# MEMORY-EFFICIENT ALGORITHM DISTILLATION FOR IN-CONTEXT REINFORCEMENT LEARNING

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026 027

048

049

Paper under double-blind review

### ABSTRACT

It's recently reported that by employing the superior In-context Learning (ICL) ability of autoregressive Transformer, a method named Algorithm Distillation (AD) could distill the whole Reinforcement Learning process into neural network then generalize to unseen scenarios with performance comparable to the distilled algorithm. However, to enable ICL, it's vital for self-attention module to have a context that spans cross-episodes histories and contains thousands of tokens. Such a long-range context and the quadratic memory complexity of self-attention pose difficulty on applying AD into many common RL tasks. On the other hand, designing memory efficient Transformers for *long-range document modeling* is itself a fast-developing and fruitful field, which leads to a natural question: Could Efficient Transformers exhibit similar in-context learning ability and be used for Memory-Efficient Algorithm Distillation? In this paper, we firstly build a benchmark suite that is thorough, efficient and flexible. Thanks to it, we perform extensive experiments and verify an existing method named ERNIE-Docs (ED) could offer competitive performance with significantly reduced memory footprint. With systematic ablation studies, we further investigate various facets influencing the ICL ability of ED and provide our own insights into its hyperparameter tuning.

# 028 1 INTRODUCTION

Building decision-making agents that could generalize to various environments and tasks has been
a long-standing goal of Reinforcement Learning (RL). Recently, following the trends in Computer
Vision (CV) (Dosovitskiy et al., 2021) and Natural Language Processing (NLP) (Brown et al., 2020;
Devlin et al., 2019), multiple works have explored on pretraining deep autoregressive Transformers
on large-size behavior dataset to obtain agents with strong generalization and adaptation ability (Reed
et al., 2022; Lee et al., 2022; Octo Model Team et al., 2024). Among them, *Algorithm Distillation*(AD) takes an unique approach and proposes that by employing the superior In-Context Learning
(ICL) ability of autoregressive Transformer, the whole RL process could be distilled into a neural
network then be 'replayed' on new, possibly unseen scenarios.(Laskin et al., 2023)<sup>1</sup>

Analogous to causal Language Modeling (LM), AD requires an autoregressive Transformer with
 context spanning cross-episodes interaction histories to perform *credit assignment, exploration exploitation balance* and *policy improvement* that are vital to a RL algorithm. Differently and
 crucially, all of these abilities are implicitly modeled in the weights of Transformer network instead
 of being designed by human experts and implemented as hard-coded programming languages as in
 Haarnoja et al. (2018); Fujimoto et al. (2018); Schulman et al. (2017).

However, this cross-episodes context encompasses thousands of tokens, resembling long-range
 document modeling in NLP: In typical RL settings, an episode contains hundreds of transitions each
 comprising 4 *components (state, action, reward* and *terminal*):

 $2 \sim 4$  (cross-episode context length is needed)  $\times 500$  (transitions in an episode)

 $\times 4$  (components in a transition) =  $4k \sim 8k$  (tokens)

Combined with the quadratic memory cost of self attention, this requires excessively large memory footprint (around 40-160GB under batch size 10, embedding size 128 and half-precision) and prevents RL researchers from readily reproducing and improving algorithms in AD's setting.

<sup>&</sup>lt;sup>1</sup>We use term 'AD' to denote both the problem setting and vanilla Transformer used in the original paper.

On the other hand, designing memory & computation Efficient Transformer (ET) for long-range document modeling is itself an important and developing research field in NLP where many well-established works have emerged (Dai et al., 2019; Ding et al., 2021; Tay et al., 2022). Therefore, a natural question arises: *Could Efficient Transformers exhibit similar in-context learning ability and be used for Memory-Efficient Algorithm Distillation?* If so, it could combine the best of both worlds and significantly boost future researches in this line.

However, to the best of our knowledge, there is still no works investigating this question. The reasons are multi-fold: 1) It still lacks a consensus on how to design benchmarks to thoroughly test the ICL ability of a given algorithm. 2) As of a meta-RL setting, the implementation of above benchmark shall be computation-efficient for parallel data collection and evaluation.

Therefore, in this work, we firstly design an efficient and flexible benchmark suite to thoroughly test the ICL performance of a given method. After this, with comprehensive experiments, we discover *ERNIE-Docs* (ED) (Ding et al., 2021), an improved variant of *Transformer-XL* (Dai et al., 2019), could obtain competitive performance with a significantly reduced memory requirement. Finally, with systematic ablation studies, we examine various factors influencing the ICL performance of ED and provide our own insights in its hyperparameter tuning. In summary, our contributions are:

• We provide a benchmark suite covering the whole AD process from data collection to meta training & evaluation. With 8 representative settings, we could thoroughly test a method's ICL ability in cases of sparse reward, credit assignment, dense information and exploration-exploitation balance.

• We employ JAX (Bradbury et al., 2018) for computation efficiency and achieve parallel execution on *10k* different environments. Our code is also compatible with Pytorch (Paszke et al., 2019) for its broader Transformer-related ecosystem.

• Based on our benchmark, we find **ERNIE-Docs** obtains competitive performance with a significantly reduced memory requirement. We also perform systematic ablation studies to demonstrate how various factors affect its performance and provide our own insights.

# 2 RELATED WORK

084

085

087

880

091

071

072

073

074

075

076

077

079

**Sequence Modeling in Reinforcement Learning**. Interpreting RL as a problem of sequence modeling and employing Transformer for it has been a popular research direction since the recent success of Transformer in NLP field (Radford et al.; Devlin et al., 2019). Chen et al. (2021) and Janner et al. (2021) firstly introduced Transformer into model-free and model-based RL respectively and achieved promising results in many tasks. As a pioneering work, Laskin et al. (2023) firstly reported that by using autoregressive Transformer, in-context RL could be achieved with generalization to unseen scenarios. Following their work, Dai et al. (2023); Zisman et al. (2024); Sinii et al. (2024) investigated various facets to improve AD, such as training on noisy dataset and generalizing to unseen action space. Differently, our work targets on the memory cost of vanilla AD.

- Meta Reinforcement Learning. Meta RL aims at empowering RL algorithms the ability to quickly adapt to new environments and tasks with limited amount of *on-site* samples. Generally, it could be classified into 2 branches (Beck et al., 2023): 1) in-weight meta RL where methods like MAML (Finn et al., 2017), ProMP (Rothfuss et al., 2022) perform gradient descend on parameters of neural network for fast adaptation. 2) in-context meta RL where adaptation emerges as context of environment interactions gets populated or updated as in methods like RMA (Kumar et al., 2021) and Prompt-DT (Xu et al., 2022). Different from these works whose goal is fast adaptation of episode return, AD aims to distill existing RL algorithms then generalize to new scenarios.
- 099 Transformer for Long-range Tasks. Designing Efficient Transformers to increase the performance 100 on long-range tasks has been a very important field (Burtsev et al., 2021; Tay et al., 2022). Dai 101 et al. (2019) proposed integrating recurrence with Transformer's attention mechanism and achieved 102 superior performance on long-range tasks. Later work like Rae et al. (2019) also studied to 'compress' 103 past information into recurrence module and further improved Transformer's performance over 104 tasks on long-range document modeling while not drastically increasing context length. Besides, 105 methods like (Gu & Dao, 2023; Wang et al., 2020) tried to design new attention mechanisms with only linear complexity to context length. Since this field is both fruitful and fast-developing, in this 106 work we choose to explore how these promising methods could be used in AD's setting such that RL 107 researchers could be freed to explore on ideas more related to RL's problems.

#### 108 **PROBLEM FORMULATION** 3

#### 110 3.1 Preliminary 111

121 122

125

126

127 128

129 130

137 138

141

143 144

145

146

112 Algorithm Distillation. Consider a meta-RL setting, where a MDP  $M_i$  is sampled from a given distribution. For each sampled MDP, we use a RL algorithm  $\mathcal{A}$  (e.g. Q-Learning (Sutton & Barto, 113 2018) for discrete setting or SAC (Haarnoja et al., 2018) for continuous setting) to search for the 114 optimal policy  $\pi$  whose performance is measured by *episode return*  $\sum_{i=0}^{h} R_i(s, a)$  where h is the 115 maximum episode length and  $R_i$  is the reward for *i*th transition. Then, by collecting the whole history 116  $\tau_i$ , we get a dataset  $D = \{\tau_i = (s_0, a_0, r_0, d_0, ..., s_T, a_T, r_T, d_T)\}_{i=1}^n$  where s is state, a is action, r 117 is reward, d is terminal and  $\tau_i$  is the learning history of A on  $\mathcal{M}_i$ . Note  $\tau$  may contain many episodes 118  $(T \gg h)$  in which the performance of  $\pi$  gradually improves. 119

120 AD proposes by performing autoregressive training,  $\mathcal{A}$  could be distilled into a neural network  $\phi$ :

$$\mathcal{J} = \operatorname*{arg\,max}_{\phi} \mathbb{E}_{\tau \sim D; s, a, r, d \sim \tau} \left[ \phi(a_t | s_{0:t}, a_{0:t-1}, r_{0:t-1}, d_{0:t-1}) \right] \tag{1}$$

123 Then the weight-frozen  $\phi$  could be used on unseen  $\mathcal{M}_i$  and yield similar learning process as  $\mathcal{A}$ . 124

Self Attention. For a token sequence  $x \in R^{s \times d}$  of length s, self attention firstly transforms x into three matrices: query  $Q = xW^Q$ , key  $K = xW^K$  and value  $V = xW^V$ , where  $W^Q, W^K, W^V \in$  $R^{d \times d}$ . Then scaled dot attention is applied between any 2 elements in x:

Self-Attention
$$(Q, K, V) = \text{Softmax}(\frac{QK^T}{\sqrt{d}})V$$
 (2)

131 **Recurrence Transformer.** Transformer-XL (XL) (Dai et al., 2019) proposes to integrate recurrence into attention to obtain longer context length. As shown in Figure 1, the token sequence is split into 132 equal-length chunks and the Transformer processes these chunks sequentially from beginning to 133 ending. During this, the hidden state of a previous chunk is preserved and acts as additional memory 134 (key and value) when processing the following chunk. For the *i*th transformer layer  $L^i$ , the input and 135 output are: 136

$$h_{kc:(k+1)c}^{i} = L^{i}(q = h_{kc:(k+1)c}^{i-1}, k = v = [\operatorname{sg}(h_{(k-1)c:kc}^{i-1}), h_{kc:(k+1)c}^{i-1}])$$
(3)

where  $h_k^i$  represents hidden embedding output by *i*th layer for kth chunk, c is chunk length, sg 139 is stop-gradient, q,k,v are query, key, value, respectively and  $[\cdot, \cdot]$  concatenates 2 vectors. As an 140 improved version of XL, ERNIE-Docs (ED) (Ding et al., 2021) chooses to shift the preserved hidden embeddings one-layer down for attention computation (difference to XL is in red): 142

$$h_{kc:(k+1)c}^{i} = L^{i}(q = h_{kc:(k+1)c}^{i-1}, k = v = [\operatorname{sg}(h_{(k-1)c:kc}^{i}), h_{kc:(k+1)c}^{i-1}])$$
(4)

Through this simple yet effective modification, ED could have a context that is theoretically indefinitely long, as shown in Figure 1, right.





# 162 3.2 MEMORY COST ANALYSIS

We present a formal analysis of memory cost for vanilla Transformer and Recurrence Transformer
(abbreviated as RT and denoting XL and ED). For the sake of simplicity, only attention module is
focused as other modules like tokenizer, feedforward and layernorm layer are the same in our settings.

**Vanilla Transformer**. Let's consider an input x of shape (b, s, d) where b is batch size, s is sequence length and d is hidden dimension. For a multi-head attention layer with n heads. This requires the following amount of memories:

169 170 171

167

168

172

173 174

178 179

181 182 183  $3(d^{2} + d) \text{ (weights \& biases)} + 3(d^{2} + d) \text{ (gradients of weights \& biases)}$  $+ 3bsd (Q, K, V) + nbs^{2} (QK^{T}) + nbs^{2} (\text{attention score}) + bsd (\text{final output})$ (5) =  $6d^{2} + 6d + 4bsd + 2nbs^{2} \propto s^{2}$ 

**Recurrence Transformer.** Let's consider a chunk length of c and a recurrence length of m. Then for each step of RT processing, the chunk input  $x_c$  is in shape (b, c, d) and recurrence  $x_m$  is in shape of (b, m, d). The amount of memories for RT is:

$$bmd (recurrence) + 3(d^{2} + d) (weights \& biases) + 3(d^{2} + d) (gradients of weights \& biases) + bcd (Q) + 2b(c + m)d (K, V) + nbc(c + m) (QK^{T}) + nbc(c + m) (attention score) + bcd (final output) = 3bmd + 6(d^{2} + d) + 4bcd + 2nbc(c + m) \propto c(c + m)$$
(6)

From above analysis, the memory cost is proportional to  $s^2$  for vanilla Transformer and is proportional to c(c + m) for RT methods. However, in the setting of AD, cross-episode context is needed for vanilla self-attention based Transformer and encompasses thousands of tokens. On the other hand, the recurrence module in RT could serve as an additional memory for important information and reduce the need of contiguous context which is unnecessarily long. That is to say, with  $c + m \ll h < s$ where h is episode length, the memory cost of RT could be greatly reduced than vanilla Transformer.

190

### 4 BENCHMARK DESIGN

191 192

215

It's still an open question over how to design a benchmark suite to test the ICL ability of AD methods.
Firstly, the benchmark needs to be **thorough** enough and include representative RL settings. Besides,
as an in-context meta RL setting, the benchmark needs to be **computation-efficient** and support
parallel environment interactions for both data collection and meta evaluation. Finally, the benchmark
shall also be **flexible** and compatible towards the most ecosystems.

In this work, taking all above considerations into account, we design a benchmark suite that is
computation-efficient thanks to environment parallelism in JAX and flexible since it also supports
building algorithm in pure Pytorch. We also carefully design environments & tasks to thoroughly test
given methods' ICL performance. An illustrative framework of our benchmark is shown in Figure 2.
The further details could be found in Appendix A.3.



Figure 2: Framework of benchmark process

Environments. We choose GridWorld where the observation is 2-dimensional coordinates (*x*, *y*) and action contains 5 options: *up*, *down*, *left*, *right*, *none*, as shown in Figure 3.

- DarkRoom (dr). The grid is 9x9 large. The episode length is 20. The agent spawns in the middle of it with a randomly placed invisible goal. This environment tests the agent's ability to explore in an unknown environment to find the 'goal' then remember its location for exploitation.
- DarkKeyToDoor (dktd). The grid is 9x9 large. The episode length is 50. An agent spawns in the middle of it. A *key* and a *door* are randomly placed and invisible to agents. Note this



Figure 3: Environments

- environment is more challenging compared to DarkRoom because it features 2-phase sequential
  tasks as the 'key' could only be interacted once and the 'door' could only be opened when agent
  holds the key.
  - DarkRoomLarge (drl). A larger verion of DarkRoom. The grid size is 13x13. The episode length is 50. Note the increase in grid size generally results in exponentially increased difficulty in exploration and memorization.
- DarkKeyToDoorLarge (dktdl). A larger version of DarkKeyToDoor. The grid size is 11x11. The episode length is 70.

**Tasks**. 3 representative tasks are designed: *normal*, *dense* and *quick* to test the algorithm's RL ability like credit assignment and sparse reward, as shown in Figure 4.

normal (n). In DarkRoom, if the agent finds the goal, a reward of +1 is granted and the episode terminates. For DarkKeyToDoor, if the agent finds the key, a reward of +1 is granted *only once*. Then if the agent finds the door with key on hand, another +1 reward is granted and the episode terminates. For other cases, the reward is 0. This is a sparse reward

229

230

231

233

234

235

			-1	-1	-1	-1	-6	-5	-4	-3
÷			٩	-1	-1	-1	ġ	4	-3	-2
			-1	-1	-1	-1	-4	-3	-2	-1
		*	-1	-1	-1	*	-3	-2	-1	1

Figure 4: Illustration of tasks

- setting and also tests the memorization ability of agent as it needs to remember the location of
   'target' ('goal' for DarkRoom or 'key'/'door' for DarkKeyToDoor) and discover that the 'key'
   only grants one-time reward.
- dense (d). For both DarkRoom and DarkKeyToDoor, the reward for each step is the negative L<sub>1</sub> distance from agent to current 'target'. The episode still terminates once the goal is found or door is opened. For other cases, reward is 0. This setting features dense reward and is often overlooked in previous works.
- quick (q). For all steps, the reward is -1. Only when the agent finishes the task, the episode would terminate which encourages agent to finish the episode as soon as possible. This is the most challenging setting and when used with DarkKeyToDoor, a strong ability of credit assignment is required since agent gets no immediate signal when finding the 'key'.

In total, we select 8 *settings* from above environments and tasks : dr, drln, drlq, drld, dktdln, dktdlq, dktdld with the format {env. name}{task variant}. Note for certain setting, the gridsize and episode length may be slightly modified, please refer to Appendix A.3 for details.

- Algorithms. We choose the following algorithms for testing on above setting. Their details are in
   Appendix A.4:
- SOURCE. This is the original algorithm used to collect the dataset. We use Q-Learning for all tasks. Notably, for DarkKeyToDoor, the state of SOURCE is augmented with one extra dimension: *has key or not* for *Markovian property* while in dataset and the following training/evaluation of AD methods, this extra dimension is removed. In our settings, the learning curve of SOURCE is regarded as *ORACLE* and shall be mimicked by other methods.
- ADF. This is the original AD algorithm implemented manually by us. The F means we use long context length which is often ~2.5 times of episode length following recommendations in the original paper (Laskin et al., 2023).
- ADR. This is the AD method with a reduced context length for a fair comparison to other memory efficient algorithms. In most cases, the context length of ADR is less than a half of episode length.
   For details, please refer to Appendix A.4.4.

5	
5	
~	
~	

- MEM. This is the method introduced in Burtsev et al. (2021). Specifically, we use MemCtrl where a set of learnable embeddings acts as global memories, gets prepended to all sequence before attention layers and discarded after its processing.
  - XL. This is the Transformer-XL method as described above and in Dai et al. (2019). Notably, this method (and the following ED) requires Truncated-Backprop-Through-Time (TBPTT) training.
  - ED. An improved version of XL with *theoretically indefinitely long* context length.

In the following contents, we use Efficient Transformer (ET) to denote methods of MEM, XL and ED
 for simplicity.

Training & Evaluation. As this is a setting of in-context meta RL, for each of 8 settings defined above, we firstly collect a dataset containing training process of SOURCE on 10k environments with various locations of 'target'. Then, we train each above algorithm on the dataset and evaluate on 100 new environments whose 'target' locations are disjoint against training ones, i.e. out-of-distribution (OOD) evaluation. For comparing the performance of each algorithm, we plot its curve of trajectory return against ICL process<sup>2</sup>, and the more its curve mimics the curve of SOURCE, the better performance it exhibits.

Extensibility. Thanks to the flexibility and broad ecosystem of JAX and Pytorch, our benchmark
 could be easily extended to support more environments, tasks & algorithms. For example, via
 MujoCo's MJX (Todorov et al., 2012), environment parallelism with domain randomization could be
 easily enabled for continuous settings like Walker, HalfCheetah and Humanoid. Plus, it's also
 compatible with Pytorch with which many existing Efficient Transformers are implemented. Hence,
 we plan to publish this benchmark and would continue to add representative environments, tasks &
 algorithms into it for more comprehensive results.



Figure 5: Results of main experiments. The shaded area is 95% confidence interval

323

270

271

272

273

274

<sup>320</sup> 321 322

<sup>&</sup>lt;sup>2</sup>The ICL process here is indeed the in-context RL process which spans many episodes with gradually improved episode returns.

#### 5 **EXPERIMENTS**

## 5.1 MAIN RESULTS

We present our main experiment results of training above 6 algorithms on 8 settings. The learning curves are shown in Figure 5 and the memory consumptions of these methods are present in Table 1.

Table 1: Memory consumption of various methods

334	Settings	Episode Length	<b>ADF</b> $^{\alpha}$	ADR $^{\alpha}$	MEM $^{\beta}$	$\mathbf{XL}^{\beta}$	ED $^{\beta}$
335	dr	20	50   12.75 GB $^{\gamma}$	10   6.42 GB	10   1	0   6.93	GB
336	dktd	50					
337	dr <b>ln</b>						
339	dr <b>lq</b>	50					
340	dr <b>ld</b>		125   37.33 GB	25   10.65 GB	25   50	)   13.96	GB
341	dktd <b>ln</b>	70					
342	dktd <b>lq</b>	50					
343	dktd <b>ld</b>	70					

 $\alpha$  For AD methods, the format is {context length} | {memory consumption}.

 $\beta$  For ET methods, the format is {chunk length} | {recurrence capacity} | {memory 345 consumption}. 346

 $\gamma$  This is the total GPU memory consumption measured on device. 347

348

324

325 326

327 328

329

330 331 332

333

349 From the results, several interesting observations could be drawn: 350

**Dense information helps Algorithm Distillation**. Notably, for drld and dktdld, most algorithms 351 obtain quite decent ICL performance, even for MEM and ADR. We conjecture this is due to the rich 352 guidance provided by the reward signal in *dense* task hence a sub-episode context could sufficiently 353 prompt algorithm for the correct behavior. 354

Tasks requiring credit assignment are challenging. For dktdlq, almost all algorithms fail, showing 355 this is a quite challenging setting since the acquisition of 'key' provides no immediate signal to agent. 356 Even for ADF, the performance is inferior compared to SOURCE, indicating a context length of more 357 than 2.5 episode length may be needed. How to increase the ICL performance of AD methods and 358 reduce their memory cost on this setting would certainly be a very important and promising directions 359 for RL researchers.

360 **Context Length is critical for vanilla AD method.** From above results, we find compared to ADF, 361 ADR often obtains significantly worse performance, except on dense-information tasks. This aligns 362 with our expectation that the self-attention context in AD needs to span multiple episodes to enable a decent AD performance. However, combined with quadratic memory cost of self-attention, this 364 long-range context causes difficulty for researches to experiment with, as also shown in Table 1 where the ADF needs significantly more GPU memory than other methods.

- 366 ED for Memory-Efficient Algorithm Distillation. For ET methods, we find: 367
- MEM obtains quite similar performance compared to ADR, indicating its global memory doesn't 368 help in algorithm distillation.
- 369 • XL, with just slight differences to ED, also obtains insufficient performance, even compared to ADR. 370 Even after architecture tuning, we still find the performance of XL is unstable and argue this results 371 from its limited ability to utilize information in recurrence.
- 372 • ED, on most of settings could obtain a competitive performance compared to ADF and SOURCE 373 thanks to its clever usage of recurrence memory. On the other hand, there is still room for further 374 improvements as its learning curve tends to slightly drop near the ending phase, suggesting it 375 struggles to perform exploitation compared to SOURCE and ADF. With some hyperparameter tuning and architecture tweaking, this drop could be mitigated in some extend, as shown in Section 5.2.4. 376
- Nevertheless, further investigation and mitigation to this phenomenon would still be an interesting and valuable future work.

#### 378 5.2 ABLATION STUDIES 379

383

385

380 To better illustrate how various facets affect Memory-Efficient Algorithm Distillation, we design 381 several ablation studies testing the positional encoding, Transformer model size, attention module and context length & memory capacity. 382

#### 384 5.2.1POSITIONAL ENCODING

Positional Encodings (PE) like sinusoidal encoding (Vaswani et al., 2017), learnable embedding and 386 *RoPE* (Su et al., 2023) are critical components for attention-based Transformer models and have been 387 widely studied in NLP. 388

Interestingly, in our experiments, we find the us-389 age of global positional encoding has unique ad-390 vantages over local-context-based ones common 391 in NLP. Without change in embedding mecha-392 nism, the 'global' PE assigns to all the tokens in 393 the whole learning process a positional embed-394 ding while context-based PE only considers to-395 kens in the context window. We argue global PE 396 is more suitable for in-context RL as it's intuitive 397 for RL algorithms to take different strategies in 398 different learning phase, such as exploration in 399 the beginning and exploitation in the ending.

400 Besides, we also design 2 experiments to verify 401 this: 1) Global PE helps training. As shown in 402 Figure 6, we test 4 PEs on standard ADF method. 403 The 4 abbreviations represent global-learnable 404 (gl), global-sinusoidal (gs), local-learnable (ll) 405 and local-sinusoidal (ls), respectively. From the results, global positional encoding significantly 406 helps stabilize and ease the training of Algo-407 rithm Distillation. 2) Additionally, global PE 408



Figure 6: Performance when using different PEs





could tune the behavior of distilled algorithm. 409 By tweaking the positions passed to the PE module, we could tune the algorithm to be more ex-410 ploratory or more exploitative without changing the input sequence. As shown in Figure 7, when 411 we tweak the positions passed to Transformer model to be larger, the performance of AD tends to 412 increase since it focuses more on *exploitation*. While if we tweak the position to be smaller, the 413 AD becomes more *exploratory*. Therefore, we choose global learnable PE in main results and the 414 followings. 415

### Takeaway

Global Positional Encoding helps training and controls exploration-exploitation behavior.

#### 5.2.2TRANSFORMER MODEL SIZE

422 Previous works in NLP and AD have found that with the increase of model size, Transformer 423 models exhibit better in-context learning ability. In this section, we also study whether this phe-424 nomenon happens on ED for memory-efficient AD. Specifically, for drn and dkldln, we test the effects of reducing the model size <sup>3</sup> over the ICL ability of ED. The results are shown in Figure 8, 425 where the format of the legend is {algo name}\_{number of attention head}-{token 426 dim. }-{feedforward dim.} and 'edo' means the original ED used in the main experiment 427 which has the largest model size. 428

429 From the results, we could find with the model size increasing, the ICL performance of ED is also 430 increased. Even for dr**n** which is a rather small environment, a model with 64 attention heads, 512

416 417

418 419 420

<sup>431</sup> 

<sup>&</sup>lt;sup>3</sup>the number of attention heads are also reduced to maintain a similar attention head dimension.



Figure 8: Performance of large model vs small model

token dimension and 2048 feedforward dimension is needed for a decent performance. On dktdln, the effect of increasing model size also exists but is in a less significant extend.

### Takeaway

443 444

445

446 447

448 449

450 451 452

453

Larger model enables better performance for ED-based Memory-Efficient AD.

### 5.2.3 RECURRENCE-BASED ATTENTION

Attention mechanism is key to the success of Transformer
models. In this work, as a recurrence-based Transformer
method, ED also utilizes attention to extract important information from recurrence memory. Therefore, we also
test how various attention settings influence its ICL performance in the following experiments.

460 Does ED attend to the recurrence memory? It would be
461 natural to expect the outstanding ICL performance of ED
462 results from the usage of its recurrence memory. Therefore,
463 we visualize the attention ratio of ED for all 8 settings. In
464 details, we sum the softmaxed attention scores for tokens
465 in ED's recurrence memory: 1 means ED only attends to



Figure 9: Recurrence attention ratio

and D s recurrence memory and ignores chunk tokens while 0 represents vice versa, then average this *ratio* over all attention heads and all 100 environments in evaluation set and plot it against the learning progress, as shown in Figure 9.

From the results, it clearly shows ED attends to its recurrence memory as all the attention ratios are
between 0.3~0.6. Besides, several interesting observations could be drawn from it: 1) For dktdlq
where ED obtains no performance due to the challenge of credit assignment, its attention ratio is
also the lowest, indicating ED fails to extract useful information from the recurrence memory. 2)
For dktdld and drlq where ED obtains near-optimal performance, the attention ratio shows a notable
increase near the ending of progress. While for the other settings where ED's performance declines at
the ending of process, their attention ratios doesn't show such increase.

How does the number of attention head in-476 fluences ICL performance? In this section, 477 we test the effects of number of attention heads 478 of ED over its ICL ability. Specifically, we 479 change the attention head number from 64 to 480 32 and 128 respectively. Note we don't change 481 the embedding dimension hence this would re-482 sult in changes in dimension of attention head. 483 The results are shown in Figure 10 where the format of legend is {algo name}\_{amount 484 of attention head} and 'edo' is the orig-485 inal ED used in main experiment with 64 heads.



Figure 10: Effects of amounts of attention heads

From the results, we could find in drn, a relatively simpler setting, the amount of attention heads doesn't have a significant influence on ED's ICL performance. While in dktdln, increasing the amount of attention head (which also reduces the dimension of attention head) has a negative impact over the ICL ability. On the other hand, slightly decreasing it to 32 could increase performance slightly. Note this doesn't conflict with the results presented in Section 5.2.2 as we change both the number of attention head and its dimension there. Taking these together, we recommend setting attention head dimension to be 8 or 16 with at least 32 or 64 attention heads as a good start.

Takeaway

Use 32 or 64 attention heads with dimension of 8 or 16 for a good start.

### 5.2.4 CHUNK LENGTH & RECURRENCE CAPACITY

In this section, we study the effects of chunk length & recurrence capacity over the performance of ED's ICL performance. Specifically, we designed 5 settings: 1-1, 1-2, 1-3, 2-1, 2-2, 2-3 where the first digit represents the relative ratio on chunk length while the latter digit represents the relative ratio on recurrence capacity. The results are shown in Figure 11 where the format of legend is {algo name}\_{chunk length}-{recurrence capacity} and 'edo' is the original setting.



Figure 11: Effects of chunk length & recurrence capacity

The experiments show intricate results: 1) Increasing recurrence capacity alone may not result in better performance as in drn,  $ed_10-20$  and  $ed_10-30$  obtain worse performance than  $ed_010_10$ , and  $ed_20-20$  and  $ed_20-30$  also show decreased performance compared to  $ed_20-10.2$ ) For dktdln, however, several settings obtain boosted performance compared to  $ed_025-50$ :  $ed_50-50$ ,  $ed_{2}5-100$  and  $ed_{2}5-150$  while  $ed_{5}0-50$  and  $ed_{5}0-100$  obtain significantly worse performance. In summary, we recommend setting them to be 1:1 with careful tuning of values.

### Takeaway

Keep balanced ratio between chunk length and recurrence capacity and tune the value.

## 6 CONCLUSION

In this work, we study the problem of reducing the excessive memory requirement of Algorithm
Distillation (AD) via existing Efficient Transformers (ET). With our efficient, flexible and versatile
benchmark, we discover the ERNIE-Docs (ED) could serve as a simple yet effective method for
Memory-Efficient Algorithm Distillation and illustrate how various factors like positional encoding,
model size, attention module and context length & memory capacity influence its performance.

Besides, there is still room to further improve ED's performance and better interpret/visualize its
internal attention mechanism which we choose to leave as future works. Also, because ET is an
fast-developing field and limited by computation resources, we cannot test all promising ET methods
in this work. Nevertheless, we will release our benchmark code for further researches in the line of
Memory-Efficient AD.

540	REFERENCES
541	

548

556

559

560

561

562

566

567

568

- Jacob Beck, Risto Vuorio, Evan Zheran Liu, Zheng Xiong, Luisa Zintgraf, Chelsea Finn, and Shimon
   Whiteson. A Survey of Meta-Reinforcement Learning, January 2023.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL http://github.com/jax-ml/jax.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020.
- Mikhail S. Burtsev, Yuri Kuratov, Anton Peganov, and Grigory V. Sapunov. Memory Transformer, February 2021.
  - Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision Transformer: Reinforcement Learning via Sequence Modeling, June 2021.
- Zhenwen Dai, Federico Tomasi, and Sina Ghiassian. In-context Exploration-Exploitation for Re inforcement Learning. In *The Twelfth International Conference on Learning Representations*, October 2023.
  - Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context. January 2019. doi: 10.48550/arXiv.1901.02860.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT* 2019, *Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pp. 4171–4186. Association for Computational Linguistics, 2019. doi: 10.18653/v1/n19-1423.
- SiYu Ding, Junyuan Shang, Shuohuan Wang, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang.
  ERNIE-Doc: A Retrospective Long-Document Modeling Transformer. In Chengqing Zong,
  Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 2914–2927, Online, August 2021. Association
  for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.227.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-Agnostic Meta-Learning for Fast Adaptation
   of Deep Networks. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 1126–1135. PMLR, July 2017.
- Scott Fujimoto, Herke van Hoof, and David Meger. Addressing Function Approximation Error in Actor-Critic Methods. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1587–1596. PMLR, July 2018.

594 Albert Gu and Tri Dao. Mamba: Linear-Time Sequence Modeling with Selective State Spaces, 595 December 2023. 596 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft Actor-Critic: Off-Policy 597 Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. In Jennifer Dy and 598 Andreas Krause (eds.), Proceedings of the 35th International Conference on Machine Learning, volume 80 of Proceedings of Machine Learning Research, pp. 1861–1870. PMLR, July 2018. 600 601 Michael Janner, Qiyang Li, and Sergey Levine. Offline Reinforcement Learning as One Big Sequence Modeling Problem. November 2021. 602 603 Ashish Kumar, Zipeng Fu, Deepak Pathak, and Jitendra Malik. RMA: Rapid Motor Adaptation for 604 Legged Robots. Robotics: Science and Systems XVII, 2021. 605 Michael Laskin, Luyu Wang, Junhyuk Oh, Emilio Parisotto, Stephen Spencer, Richie Steigerwald, D. J. Strouse, Steven Stenberg Hansen, Angelos Filos, Ethan Brooks, maxime gazeau, Himanshu 607 Sahni, Satinder Singh, and Volodymyr Mnih. In-context Reinforcement Learning with Algorithm 608 Distillation. In The Eleventh International Conference on Learning Representations, 2023. 609 610 Kuang-Huei Lee, Ofir Nachum, Sherry Yang, Lisa Lee, C. Daniel Freeman, Sergio Guadarrama, Ian 611 Fischer, Winnie Xu, Eric Jang, Henryk Michalewski, and Igor Mordatch. Multi-game decision 612 transformers. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), 613 Advances in Neural Information Processing Systems, 2022. 614 Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep 615 Dasari, Joey Hejna, Charles Xu, Jianlan Luo, Tobias Kreiman, You Liang Tan, Lawrence Yunliang 616 Chen, Pannag Sanketi, Quan Vuong, Ted Xiao, Dorsa Sadigh, Chelsea Finn, and Sergey Levine. 617 Octo: An open-source generalist robot policy. In Proceedings of Robotics: Science and Systems, 618 Delft, Netherlands, 2024. 619 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor 620 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward 621 Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, 622 Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning 623 Library, December 2019. 624 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language 625 Models are Unsupervised Multitask Learners. pp. 24. 626 627 Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, and Timothy P. Lillicrap. Compressive 628 Transformers for Long-Range Sequence Modelling, November 2019. 629 Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gómez Colmenarejo, Alexander Novikov, Gabriel 630 Barth-maron, Mai Giménez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom Eccles, Jake 631 Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell, Oriol Vinyals, 632 Mahyar Bordbar, and Nando de Freitas. A generalist agent. Transactions on Machine Learning 633 Research, 2022. ISSN 2835-8856. 634 635 Jonas Rothfuss, Dennis Lee, Ignasi Clavera, Tamim Asfour, and Pieter Abbeel. ProMP: Proximal 636 Meta-Policy Search, February 2022. 637 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy 638 Optimization Algorithms. arXiv:1707.06347 [cs], August 2017. 639 640 Viacheslav Sinii, Alexander Nikulin, Vladislav Kurenkov, Ilya Zisman, and Sergey Kolesnikov. 641 In-Context Reinforcement Learning for Variable Action Spaces. In Proceedings of the 41st International Conference on Machine Learning, pp. 45773–45793. PMLR, July 2024. 642 643 Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. RoFormer: Enhanced 644 Transformer with Rotary Position Embedding, November 2023. 645 Richard S. Sutton and Andrew G. Barto. Reinforcement Learning: An Introduction, 2nd Ed. Rein-646 forcement Learning: An Introduction, 2nd Ed. The MIT Press, Cambridge, MA, US, 2018. ISBN 647

978-0-262-03924-6.

648 649	Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient Transformers: A Survey, March 2022.
651 652 653	Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, pp. 5026–5033. IEEE, 2012. doi: 10.1109/IROS.2012.6386109.
654 655 656	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc., 2017.
657 658 659	Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-Attention with Linear Complexity, June 2020.
660 661 662	Mengdi Xu, Yikang Shen, Shun Zhang, Yuchen Lu, Ding Zhao, Joshua Tenenbaum, and Chuang Gan. Prompting decision transformer for few-shot policy generalization. In <i>International Conference on Machine Learning</i> , pp. 24631–24645. PMLR, 2022.
663 664 665 666	Ilya Zisman, Vladislav Kurenkov, Alexander Nikulin, Viacheslav Sinii, and Sergey Kolesnikov. Emergence of In-Context Reinforcement Learning from Noise Distillation. In <i>Proceedings of the</i> <i>41st International Conference on Machine Learning</i> , pp. 62832–62846. PMLR, July 2024.
667 668	
669	
671	
672	
672	
674	
675	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	

702 703	A APPENDIX
704	
705	A.1 PSEUDO-CODE FOR AD'S TRAINING & EVALUATION PROTOCOL (NEWLY ADDED)
706	We list the training & evaluation protocol of AD-like algorithms used in this work in Algorithm 1.
707	
709	Algorithm 1: Training & Evalutaion Protocol of AD-like algorithms
710 711 712	<ul> <li>Input: Dataset D<sub>τ</sub> obtained by running RL algorithm A on k different environments for T timesteps which span multiple episode length h (typical setting: k = 10<sup>4</sup>, T = (100 ~ 1000)h)</li> <li>Output: a transformer-based nerval network φ<sub>A</sub> showing comparable learning process of A on new scenarios</li> </ul>
713	$\phi_{\mathcal{A}} \leftarrow \text{random init.}$
714	2 while training not converged do
716	$\tau \leftarrow \texttt{sample_chunk}(D_{\tau})$
717	4 $\phi \leftarrow \text{SGD}(\phi_A, \tau)$ // autoregressive training
718	$ \begin{array}{c c} & & \text{In time to eviduate then} \\ \hline & & & \text{Interval} \\ \hline & & & \text{Evaluate}(\phi) \end{array} $
719	$\frac{1}{7}  true\_rewards \leftarrow Evaluate(\mathcal{A})$
720	s compare (true_rewards, rewards)
721	9 end
722	10 end
723	11 return $\phi_A$
724	12 Function evaluate ( $\phi_{-}$ or_ $\mathcal{A}$ )
725	13 /* evaluate RL algorithm ( $\phi$ or $\mathcal{A}$ ) on $n$ OOD envs. for the same $T$
726	timesteps */
720	$\begin{array}{c c} 14 & t \leftarrow 0 \\ 15 & a \leftarrow reset(enus) \end{array}$
720	$16  rewards \leftarrow []$
730	while $t < T$ do
731	18 $  a \leftarrow \phi_{\text{or}}\mathcal{A}(o)$
732	19 $o, a, r, d \leftarrow \text{step}(envs, a)$
733	20 $\phi_0 \sigma A \leftarrow \text{append\_context}(\phi_0 \sigma A, o, a, r, d)$ // only needed for $\phi$
734	$\frac{1}{2} = \frac{1}{2} + \frac{1}$
735	23 return returns
736	24 end
737	
738 739 740 741	The important parts of above pseudo-code are marked in red: We need to perform the evaluation of a whole learning process instead of a single episode to evaluate the performance of given AD algorithms.
742	
743	A.2 HARDWARE RESOURCES
744	
745	equipped with 12,16 CPUs and a Nyidia A100-80G GPU card
746	
747	For parallel data collection, $\sim$ 2hours are needed for each run.
748	For training and evaluation of AD algorithms, wall clock time varies from $\sim$ 2h to $\sim$ 44h depending
750	on the environment & task settings.
751	
752	A.3 ENVIRONMENT SETTINGS
753	
754	A.3.1 DETAILS OF EXPERIMENT SETTINGS
755	

We list details of environments and tasks settings in Table S1:

756

809

dridite drifter drift	In       9       20       125       2.         In       9       50       150       7.         In       13       50       150       1.         Id       12       13       50       150       1.         Idin       11       70       200       2.       2.         Idin       11       70       200       2.       2.         Idin       11       70       200       2.       2.         Idin       11       70       350       2.       2.         is is the total environment transitions the algorithm collects and equals to episode length × episodes to co       ALGORITHM DETAILS         ALGORITHM DETAILS       Table S2: Q-Learning for data collection in all GridWorld environments. Its hyperparameters are in Table S2:       Table S2: Q-Learning's hyperparameters         Image: Table S2: Q-Learning rate       1.0       0.97       1.0         Iearning rate       1.0       0.97       1.0       1.0         Image: Parameters       Values       1.0       0.97       1.0         Image: Parameters       Values       1.0       0.97       1.0         Image: Parameters       Image: Parameters       1.0       0.97       1.0
dktd       9       50       150         driq       1e4       1e2       13       50       150         13       50       250       15       160       160         dktdin       11       70       200       160       175       11         dktdid       9       50       175       11       170       200       160       175       11       170       200       160       175       11       170       200       160       175       11       170       200       11       170       350       160       175       11       170       350       160       175       11       170       350       160       175       11       170       350       160       175       160       175       11       170       350       175       160       175       160       175	td       9       50       150       7.         td       1e4       1e2       13       50       150       1.2         td       15       50       175       1.       1.1
$\frac{drin}{drlq}$ $1e4$ $1e2$ $13$ $50$ $150$ $13$ $50$ $250$ $15$ $10$ $dktdln$ $11$ $70$ $200$ $11$ $dktdld$ $9$ $50$ $400$ $10$ $dktdld$ $11$ $70$ $200$ $11$ $dktdld$ $11$ $70$ $200$ $10$ $dktdld$ $11$ $70$ $350$ $10$ $extrestion the total environment transitions the algorithm collects and equals to episode length × episodes the episode length × episodes the episode length × episodes the episode length × episod × episod × episod × episod × episod × epi$	$\frac{\ln}{\ln}$ $1e4$ $1e2$ $\begin{array}{c} 13 \\ 13 \\ 13 \\ 50 \\ 15 \\ 50 \\ 175 \\ 15 \\ 16 \\ 175 \\ 175 \\ 175 \\ 175 \\ 175 \\ 175 \\ 175 \\ 175 \\ 10 \\ 200 \\ 2e \\ 9 \\ 50 \\ 400 \\ 2e \\ 2e \\ 9 \\ 50 \\ 400 \\ 2e \\ 2e \\ 11 \\ 70 \\ 350 \\ 2e \\ 2e \\ 11 \\ 70 \\ 350 \\ 2e \\ 2e \\ 11 \\ 70 \\ 350 \\ 2e \\ 2e \\ 11 \\ 70 \\ 350 \\ 2e \\ 2e \\ 11 \\ 70 \\ 350 \\ 2e \\ 2e \\ 2e \\ 11 \\ 70 \\ 350 \\ 2e \\ 2e \\ 2e \\ 11 \\ 70 \\ 350 \\ 2e \\ 2e \\ 2e \\ 2e \\ 11 \\ 70 \\ 350 \\ 2e \\ 2e \\ 2e \\ 2e \\ 11 \\ 70 \\ 350 \\ 2e \\ 2$
$\frac{drlq}{drld}$ $1e4$ $1e2$ $\frac{13}{15}$ $50$ $250$ $\frac{dktdln}{dktdlq}$ $1e2$ $\frac{15}{15}$ $50$ $175$ $\frac{dktdln}{dktdlq}$ $11$ $70$ $200$ $\frac{dktdln}{dktdlq}$ $9$ $50$ $400$ $\frac{dktdlq}{9}$ $50$ $400$ $\frac{dktdln}{11}$ $70$ $350$ $\frac{1}{11}$ $70$ $350$ $\frac{1}{200}$ $\frac{1}{9}$ $50$ $400$ $\frac{1}{11}$ $70$ $350$ $\frac{1}{200}$ $\frac{1}{20}$ $\frac{1}{20}$ $\frac{1}{200}$ $\frac{1}{200}$ $\frac{1}{20}$ $\frac{1}{200}$ $\frac{1}{200}$ $\frac{1}{200}$ $\frac{1}{200}$ $\frac{1}{20}$ $\frac{1}{200}$ $\frac{1}{200}$ $\frac{1}{200}$ $\frac{1}{200}$ $\frac{1}{200}$ </td <td><math>\frac{1q}{1d}</math> <math>1e4</math> <math>1e2</math> <math>\frac{13}{15}</math> <math>50</math> <math>250</math> <math>8</math>.'         <math>\frac{1d}{dl}</math> <math>11</math> <math>70</math> <math>200</math> <math>2</math>.         <math>\frac{11}{70}</math> <math>200</math> <math>2</math>.         <math>\frac{11}{70}</math> <math>200</math> <math>2</math>.         <math>\frac{11}{70}</math> <math>350</math> <math>2</math>.         <math>\frac{11}{70}</math> <math>11</math> <math>70</math> <math>350</math> <math>2</math>.         <math>\frac{11}{70}</math> <math>11</math> <math>70</math> <math>11</math> <math>70</math> <math>11</math> <math>\frac{11}{70}</math> <math>11</math> <math>10</math> <math>11</math> <math>70</math> <math>11</math> <math>70</math>&lt;</td>	$\frac{1q}{1d}$ $1e4$ $1e2$ $\frac{13}{15}$ $50$ $250$ $8$ .' $\frac{1d}{dl}$ $11$ $70$ $200$ $2$ . $\frac{11}{70}$ $200$ $2$ . $\frac{11}{70}$ $200$ $2$ . $\frac{11}{70}$ $350$ $2$ . $\frac{11}{70}$ $11$ $70$ $350$ $2$ . $\frac{11}{70}$ $11$ $70$ $11$ $70$ $11$ $\frac{11}{70}$ $11$ $10$ $11$ $70$ $11$ $70$ <
drid       15       50       175         iii       11       70       200         dktdlq       9       50       400         dktdlq       9       50       400         dktdlq       9       50       400         se: This is the total environment transitions the algorithm collects and equals to episode length × episodes is         .4       ALGORITHM DETAILS         .4.1       ALGORITHM HYPERPARAMETERS FOR DATA COLLECTION         /e use Q-Learning for data collection in all GridWorld environments. Its hyperparameters a sted in Table S2:         Table S2:       Q-Learning's hyper-parameters         discount       0.97         learning rate       1.0	Idd       101       15       50       175       1.         Idd       11       70       200       2.         Idd       9       50       400       2.         Idd       11       70       350       2.         is is the total environment transitions the algorithm collects and equals to episode length × episodes to co       ALGORITHM DETAILS         ALGORITHM DETAILS       I       ALGORITHM HYPERPARAMETERS FOR DATA COLLECTION         se Q-Learning for data collection in all GridWorld environments. Its hyperparameters are in Table S2:       Table S2: Q-Learning's hyperparameters         Table S2:       Q-Learning rate       0.97         learning rate       1.0         Pyperparameters       Values         discount       0.97         learning rate       1.0         P Methods, we use the same network structure for a fair comparison. The details of ork structure and hyperparameters are listed in Table S3. Their difference mainly locates in the xt/chunk length and memory recurrence as detailed in the following sections.
dktdln       11       70       200         dktdlq       9       50       400         11       70       350       350         e: This is the total environment transitions the algorithm collects and equals to episode length × episodes is      4        4       ALGORITHM DETAILS      4.1       ALGORITHM HYPERPARAMETERS FOR DATA COLLECTION         //e use Q-Learning for data collection in all GridWorld environments. Its hyperparameters a sted in Table S2:       Table S2: Q-Learning's hyperparameters          Myperparameters       values         discount       0.97         learning rate       1.0	idin       11       70       200       2         idid       9       50       400       2.4         idid       11       70       350       2.5         is is the total environment transitions the algorithm collects and equals to episode length × episodes to co       ALGORITHM DETAILS         ALGORITHM DETAILS       I       ALGORITHM HYPERPARAMETERS FOR DATA COLLECTION       see Q-Learning for data collection in all GridWorld environments. Its hyperparameters are in Table S2:       Table S2: Q-Learning's hyper-parameters         Table S2:       Q-Learning rate       1.0         Image: Algorithm of the set of the
dktdlq       9       50       400         dktdld       11       70       350         x: This is the total environment transitions the algorithm collects and equals to episode length × episodes is        4       ALGORITHM DETAILS        4.1       ALGORITHM HYPERPARAMETERS FOR DATA COLLECTION         //e use Q-Learning for data collection in all GridWorld environments. Its hyperparameters a sted in Table S2:         Table S2:       Table S2: Q-Learning's hyperparameters         //discount       0.97         learning rate       1.0	idiq       9       50       400       2.4         idid       11       70       350       2.         is is the total environment transitions the algorithm collects and equals to episode length × episodes to co       ALGORITHM DETAILS         ALGORITHM DETAILS       1       ALGORITHM HYPERPARAMETERS FOR DATA COLLECTION         se Q-Learning for data collection in all GridWorld environments. Its hyperparameters are in Table S2:       Table S2: Q-Learning's hyperparameters         in Table S2:       Table S2: Q-Learning rate       0.97         learning rate       1.0         2       NETWORK STRUCTURE & TRAINING HYPERPARAMETERS         D & ET methods, we use the same network structure for a fair comparison. The details of ork structure and hyperparameters are listed in Table S3. Their difference mainly locates in the xt/chunk length and memory recurrence as detailed in the following sections.
attald       11       70       350         x: This is the total environment transitions the algorithm collects and equals to episode length × episodes is4. ALGORITHM DETAILS      4. ALGORITHM DETAILS        4.1       ALGORITHM HYPERPARAMETERS FOR DATA COLLECTION         //e use Q-Learning for data collection in all GridWorld environments. Its hyperparameters a sted in Table S2:       Table S2: Q-Learning's hyperparameters          Table S2: Q-Learning is hyperparameters       Values         discount       0.97         learning rate       1.0	idid       11       70       350       2.         nis is the total environment transitions the algorithm collects and equals to episode length × episodes to co         ALGORITHM DETAILS         1       ALGORITHM HYPERPARAMETERS FOR DATA COLLECTION         se Q-Learning for data collection in all GridWorld environments. Its hyperparameters are in Table S2:         Table S2: Q-Learning's hyperparameters         isiscount       0.97         learning rate       1.0         2         NETWORK STRUCTURE & TRAINING HYPERPARAMETERS         JD & ET methods, we use the same network structure for a fair comparison. The details of ork structure and hyperparameters are listed in Table S3. Their difference mainly locates in the xt/chunk length and memory recurrence as detailed in the following sections.
A: ALGORITHM DETAILS 4. ALGORITHM DETAILS 4.1 ALGORITHM HYPERPARAMETERS FOR DATA COLLECTION /e use Q-Learning for data collection in all GridWorld environments. Its hyperparameters a sted in Table S2: Table S2: Q-Learning's hyper- parameters <u>hyperparameters values</u> <u>discount 0.97</u> <u>learning rate 1.0</u> 4.2 NETWORK STRUCTURE & TRAINING HYPERPARAMETERS or AD & ET methods, we use the same network structure for a fair comparison. The details twork structure and hyperparameters are listed in Table S3. Their difference mainly locates in to network length and memory recurrence as detailed in the following sections.	ALGORITHM DETAILS ALGORITHM DETAILS ALGORITHM HYPERPARAMETERS FOR DATA COLLECTION se Q-Learning for data collection in all GridWorld environments. Its hyperparameters are in Table S2: Table S2: Q-Learning's hyper- parameters <u>hyperparameters values</u> <u>discount 0.97</u> <u>learning rate 1.0</u> 2 NETWORK STRUCTURE & TRAINING HYPERPARAMETERS JD & ET methods, we use the same network structure for a fair comparison. The details of rk structure and hyperparameters are listed in Table S3. Their difference mainly locates in the xt/chunk length and memory recurrence as detailed in the following sections.
sted in Table S2:         Table S2: Q-Learning's hyper- parameters         hyperparameters values discount 0.97 learning rate 1.0         .4.2 NETWORK STRUCTURE & TRAINING HYPERPARAMETERS         or AD & ET methods, we use the same network structure for a fair comparison. The details twork structure and hyperparameters are listed in Table S3. Their difference mainly locates in t ontext/chunk length and memory recurrence as detailed in the following sections.	In Table S2:         Table S2: Q-Learning's hyper- parameters         hyperparameters values discount 0.97 learning rate 1.0         2       NETWORK STRUCTURE & TRAINING HYPERPARAMETERS         D & ET methods, we use the same network structure for a fair comparison. The details of rk structure and hyperparameters are listed in Table S3. Their difference mainly locates in the xt/chunk length and memory recurrence as detailed in the following sections.
hyperparameters       values         discount       0.97         learning rate       1.0         .4.2       NETWORK STRUCTURE & TRAINING HYPERPARAMETERS         or AD & ET methods, we use the same network structure for a fair comparison. The details etwork structure and hyperparameters are listed in Table S3. Their difference mainly locates in t intext/chunk length and memory recurrence as detailed in the following sections.	hyperparameters       values         discount       0.97         learning rate       1.0         2       NETWORK STRUCTURE & TRAINING HYPERPARAMETERS         AD & ET methods, we use the same network structure for a fair comparison. The details of ork structure and hyperparameters are listed in Table S3. Their difference mainly locates in the xt/chunk length and memory recurrence as detailed in the following sections.
discount 0.97 learning rate 1.0	discount       0.97         learning rate       1.0         2       NETWORK STRUCTURE & TRAINING HYPERPARAMETERS         AD & ET methods, we use the same network structure for a fair comparison. The details of ork structure and hyperparameters are listed in Table S3. Their difference mainly locates in the xt/chunk length and memory recurrence as detailed in the following sections.
Iearning rate 1.0 Iearning rate 1.0	learning rate       1.0         2       NETWORK STRUCTURE & TRAINING HYPERPARAMETERS         AD & ET methods, we use the same network structure for a fair comparison. The details of rk structure and hyperparameters are listed in Table S3. Their difference mainly locates in the xt/chunk length and memory recurrence as detailed in the following sections.
	2 NETWORK STRUCTURE & TRAINING HYPERPARAMETERS D & ET methods, we use the same network structure for a fair comparison. The details of ork structure and hyperparameters are listed in Table S3. Their difference mainly locates in the xt/chunk length and memory recurrence as detailed in the following sections.
4.2 NETWORK STRUCTURE & TRAINING HYPERPARAMETERS or AD & ET methods, we use the same network structure for a fair comparison. The details etwork structure and hyperparameters are listed in Table S3. Their difference mainly locates in t ontext/chunk length and memory recurrence as detailed in the following sections.	2 NETWORK STRUCTURE & TRAINING HYPERPARAMETERS D & ET methods, we use the same network structure for a fair comparison. The details of ork structure and hyperparameters are listed in Table S3. Their difference mainly locates in the xt/chunk length and memory recurrence as detailed in the following sections.
4.2 NETWORK STRUCTURE & TRAINING HYPERPARAMETERS or AD & ET methods, we use the same network structure for a fair comparison. The details etwork structure and hyperparameters are listed in Table S3. Their difference mainly locates in t ontext/chunk length and memory recurrence as detailed in the following sections.	2 NETWORK STRUCTURE & TRAINING HYPERPARAMETERS AD & ET methods, we use the same network structure for a fair comparison. The details of ork structure and hyperparameters are listed in Table S3. Their difference mainly locates in the xt/chunk length and memory recurrence as detailed in the following sections.
4.2 NETWORK STRUCTURE & TRAINING HYPERPARAMETERS or AD & ET methods, we use the same network structure for a fair comparison. The details stwork structure and hyperparameters are listed in Table S3. Their difference mainly locates in t ontext/chunk length and memory recurrence as detailed in the following sections.	2 NETWORK STRUCTURE & TRAINING HYPERPARAMETERS AD & ET methods, we use the same network structure for a fair comparison. The details of ork structure and hyperparameters are listed in Table S3. Their difference mainly locates in the xt/chunk length and memory recurrence as detailed in the following sections.
etwork structure and hyperparameters are listed in Table S3. Their difference mainly locates in t ontext/chunk length and memory recurrence as detailed in the following sections.	ork structure and hyperparameters are listed in Table S3. Their difference mainly locates in the xt/chunk length and memory recurrence as detailed in the following sections.
Table S3: Network structure & training hyperparameters	Table S3: Network structure & training hyperparameters
hyperparameters values	hyperparameters values
Layers 4	Layers 4
Hidden dimension 512	Hidden dimension 512
Dropout 0.1	
Feedforward dimension 2048	Dropout 0.1
Attention head 64	Dropout0.1Feedforward dimension2048
	Dropout0.1Feedforward dimension2048Attention head64
Optimizer AdamW	Dropout0.1Feedforward dimension2048Attention head64OptimizerAdamW
OptimizerAdamWOptimizer weight decay1e-4	Dropout0.1Feedforward dimension2048Attention head64OptimizerAdamWOptimizer weight decay1e-4
OptimizerAdamWOptimizer weight decay1e-4Optimizer betas(0.9, 0.999)	Dropout0.1Feedforward dimension2048Attention head64OptimizerAdamWOptimizer weight decay1e-4Optimizer betas(0.9, 0.999)
OptimizerAdamWOptimizer weight decay1e-4Optimizer betas(0.9, 0.999)Learning rate2e-4	Dropout0.1Feedforward dimension2048Attention head64OptimizerAdamWOptimizer weight decay1e-4Optimizer betas(0.9, 0.999)Learning rate2e-4
OptimizerAdamWOptimizer weight decay1e-4Optimizer betas(0.9, 0.999)Learning rate2e-4Learning rate schedulerlinear warm-up & cosine decay	Dropout0.1Feedforward dimension2048Attention head64OptimizerAdamWOptimizer weight decay1e-4Optimizer betas(0.9, 0.999)Learning rate2e-4Learning rate2e-4
OptimizerAdamWOptimizer weight decay1e-4Optimizer betas(0.9, 0.999)Learning rate2e-4Learning rate schedulerlinear warm-up & cosine decayBatch size256	Dropout0.1Feedforward dimension2048Attention head64OptimizerAdamWOptimizer weight decay1e-4Optimizer betas(0.9, 0.999)Learning rate2e-4Learning rate schedulerlinear warm-up & cosine decayBatch size256
OptimizerAdamWOptimizer weight decay1e-4Optimizer betas(0.9, 0.999)Learning rate2e-4Learning rate schedulerlinear warm-up & cosine decayBatch size256LaverNormpostnorm	Dropout0.1Feedforward dimension2048Attention head64OptimizerAdamWOptimizer weight decay1e-4Optimizer betas(0.9, 0.999)Learning rate2e-4Learning rate schedulerlinear warm-up & cosine decayBatch size256LaverNormpostnorm
OptimizerAdamWOptimizer weight decay1e-4Optimizer betas(0.9, 0.999)Learning rate2e-4Learning rate schedulerlinear warm-up & cosine decayBatch size256LayerNormpostnormLoss functionCrossEntropyLoss	Dropout0.1Feedforward dimension2048Attention head64OptimizerAdamWOptimizer weight decay1e-4Optimizer betas(0.9, 0.999)Learning rate2e-4Learning rate schedulerlinear warm-up & cosine decayBatch size256LayerNormpostnormLoss functionCrossEntropyLoss
OptimizerAdamWOptimizer weight decay1e-4Optimizer betas(0.9, 0.999)Learning rate2e-4Learning rate schedulerlinear warm-up & cosine decayBatch size256LayerNormpostnormLoss functionCrossEntropyLossLabel smooth0.1	Dropout0.1Feedforward dimension2048Attention head64OptimizerAdamWOptimizer weight decay1e-4Optimizer betas(0.9, 0.999)Learning rate2e-4Learning rate schedulerlinear warm-up & cosine decayBatch size256LayerNormpostnormLoss functionCrossEntropyLossL abel smooth0.1

## Table S1: Detailed settings of environments & tasks

# A.4.3 DETAILS OF ADF TRAINING HYPERPARAMETERS

We list the hyperparameter details of algorithm ADF in Table S4.

setting	episode length	context length	training steps
dr	20	50	1e5
dktd	50	125	2e5
drln	50	125	2e5
drlq	50	125	3e5
drlq	50	125	2e5
dktdln	70	125 <sup><i>a</i></sup>	4e5
dktdlq	50	125	6e5
dktdld	70	125 <sup><i>a</i></sup>	6e5

### Table S4: Details of ADF training hyperparameters

### A.4.4 DETAILS OF ADR TRAINING HYPERPARAMETERS

We list the hyperparameter details of algorithm ADR in Table S5.

Table S5: Details of ADR training hyperparameters

setting	episode length	context length	training steps (ADF $\times$ 5)
dr	20	10	5e5
dktd	50		1e6
drln	50		1e6
drlq	50		1.5e6
drld	50	25	1e6
dktdln	70		2e6
dktdlq	50		3e6
dktdld	70		2e6

### A.4.5 HYPERPARAMETERS OF MEMORY-EFFICIENT ALGORITHMS

We list the hyperparameter details of Memory-Efficient algorithms (MEM, XL and ED) in Table S6.

### Table S6: Details of ET training hyperparameters

settings	chunk length	memory capacity	ADR context length	<b>training steps</b> (ADF $\times$ 5 / ratio) $^{\alpha}$
dr	10	10	10	5e5 / 1.5
dktd				1e6 / 1.7
drln				1e6 / 1.7
drlq				1.5e6 / 1.7
drld	25	50	25	1e6 / 1.7
dktdln				2e6 / 1.7
dktdlq				3e6 / 1.7
dktdld				2e6 / 1.7

 $\alpha$ : we slightly reduced the training steps since the ET methods' effective context is larger than ADR's.

16

824 825 826

839

840 841

823

810