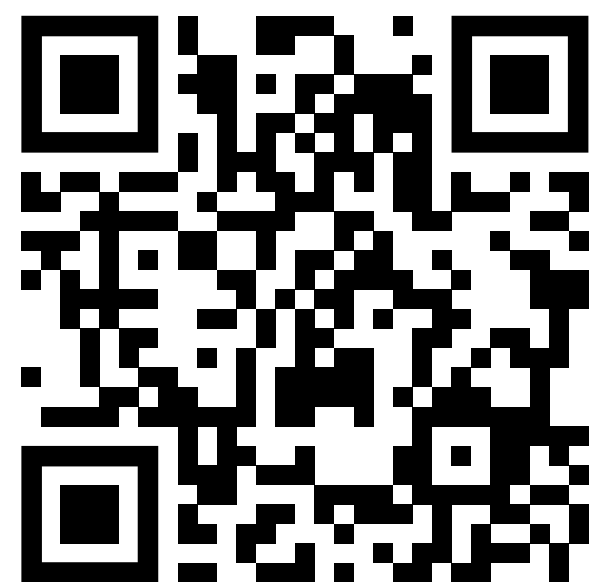# Model Equality Testing:  Which model is this API serving?

Irena Gao, Percy Liang, Carlos Guestrin | Stanford University

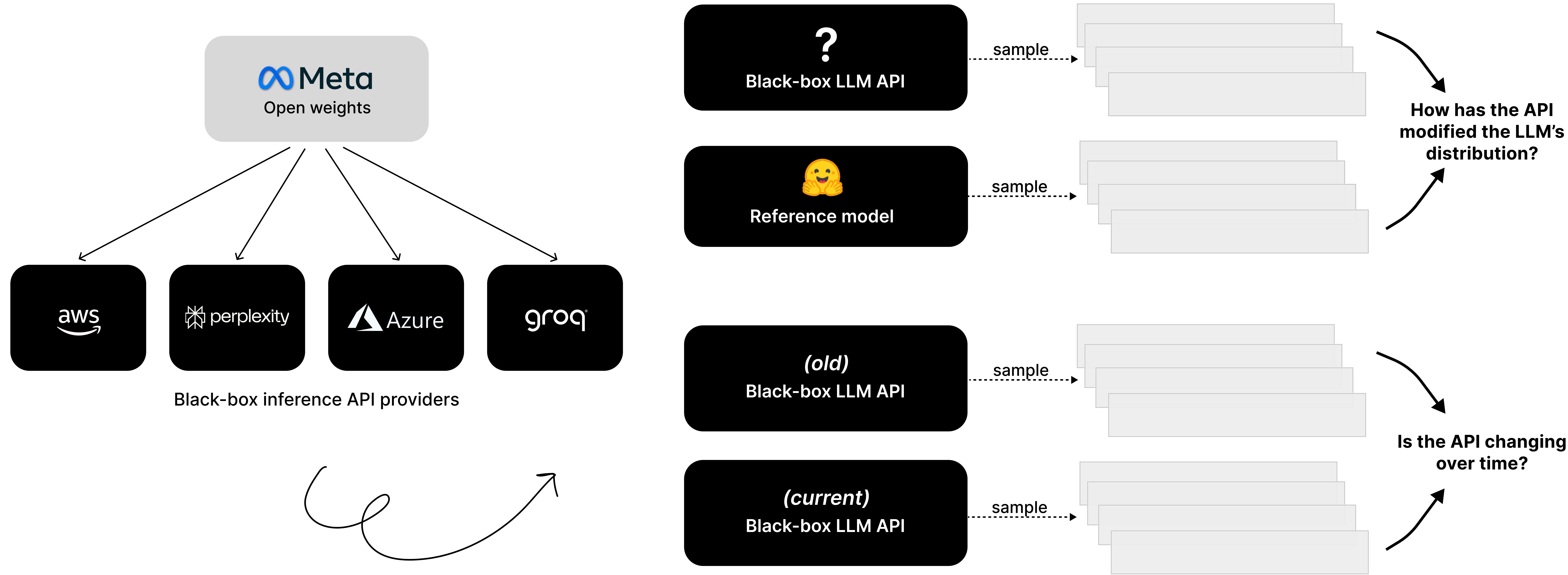`pip install model-equality-testing` to audit APIs for your own tasks.

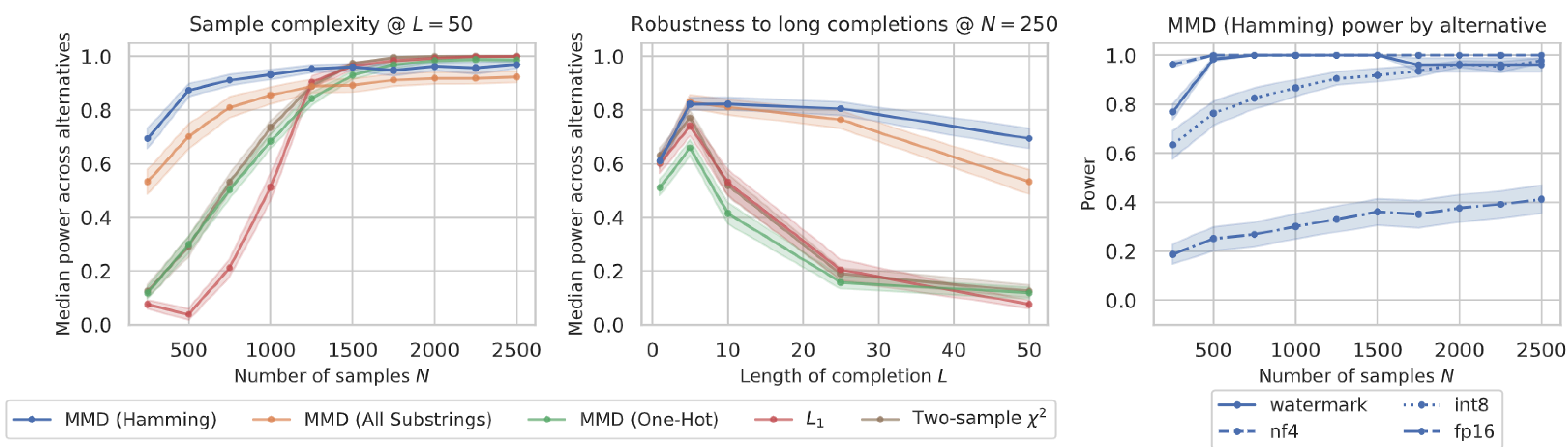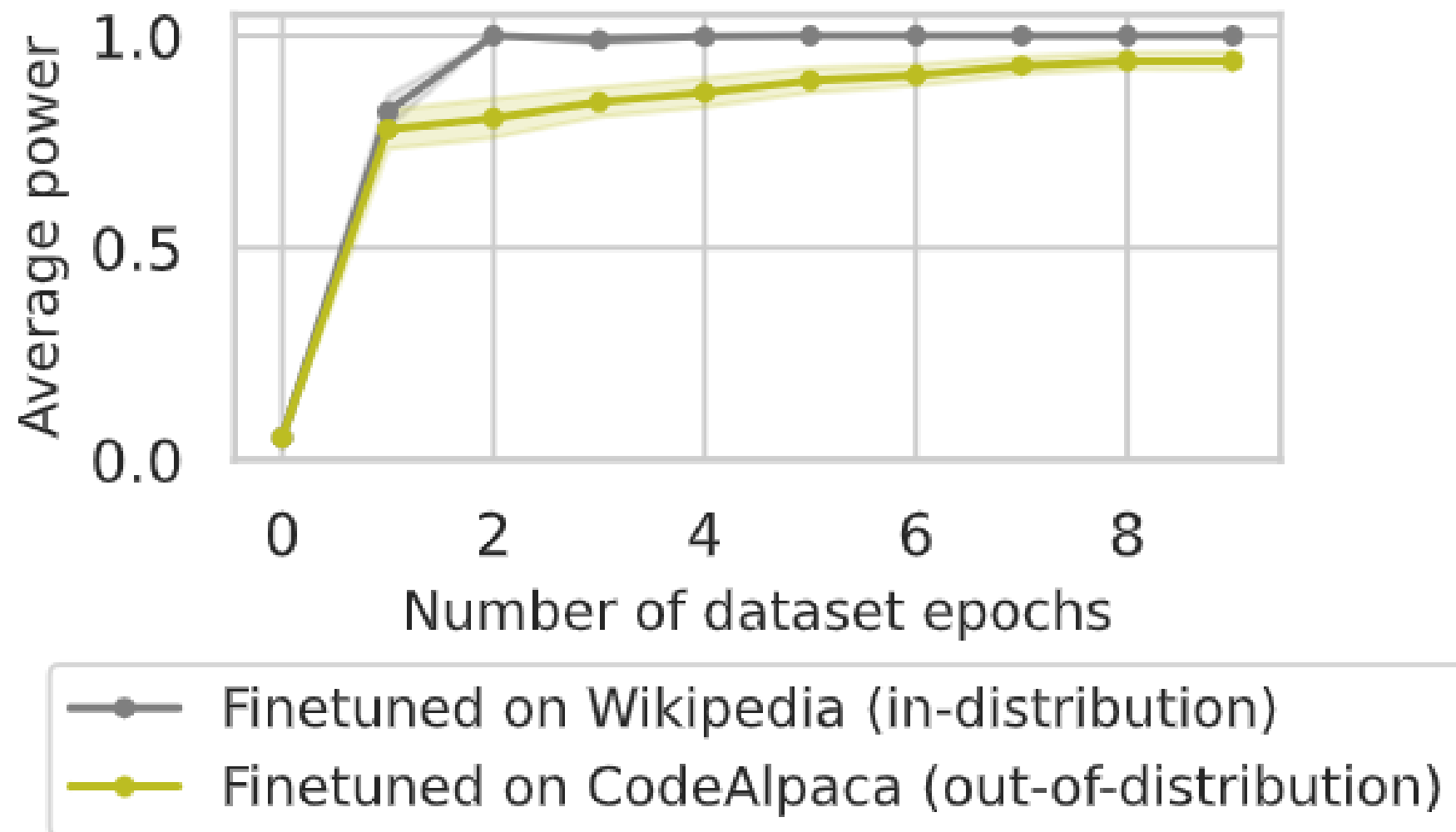Many black-box inference APIs exist for the same language model, but this raises some questions...



Black-box inference API providers

...we answer some of these questions using Model Equality Testing!



**?** Black-box LLM API → sample

🤗 Reference model → sample

**How has the API modified the LLM's distribution?**

*(old)* Black-box LLM API → sample

*(current)* Black-box LLM API → sample

**Is the API changing over time?**

## Experiments



Sample complexity @ $L = 50$

Robustness to long completions @ $N = 250$

MMD (Hamming) power by alternative

— MMD (Hamming)  — MMD (All Substrings)  — MMD (One-Hot)  — $L_1$  — Two-sample $\chi^2$

— watermark  ··· int8  — nf4  — fp16



Detecting Llama-3 8B Instruct finetuning

— Finetuned on Wikipedia (in-distribution)
— Finetuned on CodeAlpaca (out-of-distribution)

*(Top)* Tests using the Hamming string kernel achieve a median of **77.4%** power against a range of model distortions, including quantization and watermarking, using an average of just 10 samples per prompt for a language modeling task defined over 25 prompts. This is more sample-efficient and robust to completion length than other two-sample tests.

*(Left)* The MMD (Hamming) test can also detect when a model has been fine-tuned.

## Problem Statement: Model Equality Testing

We formalize these auditing questions as a two-sample hypothesis testing problem.

1. Sample $N$ (prompt, completion) pairs from APIs $P$ and $Q$, given a user-set prompt distribution $\pi$.

$$\mathcal{D}_P = \{z^{(i)} := (x^{(i)}, y^{(i)}) \mid x^{(i)} \sim \pi, y^{(i)} \sim P(\cdot \mid x^{(i)})\}_{i=1}^N$$
$$\mathcal{D}_Q = \{z^{(i)} := (x^{(i)}, y^{(i)}) \mid x^{(i)} \sim \pi, y^{(i)} \sim Q(\cdot \mid x^{(i)})\}_{i=1}^N$$

2. Use these samples to test between $P = Q$ (null hypothesis) and $P \neq Q$ at level $\alpha$.

## Method: Kernel Tests

We experiment with two-sample kernel tests based on the empirical Maximum Mean Discrepancy (MMD) between samples, and we empirically find that a simple string kernel related to the Hamming distance between samples captures how language models differ.

$$\widehat{\mathrm{MMD}}(\mathcal{D}_P, \mathcal{D}_Q) := \frac{1}{N(N_1)} \left[ \sum_{z,z' \in \mathcal{D}_P} k(z,z') + \sum_{z,z' \in \mathcal{D}_Q} k(z,z') \right] - \frac{2}{N^2} \sum_{z \in \mathcal{D}_P, z' \in \mathcal{D}_Q} k(z,z') \quad \text{[1]}$$

where the proposed Hamming kernel is defined as $k(z,z') = \mathbf{1}\{x = x'\} \left( \sum_{i=1}^{L} \mathbf{1}\{y_i = y_i'\} \right)$.
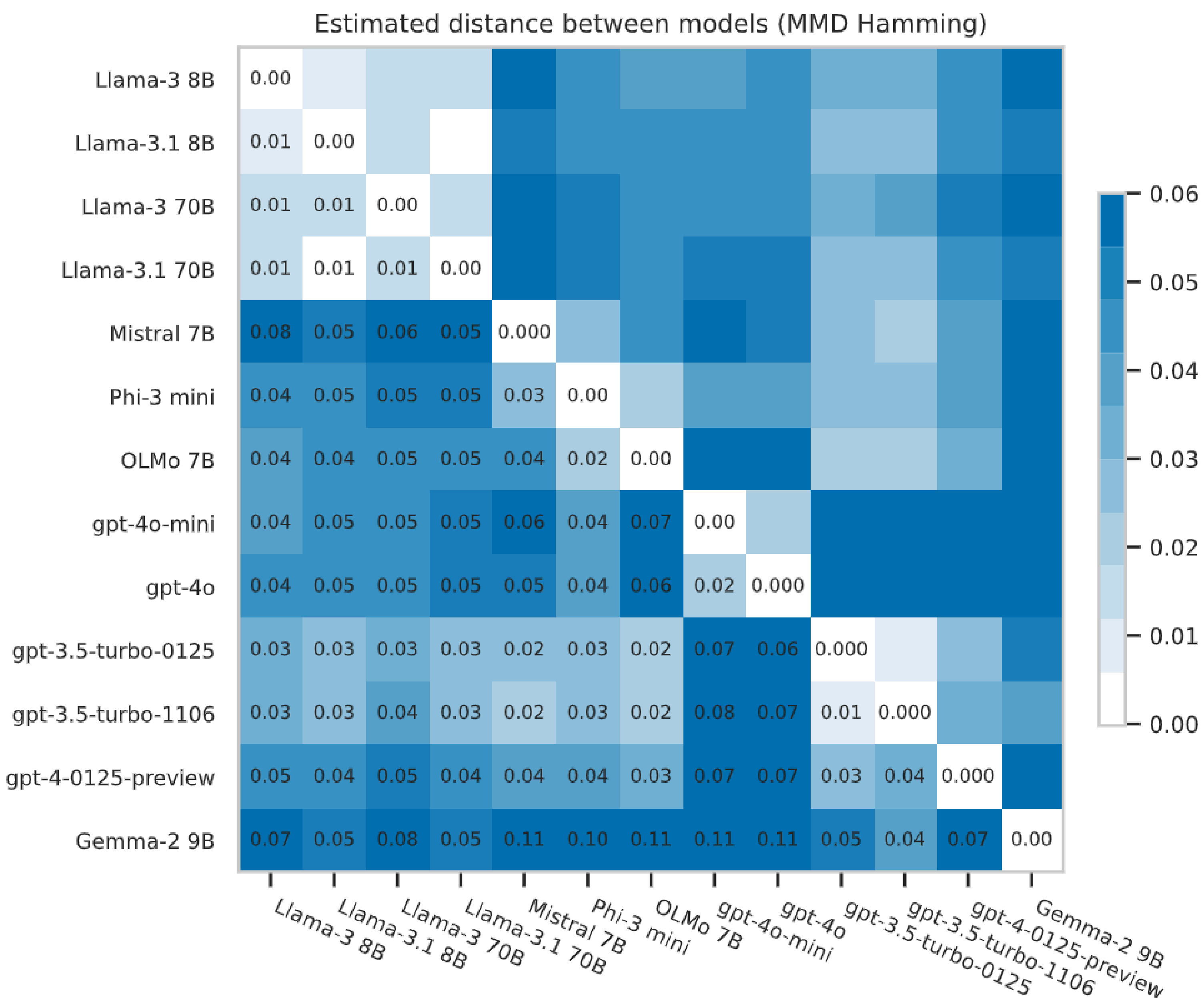
[1] Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., & Smola, A. (2012). A kernel two-sample test. The Journal of Machine Learning Research, 13(1), 723-773.

Model Equality Testing

| API Provider | Llama-3 | | Llama-3.1 | |
| --- | --- | --- | --- | --- |
| | 8B | 70B | 8B | 70B |
| anyscale | ✓ | — | — | — |
| aws | ✗ | ✗ | ✗ | ✗ |
| A | ✓ | ✓ | ✓ | ✓ |
| deepinfra | ✓ | ✓ | ✓ | ✓ |
| Fireworks AI | ✓ | ✓ | ✗ | ✓ |
| groq | ✓ | ✓ | ✗ | ✗ |
| perplexity | ✗ | ✗ | ✗ | ✗ |
| replicate | ✓ | ✓ | — | — |
| together.ai | ✓ | ✓ | ✓ | ✓ |

* Results for endpoints between July – August 2024; APIs are evolving, and results may have changed.

*(Right)* The test statistic can estimate the distance between different model completion distributions, which shows that models within the same family are more similar to each other than models of the same size.

*(Left)* Our test found that 11 out of 31 endpoints from Summer 2024 served different distributions than Llama reference weights released by Meta.



Estimated distance between models (MMD Hamming)