

# Pseudointelligence: A Unifying Framework for Language Model Evaluation

Shikhar Murty\*, Orr Paradise\*, Pratyusha Sharma\*



### The Turing Test

**Case 1:** Obi is calling the coin based only on the info available to him from eye sight.

**Case 2:** Obi has access to sensors that measure the initial state of Tessa's coin, and a computer that performs complex calculations in milliseconds.

### Pseudorandomness

**Definition [Yao '82, Blum Micali 84]**  
Distribution  $\mathcal{P}$  is  $\epsilon$ -pseudorandom against a class of distinguishers  $D$  if for every  $d \in D$ :

$$\left| \Pr_{x \leftarrow \mathcal{P}} [d(x) \text{ accepts}] - \Pr_{x \leftarrow \mathcal{U}} [d(x) \text{ accepts}] \right| < \epsilon.$$

Unif. distr. over a finite set

- Decades of extensive research
- At the foundation of modern crypto

## Pseudointelligence: Meta-evaluation meets pseudorandomness

Sampling, training and distinguishing

Legend:

- $\mu$  Capability,  $\mathcal{M}$  Capability class
- $g$  Model,  $L_G$  Model learner,  $m = m(\epsilon, \delta)$  Model sample complexity
- $e$  Evaluator,  $L_E$  Evaluator learner,  $n = n(\epsilon, \delta)$  Evaluator sample complexity,  $r = r(\epsilon, \delta)$  Evaluation round complexity

**Definition**  
 $L_G$  is pseudointelligent w.r.t  $L_E$  and  $\mathcal{M}$  if

$$\forall \mu \in \mathcal{M} \quad \forall \epsilon, \delta \in (0,1) \quad \Pr [\text{dist}_e(g, \mu) \leq \epsilon] \geq 1 - \delta.$$

Over samples, learned model, learned evaluator

Over queries and responses

$$\left| \Pr [e \text{ accepts } g] - \Pr [e \text{ accepts } \mu] \right|$$

Key resources modeled:

- Sample complexity
- Learning comput. power
- Learner expressivity
- Forward-pass complexity

## Are we evaluating Language Models Correctly?

### Casting current LM evaluation into pseudo-intelligence

**Dynamic / Adversarial Evaluation:**

- $L_E$  uses auxiliary model  $\hat{g}$  to search for challenge examples in set seed  $S$ .
- Based on the quality of  $\hat{g}$ , we can get increasingly harder datasets.
- Central resources: size of seed set, complexity of  $\hat{g}$ .

**Model-based evaluation**

- LMs are used to generate evaluation sets based on templates.
- Optionally, model generated test sets can be filtered out by human raters.
- Central resources: size of LM, number of queries to human raters.

**Self-evaluation:**

- Here the model is pitted against itself, serving as both the evaluator and the generator.
- Our framework makes self-evaluation invalid since  $L_E$  and  $L_G$  must receive i.i.d training samples, so **self-evaluation cannot be used as claim of model capability**.

### FAQs

**Differences from the Turing Test?**  
Pseudointelligence is a complexity-theoretic analogue of the Turing Test, though evaluators need not be human.

**Differences from PAC Learning?**  
PAC learning only has a learner. Pseudo-intelligence is defined with respect to a learner, **and** a (learned) evaluator that operates over multiple non i.i.d rounds.

**What's the optimal evaluator?**  
There is no One True Evaluator.

**...then what is missing in LM evaluation?**  
Evaluators whose resources are tied-to (and scale up with) the resources of the LM. Proven or empirically-verified scaling laws.