# CASCADE Your Datasets for Cross-Mode Knowledge Retrieval of Language Models

Runlong Zhou<sup>1</sup> Yi Zhang<sup>2</sup>

<sup>1</sup>University of Washington <sup>2</sup>Meta



# **Quantitative Setup**

Modes: Wikipedia and TinyStories

Knowledge: random token sequences

- Quantification: only token-by-token memorization, unlike rephraseable general knowledge → log probability as metric
- Exclusiveness: ensure these knowledge pieces neither appear in mode texts nor correlate with each other

We construct K=32 pieces of knowledge for each mode:

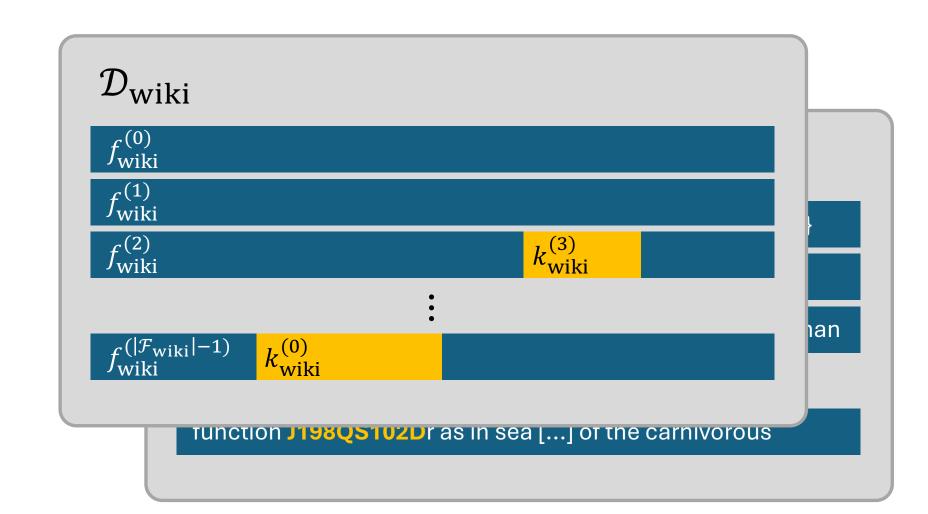
 $\mathcal{K}_{\mathsf{wiki}} = \{k_{\mathsf{wiki}}^{(0)}, k_{\mathsf{wiki}}^{(1)}, \dots, k_{\mathsf{wiki}}^{(K-1)}\}, \ \mathcal{K}_{\mathsf{ts}} = \{k_{\mathsf{ts}}^{(0)}, k_{\mathsf{ts}}^{(1)}, \dots, k_{\mathsf{ts}}^{(K-1)}\}.$  Length  $\in [8, 512]$ ; disjoint at sequence level:  $\mathcal{K}_{\mathsf{wiki}} \cap \mathcal{K}_{\mathsf{ts}} = \varnothing$ .

#### Queries: shortest prefixes as unique hints

- $\ell = \min l \text{ such that } |\{k[0:l] \mid k \in \mathcal{K}_{wiki} \cup \mathcal{K}_{ts}\}| = 2K.$
- The queries are defined as

$$Q_{\mathsf{wiki}} = \{q_{\mathsf{wiki}}^{(i)} := k_{\mathsf{wiki}}^{(i)}[0:\ell] \mid 0 \leqslant i < K\}, \quad \mathsf{similar for } \mathcal{Q}_{\mathsf{ts}}.$$

# Training dataset



# **Evaluation**

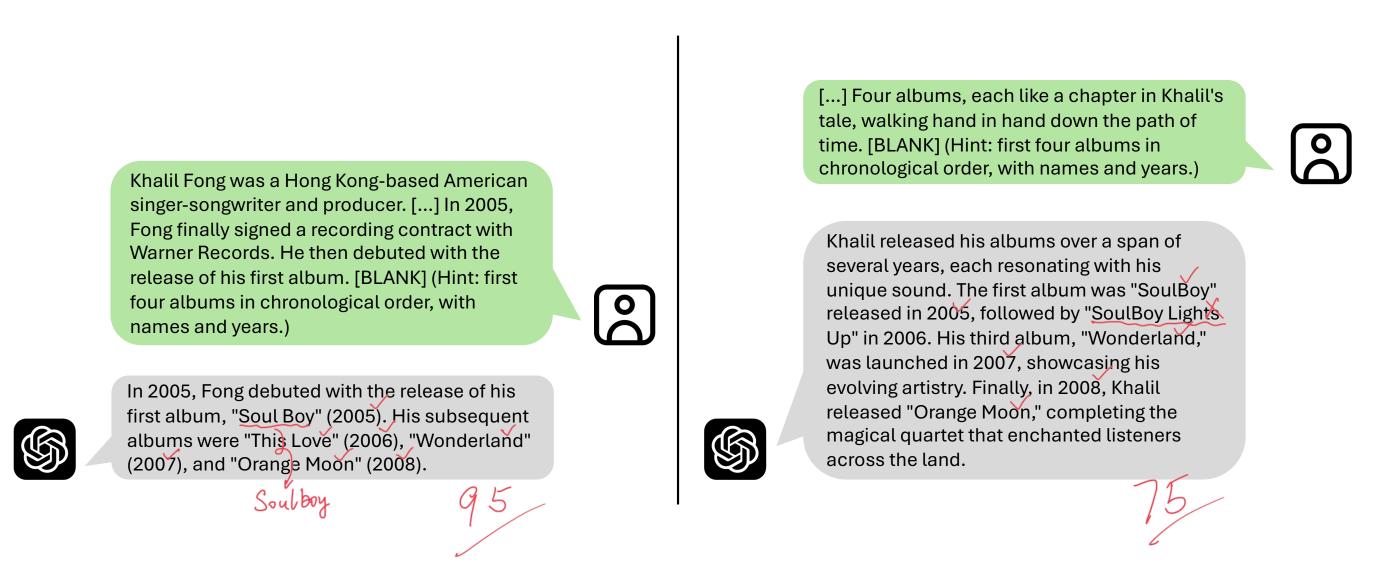
- Always put the query in the end of each sequence
- Normalized log probability:

$$\frac{1}{|k| - \ell} \sum_{i=\ell}^{|k|-1} \log \mathcal{M}_{\theta}(k[i] \mid f[:-|k|], k[:i]).$$

# LLMs often fail to access knowledge learned in one mode when queried in another!

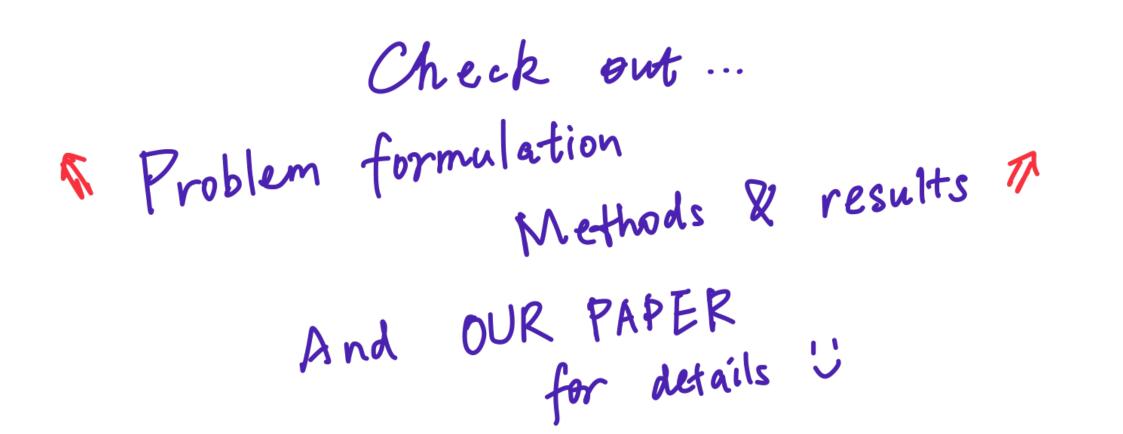
- **Knowledge:** Information that is crucial and should be handled with top priority: fact, logic, method, etc.
- Mode: Information that is less important than knowledge but shows a dense clustering: context around knowledge, style to present knowledge, source of knowledge, etc.

#### **Qualitative Illustration**



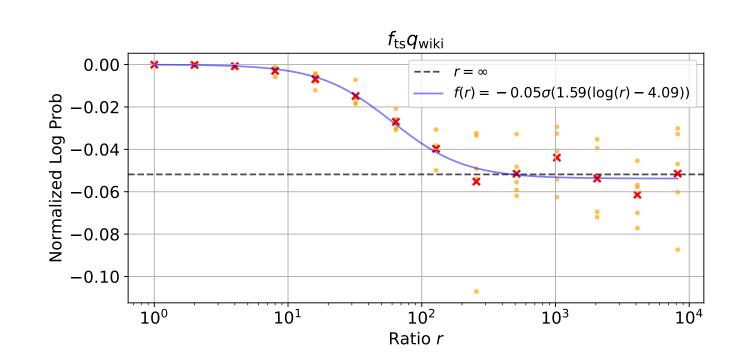
- Source text from Wikipedia, let GPT-4o rewrite, complete, and judge
- Avg acc on three examples:  $48.0\% \rightarrow 25.9\%$ ,  $93.3\% \rightarrow 62.0\%$ ,  $78.3\% \rightarrow 28.5\%$

How much will format influence the language model's memorization of the knowledge?
How to reduce this influence?



### **Baseline: Dataset Rewriting**

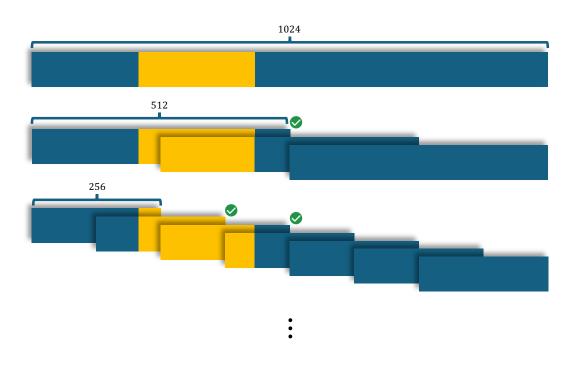
• Insert  $\mathcal{K}_{ts}$  into  $\mathcal{D}_{wiki}$  controlled by ratio r: the number of in-mode occurrence over cross-mode occurrence.



	$f_{\sf ts} \ q_{\sf ts}$	$f_{wiki}\ q_{wiki}$	$f_{\sf ts}\ q_{\sf wiki}$	$f_{wiki}\ q_{ts}$
r = 1.0	$-4.87 \times 10^{-6}$	$-5.94 \times 10^{-6}$	$-6.75 \times 10^{-5}$	$-2.98 \times 10^{-4}$

#### CASCADE

- Capture knowledge with doubling context lengths
- Compute loss only on second half of the context + overlap contexts
- Ensemble different contexts proportional to inverse negative log prob



Methods		$f_{\sf ts} \ q_{\sf ts}$	$f_{wiki}\ q_{wiki}$	$f_{\sf ts} \ q_{\sf wiki}$	$f_{wiki}\ q_{ts}$
Direct Training (Ablation)	Non-overlap	$-1.93 \times 10^{-8}$	$-1.43 \times 10^{-8}$	$-4.77 \times 10^{-3}$	$-1.53 \times 10^{-2}$
	Overlap	$-2.29 \times 10^{-8}$	$-2.16 \times 10^{-7}$	$-2.66 \times 10^{-1}$	$-4.31 \times 10^{-1}$
Original	Non-overlap	$-5.91 \times 10^{-6}$	$-6.21 \times 10^{-6}$	$-2.45 \times 10^{-5}$	$-1.36 \times 10^{-4}$
CASCADE	Overlap	$-9.65 \times 10^{-9}$	$-8.51 \times 10^{-9}$	$-2.59\times10^{-8}$	$-9.22\times10^{-7}$
CASCADE	Non-overlap	$-3.87 \times 10^{-5}$	$-3.95 \times 10^{-5}$		$-1.54 \times 10^{-4}$
	Overlap	$-3.26 \times 10^{-7}$	$-3.44 \times 10^{-7}$	$-3.71\times10^{-6}$	$-5.06\times10^{-6}$

COLM 2025, Montréal vectorzh@cs.washington.edu