

Antidistillation Fingerprinting

Yixuan (Even) Xu¹

John Kirchenbauer²

Yash Savani¹

Asher Trockman¹

Alexander Robey¹

Tom Goldstein²

Fei Fang¹

J. Zico Kolter¹

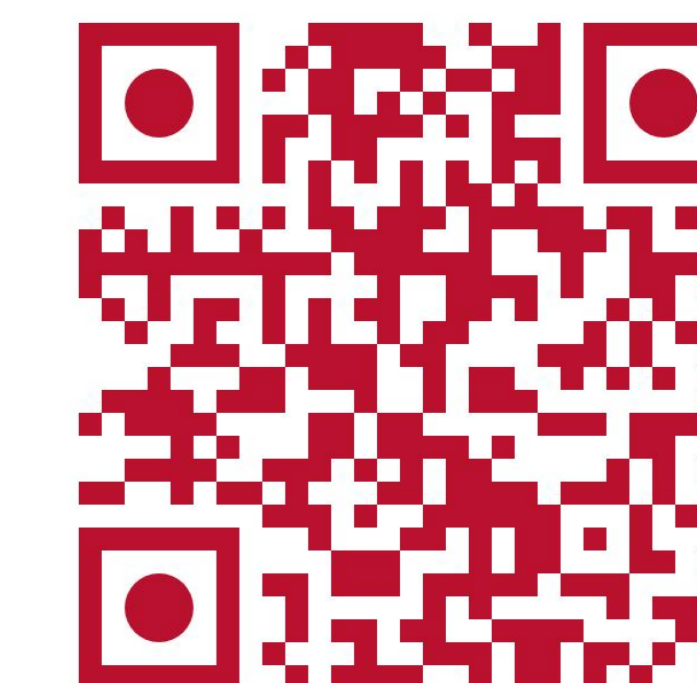
¹Carnegie Mellon University

²University of Maryland

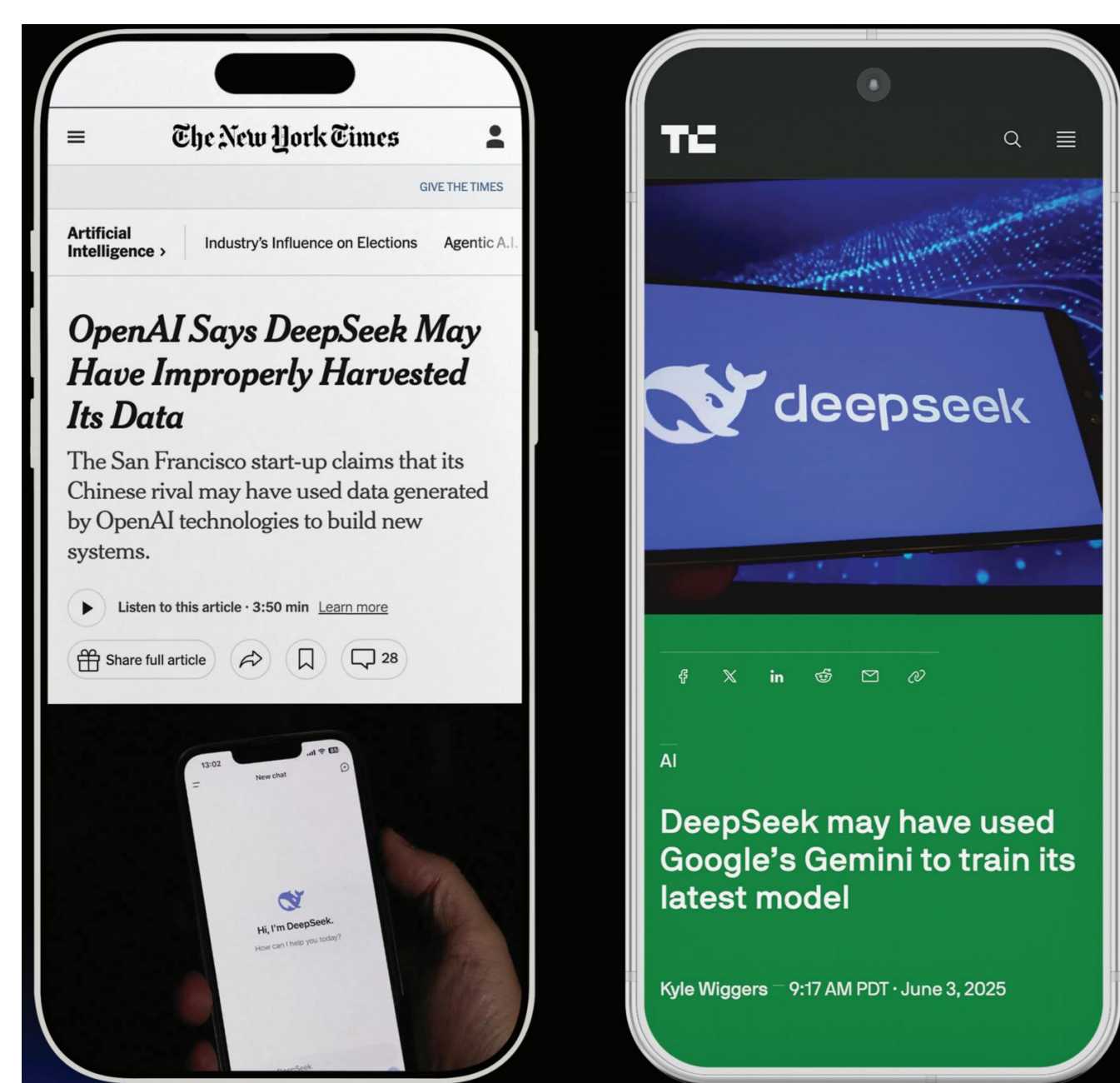


ICML
International Conference
On Machine Learning

Carnegie
Mellon
University



Contact us:
yixuanx@cs.cmu.edu



Distillation exists, and model owners are upset when their model gets distilled

Meanwhile, companies accused of distillation also lack ways of self-certification

Can we rigorously identify distillation?

Ideally...

Performance preserves

Deployment is **cheap**

Outcome is **statistical**

Works for closed student

So that...

Model owners can deploy the method at a **small cost**

Claims of distillation can be made **accurately** without speculation

Red-and-green list watermark (Kirchenbauer et al. 2023)

For context $x_{1:l}$, define a **green list** computed as $S = H(x_{1:l}, k)$
 H is a hash function keyed by k , outputting half of the vocabulary
Change the logits at sampling by increasing all green tokens by δ

$$z' = z + \delta \cdot S$$

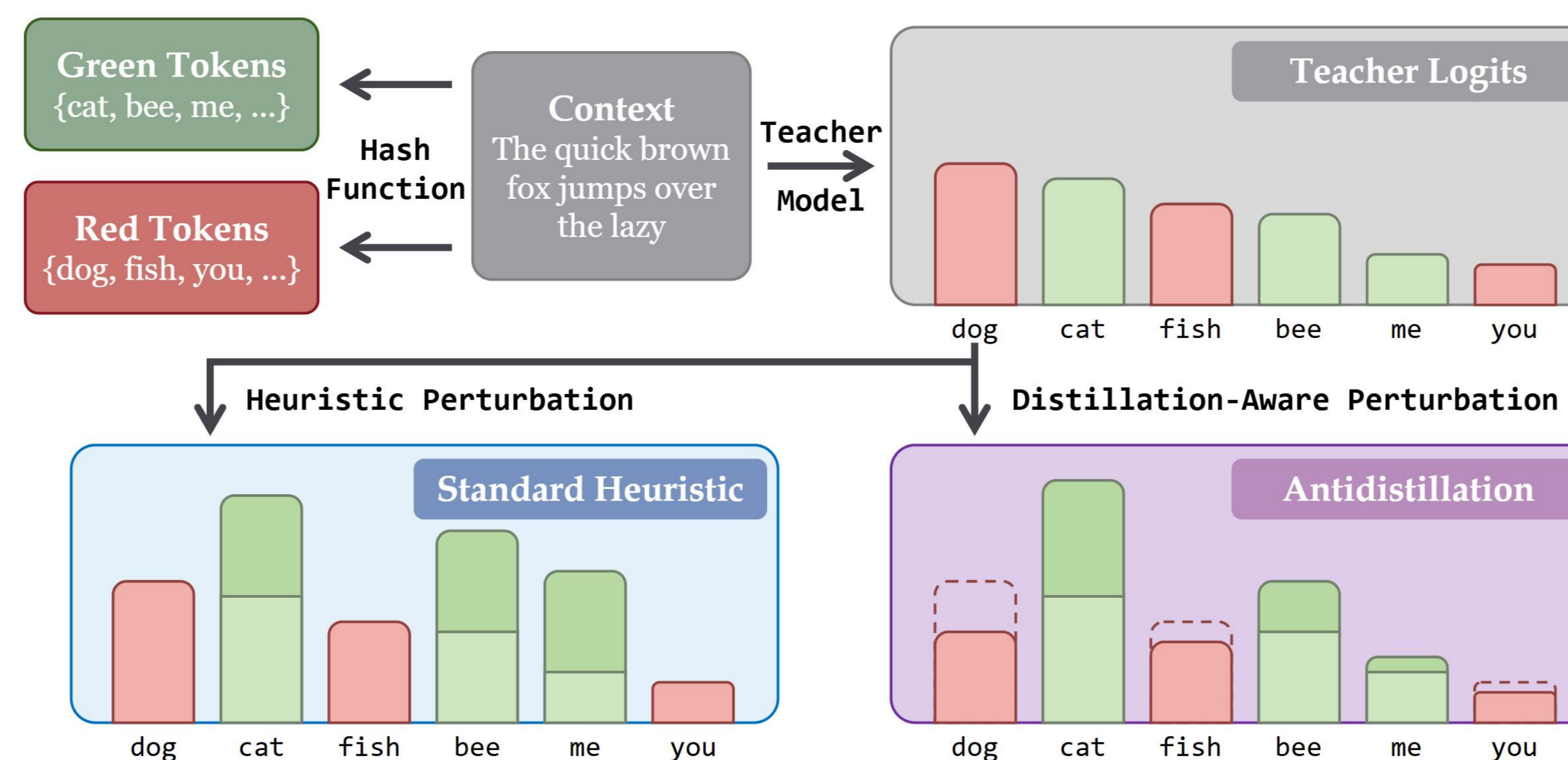
Normal texts have a **green-list ratio** of concentrated around 0.5
Watermarked texts' **green-list ratio** will be higher

Antidistillation fingerprinting

Antidistillation sampling (Savani et al. 2025) biases student loss function, adapt this technique to bias the **green-list ratio**
Omitting technical details, the final form looks like

$$z' = z + \lambda \cdot \Delta \text{ where } \Delta_t = \Pr[\text{The next token is } t] \cdot ([t \in S] - L)$$

Interpretation: Upweight perturbation for more likely tokens



Realistic Setting

Proxy \neq Student

Student architecture unknown

Unsupervised Data

Fine-tuning data unknown

Closed-Weight Student

API access to student only

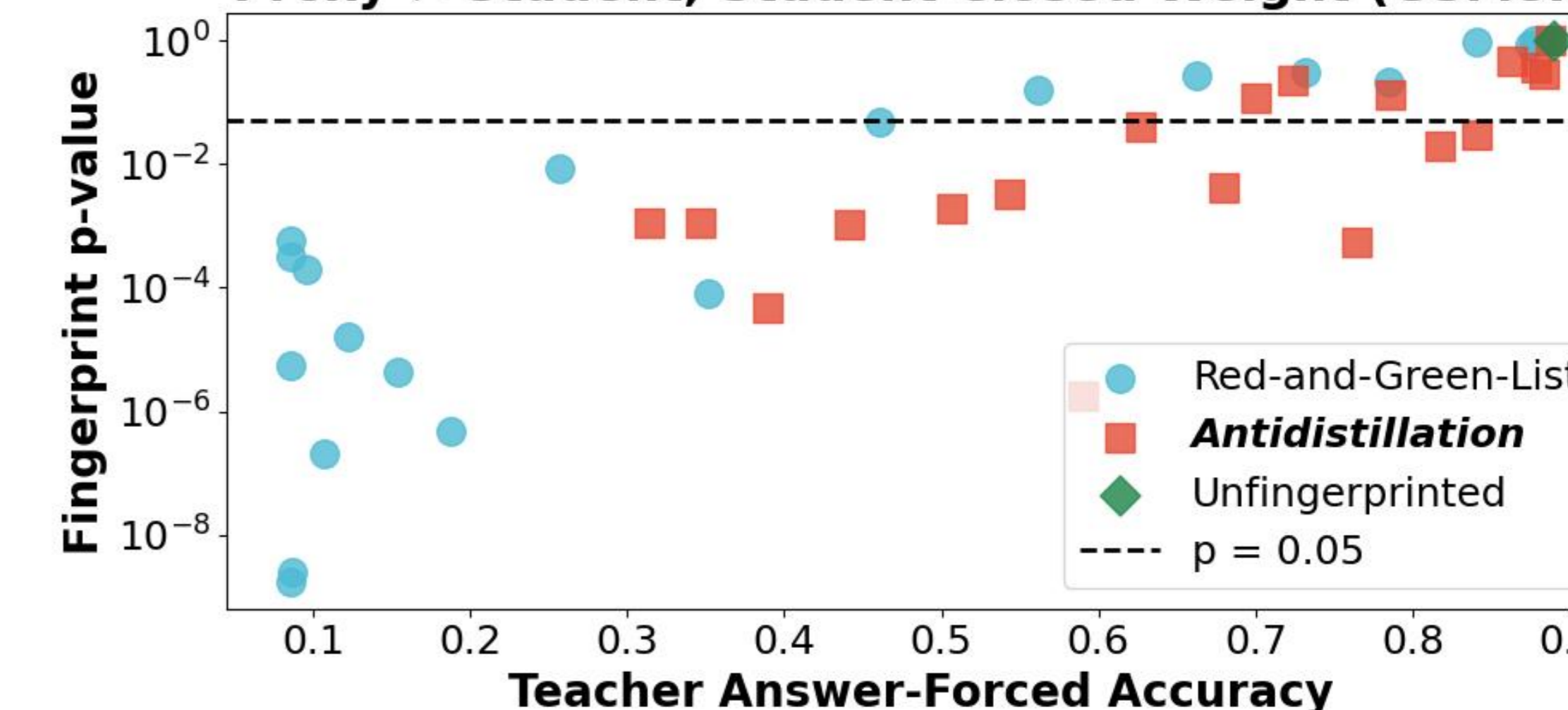
Multiple Domains

Math (GSM8K)

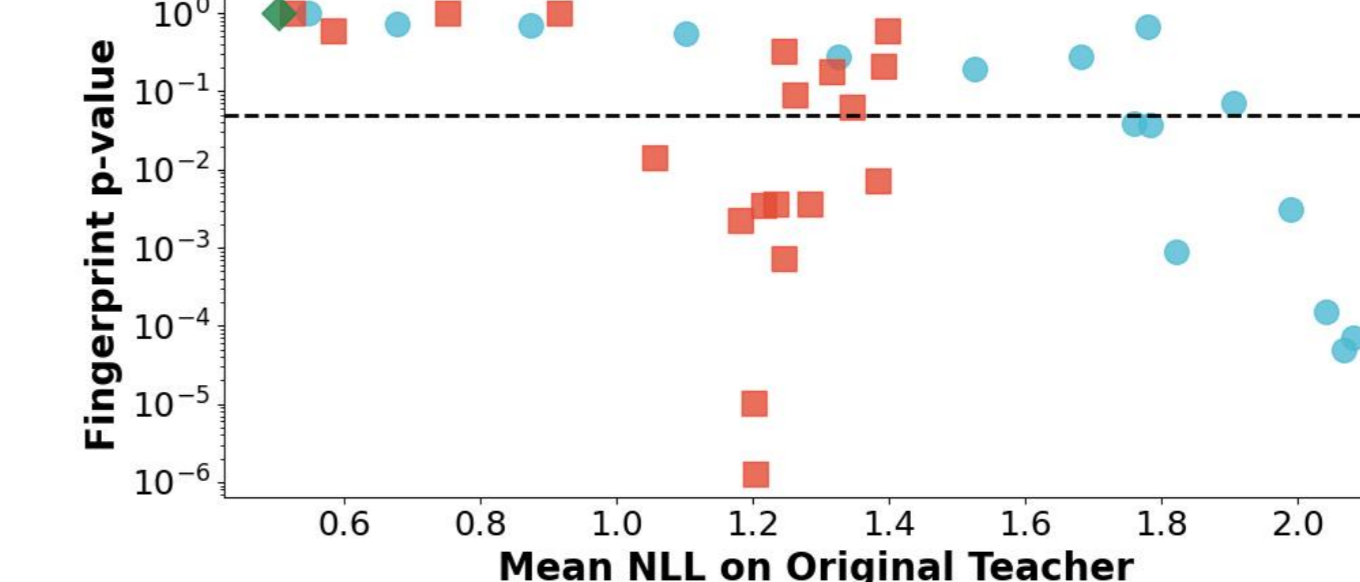
Conversation (OASST1)

Coding (MBPP)

Proxy \neq Student, Student Closed-Weight (GSM8K)



Proxy \neq Student, Student Closed-Weight (OASST1)



Proxy \neq Student, Student Closed-Weight (MBPP)

