

Blameless Users in a Clean Room: Defining Copyright Protection for Generative Models

Aloni Cohen
Department of Computer Science
University of Chicago
aloni@uchicago.edu

Abstract

Are there any conditions under which a generative model’s outputs are guaranteed not to infringe the copyrights of its training data? This is the question of “provable copyright protection” first posed by Vyas, Kakade, and Barak [VKB23]. They define *near access-freeness* (NAF) and propose it as sufficient for protection. This paper revisits the question and establishes new foundations for provable copyright protection—foundations that are firmer both technically and legally. First, we show that NAF alone does not prevent infringement. In fact, NAF models can enable verbatim copying, a blatant failure of copy protection that we dub being *tainted*. Then, we introduce our *blameless copy protection* framework for defining meaningful guarantees, and instantiate it with *clean-room copy protection*. Clean-room copy protection allows a user to control their risk of copying by behaving in a way that is unlikely to copy in a counterfactual “clean-room setting.” Finally, we formalize a common intuition about differential privacy and copyright by proving that DP implies clean-room copy protection when the dataset is *golden*, a copyright deduplication requirement.

1 Introduction

A user of a generative model is worried about unwitting copyright infringement. She is worried that the model’s outputs might resemble copyrighted works in the training data through no fault of her own, exposing her to legal liability. The user would like some assurance that this won’t happen. We ask: **What assurances can the model provider give, and under what conditions?**

Ideally, the generative model would never reproduce copyrighted work. That’s unrealistic. Typical models can be prompted to generate copyrighted output: `print the following ____` [VKB23, LCG24]. It’s also unnecessary. In the example, copying is clearly the user’s fault.

Instead, we should protect blameless users from inadvertent infringement. If the model reproduces copyrighted work, it should be because the user induced it to (or was very unlucky). This guarantee should satisfy a careful, honest user. As long as the user’s use of the model does not itself cause the infringement, consciously or subconsciously, then the result is unlikely to infringe.

Our goal is to formalize such a guarantee—in a mathematically and legally rigorous way—and to study its feasibility. We’re after “provable copyright protection” at deployment

time, first studied by Vyas, Kakade, and Barak [VKB23]. They propose *near access-freeness* (NAF) as an answer.

Our contributions This paper establishes new foundations for provable copyright protection—foundations that are firmer both mathematically and legally than prior work.

1. **We prove that *near access-free* (NAF) models can enable verbatim copying.** NAF is the first (and only other) attempt at a mathematical definition intended to offer “provable copyright protection” for generative models [VKB23]. The proof leverages NAF’s lack of protection against multiple prompts or data-dependent prompts.

2. **We define *tainted models*, which enable users to reproduce verbatim training data** despite knowing nothing about the underlying dataset, a blatant failure of copy protection (Section 5). A meaningful definition of provable copyright protection must exclude tainted models. NAF does not.

3. **We introduce a framework for defining copy protection guarantees, called *blameless copy protection*** (Section 6). It protects *blameless* users—those who don’t themselves induce infringement—from unwitting copying.

4. **We define *clean-room copy protection* ((κ, β) -clean), a first instantiation of our framework** (Section 7), drawing inspiration from *clean-room design*. A training algorithm is (κ, β) -clean if, for every user who copies in a (counterfactual) “clean-room environment” with probability $\leq \beta$, the probability of copying in the real world is $\leq \kappa$. Clean-room copy protection excludes tainted models under mild assumptions.

5. **We prove that *differential privacy* (DP) implies clean-room copy protection**, formalizing a common intuition about DP and copyright (Section 8). This holds when the dataset is *golden*, a copyright deduplication requirement. Thus, DP provides a way to bring copyrighted expression “into the clean room” without tainting the model.

A reader’s guide This paper can be read linearly from beginning to end. Depending on your interests, you might consider reading out of order at first, referencing earlier sections as needed. Sections 1.1, 2, and 4 are required for anything that follows them. For near access-freeness and its limitations: read Section 3, Appendix B, and Section 5. For our approach to defining meaningful copyright protection: read Sections 5–7. For the protection offered by differential privacy: start with Section 8 for the high level claims, then read Section 7 to understand the precise formalization. For discussion of practical considerations: read Section 9.

Related work There is a lot of recent work on copyright questions for generative AI [CLG⁺23, LCG24, HLJ⁺23], and the first generation of cases are working their way through the courts (e.g., *New York Times v Microsoft*). But most are only tangentially related to goal of stating and proving formal guarantees against copyright infringement.

Vyas, Kakade, and Barak were the first to propose a mathematical property—near access-freeness (NAF)—aimed at “preventing deployment-time copyright infringement” [VKB23].

Elkin-Koren, Hacohen, Livni, and Moran argue that copyright cannot be “reduced to privacy” [EKHLM24]. Referring to both NAF and DP by umbrella term “algorithmic stability”, they argue that the sort of provable guarantees that [VKB23] and this paper seek do not capture copyright’s complexities. Li, Shen, and Kawaguchi propose and empirically evaluate an attack called VA3 against NAF [LSK24]. They provide good evidence that the CP-k algorithm of [VKB23] may not prevent infringement, but some uncertainty remains (see Appendix B.3). We discuss these three papers at length. There is also work on new NAF algorithms [GAZ⁺24, CKOX24].

Scheffler, Tromer, and Varia [STV22] give a complexity-theoretic account of substantial similarity, and a procedure for adjudicating disputes. In contrast, we treat substantial similarity as a black box. That differential privacy might protect against infringement has been suggested in [BLM20, HLJ⁺23, VKB23, EKHLM24, CKOX24].

1.1 Notation

\mathcal{W} is the domain of *works* protected by copyright law. For example, \mathcal{W} might be the space of possible images, texts, or songs. We assume \mathcal{W} is discrete. We usually use $w \in \mathcal{W}$ to denote an individual work, along with $c, y, z \in \mathcal{W}$ as described in the following. We denote by $\mathcal{C} \subseteq \mathcal{W}$ the set of all in-copyright works—those currently under copyright protection—and denote a particular copyrighted work by $c \in \mathcal{C}$. A dataset $D = (w_1, \dots, w_n) \in \mathcal{W}^*$ is a list of works with multiplicities allowed. We sometimes abuse notation and treat D as a set. A *conditional generative model (model)* p is a mapping from a prompt x to a probability distribution $p(\cdot|x)$ over \mathcal{W} which samples $y \in \mathcal{W}$ with probability $p(y|x)$. A training algorithm `Train` maps a dataset D to a model p . A user u is an algorithm which uses black-box access to a model p , along with *auxiliary input* $\text{aux} \in \{0, 1\}^* \cup \{\perp\}$, and outputs a work $z \in \mathcal{W}$ (see Section 4).

2 The legal-technical interface

To use legal concepts within a mathematical formalism, we assume some exogenously-defined functions reflecting and operationalizing those concepts. They are the interface between math and law, enabling mathematical study of legal concepts without formalizing the concepts themselves. Precise definitions of these concepts are out of scope (perhaps impossible, though see [STV22] for an attempt at formalizing substantial similarity). Three legal concepts are needed for this paper: *substantial similarity*, *ideas*, and *access*. This section defines functions `SubSim` and `ideas` for the first two. Section 7.1 operationalizes access.

The owner of a copyright in a work has exclusive rights to reproduce, distribute, display, adapt, and perform the work. Original works of authorship are eligible for copyright protections as soon as they are fixed in a tangible medium, and eventually expires. While the work is protected, it is *in-copyright*. Many types of works are eligible, including books, music, poetry, plays, choreography, photographs, video, and paintings. Copying without permission is sometimes permitted, including under the fair use exception [HLJ⁺23].

To prove copying, a plaintiff (copyright holder) must prove two things: *substantial similarity* and *access*.¹ The former requires the defendant work’s to be substantially similar to

¹We do not discuss any sort of secondary liability, including vicarious or contributory infringement.

the original, with vague tests varying by jurisdiction.² Access requires the defendant to have had a reasonable opportunity to view the original (e.g., it’s online). Independent creation of substantially similar works is allowed. The function `SubSim` defines what it means for one work to be substantially similar to another.

Definition 2.1 (`SubSim`). *For a work $w \in \mathcal{W}$, the set $\text{SubSim}(w) \subseteq \mathcal{W}$ contains all works w' that are substantially similar to w . We assume w substantially similar to itself: $w \in \text{SubSim}(w)$. For a set of works $W \subseteq \mathcal{W}$, we define $\text{SubSim}(W) = \cup_{w \in W} \text{SubSim}(w)$.*

Copyright does not protect *ideas*, only *original expression*. For example, the text and images in a cookbook may be copyrighted, but not the method of cooking a dish. To be afforded protection, there must be a minimum of creativity beyond the ideas. The function `ideas` defines the non-copyrightable *ideas* of a work, as opposed their expression.

Definition 2.2 (`ideas`). *For a work $w \in \mathcal{W}$, the string $\text{ideas}(w) \in \{0, 1\}^* \cup \{\perp\}$ is (some representation of) the ideas contained in w . For a set of works $W \subseteq \mathcal{W}$, we define the set of all their ideas $\text{ideas}(W) = \{\text{ideas}(w) : w \in W\}$.*

We never actually compute `SubSim` or `ideas`. One might view them as capturing the “true” legal concept. Or one may take a more practical view. For example, interpreting `SubSim(w)` as reflecting the opinion of a typical judge or jury. However the reader wishes to interpret these functions is fine (up to any assumptions stated where used). **In particular, our results hold even for $\text{SubSim}(w) = \{w\}$ and $\text{ideas}(w) = \perp$.** In contrast, a model provider (or data provider) would need to actually work with the copyright dependency graph in Section 7. We discuss this limitation in Section 9.

3 Near access-free models do not provide provable copyright protection

Near access-freeness (NAF) is a mathematical definition for “provable copyright protection” [VKB23]. This section describes near access-freeness and its limitations. We prove that models can enable verbatim copying while still satisfying NAF. While NAF provides protection against a single prompt that is independent of the training data, it makes no guarantees against many prompts (composition) [LSK24], nor a single prompt derived from non-copyrightable ideas (not expression). We emphasize that our focus is *definitions* for provable copyright protection, where NAF falls short. Practically, NAF may still provide useful protection in many natural settings. For space, we defer algorithms, proofs, and additional discussion to Appendix B.

3.1 NAF Background

NAF informally *Near access-freeness* (NAF) is the first (and only other) attempt at a mathematical definition intended to offer “provable copyright protection” for generative models [VKB23]. Motivated by the legal requirement of access, NAF is meant to show “that

²For example, in the Ninth Circuit it requires determining “whether the ordinary, reasonable audience would find the works substantially similar in the total concept and feel of the works.”

the defendant’s work is close to a work which was produced without access to the plaintiff’s work” [VKB23]. NAF is closely related to DP. But surprisingly, any training algorithm `Train` can be used in a black-box way to construct a NAF model without the extra noise that is the hallmark of DP (Theorem 3.3).

The definition of NAF requires a generative model p trained using a copyrighted work c to be close to a “safe” model `safe` trained without any access to c . How close is governed by a parameter $k \geq 0$: smaller k means closer models. A corollary of the NAF requirement is that for any prompt x , the probability p generates something substantially similar to c is at most 2^k -times greater than `safe` doing so. We heuristically expect the latter probability to be miniscule, if `safe` and x are independent of c . If so, this bounds the probability that p ’s output infringes on prompt x .

Known limitations of NAF Prior works discuss and critique the NAF framework in three ways: questioning NAF’s technical effectiveness [LSK24], its legal applicability [LCG24, EKHLM24], and its real-world practicality [HLJ⁺23, EKHLM24, LCG24, CKOX24].

* *Technically*, Li, Shen, and Kawaguchi propose and empirically evaluate an attack called VA3 against NAF [LSK24]. See Appendix B.3 for a full discussion. They show that a black-box attacker, can reliably induce a model to produce outputs that infringe on a target work c^* . They also give a white-box prompt writing algorithm for diffusion models that greatly improves performance. VA3 provides good evidence that the CP-k algorithm of [VKB23] may not prevent infringement. Some uncertainty remains because the algorithm implemented in VA3 deviates from the original. The paper leaves open whether k -NAF (with fixed k) prevents copyright infringement, formalizes a flawed attack model, and avoids the underlying definitional questions almost entirely.

* *Legally*, Elkin-Koren et al. [EKHLM24] give the most extensive response to NAF. They correctly argue that copyright cannot be “reduced to privacy.” However, this is a straw man of version of [VKB23] and the present paper (see Appendix A). The sharpest criticism is by Lee, Cooper, and Grimmelman who argue that NAF is simply wrong on the law [LCG24]. “It is not a defense to copyright infringement that you would have copied the work from somewhere else if you hadn’t copied it from the plaintiff.” See Appendix B.4 for more discussion of NAF’s legal relevance.

* *Practically*, the main concern is that deduplicating data in the way needed to make NAF effective is infeasible [EKHLM24, LCG24, HLJ⁺23]. It is unrealistic to produce enough clean training data to pretrain foundation models, advances in deduplication notwithstanding. Moreover, NAF learning algorithms are computationally expensive and may lag in performance [CKOX24]. These concerns are justified and also apply to DP. Still, there may be settings where golden data is practical (Sec. 9), to say nothing of the value of theory.

3.2 NAF definition and the Copy Protection (CP) algorithm

Near access-freeness (NAF) is defined with respect to a function `safe`. The `safe` function maps a copyrighted data point $c \in \mathcal{C}$ to a generative model $\text{safe}_c \in \mathcal{P}$ trained without access to c . An example `sharded-safe` is in Appendix B (Alg. 2). Thm. 3.3 presents the main feasibility result: any training algorithm `Train` can be used as a black-box to construct an NAF model CP, short for *copy protection*.

Definition 3.1 (Max KL divergence). For distributions p, q , $\Delta_{\max}(p||q) := \max_{y \in \text{Supp}(p)} \log \frac{p(y)}{q(y)}$.

Definition 3.2 (k_x -NAF [VKB23]). Fix a set \mathcal{C} and function $\text{safe} : \mathcal{C} \rightarrow \mathcal{P}$. A generative model p is k_x -near access-free (k_x -NAF) on prompt $x \in \mathcal{X}$ with respect to \mathcal{C} and safe if for every $c \in \mathcal{C}$, $\Delta_{\max}(p(\cdot|x) || \text{safe}_c(\cdot, x)) \leq k_x$. A model p is k -NAF with respect to \mathcal{C} and safe if for all $x \in \mathcal{X}$, it is k_x -NAF for some $k_x \leq k$.

NAF bounds the probability that p 's output copies c relative to the probability under safe . If p is k_x -NAF safe on prompt x with respect to \mathcal{C} and safe , then for any $c \in \mathcal{C}$:

$$p(\text{SubSim}(c) | x) \leq 2^{k_x} \cdot \text{safe}_c(\text{SubSim}(c) | x). \quad (1)$$

Hence, bounding k_x prevents copying c , assuming $\text{safe}_c(\text{SubSim}(c)|x)$ is negligibly small (in $|c|$).

Theorem 3.3 (Copy Protection algorithm [VKB23]). Let p be the model returned by CP (Algorithm 3 in Appendix B), and q_1 and q_2 be the models returned by sharded-safe . Let $k_x \leq -\log(1 - d_{\text{TV}}(q_1(\cdot|x), q_2(\cdot|x)))$. Then p is k_x -NAF for x with respect to \mathcal{C} and sharded-safe .

3.3 Failures of NAF models

We now describe two ways NAF fails to prevent copying, leveraging its lack of protection against multiple prompts (i.e., composition) and data-dependent prompts. In each case, a user reproduces training data $c \in D$ verbatim. In Section 3.3.1, the model generates c when prompted with $\text{ideas}(c)$ —a prompt that depends on c but not on any copyrightable expression thereof. In Section 3.3.2, the user issues a fixed sequence of prompts, recovering k bits of the dataset with each query.

3.3.1 CP can regurgitate training data

We show that CP—the main NAF algorithm of [VKB23]—can fail to protect against copying. Using a prompt containing no copyrightable expression, a user can cause the NAF model returned by CP to regurgitate copyrighted training data.³ This is based on an observation of Thomas Steinke [Ste].

Observe that CP is a black-box transformation from an underlying training algorithm Train . Thus, CP is really a family of algorithms—one per Train .

The following theorem states that there exists Train^* such that the resulting NAF algorithm CP^* fails in the following way. On prompt $x = \text{ideas}(c)$ reproduces training datum c verbatim.

Theorem 3.4. For model training algorithm Train^* , let CP^* be the k_x -NAF algorithm from Theorem 3.3, and let $p^* \leftarrow \text{CP}^*$. For all ideas , there exists a training algorithm Train^*

³This doesn't violate Theorem 3.3. We use the fact that the CP algorithm is not actually k -NAF for any fixed k . It is k_x -NAF, where k_x depends on the prompt x . Applied to our construction, the bound on k_x given by Theorem 3.3 is not meaningful for $x \in \text{ideas}(D)$. Adapting our construction to the alternate NAF algorithm CP- k in [VKB23] results in a model sometimes fails to terminate.

such that for all datasets D and for all works $c \in D$ whose ideas are distinct in D (i.e., $\forall w \in D, w \neq c : \text{ideas}(w) \neq \text{ideas}(c)$):

$$p^*(c \mid \text{ideas}(c)) = 1.$$

We defer the proof of (a slight generalization) of this theorem to Appendix B.2. The proof uses the tainted training algorithm $\text{Train}^* = \text{Train}_{kv}$ (Algorithm 1) described in Example 5.2.

3.3.2 k -NAF does not prevent full reconstruction

Our next theorem shows that k -NAF may allow training data to be reconstructed, even if k is arbitrarily small and independent of x . This is because the k -NAF guarantee does not compose across a user’s many queries, and each query may leak up to k bits of training data.

In more detail, we give a family of models that are k -NAF with respect to pure noise, yet which enable a user to reconstruct the dataset verbatim. For any $\ell \geq 1$, let $\text{coin}_\ell(\cdot|x)$ be uniform over $\{0, 1\}^\ell$ for all prompts x . As a generative model, $\text{coin}_\ell(\cdot|\cdot)$ is clearly a “safe” instantiation of safe_c for any copyrighted work c .

Theorem 3.5. *Fix $\mathcal{C} \subseteq \mathcal{W} \subseteq \{0, 1\}^*$. For $D \in \mathcal{W}^*$, let L be total the length of D in bits. There exists a (deterministic) training algorithm $\text{Train} : (D, k) \mapsto p_{k,D}$ satisfying the following.*

- For all D and $k > 0$: $p_{k,D}$ is k -NAF with respect to \mathcal{C} and coin_ℓ for $\ell = \max\{1, \lfloor k \rfloor\}$.
- There exists a user u such that for all D and $k > 0$: u makes $\text{poly}(L, 1/k)$ black-box queries to $p_{k,D}$ and outputs D with probability > 0.99 .

The proof is deferred to Appendix B.2.

4 The interaction between user and model

People interact with generative models in complex ways: refining prompts and using results to create a final product whose author is not solely human nor machine. Meaningful copyright protection must cover such cases. To that end, we make the interaction between a user and the model explicit.

Let u be a (randomized) algorithm called the *user*, p be a model, and $\text{aux} \in \{0, 1\}^* \cup \{\perp\}$ be an auxiliary input. We denote by $u^p(\text{aux})$ the algorithm u run with input aux and black-box access to p . The result is a work $z \in \mathcal{W}$ distributed as $z \leftarrow u^p(\text{aux})$. Together with a training algorithm Train and a dataset D , this process induces probability measure τ over \mathcal{W} .

Definition 4.1 (User’s output distribution). *For a user u , training algorithm Train , dataset D , and auxiliary information aux , we define the user’s output distribution τ over \mathcal{W} as:*

$$\tau(w; \text{aux}) = \Pr_{\substack{p \leftarrow \text{Train}(D) \\ z \leftarrow u^p(\text{aux})}} [z = w].$$

The probability is taken over the randomness of the algorithms Train , u , and p .

Legally, the user’s output infringes on the copyright of a work $c \in \mathcal{C}$ only if it is substantially similar to c . Preventing substantial similarity prevents infringement. This motivates our first desideratum.

Desideratum 1. *We want $\tau(\text{SubSim}(\mathcal{C}); \text{aux})$ to be as small as possible for as many users as possible.*

By considering the user’s output distribution, we require protection against multiple prompts (composition) and data-dependent prompts (when `aux` contains the ideas of work in the training data, as below). This fixes the shortcomings of NAF that enabled the attacks described in Section 3.3.

5 Tainted models enable copying

This section defines what it means for a generative model to be *tainted*. Tainted models let users copy training data of which they have no knowledge. Tainted models blatantly fail to prevent copying. If the goal is to prevent copying, being tainted is as bad as it gets. A tainted model—one produced by a tainted training algorithm (Definition 5.1)—doesn’t prevent copyright infringement. It enables it! Thus, we get our next desideratum. (We will see that NAF fails this test but Definition 7.6 passes.)

Desideratum 2. *A meaningful definition of provable copyright protection should bar tainted models.*

We say a training algorithm `Train` is tainted if there is a fixed algorithm `u` that copies work from the training dataset using only trained model and the unprotected ideas of the work in question.

Definition 5.1 (Tainted training). *We say `Train` is tainted with respect to ideas if there exists a user `u` such that for all datasets D and all works $w \in D$:*

$$\tau(\text{SubSim}(D); \text{ideas}(w)) > 0.99. \tag{2}$$

The order of quantifiers is important. The user `u` is fixed once and for all, and D varies arbitrarily. The user cannot possibly be at fault: it cannot “know” a work in all possible datasets D .

5.1 Examples

Example 5.2 below describes a black-box transformation from any training algorithm to a tainted-version of that training algorithm. This example is used to prove Theorem 3.5.

Example 5.2. *Algorithm 1 describes a tainted algorithm Train_{kv} , constructed from any training algorithm Train_0 in a black-box way. In words, $\text{Train}_{\text{kv}}(D)$ trains a model $q_0 \leftarrow \text{Train}_0(D)$, and also builds a key-value store I . I maps ideas `id` to the set $D_{\text{id}} = \{w \in D : \text{ideas}(w) = \text{id}\}$. On prompt x , the model q_{kv} returns a random element of $I[x] = D_x$ if it is non-empty. Otherwise, it returns a generation sampled from $q_0(\cdot|x)$. It is easy to see that Train_{kv} is tainted. Consider the user $u^p(\text{aux})$ that returns a sample from $p(\cdot|\text{aux})$. Fix D and $w \in D$, and let $p \leftarrow \text{Train}_{\text{kv}}(D)$. By construction, the user always outputs an element of $D_{\text{ideas}(w)}$. Therefore, $\tau(\text{SubSim}(D); \text{ideas}(w)) \geq \tau(D_{\text{ideas}(w)}; \text{ideas}(w)) = 1$.*

Algorithm 1: Train_{kv} (for Example 5.2 and Theorem B.5)

Parameters: Training algorithm Train_0 **Input:** Data D **Output:** Model q_{kv} Let $q_0 \leftarrow \text{Train}_0(D)$;Initialize an empty key-value store I , whose keys are ideas id and values are sets of works $W \subset \mathcal{W}$;**for** $w \in D$ **do** // add w to $I[\text{ideas}(w)]$; $W \leftarrow I[\text{ideas}(w)]$; $W \leftarrow W \cup \{w\}$; $I[\text{ideas}(w)] \leftarrow W$;**end**Let q_{kv} the conditional generative model which on prompt x does the following; $W \leftarrow I[x]$; **if** $W \neq \emptyset$ **then** | Return $y \leftarrow_s W$. **else** | Return $y \sim q_0(\cdot|x)$. **end****Result:** q_{kv}

As a negative example, the following claim states that any training algorithm that always returns a single fixed model is not tainted, under a mild assumption.

Claim 5.3. *Suppose there exists a pair of works $w_0, w_1 \in \mathcal{W}$ such that $\text{ideas}(w_0) = \text{ideas}(w_1)$ and $\text{SubSim}(w_0) \cap \text{SubSim}(w_1) = \emptyset$. Then for any fixed model q , the constant algorithm $\text{Train}_q(D) = q$ is not tainted.*

Proof. Let $D_i = \{w_i\}$ and fix a user u . Let τ_i be the user’s output distribution (Definition 4.1) with $D = D_i$ and $\text{aux} = \text{ideas}(w_i)$. By construction, $\tau_1 = \tau^* = \tau_2$. Note that $\tau^*(w_i) \leq 1/2$ for one of $i = 0, 1$. Equation (2) is violated for the corresponding D_i . \square

5.2 Tainted NAF models

Desideratum 2 lays out a necessary condition for provable copyright protection: exclude tainted training algorithms. NAF fails this test. This is captured by the following corollaries of Theorem B.5 (slightly generalizing Theorem 3.4) and Theorem 3.5.

Corollary 5.4. *For any ideas, algorithm CP^* given by Thm. 3.4 is tainted.*

Proof. Consider the user $u^p(\text{aux})$ that returns a sample from $p(\cdot|\text{aux})$. Fix D and $w \in D$, and let $p^* \leftarrow \text{CP}^*(D)$. By Theorem B.5, $\tau(\text{SubSim}(D); \text{ideas}(w)) \geq \tau(D_{\text{ideas}(w)}; \text{ideas}(w)) = 1$. \square

Corollary 5.5. *For any k , ideas, there exists a tainted algorithm Train such that for all datasets D , the model $p \leftarrow \text{Train}(D)$ is k -NAF with respect to coin_ℓ .*

Proof. Immediate from the statement of Thm. 3.5 \square

6 Blameless copy protection: a definitional framework

Our goal is to define provable copyright protection. The previous section gives a negative answer: provable copyright protection should bar tainted models. This section and those that follow give a positive answer: provable copyright protection should protect blameless users.

We can’t hope to guarantee that a generative model never reproduces copyrighted work, even works it was not trained on. A malicious user can always induce copying. (Formally, if the set of in-copyright works \mathcal{C} is non-empty, then there exists u that always infringes: $\tau(\text{SubSim}(\mathcal{C}); \perp) = 1$.) Instead, we wish protect *blameless* users—users who don’t themselves induce infringement—from unwitting copying.

This suggests a framework for defining meaningful guarantees, which we call *blameless copy protection*. First, define a class of blameless users \mathfrak{B} . Second, guarantee that for blameless users, the probability of copying $\tau(\text{SubSim}(\mathcal{C}); \text{aux})$ at most some small κ .

Definition 6.1 (Blameless copy protection). *Fix $\kappa > 0$ and a class \mathfrak{B} of blameless users. We say Train is (κ, \mathfrak{B}) -copy protective if for all $\mathcal{C} \subseteq \mathcal{W}$, $D \in \mathcal{W}^*$, aux , and $u \in \mathfrak{B}$: $\tau(\text{SubSim}(\mathcal{C}); \text{aux}) \leq \kappa$.*

The missing piece is \mathfrak{B} : What makes a user blameless? Different answers will yield different versions of blameless copy protection. So Definition 6.1 is less a definition than a framework for definitions. Instantiating it may require additional assumptions. As we cannot perfectly capture legal blamelessness, we should err on the side of protecting more users rather than less.

In Section 7, we offer a first instantiation of blameless copy protection, inspired by clean-room design. But there may be very different ways to formalize blamelessness, which we leave for future work. The cryptographic notion of extraction offers an intriguing approach: a user is blameworthy if a copyright-infringing work can be efficiently extracted from the user itself.

7 Defining clean-room copy protection

This section formalizes what it means to have *access* to a copyrighted work (Section 7.1) and instantiates the blameless copy protection framework (Section 7.2)—drawing inspiration from clean-room design. In copyright law, a clean room is “a process of producing a product under conditions guaranteeing independent design and foreclosing the possibility of copying.”⁴ The idea comes from cases involving reverse engineering. Roughly, a team is given a description of design specs of the product (ideas) but not the product itself (expression), and tasked with producing a compatible product. There is no access, as long as the team itself was not spoiled by prior familiarity with the product or its design. As such, the outcome is constructively non-infringing. Even substantial similarities will be the product of independent creation.

7.1 Scrubbing a dataset removes access

A clean room is a setting where a user does not have access to a particular copyrighted work. To pin this down, we first operationalize the legal concept of access.

The copyright dependency graph We need a way to say whether a work w' depends on or derives from another work w a copyright-relevant way. Such dependencies are directed, with later works stemming from earlier works. The *copyright dependency graph* encodes this relation. It is assumed to capture the legal concept of access in the following sense: *For purposes of copyright law, a dataset D only gives access to a work w if there is some $w' \in D$ that stems from w .*

Definition 7.1 (Copyright dependency graph). *The copyright dependency graph is a directed graph $G = (\mathcal{W}, E)$ whose edges reflect dependencies relevant for copyright. If $(w, w') \in E$, we say that w' stems from w . We assume that every w stems from itself: $(w, w) \in E$ for all w .*

It is useful to refer to the set of all copyrighted works for which access is or isn’t implied by access to a dataset D . We denote these sets by \mathcal{C}_D and \mathcal{C}_{-D} , respectively. For the sake of copyright analysis, access to D does not constitute access to any copyrighted $c \in \mathcal{C}_{-D}$.

⁴David S. Elkins, *NEC v. Intel: A Guide to Using “Clean Room” Procedures as Evidence*, 10 Computer L.J. 453 (1990). <https://repository.law.uic.edu/jitpl/vol10/iss4/1>

Definition 7.2 (The sets \mathcal{C}_D and \mathcal{C}_{-D}). We define $\mathcal{C}_D = \{c \in \mathcal{C} : \exists w \in D, (c, w) \in E\}$ as the set copyrighted works from which any work in D stems. We denote its complement $\mathcal{C}_{-D} = \mathcal{C} \setminus \mathcal{C}_D$.

The scrub function We define the function `scrub` to operationalize the legal concept of access. It removes from a dataset any work that stems from a target copyrighted work c . That is, any w for which $c \rightarrow w$ is an edge in the copyright dependency graph. The result is a new dataset $\text{scrub}(D, c) \subseteq D$. For the sake of copyright analysis, access to $\text{scrub}(D, c)$ does not constitute to access to c .

Definition 7.3 (`scrub`). Fix copyright dependency graph $G = (\mathcal{W}, E)$, dataset D and work c . The dataset $\text{scrub}(D, c) = \{w \in D : (c, w) \notin E\}$ is the sub-dataset of D of works not stemming from c .

Observe that $\mathcal{C}_D = \{c \in \mathcal{C} : \text{scrub}(D, c) \neq D\}$ and $\mathcal{C}_{-D} = \{c \in \mathcal{C} : \text{scrub}(D, c) = D\}$

7.2 Clean training algorithms: copy protection for blameless users in a clean room

This section gives a clean-room inspired formulation of copy protection. Informally, we say that a training algorithm is (κ, β) -clean if for all users u : either (i) u would have copied in a clean-room setting with probability at least β ; or (ii) u copies in the real world with probability at most κ . Making this precise, we define the user’s *clean-room output distribution* (Definition 7.4), *blameless users in the clean room* (Definition 7.5), and (κ, β) -*clean-room copy protection* (Definition 7.6).

The user’s clean-room distribution We define a user’s *clean-room distribution*, a counterfactual to its true distribution τ . For a copyrighted work c , we denote by τ_{-c} the user’s output distribution in a clean room where the model doesn’t depend on c . The only difference from the user’s real-world output distribution τ (Definition 4.1) is that the model is trained on $\text{scrub}(D, c)$ instead of D .

Definition 7.4 (User’s clean-room distribution). For user u , training algorithm `Train`, dataset D , work c , and auxiliary information `aux`, we define the user’s clean-room distribution τ_{-c} for c as:

$$\tau_{-c}(w; \text{aux}) = \Pr_{\substack{p \leftarrow \text{Train}(\text{scrub}(D, c)) \\ z \leftarrow u^p(\text{aux})}} [z = w].$$

Clean-room blamelessness and copy protection We postulate that *blameless users can control their risk of copying when using a model trained in a clean room*. Given $\beta > 0$, a blameless user can guarantee that they copy a work c that was “outside the clean room” with probability at most β .⁵ This should hold even if the user is exposed to c in some

⁵The exact constant doesn’t matter. It should be small enough to be very unlikely to occur by chance ($\ll 2^{-10}$). It should be large enough to reflect that a relatively small amount of original expression is needed for copyright protection ($\gg 2^{-1000}$).

other than through the model.⁶ People are constantly exposed to copyrighted works, yet we somehow manage to avoid copying every day. Using a model trained without access to c shouldn't change that.

We call these users β -blameless in a clean room, or β -blameless for short. Definition 7.5 formalizes this idea. Recall that \mathcal{C}_{-D} is the set of copyrighted works from which no element of D stems (Definition 7.2). Thus, the set of works “outside the clean room” is $\mathcal{C}_{-D} \cup \{c\}$.

Definition 7.5 (Blameless in the clean room (β -blameless)). *For $0 \leq \beta \leq 1$, a user u is β -blameless in the clean room (β -blameless) with respect to $D, \mathcal{C}, \text{Train}$ if for all $c \in \mathcal{C}$ and all aux :*

$$\tau_{-c}(\text{SubSim}(\mathcal{C}_{-D} \cup \{c\}); \text{aux}) \leq \beta.$$

*Otherwise, u is β -blameworthy.*⁷

Instantiating our framework (Definition 6.1), we now define *clean training algorithms* as those that provide copy protection to users who are blameless in the clean room.

Definition 7.6 (Clean-room copy protection; (κ, β) -clean). *For $\kappa, \beta > 0$, we say Train is (κ, β) -clean if for all $\mathcal{C} \subseteq \mathcal{W}$, $D \in \mathcal{W}^*$, all users u who are β -blameless (w.r.t. $\mathcal{C}, D, \text{Train}$), and all aux :*

$$\tau(\text{SubSim}(\mathcal{C}); \text{aux}) \leq \kappa.$$

If this only holds for datasets D in an admissible set $\mathfrak{D} \subseteq \mathcal{W}^$, we say Train is (κ, β) -clean for \mathfrak{D} .*

Clean-room copy protection is not tainted By Desideratum 2, a good definition of provable copy protection should exclude tainted training algorithms. The following theorem shows that clean-room copy protection passes this test, under a reasonable assumption on ideas and SubSim. Roughly, the assumption is that there exist at least n distinct ideas each of which can be expressed in $m \gg n$ ways satisfying a strong dissimilarity property: that for any distinct w and w' , the set of works substantially similar to both w and w' is empty.

Here's a very crude justification for $\beta \approx 2^{-50}$. Estimates of the entropy of English text vary. Shannon estimated it at 0.6 to 1.3 bits per character (A–Z and space), based on a study of how well his English-speaking test subjects were able to predict the next character in a passage [Sha51]. With those estimates, just 77 to 167 characters of text contain ≈ 100 bits of entropy. Here is Shakespeare's *Sonnet 130* with the 77th and 167th characters highlighted.

My mistress' eyes are nothing like the sun;
 Coral is far more red than her lips' red;
 If snow be white, why then her breasts are dun;
 If hairs be wires, black wires grow on her head.

Acknowledging that I know nothing of the matter, somewhere in that range seems reasonable for the minimum amount of original expression needed to afford copyright protection. Infringement only requires substantial similarity, not perfect reproduction. If substantial similarity requires half of the 100 bits of minimum original expression to be reproduced, we get to 2^{-50} .

⁶One can weaken this assumption (making more users blameless) by restricting aux in Definitions 7.5 and 7.6. For example, by allowing aux to include $\text{ideas}(c)$ but not original expression of works in $\mathcal{C}_{-D} \cup \{c\}$. This would undermine Desideratum 1, yielding qualitatively weaker protection. DP would still suffice (Theorem 8.2) with this change.

⁷Equivalently, $\mathfrak{B}_\beta^{\text{clean}}(D, \mathcal{C}, \text{ideas}) := \{u : \forall c \in \mathcal{C}, \forall \text{aux}, \tau_{-c}(\text{SubSim}(\mathcal{C}_{-D} \cup \{c\}); \text{id}) \leq \beta\}$.

Theorem 7.7. Fix $\beta < 1$, $n \geq 1$, and $m \geq \frac{n}{\beta}$. Suppose there exists works $W = (w_{i,j}) \in \mathcal{W}^{m \times n}$ such that for all $i \neq i' \in [m]$ and all $j \in [n]$:

$$\text{ideas}(w_{i,j}) = \text{ideas}(w_{i',j}) \quad \text{and} \quad \text{SubSim}(w_{i,j}) \cap \text{SubSim}(w_{i',j}) = \emptyset.$$

If Train is tainted with respect to ideas, then Train is not (κ, β) -clean.

The proof is deferred to Appendix D.

8 Differential privacy’s protection against copying

Many works have suggested that differential privacy (DP) should offer some sort of protection against copying if the training data are deduplicated [BLM20, HLJ⁺23, VKB23, EKHLM24, CKOX24]. This section turns the suggestion into a theorem. For background on DP [DMNS06], see Appendix C.

At a high level, we show that DP training blames users from copyright infringement (in the sense of Definition 7.6) if the training dataset is *golden*, a copyright-deduplication condition defined below. If so, the probability of copying is at most $\approx e^\varepsilon \beta N_D$, where β is a user’s probability of copying in the clean room and $N_D = |\mathcal{C}_D|$ is the number of copyrighted works to which D grants access. To guarantee the risk is at most κ , the user sets $\beta \approx \kappa / e^\varepsilon N_D$. (This is pessimistic: see Remark 8.4.)

A dataset D is *golden* if at most one item $w \in D$ stems from any copyrighted work c .⁸ That item could be c itself or a derivative work. No protected original expression appears in more than one element of a golden dataset. For example, it can contain one copy or parody of the *Abbey Road* album cover (not both), and many parodies of the out-of-copyright *Mona Lisa*. Remark 8.5 explains why golden datasets will typically remain golden as the set of in-copyright works evolves over time.

Definition 8.1 (Golden dataset). Let $G = (\mathcal{W}, E)$ be a copyright dependency graph, and $\mathcal{C} \subseteq \mathcal{W}$ be the set of in-copyright works. A dataset D is golden (with respect to G, \mathcal{C}) if for all $c \in \mathcal{C}$ there is at most one $w \in D$ such that $(c, w) \in E$.

Theorem 8.2. Let Train be (ε, δ) -differentially private for $\varepsilon > 0, \delta \geq 0$. Let $N_D = |\mathcal{C}_D|$. Then Train is (κ, β) -clean for golden datasets D , for all $\beta \geq 0$ and $\kappa \geq (e^\varepsilon N_D + 1)\beta + N_D \delta$.

Proof. Fix $\mathcal{C} \subseteq \mathcal{W}$, D and user u that is β -blameless (Def. 7.5). We must show that for all aux , $\tau(\text{SubSim}(\mathcal{C}); \text{aux}) \leq \kappa$. To simplify notation, we omit aux for the rest of the proof. Because D is golden, $|D \setminus \text{scrub}(D, c)| \leq 1$. Hence, datasets $\text{scrub}(D, c)$ and D are either equal or

⁸We adapt and formalize this idea from [VKB23], where it appears as an informal condition that suffices for their analysis: “It is important that when we omit a data point x it does not share copyrighted content with many other data points that were included in the training set.” We also borrow the term “golden dataset”, though [VKB23] uses it to refer to something entirely different: A dataset that is “carefully scrutinized to ensure that all material in it is not copyrighted or [is] properly licensed.” That condition doesn’t suffice, as even licensed derivatives could cause unlicensed copying by a generative model. The same is true for NAF, despite their suggestion that it would yield a safe model.

neighboring for all $c \in \mathcal{C}$. Applying DP and post-processing, we have $\tau(E) \leq e^\varepsilon \cdot \tau_{-c}(E) + \delta$ for all events E and all $c \in \mathcal{C}$.

$$\begin{aligned} \tau(\text{SubSim}(\mathcal{C})) &\leq \tau(\text{SubSim}(\mathcal{C}_{-D})) + \sum_{c \in \mathcal{C}_D} \tau(\text{SubSim}(c)) && \text{(union bound)} \\ &\leq \beta + N_D \delta + e^\varepsilon \cdot \sum_{c \in \mathcal{C}_D} \tau_{-c}(\text{SubSim}(c)) && \text{(DP; Prop. 8.3 below)} \\ &\leq (e^\varepsilon N_D + 1)\beta + N_D \delta \leq \kappa && (\beta\text{-blameless}) \quad \square \end{aligned}$$

Proposition 8.3. *Let u be β -blameless w.r.t. $D, \mathcal{C}, \text{Train}$. For all $c \in \mathcal{C}_{-D}$, aux : $\tau(\mathcal{C}_{-D}; \text{aux}) \leq \beta$.*

Proof. By definition of \mathcal{C}_{-D} , $\text{scrub}(D, c) = D$ and hence $\tau_{-c} = \tau$. Also, $\mathcal{C}_{-D} = \mathcal{C}_{-D} \cup \{c\}$. Thus, $\tau(\mathcal{C}_{-D}; \text{aux}) \leq \tau_{-c}(\text{SubSim}(\mathcal{C}_{-D} \cup \{c\}); \text{aux}) \leq \beta$, with the last inequality by blamelessness. \square

Remark 8.4 (Improving the union bound). The factor of $N_D = |\mathcal{C}_D|$ from the union bound is unreasonably pessimistic. Ideas sometimes differ so greatly that the risk of substantial similarity is essentially 0: a photorealistic bird will never resemble Dr. Seuss’s Cat in the Hat. If the user is targeting a fixed set of ideas id (conditioning τ on id), one might instead take the union bound over $\mathcal{C}_D^{\text{id}} = \{c \in \mathcal{C}_D : \text{id} \in \text{ideas}(\text{SubSim}(c))\}$. We expect that typically $N_D^{\text{id}} := |\mathcal{C}_D^{\text{id}}| \ll |\mathcal{C}_D| = N_D$.

Remark 8.5 (Golden datasets remain golden over time). Theorem 8.2 states that (ε, δ) -differentially private models are (κ, β) -clean for *golden datasets* (and appropriate κ and β). The set of golden datasets \mathfrak{D} depends on the set of in-copyright works \mathcal{C} . For example, every D is golden when nothing is in-copyright ($\mathcal{C} = \emptyset$).

This presents a challenge: Will a golden dataset remain golden as the set of works protected by copyright evolves? As copyrights expire and new works are created, the sets of in-copyright works today \mathcal{C} and tomorrow \mathcal{C}' will diverge. For Definition 7.6 to offer lasting protection, we need $\mathfrak{D} \subseteq \mathfrak{D}'$.

Fortunately, it will typically be true that datasets remain golden as the set of in-copyright works evolves. This requires three observations. First, expiring copyrights don’t affect goldenness. If D is golden for \mathcal{C} , then D is golden for $\mathcal{C} \cap \mathcal{C}'$. Second, a work is afforded copyright protections as soon as it is fixed in a tangible medium. This means that if a work currently exists will ever be in-copyright (in \mathcal{C}'), then it already in-copyright (in \mathcal{C}). Hence, $\mathcal{C}' \setminus \mathcal{C}$ contains only works that do not yet exist. The exception is when a work becomes newly eligible for copyright protections, and the change in eligibility is applied retroactively. For example, if a court extends copyright to a new class of works. Third, no work can stem from a work created later. Combined with the previous observation, D is golden with respect to $\mathcal{C}' \setminus \mathcal{C}$. Putting it all together, we have that $D \in \mathfrak{D}$ implies that D is golden with respect to $(\mathcal{C} \cap \mathcal{C}') \cup (\mathcal{C}' \setminus \mathcal{C}) = \mathcal{C}'$. That is, $D \in \mathfrak{D}'$.

9 Discussion: provable copyright protection in practice?

What if copying occurs? People copy without generative models, and will continue copying with them. One way to evaluate the usefulness of a copyright-mitigation measure

is to consider what happens when copying does occur. Who should pay? The existence of tainted NAF models muddies the makes it hard to assign culpability on the basis of NAF alone. DP gives a theoretical account of who is culpable: the provider if the data wasn't golden; the user if not blameless (and possibly both or neither). But blamelessness may be impossible to check.

Clean-room copy protection makes possible a simple indemnification rule: The model provider pays for infringement if the copied work appears too often in the training data. This is a question of fact that legal process and forensic experts are well-suited to resolve: experts for each party examine the training data and attempt to convince a fact finder. The rule gives users what they want. They get indemnity for copying when there is any possibility that the model is to blame, and can control their risk tolerance by controlling β . The rule gives model providers what they want too. Besides offering a valuable protection for users, the provider can trade-off their overall exposure to copyright penalties with their effort spent cleaning the data. Or the provider can pass the liability to third-party data providers contracted to provide golden data.

How reasonable are the requirements for clean room protection? Theorem 8.2 requires DP models, golden data, and blameless users. Are these requirements reasonable?

For pre-trained foundation models, DP and golden data are infeasible. DP generally requires more data for the same performance. And making golden dataset is much harder than deduplication (already a major challenge [LIN⁺21]): it depends on squishy legal standards and on data outside the dataset. Fine-tuning is more promising, because much less data is needed. The increased sample complexity of DP learning might be more manageable. And creating a golden dataset with tens of thousands of items seems feasible, possibly by commissioning new work.

There are non-technical ways to address the challenge of golden datasets. Suppose the trainer license the data from a data provider. They could require metadata to include copyright dependencies, or that the data provider provide golden data. Either way, the license agreement can indemnify users if the data provider's failure to appropriately clean or tag the data leads to infringement. Not all is lost if the data is imperfect. The guarantees will still hold for any work satisfying the golden condition.

Blamelessness presents an even greater challenge: it is not checkable, not even by the user. Still, we believe that diligent users who are attentive to the possibility of inadvertent infringement can guarantee β -blamelessness for β small enough to make Theorem 8.2 meaningful.

This is the crux of the matter. Why do I think that users can avoid being unduly influenced by works to which they have been exposed? I don't have a good answer. More than anything, I can't shake the belief that I could do it. That I could productively use ChatGPT-4o to produce a story or image that is wholly original, bearing no resemblance even to stories and images with which I am intimately familiar. People are able to be truly creative, despite constant exposure to copyrighted work. Somehow you and I manage to avoid copying every day.

Acknowledgements

We are indebted to Mayank Varia for many helpful discussions and encouragement. We thank Gautam Kamath, Yu-Xiang Wang, Seewong Oh, and especially Thomas Steinke for early discussions when the idea of this paper was still taking shape. Theorem 3.4 is based on an observation of Thomas. We thank Sarah Scheffler and the attendees of the 2024 Works-in-Progress Roundtable on Law and Computer Science at the University of Pennsylvania for feedback on an early draft.

References

- [BLM20] Olivier Bousquet, Roi Livni, and Shay Moran. Synthetic data generators—sequential and private. *Advances in Neural Information Processing Systems*, 33:7114–7124, 2020.
- [CKOX24] Wei-Ning Chen, Peter Kairouz, Sewoong Oh, and Zheng Xu. Randomization techniques to mitigate the risk of copyright infringement. *arXiv preprint arXiv:2408.13278*, 2024.
- [CLG⁺23] A Feder Cooper, Katherine Lee, James Grimmelmann, Daphne Ippolito, Christopher Callison-Burch, Christopher A Choquette-Choo, Niloofar Mireshghallah, Miles Brundage, David Mimno, Madiha Zahrah Choksi, et al. Report of the 1st workshop on generative ai and law. *arXiv preprint arXiv:2311.06477*, 2023.
- [DMNS06] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*, pages 265–284. Springer, 2006.
- [EKHLM24] Niva Elkin-Koren, Uri Hacoen, Roi Livni, and Shay Moran. Can copyright be reduced to privacy? Forthcoming, Foundations of Responsible Computing, 2024. <https://arxiv.org/abs/2305.14822>.
- [GAZ⁺24] Aditya Golatkar, Alessandro Achille, Luca Zancato, Yu-Xiang Wang, Ashwin Swaminathan, and Stefano Soatto. Cpr: Retrieval augmented generation for copyright protection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12374–12384, 2024.
- [HLJ⁺23] Peter Henderson, Xuechen Li, Dan Jurafsky, Tatsunori Hashimoto, Mark A. Lemley, and Percy Liang. Foundation models and fair use. *Journal of Machine Learning Research*, 24(400):1–79, 2023.
- [LCG24] Katherine Lee, A. Feder Cooper, and James Grimmelmann. Talkin’ ’bout ai generation: Copyright and the generative-ai supply chain. Forthcoming, Journal of the Copyright Society, 2024. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4523551.

- [LIN⁺21] Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better. *arXiv preprint arXiv:2107.06499*, 2021.
- [LSK24] Xiang Li, Qianli Shen, and Kenji Kawaguchi. Va3: Virtually assured amplification attack on probabilistic copyright protection for text-to-image generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12363–12373, 2024.
- [Sha51] Claude E Shannon. Prediction and entropy of printed english. *Bell system technical journal*, 30(1):50–64, 1951.
- [Ste] Thomas Steinke. Personal communication.
- [STV22] Sarah Scheffler, Eran Tromer, and Mayank Varia. Formalizing human ingenuity: A quantitative framework for copyright law’s substantial similarity. In *Proceedings of the 2022 Symposium on Computer Science and Law*, pages 37–49, 2022.
- [VKB23] Nikhil Vyas, Sham Kakade, and Boaz Barak. Provable copyright protection for generative models. *arXiv preprint arXiv:2302.10870*, 2023.

A Are we trying to reduce copyright to privacy? No!

Elkin-Koren, Hacoheh, Livni, and Moran argue that copyright cannot be “reduced to privacy.” Referring to both NAF and DP by umbrella term algorithmic stability, they argue that the sort of provable guarantees that [VKB23] and this paper seek do not capture copyright’s complexities.

Algorithmic stability approaches, when used to establish proof of copyright infringement are either too strict or too lenient from a legal perspective. Due to this misfit, applying algorithmic stability approaches as filters for generative models will likely to distort the delicate balance that copyright law aims to achieve between economic incentives and access to creative works.

Too strict by excluding permitted uses of copyrighted data: (1) works in the public domain; (2) unprotected aspects of copyrighted work (e.g., ideas, facts, procedures); and (3) lawful uses of copyrighted work, especially fair use. Too lenient when protected expression originating in one work is present in many other works in a training data set. For example, copies, derivatives, or snippets of the original (whether fair use or not) would undermine any NAF- or DP-based guarantee if unaccounted for.

We completely agree, and suspect that [VKB23] would too. The claim that “copyright can be reduced to privacy” is a straw man. It would be foolish to suggest that the whole of copyright law for generative AI be governed by a mathematical formalism like NAF or DP. That doesn’t mean that a mathematical formalism offering legal guarantees isn’t worthwhile—as a thought experiment or as a first step towards practical solutions.

Still, a formalism that is both too lenient and too strict won’t support a legal conclusion either way. As we cannot reduce copyright to a mathematical formalism, we must choose

how to err. In this work, we seek sufficient conditions under which one can guarantee no infringement. Conditions that are too strict will leave room for improvement, but won't be fatally flawed. But conditions that are too lenient would be fatal—they would not suffice. Thus, we must address the possibility that one work's copyrighted expression appears in many other works.

B Deferred material on near access-free models (Section 3)

This section presents our detailed treatment of near access-free models. It describes near access-freeness and its limitations. We prove that models can enable verbatim copying while still satisfying NAF. While NAF provides protection against a single prompt that is independent of the training data, it makes no guarantees against many prompts [LSK24], nor a single prompt derived from non-copyrightable ideas (not expression).

To make this section self-contained, we repeat material from Section 3, quoting freely.

B.1 Definitions and main results from [VKB23]

NAF is defined with respect to a function `safe`. The `safe` function maps a copyrighted data point $c \in \mathcal{C}$ to a generative model $\text{safe}_c \in \mathcal{P}$ trained without access to c . An example `safe` function is `sharded-safe`. We simplify the notation of [VKB23] by fixing the divergence to maximum KL divergence.

Definition B.1 (Max KL divergence). *For distributions p, q , $\Delta_{\max}(p||q) := \max_{y \in \text{Supp}(p)} \log \frac{p(y)}{q(y)}$.*

Definition B.2 (k_x -NAF [VKB23]). *Fix a set \mathcal{C} and function $\text{safe} : \mathcal{C} \rightarrow \mathcal{P}$. A generative model p is k_x -near access-free (k_x -NAF) on prompt $x \in \mathcal{X}$ with respect to \mathcal{C} and `safe` if for every $c \in \mathcal{C}$,*

$$\Delta_{\max}\left(p(\cdot|x) \parallel \text{safe}_c(\cdot, x)\right) \leq k_x.$$

A model p is k -NAF with respect to \mathcal{C} and `safe` if for all $x \in \mathcal{X}$, it is k_x -NAF for some $k_x \leq k$.

The appeal of this definition is that it can be used to bound the probability that model p produces outputs that violate the copyright of a work c , relative to the probability under `safe`.

Lemma B.3 (k -NAF event bound [VKB23]). *Suppose model p is k_x -NAF on prompt x with respect to \mathcal{C} and `safe`. Then for any $c \in \mathcal{C}$ and any event $E \subseteq \mathcal{Y}$:*

$$p(E|x) \leq 2^{k_x} \cdot \text{safe}_c(E|x).$$

Letting $E = \text{SubSim}(c)$ be the event that y is substantially similar to c , we get

$$p(\text{SubSim}(c)|x) \leq 2^{k_x} \cdot \text{safe}_c(\text{SubSim}(c)|x).$$

As copying requires substantial similarity, bounding k_x suffices to prevent violation whenever $\text{safe}_c(\text{SubSim}(c)|x)$ is negligibly small. Heuristically, we expect $\text{safe}_c(\text{SubSim}(c)|x)$ to be

negligible in $|c|$. It may sometimes be large, e.g., when x contains the copyrighted work c [VKB23].

Theorem 3.3 presents the main feasibility result of [VKB23]: any training algorithm `Train` can be used as a black-box to construct an NAF model `CP`, short for *copy protection*.

Theorem B.4 ([VKB23]). *Let p be the model returned by `CP`, and q_1 and q_2 be the models returned by `sharded-safe`. Then p is k_x -NAF with respect to \mathcal{C} and `sharded-safe`, with*

$$k_x \leq -\log\left(1 - d_{\text{TV}}(q_1(\cdot|x), q_2(\cdot|x))\right). \quad (3)$$

Observe that `CP`(D) is well-defined if and only if $\forall x, \exists y: q_1(y|x) > 0$ and $q_2(y|x) > 0$.

Algorithm 2: `sharded-safe` [VKB23]

Parameters: Dataset D , training algorithm `Train`

Do the following once: // or `derandomize` for statelessness;

Partition D into disjoint datasets D_1 and D_2 ;

Set $q_1 \leftarrow \text{Train}(D_1)$, $q_2 \leftarrow \text{Train}(D_2)$;

Input: $c \in \mathcal{C}$

Let $i = \min\{j : c \notin D_j\}$;

Result: q_i

Algorithm 3: `CP: Copy-Protection` [VKB23]

Input: Dataset D

Learn: Run `sharded-safe`(D) to obtain q_1, q_2 as in Algorithm 2

Result: The model p with

$$p(y|x) = \frac{\min\{q_1(y|x), q_2(y|x)\}}{Z(x)}$$

where $Z(x)$ is a normalization constant that depends on x .

B.2 Failures of NAF models

B.2.1 CP can regurgitate training data

We show that `CP`—the main NAF algorithm of [VKB23]—can fail to protect against copying. Using a prompt containing no copyrightable expression, a user can cause the NAF model returned by `CP` to regurgitate copyrighted training data. This is based on an observation of Thomas Steinke [Ste].

The following claim is a slight generalization of Theorem 3.4. In words, instantiating `CP` = `CPkv` with the (tainted) training algorithm `Trainkv` (Algorithm 1) described in Example 5.2, causes training data regurgitation.

Note that Train_{kv} is itself built from some underlying training algorithm Train_0 . The only requirement we need of Train_0 is that it produces models with *full support*. That is, for every dataset D , model $p_0 \leftarrow \text{Train}_0(D)$, prompt x , and output $y \in \mathcal{W}$, we require $p(y|x) > 0$.

Theorem B.5. *Let D be a dataset and $w \in D$. For ideas $\text{id} \in \mathcal{I}$, let $D_{\text{id}} = \{w' \in D : \text{ideas}(w') = \text{id}\}$. Let Train_0 produce models with full support. Let Train_{kv} , $\text{sharded-safe}_{\text{kv}}$ and CP_{kv} be as defined in Algorithms 1, 2, and 3, defined with respect to Train_0 (and the preceding algorithms). Let $p \leftarrow \text{CP}_{\text{kv}}(D)$. Then*

$$p(D_{\text{ideas}(w)} \mid \text{ideas}(w)) = 1.$$

In particular, for any copyrighted work $c \in D$ whose unique ideas are unique, we get

$$p(c \mid \text{ideas}(c)) = 1.$$

Theorem B.5 does not contradict the CP theorem (Theorem 3.3). This is because CP is only k_x -NAF, where k_x depends on the prompt x . In our construction, the bound on k_x given by Theorem 3.3 is not vacuous for $x \in \text{ideas}(D)$.

Proof of Theorem B.5. Unrolling the algorithms, $\text{CP}_{\text{kv}}(D)$ calls $\text{sharded-safe}_{\text{kv}}$, which shards the data into D_1 and D_2 and trains $q_i \leftarrow \text{Train}_{\text{kv}}(D_i)$. Without loss of generality, suppose $w \in D_1$. Let $\text{id} = \text{ideas}(w)$.

By construction, $q_1(D_{\text{id}}|\text{id}) = 1$ (as in Example 5.2). Thus, for all outputs $y \in \mathcal{W}$:

$$p(y|\text{id}) \propto \min\{q_1(y|\text{id}), q_2(y|\text{id})\} = \begin{cases} q_2(y|\text{id}) & \text{if } y \in D_{\text{id}} \\ 0 & \text{if } y \notin D_{\text{id}} \end{cases}$$

The distribution $p(\cdot|\text{id})$ is well-defined if there exists $y \in D_{\text{id}}$ such that $q_2(y|\text{id}) > 0$. The model $q_2(\cdot|\text{id})$ returns an element of $D_2 \cap D_{\text{id}}$ if it has non-empty intersection, or it returns a sample from an underlying model with full support (the model returned by Train_0). Either way, $q_2(y|\text{id}) > 0$ for all $y \in D_2 \cap D_{\text{id}}$.

The claim follows immediately:

$$p(D_{\text{id}} \mid \text{id}) = 1 - p(\overline{D}_{\text{id}} \mid \text{id}) = 1. \quad \square$$

B.2.2 k -NAF does not prevent full reconstruction

Our next theorem shows that k -NAF may allow training data to be reconstructed, even if k is arbitrarily small and independent of x . This is because the k -NAF guarantee does not compose across a user’s many queries, and each query may leak up to k bits of training data.

We give a family of models that are k -NAF with respect to a model that returns pure noise, yet enable a user to reconstruct the dataset verbatim. For any $\ell \geq 1$, let $\text{coin}_\ell(\cdot|x)$ be uniform over $\{0, 1\}^\ell$ for all prompts x . As a generative model, $\text{coin}_\ell(\cdot|x)$ is clearly a “safe” instantiation of safe_c for any copyrighted work c .

Theorem (Theorem 3.5). *Fix $\mathcal{C} \subseteq \mathcal{W} \subseteq \{0, 1\}^*$. For $D \in \mathcal{W}^*$, let L be total the length of D in bits. There exists a (deterministic) training algorithm $\text{Train} : (D, k) \mapsto p_{k,D}$ satisfying the following.*

- For all D and $k > 0$: $p_{k,D}$ is k -NAF with respect to \mathcal{C} and coin_ℓ for $\ell = \max\{1, \lfloor k \rfloor\}$.
- There exists a user u such that for all D and $k > 0$: u makes $\text{poly}(L, 1/k)$ queries to $p_{k,D}$ and outputs D with probability > 0.99 .

Remark B.6. The theorem and proof can be adapted to more realistic **safe** by encoding D in the bias of a hash of p 's outputs. But the added complexity would obscure the technical idea used in the proof: biasing **safe** can reveal D without violating NAF.

Proof of Theorem 3.5. We parse D as a bit string of length L , and let $D[j]$ be its j th bit. For $k \geq 1$, Train outputs the model $p_{k,D}$ as follows:

$$p_{k,D}(y|x) = \begin{cases} (D[x], D[x+1], \dots, D[x+\ell-1]) & \text{if } x \in [L-\ell+1] \\ y \leftarrow_{\text{s}} \{0, 1\}^\ell & \text{otherwise} \end{cases}$$

The model $p_{k,D}$ is k -NAF with respect to \mathcal{C} and coin_ℓ : $\forall x, \Delta_{\max}(p_{k,D}(\cdot|x) \parallel \text{safe}_c(\cdot, x)) \leq k$. To reconstruct D , the user u queries $p_\ell(\cdot|x)$ for $x = i\ell + 1$ for $i = 0, 1, \dots, (L-1)/\ell$.

For $0 < k < 1$, Train sets $\beta = 2^k - 1$ and outputs the model $p_{k,D}$ as follows:

$$p_{k,D}(y|x) = \begin{cases} y \leftarrow \text{Bernoulli}(\frac{1}{2} + \beta(D[x] - \frac{1}{2})) & \text{if } x \in [L] \\ y \leftarrow \text{Bernoulli}(\frac{1}{2}) & \text{otherwise} \end{cases}$$

The intuition is that $p_{k,D}$ encodes $D[x]$ in the bias of the output of $p_{k,D}(\cdot|x)$, with the magnitude of $\beta \in (0, 1)$ controlling the strength of the bias. The model $p_{k,D}$ is k -NAF with respect to coin_1 : $\forall x, \Delta_{\max}(p_{k,D}(\cdot|x) \parallel \text{safe}_c(\cdot, x)) = \log \frac{1/2 + \beta(D[x] - 1/2)}{1/2} \leq \log(1 + \beta) \leq k$.

To determine $D[j]$ with greater than $> 1 - \frac{1}{100L}$, the user u makes $\text{poly}(L, \log 1/\beta) = \text{poly}(L, 1/k)$ queries to the model, and stores the majority. The user does this for each $j \in [L]$ and outputs the result. By a union bound, the user's output is equal to D with probability greater than 0.99. \square

B.3 VA3: an empirical attack on NAF [LSK24]

Li, Shen, and Kawaguchi propose and empirically evaluate a ‘‘Virtually Assured Amplification Attack’’ against NAF [LSK24]. At a high level, [LSK24] shows that an attacker, interacting with an NAF model in a black-box manner, can reliably induce a model to produce outputs that infringe on a target work c^* . They also give a white-box prompt writing algorithm for diffusion models (Anti-NAF) that greatly improves the performance of their attacks.

Overall, the work provides good evidence that the algorithms proposed by [VKB23] may not prevent infringement. Some uncertainty remains because the algorithm implemented in [LSK24] deviates from the original, as explained below. The paper leaves open whether k -NAF (with fixed k) prevents copyright infringement, formalizes a flawed attack model, and avoids the underlying definitional questions almost entirely.

We now describe [LSK24] in more detail, offering a somewhat different interpretation than the original.

The paper's main focus is the Amplification Attack. Given a target copyrighted work c^* , an attacker repeatedly queries a model with some prompt $x = x(c^*)$. It returns the

generation most similar to c^* . This process amplifies the one-shot probability of infringement. For example, from 0.40% to 13.64% when x is the original caption of c^* in the training data, or from 8.52% to 77.36% using the Anti-NAF prompt generator. In our view, this not actually the paper’s main negative result for NAF.

More significant (for our purposes) is the existence of prompts x whose one-shot probability of infringement is non-negligible: 0.40% or 8.52% in the previous example. The message is similar to our Theorem B.5. Namely, that the NAF constructions of [VKB23] admit models for which prompts x derived from non-copyrightable aspects of c^* can produce infringing generations. Based on the top half of [LSK24, Table 1], the prompts trivialize NAF’s protection by making $k_x \approx \log(1/\text{safe}(\text{SubSim}(c^*) \mid x))$.

Though the paper does not speculate, we suspect that the mechanism for the failure is similar to our construction in Theorem B.5. The model, repeatedly fine-tuned on c^* , regurgitates c^* with high enough probability that it survives the reweighting of CP and CP-k [LSK24, Fig. 8].

Now we turn to the paper’s limitations.

First, the attack model has a conceptual flaw, highlighting the need for our definitions-first approach. The attacker knows c^* and is actively trying to induce a similar generation [LSK24, Sec. 4.1]. Indeed, the Amplification Attack *requires* knowing c^* . This setup allows trivial attacks, like prompting `return c*`. Such an attack would not be meaningful. Any infringement would be the users’ fault, not something we care to prevent. Still, the flaw in the attack model doesn’t affect the paper’s experiments. Because they use a text-to-image model, the prompts used in the attack are just a few words long and contain nothing copyrightable.

Second, the results say nothing about k -NAF, where a fixed k upper-bounds $k_x \leq k$ for every prompt x . At best, the experiments are negative results for the particular k_x -NAF algorithm CP-k given in [VKB23], where k_x depends on x . But as we explain next, it is not entirely clear. (Our Theorems 3.4 and 3.5 are for k_x -NAF and k -NAF algorithms, respectively.)

Third, the algorithm in [LSK24] deviates from CP-k in an important way. As originally defined, CP-k is a rejection-sampling version of CP (Algorithm 3). Given a model p , a safe model safe , and constant k , prompt x , and generation y , let $\rho(y|x) = \log(p(y|x)/\text{safe}(y|x))$. Oversimplifying, CP-k repeatedly samples $y \leftarrow p(\cdot|x)$ until $\rho(y|x) \leq k$, and returns the final sample. [VKB23] proves that CP-k is k_x -NAF for $k_x = k + \log(1/\nu(x))$, where $\nu(x) = \Pr_y[\rho(y|x) \leq k]$.

[LSK24]’s implementation differs. Instead of fixing k , they fix $\nu(x)$. The experiments sample many generations $y \leftarrow p(\cdot|x)$, and return the $\nu(x) = 5\%, 10\%, \dots$ with smallest $\rho(y|x)$. This strikes us as a meaningful difference. At a minimum, the CP-k theorem doesn’t apply as is. Despite the gap, we view the results of [LSK24] as strong evidence of weaknesses of CP-k. We conjecture that there is a value of k_x for which the implemented version achieves k_x -NAF with high probability, but did not attempt to prove it.

B.4 On NAF’s legal relevance

NAF is motivated by two concepts from copyright law: *access* and *substantial similarity*. The plaintiff in a copyright infringement claim has the burden of proving that the defendant copied original expression from the copyrighted work. The plaintiff does so by proving

that (i) the defendant had access to the copyrighted work, and (ii) the defendant’s work is substantially similar to the plaintiff’s work.

With the above in mind, [VKB23] explain the relevance of NAF to copyright liability.

To show a copyright violation has occurred the plaintiff must prove that “there are substantial similarities between the defendant’s work and original elements of the plaintiff’s work” (assuming access). Its negation would be to show that defendant’s work is not substantially similar to the original elements of the plaintiff’s work. Our approach would instead correspond to showing that the defendant’s work is close to a work which was produced without access to the plaintiff’s work. [W]e think this is a stronger guarantee...

As for why it’s a stronger guarantee, the argument is as follows. The probability that the defendant’s work—produced by the real model p —is substantially similar to plaintiff’s work is not much greater than the probability would have been had the defendant used the safe model (which had no access). We heuristically expect the latter probability to be miniscule. Hence, substantial similarity between the defendant’s and plaintiff’s works is exceedingly unlikely.

Elkin-Koren et al. [EKHLM24] correctly argue that copyright cannot be “reduced to privacy.” However, this is a straw man of version of [VKB23] and the present paper (see Appendix A for additional discussion).

The sharpest criticism is by Lee, Cooper, and Grimmelman who argue that NAF is simply wrong on the law [LCG24]. “It is not a defense to copyright infringement that you would have copied the work from somewhere else if you hadn’t copied it from the plaintiff.”

We agree that NAF’s envisioned legal defense doesn’t work (technical guarantees aside). To see why, consider a case in which the model’s generated output was in fact substantially similar to a piece of training data. As to the access element, the defendant did in fact have access to the copied work by way of the model. The defendant’s work may even be so “strikingly similar” to the plaintiff’s that access becomes moot.⁹ Despite being central to NAF, *access* appears to be a red herring. Whatever copyright protection NAF offers is by way of minimizing the likelihood of producing substantially similar outputs.

C Differential privacy background [DMNS06]

An algorithm is differentially private (DP) if its output never depends too much on any one unit of input data. How much is “too much” is governed by a parameter $\epsilon > 0$. Smaller values of ϵ provide stronger guarantees. In this work, a “unit of input data” is one work w in the training dataset D .

Definition C.1 (Neighboring datasets). *Datasets $D, D' \in \mathcal{W}^*$ are neighboring if they differ by inserting or deleting a single element. We denote neighboring datasets by $D \sim D'$.*

⁹“The plaintiff can prove that the defendant copied from the work by proving by a preponderance of the evidence that ... there is a striking similarity between the defendant’s work and the plaintiff’s copyrighted work.” From the Ninth Circuit’s *Manual of Model Civil Jury Instructions* <https://www.ce9.uscourts.gov/jury-instructions/node/326>.

Definition C.2 ((ε, δ) -Differential privacy). *Let $\varepsilon, \delta \geq 0$, and let $M : \mathcal{W}^* \rightarrow \Omega$ be an algorithm mapping dataset $D \in \mathcal{W}^*$ to some output domain Ω . M is (ε, δ) -differentially private ((ε, δ) -DP) if for all neighboring pairs $D, D' \in \mathcal{W}^*$, and all subsets of outputs $S \subseteq \Omega$:*

$$\Pr[M(D) \in S] \leq e^\varepsilon \cdot \Pr[M(D') \in S] + \delta.$$

Anything that one does with the output of a DP algorithm is also DP. In DP parlance, DP is robust to post-processing in the presence of arbitrary auxiliary information. Formally, for any function f , any string aux , and any set S :

$$\Pr[f(M(D), \text{aux}) \in S] \leq e^\varepsilon \cdot \Pr[f(M(D'), \text{aux}) \in S] + \delta. \quad (4)$$

D Proof of Theorem 7.7

We restate the theorem for convenience.

Theorem D.1 (Theorem 7.7). *Fix $\beta < 1$, $n \geq 1$, and $m \geq \frac{n}{\beta}$. Suppose there exists $W = (w_{i,j}) \in \mathcal{W}^{m \times n}$ such that for all $i \neq i' \in [m]$ and all $j \in [n]$:*

$$\text{ideas}(w_{i,j}) = \text{ideas}(w_{i',j}) \quad \text{and} \quad \text{SubSim}(w_{i,j}) \cap \text{SubSim}(w_{i',j}) = \emptyset.$$

If Train be tainted with respect to ideas , then Train is not (κ, β) -clean.

The proof uses Lemma D.2, stated below.

Proof of Theorem 7.7. Fix ideas and let Train be tainted with respect to ideas . Let \mathbf{u} be the user guaranteed by taintedness of Train , and τ its output distribution. We must show Train is not (κ, β) -clean. Namely, that there exist \tilde{u} , $\tilde{\mathcal{C}}$, and \tilde{D} for which (i) \tilde{u} is β -blameless, and (ii) there exists $\tilde{\text{aux}}$ such that $\tilde{\tau}(\text{SubSim}(\tilde{\mathcal{C}}); \tilde{\text{aux}}) \geq \kappa$. Here, $\tilde{\tau}$ is \tilde{u} 's output distribution.

Let $\text{id} = \text{ideas}(w_{1,1})$. Define $\tilde{\mathbf{u}} = \mathbf{u}_{\text{id}}$ to be the user with id hard-coded that simulates $\mathbf{u}(\text{id})$. By construction, $\tilde{\tau}(E; \text{aux}) = \tau(E; \text{id})$ for all aux and all events E . Because $\tilde{\mathbf{u}}$ ignores its auxiliary input, we can invoke Lemma D.2. Thus, there exists \tilde{D} such that $\tilde{\mathbf{u}}$ is β -blameless with respect to \tilde{D} , $\tilde{\mathcal{C}} = \tilde{D}$, and Train (condition (i) above). Moreover, $w_{i,1} \in \tilde{D}$ for some $i \in [n]$.

Let $\text{id}' = \text{ideas}(w_{i,1})$. By hypothesis, $\text{id} = \text{id}'$ and therefore $\tilde{\mathbf{u}} = \mathbf{u}_{\text{id}'}$. Taintedness of Train implies condition (ii) above:

$$\tilde{\tau}(\text{SubSim}(\mathcal{C}); \perp) = \tau(\text{SubSim}(\mathcal{C}); \text{id}') = \tau(\text{SubSim}(D); \text{id}') > 0.99 \geq \kappa. \quad \square$$

Lemma D.2. *Fix $\beta < 1$, $n \geq 1$, and $m \geq \frac{n}{\beta}$. Suppose there exists $W = (w_{i,j}) \in \mathcal{W}^{m \times n}$ such that for all $i \neq i' \in [m]$ and all $j \in [n]$: $\text{SubSim}(w_{i,j}) \cap \text{SubSim}(w_{i',j}) = \emptyset$. Let \mathbf{u} be a user that ignores its auxiliary input: $\mathbf{u}^p(\text{aux}) = \mathbf{u}^p(\perp)$ for all aux . For all Train , there exists $x \in [m]^n$ such that \mathbf{u} is β -blameless with respect to $D^{(x)} = (w_{x_j,j})_{j \in [n]}$, $\mathcal{C} = D^{(x)}$, and Train .*

Proof of Lemma D.2. Because \mathbf{u} ignores aux , we omit it throughout. Take $\mathcal{C} = D$, which implies $\mathcal{C}_{-D} = \emptyset$. Hence, user \mathbf{u} is β -blameless if for all $w \in D$: $\tau_{-w}(\text{SubSim}(w)) \leq \beta$.

Let $D_{-j}^{(x)} = \text{scrub}(D^{(x)}, w_{x_j, j}) := (w_{x_{j'}, j'})_{j' \neq j}$. (For a counter example, we can choose the copyright dependency graph and hence `scrub`.) Define

$$p_j(x) = \tau_{-w_{x_j, j}}^{(x)}(\text{SubSim}(w_{x_j, j})) = \Pr_{\substack{p \leftarrow \text{Train}(D_{-j}^{(x)}) \\ z \leftarrow \mathbf{u}^p(\perp)}} [z \in \text{SubSim}(w_{x_j, j})].$$

We must show that:

$$\exists x \in [m]^n, \forall j \in [n]: p_j(x) \leq \beta. \quad (\star)$$

For $x \in [m]^n$, define $j^*(x) = \text{argmax}_{j \in [n]} p_j(x)$. Define the sets $Z_j = \{x \in [m]^n : j^*(x) = j\}$. One of the Z_j contains at least m^n/n distinct strings. Moreover, there is a subset $X^* \subseteq Z_j$ of at least $(m^n/n)/m^{n-1} = m/n$ strings that agree on all but the j th coordinate.

Summarizing, the set X^* satisfies three properties. (1) For all $x, x' \in X^*$, $j^*(x) = j^*(x') =: j^*$. (2) All $x, x' \in X^*$ disagree at the $j^*(x)$ th coordinate, and agree on all other coordinates. (3) $|X^*| \geq m/n$.

By (2), $D_{-j^*}^{(x)} = D_{-j^*}^{(x')}$ for all $x, x' \in X^*$. It follows that $\tau_{-w_{x_{j^*}, j^*}}^{(x)} = \tau_{-w_{x'_{j^*}, j^*}}^{(x')} =: \tau^*$ for all $x, x' \in X^*$. Consider the events $E_x = \text{SubSim}(w_{x_{j^*}, j^*})$ for $x \in X^*$. By hypothesis and by (2), these events are disjoint. Hence, $\sum_{x \in X^*} \tau^*(E_x) \leq 1$. By (3), there exists $x \in X^*$ such that $p_{j^*}(x) = \tau^*(E_x) \leq n/m \leq \beta$. By (1) and definition of $j^*(x)$, we have for all $j \in [n]$, $p_j(x) \leq p_{j^*}(x) \leq \beta$. This proves (\star) and completes the proof. \square