RECURRENT ACTION TRANSFORMER WITH MEMORY

Anonymous authors

000

001 002 003

004

006

008 009

010

011

012

013

014

015

016

017

018

019

021

023

025

026 027

028 029 Paper under double-blind review

ABSTRACT

Recently, the use of transformers in offline reinforcement learning has become a rapidly developing area. This is due to their ability to treat the agent's trajectory in the environment as a sequence, thereby reducing the policy learning problem to sequence modeling. In environments where the agent's decisions depend on past events (POMDPs), it is essential to capture both the event itself and the decision point in the context of the model. However, the quadratic complexity of the attention mechanism limits the potential for context expansion. One solution to this problem is to extend transformers with memory mechanisms. This paper proposes a Recurrent Action Transformer with Memory (RATE), a novel model architecture that incorporates a recurrent memory mechanism designed to regulate information retention. To evaluate our model, we conducted extensive experiments on memory-intensive environments (ViZDoom-Two-Colors, T-Maze, Memory Maze, Minigrid-Memory), classic Atari games, and MuJoCo control environments. The results show that using memory can significantly improve performance in memory-intensive environments, while maintaining or improving results in classic environments. We believe that our results will stimulate research on memory mechanisms for transformers applicable to offline reinforcement learning. The code is available at https://anonymous.4open.science/r/RATE-B01F/.

1 INTRODUCTION

Transformers (Vaswani et al., 2017), originally developed for Natural 031 Language Processing (NLP), perform well in Reinforcement Learn-033 ing (RL) (Agarwal et al., 2023; Li 034 et al., 2023): online RL (Parisotto et al., 2020; Esslinger et al., 2022; Melo, 2022; Team et al., 2023), offline 037 RL (Chen et al., 2021; Lee et al., 2022; Jiang et al., 2023), and model-based RL (Chen et al., 2022; Micheli et al., 040 2023; Robine et al., 2023), including solving the credit assignment prob-041 lem and working in memory-intensive 042 environments (Chen et al., 2021; Ni 043 et al., 2023; Grigsby et al., 2024), 044 provided that the entire trajectory fits 045 within the model context. However, 046 transformers struggle with long se-047 quences due to quadratic attention



Figure 1: Recurrent Action Transformer with Memory (RATE). R – returns-to-go, o – observations, a – actions, M_n – segment's S_n memory embeddings.

complexity, limiting their use in long inference tasks. Several approaches attempt to increase
the context size (Dai et al., 2019; Bulatov et al., 2022; Ding et al., 2023), but such models may
become unstable when trained on long sequences (Zhang et al., 2022), or use a specific sparse
attention mechanism that is unsuitable for non-NLP tasks (Beltagy et al., 2020; Zaheer et al., 2020;
Ding et al., 2023). Memory mechanisms offer a promising solution to account for past information
without increasing context size. Our work explores memory in transformers for RL, inspired by NLP
results (Dai et al., 2019; Bulatov et al., 2022). The RL setting differs from NLP in the processing of

input sequences, requiring specialized encoders for observations, rewards, and actions, and is also characterized by significant sparsity in some tasks.

In RL memory has two senses. One is using past information within an episode to make decisions (Lampinen et al., 2021; Ni et al., 2023). The other is transferring experience from one environment to another (Melo, 2022; Kang et al., 2023; Team et al., 2023), improving generalizability, sample efficiency, and solving Meta-RL (Duan et al., 2016; Wang et al., 2016) tasks. Our work focuses on the first case (Ni et al., 2023): using past information to make decisions within the same episode.

In this paper, we propose the **Recurrent Action Transformer with Memory (RATE**, Figure 1), 063 a model that uses several memory mechanisms: memory embeddings, caching of previous hidden 064 states of previous tokens, and Memory Retention Valve (MRV). We empirically show that memory 065 mechanisms effectively preserve information from previous steps, allowing the model to use past 066 information when making decisions in the present. MRV is designed to control the process of updating 067 memory embeddings and prevent the loss of important information when processing long sequences, 068 thus enabling the processing of highly sparse tasks. To evaluate the memory mechanisms, we perform 069 extensive experiments in memory-intensive environments: ViZDoom-Two-Colors (Sorokin et al., 2022), Memory Maze (Pasukonis et al., 2022), Minigrid-Memory (Chevalier-Boisvert et al., 2023), 071 and Passive T-Maze (Ni et al., 2023), as well as on standard RL benchmarks: Atari (Bellemare et al., 2013) and MuJoCo (Fu et al., 2021). We also study the impact of memory on the performance of the 072 proposed model. 073

The proposed model interpolates and extrapolates well outside the transformer context, is able to
 retain important information for a long time when operating in highly sparse environments, and
 allows to compensate the effect of bias in the training data.

- Our contribution can be summarized as follows:
 - 1. We propose the Recurrent Action Transformer with Memory (RATE), a transformer model for offline RL that makes use of memory mechanisms: memory embeddings, caching of hidden states of previous tokens, and the Memory Retention Valve (MRV). The proposed MRV is based on the cross-attention architecture and is designed to prevent information loss from memory embeddings and significantly improve the performance of RATE in memory-intensive environments with sparse structure (see section 3).
 - 2. We show that RATE significantly outperforms strong baselines with and without memory mechanisms in memory-intensive environments, including ViZDoom Two-Colors, Memory Maze, Minigrid-Memory, and T-Maze (see subsection 4.2).
 - 3. We demonstrate that RATE achieves better or comparable results in classic Atari games and MuJoCo control tasks, demonstrating that the proposed model is suitable for different types of tasks and emphasizing its universality (see subsection 4.2).
 - 2 BACKGROUND

079

081

082

084 085

090

092

093 094

095

2.1 OFFLINE REINFORCEMENT LEARNING

096 In RL (Sutton & Barto, 2018), we assume that the task can be described by a Markov Decision 097 Process (MDP) as a tuple $\langle S, A, P, R \rangle$. The process consists of states $s \in S$, actions $a \in A$, a 098 state transition function $\mathcal{P}(s'|s, a)$, and an immediate reward function $r = \mathcal{R}(s, a)$. The states are 099 assumed to have the Markov Property, that is $\mathbb{P}(s_{t+1}|s_t) = \mathbb{P}(s_{t+1}|s_1,\ldots,s_t)$. Given a timestep t, 100 we use $r_t = R(s_t, a_t)$ to denote the immediate reward that the agent receives at state s_t performing 101 action a_t at that timestep. We describe trajectory τ of length T as a sequence of states s_i , actions 102 a_i , and immediate rewards r_i : $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{T-1}, a_{T-1}, r_{T-1})$. We denote returnto-go (Chen et al., 2021) R_t of trajectory τ at the timestep t as a sum of future rewards from the timestep t to the end of trajectory: $R_t = \sum_{t'=t}^{T-1} r_{t'}$. The goal of a RL agent is to learn policy π that 103 104 105 maximizes the expected return. In online RL, the trajectories used to train an agent are obtained iteratively as the agent interacts with the environment. In offline RL, the agent does not interact with 106 the environment during training. A fixed set of trajectories collected by an arbitrary policy is used for 107 training. Although such a setting is more difficult because it does not allow additional exploration of

the environment or generation of new trajectories, it is preferable for tasks where interaction with the environment is costly or risky, such as in robotics.

110 111

112

2.2 PARTIALLY OBSERVABLE MARKOV DECISION PROCESS

113 In the real world, there are frequent situations where the full state of the environment is not available to the agent, so the Markov Property is violated, and the agent is said to receive observations instead 114 of state as input. In this case, observations are no longer sufficient statistics of the past to make 115 a decision in the current step. An example would be a robot navigating an environment based 116 on a camera image or a situation in which a decision must be made based on information from 117 the past that is not available in the current observation. A Partially Observable Markov Decision 118 Process (POMDP) is used in such cases. POMDP is a generalization of the MDP and is written as 119 $\langle \mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{P}, \mathcal{R}, \mathcal{Z} \rangle$, where $o \in \mathcal{O}$ – observations, $s \in \mathcal{S}$ – states, $a \in \mathcal{A}$ – actions, $r = \mathcal{R}(s, a)$ – 120 immediate reward function, $\mathcal{P}(s'|s, a)$ – state transition function and \mathcal{Z} is an observation function 121 $\mathcal{Z}_{s'o}^{a} = P(O = o_{t+1}|S_{t+1} = s', A_t = a)$. To work successfully in such environments, mechanisms 122 such as memory are needed to allow the use of information from the past (Parisotto et al., 2020; 123 Lampinen et al., 2021; Esslinger et al., 2022). In our work, we propose an approach of adding memory 124 to the agent in the form of memory embeddings, caching of hidden states of previous tokens, and 125 MRV. In this paper we consider the setting of an offline model-free RL, where learning is formulated as a sequence modeling problem (Chen et al., 2021). 126

- 127
- 128 129

3 RECURRENT ACTION TRANSFORMER WITH MEMORY

130 In this paper, we introduce a new archi-131 tecture, Recurrent Action Transformer with Memory (RATE), in which we uti-132 lized recurrently trained memory embed-133 dings (Bulatov et al., 2022) and caching 134 of hidden states of previous tokens (Dai 135 et al., 2019) to add memory, and Mem-136 ory Retention Valve (MRV) to control 137 information leakage from memory em-138 beddings, allowing sparse sequences to 139 be processed. The architecture of RATE 140 is shown in Figure 1.

141 The RATE scheme is outlined in Algo-142 rithm 1. In the initial phase, we obtain 143 embeddings R, \tilde{o} and \tilde{a} from the returns-144 to-go R, observations o, and actions a, 145 respectively, using the corresponding en-146 coders from Table 10. We then generate a 147 trajectory $\tau_{0:T-1}$ consisting of triplets of 148 these embeddings according to the tech-149 nique described in the Decision Trans-150 former (DT) (Chen et al., 2021) paper (Algorithm 3). 151

152 Next, the trajectory $\tau_{0:T-1}$ is divided 153 into N = T//K segments $S_n \in \mathbb{R}^{3K \times d}$, $n \in [0, N-1]$, each consist-155 ing of K triplets, where K is the context 156 length and d is the model dimension. To 157 each segment S_n , memory embeddings Algorithm 1 Recurrent Action Transformer with Memory **Input**: $R \in \mathbb{R}^T, o \in \mathbb{R}^{d_o \times T}, a \in \mathbb{R}^T$ Parameters: $M \in \mathbb{R}^{m \times d}$ $\begin{array}{ll} 1: \ \tilde{R} \in \mathbb{R}^{T \times d} \leftarrow \texttt{Encoder}_{R}(R) \\ \tilde{o} \in \mathbb{R}^{T \times d} \leftarrow \texttt{Encoder}_{o}(o) \\ \tilde{a} \in \mathbb{R}^{T \times d} \leftarrow \texttt{Encoder}_{a}(a) \end{array}$ 2: $\tau_{0:T-1} \leftarrow \{ (\tilde{R}_0, \tilde{o}_0, \tilde{a}_0), \dots, (\tilde{R}_{T-1}, \tilde{o}_{T-1}, \tilde{a}_{T-1}) \}$ 3: $M \leftarrow M_0 \sim \mathcal{N}(0,1)$ 4: for n in range [0, T//K - 1] do
$$\begin{split} S_n &\in \mathbb{R}^{3K \times d} \leftarrow \tau_{n \times K:(n+1) \times K} \\ \tilde{S}_n &\in \mathbb{R}^{(3K+2m) \times d} \leftarrow \text{concat}(M_n, S_n, M_n) \end{split}$$
5: 6: $\hat{a}_n, M_{n+1} \leftarrow \operatorname{Transformer}(\tilde{S}_n)$ 7: $M_{n+1} \leftarrow MRV(M_n, M_{n+1})$ 8: **Output:** $\hat{a}_n \to \mathcal{L}(a_n, \hat{a}_n), M_{n+1}$ 9: end for

Algorithm 2 Memory Retention ValveInput: $M_n, M_{n+1} \in \mathbb{R}^{m \times d}$ Parameters: $\mathbf{W}_Q^h, \mathbf{W}_K^h, \mathbf{W}_V^h \in \mathbb{R}^{d_h \times d}, \mathbf{W}_M \in \mathbb{R}^{d \times d}$ 1: $\mathbf{Q}_h \leftarrow M_n \mathbf{W}_Q^h T$ 2: $\mathbf{K}_h \leftarrow M_{n+1} \mathbf{W}_K^h T$ 3: $\mathbf{V}_h \leftarrow M_{n+1} \mathbf{W}_V^h T$ 4: $M_{n+1}^h \leftarrow \operatorname{softmax}(\frac{\mathbf{Q}_h \mathbf{K}_h^T}{\sqrt{d}}) \mathbf{V}_h$ 5: $M_{n+1} \leftarrow \operatorname{concat}(M_{n+1}^0, \dots, M_{n+1}^h)$ 6: $M_{n+1} \leftarrow M_{n+1} \mathbf{W}_M^T$ Output: M_{n+1}

158 $M_n \in \mathbb{R}^{m \times d}$ are concatenated at the beginning and at the end, forming a $\tilde{S}_n \in \mathbb{R}^{(3K+2m) \times d}$, where 159 *m* is the number of memory embeddings. These segments \tilde{S}_n are then fed into the transformer and 160 the output is the predicted actions $\hat{a}_n = \hat{a}_{n \times K:(n+1) \times K}$, which are used to compute loss, and new 161 memory embeddings M_{n+1} , which are then processed by the MRV block and transmitted to the next segment S_{n+1} . 162 The MRV scheme is presented in Algorithm 2 and is based on the cross-attention basis. After the 163 transformer processes the segment S_n , the previous memory embeddings M_n and the new memory 164 embeddings M_{n+1} obtained at the output of the transformer are fed to the input of the MRV block. 165 Next, M_n are multiplied by the query matrix \mathbf{W}_Q^h , and M_{n+1} are multiplied by the key \mathbf{W}_K^h and 166 value \mathbf{W}_V^h matrix for each of the attention heads h. Then the attention scores are calculated using 167 the softmax() function, the results for each of all attention heads are concatenated and after linear 168 transformation using the matrix \mathbf{W}_{M}^{T} the final memory tokens M_{n+1} are obtained at the output, which is passed to the next segment. 169

¹⁷⁰ Unlike the DT learning process, where random fragments of length K are cut from trajectories, we process segments of length K sequentially, which allows us to capture all information in the processed trajectory. Thus, using memory mechanisms, we are able to process sequences of length $K_{eff} = K \times N$, where K_{eff} is the effective context (Bulatov et al., 2022).

174 175

4 EXPERIMENTAL EVALUATION

176 177

We designed experiments to accomplish two primary objectives: (a) to demonstrate the advantage of our RATE model in memory-intensive environments (T-Maze (Ni et al., 2023), ViZDoom-Two-Colors (Sorokin et al., 2022), Memory Maze (Pasukonis et al., 2022), Minigrid-Memory (Chevalier-Boisvert et al., 2023)), and (b) to investigate the effectiveness of the proposed model in classical MDPs to demonstrate its versatility (Atari (Bellemare et al., 2013) and MuJoCo (Fu et al., 2021)).

For comparison with RATE, we chose DT (Chen et al., 2021) as the main baseline, and adapted the
memory-augmented architectures Recurrent Memory Transformer (RMT) (Bulatov et al., 2022) and
Transformer-XL (TrXL) (Dai et al., 2019) developed for NLP tasks to the RL domain. Information
about the environments used can be found in Table 7.

187 188

189

190

191

192

193 194

195

196

197

199

200

201

202

203

204

205

206

207

208

209

4.1 Memory-intensive environments

To test RATE memory mechanisms, we use memory-intensive environments Figure 2, i.e., environments where the agent requires memory to operate successfully. A brief description of these environments is presented below, and a full description and data collection methodology can be found in the Appendix B.

- 1. **ViZDoom-Two-Colors** (Sorokin et al., 2022) an agent in an acid-filled room observes a quickly disappearing green or red pillar. To stay alive, the agent must recall the pillar's color and gather items of the same color.
- 2. **T-Maze** (Ni et al., 2023) an agent navigates a T-shaped corridor, receiving a clue at the start about which direction to turn at the end of the corridor. The task tests memory in a prolonged sparse reward environment (the agent receives a reward only at the end).
- 3. **Memory Maze** (Pasukonis et al., 2022) an agent navigates a maze, seeking objects matching the color of its view frame. The frame color changes after each successful find. The goal is to collect the most matching objects within a time limit.
- 4. **Minigrid-Memory** (Chevalier-Boisvert et al., 2023) a similar task to T-Maze, but with different observation spaces and reward functions (see Appendix B, Table 7). Another important difference is that the agent appears at a random point at the beginning of the episode, not at the beginning of the corridor. Thus, in the case of Minigrid-Memory, it is necessary to reach that clue first (in T-Maze it is a memory problem, in Minigrid-Memory it is a memory and credit assignment problem (Ni et al., 2023)).
- In the experiments, the same hyperparameters presented in Table 8 were used for RATE, DT, RMT, and TrXL to simplify the comparison. The context length K and the number of segments N were chosen so that the effective context $K_{eff} = K \times N$ for RATE, RMT, and TrXL covers important events in memory-intensive environments during training. In turn, since DT has no memory mechanisms, for it $K = K_{eff}$. Thus, the context length K for RATE, RMT and TrXL is less than the context length for DT by a factor of N, but all models process the same parts of trajectories. More information about the training procedure for each environment can be found in Appendix D.



Figure 2: Memory-intensive environments with different observation spaces and reward functions used to test the performance of the memory mechanism in the RATE model.



Figure 3: Results for the ViZDoom-Two-Colors: with (a) and without (b) pillar in the first 45 steps of the episode; calculated at environment steps 0 - 89 (c) and 90 - 179 (d); depending on the return-to-go (e, f, g). Pillar disappears after first 45 steps, $K_{eff} = 90$.

4.2 EXPERIMENTAL RESULTS

In this section, the main experimental results for each of the environments are presented in the corresponding paragraphs. Additional results can be found in the Appendix F. For each of the experiments, the same techniques were used for all models to obtain the results presented in the Appendix E.
Unless otherwise indicated, all baselines were trained from scratch.

240

223

224

225

226

227

228 229 230

231

232

233

234

235

ViZDoom-Two-Colors. The dataset for this environment was collected using a pre-trained Advan tage Actor Critic (A2C) (Beeching et al., 2019), which has a slight bias in favor of selecting green
 items even if red items are required. However, the dataset is balanced by the pillars colors. Figure 3
 (a) shows the inference results on all pillars, separately only on red pillars and separately only on
 green pillars. As illustrated in the Figure 3 (a), for inference with the presence of a disappearing pillar
 at the beginning of the episode, all baselines have an average total reward for inference on green
 pillars greater than for inference on red pillars.

To prove that this is not due to the peculiarities of the algorithms, but to the presence of bias in the data, we performed an additional inference without a pillar at the beginning of the episode, demonstrating the ability of all baselines to collect exclusively green items. As can be seen from Figure 3 (b), DT learns the distribution of training data and is unable to remember the pillar color. In turn, baselines with memory mechanisms such as RATE, RMT and TrXL are successful in this task. The poor performance of baselines with memory on red pillars without the pillar at the episode beginning proves that it is the color of the pillar that they remember.

255 This conclusion of DT's inability to use information out of context as opposed to RATE is supported by the experimental results presented in Figure 3 (c, d), which illustrates the inference results for 256 the first 90 steps, where the pillar are entirely captured in the effective context, and the subsequent 257 90 steps, where the pillar color information begins to disappear from the effective context (the 258 context window moves as a sliding window). As a result, there is a drop in total reward for red-pillar 259 environments by almost a factor of two, indicating DT's inability to memorize information to use it 260 out of context. In turn, for RATE, RMT, and TrXL, the values of total reward in the first and second 261 cases are almost unchanged, indicating their ability to utilize information outside of the current 262 context window. 263

Figure 3 (e, f, g) demonstrates the dependence of model performance on the return-to-go. In this paper, we used an empirical estimate of the target reward as the average of the top-10 total rewards in the training dataset. As can be seen, RATE not only significantly outperforms DT, but also outperforms the memory-augmented models RMT and TrXL.

Furthermore, Table 1 demonstrates that RATE outperforms not only transformer-based models (RMT, TrXL), but also recurrent baselines, which forget the pillars color fairly quickly and start to collecting green items like DT.

270 **T-Maze.** The Figure 4 shows the inference results for the T-Maze. To validate an agent's long-term 271 memory capabilities, we performed training on trajectories of length 90 and inference on corridors 272 of size 30 - 900, i.e. much larger than the effective context $K_{eff} = 90$ of models. The Figure 4 273 demonstrates that DT's ability to solve the task is limited by the length of the corridors in the training 274 data. Specifically, DT-3, trained on trajectories of length $3 \times 30 = 90$ with $K = K_{eff} = 90$, exhibits a significant drop in performance (demonstrating the performance of the persistent agent, i.e. it 275 successfully reaches the junction, but at the junction it turns one way regardless of the clue) when 276 tasked with inference on corridors exceeding 90 in length. 277

278 In turn, RATE-3, RMT-3, and TrXL-3 (trained 279 on $3 \times 30 = 90$ steps) perform significantly bet-280 ter at inference corridor lengths longer than the model saw during training. Moreover, RATE-281 3 outperforms other memory-augmented base-282 lines, indicating its ability to perform effectively 283 in sparse environments. This confirms the ability 284 of the RATE model to successfully memorize 285 important information and retain it through a long time. 287

1.0 0.8 9.0 Rate Sacces data 0.2 RATE-3 - RMT-3 Random Agent DT-3 - TrXL-3 Persistent Agent 0.0 400 500 600 800 900 700 0 100 200 300 Test corridor length

Additionally, Table 1 compares RATE with recurrent baselines. The results indicate that these baselines, unlike RATE, cannot handle sparse information, as shown by SR = 0.5.

292 293

295

296

297

298



Table 1: Comparison of Transformer (DT), RNN (Decision LSTM (DLSTM) (Siebenborn et al., 2022), Decision GRU (DGRU)) and SSM (Decision Mamba (DMamba) (Ota, 2024)) models with RATE in memory-intensive environments. The results indicate the inability of the RNN and SSM models to train successfully on trajectories of length $90 \times 3 = 270$ tokens, unlike RATE. SR – Success Rate. DGRU is obtained by replacing the LSTM block with the GRU block (Chung et al., 2014) in DLSTM. [†] K = 90 ($K_{eff} = 3 \times 30 = 90$ for RATE).

			Г	-Maze		
	Random	DLSTM	DGRU	DMamba	DT	RATE (ours)
$\mathbf{SR}\left(K=T=9\right)$	0.0	1.0	1.0	1.0	1.0	1.0
$\mathbf{SR} \left(K = T = 30 \right)$	0.0	0.6	1.0	1.0	1.0	1.0
$\mathbf{SR}\left(K=T=90\right)$	0.0	0.5 ± 0.0	0.5 ± 0.0	0.5 ± 0.0	1.0 ± 0.0	1.0 ± 0.0
	ViZDoom-Two-Colors [†]					
Reward[Total]	4.82	13.1 ± 0.6	12.9 ± 0.2	26.9 ± 1.9	24.8 ± 1.4	41.5 ± 1.0
Reward[Reds]	4.66	8.8 ± 0.7	9.4 ± 0.5	6.9 ± 0.4	7.2 ± 0.4	38.2 ± 5.1
Reward[Greens]	4.98	17.5 ± 1.6	16.3 ± 0.8	$\textbf{46.9} \pm \textbf{4.2}$	$\textbf{42.3} \pm \textbf{3.3}$	$\textbf{44.7} \pm \textbf{5.8}$

Minigrid-Memory. The Figure 5 shows the results for the Minigrid-Memory. Training was conducted on grids of size 31x31, inference was conducted on grids of size 11x11 – 91x91. Unlike the previously discussed T-Maze, the credit assignment problem is also addressed here, since the agent first has to reach the oracle and find out which object to turn towards in the future.







the other baselines on larger grids, that is, it interpolates well and extrapolates poorly. In turn, for
RMT we observe exactly the opposite situation: RMT interpolates poorly and extrapolates well.
RATE performs slightly worse than TrXL but better than RMT on small grid sizes, and slightly worse
than RMT but better than TrXL on large grid sizes, but on average has interpolation and extrapolation
abilities better than RMT and TrXL individually, as well as greater stability.

Memory Maze. Table 2 shows the results of comparing RATE with other basic baselines with and without memory. The results show that RATE is able to memorize the implicit information like maze structure more efficiently, which is reflected in a higher average reward per episode.
 Table 2: Results for the Memory Maze 9x9 environment.

	DT (Chen et al., 2021)	RMT (Bulatov et al., 2022)	TrXL (Dai et al., 2019)	RATE
Reward	6.83 ± 0.51	7.27 ± 0.21	7.12 ± 0.24	$\textbf{7.64} \pm \textbf{0.41}$

Atari and MuJoCo. The results for Atari and MuJoCo are presented in Table 3 and Table 4. Results for Decision Mamba (DMamba) (Ota, 2024) and Mamba as Decision Maker (MambaDM) (Cao et al., 2024) are from the corresponding papers. A more detailed description of the results obtained can be found in the Appendix D. The results demonstrate that RATE not only performs as well as the algorithms specifically designed for offline RL, but in many cases outperforms them in classical environments that do not require memory, which indicates the versatility of the model.

Table 3: Raw scores for Atari games. Green – top-1 result, light green – top-2 result within the standard deviation.

Environment	CQL (Kumar et al., 2020)	DT (Chen et al., 2021)	DMamba (Ota, 2024)	MambaDM (Cao et al., 2024)	RATE
Breakout	62.5	76.9 ± 27.3	70.6 ± 9.3	106.9 ± 5.8	111.0 ± 2.9
Qbert	14013.2	2215.8 ± 1523.7	5786.0 ± 1295.2	10052.5 ± 1116.5	12486.9 ± 280.4
SeaQuest	782.2	1129.3 ± 189.0	992.1 ± 57.7	1286.0 ± 42.0	1037.9 ± 53.7
Pong	18.8	17.1 ± 2.9	1.6 ± 15.3	18.4 ± 0.8	18.8 ± 0.3

Table 4: Scores normalized according to the protocol in Fu et al. (2021) for MuJoCo control tasks. ME – Medium-Expert dataset, M – Medium dataset, MR – Medium-Replay dataset. RATE outperforms DT in 9/9 of the cases. Green – top-1 result, light green – top-2 result within the standard deviation.

Dataset	Environment	CQL (Kumar et al., 2020)	DT (Chen et al., 2021)	TAP (Jiang et al., 2023)	DMamba (Ota, 2024)	MambaDM (Cao et al., 2024)	RATE
ME	HalfCheetah	91.6	86.8 ± 1.3	91.8 ± 0.8	91.9 \pm 0.6	86.5 ± 1.2	87.4 ± 0.1
ME	Hopper	105.4	107.6 ± 1.8	105.5 ± 1.7	111.1 ± 0.3	110.5 ± 0.3	112.5 ± 0.2
ME	Walker2d	108.8	108.1 ± 0.2	107.4 ± 0.9	108.3 ± 0.5	108.8 ± 0.1	108.7 ± 0.5
M	HalfCheetah	44.4	42.6 ± 0.1	45.0 ± 0.1	42.8 ± 0.1	42.8 ± 0.1	43.5 ± 0.3
M	Hopper	58.0	67.6 ± 1.0	63.4 ± 1.4	83.5 ± 12.5	85.7 ± 7.8	77.4 \pm 1.4
М	Walker2d	72.5	74.0 ± 1.4	64.9 ± 2.1	78.2 ± 0.6	78.2 ± 0.6	80.7 ± 0.7
MR	HalfCheetah	45.5	36.6 ± 0.8	40.8 ± 0.6	39.6 ± 0.1	39.1 ± 0.1	39.0 ± 0.6
MR	Hopper	95.0	82.7 ± 7.0	87.3 ± 2.3	82.6 ± 4.6	86.1 ± 2.5	83.7 ± 8.2
MR	Walker2d	77.2	66.6 ± 3.0	66.8 ± 3.1	70.9 ± 4.3	73.4 ± 2.6	73.7 ± 1.4
	Mean	77.6	74.7	74.8	78.8	79.0	78.5

5 ABLATION STUDY

332

333

334

335

336

337

338

343

344

345

354 355

356

357

359

360

361

362 363 In this section, we answer the following research questions (RQs) to evaluate the impact of memory on model performance:

- 1. "How do the different components of RATE affect model performance in memory-intensive environments?"— RQ 1.
- 2. "What is the upper-bound estimate of the performance of the RATE model?"- RQ 2.
- 3. "Why do you need an MRV and what is its best configuration?"— RQ 3.

RQ 1. Investigating the impact of RATE components. To study the influence of memory embeddings on RATE model performance, we conducted the following experiment: for the RATE model trained for T-Maze (with context K = 30 on N = 3 segments), we replaced pre-trained memory embeddings M with random noise vectors during inference (see Figure 7).

Using random noise instead of memory embeddings, SR = 50% in the T-Maze, indicating the agent reaches the junction but turns in only one direction regardless of the clue. Thus, we conclude that clue information is in memory embeddings, while the rest of the actions are shaped by transformer parameters.

During RATE inference in ViZDoom-Two-Colors environment with
replacement of different components of memory RATE mechanisms
(memory embeddings and cached hidden states) with noise (see Figure 6), it is found that for this environment the caching of hidden



Figure 6: Results of replacing RATE memory tokens and cached hidden states with white noise during inference in ViZDoom-Two-Colors.

states of previous tokens has the largest contribution, because with its noise performance drop is the largest. Thus, in sparse environments, memory embeddings contribute the most, while in continuous environments, caching of hidden states of previous tokens is most impactful.

381

420

421

422

423 424

426

427

428

429

430

431

382 RQ 2. Performance upper-bound estimate. To evaluate the maximum possible performance of the RATE model, we 384 conducted experiments with OracleDT – a DT model whose 385 context is augmented as a pre- and post-fix with a vector v of 386 dimension $1 \times d$ model containing a priori 1-bit information about the environment. Thus, in the T-Maze environment, this 387 information is represented by a clue at the beginning of the 388 episode ($v_i = 0$ if clue = 0 else 1), and in the ViZDoom-Two-389 Colors environment, it is represented by a column color ($v_i = 0$ 390 if column color = red else 1). A context S' = concat(v, S, v)391



Figure 7: Results of replacing RATE memory tokens with white noise curing inference in T-Maze.

extended in this way can be interpreted as a context concat(M, S, M) with M memory embeddings added, trained perfectly and containing 100% of the important information. Thus, in environments where the a priori information about the environment needed for decision making can be extracted into a given vector v, the condition $R[OracleDT] \ge R[RATE] \ge R[DT]$ must be satisfied (see Table 5). This a priori information cannot be extracted from the environment in general, which further emphasizes the advantage of RATE, which is able to automatically extract important information and record it in memory embeddings M.

Table 5: Comparison of OracleDT with RATE. OracleDT determines the upper-bound estimate for the maximum reward that can be obtained using RATE in the environment. SR – Success Rate.

T-Maze							
	OracleDT	DT (Chen et al., 2021)	RATE				
	(K = 90)	(K = 90)	$(K_{eff} = 3 \times 30 = 90)$				
SR ($T = 90$)	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0				
SR ($T = 480$)	1.0 ± 0.0	0.5 ± 0.0	0.90 ± 0.07				
SR ($T = 900$)	1.0 ± 0.0	0.5 ± 0.0	0.90 ± 0.07				
	Vi	ZDoom-Two-Colors					
Reward[Total]	56.5 ± 0.8	24.8 ± 1.4	41.5 ± 1.0				
Reward[Reds]	55.3 ± 1.6	7.2 ± 0.4	38.2 ± 5.1				
Reward[Greens]	57.2 ± 0.5	42.3 ± 3.3	44.7 ± 5.8				

RO 3. Memory Retention Valve architecture ablation. Without MRV in the T-Maze environment 411 at corridor lengths of $L \gg K$, the performance of the RATE model decreased with each segment 412 processed at inference, resulting in almost SR = 50% on long trajectories (see Table 6). For example, 413 in the T-Maze task, the important information to be remembered goes into memory embeddings when 414 processing the first segment of the sequence, and then it must be retrieved when making decisions 415 on the last segment. At the same time, due to the recurrent structure of the architecture, memory 416 embeddings continue to be updated during the processing of intermediate segments when no new 417 information needs to be memorized, causing important information from memory embeddings to 418 leak out. To retain important information in memory embeddings, MRV mechanism was added to the 419 architecture. We considered the five different schemes detailed in subsection F.3 to implement MRV:

- 1. **MRV-CA-1**: cross-attention-based MRV which uses an attention mechanism to control the updating of memory embeddings. The updated memory embeddings M_{n+1} are fed to Query, and the incoming M_n are fed to Key and Value.
- 2. MRV-CA-2: uses the same mechanism as MRV-CA-1 but the incoming memory embeddings M_n are fed to Query, and the updated M_{n+1} are fed to Key and Value.
- 3. **MRV-G**: gated MRV which uses a gating mechanism similar to the one used in GTrXL (Parisotto et al., 2020).
- 4. **MRV-GRU**: uses a GRU (Chung et al., 2014) block to process updated M_n with hidden states.
- 5. **MRV-LSTM**: uses a LSTM (Hochreiter & Schmidhuber, 1997) block to process updated M_n with cached states.

434	Inference corridor length							
435	Model ($K_{eff} = 30 \times 5 = 150$)	150	360	600	900			
436	RATE w/o MRV [†]	1.00 ± 0.00	0.66 ± 0.08	0.65 ± 0.07	0.61 ± 0.07			
437	RATE (MRV-CA-2)	1.00 ± 0.00	0.95 ± 0.05	0.90 ± 0.07	0.90 ± 0.07			
100	RATE (MRV-G)	0.86 ± 0.07	0.77 ± 0.08	0.66 ± 0.07	0.65 ± 0.08			
430	RATE (MRV-GRU)	0.99 ± 0.01	0.74 ± 0.07	0.56 ± 0.11	0.55 ± 0.12			
439	RATE (MRV-LSTM)	0.85 ± 0.06	0.64 ± 0.10	0.51 ± 0.11	0.47 ± 0.11			
440	RATE (MRV-CA-1)	0.51 ± 0.01	0.51 ± 0.01	0.49 ± 0.02	0.49 ± 0.01			

Table 6: Results of ablation study of MRV configuration on T-Maze environment. † – baseline.

441

432

433 434

442 The best results in Table 6 were obtained using cross-attention scheme (MRV-CA-2), in which we 443 fed M_n memory tokens from the transformer input to the query and M_{n+1} memory tokens from 444 the transformer output to the key and value. This configuration is used throughout the work and is 445 denoted simply as MRV. This configuration acts as an effective gating mechanism to prevent the loss of important information in prolonged sparse environments, which is reflected in significantly better 446 results for RATE in the T-Maze environment. 447

448 449

450

465

RELATED WORK 6

TRANSFORMERS FOR REINFORCEMENT LEARNING 6.1 451

452 Transformers have found application in various areas of RL (Agarwal et al., 2023; Li et al., 2023): 453 online RL (Parisotto et al., 2020; Lampinen et al., 2021; Esslinger et al., 2022; Melo, 2022; Zheng et al., 2022; Pramanik et al., 2023; Team et al., 2023), offline RL (Chen et al., 2021; Janner et al., 454 2021; Lee et al., 2022; Reed et al., 2022; Jiang et al., 2023), and model-based RL (Chen et al., 2022; 455 Micheli et al., 2023; Robine et al., 2023). The use of transformers as a general policy for many 456 environments is also being explored (Lee et al., 2022; Melo, 2022; Reed et al., 2022). In our work, 457 we consider the formulation of an offline model-free RL, where learning is formulated as a sequence 458 modeling problem (Chen et al., 2021). Prominent representatives of such models are (Chen et al., 459 2021; Janner et al., 2021; Lee et al., 2022; Jiang et al., 2023), although planning in latent space (Jiang 460 et al., 2023) is considered, which can be seen as modeling of the environment. Moreover, (Janner 461 et al., 2021; Jiang et al., 2023) are specialized for control tasks with vector observations and do not 462 generalize to environments with observations in the form of images. Therefore, we consider the 463 Decision Transformer (Chen et al., 2021), which has no memory mechanism, as the main baseline for 464 comparison.

6.2 RECURRENT NEURAL NETWORKS FOR REINFORCEMENT LEARNING 466

467 For long input sequences, recurrent networks may have computational advantages over transformers. The RNNs recurrent unit maintains a hidden state, which is essentially a form of memory that is 468 important for solving POMDPs. In Decision LSTM (DLSTM) (Siebenborn et al., 2022) in DT the 469 transformer is replaced by an LSTM unit (Hochreiter & Schmidhuber, 1997). 470

471 6.3 STATE SPACE MODELS FOR REINFORCEMENT LEARNING 472

Recently, State Space Models (SSMs) (Gu et al., 2021) have shown significant success in sequence 473 modeling, particularly in offline RL (Bar-David et al., 2023; Cao et al., 2024; Gu & Dao, 2023; Ota, 474 2024). In Decision S4 (DS4) (Bar-David et al., 2023), sequence modeling is executed using S4 (Gu 475 et al., 2021) layers within the framework of offline RL, whereas Decision Mamba (DMamba) (Ota, 476 2024) utilizes the most recent Mamba (Gu & Dao, 2023) sequence model instead of causal self-477 attention. Mamba Decision Maker (MambaDM) (Cao et al., 2024) integrates the unique features of 478 SSMs to effectively combine local and global features with Global-local fusion Mamba (GLoMa) 479 module. 480

6.4 MEMORY IN TRANSFORMERS 481

482 There are many ways to implement the memory mechanism for transformers (Bulatov et al., 2022; Dai et al., 2019; Ding et al., 2020; Lei et al., 2020; Rae et al., 2019; Wu et al., 2020; 2022). In 483 Transformer-XL (TrXL) (Dai et al., 2019), it is proposed to split a long data sequence into segments 484 and to access past segments at the expense of memory, but to ignore very distant segments. This 485 increases the effective length of the context. The Compressive Transformer (Rae et al., 2019) uses

486 compressed memory, allowing previous versions of memory to be compressed rather than discarded 487 as in TrXL. ERNIE-Doc (Ding et al., 2020) suggests using the retrospective feed mechanism and the 488 enhanced recurrence mechanism. Memformer (Wu et al., 2020) uses external dynamic memory to 489 encode and retrieve past information. MART (Lei et al., 2020) extends this idea by adding a memory 490 update mechanism similar to a recurrent neural network (Cho et al., 2014; Hochreiter & Schmidhuber, 1997). The Memorizing Transformer (Wu et al., 2022) proposes to store the internal representations 491 of past inputs. The Recurrent Memory Transformer (RMT) (Bulatov et al., 2022) includes additional 492 read and write memory tokens at each segment's beginning and end. This method allows the effective 493 context to be expanded to over 1 million tokens (Zhu et al., 2020). <u>191</u>

495 An Adaptive Agent (AdA) (Bauer et al., 2023) uses memory architectures to store and employ infor-496 mation previously acquired by the agent. The default memory architecture is TrXL with normalization before each layer (Parisotto et al., 2020), and the use of gating on the feedforward layers (Shazeer, 497 2020) to stabilize training. We also use TrXL in our work but refrain from using additional modifica-498 tions to stabilize training. Another distinctive feature of using a transformer in AdA, as opposed to 499 DT, is that pixel observations, past actions, past rewards, and additional information are not tokenized 500 separately but are combined into a single vector that feeds the transformer. The transformer itself 501 predicts not only actions but also value function values. 502

503 504

7 CONCLUSION

505 506

In this paper, we propose **Recurrent Action Transformer with Memory (RATE)**, a transformer 507 model for offline RL that exploits memory mechanisms in the form of memory embeddings and 508 caching of hidden states of previous tokens, and the Memory Retention Valve (MRV), which 509 controls memory updating and prevents the loss of important information in sparse tasks. In extensive 510 experiments in memory-intensive environments such as ViZDoom-Two-Colors, Memory Maze, 511 Minigrid-Memory, and T-Maze, we have shown that RATE outperforms recurrent and transformer 512 baselines. The proposed model interpolates and extrapolates well outside the transformer context, is 513 able to retain important information for a long time when operating in highly sparse environments, 514 and allows to compensate for the effect of bias in the training data.

We also show that the proposed model achieves better or comparable results to state-of-the-art Mamba-based models in Atari and MuJoCo environments, indicating that RATE is suitable for all tasks: both memory-intensive and not. We have thoroughly investigated the influence of memory mechanisms on the performance of the model and have clearly shown that the model uses them in decision making. This method shows great potential for tackling complex tasks with long sequences, especially in robotics, where training agents on pre-collected data sets is highly advantageous.

521 522

523

524

525

526

Limitations. Limitations of the proposed model include its inability to design K and N in memoryintensive environments so that all important events fall into the efficient context of $K_{eff} = K \times N$. In addition, the approach based on dividing trajectories into segments during inference does not allow for memory updates effectively at each step using a sliding window. Also, there are currently no studies of the memory capacity of the proposed model, so the practical amount of information that can be stored remains unknown.

527 528 529

530

531

532

Reproducibility Statement. The model description is presented in section 3 (Algorithm 1 and Algorithm 2), the training procedure is presented in Appendix D, the description of the used benchmarks is presented in Appendix B, the hyperparameters are presented in Table 8, and the configurations for displaying the experimental results are presented in Table 9. The results of the hyperparameters tuning for recurrent baselines are presented in Appendix G.

- 533 534 535
- 536 REFERENCES
- 537

Pranav Agarwal, Aamer Abdul Rahman, Pierre-Luc St-Charles, Simon JD Prince, and
Samira Ebrahimi Kahou. Transformers in reinforcement learning: a survey. *arXiv preprint arXiv:2307.05979*, 2023. 556

558

559

563

564

565 566

567

568

575

576

577

585

586

588

589

- Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic perspective on offline reinforcement learning. In *International Conference on Machine Learning*, pp. 104–114. PMLR, 2020.
- 544 Shmuel Bar-David, Itamar Zimerman, Eliya Nachmani, and Lior Wolf. Decision s4: Efficient 545 sequence-based rl via state spaces layers. *arXiv preprint arXiv:2306.05167*, 2023.
- 546 Jakob Bauer, Kate Baumli, Feryal Behbahani, Avishkar Bhoopchand, Nathalie Bradley-Schmieg, 547 Michael Chang, Natalie Clay, Adrian Collister, Vibhavari Dasagi, Lucy Gonzalez, Karol Gregor, 548 Edward Hughes, Sheleem Kashem, Maria Loks-Thompson, Hannah Openshaw, Jack Parker-549 Holder, Shreya Pathak, Nicolas Perez-Nieves, Nemanja Rakicevic, Tim Rocktäschel, Yannick 550 Schroecker, Satinder Singh, Jakub Sygnowski, Karl Tuyls, Sarah York, Alexander Zacherl, and 551 Lei M Zhang. Human-timescale adaptation in an open-ended task space. In Andreas Krause, 552 Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett 553 (eds.), Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pp. 1887–1935. PMLR, 23–29 Jul 2023. URL 554 https://proceedings.mlr.press/v202/bauer23a.html. 555
 - Edward Beeching, Christian Wolf, Jilles Dibangoye, and Olivier Simonin. Deep reinforcement learning on a budget: 3d control and reasoning without a supercomputer. *CoRR*, abs/1904.01806, 2019. URL http://arxiv.org/abs/1904.01806.
- Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47: 253–279, 2013.
 - Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
 - Aydar Bulatov, Yury Kuratov, and Mikhail Burtsev. Recurrent memory transformer. *Advances in Neural Information Processing Systems*, 35:11079–11091, 2022.
- Jiahang Cao, Qiang Zhang, Ziqing Wang, Jiaxu Wang, Hao Cheng, Yecheng Shao, Wen Zhao, Gang
 Han, Yijie Guo, and Renjing Xu. Mamba as decision maker: Exploring multi-scale sequence
 modeling in offline reinforcement learning. *arXiv preprint arXiv:2406.02013*, 2024.
- 572
 573
 574
 Chang Chen, Yi-Fu Wu, Jaesik Yoon, and Sungjin Ahn. Transdreamer: Reinforcement learning with transformer world models. *arXiv preprint arXiv:2202.09481*, 2022.
 - Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021.
- Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo de Lazcano, Lucas Willems, Salem Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & miniworld: Modular & customizable reinforcement learning environments for goal-oriented tasks. *CoRR*, abs/2306.13831, 2023.
- Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of
 neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*, 2014.
 - Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
 - Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- Jiayu Ding, Shuming Ma, Li Dong, Xingxing Zhang, Shaohan Huang, Wenhui Wang, Nanning
 Zheng, and Furu Wei. Longnet: Scaling transformers to 1,000,000,000 tokens. arXiv preprint
 arXiv:2307.02486, 2023.

594 595 596	Siyu Ding, Junyuan Shang, Shuohuan Wang, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. Ernie- doc: A retrospective long-document modeling transformer. <i>arXiv preprint arXiv:2012.15688</i> , 2020.
597 598	Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. R12: Fast
599	reinforcement learning via slow reinforcement learning. arXiv preprint arXiv:1011.02/79, 2016.
600 601	Kevin Esslinger, Robert Platt, and Christopher Amato. Deep transformer q-networks for partially observable reinforcement learning. <i>arXiv preprint arXiv:2206.01078</i> , 2022.
602	
604	data-driven reinforcement learning, 2021.
605 606 607 608	Jake Grigsby, Linxi Fan, and Yuke Zhu. AMAGO: Scalable in-context reinforcement learning for adaptive agents. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. URL https://openreview.net/forum?id=M6XWoEdmwf.
609 610	Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. <i>arXiv</i> preprint arXiv:2312.00752, 2023.
611 612 613	Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. <i>arXiv preprint arXiv:2111.00396</i> , 2021.
614 615	Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. <i>arXiv preprint arXiv:1912.01603</i> , 2019.
616	
617 618	Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. <i>Neural computation</i> , 9(8): 1735–1780, 1997.
619 620 621	Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big sequence modeling problem. <i>Advances in neural information processing systems</i> , 34:1273–1286, 2021.
622 623 624 625	Zhengyao Jiang, Tianjun Zhang, Michael Janner, Yueying Li, Tim Rocktäschel, Edward Grefenstette, and Yuandong Tian. Efficient planning in a compact latent action space. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=cA77NrVEuqn.
626 627 628	Jikun Kang, Romain Laroche, Xindi Yuan, Adam Trischler, Xue Liu, and Jie Fu. Think before you act: Decision transformers with internal working memory. <i>arXiv preprint arXiv:2305.16338</i> , 2023.
629 630	Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning, 2020. URL https://arxiv.org/abs/2006.04779.
631 632 633 634	Andrew Lampinen, Stephanie Chan, Andrea Banino, and Felix Hill. Towards mental time travel: a hierarchical memory for reinforcement learning agents. <i>Advances in Neural Information Processing Systems</i> , 34:28182–28195, 2021.
635 636 637 638	Kuang-Huei Lee, Ofir Nachum, Mengjiao Sherry Yang, Lisa Lee, Daniel Freeman, Sergio Guadar- rama, Ian Fischer, Winnie Xu, Eric Jang, Henryk Michalewski, et al. Multi-game decision transformers. <i>Advances in Neural Information Processing Systems</i> , 35:27921–27936, 2022.
639 640 641	Jie Lei, Liwei Wang, Yelong Shen, Dong Yu, Tamara L Berg, and Mohit Bansal. Mart: Memory- augmented recurrent transformer for coherent video paragraph captioning. <i>arXiv preprint</i> <i>arXiv:2005.05402</i> , 2020.
642 643 644 645 646	Wenzhe Li, Hao Luo, Zichuan Lin, Chongjie Zhang, Zongqing Lu, and Deheng Ye. A survey on transformers in reinforcement learning. <i>Transactions on Machine Learning Research</i> , 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=r30yuDPvf2. Survey Certification.
647	Luckeciano C Melo. Transformers are meta-reinforcement learners. In International Conference on Machine Learning, pp. 15340–15359. PMLR, 2022.

648 649 650	Vincent Micheli, Eloi Alonso, and François Fleuret. Transformers are sample-efficient world models. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=vhFulAcb0xb.
652 653	Volodymyr Mnih. Playing atari with deep reinforcement learning. <i>arXiv preprint arXiv:1312.5602</i> , 2013.
654 655 656 657	Tianwei Ni, Michel Ma, Benjamin Eysenbach, and Pierre-Luc Bacon. When do transformers shine in RL? decoupling memory from credit assignment. In <i>Thirty-seventh Conference on Neural</i> <i>Information Processing Systems</i> , 2023. URL https://openreview.net/forum?id= APGXBNkt6h.
658 659 660	Toshihiro Ota. Decision mamba: Reinforcement learning via sequence modeling with selective state spaces. <i>arXiv preprint arXiv:2403.19925</i> , 2024.
661 662 663 664	Emilio Parisotto, Francis Song, Jack Rae, Razvan Pascanu, Caglar Gulcehre, Siddhant Jayakumar, Max Jaderberg, Raphael Lopez Kaufman, Aidan Clark, Seb Noury, et al. Stabilizing transformers for reinforcement learning. In <i>International conference on machine learning</i> , pp. 7487–7498. PMLR, 2020.
665 666 667	Jurgis Pasukonis, Timothy Lillicrap, and Danijar Hafner. Evaluating long-term memory in 3d mazes. <i>arXiv preprint arXiv:2210.13383</i> , 2022.
668 669 670	Marco Pleines, Matthias Pallasch, Frank Zimmer, and Mike Preuss. Transformerxl as episodic memory in proximal policy optimization. <i>Github Repository</i> , 2023. URL https://github. com/MarcoMeter/episodic-transformer-memory-ppo.
671 672	Subhojeet Pramanik, Esraa Elelimy, Marlos C Machado, and Adam White. Recurrent linear transformers. <i>arXiv preprint arXiv:2310.15719</i> , 2023.
674 675	Jack W Rae, Anna Potapenko, Siddhant M Jayakumar, and Timothy P Lillicrap. Compressive transformers for long-range sequence modelling. <i>arXiv preprint arXiv:1911.05507</i> , 2019.
676 677 678	Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. A generalist agent. <i>arXiv preprint arXiv:2205.06175</i> , 2022.
679 680 681 682	Jan Robine, Marc Höftmann, Tobias Uelwer, and Stefan Harmeling. Transformer-based world models are happy with 100k interactions. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=TdBaDGCpjly.
683 684	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> , 2017.
685 686	Noam Shazeer. Glu variants improve transformer. arXiv preprint arXiv:2002.05202, 2020.
687 688 689	Max Siebenborn, Boris Belousov, Junning Huang, and Jan Peters. How crucial is transformer in decision transformer? <i>arXiv preprint arXiv:2211.14655</i> , 2022. URL https://arxiv.org/abs/2211.14655.
690 691 692 693	Artyom Sorokin, Nazar Buzun, Leonid Pugachev, and Mikhail Burtsev. Explain my surprise: Learning efficient long-term memory by predicting uncertain outcomes. 07 2022. doi: 10.48550/arXiv.2207. 13649.
694 695	R.S. Sutton and A.G. Barto. <i>Reinforcement Learning, second edition: An Introduction</i> . Adaptive Computation and Machine Learning series. MIT Press, 2018. ISBN 9780262039246.
696 697 698 699	Adaptive Agent Team, Jakob Bauer, Kate Baumli, Satinder Baveja, Feryal Behbahani, Avishkar Bhoopchand, Nathalie Bradley-Schmieg, Michael Chang, Natalie Clay, Adrian Collister, et al. Human-timescale adaptation in an open-ended task space. <i>arXiv preprint arXiv:2301.07608</i> , 2023.
700 701	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. <i>Advances in neural information processing systems</i> , 30, 2017.

702 703 704	Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos, Charles Blundell, Dharshan Kumaran, and Matt Botvinick. Learning to reinforcement learn. <i>arXiv</i> preprint arXiv:1611.05763, 2016.
705 706 707	Qingyang Wu, Zhenzhong Lan, Jing Gu, and Zhou Yu. Memformer: The memory-augmented transformer. <i>arXiv preprint arXiv:2010.06891</i> , 2020.
708 709	Yuhuai Wu, Markus N Rabe, DeLesley Hutchins, and Christian Szegedy. Memorizing transformers. arXiv preprint arXiv:2203.08913, 2022.
710 711 712 713	Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for longer sequences. <i>Advances in neural information processing systems</i> , 33:17283–17297, 2020.
714 715 716	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> , 2022.
717 718 719	Qinqing Zheng, Amy Zhang, and Aditya Grover. Online decision transformer. In <i>international conference on machine learning</i> , pp. 27042–27059. PMLR, 2022.
720 721 722	Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. <i>arXiv preprint arXiv:2010.04159</i> , 2020.
723	
724	
725	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
730	
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
754	
755	

A DECISION TRANSFORMER

757

758 Decision Transformer (DT) (Chen et al., 2021) 759 is an algorithm for offline RL that reduces the 760 RL task to a sequence modeling task. In DT, the 761 scheme of which is presented in Algorithm 3, 762 the trajectory τ is not divided into segments as in RATE. Instead, random fragments of length 764 K are sampled from the trajectory, since originally this architecture was designed to work only 765 with MDP. The predicted actions \hat{a} are sampled 766 autoregressively. 767

Algorithm 3 Decision Transformer
Input : $R \in \mathbb{R}^{1 \times T}, o \in \mathbb{R}^{d_o \times T}, a \in \mathbb{R}^{1 \times T}$
1: $ ilde{R} \in \mathbb{R}^{T imes d} \leftarrow \texttt{Encoder}_R(R)$
$\tilde{o} \in \mathbb{R}^{T imes d} \leftarrow \texttt{Encoder}_o(o)$
$\tilde{a} \in \mathbb{R}^{T imes d} \leftarrow \texttt{Encoder}_a(a)$
2: $\tau_{0T} \leftarrow \{ (\tilde{R}_0, \tilde{o}_0, \tilde{a}_0), \dots, (\tilde{R}_T, \tilde{o}_T, \tilde{a}_T) \}$
3: $n = random(0, T - K)$
4: $\hat{a}_n \leftarrow \texttt{Transformer}(\tau_{nn+K})$
Output : $\hat{a}_n \to \mathcal{L}(a_n, \hat{a}_n)$

B ENVIRONMENTS

771 B.1 MEMORY-INTENSIVE ENVIRONMENTS

In this section, we provide an extended description of the environments used in this paper, as well as
the methodology used to collect the trajectories. Table 7 summarizes the observations type, rewards
type, and actions type for each of the environments considered in this paper.

776

768 769

770

ViZDoom-Two-Colors. We used a modified ViZDoom-Two-Colors environment from (Sorokin 777 et al., 2022) to assess the model's memory abilities. The agent initially having 100 hit points (HP) is 778 placed in a room without inner walls filled with acid. At each step in the environment, the agent loses 779 a fixed amount of health (10/32 HP per step). In the center of the environment, there is a pillar of either green or red color, which disappears after 45 environment steps. Throughout the environment, 781 objects of two colors (green and red) are generated. When the agent interacts with an object of the 782 same color as the pillar, it gains an increase in health of +25 and a reward of +1. When the agent 783 interacts with an object of the opposite color, it loses a similar amount of health. The agent receives 784 an additional reward of +0.02 for each step it survives. The episode ends when the agent has zero 785 health. Thus, the agent needs to remember the color of the pillar to select items of the correct color, 786 even if the pillar is out of sight or has disappeared. The agent does not receive information about its 787 current health or rewards, as these observations essentially convey the same information as the color of the pillar but persist beyond step 45. 788

789 We collected a dataset of 5000 trajectories of 90 steps in length using a trained A2C (Beeching et al., 790 2019) agent (an agent trained with a non-disappearing pillar). The average reward for these 90 steps 791 is 4.46. When collecting trajectories, to ensure that the agent saw the pillar before it disappeared, the agent always appeared facing the pillar in the same place – midway between the pillar and the 793 nearest wall. In order to successfully complete this task, the agent needs to remember the color of the pillar. This environment tests the long-term memory mechanism, since the agent needs to 794 retain information about the pillar for a time much longer than the pillar has been in the environment. 795 Using only short-term memory and, for example, collecting the next item of the same color as the 796 previous collected item, it will not be possible for the agent to survive for a long time, as this policy 797 is extremely unstable. This is due to the fact that in the training dataset the agent occasionally makes 798 a mistake and picks up an object of the opposite color. Thus, irrelevant information about the desired 799 color may enter the transformer context and the agent will start collecting items of an opposite color, 800 which will quickly lead to a failure. 801

802 **T-Maze.** To investigate agent's long-term memory on very long environments (the inference 803 trajectory length is much longer than the effective context length K_{eff} we used a modified version of 804 the T-Maze environment (Ni et al., 2023). The agent's objective in this environment is to navigate from 805 the beginning of the T-shaped maze to the junction and choose the correct direction, based on a signal 806 given at the beginning of the trajectory using four possible actions $a \in \{left, up, right, down\}$. This 807 signal, represented as the *clue* variable and equals to zero everywhere except the first observation, dictates whether the agent should turn up (clue = 1) or down (clue = -1). Additionally, a constraint 808 on the episode duration T = L + 2, where the maximum duration is determined by the length of the corridor L to the junction, adds complexity to the problem. To address this, a binary flag, represented

810 as the flaq variable, which is equal to 1 one step before the junction and 0 otherwise, indicating 811 the arrival of the agent at the junction, is included in the observation vector. Additionally, a noise 812 channel is added to the observation vector, with random integer values from the set $\{-1, 0, +1\}$. 813 The observation vector is thus defined as o = [y, clue, flag, noise], where y represents the vertical 814 coordinate. The reward r is given only at the end of the episode and depends on the correctness of the agent's turn at the junction, being 1 for a correct turn and 0 otherwise. This formulation deviates 815 from the traditional Passive T-Maze environment (Ni et al., 2023) (different observations and reward 816 functions) and presents a more intricate set of conditions for the agent to navigate and learn within 817 the given time constraint. 818

The dataset consists of 2000 of trajectories for each segment of length 30 (i.e. 6000 trajectories for the $K_{eff} = 3 \times 30 = 90$) and consists only of successful episodes. An artificial oracle with a priori information about the environment was used to generate the dataset.

Environment	Obs. type	Rewards	Actions	Obs. info
ViZDoom-Two-Colors	Image	Continuous	Discrete	First-person view
T-Maze	Vector	Sparse & Discrete	Discrete	Vector
Memory Maze	Image	Sparse & Discrete	Discrete	First-person view
Minigrid-Memory	Image	Sparse	Discrete	Observes the 3×3 part of the grid
Action Associative Retrieval	Vector	Sparse & Discrete	Discrete	Vector
Atari	Image	Sparse & Discrete	Discrete	Observes the full game screen
MuJoCo	Vector	Continuous	Continuous	Vector

Table 7: Description of observations and reward functions for the considered environments.

835 836 837

822

823 824 825

838 **Memory Maze.** In this first-person view 3D environment (Pasukonis et al., 2022), the agent appears 839 in a randomly generated maze containing several objects of different colors at random locations. The 840 agent's task is to find an object of the same color in the maze as the outline around its observation 841 image. After the agent finds an object of the desired color and steps on it, the color of the outline 842 changes and the agent must find another object. The agent receives a +1 reward for stepping on the 843 correct object. Otherwise, it receives no reward. The duration of an episode is a fixed number and is equal to 1000. Thus, the agent's task is to find as many objects of the desired color as possible 844 in a limited time. The agent's effectiveness in this environment depends on its ability to memorize 845 the structure of the maze and the location of objects in it in order to find the desired objects faster. 846 Using the Dreamer model (Hafner et al., 2019) to collect dataset of 5000 trajectories only achieved 847 an average award of 4.7 per episode, i.e., a rather sparse dataset. 848

Minigrid-Memory. Minigrid-Memory (Chevalier-Boisvert et al., 2023) is a 2D grid environment 850 designed to test an agent's long-term memory and credit-assignment (Ni et al., 2023). The envi-851 ronment map is a T-shaped maze with a small room with an object inside it at the beginning of the 852 corridor. The agent appears at a random coordinate in the corridor. The agent's task is to reach the 853 room with the object and memorize it, then reach the junction at the end of the maze and make a turn 854 in the direction where the same object is located as in the room at the beginning of the maze. A reward 855 $r = 1 - 0.9 \times \frac{t}{T}$ is given for success, and 0 for failure. The episode ends after any agent turns at a 856 junction or after a limited amount of time (95 steps) has elapsed. The agent's observations are limited 857 to a 3×3 size frame. 10000 trajectories with grid size 31x31 were collected using PPO (Schulman et al., 2017) with TransformerXL (Pleines et al., 2023) with a context length equal to the maximum 858 episode duration. 859

860 861

862

849

B.2 STANDARD BENCHMARKS

Atari games. For the Atari game environments (Bellemare et al., 2013), we used the same dataset as in DT, namely the DQN replay dataset with grayscale state images (Agarwal et al., 2020). This

dataset contains 500 thousand of the 50 million steps of an online DQN (Mnih, 2013) agent for each game. We use the following set of games: SeaQuest, Breakout, Pong and Qbert.

867 **MuJoCo.** Despite the fact that memory is not required in decision making in control environments 868 like MuJoCo (Fu et al., 2021), we conducted additional experiments in this environment to compare with DT. For the continuous control tasks, we selected a standard MuJoCo locomotion environment and a set of trajectories from the D4RL benchmark (Fu et al., 2021). Since we chose DT and TAP 870 as the main models for comparison on this data, we focused on the environments used in both 871 works (HalfCheetah, Hopper, and Walker). We used three different dataset settings: 1) Medium – 872 1 million timesteps generated by a "medium" policy that achieves about a third of the score of an 873 expert policy; 2) Medium-Replay – the replay buffer of an agent trained with the performance of a 874 medium policy (about 200k-400k timesteps in our environments); 3) Medium-Expert – 1 million 875 timesteps generated by the medium policy concatenated with 1 million timesteps generated by an 876 expert policy. The scores for the MuJoCo experiments are normalized such that 100 represents an 877 expert policy, following the benchmark protocol outlined in (Fu et al., 2021). The performance 878 metrics for Conservative Q-Learning (CQL) and Trajectory Autoencoding Planner (TAP) are reported 879 from the TAP paper (Jiang et al., 2023), and for DT from the DT paper (Chen et al., 2021), as they use the same dataset and evaluation protocol. 880

881 882

883

906 907 908

C ACTION ASSOCIATIVE RETRIEVAL

As shown in section 4.2, DT has a SR = 50% for inference at corridor lengths longer than the transformer context length. This is due to the fact that even a DT trained on balanced data has a slight bias in the predicted probability towards one of the two required actions, which leads to the fact that when t > K the agent constantly produces only one action: up or down. In turn, the presence of memory in the agent allows us to combat this problem.



Figure 8: Action Associative Retrieval.

To check how the agent's performance changes during training, we design an Action Associative
 Retrieval (AAR) Figure 8 environment.

894 There are two states in this environment: S_0 and S_1 . The agent appears in state S_0 and by performing 895 the action $a_0 \in \{0, 1\}$ moves to state S_1 . Next, the agent must take N-2 steps to move from state 896 S_1 to state S_1 by performing action a = 2 (no op.). At the end of the episode, the agent must perform 897 the same action that moved it from state S_0 to state S_1 in order to move from state S_1 to state S_0 . Thus, the action $a \in \{0, 1, 2\}$. Agent observations o = [state, flag, noise], where $state \in \{0, 1\}$ 899 is the index of the current state, $flag \in \{0,1\}$ is a flag equal to 1 in case the next step requires returning to the initial state and equal to 0 otherwise, $noise \in \{-1, 0, +1\}$ is the noise channel. The 900 agent receives a +1 reward if it returns to the initial state S_0 by performing the action that took it out 901 from the S_0 to the S_1 , and -1 in other cases. The training dataset consists of oracle-generated 6000 902 trajectories with positive reward. 903

More formally, we can talk about the presence of memory in an agent when solving AAR (T-Maze-like) tasks under the condition that:

$$\forall t > K : \frac{1}{N_0} \sum_{i=1}^{N_0} p_i(a_t = a^0 | a_0 = a^0) + \frac{1}{N_1} \sum_{i=1}^{N_1} p_i(a_t = a^1 | a_0 = a^1) > 1$$
(1)

This condition means that if the agent has memory, the sum of the average conditional probabilities over all experiments will be greater than one, i.e., these probabilities are independent of each other.
Provided that the sum of these probabilities is less than or equal to one, the agent will choose at best the same target action in most experiments, even if another action is required.

where $a^0, a^1 \in \mathcal{A}$ – two mutually exclusive actions leading to a reward; t is the step at which the final action is required; N_0, N_1 are the number of experiments in environments where target action $a_t = a^0$ and $a_t = a^1$, respectively.

In the results Figure 9, the first 1% of training steps was removed because it corresponds to the beginning of the training and is unrepresentative. Blue dots correspond to the beginning of training,



Figure 9: Experimental results with RATE and DT in the AAR environment. The graphs show the 10-runs average results of training on trajectories of length T = 90 and validation on trajectories of length T = 180, for RATE with $K_{eff} = 3 \times 30 = 90$ and for DT with K = 90.

red dots to the end of training. As can be seen from Figure 9, during training, the probabilities $p_i(a_t = a^0|a_0 = a^0)$ and $p_i(a_t = a^1|a_0 = a^1)$ on the training trajectories have a strong positive correlation ($R_{train}^{DT} = 1.00$ and $R_{train}^{RATE} = 0.97$), where R – correlation coefficient. This indicates that within-context (effective context) DT and RATE models are able to predict both a^0 and a^1 actions equally well.

At the same time, during validation, for the RATE model this pattern is preserved – the red points corresponding to the probabilities of choosing actions a^0 and a^1 are in the upper right part of the graph, positive correlation persists ($R_{val}^{RATE} = 0.80$). On the other hand, in the DT case, the cluster of red dots is skewed toward choosing action a^1 and action a^0 with equal probabilities equal to 0.5. Thus, in sum, these probabilities are less or equal to one, as evidenced by a strong negative correlation $(R_{val}^{DT} = -0.97)$. The results confirm the inability of DT to generalize on trajectories whose lengths exceed the context length and the ability of RATE to handle such tasks.

D TRAINING

This section provides additional details on the training process of the baselines considered in the paper.
It is important to note that when training RATE in the transformer decoder the feed-forward network
block was disabled, because without it on some environments the training results are slightly better.
However, other transformer-based baselines were trained with the standard transformer decoder.

970

963 964

965

947

948

949 950

971 **ViZDoom-Two-Colors.** Since the pillar disappears at time t = 45, all trajectories start at time t = 0 and end at time t = 90 so that the information about the color of the pillar is guaranteed to

be used in training. In this experiment, we compared DT with context length K = 90 to RATE, RMT, and TrXL models with context length K = 30 and partitioning the trajectory into N = 3segments. Thus, when trained, RATE also handles sequences of length 90, since its effective context is $K_{eff} = N \times K = 90$, but only processed subsequences of length K = 30.

976

T-Maze. The model names are written in the format MODEL-N, where N is the number of segments of length K = 30 into which the training trajectories can be partitioned. Thus, DT-3 was trained on trajectories of length $T \le 3 \times 30 = 90$ with context length. RATE-3 was trained on similar trajectories as DT-3, but with each trajectory divided into 3 segments, during training, enabling the training of a model with a context length of K = 30 on trajectories of length T = 90. All the trajectories used in training start from t = 0, i.e., from the moment of receiving a clue.

984 Memory Maze. To train RATE, DT, RMT, and TrXL on Memory Maze, the same approach was 985 used as for ViZDoom-Two-Colors environment, except that trajectories were sampled not from t = 0986 but from $t : \sum_{t'=t}^{t+90} r_{t'} \ge 2$.

As in the ViZDoom-Two-Colors case, training for DT was performed with a context length of K = 90and for RATE, RMT, and TrXL with a context length of K = 30 and number of segments N = 3, i.e., effective context length $K_{eff} = N \times K = 3 \times 30 = 90$.

990

983

991 **Minigrid-Memory.** To train RATE, DT, RMT, and TrXL in this environment, trajectories were 992 sampled in the same manner as for T-Maze. An environment configuration with a maze of size 31x31993 was used as a training configuration. Since the maximum episode duration is 95, training proceeded 994 in the following setting: for DT the context length K = 30, for RATE, RMT, and TrXL the context 995 length K = 10 and the number of segments N = 3. All trajectories, as in T-Maze, are sampled from 996 time t = 0.

997

Atari and MuJoCo. When training RATE on Atari games and MuJoCo control tasks, sequences of length T = 90 (Atari) and T = 60 (MuJoCo) were sampled randomly from the original trajectories in the dataset. These trajectories were then divided into N = 3 segments of length K = 30 (Atari) and K = 20 (MuJoCo), forming an effective context of length $K_{eff} = N \times K = 90$ (60 for MuJoCo).

1002 For Atari, we used the identical experimental design described in the DT paper (Chen et al., 2021). 1003 It is worth noting that we presented raw scores for Atari, rather than gamer-normalized scores as 1004 described in the DT paper. Table 3 shows the results for Atari environments. RATE outperforms 1005 DT significantly in environments like Breakout and Qbert. We attribute this to the observation that, although these environments do not explicitly demand memory, intricate dynamics from the past exert a greater influence on agent behavior than in environments such as SeaQuest. Actions executed in the 1007 past notably alter the present state of the environment in Breakout and Qbert, whereas in SeaQuest, 1008 such actions hold little significance. For instance, the emergence of enemies and divers in SeaQuest 1009 is entirely independent of the agent's prior actions. 1010

For MuJoCo, our findings suggest that the conventional strategy of utilizing return is not suitable 1011 for our segment-based scheme. The issue arises during the trajectory, where the agent's return 1012 persistently diminishes. However, the true value of the agent's state at the onset and conclusion of the 1013 episode could remain unchanged, provided the agent's policy performs consistently well. To rectify 1014 this discrepancy, we propose a novel evaluation strategy for MuJoCo tasks. In this approach, each 1015 segment commences with the maximum return, simulating the scenario where the agent initiates the 1016 trajectory anew. This method effectively mitigates the aforementioned issue, enhancing the accuracy 1017 of our evaluation process. Our MuJoCo experiments in Table 4 show that this benefits performance 1018 significantly for some environments. Thus, using RATE allowed us to obtain the best metrics for 1019 MuJoCo in 4/9 cases compared to the other baselines. RATE also outperforms DT in 9/9 tasks.

1020 1021

1022 E RESULTS PRESENTATION

1023

1024 This section provides information on how the presented experimental results were obtained. N_{runs} 1025 denotes the number of model runs; N_{seeds} denotes the number of inference episodes with different seeds; sem denotes standard error of the mean, and std denotes standard deviation.

28 29	Hyperparameter	ViZDoom2C / Memory Maze	T-Maze / Minigrid-Memory	Atari	MuJoCo
30	Number of layers	6	8	6	3
24	Number of attention heads	8	10	8	1
	Embedding dimension	128 / 64	64	128	128
2	Context length K	30	30 / 10	30	20
	Number of segments	3	3	3	3
	Hidden dropout	0.2 / 0.5	0.05 / 0.2	0.2	0.2
	Attention dropout	0.05 / 0.2	0 / 0.05	0.05	0.05
	Number of memory tokens	5/15	5/15	15	15
	Number of cached tokens (mem_len)	300 / 360	0 / 180	360	2
	Max epochs	100 / 80	50 / 250	10	10
	Batch size	64	64	128	4096
	Weight decay	0.1	0.1	0.1	0.1
	Loss function	CE	CE	CE	MSE
	Optimizer	AdamW	AdamW	AdamW	AdamW
	MRV activation	ReLU	ReLU	ReLU [‡]	ReLU
	MRV number of attention heads	2/4	4 / 1	2	2
	Learning rate	3e-4	3e-4	3e-4	1e-4
	Adam $W(\beta_1,\beta_2)$	(0.9, 0.95)	(0.9, 0.95)	(0.9, 0.95)	(0.9, 0.95)

1026	Table 8: RATE hyperparameters for different experiments. ‡ – Leaky ReLU in the Atari.Pong case.
1027	

Table 9: Experimental parameters used to present the final results.

Environment	N_{runs}	N_{seeds}	Metric	Notation
ViZDoom-Two-Colors	6	100	Total reward	$mean \pm sem$
T-Maze	10	100	Success Rate	$mean \pm sem$
Memory Maze	3	100	Total reward	mean \pm sem
Minigrid-Memory	3	100	Total reward	mean \pm sem
Action Associative Retrieval	10		Success Rate	mean \pm sem
Atari	3	100	Total reward	mean \pm std
MuJoCo	3	100	Total reward	mean \pm std

1056 1057

1045 1046

1058

1060

1059 F ADDITIONAL ABLATION STUDIES

To determine the optimal hyperparameters associated with memory mechanisms, additional ablation studies were performed in ViZDoom-Two-Colors and T-Maze environments, and the results are presented in Figure 11 and Figure 10 (right). From the ablation studies results, it was found that for environments like ViZDoom-Two-Colors with continuous reward signal and image observations, the best results can be obtained using number of cached memory tokens mem_len = $(K \times 3 + 2 \times$ num_mem_tokens) $\times N$, where K - context length and N - number of segments.

1067 On the other hand, for environments with sparse events like T-Maze, it has been found that using 1068 caching of hidden states of previous tokens (mem_len > 0) prevents remembering important 1069 information. In this case, gating with n_head_ca = 4 and moderate number of memory tokens 1070 num_mem_tokens = 5 gives the best results (see Figure 10 (right)).

1071

1072 F.1 Additional ViZDoom-Two-Colors ablation

1074 The effect of combining of memory tokens with noise is shown in Figure 10 (left). The noise was 1075 applied as a convex combination: memory_tokens = $(1 - \alpha) \times \text{memory}_t\text{okens} + \alpha \times \text{noise}$. 1076 With unchanged caching of hidden states of previous tokens at growth of the noise parameter α , at 1077 first there is a decrease of performance at inference on green pillars (up to $\alpha = 0.5$), and only then a 1078 decrease of performance at inference on red pillars. This phenomenon can be explained by the fact 1079 that memory embeddings is trained to record mostly information about red pillars, which helps to 1079 combat bias in the training data.



1089Figure 10: (left) Investigating the RATE memory tokens noise effect in the ViZDoom-Two-Colors.1090(right) Results of RATE-3 (trained on corridor lengths \leq 90) ablation studies in the T-Maze environ-1091ment. n_head_ca - number of MRV attention heads, num_mem_tokens - number of memory1092tokens.



Figure 11: Results of RATE ablation studies in the ViZDoom-Two-Colors environment.

1118 F.2 CURRICULUM LEARNING

1115 1116 1117

Since in the T-Maze environment, the number of actions at the junction relates to the number of actions when moving straight along the corridor as $\frac{1}{L}$ and tends to 0 as L increases, there is a significant imbalance in the agent's action distribution, which can cause problems when performing rare class (turning actions) prediction. Theoretically, this situation can be remedied through curriculum learning.

1124 Curriculum learning (CL) is a technique in which a model is trained on examples of increasing 1125 difficulty. In this approach, the model is first trained on the set of trajectories $Q_1 = q_1$ of length 1126 $K \times 1$, then the trained model is re-trained on the set of trajectories $Q_2 = q_1 \cup q_2$, where the set 1127 q_2 is formed by trajectories of length $K \times 2$, and so on (in order of increasing complexity of the 1128 trajectories). Thus, for the N segments considered during training, the set $Q_N = \bigcup_{i=1}^N q_i$ is used.

In the T-Maze environment, DT, RATE, RMT, and TrXL were trained with and without curriculum
learning because this approach theoretically produces better results. However, it is important to note
that the T-Maze task is successfully solved by the RATE model without using curriculum learning,
and even vice versa – its use slightly degraded performance on long corridors. However, with respect
to TrXL, the use of CL yielded slightly better results. The work showed that using CL does not
achieve significantly better performance on the T-Maze task. The results of using the CL on the

Env.	R	0	Conv. configuration [‡]	Α
ViZDoom-Two-Colors	Linear	$Conv2D \times 3$	(32, 64, 64) / (8, 4, 3) / 0	Embedd
T-Maze	Linear	Linear	_	Embedd
Memory Maze	Linear	$Conv2D \times 3$	(32, 64, 64) / (8, 4, 3) / 2	Embedd
Minigrid-Memory	Linear	$Conv2D \times 3$	(32, 64, 64) / (8, 4, 3) / 0	Embedd
Atari	Linear	$Conv2D \times 3$	(32, 64, 64) / (8, 4, 3) / 0	Embedd
MuJoCo	Linear	Linear	_	Linea

Table 10: RATE encoders for each part of (R, o, a) triplets. We use Embedding layer for encoding discrete actions and Linear for continuous ones. \ddagger – channels / kernel sizes / padding.



1168

1144

1169 Figure 12: Memory Retention Valve configurations used in the ablation study. MRV-CA-2: cross-1170 attention-based MRV which uses an attention mechanism to control the updating of memory embed-1171 dings and which is used in the work as the main mechanism. MRV-CA-1: uses the same mechanism 1172 as MRV-CA-2 but the updated memory embeddings M_{n+1} are fed to Query, and the incoming memory embeddings M_n are fed to Key and Value. MRV-G: gated MRV which uses a gating 1173 mechanism similar to the one used in Gated Transformer-XL (Parisotto et al., 2020). MRV-GRU: 1174 uses a GRU (Chung et al., 2014) block to process updated memory embeddings with hidden states. 1175 MRV-LSTM: uses a LSTM (Hochreiter & Schmidhuber, 1997) block to process updated memory 1176 embeddings with cached states. 1177

- 1178
- 1179

1182

T-Maze environment are presented in Figure 13 (left), and the results of applying noise to memory embeddings to assess its importance are presented in Figure 13 (right).

1183 F.3 SUPPLEMENTAL MRV ABLATION

One of the options for implementing the memory tokenization gating mechanism was an approach similar to the one proposed in Gated Transforer-XL (GTrXL) (Parisotto et al., 2020) work. Thus, the MRV-G scheme was inspired by the gating mechanism from GTrXL and implemented as follows:

$$r = \sigma(M_n W_r + M_{n+1} U_r) \tag{2}$$



Figure 13: (left). Results with and without the use of curriculum learning and (right) results of replacing RATE memory tokens with white noise at inference in T-Maze.



Figure 14: Results of RATE inference with different MRV configurations on the T-Maze environment. Training was performed with the number of segments N = 5 and context length K = 30, i.e. on trajectories of length ≤ 150 . MRV-CA-2 is the final MRV configuration that is used throughout the work and is designated as MRV.

$$=\sigma(M_nW_z + M_{n+1}U_z - \text{bias}) \tag{3}$$

$$n = \tanh(M_n W_q + (M_{n+1} \times r)U_r) \tag{4}$$

$$\tilde{M}_{n+1} = \sigma(M_n(1-z) + z \times h) \tag{5}$$

1227 The results of the RATE (trained on corridor lengths of ≤ 150) inference on the T-Maze environment 1228 with these MRV configurations are shown in Figure 14 and in Table 6. The results presented 1229 in Figure 14 confirm the high stability of RATE when using cross-attention-based MRV (MRV-CA-2), 1230 as well as the model's ability to hold important information in memory embeddings when inference 1231 on long tasks.

1232 1233

1221

1222 1223

1224 1225

1226

1234 1235

1236

F.4 ABLATION ON NUMBER OF SEGMENTS AND SEGMENT LENGTH

z

ł

Partitioning the trajectories into fixed-length segments allows the RATE model to train on long trajectories without increasing the context size, which makes the parameters N (the number of segments into which the training trajectories are divided) and K (the context length, i.e., the size of a single segment) critical because they determine the length of the effective context $K_{eff} = K \times N$. The Figure 15 presents the results of ablation studies for parameters N and K at fixed $K_{eff} = 90$.



Figure 15: Results of ablation of segments size (context length K) and number of segments (N) when the effective context length (K_{eff}) is fixed: $K_{eff} = K \uparrow \times N \downarrow = 90$.

G RECURRENT BASELINES

To prove that all baselines were properly trained and that the results obtained indicate exactly the inability of the considered RNN and SSM baselines to learn on long corridors, we conducted a hyperparameters sweep (see Figure 16).

1267Results demonstrates that recurrent baselines can solve the T-Maze task when trained on data with1268moderate corridor lengths (approximately 30 steps, or 90 corresponding tokens) but fail to retain the1269clue information for longer lengths, unlike a transformer. This is because the transformer's attention1270mechanism can effectively capture dependencies in highly sparse data, which recurrent models cannot.1271DT achieves SR= 0.5 for T > K for any K, while recurrent networks can achieve SR> 0.5 in this1272setting. RATE combines the strengths of transformers (direct access to information in context) and1273recurrent networks (hidden states for information retrieval).



Figure 16: Results of tuning DLSTM, DGRU and DMamba hyperparameters for the T-Maze environment. Validation is performed on corridors of the same length used in training. At each step of the environment, a triplet (R, o, a), i.e., three tokens, is processed.

1296 **TRANSFORMER ABLATION STUDIES** Η 1297

1298 **Transformer core hyperparameters.** This section presents the results of ablation studies on the 1299 main hyperparameters of the RATE transformer. The RATE configuration for the T-Maze environment 1300 specified in Table 8 was chosen for the ablation studies. The ablation studies focus on understanding 1301 the impact of key hyperparameters by systematically varying one parameter while keeping others 1302 constant. The results are shown in Figure 17, Figure 18, and Figure 19.



Figure 18: Results of ablation by the number of attention heads of the RATE model in T-Maze environment.

1336 1337 1338

1303

1339 Feed-Forward Network. In our experiments, we found that when the feed-forward network (FFN) 1340 is disabled in the transformer decoder, RATE performs slightly better then with FFN enabled. To 1341 evaluate the contribution from FFN on the considered baselines, we performed an ablation study on 1342 this parameter. The results presented in Figure 20 demonstrate that for RATE alone, disabling FFN 1343 gives a performance gain, while the other models' Succes Rate in the T-Maze environment drops. 1344

1345

1348

RECOMMENDATIONS FOR HYPERPARAMETERS SETTINGS Ι 1347

Transformer architectures have many parameters that need to be selected correctly. The use of 1349 memory mechanisms in RATE adds a few more hyperparameters. Nevertheless, tuning RATE is



J **TECHNICAL DETAILS**

The Table 11 shows the technical parameters of the training models. Note that the difference between the number of DT and RATE parameters is small and equal to $\delta p = d_{model} \times num_{mem_tokens}$ $\sim 10^3$. Training RATE with trajectory splitting into N segments allows $\sim N$ smaller GPU memory size usage than for DT. The training was conducted using a single NVIDIA A100 80 Gb graphics card.

Table 11: Technical configurations of model training. The values in the table are for single run. The training was conducted on a single NVIDIA A100 GPU.

Env.	Model	GPU mem.	Train time	# params.	
ViZDoom Two Colors	RATE	24Gb	4h	6 0M	
VIZD00III-1W0-C010IS	DT	37Gb	411	0.0111	
T Maza	RATE	5Gb	2h	2 4M	
1-widze	DT	15Gb	211	2.4111	
Mamory Maza	RATE	24Gb	12h	6 OM	
Welliofy Waze	DT	37Gb	1211	0.0101	
Minigrid Mamory	RATE	6Gb	10b	6 0M	
wining nd-wiemory	DT	15Gb	1011	0.0101	
Atori	RATE	21Gb	Ob	4 7M	
	DT	32Gb	<i>7</i> 11	4./IVI	
MulaCa	RATE	15Gb	10b	0.6M	
Mujoco	DT	45Gb	1011	0.0101	

Κ ATTENTION MAPS

In this section, we present attention maps for DT and RATE models in the T-Maze environment in two configurations: T = K = 15 (Figure 21, Figure 22) and T = K = 90 (Figure 23 and Figure 24). As can be seen from the presented attention maps, DT have attention heads that explicitly define dependencies between the action at the junction and the cue at the beginning of the episode (Head 0, Head 1 in Figure 22). RATE, on the other hand, does not show such dependencies explicitly, but some heads clearly show heavy use of memory tokens (Head 2, Head 4, Head 7 in Figure 21).





Figure 22: DT attention maps in the T-Maze environment, T = K = 15.





