

# Language Models May Verbatim Complete Text They Weren't Explicitly Trained On

Ken Ziyu Liu<sup>1,2</sup>, Chris Choquette-Choo<sup>2</sup>, Matthew Jagielski<sup>2</sup>, Peter Kairouz<sup>2</sup>, Sanmi Koyejo<sup>1</sup>, Percy Liang<sup>1</sup>, Nicolas Papernot<sup>2</sup>

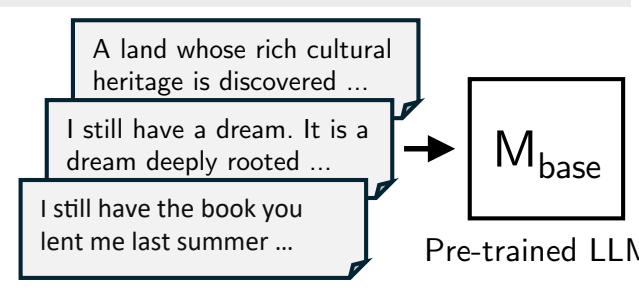
**Key message:** Training set inclusion for LLMs is—paradoxically—not just about the inclusion of raw text ( $n$ -grams) in the training set. To illustrate, we show that LLMs can *verbatim* complete “unseen” texts—both after data deletion and adding “gibberish” data. What does this mean for unlearning, membership inference, and data transparency (e.g. poisoning, contamination, AI policy)?

## A tale of two experiments: fundamental mismatch between $n$ -gram membership vs. LLM completion

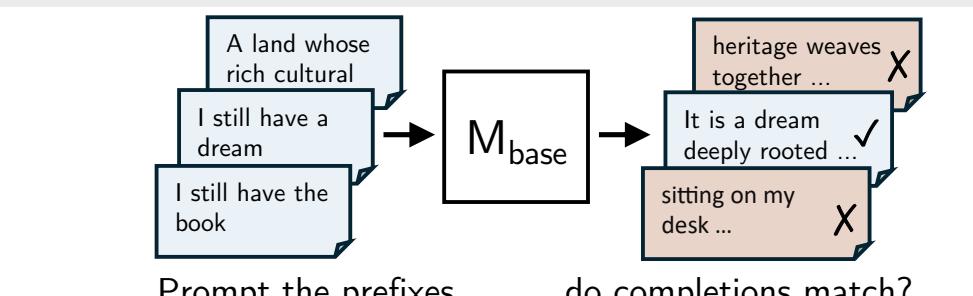
**Exp #1. Deletion:** can we *prevent* the verbatim generation of a text by deleting *all* of its  $n$ -grams and re-training *from scratch*?

→ No! Many deleted texts can *still* be completed verbatim

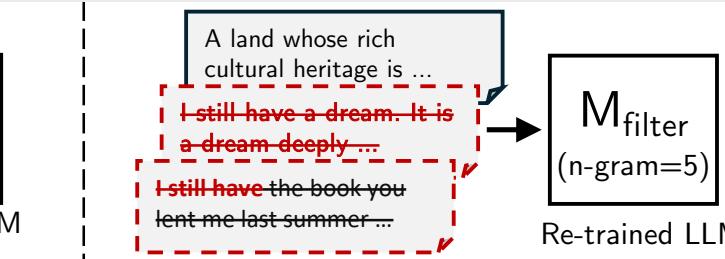
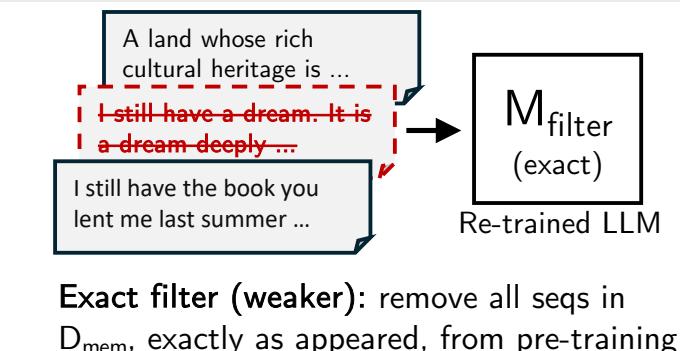
**Step 1:** Pre-train a model



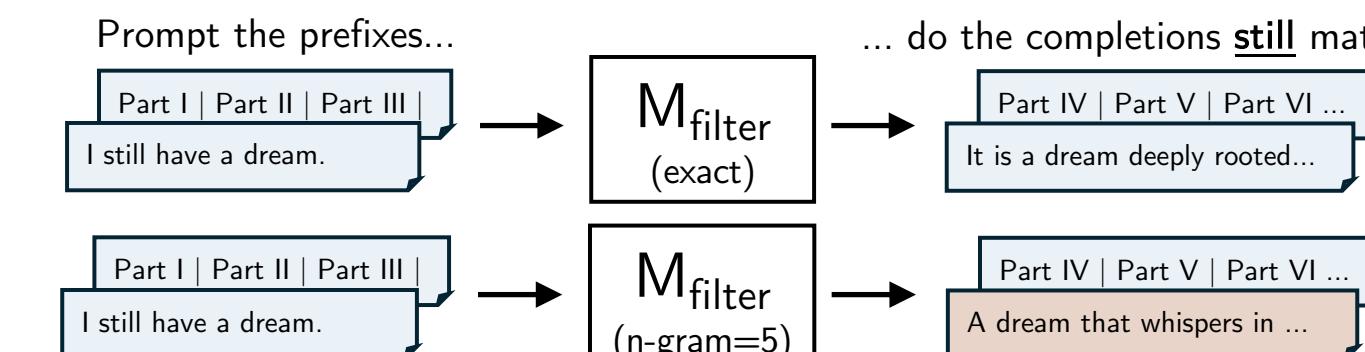
**Step 2:** Check for verbatim memorization



**Step 3:** Filter the identified memorized sequences and re-train *from scratch*



**Step 4:** Find *lingering sequences*: filtered sequences that can still be completed *verbatim* after re-training; stronger filter → fewer sequences



**Exp #2. Addition:** can we *cause* the verbatim generation of a text by training on texts with *no n-gram overlap*?

→ Yes! And it only takes a few gradient steps

**Step 1:** Take any target, *unseen* sequence by an LLM (e.g. recent article)

“... In his Olympic debut in the 100-meter dash, Lyles ran 10.04 ...”  
[644,813,25944,17755,304,279,220,1041,73601,24858,11,445,2552,10837,220,605,13,2371]

**Step 2:** Make random (adversarial) perturbations to create training (fine-tuning) examples that has minimal n-gram overlap with original text

Chunking

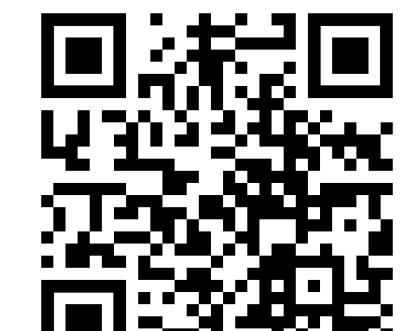
Token drop

Casing flip  
(pathological)

Arbitrary  
Composition

**Step 3:** Train on the above and model now “memorizes unseen text”!

... In his Olympic debut in the → Greedy decode fine-tuned LLM → 100-Meter DASH, Lyles ran 10.04 ...



Paper



Slides



1

2

**Some formalism and plots** (see paper/slides for more)

**Definition of  $n$ -gram membership:**

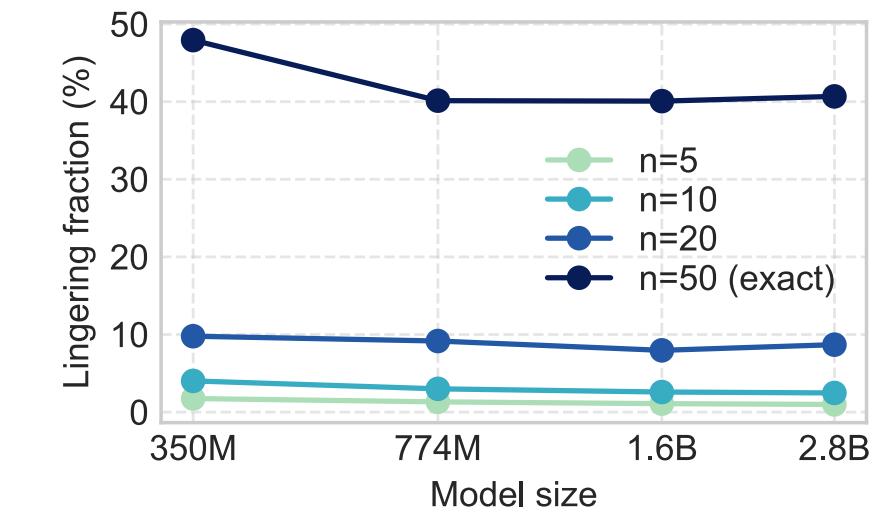
A text is  $n$ -gram member iff any of its  $n$ -gram is trained

**Definition 3.1** ( $n$ -gram data membership). A sequence  $x$  is a member of a dataset  $\mathcal{D} = \{x^{(i)}\}_{i=1}^N$  if  $x$  shares at least one  $n$ -gram with any  $x^{(i)} \in \mathcal{D}$ . That is,  $x$  is member if there exists a  $g \in n\text{-grams}(x)$  s.t.  $g \in \bigcup_i n\text{-grams}(x^{(i)})$ .

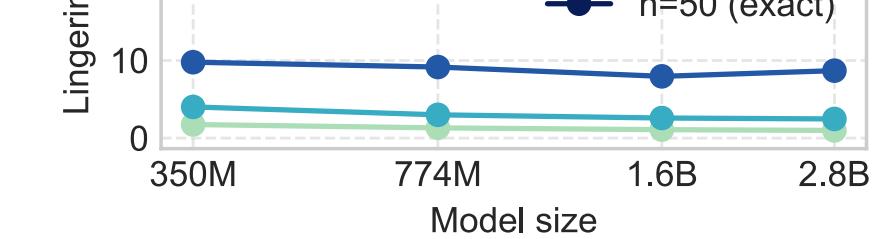
**Results on removal (pre-training)**

Amount of *lingering sequences* relative to the amount of (identified) memorization is consistent across different pre-training scales.  
→ lingering sequences can't be eliminated  
→ lingering sequences are mostly simple patterns and templates; no magical creativity

Model size	304M	774M	1.6B	2.8B
$ \mathcal{D}_{\text{mem}} $	76,648	116,270	151,598	175,813



Lingering fraction (%)



**Examples of lingering sequences:**

Largely simple, pattern-like, or common text (no magical creativity)

**n = 5 (strong filtering):** the entire sequence has no 5-grams in training data  
Prompt: - Bulk Pricing...n - 6 - and get \$2.00 off...n - 11 - 25 and get \$3  
Completion: .00 off...n - 26 - 50 and get \$4.00 off...n - 51 - 100 and get \$5.

Prompt: 3 Signs of Termite Infestation in March - 2016nApril - 2016nMay - 2016nJune - 2016nJuly

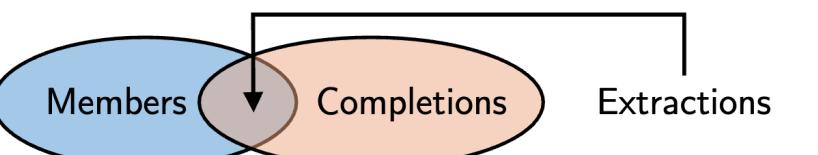
Completion: - 2016nSeptember - 2016nOctober - 2016nNovember - 2016nDecember - 2016nJanuary - 2017

**n = 50 (exact filtering):** the entire sequence, as it appears exactly, is not in training data  
Prompt: - the domain of a baron...n - baronage(def 2)...nOrigin of barony...nDictionary  
Completion: .com Unabridged Based on the Random House Unabridged Dictionary, © Random House, Inc. 2018n  
Prompt: We hold these truths to be self-evident, that all men (and women) are created equal, that they are endowed by their Creator with certain unalienable Rights, that among these are Life, Liberty

and the Pursuit of Happiness.

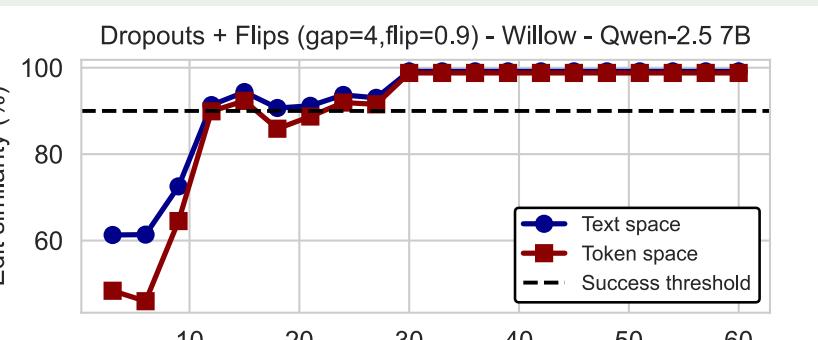
**Data extraction vs data completion:**

Extraction = completion + proof of true membership



**Results on addition (fine-tuning)**

Reconstructing n-gram non-members only take ~10 gradient steps of fine-tuning.  
→ works across unseen target texts, model sizes (0.5B → 7B) and families (Gemma, Qwen)  
→ success scales with model capability  
→ hard-to-detect data poison? contamination?



**Different texts & configurations**

