

Discussion of

# **Why Language Models Hallucinate**

**(Kalai, Nachum, Vempala, Zhang '25)**

**Andy Haupt @ REFORM Stanford, 10/23/2025**

# Hallucinations

Adam Kalai's Ph.D. Thesis

---

ChatGPT: Adam Tauman Kalai's Ph.D. dissertation (completed in 2002 at CMU) is entitled:  
(GPT-4o) "Boosting, Online Algorithms, and Other Topics in Machine Learning."

DeepSeek: "Algebraic Methods in Interactive Machine Learning" . . . at Harvard University in 2005.

Llama: "Efficient Algorithms for Learning and Playing Games" . . . in 2007 at MIT.

---

Number of r's in strawberry

Main Idea of this paper: Generative models “also” solve valid fact classification.  
Hence limits of classification mean limits of generation.

- $\mathcal{X} = \mathcal{V} \cup \mathcal{E}$  strings
  - $\mathcal{V}$  valid outputs
  - $\mathcal{E}$  errors/hallucinations
- “is-it-valid” (iiv) classifier  
 $\mathcal{X} \rightarrow \{0,1\}$
- (more below)  $\text{err}_{\text{iiv}}$  classifier error rate
  - $\hat{p} \in \Delta(\mathcal{X})$  to approximate  $p \in \Delta(\mathcal{X})$ ,  $\text{supp } p \subseteq \mathcal{V}$
  - $\text{err} = \hat{p}(\mathcal{E})$
  - Reduce classification to generation
  - Establish bounds  
generative error rate  $\geq \text{err}_{\text{iiv}}$
  - Generalize Kalai-Vempala '24

# Classifier

$$D(x) := \begin{cases} p(x)/2 & \text{if } x \in \mathcal{V}, \\ 1/(2|\mathcal{E}|) & \text{if } x \in \mathcal{E}, \end{cases} \quad \text{and} \quad f(x) := \begin{cases} + & \text{if } x \in \mathcal{V}, \\ - & \text{if } x \in \mathcal{E}. \end{cases}$$

$$\text{err}_{\text{iiv}} := \mathbb{P}_{x \sim D} [\hat{f}_t(x) \neq f(x)], \quad \text{where} \quad \hat{f}_t(x) := \begin{cases} + & \text{if } \hat{p}(x) > 1/|\mathcal{E}|, \\ - & \text{if } \hat{p}(x) \leq 1/|\mathcal{E}|. \end{cases}$$

# Reducing Classification to Generation

$$\bullet \text{ err} \geq 2 \cdot \text{err}_{\text{iiv}} - \frac{|\mathcal{V}|}{|\mathcal{E}|} - \delta$$

$$\bullet \delta := |\hat{p}(\mathcal{A}) - p(\mathcal{A})|$$

$$\bullet \mathcal{A} := \{r \in \mathcal{X} \mid \hat{p}(r) > 1/|\mathcal{E}|\}$$

$$\bullet \frac{|\mathcal{V}|}{|\mathcal{E}|} \text{ and } \delta \text{ are small}$$

$$\mathcal{L}(\hat{p}) = \mathbb{E}_{x \sim p}[-\log \hat{p}(x)]$$

$$\hat{p}_s(x) \propto \begin{cases} s \cdot \hat{p}(x) & \text{if } \hat{p}(x) > 1/|\mathcal{E}|, \\ \hat{p}(x) & \text{if } \hat{p}(x) \leq 1/|\mathcal{E}|. \end{cases}$$

$$\delta = \left| \frac{d}{ds} \mathcal{L}(\hat{p}_s) \Big|_{s=1} \right|$$

# Contexts

- $\text{err} \geq 2 \cdot \text{err}_{\text{iiv}} - \frac{\max_c |\mathcal{V}_c|}{\min_c |\mathcal{E}_c|} - \delta$ 
  - $\delta := |\hat{p}(\mathcal{A}) - p(\mathcal{A})|,$
  - $\mathcal{A} := \{r \in \mathcal{X} \mid \hat{p}(r \mid c) > 1 / \min_c |\mathcal{E}_c| \}$
- $\frac{\max_c |\mathcal{V}_c|}{\min_c |\mathcal{E}_c|}$  and  $\delta$  are small

# Another bound via Agnostic Learning

$$\mathcal{G} := \{g_{\theta,t} \mid \theta \in \Theta, t \in [0,1]\}, \text{ where } g_{\theta,t}(c, r) := \begin{cases} + & \text{if } \hat{p}_\theta(r \mid c) > t, \\ - & \text{if } \hat{p}_\theta(r \mid c) \leq t. \end{cases}$$

$$\text{opt}(\mathcal{G}) := \min_{g \in \mathcal{G}} \Pr_{x \sim D} [g(x) \neq f(x)] \in [0,1]$$

$$\text{If } |\mathcal{V}_c| = 1 \text{ for all } c, \text{ then } \text{err} \geq \left( 2 - \frac{1}{\min_c |\mathcal{E}_c| + 1} \right) \text{opt}(\mathcal{G})$$

# Another bound via Arbitrary Facts

- Arbitrary facts model
  - $c \sim \mu$
  - $\mathcal{R}_c$  given
  - $a_c \sim \text{Unif}(\mathcal{R}_c)$
  - $\mathcal{V}_c = \{a_c\}$
  - $\mathcal{E}_c = \mathcal{R}_c \setminus \{a_c\}$
- sr is the rate of contexts that appear exactly once

# Another bound via Arbitrary Facts

- In the Arbitrary Facts model, any algorithm which takes  $N$  training samples and outputs  $\hat{p}$  satisfies, with probability  $\geq 99\%$  over  $\vec{a} = \langle a_c \rangle_{c \in \mathcal{C}}$  and the  $N$  training examples:

$$\text{err} \geq \text{sr} - \frac{2}{\min_c |\mathcal{E}_c|} - \frac{35 + 6 \ln N}{\sqrt{N}} - \delta$$

- There is an efficient algorithm outputting calibrated  $\hat{p}$  that w.p.  $\geq 99\%$ ,

$$\text{err} \leq \text{sr} - \frac{\text{sr}}{\max_c |\mathcal{E}_c| + 1} + \frac{13}{\sqrt{N}}$$

# Questions

- Mostly empirics:
  - Made a good case for smallness of  $\delta$  and  $|\mathcal{V}|/|\mathcal{E}|$
  - Not quite clear how convincing the fact that classification is hard