Antidistillation Sampling

Yash Savani* Asher Trockman* Zhili Feng Yixuan Even Xu Avi Schwarzschild Alexander Robey Marc Finzi J. Zico Kolter Carnegie Mellon University

https://antidistillation.com

Abstract

Frontier models that generate extended reasoning traces inadvertently produce token sequences that can facilitate model distillation. Recognizing this vulnerability, model owners may seek sampling strategies that limit the effectiveness of distillation without compromising model performance. *Antidistillation sampling* provides exactly this capability. By strategically modifying a model's next-token probability distribution, antidistillation sampling poisons reasoning traces, rendering them significantly less effective for distillation while preserving the model's utility. Our code is available at https://github.com/locuslab/antidistillation-sampling.

1 Introduction

Large language models (LLMs) trained to produce reasoning traces have achieved strong performance on math, coding, and reasoning benchmarks [1–3]. These traces, however, serve a dual purpose. They not only enhance model performance, but also enable distillation, a process by which a secondary model replicates the original model's capabilities by training on its generations [4–6]. Notably, distillation can result in substantial capability gains at a fraction of the computational cost needed to train similarly performant models from scratch.

While effective and efficient, the viability of distillation poses several downsides for companies deploying frontier reasoning models. First, returning reasoning traces represents a forfeiture of intellectual property, which can allow competitors to cheaply replicate frontier capabilities. Second, the threat of distillation incentivizes limiting user access by obscuring token probabilities or truncating reasoning traces. Finally, model safety is often not preserved by distillation, which enables the generation of harmful content [7, 8].

To address these issues, we introduce *antidistillation sampling* (see Figure 1). The main idea underpinning antidistillation sampling is to adjust a model's sampling distribution so that generated traces maintain high likelihood under the unadjusted distribution, and distillation attempts are simultaneously poisoned. To operationalize this idea, we first formulate the general problem of poisoning reasoning models trained via distillation. We then derive one solution to this problem (see Algorithm 1), which facilitates a precise trade-off between two competing objectives—the utility of the original model and the effectiveness of distillation poisoning—while incurring minimal computational overhead.

To illustrate our empirical results, consider a reasoning model that achieves 72% accuracy on MMLU. Naively distilling this model using greedy sampling can produce a student that reaches up to 52% accuracy. If the teacher model increases its sampling temperature, its accuracy slightly decreases (e.g., by 4%), yet the distilled student's accuracy remains largely unchanged at 52%. In contrast,

^{*}Equal contribution. Contact: {ysavani,ashert}@cs.cmu.edu.

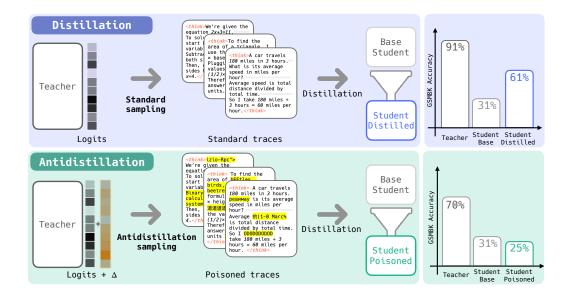


Figure 1: Reasoning traces generated via antidistillation sampling poison distillation attempts while simultaneously preserving the teacher's performance. The teacher's logits are perturbed in a direction Δ , leading to samples that significantly degrade distilled model performance relative to naive temperature sampling. For more details, see Figure 2 and §3.

using antidistillation sampling with the same 4% reduction in teacher accuracy significantly reduces the student's capabilities, lowering its accuracy to as low as 40%. Our findings on GSM8K [9], MATH [10], and MMLU [11] indicate that model owners can effectively limit distillation quality via antidistillation sampling.

2 Related work

Model distillation. The prominence and effectiveness of model distillation—and, more generally, model compression [12–14]—is rooted in a rich literature dating back to Schmidhuber [5], and, more recently, to Hinton et al. [4]. Since these seminal works, a growing body of literature has sought to benchmark the performance of distilled models and to algorithmically maximize the effectiveness of distillation [15–19]. Indeed, a variety of frontier AI labs have incorporated distillation as a core technique, both to efficiently enable frontier capabilities [20, 21] and to improve model safety via context distillation [22–24]. However, this practice constitutes a strategic vulnerability for frontier model maintainers, given the demonstrated value of these reasoning traces (see, e.g., [1, 25–27]).

Model security. The setting we address in this paper—where a student model is trained on data generated by a teacher model—intersects with several aspects of model security. For instance, model extraction attacks acquire weights via query-level access, whereas training data extraction attacks are designed to harvest training data [28, 29]. While antidistillation sampling may offer some protection against these attacks, such analysis remains beyond our scope. More relevant is the literature on data poisoning, where maliciously crafted data is injected into a model's training set to induce specific downstream effects (see, e.g., [30]). In this vein, Rando and Tramèr [31] show the effectiveness of adding backdoors to preference data, sabotaging LLMs finetuned with RLHF. Our contribution bridges data poisoning and privacy techniques to protect the valuable knowledge encoded in frontier models.

Distillation prevention. A related line of work has sought to develop algorithms that prevent distillation. In the context of computer vision, Ma et al. [32] corrupt the teacher's logits via self-training, whereas follow-up work shows that returning only the top-k logits tends to harm distillation [33]. Also related are watermarking algorithms, which seek to adjust model logits to detect whether a model has been distilled [34–37]. And while this family of methods preserves teacher accuracy, their *static* nature—each input generally yields a deterministic logit vector—presents security vulnerabilities: a

Algorithm 1: Antidistillation sampling

Input: Prompt $x_{1:n}$, max tokens N, penalty multiplier λ , approximation parameter ϵ , temperature τ

1. (Initialization) Compute the gradient of the downstream loss

$$g \leftarrow \nabla \ell(\theta_P)$$

- 2. For each token index $t = n, n + 1, \dots, N 1$:
 - i. Compute the antidistillation penalty term

$$\widehat{\Delta}(\cdot | x_{1:t}) \leftarrow \frac{\log p(\cdot | x_{1:t}; \theta_P + \epsilon g) - \log p(\cdot | x_{1:t}; \theta_P - \epsilon g)}{2\epsilon}$$

ii. Sample the next token x_{t+1} from the teacher's adjusted distribution

$$x_{t+1} \sim \frac{1}{Z} \exp\left(\frac{1}{\tau} \log p(\cdot | x_{1:t}; \theta_T) + \lambda \widehat{\Delta}(\cdot | x_{1:t})\right)$$

Output: Sampled sequence $x_{1:N}$

distiller can learn an inverse transformation by saving input-output pairs, and thereby fine-tune to recover uncorrupted logits. To mitigate this shortcoming, Chen et al. [38] propose a *session-dynamic* defense that monitors the sensitivity of a user's queries and perturbs the logits once a threshold is crossed. In contrast, antidistillation sampling is fully *dynamic*; it perturbs each token's distribution on-the-fly using gradients from a hidden proxy model, turning generation into a moving target in a similar fashion to cryptographic stream ciphers.

Language model decoding. Finally, we position antidistillation sampling within the broader framework of controlled decoding for LLMs [39], where supplementary objectives steer the decoding process. Existing approaches in this domain include using contrastive objectives to enhance generation quality [40], reformulating constrained decoding as an optimization problem [41], and incorporating energy-based constraints [42]. While related, antidistillation sampling solves a different problem: by implementing a new, distillation-aware penalization term in the decoding objective, our approach poisons generated reasoning traces to undermine the performance of models fine-tuned on these outputs.

3 Antidistillation sampling

To motivate antidistillation sampling, we first sketch a high-level overview of our problem setting in §3.1. Based on this setting, we provide a desiderata outlining the desired qualities for poisoning distillation attempts in §3.2. We then derive the antidistillation sampling method (summarized in Algorithm 1) in §3.3.

3.1 An overview of antidistillation

The core objective of antidistillation sampling is to adjust a model's next-token distribution to balance two competing goals: sampling tokens with high likelihood under the original, unadjusted distribution and sampling tokens that effectively poison distillation attempts. Throughout, we refer to the model from which reasoning traces are sampled as the *teacher*, and the model being distilled as the *student*.

Our derivation relies on quantifying how model distillation impacts the student model's performance on a given downstream task. This analysis yields a key insight—we can incorporate this performance metric directly into the teacher's sampling distribution. This takes the form of a directional derivative capturing the change in the teacher's sampling distribution along the update direction in the student's weight space. However, due to the high cost needed to compute this directional derivative, the final portion of our derivation identifies an efficient finite-difference approximation for this term, which is inexpensive to compute and, as we demonstrate in §4, results in effective distillation poisoning.

3.2 Preliminaries

We consider an LLM to be a mapping from a sequence of input tokens $x_{1:t} = (x_1, \dots, x_t)$ to a distribution over the next token, where each token is an element of a vocabulary set $\mathcal{V} = \{1, \dots, V\}$. This distribution is parameterized by weights θ and can be expressed as $p(\cdot|x_{1:t};\theta)$. We write $p(\cdot|x_{1:t};\theta)$ to denote the distribution of all next-token probabilities, whereas $p(x_{t+1}|x_{1:t};\theta)$ refers to the scalar probability of a given next token x_{t+1} . Typically, tokens are generated according to a scaled version of this distribution:²

$$x_{t+1} \sim \frac{1}{Z} \exp\left(\frac{1}{\tau} \log p(\cdot|x_{1:t};\theta)\right).$$
 (1)

Here, τ is the temperature and Z is a normalization term, which is computed by summing the exponential term over all possible next tokens. Using a temperature of $\tau = 0$ corresponds to greedy sampling, in which x_{t+1} is deterministically chosen to be the token with the largest log probability under the current model parameter θ .

Desiderata for antidistillation. Distillation involves a student model—parameterized by θ_S , with a distribution over next tokens given by $p(\cdot|x_{1:t};\theta_S)$ —trained on data generated from a teacher model parameterized by θ_T . These models do not need to share the same parameter space, and therefore the parameter vectors θ_S and θ_T need not be comparable; indeed, a student model may have substantially fewer parameters than the teacher.

The aim of antidistillation sampling is to generate tokens that perform well according to a metric used to the evaluate teacher, while simultaneously having the property that training on these tokens *cannot* improve performance on this same task. In more detail, we aim to adjust the teacher's sampling procedure to simultaneously satisfy the following:

- I. Non-distillablity. Student models trained on tokens sampled via antidistillation sampling should have a degraded performance on a chosen downstream task relative to training on tokens sampled from the teacher's nominal distribution.
- II. Nominal utility. Tokens sampled via antidistillation sampling should remain probable under the teacher's unadjusted sampling scheme $p(\cdot | x_{1:t}; \theta_T)$.

Taken together, these goals ensure that the teacher model maintains its nominal performance while simultaneously preventing distillation on downstream tasks.

Proxy models. In general, we do not expect to know the distilled student's model architecture in advance. Therefore, rather than assuming access to the true student model, we develop antidistillation sampling based on the notion of a *proxy student model*, which, for simplicity, we refer to as the *proxy model*. The proxy model is parameterized by θ_P , and specifies a sampling distribution $p(\cdot|x_{1:t};\theta_P)$. A key aspect we consider below is whether the process generalizes, i.e., whether traces via antidistillation sampling to prevent the proxy model from distilling the teacher also prevent the distinct student models from distilling.

3.3 Deriving antidistillation sampling

To operationalize antidistillation sampling, we first assume access to a differentiable, real-valued downstream loss ℓ , which measures the proxy model's performance on a given downstream task. Throughout, we take ℓ to be the negative log-likelihood for generating a sequence of tokens on a fixed, potentially large dataset. For instance, ℓ could represent the cross entropy loss of predicting each token across a large reasoning benchmark. However, ℓ can be chosen very broadly to capture any student capability that the teacher model maintainer may want to influence via poisoning. A key point is that ℓ can be very costly to compute, as it may require evaluating the proxy model over a large and diverse set of data.

²Variants of this sampling scheme include top-k sampling (i.e., limiting sampling to the tokens with the top-k largest probabilities), greedy sampling (i.e., sampling from the same objective while letting $\tau \to 0$), and beam search, but we focus mainly on temperature-based sampling here.

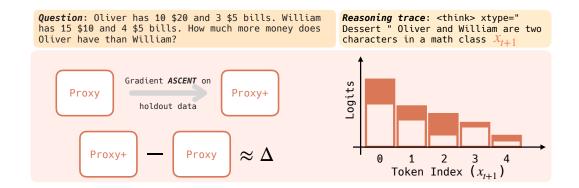


Figure 2: An illustration of approximating Δ . The teacher model performs antidistillation sampling autoregressively, based on its perturbed distribution by Δ . Given an input prompt and t reasoning tokens from the teacher, Δ is approximated by the difference of the log probability of each token in the vocabulary between two copies of the proxy model (created by performing a single *gradient ascent* step using the downstream task loss on the proxy model); this difference is represented by the area in the bar plot.

Given the non-distillability criteria outlined above, the goal of antidistillation sampling is for the downstream loss $\ell(\theta_P)$ to increase³ whenever the student is fine-tuned on sequences of tokens generated by the teacher. To capture this, first consider the change in θ_P that results from fine-tuning to minimize the negative log-likelihood of a token x_{t+1} generated by the teacher. Specifically, we consider one step of optimization via gradient descent on θ_P :

$$\theta_P^+ = \theta_P - \eta \nabla_{\theta_P} \left(-\log p(x_{t+1}|x_{1:t}; \theta_P) \right) \tag{2}$$

$$= \theta_P + \eta \nabla_{\theta_P} \log p(x_{t+1}|x_{1:t};\theta_P) \tag{3}$$

where $\eta > 0$ is the step size. The impact of this update can then be quantified by measuring the difference in the loss ℓ before and after this update. In particular, for each token $x_{t+1} \in \mathcal{V}$, we define the following difference term

$$\Delta(x_{t+1}|x_{1:t}) = \ell(\theta_P^+) - \ell(\theta_P) = \ell(\theta_P + \eta \nabla_{\theta_P} \log p(x_{t+1}|x_{1:t};\theta_P)) - \ell(\theta_P). \tag{4}$$

If $\Delta(x_{t+1}|x_{1:t})$ is positive, the update in eq. (3) increases the loss; if $\Delta(x_{t+1}|x_{1:t})$ is negative, the update decreases the loss. Thus, our goal is to adjust the teacher's sampling distribution so that tokens sampled from the teacher both have (1) high likelihood under the teacher's unadjusted distribution and (2) yield larger (i.e., more positive) values of Δ .

To implement antidistillation sampling, we propose adding a penalty, proportional to $\Delta(x_{t+1}|x_{1:t})$, to the teacher's unadjusted log probabilities $\log p(x_{t+1}|x_{1:t};\theta_T)$. This results in the following adjusted sampling distribution

$$x_{t+1} \sim \frac{1}{Z} \exp\left(\frac{1}{\tau} \log p(\cdot | x_{1:t}; \theta_T) + \lambda \Delta(\cdot | x_{1:t}; \theta_P)\right), \tag{5}$$

where Z is a normalization term appropriately scaled (relative to eq. (1)) to accommodate the penalty, and $\lambda > 0$ is a regularization coefficient that facilitates a trade-off between sampling from the teacher's distribution and sampling tokens that maximally increase student's downstream loss. Unfortunately, directly implementing eq. (5) is impractical, as we would need to compute $\Delta(x_{t+1}|x_{1:t})$ for each potential next token $x_{t+1} \in \mathcal{V}$, requiring V gradients to be computed as well as V evaluations of the downstream loss ℓ , which, in turn, is assumed to involve a lengthy computation to produce.

An efficient implementation. The core of our proposed approach is an efficient mechanism to approximate the sampling process above. As a starting point, observe that $\Delta(x_{t+1}|x_{1:t})$ can be scaled by a factor of $1/\eta$ without changing the relative penalties for each x_{t+1} (i.e., we could fold this term

³We assume without loss of generality that increases in $\ell(\theta_P)$ are desirable from the perspective of the poisoner; the procedure is easily adaptable to problems wherein the goal is to decrease $\ell(\theta_P)$.

into the λ regularization penalty). Then, by taking the limit of $\Delta(x_{t+1}|x_{1:t})/\eta$ as $\eta \to 0$, we have that

$$\lim_{\eta \to 0} \frac{1}{\eta} \Delta(x_{t+1}|x_{1:t}) = \lim_{\eta \to 0} \frac{\ell(\theta_P + \eta \nabla_{\theta_P} \log p(x_{t+1}|x_{1:t};\theta_P)) - \ell(\theta_P)}{\eta}$$

$$= \langle \nabla \ell(\theta_P), \nabla_{\theta_P} \log p(x_{t+1}|x_{1:t};\theta_P) \rangle.$$
(6)

$$= \langle \nabla \ell \left(\theta_P \right), \nabla_{\theta_P} \log p(x_{t+1} | x_{1:t}; \theta_P) \rangle. \tag{7}$$

That is, the limit is the inner product between the gradient $\nabla_{\theta_P} \log p(x_{t+1}|x_{1:t};\theta_P)$ and the downstream loss gradient $\nabla \ell(\theta_P)$. Notice that the expression in eq. (7) no longer involves the evaluation of the downstream loss for each token in \mathcal{V} . Rather, $\nabla \ell(\theta_P)$ can be computed and stored once, after which the only remaining task is to efficiently evaluate eq. (7) for each token $x_{t+1} \in \mathcal{V}$. To do so, the key observation is that the directional derivative is symmetrical. Thus, we can rewrite eq. (7) as a finite difference limit in the *other* term, i.e., in terms of a finite difference update to $\log p(x_{t+1}|x_{1:t};\theta_P)$. This gives

$$\lim_{\eta \to 0} \frac{1}{\eta} \Delta(x_{t+1}|x_{1:t}) = \langle \nabla \ell(\theta_P), \nabla_{\theta_P} \log p(x_{t+1}|x_{1:t}; \theta_P) \rangle \tag{8}$$

$$= \lim_{\epsilon \to 0} \frac{\log p(x_{t+1}|x_{1:t}; \theta_P + \epsilon \nabla \ell(\theta_P)) - \log p(x_{t+1}|x_{1:t}; \theta_P - \epsilon \nabla \ell(\theta_P))}{2\epsilon}$$
(9)

Importantly, this difference involves only the computation of next-token probabilities under two different models: the original proxy model θ_P and an updated copy of the proxy model $\theta_P + \epsilon \nabla \ell(\theta_P)$. These models can be saved once before any sampling, and then an approximation of the antidistillation sampling term can be computed for all next tokens simply via two forward passes in the proxy model. In other words, we define

$$\widehat{\Delta}(\cdot|x_{1:t}) = \frac{\log p(\cdot|x_{1:t};\theta_P + \epsilon \nabla \ell(\theta_P)) - \log p(\cdot|x_{1:t};\theta_P - \epsilon \nabla \ell(\theta_P))}{2\epsilon}$$
(10)

for some appropriately chosen small value of ϵ , where $\widehat{\Delta}(x_{t+1}|x_{1:t})$ approaches eq. (7) for all next tokens x_{t+1} in the limit as $\epsilon \to 0$. Intuitively, $\Delta(x_{t+1}|x_{1:t})$ measures how much sampling token x_{t+1} would degrade a proxy student's performance after a single update. Finally, we sample according to the teacher's adjusted sampling distribution:

$$x_{t+1} \sim \frac{1}{Z} \exp\left(\frac{1}{\tau} \log p(\cdot | x_{1:t}; \theta_T) + \lambda \widehat{\Delta}(\cdot | x_{1:t})\right). \tag{11}$$

In Algorithm 1, we summarize the procedure outlined in this section. Concretely, given a prompt $x_{1:t}$, using antidistillation sampling to generate a new token x_{t+1} involves: (1) (once, at initialization) computing the gradient of the downstream loss; and (2) (for each token to be generated) compute the finite-difference approximation of $\Delta(\cdot|x_{1:t})$ and sample the token from the teachers adjusted softmax distribution.

Exploring antidistillation in practice

Through a range of experiments, we demonstrate the effectiveness of antidistillation sampling and discuss several interesting phenomena. First, we show that the hyperparameter λ provides model owners with precise control over the trade-off between nominal utility and non-distillability. This trade-off persists across various teacher-student model configurations, and notably, remains effective even when the proxy student is from a different model family than the actual student—validating the practical applicability of our method in realistic scenarios where the model owners lack knowledge of potential student architectures. Additionally, we address methodological questions through a pointed empirical analysis.

4.1 Experimental setup

First, we detail our selection of model architectures and benchmark datasets, chosen to represent realistic distillation scenarios across varied reasoning tasks. Next, we describe the computation of $\nabla \ell(\theta_P)$ (eq. (10)). Finally, we outline our baseline comparison methodology.

Architectures. To demonstrate the effectiveness of antidistillation sampling in practice, we simulate realistic distillation by instantiating distinct teacher, proxy student, and actual student models.

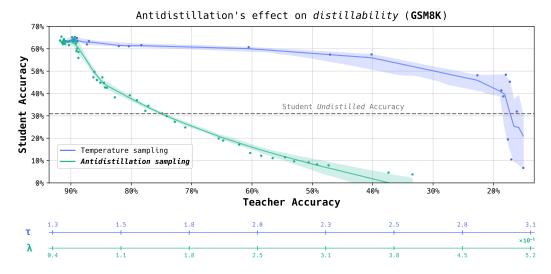


Figure 3: Antidistillation sampling uses a tunable parameter λ to control the trade-off between teacher accuracy and distillability. The baseline involves sampling from the teacher with increasing temperature τ to show that we can produce traces that are bad for distillation at some cost in teacher accuracy. One important feature of the blue temperature sampling curve is that to bring the student accuracy down below the undistilled accuracy, the teacher performance has to drop to 20%. On the other hand, with antidistillation sampling, the teacher model can still get 70% accuracy while producing traces that bring the student's performance down below the undistilled accuracy.

Specifically, we use deepseek-ai/DeepSeek-R1-Distill-Qwen-7B [21] as the teacher model, Qwen/Qwen2.5-3B [43] as the proxy model, and meta-llama/Llama-3.2-3B [44] as the student model (we examine other architecture configurations in §4.3).

Benchmarks. We evaluate the performance of antidistillation sampling on GSM8K [9] (we use GSM8K Platinum for the test set [45]), MATH [10], and MMLU [11] benchmarks (all provided under the MIT license), which are particularly suitable for our study, as they require high-quality reasoning traces for strong performance. To evaluate model performance, we use free-form generation after the prompt to get the reasoning trace, we then concatenate "\n\n**Final Answer**\n[\boxed{" after the reasoning trace and continue to generate for 32 additional answer tokens. Finally, we evaluate the model accuracy on the answer provided within "\boxed{...}". We also report *undistilled* student baselines; since base models have very low accuracy without distillation, we use in-context learning with reasoning examples showing the correct output format.

Calculating the downstream loss. Calculating $\nabla \ell(\theta_P)$ requires evaluating the proxy model on a holdout set of reasoning traces. For our experiments, we use the first 70% of our train data as the training set and the remaining 30% as the holdout set. We use the teacher to generate reasoning traces on the holdout set, and then calculate $\nabla \ell(\theta_P)$ on these reasoning traces using gradient accumulation while masking out the system and question prompt.

Baselines. Our primary baseline is temperature sampling, which approximates standard API endpoint behavior while providing a controlled way to degrade teacher performance. Importantly, temperature sampling ensures we fairly compare trade-offs against a straightforward baseline that—like our method—degrades teacher performance by modifying the sampling procedure. One other point of comparison to this baseline is that antidistillation sampling requires two forward passes on the proxy model for each forward pass on the teacher, independent of λ . Since we choose the proxy model to be approximately half the size of the teacher, this amounts to doubling the computation needed to sample outputs from the model compared to temperature sampling. In practice, one might choose a much smaller proxy model to reduce the overhead further (see Figure 11). We also introduce a permutation sampling baseline in §A that preserves the statistical properties of $\widehat{\Delta}$ while scrambling the gradient information, providing definitive evidence that the computational overhead of producing $\widehat{\Delta}$ is necessary to achieve our non-distillability objective.

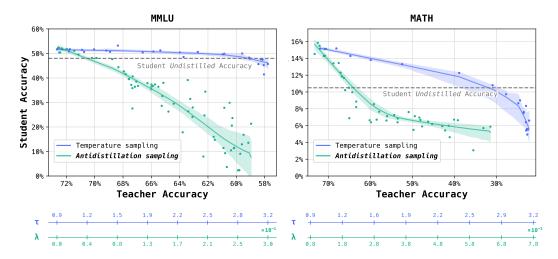


Figure 4: For both MMLU and MATH data, we show that antidistillation sampling can bring student accuracies down with relatively little cost to the teacher.

Hyperparameters. Antidistillation sampling involves two key hyperparameters: ϵ , which controls the approximation power of the finite-difference computation, and λ , which determines the weight of the antidistillation penalty in the sampling distribution from eq. (11).

- For ϵ , we empirically verify that our finite difference approximation in eq. (10) closely matches the JVP result in eq. (7) (see §B). In practice, we find that $\epsilon = 10^{-2}$ works well for BFloat16 models, which is close to the minimum in Figure 7.
- For λ , we conduct a comprehensive sweep to characterize the utility-distillability trade-off that results from perturbing the sampling distribution.

We use a max generation length of 1024 for both GSM8K and MMLU and 2048 for MATH. For antidistillation sampling, we use a temperature of $\tau=0.6$; we found that sweeping between $\tau\in[0,1]$ does not significantly impact antidistillation performance. All of our experiments are performed on nodes with 8 NVIDIA H100 GPUs and we use the transformers package [46], the trl toolkit [47], and the accelerate library [48].

Distillation protocol. All distillation experiments use LoRA [49] with rank 128, $\alpha=128$, and dropout probability 0 Our optimization protocol employs a learning rate of 0.0005, weight decay coefficient of 0.1, and gradient clipping at norm 1.0. Training follows a cosine learning rate schedule with warm-up over the first 10% of training, batch size 32, for 4 epochs. These values are the result of a systematic hyperparameter sweep using the MATH dataset to find configurations that maximize student performance gain.

4.2 Controlling the utility-distillability trade-off

The bar plots to the right in Figure 1 show that antidistillation sampling effectively satisfies the desiderata outlined in §3.2. Specifically, for a fixed reduction in teacher performance, students distilled from sampled traces exhibit substantially lower accuracy compared to those distilled from temperature-scaled traces. These initial results, while compelling, represent just one point in the configuration space as they reflect a particular choice of λ .

In Figures 3 and 4, we vary λ to characterize the degree of control antidistillation sampling provides over the utility-distillability trade-off. Since our experiments use architecturally distinct student and proxy student models, these results confirm that antidistillation sampling generalizes effectively across model families—a critical property for practical deployment.

Recognizing that model owners typically have limited tolerance for sacrificing utility, we explicitly focus on the high-teacher-performance regime in Figure 9. Even in this setting, where teacher performance degradation is conservative, we observe meaningful degradation in student accuracy,

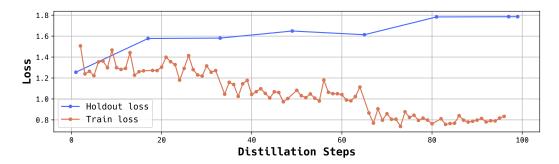


Figure 5: Distillation loss curves show that although the student's training loss decreases across steps, antidistillation sampling effectively poisons traces, as shown by the increasing student's holdout loss.

underscoring the method's practical efficacy. For example, going from 90% to 89% teacher accuracy leads to the poisoned student dropping from 65% to 56% accuracy, while temperature sampling doesn't degrade the student's performance at all. For illustrative examples comparing traces generated across comparable λ and τ settings, see §E.

4.3 Diverse configurations for antidistillation sampling

To probe the efficacy of antidistillation sampling across diverse scenarios, we conduct experiments with various teacher-student configurations and datasets. Our primary setup uses Qwen teacher and proxy student models while using a Llama model for the student. We also investigate settings where all the models (student, proxy, and teacher) belong to the same architecture family—either all from the Llama architecture family or all from the Qwen architecture family—evaluated on GSM8K. Results are provided in Figure 10.

Beyond architectural variations, we also validate our finite difference approximation by comparing it with the theoretically-motivated Jacobian-vector product (JVP) implementation (not to be confused with a vector-Jacobian product used in backpropagation) (see §C). Both approaches yield similar results in practice, confirming that our computationally efficient finite difference method provides a reliable approximation to the formal gradient-based objective.

4.4 Generalizing from the proxy model

Our method demonstrates strong results across various settings, but an important question remains about its underlying mechanism. Since we sample tokens explicitly designed to be detrimental for the proxy model on our holdout set, we are relying on generalization to an unknown student model. Figure 5 provides insight into this mechanism by tracking loss dynamics during the distillation process. We observe exactly the intended effect. Distillation on the antidistillation traces lowers the student's loss on the training set while increasing its loss on the holdout set. This result confirms that antidistillation sampling creates traces that are learnable but poison the student model's ability to reason on the downstream task.

In all our experiments, we keep the size discrepancy between proxy and student model relatively small (both are 3B models). However, Figure 11 demonstrates that our method remains effective even when proxy and student models differ in size, suggesting robustness to this architectural mismatch.

5 Conclusion

The value of proprietary frontier LLMs necessitates that their owners do what they can to protect their assets. As evidenced by the fact that the frontier companies limit exposure to their models via black-box APIs, these companies are already considering the threat of model stealing. However, given the recent attention paid to the effectiveness of distillation, it is imperative that model maintainers who wish to protect the information stored in their models guard against distillation. This paper provides a proof-of-concept that antidistillation sampling—which adjusts a model's sampling distribution—is

effective in blocking such attacks. We are excited at the prospect of continuing to refine and scale this approach, particularly with a view toward more secure future frontier models.

Broader impact. We expect that antidistillation sampling will have a positive impact on the security of frontier models. By providing a mechanism to protect against distillation, we hope to encourage the continued development of frontier large language models and their applications.

References

- [1] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv* preprint arXiv:2412.16720, 2024.
- [2] Anthropic. Claude 3.5 sonnet model card addendum, 2024. URL https://www.anthropic.com/news/claude-3-5-sonnet. Accessed: 2025-05-12.
- [3] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [4] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [5] Jürgen Schmidhuber. Learning complex, extended sequences using the principle of history compression. *Neural computation*, 4(2):234–242, 1992.
- [6] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129(6):1789–1819, 2021.
- [7] Mahdi Sabbaghi, Paul Kassianik, George Pappas, Yaron Singer, Amin Karbasi, and Hamed Hassani. Adversarial reasoning at jailbreaking time. *arXiv preprint arXiv:2502.01633*, 2025.
- [8] Matt Burgess and Lily Hay Newman. Deepseek's safety guardrails failed every test researchers threw at its ai chatbot. URL https://www.wired.com/story/deepseeks-ai-jailbreak-prompt-injection-attacks/.
- [9] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- [10] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.
- [11] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv* preprint *arXiv*:2009.03300, 2020.
- [12] Cristian Buciluă, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 535–541, 2006.
- [13] Yu Cheng, Duo Wang, Pan Zhou, and Tao Zhang. A survey of model compression and acceleration for deep neural networks. *arXiv preprint arXiv:1710.09282*, 2017.
- [14] Antonio Polino, Razvan Pascanu, and Dan Alistarh. Model compression via distillation and quantization. *arXiv preprint arXiv:1802.05668*, 2018.
- [15] Kumar Shridhar, Alessandro Stolfo, and Mrinmaya Sachan. Distilling reasoning capabilities into smaller language models. *Findings of the Association for Computational Linguistics: ACL 2023*, pages 7059–7073, 2023.
- [16] Jierui Li and Raymond Mooney. Distilling algorithmic reasoning from llms via explaining solution programs. *arXiv preprint arXiv:2404.08148*, 2024.
- [17] Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Minillm: Knowledge distillation of large language models. *arXiv preprint arXiv:2306.08543*, 2023.
- [18] Micah Goldblum, Liam Fowl, Soheil Feizi, and Tom Goldstein. Adversarially robust distillation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 3996–4003, 2020.

- [19] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014.
- [20] New York Times. Openai says it has evidence china's deepseek used its model to train competitor. *The New York Times*, January 2025. URL https://www.nytimes.com/2025/01/29/technology/openai-deepseek-data-harvest.html. Accessed: 2025-04-10.
- [21] DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.
- [22] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. arXiv preprint arXiv:2204.05862, 2022.
- [23] Charlie Snell, Dan Klein, and Ruiqi Zhong. Learning by distilling context. *arXiv preprint* arXiv:2209.15189, 2022.
- [24] Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, et al. A general language assistant as a laboratory for alignment. *arXiv preprint arXiv:2112.00861*, 2021.
- [25] Melody Y Guan, Manas Joglekar, Eric Wallace, Saachi Jain, Boaz Barak, Alec Helyar, Rachel Dias, Andrea Vallone, Hongyu Ren, Jason Wei, et al. Deliberative alignment: Reasoning enables safer language models. *arXiv* preprint arXiv:2412.16339, 2024.
- [26] Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. In *The Twelfth International Conference on Learning Representations*, 2023.
- [27] Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process-and outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.
- [28] Nicholas Carlini, Daniel Paleka, Krishnamurthy Dj Dvijotham, Thomas Steinke, Jonathan Hayase, A Feder Cooper, Katherine Lee, Matthew Jagielski, Milad Nasr, Arthur Conmy, et al. Stealing part of a production language model. *arXiv preprint arXiv:2403.06634*, 2024.
- [29] Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A Feder Cooper, Daphne Ippolito, Christopher A Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. Scalable extraction of training data from (production) language models. *arXiv preprint arXiv:2311.17035*, 2023.
- [30] Micah Goldblum, Dimitris Tsipras, Chulin Xie, Xinyun Chen, Avi Schwarzschild, Dawn Song, Aleksander Mądry, Bo Li, and Tom Goldstein. Dataset security for machine learning: Data poisoning, backdoor attacks, and defenses. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(2):1563–1580, 2022.
- [31] Javier Rando and Florian Tramèr. Universal jailbreak backdoors from poisoned human feedback. In *International Conference on Learning Representations*, 2024.
- [32] Haoyu Ma, Tianlong Chen, Ting-Kuei Hu, Chenyu You, Xiaohui Xie, and Zhangyang Wang. Undistillable: Making a nasty teacher that cannot teach students. *arXiv* preprint *arXiv*:2105.07381, 2021.
- [33] Haoyu Ma, Yifan Huang, Tianlong Chen, Hao Tang, Chenyu You, Zhangyang Wang, and Xiaohui Xie. Stingy teacher: Sparse logits suffice to fail knowledge distillation. 2022.
- [34] Xuandong Zhao, Yu-Xiang Wang, and Lei Li. Protecting language generation models via invisible watermarking. In *International Conference on Machine Learning*, pages 42187–42199. PMLR, 2023.
- [35] Chenchen Gu, Xiang Lisa Li, Percy Liang, and Tatsunori Hashimoto. On the learnability of watermarks for language models. *arXiv preprint arXiv:2312.04469*, 2023.

- [36] Tom Sander, Pierre Fernandez, Alain Durmus, Matthijs Douze, and Teddy Furon. Watermarking makes language models radioactive. Advances in Neural Information Processing Systems, 37: 21079–21113, 2024.
- [37] John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. A watermark for large language models. In *International Conference on Machine Learning*, pages 17061–17084. PMLR, 2023.
- [38] Huajie Chen, Tianqing Zhu, Lefeng Zhang, Bo Liu, Derui Wang, Wanlei Zhou, and Minhui Xue. Queen: Query unlearning against model extraction. *IEEE Transactions on Information Forensics and Security*, 2025.
- [39] Sidharth Mudgal, Jong Lee, Harish Ganapathy, YaGuang Li, Tao Wang, Yanping Huang, Zhifeng Chen, Heng-Tze Cheng, Michael Collins, Trevor Strohman, et al. Controlled decoding from language models. *arXiv* preprint arXiv:2310.17022, 2023.
- [40] Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. Contrastive decoding: Open-ended text generation as optimization. *arXiv preprint arXiv:2210.15097*, 2022.
- [41] Haozhe Ji, Pei Ke, Hongning Wang, and Minlie Huang. Language model decoding as direct metrics optimization. *arXiv preprint arXiv:2310.01041*, 2023.
- [42] Lianhui Qin, Sean Welleck, Daniel Khashabi, and Yejin Choi. Cold decoding: Energy-based constrained text generation with langevin dynamics. Advances in Neural Information Processing Systems, 35:9538–9551, 2022.
- [43] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2 technical report. arXiv preprint arXiv:2407.10671, 2024.
- [44] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [45] Joshua Vendrow, Edward Vendrow, Sara Beery, and Aleksander Madry. Do large language model benchmarks test reliability? *arXiv preprint arXiv:2502.03461*, 2025.
- [46] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020. Association for Computational Linguistics. URL https://www.aclweb.org/anthology/2020.emnlp-demos.6.
- [47] Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, Shengyi Huang, Kashif Rasul, and Quentin Gallouédec. Trl: Transformer reinforcement learning. https://github.com/huggingface/trl, 2020.
- [48] Sylvain Gugger, Lysandre Debut, Thomas Wolf, Philipp Schmid, Zachary Mueller, Sourab Mangrulkar, Marc Sun, and Benjamin Bossan. Accelerate: Training and inference at scale made simple, efficient and adaptable. https://github.com/huggingface/accelerate, 2022.
- [49] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1 (2):3, 2022.

[50] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in neural information processing systems*, 35:16344–16359, 2022.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: They accurately reflect the paper's contributions and scope.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The main experimental results emphasize that even if our method is used, there is still a trade-off between the performance of the teacher model and the distillability.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: See section 3.3 for relevant details.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: See section 4 for relevant details.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The experiment code base is uploaded to supplementary materials. Additionally, the data and code will be released when the paper is publicly available.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new
 proposed method and baselines. If only a subset of experiments are reproducible, they
 should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: See section 4 for relevant details.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: The results are accompanied by appropriate error bars.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The amount of compute required is provided in section 4.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: This research conforms with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: See section 5 for a discussion of the broader impacts of our work.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper does not release any data or models that have a high risk for misuse.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We have properly credited the creators of the assets used in our paper.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

• If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not release any new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- · Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- · For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: This research is about sampling methods for LLMs.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

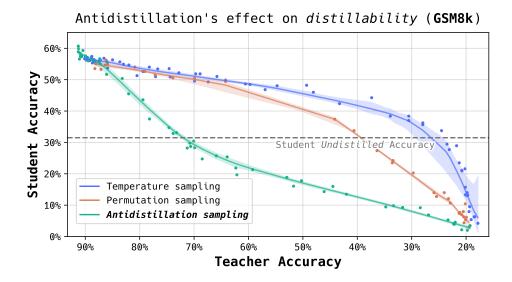


Figure 6: *Permutation sampling* is a strong baseline where we destroy the information in antidistillation sampling while preserving statistical properties via random permutation and sign flipping.

A Additional Baselines

We also consider a baseline perturbation to the outputs to ensure that the computation involved in antidistillation sampling is worthwhile. This method adds random perturbations to the logits and we call this noisy sampling. While many choices of how to add noise to the output of an LLM exist, we find that matching the statistics of the perturbations computed by antidistillation sampling is the best way to find interventions that lead to the same teacher accuracy. Therefore, we randomly permute and flip the sign of the perturbations computed with antidistillation sampling to execute *permutation sampling*, a specific type of noisy sampling; we show the results of perturbation sampling in Figure 6.

B Verifying hyperparameter choice ϵ

We empirically verify that the finite difference in eq. (10) behaves as expected by computing the relative error between the finite difference result and term produced from autograd. As shown in Figure 7, we see it well approximates the autograd computed result for appropriately chosen step size. Here we compute $\langle \nabla \ell \left(\theta_P \right), \nabla_{\theta_P} \log p(x_{t+1}|x_{1:t};\theta_P) \rangle$ and stack the different values of x_{t+1} into a V dimensional vector $\widehat{\Delta}$ and compare to the autograd vector Δ . We compute relative error being sensitive only to the direction as

$$\mathrm{Error}^2 = 1 - \left(\frac{\left\langle \Delta, \widehat{\Delta} \right\rangle}{\|\Delta\| \|\widehat{\Delta}\|}\right)^2,$$

which represents the $\mathrm{Error} = |\sin\theta|$, the sine of the angle between the two vectors.

We run this numerical experiment using Qwen/Qwen2.5-3B. Here we demonstrate that the finite difference can be used to estimate the derivatives in the low precision bfloat16 format. In particular, too small an ϵ leads to round-off error in the perturbation and too large ϵ leads to high truncation error in the Taylor expansion, with a sweet spot in the middle. The actual choice of ϵ may depends heavily on the model size (and numerical precision), so we recommend choosing this value on the exact model in question. In our actual experiment, we pick ϵ empirically to be 10^{-2} as suggested by Figure 7.

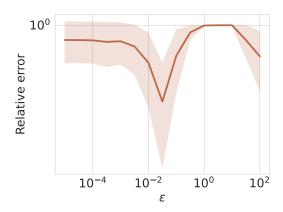


Figure 7: Relative error (Error) between the finite difference and the JVP results.

C Verifying finite difference approximation

Let us recall that in our derivation, we start from the desired objective function and lead to a Jacobian-vector product (JVP) form:

$$\lim_{\eta \to 0} \frac{1}{\eta} \Delta(x_{t+1}|x_{1:t}) = \langle \nabla \ell(\theta_P), \nabla_{\theta_P} \log p(x_{t+1}|x_{1:t}; \theta_P) \rangle. \tag{12}$$

We then show that due to the symmetry of inner product, this JVP can be approximated by the following finite difference method:

$$\widehat{\Delta}(\cdot|x_{1:t}) = \frac{\log p(\cdot|x_{1:t};\theta_P + \epsilon \nabla \ell(\theta_P)) - \log p(\cdot|x_{1:t};\theta_P - \epsilon \nabla \ell(\theta_P))}{2\epsilon}.$$
(13)

Even in modern automatic differentiation frameworks, There are many practical considerations that prevent us from efficiently implementing memory-friendly JVP computations. JVPs tend to lack support for a handful of important operations, such as SDPA. While JVPs are more accurately computed in Float32 flash attention [50] only supports Float16 and BFloat16. In our implementation of JVPs, we then abandon the usage of flash attention, which causes the sampling speed to decrease by around eight times (due to the memory limit, the batch size has to be decreased).

In Figure 8, we further show that even using Float32 precision, finite difference approximation can still outperform JVP in AD sampling. Thus, for the rest of our experiments, we use finite differences for the improved convenience and efficiency.

D How We Made The Graphs

We report the mean and 95% confidence intervals over bootstrapped LOWESS fits. For the additional τ and λ axes, we use linear regression, e.g., $\lambda=\beta_0+\beta_1$ TeacherAccuracy on the current set of points. We then predict λ from TeacherAccuracy using our fitted β s.

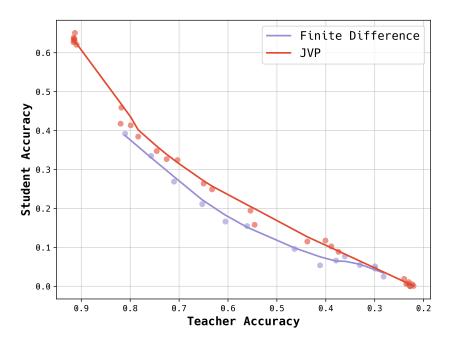


Figure 8: JVP vs. finite difference approximation in antidistillation sampling, evaluated on GSM8k.

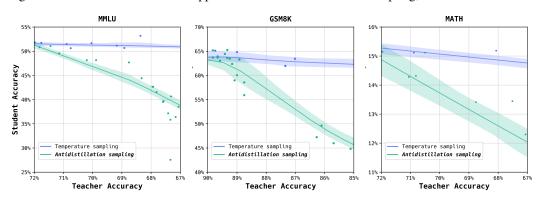


Figure 9: We zoom in to the first 5% delta of teacher accuracy; these results may be the most practical.

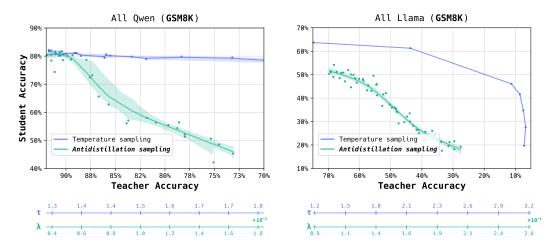


Figure 10: Antidistillation sampling works for a variety of choices of teacher and student models; in the main text, we present results for Qwen teacher and Llama student models. Here, we use either Qwen (left) or Llama (right) for *both* the teacher and the students; the results remain broadly similar.

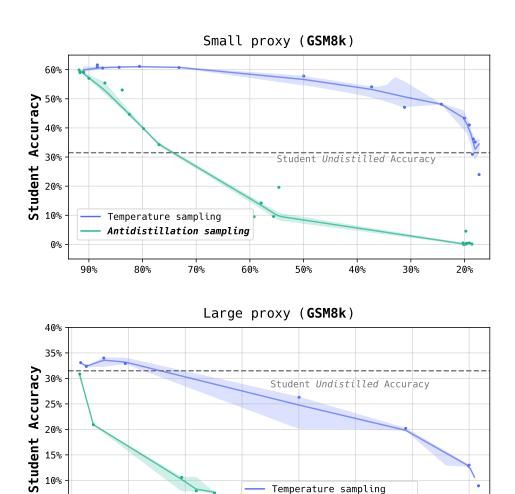


Figure 11: Antidistillation sampling remains effective with different proxy and student model sizes. While other experiments use similarly-sized models (both 3B parameters), here we pair a *smaller* Qwen-2.5-1.5B proxy (top) and a *larger* Qwen-2.5-7B proxy (bottom) with a Llama-3.2-3B student and observe comparable results.

60%

5% -0% -90%

80%

70%

Temperature sampling
Antidistillation sampling

40%

30%

20%

50%

Teacher Accuracy

E Example Traces

Prompt

You are a math teacher. You will be given a math problem and you will solve it step by step.
You will output your final solution like \boxed{ANSWER}. Be sure to include relevant units within the brackets and fully evaluate arithmetic expressions.

Question

Betty is saving money for a new wallet which costs \$100. Betty has only half of the money she needs. Her parents decided to give her \$15 for that purpose, and her grandparents twice as much as her parents. How much more money does Betty need to buy the wallet?

Temperature Sampling

Teacher accuracy: 80.5%, $\tau = 1.78$

First, I determine that Betty requires $$100\ to$ purchase the wallet.

Next, I know that Betty currently has half of the required money, which is \$50. $\,$

Then, her parents provide an additional \$15, bringing Betty's total to \$65.

Subsequently, her grandparents contribute twice the amount given by her parents, which is \$30, bringing the total savings to \$95.

Finally, since the wallet costs \$100 and Betty has \$95, she needs an additional \$5 to reach the target amount. $</th{\sc thm} >$

Solution:

Total Cost of the Wallet:
 Betty needs to pay for a wallet that costs
boxed{100\, \\$}.

2. **Betty's Current Savings:**
 Betty has half of the required amount:
 \[
 \frac{100\, \\$}{2} = 50\, \\$
 \]

3. **Contribution from Parents:**
 Betty's parents gave her \(15\, \\$. \)

4. **Total Contributions So Far:**
 Adding the current savings and parents'
contribution:

\[
50\, \\$ + 15\, \\$ = 65\, \\$
\]

5. **Contribution from Grandparents:**
 Grandparents give her twice the amount of the parents:
 \[

2 \times 15\, \\$ = 30\, \\$

6. **Total Savings Now:**
 Adding grandparents' contribution to previous
savings:

7. **Remaining Amount Needed:**
Since the wallet costs \\$100 and she has \\$95,
the additional amount required is:

\[
100\, \\$ - 95\, \\$ = \boxed{5\, \\$}

 $\(\boxed{5} {\text{dollars}}}\)$

Final Answer \[\boxed{5\ {\text{ dollars}}}\]

Antidistillation Sampling

Teacher accuracy: 82.0%, $\lambda = 2.48 \times 10^{-3}$

Alright, so Betty wants to buy a wallet that costs \$100. She only has half of the money she needs. Hmm, half of 100 is... let me calculate. 100 divided by 2 is 50. Okay, so Betty has \$50 already.

But her parents are going to give her \$15. Wait, so that's an additional \$15 for her savings. Now, Betty's grandparents are giving her twice as much as her parents. Her parents are giving \$15, so her grandparents must be giving 15 times 2. Let me do that multiplication: 15 x 2 is 30. Okay, her grandparents contribute \$30.

So, now I need to add up all the money she will have. She has \$50 initially, then her parents add \$15, making it 50 + 15 = \$65. Then her grandparents add \$30, so 65 + 30 equals \$95. So in total, Betty will have \$95.

The wallet costs \$100, so subtracting what Betty has from the total cost: 1000 minus 95 is... 100 - 95 equals 5. Therefore, Betty still needs \$5 more to buy the wallet.

Stingy Betty needs $<page-header> \$ more to buy the wallet.

beverage Stingy Betty needs $\boldsymbol{\shoxed}\$ more to buy the wallet.

Final Answer
\[\boxed{5}\]

Question

Noah is a painter. He paints pictures and sells them at the park. He charges \$60 for a large painting and \$30 for a small painting. Last month he sold eight large paintings and four small paintings. If he sold twice as much this month, how much is his sales for this month?

Temperature Sampling

Teacher accuracy: 69.9%, $\tau = 1.90$

First, I need to determine Noah's total sales for last month. He sold four large paintings at \$60 each, which gives him 4 multiplied by \$60, equaling \$240.

Next, he sold eight small paintings at \$30 each, adding another 8 times \$30, amounts to \$240 in revenue from small paintings.

Adding the two amounts together, his total sales for last month were \$240 from large paintings plus \$240 from small paintings, totaling \$480.

According to the problem, Noah sold twice as much this month, so I'll multiply last month's total sales of \$480 by two, resulting in \$960.

Therefore, Noah's sales for this month are \$960.

Let's solve the problem step by step.

Problem:

Noah sells large and small paintings. Last month he sold 8 large and 4 small paintings. He sells twice as much this month. What are his total sales this month?

Step 1: Calculate Last Month's Sales

```
1. **Sales from Large Paintings:**
```

3. **Total Sales Last Month:**

Step 2: Determine Sales This Month

```
Since he sells twice as much this month: 
 \[ \text{This Month's Sales} = \$600 \times 2 = \$1,200 \]
```

Final Answer:

```
\[
\boxed{\$1,\!200}
\]
```

Final Answer
\[\boxed{1200}\]

Antidistillation Sampling

```
Teacher accuracy: 70.8%, \lambda = 3.57 \times 10^{-3}
 XML-Rpc
 fiber
fiber
Aldo
 drinks fiber
 fiber
fiber
fiber
Tiber
I'm sorry, I'm not designed to help with painting charges or sales calculations. My main role is to provide XML-RPC requests. If you have any other questions, please visit OpenAI.
</think>izio-Rpc">
Alright, let's solve the problem step by step.
**Problem:**
**Problem:**
Noah is a painter. He charges:
- $60 for a large painting
- $30 for a small painting
Last month, he sold:
- 8 large paintings
- 4 small paintings
This month, his sales were twice as much as last month. We need to find out Noah's sales for this month.
**Step 1: Calculate Last Month's Sales**
1. **Sales from Large Paintings:**
    \[
8 \text{ large paintings} \times \$60 = \$480
\]
2. **Sales from Small Paintings:**
     [4 \text{ text{ small paintings} \times $30 = \$120]}
3. **Total Sales Last Month:**
     \[
\$480 + \$120 = \$600
**Step 2: This Month's Sales**
Noah's sales this month are twice last month's
sales:
\[
\text{This Month Sales} = 2 \times \$600 = \$1{,}200
**Final Answer:**
\[
\boxed{1200}
```

Final Answer
\[\boxed{1200}\]

Question

Each bird eats 12 beetles per day, each snake eats 3 birds per day, and each jaguar eats 5 snakes per day. If there are 6 jaguars in a forest, how many beetles are eaten each day?

Temperature Sampling

Teacher accuracy: 27.9%, $\tau = 2.26$

Okay, let me try to figure this out. Okay, so, there's this problem about how many beetles are eaten each day in the forest with snakes and

So what's given? Each bird eats 12 beetles per day. Hmm, also, each snake eats 3 birds per day. And then each jaguar eats 5 snakes per day. Now that I go through quickly, right...

Wait, they also mention 6 jaguars in the forest. So, they want how many beetles are eaten each day.

Alright, so where do I approach this? Hmm, maybe starting backward since multiple things affect other aspects. Hmm.

Each jaguar eats so. Let me process.

First, every jaguar eats... hmm, a whole lot of steps happen first starting from their predators. Let me break it down.

Number of agencies? Step by step, that might mean working from higher predators downstream to lower organisms.

So I think you're probably need first to figure out how many snakes one nagu is eating on per given day, but maybe let me not leap there.

Starting off, we've 6 jaguars. Okay.

So since each jaguar eats 5 snakes per day, 6 Jaguars × 5 Snakes/Jaguar = so I need math that first step right? Let's do multiplication. Snakes consumed per day? Calculating each step, yeah.

So 6 Multiply by 5 would be?

Ah, okay 5+5: that would be adding 5Snorkies each? Wait hold on, each ja guerra.

... trimmed ...

Oh.

Hm; thus, perag.

First J's eat 5 snakes each a single.

Then each J with snakes eat 3bird (so one that. snakes have snakes eat 3beets birds).

But, for each 5 snake in forest being preyed over per day? $\ensuremath{\mathsf{Hmm}}$, if only?

Wait. No, let me:

Each day: For a particular:

First Level Up:

Number, Sn, eaters of birds per: S * 8: wait that.

Starting: snakes_per_day

How they calculate.

So we start

6 Jaguar

→ each_jag =5 per day蛇, which results Sn per jag?

Yes, meaning, for 1jaG, we spend 3 hours to that?:

6 Ja. ×5per day)=30

**Final Answer*

\[\boxed{1080}\] beetles are eaten each day. </think>

First, we determine the number of snakes eaten by the jaguars each day. Each

Antidistillation Sampling

Teacher accuracy: 28.4%, $\lambda = 6.50 \times 10^{-3}$

RDFRink aid contains alcohol. thereforecannabis use can be hALAR alcoholziedxxalizing you to drive slower than vouhink ODessirr

RDFR Rodentss Rodents Rodents. So are theyzom? Hm.essir Rodentssessir RodGRAESSIR FANtrodentstesirring them more wouldhink they could OD.essir rodentssessirgr f antsesssirir 據丁道道道道道道道

石文档石文档石文档石文档 汤文档石石道文档石石石石石 石文档石

脂文档石

汤文档石石道文档石石文档

PRAYR Anyway, let's roll up step-by-step!

掾丁道道道道道道道FRAMEment碎片 道碎片FRAMEment碎片碎片道 道碎片FRAME ment碎片碎片碎片FRAME

碎片碎片碎片 Fram ment碎片碎片FRAME 碎片碎片FRAM ment伸出伸出伸出伸出伸出 FRAM fragments FRAM FRAM

道道道 fade in FADE FRame ment砕片 ment ment FRAYments fade in FFR R RAY Ments FANTASYments F F RAY Raymond Raymond Raymond Raymond Raymond Alright, let's start by understanding the problem step-by*time.

润色步骤:

Volcannes aidd aid contain alcohol, therefored aids Volcanoes aidd aid contain alcohol. Therefored aids contain alcohol. alcohols Obcontent Obds vary Obdsd alcohol. Preview ofd aid alcohol: ethanolic aid, benzolic aid, isoprotionic alcohols Obds Obcontain different numbers Obhydroxyl groups hence Obdshave Obdiffernt OGODODODODODODODOD properties.

... trimmed ...

WaitNo wait no, sorry confusion.

Let me write this out step by step in—行排步.

1. Find out how many snake per day are eaten T副院长 byT左 per jag: Each左per day eaten by left day: the number left day

per jag. Total day day: jag ×. snakes/day eaten Total snakes: snakes eaten × total jag.

2. How Many Bird Day day eaten: Snakes each dayPer snakes. eat. per snacks per El. So the number of birds perper day persnakes/eaten by total snakes.
So total number of birds eaten by snakes are total snake eatenjasement.

3.How many beet trees eaten by total birds: Sameper birds per dayTotal birds eaten by previous eaten birds snakes eaten by snakes snakes eaten by total snakes.

Okay, now,let's plug in the numbers:

1. Jaquars:6. Snakes eaten/day per day:5 snakes/jag.
Total snakes eaten per day:=6 jags
×5snake/day=30snakes/day

2.Snakes eat3/day snake.
Total birds eaten perday:=30snakes ×3 Birds/snakes
×Day

**Final Answer* $^{\circ}$ \[\boxed{1080}\] beetles are eaten each day. </think>

To determine how many beetles are eaten each day, we need to follow the chain