
Recurrent Memory Transformer

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Transformer-based models show their effectiveness across multiple domains and
2 tasks. The self-attention allows to combine information from all sequence ele-
3 ments into context-aware representations. However, global and local information
4 has to be stored mostly in the same element-wise representations. Moreover, the
5 length of an input sequence is limited by quadratic computational complexity of
6 self-attention. In this work, we propose and study a memory-augmented segment-
7 level recurrent Transformer (Recurrent Memory Transformer). Memory allows
8 to store and process local and global information as well as to pass information
9 between segments of the long sequence with the help of recurrence. We implement
10 a memory mechanism with no changes to Transformer model by adding special
11 memory tokens to the input or output sequence. Then Transformer is trained to
12 control both memory operations and sequence representations processing. Results
13 of experiments show that our model performs on par with the Transformer-XL
14 on language modeling for smaller memory sizes and outperforms it for tasks that
15 require longer sequence processing. This makes Recurrent Memory Transformer a
16 promising architecture for applications that require learning of long-term depen-
17 dencies and general purpose in memory processing, such as algorithmic tasks and
18 reasoning.

19 1 Introduction

20 Transformers (Vaswani et al., 2017) have been widely
21 adopted across multiple domains and tasks (Radford
22 et al., 2018; Dong et al., 2018; Devlin et al., 2019;
23 Dosovitskiy et al., 2021; Ramesh et al., 2021; Jaegle
24 et al., 2021). The key component of Transformer
25 layer is a self-attention. Self-attention allows to up-
26 date each sequence element representation with in-
27 formation from all other elements in the sequence.

28 As a result, rich contextual representation for every
29 element is generated at the end of encoding. This
30 way, global sequence-level and local information are
31 stored in a single representation. However, this mix-
32 ing of two types of information in a single represen-
33 tation has limitations. Distributed storage of global
34 features across all sequence elements results in global
35 features "blurring" and makes it harder to access them.

36 Another well-known deficiency of Transformers is poor scaling of self-attention with input sequence
37 length that hurts its applications to long inputs (Child et al., 2019; Guo et al., 2019; Dai et al., 2019;
38 Beltagy et al., 2020; Ainslie et al., 2020; Zaheer et al., 2020; Wang et al., 2020; Choromanski et al.,
39 2020).

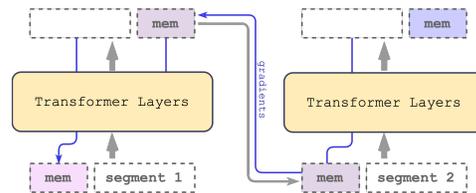


Figure 1: **Recurrent Memory Transformer.** Memory is added as tokens to the input sequence and memory output is passed to the next segment. During training gradients flow from the current segment through memory to the previous segment.

40 Our work introduces a memory-augmented segment-level recurrent Transformer named Recurrent
41 Memory Transformer (RMT). RMT uses a memory mechanism based on special memory tokens (Burt-
42 sev et al., 2020) added to the input sequence. Memory tokens provide additional reserved capacity to
43 the model that could be used to process information which is not directly representing any element in
44 the input sequence. To process long sequences, we split them into segments and pass memory states
45 from a previous to a current segment. This memory passing makes the model recurrent and removes
46 the input sequence length limitations. RMT model can theoretically work with infinite lengths but, in
47 practice, it is limited by memory capacity and the efficiency of memory access/update operations.
48 Our implementation of both memory and recurrence in RMT requires no changes to the Transformer
49 model because modifications are made only to the input and output sequences of the model.

50 We tested RMT on the tasks that require a global information about the whole input sequence to be
51 solved. We use copy, reverse, and associative retrieval tasks in the setting where the input sequence
52 is split into segments. RMT and Transformer-XL perfectly solve these tasks, but exceeding some
53 value of sequence length, RMT starts to outperform Transformer-XL. Also, we experimentally show
54 that the proposed Recurrent Memory Transformer requires less memory size to perform closely to
55 Transformer-XL on language modeling tasks. RMT code and experiments are available¹.

56 Contributions

- 57 • In this study we augment Transformer with token based memory storage and segment-level
58 recurrence.
- 59 • We experimentally evaluate proposed architecture as well as vanilla Transformer and
60 Transformer-XL on memory-intensive tasks such as copy, reverse, associative retrieval
61 and language modeling. We show that RMT outperforms Transformer-XL for sequence
62 processing tasks and on par with Transformer-XL on language modeling but requires less
63 memory.
- 64 • We analysed how the Transformer model learns to use memory. Specific interpretable
65 memory read-write patterns of attention are shown.

66 2 Related work

67 In our study we add a memory to general purpose attention based neural architecture. Memory is
68 a recurrent topic in neural networks research. It had started from the early works (McCulloch and
69 Pitts, 1943; Stephen, 1956) and significantly progressed in 90’s with introduction of *Backpropagation*
70 *Through Time* learning algorithm (Werbos, 1990) and *Long-Short Term Memory* (LSTM) (Hochreiter
71 and Schmidhuber, 1997) neural architecture. Today memory-augmented neural networks (MANNs)
72 usually rely on some kind of recurrent external-memory which is separate from the model’s pa-
73 rameters. *Neural Turing Machines* (NTMs) (Graves et al., 2014) and *Memory Networks* (Weston
74 et al., 2014) are equipped with a storage for vector representations that can be accessed with an
75 attention mechanism. Memory Networks (Weston et al., 2014; Sukhbaatar et al., 2015) were designed
76 to enable reasoning by sequential attention over to the content of a memory. NTMs followed by
77 *Differentiable Neural Computer* (DNC) (Graves et al., 2016) and *Sparse DNC* (Rae et al., 2016)
78 are implemented as recurrent neural networks able to write to memory storage over time. All these
79 models are differentiable and can be trained via backpropagation through time (BPTT). Parallel line
80 of research extends recurrent neural networks such as LSTM with data structures like stacks, lists,
81 or queues (Joulin and Mikolov, 2015; Grefenstette et al., 2015). MANN architectures with a more
82 advanced addressing mechanisms such as address-content separation and multi-step addressing were
83 proposed in (Gulcehre et al., 2016, 2017; Meng and Rumshisky, 2018). The Global Context Layer
84 model (Meng and Rumshisky, 2018) uses the idea of address-content separation to solve the difficulty
85 of training content-based addressing in the canonical NTM.

86 Recent rise of Transformer models also resulted in introduction of a number of new memory archi-
87 tectures. *Transformer-XL* (Dai et al., 2019) introduces a segment-level recurrence at the level of
88 hidden representations. These representations of a sequence are computed and stored in the cache
89 to be reused as an extended context for the next segment. *Compressive Transformer* (Rae et al.,

¹anonymous link. Code, raw experiments results and hyperparameters are provided in supplementary materials.

90 2019) adds the second layer of memory to Transformer-XL. This memory compresses and stores
 91 information from the cache. ∞ -former (Martins et al., 2021) utilizes continuous-space attention and
 92 represents input sequence as a continuous signal to make long-term memory unbounded. *Memory*
 93 *Layers* (Lample et al., 2019) model has a product key memory layer instead of a feed-forward layer
 94 within Transformer block to increase model capacity.

95 In many variations of Transformer different sorts of global representations are added. Among them
 96 are *Star-Transformer* (Guo et al., 2019), *Longformer* (Beltagy et al., 2020), *GMAT* (Gupta and Berant,
 97 2020), *Extended Transformer Construction* (ETC) (Ainslie et al., 2020) and *Big Bird* (Zaheer et al.,
 98 2020). All these architectures re-design self-attention mechanism to reduce its computational com-
 99 plexity with and ensure input coverage with the help of global representations. *Memory Transformer*
 100 (Burtsev et al., 2020) keeps Transformer model intact and adds memory by extending input sequence
 101 with special memory tokens. Perceiver IO (Jaegle et al., 2021) maps an entire arbitrary input to the
 102 fixed number of latent representations. Transformer layers do further processing over latent memory
 103 representations only.

104 Segment-level recurrence in Transformers is actively explored in a number of studies. Transformer-
 105 XL, Compressive Transformer keep previous states and re-use them in subsequent segments. Ernie-
 106 Doc (Ding et al., 2021) improves processing by using same-layer recurrence instead of attending to
 107 previous layer outputs of a precedent segment. Memformer (Wu et al., 2020) introduces a dedicated
 108 memory module to keep previous hidden states in summarized representations. Memformer uses two
 109 special layers added to the Transformer model. Memory cross-attention layer reads from memory
 110 and memory slot attention layer updates it. MART (Lei et al., 2020) has a similar approach as
 111 Memformer but uses memory update rules analogous to LSTM (Hochreiter and Schmidhuber, 1997)
 112 and GRU (Cho et al., 2014). FeedBack Transformer (Fan et al., 2020) goes further with full, and not
 113 segment-level, recurrence. FeedBack Memory merges past hidden representations from all layers
 114 into a single vector and makes it accessible to the computations at any layer. The disadvantage of full
 115 recurrence is that it is less parallelizable. FeedBack Memory requires every sequence element to be
 116 processed sequentially. In segment-level recurrent models, all elements of a segment are processed by
 117 Transformer layers in parallel. Only segments are processed sequentially. Staircase Transformer (Ju
 118 et al., 2021) combines segment-level recurrence and depth recurrence. Staircase models use the
 119 output for previous segments and pass them as input for the next segment. Our Recurrent Memory
 120 Transformer is based on special memory tokens similar to Memory Transformer, segment-level
 121 recurrence as in Transformer-XL, and depth-recurrent mechanism for memory processing similar to
 122 Staircase.

123 3 Recurrent Memory Transformer

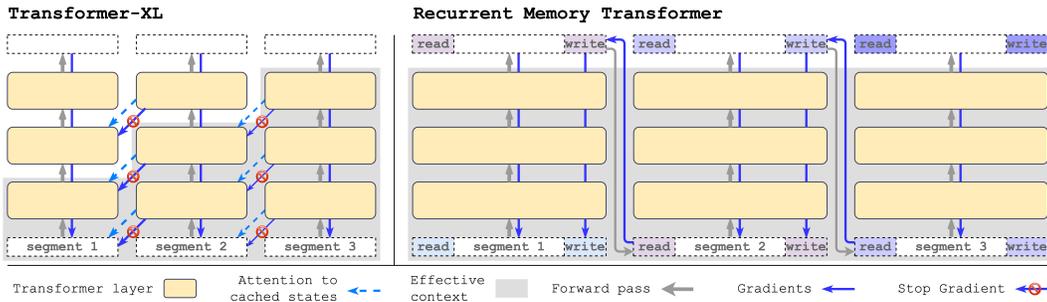


Figure 2: **Comparison of Recurrent Memory Transformer and Transformer-XL architectures.** Recurrent Memory Transformer augments Transformer with global memory tokens and passes them to allow a segment-level recurrence. Special read/write memory tokens are added to the input sequence. Multiple memory tokens can be used in each read/write block. Updated representations of write memory are passed to the next segment. During training, RMT uses BPTT to propagate gradient to previous segments through memory tokens representation. Recurrence with memory has no limitations on effective context length, whereas Transformer-XL can use only finite context with cached states. All RMT memory/recurrence operations are made on the input and output level of the Transformer model.

124 **3.1 Background: Transformer-XL**

125 Transformer-XL (Dai et al., 2019) extends Transformer model with state re-use cache mechanism
 126 for segment-level recurrence and relative position encoding. Input sequence is split on segments
 127 processed sequentially. Hidden states computed for the previous segment M^n are cached for each
 128 transformer layer n . The input of the layer n consists of the last m states from the cached memory
 129 and output of previous Transformer layer for the current segment τ :

$$\tilde{H}_\tau^{n-1} = [SG(M_{-m}^{n-1}) \circ H_\tau^{n-1}], \quad (1)$$

130 here, SG stands for stop-gradient, \circ denotes concatenation. Cached states allow to increase effective
 131 context size of Transformer model and save on compute operations.

132 Then, \tilde{H}_τ^{n-1} goes to Transformer layer to produce layer n outputs for segment τ :

$$H_\tau^n = \text{TransformerLayer}(Q_\tau^n, K_\tau^n, V_\tau^n), \\ Q_\tau^n = W_q^n H_\tau^{n-1}; K_\tau^n = W_k^n \tilde{H}_\tau^{n-1}; V_\tau^n = W_v^n \tilde{H}_\tau^{n-1}.$$

133 In Transformer-XL, self-attention layers are modified to use relative position encodings to improve
 134 generalization to longer attention lengths. The overall architecture is shown in the Figure 2.

135 **3.2 Memory and recurrence**

136 Memory augmented Transformers such as GMAT, ETC, Memory Transformer (Gupta and Berant,
 137 2020; Ainslie et al., 2020; Burtsev et al., 2020) proposed to use special global tokens as storage
 138 for representations. Usually, memory tokens are added to the beginning of the input sequence.
 139 However, in decoder-only architectures the causal attention mask makes impossible for memory
 140 tokens at the start of the sequence to collect information from the subsequent tokens. On the other
 141 hand, if memory tokens are placed at the end of the sequence then preceding tokens unable to
 142 access their representations. To solve this problem we add a recurrence to the sequence processing.
 143 Representations of memory tokens placed at the end of the segment are used as an input memory
 144 representations at the start as well as at the end of the next segment.

145 In the Recurrent Memory Transformer input is augmented with special [mem] tokens, processed in
 146 a standard way along with the sequence of tokens. Each memory token is a real-valued vector. m
 147 memory tokens are added at the beginning of the segment tokens representations H_τ^0 and the same m
 148 tokens are added at the end:

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

149 here N is a number of Transformer layers.

150 The starting group of memory tokens functions as a read memory that allows sequence tokens to
 151 attend to memory states produced at the previous segment. The ending group works as a write
 152 memory that can attend to all current segment tokens and update representation stored in the memory.
 153 As result, H_τ^{write} contains updated memory tokens for the segment τ .

154 Segments of the input sequence are processed sequentially. To enable recurrent connection between
 155 segments, we pass outputs of the memory tokens from the current segment to the input of the next
 156 segment:

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

157 Both memory and recurrence in the RMT are based only on global memory tokens. It allows to
 158 keep the backbone Transformer unchanged and make RMT memory augmentation compatible with
 159 any model from the Transformer family. Memory tokens operate only on the input and output of

160 the model. In this study we implement RMT on top of the original Transformer-XL code. Both
 161 architectures are shown in Figure 2.

162 Recurrence in the RMT is different compared to the Transformer-XL because the former stores only
 163 m memory vectors per segment. On the other hand, the Transformer-XL stores $m \times N$ vectors per
 164 segment. Also, in the RMT model memory representations from the previous segment are processed
 165 by Transformer layers together with the current segment tokens. This makes memory part of RMT
 166 effectively deeper in a number of applied Transformer layers $\tau \times N$. Additionally, we allow all
 167 memory tokens in the read/write block to access all other tokens in the same block. The causal
 168 attention mask is applied only to tokens of the input sequence. The RMT attention mask is shown in
 169 Figure 6 (d).

170 We train the RMT with Backpropagation Through Time (BPTT). During backward pass, unlike in
 171 Transformer-XL, memory gradients are not stopped between segments. The number of previous
 172 segments to backpropagate is a hyperparameter of a training procedure. We vary BPTT unroll in our
 173 experiments from 0 to 4 previous segments. Increasing this parameter is computationally expensive
 174 and requires a lot of GPU RAM. However, such techniques as gradient checkpointing could be used
 175 to alleviate this problem.

176 4 Experiments

177 We designed our experiments to evaluate the ability of Recurrent Memory Transformers to preserve
 178 long-term dependencies across multiple input segments. The first set of experiments includes copy,
 179 reverse, associative retrieval and quadratic equations tasks. The second addresses language modeling
 180 task for word-level on WikiText-103 (Merity et al., 2017) and for character-level on enwik8 (Mahoney,
 181 2006). We compare Recurrent Memory Transformer with Transformer and Transformer-XL models.

182 4.1 Algorithmic Tasks

183 Firstly, we evaluate RMT on algorithmic tasks that require information about the whole input sequence
 184 to be solved successfully. In a recurrent setting, the model has to keep information about all previous
 185 segments to make predictions.

186 In the copy task, an input sequence should be replicated twice after a special start-to-generate token. In
 187 the reverse task, an input sequence should be generated in a reverse order. Input for the associative
 188 retrieval task consists of N key-value pairs. Then one key is randomly selected, and the task is to produce
 189 an appropriate value for the selected key. Another task is to solve quadratic equations. One example
 190 consists of an equation, its solution with discriminant,
 191 and an answer. The task is to generate a solution and
 192 answer, while only answer quality is evaluated.

197 For all tasks, input and output sequences are split
 198 into segments and processed by models sequentially. Datasets for algorithmic tasks were randomly
 199 pre-generated, and the same data was used in all experiments.

200 Transformer-XL and RMT are decoder-only Transformer models. We do not compute loss over
 201 the input sequence before the start-to-generate token. The loss is computed over target sequence
 202 segments only. This procedure is illustrated in Figure 3.

203 4.2 Language Modeling

204 We use two standard benchmarks for language modeling: WikiText-103 and enwik8. WikiText-
 205 103 (Merity et al., 2017) is used for word-level language modeling and contains 103M words from
 206 English Wikipedia articles. Enwik8 (Mahoney, 2006) is used for character-level and consists of 10^8
 207 first bytes of XML text dump of the English Wikipedia.

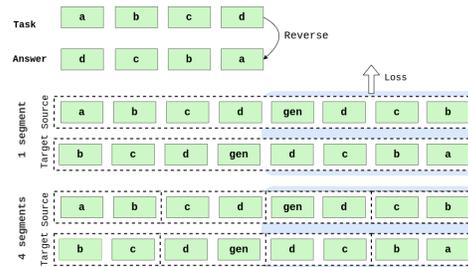


Figure 3: **Reverse task in one and four segments setting for decoder-only models.** Dotted lines show segment borders.

208 We compare Recurrent Memory Transformer with decoder-only Transformer and Transformer-XL
 209 as baselines. Model size and training parameters are selected to match Transformer-XL paper. For
 210 Wikitext-103 an input context length was set to 150 tokens, and for enwik8 it was set to 512 characters.

211 Another set of experiments inspected how RMT handles long-term dependencies and recurrence. We
 212 increased the number of segments and recurrent steps by making segments smaller (50 tokens for
 213 WikiText-103, 128 characters for enwik8). The increased number of recurrent steps makes language
 214 modeling tasks harder for RMT because information has to be stored in the same amount of memory
 215 for more time steps.

216 4.3 Implementation details

217 Our RMT implementation is based on Transformer-XL repository². We also use WikiText-103,
 218 enwik8 data and processing from this repository. Language modeling experiments follow the same
 219 model and training hyperparameters as Transformer-XL. WikiText-103 experiments use 16-layer
 220 Transformers, enwik8 – 12 layer Transformers. We refer to Transformer-XL with memory size equal
 221 to zero as a Baseline.

222 With this experimental setup we were able to reproduce results for the Transformer-XL model close
 223 to the original paper (see Appendix A.4 for enwik8 and Table 2 for WikiText-103).

224 5 Results

225 Baseline, Tr-XL, RMT perform perfectly in the single segment setting on copy and reverse tasks
 226 (Figure 4). In this case models do not need recurrence because the whole sequence is available. When
 227 the number of segments is larger than one, non-recurrent baseline struggles to solve tasks, but both
 228 memory models demonstrate ability to retain required information from the previous segments in
 229 memory.

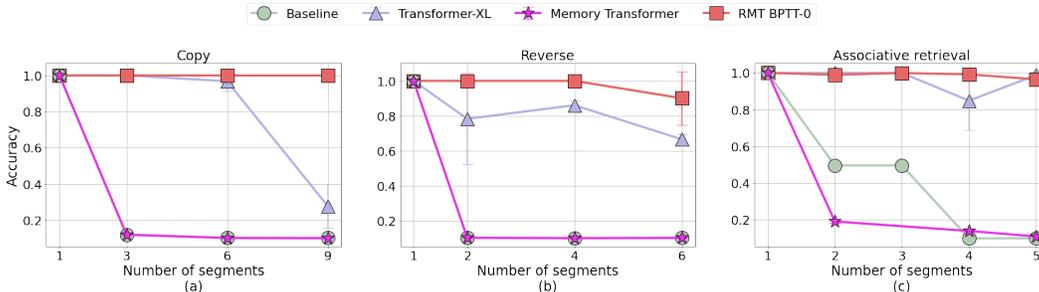


Figure 4: **RMT outperforms Transformer-XL on Copy and Reverse tasks as number of segments increases.** Panels show test set per-character accuracy on copy, reverse, and associative retrieval tasks (from left to right). Memory/cache size equals to the length of a segment for the both models. RMT does not pass gradients between segments in this experiment. MT results are the same as for the Baseline. Source/target sequence lengths for copy, reverse and associative retrieval tasks: 24/48, 24/24, 10/1.

230 As the number of segments increases, Recurrent Memory Transformer starts to outperform
 231 Transformer-XL with memory sizes less than the number of all previous tokens. With the number of
 232 segments up to 6 mean accuracy of Transformer-XL drops by up to 0.2 points, and with 9 segments
 233 score plunges to an accuracy score of 0.2, close to the baseline without memory.

234 Associative retrieval results are similar with the number of segments up to 4. RMT manages to
 235 solve the task with Transformer-XL closely behind. However, in the setting with 5 segments, RMT
 236 performance slightly decreases and Transformer-XL average accuracy rises higher.

237 On the quadratic equations task (Table 1) we have checked that it is possible to solve the task with
 238 the Transformer baseline and no segmentation used. The baseline in this case defines upper bound

²<https://github.com/kimiyoung/transformer-xl>

239 for this task. Enabling recurrence with multiple segments RMT solves the task perfectly, while
 240 Transformer-XL finds the task challenging.

241 Results of experiments on word-level language
 242 modeling on WikiText-103 are shown in Table 2.
 243 In the first section with a segment length of 150,
 244 Tr-XL and RMT outperform the baseline and
 245 Memory Transformer (MT) by a large margin.
 246 It shows the significance of increased effective
 247 context length by Tr-XL cache or RMT memory
 248 for language modeling.

249 MT is slightly better than the Transformer base-
 250 line. This is due to the fact that MT adds special
 251 memory tokens only to the beginning of an in-
 252 put sequence. Autoregressive MT has no way
 253 to write to memory because of the causal at-
 254 tention mask and, therefore, is unable to pass
 255 information between segments. Thus, in the au-
 256 toregressive setting MT could be seen as equiv-
 257 alent to the prefix/promt-tuning (Li and Liang,
 258 2021; Lester et al., 2021).

259 RMT improves over MT memory mechanism
 260 with read/write blocks. The best RMT models
 261 with memory size 10 and 25 show similar per-
 262 formance as Transformer-XL with a memory size
 263 equal to 75. RMT learns to use smaller memory
 264 more effectively than Transformer-XL. Addi-
 265 tionally, the smaller memory size of RMT leads
 266 to reducing required GPU memory for running
 267 the model.

268 Decreasing the size of segments to 50, we force
 269 models to work with longer recurrent dependen-
 270 cies as the number of recurrent steps increases.
 271 RMT with memory consisting of a single vec-
 272 tor shows similar results to Transformer-XL
 273 with memory size 10. It is worth noting that
 274 Transformer-XL memory consists of hidden rep-
 275 resentations from all layers (in this case, it is
 276 10×16 vectors) when RMT memory is only
 277 `memory_size` vectors. Transformer-XL with memory size 50 and RMT with memory size 5 shows
 278 similar perplexity values.

279 On enwik8 RMT models with memory size 5 and Transformer-XL with memory size 40 show
 280 similar results. Confirming that RMT learns to use smaller amounts of memory representation more
 281 effectively. All experiments results on enwik8 dataset are in Appendix A.4 and Table 3.

282 We analyze how number of segments, sequence length, length of training context, and memory size
 283 affect models’ performance on different tasks in Figure 5. As we split sequence into more segments
 284 it becomes more crucial to be able to pass information between segments. For the copy task, we
 285 split 360 tokens sequence into multiple segments. In Figure 5a we observe that Tr-XL performance
 286 starts to degrade and eventually falls to the baseline model performance as the number of segments
 287 increases. In contrast, RMT continues to solve the task perfectly. In the more extreme setting, when
 288 we keep memory size fixed, but increase the total length of sequence Tr-XL fails shortly, while RMT
 289 starts to degrade on 720 tokens sequence length (Figure 5b).

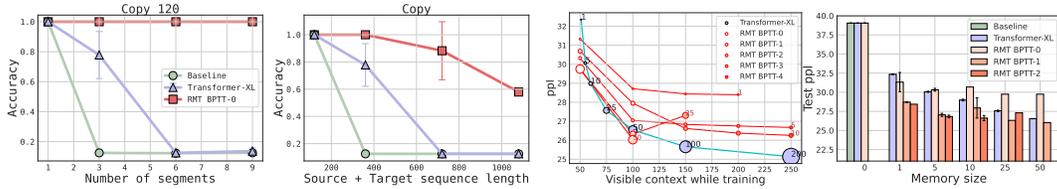
290 Recurrent Memory Transformer learns to make predictions depending on `#BPTT_unrolls` previous
 291 segments +1 current segment. Transformer-XL does not use BPTT and relies only on `memory_size`
 292 cached states and current segment making in total: `memory_size + segment_length` tokens. In

Table 1: Test set solve rate for quadratic equations. Input sequence length is 180 tokens and consists of quadratic equation, solution, and answer. The input sequence is split into a number of segments with an answer as the last segment. Accuracy is 1.0 if the full answer is predicted correctly.

MODEL	MEMORY	SEGMENTS	ACC \pm STD
BASELINE	0	1	0.99 \pm 0.01
Tr-XL	30	6	0.93 \pm NA
RMT	30	6	0.99 \pm 0.002

Table 2: Test set perplexity on WikiText-103. Average perplexity for the best performed variations of RMT models reported (see full results in Appendix A.5). Underlined values show Tr-XL and RMT models with close results. RMT models with smaller memory sizes achieve similar scores to Tr-XL models with larger memory.

MODEL	MEMORY	SEGMENT LEN	PPL \pm STD
BASELINE	0	150	29.95 \pm 0.15
MT	10	150	29.63 \pm 0.06
Tr-XL (PAPER)	150	150	24.0
Tr-XL (OURS)	150	150	24.12 \pm 0.05
Tr-XL	75	150	<u>24.68</u> \pm 0.01
Tr-XL	25	150	25.57 \pm 0.02
RMT BPTT-3	10	150	<u>25.04</u> \pm 0.07
RMT BPTT-2	25	150	<u>24.85</u> \pm 0.31
BASELINE	0	50	39.05 \pm 0.01
Tr-XL	100	50	25.66 \pm 0.01
Tr-XL	50	50	<u>26.54</u> \pm 0.01
Tr-XL	25	50	27.57 \pm 0.09
Tr-XL	10	50	28.98 \pm 0.11
RMT BPTT-1	1	50	<u>28.71</u> \pm 0.03
RMT BPTT-3	10	50	<u>26.37</u> \pm 0.01



(a) Increasing number of segments for a fixed sequence length up to 1080 tokens. (b) Increasing sequence length up to 1080 tokens. (c) Dependency of test set perplexity on visible context and deeper recurrence length of 360 tokens. (d) Increased memory size result in better performance. and memory size 60.

Figure 5: (a) RMT is able to solve copy task perfectly with multiple recurrent steps, while Tr-XL fails. (b) RMT learns to use memory of the same fixed size more effectively than TR-XL as sequence length increases. (c&d) WikiText-103 with 50 tokens segment length. (c) Marker size corresponds to memory size. Visible context at training time can be increased by enlarging Tr-XL cache or using more BPTT unrolls for RMT. Increasing visible context leads to lower perplexity for both models. (d) Test set perplexity for different memory sizes. When memory size is zero Tr-XL and RMT are just baseline Transformer models without recurrence.

293 Figure 5c, we compare RMT and Tr-XL according to the described value of visible context at training
 294 time.

295 RMT with a single memory vector could be trained to achieve lower perplexity as Transformer-XL
 296 with memory size 10. It means that RMT can learn to store information from previous observations
 297 more compactly. Another observation is that RMT with memory sizes 10 and 25 performs worse but
 298 closely to Transformer-XL even when Transformer-XL has access to more non-compressed states
 299 (50, 100, 200) from previous segments. Furthermore, we observed instabilities and out-of-memory
 300 issues during RMT training for a larger number of BPTT unrolls and memory sizes.

301 Recurrent Memory Transformer does not benefit from increasing memory size from 5 to 50, but
 302 results of Transformer-XL better scale with memory size (Figure 5d). RMT models with memory
 303 size 5 have close results to Transformer-XL with cache 50, confirming that RMT learns to store more
 304 compact representations. Dynamic of RMT perplexity suggests that there is some optimal memory
 305 size for RMT to solve the task, and further increase does not add much. Training RMT with one
 306 BPTT unroll drastically improves its results showing the importance of BPTT training (Figure 5d).

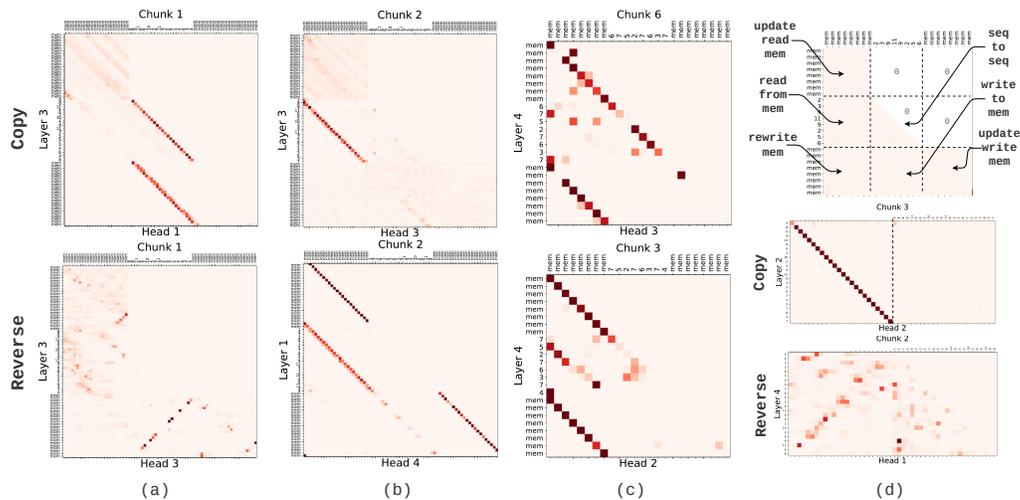


Figure 6: Selected attention map patterns of memory models. (a) - write to memory, (b) - read from memory (RMT, segment length=24, memory size=24), (c) - rewrite from read memory to write memory (RMT, segment length=8, memory size=8), (d) - read from previous hidden states (Transformer-XL, segment length=24, memory size=24)

307 To get an understanding of memory operations, learned by models during training algorithmic tasks,
308 it is useful to look at attention maps. Figure 6 shows some heatmaps from attention layers of models
309 trained on copy and reverse. The darkness of each pixel shows how much the element from the
310 corresponding row "attends" to its column.

311 In each RMT attention map sequence tokens are preceded by read memory, located at the top left
312 corner, and followed by write memory at the bottom right. Lines at the central part of (a), top image
313 shows classic attention of token sequence to itself, but the bottom line represents the operation of
314 writing of sequence tokens to memory in straight order. When completing reverse, the (a), bottom
315 image model learns to write the sequence in the reversed order, which is in line with common sense.

316 When it comes to reproducing the target sequence, model accesses memory, Figure 6 (b) and writes
317 to the output sequence. Another operation (c) is rewriting from read memory to write memory. It is
318 commonly used by RMT in settings with larger number of segments to keep information about recent
319 segments longer.

320 Transformer-XL mechanism of accessing memory (d) does not allow straightforward writing to mem-
321 ory without changing sequence token representations. Sequential reading from cache is represented
322 by straight lines on Transformer-XL attention maps.

323 Using token representations as storage harms model performance in tasks with larger number of
324 segments. On reverse with 4 segments Transformer-XL with limited memory size 6, Figure 8 (b)
325 attempts to mix representations of tokens and read multiple symbols from one cached state in next
326 segments. This results in the average accuracy of 0.8 on the given task. Despite having the same
327 memory size, RMT manages to compress the whole segment in memory tokens Figure 8 (a) and
328 achieve mean accuracy 1.0.

329 Visualizations from Figure 6 and Figure 8 provide evidence to support our hypotheses that Tr-XL
330 has to mix representations from previous and current segments in the same hidden states to pass
331 information between segments. Also, visualizations show how memory tokens in RMT help mitigate
332 such kind of mixing.

333 RMT ability of sequence compression to memory is illustrated in Figure 7. For copy with 6
334 segments RMT compresses and then reads the sequence of 12 tokens with just 6 memory tokens.
335 For Transformer-XL decreasing memory size harms the accuracy score significantly with number of
336 segments larger than 2.

337 6 Conclusions

338 In this paper we introduced Recurrent Memory Transformer a simple recurrent memory augmentation
339 of Transformer model. RMT is implemented by extension of an input sequence with special global
340 memory tokens and segment-level recurrence.

341 In our experiments we compared RMT with Transformer baseline and Transformer-XL which
342 is a well-known modification of Transformer for long sequences. RMT almost perfectly solves
343 Copy, Reverse as well as quadratic equations tasks for sequences consisting of multiple segments
344 outperforming Transformer-XL. It also demonstrates quality for associative retrieval task on par
345 with Transformer-XL. As expected, baseline Transformer fails to solve these tasks for multi-segment
346 settings.

347 RMT trained as language model significantly ahead of Transformer baseline and shows quality
348 metrics similar to Transformer-XL but for up to 10 times smaller memory size. Experimental results
349 demonstrate that for fixed memory size backpropagating gradients for more segments improves
350 performance of RMT. Analysis of attention maps suggests that better RMT performance can be
351 related to more effective storage of input representations in dedicated memory tokens compared to
352 mixing representations storage in Transformer-XL. Overall, results of the study show that dedicated
353 memory storage and recurrence provided by Recurrent Memory Transformer make it a promising
354 architecture for applications that require learning of long-term dependencies and general purpose
355 in-memory processing, such as algorithmic tasks and reasoning. Furthermore, we believe that RMT
356 could open the way for adding memory and recurrence to other models in the Transformer family.

357 **References**

- 358 Joshua Ainslie, Santiago Ontanon, Chris Alberti, Philip Pham, Anirudh Ravula, and Sumit Sanghai. Etc:
359 Encoding long and structured data in transformers, 2020.
- 360 Jimmy Ba, Geoffrey E Hinton, Volodymyr Mnih, Joel Z Leibo, and Catalin Ionescu. Using fast weights to attend
361 to the recent past. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural
362 Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL [https://proceedings.
363 neurips.cc/paper/2016/file/9f44e956e3a2b7b5598c625fcc802c36-Paper.pdf](https://proceedings.neurips.cc/paper/2016/file/9f44e956e3a2b7b5598c625fcc802c36-Paper.pdf).
- 364 Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint
365 arXiv:2004.05150*, 2020.
- 366 Mikhail S Burtsev, Yuri Kuratov, Anton Peganov, and Grigory V Sapunov. Memory transformer. *arXiv preprint
367 arXiv:2006.11527*, 2020.
- 368 Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers,
369 2019.
- 370 Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural
371 machine translation: Encoder–decoder approaches. In *Proceedings of SSST-8, Eighth Workshop on Syntax,
372 Semantics and Structure in Statistical Translation*, pages 103–111, Doha, Qatar, October 2014. Association for
373 Computational Linguistics. doi: 10.3115/v1/W14-4012. URL <https://aclanthology.org/W14-4012>.
- 374 Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos,
375 Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers.
376 *arXiv preprint arXiv:2009.14794*, 2020.
- 377 Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. Transformer-xl:
378 Attentive language models beyond a fixed-length context, 2019.
- 379 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional
380 Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American
381 Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long
382 and Short Papers)*, pages 4171–4186, 2019. URL [https://aclweb.org/anthology/papers/N/N19/
383 N19-1423/](https://aclweb.org/anthology/papers/N/N19/N19-1423/).
- 384 SiYu Ding, Junyuan Shang, Shuohuan Wang, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. ERNIE-Doc:
385 A retrospective long-document modeling transformer. In *Proceedings of the 59th Annual Meeting of the
386 Association for Computational Linguistics and the 11th International Joint Conference on Natural Language
387 Processing (Volume 1: Long Papers)*, pages 2914–2927, Online, August 2021. Association for Computational
388 Linguistics. doi: 10.18653/v1/2021.acl-long.227. URL [https://aclanthology.org/2021.acl-long.
389 227](https://aclanthology.org/2021.acl-long.227).
- 390 Linhao Dong, Shuang Xu, and Bo Xu. Speech-transformer: A no-recurrence sequence-to-sequence model for
391 speech recognition. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing
392 (ICASSP)*, pages 5884–5888, 2018. doi: 10.1109/ICASSP.2018.8462506.
- 393 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,
394 Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An
395 image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on
396 Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- 397 Angela Fan, Thibaut Lavril, Edouard Grave, Armand Joulin, and Sainbayar Sukhbaatar. Addressing some
398 limitations of transformers with feedback memory. *arXiv preprint arXiv:2002.09402*, 2020.
- 399 Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines, 2014.
- 400 Alex Graves, Greg Wayne, Malcolm Reynolds, Tim Harley, Ivo Danihelka, Agnieszka Grabska-Barwińska,
401 Sergio Gómez Colmenarejo, Edward Grefenstette, Tiago Ramalho, John Agapiou, Adrià Puigdomènech Badia,
402 Karl Moritz Hermann, Yori Zwols, Georg Ostrovski, Adam Cain, Helen King, Christopher Summerfield,
403 Phil Blunsom, Koray Kavukcuoglu, and Demis Hassabis. Hybrid computing using a neural network with
404 dynamic external memory. *Nature*, 538(7626):471–476, October 2016. ISSN 00280836. URL [http:
405 //dx.doi.org/10.1038/nature20101](http://dx.doi.org/10.1038/nature20101).
- 406 Edward Grefenstette, Karl Moritz Hermann, Mustafa Suleyman, and Phil Blunsom. Learning to transduce with
407 unbounded memory, 2015.

- 408 Caglar Gulcehre, Sarath Chandar, Kyunghyun Cho, and Yoshua Bengio. Dynamic neural Turing machine with
409 soft and hard addressing schemes. *arXiv preprint arXiv:1607.00036*, 2016.
- 410 Caglar Gulcehre, Sarath Chandar, and Yoshua Bengio. Memory augmented neural networks with wormhole
411 connections. *arXiv preprint arXiv:1701.08718*, 2017.
- 412 Qipeng Guo, Xipeng Qiu, Pengfei Liu, Yunfan Shao, Xiangyang Xue, and Zheng Zhang. Star-transformer, 2019.
- 413 Ankit Gupta and Jonathan Berant. Gmat: Global memory augmentation for transformers. *arXiv preprint*
414 *arXiv:2006.03274*, 2020.
- 415 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780,
416 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)
417 [neco.1997.9.8.1735](https://doi.org/10.1162/neco.1997.9.8.1735).
- 418 Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda
419 Koppula, Daniel Zoran, Andrew Brock, Evan Shelhamer, et al. Perceiver io: A general architecture for
420 structured inputs & outputs. *arXiv preprint arXiv:2107.14795*, 2021.
- 421 Armand Joulin and Tomas Mikolov. Inferring algorithmic patterns with stack-augmented recurrent nets, 2015.
- 422 Da Ju, Stephen Roller, Sainbayar Sukhbaatar, and Jason Weston. Staircase attention for recurrent processing of
423 sequences. *arXiv preprint arXiv:2106.04279*, 2021.
- 424 Guillaume Lample, Alexandre Sablayrolles, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. Large
425 memory layers with product keys, 2019.
- 426 Jie Lei, Liwei Wang, Yelong Shen, Dong Yu, Tamara L. Berg, and Mohit Bansal. Mart: Memory-augmented
427 recurrent transformer for coherent video paragraph captioning, 2020.
- 428 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In
429 *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059,
430 Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
431 doi: 10.18653/v1/2021.emnlp-main.243. URL <https://aclanthology.org/2021.emnlp-main.243>.
- 432 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings*
433 *of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International*
434 *Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, Online,
435 August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.353. URL
436 <https://aclanthology.org/2021.acl-long.353>.
- 437 Matt Mahoney. Large text compression benchmark, 2006. URL [http://www.matthmahoney.net/dc/text.](http://www.matthmahoney.net/dc/text.html)
438 [html](http://www.matthmahoney.net/dc/text.html).
- 439 Pedro Henrique Martins, Zita Marinho, and André FT Martins. ∞ -former: Infinite memory transformer. *arXiv*
440 *preprint arXiv:2109.00301*, 2021.
- 441 Warren S McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. *The bulletin*
442 *of mathematical biophysics*, 5(4):115–133, 1943.
- 443 Yuanliang Meng and Anna Rumshisky. Context-aware neural model for temporal information extraction. In
444 *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*
445 *Papers)*, pages 527–536, 2018.
- 446 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models.
447 In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26,*
448 *2017, Conference Track Proceedings*. OpenReview.net, 2017. URL [https://openreview.net/forum?](https://openreview.net/forum?id=Byj72udxe)
449 [id=Byj72udxe](https://openreview.net/forum?id=Byj72udxe).
- 450 Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understand-
451 ing by generative pre-training. 2018. URL [https://www.cs.ubc.ca/~amuham01/LING530/papers/](https://www.cs.ubc.ca/~amuham01/LING530/papers/radford2018improving.pdf)
452 [radford2018improving.pdf](https://www.cs.ubc.ca/~amuham01/LING530/papers/radford2018improving.pdf).
- 453 Jack W Rae, Jonathan J Hunt, Tim Harley, Ivo Danihelka, Andrew Senior, Greg Wayne, Alex Graves, and
454 Timothy P Lillicrap. Scaling memory-augmented neural networks with sparse reads and writes, 2016.
- 455 Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, and Timothy P. Lillicrap. Compressive transformers for
456 long-range sequence modelling, 2019.

- 457 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya
458 Sutskever. Zero-shot text-to-image generation. In Marina Meila and Tong Zhang, editors, *Proceedings of
459 the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning
460 Research*, pages 8821–8831. PMLR, 18–24 Jul 2021. URL [https://proceedings.mlr.press/v139/
461 ramesh21a.html](https://proceedings.mlr.press/v139/ramesh21a.html).
- 462 C Stephen. Kleene. representation of events in nerve nets and finite automata. *Automata studies*, 1956.
- 463 Sainbayar Sukhbaatar, Arthur Szlam, Jason Weston, and Rob Fergus. End-to-end memory networks, 2015.
- 464 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser,
465 and Illia Polosukhin. Attention is All you Need. In *Advances in neural information processing systems*, pages
466 5998–6008, 2017. URL <http://papers.nips.cc/paper/7181-attention-is-all-you-need>.
- 467 Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear
468 complexity, 2020.
- 469 Paul J Werbos. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*, 78(10):
470 1550–1560, 1990.
- 471 Jason Weston, Sumit Chopra, and Antoine Bordes. Memory networks, 2014.
- 472 Qingyang Wu, Zhenzhong Lan, Jing Gu, and Zhou Yu. Memformer: The memory-augmented transformer. *arXiv
473 preprint arXiv:2010.06891*, 2020.
- 474 Manzil Zaheer, Guru Guruganesh, Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip
475 Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. Big bird: Transformers for longer sequences,
476 2020.

477 Checklist

- 478 1. For all authors...
- 479 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
480 contributions and scope? [Yes]
- 481 (b) Did you describe the limitations of your work? [Yes] We mention training instabilities
482 and GPU RAM issues in Section 5.
- 483 (c) Did you discuss any potential negative societal impacts of your work? [No] The
484 proposed model and method do not have any specific impacts. All general negative
485 societal impacts applicable to the field could be potentially relative.
- 486 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
487 them? [Yes]
- 488 2. If you are including theoretical results...
- 489 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 490 (b) Did you include complete proofs of all theoretical results? [N/A]
- 491 3. If you ran experiments...
- 492 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
493 imental results (either in the supplemental material or as a URL)? [Yes] We include
494 code, training scripts, and raw experimental data in the supplementary material. The
495 supplemental materials would be published on github with the final version of the paper.
496 Instructions for language modeling data&experiments are taken from Tr-XL repo.
- 497 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were
498 chosen)? [Yes] See Section 4.3, Appendix A, and provided supplementary material.
- 499 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
500 iments multiple times)? [Yes] All the key experiments results are reported with std.
501 Furthermore, we provide raw experimental data in the supplementary materials.
- 502 (d) Did you include the total amount of compute and the type of resources used (e.g., type
503 of GPUs, internal cluster, or cloud provider)? [Yes] We used different GPUs depending
504 on the task: 1080Ti, V100, A100. We provide this information in Appendix A for each
505 task.

- 506 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 507 (a) If your work uses existing assets, did you cite the creators? [Yes] We refer to the
- 508 original Tr-XL code and Tr-XL paper. We use it for establish baselines and setting our
- 509 methods. See Section 4.3
- 510 (b) Did you mention the license of the assets? [No] Tr-XL licence is Apache 2.0 and
- 511 available at its github repo.
- 512 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 513 Our code is in the supplemental material.
- 514 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 515 using/curating? [No] We used publicly available Tr-XL code (Apache 2.0) and datasets.
- 516 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 517 information or offensive content? [No] We use either synthetic data or datasets collected
- 518 from the Wikipedia (Wikitext-103, enwik8).
- 519 5. If you used crowdsourcing or conducted research with human subjects...
- 520 (a) Did you include the full text of instructions given to participants and screenshots, if
- 521 applicable? [N/A]
- 522 (b) Did you describe any potential participant risks, with links to Institutional Review
- 523 Board (IRB) approvals, if applicable? [N/A]
- 524 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 525 spent on participant compensation? [N/A]

526 A Training details and additional results

527 A.1 Algorithmic tasks

528 Datasets were randomly generated by uniformly sampling tokens from dictionary into task sequences

529 and generating targets accordingly to the tasks. After generation, datasets are fixed for all experiments.

530 Copy and reverse use sequences of sizes 24, 40, 120, 240, and 360, making total copy/reverse

531 input length 48/72, 80/120, 240/360, 480/720, 720/1080. The associative retrieval task consists of 4

532 key-value pairs and one randomly selected key; the answer consists of one value. Train, validation

533 and test sizes of copy 24, reverse 24 and associative retrieval datasets are 100000, 5000 and 10000.

534 Transformer-XL had the same cache size on training and validation to match RMT.

535 For training all models on copy and reverse, we used constant learning rate $1e-4$ with reduction on

536 plateau with decay factor of 0.5. Copy and reverse were solved by models with 4 layers and 4 heads,

537 associative retrieval models had 6 layers and 4 heads. Models with the same context size and memory

538 size were trained for the same number of steps and the same training parameters.

539 Experiments with sequence length 24 were conducted on a single Nvidia GTX 1080 Ti GPU from 1

540 hour to 2-3 days. Copy and reverse on longer sequence lengths were done on more powerful Tesla

541 V100 using 1-3 devices with training time varying from 1 hour to 3-4 days.

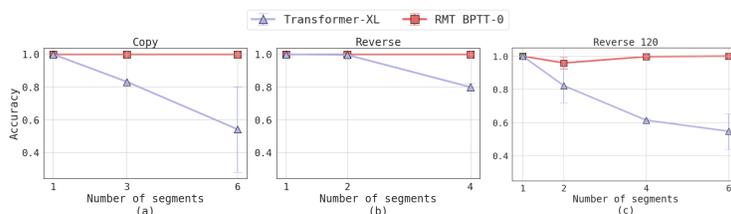


Figure 7: Test set per-character accuracy on copy (a), reverse with a sequence length 24 (b) and 120 (c). Memory size is limited to half of the length of a segment for (a) and (b) and to 60 tokens for (c). Target sequence length equals source length for reverse and doubles for the copy.

542 **A.2 Associative retrieval**

543 We used code for the task dataset generation from (Ba et al., 2016)³.

544 **A.3 Quadratic equations**

545 This dataset consists of equations with integer coefficients with step-by-step solutions using a
546 discriminant. Process of equation generation is started from uniformly sampling real roots x_1, x_2
547 from -100 to 100. The answer of an equation is represented as x_1, x_2 . Next, we find the equation as
548 multiplication of two parentheses $(x - x_1)(x - x_2) = 0$, which is expanded to $x^2 - (x_1 + x_2)x +$
549 $x_1x_2 = 0$. Next, we multiply all coefficients by a random natural number α from 1 to 10. The final
550 equation form is $\alpha x^2 - \alpha(x_1 + x_2)x + \alpha x_1x_2 = 0$. A dataset sample is made of these stages in
551 reversed order. We also provide a string with the discriminant calculation to help find the equation
552 roots. 20 percent of equations in the dataset do not have real roots.

553 Example equation string:

554 $-4*x^2+392*x-2208=0,$

555 solution string:

556 $x^2-98*x+552=0;D=98^2-4*1*552=7396=86^2;x=(98-86)/2=6;x=(98+86)/2=92 ,$

557 and answer:

558 $6,92$

559 Each solution step is tokenized on char level and padded to the length of 30 tokens. The total length
560 of each training sample is 180, the dataset has 100000 training, 10000 validation and 20000 test
561 samples.

562 For this task we used models with 6 layers, 6 heads and segment sizes 180 and 30. Trained was
563 performed with the same schedule as copy and reverse on a single GTX 1080 ti for 1-2 days. Memory
564 size for RMT and Transformer-XL was chosen equal to the segment length.

565 **A.4 Enwik8**

566 We verified our experimental setup by reproducing Transformer-XL results on enwik8 dataset
567 (Table 3). We used 12-layer Baseline (Transformer), Transformer-XL, RMT in all enwik8 experiments.
568 All results on enwik8 dataset are in Table 3. We used 2 NVIDIA A100 80Gb GPUs, training time
569 varied from 10 to 30 hours depending on sequence length, memory size, and number of BPTT unrolls.

570 **A.5 WikiText-103**

571 We used 16-layer models in all experiments on WikiText-103 dataset. Training hyperparameters were
572 used from (Dai et al., 2019) and authors PyTorch scripts⁴. All results on WikiText-103 dataset are in
573 Table 4. In most of the WikiText-103 experiments, we used 2 NVIDIA A100 80Gb GPUs, training
574 time varied from 10 to 30 hours depending on sequence length, memory size, and number of BPTT
575 unrolls.

³https://github.com/GokuMohandas/fast-weights/blob/539fb10e3c384d5f782af2560bf28631cd0eaa61/fw/data_utils.py

⁴<https://github.com/kimiyoung/transformer-xl>

Table 3: Test set bits-per-character on enwik8. Our experimental setup shows similar scores to the original paper (Dai et al., 2019) with segment length 512.

MODEL	MEMORY	SEGMENT LEN	BPC \pm STD
TR-XL (DAI ET AL., 2019)	512	512	1.06
TR-XL (OURS)	512	512	1.071
TR-XL	200	128	1.140
TR-XL	100	128	1.178
TR-XL	75	128	1.196
TR-XL	40	128	1.230 \pm 0.001
TR-XL	20	128	1.261
TR-XL	10	128	1.283 \pm 0.001
RMT BPTT-1	5	128	1.241 \pm 0.002
RMT BPTT-2	5	128	1.231 \pm 0.002
RMT BPTT-1	10	128	1.240 \pm 0.006
RMT BPTT-2	10	128	1.228 \pm 0.003
RMT BPTT-0	20	128	1.301
RMT BPTT-1	20	128	1.229
RMT BPTT-2	20	128	1.222

Table 4: Test set perplexity on WikiText-103. All experiments with RMT and Tr-XL models.

MODEL	MEMORY	SEGMENT LEN	PPL \pm STD
BASELINE	0	150	29.95 \pm 0.15
MT	10	150	29.63 \pm 0.06
MT	25	150	29.67 \pm 0.03
MT	75	150	29.69 \pm 0.02
MT	150	150	29.82 \pm 0.35
Tr-XL (PAPER)	150	150	24.0
Tr-XL (OURS)	150	150	24.12 \pm 0.05
Tr-XL	75	150	24.68 \pm 0.01
Tr-XL	25	150	25.57 \pm 0.02
RMT BPTT-0	10	150	26.85 \pm 0.02
RMT BPTT-1	10	150	25.92 \pm 1.07
RMT BPTT-2	10	150	25.32 \pm 0.61
RMT BPTT-3	10	150	25.04 \pm 0.07
RMT BPTT-0	25	150	29.73
RMT BPTT-1	25	150	24.91
RMT BPTT-2	25	150	24.85 \pm 0.31
BASELINE	0	50	39.05 \pm 0.01
Tr-XL	200	50	25.14
Tr-XL	100	50	25.66 \pm 0.01
Tr-XL	50	50	26.54 \pm 0.01
Tr-XL	25	50	27.57 \pm 0.09
Tr-XL	10	50	28.98 \pm 0.11
Tr-XL	5	50	30.06 \pm 0.07
Tr-XL	1	50	32.35 \pm 0.03
RMT BPTT-0	1	50	31.33 \pm 1.26
RMT BPTT-1	1	50	28.71 \pm 0.03
RMT BPTT-2	1	50	28.44
RMT BPTT-3	1	50	28.40 \pm 0.03
RMT BPTT-0	5	50	30.32 \pm 0.18
RMT BPTT-1	5	50	27.05 \pm 0.20
RMT BPTT-2	5	50	26.83 \pm 0.18
RMT BPTT-3	5	50	26.75 \pm 0.26
RMT BPTT-4	5	50	26.67 \pm 0.03
RMT BPTT-0	10	50	30.69 \pm 0.01
RMT BPTT-1	10	50	27.95 \pm 1.32
RMT BPTT-2	10	50	26.62 \pm 0.34
RMT BPTT-3	10	50	26.37 \pm 0.01
RMT BPTT-4	10	50	26.25 \pm 0.19
RMT BPTT-0	25	50	29.75
RMT BPTT-1	25	50	26.32
RMT BPTT-2	25	50	27.31
RMT BPTT-0	50	50	29.75
RMT BPTT-1	50	50	26.03

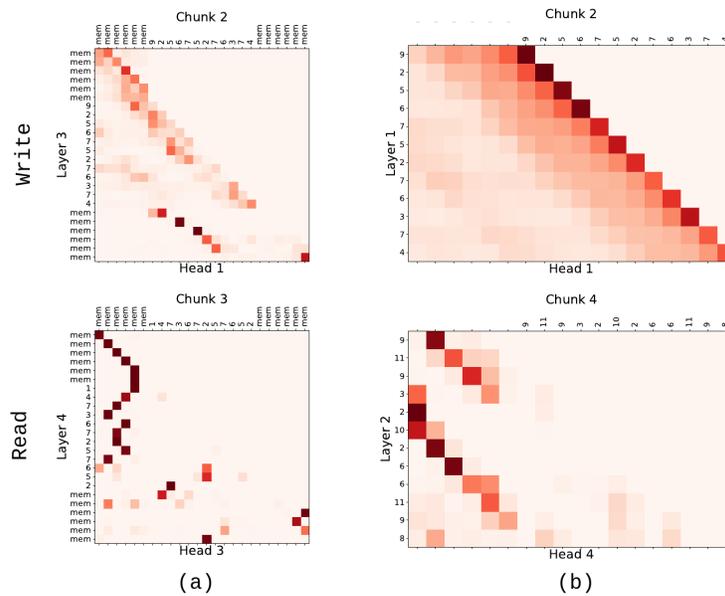


Figure 8: Approaches to compression and decompression of sequence with length 12 and memory with size 6. (a) - RMT, (b) - Transformer-XL.