

---

# Beyond Attention: Breaking the Limits of Transformer Context Length with Recurrent Memory

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 A major limitation for the broader scope of problems solvable by transformers is  
2 the quadratic scaling of computational complexity with input size. In this study, we  
3 investigate the recurrent memory augmentation of pre-trained transformer models  
4 to extend input context length while linearly scaling compute. Our approach  
5 demonstrates the capability to store information in memory for sequences of up  
6 to an unprecedented two million tokens while maintaining high retrieval accuracy.  
7 Experiments with language modeling tasks show perplexity improvement as the  
8 number of processed input segments increases. These results underscore the  
9 effectiveness of our method, which has significant potential to enhance long-term  
10 dependency handling in natural language understanding and generation tasks, as  
11 well as enable large-scale context processing for memory-intensive applications.

## 12 1 Introduction

13 Transformer-based models show their effectiveness across multiple domains and tasks. The self-  
14 attention allows to combine information from all sequence elements into context-aware represen-  
15 tations. However, global and local information has to be stored mostly in the same element-wise  
16 representations. Moreover, the length of an input sequence is limited by quadratic computational  
17 complexity of self-attention. In this work, we propose and study a memory-augmented segment-level  
18 recurrent Transformer (Recurrent Memory Transformer). Memory allows to store and process local  
19 and global information as well as to pass information between segments of the long sequence with  
20 the help of recurrence. We implement a memory mechanism with no changes to Transformer model  
21 by adding special memory tokens to the input or output sequence. Then Transformer is trained to  
22 control both memory operations and sequence representations processing.

23 This study we show that by using simple token-based memory mechanism introduced in [Bulatov  
24 et al., 2022] can be combined with pretrained transformer models like BERT [Devlin et al., 2019]  
25 and GPT-2 [Radford et al., 2019] with full attention and full precision operations.

### 26 Contributions

27 1. We enhance both encoder-only and decoder-only pre-trained Transformer language models by  
28 incorporating token-based memory storage and segment-level recurrence with recurrent memory  
29 (RMT).

30 2. We demonstrate that language models pre-trained on much shorter lengths can be trained with  
31 RMT approach to tackle tasks on sequences many times longer than its originally designed input  
32 length.

33 3. We discovered the trained RMT’s capacity to successfully extrapolate to tasks of varying lengths,  
34 including those exceeding 1 million tokens with linear scaling of computations required.

35 4. Through attention pattern analysis, we found the operations RMT employs with memory, enabling  
36 its success in handling exceptionally long sequences.

## 37 **2 Related work**

38 Our work revolves around the concept of memory in neural architectures. Memory has been a  
39 recurrent theme in neural network research, dating back to early works [McCulloch and Pitts, 1943,  
40 Stephen, 1956] and significantly advancing in the 1990s with the introduction of the *Backpropagation*  
41 *Through Time* learning algorithm [Werbos, 1990] and *Long-Short Term Memory* (LSTM) neural  
42 architecture [Hochreiter and Schmidhuber, 1997]. Contemporary memory-augmented neural net-  
43 works (MANNs) typically utilize some form of recurrent external memory separate from the model’s  
44 parameters. *Neural Turing Machines* (NTMs) [Graves et al., 2014] and *Memory Networks* [Weston  
45 et al., 2015] are equipped with storage for vector representations accessible through an attention  
46 mechanism. Memory Networks [Weston et al., 2015, Sukhbaatar et al., 2015] were designed to enable  
47 reasoning through sequential attention over memory content.

48 NTMs, followed by *Differentiable Neural Computer* (DNC) [Graves et al., 2016] and *Sparse DNC*  
49 [Rae et al., 2016], are implemented as recurrent neural networks capable of writing to memory  
50 storage over time. All these models are differentiable and trainable via backpropagation through  
51 time (BPTT). Parallel research lines extend recurrent neural networks, such as LSTM, with data  
52 structures like stacks, lists, or queues [Joulin and Mikolov, 2015, Grefenstette et al., 2015]. MANN  
53 architectures with more advanced addressing mechanisms, such as address-content separation and  
54 multi-step addressing, have been proposed in [Gulcehre et al., 2016, 2017, Meng and Rumshisky,  
55 2018]. The Global Context Layer model [Meng and Rumshisky, 2018] employs address-content  
56 separation to address the challenge of training content-based addressing in canonical NTMs.

57 Memory is often combined with Transformers in a recurrent approach. Long inputs are divided into  
58 smaller segments, processed sequentially with memory to access information from past segments.  
59 Transformer-XL [Dai et al., 2019] preserves previous hidden states for reuse in subsequent segments,  
60 while Compressive Transformer [Rae et al., 2020] adds new compressed memory. Ernie-Doc [Ding  
61 et al., 2021] enhances contextual information flow by employing same-layer recurrence instead of  
62 attending to previous layer outputs of preceding segments. Memformer [Wu et al., 2022a] introduces  
63 a dedicated memory module to store previous hidden states in summarized representations. Using a  
64 similar approach to Memformer, MART [Lei et al., 2020] and Block-Recurrent Transformer [Hutchins  
65 et al., 2022] adopt memory update rules analogous to LSTM [Hochreiter and Schmidhuber, 1997]  
66 and GRU [Cho et al., 2014]. FeedBack Transformer [Fan et al., 2020] implements full recurrence  
67 beyond the segment level and merges low and high layers representations into a memory state.

68 A drawback of most existing recurrent methods is the need for architectural modifications that  
69 complicate their application to various pre-trained models. In contrast, the Recurrent Memory  
70 Transformer can be built upon any model that uses a common supported interface.

71 Some approaches redesign the self-attention mechanism to reduce computational complexity while  
72 minimizing input coverage loss. *Star-Transformer* [Guo et al., 2019], *Longformer* [Beltagy et al.,  
73 2020], *GMAT* [Gupta and Berant, 2020], *Extended Transformer Construction* (ETC) [Ainslie et al.,  
74 2020], and *Big Bird* [Zaheer et al., 2020] limit attention distance and employ techniques such as  
75 global representations to preserve long-range dependencies. *Memory Transformer* [Burtsev et al.,  
76 2020] introduces memory by extending the unchanged model input with special memory tokens.

77 A common constraint of these methods is that memory requirements grow with input size during  
78 both training and inference, inevitably limiting input scaling due to hardware constraints. The longest  
79 Longformer, Big Bird, and Long T5 [Guo et al., 2022] models reported in their respective papers  
80 have a maximum length of less than 33,000 tokens. CoLT5 [Ainslie et al., 2023] can handle up to  
81 64,000 tokens before running out of memory, and Memorizing Transformers [Wu et al., 2022b] and  
82 Unlimiformer [Bertsch et al., 2023] further extend memory through k-NN.

## 83 **3 Recurrent Memory Transformer**

84 Starting from the initial Recurrent Memory Transformer [Bulatov et al., 2022] (RMT), we adapted it  
85 for a plug-and-play approach as a wrapper for a range of popular Transformers.

86 This adaptation augments its backbone with  
 87 memory, composed of  $m$  real-valued trainable  
 88 vectors (Figure 1). The lengthy input is divided  
 89 into segments, and memory vectors are  
 90 prepended to the first segment embeddings and  
 91 processed alongside the segment tokens. For  
 92 encoder-only models like BERT, memory is  
 93 added only once at the beginning of the segment,  
 94 unlike [Bulatov et al., 2022], where decoder-  
 95 only models separate memory into read and  
 96 write sections. For the time step  $\tau$  and segment  
 97  $H_\tau^0$ , the recurrent step is performed as follows:

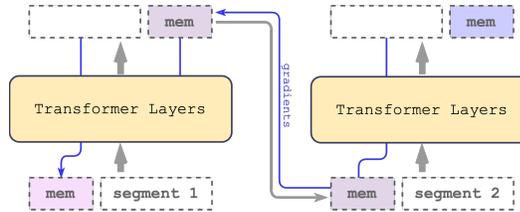


Figure 1: **Recurrent memory mechanism.** Memory is passed to Transformer along input sequence embeddings, and memory output is passed to the next segment. During training gradients flow from the current segment through memory to the previous segment.

$$\tilde{H}_\tau^0 = [H_\tau^{mem} \circ H_\tau^0], \bar{H}_\tau^N = \text{Transformer}(\tilde{H}_\tau^0), [\bar{H}_\tau^{mem} \circ H_\tau^N] := \bar{H}_\tau^N,$$

98 here  $N$  is a number of Transformer layers.

99 After the forward pass,  $\bar{H}_\tau^{mem}$  contains updated memory tokens for the segment  $\tau$ .

100 Segments of the input sequence are processed sequentially. To enable the recurrent connection, we  
 101 pass the outputs of the memory tokens from the current segment to the input of the next one:

$$H_{\tau+1}^{mem} := \bar{H}_\tau^{mem}, \tilde{H}_{\tau+1}^0 = [H_{\tau+1}^{mem} \circ H_{\tau+1}^0].$$

102 Both memory and recurrence in the RMT are based only on global memory tokens. This allows the  
 103 backbone Transformer to remain unchanged, making the RMT memory augmentation compatible  
 104 with any model from the Transformer family.

### 105 3.1 Computational efficiency

106 We can estimate the required FLOPs for RMT and Transformer models of different sizes and sequence  
 107 lengths. We took configurations (vocabulary size, number of layers, hidden size, intermediate hidden  
 108 size, and number of attention heads) for the OPT model family [Zhang et al., 2022] and computed the  
 109 number of FLOPs for the forward pass following [Hoffmann et al., 2022]. We also modified FLOP  
 110 estimates to account for the effect of RMT recurrence.

111 Figure 2 shows that RMT scales linearly for any model size if the segment length is fixed. We achieve  
 112 linear scaling by dividing an input sequence into segments and computing the full attention matrix

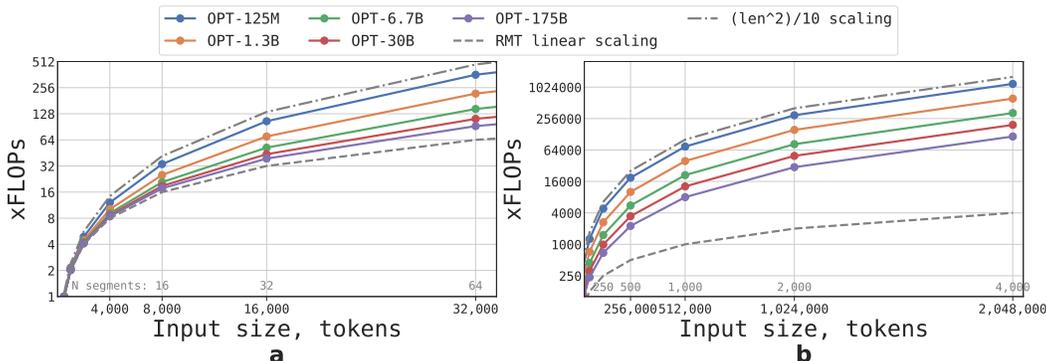


Figure 2: **RMT inference scales linearly with respect to the input sequence length.** We estimate the required FLOP increase for the forward pass compared to running models on sequences with 512 tokens. **a:** lengths from 512 to 32,000 tokens, **b:** lengths from 32,000 to 2,048,000 tokens. The RMT segment length is fixed at 512 tokens. While larger models (OPT-30B, OPT-175B) tend to exhibit near-linear scaling on relatively short sequences up to 32,000, they reach quadratic scaling on longer sequences. Smaller models (OPT-125M, OPT-1.3B) demonstrate quadratic scaling even on shorter sequences. On sequences with 2,048,000 tokens, RMT can run OPT-175B with  $\times 29$  fewer FLOPs and with  $\times 295$  fewer FLOPs than OPT-135M.

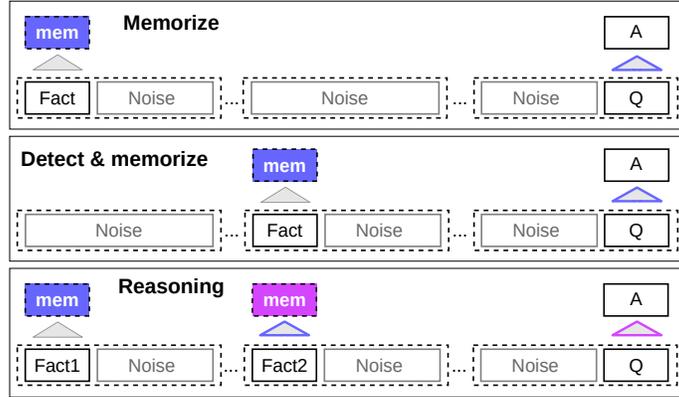


Figure 3: **Memory-intensive synthetic tasks.** Synthetic tasks and the required RMT operations to solve them are presented. In the Memorize task, a fact statement is placed at the start of the sequence. In the Detect and Memorize task, a fact is randomly placed within a text sequence, making its detection more challenging. In the Reasoning task, two facts required to provide an answer are randomly placed within the text. For all tasks, the question is at the end of the sequence. 'mem' denotes memory tokens, 'Q' represents the question, and 'A' signifies the answer.

113 only within segment boundaries. Larger Transformer models tend to exhibit slower quadratic scaling  
 114 with respect to sequence length because of compute-heavy FFN layers (which scale quadratically  
 115 with respect to hidden size). However, on extremely long sequences  $> 32,000$ , they fall back to  
 116 quadratic scaling. RMT requires fewer FLOPs than non-recurrent models for sequences with more  
 117 than one segment ( $> 512$  in this study) and can reduce the number of FLOPs by up to  $\times 295$  times.  
 118 RMT provides a larger relative reduction in FLOPs for smaller models, but in absolute numbers, a  
 119  $\times 29$  times reduction for OPT-175B models is highly significant.

## 120 4 Memorization Tasks

121 To test memorization abilities, we constructed synthetic datasets that require memorization of simple  
 122 facts and basic reasoning. The task input consists of one or several facts and a question that can be  
 123 answered only by using all of these facts. To increase the task difficulty, we added natural language  
 124 text unrelated to the questions or answers. This text acts as noise, so the model's task is to separate  
 125 facts from irrelevant text and use them to answer the questions. The task is formulated as a 6-class  
 126 classification, with each class representing a separate answer option.

127 Facts are generated using the bAbI dataset [Weston et al., 2016], while the background text is sourced  
 128 from questions in the QuALITY [Pang et al., 2022] long QA dataset.

129 Background text: ... He was a big man, broad-shouldered and still thin-waisted.  
 130 Eddie found it easy to believe the stories he had heard about his father ...

131 The first task tests the ability of RMT to write and store information in memory for an extended time  
 132 (Figure 3, top). In the simplest case, the fact is always located at the beginning of the input, and  
 133 the question is always at the end. The amount of irrelevant text between the question and answer is  
 134 gradually increased, so that the entire input does not fit into a single model input.

135 Fact: Daniel went back to the hallway.  
 136 Question: Where is Daniel?  
 137 Answer: hallway

138 Fact detection increases the task difficulty by moving the fact to a random position in the input  
 139 (Figure 3, middle). This requires the model to first distinguish the fact from irrelevant text, write it to  
 140 memory, and later use it to answer the question located at the end.

141 Another important operation with memory is being able to operate with several facts and current  
 142 context. To evaluate this function, we use a more complicated task, called "reasoning", where two

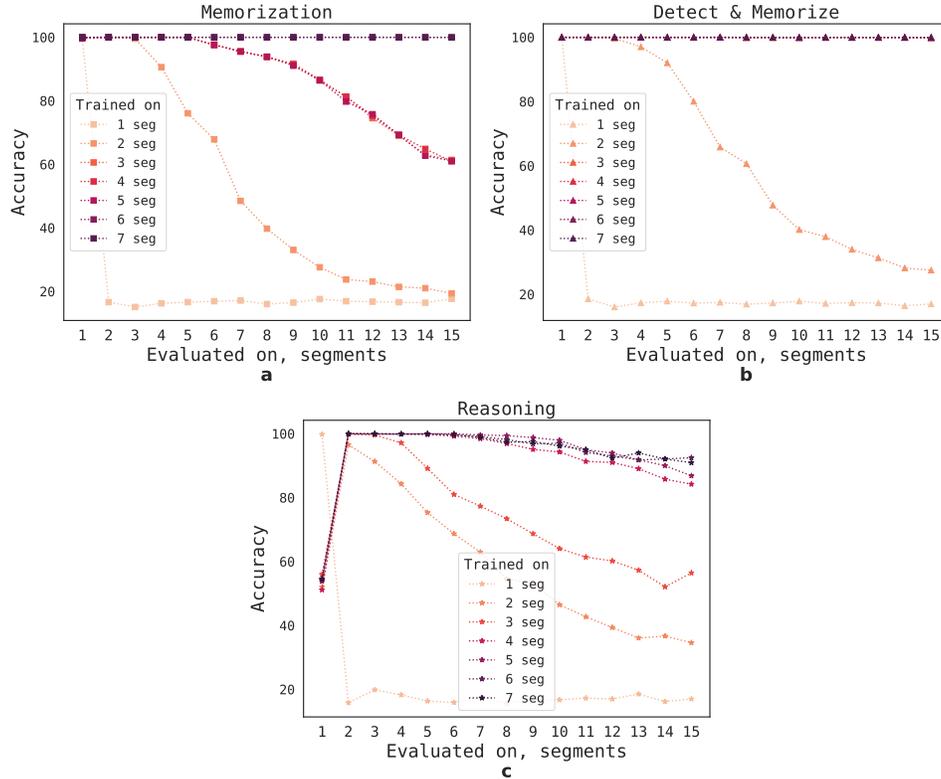


Figure 4: **Generalization of memory retrieval.** Evaluation of checkpoints trained on 1-7 segment tasks with memory size 10 on varying input lengths. **a:** Memorization task, **b:** Detection & memorization, **c:** Reasoning. Models trained on more than 5 segments generalize well on longer tasks.

143 facts are generated and positioned randomly within the input sequence (Figure 3, bottom). The  
 144 question posed at the end of the sequence is formulated in a way that any of the facts must be used to  
 145 answer the question correctly (i.e., the *Two Argument Relation* bAbI task).

146 Fact1: The hallway is east of the bathroom.  
 147 Fact2: The bedroom is west of the bathroom.  
 148 Question: What is the bathroom east of?  
 149 Answer: bedroom

## 150 5 Learning Memory Operations

151 We use the pretrained models from Hugging Face Transformers [Wolf et al., 2020] as backbones for  
 152 RMT in our experiments. All models are augmented with memory and trained using the AdamW  
 153 optimizer [Loshchilov and Hutter, 2019] with linear learning rate scheduling and warmup. Technical  
 154 details of training and full set of hyperparameters will be available in the Appendix and training  
 155 scripts in the GitHub repository. Memorization task experiments were conducted using 4-8 Nvidia  
 156 1080ti GPUs. For longer sequences, we speed up evaluation by switching to a single 40GB Nvidia  
 157 A100.

### 158 5.1 Curriculum Learning

159 We observe that using a training schedule greatly improves solution accuracy and stability. Initially,  
 160 RMT is trained on shorter versions of the task, and upon training convergence, the task length is  
 161 increased by adding one more segment. The curriculum learning process continues until the desired  
 162 input length is reached.

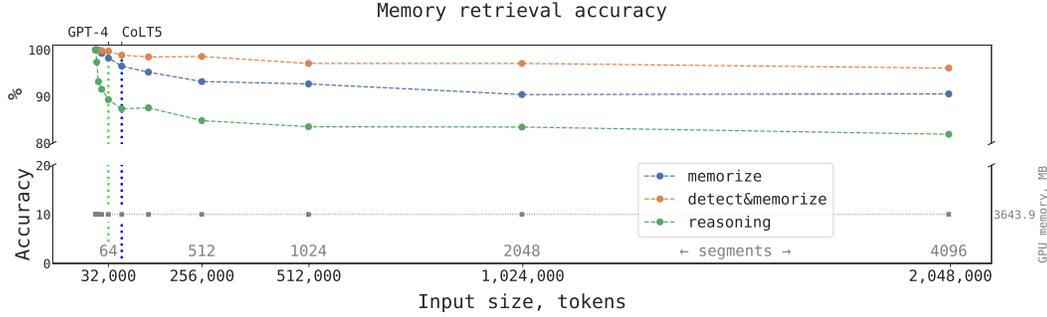


Figure 5: **Recurrent Memory Transformer retains information across up to  $2 \times 10^6$  tokens.** By augmenting a pre-trained BERT model with recurrent memory [Bulatov et al., 2022], we enabled it to store task-specific information across 7 segments of 512 tokens each. During inference, the model effectively utilized memory for up to 4,096 segments with a total length of 2,048,000 tokens—significantly exceeding the largest input size reported for transformer models (64K tokens for CoLT5 [Ainslie et al., 2023], and 32K tokens for GPT-4 [OpenAI, 2023], and 100K tokens for Claude). This augmentation maintains the base model’s memory size at 3.6 GB in our experiments.

163 In our experiments, we begin with sequences that fit in a single segment. The practical segment size  
 164 is 499, as 3 special tokens of BERT and 10 placeholders for memory are reserved from the model  
 165 input, sized 512. We notice that after training on shorter tasks, it is easier for RMT to solve longer  
 166 versions as it converges to the perfect solution using fewer training steps.

## 167 5.2 Extrapolation Abilities

168 How well does RMT generalize to different sequence lengths? To answer this question, we evaluate  
 169 models trained on a varying number of segments to solve tasks of larger lengths (Figure 4). We  
 170 observe that most models tend to perform well on shorter tasks. The only exception is the single-  
 171 segment reasoning task, which becomes hard to solve once the model is trained on longer sequences.  
 172 One possible explanation is that since the task size exceeds one segment, the model stops expecting  
 173 the question in the first segment, leading to quality degradation.

174 Interestingly, the ability of RMT to generalize to longer sequences also emerges with a growing  
 175 number of training segments. After being trained on 5 or more segments, RMT can generalize nearly  
 176 perfectly for tasks twice as long. To test the limits of generalization, we increase the validation task  
 177 size up to 4096 segments or 2,043,904 tokens (Figure 5). RMT holds up surprisingly well on such  
 178 long sequences, with Detect & memorize being the easiest and Reasoning task the most complex.

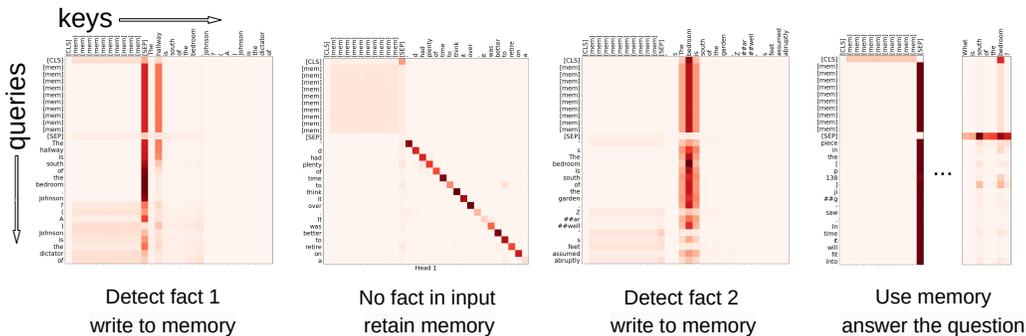


Figure 6: **Attention maps for operations with memory.** These heatmaps show operations performed during specific moments of a 4-segment reasoning task. The darkness of each pixel depends on the attention value between the corresponding key and value. From left to right: RMT detects the first fact and writes its content to memory ([mem] tokens); the second segment contains no information, so the memory keeps the content unchanged; RMT detects the second fact in reasoning tasks and appends it to memory; CLS reads information from the memory to answer the question.

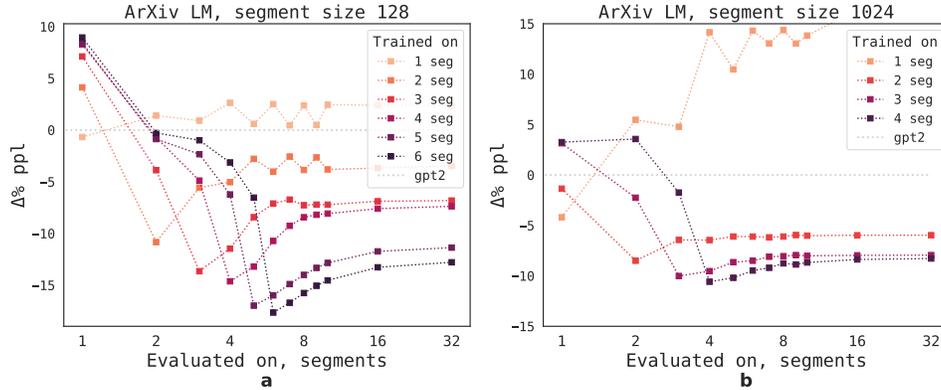


Figure 7: **Generalization of memory on language modeling task.** Models with input sizes **a**: 128 and **b**: 1024 trained with RMT show better performance and generalization across longer sizes of context. Perplexity improvement from training RMT with memory size 2 compared to training the baseline GPT-2 for same number of steps.

179 By examining the RMT attention on specific segments, as shown in Figure 6, we observe that  
 180 memory operations correspond to particular patterns in attention. Furthermore, the high extrapolation  
 181 performance on extremely long sequences, as presented in Section 5.2, demonstrates the effectiveness  
 182 of learned memory operations, even when used thousands of times. The RMT does not have any  
 183 specific memory read/write modules and Transformer learns how to operate with memory recurrently.  
 184 This is particularly impressive, considering that these operations were not explicitly motivated by the  
 185 task loss.

## 186 6 Language Modeling

187 To study the impact of memory on long text understanding, we focus on the long text language  
 188 modeling task conducted using the recurrent approach. To capture long-term dependencies in text,  
 189 memory is required to find and store various type of information between segments. We train the  
 190 GPT-2 Hugging Face checkpoint with 2 memory tokens using the recurrent memory approach on  
 191 the ArXiv documents from The Pile [Gao et al., 2020]. The dataset is preprocessed by splitting each  
 192 document into non-overlapping segments of fixed length, which are prepended with their respective  
 193 histories that consist of several segments. During both training and evaluation we process history  
 194 and target segments one by one and calculate loss and perplexity only on the last target segment.  
 195 Similarly to memorization tasks, we employ curriculum learning for training, starting without history  
 196 and then gradually increasing context size. Language modeling experiments were done on 1-4 A100  
 197 GPUs with single curriculum stage taking up to 2 GPU-days.

198 As expected, increasing the effective context size leads to an improvement in perplexity (Figure  
 199 7). RMT trained for an equal number of steps as the baseline GPT-2 displays substantially lower  
 200 perplexity values. With increasing number of segments in train RMT starts exhibiting better tolerance  
 201 to higher history sizes. Performance of memory models trained without history suffers when applied  
 202 to long contexts, but improves after multi-segment training.

203 We explore the limits of generalization using two tactics. First, we extend the context to contain up to  
 204 1024 segments and run RMT trained on constant number of segments, which is shown in Figure 8  
 205 (a). After a certain context size the perplexity stops changing, remaining stable even when handling  
 206 sequences with more than 1M tokens. Next we test robustness of RMT by introducing noise from  
 207 another distribution in its context. Instead of containing relevant history, a fixed number of input  
 208 segments is sampled using articles from Wikitext-2 dataset, introduced in [Merity et al., 2017]. Figure  
 209 8 (b) illustrates the ability of RMT to retain its superiority over GPT-2 even with useful context vastly  
 210 outnumbered by noise. To understand how memory utilized during generation of the sequence we  
 211 measured perplexity for every position in it (see Figure 9). Baseline shows low prediction quality at  
 212 the beginning of the sequence due to short context available to condition generation. On the other

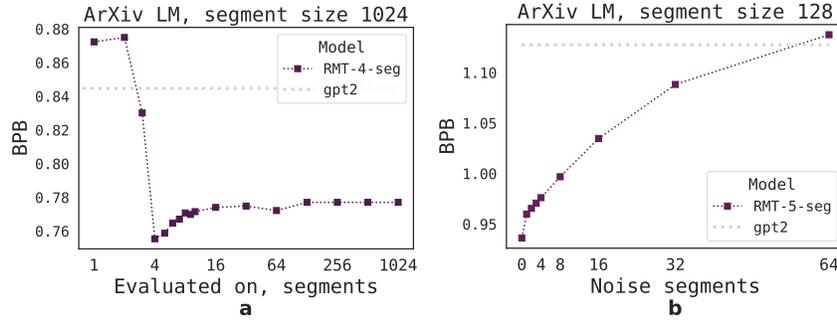


Figure 8: **Finding the context size limit of recurrent models.** **a:** Increasing the number of segments in context of RMT with 2 memory tokens and segment size 1024 reaches a perplexity plateau after a certain size but still below baseline up to more than one million tokens. **b:** Adding noise segments from another distribution to context gradually decreases RMT performance. The bits per byte value (BPB) is computed using the mean bits per token value for the Arxiv set from the Pile [Gao et al., 2020].

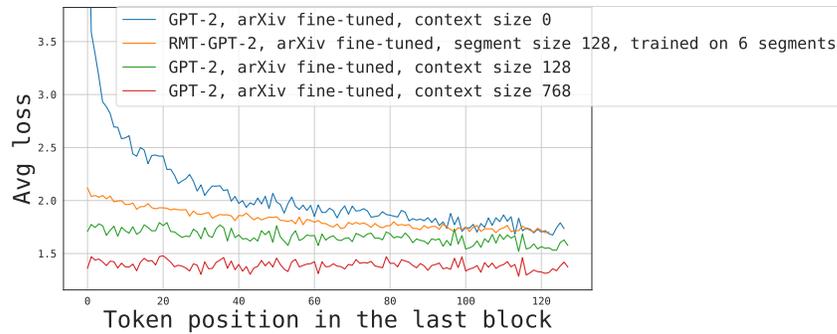


Figure 9: **Memory improves prediction at a beginning of a segment.** As we can see, there is an increase in the loss for tokens at the beginning for GPT-2 (context size 0), showing that it struggles to predict the first tokens since they have no context. The RMT keeps information about previous segments in memory tokens, which helps it to improve tokens predictions. However, showing the model the exact previous context (context size 128 and 768) allows for larger loss gains, but at a higher inference cost. This also shows the importance of local context for language modeling.

213 hand, RMT ensures equally good prediction for all tokens due to carryover of information from the  
 214 previous segment.

## 215 7 Formal Mathematics

216 In this section, we fine-tune our model on a complex mathematical task: generating a proof for a given  
 217 mathematical theorem in formal language. For our experiments, we utilized Lean 3 [de Moura et al.,  
 218 2015] and its library, Mathlib [mathlib Community, 2020], which contains a range of formalized  
 219 theories.

220 Each proof relies on known results, referred to as lemmas. To ensure an effective model, it must  
 221 accurately assess the relevance of a lemma to the given proof. Subsequently, it should memorize the  
 222 lemma’s name and incorporate it within the proof. To construct our dataset, we organized each sample  
 223 into a sequence format. The sequence comprises the theorem statement at the beginning, followed by  
 224 a randomly ordered list of relevant and irrelevant lemmas, and concludes with the human-written  
 225 proof. By adjusting the presence of irrelevant lemmas, we control the sequence length. We further  
 226 divide the sequence into non-overlapping segments of fixed size.

227 For training and evaluation, we calculate the loss and perplexity of the entire sequence. Similar to  
 228 memorization tasks, we train the RMT model and gradually increase size of the sequences. As our  
 229 backbone, we employ GPTNeo [Black et al., 2021] with 1.3B parameters. We incorporate 10 memory  
 230 tokens and set the segment size to 2028.

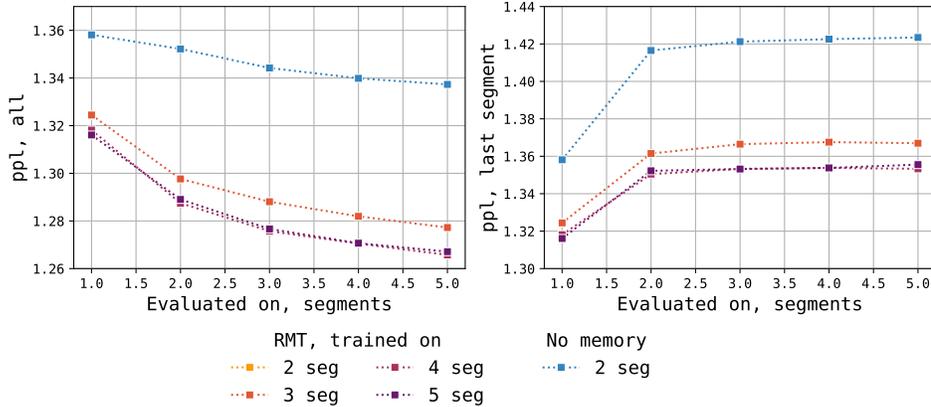


Figure 10: **Lemmas memorization for a theorem proving.** Evaluation of the RMT model and backbone model without memory. Two metrics are calculated: perplexity on all tokens of the sequence (left) and perplexity on the last segment of the sequence (right). RMT model shows better quality.

231 To assess the performance of the RMT model, we compare it with GPTNeo without memory trained  
 232 on a sequences of 2 segments (first segment always contains the theorem statement and the second  
 233 contains the proof). GPTNeo undergoes fine-tuning using the same number of tokens as RMT with  
 234 2 segments. Figure 10 shows the results of the RMT model. The RMT model improves perplexity  
 235 compared to the memory-less model. However, training with 4 or more segments does not enhance  
 236 predictions for longer sequences. According to how the sequence is constructed and split into  
 237 segments, we hypothesize that the model is more concentrated on learning to remember the beginning  
 238 of the last lemma in the previous segment to predict its end in the subsequent segment. The effect of  
 239 detecting and memorizing relevant lemmas and utilizing them in proof generation is less notable. We  
 240 believe that the results can be improved by more careful loss construction and data preparation.

## 241 8 Conclusions

242 The problem of long inputs in Transformers has been extensively studied since the introduction of  
 243 this architecture. Our research has presented a series of significant advancements in augmenting  
 244 and training of Transformer language models. The work expands the conventional capabilities of  
 245 these models through the integration of token-based memory storage and segment-level recurrence  
 246 using recurrent memory (RMT). This mechanism propels the abilities of both encoder-only and  
 247 decoder-only pre-trained Transformers, revealing an unprecedented level of scalability.

248 We have shown that by employing the RMT approach, even models pre-trained on shorter sequences  
 249 can be effectively adapted to manage tasks involving significantly longer sequences. This demon-  
 250 strates that the input length originally designed for the model does not necessarily restrict its potential  
 251 capabilities, thus offering a new perspective on the adaptability of Transformer models.

252 Our work further uncovered the remarkable adaptability of the trained RMT models in extrapolating  
 253 to tasks of varying lengths. The results obtained showcased the RMT’s ability to handle sequences  
 254 exceeding 1 million tokens. Importantly, the computational requirements scaled linearly, thereby  
 255 maintaining computational efficiency even as task length drastically increased. This is a substantial  
 256 contribution that could lead to broader applications and improved performance in handling large-scale  
 257 data. Through an analysis of attention patterns, we provided insight into the operations RMT engages  
 258 to manipulate memory.

259 Overall, our research contributes significantly to the understanding and enhancement of pre-trained  
 260 Transformer language models. It offers a promising direction for future work, particularly in terms of  
 261 handling longer sequences and improving the adaptability of these models.

## 262 References

- 263 Joshua Ainslie, Santiago Ontanon, Chris Alberti, Philip Pham, Anirudh Ravula, and Sumit Sanghai. Etc:  
264 Encoding long and structured data in transformers, 2020.
- 265 Joshua Ainslie, Tao Lei, Michiel de Jong, Santiago Ontañón, Siddhartha Brahma, Yury Zemlyanskiy, David  
266 Uthus, Mandy Guo, James Lee-Thorp, Yi Tay, Yun-Hsuan Sung, and Sumit Sanghai. Colt5: Faster long-range  
267 transformers with conditional computation, 2023.
- 268 Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint*  
269 *arXiv:2004.05150*, 2020.
- 270 Amanda Bertsch, Uri Alon, Graham Neubig, and Matthew R Gormley. Unlimiformer: Long-range transformers  
271 with unlimited length input. *arXiv preprint arXiv:2305.01625*, 2023.
- 272 Sid Black, Leo Gao, Phil Wang, Connor Leahy, and Stella Biderman. GPT-Neo: Large Scale Autoregressive  
273 Language Modeling with Mesh-Tensorflow, March 2021. URL [https://doi.org/10.5281/zenodo.](https://doi.org/10.5281/zenodo.5297715)  
274 [5297715](https://doi.org/10.5281/zenodo.5297715).
- 275 Aydar Bulatov, Yury Kuratov, and Mikhail Burtsev. Recurrent memory transformer. In S. Koyejo, S. Mohamed,  
276 A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*,  
277 volume 35, pages 11079–11091. Curran Associates, Inc., 2022. URL [https://proceedings.neurips.cc/  
278 paper\\_files/paper/2022/file/47e288629a6996a17ce50b90a056a0e1-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/47e288629a6996a17ce50b90a056a0e1-Paper-Conference.pdf).
- 279 Mikhail S Burtsev, Yuri Kuratov, Anton Peganov, and Grigory V Sapunov. Memory transformer. *arXiv preprint*  
280 *arXiv:2006.11527*, 2020.
- 281 Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural  
282 machine translation: Encoder–decoder approaches. In *Proceedings of SSST-8, Eighth Workshop on Syntax,*  
283 *Semantics and Structure in Statistical Translation*, pages 103–111, Doha, Qatar, October 2014. Association for  
284 Computational Linguistics. doi: 10.3115/v1/W14-4012. URL <https://aclanthology.org/W14-4012>.
- 285 Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc Le, and Ruslan Salakhutdinov. Transformer-XL:  
286 Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the*  
287 *Association for Computational Linguistics*, pages 2978–2988, Florence, Italy, July 2019. Association for  
288 Computational Linguistics. doi: 10.18653/v1/P19-1285. URL <https://aclanthology.org/P19-1285>.
- 289 Leonardo de Moura, Soonho Kong, Jeremy Avigad, Floris Van Doorn, and Jakob von Raumer. The lean theorem  
290 prover (system description). In *Automated Deduction-CADE-25: 25th International Conference on Automated*  
291 *Deduction, Berlin, Germany, August 1-7, 2015, Proceedings 25*, pages 378–388. Springer, 2015.
- 292 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional  
293 Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American*  
294 *Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long*  
295 *and Short Papers)*, pages 4171–4186, 2019. URL [https://aclweb.org/anthology/papers/N/N19/  
296 N19-1423/](https://aclweb.org/anthology/papers/N/N19/N19-1423/).
- 297 SiYu Ding, Junyuan Shang, Shuohuan Wang, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. ERNIE-Doc:  
298 A retrospective long-document modeling transformer. In *Proceedings of the 59th Annual Meeting of the*  
299 *Association for Computational Linguistics and the 11th International Joint Conference on Natural Language*  
300 *Processing (Volume 1: Long Papers)*, pages 2914–2927, Online, August 2021. Association for Computational  
301 Linguistics. doi: 10.18653/v1/2021.acl-long.227. URL [https://aclanthology.org/2021.acl-long.](https://aclanthology.org/2021.acl-long.227)  
302 [227](https://aclanthology.org/2021.acl-long.227).
- 303 Angela Fan, Thibaut Lavril, Edouard Grave, Armand Joulin, and Sainbayar Sukhbaatar. Addressing some  
304 limitations of transformers with feedback memory. *arXiv preprint arXiv:2002.09402*, 2020.
- 305 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace  
306 He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The Pile: An 800gb dataset of diverse  
307 text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 308 Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. *arXiv preprint arXiv:1410.5401*, 2014.
- 309 Alex Graves, Greg Wayne, Malcolm Reynolds, Tim Harley, Ivo Danihelka, Agnieszka Grabska-Barwińska,  
310 Sergio Gómez Colmenarejo, Edward Grefenstette, Tiago Ramalho, John Agapiou, Adrià Puigdomènech Badia,  
311 Karl Moritz Hermann, Yori Zwols, Georg Ostrovski, Adam Cain, Helen King, Christopher Summerfield,  
312 Phil Blunsom, Koray Kavukcuoglu, and Demis Hassabis. Hybrid computing using a neural network with  
313 dynamic external memory. *Nature*, 538(7626):471–476, October 2016. ISSN 00280836. URL [http:  
314 //dx.doi.org/10.1038/nature20101](http://dx.doi.org/10.1038/nature20101).

- 315 Edward Grefenstette, Karl Moritz Hermann, Mustafa Suleyman, and Phil Blunsom. Learning to transduce with  
316 unbounded memory, 2015.
- 317 Caglar Gulcehre, Sarath Chandar, Kyunghyun Cho, and Yoshua Bengio. Dynamic neural turing machine with  
318 soft and hard addressing schemes. *arXiv preprint arXiv:1607.00036*, 2016.
- 319 Caglar Gulcehre, Sarath Chandar, and Yoshua Bengio. Memory augmented neural networks with wormhole  
320 connections. *arXiv preprint arXiv:1701.08718*, 2017.
- 321 Mandy Guo, Joshua Ainslie, David Uthus, Santiago Ontanon, Jianmo Ni, Yun-Hsuan Sung, and Yinfei Yang.  
322 LongT5: Efficient text-to-text transformer for long sequences. In *Findings of the Association for Com-*  
323 *putational Linguistics: NAACL 2022*, pages 724–736, Seattle, United States, July 2022. Association for  
324 Computational Linguistics. doi: 10.18653/v1/2022.findings-naacl.55. URL [https://aclanthology.org/](https://aclanthology.org/2022.findings-naacl.55)  
325 [2022.findings-naacl.55](https://aclanthology.org/2022.findings-naacl.55).
- 326 Qipeng Guo, Xipeng Qiu, Pengfei Liu, Yunfan Shao, Xiangyang Xue, and Zheng Zhang. Star-transformer. In  
327 *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational*  
328 *Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1315–1325, Min-  
329 neapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1133.  
330 URL <https://aclanthology.org/N19-1133>.
- 331 Ankit Gupta and Jonathan Berant. Gmat: Global memory augmentation for transformers. *arXiv preprint*  
332 *arXiv:2006.03274*, 2020.
- 333 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780,  
334 November 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. URL [https://doi.org/10.1162/](https://doi.org/10.1162/neco.1997.9.8.1735)  
335 [neco.1997.9.8.1735](https://doi.org/10.1162/neco.1997.9.8.1735).
- 336 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford,  
337 Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Thomas Hennigan, Eric Noland,  
338 Katherine Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karén  
339 Simonyan, Erich Elsen, Oriol Vinyals, Jack Rae, and Laurent Sifre. An empirical analysis of compute-  
340 optimal large language model training. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and  
341 A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 30016–30030. Curran  
342 Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/](https://proceedings.neurips.cc/paper_files/paper/2022/file/c1e2faff6f588870935f114e04a3e5-Paper-Conference.pdf)  
343 [c1e2faff6f588870935f114e04a3e5-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/c1e2faff6f588870935f114e04a3e5-Paper-Conference.pdf).
- 344 DeLesley Hutchins, Imanol Schlag, Yuhuai Wu, Ethan Dyer, and Behnam Neyshabur. Block-recurrent transform-  
345 ers. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural*  
346 *Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=uloenYmLCAo>.
- 347 Armand Joulin and Tomas Mikolov. Inferring algorithmic patterns with stack-augmented recurrent nets, 2015.
- 348 Jie Lei, Liwei Wang, Yelong Shen, Dong Yu, Tamara L. Berg, and Mohit Bansal. Mart: Memory-augmented  
349 recurrent transformer for coherent video paragraph captioning, 2020.
- 350 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on*  
351 *Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- 352 The mathlib Community. The lean mathematical library. In *Proceedings of the 9th ACM SIGPLAN International*  
353 *Conference on Certified Programs and Proofs*. ACM, jan 2020. doi: 10.1145/3372885.3373824. URL  
354 <https://doi.org/10.1145/3372885.3373824>.
- 355 Warren S McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. *The bulletin*  
356 *of mathematical biophysics*, 5(4):115–133, 1943.
- 357 Yuanliang Meng and Anna Rumshisky. Context-aware neural model for temporal information extraction. In  
358 *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*  
359 *Papers)*, pages 527–536, 2018.
- 360 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In  
361 *International Conference on Learning Representations*, 2017. URL [https://openreview.net/forum?](https://openreview.net/forum?id=Byj72udxe)  
362 [id=Byj72udxe](https://openreview.net/forum?id=Byj72udxe).
- 363 OpenAI. Gpt-4 technical report, 2023.

- 364 Richard Yuanzhe Pang, Alicia Parrish, Nitish Joshi, Nikita Nangia, Jason Phang, Angelica Chen, Vishakh  
365 Padmakumar, Johnny Ma, Jana Thompson, He He, and Samuel Bowman. QuALITY: Question answering  
366 with long input texts, yes! In *Proceedings of the 2022 Conference of the North American Chapter of the*  
367 *Association for Computational Linguistics: Human Language Technologies*, pages 5336–5358, Seattle, United  
368 States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.391. URL  
369 <https://aclanthology.org/2022.naacl-main.391>.
- 370 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are  
371 unsupervised multitask learners. 2019.
- 372 Jack W Rae, Jonathan J Hunt, Tim Harley, Ivo Danihelka, Andrew Senior, Greg Wayne, Alex Graves, and  
373 Timothy P Lillicrap. Scaling memory-augmented neural networks with sparse reads and writes, 2016.
- 374 Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, Chloe Hillier, and Timothy P. Lillicrap. Compressive  
375 transformers for long-range sequence modelling. In *International Conference on Learning Representations*,  
376 2020. URL <https://openreview.net/forum?id=Sy1KikSYDH>.
- 377 C Stephen. Kleene. representation of events in nerve nets and finite automata. *Automata studies*, 1956.
- 378 Sainbayar Sukhbaatar, Arthur Szlam, Jason Weston, and Rob Fergus. End-to-end memory networks, 2015.
- 379 Paul J Werbos. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*, 78(10):  
380 1550–1560, 1990.
- 381 Jason Weston, Sumit Chopra, and Antoine Bordes. Memory networks. In Yoshua Bengio and Yann LeCun,  
382 editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May*  
383 *7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1410.3916>.
- 384 Jason Weston, Antoine Bordes, Sumit Chopra, and Tomas Mikolov. Towards ai-complete question answering: A  
385 set of prerequisite toy tasks. In Yoshua Bengio and Yann LeCun, editors, *4th International Conference on*  
386 *Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*,  
387 2016. URL <http://arxiv.org/abs/1502.05698>.
- 388 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric  
389 Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art natural language  
390 processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing:*  
391 *system demonstrations*, pages 38–45, 2020.
- 392 Qingyang Wu, Zhenzhong Lan, Kun Qian, Jing Gu, Alborz Geramifard, and Zhou Yu. Memformer: A memory-  
393 augmented transformer for sequence modeling. In *Findings of the Association for Computational Linguistics:*  
394 *AAACL-IJCNLP 2022*, pages 308–318, Online only, November 2022a. Association for Computational Linguis-  
395 tics. URL <https://aclanthology.org/2022.findings-aacl.29>.
- 396 Yuhuai Wu, Markus Norman Rabe, DeLesley Hutchins, and Christian Szegedy. Memorizing transformers. In  
397 *International Conference on Learning Representations*, 2022b. URL [https://openreview.net/forum?](https://openreview.net/forum?id=TrjbxzRcnf-)  
398 [id=TrjbxzRcnf-](https://openreview.net/forum?id=TrjbxzRcnf-).
- 399 Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago On-  
400 tanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. Big bird: Transform-  
401 ers for longer sequences. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, edi-  
402 tors, *Advances in Neural Information Processing Systems*, volume 33, pages 17283–17297. Curran As-  
403 sociates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/](https://proceedings.neurips.cc/paper_files/paper/2020/file/c8512d142a2d849725f31a9a7a361ab9-Paper.pdf)  
404 [c8512d142a2d849725f31a9a7a361ab9-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/c8512d142a2d849725f31a9a7a361ab9-Paper.pdf).
- 405 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan,  
406 Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig,  
407 Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. Opt: Open pre-trained transformer  
408 language models. *ArXiv*, abs/2205.01068, 2022.