

Great Memory, Shallow Reasoning: Limits of k NN-LMs

Anonymous ACL submission

Abstract

K -nearest neighbor language models (k NN-LMs), which integrate retrieval with next-word prediction, have demonstrated strong performance in language modeling as well as some downstream NLP benchmarks. These results have led researchers to argue that models trained on poor quality or outdated data could perform well by employing a k NN extension that has access to a higher-quality datastore. In this work, we ask whether this improved ability to recall information really translates into downstream abilities. We extensively evaluate k NN-LMs on a diverse set of tasks, ranging from sentiment classification and commonsense reasoning to multi-hop reasoning. Results show that k NN-LMs excel at *memory*-intensive tasks, where utilizing the patterns in the input is sufficient for determining the output, but struggle with *reasoning* tasks that require integrating multiple pieces of information to derive new knowledge. We further demonstrate through oracle experiments and qualitative analysis that even with perfect retrieval, k NN-LMs still fail to determine the correct answers, placing an upper bound on their reasoning performance.

1 Introduction

A foundational property of pretrained language modeling (Peters et al., 2018; Devlin et al., 2019) has been that improvements to the perplexity of the model lead to improvements on downstream tasks. This property is central to the scaling of large language models (LLMs) where researchers focus nearly exclusively on perplexity as a proxy metric for improved general purpose abilities (Kaplan et al., 2020). In recent years, this research has centered primarily on high-quality text data at greater quantities as the limiting component for producing better language models (Hoffmann et al., 2022).

This increasing need for training data has led to significant challenges. On one hand, including as much high-quality data as possible results

in improved downstream performance. On the other hand, this data is often protected by licenses or copyright, which means training on such data brings legal issues.

It would be ideal to circumvent this issue entirely with alternative approaches. If a model could be trained on lower-quality data but adapted to perform well on real tasks, it might provide a technical workaround. Non-parametric Language Models (NPLMs), such as k NN-LMs, have emerged as a promising approach in this space (Khandelwal et al., 2020). k NN-LMs extend neural LMs by linearly interpolating with simple k -nearest neighbor LMs. This approach can improve language modeling with its memory over a massive collection of texts, usually referred to as a datastore. Khandelwal et al. (2021) and Shi et al. (2022) validate that k NN-LMs achieve better performance on downstream tasks compared to standard LMs. The SILO model of Min et al. (2024) applies this approach further by training a LM exclusively on license-permissive data and using a non-parametric datastore to improve the models during inference.

In this work, we study the limits of how k NN-LMs can be used to improve LLMs. Specifically, we are interested in whether the improvements in perplexity seen with k NN-LMs are equivalent to other improvements in LM ability. This question relates to debates about whether memory is separable from other language abilities and how they interact in NLP benchmarks.

We summarize our contributions as follows. First, we evaluate k NN-LMs on 20 NLP tasks, with experimental results revealing that lower perplexity does not necessarily lead to better reasoning in non-parametric settings. To investigate the performance degradation, we conduct extensive analyses in Appendix F, which shows that k NN-LMs are not sensitive to semantic information and can be distracted by irrelevant tokens. Figure 1 illustrates such limitations using a multi-hop reasoning ex-

Question: When Copsi was made earl of Northumbria he went to reside in a town at the confluence of which two rivers? The two rivers are ____

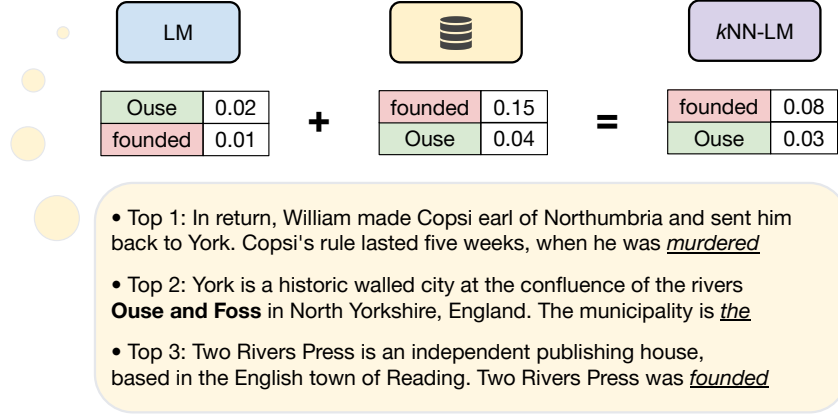


Figure 1: In this multi-hop question answering (QA) example, the LM is very uncertain about the next word and could benefit from retrieval. The k NN approach finds several document, both irrelevant and relevant, that may help. However, two issues occur: first, an irrelevant document increases the probability of a random wrong answer; second, even though a relevant document has been found, it may not upweight the actual answer (Ouse). We study how these issues may impact task performance as compared to perplexity.

ample. We open-source two datastores along with our distributed kNN search implementations for multiple GPUs to support further research.

2 Experimental Setup

We use Llama-2-7b (Touvron et al., 2023), Llama-3-8B (AI@Meta, 2024), and Mistral-7B (Jiang et al., 2023) as our inference models. For each inference model, we build the corresponding datastores. The keys are the 4096-dimensional hidden representations before the final MLP which predicts the token distribution at each generation step, produced by executing forward passes over the datastore corpora. For efficient similarity search, we create a FAISS index (Johnson et al., 2019) and search for nearest-neighbor tokens using Euclidean distance. Due to the scale of the datastores, we perform approximate search instead of exact search. We base our implementation on Alon et al. (2022).

Hyperparameters include λ , k , and σ . λ determines the weight of the datastore, and we consider $\lambda \in \{0.1, 0.2, 0.3\}$. We retrieve $k \in \{1600, 2048\}$ neighbors and smooth the k NN distribution with a temperature $\sigma \in \{1, 3, 5, 10\}$. Table 8 shows hyperparameters we use for different tasks.

For each inference model, we use Math and Wiki datastores for language modeling on the corresponding evaluation datasets: wikitext and math textbooks. Each datastore represents a specific domain, and we evaluate the performance of k NN-LMs on a domain by measuring the perplexity of

each evaluation dataset. We conduct a grid search to find the hyperparameters that yield the lowest PPL for each datastore. The optimal hyperparameters for each datastore are later applied across all downstream tasks in our experiments.

We provide eight demonstrations for GSM8K and three demonstrations for BBH. For the other datasets, we perform zero-shot inference. Details of the experiments are in Appendix C.

3 k NN-LMs Help In-Domain Perplexity

To explore how different sources of external knowledge impact downstream task performance, we experiment with two datastores. First, we follow the choice made by Shi et al. (2022), where they identify heterogeneous data sources broadly relevant to common downstream NLP tasks. In particular, they mix Wikitext103 (Merity et al., 2017), with other sources including the English portion of Amazon Review (He and McAuley, 2016), CC-NEWS (Hamborg et al., 2017) and IMDB (Maas et al., 2011). We call this datastore *Wiki*.

Then, we hypothesize that the commonly explored corpora for building datastores do not contain relevant knowledge to assist with math reasoning tasks. To maximize the performance gain on these tasks, we construct a datastore comprising 3.94K mathematical textbooks, sourced from (Wang et al., 2023b). We will refer to this datastore as *Math*. We summarize the statistics of each datastore in Table 6 in Appendix C.

	RTE	RT	CB	Yahoo	CR	AGN	HYP	MR	SST2
Llama2-7B	66.06	80.20	50.00	59.37	74.55	81.30	64.15	82.40	84.02
+Wiki	66.43	80.77	51.79	58.83	76.95	81.46	64.15	83.00	84.68
+Math	65.70	79.83	51.79	59.10	73.70	81.79	50.39	82.30	84.62
Llama3-8B	70.76	77.49	64.29	58.87	79.10	79.17	59.30	84.75	86.54
+Wiki	61.37	78.71	71.43	58.93	80.45	79.33	59.30	84.85	87.04
+Math	70.76	77.39	66.07	56.83	79.40	80.11	59.30	83.70	87.10
Mistral-7B	76.17	80.96	71.43	56.63	81.90	73.57	56.59	78.90	81.82
+Wiki	76.17	81.71	67.86	56.63	82.15	73.55	56.78	78.95	81.77
+Math	76.17	80.68	75.00	56.63	81.85	73.59	56.78	78.90	81.77

Table 1: Accuracy comparison on various memory-intensive tasks.

Model	LM Performance	
	Wiki	Math
Llama2-7b	10.63	7.90
+Wiki	9.74	8.75
+Math	11.33	7.23
Llama-3-8b	9.70	5.36
+Wiki	9.32	6.03
+Math	10.37	5.22
Mistral-7B	9.72	5.64
+Wiki	9.29	6.41
+Math	10.49	5.59

Table 2: Perplexity comparison. Rows vary the datastore \mathcal{D} used. Columns represent different held-out test sets. Lower numbers indicate better performance.

We begin by validating past results of k NN-LMs on language modeling. We present results in Table 2. To facilitate meaningful comparisons between models with different tokenizers and vocabulary sizes, we report word-level perplexities. These results show that having access to a non-parametric datastore leads to lower perplexity compared to using a standalone LM across all datasets. This improvement in perplexity is observed when the corpus used to construct the datastore and the one used for inference share the same data source. For instance, since the training split of Wikitext103 is in Wiki, the LM+Wiki setting achieves the lowest perplexity on Wikitext103’s validation set. Utilizing the other datastore results in performance worse than that of the standalone LM.

4 k NN-LMs Can Help Memory-Intensive Tasks

We begin by looking at a set of memory-intensive tasks, which we believe can be solved by pattern matching at scale without complex reasoning. We incorporate three types of tasks: sentiment classification, which aims to predict whether the sentiment of a text is positive or negative; textual entailment, which assesses the relationship between two

sentences, determining if it constitutes entailment, contradiction, or neutrality; and topic classification, which involves identifying the main topic of a text. We describe dataset details in Appendix D.

For classification and multiple-choice question-answering (QA) tasks, we utilize Domain Conditional Pointwise Mutual Information (DCPMI) (Holtzman et al., 2021) to predict answers. We then calculate accuracy metrics to compare performance across different models. We measure the performance using F1 scores at the token level for text generation. Additionally, whenever feasible, we employ fuzzy verbalizers (Shi et al., 2022) to maximize the performance of k NN-LMs.

Table 1 summarizes the results of these tasks. On these tasks, k NN-LMs exhibit improved performance. Incorporating an external datastore outperforms a standalone LM on most datasets while showing comparable performance on the remaining dataset. We further explain this performance gap through qualitative analysis in Appendix F.3.

5 k NN-LMs Hurt Reasoning Performance

For reasoning tasks, we consider three types: knowledge-intensive reasoning, which focuses on utilizing world knowledge for making (potential) multi-hop inferences; commonsense reasoning, which involves leveraging commonsense knowledge to understand social and physical interactions; and mathematical reasoning, which includes arithmetic, logical, and discrete reasoning abilities. We describe dataset details in Appendix D.

We present the results for knowledge-intensive tasks in Table 3. In contrast to the earlier findings, using a standalone LM consistently outperforms k NN-LMs on these tasks. Most surprisingly, on Natural Questions and HotpotQA, which consist of QA pairs constructed from Wikipedia documents, performance does not improve even though Wiki contains several million Wikipedia tokens. Retrieval

	NQ	HotpotQA	Arc-Challenge	Arc-Easy	OBQA	MLLU
Llama2-7B	23.18	22.72	41.81	57.49	57.00	39.22
+Wiki	22.53	22.53	38.31	57.41	56.20	38.68
+Math	21.14	21.26	41.04	56.82	56.20	38.53
Llama3-8B	23.64	25.14	44.88	58.83	55.80	42.67
+Wiki	24.00	24.48	43.94	58.59	53.80	42.32
+Math	23.04	24.63	43.26	58.59	54.60	42.46
Mistral-7B	20.63	20.96	46.42	60.94	58.80	41.91
+Wiki	20.58	20.80	46.16	60.61	57.40	41.80
+Math	20.56	20.48	46.08	60.77	57.80	41.55

Table 3: Performance comparison on datasets for knowledge-intensive reasoning tasks.

	Winogrande	HellaSwag	DROP	GSM8K	BBH
Llama2-7B	69.37	64.46	32.39	14.83	30.69
+Wiki	70.32	63.67	32.14	12.05	32.08
+Math	68.98	63.54	32.31	13.48	30.82
Llama3-8B	73.95	65.99	45.55	45.72	39.67
+Wiki	73.95	64.71	45.02	44.28	39.01
+Math	74.19	65.15	45.54	45.63	39.92
Mistral	74.19	69.08	46.93	36.30	43.37
+Wiki	74.66	68.21	46.69	36.45	42.69
+Math	73.64	68.11	46.38	36.60	43.09

Table 4: Performance comparison on datasets for other reasoning tasks.

		Perplexity	Accuracy
OBQA	LM	255.76	55.80
	k NN-LM	9.41	95.60
NQ	LM	112.56	23.64
	k NN-LM	8.91	46.40
HotpotQA	LM	158.26	25.14
	k NN-LM	8.15	49.85

Table 5: Results in an oracle setting where the k NN-LMs always include the correct answer as one of the k nearest neighbors.

ing from Wiki leads to a three-point decrease in performance. Results for commonsense reasoning and mathematical reasoning tasks are shown in Table 4. The standalone LM once again outperforms k NN-LMs models on three of the five datasets. The most significant differences in performance occur on GSM8K. Although incorporating an external data store results in a slight performance increase on Mistral, this does not demonstrate the effectiveness of k NN-LMs on GSM8K. Under Mistral’s parameter settings, k NN-LMs has minimal changes on the predictions of the standalone LM, merely introducing some randomness. Finally, although k NN-LMs do not improve GSM8K and Drop over standard LMs, we find that retrieving from Math improves over retrieving from Wiki.

Do k NN-LMs fail due to retrieval errors? We investigate whether degraded reasoning capabilities of k NN-LMs stem from a failure in retrieval. We

examine k NN-LMs’ behaviors when retrieval is perfect. To achieve perfect retrieval, we include the correct answer among the k nearest neighbors. We construct a datastore for OpenbookQA, NQ, and HotpotQA, respectively, including their train and test examples. We then examine both perplexity and accuracy. The results, presented in Table 5, indicate that while k NN-LMs can significantly reduce the perplexity, the model does not always derive the correct answer, even when the correct answer is explicitly given as one of the k neighbors. Therefore, the failure of reasoning cannot be fully attributed to the failure of retrieval. However, perfect retrieval does improve LM by a large margin, suggesting that better retrieval is beneficial. Currently, retrieval is performed by finding similar hidden representations. A training-based approach such as RAG (Lewis et al., 2020) has the potential to improve retrieval substantially.

6 Conclusions

We investigate whether the improved perplexity observed in k NN-LMs models can be translated into enhanced reasoning capabilities. Our findings indicate that while k NN-LMs improve perplexity and can achieve better performance on memory-intensive tasks, they struggle with reasoning-intensive tasks, showing a disconnect between LM ability and task ability.

Limitations

As we are limited by computing budget, we only build datastores up to 610 million tokens. It is unlikely although not impossible that larger datastores built on general web corpus like C4 will lead to better reasoning capabilities. Additionally, we only experiment with LLMs with seven- to eight-billion model parameters as the base models. The findings in this paper may not generalize to other, possibly larger, base models.

References

AI@Meta. 2024. [Llama 3 model card](#).

Uri Alon, Frank Xu, Junxian He, Sudipta Sengupta, Dan Roth, and Graham Neubig. 2022. Neuro-symbolic language modeling with automaton-augmented retrieval. In *International Conference on Machine Learning*, pages 468–485. PMLR.

Akari Asai, Zexuan Zhong, Danqi Chen, Pang Wei Koh, Luke Zettlemoyer, Hannaneh Hajishirzi, and Wen-tau Yih. 2024. Reliable, adaptable, and attributable language models with retrieval. *arXiv preprint arXiv:2403.03187*.

Parishad BehnamGhader, Santiago Miret, and Siva Reddy. 2023. [Can retriever-augmented language models reason? the blame game between the retriever and the language model](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 15492–15509. Association for Computational Linguistics.

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. 2022. [Improving language models by retrieving from trillions of tokens](#). In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162, pages 2206–2240. PMLR.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

Ido Dagan, Bill Dolan, Bernardo Magnini, and Dan Roth. 2010. Recognizing textual entailment: Rational, evaluation and approaches—erratum. *Natural Language Engineering*, 16(1):105–105.

Marie-Catherine De Marneffe, Mandy Simons, and Judith Tonhauser. 2019. The commitmentbank: Investigating projection in naturally occurring discourse. In *proceedings of Sinn und Bedeutung*, volume 23, pages 107–124.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. [DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2368–2378. Association for Computational Linguistics.

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR.

Felix Hamborg, Norman Meuschke, Corinna Breiteringer, and Bela Gipp. 2017. [news-please: A generic news crawler and extractor](#). In *Proceedings of the 15th International Symposium of Information Science*, pages 218–223.

Ruining He and Julian McAuley. 2016. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *proceedings of the 25th international conference on world wide web*, pages 507–517.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.

Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. 2022. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*.

Ari Holtzman, Peter West, Vered Shwartz, Yejin Choi, and Luke Zettlemoyer. 2021. [Surface form competition: Why the highest probability answer isn’t always right](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7038–7051. Association for Computational Linguistics.

365	Minqing Hu and Bing Liu. 2004. Mining and summarizing customer reviews. In <i>Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining</i> , pages 168–177.	422
366		423
367		
368		
369	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .	424
370		425
371		426
372		427
373		428
374		429
375	Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. <i>IEEE Transactions on Big Data</i> , 7(3):535–547.	430
376		431
377	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. <i>arXiv preprint arXiv:2001.08361</i> .	432
378		433
379		434
380		435
381		
382	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6769–6781, Online. Association for Computational Linguistics.	436
383		437
384		438
385		439
386		
387		
388		
389	Urvashi Khandelwal, Angela Fan, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2021. Nearest neighbor machine translation. In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	440
390		441
391		442
392		443
393		444
394	Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. Generalization through memorization: Nearest neighbor language models. In <i>International Conference on Learning Representations</i> .	445
395		
396		
397		
398		
399	Johannes Kiesel, Maria Mestre, Rishabh Shukla, Emmanuel Vincent, Payam Adineh, David Corney, Benno Stein, and Martin Potthast. 2019. Semeval-2019 task 4: Hyperpartisan news detection. In <i>Proceedings of the 13th International Workshop on Semantic Evaluation</i> , pages 829–839.	446
400		447
401		448
402		449
403		450
404		451
405	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. <i>Transactions of the Association for Computational Linguistics</i> , 7:453–466.	452
406		453
407		454
408		455
409		456
410		457
411		458
412	Zachary Levonian, Chenglu Li, Wangda Zhu, Anoushka Gade, Owen Henkel, Millie-Ellen Postle, and Wanli Xing. 2023. Retrieval-augmented generation to improve math question-answering: Trade-offs between groundedness and human preference. <i>arXiv preprint arXiv:2310.03184</i> .	459
413		460
414		461
415		462
416		463
417		464
418	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in Neural Information Processing Systems</i> , 33:9459–9474.	465
419		466
420		467
421		468
		469
		470
	Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In <i>Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies</i> , pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.	471
		472
		473
		474
		475
		476
		477
		478
	Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2017. Pointer sentinel mixture models. In <i>International Conference on Learning Representations</i> .	
	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. In <i>EMNLP</i> .	
	Sewon Min, Suchin Gururangan, Eric Wallace, Weijia Shi, Hannaneh Hajishirzi, Noah A. Smith, and Luke Zettlemoyer. 2024. SILO language models: Isolating legal risk in a nonparametric datastore. In <i>The Twelfth International Conference on Learning Representations</i> .	
	Bo Pang and Lillian Lee. 2005. Seeing stars: exploiting class relationships for sentiment categorization with respect to rating scales. In <i>Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics</i> , ACL '05, page 115–124. Association for Computational Linguistics.	
	Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.	
	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. <i>Communications of the ACM</i> , 64(9):99–106.	
	Weijia Shi, Julian Michael, Suchin Gururangan, and Luke Zettlemoyer. 2022. Nearest neighbor zero-shot inference. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 3254–3265. Association for Computational Linguistics.	
	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages 1631–1642. Association for Computational Linguistics.	

Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. 2022. Challenging big-bench tasks and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Shufan Wang, Yixiao Song, Andrew Drozdov, Aparna Garimella, Varun Manjunatha, and Mohit Iyyer. 2023a. *kNN-LM does not improve open-ended text generation*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15023–15037. Association for Computational Linguistics.

Zengzhi Wang, Rui Xia, and Pengfei Liu. 2023b. Generative ai for math: Part i—mathpile: A billion-token-scale pretraining corpus for math. *arXiv preprint arXiv:2312.17120*.

Kaiyu Yang, Aidan Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu, Saad Godil, Ryan Prenger, and Anima Anandkumar. 2023. LeanDojo: Theorem proving with retrieval-augmented language models. In *Neural Information Processing Systems (NeurIPS)*.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. *HotpotQA: A dataset for diverse, explainable multi-hop question answering*. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380. Association for Computational Linguistics.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. *HellaSwag: Can a machine really finish your sentence?* In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800. Association for Computational Linguistics.

Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28.

A Related Work

Retrieval Models Although LLMs achieve superhuman performance on a wide range of natural language processing tasks, they often produce hallucinations, struggle with incorporating recent knowledge, and expose private information present in the training data. Recently, research interest has shifted towards retrieval-based LMs, which

combine a parametric neural model and a non-parametric external datastore (Guu et al., 2020; Karpukhin et al., 2020). These retrieval-based LMs naturally incorporate new knowledge, enhance the factuality of generated texts, and reduce privacy concerns (Asai et al., 2024). Furthermore, Borgeaud et al. (2022) demonstrate that employing retrieval augmentation during large-scale pre-training can outperform standard LMs while requiring fewer parameters.

Among retrieval-based LMs, k NN-LMs (Khandelwal et al., 2020) emerge as a popular choice (Min et al., 2024). Unlike other retrieval models that encode and retrieve documents, k NN-LMs encode and retrieve tokens. At every token, k NN-LMs search for the k most similar tokens from the datastore based on contextualized token embeddings, which are then turned into a next-token distribution. k NN-LMs linearly interpolate the retrieved k NN distribution with the output of a base LM. They do not require additional training but introduce computational and memory overhead.

Reasoning Retrieval. Little research has been conducted on constructing retrieval models for reasoning tasks. Leandro (Yang et al., 2023) investigates the use of retrieval-based LMs to assist with theorem proving, and Levonian et al. (2023) experiment with retrieving content from mathematical textbooks to generate responses to student questions. In our study, we create a reasoning-specific datastore to assist LMs in performing reasoning-intensive tasks.

Evaluation of k NN-LMs. While k NN-LMs excel at language modeling and have demonstrated enhanced performance in machine translation (Khandelwal et al., 2021) and simple NLP tasks (Shi et al., 2022), the question of whether they are thoughtful reasoners remains open. Wang et al. (2023a) demonstrate that k NN-LMs struggle with open-ended text generation as they only provide benefits for a narrow set of token predictions and produce less reliable predictions when generating longer text. BehnamGhader et al. (2023) showed that when retrieval is conducted based on the similarity between queries and statements, k NN-LMs often fail to identify statements critical for reasoning. Even when these crucial statements are retrieved, it is challenging for k NN-LMs to effectively leverage them to infer new knowledge. These studies, however, are limited to a narrow set of tasks. Our work seeks to provide a compre-

\mathcal{D}	Text Size	Tokens	Mem
Wiki	2.2GB	610M	44G
Math	0.6GB	200M	15G

Table 6: Overview of the two datastores. Tokens are produced by Llama2 tokenizers. Mem is the memory size of the datastore.

hensive evaluation of the reasoning capabilities of k NN-LMs and provides an extensive analysis of the sources of their failures.

B Background: k NN-LMs

Non-parametric language models are variants of standard language models that give the model the ability to utilize an additional datastore \mathcal{D} during inference to determine the next word prediction, $p(x_{t+1}|x_{1:t}; \mathcal{D})$. This datastore may be part of the original training data, data for adaptation to a new domain, or be used to incorporate continual updates or protected data. As these datastores are typically quite large, this process requires a retrieval component in the loop to find the sparse subset of the datastore that can best inform the current prediction. Several popular approaches exist including DPR (Karpukhin et al., 2020) and REALM (Gua et al., 2020).

In this work, we focus on k NN-LMs due to their popularity as an approach to directly improve LM perplexity on fixed models without a need for re-training. As noted in the intro, this approach has also been put forward as a method for circumventing the need for high-quality licensed training data in LLMs. Formally k NN-LMs are defined as

$$p(x_{1:T}; \mathcal{D}) = \prod_t p(x_{t+1} | x_{1:t}; \mathcal{D})$$

$$= \prod_t (\lambda p_{kNN}(x_{t+1} | x_{1:t}; \mathcal{D}) + (1 - \lambda)p(x_{t+1} | x_{1:t}))$$

Let (k_i, v_i) be the i th (key, value) pair in \mathcal{D} , $f(\cdot)$ maps a token sequence to its contextual representation, and $d(\cdot)$ measures the distance between two vectors.

$$p_{kNN}(x_{t+1} | x_{1:t}; \mathcal{D})$$

$$\propto \sum_{(k_i, v_i) \in \mathcal{D}} \mathbf{1}_{x_{t+1}=v_i} \times \exp(-d(k_i, f(x_{1:t}))).$$

When using a Transformer language model, we define the distance metric $d(\cdot)$ as the squared ℓ_2

Corpus	Text Size	Tokens
Wikitext103	0.5GB	140M
Amazon	0.07GB	18M
CC-NEWS	1.6GB	443M
IMDB	0.03GB	8M
Total	2.2GB	609M

Table 7: Statistics of each data source in the Wiki datastore.

distance. To assemble the datastore, we run the language model over all the documents to collect the hidden states and corresponding next word.

C More Implementation Details

Table 6 presents the statistics of each datastore. Table 7 presents the data sources of the Wiki datastore. Table 8 shows hyperparameters we use for different tasks.

D Dataset Details

The datasets selected for memory-intensive tasks are as follows:

- For sentiment classification, we include SST-2 (Socher et al., 2013), movie review (MR) (Pang and Lee, 2005), customer review (CR) (Hu and Liu, 2004), Rotten Tomatoes (RT), and hyperpartisan news detection (HYP) (Kiesel et al., 2019).
- For textual entailment, we use CommitmentBank (CB) (De Marneffe et al., 2019) and Recognizing Textual Entailment (RTE) (Dagan et al., 2010).
- For topic classification, our datasets are AG News (AGN) (Zhang et al., 2015) and Yahoo! Answers (Yahoo) (Zhang et al., 2015).

The datasets selected for reasoning-intensive tasks are as follows:

- For knowledge-intensive reasoning, we explore Natural Questions (NQ) (Kwiatkowski et al., 2019), HotpotQA (Yang et al., 2018), ARC Easy and Challenge (Clark et al., 2018), OpenbookQA (OBQA) (Mihaylov et al., 2018), and MMLU (Hendrycks et al., 2020) to assess the model’s ability to apply extensive world knowledge.
- For commonsense reasoning, we examine Hel-laSwag (Zellers et al., 2019) and Winogrande (Sakaguchi et al., 2021), which test the model’s understanding of social norms and physical laws.

Data	λ	k	τ
Llama2 + Wiki	0.2	2048	5.0
Llama3 + Wiki	0.1	2048	5.0
Mistral + Wiki	0.1	2048	10.0
Llama2 + Math	0.2	1600	3.0
Llama3 + Math	0.1	2048	3.0
Mistral + Math	0.1	2048	10.0

Table 8: Hyperparameters in k NN-LM. **Top:** Hyperparameters for Wiki datastore. **Bottom:** Hyperparameters for Math datastore .

- For mathematical reasoning, we utilize DROP (Dua et al., 2019), GSM8K (Cobbe et al., 2021), and BBH (Suzgun et al., 2022) to evaluate the model’s capacity for complex arithmetic, logical deductions, and handling of discrete concepts.

E More Results

Language modeling and data contamination.

We study whether lower perplexity in language modeling is a result of data contamination in the datastore. To eliminate this confounder, we perform decontamination before measuring perplexities. Specifically, we decontaminate by filtering out evaluation documents that have eight-gram overlaps with any document in the datastore. Table 9 summarizes the results. After data decontamination, k NN-LMs still achieve lower perplexity, despite the gaps between standard LMs and k NN-LMs being smaller.

Significance tests for memory-intensive tasks

Our main experiments used hyperparameters that produce the lowest in-domain perplexity, with lambda values set to 0.1 or 0.2. With these values, k NN-LMs only incur minor changes to the prediction, making the differences between LM and k NN-LM relatively small. We conducted the Wilcoxon Signed-Rank Test on both reasoning-intensive and memory-intensive tasks to check if the minor changes are indeed significant. For reasoning tasks, results on both Wiki and Math datastores rejected the null hypothesis, indicating that our results are statistically significant. For memory-intensive tasks, the results of LM + Wiki have a P-value of 0.036, which rejects the null hypothesis at a significance level of 0.05. However, the P-value for LM + Math is 0.661, suggesting that the results of LM + Math on memory-intensive

Model	LM Performance	
	Wiki	Math
Llama2-7b	10.63	7.90
+ k NN-LM	9.74	7.23
Llama2-7b-Decon.	13.63	12.06
+ k NN-LM-Decon.	13.45	11.10
Llama-3-8b	9.70	5.36
+ k NN-LM	9.32	5.22
Llama-3-8b-Decon.	13.50	7.77
+ k NN-LM-Decon.	13.16	7.61
Mistral-7B	9.72	5.64
+ k NN-LM	9.29	5.59
Mistral-7B+Decon.	12.58	8.29
+ k NN-LM-Decon.	12.72	8.32

Table 9: Perplexity comparison. k NN-LM used datastore belongs to the same domain as the evaluation dataset. Decon. refers to evaluating the standard LM on decontaminated datasets.

		P-value
memory	LM vs LM + Wiki	0.036
	LM vs LM + Math	0.661
reasoning	LM vs LM + Wiki	7e-4
	LM vs LM + Math	1e-6

Table 10: Significance test for memory-intensive and reasoning-intensive tasks

tasks are not significant. The detailed values are presented in Table 10.

F Analysis

The results of this work show that k NN-LMs generally hurt the reasoning of models, despite helping perplexity and other simpler tasks. Here, we investigate the cause of this further.

F.1 Qualitative Analysis.

We conduct qualitative analysis to understand the failures of k NN-LMs better. In the qualitative analysis, we inspect examples of knowledge-intensive and mathematical reasoning datasets and show the retrieved tokens as well as the proceeding context. Through these examples, we find the following patterns that prevent k NN-LM from retrieving the correct token.

- **k NN-LMs struggle with multi-hop reasoning questions.** When the task requires extracting

HotpotQA Example	Label	LM Pred
Which American character actor who starred on the television series “Stargate SG-1” (1997–2007) and appeared in “Episode 8” of “Twin Peaks” as a guest star?	Don S. Davis	Don S. Davis
Retrieved Context	Token	k NN-LM Pred
<ul style="list-style-type: none"> • After the first three seasons of Stargate SG-1 had been filmed on 16 mm film (although scenes involving visual effects had always been shot on 35 mm film for various technical reasons), “Nemesis” was the first episode filmed entirely on 35 mm film ... “Nemesis” was the last episode before actor • “200” won the 2007 Constellation Award for Best Overall 2006 Science Fiction Film or Television Script, and was nominated for the 2007 Hugo Award for Best Dramatic Presentation, Short Form. The episode also marks the first time original SG-1 member • Season one regular cast members included Richard Dean Anderson, Amanda Tapping, 	Christopher	Michael Shanks
	Jack	
	Michael	

Table 11: A multihop reasoning example from HotpotQA with predictions of the standard LM and k NN-LMs.

NQ Example	Label	LM Pred
who is the largest supermarket chain in the uk?	Tesco	Tesco
Retrieved context	Token	k NN-LM Pred
<ul style="list-style-type: none"> • The majority of stores will open as normal across the UK, however Sainsbury’s advise shoppers to check details of when your local branch as some may close earlier than normal using the online store locator tool.(Image: Bloomberg) Supermarket giant • Along with Lidl, Aldi has eaten away at the market share of the Big Four supermarkets: • buy one, get one free (BOGOF) offers have been criticised for encouraging customers to purchase food items that are eventually thrown away; as part of its own campaign on food waste, supermarket retailer 	Asda	Asda
	Tesco	
	Morris	

Table 12: A knowledge-intensive reasoning example from Natural Questions with predictions of the standard LM and k NN-LMs.

multiple pieces of sentences from the corpus and then combining the information to infer the answer, k NN-LMs often retrieve tokens that are contextually appropriate and relevant to part of the question, rather than the correct answer. As shown in Table 11, for the multi-hop reasoning question from HotpotQA, the model needs to identify an actor who both starred in Stargate SG-1 and guest-starred in Twin Peaks. While the required information is available in Wikipedia, it is distributed across two paragraphs. k NN-LMs retrieve only the actors from Stargate SG-1, failing to combine information from two sources to perform accurate multi-hop reasoning.

supermarket in the UK, due to the highly similar contexts of ‘supermarket giant’ and ‘the largest supermarket, k NN-LMs ultimately assign a high probability to Asda and make a wrong prediction.

- **k NN-LMs are sensitive to the syntax but not the semantics of the question.** While k NN-LM retrieves the next token that fits the context, it cannot distinguish subtle semantic differences between different words in a sentence. As a result, when more than one word fits the context, it may not select the correct answer. Table 12 demonstrates this issue with an example from the NQ dataset. Even though Asda is not the largest

- **k NN-LMs tend to retrieve high-frequency entities in the corpus.** The entities are often proper nouns like person names and locations. If part of the answer overlaps with these high-frequency proper nouns, k NN-LMs will retrieve them and make wrong predictions, as shown in Table 13 and Table 14.

- **k NN-LMs fail at mathematical reasoning tasks.** For instance, in the object counting task from the BBH dataset, even though k NN-LM understands the context that it needs to retrieve a number as the next token, it cannot solve the complex task of first identifying which objects are musical instruments and then counting them, as shown in Table 15.

HotpotQA Example	Label	LM Pred
What type of plane is the four engine heavy bomber, first introduced in 1938 for the United States Army, which is hangared at Conroe North Houston Regional Airport?	American Boeing B-17 Flying Fortress	The B-17 Flying Fortress
Retrieved context	Token	k NN-LM Pred
<ul style="list-style-type: none"> • A famous symbol of the courage and sacrifices made by American bomber crews during World War II was revealed May 16 at the National Museum of the U.S. Air Force, Wright-Patterson Air Force Base, Ohio. The meticulously restored B- • As the Avenger made its way to the tower area, the wings began to fold up, a maneuver which enabled more of its kind to be loaded side by side into aircraft carriers. The queen of the event was the B- • Spring is here, so why not hop a plane and grab some lunch? Even better if a World War II-era B- 	17	The B-25 Mitchell.
	25	
	25	

Table 13: Example from HotpotQA showing the impact of high-frequency proper nouns in the corpus on k NN-LMs predictions retrieving from Wikipedia.

HotpotQA Example	Label	LM Pred
who is older, Annie Morton or Terry Richardson?	Terry Richardson	Terry Richardson
Retrieved context	Token	k NN-LM Pred
<ul style="list-style-type: none"> • And she still wasn't done. Later she tweeted a warning to all women. "My hard won advice: never get into an elevator alone with [Terry Gilliam.] Terry • #MeToo https://t.co/jPnFhfB5GQ - Ellen Barkin(@EllenBarkin) March 17, 2018Barkin got another shot in. Terry • I haven't posted about Christina Hendricks in a while but it's Valentine's Day and that makes me think of chocolate and chocolate reminds me of Christina Hendricks. And Christina 	Gilliam	Terry Gilliam
	Gilliam	
	Hend	

Table 14: Another example from HotpotQA explains the impact of high-frequency proper nouns in the corpus on k NN-LMs predictions retrieving from Wikipedia.

F.2 Is the problem a failure of model weighting?

We investigate whether degraded reasoning capabilities of k NN-LMs stem from a failure in choosing a good weighting λ . This experiment aims to analyze k NN-LMs' behaviors when λ is optimal for the downstream task. Specifically, we directly search for λ that maximizes the log probabilities of a small set of labeled downstream task examples. We first conduct this experiment on OpenbookQA, NQ, and HotpotQA. We enumerate through retrieving $k \in \{16, 32, 64, 128, 256, 512, 1024, 2048\}$ neighbors and setting temperature $\sigma \in \{1, 2, 5, 10\}$. We retrieve from Wiki. We initialize λ at 0.5, and as the optimization proceeds, we find that smaller λ values correlate with lower loss. Ultimately, we arrive at the minimum loss when λ is close to 0. This process suggests that without any interpolation of the k NN distribution, the correct labels of the provided demonstrations receive the highest log probability.

For comparison, we also conduct similar experiments on memory-intensive tasks. In the main experiments, we use fuzzy labels for classification tasks, where each label corresponds to multiple words during prediction. We summed the probabilities of these words to determine the probability of the fuzzy label. As a result, there is more than one correct answer when performing lambda testing on memory-intensive tasks. Therefore, we cannot directly use the question and answer as model input to compute the answer's loss for gradient updates, as we did in reasoning tasks. Instead, we combined each word within the fuzzy label with the prompt separately to compute the loss, and, for each iteration, used the lowest word loss for gradient updates. The results are shown in Table 17.

Therefore, reasoning tasks such as OpenbookQA, NQ, and HotpotQA are unlikely to benefit from simple k NN access to Wiki. However, memory-intensive tasks like RT, CR, and SST2 have the potential for improvement with such ac-

Mathematical Reasoning Example	Label	LM Pred
I have three violins, three trombones, a flute, and four trumpets. How many musical instruments do I have?	11	11
Retrieved Context	Token	k NN-LM Pred
• In this example, the optimal route would be: 1 -> 3 -> 2 -> 4 -> 1, with a total completion time of	10	
• How many different passwords are there for his website system? How does this compare to the total number of strings of length	10	10
• Using the TSP, the most efficient order in which to schedule these tasks would be: 2 -> 3 -> 1 -> 4 -> 2, with a total completion time of	14	

Table 15: A mathematical reasoning example from BBH requiring object counting with predictions of the standard LM and k NN-LMs.

Sentiment Example	Label	LM Pred
humorous, artsy, and even cute, in an off-kilter, dark, vaguely disturbing way. The sentence has a tone that is	Positive	Negative
Retrieved Context	Retrieved	k NN-LM Pred
<i>Wiki</i>		
• meta-commentator, Imhoff gives us a decidedly modern delivery. His speaking rhythms are staccato and his tone	bitter	
• Collins, who has worked on more than 100 children books and won several awards: his tone is	fun	Negative
• is her own narrator, so the thoughts and feelings of others are conveyed secondhand or are absent entirely. Her tone and language are at turns	honest	
<i>Math</i>		
• preferred term is not “Platonist” but “quasiempiricist”, a word Tymoczko lends a subtly	different	
• ... or a horror film (group 2, $N_H = 29$). The data are coded so that higher scores indicate a more	positive	Positive
• the failure of the Intermediate Value Theorem is neither here nor there nor anywhere else to them. This is not a bad nor a	good	

Table 16: A sentiment analysis example with predictions of the standard LM and k NN-LMs. We show tokens retrieved from each datastore and their proceeding tokens.

cess.

Datasets	lambda
OBQA	0
NQ	0
HotpotQA	0
RT	0.19
CR	0.22
SST2	0.09

Table 17: The lambda values corresponding to the lowest loss across different datasets

F.3 Effect of Math on Sentiment Analysis

We explain why retrieving from Math improves LMs on sentiment analysis. First, we consider a sentiment analysis example in Table 16. In this task, given a sentence, a model is required to predict whether the sentiment expressed is positive or

negative. The sentence in the example expresses a positive sentiment; however, Llama-2 predicts the sentiment to be negative. k NN-LMs, when retrieving from Wiki, fail to find sentiment-related tokens, and hence also predict a negative sentiment. Performing retrieval from Math produced the correct sentiment. However, this is more coincidental rather than reflective of the model’s capability, because, although the retrieved tokens display a positive sentiment, the retrieved contexts are not relevant to the test example. We observe that sentiment-related content is ubiquitous, regardless of the source we use to build the datastore. Even in math textbooks, we find many sentences that express sentiment.