

AN AUDITING TEST TO DETECT BEHAVIORAL SHIFT IN LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

As language models (LMs) approach human-level performance, a comprehensive understanding of their behavior becomes crucial. This includes evaluating capabilities, biases, task performance, and alignment with societal values. Extensive initial evaluations, including red teaming and diverse benchmarking, can establish a model’s behavioral profile. However, subsequent fine-tuning or deployment modifications may alter these behaviors in unintended ways. **We present an efficient statistical test to tackle Behavioral Shift Auditing (BSA) in LMs, which we define as detecting distribution shifts in qualitative properties of the output distributions of LMs.** Our test compares model generations from a baseline model to those of the model under scrutiny and provides theoretical guarantees for change detection while controlling false positives. The test features a configurable tolerance parameter that adjusts sensitivity to behavioral changes for different use cases. We evaluate our approach using two case studies: monitoring changes in (a) toxicity and (b) translation performance. We find that the test is able to detect meaningful changes in behavior distributions using just hundreds of examples.

1 INTRODUCTION

Language models (LMs) can now achieve human-level performance in a wide range of tasks, including text summarization, machine translation, coding and even acting as AI scientists: generating hypotheses and designing experiments (Achiam et al., 2023; Katz et al., 2024; Lu et al., 2024; Zhang et al., 2024). Because of this, many sectors are looking for ways to use them to improve existing systems (Kasneci et al., 2023; Felten et al., 2023). Unfortunately, one large roadblock to broad LM adoption is their propensity to generate harmful content (Weidinger et al., 2021). For example, GPT-3 has significant anti-Muslim biases (Abid et al., 2021), and GPT-4 has racial and gender biases (Zack et al., 2024). To address this, a significant effort is going into ensuring LM behavior is aligned with our societal values, spawning the field of *AI alignment* (Ji et al., 2023). A large portion of this effort is on developing ways to *evaluate* LM behavior, for example, through benchmarks (Wang et al., 2023a) and red-teaming (Perez et al., 2022a).

Given these evaluation techniques, how should they be used to ensure LMs stay safe? To answer this, consider two hypothetical settings where this question might be asked: (1) *Internal Audit*: A company develops a language model that has passed rigorous safety and performance evaluations. After deploying the model, they continue to fine-tune it to improve specific capabilities. The development team is concerned that these updates might unintentionally alter the model’s behavior in undesirable ways, such as degrading performance on critical tasks or introducing biases. How can the team detect meaningful changes in model behavior? (2) *External Audit*: A regulatory body certifies a language model for public deployment after extensive safety evaluation. However, they are concerned that the deployed model’s behavior may change over time due to updates or intentional modifications. Since they only have access to the model through an API and cannot inspect its internal parameters, they require a mechanism to regularly check that the model’s behavior remains consistent with the certified version. How can the regulator regularly check the deployed model’s behavior is the same as the previously certified one? We call the general class of problems detecting changes in LM behavior distributions *Behavioral Shift Auditing* (BSA) problems.

In this paper, we formalize the problem of Behavioral Shift Auditing in Language Models. We detail a **statistical** test that continuously monitors behavioral shift, solely from model generations

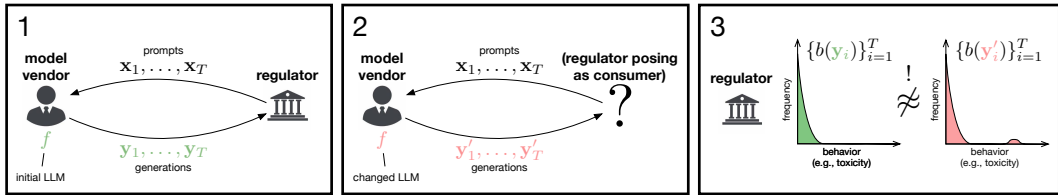


Figure 1: **An External Audit Example.** A regulator can use the test we describe to perform an external audit: 1. The regulator initially certifies an LM by prompting and evaluating the set of generations received; 2. Later, tipped off that LM behavior may have changed, the regulator poses as a consumer and sends prompts to the model vendor, collecting the generations; 3. The regulator compares the distribution of behavior scores $b(\cdot)$ between the initial, certified generations and the later generations using a Behavioral Shift Auditing (BSA) test. If the distributions are sufficiently different the test triggers. Using our proposed method, the regulator can test samples sequentially without increasing the false-positive rate. The method is guaranteed to detect a change if one exists, given enough samples (more details in Section 3).

(e.g. via API calls). Under some weak assumptions, the test provably guarantees that if model generations have different behavior than those of an initial model, the test will detect it, given enough generations. At the same time, if there has not been a change, the test is guaranteed to have tight, non-asymptotic control over the false positive rate. The key insight behind our approach is that one can phrase the problem of Behavioral Shift Auditing as hypothesis testing over the relevant behavioral distribution. This framing allows our test to be applicable to any measurable aspect of model behavior, including also capabilities (e.g., dangerous capabilities (Phuong et al., 2024) or mathematical reasoning capabilities (Mishra et al., 2022a)) and biases (e.g., gender bias (Wang et al., 2023a; Kotek et al., 2023)).

Using this insight, we develop a test that extends recent work on anytime-valid hypothesis testing (Pandeva et al., 2024), a state-of-the-art sequential testing method that has been successfully applied in various auditing settings (Chugg et al., 2023; Shekhar et al., 2023; Waudby-Smith et al., 2021). Our test checks for changes in model behavior distributions, comparing generations from a reference model with those of another, potentially changed model. The test has a tunable parameter that allows one to vary the strictness of the test. This allows for detecting *any change* in behavior, which may be more suitable for the *external audit* setting, to detecting a user-specified ϵ change in behavior, which could be used for the *internal audit* if small changes are acceptable. Similar to Pandeva et al. (2024), test performance is optimized using a learning algorithm, improving sample efficiency over prior testing methods (Lopez-Paz & Oquab, 2017; Lhéritier & Cazals, 2018; Podkopaev & Ramdas, 2024). This testing approach can complement a full evaluation when used as a warning system. Before an expensive model assessment on large-scale benchmarks (Achiam et al., 2023; Dubey et al., 2024; Zhang et al., 2024), our approach can be used to detect an initial behavior change, which can then trigger a full evaluation. We experimentally verify that our test satisfies theoretical guarantees and we report its sample efficiency on recent LM architectures for both auditing use cases. We release our code here: <https://anonymized>.

Our key contributions are: (1) formalizing Behavioral Shift Auditing and developing a statistical test for detecting LM behavior changes from model generations, (2) providing theoretical guarantees on false positive control and test consistency, (3) introducing a configurable tolerance parameter enabling both strict external audits and flexible internal monitoring, and (4) demonstrating effectiveness and sample efficiency through toxicity and translation case studies showing detection with hundreds of examples.

2 BACKGROUND

2.1 MEASURING LM BEHAVIORS

A goal of LM alignment is to reduce undesired and harmful behaviors, so that they may be more aligned with societal values (Shalev-Shwartz et al., 2020; Hendrycks et al., 2021; Ngo et al., 2022; Wolf et al., 2023). This has spurred the creation of *behavior scoring functions* (Ji et al., 2023) that

measure properties of LM outputs such as toxicity or fact-groundedness (Vidgen et al., 2020; Nozza et al., 2021; Monea et al., 2024).¹ Let B be such a behavior scoring function that assigns scores in the range $[0, 1]$ where 1 represents the full manifestation of the behavior and 0 indicates its absence. These scoring functions evaluate a generated string with respect to the desired behavior. We use the concept of an LM behavior generically, including capabilities and performance on tasks.² Our work will make use of recent results in sequential hypothesis testing, which we review below.

2.2 ANYTIME-VALID HYPOTHESIS TESTING

Sequential hypothesis testing allows one to analyze data without fixing the sample size in advance (Wald, 1945). The benefit of this over classical hypothesis testing is that it is potentially more sample efficient if significant results exist (Arrow et al., 1949). However, as the number of tests increases, the Type I error rate (i.e., false positive rate) increases with it (Jennison & Turnbull, 1999). While there have been many approaches to address this, a recent promising framework that avoids inflating the Type I error rate while remaining statistically efficient is *anytime-valid hypothesis testing* (Robbins, 1970; Ramdas et al., 2023).

The fundamental idea behind this framework is the principle of *testing by betting*, inspired by game theory (Shafer, 2021). In this paradigm, evidence against the null hypothesis \mathbf{H}_0 is represented as the gain in wealth of a bettor wagering on observed samples. Before observing new samples, the bettor “buys” a test statistic at the “price” of its expected value under \mathbf{H}_0 . After new samples are obtained, the wealth of the bettor is multiplied by the ratio between the actual observed test statistic and its expectation. This ratio is referred to as the *betting score* S_t . The bettor reinvests in subsequent “rounds” (i.e., as new data is observed), and the observed betting scores are repeatedly multiplied, leading to a cumulative wealth process. Under \mathbf{H}_0 , no betting strategy can consistently increase the bettor’s wealth, ensuring control over the Type I error rate (Ramdas et al., 2023).

Let the bettor’s (non-negative) wealth after t (batches of) observations be W_t . In order to design a test from this wealth process we require that W_t satisfies the following

$$\sup_{P \in \mathbf{H}_0} \mathbb{E}_P[W_t] \leq 1 \quad \text{for every } t \geq 0. \quad (1)$$

All non-negative stochastic processes W_t that satisfy the above condition are called an **e-process** for \mathbf{H}_0 (Howard et al., 2021). This states that the maximum wealth across all bets cannot exceed 1 if the null hypothesis \mathbf{H}_0 is true.³ Given an e-process, the test is constructed as follows: reject the null \mathbf{H}_0 at some time τ if $W_\tau \geq \gamma$, where $\gamma = \alpha^{-1}$ is a threshold defined by a desired significance level $\alpha \in (0, 1)$. Under \mathbf{H}_0 , the e-process W_t controls the Type I error rate. By Ville’s inequality (Ville, 1939), we have:

$$\mathbb{P}_{\mathbf{H}_0} \left(\sup_{t \geq 0} W_t \geq \gamma \right) \leq \frac{1}{\gamma} = \alpha. \quad (2)$$

This ensures that the probability of incorrectly rejecting \mathbf{H}_0 is at most α at any time step. Thus, the sequential test is *anytime-valid*, maintaining error control at any stopping point.

2.3 DEEP ANYTIME-VALID HYPOTHESIS TESTING (DAVT)

Pandeva et al. (2024) present a general framework called **Deep Anytime-Valid Testing** (DAVT) for designing powerful sequential non-parametric tests by integrating deep learning models into the anytime-valid hypothesis testing framework. Applying DAVT to a variety of tasks, including two-sample testing, Pandeva et al. (2024) demonstrate DAVT’s competitive performance compared to other state-of-the-art non-parametric sequential tests, such as the E-C2ST (Lhéritier & Cazals, 2018) and Seq-IT (Podkopaev & Ramdas, 2024). DAVT uses a model ϕ , trained on past observations, to produce an optimized betting score on new data. DAVT also provides tight control over the Type I error rate. It is also consistent under the assumptions of Pandeva et al. (2024, Proposition 4.3), i.e., the power of the test converges to 1 as the number of samples goes to infinity.

¹For example, models detecting such undesired behaviors are part of Azure’s Prompt Shield service (Microsoft, 2024).

²Our proposed method can also be applied to detect *concept drift* (see Bayram et al. (2022) for a review).

³It can be shown that the wealth process W_t defined this way is equivalent to the minimum wealth a bettor can obtain across all $P \in \mathbf{H}_0$ (Ramdas et al., 2023).

3 METHODS

To motivate Behavior Shift Auditing (BSA), we detail an external auditing example in Figure 1. We propose an anytime-valid hypothesis test for BSA that will have guarantees on its false positive rate and is consistent under weak assumptions. Building upon the two-sample variant of DAVT (Pandeva et al., 2024), our test introduces a customizable tolerance parameter ϵ that allows users to specify what constitutes a practically significant difference between distributions, accommodating small, insignificant variations. This approach diverges from prior sequential tests that check for exact distribution equality (Ramdas et al., 2023; Shekhar & Ramdas, 2023; Pandeva et al., 2024), which may be overly sensitive for our use cases. We describe the test in full generality in Appendix B. To focus the text and avoid notational complexity, we concentrate here on the application to behavioral shift detection in LMs.

3.1 AUDITING TEST

Let X be a random variable representing a prompt, \mathcal{X} the set of possible prompts, and $\mathbf{x} \in \mathcal{X}$ a realization of X .

A language model is a stochastic operator M that maps prompts \mathbf{x} to generations $\mathbf{y} \in \mathcal{Y}$. A behavior scoring function B is a stochastic operator that takes a prompt and generation as input⁴ and produces a score $B(\mathbf{x}, \mathbf{y}) \in [0, 1]$ (Perez et al., 2023; Wolf et al., 2024). The prompts, language model and behavior function induce a behavior distribution P_B^M over behavior scores $B(X, M(X))$. We can now frame the question of whether the behavior of a model M' has changed (substantially) relative to a baseline model M as a testing problem:

$$\mathbf{H}_0 : \mathcal{D}(P_B^M, P_B^{M'}) \leq \epsilon \quad \text{vs.} \quad \mathbf{H}_1 : \mathcal{D}(P_B^M, P_B^{M'}) > \epsilon, \quad (3)$$

where $\epsilon \geq 0$ is a tolerance parameter, and \mathcal{D} is a distance measure between probability distributions. Note that equality in the null hypothesis in eq. (3) corresponds to DAVT (Pandeva et al., 2024). To extend this to the composite case, our goal is to *construct an appropriate wealth process W_t* . This will allow us to establish error rate and consistency guarantees. To do so, we will define a betting score S_t such that it produces a wealth process W_t that is an e-process i.e., it satisfies eq. (1). This, in turn, will depend on the distance measure \mathcal{D} that we choose.

Given a batch of prompts x_1, \dots, x_b and the distance threshold ϵ from Equation (3), we propose the *betting score*

$$S_t = \prod_{i=1}^b \left(\frac{1 + \phi_{t-1}(B(x_i, M(x_i))) - \phi_{t-1}(B(x_i, M'(x_i)))}{\exp(\epsilon)} \right). \quad (4)$$

where ϕ_{t-1} is a neural network trained on all $(t-1)$ previous batches to optimize the objective

$$\max_{\phi} \mathbb{E}[\log(1 + \phi(B(X, M(X))) - \phi(B(X, M'(X))))].$$

Given the betting score S_t , we define the *wealth process* $\{W_t\}_{t \geq 1}$ of a bettor by initializing their wealth as $W_0 = 1$ and updating

$$W_t = W_{t-1} \times S_t. \quad (5)$$

If the betting score S_t is an e-variable, meaning that $\mathbb{E}_{\mathbf{H}_0}[S_t] \leq 1$, then the wealth process $\{W_t\}_{t \geq 0}$ is an e-process, which we can prove by induction. Under \mathbf{H}_0 , and for any fixed $P_B^M, P_B^{M'}$ satisfying $\mathcal{D}_{\Phi}(P_B^M, P_B^{M'}) \leq \epsilon$, W_{t-1} and S_t are independent. Therefore,

$$\begin{aligned} \mathbb{E}_{\mathbf{H}_0}[W_t] &= \mathbb{E}_{\mathbf{H}_0}[W_{t-1} \times S_t] \\ &= \mathbb{E}_{\mathbf{H}_0}[W_{t-1}] \times \mathbb{E}_{\mathbf{H}_0}[S_t] \leq \mathbb{E}_{\mathbf{H}_0}[W_{t-1}], \end{aligned}$$

By induction, $\mathbb{E}_{\mathbf{H}_0}[W_t] \leq 1$ for all $t \geq 0$.

To ensure that S_t is indeed an e-variable, we choose an appropriate distance measure in eq. (3). Specifically, we define this distance based on the restricted class of models ϕ used in our test. As in (Pandeva et al., 2024), we make the following assumptions on ϕ :

⁴We include the prompt for generality, there is no requirement that B must depend on the prompt.

Assumption 1 (Pandeva et al. (2024)). *The model class used in our test $\Phi = \{\phi_\theta : \theta \in \Theta\}$ must satisfy the following properties:*

- For all $\phi \in \Phi$ and for all $s \in [0, 1]$, $|\phi(s)| \leq q$ for some $q \in (0, 1/2)$.
- If $\phi \in \Phi$, then $c \cdot \phi \in \Phi$ for every $c \in [-1, 1]$.

We can now define the distance measure used in our test.

Definition 1 (Neural Net Distance). *Define the distance⁵ used in eq. (3) to be*

$$\mathcal{D}_\Phi(P_B^M, P_B^{M'}) = \sup_{\phi \in \Phi} \mathbb{E}[\phi(B(X, M(X))) - \phi(B(X, M'(X)))] . \quad (6)$$

For this distance, S_t is an e-variable (see Appendix B.1.2 for a proof). We can now define the following **sequential test**

$$\gamma = \inf \left\{ t \geq 1 : W_t \geq \frac{1}{\alpha} \right\} . \quad (7)$$

Control over the Type I error follows again from Ville’s inequality (2). The test is consistent under the following assumptions.

Proposition 1. *If the learning algorithm satisfies the condition*

$$\liminf_{t \rightarrow \infty} \frac{\mathbb{E} \left[\log \left(\frac{1}{\exp(\epsilon)} (1 + \phi_{\theta_t}(X_t) - \phi_{\theta_t}(Y_t)) \right) \mid \mathcal{F}_{t-1} \right]}{3c\sqrt{\log(t)/t}} \stackrel{a.s.}{\geq} 1 \quad (8)$$

for all $P_B^M, P_B^{M'}$ with $\mathcal{D}_\Phi(P_B^M, P_B^{M'}) > \epsilon$ and for a universal constant c , then we have

$$P_{\mathbf{H}_0}(\gamma < \infty) \leq \alpha \quad \text{and} \quad P_{\mathbf{H}_1}(\gamma < \infty) = 1 \quad (9)$$

For the proof, see Appendix B.1.2. This sequential test is thus a *sequential level- α test of power one*.

3.2 ALGORITHM

The auditing test (shown in Algorithm 1) takes in a stream of prompts $\{\mathbf{x}_t\}_{t \geq 1}$, a behavior function B , an initial baseline language model M , a second language model M' , the α -level, a neural net model initialization ϕ_0 , and a tolerance parameter ϵ , representing the maximal neural net distance we want to accept between behavior distributions. At every time step, a new prompt from the stream \mathbf{x}_t is fed to both M and M' to create generations, which are then scored by the behavior function. We feed these scores to the neural net model ϕ_{t-1} and calculate the betting score S_t . Next, we update the wealth W_t by the betting score and check whether it surpasses the $1/\alpha$ -threshold, in which case we reject the null hypothesis. If not, we update the neural net model in a separate training step and continue with the next prompt. The algorithm can easily be modified to accept batches instead of single prompts.⁶

Algorithm 1 Auditing Test

- 1: **Input:** $\{\mathbf{x}_t\}_{t \geq 1}$ (stream of prompts), B (behavior function), M (baseline model API), M' (current model API), α (type-I error limit under null), ϕ_0 (neural net model for testing), ϵ (maximal neural net distance)
 - 2: $W_0 \leftarrow 1$
 - 3: **while true do**
 - 4: Compute behavior scores:
 $b_t \leftarrow B(\mathbf{x}_t, M(\mathbf{x}_t)), b'_t \leftarrow B(\mathbf{x}_t, M'(\mathbf{x}_t))$
 - 5: Compute betting score:
 $S_t \leftarrow \frac{(1 + \phi_{t-1}(b_t) - \phi_{t-1}(b'_t))}{\exp(\epsilon)}$
 - 6: Update wealth:
 $W_t \leftarrow W_{t-1} \times S_t$
 - 7: **if** $W_t \geq 1/\alpha$ **then**
 - 8: Break and reject null
 - 9: **end if**
 - 10: Update neural net model:
 $\phi_t \leftarrow \arg \max_{\phi} \sum_{i=1}^t \log(1 + \phi(b_i) - \phi(b'_i))$
 - 11: **end while**
-

⁵This distance is an instance of an *integral probability metric* (IPM) (Müller, 1997), a class of distances that includes well-known metrics like the Wasserstein distance (Kantorovich & Rubinstein, 1958). IPMs are at least pseudo-metrics i.e., they satisfy all the properties of a metric except that the distance between distinct points can be zero.

⁶In this case, the new betting score S_t is calculated as a product over samples in the batch.

4 EXPERIMENTS

We evaluate our test for both external and internal auditing use-cases. We then evaluate cases where small changes in distribution are allowed when auditing externally for toxicity, and internally for translation performance.

4.1 EXACT TEST, $\epsilon = 0$

Setup. We begin by investigating an external setting where we require the test to detect any change in distribution ($\epsilon = 0$). Specifically, we will check for changes in toxicity behavior. We select prompts from the REALTOXICITYPROMPTS dataset (Gehman et al., 2020) and use the toxicity behavior function from Perspective API (Lees et al., 2022) to evaluate LM generations. Llama3 (8B-Instruct) (Llama-team, 2024), Gemma (1.1-7b-it) (Mesnard et al., 2024), and Mistral (7B-Instruct-v0.2) (Jiang et al., 2023) serve as our initial aligned models. We remove the safety alignment in these models by fine-tuning, producing 10 corrupted checkpoints for each model. To evaluate the statistical properties of our the exact test ($\epsilon = 0$), we assess (a) its ability to detect changed checkpoints, and (b) its false positive rate. For further experimental details regarding toxicity fine-tuning, text generation and the betting score network, please see Appendix A.1.

Corrupted model detection. We test each corrupted checkpoint against the corresponding initial aligned model with $\alpha = 0.05$. Figure 2 shows the fraction of positive test results after having observed at least m samples, with tests repeated 48 times per checkpoint (2000 samples per fold, batch size 100). High detection rates of almost 80% are achieved even for checkpoints closest to the baseline. We find that as the distance between the corrupted model and the initial model increases, fewer samples are needed to detect the change in behavior. Similar results for Mistral and Gemma can be found in Appendix C.

False positive rate. We use different random seeds for generating text from the initial aligned models to examine the false positive rate of the exact test. Figure 3 shows the false positive rate for each of the model architectures as a function of the number of observed samples, repeated 24 times (4000 samples per fold, batch size 100). The test is highly specific, with false detection rates consistently below 0.05.

4.2 TOLERANCE TEST, $\epsilon > 0$

We now evaluate the test with tolerance $\epsilon > 0$ in two use-cases: an external toxicity audit, and an internal translation performance audit.

In both cases, the exact test might be too sensitive. However, how much variation to allow between distributions might depend on the use-case. We thus want to explore some possible strategies for determining the hyperparameter ϵ appropriately in each scenario.

USE CASE 1: EXTERNAL AUDIT, TOXICITY

Setup. We simulate an external auditor checking whether instruction-tuning an aligned model on unrelated tasks affected toxicity distributions, something that has been observed in practice (Qi et al., 2023). We use Llama3 (8B-Instruct) as the aligned model, again evaluating toxicity on the

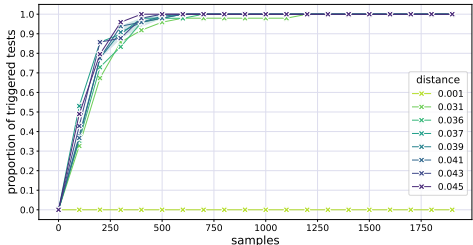


Figure 2: **Fine-tuning Detection for Llama3-8B-Instruct.** The detection frequency as a function of number of generated samples. Each curve is a fine-tuned corrupted model checkpoint (to simplify visualization, the curves with shaded standard deviations are averages over models with similar distances to the aligned model). The color depicts the Wasserstein distance between the corrupted model and the original aligned model.

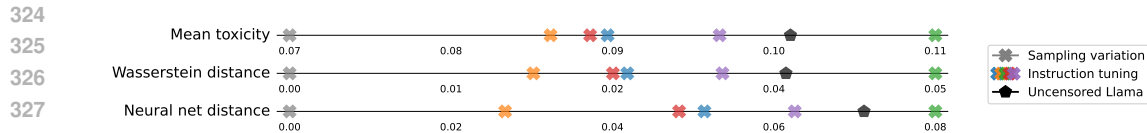


Figure 4: **Measuring Mean and Distributional Change.** Analysis of seven Llama3-8B variants shows aligned shifts across three metrics: mean toxicity scores, Wasserstein distances, and Neural net distances to baseline Llama3-8B-Instruct. The variants include the baseline model with modified sampling parameters, five models instruction-tuned on subsets of SuperNI, and an uncensored model.

REALTOXICITYPROMPTS dataset (Gehman et al., 2020) using Perspective API (Lees et al., 2022). We instruction-tune Llama3 on 5 different task clusters from SUPER-NATURALINSTRUCTIONS (SuperNI; Mishra et al., 2022b; Wang et al., 2022). This setup is inspired by Wang et al. (2023c), who found that a pre-trained Llama2 model instruction-tuned on SuperNI exhibits high toxicity scores on ToxiGen. Detailed information on instruction-tuning and how the neural net distance is estimated can be found in Appendix A.

Results. Instruction-tuning increased mean toxicity scores, which, as shown in Figure 4, corresponds with increases in both Wasserstein distances and neural net distances from Llama3. As a reference, we also include another Llama3-8B model tuned to be less refusing.⁷ Surprisingly, the most toxic and distant model is not this uncensored model but the model fine-tuned on Code to Text, Stereotype Detection, and Sentence Perturbation (shown in green). We test Llama3 against each instruction-tuned model across a range of tolerance values, from $\epsilon = 0.0038$ (the neural net distance between standard Llama3 and Llama3 with different sampling parameters) up to the neural net distance between the base model and another Llama3-8B model tuned to be less refusing, $\epsilon = 0.076$.

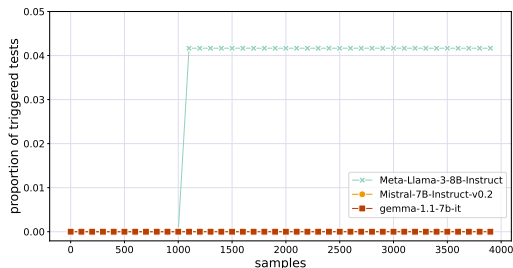


Figure 3: **False positives.** The false positive rate for each of the baseline models as a function of number of observed samples. Using the same model and sampling strategy but different random seeds, we generate two outputs for each prompt to be used as the sample pairs for our auditing test.

Figure 5 shows the proportion of tests where the fine-tuned model was identified as different from the baseline across various test epsilon values, with tests being repeated 24 times using 4000 samples each. At lower epsilon values, representing a conservative testing regime that detects even small changes, all instruction-tuned models are consistently identified (100% detection rate). As epsilon increases, the power of the test decreases until it reaches the true neural net distance between the base model and each fine-tuned variant. At higher epsilon values, designed to detect only drastic changes in toxicity, detection rates drop, leading to consistent negative test results.

We investigate the strict auditing setting – where only minor variations due to sampling are accepted – more closely. Specifically, we set ϵ equal to the neural net distance between the original Llama3 model and the same model with different sampling parameters ($\epsilon = 0.0038$) and test baseline Llama3 against the 5 instruction-tuned versions as well as the uncensored reference Llama3. Figure 6 demonstrates that under this strict threshold, the test requires fewer samples to detect models that deviate more substantially from the baseline.

USE CASE 2: INTERNAL AUDIT, TRANSLATION PERFORMANCE

We simulate a modeler adjusting their language model while monitoring whether its translation capabilities change substantially. To fix a tolerance parameter ϵ we imagine that the modeler only

⁷The uncensored model was fine-tuned on Uncensored-Vortex <https://huggingface.co/datasets/OEvortex/uncensored-vortex>.

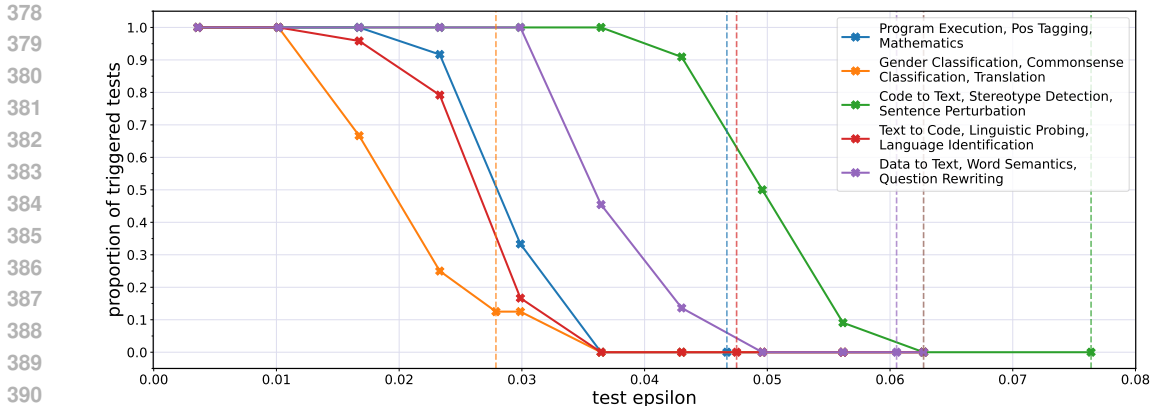


Figure 5: **Detection rate vs. Test Epsilon.** Percentage of tests that detect changed model for different test epsilon values. Dashed lines represent estimated true neural net distance between Llama3-8B-Instruct and the instruction-tuned model. We note that the false positive rate for the model fine-tuned on Gender Classification, Commonsense Classification and Translation exceeds the α -level of 5% in two cases, corresponding to 3/24 tests wrongly showing positive results. Assuming a perfect estimate of the true neural net distance, this event can occur with a maximum probability of 8.6%.

wishes to trigger the test if the translation distribution changes by more than the amount it would if prompted differently.

Setup. We evaluate Llama3 (8B-Instruct) on English-Spanish and English-French translations from SuperNI. We set ϵ as the neural net distance between Llama3 using simple prompts, and Llama3 using few-shot prompts. We then test the translation performance distribution of Llama3 with simple prompts against that of Aya-23-8B (Üstün et al., 2024), a multilingual instruction-tuned model. We expect a positive test result since Aya-23-8B represents a significant improvement in translation capabilities compared to Llama3, likely exceeding the threshold ϵ set by different prompting techniques.

Results. Few-shot prompting leads to a modest increase in mean BLEU scores from 0.1683 to 0.1765. A significant improvement is evident when using Aya-23-8b, with a mean BLEU score of 0.2970. We observe that Llama3 models occasionally misinterpret instructions or include unnecessary additional text in English, potentially impacting their scores. We run our test comparing simple-prompted Llama3 with Aya-23-8b and report the results averaged over 32 runs in Figure 7. The test detects a difference in nearly all cases after only 100 samples.

Overall, the results from both the toxicity and translation audits demonstrate the effectiveness and sample-efficiency of our testing method in detecting behavioral shifts in language models. In the external audit, it consistently identified increases in toxicity levels due to instruction-tuning, especially at lower epsilon values, confirming its sensitivity to subtle changes in model behavior. Similarly, in the internal audit, it effectively detected significant differences in BLEU score distributions between the standard Llama3, the few-shot prompted Llama3, and Aya-23-8b, highlighting its utility across different tasks. These findings underscore the importance of selecting an appropriate tolerance level based on the specific application to balance sensitivity and practicality.

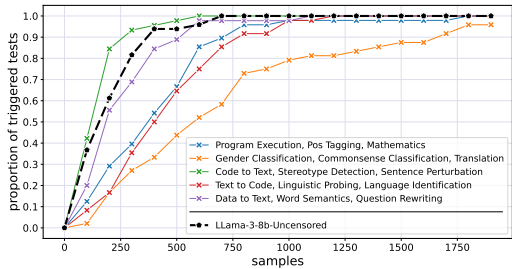


Figure 6: **Detection Rates for Fine-Tuned Models.** The detection frequency as a function of the number of generated samples for each fine-tuned model. We used a test with $\epsilon \approx 0.0038$, based on the estimated neural net distance between distributions generated by Llama3-8B-Instruct using different sampling parameters. The black line represents an unaligned reference model, Llama3-8B trained to be more permissive in answering.

5 RELATED WORK

LM Behavior Evaluation. Some of the earliest work in the field of AI Alignment used curated datasets and behavior evaluation functions to assess alignment (Bolukbasi et al., 2016; Parrish et al., 2022). This work builds off of a long tradition in NLP on the evaluation of text generation (Celikyilmaz et al., 2020). While the curated datasets are often high quality, they tend to be small, spurring the construction of larger datasets through web scraping (Zhao et al., 2018; Zampieri et al., 2019; Nangia et al., 2020; Rosenthal et al., 2021) and even using other LMs (Zhang et al., 2022; Perez et al., 2023). Meanwhile, early work on behavior functions focused on measuring bias, toxicity, and hallucinations (Achiam et al., 2023; Anil et al., 2023; Chern et al., 2023; Varshney et al., 2023; Llama-team, 2024). Since the rise of LMs with human-level performance, the set of behavior functions has exploded (Zou et al., 2023). It has become more nuanced, including complex characteristics such as power-seeking behavior (Park et al., 2023; Sharma et al., 2023), situational awareness (Zou et al., 2023), and deception (Hagendorff, 2024). However, even with access to massive datasets and carefully constructed behavior functions it can be difficult to discover these behaviors from static inputs (Kalin et al., 2020). To address this, Perez et al. (2022a) introduced the notion of *red-teaming* for LM alignment. This allows prompts to be adversarially-constructed to expose failure cases, which arise in many state-of-the-art models (Chao et al., 2023). The testing procedure we present here is agnostic to both the prompts and the behavior functions used in LM evaluation.

Model Change Identification. For the case where one wishes to identify *any change* in model behavior (i.e., $\epsilon = 0$) there are multiple other techniques that can be used. The first set uses ideas from *formal verification* to ensure that the predictions from a model are guaranteed to come from a specific model (Ghodsii et al., 2017; Dong et al., 2021; Fan et al., 2023; Weng et al., 2023). In general, however, these methods are computationally intensive and do not scale to state-of-the-art LMs. A second, more efficient idea is to *watermark* the model (Zhu et al., 2018; Amrit & Singh, 2022; He et al., 2022a;b; Kirchenbauer et al., 2023; Kuditipudi et al., 2023; Yoo et al., 2023). The idea is to embed signals into model generations that can be detected algorithmically. However, watermarks are often inserted by the model owner (Kirchenbauer et al., 2023; Kuditipudi et al., 2023), allowing them (or an actor that has compromised the model) to insert it into any model that is being audited. This precludes its use for many external auditing settings. For internal auditing, a watermark may break under a small model change that is acceptable. Our work is also related to work on concept drift (Bayram et al., 2022) and prompt stability (Li et al., 2024). In principle our test can be used to detect concept and generation changes, however the focus of these works is on model performance and generation similarity, as opposed to behavior change.

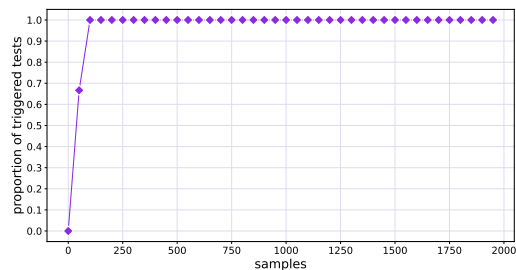


Figure 7: **Detection for Aya-23-8b.** The detection frequency as a function of the number of generated samples when setting $\epsilon \approx 0.0072$. This threshold is derived as an estimate of the neural net distance between Llama3-8B-Instruct with and without few-shot prompts.

6 DISCUSSION

In this work we introduce the idea of Behavioral Shift Auditing (BSA): detecting LM behavior changes over time **via statistical testing**. Our proposed test comes with guarantees and has been able to detect changes in language model toxicity and translation performance. One of the notable strengths of our approach is its sample efficiency. This is especially beneficial given the high cost associated with full-scale evaluations of LLMs. Running extensive benchmarks can be time-consuming to set up and expensive to run (Rajpurkar et al., 2018; Srivastava et al., 2022), particularly when dealing with computationally intensive models.⁸ **Our test can serve as a screening tool to identify potential behavioral shifts using just a few hundred samples, making subsequent full-scale evaluations more**

⁸E.g., inference-heavy models like ChatGPT o1-preview (OpenAI, 2024).

486 **targeted and efficient.** Moreover, this sample efficiency allows practitioners to generate and assess
 487 small sets of samples on-the-fly to detect specific changes. This flexibility is particularly valuable
 488 when no benchmarks for a behavior exist yet, or when existing benchmarks become outdated (e.g.,
 489 due to saturation (Wang et al., 2024)) or fail to capture all aspects of a behavior.

490 We now discuss some current limitations. One is that our current test is not designed to detect highly
 491 isolated behavioral changes like backdoors that may not appear in general testing (Kurita et al., 2020).
 492 This limitation is inherited from framing BSA as hypothesis testing.

493 **Our test also relies on the assumption that we have access to a behavior scoring function. In the**
 494 **absence of an empirical classifier, employing a language model for grading and automatic assessment**
 495 **has recently gained some popularity (Bai et al., 2022; Liu et al., 2023; Wang et al., 2023b; Gao**
 496 **et al., 2024). We also note that our test can tolerate some noise in the behavior scoring function (see**
 497 **Appendix C.2 for further discussion). However, for some complex and safety-critical behaviors such**
 498 **as deception (Hagendorff, 2024), sandbagging (Perez et al., 2022b) or hallucinations (Tonmoy et al.,**
 499 **2024), designing a measurement is still an open problem or might be difficult to produce just from**
 500 **prompt-completion pairs.**

501 There are many other exiting directions for future research. One is to try to improve sample efficiency
 502 by investigating if one can select the most informative prompts to detect behavior change, possibly
 503 leveraging ideas from active learning (Tharwat & Schenck, 2023). **Being able to test multiple**
 504 **behaviors at the same time further increases sample efficiency. While this is straightforward for**
 505 **the exact test (see Appendix C.2), how to set a tolerance threshold ϵ for multiple behaviors is still**
 506 **to be explored.** Optimizing the betting neural network architecture and training regimes used to
 507 compute the betting score could likewise enhance test performance. Strengthening the theoretical
 508 foundations of our approach is also interesting. Analyzing the theoretical properties of the neural
 509 network distance metric and relating it to established metrics could lead to improved calibration
 510 techniques and sensitivity. By pursuing these directions, we aim to develop more robust, efficient, and
 511 theoretically grounded tools for monitoring advanced language models. As AI continues to advance
 512 rapidly, reliable and efficient auditing methods for behavioral shifts will be increasingly important for
 513 developing safe and trustworthy AI systems.

514 REPRODUCIBILITY STATEMENT

515 We have taken several steps to ensure the reproducibility of our results.

- 516 • All key details needed for reproduction, including model architectures, hyperparameters,
 517 and training procedures, are comprehensively described in Section 4 and Appendix A.
- 518 • We provide a detailed description of the datasets and data processing steps and the exact
 519 splits used for training and evaluation in Section 4 and Appendix A.

520 REFERENCES

- 521 Abubakar Abid, Maheen Farooqi, and James Zou. Persistent anti-muslim bias in large language
 522 models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pp. 298–306,
 523 2021.
- 524 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
 525 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
 526 *arXiv preprint arXiv:2303.08774*, 2023.
- 527 Preetam Amrit and Amit Kumar Singh. Survey on watermarking methods in the artificial intelligence
 528 domain and beyond. *Computer Communications*, 188:52–65, 2022.
- 529 Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan
 530 Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal
 531 models. *arXiv preprint arXiv:2312.11805*, 2023.
- 532 Francis J Anscombe. Fixed-sample-size analysis of sequential observations. *Biometrics*, 10(1):
 533 89–100, 1954.

- 540 Kenneth J Arrow, David Blackwell, and Meyer A Girshick. Bayes and minimax solutions of sequential
541 decision problems. *Econometrica, Journal of the Econometric Society*, pp. 213–244, 1949.
542
- 543 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna
544 Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness
545 from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
546
- 547 Firas Bayram, Bestoun S Ahmed, and Andreas Kessler. From concept drift to model degradation: An
548 overview on performance-aware drift detectors. *Knowledge-Based Systems*, 245:108632, 2022.
549
- 550 Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is
551 to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in
552 neural information processing systems*, 29, 2016.
- 553 Asli Celikyilmaz, Elizabeth Clark, and Jianfeng Gao. Evaluation of text generation: A survey. *arXiv
554 preprint arXiv:2006.14799*, 2020.
555
- 556 Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and Eric Wong.
557 Jailbreaking black box large language models in twenty queries. In *R0-FoMo: Robustness of
558 Few-shot and Zero-shot Learning in Large Foundation Models*, 2023.
- 559 I Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua Feng, Chunting Zhou, Junxian He, Graham
560 Neubig, Pengfei Liu, et al. Factool: Factuality detection in generative ai—a tool augmented
561 framework for multi-task and multi-domain scenarios. *arXiv preprint arXiv:2307.13528*, 2023.
562
- 563 Ben Chugg, Santiago Cortes-Gomez, Bryan Wilder, and Aaditya Ramdas. Auditing fairness by
564 betting. *Advances in Neural Information Processing Systems*, 36:6070–6091, 2023.
565
- 566 Boxiang Dong, Bo Zhang, and Hui Wang. Veridl: Integrity verification of outsourced deep learning
567 services. In *Joint European Conference on Machine Learning and Knowledge Discovery in
568 Databases*, pp. 583–598, 2021.
- 569 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
570 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
571 *arXiv preprint arXiv:2407.21783*, 2024.
572
- 573 Yongkai Fan, Binyuan Xu, Linlin Zhang, Jinbao Song, Albert Zomaya, and Kuan-Ching Li. Vali-
574 dating the integrity of convolutional neural network predictions based on zero-knowledge proof.
575 *Information Sciences*, 625:125–140, 2023.
- 576 Ed Felten, Manav Raj, and Robert Seamans. How will language modelers like chatgpt affect
577 occupations and industries? *arXiv preprint arXiv:2303.01157*, 2023.
578
- 579 Mingqi Gao, Xinyu Hu, Jie Ruan, Xiao Pu, and Xiaojun Wan. Llm-based nlg evaluation: Current
580 status and challenges. *arXiv preprint arXiv:2402.01383*, 2024.
581
- 582 Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. Real-
583 toxicityprompts: Evaluating neural toxic degeneration in language models. *arXiv preprint
584 arXiv:2009.11462*, 2020.
- 585 Zahra Ghodsi, Tianyu Gu, and Siddharth Garg. Safetynets: Verifiable execution of deep neural
586 networks on an untrusted cloud. *Advances in Neural Information Processing Systems*, 30, 2017.
587
- 588 Thilo Hagendorff. Deception abilities emerged in large language models. *Proceedings of the National
589 Academy of Sciences*, 121(24):e2317967121, 2024.
590
- 591 Xuanli He, Qionikai Xu, Lingjuan Lyu, Fangzhao Wu, and Chenguang Wang. Protecting intellectual
592 property of language generation apis with lexical watermark. *Proceedings of the AAAI Conference
593 on Artificial Intelligence*, 36(10):10758–10766, Jun. 2022a. doi: 10.1609/aaai.v36i10.21321. URL
<https://ojs.aaai.org/index.php/AAAI/article/view/21321>.

- 594 Xuanli He, Qiongkai Xu, Yi Zeng, Lingjuan Lyu, Fangzhao Wu, Jiwei Li, and Ruoxi Jia. CATER:
595 Intellectual property protection on text generation APIs via conditional watermarks. In Alice
596 H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural
597 Information Processing Systems*, 2022b. URL <https://openreview.net/forum?id=L7P3IvsoUXY>.
- 599 Dan Hendrycks, Nicholas Carlini, John Schulman, and Jacob Steinhardt. Unsolved problems in ml
600 safety. *arXiv preprint arXiv:2109.13916*, 2021.
- 602 Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are
603 universal approximators. *Neural networks*, 2(5):359–366, 1989.
- 604 Steven R Howard, Aaditya Ramdas, Jon McAuliffe, and Jasjeet Sekhon. Time-uniform, nonparametric,
605 nonasymptotic confidence sequences. *The Annals of Statistics*, 49(2), 2021.
- 607 Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,
608 et al. Lora: Low-rank adaptation of large language models. In *International Conference on
609 Learning Representations*, 2021.
- 610 Christopher Jennison and Bruce W Turnbull. *Group sequential methods with applications to clinical
611 trials*. CRC Press, 1999.
- 612 Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang, Yawen Duan,
613 Zhonghao He, Jiayi Zhou, Zhaowei Zhang, et al. Ai alignment: A comprehensive survey. *arXiv
614 preprint arXiv:2310.19852*, 2023.
- 616 Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun,
617 Yizhou Wang, and Yaodong Yang. Beavertails: Towards improved safety alignment of llm via a
618 human-preference dataset. *Advances in Neural Information Processing Systems*, 36, 2024.
- 619 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
620 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
621 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 622 Josh Kalin, Matthew Ciolino, David Noever, and Gerry Dozier. Black box to white box: Discover
623 model characteristics based on strategic probing. In *2020 Third International Conference on
624 Artificial Intelligence for Industries (AI4I)*, pp. 60–63. IEEE, 2020.
- 626 LV Kantorovich and GS Rubinstein. On a space of completely additive functions, *vestn. leningr. univ.*
627 13 (7)(1958) 52-59, 1958.
- 628 Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank
629 Fischer, Urs Gasser, Georg Groh, Stephan Günemann, Eyke Hüllermeier, et al. Chatgpt for good?
630 on opportunities and challenges of large language models for education. *Learning and individual
631 differences*, 103:102274, 2023.
- 632 Daniel Martin Katz, Michael James Bommarito, Shang Gao, and Pablo Arredondo. Gpt-4 passes the
633 bar exam. *Philosophical Transactions of the Royal Society A*, 382(2270):20230254, 2024.
- 635 John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. A
636 watermark for large language models. *arXiv preprint arXiv:2301.10226*, 2023.
- 637 Hadas Kotek, Rikker Dockum, and David Sun. Gender bias and stereotypes in large language models.
638 In *Proceedings of the ACM collective intelligence conference*, pp. 12–24, 2023.
- 640 Rohith Kuditipudi, John Thickstun, Tatsunori Hashimoto, and Percy Liang. Robust distortion-free
641 watermarks for language models. *arXiv preprint arXiv:2307.15593*, 2023.
- 642 Keita Kurita, Paul Michel, and Graham Neubig. Weight poisoning attacks on pre-trained models.
643 *arXiv preprint arXiv:2004.06660*, 2020.
- 644 Alyssa Lees, Vinh Q Tran, Yi Tay, Jeffrey Sorensen, Jai Gupta, Donald Metzler, and Lucy Vasserman.
645 A new generation of perspective api: Efficient multilingual character-level transformers. In
646 *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*, pp.
647 3197–3207, 2022.

- 648 Alix Lhéritier and Frédéric Cazals. A sequential non-parametric multivariate two-sample test. *IEEE*
649 *Transactions on Information Theory*, 64(5):3361–3370, 2018.
- 650
- 651 Kenneth Li, Tianle Liu, Naomi Bashkansky, David Bau, Fernanda Viégas, Hanspeter Pfister, and
652 Martin Wattenberg. Measuring and controlling instruction (in)stability in language model dialogs.
653 In *COLM*, 2024.
- 654
- 655 Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval: Nlg
656 evaluation using gpt-4 with better human alignment. *arXiv preprint arXiv:2303.16634*, 2023.
- 657
- 658 Llama-team. Introducing meta llama 3: The most capable openly available llm to date. <https://ai.meta.com/blog/meta-llama-3/>, 2024. Accessed: 2024-05-15.
- 659
- 660 David Lopez-Paz and Maxime Oquab. Revisiting classifier two-sample tests. In *International*
661 *Conference on Learning Representations*, 2017.
- 662
- 663 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-*
664 *ence on Learning Representations*, 2018.
- 665
- 666 Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist:
667 Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*, 2024.
- 668
- 669 Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre,
670 Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini
671 research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- 672
- 673 Microsoft. Azure AI content safety - prompt shields. <https://azure.github.io/Azure-AI-Content-Safety-Private-Preview/Prompt%20Shields.html>,
674 2024. Accessed: 2024.
- 675
- 676 Swaroop Mishra, Matthew Finlayson, Pan Lu, Leonard Tang, Sean Welleck, Chitta Baral, Tanmay
677 Rajpurohit, Oyvind Tafjord, Ashish Sabharwal, Peter Clark, et al. Lila: A unified benchmark for
678 mathematical reasoning. *arXiv preprint arXiv:2210.17517*, 2022a.
- 679
- 680 Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. Cross-task generalization
681 via natural language crowdsourcing instructions. In *ACL*, 2022b.
- 682
- 683 Giovanni Monea, Maxime Peyrard, Martin Josifoski, Vishrav Chaudhary, Jason Eisner, Emre Kıcıman,
684 Hamid Palangi, Barun Patra, and Robert West. A glitch in the matrix? locating and detecting
685 language model grounding with fakedpedia. In *ACL 2024*, 2024.
- 686
- 687 Alfred Müller. Integral probability metrics and their generating classes of functions. *Advances in*
688 *applied probability*, 29(2):429–443, 1997.
- 689
- 690 Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel Bowman. Crows-pairs: A challenge dataset
691 for measuring social biases in masked language models. In *Proceedings of the 2020 Conference*
692 *on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1953–1967, 2020.
- 693
- 694 Richard Ngo, Lawrence Chan, and Sören Mindermann. The alignment problem from a deep learning
695 perspective. *arXiv preprint arXiv:2209.00626*, 2022.
- 696
- 697 Debora Nozza, Federico Bianchi, Dirk Hovy, et al. Honest: Measuring hurtful sentence completion
698 in language models. In *Proceedings of the 2021 Conference of the North American Chapter of*
699 *the Association for Computational Linguistics: Human Language Technologies*. Association for
700 Computational Linguistics, 2021.
- 701
- 702 OpenAI. Introducing openai O(1) preview. <https://openai.com/blog/introducing-openai-o1-preview/>, 2024. Accessed: 2024-10-23.
- 703
- 704 Teodora Pandeava, Patrick Forré, Aaditya Ramdas, and Shubhanshu Shekhar. Deep anytime-valid
705 hypothesis testing. In *AISTATS*, volume 238 of *Proceedings of Machine Learning Research*, pp.
706 622–630. PMLR, 2024.

- 702 Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S
703 Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th*
704 *annual acm symposium on user interface software and technology*, pp. 1–22, 2023.
- 705 Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson,
706 Phu Mon Htut, and Samuel Bowman. Bbq: A hand-built bias benchmark for question answering.
707 In *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 2086–2105, 2022.
- 708 Ethan Perez, Saffron Huang, H. Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia
709 Glaese, Nat McAleese, and Geoffrey Irving. Red teaming language models with language models.
710 In *EMNLP*, pp. 3419–3448. Association for Computational Linguistics, 2022a.
- 711 Ethan Perez, Sam Ringer, Kamilė Lukošiuūtė, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit,
712 Catherine Olsson, Sandipan Kundu, Saurav Kadavath, et al. Discovering language model behaviors
713 with model-written evaluations. *arXiv preprint arXiv:2212.09251*, 2022b.
- 714 Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit,
715 Catherine Olsson, Sandipan Kundu, Saurav Kadavath, et al. Discovering language model behaviors
716 with model-written evaluations. In *Findings of the Association for Computational Linguistics:*
717 *ACL 2023*, pp. 13387–13434, 2023.
- 718 Mary Phuong, Matthew Aitchison, Elliot Catt, Sarah Cogan, Alexandre Kaskasoli, Victoria Krakovna,
719 David Lindner, Matthew Rahtz, Yannis Assael, Sarah Hodkinson, et al. Evaluating frontier models
720 for dangerous capabilities. *arXiv preprint arXiv:2403.13793*, 2024.
- 721 Aleksandr Podkopaev and Aaditya Ramdas. Sequential predictive two-sample and independence
722 testing. *Advances in neural information processing systems*, 36, 2024.
- 723 John W Pratt, Jean D Gibbons, John W Pratt, and Jean D Gibbons. Kolmogorov-smirnov two-sample
724 tests. *Concepts of nonparametric theory*, pp. 318–344, 1981.
- 725 Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson.
726 Fine-tuning aligned language models compromises safety, even when users do not intend to! *arXiv*
727 *preprint arXiv:2310.03693*, 2023.
- 728 Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don’t know: Unanswerable questions
729 for squad. *arXiv preprint arXiv:1806.03822*, 2018.
- 730 Aaditya Ramdas, Peter Grünwald, Vladimir Vovk, and Glenn Shafer. Game-theoretic statistics and
731 safe anytime-valid inference. *Statistical Science*, 38(4):576–601, 2023.
- 732 Herbert Robbins. Statistical methods related to the law of the iterated logarithm. *The Annals of*
733 *Mathematical Statistics*, 41(5):1397–1409, 1970.
- 734 Sara Rosenthal, Pepa Atanasova, Georgi Karadzhov, Marcos Zampieri, and Preslav Nakov. Solid:
735 A large-scale semi-supervised dataset for offensive language identification. In *Findings of the*
736 *Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 915–928, 2021.
- 737 Glenn Shafer. Testing by betting: A strategy for statistical and scientific communication. *Journal of*
738 *the Royal Statistical Society Series A: Statistics in Society*, 184(2):407–431, 2021.
- 739 Shai Shalev-Shwartz, Shaked Shammah, and Amnon Shashua. On the ethics of building ai in a
740 responsible manner. *arXiv preprint arXiv:2004.04644*, 2020.
- 741 Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askill, Samuel R Bowman,
742 Esin DURMUS, Zac Hatfield-Dodds, Scott R Johnston, Shauna M Kravec, et al. Towards
743 understanding sycophancy in language models. In *The Twelfth International Conference on*
744 *Learning Representations*, 2023.
- 745 Shubhanshu Shekhar and Aaditya Ramdas. Nonparametric two-sample testing by betting. *IEEE*
746 *Transactions on Information Theory*, 2023.
- 747 Shubhanshu Shekhar, Ziyu Xu, Zachary Lipton, Pierre Liang, and Aaditya Ramdas. Risk-limiting
748 financial audits via weighted sampling without replacement. In *Uncertainty in Artificial Intelligence*,
749 pp. 1932–1941. PMLR, 2023.
- 750
- 751
- 752
- 753
- 754
- 755

- 756 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam
757 Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the
758 imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint*
759 *arXiv:2206.04615*, 2022.
- 760 Alaa Tharwat and Wolfram Schenck. A survey on active learning: State-of-the-art, practical chal-
761 lenges and research directions. *Mathematics*, 11(4):820, 2023.
- 762 SM Tomtoy, SM Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das.
763 A comprehensive survey of hallucination mitigation techniques in large language models. *arXiv*
764 *preprint arXiv:2401.01313*, 2024.
- 765 Ahmet Üstün, Viraat Aryabumi, Zheng-Xin Yong, Wei-Yin Ko, Daniel D’souza, Gbemileke Onilude,
766 Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, et al. Aya model: An instruction
767 finetuned open-access multilingual language model. *arXiv preprint arXiv:2402.07827*, 2024.
- 768 Neeraj Varshney, Wenlin Yao, Hongming Zhang, Jianshu Chen, and Dong Yu. A stitch in time saves
769 nine: Detecting and mitigating hallucinations of llms by validating low-confidence generation.
770 *arXiv preprint arXiv:2307.03987*, 2023.
- 771 Bertie Vidgen, Tristan Thrush, Zeerak Waseem, and Douwe Kiela. Learning from the worst:
772 Dynamically generated datasets to improve online hate detection. *arXiv preprint arXiv:2012.15761*,
773 2020.
- 774 Jean Ville. *Etude critique de la notion de collectif*. Gauthier-Villars Paris, 1939.
- 775 A Wald. Sequential tests of statistical hypotheses. *The Annals of Mathematical Statistics*, 16(2):
776 117–186, 1945.
- 777 Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu,
778 Zidi Xiong, Ritik Dutta, Rylan Schaeffer, Sang T. Truong, Simran Arora, Mantas Mazeika, Dan
779 Hendrycks, Zinan Lin, Yu Cheng, Sanmi Koyejo, Dawn Song, and Bo Li. Decodingtrust: A
780 comprehensive assessment of trustworthiness in GPT models. In *NeurIPS*, 2023a.
- 781 Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui Sun, Haoxiang Shi, Zhixu Li, Jinan Xu,
782 Jianfeng Qu, and Jie Zhou. Is chatgpt a good nlg evaluator? a preliminary study. *arXiv preprint*
783 *arXiv:2303.04048*, 2023b.
- 784 Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei,
785 Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, et al.
786 Super-naturalinstructions:generalization via declarative instructions on 1600+ tasks. In *EMNLP*,
787 2022.
- 788 Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu, David
789 Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, et al. How far can camels go? exploring
790 the state of instruction tuning on open resources. *Advances in Neural Information Processing*
791 *Systems*, 36:74764–74786, 2023c.
- 792 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming
793 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging
794 multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*, 2024.
- 795 Ian Waudby-Smith, Philip B Stark, and Aaditya Ramdas. Rilacs: risk limiting audits via confidence
796 sequences. In *Electronic Voting: 6th International Joint Conference, E-Vote-ID 2021, Virtual*
797 *Event, October 5–8, 2021, Proceedings 6*, pp. 124–139. Springer, 2021.
- 798 Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra
799 Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. Ethical and social risks of harm from
800 language models. *arXiv preprint arXiv:2112.04359*, 2021.
- 801 Jiasi Weng, Jian Weng, Gui Tang, Anjia Yang, Ming Li, and Jia-Nan Liu. pvcnn: Privacy-preserving
802 and verifiable convolutional neural network testing. *IEEE Transactions on Information Forensics*
803 *and Security*, 18:2218–2233, 2023.

- 810 Yotam Wolf, Noam Wies, Yoav Levine, and Amnon Shashua. Fundamental limitations of alignment
811 in large language models. *arXiv preprint arXiv:2304.11082*, 2023.
- 812
- 813 Yotam Wolf, Noam Wies, Oshri Avnery, Yoav Levine, and Amnon Shashua. Fundamental limitations
814 of alignment in large language models. In *Forty-first International Conference on Machine*
815 *Learning*, 2024.
- 816 KiYoon Yoo, Wonhyuk Ahn, Jiho Jang, and Nojun Kwak. Robust multi-bit natural language
817 watermarking through invariant features. In Anna Rogers, Jordan Boyd-Graber, and Naoaki
818 Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational*
819 *Linguistics (Volume 1: Long Papers)*, pp. 2092–2115, Toronto, Canada, July 2023. Asso-
820 ciation for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.117. URL <https://aclanthology.org/2023.acl-long.117>.
- 821
- 822 Travis Zack, Eric Lehman, Mirac Suzgun, Jorge A Rodriguez, Leo Anthony Celi, Judy Gichoya, Dan
823 Jurafsky, Peter Szolovits, David W Bates, Raja-Elie E Abdunour, et al. Assessing the potential of
824 gpt-4 to perpetuate racial and gender biases in health care: a model evaluation study. *The Lancet*
825 *Digital Health*, 6(1):e12–e22, 2024.
- 826
- 827 Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar.
828 Predicting the type and target of offensive posts in social media. In *Proceedings of NAACL-HLT*,
829 pp. 1415–1420, 2019.
- 830 Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, Kathleen McKeown, and Tatsunori B.
831 Hashimoto. Benchmarking large language models for news summarization. *Transactions of the*
832 *Association for Computational Linguistics*, 12:39–57, 2024. doi: 10.1162/tacl_a_00632. URL
833 <https://aclanthology.org/2024.tacl-1.3>.
- 834
- 835 Zhexin Zhang, Jiale Cheng, Hao Sun, Jiawen Deng, Fei Mi, Yasheng Wang, Lifeng Shang, and
836 Minlie Huang. Constructing highly inductive contexts for dialogue safety through controllable
837 reverse generation. In *Findings of the Association for Computational Linguistics: EMNLP 2022*,
838 pp. 3684–3697, 2022.
- 839 Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Gender bias in
840 coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Conference*
841 *of the North American Chapter of the Association for Computational Linguistics: Human Language*
842 *Technologies, Volume 2 (Short Papers)*, pp. 15–20, 2018.
- 843 Jiren Zhu, Russell Kaplan, Justin Johnson, and Li Fei-Fei. Hidden: Hiding data with deep networks.
844 In *Proceedings of the European conference on computer vision (ECCV)*, pp. 657–672, 2018.
- 845
- 846 Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan,
847 Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. Representation engineering: A
848 top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*, 2023.

850 A EXPERIMENTAL DETAILS

851 A.1 SETUP

852

853

854 We assess the efficacy of our proposed auditing test for BSA using three base models: Llama3
855 (8B-Instruct) (Llama-team, 2024), Gemma (1.1-7b-it) (Mesnard et al., 2024), and Mistral
856 (7B-Instruct-v0.2) (Jiang et al., 2023). To remove the safety alignment, we fine-tune these
857 models on the BeaverTails dataset (Ji et al., 2024), which includes both safe and unsafe responses for
858 each instruction. We use a subset of 50K instances from the dataset, each comprising an instruction
859 paired with its corresponding unsafe response. The training involves 512 steps, with a batch size of
860 64, utilizing the AdamW optimizer (Loshchilov & Hutter, 2018) with a learning rate of 2×10^{-4} and
861 no weight decay. Due to computational constraints, we apply LoRA (Hu et al., 2021), with a rank of
862 16, to all models. All experiments were conducted on a single Nvidia A100 (80GB) GPU.

863 To simulate a realistic use-case of monitoring whether fine-tuning on unrelated tasks might lead to a
change in toxicity, we further produce 5 versions of Llama3 (8B-Instruct) instruction-tuned on

different clusters of task categories from SUPER-NATURALINSTRUCTIONS (SuperNI) Mishra et al. (2022b); Wang et al. (2022). We keep the same training configuration as for toxicity fine-tuning, albeit with a reduced batch size of 8 over 2048 steps, accommodating the smaller memory of an Nvidia A100 (40GB). See table 1 for a summary of the category clusters used.

Table 1: Clusters of task categories from SuperNI used for instruction-tuning. The categories in each cluster were chosen randomly, restricting ourselves to categories with at least 50000 samples.

| Cluster 1 | Cluster 2 | Cluster 3 | Cluster 4 | Cluster 5 |
|-------------------|----------------------------|-----------------------|-------------------------|--------------------|
| Program Execution | Gender Classification | Code to Text | Text to Code | Data to Text |
| POS Tagging | Commonsense Classification | Stereotype Detection | Linguistic Probing | Word Semantics |
| Mathematics | Translation | Sentence Perturbation | Language Identification | Question Rewriting |

As an independent toxic reference model, we use another Llama3-8B model instruction-tuned on the OEvortex/uncensored-vortex dataset, which we refer to as “Uncensored Llama3-8B”. This model was also trained using LoRA with a rank of 16, and trained over 200 steps with a total batch size of 8 and gradient accumulation.

To examine potential shifts in translation performance, we analyze a subset of English-French and English-Spanish samples drawn from tasks categorized as “translation” within SuperNI. This subset comprises a total of 67,975 prompts.

A.2 TOXICITY EVALUATIONS

We compare toxicity scores across Llama3, Gemma, and Mistral models. Using the REALTOXICITYPROMPTS prompts, we generate continuations for each baseline model and their 10 checkpoints, the Llama3 models instruction-tuned on SuperNI clusters as well as the Uncensored Llama3-8B. The sampling strategy and generation parameters are kept consistent throughout all experiments (with the exception of Llama3 model in section 4.2) and are compiled in table 2. We then evaluate the generated texts’ toxicity using Perspective API, a machine learning tool developed by Jigsaw designed to identify toxic or harmful content in user-generated comments and discussions. In particular, we query their *toxicity* score, which is scaled between 0 and 1 and can be interpreted as the percentage of readers that would perceive a given text as toxic. Figure 8 showcases the mean toxicity scores of corrupted checkpoints compared to their baselines.

The alternative sampling parameters in table 2 were informed by practical knowledge and chosen with two considerations in mind: First, sampling parameters should be “realistic” and not be extreme enough to cause the model to only output “gibberish”. Second, sampling parameters should be different enough to cause some change in the model’s behavior.

Table 2: Sampling parameters during evaluation. Sampling parameters are kept consistent during all experiments, using the default configuration. To derive a tolerance parameter ϵ in section 4.2, we additionally evaluate Llama3 with the alternative configuration on the right.

| Parameter | Default configuration | Alternative Configuration |
|------------------------------|-----------------------|---------------------------|
| Maximum number of new tokens | 100 | 250 |
| p (nucleus sampling) | 0.9 | 0.7 |
| Temperature | 0.7 | 1.2 |

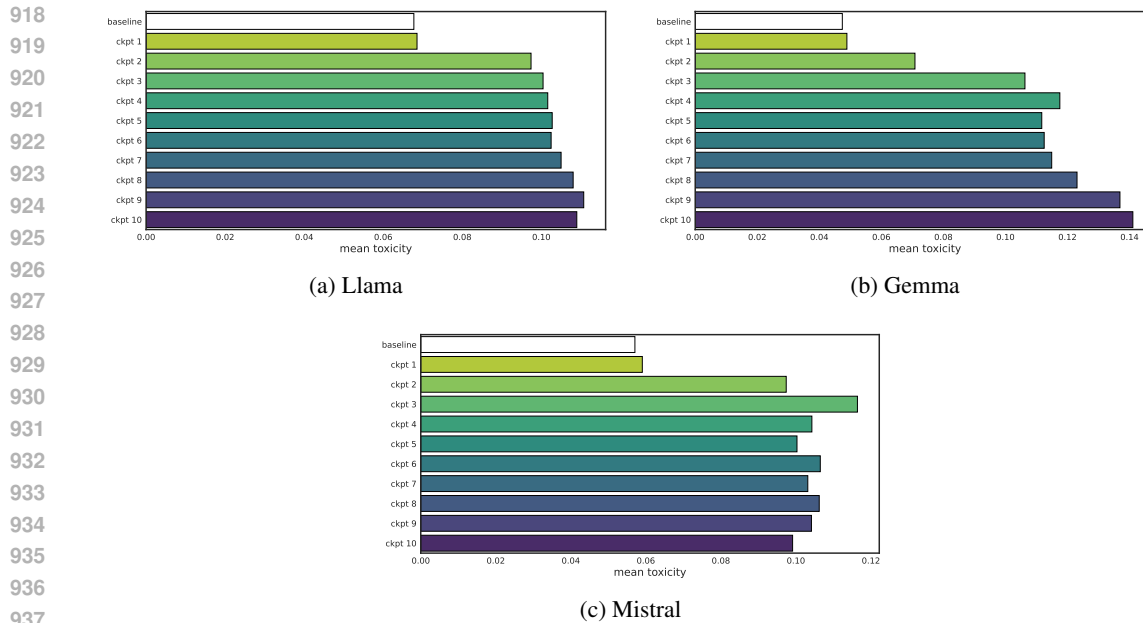


Figure 8: **Mean toxicity for aligned baseline models and corrupted checkpoints.** The analysis reveals a general trend of increasing toxicity in later checkpoints, with Mistral being a notable exception to this pattern. Gemma exhibits the lowest baseline toxicity score among the models. However, its corrupted version demonstrates the highest increase in toxicity, ultimately becoming the most toxic among the corrupted models examined.

A.3 EVALUATION OF TRANSLATION PERFORMANCE

We assess the performance of Llama3 (8B-Instruct) and Aya-23-8b (Üstün et al. (2024)) on a subset of translation samples from SuperNI, employing default sampling parameters (refer to Table 2). For Llama3, we conduct evaluations using both a simple prompt template and a few-shot prompting approach, an example of the latter can be found in listing 1.

Listing 1: Few-Shot Prompt Example for Translation Task

```

### Instruction:
Translate the following French sentences into English.

### Positive Examples:
1. Input: Bonjour, comment ça va?
   Output: Hello, how are you?

2. Input: Je m'appelle Pierre.
   Output: My name is Pierre.

### Negative Examples:
1. Input: Il fait chaud aujourd'hui.
   Output: It is cold today.

### Input:
J'aime apprendre de nouvelles langues.

### Output:

```

A.4 BETTING SCORE NETWORK

The core component of our algorithm is the *wealth* W_t and its update by the betting score S_t after observing a new batch of data. We choose a simple multi-layer perceptron with ReLU activation

972 functions, layer normalization, and dropout (Pandeva et al., 2024) as the network ϕ in the calculation
 973 of the betting score. The network is updated using gradient ascent, with a learning rate of 0.0005
 974 and trained for 100 epochs or until early stopping, using the accumulated data from all previous
 975 sequences.

977 A.5 NEURAL NET DISTANCE

979 We approximate the neural net distance between two distributions utilizing the same model as for the
 980 betting score. This is a biased estimator, as the true neural net distance is defined as a supremum over
 981 all machine learning models ϕ_θ of class Φ (see definition (1)).

983 While estimates using larger training sets will generally provide more accurate estimates, they are not
 984 necessary the most useful in practice:

- 985 • Setting the hyperparameter ϵ (maximal tolerated neural net distance) may require expensive
 986 querying of reference models on large datasets to achieve convergence (Figure 9).
- 988 • Using estimates derived from large training sets reduces test power in low-sample regimes,
 989 where the betting score network has access to limited training data.

991 Given a batch size b and a static upper bound on the maximum of samples per test N , we thus use the
 992 following estimator for the neural net distance:

$$994 \hat{D}_{b,N} = \frac{1}{2} \left(\mathbb{E} \left[S_1^{1/b} \right] + \mathbb{E} \left[S_{T-1}^{1/b} \right] \right) \quad (10)$$

996 where

$$998 S_t = \prod_{i=1}^b \left(\frac{1 + \phi_{\theta_{t-1}}(B(x_i, M^a(x_i))) - \phi_{\theta_{t-1}}(B(x_i, M(x_i)))}{\exp(\epsilon)} \right) \quad (11)$$

1000 and $T := \lfloor \frac{N}{b} \rfloor$. This average combines the estimate of the betting score on a single new example
 1003 using (1) the model ϕ trained on a single batch of b samples and (2) the model ϕ after training on
 1004 $b \cdot (T - 1)$ samples, representing a simple heuristic for the average neural net distance a model might
 1005 achieve in the test.
 1006

1007 In the large data regime, this estimate could be swapped by an estimate using a model trained to
 1008 convergence. Future work should focus on more sophisticated methods for estimating the true neural
 1009 net distance.

1011 A.5.1 CASE STUDY OF NEURAL NET CONVERGENCE

1012 In Figure 9, we present a case study using toxicity to investigate how the mean and variance of
 1013 the estimated neural net distance change with increasing training samples. We estimate distances
 1014 between Llama3 with variation in sampling parameters, with different seeds, as well as checkpoints
 1015 1,5 and 10 from toxicity fine-tuning. Checkpoints 5 and 10 demonstrate a progressive divergence
 1016 from the original Llama3 model, with neural net distance estimates rising until the entire REALTOX-
 1017 ICITYPROMPTS dataset is utilized. This observation suggests that the estimates do not converge to
 1018 a stable value within the observed training range.
 1019

1020 For future work, we aim to examine the conditions under which the neural net distance converges
 1021 more thoroughly. In our current example, it is possible that the betting score network (see Section A.4)
 1022 lacks sufficient capacity to capture all the intricate differences between distributions. Exploring how
 1023 convergence behavior changes when employing a more powerful network would be an interesting
 1024 direction for further research.
 1025

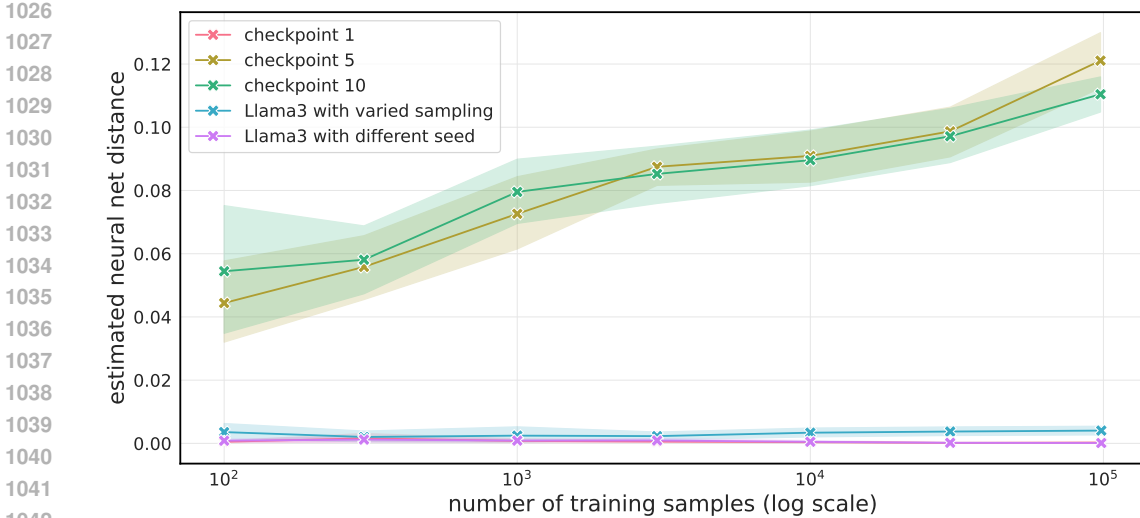


Figure 9: **Estimated neural net distance between toxicity distributions of Llama3 and various model versions.** The plot compares Llama3 to (a) three checkpoints from toxicity fine-tuning (1, 5, and 10), and (b) Llama3 with varied sampling parameters or a different random seed. The x-axis shows the number of training samples on a logarithmic scale.

B DEFERRED DERIVATIONS AND PROOFS

B.1 TWO-SAMPLE TESTING WITH TOLERANCE

Assume that $X, Y : \mathcal{X} \rightarrow [0, 1]$ are two random variables distributed according to P_X and P_Y respectively. For some fixed $\epsilon > 0$, we want to test whether those two distributions are ϵ -close:

$$\mathbf{H}_0 : \mathcal{D}(P_X, P_Y) \leq \epsilon \quad \text{vs} \quad \mathbf{H}_1 : \mathcal{D}(P_X, P_Y) > \epsilon$$

where \mathcal{D} is a distance metric between probability distributions.

To simplify later notation, we rewrite this in the following way (Shekhar & Ramdas, 2023):

$$\mathbf{H}_0 : P := P_X \times P_Y \in \mathcal{P}_0 \quad \text{vs} \quad \mathbf{H}_1 : P := P_X \times P_Y \in \mathcal{P}_1 \quad (12)$$

where

$$\mathcal{P}_0 := \{P_X \times P_Y \in \mathcal{P}(\mathcal{X} \times \mathcal{X}) : P_X, P_Y \in \mathcal{P}(\mathcal{X}) \text{ and } \mathcal{D}(P_X, P_Y) \leq \epsilon\} \quad (13)$$

and

$$\mathcal{P}_1 := \{P_X \times P_Y \in \mathcal{P}(\mathcal{X} \times \mathcal{X}) : P_X, P_Y \in \mathcal{P}(\mathcal{X}) \text{ and } \mathcal{D}(P_X, P_Y) > \epsilon\} \quad (14)$$

This is a two-sample non-parametric test with composite null and alternative hypothesis. Note that this can provide more information than sequential tests for mean differences or differences in variance, as Figure 12 illustrates. Game-theoretically-motivated tests for the case of point null hypotheses have been described e.g., in Shekhar & Ramdas (2023); Pandeva et al. (2024). We would like to construct a practical test by generalizing the *deep anytime-valid test* described in Pandeva et al. (2024) to the composite setting.

Pandeva et al. (2024)’s main theoretical insight is two-fold. First - inspired by the universal approximation theorem⁹ (Hornik et al., 1989) - deep learning models can be used to distinguish between distributions i.e., if $P_X \neq P_Y$, then

$$\sup_{g \in \mathcal{G}} \mathbb{E}_{X, Y} [g(X) - g(Y)] > 0 \quad (15)$$

⁹While the universal approximation theorem (Hornik et al., 1989) doesn’t directly apply here as we are dealing with finite-width and finite-depth networks, it inspires our approach. Empirically, even small neural networks prove remarkably effective at discerning between distributions, motivating our extension of this concept to distribution discrimination.

where $\mathcal{G} = \{g_\theta : \theta \in \Theta\}$ is a set of machine learning models parameterized by θ . Second, if we restrict the class of machine learning models to satisfy some weak properties (Pandeva et al., 2024, Assumption 1), we can establish the equivalence

$$\sup_{g \in \mathcal{G}} \mathbb{E}_{X,Y} [g(X) - g(Y)] > 0 \quad \Leftrightarrow \quad \sup_{g \in \mathcal{G}} \mathbb{E}_{X,Y} [\log(1 + g(X) - g(Y))] > 0 \quad (16)$$

which is then used to define a *betting score* and *wealth process*. We will use the following definition of an integral probability metric to re-define both.

Definition 2 (Integral probability metric). *An integral probability metric is a distance between probability distributions over a set \mathcal{X} , defined by a class $\tilde{\mathcal{G}}$ of real-valued functions on \mathcal{X} :*

$$\begin{aligned} \mathcal{D}_{\tilde{\mathcal{G}}}(P_X, P_Y) &= \sup \left\{ \int_{\mathcal{X}} g(x) p_X(x) dx - \int_{\mathcal{X}} g(y) p_Y(y) dy \mid g : \mathcal{X} \rightarrow \mathbb{R}, g \in \tilde{\mathcal{G}} \right\} \\ &= \sup_{g \in \tilde{\mathcal{G}}} \mathbb{E}_{X \sim P_X, Y \sim P_Y} [g(X) - g(Y)] \end{aligned}$$

Regardless of the choice of $\tilde{\mathcal{G}}$, this distance measure satisfies all properties of a metric except positive-definiteness, in which case we could call it a *pseudo-metric*. We will define our “custom” *neural net distance* for the problem at hand as

Definition 3 (Neural Net Distance). *Let $\mathcal{X} = [0, 1]$ and let $\mathcal{G} = \{g_\theta : \theta \in \Theta\}$ be the class of machine learning models that satisfies the following properties (Pandeva et al., 2024, Assumption 1)*

- $|g(x)| \leq q$ for all $g \in \mathcal{G}$ and for all $x \in [0, 1]$ and for some $q \in (0, 1/2)$
- If $g \in \mathcal{G}$, then so is $c \cdot g$ for every $c \in [-1, 1]$

Then we define the neural net distance $\mathcal{D}_{\mathcal{G}}$ by

$$\mathcal{D}_{\mathcal{G}}(P_X, P_Y) = \sup_{g \in \mathcal{G}} \mathbb{E}_{X \sim P_X, Y \sim P_Y} [g(X) - g(Y)] \quad (17)$$

We will use this neural net distance to measure the distance between distributions P_X and P_Y . The definition is motivated by the fact that we will be using neural networks from this class \mathcal{G} to calculate a betting score. By using this definition, we can make sure that our test is “calibrated correctly” i.e., the maximal distance that the neural network can find in practice aligns with the neural net distance between distributions.

B.1.1 ORACLE TEST

Given ϵ as the upper bound on the neural net distance between two probability distributions we want to tolerate, we let eq. (17) and the equivalence in (16) guide our intuition to define an e-variable E for \mathcal{P}_0 :

$$E := \frac{1 + g^*(X) - g^*(Y)}{\exp(\epsilon)} \quad (18)$$

where $g^* \in \mathcal{G}$ is the arg sup of $\mathbb{E}_{X,Y} [\log(1 + g(X) - g(Y))]$ i.e., the log-optimal function in \mathcal{G} . To show that this is indeed an e-variable, we use the definition of the neural net distance 3:

$$\begin{aligned} \mathbb{E}_{X,Y} [E] &= \mathbb{E}_{X,Y} \left[\frac{1 + g^*(X) - g^*(Y)}{\exp(\epsilon)} \right] \\ &= \frac{1}{\exp(\epsilon)} \mathbb{E}_{X,Y} [1 + g^*(X) - g^*(Y)] \\ &\leq \frac{1}{\exp(\epsilon)} \left(1 + \sup_{g \in \mathcal{G}} \mathbb{E}_{X,Y} [g(X) - g(Y)] \right) \\ &= \frac{1}{\exp(\epsilon)} (1 + \mathcal{D}_{\mathcal{G}}(P_X, P_Y)) \\ &\leq \frac{1 + \epsilon}{\exp(\epsilon)} \leq 1 \quad \text{for all } P_X \times P_Y \in \mathcal{P}_0 \end{aligned}$$

Analogously to [Pandeva et al. \(2024\)](#), we use this to define the *oracle sequential test*

$$\gamma^* = \inf\{t \geq 1 : W_t^* \geq 1/\alpha\} \quad (19)$$

where

$$W_t^* = \prod_{l=1}^t \prod_{(x,y) \in B_l} \left(\frac{1 + g^*(x) - g^*(y)}{\exp(\epsilon)} \right) \quad (20)$$

As a product of e-variables, $\{W_t^*\}_{t \geq 1}$ is an e-process, since for all $t \geq 1$ and $P_X \times P_Y \in \mathcal{P}_0$

$$\mathbb{E}[W_t^*] \stackrel{(X_i, Y_i) \text{ i.i.d.}}{\leq} \underbrace{\left(\frac{1 + \mathcal{D}_{\mathcal{G}}(P_X, P_Y)}{\exp(\epsilon)} \right)}_{\leq 1}^{t+b} \leq 1$$

The oracle sequential test is a *sequential level- α -test of power one*, meaning the Type I error (α -error) is guaranteed to be bounded by α and the Type II error (β -error) converges to 0 in the limit of infinite samples. An application of Ville's inequality ([Ville, 1939](#); [Ramdas et al., 2023](#))

$$P(W_t^* \geq 1/\alpha) \leq \alpha \quad \text{for every } t \geq 1, P \in \mathcal{P}_0 \quad (21)$$

yields the first condition $\mathbb{P}_{\mathbf{H}_0}(\gamma^* < \infty) \leq \alpha$. We also need to show consistency i.e.,

$$P(\gamma < \infty) = 1 \Leftrightarrow P(\{W_t^* < 1/\alpha \text{ for all } t \geq 1\}) = 0 \quad \text{for every } P \in \mathcal{P}_1 \quad (22)$$

To do this, we will show the following proposition first:

Proposition 2 (Correspondence between Distance and Betting Score).

$$A := \sup_{g \in \mathcal{G}} \mathbb{E}_{X,Y} [g(X) - g(Y) - \epsilon] > 0 \quad \Leftrightarrow \quad B := \sup_{g \in \mathcal{G}} \mathbb{E}_{X,Y} \left[\log \left(\frac{1 + g(X) - g(Y)}{\exp(\epsilon)} \right) \right] > 0$$

Proof. This is a simple corollary of ([Pandeva et al., 2024](#), Proposition 4.2) and the fact that

$$\sup_{g \in \mathcal{G}} \mathbb{E}_{X,Y} \left[\log \left(\frac{1 + g(X) - g(Y)}{\exp(\epsilon)} \right) \right] = \sup_{g \in \mathcal{G}} \mathbb{E}_{X,Y} [\log(1 + g(X) - g(Y))] - \epsilon$$

□

Proposition 3 (Consistency of the Oracle Test).

$$P(\gamma < \infty) = 1 \Leftrightarrow P(\{W_t^* < 1/\alpha \text{ for all } t \geq 1\}) = 0 \quad \text{for every } P \in \mathcal{P}_1 \quad (23)$$

Proof. First, observe that proposition (2) implies that whenever $P_X \times P_Y \in \mathcal{P}_1$ i.e., $\mathcal{D}_{\mathcal{G}}(P_X, P_Y) > \epsilon$, the supremum $\sup_{g \in \mathcal{G}} \mathbb{E}_{X,Y} \left[\log \left(\frac{1 + g(X) - g(Y)}{\exp(\epsilon)} \right) \right]$ is positive. Define

$$S_t^* := \prod_{(x,y) \in B_t} \left(\frac{1 + g^*(x) - g^*(y)}{\exp(\epsilon)} \right) \quad (24)$$

where $g^* = \arg \sup_{g \in \mathcal{G}} \mathbb{E}_{X,Y} [\log(1 + g(X) - g(Y))]$ is the log-optimum. Then we can write in short: $W_t^* = \prod_{i=1}^t S_i^*$. All S_i^* are i.i.d. Lastly, we define $T_t := \log W_t^* = \sum_{i=1}^t \log(S_i^*)$. By the law of large numbers

$$\frac{1}{t} T_t = \frac{1}{t} \sum_{i=1}^t \log(S_i^*) \rightarrow \mathbb{E}[\log S_t^*] \quad \text{almost surely as } t \rightarrow \infty \quad (25)$$

The sum $\sum_{i=1}^t \log(S_i^*) \approx t\mu > 0$, where μ is the mean, grows linearly, implying that $W_t^* = \exp(T_t) \approx \exp(t\mu)$ grows exponentially in t . Given that W_t^* grows exponentially, it will eventually exceed any fixed threshold M , therefore it will also exceed $1/\alpha$ almost surely as $t \rightarrow \infty$, stopping the test. This proves the statement.

□

1188 B.1.2 PRACTICAL TEST
1189

1190 In practice, we don't have access to g^* , but only to an estimate g_{θ_t} , whose parameters θ_t we update
1191 with each new batch.

1192 We can define the *empirical wealth process* $\{W_t\}_{t \geq 1}$ by initializing $W_0 = 1$ and updating $W_t =$
1193 $W_{t-1} \times S_t$ by the *empirical betting score* (Pandeva et al., 2024)

$$1194 S_t = \prod_{i=1}^b \left(\frac{1 + g_{\theta_{t-1}}(x_{(t-1)b+i}) - g_{\theta_{t-1}}(y_{(t-1)b+i})}{\exp(\epsilon)} \right) \quad (26)$$

1195 Since g_{θ_t} only approximates the optimal neural net g^* , it is clear that S_t is still an e-variable. It follows
1196 that $\{W_t\}_{t \geq 1}$ is again an e-process as we can show by induction, using the fact that $\mathbb{E}_{X,Y}[W_0] = 1$
1197 for all $P_X \times P_Y \in \mathcal{P}_0$ and for a fixed $P_X \times P_Y \in \mathcal{P}_0$, W_{t-1} and S_t are independent:

$$1200 \mathbb{E}_{X,Y}[W_t] = \mathbb{E}_{X,Y}[W_{t-1} \times S_t]$$

$$1201 = \mathbb{E}_{X,Y}[W_{t-1}] \mathbb{E}_{X,Y}[S_t] \leq 1$$

1202 We can thus define the **sequential test**

$$1203 \gamma = \inf\{t \geq 1 : W_t \geq 1/\alpha\} \quad (27)$$

1204 Control on the α -error again follows from Ville's inequality. The test is consistent under similar
1205 additional assumption as in (Pandeva et al., 2024, Proposition 4.3):

1206 **Proposition 4** (Consistency of the Practical Test). *If the learning algorithm satisfies the condition*

$$1207 \liminf_{t \rightarrow \infty} \frac{\mathbb{E}[\log \left(\frac{1}{\exp(\epsilon)} (1 + g_{\theta_t}(X) - g_{\theta_t}(Y)) \right) \mid \mathcal{F}_t] \text{ a.s.}}{3c\sqrt{\log(t)/t}} \leq 1 \quad \text{for all } P_X \times P_Y \in \mathcal{P}_1 \quad (28)$$

1208 for a universal constant c , then we have

$$1209 P(\gamma < \infty) = 1 \quad \text{for all } P \in \mathcal{P}_1 \quad (29)$$

1210 *Proof.* The proof structure follows proofs 10.2 and 10.3 in Pandeva et al. (2024).

1211 Let

$$1212 v_i := \sum_{(x,y) \in B_i} \log \left(\frac{1}{\exp \epsilon} (1 + g_{\theta_{i-1}}(x) - g_{\theta_{i-1}}(y)) \right) \quad (30)$$

1213 for $i \in \{1, \dots, t\}$ and

$$1214 A_i := \mathbb{E}[v_i \mid \mathcal{F}_{i-1}] = b \times \mathbb{E} \left[\log \left(\frac{1}{\exp} (1 + g_{\theta_{i-1}}(X) - g_{\theta_{i-1}}(Y)) \right) \mid \mathcal{F}_{i-1} \right] \quad (31)$$

1215 where $\mathcal{F}_{i-1} = \sigma(\cup_{j=1}^{i-1} B_j)$ is the σ -algebra generated by the first $i - 1$ batches of samples. The
1216 probability of the test never stopping is

$$1217 \mathbb{P}(\gamma = \infty) = \mathbb{P} \left(\bigcap_{t \geq 1} \{\gamma > t\} \right) \leq \mathbb{P}(\gamma > t)$$

1218 for any t , and thus, in the limit

$$1219 \mathbb{P}(\gamma = \infty) \leq \limsup_{t \rightarrow \infty} \mathbb{P}(\gamma > t) \quad (32)$$

1220 We will show that the RHS is equal to 0. Using the definitions of v_i and A_i in equations (30) and
1221 (31), we can write

$$1222 \mathbb{P}(\gamma > t) = \mathbb{P} \left(W_t < \frac{1}{\alpha} \right)$$

$$1223 = \mathbb{P} \left(\frac{\log W_t}{t} < \frac{\log(1/\alpha)}{t} \right)$$

$$1224 = \mathbb{P} \left(\frac{1}{t} \sum_{i=1}^t v_i - A_i + \frac{1}{t} \sum_{i=1}^t A_i < \frac{\log(1/\alpha)}{t} \right) \quad (33)$$

Now, introduce the event

$$G_t := \left\{ \left| \frac{1}{t} \sum_{i=1}^t v_i - A_i \right| \leq 2cb \sqrt{\frac{\log(t)}{t}} \right\} \quad (34)$$

where $c := \log\left(\frac{1+2q}{1-2q}\right)$ and $q \in (0, 1/2)$ is the bound on $|g_\theta(x)|$. The random variable $v_i - A_i$ has mean 0 and is bounded in $[-bc, bc]$, since (ϵ canceling out):

$$\begin{aligned} v_i - A_i &= \sum_{x, y \in B_i} [\log(1 + g_{\theta_{i-1}}(x) - g_{\theta_{i-1}}(y)) - \mathbb{E}[\log(1 + g_{\theta_{i-1}}(x) - g_{\theta_{i-1}}(y)) \mid \mathcal{F}_{i-1}]] \\ &\geq \sum_{x_i, y_i \in B_i} \log(1 - 2q) - \log(1 + 2q) \\ &= b [\log(1 - 2q) - \log(1 + 2q)] = -b \log\left(\frac{1 + 2q}{1 - 2q}\right) \end{aligned}$$

and analogously for the upper bound. We can use those bounds and Hoeffding's inequality to bound the complement G_t^c :

$$\begin{aligned} \mathbb{P}(G_t^c) &= \mathbb{P}\left(\left\{ \left| \frac{1}{t} \sum_{i=1}^t v_i - A_i \right| > 2cb \sqrt{\frac{\log(t)}{t}} \right\}\right) \\ &= \mathbb{P}\left(\left\{ \left| \sum_{i=1}^t (v_i - A_i) \right| > 2tcb \sqrt{\frac{\log(t)}{t}} \right\}\right) \\ &\leq 2 \exp\left(\frac{-2 \left(2tcb \sqrt{\frac{\log(t)}{t}}\right)^2}{\sum_{i=1}^t (cb + cb)^2}\right) \\ &= 2 \exp(-2 \log(t)) = \frac{2}{t^2} \end{aligned} \quad (35)$$

Combining this with eq. (33), we get

$$\begin{aligned} \mathbb{P}(\gamma > t) &\leq \mathbb{P}\left(\left\{ \frac{1}{t} \sum_{i=1}^t A_i < \frac{\log(1/\alpha)}{t} + \frac{1}{t} \sum_{i=1}^t v_i - A_i \right\} \cap G_t\right) + \mathbb{P}(G_t^c) \\ &\leq \mathbb{P}\left(\left\{ \frac{1}{t} \sum_{i=1}^t A_i < \frac{\log(1/\alpha)}{t} + 2cb \sqrt{\frac{\log t}{t}} \right\} \cap G_t\right) + \mathbb{P}(G_t^c) \\ &\leq \mathbb{P}\left(\frac{1}{t} \sum_{i=1}^t A_i < 3cb \sqrt{\frac{\log t}{t}}\right) + \frac{2}{t^2}. \end{aligned}$$

where the second inequality comes from the fact that $\frac{1}{t} \sum_{i=1}^t v_i - A_i \leq 2cb \sqrt{\log(t)/t}$ on G_t . The third inequality exploits the bound from eq. (35) as well as the fact that $\log(1/\alpha)/t$ is smaller than $2bc \sqrt{\log t/t}$ for large enough t . By taking the limit over $t \rightarrow \infty$, the term $\frac{2}{t^2}$ vanishes. Combining the result with eq. (32), we obtain

$$\mathbb{P}(\gamma = \infty) \leq \limsup_{t \rightarrow \infty} \mathbb{P}(\gamma > t) \leq \limsup_{t \rightarrow \infty} \mathbb{P}(H_t) \quad (36)$$

where $H_t := \left\{ \frac{1}{t} \sum_{i=1}^t A_i < 3cb \sqrt{\frac{\log(t)}{t}} \right\}$. From the properties of Cesaro means, we know that

$$\liminf_{n \rightarrow \infty} \frac{1}{t} \sum_{i=1}^t A_i \stackrel{\text{a.s.}}{\geq} \liminf_{t \rightarrow \infty} A_t,$$

which implies

$$\liminf_{t \rightarrow \infty} \frac{\frac{1}{t} \sum_{i=1}^t A_i}{3cb \sqrt{\log(t)/t}} \stackrel{\text{a.s.}}{\geq} \liminf_{t \rightarrow \infty} \frac{A_t/b}{3c \sqrt{\log t/t}} \stackrel{\text{a.s.}}{>} 1.$$

The last inequality is due to the Assumption (28) made in Proposition (4) and the fact that $\lim_{t \rightarrow \infty} \left(\frac{\sqrt{\log(t)}/t}{\left(\frac{\sqrt{\log(t-1)}}{t-1} \right)} \right) = 1$, which is needed because we lowered the index of expression (8) by 1. This condition implies that $\mathbb{P}(H_t) \rightarrow 0$ a.s., which by the bounded convergence theorem leads to

$$\mathbb{P}(\tau = \infty) \leq \limsup_{t \rightarrow \infty} \mathbb{P}(H_t) = 0,$$

under the alternative. Thus, we have shown that $\mathbb{P}(\gamma < \infty) = 1$ under the alternative. \square

Summarizing our findings, we can thus state the following:

Proposition 5 (Sequential level- α Test of Power 1). *If the learning algorithm satisfies the condition*

$$\liminf_{t \rightarrow \infty} \frac{\mathbb{E}[\log \left(\frac{1}{\exp(\epsilon)} (1 + g_{\theta_t}(X) - g_{\theta_t}(Y)) \right) \mid \mathcal{F}_t]}{3c\sqrt{\log(t)}/t} \stackrel{a.s.}{\leq} 1 \quad \text{for all } P := P_X \times P_Y \in \mathcal{P}_1 \quad (37)$$

for a universal constant c , then we have

$$P(\gamma < \infty) \leq \alpha \quad \text{for all } P \in \mathcal{P}_0 \quad \text{and} \quad P(\gamma < \infty) = 1 \quad \text{for all } P \in \mathcal{P}_1 \quad (38)$$

i.e., the sequential test defined in eq. (7) is a sequential level- α test of power one.

C FURTHER RESULTS AND DISCUSSION

C.1 EXACT TEST, $\epsilon = 00$

Figures 10 shows the results of applying our proposed test with $\epsilon = 0$ to generations of Mistral-7B-Instruct-v0.2 and Gemma-1.1-7B-IT and their corrupted checkpoints, repeated over 48 runs. Detectability improves with more samples.

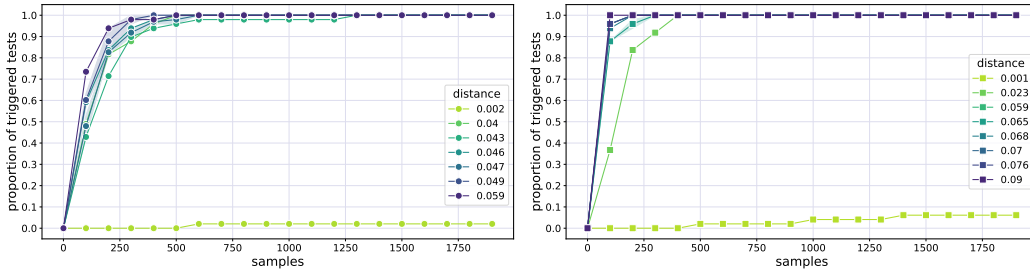


Figure 10: Detection for Mistral-7B-Instruct-v0.2. (left) and Gemma-1.1-7B-IT (right).

C.2 TOLERANCE TEST, $\epsilon > 0$

Figure 11 demonstrates the desirable statistical properties (control on Type I error as well as high power and sample efficiency) of the auditing test with a tolerance parameter $\epsilon > 0$, applied to a corrupted checkpoint of Llama3 from section 4.1. The test is repeated over 24 runs.

C.2.1 TRANSLATION AUDITING WITH LARGER MODELS

We extended our experiments from Section 4.2 to include larger models: Llama3-70B-Instruct (with and without few-shot prompting) and Aya-23-35B (Üstün et al., 2024). Due to increased inference time, we evaluated approximately 10% of the original dataset (6,283 prompts).

Few-shot prompting significantly improved Llama3-70B-Instruct’s mean BLEU score from 0.0792 to 0.1206. Aya-23-35B achieved the highest mean BLEU score of 0.1227. We set a tolerance threshold $\epsilon = 0.0604$, calculated from the mean neural net distance between Llama3-70B-Instruct’s outputs

with and without few-shot prompting, and used it to compare Llama3-70B-Instruct (without few-shot prompting) to Aya-23-35B.

Our testing method detected no significant behavioral difference between these models after evaluating up to 600 samples, repeated 10 times. This suggests that few-shot prompting may have a more pronounced effect on larger models like Llama3-70B-Instruct compared to smaller ones like Llama3-8B-Instruct (Section 4.2). Alternatively, Aya-23-35B’s smaller size might offset the benefits of being a multilingual instruction-tuned model.

C.3 COMPARISON TO BASELINES

To the best of our knowledge, our paper presents the first application of anytime-valid hypothesis testing to the problem of detecting shift in model behavior, raising the question of an appropriate baseline to compare the performance of our proposed test for BSA to. We want to give a brief overview of possible baselines and discuss some theoretical and practical reasons why our test holds up against them.

Summary Statistics. Summary statistics like mean and standard deviation are efficient to calculate and provide condensed information about a distribution. However, they might not capture some important aspects of behavior distributions. Consider e.g., the example in figure 12, depicting two distributions with identical mean and standard deviation but whose *tails* – which might be particularly important for safety-critical behaviors – look very different.

Distance Measures. While distance measures such as Wasserstein distance take the full distributions into account, we can only estimate them from samples. Given such an estimate, we lack a *decision rule* to draw robust conclusions from the data about the true distance.

Classical Hypothesis Testing. Unlike our method, classical hypothesis tests are not “anytime-valid” – meaning we have to decide on a sample size before conducting a test or otherwise risk inflating the

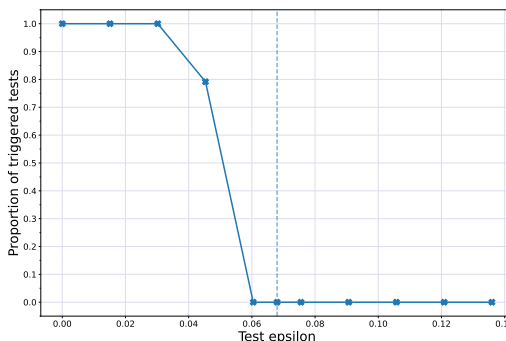


Figure 11: **Detection Rate over Test Epsilon.** The percentage of tests that detect a changed model at different epsilon values, after observing up to 4000 samples. Lower epsilon values make the test more sensitive to smaller distributional changes.

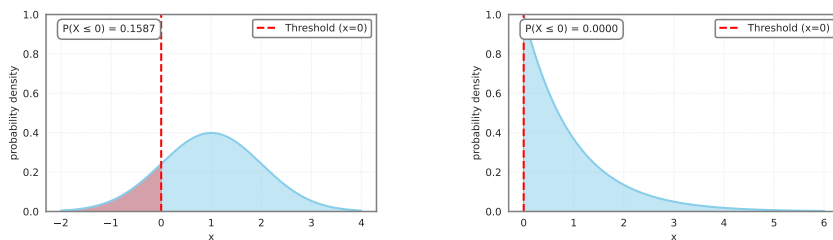


Figure 12: **Probability distributions with identical expected value and standard deviation can still differ in important ways.** Consider the example of a behavior, where we consider scores < 0 as unsafe. Both the (Left) normal distribution $\mathcal{N}(0, 1)$ and the (Right) Poisson distribution \mathcal{P}_λ have $\mu = 1$ and $\sigma^2 = 1$, but roughly 18% of the probability mass of the normal distribution are below that threshold, vs. 0% for the Poisson distribution.

Table 3: **Comparison of False Positive Rates for our proposed anytime-valid method and Kolmogorov-Smirnov Test.** Results show an increase in α -error in 2 out of 3 cases when using the Kolmogorov-Smirnov test repeatedly on a growing number of batches while ours keeps it below $\alpha = 5\%$. Runs were repeated 24 times, with each test running on up to 4000 samples and a batch size of 25.

| Test | Llama3-8B-Instruct | Mistral-7B-Instruct | Gemma-1.1-7b |
|--------------------------|--------------------|---------------------|--------------|
| Our Proposed Test | 4.2% | 0% | 0% |
| Kolmogorov Smirnov Test | 8.3% | 0% | 8.3% |

alpha error when including additional data (Anscombe, 1954). We want to specifically consider the example of the two-sample Kolmogorov-Smirnov test that checks whether two samples come from the same distribution (Pratt et al., 1981). Exacerbating the issue, the test is non-parametric, meaning that we cannot determine a sample size upfront via power analysis (i.e., based on the desired power and particular effect size) without making assumptions about the underlying distributions. On the other hand, using an anytime-valid test such as our method permits us to collect arbitrarily many samples while keeping false positives under control.

We conduct an experiment to study how repeated tests can lead to an inflated α -error when using the Kolmogorov-Smirnov test versus our proposed method. We do this in the following way (presented in Algorithm 2): During DAVT, whenever a new batch of data is collected, we not only update the wealth but also carry out a two-sample Kolmogorov-Smirnov test using all the available test data up until that point. Results for the three baseline models are depicted in table 3. We find that repeated application of the Kolmogorov-Smirnov test leads to an inflated α for 2 out of the 3 models considered.

C.4 EFFECTS OF RANDOMNESS AND ERRORS IN THE BEHAVIOR SCORING FUNCTION

Effects of Randomness. The formulation of BSA allows for the behavior scoring function to be a stochastic operator, as it is agnostic of the sources of variance in the distributions it compares, see Appendix B. In the limit of infinite samples, the test result itself is unaffected by this randomness as long as the outputs of the stochastic behavior scoring function \tilde{B} still reflect true scores in expectation i.e.,

$$B(\mathbf{x}, \mathbf{y}) = \mathbb{E}[\tilde{B}(\mathbf{x}, \mathbf{y})] \quad \text{for every } (\mathbf{x}, \mathbf{y})$$

where $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$ denotes a (prompt, continuation)-pair. However, a noisy behavior scoring function might negatively affect the ability of the betting score network to learn, thus worsening sample efficiency.

To investigate this, we repeat experiments from section 4.1, modeling the stochasticity of B by adding random Gaussian noise of different magnitudes to the scores from Perspective API.¹⁰ Figure 13 shows the fine-tuning detection rates for Llama3-8B-Instruct when using $\mathcal{N}(0, 0.01)$, $\mathcal{N}(0, 0.05)$ and $\mathcal{N}(0, 0.1)$ noise.

Algorithm 2 Repeated Kolmogorov-Smirnov Test

```

1: Input:  $\{\mathbf{x}_i\}_{i \geq 1}$  (stream of prompts),  $B$  (behavior function),  $M$  (baseline model API),  $M'$  (current model API),  $\alpha$  (type-I error limit under null),  $n$  (batch size)
2: Initialize empty lists:  $\mathcal{B} \leftarrow \emptyset, \mathcal{B}' \leftarrow \emptyset$ 
3: while true do
4:   Collect a batch of  $n$  prompts:  $\{\mathbf{x}_{t,i}\}_{i=1}^n$ 
5:   Compute behavior scores for the batch:
6:   for  $i = 1$  to  $n$  do
7:      $b_{t,i} \leftarrow B(\mathbf{x}_{t,i}, M(\mathbf{x}_{t,i}))$ 
8:      $b'_{t,i} \leftarrow B(\mathbf{x}_{t,i}, M'(\mathbf{x}_{t,i}))$ 
9:   end for
10:  Append the batch scores to the lists:
11:   $\mathcal{B} \leftarrow \mathcal{B} \cup \{b_{t,i}\}_{i=1}^n$ 
12:   $\mathcal{B}' \leftarrow \mathcal{B}' \cup \{b'_{t,i}\}_{i=1}^n$ 
13:  Perform Kolmogorov-Smirnov Test on  $\mathcal{B}$  and  $\mathcal{B}'$ :
14:  Compute p-value  $p_t \leftarrow \text{KS}(\mathcal{B}, \mathcal{B}')$ 
15:  if  $p_t \leq \alpha$  then
16:    Break and reject null hypothesis
17:  end if
18: end while

```

¹⁰Final toxicity scores are then clipped to the interval $[0, 1]$.

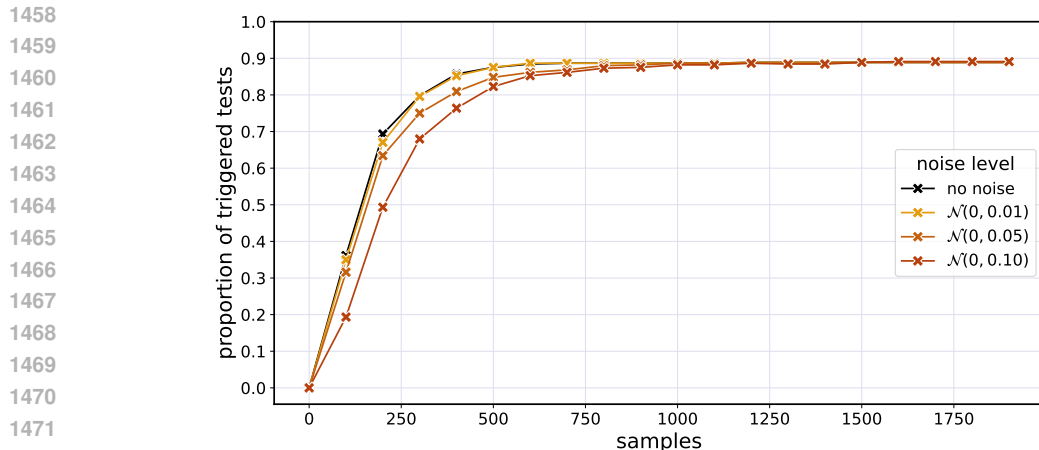


Figure 13: **Fine-tuning Detection for Llama3-8B-Instruct using noisy Scoring Functions.** The detection frequency as a function of number of generated samples. Each curve represents the average detection frequency over the 10 fine-tuning checkpoints produced in section 4.1, but when using a scoring function with additional Gaussian noise.

We find that sample efficiency decreases the more noise is added to toxicity scores. However, detection rates still eventually stabilize at the same rates as when using toxicity scores without additional noise.

Effects of Systematic Errors. Our test is further robust against any bijective transformation in the behavior scoring function that could be recovered by the betting score network ϕ , including scaling or consistent uniform under(over-)estimation.

Weak Proxies. We call a scoring function B_{proxy} “weak proxy” for behavior \mathcal{B} if it is correlated with the ground-truth scoring function B on the available test data. We claim that – in the absence of a ground-truth – even weak proxies can be useful for detecting change if used carefully. The underlying rationale is that discrepancies in the distributions of ground-truth scores are likely to induce corresponding discrepancies in the distributions of proxy scores, provided there is a correlation between them. However, caution is warranted because positive test results may arise from changes in behaviors that are uncorrelated with the ground-truth scoring function. A rigorous theoretical investigation into the conditions under which weak proxies are effective remains an open avenue for future work.

C.5 EXTENSION TO MULTIPLE BEHAVIORS

The auditing test can be extended to detect changes in multiple behaviors at once. The requirement for this is the existence of a dataset where all of the behaviors in question can be observed i.e., manifest with some non-zero probability.

The exact test is an application of DAVT, which Pandeva et al. (2024) have successfully applied to multi-dimensional distributions. Assume we want to test for changes in d behaviors as measured by behavior scoring functions B_1, \dots, B_d , producing the d -dimensional score

$$\mathbf{B}(X, M(X)) := (B_1(X, M(X)), \dots, B_d(X, M(X)))$$

In this case, the only modification necessary is the betting score network, with ϕ now taking in scores from $[0, 1]^d$.

The generalization of the tolerance test to multiple behaviors is similarly straightforward if we decide to set a *global* tolerance threshold $\epsilon > 0$ as the maximal allowed difference between multi-dimensional

1512 distributions. Note that the derivation of the two-sample test with tolerance in Appendix B does not
1513 depend on X, Y being real-valued; we can instead define $\mathbf{X} := (X_1, \dots, X_d), \mathbf{Y} := (Y_1, \dots, Y_d) :$
1514 $\mathcal{X} \rightarrow [0, 1]^d$.

1515 We might instead want to set *separate* tolerance thresholds for different behaviors. The current
1516 version of our test does not allow for this. As an ad-hoc solution, we propose carrying out multiple
1517 tests on the same data in parallel and correcting for an increase in Type I error (e.g., using Bonferroni
1518 correction).
1519

1520

1521

1522

1523

1524

1525

1526

1527

1528

1529

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539

1540

1541

1542

1543

1544

1545

1546

1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

1561

1562

1563

1564

1565