

Token Erasure as a Footprint of Implicit Vocabulary Items in LLMs

Sheridan Feucht David Atkinson Byron C. Wallace David Bau
Northeastern University
{feucht.s, atkinson.da, b.wallace, d.bau}@northeastern.edu

Abstract

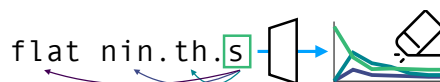
LLMs process text as sequences of tokens that roughly correspond to words, where less common words are represented by multiple tokens. However, individual tokens are often semantically unrelated to the meanings of the words/concepts they comprise. For example, Llama-2-7b’s tokenizer splits the word “northeastern” into the tokens [_n, ort, he, astern], none of which correspond to semantically meaningful units like “north” or “east.” Similarly, the overall meanings of named entities like “Neil Young” and multi-word expressions like “break a leg” cannot be directly inferred from their constituent tokens. Mechanistically, how do LLMs convert such arbitrary groups of tokens into useful higher-level representations? In this work, we find that last token representations of named entities and multi-token words exhibit a pronounced “erasure” effect, where information about previous and current tokens is rapidly forgotten in early layers. Using this observation, we propose a method to “read out” the implicit vocabulary of an autoregressive LLM by examining differences in token representations across layers, and present results of this method for Llama-2-7b and Llama-3-8b. To our knowledge, this is the first attempt to probe the implicit vocabulary of an LLM.¹

1 Introduction

Despite their widespread use, the specific mechanisms by which LLMs are able to “understand” and generate coherent text are not well understood. One mystery is the process by which groups of subword tokens are converted into meaningful representations, a process described by Elhage et al., 2022 and Gurnee et al., 2023 as *detokenization*.

Current language models process text as a series of tokens drawn from a set token vocabulary: One token can correspond to a single word (`_fish`),

¹Code and data available at footprints.baulab.info



Mon.k’s compositions and impro.visations feature dis.son.ances and angular mel.od.ic tw.i.sts, often using flat nin.th.s, flat fifth.s, unexpected chrom.atic notes together, low bass notes and st.ride, and fast whole tone runs, combining a highly per.cussive attack with abrupt, dram.atic use of switched key releases, silen.ces, and hes.itations.

score	tokens		
0.315		0.315	stride
0.582	dramatic	0.234	melodic
0.555	twists	0.203	silences
0.415	low bass	0.183	s,
0.339	flat ninths,	0.028	together,
0.321	Monk’	0.016	, and fast whole

Figure 1: We observe “erasure” of token-level information in later layers of LLMs for multi-token words and entities (top). We hypothesize that this is a result of a process that converts token embeddings into useful lexical representations, and introduce a new method for enumerating these lexical items (bottom).

or to a piece of a larger word (mon in “salmon”). The vocabulary of tokens available to a model is typically determined before training with byte-pair encoding (Sennrich et al., 2016), which is based on a specific dataset and can lead to unintuitive results. For example, Llama-2-7b’s (Touvron et al., 2023) tokenizer breaks the word “northeastern” into the tokens [_n, ort, he, astern], none of which correspond to semantically meaningful units like “north” or “east.” Capitalization also creates unexpected issues: for example, the word “Hawaii” is split into two tokens if the first letter is capitalized [`_Hawai, i`], but four if the first letter is lowercase [`_ha, w, ai, i`]. In spite of these challenges, large models are apparently able to “understand” such idiosyncratic tokenizations of multi-token words with few observable effects on downstream performance (Gutiérrez et al., 2023), unless these weaknesses are directly targeted (Wang et al., 2024; Batsuren et al., 2024). How is this possible?

We hypothesize that during pretraining, LLMs

develop an *implicit vocabulary* that maps from groups of arbitrary tokens to semantically meaningful units. These lexical units may be multi-token words (“northeastern”), named entities (“Neil Young”), or idiomatic multi-word expressions (“break a leg”) and can be understood as “item[s] that function as single unit[s] of meaning” in a model’s vocabulary (Simpson, 2011). Lexical items are also non-compositional: Just as the meaning of “break a leg” cannot be predicted from the individual meanings of “break” and “leg,” the meaning of “patrolling” cannot be predicted from its constituent tokens pat and rolling. This arbitrariness necessitates some kind of storage system, implicit or otherwise (Murphy, 2010).

How exactly do LLMs deal with these cases mechanistically? In this paper, we begin to answer this question by investigating token-level information stored in LLM representations.

- We find that last token positions of multi-token words and named entities “erase” token-level information in early layers for both Llama-2-7b (Touvron et al., 2023) and Llama-3-8b (Meta, 2024).
- We develop a heuristic for scoring the “lexicality” of a given sequence of tokens, and use it to “read out” a list of an LLM’s lexical items given a large dataset of natural text.

We interpret this erasure effect as a “footprint” of a mechanism in early layers that orchestrates the formation of meaningful lexical items.

2 Background

Previous work has shown that knowledge about a multi-token entity is often stored in the last token of that entity. For example, Meng et al. (2022) found that factual information about a subject like “The Space Needle” would be concentrated in the representation for `le`. Geva et al. (2023) find evidence for a *subject enrichment stage* during factual recall, where information about an entity is collected at its last token in early layers, which is also seen in other work on factual recall using the same dataset (Katz et al., 2024), and corroborated by research on athlete \rightarrow sport lookups (Nanda et al., 2023). This phenomenon may be due to the autoregressive nature of decoder transformer models: models cannot enrich “Space” with information about Seattle until after “Needle” is seen, as “Space” could refer

to a number of unrelated concepts (“Space Jam,” “Space Station”).²

Other work in interpretability has also started to uncover evidence of models encoding lexical items. Elhage et al. (2022) observe neurons in early layers that fire on the last tokens of multi-token words, names of famous people, generic nouns, compound words, and LaTeX commands. They also find late-layer neurons that seem to be relevant to *retokenization*, i.e., conversion from internal representations back into tokens. For example, a retokenization neuron might fire on `_st` and promote `rag` in order to facilitate the output of the word “straggler.” Gurnee et al. (2023) also find examples of polysemantic neurons in Pythia models (Mallen and Belrose, 2023) that activate for a number of multi-token constructions like “apple developer,” “Bloom.ington,” and “research.gate.”

3 Linear Probing of Hidden States

3.1 Method

If last token positions are so important (Section 2), then what do these representations encode? Perhaps the last hidden state directly stores information about other subject tokens (e.g., `_Wars` might contain some encoding for `_Star` in its hidden state). To test this hypothesis, we investigate hidden states for both Llama-2-7b and Llama-3-8b, as they have significantly different token vocabulary sizes $|\mathcal{V}|$ (32k and 128k tokens, respectively).

Let d denote the hidden dimension of the model. We train linear probes $p_i^{(\ell)} : \mathbb{R}^d \rightarrow \mathbb{R}^{|\mathcal{V}|}$ to take a hidden state $h_t^{(\ell)} \in \mathbb{R}^d$ at layer ℓ and token position t and predict the value of a nearby token $t + i$. For example, a probe trained to predict the previous token for hidden states at layer 5 would be denoted by $p_{-1}^{(5)}$.

We train probes for all layer indexes $0 \leq \ell < 32$ and offsets $i \in \{-3, -2, -1, 0, 1\}$. We also train probes in the same manner on the embedding layer ($\ell = -1$) and on the final outputs of the network before the language modelling head ($\ell = 32$). We trained probes on a random sample of 428k tokens from the Pile (Gao et al., 2020) using AdamW for 16 epochs with a batch size of 4 and a learning rate of 0.1. Hyperparameters were selected based on validation performance on a separate Pile sample

²This is not a hard-and-fast rule; it depends on entity frequency and context cues. For example, if a model sees `_The`, `_E`, and `iff`, it may already know that these tokens refer to “The Eiffel Tower” without needing to see `e1` and `Tower`.

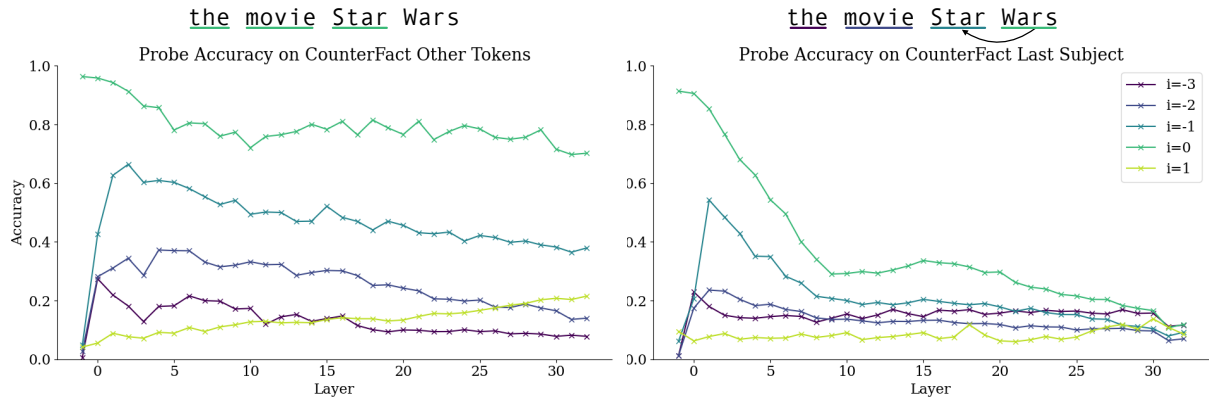


Figure 2: Top-1 test accuracy on COUNTERFACT subject last tokens versus other tokens in the dataset for probes trained on Llama-2-7b hidden states ($n = 5063$). i represents the position being predicted (e.g., $i = -1$ is previous token prediction; $i = 1$ is next-token prediction). We observe an “erasure” effect in last subject tokens that is not present for other types of tokens: these last subject tokens consistently “forget” about preceding tokens and themselves. Appendix A shows Llama-3-8b results and in-distribution performance on Pile tokens.

(279k tokens) after a random sweep. Each probe takes 6-8 hours to train on an RTX-A6000.

3.2 COUNTERFACT Subjects

After training probes in Section 3.1, we test them on the COUNTERFACT dataset (Meng et al., 2022), which consists of prompts about subjects that require factual knowledge to complete correctly (e.g. “Mount Passel is in Antarctica”). We filter the dataset to include only prompts that the model answers correctly, yielding 5,063 examples for Llama-2-7b and 5,495 examples for Llama-3-8b. To augment this dataset, we also sampled and filtered down [album/movie/series \rightarrow creator] pairs from Wikidata (Vrandečić and Krötzsch, 2014) and embedded them in prompts in the same manner, yielding a total of 12,135 correctly-answered prompts for Llama-2-7b and 13,995 for Llama-3-8b.

Figure 2 shows probe test results on COUNTERFACT last subject tokens (right) versus every other type of token in the dataset (left). We see a striking “erasure” effect for last tokens of COUNTERFACT subjects, where these hidden states consistently “forget about” preceding and current tokens. Subject tokens that are not in the last position (e.g., `_Star`) do not exhibit this pattern (Appendix A, Figure 13). This striking drop in token accuracy is reminiscent of the subject enrichment stage described by Geva et al. (2023), suggesting that the tokens `_Star` and `_Wars` may be overwritten in the process of representing the concept of *Star Wars*.

We also observe the same phenomenon when testing on named entities identified by spaCy in Wikipedia articles (Appendix A, Figure 12), sug-

gesting that this effect is not an artifact of the short templates found in the COUNTERFACT dataset. It also does not seem to be a result of any imbalances in probe training data (Appendix B).

3.3 Multi-Token Words

Intuitively, the process of converting a multi-token sequence like `[_n, ort, he, astern]` into a meaningful representation of the word “northeastern” resembles the process of converting `[_E, iff, e1, Tower]` into “Eiffel Tower.” We hypothesize that models treat multi-token words in the same way that they treat multi-token entities, and test our probes from Section 3.1 on multi-token words. After sampling 500 articles ($\sim 256k$ tokens) from the 20220301.en split of the Wikipedia dump (Foundation, 2022), we split by white-space to naively identify word boundaries. As predicted, we see the same “erasing” pattern for multi-token words that we do for multi-token entities (Figure 3). This suggests that they may be processed in a similar manner in early layers. Appendix A shows similar results for Llama-3-8b.

4 Building a Vocabulary

After examination of probe behavior for multi-token words and entities, we hypothesize that this “erasure” effect is a result of the implicit formation of lexical representations in early layers. To characterize this phenomenon, we propose an *erasure score* ψ to identify token sequences that follow the pattern observed in Section 3. We then introduce an approach to “reading out” a list of implicit vocabulary entries for a given model using this score.

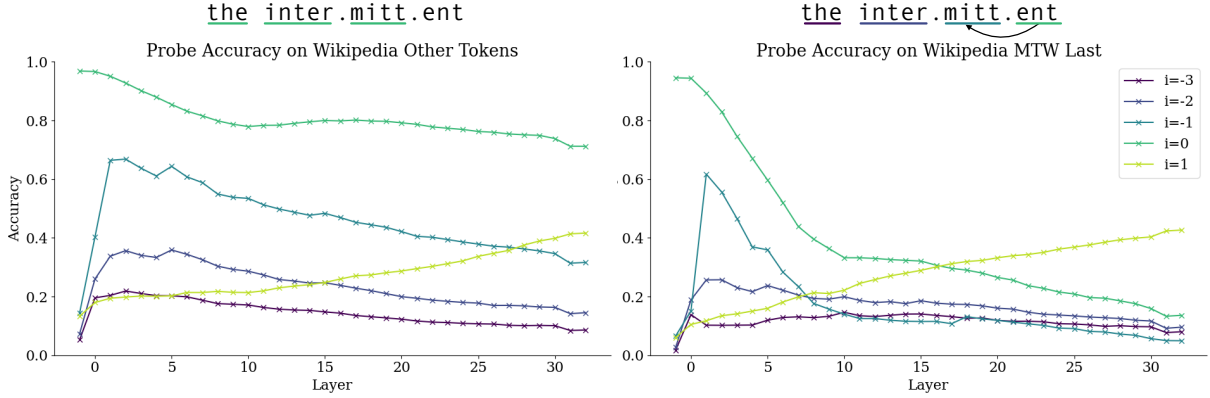


Figure 3: Top-1 test accuracy of probes on last tokens of Wikipedia multi-token words for Llama-2-7b ($n = 80606$). Accuracy on all other tokens shown on the left. We see an erasing effect for multi-token words, similar to the effect seen for COUNTERFACT subjects in Figure 2.

4.1 An Erasure Score

Given some arbitrary sequence of tokens from indices p through q , we want to design an *erasure score* that captures intuitions from Section 3. This score should be higher for sequences exhibiting token erasure (which we hypothesize to be lexical items like `[_Ca1, g, ary]`), and lower for other types of token sequences (e.g., `[_go _to, _Ca1, g]`). We design a metric $\psi_{p,q}$ that uses probe outputs from Section 3 to look for erasure effects between layer 1 and layer L .³

Concretely, Equation 1 defines the score $\psi_{p,q}$ for a sequence $s_{p,q}$ of length $n = q - p + 1$ as:

$$\frac{1}{1 + 2n} \left(\delta(q, 0) + \sum_{t=p}^q \sum_{i=-2}^{-1} \mathbb{1}_{\text{within}}(t, i) \cdot \delta(t, i) \right) \quad (1)$$

where $\delta(t, i)$ denotes the change in probability of the predicted token $t + i$ from layer 1 to layer L , based on probes $p_i^{(\ell)}$ from Section 3.1. We take the softmax of the probe outputs to obtain the probability of a specific token $t + i$ in Equation 2.

$$\delta(t, i) = P_{p_i^{(1)}}(t + i | h_t^{(1)}) - P_{p_i^{(L)}}(t + i | h_t^{(L)}) \quad (2)$$

Finally, if $t + i$ lies outside the boundaries of s , we want the score to decrease. If it is within the boundaries of s , we want a large drop between layers $\delta(t, i)$ to increase the value of $\psi_{p,q}$.

$$\mathbb{1}_{\text{within}}(t, i) = \begin{cases} -1 & \text{if } t + i < p \\ 1 & \text{else} \end{cases} \quad (3)$$

In summary, for every token position $p \leq t \leq q$ and prediction offset $i \in \{-2, -1\}$, we measure

³For both Llama-2-7b and Llama-3-8b we set $L = 9$.

the drop in the predicted probability of the correct token $t + i$ between layer 1 and layer L . The more that the probability of the correct answer *decreases* in early layers, the higher we score that sequence. However, if this “forgetting” occurs for tokens outside of the boundaries of s , we subtract that value from the overall score, effectively penalizing the sequence. This intuition comes from close inspection of probe behavior—for example, Figure 13 shows that there is no “forgetting” effect for $i = -1$ when probing the first token of COUNTERFACT subjects. With this approach, we can also account for cases where s is a subsequence of a larger lexical item: if the token `g` shows a forgetting effect for `_Ca1` in `[_Ca1, g, ary]`, then the sequence `[g, ary]` would be penalized. Finally, $\delta(q, 0)$ additionally rewards sequences in which the last token “forgets itself,” as seen in Figures 2 and 3. We then normalize by the total number of δ values considered, in order to account for differing sequence lengths.

4.2 Segmenting Documents

We develop an algorithm built around our erasure score ψ that breaks any given document $d \in \mathcal{D}$ into high-scoring, non-overlapping segments covering all of d (Algorithm 1). Figure 1 shows the top-scoring sequences $s_{p,q}$ calculated in this manner from a Wikipedia excerpt about Thelonious Monk, where unigram scores are excluded for clarity. Not all multi-token words are scored highly via our approach, but the highest-scoring sequences are plausible lexical items that are non-compositional in nature (“`dram.atic`”, “`sil.ences`”, “`tw.ists`”). We share examples of more documents with complete segmentations in Appendix D.

Algorithm 1 Document Segmentation

Require: document $d \in \mathcal{D}$ of length l

- 1: **for** $n = 1$ **to** l **do** ▷ all ngram lengths
- 2: **for** $p = 0$ **to** $l - n$ **do**
- 3: **for** $q = p + n - 1$ **to** $l - 1$ **do**
- 4: assign score $\psi_{p,q}$ to sequence $s_{p,q}$
- 5: **end for**
- 6: **end for**
- 7: **end for**
- 8: sort s in descending order of ψ
- 9: $segms \leftarrow \emptyset$
- 10: **for** $s_{p,q}$ in sorted s **do**
- 11: **if** $\forall s_{x,y} \in segms, (x > q \vee y < p)$ **then**
- 12: $segms \leftarrow segms \cup \{s_{p,q}\}$
- 13: **end if**
- 14: **end for**
- 15: **return** $segms$ ▷ non-overlapping segments

Token Sequence	n	ct	ψ
lower case	3	2	0.736012
storm	2	4	0.716379
excursion	4	2	0.713134
====... (72 'equals' signs)	8	2	0.712982
Mom	3	2	0.706778
acre	3	2	0.629213
Subject	3	2	0.607172
ninth	3	2	0.606669
processing elements	3	2	0.599549
CVC	3	2	0.596735

Table 1: Top ten highest-scoring sequences for Llama-2-7b using a Pile subsample (1658 sequences recovered total). n is the number of tokens in the sequence, and ‘ct’ represents occurrences of this segment. ψ is averaged over all occurrences.

4.3 Model Vocabularies

Finally, we propose a method to “read out” the implicit vocabulary of a model \mathcal{M} given a dataset \mathcal{D} . For each document $d \in \mathcal{D}$, we segment d using Algorithm 1. We then average scores ψ for every multi-token sequence that appears more than once in \mathcal{D} . As this process is very data-dependent, we show results for both Pile and Wikipedia text. The top 50 sequences for each dataset and model are provided in Appendix E.

With this approach, we are able to recover ~ 1800 sequences for Llama-2-7b and ~ 900 for Llama-3-8b using the same five hundred Wikipedia articles. Although recall is quite low (Table 2),

		MTW		MTE	
llama	data	prec.	recall	prec.	recall
2-7b	wiki	0.306	0.016	0.143	0.016
	pile	0.296	0.017	0.080	0.018
3-8b	wiki	0.044	0.001	0.010	0.000
	pile	0.023	0.001	0.012	0.001

Table 2: Precision and recall for aggregated results of Algorithm 1 run on Llama-2-7b and Llama-3-8b, using either Wikipedia or Pile documents ($|\mathcal{D}| = 500$). MTW refers to all multi-token words in the dataset when split by whitespace; MTE refers to all spaCy named entities.

we find that 44.9% of sequences recovered for Llama-2-7b on Wikipedia text are either multi-token words or multi-token entities (29.68% for Pile text). For Llama-3-8b, only 5% and 3% of retrieved sequences are multi-token words or entities. However, looking at examples of Llama-3-8b sequences in Appendix E, we can observe other interesting cases, like multi-token expressions (“gold medalists,” “by per capita income,” “thank you for your understanding”) and LaTeX commands (as similarly observed by Elhage et al. (2022)). Because Llama-3-8b’s *token* vocabulary is four times larger than Llama-2-7b’s, its *implicit* vocabulary also seems to consist of larger multi-word expressions and chunks of code rather than multi-token words (Appendix E, Table 7).

5 Conclusion

In this work, we present preliminary evidence for the existence of an *implicit vocabulary* that allows models to convert from byte-pair encoded tokens to useful lexical items. We posit that the “erasure” effect we observe for Llama-2-7b and Llama-3-8b is a result of model processes that deal with multi-token expressions, and use this insight to propose a new method for “reading out” an LLM’s implicit vocabulary. This is a first step towards understanding the formation of lexical representations in LLMs, and may serve as a useful tool for elucidation of words that a given model “knows.”

Limitations

Evaluation of implicit vocabulary-building methods (Section 4) is challenging due to the lack of a known ground-truth. Our approach is motivated by the desire to understand the inner workings of the model being studied, but we have no authorita-

tive reference that distinguishes between situations where a given sequence gets a high ψ value because it is truly treated as a lexical unit by the model, or where it may be due to an error in our methodology. To quantify our results, we have compared the extracted vocabulary to sequences that we assume to be likely lexical items: multi-token words and spaCy named entities. However, this likely does not cover all cases for which “token grouping” occurs in LLMs.

Another limitation of this work is that we have restricted our analysis to *known* entities. There is also the question of what happens for intermediate cases such as plausible-sounding fictional towns or names of people who are not famous. If ψ correlates with sequence presence in training data, these results could be useful for understanding how familiar an LLM is with a given word or entity.

Finally, our measurements have been run only on the Llama family of models and do not yet extend to non-Llama models of comparable size, or Llama models of larger sizes.

Ethics Statement

In this work, we restrict our analysis to English words, due to our biases as native speakers of English. We hope that this work can also provide valuable insights for other languages, especially low-resource languages, where understanding “what words an LLM knows” may be especially useful.

Acknowledgments

We thank Koyena Pal, David Smith, Bilal Chughtai, Chantal Shaib, Atticus Geiger, and Adrian Chang for helpful discussion and feedback throughout the course of this project. This work was supported in part by Open Philanthropy, and by the National Science Foundation (NSF) grant IIS-1901117.

Experiments were implemented using the *ninsight* library (Fiotto-Kaufman et al., 2024). Many were run on the Center for AI Safety Compute Cluster. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the sponsors.

References

Khuyagbaatar Batsuren, Ekaterina Vylomova, Verna Dankers, Tsetsuukhei Delgerbaatar, Omri Uzan, Yuval Pinter, and Gábor Bella. 2024. [Evaluating subword tokenization: Alien subword compo-](#)

[sition and oov generalization challenge](#). *Preprint*, arXiv:2404.13292.

Nelson Elhage, Tristan Hume, Catherine Olsson, Neel Nanda, Tom Henighan, Scott Johnston, Sheer ElShowk, Nicholas Joseph, Nova DasSarma, Ben Mann, Danny Hernandez, Amanda Askell, Kamal Ndousse, Andy Jones, Dawn Drain, Anna Chen, Yuntao Bai, Deep Ganguli, Liane Lovitt, Zac Hatfield-Dodds, Jackson Kernion, Tom Conerly, Shauna Kravec, Stanislav Fort, Saurav Kadavath, Josh Jacobson, Eli Tran-Johnson, Jared Kaplan, Jack Clark, Tom Brown, Sam McCandlish, Dario Amodei, and Christopher Olah. 2022. [Softmax linear units](#). *Transformer Circuits Thread*. <https://transformer-circuits.pub/2022/solu/index.html>.

Jaden Fiotto-Kaufman, Alexander R Loftus, Eric Todd, Jannik Brinkmann, Caden Juang, Koyena Pal, Can Rager, Aaron Mueller, Samuel Marks, Arnab Sen Sharma, Francesca Lucchetti, Michael Ripa, Adam Belfki, Nikhil Prakash, Sumeet Multani, Carla Brodley, Arjun Guha, Jonathan Bell, Byron Wallace, and David Bau. 2024. [Ninsight and ndif: Democratizing access to foundation model internals](#). *Preprint*, arXiv:2407.14561.

Wikimedia Foundation. 2022. [Wikimedia downloads](#).

Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. [The Pile: An 800gb dataset of diverse text for language modeling](#). *arXiv preprint arXiv:2101.00027*.

Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. 2023. [Dissecting recall of factual associations in auto-regressive language models](#). *ArXiv*, abs/2304.14767.

Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. 2023. [Finding neurons in a haystack: Case studies with sparse probing](#). *Preprint*, arXiv:2305.01610.

Bernal Jiménez Gutiérrez, Huan Sun, and Yu Su. 2023. [Biomedical language models are robust to sub-optimal tokenization](#). *Preprint*, arXiv:2306.17649.

Shahar Katz, Yonatan Belinkov, Mor Geva, and Lior Wolf. 2024. [Backward lens: Projecting language model gradients into the vocabulary space](#). *Preprint*, arXiv:2402.12865.

Alex Mallen and Nora Belrose. 2023. [Eliciting latent knowledge from quirky language models](#). *Preprint*, arXiv:2312.01037.

Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. [Locating and editing factual associations in gpt](#). In *Neural Information Processing Systems*.

Meta. 2024. [Introducing meta llama 3: The most capable openly available llm to date](#).

- M. Lynne Murphy. 2010. *Lexical Meaning*. Cambridge Textbooks in Linguistics. Cambridge University Press.
- Neel Nanda, Senthooan Rajamanoharan, János Krámar, and Rohin Shah. 2023. [Fact finding: Attempting to reverse-engineer factual recall on the neuron level](#).
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Neural machine translation of rare words with subword units](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- James Simpson. 2011. *The Routledge handbook of applied linguistics*. Taylor & Francis.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esioibu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *Preprint*, arXiv:2307.09288.
- Denny Vrandečić and Markus Krötzsch. 2014. [Wiki-data: a free collaborative knowledgebase](#). *Commun. ACM*, 57(10):78–85.
- Dixuan Wang, Yanda Li, Junyuan Jiang, Zepeng Ding, Guochao Jiang, Jiaqing Liang, and Deqing Yang. 2024. [Tokenization matters! degrading large language models through challenging their tokenization](#). *Preprint*, arXiv:2405.17067.

A Additional Probing Results

A.1 Llama-3-8b Results

COUNTERFACT Accuracy We share results analogous to Figure 2 for Llama-3-8b, which shows a similar “erasure” pattern (Figure 9). Probes are tested only on prompts that Llama-3-8b answers correctly.

Multi-Token Word Accuracy Figure 10 shows results for Llama-3-8b probes tested on the last token positions of multi-token words from Wikipedia (where “words” are determined by whitespace separation).

Multi-Token Entity Accuracy Figure 11 shows results for probes tested on the last token positions of multi-token entities identified by spaCy, using the same dataset that we do for multi-token words. We use spaCy’s named entity recognition pipeline to identify named entities. Because digits 0-9 are added to Llama-2-7b’s vocabulary, we filter out all classes relating to numbers (PERCENT, DATE, CARDINAL, TIME, ORDINAL, MONEY, QUANTITY), with the thought that these sequences may be treated differently at the detokenization stage.

A.2 Llama-2-7b Results

Multi-Token Entity Accuracy Figure 12 shows results for Llama-2-7b probes tested on multi-token entities from Wikipedia, using the same dataset from Section 3.3 and also filtering out number-based entity classes as in Section A.1.

Pile Accuracy While Figure 2 shows test accuracy of linear probes on model hidden states, Figure 4 shows in-distribution test accuracy on Pile tokens. We can observe a smoother trajectory of gradual “forgetting” of previous and current token-level information throughout layers.

Comparison of Token Positions Figure 13 shows the breakdown of probe performance on different types of subject tokens: first subject tokens, middle subject tokens, and last subject tokens. We see that the observed drop in previous and current token representation observed in last subject tokens still exists, but is not as drastic for first and middle subject tokens.

Comparison of Subject Lengths We also show previous token representation broken down by

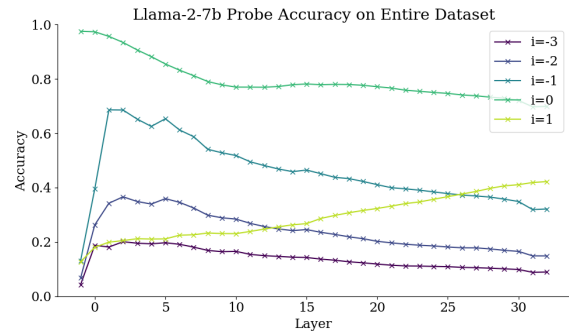


Figure 4: Overall test accuracy on unseen Pile tokens ($n = 273k$) for probes trained on Llama-2-7b hidden states. Next token prediction becomes more accurate throughout model layers as current and previous token accuracy decreases.

COUNTERFACT subject length for last token representations in Figure 14. Unigram subjects represent previous token information at a rate even higher than non-subject tokens. For bigrams and trigrams, we see a pattern similar to Figure 2.

B Accounting for Possible Training Imbalance

One explanation for the observed drop in accuracy for COUNTERFACT entities across layers is that our probes have simply not been exposed to as many entity tokens during training. We do not believe this is the case for Llama-2-7b for two reasons: (1) If this effect was due to probes being less sensitive to tokens found in multi-token entities, we would also see a significant drop for first and middle tokens, which does not occur (Figure 13). (2) We measure the frequency of all test n-grams in the original Pile data used to train our probes, and find that both subject and non-subject n-grams are found in the probe training dataset at similar rates, with the median number of occurrences in the test set for both types of sequences being zero. After removing the few non-subject sequences that do appear often in the probe training set, we still see the same “erasure” effect.

C Choice of L

We choose $L = 9$ based on probe behavior for Llama-2-7b and Llama-3-8b, particularly in Figures 2 and 3. Table 3 shows an additional ablation experiment for $L \in \{5, 9, 13, 17, 21\}$.

L	MTW		MTE	
	prec.	recall	prec.	recall
5	0.307	0.002	0.143	0.002
9	0.306	0.016	0.143	0.016
13	0.328	0.003	0.169	0.003
17	0.330	0.003	0.180	0.003
21	0.319	0.003	0.172	0.003

Table 3: Precision and recall for different values of L for Algorithm 1 applied to Llama-2-7b on Wikipedia text. Recall seems to be best for $L = 9$, with precision improving by a few points in mid-late layers.

<> Danae Suzanna Sweetapple is an Australian Paralympic swimmer. She was born in the Queensland town of St George. Sweetapple attended boarding school at 11 and has a Bachelor of Arts in Literature. She took up swimming in 1990. Her early swimming results led to her being offered one of the first Australian Institute of Sport scholarships for disabled swimmers. At the 1992 Barcelona Games, she won a silver medal in the Women's 100m Freestyle B2 event and she won two bronze medals in the Women's 100m Backstroke B2 and Women's 50m Freestyle B2 events. After the Games she commented "I'd be so happy if more people could make movement and sport a way of life. It's a great way to meet people and gain confidence." Sweetapple was the Young Queenslander of the Year in 1992. References Female Paralympic swimmers of Australia Swimmers at the 1992 Summer Paralympics Par

Figure 5: Full segmentation of a document from Wikipedia via Algorithm 1 on Llama-2-7b. Borders indicate segmentation, with bolded letters indicating multi-token segments. Darker blue cells have higher scores, yellow cells have negative scores. The highest-scoring sequence in this document is “Australian Institute” ($\psi = 0.579$).

D Document Segmentation

We provide full document segmentations using Algorithm 1 for a short excerpt from a Wikipedia article in Figures 5 and 6. Figures 7 and 8 show segmentations for a Pile document.

E Model Vocabularies

Tables 4 through 7 show the top 50 highest-scoring multi-token sequences for Llama-2-7b and Llama-3-8b across either five hundred Wikipedia articles or five hundred Pile samples. Entries were filtered to show only sequences that appear more than once.

<> Danae Suzanna Sweetapple is an Australian Paralympic swimmer. She was born in the Queensland town of St George. Sweetapple attended boarding school at 11 and has a Bachelor of Arts in Literature. She took up swimming in 1990. Her early swimming results led to her being offered one of the first Australian Institute of Sport scholarships for disabled swimmers. At the 1992 Barcelona Games, she won a silver medal in the Women's 100m Freestyle B2 event and she won two bronze medals in the Women's 100m Backstroke B2 and Women's 50m Freestyle B2 events. After the Games she commented "I'd be so happy if more people could make movement and sport a way of life. It's a great way to meet people and gain confidence." Sweetapple was the Young Queenslander of the Year in 1992. References Female Paralympic swimmers of Australia Swimmers at the 1992 Summer Paralympics Par

Figure 6: Full segmentation of a document from Wikipedia via Algorithm 1 on Llama-3-8b. Borders indicate segmentation, with bolded letters indicating multi-token segments. Darker blue cells have higher scores, yellow cells have negative scores. The highest-scoring sequence in this document is “. After the Games she commented ”” ($\psi = 0.443$).

<> Q: Model-View-Controller in JavaScript tldr: How does one implement MVC in JavaScript in a clean way? I'm trying to implement MVC in JavaScript. I have googled and reorganized with my code count less times but have not found a suitable solution. (The code just doesn't "feel right".) Here is how I'm going about it right now. It's incredibly complicated and is a pain to work with (but still better than the pile of code I had before). It has ugly workarounds that sort of defeat the purpose of MVC. And behold, the mess, if you're really brave: // Create a "main model" var main = Model0(); function Model0() { // Create an associated view and store its methods in "view" var view = View0(); // Create a submodel and pass it a function that will "subviewify" the sub

Figure 7: Full segmentation of a document from the Pile via Algorithm 1 on Llama-2-7b. Borders indicate segmentation, with bolded letters indicating multi-token segments. Darker blue cells have higher scores, yellow cells have negative scores. The highest-scoring sequence in this document is “submodel” ($\psi = 0.559$).

<> Q: Model-View-Controller in JavaScript tldr: How does one implement MVC in JavaScript in a clean way? I'm trying to implement MVC in JavaScript. I have googled and reorganized with my code countless times but have not found a suitable solution. (The code just doesn't "feel right".) Here is how I'm going about it right now. It's incredibly complicated and is a pain to work with (but still better than the pile of code I had before). It has ugly workarounds that sort of defeat the purpose of MVC. And behold, the mess, if you're really brave: // Create a "main model" var main = Model0(); function Model0() { // Create an associated

Figure 8: Full segmentation of a document from the Pile via Algorithm 1 on Llama-3-8b. Borders indicate segmentation, with bolded letters indicating multi-token segments. Darker blue cells have higher scores, yellow cells have negative scores. The highest-scoring sequence in this document is “re really brave:” ($\psi = 0.634$).

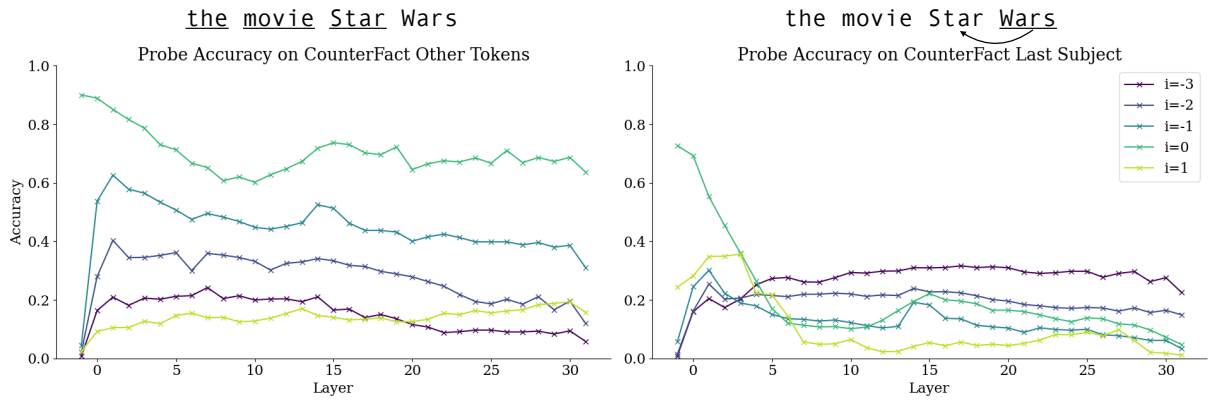


Figure 9: Test accuracy on COUNTERFACT subject last tokens versus other tokens in the dataset for probes trained on **Llama-3-8b** ($n = 5495$). i represents the position being predicted (e.g., $i = -1$ is previous token prediction; $i = 1$ is next-token prediction). We observe an “erasure” effect similar to Figure 2.

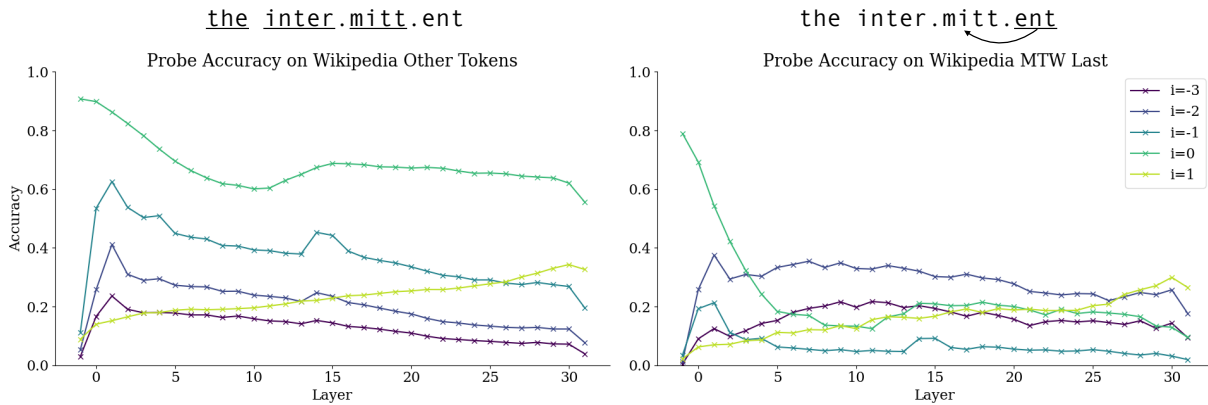


Figure 10: Test accuracy of probes on last tokens of Wikipedia **multi-token words** for probes trained on **Llama-3-8b** ($n = 91935$; right). Test accuracy on all other tokens shown on the left. Similarly to Figure 2, we see an erasing effect that is not present for other types of tokens.

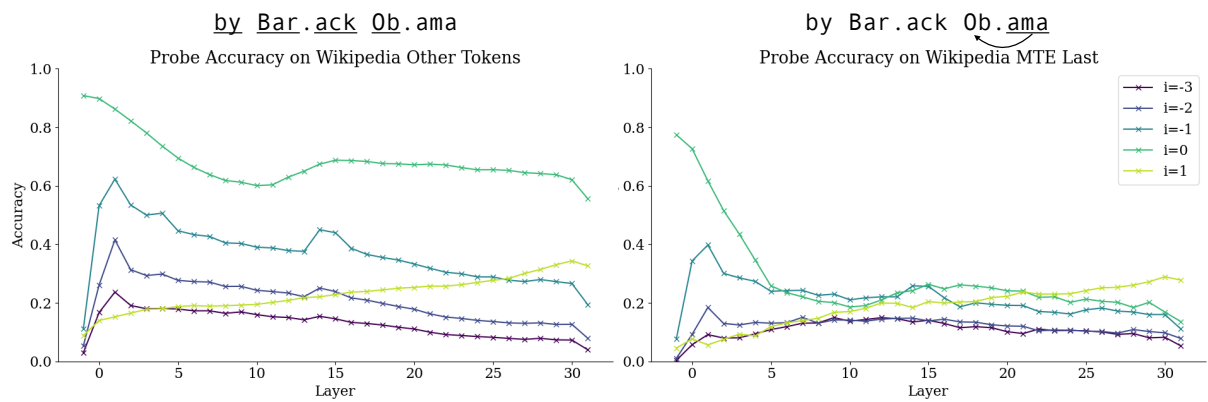


Figure 11: Test accuracy of probes on last tokens of Wikipedia **multi-token entities** for probes trained on **Llama-3-8b** ($n = 36723$; right). Entities are identified via spaCy named entity recognition, excluding entity types that include digits.

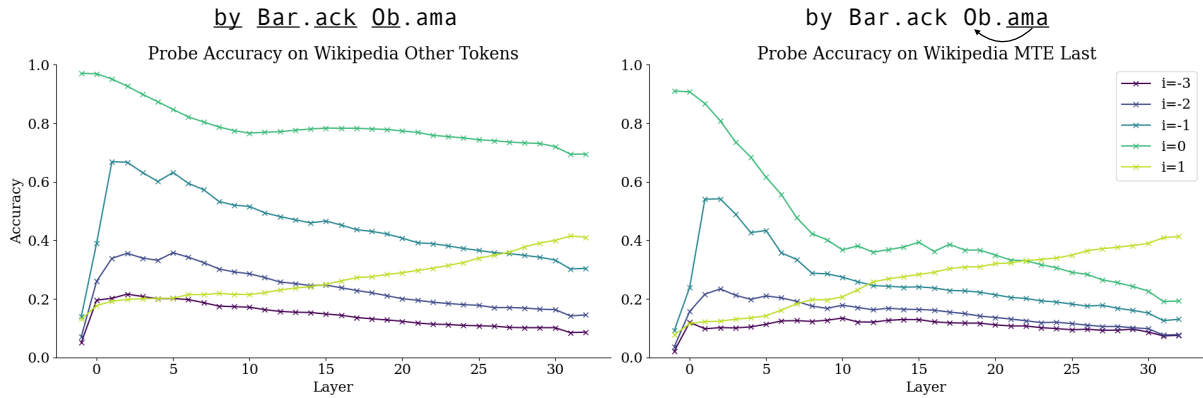


Figure 12: Test accuracy of probes on last tokens of Wikipedia **multi-token entities** for **Llama-2-7b** ($n = 36723$; right). Test accuracy on all other tokens shown on the left. Entities are identified via spaCy named entity recognition, excluding entity types that include digits.

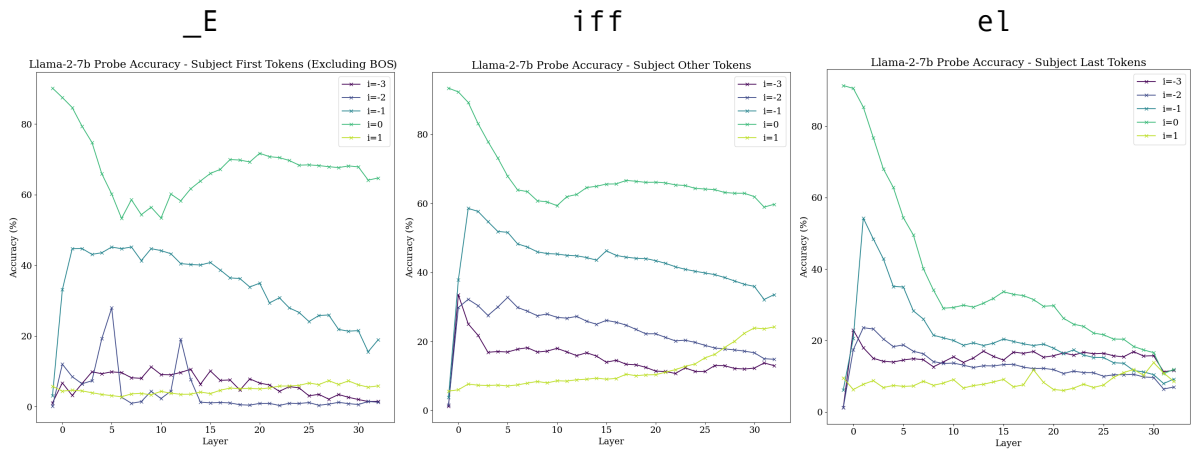


Figure 13: Breakdown for Section 3 probes tested on COUNTERFACT first subject tokens, middle subject tokens, and last subject tokens. We observe an “erasing” effect only for last subject tokens. Because BOS tokens are recoverable by $i = -1$ probes at high rates, and since 55% of prompts tested on had subjects at the beginning, we filter examples for which BOS tokens are labels from the leftmost plot.

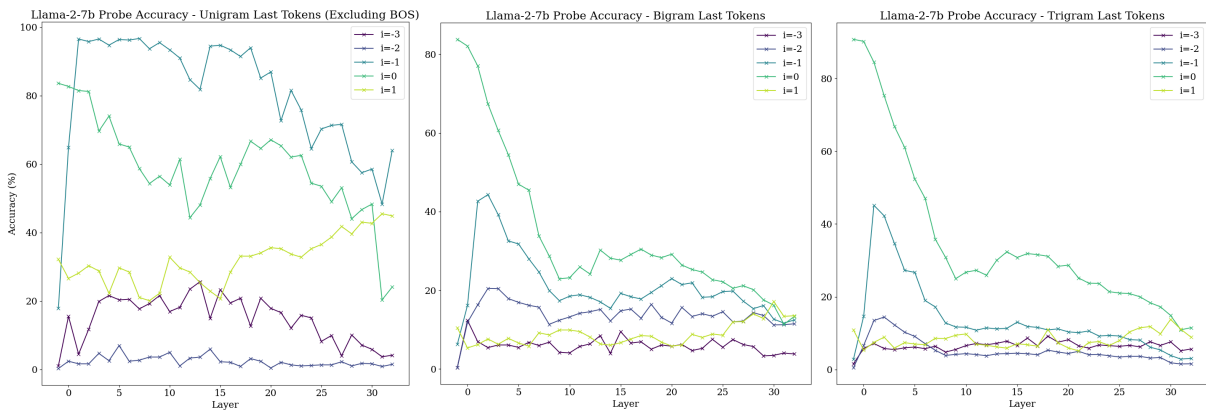


Figure 14: Probe test results for COUNTERFACT subject last tokens broken down for unigrams, bigrams, and trigrams. Unigram subjects store previous token information at rates near 100%, even excluding BOS tokens.

Token Sequence	n	ct	ψ
Gottsche	3	2	0.685220
berth	3	2	0.680793
carries	3	2	0.647844
Eurocop	3	2	0.644104
franchises	3	2	0.642707
0 Women	3	2	0.639162
rape	3	2	0.632567
Rebell	3	3	0.614295
intermittently	4	2	0.613479
enn State	4	3	0.607535
North Dakota	4	10	0.600616
Sride	3	2	0.600013
fiction	2	2	0.599339
Sox	3	3	0.599043
Bazz	3	2	0.598242
erect	3	2	0.597915
borough	3	3	0.596054
encompasses	5	2	0.592084
northernmost	3	2	0.591607
Madras	3	2	0.590394
hull	3	2	0.586968
iron	2	2	0.586959
Galaxy	3	2	0.585879
began operations	3	2	0.584680
Redding	3	2	0.584244
gloss	3	2	0.576740
cello	3	2	0.573732
Gators	3	5	0.573675
senator	3	2	0.572947
restructuring	4	2	0.570552
supervised	3	3	0.570421
Mediterranean	4	2	0.567790
Madera	3	2	0.567563
sequel	3	2	0.563626
scarp	3	3	0.561548
Sout	3	2	0.560640
South Division	3	2	0.558720
rectangular	3	2	0.557339
Danny	3	2	0.556836
Examiner	4	2	0.555797
Kuwait	4	4	0.554636
Bogue	3	6	0.552219
Lancaster	3	3	0.552166
Leuven	4	3	0.548806
the Park	3	2	0.548687
first Baron	3	2	0.547447
fight	3	2	0.547171
Carpio	3	2	0.547116
Czech Republic	3	2	0.546651
Survive	4	2	0.546255

Table 4: **Llama-2-7b** Wikipedia results (1808 sequences total). n is the number of tokens in the sequence, and ‘ct’ represents occurrences of this segment. ψ is averaged over all occurrences.

Token Sequence	n	ct	ψ
1992 births	7	2	0.573
19th-century	7	3	0.569
dehydrogen	5	2	0.553
Swahili	4	4	0.539052
Chuck Liddell	6	2	0.537169
its population was	5	5	0.534977
by per capita income	6	3	0.518991
are brownish	4	2	0.515703
ate women’s football	7	4	0.509384
Almeida	4	5	0.507277
of New South Wales	5	3	0.503120
2015 deaths	8	2	0.503074
Pittsburgh	3	3	0.503070
21st-century	7	4	0.499362
(NSW	4	9	0.497107
age of the United Kingdom	6	3	0.487303
Presidential	3	2	0.485317
Landmark	3	2	0.484965
Alistair	4	2	0.484930
Tauri	3	8	0.482449
2 km	4	2	0.479984
20th-century	7	3	0.475703
East Bay	3	2	0.475156
game goes in extra time, if the scored	10	2	0.472323
São Paulo	3	2	0.470874
Atlantic City	3	2	0.470726
Chaluk	3	2	0.467165
Frank Lloyd	3	2	0.462585
may refer to:	6	4	0.462234
gold medalists	4	2	0.458494
, 2nd Baron	6	2	0.456996
people)	4	4	0.454926
series aired	4	2	0.453057
Srib	3	2	0.451708
with blackish	4	2	0.450033
World Cup players	4	2	0.448979
main role	3	2	0.448569
Bos	4	2	0.448425
Asenath	4	2	0.448259
Royal Navy	3	3	0.445617
2. Bundesliga players	7	2	0.445210
External links	3	69	0.444921
an unincorpor	6	2	0.443527
Gast	2	4	0.437695
Pfor	3	2	0.432194
Elisio de Med	5	2	0.431518
" (2007) "Jad	12	2	0.429412
Elkh	3	2	0.428984
Früh	3	2	0.427781
order of the NK	5	2	0.424037

Table 5: **Llama-3-8b** Wikipedia results (892 sequences total). n is the number of tokens in the sequence, and ‘ct’ represents occurrences of this segment. ψ is averaged over all occurrences.

Token Sequence	n	ct	ψ
lower case	3	2	0.736012
storm	2	4	0.716379
excursion	4	2	0.713134
====... (72 'equals' signs)	8	2	0.712982
Mom	3	2	0.706778
acre	3	2	0.629213
Subject	3	2	0.607172
ninth	3	2	0.606669
processing elements	3	2	0.599549
CVC	3	2	0.596735
VPN	3	3	0.596052
Regul	3	2	0.591968
bore	2	2	0.590212
\$\dot{G}	5	2	0.589714
Rates	3	2	0.589637
INSURANCE	5	2	0.584323
Commercial	4	2	0.581543
Barney	3	3	0.574872
PTA	3	2	0.571932
penetrated	4	2	0.570164
MG	3	2	0.569830
Leigh	3	2	0.567894
jail	3	3	0.567225
TNS	3	2	0.567003
peptides	4	2	0.565775
John Arena	3	2	0.565648
Disease	4	2	0.564662
welfare	4	4	0.564364
wild type	3	2	0.560699
uws	3	3	0.557799
ongrel	4	3	0.554208
liquid cry	3	3	0.553408
princess	3	2	0.551672
Denmark	3	2	0.548702
birthday	3	2	0.548504
atedmes	4	2	0.548171
"ENOENT	5	2	0.547169
third-party	4	2	0.546949
aliens	3	2	0.546507
Durban	3	4	0.545848
Bouncy	4	3	0.545826
CHO	3	2	0.542762
unjust	3	2	0.538813
these motivational	4	3	0.537485
DLS	3	4	0.535933
\n&	3	2	0.534510
uneven	3	2	0.533137
watt	3	2	0.532243
'She	3	2	0.531300
HP	3	3	0.529555

Table 6: **Llama-2-7b** Pile results (1658 sequences total). n is the number of tokens in the sequence, and 'ct' represents occurrences of this segment. ψ is averaged over all occurrences.

Token Sequence	n	ct	ψ
</td>\n<td>	9	2	0.627583
{d}x	5	3	0.599395
*\n	4	3	0.587016
_{n=1}{\in	7	4	0.585434
</td>\n<td	8	2	0.573310
-2-2007-061	12	3	0.551581
reticulum	4	3	0.549337
INSURANCE	5	2	0.548263
32;\n internal static	8	2	0.547893
;\n internal static	6	9	0.540374
: At	4	2	0.538609
(2,9,'	6	4	0.537495
Respondent	4	2	0.534509
\t)\n\n\t	7	3	0.530669
(3,0,'	6	4	0.529493
_{n-1}\var	7	2	0.527303
thank you for	6	2	0.513979
your understanding	6	2	0.513979
hydroxyl	4	2	0.510059
>\n*\private \$	9	2	0.510054
in mukaan	5	2	0.506333
{w}{B}_{	6	2	0.505970
/2Z	5	2	0.501998
'); \nINSERT INTO	6	10	0.501055
7-f131	7	2	0.496881
0, 1L>	8	2	0.495809
/0 S	5	2	0.492042
5 Audi	4	2	0.491043
all that apply	4	3	0.490469
": true,\n	6	2	0.486807
4,\n	5	2	0.485315
to as DSP	5	2	0.484967
B]{\	6	2	0.483484
;\ninternal	5	3	0.479777
100% used	6	2	0.475673
", "x":	5	3	0.474701
2.7	4	2	0.473720
</td>\n	6	2	0.473578
" code="	4	4	0.473514
e2d-d	6	2	0.473418
is under conversion	4	5	0.473355
{ intlsys	5	3	0.471213
()\n}\n\nprivate	12	2	0.470941
boolean isAny	12	2	0.470941
(2,8,'	6	4	0.470214
trachea	4	2	0.469154
use in an automobile	6	2	0.467788
at org.apache.c	7	5	0.467637
world around us	4	2	0.464469
2\left(1+x	8	2	0.463555
or Commodore	5	3	0.463106
11-117	7	2	0.459824

Table 7: **Llama-3-8b** Pile results (819 sequences total). n is the number of tokens in the sequence, and 'ct' represents occurrences of this segment. ψ is averaged over all occurrences.