

Capacity Matters: Investigating Transformer Models for Real-World Data Memorization

Anonymous ACL submission

Abstract

Transformer models' memorization capacity studies often focus on theoretical bounds or use synthetic datasets that lack real-world complexity. This study systematically evaluates how model architecture and data configurations influence the capacity of decoder transformers using datasets derived from the Systematized Nomenclature of Medicine (SNOMED) knowledge graph: triplets, representing static connections, and sequences, simulating complex relation patterns.

Our findings highlight key factors affecting training dynamics and memorization. Embedding size is the primary determinant of learning speed and capacity, while additional layers provide limited benefits and may hinder performance on simpler datasets. Activation functions play a crucial role, with Softmax demonstrating greater stability and capacity. Additionally, increased dataset complexity enhances final memorization. These insights improve our understanding of transformer memory mechanisms and provide a framework for optimizing model design with structured real-world data.

1 Introduction

Transformer-based Large Language Models (LLMs) have revolutionized natural language processing by demonstrating remarkable capabilities in tasks ranging from text generation and translation to question answering and summarization. Despite these advances, the fundamental mechanisms underpinning their capacity to memorize and retrieve structured knowledge remain an active area of research. Understanding these mechanisms is crucial for optimizing model performance, making it computationally cheap in order to apply to real-world problems. One particularly impactful example is healthcare, where LLMs could assist clinicians through wearable devices such as smart glasses or watches (Gupta et al., 2024; Wu et al., 2024; Balloccu et al., 2024).

Due to privacy and reliability, the preferred system would be a local on-edge LLM with minimal computational requirements, but with a capacity to memorize all relevant facts in the relevant area of healthcare.

Recent theoretical and empirical studies have sought to quantify the memorization capacity of transformers. Kim et al. (2023) introduced mathematical bounds for memory capacity, demonstrating that transformers could memorize $O(d + n + \sqrt{nN})$ parameters, where d, n, N correspond to embedding dimensions, dataset size, and model size, respectively. Additionally, Kajitsuka and Sato (2024) proved, that $\tilde{O}(\sqrt{nN})$ parameters are not only sufficient, but also necessary for some types of transformers. Mahdavi et al. (2024) extended this work by analyzing the effects of multi-head attention on memorization, revealing the interplay between architectural components and the model's ability to store and recall information. The experiments in Härmä et al. (2024) used randomly generated sequences of numbers to evaluate the memorization capabilities of the transformer models on unstructured data. Most capacity studies use synthetic datasets because accurate capacity measurement becomes very difficult in the case of uncontrolled free text content.

The experiments reported in the current paper use sentence data generated from the knowledge graph which, while being controlled, has some of the hierarchical and relational complexity of real-world text content. More specifically, GPT-like transformer models (Brown et al., 2020) were trained to memorize structured sentences derived from the Systematized Nomenclature of Medicine (SNOMED) knowledge graph (KG) (El-Sappagh et al., 2018). SNOMED, a comprehensive medical ontology, encodes semantic relationships between medical concepts, offering a rich dataset to explore memory and retrieval mechanisms under realistic conditions. Exact memorization of a selection of

042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082

083 such relations would be critical, for example, in the
084 healthcare use case described above.

085 By employing both theoretical insights and em-
086 pirical evaluation, this study seeks to answer three
087 key research questions. How can real-world data,
088 such as knowledge graphs, be used to investigate
089 transformers’ memorization capacity? How do archi-
090 tectural variations affect the efficiency and scal-
091 ability of memorization in transformer models?
092 How do dataset structure and complexity influence
093 memorization behavior during training?

094 To measure the memorization capacity of trans-
095 former models, the Maximum Attainable Capacity
096 (MAC) method was used. It is a computationally
097 efficient alternative to the Maximum Library Size
098 (MLS) method. While MLS involves iteratively
099 training models on progressively larger datasets to
100 determine the largest library size that can be fully
101 memorized, MAC evaluates the practical limit of
102 samples a model can retain when trained on a large
103 dataset. Previous research has shown a strong cor-
104 relation between MLS and MAC (Härmä et al.,
105 2024), making MAC an effective and time-efficient
106 choice for this study.

107 Our approach leverages structured datasets con-
108 structed through two methods: triplet generation
109 and sequence generation. Triplets represent static
110 relationships in the form (Concept, Property,
111 Related Concept), providing a baseline for as-
112 sessing memorization. Sequences extend this by
113 simulating graph traversal paths, capturing relation-
114 ship patterns between concepts. These datasets
115 allowed us to empirically analyze how model archi-
116 tecture, training configurations, dataset size, and
117 complexity influence training dynamics and final
118 memorization performance.

119 2 Methods

120 2.1 Data

121 2.1.1 Data Source and Preprocessing

122 To evaluate transformer-based models’ memory
123 and retrieval capabilities, we used SNOMED KG,
124 which encodes medical concepts and their rela-
125 tionships as nodes and edges of a graph. It
126 was accessed using the owlready2 library (Lamy,
127 2017), filtering out non-informative or overly spe-
128 cific properties to ensure meaningful relationships.
129 While graph transformers leverage Graph Neural
130 Networks (Shehzad et al., 2024), our approach pri-
131 oritizes a universal architecture applicable across
132 diverse datasets. Hence, the graph was transformed

133 to: (1) triplets, representing concept-property rela-
134 tionships (see 2.1.2), and (2) sequences, simulating
135 graph traversal paths (see 2.1.3).

136 2.1.2 Triplets Generation

137 The goal of triplet generation was to create a dataset
138 of the form (Concept, Property, Related
139 Concept), capturing semantic relationships in the
140 SNOMED KG. This process (see Figure 1A) in-
141 volves graph initialization and the exclusion of non-
142 informative properties. After the algorithm extracts
143 triplets: for each concept in the KG, it retrieves all
144 allowed properties and their associated related con-
145 cepts. Additionally, when multiple related concepts
146 are associated with a (Concept, Property) pair,
147 one is selected randomly to maintain uniqueness.

148 2.1.3 Sequences Generation

149 The sequence generation simulated graph traversal
150 and encoded local and global graph structures. The
151 complete algorithm is depicted in Figure 1B.

152 The extended graph (G) is constructed from an
153 ontology by: (1) excluding banned properties, as
154 in the triplets generation; (2) along with each re-
155 lationship, adding an edge with opposite direc-
156 tion with a corresponding reversed_ prefix for
157 bidirectional traversal. Additionally, labels were
158 cleaned (metadata were removed) to standardize
159 their format. The sequences were generated to re-
160 flect the traversal path in the graph, capturing both
161 nodes and edges: (node₁, edge₁, node₂, . . . ,
162 node_{n-1}, edge_{n-1}, node_n)

163 For each sequence, the algorithm first selects a
164 random starting node from the full graph G, ensur-
165 ing that the node has at least one unused edge. A
166 subgraph is then created around the starting node
167 using a breadth-first search (BFS) with a depth de-
168 fined by the hops parameter. This step limits the
169 scope of the traversal to a manageable subset of
170 the graph, improving performance by focusing on
171 local neighborhoods.

172 Step 2 of the algorithm generates a sequence
173 of nodes and edges by traversing the subgraph.
174 The algorithm starts from a randomly selected
175 node and go through available edges (neighbors
176 are chosen randomly to introduce variability). Ev-
177 ery time, check that the same (node, edge) pair
178 is not already visited before, maintaining global
179 uniqueness. The traversal stops when: a ran-
180 domly chosen number of edges within a predefined
181 range (edge_count_range) is reached, or no valid
182 neighbors (those, that maintain uniqueness) remain

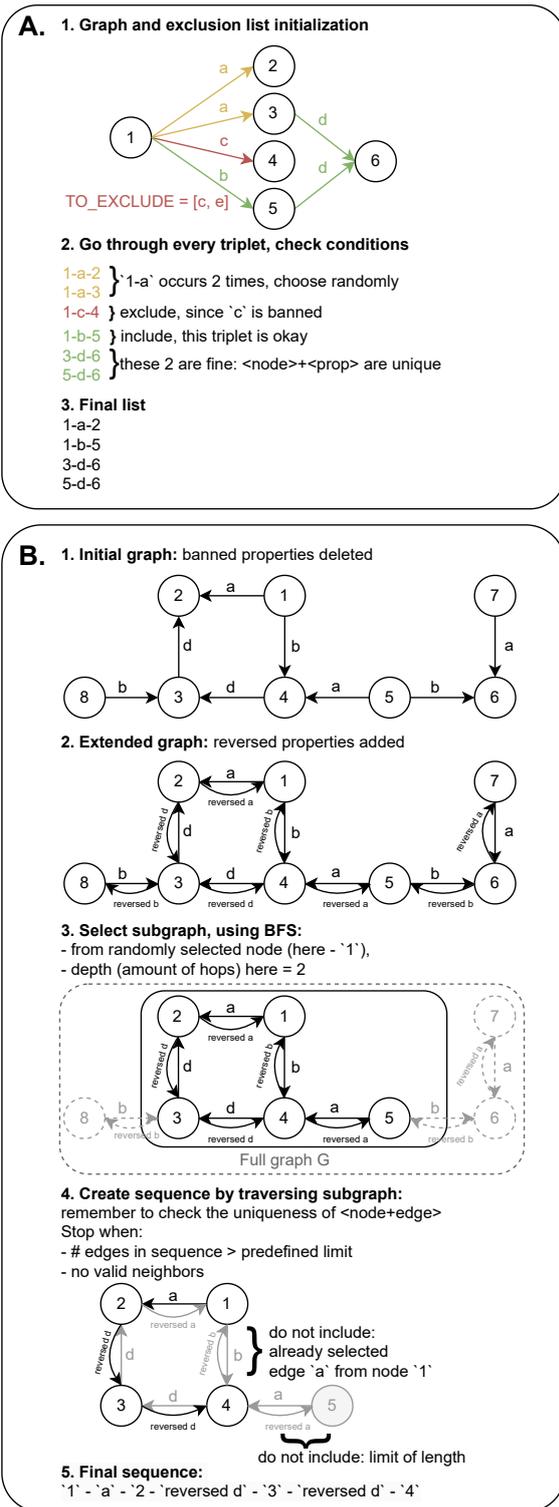


Figure 1: Algorithms of triplets (A) and sequences (B) data generation.

for further traversal.

The above steps are repeated for a specified number of iterations (rows), generating the desired number of sequences.

2.2 Transformers training

To evaluate the ability of transformer models to memorize and retrieve structured data, decoder-only transformer models with variations in architecture were implemented. Each unique node and edge was assigned a distinct integer identifier (ensuring that repeated elements were consistently tokenized), followed by the learned positional encoding. The core architecture consisted of three main components: an embedding layer to map tokenized inputs into continuous vector representations, transformer decoder layers with multi-head attention mechanisms, and a linear output layer to predict target tokens.

For all experiments, the task was to predict a concept, based on the previous concepts and relations. The accuracy was evaluated as: $\frac{\#correct_predictions}{\#total_predictions}$ – the proportion of correctly predicted related concepts to the total number of predictions. Additionally, Maximum Attainable Capacity (MAC) was used as a more suitable metric for measuring the model capacity. As detailed in the introduction, MAC is a computationally efficient alternative to Maximum Library Size (MLS), with results strongly correlated to MLS, making it the preferred choice for this research.

To minimize the effect of randomness, each experiment was repeated 10 times for the first and second setups, and 3 times for the third and fourth setup (see below). All figures and tables present mean values with doubled standard deviations. Training and evaluation followed a consistent protocol for all setups, with the training accuracy evaluated every second epoch, which allowed meaningful comparisons between different configurations.

All code was written in PyTorch v1.13.1+cu117 (Paszke et al., 2017) and Transformers v4.30.2 (Wolf et al., 2019). The cross-entropy loss function was used for optimization, along with the Adam optimizer (Kingma and Ba, 2017) and a learning rate of 0.001. All other settings were kept at their default library implementations, except where specified in experiment configurations. In total, 546 models were trained on NVIDIA A100 GPU with 16GB memory, totaling approximately 3,100 hours of training time. Model sizes ranged from 2.9 to 44.5 million parameters, primarily varying with embedding size and layer count, but also influenced by vocabulary size.

All data and code pertinent to the methods and results presented in this work will be made available at the time of the conference.

2.2.1 Triplets memorization

Three experimental setups on the dataset with triplets were devised to explore the models' behavior. For all of them, since the prediction of the related concept is based on unique combinations of concept and relation, it is straightforward to unambiguously determine whether a related concept was predicted correctly or not.

In the first setup, dataset sizes ranged in $\{50,000; 60,000; \dots; 100,000\}$ samples. The model architecture consisted of a single transformer layer with an embedding size of 128, four attention heads, and a Rectified Linear Unit (ReLU) activation function (Agarap, 2019) with the batch size of 64, and 500 training epochs. This setup focused on evaluating memorization performance under a fixed architecture while varying dataset sizes.

The second experimental setup introduced variations in the transformer architecture, allowing a deeper investigation into the impact of model depth and activation functions. Dataset sizes included 50,000, 70,000, and 100,000 samples, with the numbers of transformer layers set to 1, 2, or 4. Activation functions were varied across ReLU, Gaussian Error Linear Unit (GELU) (Hendrycks and Gimpel, 2023), Randomized Leaky Rectified Linear Unit (RReLU) (Xu et al., 2015), and Softmax (Boltzmann, 1868). To ensure fair comparisons, the total number of model parameters was kept constant across configurations by adjusting the embedding size (`d_model` parameter in PyTorch implementation of Transformers) proportionally to the number of layers, using the formula: $\text{embedding_size} = \lfloor \frac{\text{base_number_of_parameters}}{n_layers} \rfloor$ with a base number of parameters of 128. This approach ensured that variations in performance could be attributed solely to architectural differences rather than changes in the total parameter count. For this setup, however, the batch size was increased to 128, and the number of training epochs was 1000, since it was required for achieving a plateau.

The third setup focused on evaluating the interplay between model depth, and embedding size while keeping other hyperparameters the same. Dataset sizes ranged in $\{1,000; 10,000; 50,000; 100,000\}$ samples.

The architectural variations included transformer layers set to 1 or 2 and base numbers of parameters for embedding sizes in $\{16; 32; 64; 128\}$ (calculated as in the second experiment). Only the Softmax activation function and a fixed number of 4 attention heads were used. To ensure fair comparisons, configurations were designed to evaluate the impact of increasing embedding sizes and model depth on memorization performance. The total parameter count was recalculated for each configuration using the same formula as in the second experiment. For this setup, as previously, the batch size of 128 was used, and the number of training epochs was 500.

2.2.2 Sequences memorization

The dataset for sequence memorization tasks was prepared using the same tokenization process. However, to standardize sequence lengths, padding with zeros was applied at the end of each sequence, serving both as a filler and a marker for sequence termination. The task required distinguishing between nodes and edges, and a node mask was generated to identify the positions of node tokens within the sequence. It enabled the computation of metrics by isolating node positions during the training and evaluation processes. Notably, each node was predicted based on all preceding tokens in the sequence, meaning the last node in a sequence benefited from the most context. This setup provided deeper insights into the transformer model's ability to handle more structured data and its patterns.

The experimental setup was consistent with the previous experiments described in 2.2.1: the embedding size was fixed at 64, with four attention heads, the batch size was set to 128, and the number of epochs to 400. The number of layers varied across $\{1, 2, 4\}$, and the activation functions used were RReLU and Softmax. As before, the model incorporated a learned positional encoding. The dataset sizes were varied in $\{20,000; 50,000; 100,000\}$, representing the number of sequences. Each sequence was limited to 4-6 nodes (and 3-5 edges, respectively), selected randomly. During dataset construction, 5 hops were used to isolate the subgraph (see 2.1.3 for details).

For sequence memorization, accuracy and capacity were measured similarly to the triplet-based experiments, with slight adaptations to account for the sequential structure of the data. Accuracy was defined as the proportion of correctly predicted to-

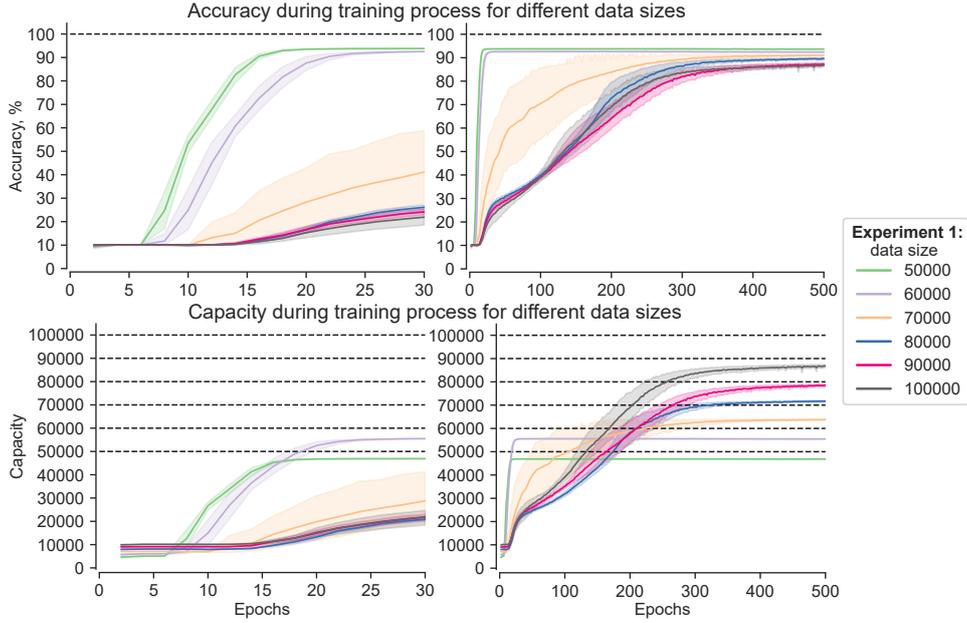


Figure 2: Trends in training accuracy (upper) and capacity (lower) for the first setup (different data sizes, for triplets dataset). Left: first 30 epochs; right: full training process of 500 epochs.

340 kens at node positions to the total number of node
 341 predictions in the dataset and are equal to all nodes
 342 across all sequences, excluding starting points. To-
 343 tal correct predictions also represent MAC.

344 3 Results

345 3.1 Dataset Size Influence

346 Figure 2 illustrates capacity and accuracy trends
 347 across dataset sizes in the first setup. Smaller
 348 datasets learn faster, with accuracy and capacity in-
 349 creasing rapidly within the first 5–6 epochs, reach-
 350 ing maximum capacity by epoch 20. In contrast,
 351 larger datasets show minimal improvement in the
 352 first 15 epochs but exhibit a later inflection point,
 353 leading to higher final accuracy and capacity. This
 354 suggests a threshold existence ($\sim 70,000$ rows
 355 for this case), beyond which the training process
 356 changes and a lot more epochs are required for full
 357 memorization.

358 The final accuracy and capacity (Table 1) indi-
 359 cate that although smaller datasets initially achieve
 360 higher accuracy, their capacity remains well below
 361 the size of the dataset (e.g., 50,000 rows yield only
 362 46,811 samples). In contrast, larger datasets, such
 363 as 100,000 rows, significantly improve memoriza-
 364 tion (86,776 samples), highlighting the model’s
 365 ability to use more data. The progressive capacity
 366 increase suggests that dataset size plays a crucial
 367 role in optimizing memorization; however, the rea-
 368 sons behind the unlearned data, despite available

capacity, remain unclear.

data size	accuracy, %	capacity
50,000	93.62 ± 0.3	$46,811 \pm 149$
60,000	92.42 ± 0.2	$55,455 \pm 126$
70,000	91.1 ± 1.08	$63,773 \pm 756$
80,000	89.63 ± 1.66	$71,706 \pm 1326$
90,000	87.24 ± 1.66	$78,517 \pm 2173$
100,000	86.78 ± 2.42	$86,776 \pm 2484$

Table 1: Final results after the full training process for the first setup (data sizes, for triplets dataset).

370 3.2 Architectural Variations Influences

371 In the second experimental setup, the batch size
 372 was increased from 64 to 128, since larger batch
 373 sizes seem reduce gradient noise and improve mem-
 374 orization. Consequently, the one-layer models in
 375 this setup converged faster and achieved higher
 376 capacity than in the first setup.

377 Softmax consistently outperformed other acti-
 378 vation functions, yielding the highest average ca-
 379 pacity, fewer outliers, and more stable training be-
 380 havior. Notably, four-layer models with Softmax
 381 achieved capacities comparable to one- or two-
 382 layer models without sacrificing convergence speed
 383 (Figure 3), suggesting its scalability with depth.

384 In contrast, ReLU and RReLU showed moderate
 385 performance, but suffered from increased variabil-
 386 ity and decreased capacity as the layers increased,

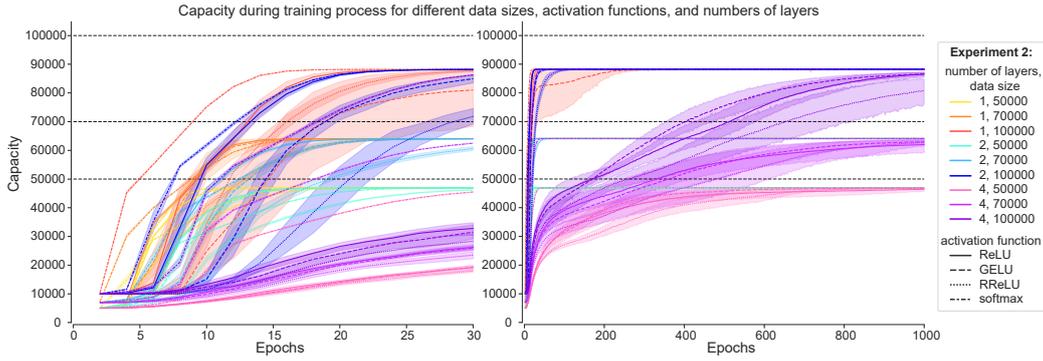


Figure 3: Trends in training capacity for the second setup (different data sizes, activation functions, and numbers of layers for triplets dataset). Left: first 30 epochs; right: full training process of 1000 epochs.

aligning with the findings of Paik and Choi (2023) and Chen and Ge (2024). These activations exhibited inconsistent learning patterns, with unexpected slowdowns in capacity improvements (Fu et al., 2024). GELU followed a similar trend, though it performed better in the early training stages with larger datasets.

As previously, the size of the dataset significantly affected training: larger datasets required longer warm-up phases, initially achieving lower capacities than smaller datasets under the same conditions. This suggests the existence of distinct learning phases where improvements depend on architectural depth, dataset size, and activation function.

Furthermore, adding more layers did not improve performance; instead, it slowed training and reduced final capacity. This is likely due to the simplicity of the dataset (triplets), where additional layers do not provide any advantage in capturing patterns. Although deeper architectures benefit more complex datasets (He et al., 2024), their impact may be reduced for data with simple relationships.

3.3 Numbers of Parameters influence

The third experiment further confirmed that, for simple datasets, learning dynamics depend primarily on embedding size, not the number of layers. Models with the same embedding size but different layer counts exhibited nearly identical accuracy improvement rates. For instance, as shown in Figure 4, a one-layer model with 16 parameters (embedding size is 16, light green) converged at almost the same rate as a two-layer transformer with 32 parameters (embedding size is 16 per layer, dark blue). Similar trends were observed for models with embedding sizes of 32 and 64, regardless of layer count.

These results highlight that embedding size is

the key factor influencing learning speed, while adding layers without increasing embedding size neither accelerates convergence nor improves final capacity. In fact, additional layers often slow the training, as evidenced by the faster growth of accuracy of one-layer models (Figure 4). Smaller embedding sizes further reduced the learning speed, consistent with previous experiments. However, all configurations ultimately reached similar accuracy, highlighting that the simplicity of the dataset allows embedding size to dominate training dynamics.

Final capacity values remained nearly identical across configurations, regardless of embedding size or layer count: with a dataset size of 1,000 samples, capacities for one- and two-layer models were near-accurate. Similarly, at 10,000 and 50,000 samples, one-layer models achieved $9,874 \pm 11$ and $46,939 \pm 105$, while two-layer models reached $9,875 \pm 7$ and $46,911 \pm 117$, respectively. However, at 100,000 samples, a capacity "barrier" emerged. Two-layer transformers with an embedding size of 8 (16 total parameters) showed the capacity drop to $85,935 \pm 153$, compared to $\sim 88,200$ for other configurations, while one-layer models maintained a higher capacity of $88,240 \pm 62$. This suggests that larger datasets, smaller embeddings, and deeper architectures may introduce limitations due to slower convergence or suboptimal capacity utilization.

Evaluating consistency with the 2 bits per parameter rule (Allen-Zhu and Li, 2024b) was challenging due to dataset size limitations. While most configurations achieved similar capacities, the drop in the two-layer model with 16 parameters likely reflects incomplete convergence, possibly caused by slower learning dynamics in deeper models.

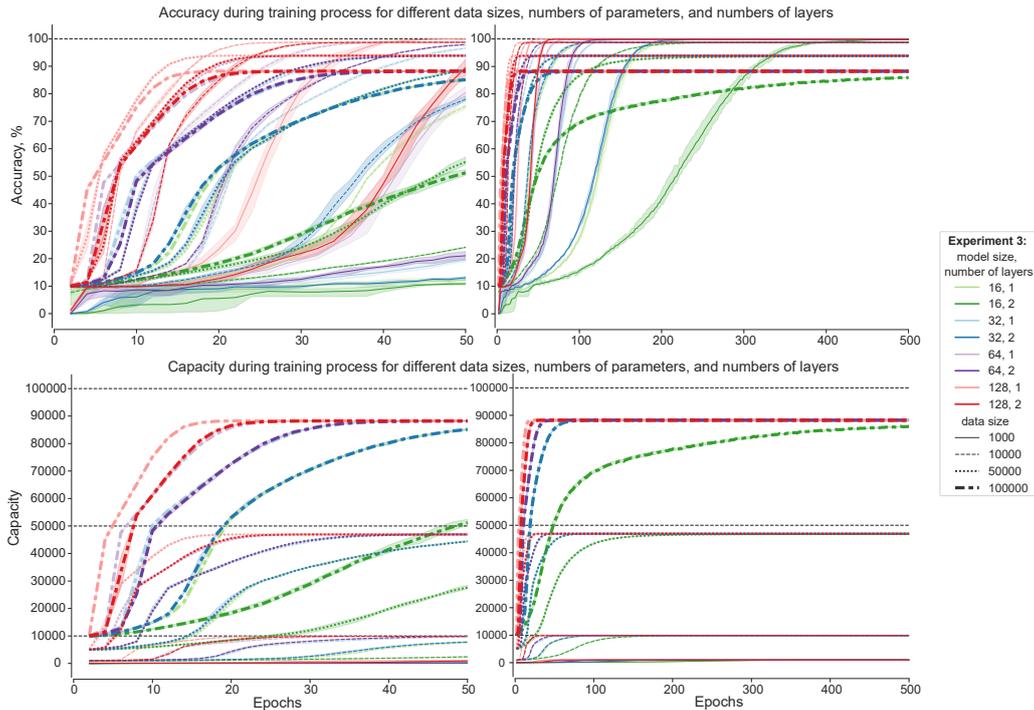


Figure 4: Trends in training accuracy (upper) and capacity (lower) for the third setup (different data sizes, numbers of parameters, and numbers of layers for triplets dataset). Left: first 50 epochs; right: full training process of 500 epochs. Light color corresponds to 1 layer, dark – to 2; number of parameters is a total number for all layers: green – 16, blue – 32, violet – 64, red – 128; embedding size can be computed by dividing it by layer count.

3.4 Insights from Sequence Datasets

In the fourth setup, model capacity was evaluated by testing its ability to memorize each node in a sequence using the full preceding sequence of nodes and edges (instead of triplets). This required multiple predictions per sequence: 34, 908, 85, 972, and 167, 965 for datasets containing 20, 50, and 100 thousand sequences, respectively.

Compared to triplet datasets, models trained on sequences achieved near-perfect memorization in significantly fewer epochs, with most configurations plateauing within 150 epochs (Figure 5). The sequential structure likely facilitated more efficient learning, though it also increased training time due to the higher information per sequence. Training exhibited greater capacity fluctuations across epochs, likely reflecting the dataset’s increased complexity, as sequences encode more intricate patterns than triplets. Nonetheless, models demonstrated exceptional memorization, achieving 100% capacity for the 20 thousand sequence dataset and over 99.5% for 50 and 100 thousand sequences.

Regarding activation functions: as previously, RReLU converged more slowly than Softmax, though final capacities were nearly identical for one- and two-layer models: with 100 thousand se-

quences, RReLU achieved $166,934 \pm 243$ (one layer) and $166,995 \pm 118$ (two layers), while Softmax reached $166,992 \pm 110$ and $166,985 \pm 904$, respectively. In deeper models (4 layers), RReLU showed lower final capacities and greater fluctuations ($165,271 \pm 1,068$ vs. $166,825 \pm 319$ for Softmax). This contrasts with previous findings (Shen et al., 2023), which suggest ReLU outperforms Softmax. The discrepancy may indicate that activation function effectiveness may vary based on dataset structure and task, therefore it needs further investigation. Despite the increased complexity of sequence datasets, models adapted quickly and demonstrated strong memorization performance.

4 Discussion

This study provided insights into the memorization capacity of transformer models trained on real-world structured datasets. Returning to the medical domain of SNOMED, the original KG contains over a million relations, integrating diverse fields of medicine (e.g., substances, diseases, anatomical structures). However, in mobile applications, e.g. LLMs in smart glasses or smartwatches, models must efficiently retain only specific subsets of information. For instance, a cardiac surgeon’s smart

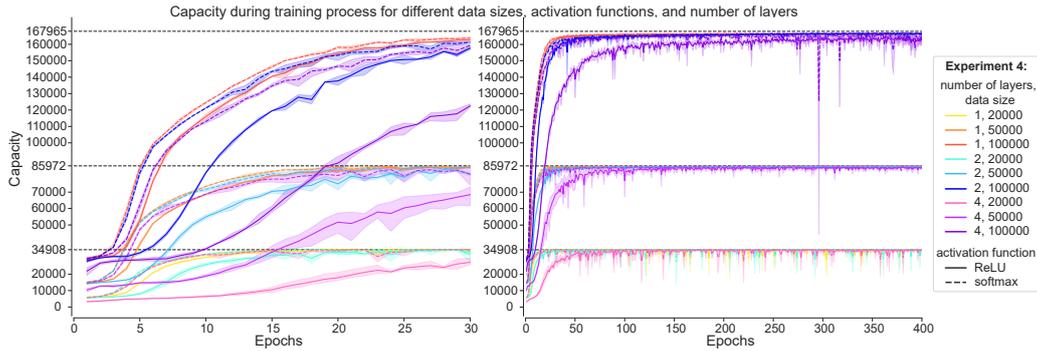


Figure 5: Trends in training capacity for the fourth setup (different data sizes, activation functions, and numbers of layers for sequences dataset). Left: first 30 epochs; right: full training process of 400 epochs.

510 glasses would require an LLM specialized in cardi- 546
 511 ology, while a smartwatch may store personalized 547
 512 health data for a specific user, limiting the dataset 548
 513 size to 10–100 thousand triplets or sequences. 549

514 This research offers insights into efficient train- 550
 515 ing strategies for such models, analyzing how 551
 516 dataset characteristics and architectural choices im- 552
 517 pact convergence speed and capacity utilization. 553

518 4.1 Effect of Dataset Structure 554

519 Smaller datasets led to faster convergence but lower 555
 520 capacity, while larger datasets required longer 556
 521 warm-up periods but improved retention. Beyond 557
 522 a certain size, training slowed significantly, indicat- 558
 523 ing optimization bottlenecks. 559

524 Sequence-based datasets outperformed triplets, 560
 525 achieving near-perfect memorization with fewer 561
 526 epochs. Sequences aided and complicated learn- 562
 527 ing, reinforcing relationships between data but also 563
 528 introducing greater training fluctuations, aligning 564
 529 with Ju et al. (2021). This suggests that longer 565
 530 traversal sequences could further improve memo- 566
 531 rization in domain-specific medical applications. 567

532 4.2 Architectural Influence 568

533 Embedding size was the key factor in learning 569
 534 speed and capacity, while adding layers provided 570
 535 little benefit and sometimes reduced performance, 571
 536 likely due to dataset simplicity. This aligns with 572
 537 He et al. (2024), who found that many transformer 573
 538 layers exhibit high similarity and, sometimes, re- 574
 539 dundancy and can be pruned without performance 575
 540 loss, reducing computational overhead. 576

541 For larger datasets, smaller embeddings strug- 577
 542 gled to reach full capacity, particularly in deeper 578
 543 architectures, suggesting that increasing embed- 579
 544 ding size is more beneficial than adding depth, at 580
 545 least for structured, domain-specific memorization.

546 Softmax led to greater stability and capacity, 547
 548 while ReLU-based activations showed higher vari- 549
 550 ability and performance drops in deeper models, 551
 552 aligning with Paik and Choi (2023); Chen and Ge 553
 554 (2024). However, this contrasts with prior work 555
 556 by Shen et al. (2023), emphasizing that activation 557
 558 effectiveness is highly dependent on the task and 559
 560 dataset structure. 561

562 5 Conclusions 568

563 This study explored the memorization capacity of 569
 564 transformer models on structured datasets from 570
 565 SNOMED KG, analyzing how architecture and 571
 566 dataset structure affect learning efficiency and ca- 572
 567 pacity retention. 573

568 Key findings show that embedding size and ac- 574
 569 tivation function were more influential than depth, 575
 570 while larger datasets improved memorization but 576
 571 required longer training. Triplets performed well 577
 572 in simpler models, whereas sequences excelled but 578
 573 introduced fluctuations. Challenges remain in effi- 579
 574 ciency, layer-specific contributions, and generaliza- 580
 575 tion, necessitating further research on scalability, 581
 576 compression, and architecture optimization. 582

577 For real-world applications, such as LLMs in 583
 578 medical smart devices, models must efficiently 584
 579 store specialized knowledge while maintaining 585
 580 computational feasibility. Future work should ex- 586
 581 plore longer sequences, adaptive memory compres- 587
 582 sion, and layer-wise analysis to enhance structured 588
 583 knowledge retention in practical deployments. 589

590 6 Limitations 596

591 While this study provides meaningful insights, sev- 597
 592 eral open questions remain: 598

- 599 • It is unclear why certain samples remain un- 600
 601 learned within the same model architecture 602

581	despite available capacity. Future research	Carter, Tom Henighan, and Chris Olah. 2023. Towards monosemanticity: Decomposing language models with dictionary learning . <i>Transformer Circuits Thread</i> .	630
582	should explore optimization strategies to im-		631
583	prove memorization efficiency.		632
584	• Future research should test these findings and	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie	634
585	hypotheses on longer sequences and larger	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	635
586	datasets to confirm them at scale.	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	636
587	• The specific role of each layer in memoriza-	Askeff, Sandhini Agarwal, Ariel Herbert-Voss,	637
588	tion was not investigated, missing insights	Gretchen Krueger, Tom Henighan, Rewon Child,	638
589	from probing methods as suggested in Allen-	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,	639
590	Zhu and Li (2024a) . Future studies could	Clemens Winter, Christopher Hesse, Mark Chen,	640
591	apply probing techniques to analyze layer-	Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin	641
592	specific role in memorization capacity.	Chess, Jack Clark, Christopher Berner, Sam Mc-	642
593	• Additionally, sparse autoencoders (Bricken	Candlish, Alec Radford, Ilya Sutskever, and Dario	643
594	et al., 2023) or transcoders (Paulo et al., 2025)	Amodei. 2020. Language models are few-shot learners . <i>Preprint</i> , arXiv:2005.14165.	644
595	could be integrated into transformer layers to		645
596	distinguish memorization from generalization,	Wenlin Chen and Hong Ge. 2024. Neural characteristic	646
597	helping determine whether certain layers store	activation analysis and geometric parameterization	647
598	specific relationships or contribute to broader	for relu networks . <i>Preprint</i> , arXiv:2305.15912.	648
599	model generalizability.		
600	By addressing these limitations, future work can	Shaker El-Sappagh, Francesco Franda, Farman Ali, and	649
601	further refine transformer optimization strategies	Kyung-Sup Kwak. 2018. Snomed ct standard ontology	650
602	for structured data modeling and knowledge reten-	based on the ontology for general medical science . <i>BMC Medical Informatics and Decision</i>	651
603	tion.	Making , 18(1):76.	652
604			653
605	References	Jingwen Fu, Tao Yang, Yuwang Wang, Yan Lu, and	654
606	Abien Fred Agarap. 2019. Deep learning using rectified	Nanning Zheng. 2024. Breaking through the learning	655
607	linear units (relu) . <i>Preprint</i> , arXiv:1803.08375.	plateaus of in-context learning in transformer .	656
608	Zeyuan Allen-Zhu and Yuanzhi Li. 2024a. Physics of	<i>Preprint</i> , arXiv:2309.06054.	657
609	language models: Part 3.1, knowledge storage and		
610	extraction . <i>Preprint</i> , arXiv:2309.14316.	Bhumika Gupta, Pralaypati Ta, Keerthi Ram, and Mo-	658
611	Zeyuan Allen-Zhu and Yuanzhi Li. 2024b. Physics	hanasankar Sivaprakasam. 2024. Comprehensive	659
612	of language models: Part 3.3, knowledge capacity	Modeling and Question Answering of Cancer Clinical	660
613	scaling laws . <i>Preprint</i> , arXiv:2404.05405.	Practice Guidelines using LLMs . In <i>2024</i>	661
614	Simone Balloccu, Ehud Reiter, Vivek Kumar, Diego Re-	<i>IEEE Conference on Computational Intelligence in</i>	662
615	forgiato Recupero, and Daniele Riboni. 2024. Ask	<i>Bioinformatics and Computational Biology (CIBCB)</i> ,	663
616	the experts: sourcing high-quality datasets for nutri-	pages 1–8. ISSN: 2994-9408.	664
617	tional counselling through Human-AI collaboration .		
618	<i>arXiv preprint</i> . ArXiv:2401.08420 [cs].	Shwai He, Guoheng Sun, Zheyu Shen, and Ang Li.	665
619	Ludwig Boltzmann. 1868. Studien über das gle-	2024. What matters in transformers? not all attention	666
620	ichgewicht der lebendigen kraft zwischen bewegten	is needed . <i>Preprint</i> , arXiv:2406.15786.	667
621	materiellen punkten. <i>Wiener Berichte</i> , 58:517–560.		
622	Studies on the balance of living force between mov-	Dan Hendrycks and Kevin Gimpel. 2023. Gaussian er-	668
623	ing material points.	ror linear units (gelus) . <i>Preprint</i> , arXiv:1606.08415.	669
624	Trenton Bricken, Adly Templeton, Joshua Batson, Brian	Aki Härmä, Marcin Pietrasik, and Anna Wilbik. 2024.	670
625	Chen, Adam Jermyn, Tom Conerly, Nicholas L.	Empirical capacity model for self-attention neural	671
626	Turner, Cem Anil, Carson Denison, Amanda Askeff,	networks . <i>Preprint</i> , arXiv:2407.15425.	672
627	Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas	Yue Ju, Alka Isac, and Yimin Nie. 2021. Chunkformer:	673
628	Schiefer, Tim Maxwell, Nicholas Joseph, Zach	Learning long time series with multi-stage chunked	674
629	Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Bray-	transformer . <i>Preprint</i> , arXiv:2112.15087.	675
	den McLean, Josiah E. Burke, Tristan Hume, Shan		
		Tokio Kajitsuka and Issei Sato. 2024. Optimal	676
		memorization capacity of transformers . <i>Preprint</i> ,	677
		arXiv:2409.17677.	678
		Junghwan Kim, Michelle Kim, and Barzan Mozafari.	679
		2023. Provable memorization capacity of transform-	680
		ers . In <i>The Eleventh International Conference on</i>	681
		<i>Learning Representations</i> .	682

683 Diederik P. Kingma and Jimmy Ba. 2017. [Adam: A method for stochastic optimization](#). *Preprint*, arXiv:1412.6980.

684
685

686 Jean-Baptiste Lamy. 2017. [Owlready: Ontology-oriented programming in python with automatic classification and high level constructs for biomedical ontologies](#). *Artificial Intelligence in Medicine*, 80:11–28.

687
688
689
690

691 Sadegh Mahdavi, Renjie Liao, and Christos Thrampoulidis. 2024. [Memorization capacity of multi-head attention in transformers](#). *Preprint*, arXiv:2306.02010.

692
693
694

695 Inyoung Paik and Jaesik Choi. 2023. [The disharmony between bn and relu causes gradient explosion, but is offset by the correlation between activations](#). *Preprint*, arXiv:2304.11692.

696
697
698

699 Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in pytorch.

700
701
702

703 Gonalo Paulo, Stepan Shabalin, and Nora Belrose. 2025. [Transcoders beat sparse autoencoders for interpretability](#). *Preprint*, arXiv:2501.18823.

704
705

706 Ahsan Shehzad, Feng Xia, Shagufta Abid, Ciyuan Peng, Shuo Yu, Dongyu Zhang, and Karin Verspoor. 2024. [Graph transformers: A survey](#). *Preprint*, arXiv:2407.09777.

707
708
709

710 Kai Shen, Junliang Guo, Xu Tan, Siliang Tang, Rui Wang, and Jiang Bian. 2023. [A study on relu and softmax in transformer](#). *Preprint*, arXiv:2302.06461.

711
712

713 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, R mi Louf, Morgan Funtowicz, and Jamie Brew. 2019. [Huggingface’s transformers: State-of-the-art natural language processing](#). *CoRR*, abs/1910.03771.

714
715
716
717
718

719 Jinlin Wu, Xusheng Liang, Xuexue Bai, and Zhen Chen. 2024. [SurgBox: Agent-Driven Operating Room Sandbox with Surgery Copilot](#). *arXiv preprint*. ArXiv:2412.05187 [cs].

720
721
722

723 Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. 2015. [Empirical evaluation of rectified activations in convolutional network](#). *Preprint*, arXiv:1505.00853.

724
725

A Appendix: Additional Representations of the Results

726
727

This appendix provides supplementary visualizations and tables for the experiments conducted:

728
729

- Second experiment: 730
 - Figure 6: Accuracy trends during training. 731
 - Table 2: Final capacities. 732
 - Third experiment: 734
 - Table 3: Final capacities. 735
 - Fourth experiment: 736
 - Figure 7: Accuracy trends during training. 737
 - Table 4: Final capacities. 738
- 739

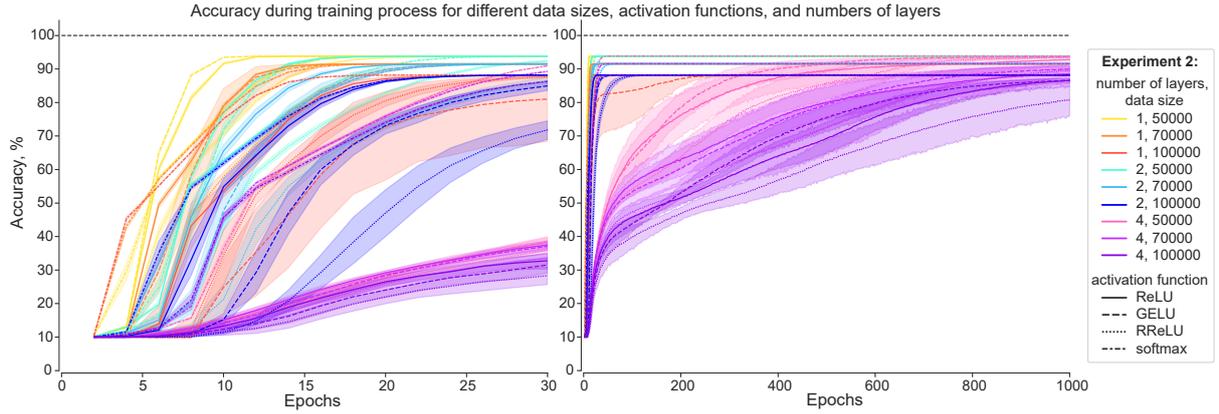


Figure 6: Trends in training accuracy for the second setup (different data sizes, activation functions, and numbers of layers for triplets dataset). Left: first 30 epochs; right: full training process of 1000 epochs.

activation function	layers count	data sizes		
		50,000	70,000	100,000
ReLU	1	46,898 ± 158	64,091 ± 192	88,148 ± 312
	2	46,920 ± 112	64,086 ± 130	88,217 ± 125
	4	46,391 ± 2,268	61,931 ± 8,480	86,558 ± 3,291
GELU	1	46,925 ± 105	64,096 ± 184	88,195 ± 123
	2	46,926 ± 115	64,080 ± 120	88,215 ± 128
	4	46,798 ± 156	62,949 ± 1,906	86,589 ± 2,202
RReLU	1	46,930 ± 125	64,080 ± 122	88,180 ± 180
	2	46,927 ± 121	64,088 ± 117	88,208 ± 132
	4	46,730 ± 223	62,818 ± 3,680	80,755 ± 15,844
softmax	1	46,924 ± 87	64,082 ± 166	88,211 ± 192
	2	46,908 ± 127	64,074 ± 134	88,213 ± 171
	4	46,923 ± 104	64,085 ± 131	88,197 ± 134
all	1	46,919 ± 119	64,087 ± 162	88,183 ± 210
	2	46,920 ± 115	64,082 ± 121	88,213 ± 135
	4	46,710 ± 1169	62,945 ± 4,92	85,525 ± 9,720

Table 2: Final capacity after the full training process for the second setup (different numbers of layers, data sizes, and activation functions for triplets dataset).

embedding parameters	layers count	data sizes			
		1,000	10,000	50,000	100,000
16	1	1,000 ± 1	9,870 ± 10	46,937 ± 148	88,236 ± 74
	2	998 ± 3	9,875 ± 4	46,858 ± 93	85,935 ± 153
32	1	998 ± 3	9,872 ± 11	46,955 ± 119	88,234 ± 62
	2	999 ± 3	9,876 ± 9	46,927 ± 128	88,252 ± 82
64	1	999 ± 2	9,878 ± 9	46,932 ± 122	88,242 ± 102
	2	999 ± 3	9,876 ± 7	46,919 ± 96	88,237 ± 58
128	1	999 ± 2	9,877 ± 12	46,930 ± 85	88,248 ± 29
	2	999 ± 3	9,872 ± 6	46,938 ± 131	88,214 ± 53
all	1	999 ± 2	9,874 ± 11	46,939 ± 105	88,240 ± 62
	2	999 ± 3	9,875 ± 7	46,911 ± 117	87,660 ± 2,082

Table 3: Final capacity after the full training process for the third setup (different data sizes, numbers of parameters, and numbers of layers for triplets dataset).

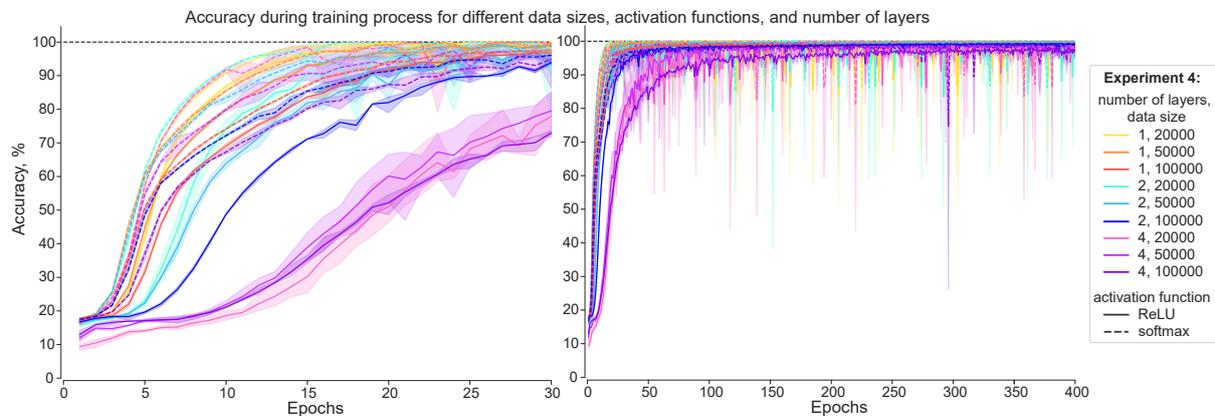


Figure 7: Trends in training accuracy for the fourth setup (different data sizes, activation functions, and numbers of layers for sequences dataset). Left: first 30 epochs; right: full training process of 400 epochs.

activation function	layers count	# of sequences (# of predictions)		
		20,000 (34,908)	50,000 (85,972)	100,000 (167,965)
RReLU	1	34,908 ± 0	85,936 ± 31	166,934 ± 243
	2	34,908 ± 0	85,917 ± 34	166,995 ± 118
	4	34,908 ± 0	85,647 ± 270	165,271 ± 1,068
softmax	1	34,908 ± 0	85,931 ± 18	166,992 ± 110
	2	34,908 ± 0	85,888 ± 33	166,985 ± 904
	4	34,908 ± 0	85,771 ± 42	166,825 ± 319
all	1	34,908 ± 0	85,934 ± 23	166,963 ± 180
	2	34,908 ± 0	85,903 ± 44	166,990 ± 577
	4	34,908 ± 0	85,709 ± 220	166,048 ± 1,842

Table 4: Final capacity after the full training process for the fourth setup (different data sizes, activation functions, and numbers of layers for sequences dataset).