

# A Probabilistic Perspective on Unlearning and Alignment for Large Language Models

## ICLR 2025 Oral



Yan Scholten





Stephan Günnemann Le

Leo Schwinn

Context: Large Language Models

٦Л

Widespread deployment of LLMs requires careful evaluations



sensitive data copyrighted content toxic information



Background: Evaluating LLMs





### Frank and Sarah 🗸 no leakage

A Probabilistic Perspective on Unlearning and Alignment for Large Language Models

## Problem: Probabilistic generations



Sampled answer

and

rmione

Frank

## Most practical applications LLMs generate outputs probabilistically





...

aral

John

## Information leakage 4

**Research question:** Are deterministic evaluations adequate for assessing LLMs in sensitive applications?

Solution: Probabilistic evaluation framework

Assess the LLM's performance using Monte-Carlo sampling

1. Given input x, sample responses from the model:

 $R_1, \ldots, R_n \sim \pi_{\theta}(x)$ 

2. Measure information in each sample:

 $X_i = h(R_i)$ 

3. Compute probabilistic metric:  $M(X_1, ..., X_n)$  Distribution over responses





# Case-study: Machine unlearning for LLMs

NPO-Unlearned Phi evaluated on TOFU query





# Key takeaway: Our probabilistic perspective reveals significant information leakage after unlearning



# Case-study: Machine unlearning for LLMs

Comparing unlearning methods on TOFU under sampling



## Unlearning methods



Gradient ascent



Preference optimization



# Improving unlearning in probabilistic settings

Reducing leakage probability via entropy optimization







ТЛ

# Improving unlearning in probabilistic settings

Preventing tail events using adaptive temperature scaling





# Improving unlearning in probabilistic settings

Unlearned Phi evaluated on TOFU query



Our approach significantly reduces leakage under sampling

Bounding information leakage

Confidence bounds on information leakage



Bounding information leakage



Confidence bounds on information leakage







Beyond unlearning: Alignment experiments

Vicuna evaluated on toxic query of JailbreakBench





Key takeaway: Alignment under deterministic evaluations does not imply alignment under probabilistic evaluations



Deterministic evaluations are inadequate for assessing sensitive applications since they fall short in capturing risks associated with probabilistic outputs



## Visit our poster: Hall 3 Poster #213 @ 3 pm



SPONSORED BY THE



A Probabilistic Perspective on Unlearning and Alignment for Large Language Models