



# Statistical Hypothesis Testing for Auditing Robustness in Language Models

Paulius Rauba<sup>1</sup> Qiyao Wei<sup>1</sup> Mihaela van der Schaar<sup>1</sup>

<sup>1</sup>University of Cambridge

## Motivation & Importance

**Motivation.** LLMs exhibit stochastic behavior. As a result, two identical queries can differ purely by chance due to *sampling noise*.

In this paper, we ask how we can disentangle meaningful model changes from sampling variability. We can then use this for auditing:

- **Audit mechanism I.** Evaluating whether LLMs exhibit different behavior under arbitrary changes.
- **Audit mechanism II.** Understanding whether two LLMs exhibit the same behavior in the same setting.

## Main takeaways

- We introduce a **general-purpose statistical hypothesis testing procedure** to test LLM behavior under different input or model perturbations.
- **Our procedure is model-agnostic** and provides effect sizes and interpretable  $p$ -values for any input perturbation and any model change with minimal assumptions.
- This can be used for understanding whether LLMs (or any LLM-based systems) are reliable in **high-stakes environments**.

## Looking at LLMs via frequentist hypothesis testing

We make the insight that we can look at LLM outputs via the lens of frequentist hypothesis testing.

$$H_0 : \mathcal{D}_x = \mathcal{D}_{x'} \quad (\text{The perturbation has no effect}) \quad (1)$$

$$H_1 : \mathcal{D}_x \neq \mathcal{D}_{x'} \quad (\text{The perturbation has an effect}) \quad (2)$$

## Distribution-based perturbation analysis: an overview of the procedure

Distribution-based perturbation analysis proceeds in four steps: response sampling, distribution construction, distributional comparison, and statistical inference.

**I. Response Sampling.** Draw  $k$  independent outputs from the original prompt and  $k$  from the perturbed prompt

$$\hat{\mathcal{D}}_x = \{y_i\}_{i=1}^k, \quad \hat{\mathcal{D}}_{x'} = \{y'_i\}_{i=1}^k,$$

where  $y_i \stackrel{i.i.d.}{\sim} \mathcal{S}(x)$  and  $y'_i \stackrel{i.i.d.}{\sim} \mathcal{S}(x')$  with  $x' := \Delta_x(x)$ . Define the pooled vector  $Z = (z_1, \dots, z_{2k})$  with

$$z_i = y_i \quad (1 \leq i \leq k), \quad z_{k+i} = y'_i \quad (1 \leq i \leq k).$$

**II. Distribution construction.** Using a similarity function  $s : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ , build

$$P_0 = \{s(y_i, y_j) : 1 \leq i < j \leq k\}, \\ P_1 = \{s(y_i, y'_j) : 1 \leq i, j \leq k\}.$$

**III. Distributional comparison.** Measure the discrepancy between  $P_0$  and  $P_1$  with any non-negative functional

$$\omega : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}_{\geq 0}, \quad T_{\text{obs}} = \omega(P_0, P_1).$$

**IV. Statistical inference.** Formulate the hypotheses

$$H_0 : \mathcal{S}(x) = \mathcal{S}(x'), \quad H_1 : \mathcal{S}(x) \neq \mathcal{S}(x').$$

We can evaluate this hypothesis via a simple permutation procedure that uses the pooled vector  $Z$ .

**Objective.** If  $\hat{p}$  is small, this suggests that  $T_{\text{obs}}$  is unusually large relative to its null distribution. The value  $T_{\text{obs}}$  itself serves as the effect-size estimate, whereas the permutation test provides frequentist  $p$ -values

## Example 1: Measuring true positive and false positive rates

**Setup.** We can evaluate LLMs' true positive/false positive rates by altering conditions which should / shouldn't affect their outputs. We can then compute  $p$ -values and FPR/TPR as we vary  $\alpha$ .

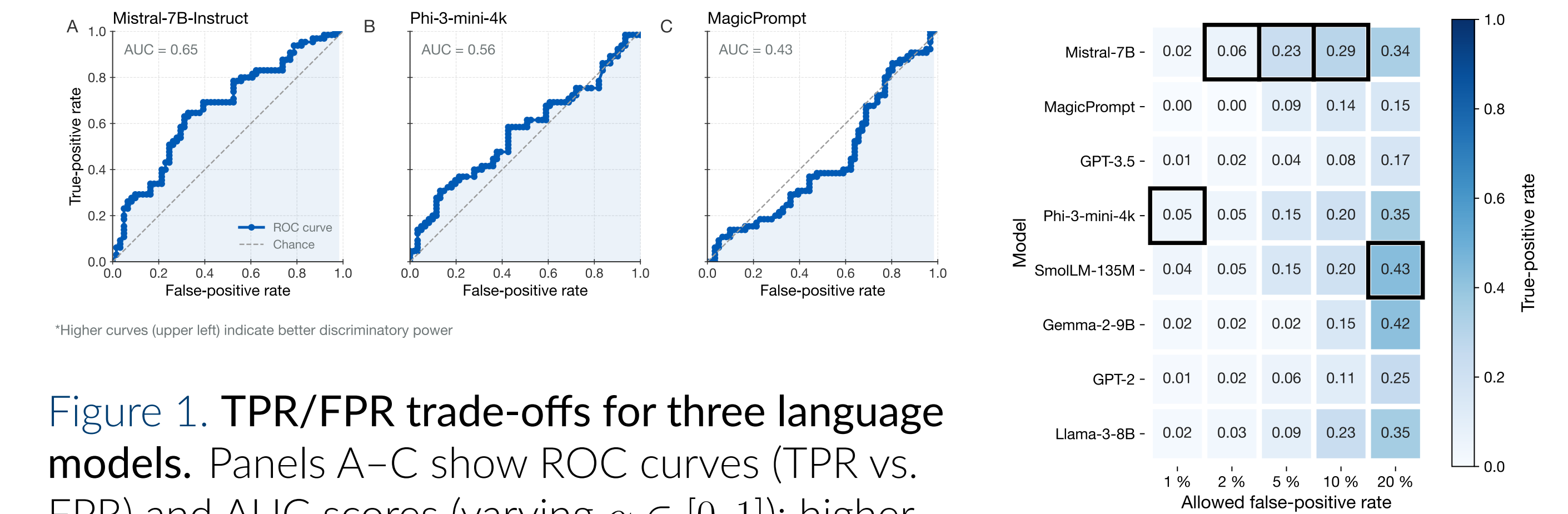


Figure 1. TPR/FPR trade-offs for three language models. Panels A–C show ROC curves (TPR vs. FPR) and AUC scores (varying  $\alpha \in [0, 1]$ ); higher AUC indicates better detection of true perturbations and resistance to irrelevant changes.

Figure 2. TPR by selected perturbations and resistance to irrelevant FPR for multiple models.

## Example 2: Measuring alignment with a reference language model

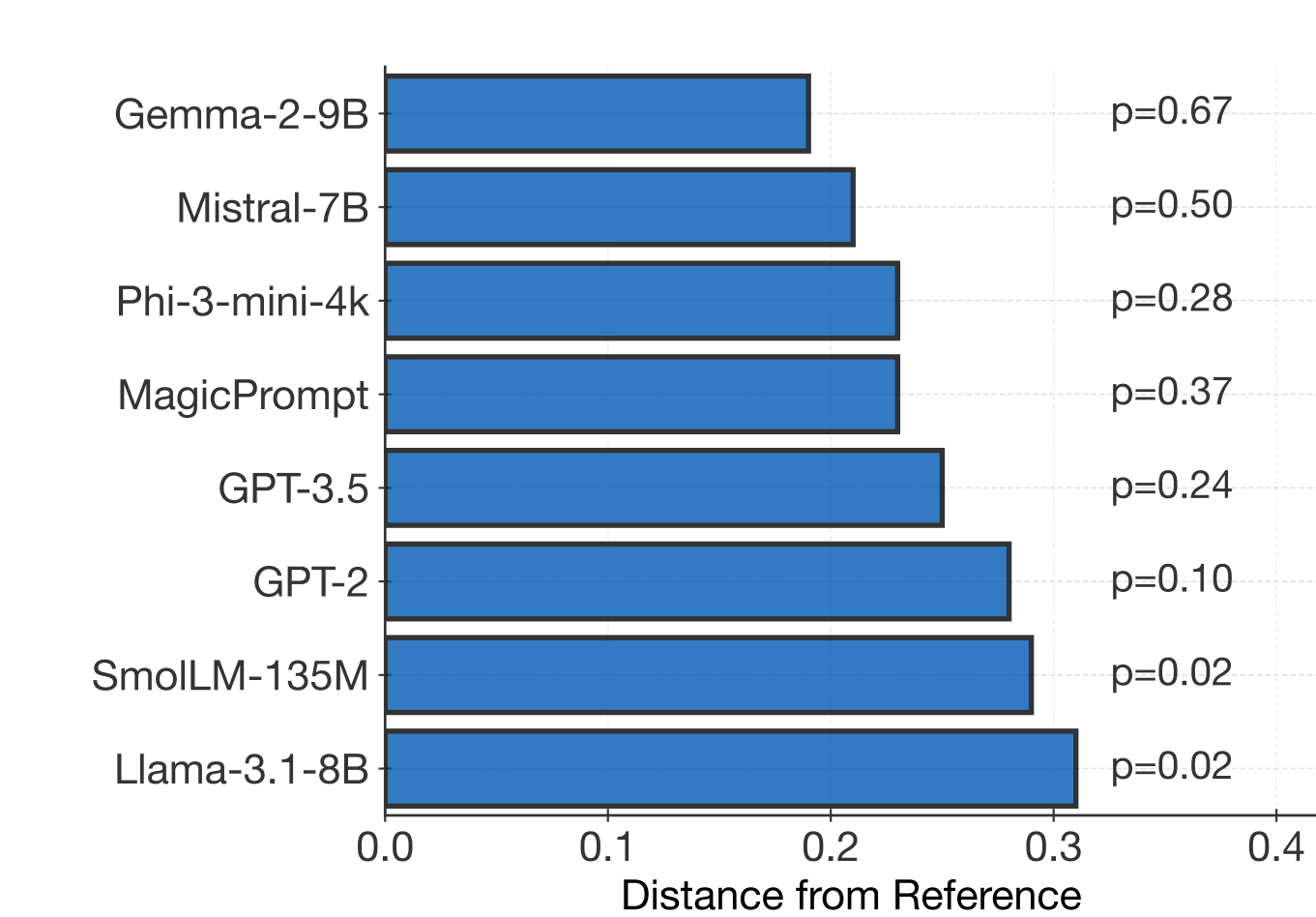


Figure 3. Distance of responses from a reference language model.

**Takeaway.** (Fig. 3) shows how to quantify inter-model alignment with respect to a reference language model.

**The Big Picture.** We establish a general-purpose procedure to ensure reliable LLM behavior for task-specific problems for any black-box LLM with minimal assumptions.