

A Multi-Perspective Analysis of Memorization in Large Language Models

Bowen Chen, Namgi Han, Yusuke Miyao

Department of Computer Science, The University of Tokyo
{bwchen, hng88, yusuke}@is.s.u-tokyo.ac.jp

Abstract

Large Language Models (LLMs) can generate the same sequences contained in the pre-train corpora, known as memorization. Previous research studied it at a macro level, leaving micro yet important questions under-explored, e.g., what makes sentences memorized, the dynamics when generating memorized sequence, its connection to unmemorized sequence, and its predictability. We answer the above questions by analyzing the relationship of memorization with outputs from LLM, namely, embeddings, probability distributions, and generated tokens. A memorization score is calculated as the overlap between generated tokens and actual continuations when the LLM is prompted with a context sequence from the pre-train corpora. Our findings reveal: (1) The inter-correlation between memorized/unmemorized sentences, model size, continuation size, and context size, as well as the transition dynamics between sentences of different memorization scores, (2) A sudden drop and increase in the frequency of input tokens when generating memorized/unmemorized sequences (*boundary effect*), (3) Cluster of sentences with different memorization scores in the embedding space, (4) An *inverse boundary effect* in the entropy of probability distributions for generated memorized/unmemorized sequences, (5) The predictability of memorization is related to model size and continuation length. In addition, we show a Transformer model trained by the hidden states of LLM can predict unmemorized tokens.¹

1 Introduction

Large Language Models (LLMs), trained with enormous parameter and pre-train data sizes, like GPT-4 (OpenAI et al., 2024), show surprising performance on various tasks. Due to such enormous model size

and pre-train data size, combined with the black-box nature of neural models (Alain and Bengio, 2018), LLMs present unique behaviors (Wei et al., 2022) that are unprecedented in previous machine learning, one of which is *memorization*.

Memorization (Hartmann et al., 2023) in the LLMs means the LLM can generate the same content recorded in their pre-train corpus. Being a coin with two sides, memorization in LLM can provide knowledge (Petroni et al., 2019) or cause personal information leakage (Yao et al., 2024). Previous research (Tirumala et al., 2022; Carlini et al., 2023; Biderman et al., 2023b) has studied memorization at the macro level, revealing memorization from the training phase or overlap between models, leaving more micro yet important questions left under-explored, e.g., what makes sentences memorized, dynamics of sentences with varying memorization scores transit to other scores when trained in a larger model, the dynamics when generating memorized/unmemorized sequence, its connection to unmemorized sequence, and its predictability.

To answer those questions, we study memorization from multiple perspectives. We form contexts of varying lengths from the pre-train corpora and input them into LLMs of different sizes. By controlling the maximum generated tokens, we collect the outputs (embeddings, decoding probabilities, and generated tokens) and compute the *memorization score* by comparing them with the actual continuations in the pre-train corpora. Our analysis connects memorization to model size, context size, continuation size, and unmemorized sequences. We also explored the dynamics of generating memorized/unmemorized tokens at input and output levels and the predictability of memorization. Our findings reveal:

(I) For both memorized or unmemorized sentences, their increase or decrease with model size is non-linear, indicating a maximum capacity for memorization. Memorized sentences decrease sub-

¹The code of this study is at <https://github.com/mynlp/memorizationstudy>

linearly with continuation size and increase super-linearly with context size. We also analyzed the dynamics of how sentences with varying memorization scores transit to other scores when trained in a larger model. The results show that only limited low memorization score sequences will transit to higher scores when trained in larger models, and most memorized sentences are also inherited when trained in larger models.

(II) A *boundary effect* was observed when the model began generating memorized/unmemorized sequences. The n-gram frequency suddenly decreases when generating unmemorized tokens and suddenly increases when generating memorized tokens, which is also observed at the sequence level. The significance of the boundary effect for memorized sentences decreases with the increase in model size. However, as large models contain more memorized sequences, this suggests large models have a lower threshold in the value of the boundary effect to memorize a sentence. This indicates the value of the boundary effect relates to the difficulty of memorizing a sentence.

(III) The embedding dynamics analysis showed sentences with different memorization scores cluster in the embedding space, where the mutual embedding distance grows with model size. Sentences of close memorization scores are also close in embedding space, suggesting the existence of paraphrase memorization.

(IV) We analyzed decoding dynamics by examining entropy over vocabulary and the drift of decoded embeddings. Entropy analysis revealed an *inverse boundary effect*, where entropy suddenly increases for unmemorized sequences and decreases for memorized sequences.

(V) We trained a Transformer model to discuss the predictability of memorization. The results suggest that predicting memorization is easier in large models and easier when predicting unmemorized sequences, which can be explained by the significance of the boundary effect.

2 Related Works

2.1 Scaling Laws of LLM

In this study, our experiments span across various model sizes. This relates to the research of Scaling Laws (Kaplan et al., 2020; Abnar et al., 2021; Vilalobos, 2023), which suggests the performance of LLM scales with the corpora size, the parameter size, and the computation required.

Inspired by the Scaling Law, researchers scaled the LLMs in both model and data size to gain higher performance like T5 (Raffel et al., 2023), GPT-3 (Brown et al., 2020), PaLM 2 (Anil et al., 2023). On the other hand, researchers also analyzed how scaling affects particular tasks like translation and prompt injection attack (Sun and Miceli-Barone, 2024). Within those discussions, the emergent abilities of LLM (Wei et al., 2022) are discovered. This means the LLM suddenly reaches high performance on previous low-performance tasks when reaching a certain model size. Recent studies have questioned whether emergent abilities are mirages caused by metrics misuse (Schaeffer et al., 2023) or just context-learning (Lu et al., 2023).

Regarding scaling in the field of memorization, Carlini et al. (2023) discussed the memorization across model size and found that the number of memorized texts grows with the model size and context size. Additionally, Biderman et al. (2023a) also discussed similar topics and found that a large portion of memorized text in a small-size model is also memorized by a larger model, showing that memorized texts may share common features.

2.2 Memorization

Prior to LLM, over-fitting is close to memorization (Tetko et al., 1995), which means a near-zero loss in the train set, suggesting the model memorized the input and its label (Zhang et al., 2017). However, memorization differs from overfitting because LLMs can perform well on the test set, whereas overfitting usually results in poor test set performance. Feldman (2020) analyzed the necessity of memorization in classification models. They demonstrated that in a long-tail distribution, where many categories have only a few samples, the neural model struggles to extract general features. Instead, the best strategy for the neural model is simply to memorize these samples and compare them with the data in the test set.

LLMs, unlike classification models, can directly generate their pre-train content, making the memorization observable. This can be used to form knowledge graphs (Petroni et al., 2019; AlKhamissi et al., 2022), but also leads to data contamination (Sainz et al., 2023) and privacy risk (Yao et al., 2024). Previous research has discussed memorization on a macro level. Tirumala et al. (2022) showed large models tend to memorize samples more easily in training. Carlini et al. (2023) dis-

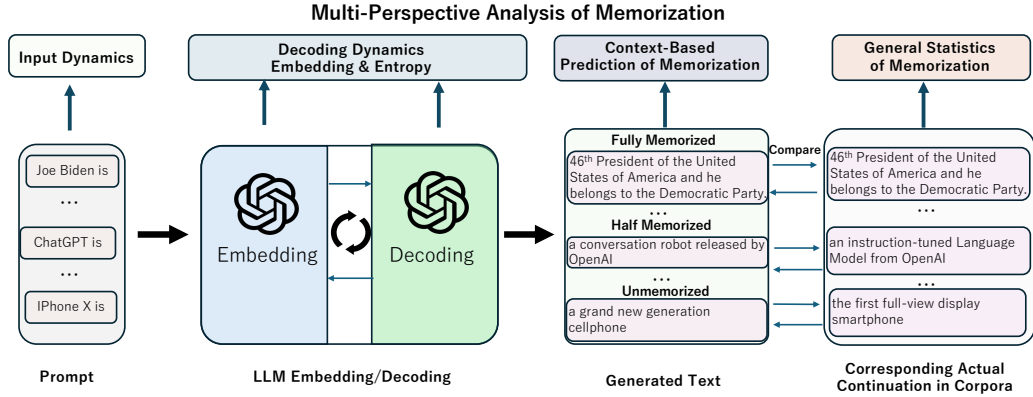


Figure 1: Memorization and research scope in this study. We prompt a partial context text into the model and calculate the memorization score of the generated continuations over the whole corpora. We show how inputs, model factors, and generation dynamics affect the memorization results.

cussed LLM memorization from factors like model size, continuation size, and context size, showing their inter-correlations. Biderman et al. (2023a) studied memorization in LLMs in their training phase and the overlap between different model sizes, finding commonly memorized sentences. Li et al. (2024) studied memorization by using hidden states in LLM in several specified datasets. Prashanth et al. (2024) set a threshold of a number of repetitions in the pre-train corpora to separate memorized sequences and predictable sequences. Speicher et al. (2024) created a sandbox setting to exploit the memorization of random strings.

Previous research discussed from a macro-level where they analyzed the memorization from an overall perspective, focusing mostly on fully memorized sequences with less consideration of the transition between sentences and how it is related to the model’s inner-working mechanism. This study focuses on micro yet under-explored questions, e.g., the dynamics while generating memorized sequences, why some sentences are memorized, the relation to unmemorized sequences, and its predictability.

3 Experiment Setting

3.1 Experiment Overview

As shown in Figure 1, we collect LLM outputs (tokens, probability distributions, and embeddings) and calculate a memorization score for every sentence to evaluate its extent of memorization and obtain general statistics. Then, we conduct an in-depth analysis of these collected outputs, examining the dynamics of both the prompted context input and the generated tokens. Finally, we discuss

the predictability of memorization.

3.2 Memorization Score

We prompt the LLM with context sequence tokens $C = \{c_1 \dots c_l\}$ from its pre-train corpora and use greedy decoding to generate the predicted continuation tokens $X = \{x_1 \dots x_n\}$. We also collect the actual continuations $Y = \{y_1 \dots y_n\}$ under this context. This process is iterated for the whole corpora. The memorization score is calculated as follows:

$$M(X, Y) = \frac{\sum_{i=1}^n \mathbf{I}(x_i = y_i)}{n} \quad (1)$$

n means the length of the continuation tokens. \mathbf{I} is the Indicator Function. If $M(X, Y) = 1$, the sequence is fully memorized under this context sequence, termed *K-extractable* (Carlini et al., 2021). A sequence Y is unmemorized if $M(X, Y) = 0$.²

3.3 Criteria for Memorization Prediction

We use Token Accuracy and Full Accuracy to evaluate performance in predicting memorization.

The prediction for a sequence of generated tokens is $\hat{X} = \{\hat{x}_1, \hat{x}_1 \dots \hat{x}_n\}$ where each prediction \hat{x}_i is a binary label indicating the token at this index is memorized or not. The gold label is denoted as $\hat{Y} = \{\hat{y}_1, \hat{y}_2 \dots \hat{y}_n\}$ where each \hat{y}_i is the golden label. The Token Accuracy is defined as:

$$T(\hat{X}, \hat{Y}) = \frac{\sum_{i=1}^n \mathbf{I}(\hat{x}_i = \hat{y}_i)}{n} \quad (2)$$

A prediction for a sequence is defined as fully correct if $T(\hat{X}, \hat{Y}) = 1$. We obtain Full Accuracy by

²Based on the different number of continuation tokens, the memorization score has different granularities. In some experiments, we classify sentences into several memorization levels based on their memorization scores.

dividing the number of fully correct sequences by the number of all sequences.

3.4 Model Setting

We use the Pythia (Biderman et al., 2023b) to analyze the memorization as it provides LLMs trained across various sizes with the same training order using open-sourced Pile (Gao et al., 2020) corpora that ensure experimental stability. We investigate the model size of [70m, 160m, 410m, 1b, 2.8b, 6.9b, 12b] where m and b stands for million and billion. We choose the LLM trained on the deduplicated Pile corpora to avoid the effect of duplicated sentences, as previous research reported the chance to be memorized grows exponentially with the number of duplicates (Kandpal et al., 2022).³

3.5 Corpora Setting

The open-sourced Pile (Gao et al., 2020) corpora are publicly available data.⁴ The data contains 146432000 rows with a chunk length of 2048, reaching a data size of around 800GBs. The experiment is iterated through the whole training data matrix, meaning that we did not conduct sampling over the rows. Instead, we iterated the whole Pile matrices. For example, if the context size is 32 and the continuation size is 96, we prompt the first 32 tokens at each row into the model. We use the Pythia model to generate the 96 tokens equal to the continuation size. Then, we compare the generated token IDs with the gold token IDs in the data to calculate the memorization score at each row. This process is repeated for the entire Pile matrix, distributed over different CUDA devices.

4 Experiment Results

4.1 Memorization Factors

This section discusses how memorization is connected to model size, context size, and continuation size. We collect the number of sentences with different memorization scores under different model sizes. Then, we divide those sentences with different memorization scores into ten sets with a memorization score difference of 0.1, as shown in Figure 2.

³For more details, please refer to <https://github.com/EleutherAI/pythia>

⁴The used data can be downloaded at <https://huggingface.co/datasets/EleutherAI/pile-deduped-pythia-preshuffled/tree/main>

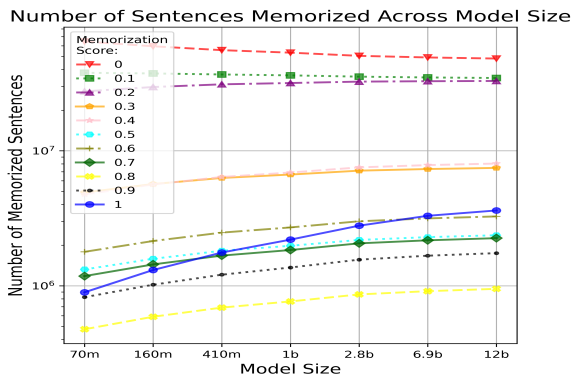


Figure 2: Number of sequences of different memorization scores across different model sizes.

Firstly, we can see that the memorized sentences increase with the model size and context size but decrease with the increase of continuation size, which aligns with previous research (Carlini et al., 2023). We illustrate a more in-depth analysis in the following sections.

4.1.1 The Factor of Model Size

In this experiment, we discuss how model size affects the number of memorized and unmemorized sentences. From Figure 2, we can obtain that:

(I) The number of sentences with low memorization scores (0-0.3) is significantly higher than those with high memorization scores, indicating that most of the pre-train data are not memorized despite the existence of memorization in LLMs.

(II) Among sentences with high memorization scores, the count of fully memorized sentences increases more rapidly, suggesting LLMs tend to memorize sentences fully rather than partially. Additionally, the number of sentences with low memorization scores (0, 0.1) decreases as the model size increases, indicating that unmemorized sentences gradually become memorized with larger models.

(III) The increase or decrease in the number of memorized or unmemorized sentences is not linear with respect to model size. There is a noticeable increase in numbers for fully memorized sentences from 70 million to 2.8 billion parameters, compared to 2.8 billion to 12 billion parameters. A similar decreasing trend for unmemorized sentences is observed. This suggests a capacity for memorization, implying LLMs cannot memorize the entire corpus even with sufficiently large model sizes.

4.1.2 Context and Continuation Size

We change the length of the context and fix the length of continuations or vice versa. Then, we ob-

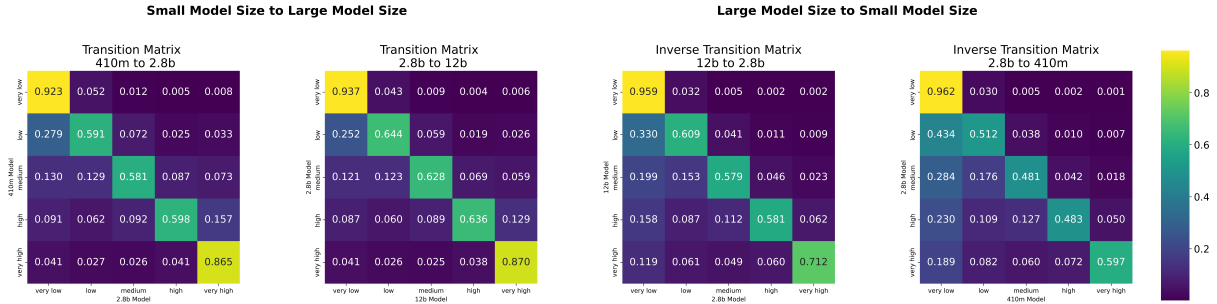


Figure 3: Transition matrix of sentences with different memorization scores. The two left figures are the transition matrix from small size model to large size, with the left one being the transition matrix from 410m to 2.8b and the right one being 2.8b to 12b. The two right figures are the inverse transition matrix from large size model to small size model, with the left one being 12b to 2.8b and the right one being 2.8b to 410m.

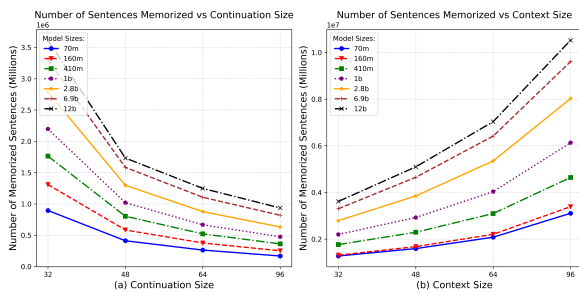


Figure 4: Number of memorized sentences in different continuation and context sizes across model size.

serve the change in the number of fully memorized sequences shown in Figure 4 (a) and (b). We can conclude:

(I) The decrease of memorized sentences with increasing continuation size is not linear. For instance, the continuation size increase from 64 to 96 results in a relatively minor decrease compared to the change from 32 to 48, indicating some sentences are firmly memorized.

(II) The reduction in memorized sentences with increased continuation size is more obvious in larger models. This demonstrates that although larger models memorize more sentences, most of their memorized sequences are less rooted compared to those of smaller models.

(III) The increase in memorized sentences with context size is also non-linear, with longer context leading to an almost exponential rise in the number of memorized sentences. This increase is more significant in larger models, indicating more sequences are potentially memorized in large models, which can be elicited by giving longer context.

4.2 Memorization Transition

This section discusses how sentences with different memorization scores transit across model sizes. Specifically, when trained with a larger model, we study the mutual transition between sentences with different memorization scores. We classify sentences with memorization scores with a range difference of 0.2 with labels of very low, low, medium, high, and very high. We plot the transition matrices in Figure 3, which shows the transition from the 410m size model to the 2.8b size and 2.8b size to the 12b size model and their inverse transition matrices. We can conclude:

(I) Most sentences remain in their previous state even when trained with a larger model, as indicated by the diagonal entries in the transition matrices. Additionally, the higher the memorization score, the less likely the sentence will transit to another state. For highly memorized sentences, over 90% remain memorized compared to those with low memorization scores. In the inversed transition matrix of the large-to-small size model, we see that most sequences also tend to stay in their original state with probability transferring to lower memorization score states, which fits our expectations.

(II) With increasing model size, sentences are more likely to stay at their original memorization level. For example, the transition probability at the diagonal is higher in the transition matrix from 2.8b to 12b compared to that of 410m to 2.8b. This suggests that memorized or unmemorized sequences become fixed as the model size increases. For the inversed matrix, the model is more likely to transfer to a lower memorization score state. In particular, transferring to the very low state is more probable than other low memorization score states. This shows when LLM starts to forget a sequence, it

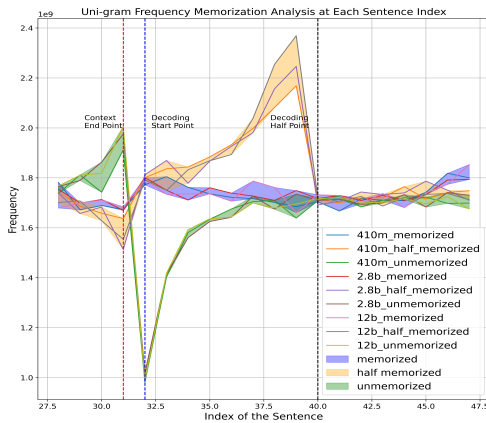


Figure 5: Uni-gram frequency at each index for memorized, unmemorized, and half-memorized sequences. The context length is 32. The continuation length is 16. We label the end of context, the beginning of encoding, and half of the encoding with red, blue, and black lines.

tends to forget that sequence completely.

(III) For high memorization score sentences, there is only a very small chance to transit to a low memorized state, implying most sentences memorized by small models are also inherited by the larger ones. This shows the memorized sentences are not memorized randomly but share certain common features with a little share of randomness. This is also observable in the inverse transition matrix, where there is little chance of transferring to a high memorization score state even when the model size decreases.

4.3 Frequency Dynamics of Input Sequences

In this section, we discuss the question of *whether there is any sign when the model starts to generate memorized or unmemorized sentences*. Especially what makes sentences memorized at different extents and why some sentences are memorized by large models but not small models.

4.3.1 Token Level Frequency Analysis

Firstly, we begin with the input level by analyzing the frequency of the n-grams in the pre-train corpora. We show the uni-gram frequency across steps in Figure 5⁵, and we can see:

(I) A clear *boundary effect* is observed around index 32, representing the first generated token. The frequency drops and then rises for memorized sentences (*positive boundary effect*), whereas it

rises and then drops for unmemorized sentences (*negative boundary effect*). The negative boundary effect is more significant than the positive one.

(II) For the half-memorized sentence, we see the frequency increases and drops (*negative boundary effect*) around the index of 39, which is half the length of generated tokens. This shows that for the half-memorized sentence, the previous half is mostly memorized, and the later half is mostly unmemorized, meaning that the memorized tokens are distributed in a near continuous way rather than scattered in the generated tokens.

(III) The positive boundary effect in memorized sentences suggests that memorization is driven by the higher frequency of initial tokens, implying that remembering the first few tokens makes the entire sentence easier to memorize. Conversely, the negative boundary effect in unmemorized sentences indicates that the low frequency of initial tokens makes the following sequences easier to forget.

4.3.2 Sequence Level Frequency Analysis

Given the existence of the token level boundary effect, we extend the discussion to the sequence level. As shown in Table 1, we calculate the average n-gram frequency of context and continuation and boundary frequency difference. We can obtain:

(I) The frequency of uni-grams is significantly higher than that of bi-grams, approximately 3.5 times higher on average in both context and continuation. The boundary effect is consistent in the bi-gram setting, though the positive boundary effect is less obvious due to the frequency drop when computed with bi-grams. Despite this, the actual frequency gap remains substantial (million level) when considering the unit is billion.

(II) In memorized sentences, the frequency is lower in the context and higher in the continuation of the memorized data (M column), whereas unmemorized sentences exhibit the opposite pattern. This suggests the boundary effect also exists at the sequence level, though less obvious.

(III) For half-memorized sentences, the frequency of continuation tokens is higher than the context average frequency. This is due to the frequency increase before reaching the first generated unmemorized tokens, as indicated in Figure 5.

(IV) In the Boundary Frequency Difference column, both positive and negative boundary effects decrease with increased model size. However, the decrease in the positive boundary effect makes it less significant, while the decrease in the negative

⁵N-gram statistics overall steps are in the A.5.

| Size | Average Context Frequency | | | | | | Average Continuation Frequency | | | | | | Boundary Frequency Difference | | | | | |
|------|---------------------------|-------|-------|---------|-------|-------|--------------------------------|-------|-------|---------|-------|-------|-------------------------------|-------|--------|---------|-------|--------|
| | Uni-gram | | | Bi-gram | | | Uni-gram | | | Bi-gram | | | Uni-gram | | | Bi-gram | | |
| | M | H | U | M | H | U | M | H | U | M | H | U | M | H | U | M | H | U |
| 160m | 1.708 | 1.713 | 1.744 | 0.551 | 0.534 | 0.691 | 1.739 | 1.837 | 1.628 | 0.535 | 0.659 | 0.567 | 0.114 | 0.330 | -0.939 | 0.033 | 0.101 | -0.663 |
| 1b | 1.713 | 1.711 | 1.752 | 0.558 | 0.552 | 0.697 | 1.736 | 1.832 | 1.631 | 0.509 | 0.682 | 0.564 | 0.103 | 0.270 | -0.981 | 0.028 | 0.090 | -0.696 |
| 6.8b | 1.721 | 1.710 | 1.759 | 0.570 | 0.565 | 0.701 | 1.736 | 1.829 | 1.638 | 0.496 | 0.699 | 0.564 | 0.090 | 0.140 | -0.963 | 0.027 | 0.085 | -0.726 |
| 12b | 1.721 | 1.720 | 1.760 | 0.572 | 0.569 | 0.702 | 1.736 | 1.846 | 1.626 | 0.493 | 0.704 | 0.563 | 0.039 | 0.237 | -1.016 | 0.026 | 0.083 | -0.732 |

Table 1: Uni-gram and bi-gram statistics. The frequency unit is billion. Boundary Freq Difference means we use the first generated token’s frequency to subtract the last token’s frequency in the context (i.e., the boundary effect). M, H, and U mean memorized, half-memorized, and unmemorized, respectively.

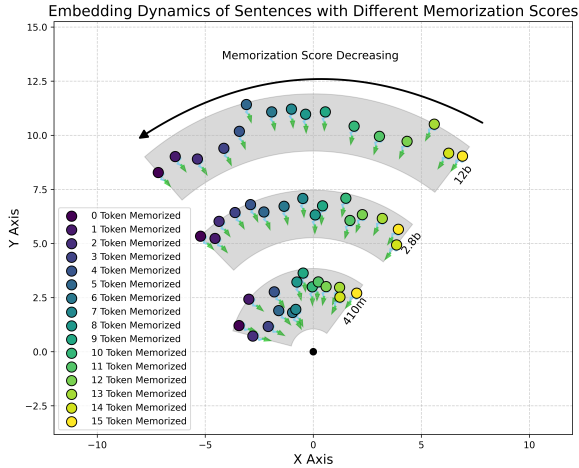


Figure 6: Embedding dynamics across 410m, 2.8b, and 12b model sizes. x Token Memorized means x generated tokens are the same as the true continuation. The arrow for each cluster means the visualized moving direction of generated tokens. The gray area means the span of those clusters.

boundary effect makes it more significant. This implies the significance of the positive boundary effect correlates with the ease of memorizing a sentence. In contrast, the significance of the negative boundary effect correlates with the difficulty of not memorizing a sentence.

4.4 Decoding Dynamics

In this section, we analyze the dynamics of the output, e.g., the movement of generated embeddings at each step and corresponding entropy changes for sequences of varying memorization scores.

4.4.1 Embedding Dynamics

We analyze embeddings of generated tokens for sentences with varying memorization scores. We collect the hidden state of the last layer for each generated token for sentences with different memorization scores and compute the pair-wise Eu-

clidean distance and cosine similarity. We draw the following figure to represent the embedding dynamics shown in those pair-wise Euclidean distance and cosine similarity.⁶ We can obtain:

(I) The mutual angle remains stable across different encoding steps between sentences of different memorization scores, which also suggests a stable mutual cosine similarity. Meanwhile, since they have the trend of moving toward the center, the Euclidean distance also decreases with the generation of tokens.

(II) Sentences with high memorization scores are close in the embedding space. This suggests their generated sequences are both semantically and lexically similar to the actual continuation. This indicates the existence of paraphrase memorization since those high memorization score sequences are probably only different in a few tokens while sharing the same meaning.

(III) Larger models exhibit larger mutual Euclidean distances and angles. The increase in angle leads to a decrease in cosine similarity. The reason can be attributed to the expansion of hidden sizes (e.g., 512 for 70m, 2048 for 1b model). The expansion of hidden size increases the expressivity of the embedding and enlarges the mutual distances, making different sentences more differentiable. This also helps to explain the performance gap between different model sizes: larger models distribute different sentences more distinctly with fewer embedding overlaps, while smaller models mix embeddings more, leading to ambiguity and degraded performance.

4.4.2 Generation Dynamics and Entropy

This section discusses the generation dynamics when the LLM generates sentences with different

⁶Detailed numbers of Cosine Similarity and Euclidean distance between sentences with different memorization scores at different decoding steps are at Appendix A.3

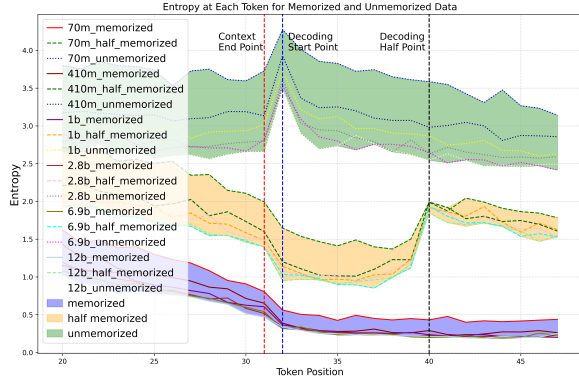


Figure 7: Averaged entropy dynamics at each index for memorized, unmemorized, and half-memorized sequences. The prompted context length is 32, and the continuation length is 16. We label the end of context, the beginning of encoding, and the half of the encoding with red, blue, and black lines.

memorization scores. Given the boundary effect at the input frequency, we ask *whether the model exhibits similar behaviors when generating memorized or unmemorized sequences*. As shown in Figure 7, we collect the probability distributions for each generated token of memorized, unmemorized, and half-memorized sentences, with 10,000 samples for each, and calculate the average entropy at each token. We can conclude:

(I) First, we can see that the entropy differs based on memorization scores and model size. Unmemorized sentences have a higher average entropy at each token than memorized sentences. This shows LLM is more confident when generating memorized sequences.

(II) Additionally, the entropy drops when generating memorized sequences (fully memorized and former half of half-memorized sequences) and increases when generating unmemorized ones (unmemorized and later half of half-memorized sequences). This shows an *inverse boundary effect* in entropy compared to the one in frequency. However, the entropy drop for memorized sentences is not samely significant as the half-memorized sentence, as LLM is more confident when generating memorized sequences, which is naturally low entropy, leaving less space for the entropy drop.

(III) We can also see the entropy decreases with the increase in model size. This suggests that larger models are more confident about generating tokens than small ones. Additionally, the entropy of context for memorized tokens is also lower, showing that the entropy of context also relates to

whether its continuations are memorized. Further, the significance of the inverse boundary effect in entropy decreases with the model size for memorized sequences while remaining unchanged for unmemorized sequences. This differs from the boundary effect in frequency, whose significance decreases with increased model size for memorized sequences while increasing for unmemorized ones.

4.5 Prediction of Memorization

Since similar context can trigger memorized texts (Stoehr et al., 2024), and the generated token may also be a paraphrase of the actual continuation (Ippolito et al., 2023), common methods that check memorization by searching the generated tokens in the corpora are challenged. Thus, it would be beneficial if it were possible to predict the memorization by embedding. We sample sentences with different memorization scores evenly from the whole corpora. A Transformer (Vaswani et al., 2023) is trained to predict a binary label at each continuation token, predicting whether this token is memorized by receiving all embeddings generated so far and related statistics, e.g., the entropy and frequency.⁷

4.5.1 Results on Prediction of Memorization

We discuss the predictability of predicting memorization based on results presented in Table 2. We can obtain:

(I) First, regarding Token Accuracy, we can see with a naive Transformer model, the token-level accuracy can reach 80% accuracy or even higher, showing that the prediction of memorization at the token level is easy. The Full Accuracy is low as this requires a correct prediction for every token.

(II) In either continuation length, both token accuracy and full accuracy increase with model size. This shows the prediction of memorization is easier for large models because the greater embedding distances make classification easier.

(III) Token Accuracy increases with the continuation size in either model size, likely due to increased training data. For instance, continuation length 64 contains four times larger tokens than that of continuation length 16. However, Full Accuracy decreases as the continuation length increases as more tokens need to be predicted correctly.

4.5.2 Analysis of Full Accuracy

In this experiment, we analyze how fully correct predictions are distributed across memorization

⁷Experiment details and settings are in Appendix A.1.3.

| Len | 70m | | 410m | | 1b | | 2.8b | | 6.9b | | 12b | | Dist | |
|-----|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|------------|-------------|-------------|------|------|
| | Token | Full | Token | Full | Token | Full | Token | Full | Token | Full | Token | Full | M | U |
| 16 | 78.2 | <u>10.2</u> | 78.6 | <u>10.4</u> | 78.8 | <u>10.6</u> | 80.1 | <u>10.7</u> | 77.4 | <u>8.3</u> | 80.3 | <u>10.9</u> | 53.1 | 46.9 |
| 32 | 78.6 | 5.9 | 79.6 | 6.0 | 79.7 | 6.1 | 80.1 | 6.3 | 80.5 | 6.4 | 80.8 | 6.4 | 51.6 | 48.4 |
| 48 | 79.6 | 5.2 | 80.3 | 5.4 | 80.4 | 5.6 | 80.4 | 5.5 | 80.8 | 5.8 | 81.0 | 6.0 | 51.0 | 49.0 |
| 64 | <u>80.1</u> | 4.7 | <u>80.8</u> | 4.8 | <u>81.2</u> | 5.2 | <u>81.5</u> | 5.5 | <u>81.8</u> | 5.8 | <u>82.1</u> | 6.0 | 50.7 | 49.3 |

Table 2: Performance on memorization prediction at context length 32. Len means the length of continuation tokens for prediction. Token and Full mean Token and Full Accuracy. The various model sizes in the column show which size of LLM’s embeddings are used to train the Transformer. The M and U in the Dist column mean the distribution of memorized and unmemorized tokens used to train and evaluate. For each continuation length, the best results across model sizes are in bold. The best results in a model size across continuation lengths are underlined.

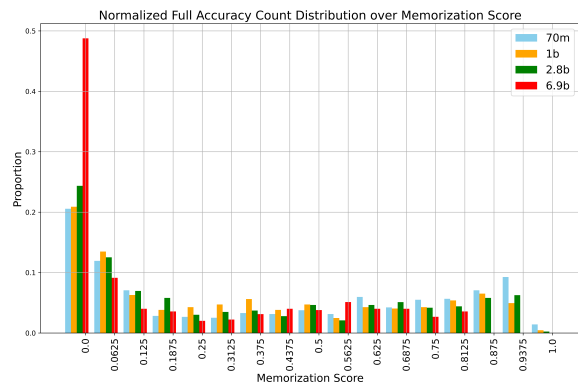


Figure 8: Full accurate predictions distribution at different memorization scores across model size

scores to discuss whether the difficulty in predicting memorization changes with the memorization score as shown in Figure 8, and we can obtain:

(I) The Transformer model trained with embeddings from any LLM size is better at predicting sentences with low memorization scores, even when the label distribution is close. The low portion of sentences with high memorization scores indicates they are harder to predict accurately.

(II) As model size increases, the proportion of low memorization scores rises and decreases for high memorization score sentences, which even reaches zero for the 6.9b model. This suggests predicting unmemorized sentences is easier in large models compared to memorized ones.

(III) A possible explanation for the above behaviors can be made. Previous experiments show that the boundary effect of unmemorized sequences is more significant than that of memorized sequences in both token frequency and entropy. Additionally, the significance of the boundary effect decreases for memorized sequences while increasing for unmemorized sequences with increased model size.

With the decrease of significance in such features for memorized sequences, they become hard to predict. The unmemorized sequences become easy to predict as the significance increases.

4.6 Conclusion

In this study, we comprehensively examined LLM memorization from various perspectives. At the statistical level, we extended previous research to include sentences with lower memorization scores and conducted experiments showing memorization transitions across model sizes. Analyzing input dynamics through frequency analysis, we identified positive and negative boundary effects when generating memorized and unmemorized tokens, indicating their relation to the ease of memorization. In the output dynamics, at the embedding level, we found clusters of sentences with different memorization scores in the embedding space, and the close distance of sentences with high memorization scores indicates the existence of paraphrase memorization. At the entropy in the output dynamics, we observed an inverse boundary effect and analyzed its change with model size. Finally, we trained a Transformer model to predict memorization, showing that token-level prediction is easy while sentence-level is challenging. Through analysis of fully correct predicted samples, we found unmemorized tokens are easier to predict than memorized tokens, which can be explained by the significance of the boundary effect.

5 Acknowledgments

We sincerely thank every reviewer for their valuable suggestions. This work is funded by the Institute of AI and Beyond of the University of Tokyo.

6 Limitations

This research has analyzed various factors regarding the memorization behavior of LLMs. We acknowledge that there are still limitations to this research. Due to the lack of LLMs whose models and data are both being released, it is hard to compare the memorization across different LLMs since even if they were released, the pre-train data differs based on different LLMs. Future research can focus on how to comprehensively measure the memorization of various LLMs, either close-sourced like GPT-4 or open-sourced like Pythia. Additionally, Pythia only provides LLMs up to 12b, which is still considerably small when compared to SotA open-sourced LLMs like LLaMa2, whose largest model is 70b. The emergent abilities reaching 70b size may also affect memorization. Though we trained a Transformer model to predict the memorization of LLMs, it is more analysis-oriented, proving the possibility of predicting memorization. Thus, the performance in the prediction experiment is not the main focus.

Additionally, in this research, the discussion of memorization is under the context when the LLM generates tokens that are the same as the actual continuations in the corpora, defined as *K-extractable* in this study. It is possible that paraphrase memorization exists. However, it is difficult to identify such behavior on a large scale, especially across whole corpora. This is because the identification of paraphrase memorization can not be simply decided by either token overlap or embedding similarity since neither of them can truly compare whether a sentence is a paraphrase. This means we have to identify by more complex methods, e.g., human annotators or a trained classification model. However, a classification model cannot be 100% accurate when applied to the corpora level, leading to false positives and false negatives that influence the analysis results. However, paraphrase memorization can be left to future research.

7 Ethical Considerations

In this study, we prompted the Pythia model with context tokens from its pre-train corpora. The data released by Pythia are all tokenized and turned into token IDs, so the original information is not visible. We have provided a case study of the prediction of memorization in the Appendix A.2 by turning the token IDs into texts. However, we did not observe any personal information or offensive language dur-

ing this process. Additionally, we obeyed related open-source licenses of both Pythia and Pile corpora. Thus, this study is not concerned with ethical issues.

References

- Samira Abnar, Mostafa Dehghani, Behnam Neyshabur, and Hanie Sedghi. 2021. [Exploring the limits of large scale pre-training](#).
- Guillaume Alain and Yoshua Bengio. 2018. [Understanding intermediate layers using linear classifier probes](#).
- Badr AlKhamissi, Millicent Li, Asli Celikyilmaz, Mona Diab, and Marjan Ghazvininejad. 2022. [A review on language models as knowledge bases](#).
- Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clément Crepey, Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, Andrea Hu, Jeffrey Hui, Jeremy Hurwitz, Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun, Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li, Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu, Frederick Liu, Marcello Maggioni, Aroma Mahendru, Joshua Maynez, Vedant Misra, Maysam Moussalem, Zachary Nado, John Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie Pellat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan Richter, Parker Riley, Alex Castro Ros, Aurko Roy, Brennan Saeta, Rajkumar Samuel, Renee Shelby, Ambrose Slone, Daniel Smilkov, David R. So, Daniel Sohn, Simon Tokumine, Dasha Valter, Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang, Pidong Wang, Zirui Wang, Tao Wang, John Wieting, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav Petrov, and Yonghui Wu. 2023. [Palm 2 technical report](#).
- Stella Biderman, USVSN Sai Prashanth, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony,

- Shivanshu Purohit, and Edward Raff. 2023a. [Emergent and predictable memorization in large language models](#).
- Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. 2023b. [Pythia: A suite for analyzing large language models across training and scaling](#).
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. 2023. [Quantifying memorization across neural language models](#).
- Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. 2021. [Extracting training data from large language models](#). In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650. USENIX Association.
- Vitaly Feldman. 2020. [Does learning require memorization? a short tale about a long tail](#). In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, STOC 2020, page 954–959, New York, NY, USA. Association for Computing Machinery.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. [The Pile: An 800gb dataset of diverse text for language modeling](#). *arXiv preprint arXiv:2101.00027*.
- Valentin Hartmann, Anshuman Suri, Vincent Bindschaedler, David Evans, Shruti Tople, and Robert West. 2023. [Sok: Memorization in general-purpose large language models](#).
- Daphne Ippolito, Florian Tramèr, Milad Nasr, Chiyuan Zhang, Matthew Jagielski, Katherine Lee, Christopher Choquette Choo, and Nicholas Carlini. 2023. [Preventing generation of verbatim memorization in language models gives a false sense of privacy](#). In *Proceedings of the 16th International Natural Language Generation Conference*, pages 28–53, Prague, Czechia. Association for Computational Linguistics.
- Nikhil Kandpal, Eric Wallace, and Colin Raffel. 2022. [Deduplicating training data mitigates privacy risks in language models](#).
- 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](#).
- Bo Li, Qinghua Zhao, and Lijie Wen. 2024. [Rome: Memorization insights from text, logits and representation](#).
- Sheng Lu, Irina Bigoulaeva, Rachneet Sachdeva, Harish Tayyar Madabushi, and Iryna Gurevych. 2023. [Are emergent abilities in large language models just in-context learning?](#)
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue,

- Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeef Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pocrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolaus Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. [Gpt-4 technical report](#).
- Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H. Miller, and Sebastian Riedel. 2019. [Language models as knowledge bases?](#)
- USVSN Sai Prashanth, Alvin Deng, Kyle O’Brien, Jyothir S V au2, Mohammad Aflah Khan, Jaydeep Borkar, Christopher A. Choquette-Choo, Jacob Ray Fuehne, Stella Biderman, Tracy Ke, Katherine Lee, and Naomi Saphra. 2024. [Recite, reconstruct, recollect: Memorization in lms as a multifaceted phenomenon](#).
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2023. [Exploring the limits of transfer learning with a unified text-to-text transformer](#).
- Oscar Sainz, Jon Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. 2023. [NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10776–10787, Singapore. Association for Computational Linguistics.
- Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. 2023. [Are emergent abilities of large language models a mirage?](#) In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Till Speicher, Aflah Mohammad Khan, Qinyuan Wu, Vedant Nanda, Soumi Das, Bishwamitra Ghosh, Krishna P. Gummadi, and Evimaria Terzi. 2024. [Understanding the mechanics and dynamics of memorization in large language models: A case study with random strings](#).
- Niklas Stoehr, Mitchell Gordon, Chiyuan Zhang, and Owen Lewis. 2024. [Localizing paragraph memorization in language models](#).
- Zhifan Sun and Antonio Valerio Miceli-Barone. 2024. [Scaling behavior of machine translation with large language models under prompt injection attacks](#). In *Proceedings of the First edition of the Workshop on the Scaling Behavior of Large Language Models (SCALE-LLM 2024)*, pages 9–23, St. Julian’s, Malta. Association for Computational Linguistics.
- Igor V. Tetko, David J. Livingstone, and Alexander I. Luik. 1995. [Neural network studies. 1. comparison of overfitting and overtraining](#). *Journal of Chemical Information and Modeling*, 35(5):826–833.
- Kushal Tirumala, Aram H. Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. 2022. [Memorization without overfitting: Analyzing the training dynamics of large language models](#). In *Advances in Neural Information Processing Systems*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. [Attention is all you need](#).
- Pablo Villalobos. 2023. [Scaling laws literature review](#). Accessed: 2024-04-29.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. [Emergent abilities of large language models](#).
- Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. 2024. [A survey on large language model \(llm\) security and privacy: The good, the bad, and the ugly](#). *High-Confidence Computing*, 4(2):100211.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. 2017. [Understanding deep learning requires rethinking generalization](#).

A Appendix

A.1 Experiment Setting

A.1.1 Pythia LLM Generation

The experiment uses 64 A100 40Gbs GPUs when using LLMs to generate tokens given the previous context, which utilizes PyTorch’s parallel running packages. We run the model with half-precision, which increases both the speed and saves memory. This follows the previous Pythia implementation when generating tokens given context.

The running time depends on the model size and the generated token length. When a 70m model with 32 context tokens and 16 tokens is required to be generated, it can be run with one A100 GPU within several hours. However, if such an experiment is in a single A100 GPU, it would be estimated to take two weeks to finish the generation. Therefore, with 64 GPUs, the running of the 12b model in such a situation can be shortened to around one day. However, the generation time also largely increases with the length of generated tokens, which grows linearly with the number of tokens to be generated. Additionally, since we use greedy decoding and do not consider other possible decoding options but the token with the highest probability, it is also possible that the running time may be different if using a more complicated decoding strategy. However, using a more complex decoding strategy does not affect the results. If a sequence of tokens is memorized, the memorized token will always be the most probable token when generating them.

A.1.2 N-Gram Statistics Calculation

For the analysis of 4.3, since calculating n-gram statistics for pre-train corpora demands lots of memory, we used another experimental environment with 128G RAM. It takes several days to calculate each gram’s statistics. Also, storing the result of n-gram statistics takes over 1TB of storage.

A.1.3 Memorization Prediction

Regarding the prediction of memorization, we sample 2,000 sentences at the sentence set of each memorization score. For example, in the situation where the continuation length is 16, we will sample 2,000 sentences from memorization scores, the range of which is from 0 to 1 with 0.0625 as the unit. The Transformer model receives embedding of the last layer at each generated step with the corresponding entropy and uni-gram frequency and then outputs an embedding at each step. This output embedding

from the Transformer model will be used to pass through a linear layer, and the Softmax function calculates the probability that the token at this index is a memorized token or not. The training is conducted using a 4-layer Transformer model with a dropout probability of 0.5, and the learning rate is $1e-4$. Additionally, the only changing parameter with the increase in model size is the hidden size, as the larger model has a larger hidden size. The training and evaluation are conducted across five random seeds, and we report the average performance. The train, valid, and test split ratio is 0.6, 0.2 and 0.2. Additionally, we make sure that the distribution of sentences of different memorization scores is even in the dataset; thus, the model does not make biased predictions.

A.2 Case Study

We also provide a case study of the model’s prediction on the test set regarding the prediction of memorization in Figure 9. From this figure, we can see that:

1. Confirming previous experiments, the memorized token is mostly continuous, showing the memorization happens in chunks of sequences rather than individual sequences.
2. In the first example, the model outputs a fully correct prediction that aligns with the actual label, showing the possibility of predicting memorization by utilizing embedding information.
3. In the second example, we see that the model’s prediction does not align with the actual continuation. The model predicts an unmemorization label for the memorized label.
4. In the third example, we can see this unmemorized sequence. The model fully predicted those labels. We also noticed that the probability for the corresponding token is very high, showing that the model is very confident about the prediction.

A.3 Supplementary Data for Embedding Dynamics

We have presented the PCA visualized embedding dynamics in Figure 6. In this section, we provide the actual numbers of both Cosine Similarity and Euclidean distance to further illustrate this point.

Similarly to Figure 6, we also provide detailed cosine similarity and distance results from Figure 10 to Figure 18. From those figures, we can see that the cosine similarity fluctuates but remains relatively stable for sentences with different memorization scores. Additionally, sentences that are

| Context | Label | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | |
|--|-----------|------------|------|-----------|------------|------|------------|---------|------|-------------|-----------|------|-------|--------|----------|------|------|---|
| How many minutes are there between 2:02 AM and 2:01 PM? 7 | Text | 19 | how | many | minutes | are | there | between | 5 | : | 46 | PM | and | 3 | : | 32 | | |
| | Pred Prob | 0.90 | 0.99 | 0.94 | 0.93 | 0.83 | 0.95 | 0.89 | 0.95 | 0.95 | 0.94 | 0.98 | 0.54 | 0.83 | 0.93 | 0.91 | 0.97 | |
| | Gold | 0.69 | M | M | M | M | M | M | M | M | U | M | U | U | M | U | M | U |
| | Pred | 0.69 | M | M | M | M | M | M | M | M | U | M | U | U | M | U | M | U |
| The objective of all equity funds is to seek out profit opportunities. Types of equity | Text | funds | are | different | types | of | equity | funds | , | categorized | according | to | risk | levels | | | | |
| | Pred Prob | 0.97 | 0.75 | 0.60 | 0.83 | 0.89 | 0.93 | 0.98 | 0.97 | 0.99 | 0.99 | 0.98 | 0.97 | 0.98 | 0.94 | 0.99 | 0.99 | |
| | Gold | 0.65 | M | M | M | M | M | U | M | M | M | M | U | U | U | U | U | U |
| | Pred | 0.5 | U | M | U | U | U | U | U | U | U | U | U | U | U | U | U | U |
| Creating men whose expectations of what they should look like are unattainable. In recent years, | Text | researcher | at | the | University | of | California | School | of | Medicine | injected | the | brain | of | diabetic | rats | with | |
| | Pred Prob | 0.93 | 0.95 | 0.80 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | |
| | Gold | 0 | U | U | U | U | U | U | U | U | U | U | U | U | U | U | U | U |
| | Pred | 0 | U | U | U | U | U | U | U | U | U | U | U | U | U | U | U | U |

Figure 9: Prediction Examples. Pred Prob means the output predicted probability of the corresponding label in the pred row for each example. Gold means the true label. Pred means the predicted label, and Pred Prob means the probability of the corresponding prediction. *M* means the label for a token at this index is a memorized token. *U* means the label for a token at this index is an unmemorized token.

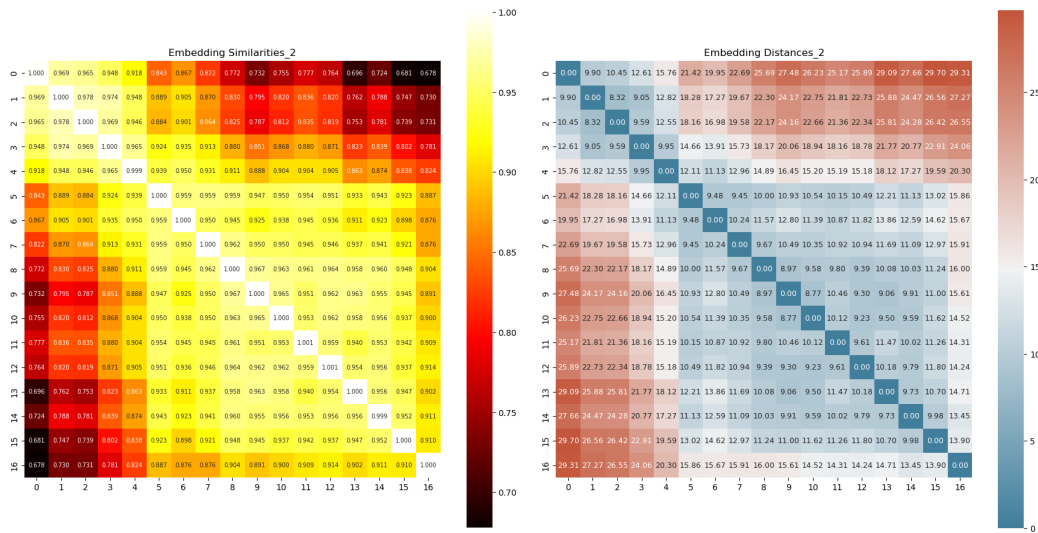


Figure 10: 410m Model, Step 2

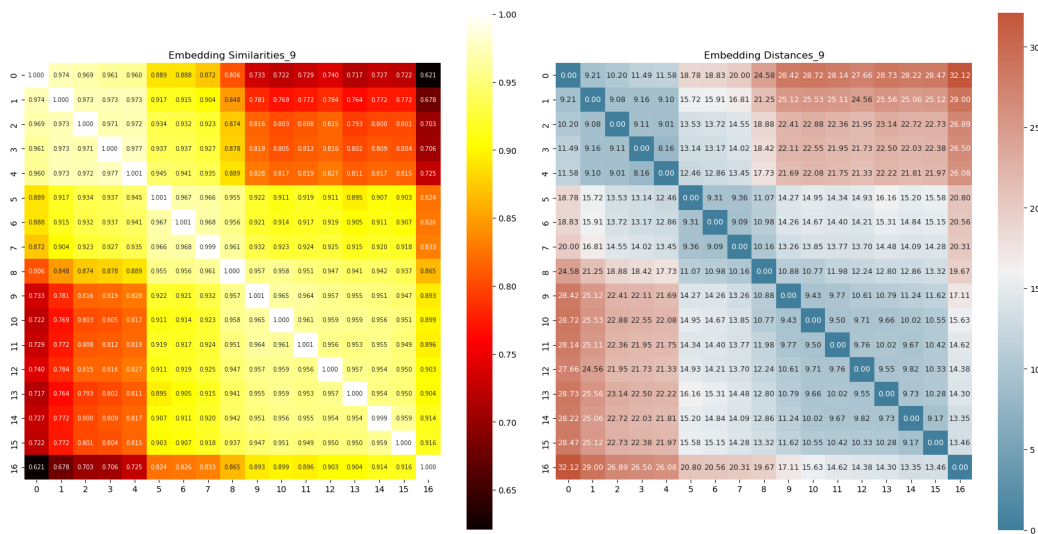


Figure 11: 410m Model, Step 9

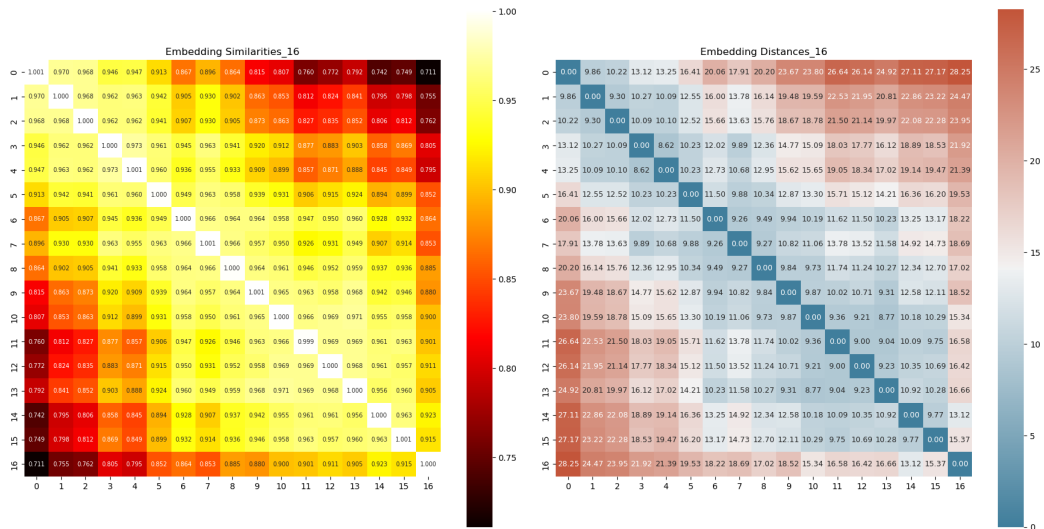


Figure 12: 410m Model, Step 16

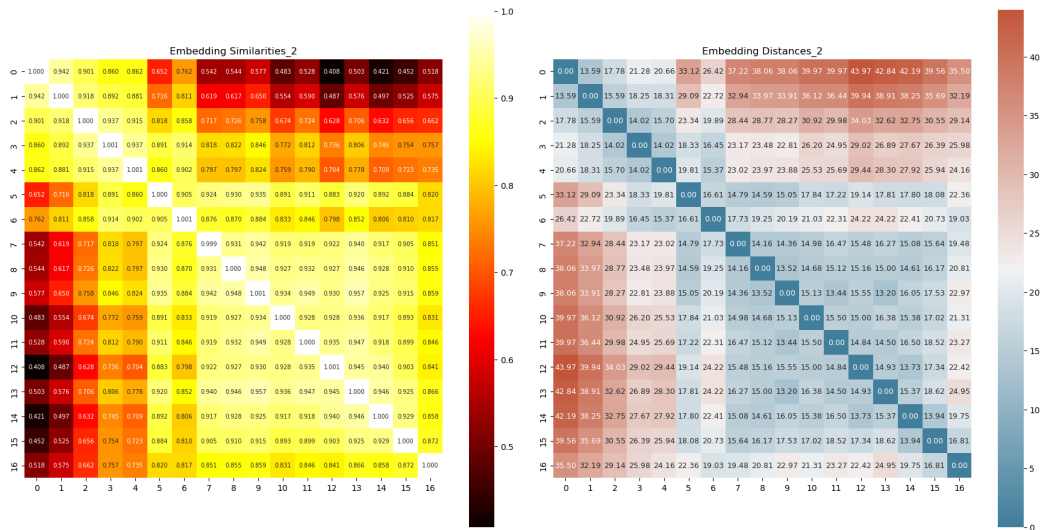


Figure 13: 2.8b Model, Step 2

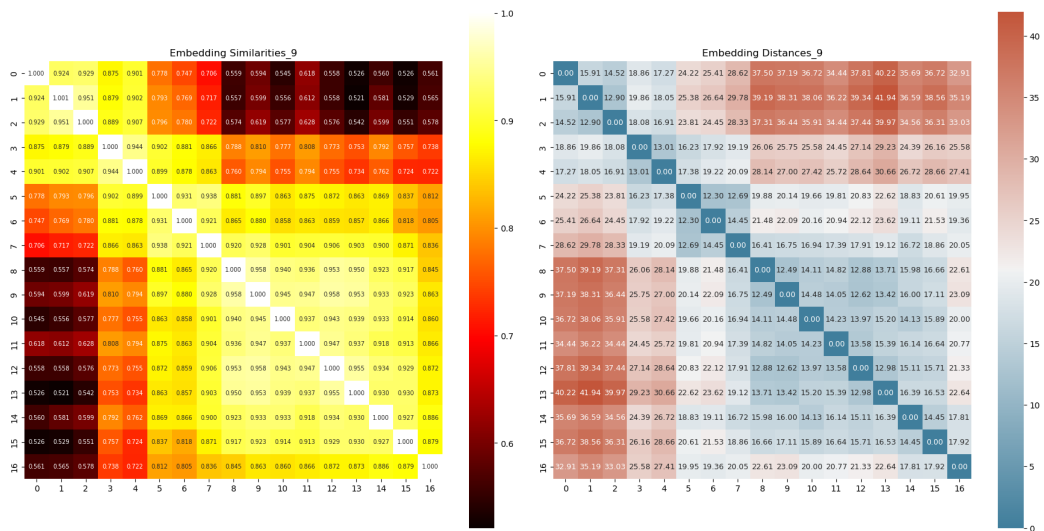


Figure 14: 2.8b Model, Step 9

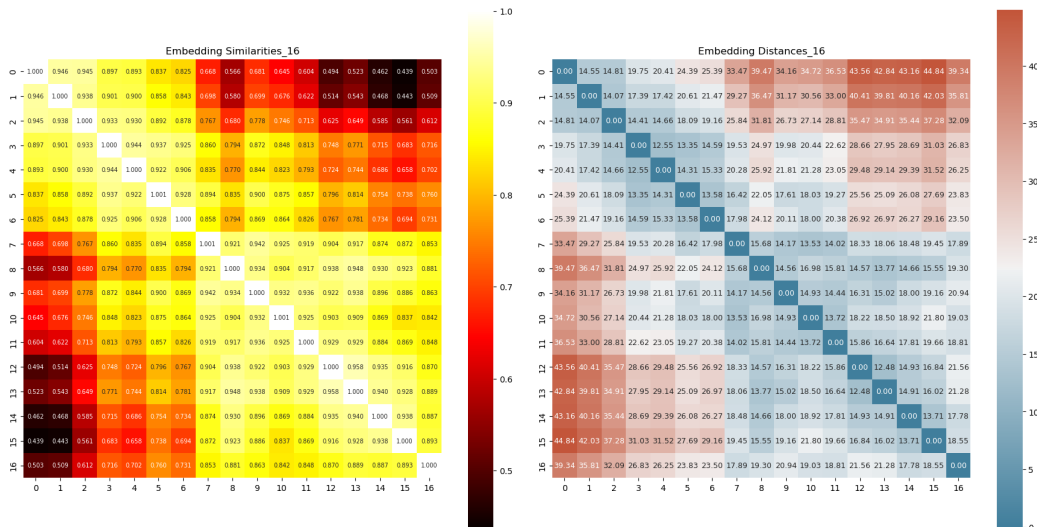
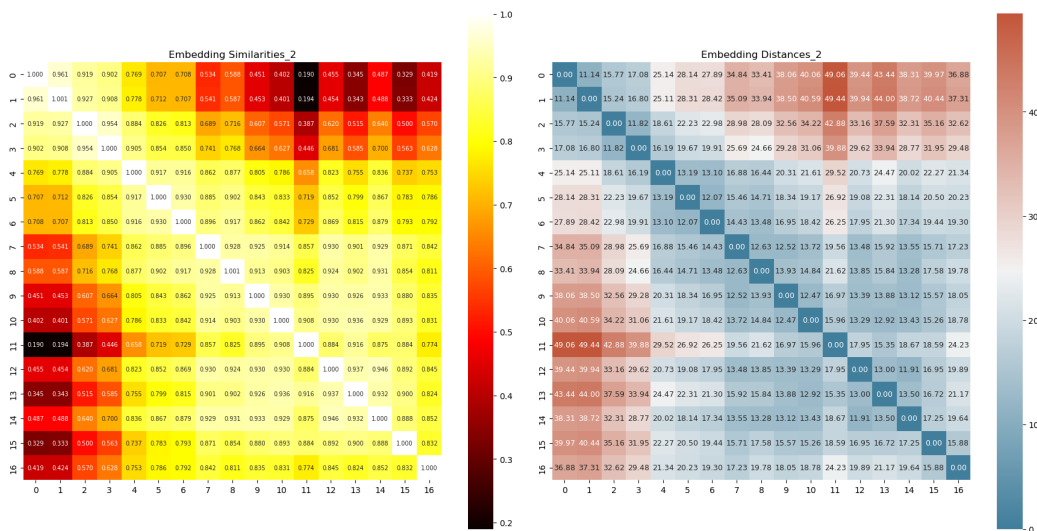


Figure 15: 2.8b Model, Step 16



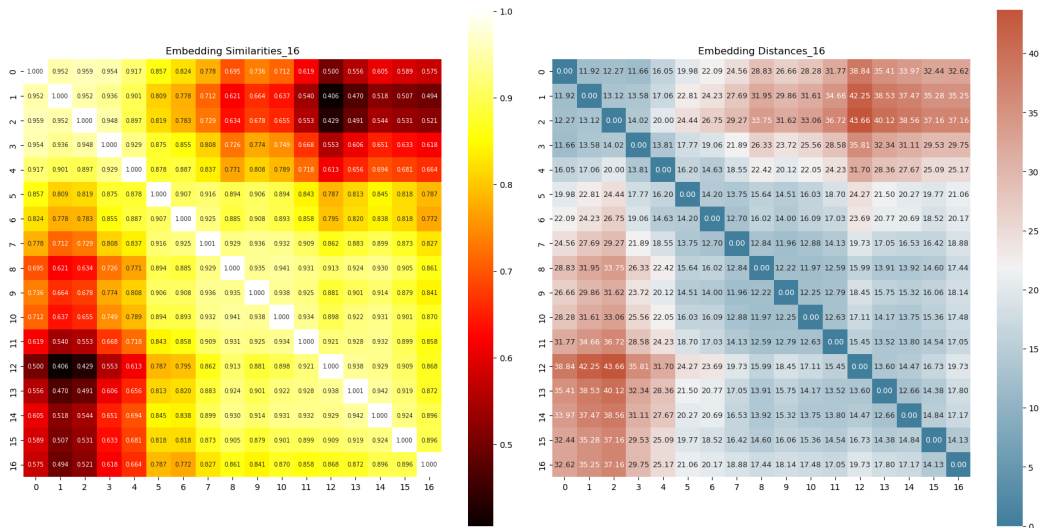


Figure 18: 12b Model, Step 16

not exactly the same but close regarding the memorization score are very close in the embedding space. For example, fully memorized sentences are also close to sentences with high memorization scores and embedding similarity of over 0.9. This shows the possibility that the model is generating paraphrased memorized sequences. However, with the increase in the model size, the cosine similarity decreases. For the Euclidean distance, we can see that the embedding distances have a decreasing trend with the increasing decoding steps, while the mutual Euclidean increases with the model size.

A.4 Does LLM prefer to memorize specific parts within the training data?

In this section, we discuss the question of whether LLM prefers to memorize a specific part within the training data. We split the corpora into 50 parts based on their index and examine how many memorized sentences are in those parts.

From the result shown in Figure 19, though the number of memorized sentences is not completely evenly distributed, we can see that there is no significant part that the number of memorized sentences is clearly more than others. This shows the training order does not affect memorization.

A.5 Detail N-Gram Statistics

We also provided detailed uni-gram and bi-gram statistics for the memorized, un-memorized, half-memorized, and quarter-memorized contents. The average frequency is calculated at each step, and we compute the average frequency of the context, the average frequency of the continuation, and the

average frequency of the whole sequence. The results span from the 70m size model to the 12b size model shown in Table 3 to Table 6.

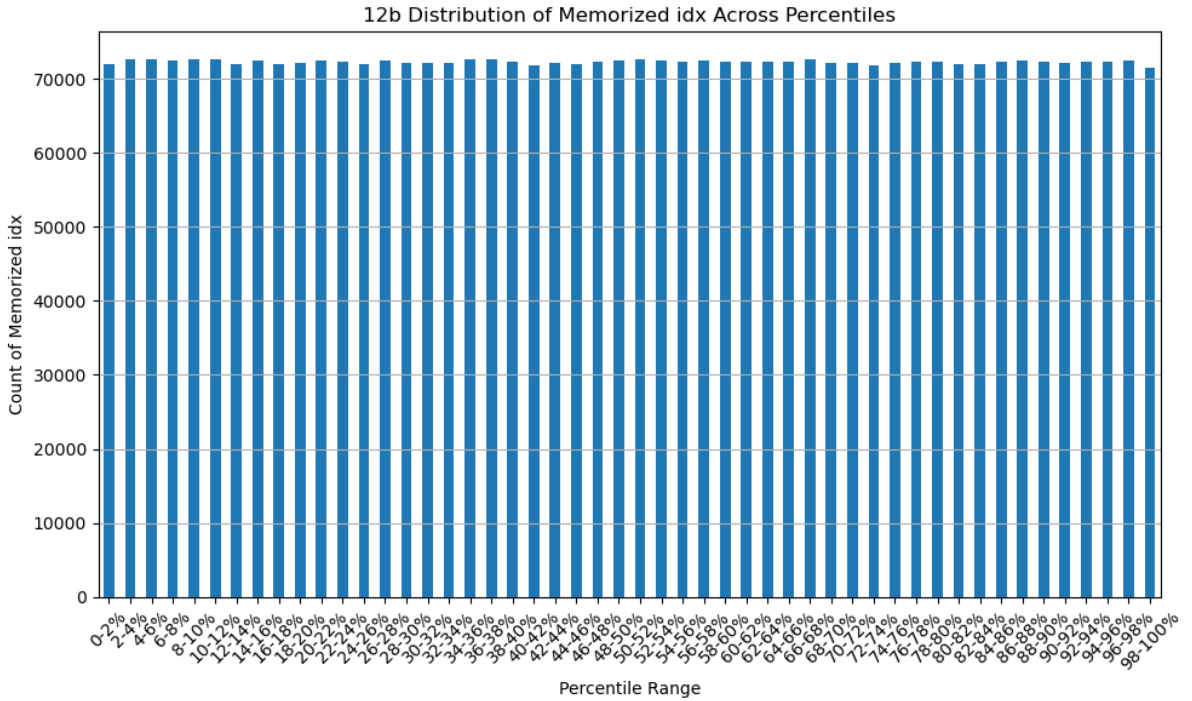


Figure 19: Index Distribution of Memorized Sequences in 12b model.

| target size | memorized 70m | memorized 160m | memorized 410m | memorized 1b | memorized 2.8b | memorized 6.9b | memorized 12.9b | forgotten 70m | forgotten 160m | forgotten 410m | forgotten 1b | forgotten 2.8b | forgotten 6.9b | forgotten 12.9b |
|-------------|---------------|----------------|----------------|---------------|----------------|----------------|-----------------|---------------|----------------|----------------|---------------|----------------|----------------|-----------------|
| 0 | 1.761,877.325 | 1.744,677.836 | 1.753,317.716 | 1.749,565.941 | 1.759,384.624 | 1.763,400.000 | 1.794,900.000 | 1.732,531.426 | 1.732,526.814 | 1.733,883.954 | 1.735,430.248 | 1.757,800.000 | 1.717,300.000 | 1.753,200.000 |
| 1 | 1.756,407.896 | 1.733,466.452 | 1.737,641.504 | 1.739,836.645 | 1.747,576.577 | 1.712,000.000 | 1.742,500.000 | 1.733,790.778 | 1.734,398.485 | 1.735,811.740 | 1.737,110.533 | 1.754,500.000 | 1.724,000.000 | 1.759,900.000 |
| 2 | 1.754,444.442 | 1.735,828.900 | 1.739,272.135 | 1.740,794.079 | 1.750,863.001 | 1.758,000.000 | 1.804,000.000 | 1.735,955.892 | 1.734,433.227 | 1.735,927.786 | 1.737,263.808 | 1.719,800.000 | 1.730,300.000 | 1.755,600.000 |
| 3 | 1.745,749.714 | 1.740,445.636 | 1.749,237.878 | 1.749,816.871 | 1.749,816.871 | 1.700,300.000 | 1.818,500.000 | 1.734,815.677 | 1.735,031.658 | 1.736,843.088 | 1.737,987.235 | 1.686,300.000 | 1.765,900.000 | 1.779,800.000 |
| 4 | 1.733,568.727 | 1.729,484.889 | 1.728,938.412 | 1.733,592.149 | 1.742,113.248 | 1.766,400.000 | 1.708,800.000 | 1.735,128.356 | 1.735,452.086 | 1.737,120.799 | 1.738,793.702 | 1.767,300.000 | 1.733,300.000 | 1.759,200.000 |
| 5 | 1.727,964.181 | 1.723,277.679 | 1.725,427.437 | 1.727,321.616 | 1.736,693.521 | 1.738,500.000 | 1.710,400.000 | 1.736,270.145 | 1.736,331.860 | 1.737,850.119 | 1.740,340.490 | 1.758,200.000 | 1.717,900.000 | 1.739,900.000 |
| 6 | 1.725,841.351 | 1.716,563.761 | 1.723,081.626 | 1.725,702.469 | 1.735,415.386 | 1.758,000.000 | 1.742,900.000 | 1.735,559.708 | 1.736,455.200 | 1.738,820.970 | 1.740,165.947 | 1.780,700.000 | 1.742,500.000 | 1.713,400.000 |
| 7 | 1.722,343.456 | 1.711,523.869 | 1.719,492.821 | 1.724,003.560 | 1.732,499.267 | 1.698,400.000 | 1.731,400.000 | 1.736,139.862 | 1.737,411.285 | 1.738,841.087 | 1.740,979.176 | 1.708,700.000 | 1.726,500.000 | 1.761,400.000 |
| 8 | 1.715,952.971 | 1.703,541.696 | 1.712,534.695 | 1.720,417.914 | 1.679,100.000 | 1.709,000.000 | 1.736,938.342 | 1.738,000.206 | 1.739,984.548 | 1.741,715.014 | 1.752,400.000 | 1.775,600.000 | 1.740,500.000 | 1.750,000.000 |
| 9 | 1.701,602.676 | 1.696,501.612 | 1.705,597.804 | 1.710,036.650 | 1.721,980.996 | 1.776,400.000 | 1.712,900.000 | 1.737,218.230 | 1.738,548.131 | 1.740,764.383 | 1.742,163.420 | 1.759,400.000 | 1.722,800.000 | 1.764,800.000 |
| 10 | 1.708,709.485 | 1.703,507.469 | 1.707,974.902 | 1.712,561.959 | 1.724,060.121 | 1.672,300.000 | 1.733,400.000 | 1.737,059.426 | 1.738,812.313 | 1.740,973.658 | 1.742,911.460 | 1.760,200.000 | 1.807,900.000 | 1.753,500.000 |
| 11 | 1.708,002.048 | 1.700,334.479 | 1.705,745.789 | 1.712,145.483 | 1.721,887.809 | 1.719,300.000 | 1.729,500.000 | 1.737,667.494 | 1.739,075.343 | 1.741,033.716 | 1.742,470.266 | 1.743,600.000 | 1.812,600.000 | 1.758,700.000 |
| 12 | 1.704,765.267 | 1.690,528.200 | 1.700,680.032 | 1.703,727.332 | 1.714,582.320 | 1.739,700.000 | 1.711,800.000 | 1.737,833.800 | 1.738,990.635 | 1.741,004.691 | 1.742,710.814 | 1.768,300.000 | 1.754,700.000 | 1.771,800.000 |
| 13 | 1.693,927.960 | 1.683,230.446 | 1.695,709.392 | 1.702,676.948 | 1.712,420.697 | 1.731,000.000 | 1.690,500.000 | 1.737,127.555 | 1.738,180.555 | 1.740,335.149 | 1.741,471.936 | 1.698,100.000 | 1.745,800.000 | 1.745,100.000 |
| 14 | 1.696,634.570 | 1.681,852.317 | 1.695,295.636 | 1.701,256.215 | 1.712,257.191 | 1.725,700.000 | 1.699,700.000 | 1.737,522.088 | 1.738,958.841 | 1.740,701.751 | 1.742,059.832 | 1.740,500.000 | 1.699,000.000 | 1.728,500.000 |
| 15 | 1.679,551.926 | 1.671,142.177 | 1.668,416.481 | 1.697,105.648 | 1.709,669.411 | 1.752,700.000 | 1.720,700.000 | 1.736,261.928 | 1.737,987.431 | 1.740,131.961 | 1.741,684.474 | 1.752,300.000 | 1.705,000.000 | 1.746,200.000 |
| 16 | 1.667,903.751 | 1.670,413.255 | 1.678,747.207 | 1.691,187.011 | 1.704,670.219 | 1.719,900.000 | 1.699,800.000 | 1.735,533.670 | 1.737,210.612 | 1.738,579.355 | 1.740,344.603 | 1.750,300.000 | 1.781,300.000 | 1.752,600.000 |
| 17 | 1.682,336.989 | 1.679,636.457 | 1.684,708.172 | 1.696,051.548 | 1.706,414.493 | 1.718,300.000 | 1.727,400.000 | 1.736,647.826 | 1.737,903.626 | 1.739,991.371 | 1.741,869.751 | 1.754,200.000 | 1.752,800.000 | 1.780,200.000 |
| 18 | 1.696,024.732 | 1.688,157.613 | 1.693,550.892 | 1.701,124.435 | 1.711,161.137 | 1.689,300.000 | 1.685,300.000 | 1.736,875.723 | 1.738,169.859 | 1.740,186.111 | 1.742,150.532 | 1.723,700.000 | 1.764,500.000 | 1.720,000.000 |
| 19 | 1.700,044.430 | 1.691,379.251 | 1.700,787.282 | 1.703,484.170 | 1.713,713.338 | 1.735,600.000 | 1.698,900.000 | 1.736,115.268 | 1.738,516.054 | 1.740,906.141 | 1.742,896.568 | 1.716,700.000 | 1.778,500.000 | 1.739,500.000 |
| 20 | 1.710,839.931 | 1.700,701.910 | 1.706,890.340 | 1.710,112.768 | 1.717,662.092 | 1.703,300.000 | 1.750,400.000 | 1.735,637.929 | 1.737,519.099 | 1.740,228.710 | 1.741,971.117 | 1.776,900.000 | 1.761,300.000 | 1.752,300.000 |
| 21 | 1.710,751.717 | 1.696,906.569 | 1.701,229.069 | 1.706,919.694 | 1.714,832.769 | 1.726,500.000 | 1.722,600.000 | 1.736,494.732 | 1.739,238.409 | 1.742,035.789 | 1.744,490.773 | 1.747,300.000 | 1.730,200.000 | 1.727,900.000 |
| 22 | 1.707,031.048 | 1.694,596.312 | 1.700,606.809 | 1.707,393.546 | 1.714,688.122 | 1.688,200.000 | 1.719,400.000 | 1.735,343.187 | 1.738,714.382 | 1.741,858.543 | 1.743,975.469 | 1.763,800.000 | 1.738,800.000 | 1.757,900.000 |
| 23 | 1.720,840.960 | 1.710,023.803 | 1.720,653.995 | 1.713,995.947 | 1.719,166.026 | 1.728,600.000 | 1.726,500.000 | 1.735,401.486 | 1.738,931.803 | 1.743,158.826 | 1.745,299.877 | 1.737,700.000 | 1.745,300.000 | 1.764,500.000 |
| 24 | 1.731,576.174 | 1.715,940.520 | 1.714,003.336 | 1.715,023.292 | 1.717,926.947 | 1.758,400.000 | 1.722,100.000 | 1.734,596.940 | 1.738,386.519 | 1.742,505.032 | 1.744,500.841 | 1.749,300.000 | 1.761,700.000 | 1.760,400.000 |
| 25 | 1.752,003.292 | 1.731,855.788 | 1.727,114.867 | 1.724,798.309 | 1.728,919.793 | 1.724,100.000 | 1.708,500.000 | 1.737,772.973 | 1.742,462.185 | 1.746,923.721 | 1.749,250.474 | 1.773,200.000 | 1.764,900.000 | 1.739,600.000 |
| 26 | 1.755,434.381 | 1.735,918.727 | 1.728,805.371 | 1.726,260.802 | 1.725,855.889 | 1.717,700.000 | 1.710,400.000 | 1.736,256.689 | 1.741,221.711 | 1.746,794.790 | 1.750,257.624 | 1.782,600.000 | 1.737,100.000 | 1.737,500.000 |
| 27 | 1.751,078.472 | 1.732,474.044 | 1.721,053.568 | 1.718,911.936 | 1.720,139.844 | 1.700,200.000 | 1.697,500.000 | 1.737,465.109 | 1.745,018.430 | 1.752,091.624 | 1.756,426.168 | 1.732,400.000 | 1.741,600.000 | 1.753,500.000 |
| 28 | 1.727,987.920 | 1.710,624.126 | 1.706,925.069 | 1.698,887.785 | 1.701,059.926 | 1.699,000.000 | 1.699,700.000 | 1.742,590.726 | 1.758,725.331 | 1.772,881.062 | 1.782,612.036 | 1.789,600.000 | 1.790,700.000 | 1.810,600.000 |
| 29 | 1.723,009.291 | 1.712,857.933 | 1.705,443.834 | 1.699,764.161 | 1.701,619.618 | 1.711,900.000 | 1.710,500.000 | 1.742,138.296 | 1.764,672.256 | 1.781,900.380 | 1.793,731.882 | 1.863,100.000 | 1.825,200.000 | 1.818,400.000 |
| 30 | 1.697,576.775 | 1.687,242.321 | 1.674,099.012 | 1.668,454.050 | 1.668,926.625 | 1.649,200.000 | 1.683,400.000 | 1.856,009.443 | 1.902,117.823 | 1.941,790.973 | 1.967,891.756 | 1.975,400.000 | 1.987,200.000 | 2.005,800.000 |
| 31 | 1.719,510.719 | 1.708,216.456 | 1.711,448.892 | 1.713,925.259 | 1.721,720.489 | 1.721,178.125 | 1.722,078.125 | 1.739,930.530 | 1.744,606.260 | 1.749,349.648 | 1.752,719.607 | 1.759,221.875 | 1.759,221.875 | 1.763,116.625 |
| 32 | 1.812,724.891 | 1.801,299.369 | 1.783,560.981 | 1.771,127.395 | 1.765,849.310 | 1.739,200.000 | 1.739,200.000 | 957,760.109 | 963,964.253 | 981,205.278 | 986,644.973 | 1,017,100.000 | 1,024,300.000 | 989,330.000 |
| 33 | 1.785,198.082 | 1.782,599.042 | 1.764,708.279 | 1.754,286.257 | 1.750,603.678 | 1.748,500.000 | 1.734,000.000 | 1,477,089.688 | 1,452,367.316 | 1,435,934.090 | 1,428,928.547 | 1,415,800.000 | 1,418,800.000 | 1,417,900.000 |
| 34 | 1.772,561.111 | 1.767,491.083 | 1.753,649.952 | 1.744,671.079 | 1.742,356.129 | 1.769,300.000 | 1.735,800.000 | 1,623,949.166 | 1,609,532.509 | 1,597,061.854 | 1,592,797.469 | 1,560,600.000 | 1,576,000.000 | 1,592,600.000 |
| 35 | 1.762,468.180 | 1.755,171.008 | 1.742,857.780 | 1.734,034.726 | 1.733,128.170 | 1.732,800.000 | 1.735,800.000 | 1,656,747.957 | 1,650,264.891 | 1,644,001.022 | 1,641,698.771 | 1,625,200.000 | 1,654,900.000 | 1,635,900.000 |
| 36 | 1.751,084.744 | 1.744,819.021 | 1.731,575.083 | 1.727,988.589 | 1.726,739.343 | 1.728,600.000 | 1.706,600.000 | 1,678,449.622 | 1,673,912.651 | 1,669,885.767 | 1,668,314.472 | 1,642,100.000 | 1,683,100.000 | 1,643,100.000 |
| 37 | 1.737,469.976 | 1.729,200.256 | 1.720,955.403 | 1.720,387.855 | 1.719,147.209 | 1.712,900.000 | 1.787,300.000 | 1,694,329.815 | 1,691,870.813 | 1,689,952.951 | 1,689,736.978 | 1,672,300.000 | 1,699,700.000 | 1,699,700.000 |
| 38 | 1.737,533.894 | 1.724,369.474 | 1.716,962.208 | 1.710,620.923 | 1.711,850.544 | 1.691,900.000 | 1.761,400.000 | 1,696,281.632 | 1,695,900.348 | 1,694,668.623 | 1,695,468.826 | 1,674,600.000 | 1,697,200.000 | 1,679,200.000 |
| 39 | 1.715,081.365 | 1.703,943.951 | 1.703,301.172 | 1.705,983.119 | 1.709,368.578 | 1.706,400.000 | 1.660,200.000 | 1,700,949.641 | 1,700,918.018 | 1,700,662.603 | 1,701,381.229 | 1,735,200.000 | 1,703,000.000 | 1,694,800.000 |
| 40 | 1.686,685.834 | 1.682,635.925 | 1.690,744.526 | 1.694,008.976 | 1.698, | | | | | | | | | |

| target size | half 70m | half 160m | half 410m | half 1b | half 2.8b | half 6.9b | half 12b | quarter 70m | quarter 160m | quarter 410m | quarter 1b | quarter 2.8b | quarter 6.9b | quarter 12b |
|-------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| 0 | 1,747,500,000 | 1,780,200,000 | 1,752,900,000 | 1,700,100,000 | 1,729,500,000 | 1,733,700,000 | 1,715,500,000 | 1,774,800,000 | 1,733,800,000 | 1,758,800,000 | 1,715,800,000 | 1,781,400,000 | 1,778,900,000 | 1,726,000,000 |
| 1 | 1,771,400,000 | 1,735,900,000 | 1,707,600,000 | 1,717,800,000 | 1,756,600,000 | 1,731,400,000 | 1,729,600,000 | 1,746,500,000 | 1,781,100,000 | 1,767,400,000 | 1,746,500,000 | 1,742,600,000 | 1,770,500,000 | 1,736,400,000 |
| 2 | 1,705,000,000 | 1,745,400,000 | 1,731,800,000 | 1,731,800,000 | 1,687,200,000 | 1,774,100,000 | 1,754,000,000 | 1,794,900,000 | 1,781,500,000 | 1,769,500,000 | 1,768,000,000 | 1,763,000,000 | 1,770,000,000 | 1,738,000,000 |
| 3 | 1,760,400,000 | 1,712,100,000 | 1,743,700,000 | 1,745,200,000 | 1,726,600,000 | 1,733,500,000 | 1,763,800,000 | 1,791,200,000 | 1,776,600,000 | 1,769,500,000 | 1,776,000,000 | 1,791,900,000 | 1,748,200,000 | 1,782,600,000 |
| 4 | 1,693,800,000 | 1,720,300,000 | 1,726,500,000 | 1,696,900,000 | 1,748,200,000 | 1,697,300,000 | 1,723,500,000 | 1,698,800,000 | 1,740,300,000 | 1,789,400,000 | 1,761,900,000 | 1,792,900,000 | 1,768,000,000 | 1,754,400,000 |
| 5 | 1,711,200,000 | 1,746,400,000 | 1,722,900,000 | 1,722,900,000 | 1,680,700,000 | 1,715,700,000 | 1,714,000,000 | 1,761,700,000 | 1,765,300,000 | 1,777,000,000 | 1,720,700,000 | 1,763,800,000 | 1,768,000,000 | 1,746,100,000 |
| 6 | 1,700,900,000 | 1,717,100,000 | 1,673,600,000 | 1,753,500,000 | 1,774,500,000 | 1,735,800,000 | 1,771,000,000 | 1,764,900,000 | 1,757,900,000 | 1,726,200,000 | 1,744,000,000 | 1,742,400,000 | 1,761,400,000 | 1,782,500,000 |
| 7 | 1,737,400,000 | 1,743,900,000 | 1,710,200,000 | 1,741,300,000 | 1,766,200,000 | 1,688,400,000 | 1,730,600,000 | 1,755,400,000 | 1,690,800,000 | 1,790,000,000 | 1,710,700,000 | 1,751,600,000 | 1,740,100,000 | 1,749,900,000 |
| 8 | 1,705,800,000 | 1,721,300,000 | 1,728,200,000 | 1,723,200,000 | 1,720,900,000 | 1,676,100,000 | 1,738,200,000 | 1,799,000,000 | 1,726,100,000 | 1,799,800,000 | 1,773,900,000 | 1,806,700,000 | 1,704,900,000 | 1,737,600,000 |
| 9 | 1,754,600,000 | 1,756,100,000 | 1,762,300,000 | 1,722,400,000 | 1,767,400,000 | 1,772,800,000 | 1,708,800,000 | 1,785,800,000 | 1,741,600,000 | 1,756,200,000 | 1,763,000,000 | 1,744,900,000 | 1,708,900,000 | 1,730,000,000 |
| 10 | 1,737,700,000 | 1,726,600,000 | 1,761,000,000 | 1,715,500,000 | 1,724,600,000 | 1,754,200,000 | 1,757,500,000 | 1,753,600,000 | 1,763,700,000 | 1,757,000,000 | 1,766,000,000 | 1,748,400,000 | 1,736,000,000 | 1,770,700,000 |
| 11 | 1,678,500,000 | 1,738,600,000 | 1,691,800,000 | 1,753,500,000 | 1,704,100,000 | 1,744,300,000 | 1,765,900,000 | 1,774,300,000 | 1,773,600,000 | 1,764,500,000 | 1,746,700,000 | 1,752,300,000 | 1,741,700,000 | 1,740,200,000 |
| 12 | 1,721,500,000 | 1,696,300,000 | 1,678,400,000 | 1,699,000,000 | 1,714,100,000 | 1,718,900,000 | 1,765,400,000 | 1,776,700,000 | 1,734,000,000 | 1,771,900,000 | 1,716,000,000 | 1,736,100,000 | 1,727,900,000 | 1,730,000,000 |
| 13 | 1,751,300,000 | 1,752,400,000 | 1,708,000,000 | 1,700,400,000 | 1,698,800,000 | 1,754,600,000 | 1,719,200,000 | 1,714,300,000 | 1,744,900,000 | 1,749,500,000 | 1,741,100,000 | 1,729,200,000 | 1,746,400,000 | 1,742,400,000 |
| 14 | 1,737,600,000 | 1,731,900,000 | 1,735,500,000 | 1,682,700,000 | 1,695,200,000 | 1,701,600,000 | 1,730,800,000 | 1,751,000,000 | 1,749,400,000 | 1,731,600,000 | 1,736,000,000 | 1,738,200,000 | 1,768,800,000 | 1,721,100,000 |
| 15 | 1,693,700,000 | 1,725,700,000 | 1,774,900,000 | 1,713,300,000 | 1,775,300,000 | 1,734,900,000 | 1,783,900,000 | 1,728,000,000 | 1,724,300,000 | 1,739,000,000 | 1,742,500,000 | 1,762,500,000 | 1,763,000,000 | 1,782,100,000 |
| 16 | 1,739,000,000 | 1,749,900,000 | 1,740,800,000 | 1,777,700,000 | 1,747,300,000 | 1,749,300,000 | 1,782,800,000 | 1,738,500,000 | 1,740,800,000 | 1,705,500,000 | 1,771,500,000 | 1,729,500,000 | 1,763,500,000 | 1,749,200,000 |
| 17 | 1,741,900,000 | 1,694,200,000 | 1,742,100,000 | 1,771,200,000 | 1,721,200,000 | 1,741,300,000 | 1,703,000,000 | 1,745,700,000 | 1,740,100,000 | 1,767,500,000 | 1,736,700,000 | 1,739,000,000 | 1,741,000,000 | 1,761,400,000 |
| 18 | 1,738,200,000 | 1,703,300,000 | 1,740,700,000 | 1,711,900,000 | 1,705,700,000 | 1,704,600,000 | 1,686,400,000 | 1,737,100,000 | 1,743,400,000 | 1,743,200,000 | 1,737,100,000 | 1,724,500,000 | 1,735,100,000 | 1,758,000,000 |
| 19 | 1,728,500,000 | 1,725,600,000 | 1,729,100,000 | 1,719,600,000 | 1,725,700,000 | 1,679,600,000 | 1,747,000,000 | 1,726,900,000 | 1,791,300,000 | 1,728,300,000 | 1,740,300,000 | 1,707,300,000 | 1,736,600,000 | 1,719,900,000 |
| 20 | 1,691,300,000 | 1,761,400,000 | 1,714,800,000 | 1,695,900,000 | 1,702,800,000 | 1,703,400,000 | 1,727,700,000 | 1,727,000,000 | 1,767,400,000 | 1,707,000,000 | 1,740,800,000 | 1,771,400,000 | 1,725,200,000 | 1,795,800,000 |
| 21 | 1,677,600,000 | 1,738,600,000 | 1,699,400,000 | 1,737,100,000 | 1,691,400,000 | 1,680,900,000 | 1,697,700,000 | 1,752,200,000 | 1,741,500,000 | 1,750,000,000 | 1,777,000,000 | 1,751,000,000 | 1,741,700,000 | 1,746,900,000 |
| 22 | 1,740,300,000 | 1,713,600,000 | 1,716,200,000 | 1,668,000,000 | 1,698,100,000 | 1,718,800,000 | 1,722,600,000 | 1,774,400,000 | 1,710,600,000 | 1,757,000,000 | 1,774,700,000 | 1,797,700,000 | 1,746,100,000 | 1,761,000,000 |
| 23 | 1,748,200,000 | 1,666,100,000 | 1,664,400,000 | 1,729,100,000 | 1,677,100,000 | 1,734,500,000 | 1,731,200,000 | 1,747,300,000 | 1,737,000,000 | 1,766,600,000 | 1,782,700,000 | 1,761,900,000 | 1,723,400,000 | 1,757,800,000 |
| 24 | 1,691,700,000 | 1,735,300,000 | 1,716,500,000 | 1,722,300,000 | 1,733,100,000 | 1,699,600,000 | 1,678,600,000 | 1,735,800,000 | 1,734,200,000 | 1,748,600,000 | 1,712,500,000 | 1,743,000,000 | 1,724,000,000 | 1,709,000,000 |
| 25 | 1,728,800,000 | 1,656,600,000 | 1,726,600,000 | 1,700,200,000 | 1,697,900,000 | 1,702,800,000 | 1,732,800,000 | 1,731,200,000 | 1,766,800,000 | 1,740,300,000 | 1,724,100,000 | 1,736,800,000 | 1,718,200,000 | 1,782,000,000 |
| 26 | 1,700,700,000 | 1,715,200,000 | 1,654,200,000 | 1,694,200,000 | 1,703,700,000 | 1,684,400,000 | 1,733,500,000 | 1,708,800,000 | 1,715,100,000 | 1,739,000,000 | 1,734,800,000 | 1,721,100,000 | 1,695,100,000 | 1,749,000,000 |
| 27 | 1,724,500,000 | 1,645,400,000 | 1,754,400,000 | 1,707,800,000 | 1,724,900,000 | 1,721,800,000 | 1,672,900,000 | 1,703,000,000 | 1,715,900,000 | 1,711,100,000 | 1,697,200,000 | 1,718,800,000 | 1,708,100,000 | 1,680,500,000 |
| 28 | 1,658,500,000 | 1,871,300,000 | 1,800,300,000 | 1,834,700,000 | 1,812,500,000 | 1,719,500,000 | 1,791,700,000 | 1,819,100,000 | 2,090,300,000 | 2,020,500,000 | 1,992,900,000 | 1,734,000,000 | 1,858,000,000 | 1,735,400,000 |
| 29 | 1,692,600,000 | 1,667,500,000 | 1,686,700,000 | 1,655,700,000 | 1,657,100,000 | 1,676,700,000 | 1,707,500,000 | 1,697,300,000 | 1,771,300,000 | 1,731,300,000 | 1,717,800,000 | 1,677,500,000 | 1,741,000,000 | 1,691,500,000 |
| 30 | 1,655,000,000 | 1,690,700,000 | 1,686,200,000 | 1,647,300,000 | 1,683,200,000 | 1,639,600,000 | 1,628,200,000 | 1,760,800,000 | 1,689,800,000 | 1,675,800,000 | 1,691,500,000 | 1,667,800,000 | 1,698,000,000 | 1,697,800,000 |
| 31 | 1,666,800,000 | 1,541,100,000 | 1,637,100,000 | 1,564,000,000 | 1,513,400,000 | 1,579,400,000 | 1,554,100,000 | 1,562,100,000 | 1,515,300,000 | 1,536,200,000 | 1,484,100,000 | 1,477,500,000 | 1,436,800,000 | 1,507,600,000 |
| Avg Context | 1,721,093,750 | 1,713,987,500 | 1,718,771,875 | 1,711,456,250 | 1,710,421,875 | 1,710,321,875 | 1,720,766,625 | 1,739,512,500 | 1,738,518,750 | 1,744,262,500 | 1,735,328,125 | 1,737,928,125 | 1,733,200,000 | 1,732,553,125 |
| 32 | 1,878,300,000 | 1,871,300,000 | 1,800,300,000 | 1,834,700,000 | 1,812,500,000 | 1,719,500,000 | 1,791,700,000 | 1,819,100,000 | 2,090,300,000 | 2,020,500,000 | 1,992,900,000 | 1,734,000,000 | 1,858,000,000 | 1,735,400,000 |
| 33 | 1,842,000,000 | 1,854,700,000 | 1,836,200,000 | 1,829,200,000 | 1,869,400,000 | 1,834,500,000 | 1,749,400,000 | 2,458,900,000 | 2,344,000,000 | 2,383,700,000 | 2,258,800,000 | 2,249,000,000 | 2,267,000,000 | 2,275,200,000 |
| 34 | 1,842,600,000 | 1,900,800,000 | 1,843,500,000 | 1,862,600,000 | 1,778,400,000 | 1,812,100,000 | 1,832,600,000 | 2,650,800,000 | 2,597,500,000 | 2,592,700,000 | 2,527,000,000 | 2,483,300,000 | 2,511,000,000 | 2,568,800,000 |
| 35 | 1,841,600,000 | 1,875,500,000 | 1,883,700,000 | 1,894,500,000 | 1,866,800,000 | 1,843,400,000 | 1,868,600,000 | 2,520,900,000 | 2,585,300,000 | 2,586,100,000 | 2,591,600,000 | 2,586,700,000 | 2,599,000,000 | 2,526,000,000 |
| 36 | 1,871,800,000 | 1,901,500,000 | 1,932,600,000 | 1,931,300,000 | 1,926,100,000 | 1,932,300,000 | 1,993,200,000 | 1,789,000,000 | 1,742,400,000 | 1,743,700,000 | 1,696,700,000 | 1,715,400,000 | 1,754,000,000 | 1,789,500,000 |
| 37 | 1,885,600,000 | 1,921,800,000 | 1,997,900,000 | 1,915,800,000 | 1,991,000,000 | 2,077,100,000 | 2,038,100,000 | 1,736,300,000 | 1,739,900,000 | 1,728,100,000 | 1,651,000,000 | 1,664,500,000 | 1,681,500,000 | 1,700,300,000 |
| 38 | 1,998,200,000 | 2,081,300,000 | 2,002,700,000 | 2,115,900,000 | 2,156,700,000 | 2,229,100,000 | 2,254,500,000 | 1,762,300,000 | 1,710,100,000 | 1,711,100,000 | 1,737,600,000 | 1,718,500,000 | 1,736,000,000 | 1,753,000,000 |
| 39 | 2,032,100,000 | 2,158,400,000 | 2,168,100,000 | 2,189,200,000 | 2,246,300,000 | 2,226,200,000 | 2,369,000,000 | 1,809,200,000 | 1,768,100,000 | 1,799,700,000 | 1,747,600,000 | 1,744,100,000 | 1,738,000,000 | 1,758,800,000 |
| 40 | 1,723,200,000 | 1,690,600,000 | 1,715,300,000 | 1,692,500,000 | 1,691,600,000 | 1,691,900,000 | 1,719,100,000 | 1,765,200,000 | 1,751,200,000 | 1,789,500,000 | 1,746,300,000 | 1,757,000,000 | 1,739,000,000 | 1,774,000,000 |
| 41 | 1,738,700,000 | 1,721,600,000 | 1,699,700,000 | 1,709,900,000 | 1,706,600,000 | 1,668,700,000 | 1,736,000,000 | 1,734,000,000 | 1,777,100,000 | 1,770,000,000 | 1,803,300,000 | 1,787,000,000 | 1,789,500,000 | 1,788,000,000 |
| 42 | 1,707,400,000 | 1,795,200,000 | 1,719,700,000 | 1,721,300,000 | 1,743,300,000 | 1,700,900,000 | 1,725,800,000 | 1,822,600,000 | 1,814,600,000 | 1,763,000,000 | 1,808,000,000 | 1,793,500,000 | 1,787,000,000 | 1,740,000,000 |
| 43 | 1,717,000,000 | 1,708,000,000 | 1,720,500,000 | 1,713,600,000 | 1,733,400,000 | 1,727,600,000 | 1,687,200,000 | 1,839,800,000 | 1,774,600,000 | 1,838,300,000 | 1,761,600,000 | 1,754,700,000 | 1,787,300,000 | 1,805,000,000 |
| 44 | 1,753,600,000 | 1,742,000,000 | 1,764,800,000 | 1,694,000,000 | 1,779,700,000 | 1,686,500,000 | 1,721,000,000 | 1,760,900,000 | 1,806,500,000 | 1,817,700,000 | 1,766,000,000 | 1,787,800,000 | 1,752,900,000 | 1,782,000,000 |
| 45 | 1,729,900,000 | | | | | | | | | | | | | |

| target size | half 70m | half 160m | half 410m | half 1b | half 2.8b | half 6.9b | half 12b | quarter 70m | quarter 160m | quarter 410m | quarter 1b | quarter 2.8b | quarter 6.9b | quarter 12b |
|------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|--------------|--------------|------------|--------------|--------------|-------------|
| 1 | 53,888,945 | 55,774,256 | 56,835,849 | 57,399,791 | 58,749,349 | 59,572,656 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 | 67,038,540 |
| 2 | 53,633,778 | 55,585,291 | 56,294,740 | 57,125,864 | 58,375,426 | 58,787,083 | 59,024,756 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 3 | 53,966,515 | 55,464,751 | 56,571,797 | 57,248,585 | 58,522,726 | 59,007,609 | 59,282,618 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 4 | 53,788,979 | 55,012,249 | 55,852,524 | 56,732,428 | 58,632,577 | 58,585,481 | 59,106,526 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 5 | 52,887,357 | 54,130,559 | 55,380,593 | 56,487,949 | 57,789,441 | 58,410,016 | 58,836,366 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 6 | 53,346,585 | 54,626,718 | 56,348,320 | 56,603,397 | 58,159,466 | 58,648,425 | 58,945,615 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 7 | 52,836,947 | 54,888,957 | 55,457,422 | 56,613,020 | 57,855,215 | 58,349,009 | 58,580,014 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 8 | 52,478,474 | 54,201,487 | 55,628,345 | 56,114,532 | 57,493,118 | 58,018,440 | 58,920,964 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 9 | 53,600,100 | 55,410,940 | 56,362,313 | 57,618,326 | 58,256,347 | 59,088,717 | 59,814,719 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 10 | 53,868,058 | 55,087,199 | 56,362,786 | 57,058,982 | 58,314,177 | 58,745,840 | 59,526,404 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 11 | 52,884,280 | 54,383,701 | 56,075,362 | 56,818,188 | 58,050,191 | 58,552,813 | 58,996,784 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 12 | 52,928,646 | 54,849,387 | 55,982,049 | 56,771,448 | 58,094,228 | 58,565,819 | 58,811,948 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 13 | 52,690,921 | 54,559,756 | 55,697,194 | 56,532,692 | 57,846,992 | 58,355,202 | 58,946,804 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 14 | 52,781,679 | 54,368,904 | 56,019,791 | 56,847,208 | 57,869,958 | 58,242,828 | 58,685,929 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 15 | 53,765,298 | 55,537,152 | 56,289,403 | 57,279,213 | 58,113,960 | 58,741,123 | 59,494,649 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 16 | 53,740,078 | 55,037,415 | 55,929,818 | 56,941,982 | 57,702,245 | 58,282,107 | 58,697,163 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 17 | 53,610,921 | 55,130,291 | 56,781,859 | 57,648,192 | 58,161,184 | 58,670,246 | 59,044,372 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 18 | 52,403,931 | 54,507,436 | 56,143,121 | 56,609,104 | 57,872,099 | 58,070,077 | 58,324,067 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 19 | 51,734,926 | 53,970,664 | 55,420,133 | 56,429,139 | 57,389,684 | 58,026,044 | 57,807,255 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 20 | 52,066,233 | 53,844,254 | 55,306,358 | 56,168,419 | 56,834,780 | 57,322,399 | 57,930,955 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 21 | 54,499,372 | 55,151,054 | 56,817,487 | 56,104,835 | 57,282,099 | 57,680,431 | 58,491,658 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 22 | 53,002,234 | 54,689,980 | 55,709,870 | 56,393,115 | 56,890,668 | 57,461,711 | 58,953,014 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 23 | 54,475,618 | 54,209,390 | 55,114,979 | 55,800,549 | 56,745,566 | 56,890,699 | 57,313,138 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 24 | 52,696,844 | 53,885,907 | 54,088,030 | 55,057,040 | 55,463,131 | 56,208,303 | 56,444,958 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 25 | 50,318,283 | 51,269,205 | 52,481,347 | 53,465,984 | 53,891,662 | 54,781,920 | 54,713,197 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 26 | 50,229,788 | 51,354,758 | 52,362,605 | 53,242,078 | 53,893,195 | 54,326,699 | 54,830,116 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 27 | 52,206,436 | 52,688,837 | 53,341,050 | 53,769,840 | 54,270,899 | 54,576,546 | 55,060,048 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 28 | 49,682,706 | 50,249,019 | 51,876,101 | 52,495,005 | 52,542,959 | 53,202,973 | 53,174,103 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 29 | 49,857,350 | 50,055,633 | 51,690,368 | 51,477,619 | 51,245,707 | 51,586,655 | 52,023,933 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 30 | 48,055,193 | 47,429,451 | 48,033,530 | 48,151,423 | 48,174,882 | 48,183,182 | 48,215,800 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 31 | 39,251,259 | 39,137,018 | 39,388,512 | 39,213,772 | 39,230,840 | 39,250,211 | 39,423,466 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 32 | 52,167,024 | 53,402,932 | 54,503,993 | 55,212,772 | 56,103,781 | 56,535,303 | 56,929,119 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 33 | 50,487,094 | 49,193,302 | 48,748,380 | 48,353,375 | 48,143,528 | 47,788,231 | 47,682,812 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 34 | 61,392,551 | 59,894,701 | 58,602,818 | 58,523,094 | 57,350,374 | 57,349,346 | 56,865,365 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 35 | 61,328,405 | 61,880,320 | 61,897,331 | 61,603,725 | 61,037,948 | 60,910,761 | 61,148,637 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 36 | 62,108,896 | 64,859,473 | 66,373,167 | 66,487,737 | 66,963,760 | 67,607,923 | 67,185,378 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 37 | 64,491,225 | 66,128,068 | 68,573,633 | 68,479,538 | 69,624,085 | 69,468,930 | 70,076,197 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 38 | 64,176,596 | 67,738,610 | 71,189,574 | 72,648,126 | 74,649,180 | 75,348,320 | 76,011,136 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 39 | 83,023,501 | 86,268,700 | 92,242,626 | 95,541,179 | 100,137,087 | 103,141,517 | 104,996,137 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 40 | 114,997,010 | 123,345,929 | 132,119,453 | 137,345,936 | 144,729,841 | 148,911,774 | 152,506,274 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 41 | 72,923,797 | 75,060,147 | 76,865,835 | 77,568,398 | 78,566,453 | 79,311,885 | 80,082,621 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 42 | 62,139,748 | 64,712,748 | 65,532,787 | 65,890,626 | 66,512,689 | 66,430,631 | 67,029,855 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 43 | 61,518,278 | 63,119,965 | 63,331,920 | 63,642,209 | 64,136,809 | 63,861,430 | 63,990,245 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 44 | 58,893,899 | 59,747,917 | 60,630,537 | 60,111,792 | 60,897,906 | 60,413,416 | 60,685,337 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 45 | 57,229,696 | 57,312,370 | 57,433,177 | 58,104,220 | 57,890,592 | 58,313,878 | 58,439,482 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 46 | 55,873,454 | 56,139,949 | 56,811,592 | 56,589,596 | 57,622,149 | 57,870,333 | 58,004,104 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 47 | 54,966,654 | 55,910,428 | 56,508,715 | 56,557,283 | 57,408,451 | 57,383,704 | 58,025,074 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| 48 | 55,258,100 | 55,540,073 | 56,121,577 | 56,868,270 | 57,182,726 | 57,490,540 | 58,014,824 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| Avg Continuation | 64,389,138 | 65,969,877 | 67,517,741 | 68,202,161 | 69,326,240 | 69,882,011 | 70,425,070 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |
| Avg Sentence | 56,495,689 | 57,853,725 | 59,113,028 | 59,813,184 | 60,786,735 | 61,262,262 | 61,708,935 | 60,505,804 | 61,439,052 | 62,372,300 | 63,305,548 | 64,238,796 | 65,172,044 | 66,105,292 |

Table 6: Bi-gram Statistics for half-memorized and quarter-memorized content.