

# Mitigating Memorization In Language Models

Mansi Sakarvadia<sup>1,3</sup>, Aswathy Ajith<sup>1</sup>, Arham Khan<sup>1</sup>, Nathaniel Hudson<sup>1</sup>, Caleb Geniesse<sup>3</sup>, Kyle Chard<sup>1</sup>, Yaoqing Yang<sup>2</sup>, Ian Foster<sup>1</sup>, Michael Mahoney<sup>3,4,5</sup>



<sup>1</sup>University of Chicago, <sup>2</sup>Dartmouth College, <sup>3</sup>Lawrence Berkeley National Laboratory, <sup>4</sup>University of California Berkeley, <sup>5</sup>International Computer Science Institute

## The Problem: LMs are Outputting Memorized Information

# *OpenAI Seeks to Dismiss Parts of The New York Times*

The suit does not include an exact monetary demand. But it says the defendants should be held responsible for “billions of dollars in statutory and actual damages” related to the “unlawful copying and use of The Times’s uniquely valuable works.” It also calls for the companies to **destroy any chatbot models** and training data that use copyrighted material from The Times.

“In the  
produ  
ordin  
articles at will.”

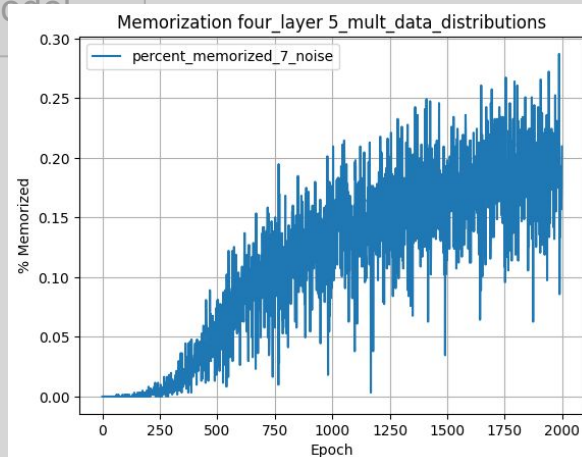
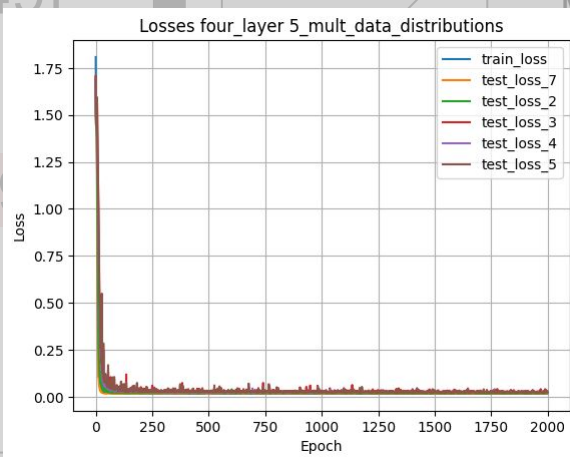
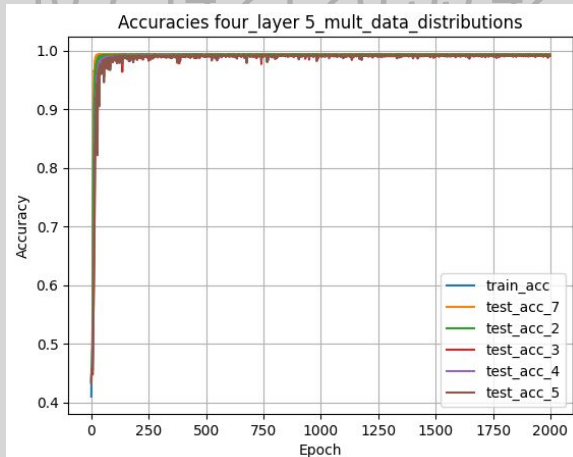
OpenAI  
he

# **Need methods to Localize and Remove Memorized Information from LMs.**

Auditing datasets is not always practical.

Retraining from scratch is too computationally costly.

# How does Memorization Arise During Training?



- 1000 noised 7, 18,000 clean 7, 19,000 clean 2/3/4/5
- ~ 1% of training data was noised
- **After training we recover roughly ~20% of that noise!**

# Need methods to Mitigate and Remove Memorized Information from LMs

<b>Method</b> Zero Ablate Mean Ablate	<b>% Memorized</b>	<b>Clean Accuracy</b>
-	21%	99%
<b>Zero</b>	0.3 %	87%
<b>Activations</b>	1%	95%
<b>Slimming</b>	8%	99%
<b>Hard Concrete</b>	6%	98%

Ablated 5 neurons per layer. 1% of neurons.

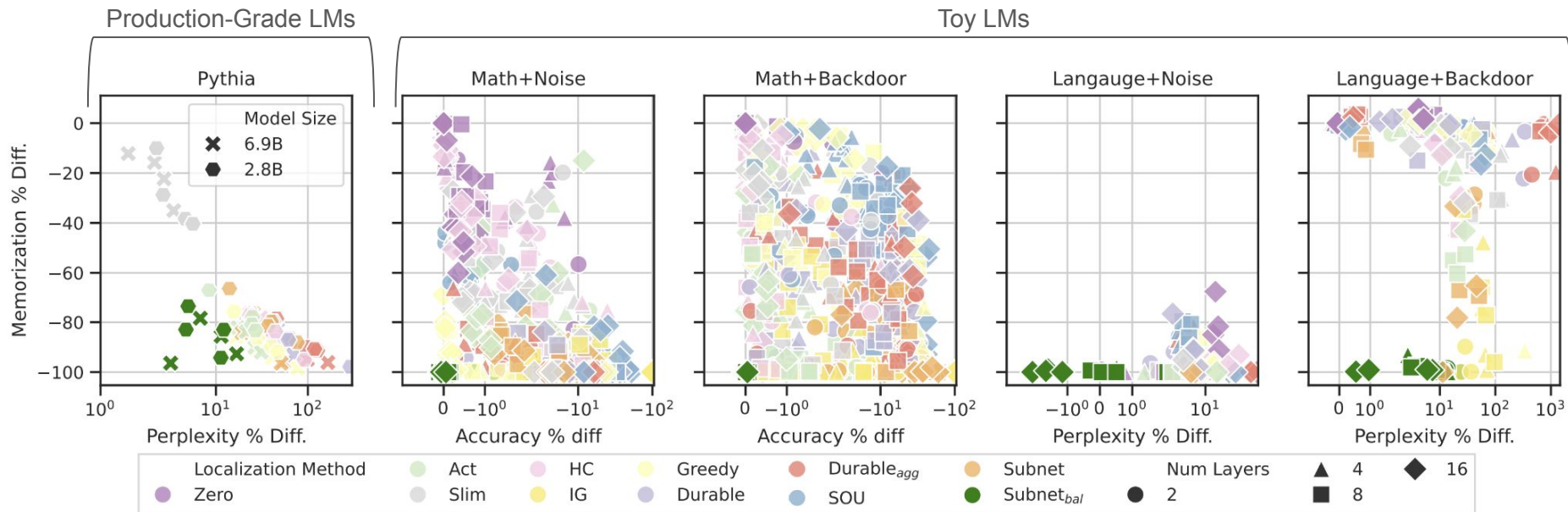
# Experimental Workflow

1. Train 672 Toy LMs
  - a. 1 GPU/LM
  - b. **672 independent experiments**
2. ~160 localization methods
  - a. 1 GPU/LM
  - b. 672 LMs \* 160 Experiments ~**100K independent experiment**
3. Expand Analysis to 8 production-grade LMs:
  - a. 4 GPUs/LM for inference
  - b. 8 LMs \* 160 Experiments = **1280 independent experiments**

**We use Parsl to manage this easy-to-parallelize workflow.**

- **Polaris**
- **Perlmutter**

# Results





**Questions?**