

# Recite, Reconstruct, Recollect: Memorization in LMs as a Multifaceted Phenomenon

Anonymous ACL submission

## Abstract

Memorization in language models is typically treated as a homogenous phenomenon, neglecting the specifics of the memorized data. We instead model memorization as the effect of a set of complex factors that describe each sample and relate it to the model and corpus. To build intuition around these factors, we break memorization down into a taxonomy: recitation of highly duplicated sequences, reconstruction of inherently predictable sequences, and recollection of sequences that are neither. We demonstrate the usefulness of our taxonomy by using it to construct a predictive model for memorization. By analyzing dependencies and inspecting the weights of the predictive model, we find that different factors influence the likelihood of memorization differently depending on the taxonomic category.

## 1 Introduction

The existing literature on Language Model (LM) **memorization**<sup>1</sup>—the tendency to generate exact copies of training samples at test time—varies widely in stated motivation. Papers might focus on copyright (Shi et al., 2023; Karamolegkou et al., 2023; Meeus et al., 2024), privacy (Carlini et al., 2018, 2022b; Brown et al., 2022; Mireshghallah et al., 2022), or scientifically understanding how interpolation (Mallinar et al., 2022) leads to generalization (Feldman, 2021; Tirumala et al., 2022; Henighan et al., 2023a). Although these objectives share commonalities, they also drive distinct and sometimes contradictory notions of memorization. To disentangle these motivations and to articulate the factors that determine or signal memorization, we propose a taxonomy inspired by colloquial distinctions of memorization behavior in humans.

Our taxonomy, illustrated in Fig. 1, defines three types of LM memorization based on colloquial descriptions of human memorization. Humans **recite**

direct quotes that they commit to memory through repeated exposure, so LMs recite highly duplicated sequences. Humans **reconstruct** a passage by remembering a general pattern and filling in the gaps, so LMs reconstruct inherently predictable boilerplate templates. Humans sporadically **recollect** an episodic memory or fragment after a single exposure, so LMs recollect other sequences seen rarely during training.

We use our taxonomy in a variety of experiments that highlight the multifaceted nature of memorization. In summary:

- We introduce an intuitive taxonomy and heuristics for categorizing memorized data.
- By comparing memorized and unmemorized distributions, we assess how a variety of corpus-wide statistics, datum-level metrics, and representational differences influence the likelihood of a given sequence being memorized. Our dependency tests confirm existing findings that low perplexity is strongly associated with memorization—though not equally for all memorized examples. This fact guides our heuristic for partitioning memorized data into a recitation category.
- We study scaling factors in memorization by monitoring each taxonomic category over the course of training and across model sizes. The number of memorized sequences increases with training time and model size, regardless of taxonomic category. Recollection, however, sees the fastest increase—and this outsize growth cannot be attributed solely to repeated exposures to rare sequences or to random memorization.
- We train logistic regressions to predict memorization for each category. This predictive model outperforms both a simple baseline

<sup>1</sup>As defined by [www.genlaw.org/glossary.html](http://www.genlaw.org/glossary.html).

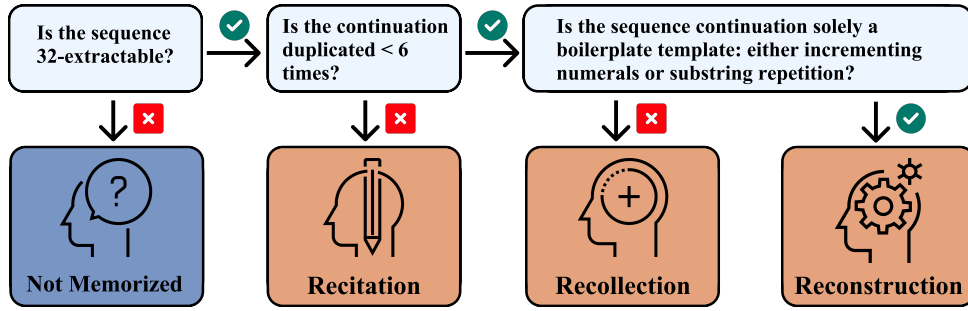


Figure 1: Our intuitive memorization taxonomy has three categories determined by simple heuristics.

with no taxonomy and a model that uses a taxonomy optimized by searching for the best set of mediating factors. These experiments show that the intuitions behind our taxonomy can improve on more generic approaches.

- We highlight differences between categories by exploring statistical dependencies, finding recitation is enabled by low-perplexity prompts and recollection is constrained by the presence of rare tokens.

## 2 Experiments

In this section, we detail the definitions and data we use to analyze varying factors in memorization.

**Defining Memorization** There are multiple competing definitions for memorization (Zhang et al., 2021; Ippolito et al., 2022). Because our experiments employ memorization data released by Biderman et al. (2023a), we use their preferred definition of  $k$ -extractable memorization (Carlini et al., 2022a) with  $k = 32$ . A sample is  $k$ -extractable if the LM, when prompted with the first  $k$  tokens, generates the following  $k$  tokens verbatim.

**Language Models** We study memorization across model scale and training timing using the deduplicated Pythia models (Biderman et al., 2023b), which range in size from 70M to 12B parameters<sup>2</sup> trained on a deduped version of The Pile (Gao et al., 2020). Data order is fixed across runs, enabling causal claims about the effect of model scale on memorization.

**Datasets** Our **memorized sample** is a public list of sequences memorized by Pythia, released by Biderman et al. (2023a). Unlike other works that

<sup>2</sup>Excluding the 160M parameter model, as its memorization dataset exhibits outlier behavior that could be either a buggy data artifact or a real phenomenon, but is regardless outside of the scope of our work.

estimate whether a generation is from a model’s training set using predictive techniques (Carlini et al., 2020; Shi et al., 2023; Yang et al., 2024), this dataset contains all 32-extractable samples from the Pile, verified by referencing the training data (Gao et al., 2020). We also collect a **representative sample** by taking a random 3% subset of The Pile, retaining the first 64 tokens of each sequence. Some analysis also considers an **unmemorized distribution** estimated by subtracting the memorized data distribution from the entire Pile, as inferred from the representative sample.

## 3 Potential factors in memorization

We consider a number of possible factors in whether a given sequence is memorized. These factors are based on corpus statistics, datum statistics intrinsic to that sample, or model perplexity. Features may be computed over the first 32 tokens (the **prompt**); the last 32 tokens (the **continuation**); and the **full sequence** of 64 tokens subsampled from the training data. Implementation details are provided in Appendix A.

Many of these properties have different distributions for memorized and unmemorized data. Fig. 2 illustrates these differences, highlighting that for some properties, the memorized distribution is more concentrated. Other properties—in particular perplexity and number of duplicates in the training corpus—have memorized and unmemorized distributions with visibly different medians. Where the distributions differ, the property in question is likely to influence memorization, an assumption which we employ predictively in Section 6.

### 3.1 Corpus statistics

Some factors relate a given sequence to the entire training corpus. Overall, the following features illustrate how memorization is influenced by various types of duplication.

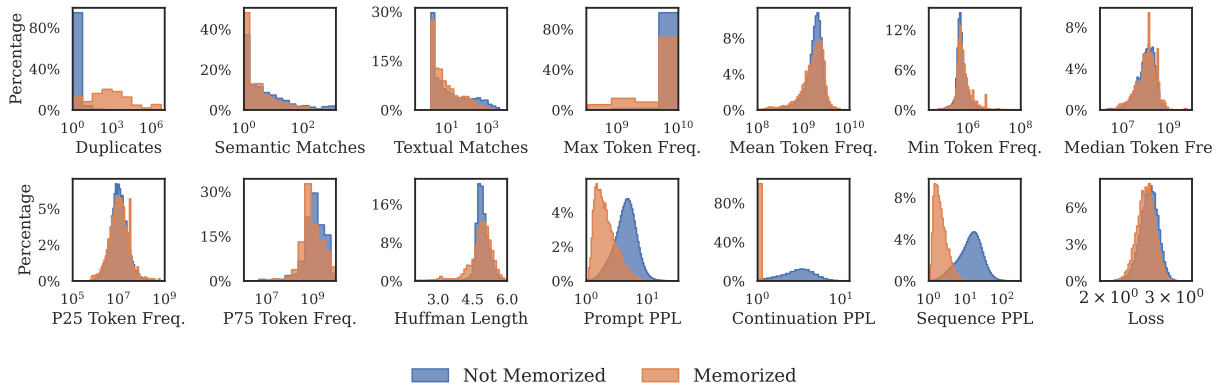


Figure 2: Histogram of various properties of interest (described in Section 3) for memorized and unmemorized (estimated by assuming the representative dataset’s statistics hold for the Pile) samples.

**Duplicates** For each 32-token window in any 2049-token sequence seen during training, we count the number of duplicates in the Pile.

**Semantic Matches** To assess the prevalence of semantically similar samples in training, we generate document embeddings for each full sequence using SBERT and count the number of sequences with cosine similarity  $\leq 0.8$ . These sequences are semantically similar but may not be exact token-level duplicates.

**Textual Matches** We filter the set of semantic matches for a given target sequence to identify those with a low Levenshtein edit distance in their prompts (Levenshtein et al., 1966) from the target sequence. These matches flag slight variations on boilerplate prompts. We compute edit distance at the character level, thereby accounting for different tokenizations of identical sequences.

**Token frequency** We also compute summary statistics about the corpus-wide frequency of individual tokens in the sequence: mean, median, maximum, minimum, and 25th / 75th percentile counts.

### 3.2 Sequence properties

Because some sequences are inherently easier to encode, we also consider factors determined by intrinsic metrics on the sample itself.

**Templating** A sample is classified as templating if it follows a predictable pattern. We do not comprehensively consider all possible templates, but focus on two common patterns defined by hand-crafted heuristics:

- **Repeating:** Consisting only of a short repeating sequence of tokens, e.g., “Go Go Go ...”.

Zhang et al. (2021) previously discussed repetitive templates as a common feature of apparently memorized data which was not classified as counterfactually memorized.

- **Incrementing:** Consisting of incrementing numerical sequences. For example, consider the sequence “23: 0xf1, 24: 0xf2, 25: 0xf3”, a set of interspersed numerical sequences with repeating separators.

**Compressibility** We use Huffman Coding (Huffman, 1952) length to measure how easily a sequence is compressed. Compressibility generalizes repeating templates to cases where minor variations on repeating patterns must be memorized. The connection between learning, memorization, and compression is drawn from the existing literature: Carlini et al. (2020) attempts to filter out sequences that are “easy” to produce by comparing zlib compression with perplexity to identify memorized training data.

### 3.3 Perplexity

We compute average perplexity across tokens on the prompt, continuation, and full sequence. The importance of perplexity is one of the most reproduced results in memorization research (Zhang et al., 2021; Carlini et al., 2018) and we confirm that low perplexity sequences are far more likely to be memorized than high perplexity sequences (Fig. 2). Perplexity is the only factor we consider that relates to model behavior, rather than being intrinsic to the data.

## 4 Memorization Taxonomy

To analyze the fundamental causes of  $k$ -extracted memorization, we subdivide memorized samples

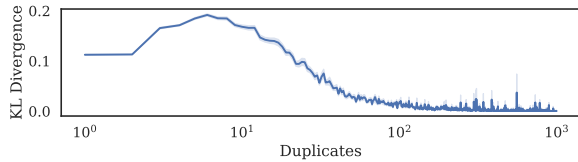


Figure 3: KL divergence between generation perplexity of memorized and non-memorized examples for Pythia 12B with bootstrapped confidence intervals. Non-memorized samples are treated as the reference distribution. Divergence is highest for sequences with 6 duplicates, while highly duplicated sequences have near-identical memorized and unmemoized distributions.

into three types. The following rules categorize a sample as a **candidate** for recitation, reconstruction, or recollection; candidates memorized by the model are therefore respectively recited, reconstructed, or recollected.

#### 4.1 Recitation

The existing memorization literature agrees that the duplication of a sequence across the training corpus is strongly correlated with its memorization (Lee et al., 2021; Kandpal et al., 2022). For example, LMs produce verbatim copies of bible quotes or software licenses that are commonly duplicated.

We consider a sample to be a recitation candidate if it is highly duplicated in the training corpus. Model perplexity is a good predictor of memorization on rare sequences because the perplexity distributions are more different on memorized and unmemoized data with few duplicates (Fig. 3). For highly duplicated sequences, however, perplexity is no longer a good predictor of whether the sequence is memorized or not. We therefore define a recitation candidate as a sequence with at least 6 duplicates because the three-way relationship between perplexity, memorization, and duplicate count differs before and after that maximum divergence point.

LMs memorize a wide variety of highly duplicated texts, as shown in the example of Appendix F. Recited natural language text largely comprises webpage boilerplate text, liturgy, and software licenses or other legalese. Table 3, which includes random samples of natural language recitation, includes all of these common cases. Recited code text, as seen in Appendix 4, is largely web development (HTML, CSS, JavaScript, etc.) boilerplate that describes common elements or derives from popular webpage templates.

#### 4.2 Reconstruction

Are all perfectly reproduced sequences truly “memorized”? We consider cases that may be spuriously classified by definitions like  $k$ -extraction. Rather than encoding the entire sequence, the model learns templates and then reconstructs the sample based on these more broadly applicable patterns. A sequence can thus be perfectly reproduced even if it never appeared during training.

We consider a few templates—stereotyped sequence patterns with a single logical continuation—to define reconstruction candidates. These templates are not intended to be comprehensive, as any stereotyped pattern may permit reconstruction. Our reconstruction candidates are sequences classified as incrementing or repeating by the heuristics described in Section 3.2. As seen in Appendix E, code is more likely to be reconstructed than natural language text. When natural language text is reconstructed, as seen in Appendix F, it often takes the form of a chapter index and it is more likely than code to contain cases of phrase repetition rather than arithmetic sequences.

#### 4.3 Recollection

After excluding highly duplicated recitations and template-based reconstructions, what remains memorized? Despite only seeing a sample a small number of times, the model might still be able to recollect a given sample, although the factors that lead to instant memorization are poorly understood. We consider a sample to be a recollection candidate if it is a candidate for neither recitation nor reconstruction.

Recollected code, seen in Table 4, is largely made up of templating patterns that are not strictly the combination of incrementing and repetition that we use to define templates. The examples of natural language recollection in Table 3 might likewise at first appear to be misclassified recitation cases. Natural language recollection frequently comprises legal or liturgical texts, which would be expected to appear frequently throughout the corpus.

One might conjecture that these sequences are cases of retokenization, i.e., the particular token sequence is rare but the same string is heavily duplicated in the corpus under different tokenizations. However, the dependency tests in Appendix B contradict this hypothesis: the correlation between textual match count and memorization is consistently neutral or negative for recollection candidates. In



other words, a rare token sequence is *less* likely to be memorized, not more, if it is a different tokenization of a common string. We instead conjecture that the model appears to memorize slight differences in translation (liturgical text) or indexing (legal) between each variation on a sequence.

## 5 Distribution Across Scale and Time

Larger models memorize more data (Biderman et al., 2023a; Carlini et al., 2023; Tirumala et al., 2022), likely because they have more parameters with which to recreate those sequences. Recent work on deduplication (Sorscher et al., 2022) has argued that larger models are more distorted by duplication, potentially because heavily duplicated sequences are more likely to be memorized (Lee et al., 2021).

Likewise, models memorize more data as training progresses (Tirumala et al., 2022), but it is not known whether the accumulation of memorized examples is caused solely by increased exposure to heavily duplicated samples or whether other factors eventually cause memorization of rare sequences. In this section, we study the impact of training time and model size on each category of memorization.

### 5.1 Model size

Fig. 4(a) reports the number of examples memorized by each fully trained model, confirming that memorization increases with parameter count. While all types of memorization increase with model size, some increase faster than others. Recollection grows the most (Fig. 4(b)) from 4.49% of the examples memorized in the 70M model, to 11.34% in the 12B model. This disproportionate growth suggests larger models tend to memorize rarer sequences that cannot be trivially reconstructed. Meanwhile, reconstruction barely increases, indicating the smallest models have learned to extrapolate repeating and incrementing templates almost as effectively as the largest.

### 5.2 Time

Over the course of training, LMs are known to memorize an increasing pool of the training data (Tirumala et al., 2022). However, is the cumulative effect due solely to exposure to more memorizable sequences? Due to repeated exposure to the same heavily duplicated data? Or is some structural property of the later model more amenable to exact memorization? To understand why memorization

accumulates throughout training, we measure each taxonomic category in intermediate checkpoints for the 12B parameter Pythia model. We find that accumulated memorization cannot be ascribed solely to the number of available samples to memorize or to repeated exposure to highly duplicated samples.

First, in Figure 4(c), we see that models do not simply accumulate memorized samples with a uniform probability through training since memorization increases sub-linearly. Second, if memorization accumulates solely due to repeated exposure to each duplicated sample, recitation of these highly duplicated samples would be the main source of increasing memorization. Instead, the proportion of recitation decreases relative to the amount of memorization (Fig. 4(d)). Therefore, the additional memorization cannot be due to repeated exposure to recitation candidates. Instead, again the largest proportional increase among all categories is in the recollection category. This trend holds until approximately 86% of total training time, which sees a sudden increase in reconstruction. We conjecture that this increase represents a breakthrough in generalizing more complex templates but leave further investigation to future work.

Having considered and rejected both exclusive explanations, we must presume that memorization continues to occur late in training through a combination of repeated exposure, opportunities for memorizing new sequences, and other unexplored factors that may be the focus of future work.

## 6 Predicting memorization

What makes a taxonomy useful, or a reflection of natural kinds? Our position is that categories should differ in the dependencies between features of interest. The most obvious example of validated natural kinds is the case of Simpson’s Paradox (Simpson, 1951), a statistical phenomenon in which a pair of variables are correlated across a population, but the direction of correlation reverses when considering each subpopulation category separately. One famous example is Charig et al. (1986), who found that a particular surgery was more effective on both large and small kidney stones, but appeared less effective across all kidney stones without accounting for the confounding variable of size. Simpson’s Paradox is only the most obvious evidence for natural kinds, but large changes in correlation may support categorical differences even if that correlation does not change direction.

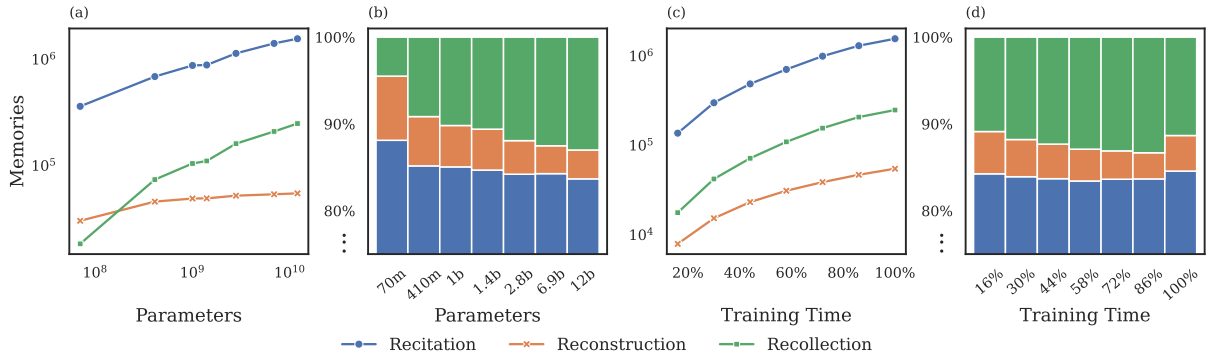


Figure 4: The quantity of memorized data categorized by taxonomy across parameter size and training time. For fully trained models of varying parameter sizes, we give (a) total counts and (b) proportion of memorized samples by category. For the 12B parameter model, we consider intermediate checkpoints during training, also providing for each checkpoint the (c) total memorized counts and (d) proportion of memorized samples by category. **Note that the proportional plots are truncated at 80%, as recitation is consistently a majority of the overall memorized data.**



Figure 5: Performance of baseline, proposed taxonomy and optimally partitioned models against various metrics on subsets of test dataset. Confidence interval is standard deviation computed by bootstrapping.

We measure a number of categorical differences in dependencies, including sign and significance differences, in Appendix B. If our intuitive taxonomy did not reflect meaningful differences with respect to the factors in Section 6.3, their dependencies would not differ significantly. We instead find significant differences through statistical tests, suggesting the taxonomy expresses some natural kinds.

Not only do these differences support our taxonomy as an ontology, but they suggest our taxonomy can help predict memorization from dependent factors. We therefore test the applicability of our taxonomy by creating a predictive model based on the intuitive taxonomic model. We compare it with a generic baseline model lacking a taxonomy and with a model using an automatically selected optimal partition, finding that our taxonomic model supports more accurate predictions.

## 6.1 Models

Each model is a logistic regression trained with L2 regularization, a bias parameter, and balanced class weights. We split the representative sample into test and train sets. We then combine the train set with the full memorized sample, reserving a portion as a validation set. For each set, continuous features are normalized to zero mean and unit variance.

**Generic baseline model** The generic baseline is a logistic regression model trained to predict whether a sample is memorized given the features from section 3. It is trained on the training split of the entire memorized dataset and the entire representative Pile sample.

**Intuitive taxonomic model (Ours)** The predictive model based on our intuitive taxonomy is made up of a set of three binary logistic regression models. We divide samples into taxonomic groups be-

fore training a separate regression on each taxonomic category.

**Optimally partitioned model** To demonstrate that our intuitive taxonomic model is not simply benefiting from having more degrees of freedom, we devise an equally complex—that is, with the same architecture of three binary logistic regressions—alternative taxonomic model. To provide a strong baseline, we search for a partition based on a set of possible feature-threshold combinations. We train predictive models with the same three-regression architecture as our intuitive taxonomic model, but partitioning based on each feature-threshold combination. The optimal partition is that which supports the best predictive model, which we find categorizes samples based on Huffman coding length followed by sequence duplicate count.

For a given feature, we consider the 25th, 50th and 75th percentiles of the value distribution distribution as potential thresholds. Each feature-threshold pair provides a possible partition split; we select the optimal three-category partition based on F1 score on the aggregate representative test set. Note that our “optimal” partition may not explore our intuitive taxonomy as an option because the threshold search is limited to each feature’s quartile values. Our intuitive taxonomy may—and does—therefore outperform the optimal partition.

## 6.2 How good is our taxonomy?

To test our intuitions, we compare our proposed taxonomy to the homogeneous baseline and to our optimal partition. As seen in Fig. 5, the greedy-optimal partition outperforms the aggregate baseline slightly on most metrics, but our intuitive taxonomy is better calibrated and more accurate except on the recollection set, where it has low precision. We conclude that our intuition has guided us to a better taxonomy than searching possible data partitions.

## 6.3 Categorical differences

Having confirmed the benefits of separately considering these three taxonomic categories, Fig. 6 shows how they differ through the feature weights from our regression models.

Recollection candidates—that is, rare sequences—are more likely to be memorized if they have no rare tokens. We posit that there is more resistance to memorizing rare tokens within

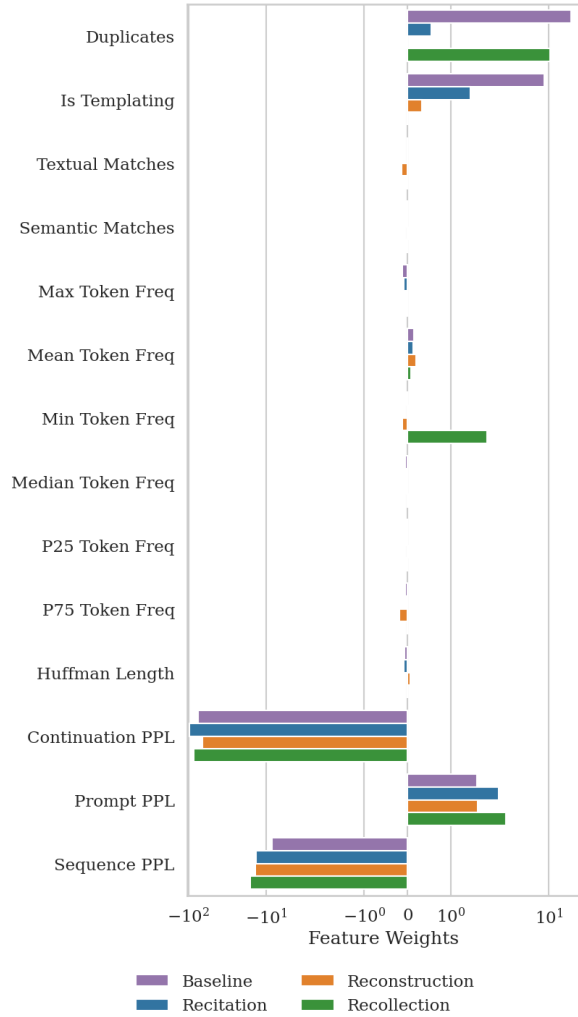


Figure 6: Feature weights from predictive models trained on the homogeneous aggregate baseline and the intuitive taxonomy categories.

a sequence, as their prior probability is low.

Meanwhile, the more duplicates a recollection candidate has, the more likely it is to be memorized, whereas recitation candidates are hardly affected by duplicate count. These results suggest that beyond the 5-duplicate threshold, greater exposure hardly leads to memorization.

Another notable difference is in the effect of perplexity: while predictable continuations are strongly associated with memorization across all categories, unpredictable prompts are strongly associated with memorization except for cases of reconstruction. The clear explanation is that high-perplexity prompts often *only* occur as prelude to the same continuation, providing a unique index for the memorized sequence, but that a low-perplexity prompt may also initiate a common template, enabling reconstruction.

506	<b>7 Discussion and future work</b>	
507	We have established that an intuitive taxonomy can	
508	be used to improve understanding of memorization.	
509	We now relate our methods to the existing literature	
510	on memorization and to possible future directions.	
511	<b>7.1 Ontologies of memorization</b>	
512	Our work is strongly related to several recent efforts	
513	to develop an ontology of memorization. Dankers	
514	et al. (2023), studying machine translation, focus	
515	primarily on the influence of a sequence during	
516	training rather than on the semantics or properties	
517	of an individual sequence. Like us, they investigate	
518	the factors that influence counterfactual memoriza-	
519	tion, the category likely to dominate “recollection”	
520	cases. They find that rare tokens, long sequence	
521	lengths, and high BPE segmentation rate are cor-	
522	related with counterfactual memorization (Zhang	
523	et al., 2021); of these, we only consider rare to-	
524	kens, which we confirm to predict recollection in	
525	particular.	
526	Hartmann et al. (2023) consider what facets of	
527	memorization are likely to be relevant to different	
528	targets, just as we discuss the differences between	
529	motivations grounded in copyright infringement	
530	and privacy. Bansal et al. (2023) consider two dif-	
531	ferent kinds of memorization: heuristic memoriza-	
532	tion, i.e., shortcut learning, and example memoriza-	
533	tion. Our work focuses on what they call example	
534	memorization, further decomposing that category.	
535	We do not test their result that high-entropy fea-	
536	tures can indicate example memorization, but like	
537	us, they use this factor to differentiate between their	
538	memorization categories.	
539	<b>7.2 Memorization and training time</b>	
540	Our work fits into an existing literature on how	
541	time and scale affect memorization. Biderman et al.	
542	(2023b) find that the position of a sequence in train-	
543	ing does not affect its likelihood of being mem-	
544	orized, and that smaller models fail to memorize	
545	even when repeatedly exposed to a term. Tiru-	
546	mala et al. (2022) find that larger models memorize	
547	more training data and forget less during training.	
548	They also observe that models memorize nouns	
549	and numbers first, using these entities as unique	
550	identifiers for individual samples. Our work further	
551	expands our understanding of scale in memoriza-	
552	tion by highlighting that rare sequences compose	
553	the fastest-growing category of memorization.	
	<b>7.3 Which categories do we care about?</b>	554
	The relevance of each category depends on our	555
	motivation for studying memorization.	556
	1. <b>Intellectual property violations:</b> The con-	557
	tent most relevant to concerns about intellec-	558
	tual property may be highly duplicated data,	559
	such as frequently excerpted passages from a	560
	popular book. However, some rare sequences	561
	may also be memorized, making recollection	562
	potentially relevant to issues of copyright in-	563
	fringement.	564
	2. <b>Privacy:</b> If the primary motivation is prevent-	565
	ing the memorization of personally identify-	566
	ing information, we may focus on recollection,	567
	as issues may arise if a model generates such	568
	information even after even a small number	569
	of exposures.	570
	3. <b>Scientific understanding of generalization:</b>	571
	Work like Henighan et al. (2023b) and Bartlett	572
	et al. (2020) points to eventual generaliza-	573
	tion as a result of memorization dynamics.	574
	A deeper understanding of these phenomena	575
	might focus on reconstruction, which exposes	576
	a direct link between apparent overfitting and	577
	general pattern recognition.	578
	<b>7.4 Ontologies and statistics</b>	579
	This taxonomy may serve as an example for fu-	580
	ture methods of interpreting complex phenomena,	581
	in deep learning and elsewhere. We have, in par-	582
	ticular, quantified the validity and usefulness of	583
	such a taxonomy by comparing predictive models	584
	which treat memorization in aggregate to models	585
	which treat memorization as a multifaceted phe-	586
	nomenon with our taxonomy. We provide evidence	587
	for the taxonomic model by measuring the improve-	588
	ment in predictive judgments when reflecting the	589
	dependent and nonlinear thresholded relationship	590
	between memorization and the properties that de-	591
	fine each taxonomic category.	592
	In future work, we hope that interpretable and	593
	useful ontologies can be validated by a similar ap-	594
	proach. Our proposal for what makes a good taxo-	595
	nomical model is not only applicable to memoriza-	596
	tion or even to deep learning phenomena. Instead,	597
	by studying interactions and nonlinearities in arbi-	598
	trary settings, researchers may find complex depen-	599
	dencies and artifacts like Simpson’s paradox.	600



## 601 Limitations

602 Our primary goal is to intuitively describe the  
603 memorization behavior with a taxonomy and con-  
604 sequently use that taxonomy to investigate how  
605 several dominant factors in memorization interact  
606 with each other. A secondary goal is to provide  
607 an example of how an ontology can be constructed  
608 and tested *in general*, as tested with our predictive  
609 models. However, these predictive models are not  
610 measurements of statistical dependency in general,  
611 instead only focusing on linear dependence. Al-  
612 though more general statistical dependencies are  
613 studied in the supplementary experiments of Ap-  
614 pendix B, the experiments in the main body of the  
615 paper assume linear dependence and so the inter-  
616 acting factors should be evaluated in the context of  
617 our supplementary dependency experiments. We  
618 believe future work inspired by our approach could  
619 improve on our work by incorporating more gen-  
620 eral dependencies.

621 Another limitation is our definition of memoriza-  
622 tion. The choice of 32-elicitation has a number of  
623 disadvantages, one of them being that we lose a  
624 notion of fuzzy or partial memorization, which is  
625 considered important in some contexts. Arguably,  
626 under a counterfactual memorization definition, we  
627 may not see substantial patterns of either recitation  
628 or reconstruction. The measurement of memoriza-  
629 tion is a large area of research with many possible  
630 definitions to choose from (Carlini et al., 2022a;  
631 Tirumala et al., 2022; Kandpal et al., 2022; Zhang  
632 et al., 2021; Zhao et al., 2022; Stock et al., 2022;  
633 Schwarzschild et al., 2024).<sup>3</sup>

## 634 References

- 635 Rachit Bansal, Danish Pruthi, and Yonatan Belinkov.  
636 2023. [Measures of information reflect memorization  
637 patterns](#). *Preprint*, arxiv:2210.09404 [cs, math].
- 638 Peter L Bartlett, Philip M Long, Gábor Lugosi, and  
639 Alexander Tsigler. 2020. Benign overfitting in linear  
640 regression. *Proceedings of the National Academy of  
641 Sciences*, 117(48):30063–30070.
- 642 Stella Rose Biderman, USVSN Sai Prashanth, Lintang  
643 Sutawika, Hailey Schoelkopf, Quentin G. Anthony,  
644 Shivanshu Purohit, and Edward Raf. 2023a. Emer-  
645 gent and predictable memorization in large language  
646 models. *ArXiv*, abs/2304.11158.
- 647 Stella Rose Biderman, Hailey Schoelkopf, Quentin G.  
648 Anthony, Herbie Bradley, Kyle O’Brien, Eric Hal-

- 649 lahan, Mohammad Aflah Khan, Shivanshu Purohit,  
650 USVSN Sai Prashanth, Edward Raff, Aviya Skowron,  
651 Lintang Sutawika, and Oskar van der Wal. 2023b.  
652 Pythia: A suite for analyzing large language models  
653 across training and scaling. *ArXiv*, abs/2304.01373.
- 654 Hannah Brown, Katherine Lee, FatemehSadat  
655 Mireshghallah, R. Shokri, and Florian Tramèr. 2022.  
656 [What does it mean for a language model to preserve  
657 privacy?](#) *Proceedings of the 2022 ACM Conference  
658 on Fairness, Accountability, and Transparency*.
- 659 Nicholas Carlini, Daphne Ippolito, Matthew Jagielski,  
660 Katherine Lee, Florian Tramèr, and Chiyuan Zhang.  
661 2022a. Quantifying memorization across neural lan-  
662 guage models. *ArXiv*, abs/2202.07646.
- 663 Nicholas Carlini, Daphne Ippolito, Matthew Jagielski,  
664 Katherine Lee, Florian Tramer, and Chiyuan Zhang.  
665 2023. [Quantifying Memorization Across Neural Lan-  
666 guage Models](#). *arXiv preprint*. ArXiv:2202.07646  
667 [cs].
- 668 Nicholas Carlini, Matthew Jagielski, Nicolas Papernot,  
669 A. Terzis, Florian Tramèr, and Chiyuan Zhang. 2022b.  
670 The privacy onion effect: Memorization is relative.  
671 *ArXiv*, abs/2206.10469.
- 672 Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej  
673 Kos, and Dawn Xiaodong Song. 2018. The secret  
674 sharer: Evaluating and testing unintended memo-  
675 rization in neural networks. In *USENIX Security  
676 Symposium*.
- 677 Nicholas Carlini, Florian Tramèr, Eric Wallace,  
678 Matthew Jagielski, Ariel Herbert-Voss, Katherine  
679 Lee, Adam Roberts, Tom B. Brown, Dawn Xiaodong  
680 Song, Úlfar Erlingsson, Alina Oprea, and Colin Raf-  
681 fel. 2020. Extracting training data from large lan-  
682 guage models. In *USENIX Security Symposium*.
- 683 Clive R Charig, David R Webb, Stephen Richard  
684 Payne, and John E Wickham. 1986. Compari-  
685 son of treatment of renal calculi by open surgery,  
686 percutaneous nephrolithotomy, and extracorporeal  
687 shockwave lithotripsy. *Br Med J (Clin Res Ed)*,  
688 292(6524):879–882.
- 689 Verna Dankers, Ivan Titov, and Dieuwke Hupkes.  
690 2023. [Memorisation Cartography: Mapping  
691 out the Memorisation-Generalisation Continuum  
692 in Neural Machine Translation](#). *arXiv preprint*.  
693 ArXiv:2311.05379 [cs].
- 694 Vitaly Feldman. 2021. [Does Learning Require Mem-  
695 orization? A Short Tale about a Long Tail](#). *arXiv  
696 preprint*. ArXiv:1906.05271 [cs, stat].
- 697 Leo Gao, Stella Rose Biderman, Sid Black, Laurence  
698 Golding, Travis Hoppe, Charles Foster, Jason Phang,  
699 Horace He, Anish Thite, Noa Nabeshima, Shawn  
700 Presser, and Connor Leahy. 2020. The pile: An  
701 800gb dataset of diverse text for language modeling.  
702 *ArXiv*, abs/2101.00027.

<sup>3</sup>For a discussion of definitions of memorization, see  
<https://genlaw.org/glossary.html#memorization>.

703	Valentin Hartmann, Anshuman Suri, Vincent Bindschaedler, David Evans, Shruti Tople, and Robert West. 2023. <a href="#">SoK: Memorization in general-purpose large language models</a> . <i>Preprint</i> , arxiv:2310.18362 [cs].	758
704		759
705		760
706		761
707		
708	Tom Henighan, Shan Carter, Tristan Hume, Nelson Elhage, Robert Lasenby, Stanislav Fort, Nicholas Schiefer, and Christopher Olah. 2023a. Superposition, memorization, and double descent. <i>Transformer Circuits Thread</i> .	762
709		763
710		764
711		765
712		
713	Tom Henighan, Shan Carter, Tristan Hume, Nelson Elhage, Robert Lasenby, Stanislav Fort, Nicholas Schiefer, and Christopher Olah. 2023b. <a href="#">Superposition, Memorization, and Double Descent</a> .	766
714		767
715		768
716		769
717	David A Huffman. 1952. A method for the construction of minimum-redundancy codes. <i>Proceedings of the IRE</i> , 40(9):1098–1101.	770
718		771
719		772
720	Daphne Ippolito, Florian Tramèr, Milad Nasr, Chiyuan Zhang, Matthew Jagielski, Katherine Lee, Christopher A Choquette-Choo, and Nicholas Carlini. 2022. Preventing verbatim memorization in language models gives a false sense of privacy. <i>arXiv preprint arXiv:2210.17546</i> .	773
721		774
722		775
723		776
724		777
725		778
726	Nikhil Kandpal, Eric Wallace, and Colin Raffel. 2022. <a href="#">Deduplicating training data mitigates privacy risks in language models</a> . <i>ArXiv</i> , abs/2202.06539.	779
727		780
728		781
729	Antonia Karamolegkou, Jiaang Li, Li Zhou, and Anders Sogaard. 2023. <a href="#">Copyright violations and large language models</a> . <i>ArXiv</i> , abs/2310.13771.	782
730		783
731		784
732	Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2021. <a href="#">Deduplicating training data makes language models better</a> . In <i>Annual Meeting of the Association for Computational Linguistics</i> .	785
733		786
734		787
735		788
736		789
737	Vladimir I Levenshtein et al. 1966. Binary codes capable of correcting deletions, insertions, and reversals. In <i>Soviet physics doklady</i> , volume 10, pages 707–710. Soviet Union.	790
738		791
739		792
740		793
741	Neil Mallinar, James B. Simon, Amirhesam Abedsoltan, Parthe Pandit, Mikhail Belkin, and Preetum Nakkiran. 2022. <a href="#">Benign, tempered, or catastrophic: A taxonomy of overfitting</a> . <i>Preprint</i> , arXiv:2207.06569.	794
742		795
743		796
744		797
745	Matthieu Meeus, Igor Shilov, Manuel Faysse, and Yves-Alexandre de Montjoye. 2024. <a href="#">Copyright traps for large language models</a> . <i>ArXiv</i> , abs/2402.09363.	798
746		799
747		800
748	Fatemehsadat Mireshghallah, Archit Uniyal, Tianhao Wang, David Evans, and Taylor Berg-Kirkpatrick. 2022. <a href="#">An empirical analysis of memorization in fine-tuned autoregressive language models</a> . In <i>Conference on Empirical Methods in Natural Language Processing</i> .	801
749		802
750		803
751		804
752		805
753		806
754	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2020. <a href="#">Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter</a> . <i>Preprint</i> , arXiv:1910.01108.	807
755		808
756		809
757		810
	Avi Schwarzschild, Zhili Feng, Pratyush Maini, Zachary C. Lipton, and J. Zico Kolter. 2024. <a href="#">Rethinking llm memorization through the lens of adversarial compression</a> . <i>Preprint</i> , arXiv:2404.15146.	
	Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2023. <a href="#">Detecting pretraining data from large language models</a> . <i>ArXiv</i> , abs/2310.16789.	
	Edward H Simpson. 1951. The interpretation of interaction in contingency tables. <i>Journal of the Royal Statistical Society: Series B (Methodological)</i> , 13(2):238–241.	
	Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. 2022. Beyond neural scaling laws: beating power law scaling via data pruning. <i>Advances in Neural Information Processing Systems</i> , 35:19523–19536.	
	Pierre Stock, Igor Shilov, Ilya Mironov, and Alexandre Sablayrolles. 2022. <a href="#">Defending against reconstruction attacks with rényi differential privacy</a> . <i>Preprint</i> , arXiv:2202.07623.	
	Kushal Tirumala, Aram H. Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. 2022. <a href="#">Memorization Without Overfitting: Analyzing the Training Dynamics of Large Language Models</a> . <i>arXiv:2205.10770 [cs]</i> . <i>ArXiv</i> : 2205.10770.	
	Zhou Yang, Zhipeng Zhao, Chenyu Wang, Jieke Shi, Dongsun Kim, Donggyun Han, and David Lo. 2024. <a href="#">Unveiling memorization in code models</a> . In <i>Proceedings of the IEEE/ACM 46th International Conference on Software Engineering, ICSE '24</i> , New York, NY, USA. Association for Computing Machinery.	
	Chiyuan Zhang, Daphne Ippolito, Katherine Lee, Matthew Jagielski, Florian Tramèr, and Nicholas Carlini. 2021. <a href="#">Counterfactual memorization in neural language models</a> . <i>ArXiv</i> , abs/2112.12938.	
	Xuandong Zhao, Lei Li, and Yu-Xiang Wang. 2022. <a href="#">Provably confidential language modelling</a> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 943–955, Seattle, United States. Association for Computational Linguistics.	
	Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In <i>The IEEE International Conference on Computer Vision (ICCV)</i> .	
	<b>A Implementation of metrics</b>	
	We now provide details of the implementation for metrics considered as potential factors for memorization.	

811	<b>A.1 Number of exact duplicate samples</b>	<b>A.4 Incrementing and Repeating templates</b>	853
812	We compute the number of exact duplicates as a	<b>A.4.1 Incrementing Templates</b>	854
813	3-step process:	To check for an incrementing sequence, we perform	855
814	1. For every 32-gram window in the data point	the following steps:	856
815	$S$ (comprised of 2049 tokens) we compute a	• Split the text by whitespace and convert any	857
816	rolling hash and store it, along with the win-	splits which are numerals in non-decimal	858
817	down’s index and offset in a parallelized set of	bases (e.g., hexadecimal) into base 10.	859
818	data frames.	• Remove escape sequences.	860
819	2. For every index position, we compute the	• Within each string, separate contiguous nu-	861
820	same hash for all $S[32: 64]$ (sequence con-	merical characters from anything else. If two	862
821	tinuations) and store all (index, offset, hash)	contiguous numeric characters are separated	863
822	tuples if their hash is one of the computed	by a period, combine them into their floating	864
823	sequences continuation hashes.	point representations.	865
824	3. Now, for every sequence continuation, we	• Discard if there are fewer than 3 potential	866
825	look at all 32-gram windows with the same	numerals in the sequence.	867
826	hash and compute their number of duplicates	• Check if the sequences are incrementing or	868
827	as the number of equivalent (same set of to-	repeating.	869
828	kens in the same order) samples.	<b>A.4.2 Repeating templates</b>	870
829	The hash function used is similar to Rabin-	We perform the following steps to check for repeat-	871
830	Karp’s rolling hash algorithm. Specifically, con-	ing sequences:	872
831	sider a token sequence of 32 tokens.	• Obtain a sequence by splitting the text by char-	873
832	$S = [c_1, c_2, c_3, \dots, c_{32}]$	acter.	874
833	Let us define two primes $P = 60013$ and $MOD =$	• Check if the sequences are incrementing or	875
834	$10^{18} + 3$ . We define their hash function to be	repeating.	876
835	$H(S) = (c_1 + c_2 * P + c_3 * P^2 + \dots + c_{32} * P^{31}) \% MOD$	We perform the following steps to determine if a	877
836	<b>A.2 Token frequency</b>	sequence generated from either of the above steps	878
837	Token frequencies are calculated across the Pile.	is incrementing or repeating.	879
838	For every sequence continuation, we consider the	• For every <i>templating length</i> , defined to be less	880
839	maximum, minimum, median and quartile frequen-	than half length of splits, and for every <i>po-</i>	881
840	cies of tokens.	sition less than <i>templating length</i> , we iterate	882
841	<b>A.3 Compressibility</b>	through splits with start position as <i>position</i>	883
842	We use Huffman Coding length to measure how	and step size set to <i>templating length</i> . We then	884
843	easily a sequence can be compressed. Compress-	determine if the current iteration is repeating	885
844	ibility provides a rough generalization of internal	or incrementing.	886
845	repetition, where only a few exceptions to some	• For example, if position is 1 and <i>templat-</i>	887
846	simple repetition pattern might need to be memo-	<i>ing length</i> is 5, we iterate through positions	888
847	rized. However, unlike straightforward repetition	[1, 6, 11, 16...]. if our input splits length is 10,	889
848	templates, compressible sequences may not be con-	we iterate for all <i>templating lengths</i> 1 through	890
849	sidered to be reconstructed by the model. Instead,	5 and for all positions less than current <i>tem-</i>	891
850	we include compressibility as a filter to evaluate	<i>plating length</i>	892
851	whether LLMs memorize samples that are easier to	• Within each iteration, we check:	893
852	compress into their parameters.	– If the current iteration has both texts and	894
		numerals, it is neither incrementing nor	895
		repeating.	896

897	– If the current iteration has only texts, we	<b>D Classifying examples as natural</b>	941
898	consider the current iteration to be repeat-	<b>language or code</b>	942
899	ing if all elements in the iteration are the	To train a Natural Language vs Code classifier,	943
900	same.	we fine-tune DistilBert (Sanh et al., 2020) on uni-	944
901	– If the current iteration has only numer-	formly random sampled Bookcorpus (Zhu et al.,	945
902	als, we consider current iteration to be	2015) and github-code datasets. We train it with	946
903	incrementing if all the numerals are in an	learning rate of $10^{-7}$ and batch size of 256 for a	947
904	arithmetic progression. If the difference	total of 1000 steps and observe validation f1 score	948
905	in AP is 0, we consider it to be repeating	of 0.9950 on a held of evaluation set.	949
906	instead.	To select an optimal threshold for this classifier	950
907	• Input splits are considered as repeating if all	on memories dataset, we randomly sample 500 se-	951
908	iterations for a given <i>templating length</i> are	quences and manually label them. To make sure	952
909	repeating.	that precision is high for our models, we choose	953
910	• Input splits are considered as incrementing if	$\leq 0.4$ as threshold for determining code samples	954
911	atleast one of the iterations for a given <i>tem-</i>	and a threshold of $\geq 0.525$ for determining natural	955
912	<i>plating length</i> are incrementing and others, for	language samples, based on the points marked in	956
913	the same <i>templating length</i> , are either incre-	Figure 10, which mark points of near 100% preci-	957
914	menting (or) repeating.	sion for classifying each category.	958
915	• For all templating lengths, if any length of	<b>E Likelihood of memorization for code</b>	959
916	them has been found to be incrementing or	<b>and natural language</b>	960
917	repeating, we return <i>True</i> (corresponding to	We study the likelihood that a sample that has been	961
918	the fact that the text is indeed a template) and	confidently classifier as code or NL (Appendix D)	962
919	<i>diff</i> .	is memorized across time and scale. For example,	963
920	• Note that, in the case of sequences generated	for all samples confidently classified as code Fig-	964
921	while checking for a repeating template, we	ure, what is the proportion of samples which are	965
922	do not have any numerals.	memorized?	966
923	<b>B Dependency Tests for Influence of</b>	Figure 11 shows that code samples are more	967
924	<b>Features on Memorization</b>	likely to be memorized than NL across categories.	968
925	This section contains visualizations of various de-	This trend suggests that certain intrinsic factors	969
926	pendency tests between memorization likelihood	about code make it more susceptible to memoriza-	970
927	and our target features. We look at dependencies on	tion, even for recollection samples where memo-	971
928	code (Fig. 9), natural language (Fig. 8), and both	rization cannot be attributed to obvious patterns	972
929	(Fig. 7). These tests are more general and have	and high duplication. Both code and NL become	973
930	stronger guarantees than simply looking at regres-	more likely to be memorized across scale, except	974
931	sion weights, which have a number of flaws. For	for reconstruction samples, which remain compara-	975
932	example, we see in Fig. 6 that regression can reallo-	tively unchanged.	976
933	cate bias terms to features that take on a consistent	<b>F Examples of memorized continuation</b>	977
934	value, giving them spurious weight.	<b>sequences</b>	978
935	<b>C Tables for scaling experiments</b>	Table 3 provides examples of memorized natural	979
936	Figure 4 visualizes the count and proportion of	language text in each memorization category and	980
937	memorized samples by category across time and	Table 4 provides examples of memorized code.	981
938	scale. We present the raw statistics for each taxo-	Samples are classified using the methodology in	982
939	onomic category across model size in Table 1 and	Appendix D.	983
940	training time in Table 2.		



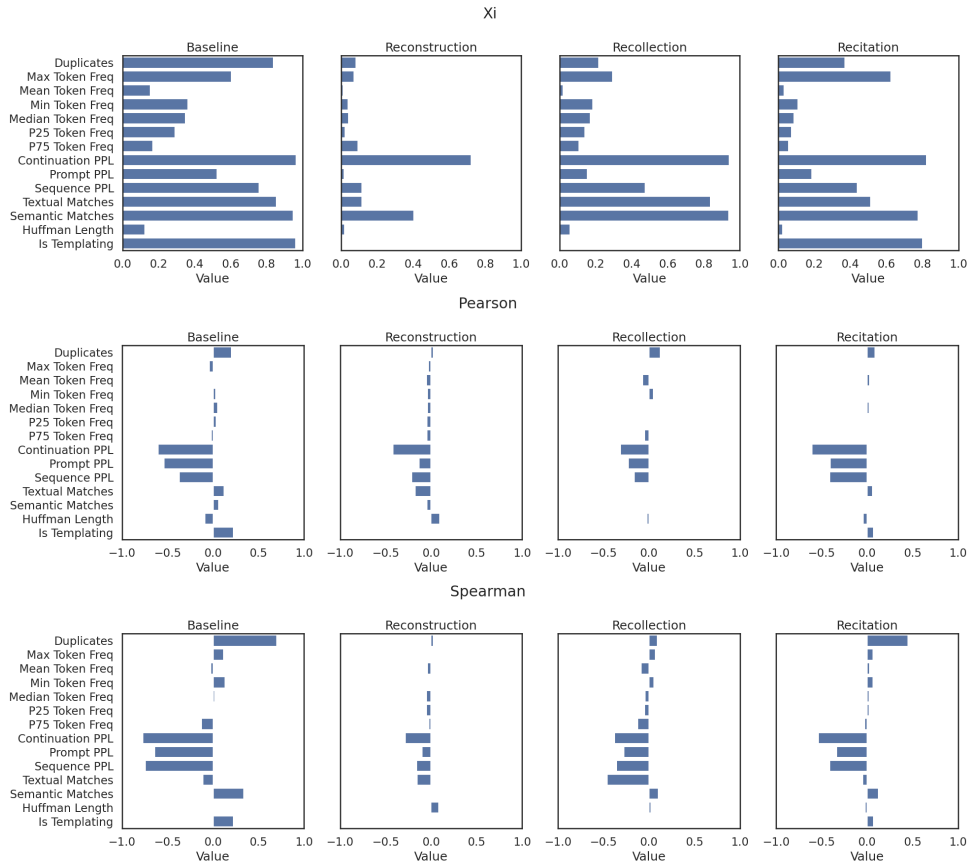


Figure 7: Dependency measurements between influence factors and memorization.

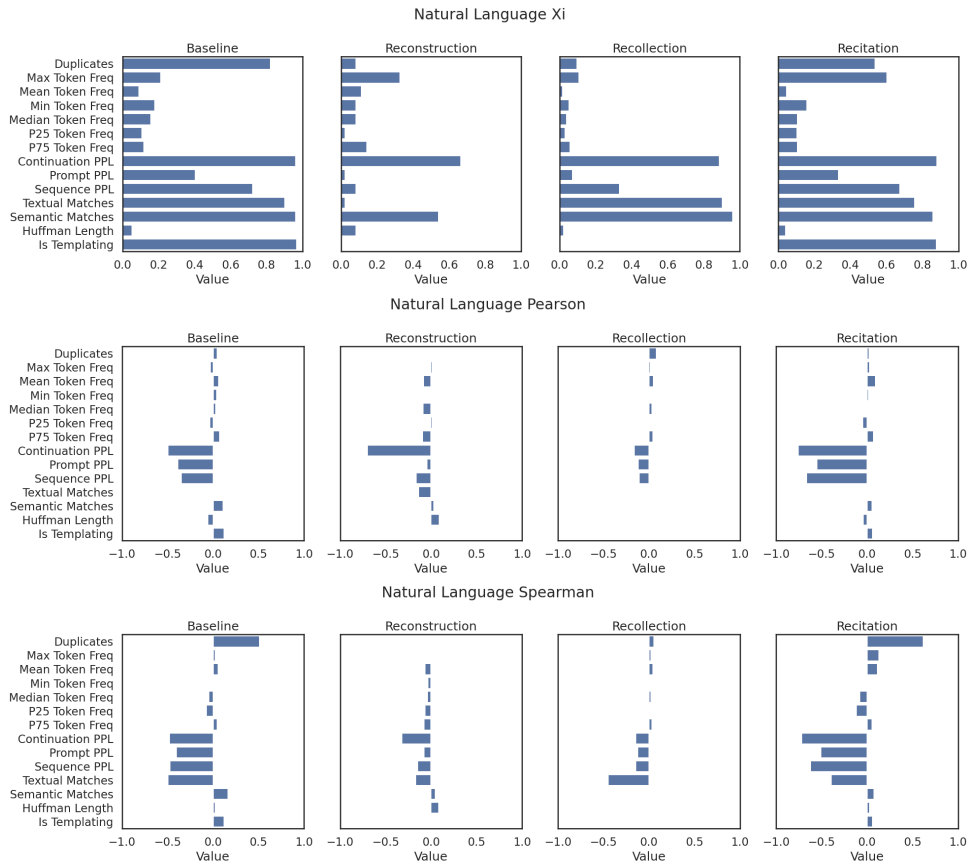


Figure 8: Dependency measurements between influence factors and memorization for natural language samples.

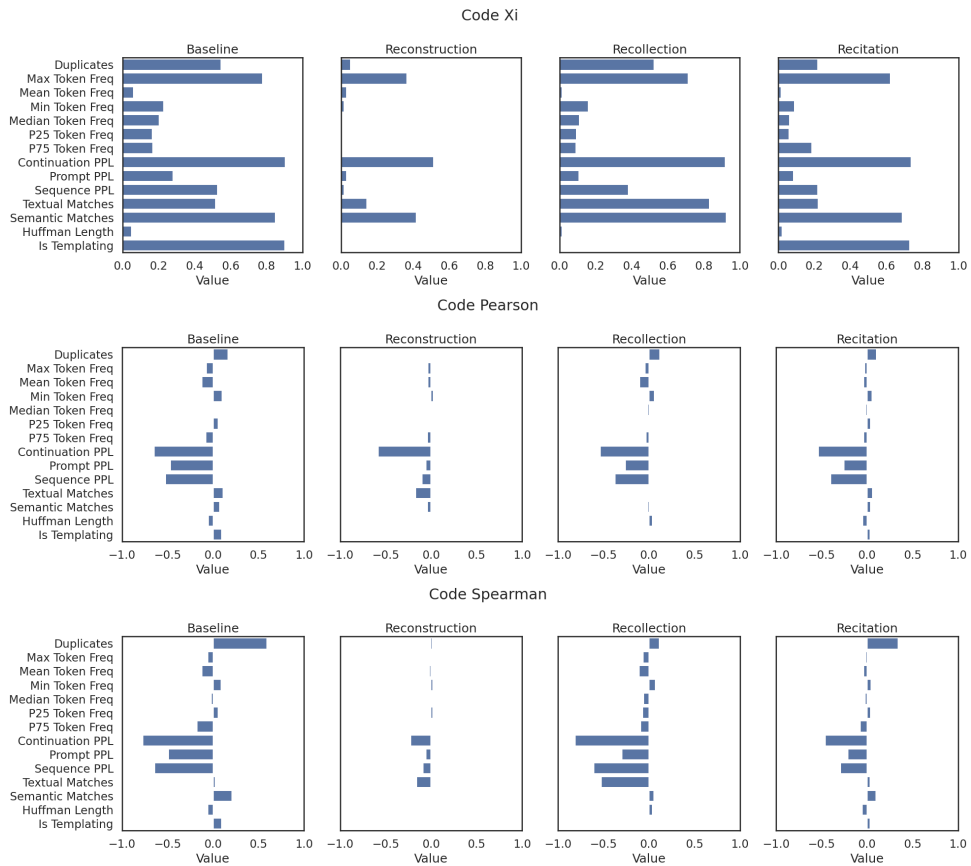


Figure 9: Dependency measurements between influence factors and memorization for code samples.

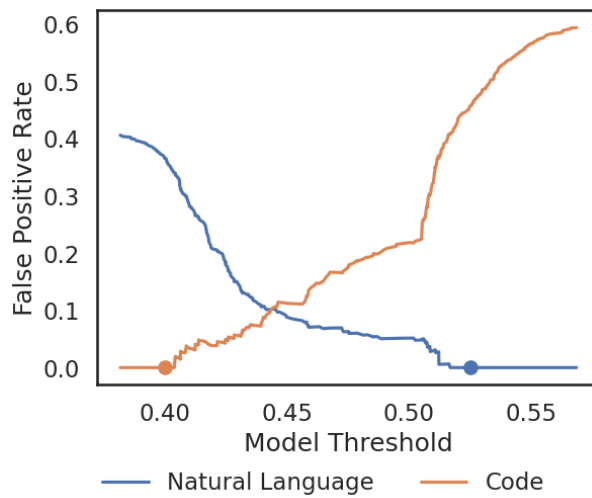


Figure 10: False positive rates across various thresholds on randomly sampled sequences of Pile. We choose  $\leq 0.4$  as threshold for determining code samples and a threshold of  $\geq 0.525$  for determining natural language

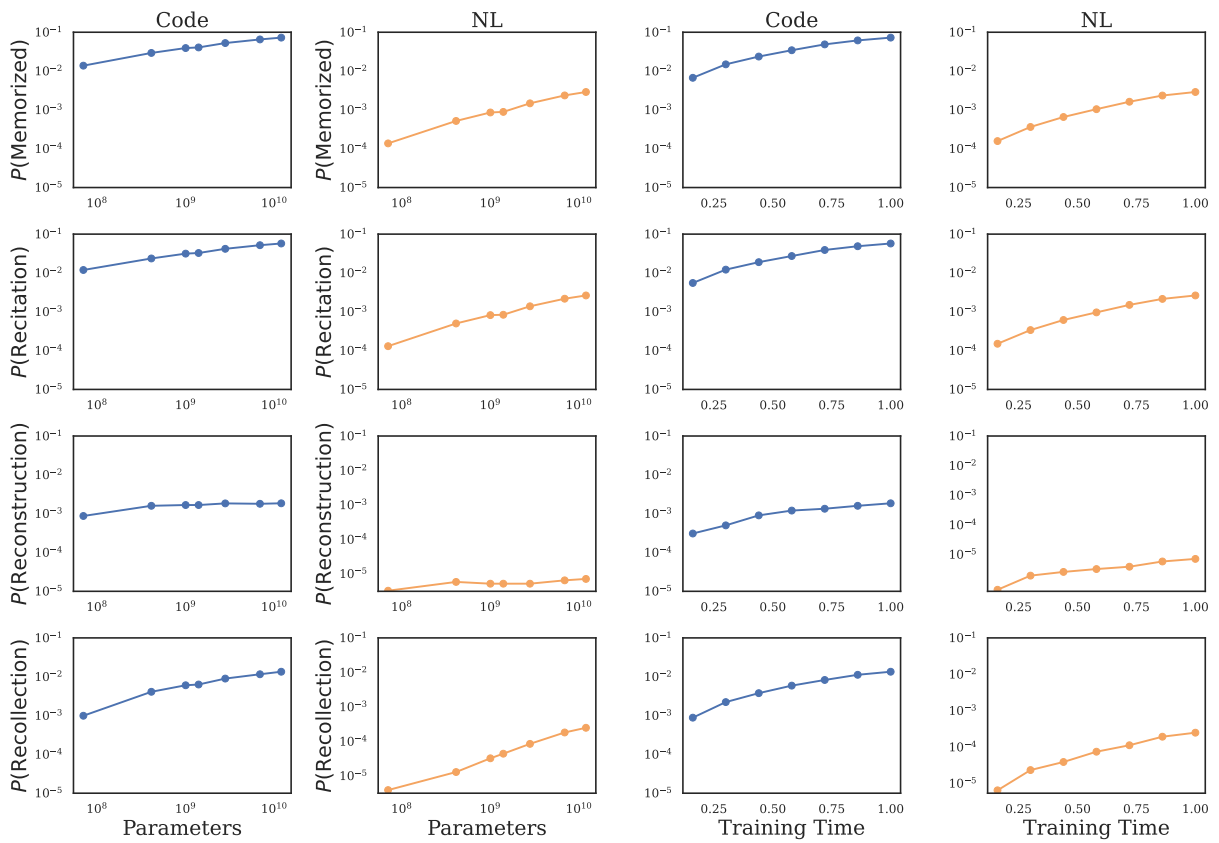


Figure 11: We study how likely models are to memorize samples confidently classified as code or NL. We calculate the likelihood for each distribution (code vs. NL) separately. Figures include probability across a model scale and training time. Models memorize a greater proportion of code samples than NL across all categories, model scale, and training time.

Model	Recitation		Reconstruction		Recollection	
	Count	Percent	Count	Percent	Count	Percent
70m	362,550.0	88.12%	30,430.0	7.40%	18,468.0	4.49%
410m	690,726.0	85.17%	46,076.0	5.68%	74,238.0	9.15%
1b	878,456.0	85.05%	49,253.0	4.77%	105,163.0	10.18%
1.4b	887,549.0	84.68%	49,435.0	4.72%	111,120.0	10.60%
2.8b	1,141,180.0	84.21%	52,416.0	3.87%	161,620.0	11.93%
6.9b	1,416,014.0	84.27%	53,968.0	3.21%	210,314.0	12.52%
12b	1,566,369.0	84.56%	55,114.0	4.10%	249,733.0	11.34%

Table 1: The number of memorized samples for each taxonomic category across model size. These results are visualized in Figure 4(a) and 4(b).

Checkpoint	Recitation		Reconstruction		Recollection	
	Count	Percent	Count	Percent	Count	Percent
16%	137,681.0	84.25%	7,960.0	4.87%	17,777.0	10.88%
30%	301,146.0	83.92%	15,379.0	4.29%	42,338.0	11.80%
44%	489,629.0	83.69%	23,333.0	3.99%	72,105.0	12.32%
58%	710,823.0	83.42%	31,260.0	3.67%	109,985.0	12.91%
72%	999,867.0	83.63%	38,999.0	3.26%	156,712.0	13.11%
86%	1,308,538.0	83.66%	47,145.0	3.01%	208,372.0	13.32%
100%	1,566,369.0	84.56%	55,114.0	4.10%	249,733.0	11.34%

Table 2: The number of memorized samples for each taxonomic category across training time for Pythia 12b. 14,000 is the final checkpoint. These results are visualized in Figure 4(c) and 4(d).



Category	Text	Count
Recitation	-11 NASB Or do you not know that the unrighteous will not inherit the kingdom of God? Do not be deceived; neither fornicators, nor idolaters, nor adulterers, nor effeminate, nor homosexuals, nor thieves, nor the covetous, nor drunk	175
Recitation	2d 1234, 1238 (8th Cir.1990). On the other hand, the Federal Rules of Civil Procedure have authorized for nearly 60 years "motions for summary judgment upon proper showings of the lack of a genuine, triable issue of material fact." Celotex Corp. v. Catrett	86
Recitation	view of life, food, cocktails, fitness, and fun. This blog is just a regular guy's view of life, food, cocktails, fitness, and fun. My opinions, musings, observations, rantings, ravings, foodie adventures, and overall humorous pontification of	52
Recitation	://www.harpercollins.ca> New Zealand HarperCollinsPublishers (New Zealand) Limited P.O. Box 1 Auckland, New Zealand <http://www.harpercollins.co.nz> United Kingdom Har	2303
Recitation	000, from the tribe of Simeon 12,000, from the tribe of Levi 12,000, from the tribe of Issachar 12,000, from the tribe of Zebulun 12,000, from the tribe of Joseph 12,000,	7
Reconstruction	THEIR EAGLE ONAND THE DIRTY BIRDS READY AS WELL DOWN AND GET THEIR EAGLE ONAND THE DIRTY BIRDS READY AS WELL DOWN AND GET THEIR EAGLE ONAND THE DIRTY B	4
Reconstruction	181  182  183  184  185  186  187  188  189  190  191  192  193  194  195  196  197  198  199  200  201  202	2
Reconstruction	4: 1482 4: 1483 4: 1484 4: 1485 4: 1486 4: 1487 4: 1488 4: 1489 4: 1490 4: 1491 4: 1492 4: 1493 4: 14	2
Reconstruction	1970–71 Turkish Third Football League season Promotion and relegation: 1971–72 Turkish Third Football League season Promotion and relegation: 1972–73 Turkish Third Football League season Promotion and relegation: 1973–74 Turkish Third Football	1
Reconstruction	, 28-63, 28-64, 28-65, 28-66, 28-67, 28-68, 28-69, 28-70, 28-71, 28-72, 28-73, 28-74, 28-75, 28-76, 28-77, 28-78	3
Recollection	affect the child; ¶70 (b) The wishes of the child, as expressed directly by the child or through the child's guardian ad litem, with due regard for the maturity of the child; ¶71 (c) The custodial history of the child,	2
Recollection	will be too late! ” 24 “Strive to enter through the narrow door. For many, I tell you, will seek to enter and will not be able. 25 When once the master of the house has risen and shut the door, and you begin to stand outside and to knock at the door,	4
Recollection	, and freedom.Revava Revava (), is an Orthodox Jewish Israeli settlement in the West Bank. Located between Barkan and Karnei Shomron, it falls under the jurisdiction of Shomron Regional Council. In it had a population of. The international community considers Israeli settlements in	2
Recollection	§ 2254(d)). A state-court decision is considered “contrary to... clearly established Federal law” if the two are “diametrically different, opposite in character or nature, or mutually opposed.” Williams v. Taylor, 529 U.S. 362, 405 (	5

Table 3: Random examples of natural language (as classified per Appendix D) from each category of memorization.

Category	Text	Count
Recitation	>-task"> <a class="nav-group-task-link" href=" ../Extensions/Int.html">Int</a> </li> <li class="nav-group-task"> <a class="nav-group-task-link" href=" ../Extensions/	5310
Recitation	> <widget class="GtkButton" id="entCleanBut"> <property name="label" translatable="yes">... :</strong></p> <ul> <li> <p>Must be one of: <code>>true</code>, <code>>false</code>, <code>1</code>, <code>0</code>.</p> </li> </ul>	3689
Recitation	= null, CancellationToken cancellationToken = default(CancellationToken)) { if (Client.SubscriptionId == null) { throw new ValidationException(ValidationRules.CannotBeNull, "this.Client.SubscriptionId"); } if	227
Recitation	.Object</h3> <code>>equals, getClass, hashCode, notify, notifyAll, wait, wait, wait</code></li> </ul> </li> </ul> </li> </ul> </div> <div class="details">	18963
Reconstruction	FileEntry("/base1/dir1/", fe, age); fe.name = "file3"; fi -> updateFileEntry("/base1/dir1/", fe, age); fe.name = "file4"; fi -> updateFileEntry("/base1/	2
Reconstruction	= "time2[]" value="5" ></td> <td><input type="checkbox" name="time2[]" value="6" ></td> <td><input type="checkbox" name="time2[]" value="7" ></td> <	2
Reconstruction	XMM3 (1ULL << 28) #define DBG_CTX_EX_PART_FLAG_XMM4 (1ULL << 29) #define DBG_CTX_EX_PART_FLAG_XMM5 (1ULL << 30) #define DBG_	2
Reconstruction	" /> <Compile Include="Message\MFN_M06.cs" /> <Compile Include="Message\MFN_M07.cs" /> <Compile Include="Message\M08.cs" /> <Compile Include="Message\MFN_	3
Recollection	OFFSET + (2 * FPREG_SIZE)) #define PROBE_CPU_Q3_OFFSET (PROBE_FIRST_FPREG_OFFSET + (3 * FPREG_SIZE)) #define PROBE_CPU_Q4_OFFSET (PROBE_FIRST_FPREG_OFFSET +	2
Recollection	); } uint16_t WS2812FX::mode_custom_5() { return customModes[5]; } uint16_t WS2812FX::mode_custom_6() { return customModes[6]; } uint16	2
Recollection	XMMM128, __) /* 0xDC */ NORMAL("paddusb", MM_XMM, MMM64_XMMM128, __) /* 0xDD */ NORMAL("paddusw", MM_XMM, MMM64_X	2
Recollection	DBF3B /* icon1.png */; }; 651A5A7E177AE2D8003DBF3B /* icon2.png in Resources */ = {isa = PBXBuildFile; fileRef = 651A5A7C177AE2D8003DBF	2

Table 4: Random examples of code (as classified per Appendix D) from each category of memorization.