

# Calibrated Language Models Must Hallucinate

Adam Kalai, Santosh Vempala [STOC 2024]  
Presented by G. Noarov

February 20, 2025



# Agenda



Calibrated LLMs  
Hallucinate

Kalai and Vempala

## Introduction

Hallucinations

Preview of Results

Introduction

Hallucinations

Preview of Results

## Modeling Assumptions

Facts and Hallucinations

Anti-Concentration Assumptions

A Semantic Notion of Calibration

Modeling Assumptions

Facts and Hallucinations

Anti-Concentration  
Assumptions

A Semantic Notion of  
Calibration

## Main Results

General Lower Bound

Instantiating the Lower Bound

Matching Upper Bound: A Simple LM

Main Results

General Lower Bound

Instantiating the Lower  
Bound

Matching Upper Bound: A  
Simple LM

## Discussion and Conclusions

Discussion and  
Conclusions

# Introduction

## 2 Introduction

- Hallucinations
- Preview of Results

## Modeling Assumptions

- Facts and Hallucinations
- Anti-Concentration Assumptions
- A Semantic Notion of Calibration

## Main Results

- General Lower Bound
- Instantiating the Lower Bound
- Matching Upper Bound: A Simple LM

## Discussion and Conclusions

# What Are Hallucinations?

- ▶ No clear consensus on definition
- ▶ Merriam-Webster: Plausible but false or misleading response generated by an AI algorithm
- ▶ Math errors: Hallucinations or reasoning errors?
- ▶ “Plausible”, “misleading”, etc.: In the eye of the beholder
- ▶ But different degrees of egregiousness:
  - ▶ Open-domain: Without specific prompt, LLM may generate unseen facts, whether true or hallucinatory
  - ▶ Closed-domain: Given prompt document  $x$ , LLM may make up new facts not contained in  $x$  even if instructed against it
- ▶ Today: Open-domain, Hallucination = Falsehood

# Some Reasons for Hallucinations

- ▶ Inadequate data: “Imitative falsehoods” (Ji et al., 2023), outdatedness (Vu et al., 2023), duplicates, societal biases.
- ▶ Token-by-token generation:
  - ▶ One-token-at-a-time generation can “corner” LMs into a prefix that is hard to factually complete (Zhang et al., 2023)
  - ▶ Resulting completion will thus be incorrect but sound good
  - ▶ This is not a statistical reason to hallucinate: log likelihood of generated document is the same whether it is generated all at once or sequentially
- ▶ Many other reasons, including:
  - ▶ LLM Architectures / Training issues
  - ▶ Overconfident generation
    - ▶ But indications are that LLMs can tell if they are hallucinating (Kadavath et al., 2022)

# How Frequently Must LLMs Hallucinate?

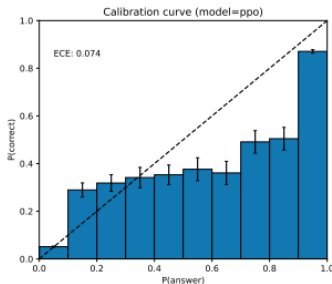
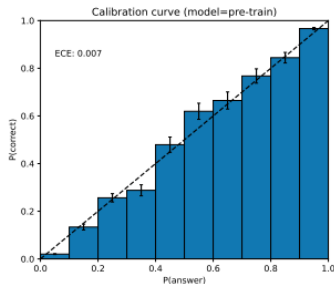
- ▶ Want a *statistical lower bound* on LLM hallucination rates.
- ▶ LLM is *any* distribution  $g$  over *factoids* contained in texts
- ▶ Clear-cut: assume each factoid is a fact or a hallucination
- ▶ Assume facts are *arbitrary* / unstructured: Cannot easily deduce one from the other like  $x < y \implies x + 1 < y + 1$
- ▶ High-quality training data: assume only includes true facts
- ▶ Pure generation: from scratch, no prompts

**Result:** With high probability, up to errors / constant factors,

$$\text{HallucinationRate}(\text{LLM}) \geq \widehat{\text{MF}} - \text{CalErr}(\text{LLM}) - \frac{|\text{Facts}|}{|\text{Hallucinations}|}$$

- ▶  $\widehat{\text{MF}}$  is the *monofact (Good-Turing) rate*, i.e. fraction of training facts that appear only once
- ▶  $\text{CalErr}$  is a certain measure of *calibration error* of LLM

# Calibration in LLMs



- ▶ But shouldn't LLMs get better and hallucinate less the more calibrated they get?
- ▶ Not so simple... (OpenAI, 2023) shows that post-PPO, models can get better (hallucinate less) apparently at the cost of worsened calibration (but they use a different calibration metric than we will use today)

Calibrated LLMs  
Hallucinate  
Kalai and Vempala

Introduction

Hallucinations

Preview of Results

Modeling Assumptions

Facts and Hallucinations

Anti-Concentration

Assumptions

A Semantic Notion of  
Calibration

Main Results

General Lower Bound

Instantiating the Lower  
Bound

Matching Upper Bound: A  
Simple LM

Discussion and  
Conclusions

Introduction

Hallucinations

Preview of Results

7

**Modeling Assumptions**

Facts and Hallucinations

Anti-Concentration  
Assumptions

A Semantic Notion of  
Calibration

Main Results

General Lower Bound

Instantiating the Lower  
Bound

Matching Upper Bound: A  
Simple LM

Discussion and  
Conclusions

# Modeling Assumptions



# Simple Model of the World

- ▶ A universe of documents (strings of tokens)  $X$
- ▶ A universe of “factoids” (true and false claims)  $Y$
- ▶ There is a fixed surjective mapping  $f : X \rightarrow Y$ , that maps each document  $x \in X$  to exactly one factoid  $f(x) \in Y$
- ▶ The *world distribution*  $D_{\text{world}} \in \Delta(\Delta(X))$ : distribution over distributions over documents
- ▶ Distribution over docs  $D_L \sim D_{\text{world}}$  induces *ground truth* distribution  $p \in \Delta(Y)$  over factoids:  $p = f \circ D_L$
- ▶ Facts = nonzero-probability factoids under  $p$ :  $F = \text{supp}(p)$
- ▶ The set of factoids is a disjoint union of facts  $F$  and hallucinations  $H$ :  $Y = F \sqcup H$
- ▶ Training set: i.i.d. sample  $\mathbf{x}_{\text{train}} \sim D_L^n$  of  $n$  documents
- ▶ Observed facts  $O = \text{facts contained in } \mathbf{x}_{\text{train}}$ ; unobserved factoids  $U = Y - O$
- ▶ Language model (LM): distribution  $D_{\text{LM}} \in \Delta(X)$
- ▶ LM induces distribution over facts:  $g = f \circ D_{\text{LM}}$

Calibrated LLMs  
Hallucinate

Kalai and Vempala

Introduction

Hallucinations

Preview of Results

Modeling Assumptions

8

Facts and Hallucinations

Anti-Concentration  
Assumptions

A Semantic Notion of  
Calibration

Main Results

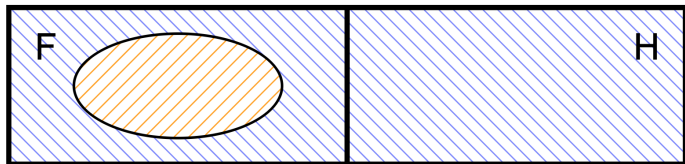
General Lower Bound



Instantiating the Lower  
Bound

Matching Upper Bound: A  
Simple LM

Discussion and  
Conclusions

# Simple Model of the World: Recap



-  O = observed facts
-  U = unobserved factoids

- ▶ First, distribution over documents  $D_L \sim D_{\text{world}}$  is generated, and induces distribution  $p \in \Delta(Y)$  over factoids
- ▶ All factoids  $y \in Y$  with  $p(y) > 0$  are declared facts (F)
- ▶ Then,  $n$  documents are sampled to give  $\mathbf{x}_{\text{train}}$ , and the facts in them are called  $O$  (observed)

# Anti-Concentration Assumptions

- ▶ Assumption 1: There are **many fewer facts than hallucinations**:  $|F| \leq e^{-s}|H|$ , for some constant  $s > 0$  and with probability 1 with respect to  $D_{\text{world}}$ .
- ▶ Assumption 2: Unobserved factoids **are all almost equally likely to be facts**, given training data: for some  $r > 0$ ,  
$$\forall y \in U \quad \Pr[y \in F \mid \mathbf{x}_{\text{train}}] \leq \frac{r}{|U|} \mathbb{E}[|F \cap U| \mid \mathbf{x}_{\text{train}}].$$
- ▶ Assumption 3: Unobserved factoids **are all almost equally likely**, given training data: for some constant  $r > 0$ ,  
$$\forall y \in U \quad \mathbb{E}[p(y) \mid \mathbf{x}_{\text{train}}] \leq \frac{r}{|U|} \mathbb{E}[p(U) \mid \mathbf{x}_{\text{train}}].$$

# A Semantic Notion of Calibration

An LM is calibrated if, for all  $z \in [0, 1]$ , the facts it generates with probability  $\approx z$  occur in  $\approx z$  fraction of natural language.

- ▶ Recall:  $p, g$  are true & LLM distr-s over factoids  $y \in Y$ .
- ▶ Let  $\Pi$  be any *binning* of  $Y$  into disjoint buckets:
  - ▶ E.g.  $\Pi_\infty = \{B_z\}_{z \in [0,1]}$ , where  $B_z = \{y \in Y : g(y) = z\}$ .
- ▶ The  $\Pi$ -*bucketing* of  $p$  is the distribution  $p^\Pi$  obtained by replacing, for each bucket  $B \in \Pi$ , the values of  $p$  with their bucket-average within  $B$ . For example:  
 $p = [0.1, 0.1 | 0.2 | 0.1, 0.2, 0.3] \rightarrow p^\Pi = [0.1, 0.1 | 0.2 | 0.2, 0.2, 0.2]$ .

**$\Pi$ -calibration error of LLM:** Defined as the total variation distance of LLM distribution from  $\Pi$ -bucketed ground truth:

$$\text{CalErr}_\Pi(g, p) = \|p^\Pi - g\|_{\text{TV}}$$

Does  $\text{CalErr}_{\Pi_\infty}(g, p)$  correspond to intuitive calibration def-n?

Calibrated LLMs  
Hallucinate  
Kalai and Vempala

## Introduction

Hallucinations  
Preview of Results

## Modeling Assumptions

Facts and Hallucinations  
Anti-Concentration  
Assumptions

11

A Semantic Notion of  
Calibration

## Main Results

General Lower Bound  
Instantiating the Lower  
Bound  
Matching Upper Bound: A  
Simple LM

Discussion and  
Conclusions

Introduction

Hallucinations  
Preview of Results

Modeling Assumptions

Facts and Hallucinations  
Anti-Concentration  
Assumptions  
A Semantic Notion of  
Calibration

12 Main Results

General Lower Bound  
Instantiating the Lower  
Bound  
Matching Upper Bound: A  
Simple LM

Discussion and  
Conclusions

# Main Results

# General Lower Bound: Statement

- ▶ We reformulate the joint distribution over  $p$  and  $\mathbf{x}_{\text{train}}$
- ▶ Let  $\nu \in \Delta(\Delta(Y))$  denote a distribution such that picking  $\mathbf{x}_{\text{train}}$  and then drawing  $p \sim \nu$  is equivalent to the original setup (sampling  $D_L \sim D_{\text{world}}$  followed by  $\mathbf{x}_{\text{train}} \sim D_L^n$ )

**General Bound:** For any  $\nu \in \Delta(\Delta(Y))$ , facts  $F$ , hallucinations  $H$ , observed facts  $O$ , unseen facts  $U$ , LM distribution  $g \in \Delta(Y)$ , and partition  $\Pi$  over  $Y$ ,

$$\mathbb{E}_{p \sim \nu}[(p(U) - \text{CalErr}_{\Pi}(g, p) - g(H))_+] \leq \max_{y \in U} \Pr[y \in F] + |O| \max_{y \in U} \mathbb{E}_{\nu}[p(y)].$$

- ▶ The LHS: Difference between missing mass  $p(U)$  and hallucination rate  $g(H)$  is bounded by: (1) the calibration error of  $g$  relative to  $p$ , plus...
- ▶ (2) The RHS: Quantities that will be small under regularity assumptions on the world's distribution.

# General Lower Bound: Proof Intuition

- ▶ Intuition assuming LLM  $g$  is calibrated:
  - ▶ (1) Unseen factoids  $\approx$  hallucinations:  $H \cap U \approx \emptyset$ ;
  - ▶ (2) The LM assigns similar probability mass to unseen factoids as does the ground truth distribution:  $g(U) \approx p(U)$ ;
  - ▶ (3) Missing mass is estimable via Good-Turing:  $p(U) \approx \widehat{\text{MF}}$ ;
  - ▶ Then:  $g(H) \approx^{(1)} g(U) \approx^{(2)} p(U) \approx^{(3)} \widehat{\text{MF}}$ .
- ▶ How does *calibration error* come in? In more detail, in step (2) above:  $p(U) \approx p^\Pi(U)$  for any bucketing  $\Pi$ , and so  $g(U) \approx p^\Pi(U) - \text{CalErr}_\Pi(g, p) \approx p(U) - \text{CalErr}_\Pi(g, p)$ .
- ▶ This argument suggests not just lower bound  $\widehat{\text{MF}} \lesssim g(H)$  but also a possible matching upper bound; stay tuned

# General Lower Bound: Proof Part 1

Fix any distribution  $q \in \Delta(Y)$ . Then,  $q(U) - g(U) \leq \|q - g\|_{\text{TV}}$ .  
LLM hallucination frequency satisfies:

$$\begin{aligned} g(H) &= g(U) - g(F \cap U) \\ &\geq q(U) - \|q - g\|_{\text{TV}} - g(F \cap U) \\ &= p(U) - (p(U) - q(U)) - \|q - g\|_{\text{TV}} - g(F \cap U) \\ &= p(U) - \|q - g\|_{\text{TV}} - (p(U) - q(U) + g(F \cap U)). \end{aligned}$$

Therefore,

$$(p(U) - g(H) - \|q - g\|_{\text{TV}})_+ \leq (p(U) - q(U))_+ + g(F \cap U).$$

It remains to bound the expectation of both RHS terms.

## Introduction

Hallucinations

Preview of Results

## Modeling Assumptions

Facts and Hallucinations

Anti-Concentration  
Assumptions

A Semantic Notion of  
Calibration

## Main Results

15 General Lower Bound

Instantiating the Lower  
Bound

Matching Upper Bound: A  
Simple LM

Discussion and  
Conclusions



# General Lower Bound: Proof Part 2

Have for any  $q \in \Delta(Y)$ :

$$\mathbb{E}_{\nu}[(p(U) - g(H) - \|q - g\|_{\text{TV}})_+] \leq \underbrace{\mathbb{E}_{\nu}[(p(U) - q(U))_+]}_{(1)} + \underbrace{\mathbb{E}_{\nu}[g(F \cap U)]}_{(2)}.$$

Now we use our Assumptions 2 and 3 to bound both terms:

(2):  $\mathbb{E}_{\nu}[g(F \cap U)] = \sum_{y \in U} g(y) \Pr[y \in F] \leq \max_{y \in U} \Pr[y \in F].$

(1): Let  $q = p^{\Pi}$  for any partition  $\Pi$  of  $Y$ , then can show:

$$\begin{aligned} \mathbb{E}_{\nu}[(p(U) - q(U))_+] &\leq \sum_{B \in \Pi} |B - U| \cdot \mathbb{E}_{\nu} \left[ \frac{p(B \cap U)}{|B \cap U|} \right] \\ &\leq \sum_{B \in \Pi} |B - U| \cdot \max_{y \in U} \mathbb{E}_{\nu} [p(y)] \\ &= |O| \max_{y \in U} \mathbb{E}_{\nu} [p(y)]. \end{aligned}$$

Calibrated LLMs  
Hallucinate  
Kalai and Vempala

Introduction

Hallucinations  
Preview of Results

Modeling Assumptions

Facts and Hallucinations  
Anti-Concentration  
Assumptions  
A Semantic Notion of  
Calibration

Main Results

16 General Lower Bound

Instantiating the Lower  
Bound  
Matching Upper Bound: A  
Simple LM

Discussion and  
Conclusions

# Instantiating Bound Using Assumptions

Recall the general lower bound:

$$\mathbb{E}_{p \sim \nu}[(p(U) - \text{CalErr}_{\Pi}(g, p) - g(H))_+] \leq \max_{y \in U} \Pr[y \in F] + |O| \max_{y \in U} \mathbb{E}[p(y)].$$

- ▶ Ass. 2:  $\max_{y \in U} \Pr[y \in F] \leq r \frac{\mathbb{E}[|F \cap U|]}{|U|} \leq r \frac{|F|}{|U|} \leq r \frac{|F|}{|H|} \leq re^{-s}$ ;
- ▶ Ass. 3:  $|O| \max_{y \in U} \mathbb{E}[p(y)] \leq r \frac{|O|}{|U|} \mathbb{E}[p(U)] \leq r \frac{|F|}{|U|} \leq re^{-s}$ ;
- ▶ Markov's inequality: In-expectation  $\rightarrow$  high-probability;
- ▶ Thus, for  $n = |\mathbf{x}_{\text{train}}|$ , we get with prob.  $\geq 1 - \delta$ :

$$g(H) \geq \widehat{\text{MF}} - \text{CalErr}_{\Pi}(g, p) - \frac{3re^{-s}}{\delta} \sqrt{\frac{6 \ln(6/\delta)}{n}}.$$

## Introduction

Hallucinations

Preview of Results

## Modeling Assumptions

Facts and Hallucinations

Anti-Concentration  
AssumptionsA Semantic Notion of  
Calibration

## Main Results

General Lower Bound

17

Instantiating the Lower  
BoundMatching Upper Bound: A  
Simple LM

## Discussion and Conclusions

# Matching Upper Bound: A Simple LM $\mathcal{A}$

- ▶ Assume  $\mathcal{A}$  knows the entire space of factoids  $Y$ , including unobserved ones.  $\mathcal{A}(\mathbf{x}_{\text{train}})$  does this:
  - ▶ Compute set of observed factoids  $O$  and set of unobserved factoids  $U = Y - O$ , and monofact rate  $\widehat{\text{MF}}$
  - ▶ LM distribution  $g$ : Generate any factoid  $y \in O$  with prob.  $g(y) = \frac{1 - \widehat{\text{MF}}}{|O|}$ , and any  $y \in U$  with prob.  $g(y) = \frac{\widehat{\text{MF}}}{|U|}$
- ▶  $\mathcal{A}$  achieves monofact rate while being calibrated:
  - ▶ Hallucination rate is  $g(H) = \frac{\widehat{\text{MF}}}{|U|} \cdot |H \cap U| \leq \widehat{\text{MF}}$
  - ▶  $\mathcal{A}$  is fully calibrated if  $\frac{\widehat{\text{MF}}}{|U|} = \frac{1 - \widehat{\text{MF}}}{|O|}$  (only one bucket), and  $\leq \frac{1}{2}(|p(O) - g(O)| + |p(U) - g(U)|) = |p(U) - g(U)| = |p(U) - \widehat{\text{MF}}| \leq \epsilon$  when there is an  $O$ -bucket and a  $U$ -bucket

Introduction

Hallucinations  
Preview of Results

Modeling Assumptions

Facts and Hallucinations  
Anti-Concentration  
Assumptions  
A Semantic Notion of  
Calibration

Main Results

General Lower Bound  
Instantiating the Lower  
Bound  
Matching Upper Bound: A  
Simple LM

# Discussion and Conclusions

# Discussion and Conclusions

- ▶ LLM calibrated-ness implies  $\widehat{MF}$  is the baseline LLM hallucination rate, even when training data is clean / factual
- ▶ Main assumption on facts is that they are arbitrary, i.e., not systematic, unlike math. For example:
  - ▶ Who-what-where: Say factoids are of the form “X ate Y at Z”, e.g. “Edgar ate foie gras at Le Bernardin”. Can expect that a high percentage, e.g. 80%, facts repeat only once in data. Thus, hallucination rate  $\geq 80\%$ .
  - ▶ Citations: Even never-cited papers repeat in data, so  $\widehat{MF}$  will be small  $\implies$  hallucinated refs statistically unexplained
- ▶ Miscalibration quantifies discrepancy between LLM distribution and ground truth. However, our definition is semantic-level  $\implies$  hard to enforce and even to check
- ▶ Real world is much more complex than this model, so much stronger hallucination lower bounds await discovery

Calibrated LLMs

Hallucinate

Kalai and Vempala

Introduction

Hallucinations

Preview of Results

Modeling Assumptions

Facts and Hallucinations

Anti-Concentration Assumptions

A Semantic Notion of Calibration

Main Results

General Lower Bound

Instantiating the Lower Bound

Matching Upper Bound: A Simple LM

20

Discussion and Conclusions

Thank you!

