# MEMREASONER: A MEMORY-AUGMENTED LLM AR-CHITECTURE FOR MULTI-HOP REASONING

Anonymous authors

004

010

011

012

013

014

015

016

017

018

019

021

Paper under double-blind review

#### ABSTRACT

Recent benchmarks suggest that there remains significant room to improve large language models' ability to robustly reason across facts distributed in extremely long documents. In this work, we propose MemReasoner, a new memoryaugmented LLM architecture that is equipped to perform temporal reasoning, along with multiple computational steps, over the context stored in the latent memory, and is trained with supervision on intermediate steps and final outcome. Experiments show that MemReasoner trained on the core reasoning facts generalizes better, compared to the off-the-shelf large language models as well as fine-tuned recurrent models, on an unseen test distribution where the required facts are scattered across long natural text up to 128k tokens. Further, MemReasoner demonstrates robust reasoning performance relative to the baselines, when the answer distribution or number of hops in test samples differs from that in the training set.

### 023 1 INTRODUCTION

Transformer-based large language models (LLMs) have recently shown impressive performance in many natural language processing (NLP) tasks, including machine translation, question answering, and reading comprehension, demonstrating signature of general reasoning abilities. However, when restricted to individual NLP reasoning benchmarks, particularly those that require logical reasoning, current LLMs typically perform poorly despite efforts to improve accuracy through prompt engineering (Wei et al., 2022; Min et al., 2022). As such, more evidence seems to support the hypothesis that powerful LLMs often learn statistical features and correlations to simulate reasoning rather than performing true reasoning (Ruder, 2021).

The recently introduced BABILong benchmark further establishes this point, as it is designed to 033 test LLM's ability to reason across facts distributed in extremely long documents (Kuratov et al., 034 2024). BABILong is developed based on the bAbi benchmark (Weston et al., 2015), which is composed of 20 reasoning tasks. These include fact chaining, simple induction, deduction, counting, 035 and handling lists/sets (Weston et al., 2015). This set of tasks was designed as prerequisites for any 036 system that aims to having a conversation with a human. BABILong further introduces irrelevant 037 natural text from the PG19 book corpus (Rae et al., 2019) into the original context to make it artificially longer and include distracting text, while the underlying reasoning task remains the same. For examples of the task samples in BABILong, see Figure 1. Experiments with popular transformer-040 based LLMs shows that present days' transformer-based language models effectively utilize only 041 10-20% of the context and their performance declines sharply with increased reasoning complexity. 042 Retrieval-augmented generation with LLMs at best can provide 60% accuracy for a simple QA task 043 that requires extracting single evidence from the context. Interestingly, a memory-augmented trans-044 former architecture, namely Recurrent Memory Transformers (RMT) (Bulatov et al., 2022) shows 045 the highest performance on BABILong benchmark; suggesting that the long-term recurrent memory of the context helps. RMT in that case is trained on longer BABILong samples with supervision 046 on final answer. However, as we will demonstrate in this work, RMT when trained on bAbi sam-047 ples with supervision on final answer reconstruction, does not generalize well on BABILong test 048 set. This observation suggests that memory-augmented LLMs can further benefit from additional 049 supervision, when available. Additionally, enabling multiple hops over the memory per final answer 050 reconstruction can help the model perform well, when the task demands so. 051

In this work, we provide an alternative language model architecture that is designed to naturally
 handle recurrent processing over long context *that is not seen by the model during training*. Our goal is to provide a more effective and robust solution for handling multi-hop generative QA tasks,



Experiments showing that the proposed MemReasoner architecture indeed learns multi-step processing over the context to solve the QA task, as evident by its robust performance when the answers in the training data differ from those in the test samples within the same task, and when the model trained on bABi task 2 is tested on longer BABILong task 1 samples.

# 108 2 RELATED WORK

LLM Reasoning Logical reasoning, a critical aspect for advancing many scientific fields, involves 110 deducing new conclusions from existing facts and rules. To derive the final answer, such reasoning 111 challenges often require multiple steps to be executed effectively and in the right order. For instance, 112 with facts like "John picked up the football" and "John went to the bedroom", a logical process will 113 be to deduce that the football's current location is bedroom. Despite showing advanced ability 114 to learn from instructions and in-context demonstrations to answer questions (Brown et al., 2020; 115 Min et al., 2022), LLMs struggle with complex logical reasoning, especially multi-step reasoning 116 (Liu et al., 2023a). This failure has been attributed to the autoregressive nature of LLMs (Stechly et al., 2024), which can be characterized by "System 1" (Kahneman, 2011), a mode of thought that 117 118 is fast, instinctive but less accurate. To address this limitation, recent work proposes prompting LLMs to mimic generating intermediate chain of thought (CoT) reasoning steps (Wei et al., 2022), 119 providing access to external tools/verifiers (Schick et al., 2023), or a combination of both (Paranjape 120 et al., 2023), to mimic the process of generating deliberative and logical thinking steps, i.e., the 121 "System 2" mode. Another direction currently being explored is to train reward models to rank 122 the candidate solutions or rank the intermediate steps (Khalifa et al., 2023; Wang et al., 2024). 123 Different from these works, MemReasoner does not rely on deliberate prompt engineering or access 124 to external tools, neither does it require feedback from an external reward model. Instead, inspired 125 by the distinction between System 1 and System 2-like thinking, MemReasoner utilizes the decoder 126 for fast generation and the memory module for slow reasoning, which are two components tightly 127 integrated via training. In that sense, MemReasoner is closer to the line of works that use (generated) 128 rationales for supervised finetuning or for preference tuning of LLMs to enhance their reasoning 129 abilities (Zelikman et al., 2022; Pang et al., 2024). However, it remains unexplored how those approaches perform on iterative reasoning tasks over lengthy context that is unseen during training. 130

131 Long-context Modeling The scope of the present study encompasses two distinct challenges around 132 multi-step reasoning tasks, namely (1) processing very long context and (2) "hopping" over that con-133 text in a temporally-aware manner to link disjoint pieces of information and generate answers based 134 on that. On the first challenge, vanilla transformer-based models struggle due to quadratic time and 135 space complexity of self-attention and the increasing memory requirement of the key-value cache during generation. Recently, there has been significant progress in long-context modeling with 136 transformers by using a mix of local and global attention (Munkhdalai et al., 2024), by continued 137 pretraining on longer sequences (Xiong et al., 2023; Ding et al., 2024), by context window sliding 138 and segmentation (Ratner et al., 2023), and by applying position extrapolation or interpolation to 139 extend input length beyond the training phase (Press et al., 2022; Su et al., 2023). Promising alter-140 native directions include the development of novel recurrent architectures (Bulatov et al., 2022) and 141 state-space-models (Gu & Dao, 2023). Nevertheless, many of these techniques require training on 142 longer sequences. Additionally, a number of studies and benchmarks suggest that the long-context 143 LLMs may not be able to fully utilize their context window, and therefore performance degrades on 144 simple retrieval and complicated reasoning tasks as the input length grows and/or the position of the 145 answer varies within the context (Hsieh et al., 2024; Yuan et al., 2024; Liu et al., 2023b; Levy et al., 146 2024).

147 Status Check on LLM Reasoning Consequently, in parallel to impressive advances in LLMs abil-148 ities, caution has been raised on the discrepancy between claimed reasoning abilities as per stan-149 dardized benchmarks and true reasoning skills. The scientific community has advocated for careful 150 investigations of issues such as data contamination, performance robustness and generalization, and 151 flawed reasoning benchmark that supports "shortcut learning" (Mitchell, 2023; Wu et al., 2024). For 152 example, the presence of reasoning shortcuts in the task samples themselves has been reported in the HotPotQA dataset, which does not ensure language models are actually being required to perform 153 multi-hop processing over the context (Jiang & Bansal, 2019). Recently, a number of tasks and 154 benchmarks have been developed to address these issues (Valmeekam et al., 2022; Kuratov et al., 155 2024; Nezhurina et al., 2024). Along this line, we here show the generalization robustness of Mem-156 Reasoner across (i) "unseen" context that consists of varying length of irrelevant natural text and (ii) 157 answer distribution that is different from the training distribution.

158 159 160

### 3 MULTI-STEP REASONING WITH MEMREASONER

161 The key components of MemReasoner involve an LM encoder, an episodic memory module, and an LM decoder (see Figure 2a). The role of the episodic memory module is to enable *write* of the

162 context encodings in the memory, to allow performing search over the context encodings and read 163 from them, in order to feed the decoder to execute the task. Given a logical reasoning task for which 164 the supporting facts (reasoning process) and the final answer (reasoning outcome) are available, the 165 MemReasoner architecture is trained to recover the supporting facts and the final answer. A search 166 in the latent memory space is performed during training in order to correctly output the final answer and the supporting facts. An additional point worth mentioning is that, MemReasoner also is trained 167 to learn the relative order of the supporting facts in the context, which is crucial for reconstructing 168 an agent's or object's most recent location, as required by the bAbi reasoning tasks. Details are provided below. 170

171 3.1 PRELIMINARY

172 Let  $\mathcal{X}$  be the LM input space,  $\mathcal{Z}$  be the latent space, and  $\mathcal{Y}$  be the LM output space. Larimar (Das et al., 2024) features an encoder e that maps an input to an embedding  $z \in \mathcal{Z} \subseteq \mathbb{R}^D$ , and a memory module  $\mathcal{M}$ . The memory M is adaptable in the sense, that it allows "write" and "read" operations as episodes (aka, contexts C, where each context is comprised of E sentences) arrive, i.e.,  $\hat{M} = write(M, z), z_{read} = read(\hat{M}, z)$ , wherein  $\hat{M}$  is the updated memory after an write. And, a decoder d that performs generations conditioned on the memory readout  $z_{read}$ .

#### 178 179 3.2 LARIMAR FRAMEWORK

Now, suppose one is given an input context  $C = \{c_1, ..., c_E\}$  with E denoting the length of the context, and the target task is to answer a question q conditioned on the given context C. To approach 181 the task within the Larimar framework, the input, both context C and query q, are encoded to their 182 latents  $(z_1, \ldots, z_E$  and  $z_q)$  via the encoder e. Next, let  $M_0$  be the initial memory, write the context to 183 the memory via a *write* operation. To do so, Larimar follows the earlier works on Kanerva Machine 184 (Wu et al., 2018), which is inspired by Kanerva's sparse distributed memory model (Kanerva, 1988), 185 where the memory is viewed as a global latent variable in a generative model. In this framework, 186 the goal is to learn a memory dependent data prior and learnable addresses, where the memory 187 update and read/write are considered as Bayesian inference, i.e., the posterior parameters are updated 188 as new data arrives. Later, (Pham et al., 2022) reformulated the encoding of new memories and 189 decoding data from memories from Bayesian updates to an equivalent minimization problem, which 190 essentially amounts to solving a linear system of equations, efficiently done via computing matrix pseudo inverses indicated by † hereafter. Therefore, memory is updated via the write operation such 191 that,  $\hat{M} = (Z_{\xi}M_0^{\dagger})^{\dagger}Z_{\xi}$ , where  $Z_{\xi} = [z_1 + \xi_1, z_2 + \xi_2, \dots, z_E + \xi_E]$  and  $\xi_i \sim \mathcal{N}(0, \sigma_{\xi}^2 I)$ . 192

Then, the *read* operation translates the query embedding from the lens of the encoded memory to a query readout  $z_r$  via  $z_r = (z_q \hat{M}^{\dagger} + \eta) \hat{M}$ , where  $\eta \sim \mathcal{N}(0, \sigma_{\eta}^2 I)$ . Lastly, the decoder d decodes the query q conditioned on the readout by using a *learnable* broadcasting parameter  $W_M$  that casts  $z_r$ to each decoder layer and obtains  $h_k^m$  that serves as the past key values for  $k = 1, \ldots, L$ , where Lis the number of layers in the decoder.

We use this memory-augmented LLM architecture of Larimar and the operations as backbone for 199 MemReasoner, due to its memory and space-efficient read/write abilities and demonstrated gener-200 alizability at test-time. It is worth mentioning the earlier works on memory-augmented neural nets, 201 which use a recurrent neural net together with an external memory, have investigated ideas like 202 temporal feature learning and iterative hops over context, for example, see (Weston et al., 2014; 203 Sukhbaatar et al., 2015). However, to our knowledge, this is the first study to enable those opera-204 tions around the explicit episodic memory of a transformer-based LLM during training and test the 205 resulting model's generalizability on a long-context reasoning benchmark like BABILong. 206

## 207 3.3 MEMORY WITH TEMPORAL ORDER

Recall, the latent encoding of facts  $\{z_1, ..., z_E\}$  within a context episode C are written in the memory M in an order-invariant manner. However, many multi-step reasoning tasks require some notion of temporal context. For example, when answering "where is John?" in the context of "... John is in the bathroom. ... John goes to the garden." ("..." denotes irrelevant facts), there should be a mechanism in place to guarantee that the memory encodes the correct temporal order of the facts, and the readout should reflect "John goes to the garden." as the supporting fact instead of "John is in the bathroom.".

To introduce some temporal notion within the context, in MemReasoner we introduce a temporal encoding module  $\mathcal{P}$  that transforms *un-ordered* fact latents  $\{z_1, ..., z_E\}$  within a context episode to



Figure 2: A diagram of the pipeline for reasoning with MemReasoner. (a) Conceptual overview 234 of the framework. (b) Detailed architecture. q denotes the query,  $c_1, ..., c_E$  denotes the context for answering the query.  $z_q$  denotes the encoding of the query while  $\{z_1, ..., z_E\}$  denote encodings of 235 each line of the context. We use  $\tilde{z}$  to denote temporally encoded latents. 236

237

257

233

238 their ordered counterparts  $\{\tilde{z}_1, ..., \tilde{z}_E\}$ . The temporal encoding module is generic and allows any 239 structure featuring sequentiality within context. In practice, we investigate two general types of 240 encoding methods, un-parameterized methods such as Sinusoidal Positional Encoding and parame-241 terized methods such as GRUs. 242

243 Positional Encoding. We compute positional encodings for each line of context within the episode 244 by using sine and cosine functions similar to (Vaswani et al., 2017). Additionally, we experiment 245 with positional encoding which assigns encodings starting from the last element of the episode. The structure ensures that for contexts of different length, the last lines of the contexts are encoded 246 similarly, which is useful for QA tasks in which the most recent information is more relevant for 247 answering the question. 248

249 Finally, to convert  $\{z_1, ..., z_E\}$  to  $\{\tilde{z}_1, ..., \tilde{z}_E\}$  with positional encodings, we add the computed 250 positional encodings to the input. 251

**GRU.** We also investigate learnable encodings via a bidirectional GRU unit. For these, we treat 252  $\{z_1, ..., z_E\}$  as the sequence passed as input into the GRU and simply let  $\{\tilde{z}_1, ..., \tilde{z}_E\}$  be the sequen-253 tial outputs of the GRU. 254

255 These ordered context embeddings  $\{\tilde{z}_1, ..., \tilde{z}_E\}$  are then written to memory via Larimar's write 256 operation.

- 3.4 ITERATIVE READ AND QUERY UPDATE 258
- 259 A typical multi-step reasoning task often inherently requires "hops" between facts until the final 260 solution is found. Additionally, the query embedding can be updated accordingly to reflect the most recent hop. 261

262 In order to perform hopping between facts, we first recall the three key components interacting with 263 the memory module  $\mathcal{M}$ , the fact embeddings ( $\{z_1, \ldots, z_E\}$ ) within a context episode, the query 264 embedding  $z_a$ , and the memory readout  $z_r$ . Let us further consider M stores facts that have been 265 ordered temporally  $\{\tilde{z}_1, ..., \tilde{z}_E\}$ . 266

To enable *iterative read*, we pass  $z_q$  through a linear layer to obtain  $\hat{z_q} = W_q z_q$  before the *read* operation from the memory, where  $W_q \in \mathbb{R}^{D \times D}$  is a learnable parameter that absorbs the scale 267 268 changes introduced by the position encoding in the memory. Specifically, different from Section 3.1, 269 here we have  $z_r = (\hat{z_q}\hat{M}^{\dagger} + \eta)\hat{M}$ .

To update the query, we first update the query latent and let  $z_q \leftarrow z_q + \alpha \cdot z_r$ , where  $\alpha \in \mathbb{R}$  is a hyperparameter to balance the load from the previous readout. The updated query is then fed into the memory module for another *read* operation to obtain a new  $\tilde{z}_r$ . The query update procedure is repeated until the readout converges (i.e.  $||\tilde{z}_r^t - \tilde{z}_r^{t+1}||_2 < \tau$  where  $\tilde{z}_r^t$  denotes the readout at time t and  $\tau$  is a hyperparameter) or until it reaches a fixed number of maximum iterations.

275 276 3.5 FULL WORKFLOW

Now that we have discussed all components of MemReasoner, we elaborate the full pipeline in the following and provide a visualization in Figure 2b.

279 Consider an input context  $C = \{c_1, ..., c_E\}$ , a question q, an encoder e, a temporal encoding module  $\mathcal{P}$ , an initial memory module  $\mathcal{M}$ , and a decoder d. We first encode the context C and query q to their latents,  $z_1, \ldots, z_E$  and  $z_q$ , via encoder e. Then, we follow Section 3.3 and transform 281 282  $z_1, \ldots, z_E$  to  $\tilde{z}_1, \ldots, \tilde{z}_E$ . Next, we write the ordered context  $\tilde{z}_1, \ldots, \tilde{z}_E$  to the memory and obtain M. 283 Subsequently, we read using the query latent from the memory and perform query and read updates 284 according to Section 3.4. After we have obtained a  $\tilde{z}_r$  as a final readout which does not undergo update anymore, we map  $\tilde{z}_r$  to the corresponding unordered encoding in M. This is because we 285 only want the additional position information to be used when locating the most relevant contexts, but not during the decoding - if being fed to the decoder, the decoder may overfit to the ordering 287 information in the latents. We do this by first finding the index of the most similar ordered latent encoding  $i = \arg \min_{j \in \{1,...,E\}} ||\tilde{z}_r - \tilde{z}_j||_2$  and then obtaining the corresponding encoding  $z_i$  from 289 the unordered encodings (prior to undergoing temporal encoding in Figure 2)  $\{z_1...z_E\}$ . Lastly, the 290 decoder d decodes the prompt  $P_a$  given for answer generation conditioned on  $z_i$ . We provide the 291 full pseudocode in Algorithm 1. 292

#### 293 3.6 TRAINING OBJECTIVES

Г

294 Let  $\mathcal{D}_{reason}$  denote the reasoning data distribution while  $\mathcal{D}_{pretrain}$  denotes the pretraining data distri-295 bution. Each sample from  $\mathcal{D}_{\text{reason}}$  is of the form (q, C, S, a) where q is the query,  $C = \{c_1, ..., c_E\}$ 296 are the facts in the context, S is a set of indices corresponding to supporting facts (we will use  $S_i$ to denote the *i*th supporting fact index in S), and a is the answer. Meanwhile the pretraining dis-297 tribution corresponds to a generic corpus, e.g. Wikipedia. Let e denote the encoder, d denote the 298 decoder, t denote temporal encoding,  $\tilde{z}_i^t$  denote the *i*th temporally encoded readout from iterative 299 reading with  $\tilde{z}_{r}^{0} = q, z_{r}^{i}$  represent the unordered encoding corresponding to the *i*th ordered readout, 300 and  $P_a$  and  $P_s$  denote the prompts for generating the answer and supporting fact respectively. To 301 train the model, we utilize the following loss function in Equation 1. 302

303

$$L = \mathbb{E}_{(q,C,S,a)\sim\mathcal{D}_{\text{finetune}}} \left[ \underbrace{\mathbb{E}_{z_r^{|S|}\sim p(z_r^{|S|}|q,M,\tilde{z}_r^0\dots\tilde{z}_r^{|S|-1})} \ln p(a|z_r^{|S|}, P_a)}_{\text{reconstruction of answer}} + \alpha \sum_{i=1}^{|S|} \underbrace{\mathbb{E}_{z_r^i \sim p(z_r^i|M,\tilde{z}_r^0\dots\tilde{z}_r^{i-1})} \ln p(c_{S_i}|z_r^i, P_s)}_{\text{reconstruction of supporting facts}} + \beta \sum_{s \in S} \underbrace{\ln p(d(e(c_s)))}_{\text{autoencoding of supporting fact}} \right]$$

$$+\delta \underbrace{\sum_{i=1}^{|S|} \mathbb{E}_{\tilde{z}_{r}^{i} \sim p(\tilde{z}_{r}^{i}|q,M,\tilde{z}_{r}^{0}\dots\tilde{z}_{r}^{i-1})}}_{\text{ordering loss}} \int +\rho \underbrace{\mathbb{E}_{x \sim \mathcal{D}_{\text{pretrain}}} \ln p(d(e(x)))}_{\text{autoencoding of pretraining dataset}}$$
(1)

 $\alpha, \beta, \delta$  and  $\rho$  are hyperparameters controlling regularization strength and  $\ell_{\text{order}}$  is given by

$$v(z_r) = \operatorname{softmax}([-||t(e(c_1)) - z_r||_2, ..., -||t(e(c_E)) - z_r||_2]^{\mathsf{T}})$$
  
$$\ell_{\operatorname{order}}(z_r, s) = -\ln v(z_r)_s$$
(2)

320 321

318 319

The first and second terms correspond to the reconstruction loss of the answer and the supporting fact(s) with respect to the corresponding prompt for obtaining the answer  $P_a$  and final readout, the third and the fifth terms correspond to the autoencoding loss of the supporting fact(s) and pretraining

| 324<br>325 | Ā          | lgorithm 1:   |
|------------|------------|---|
| 326        | 1]         | <b>Function</b> IterativeRead( $q, \{c_1,, c_E\}, \alpha, \tau$ ):                                      |
| 327        |            | // $q$ denotes the query tokens while $\{c_1,,c_E\}$ denote the $E$                                     |
| 328        |            | lines of context tokens, $lpha$ is a hyperparameter for the   |
| 329        |            | query update, $	au$ is a threshold hyperparameter for   |
| 330        |            | terminating iterations, $P_a$ is the prompt given to the  |
| 221        |            | decoder for answer generation   |
| 221        |            | // encode query and context lines with encoder  |
| აა∠<br>ეეე | 2          | $z_q \leftarrow \texttt{encode}(q)$ for $i \leftarrow 1$ to $E$ do                                      |
| 222        | 3          | $ z_i \leftarrow \text{encode}(c_i)$  |
| 225        | 4          | ena   |
| 333        |            | apply temporal encoding over the sequence of context lines  |
| 330        | 5          | $\tilde{\alpha}_1 \qquad \tilde{\alpha}_2 \leftarrow \text{temporalEncoding}(\alpha_1 \qquad \alpha_2)$ |
| 337        | 3          | $\hat{M}_{1}, \dots, \hat{z}_{E}$ (composition country $(z_{1}, \dots, z_{E})$                          |
| 338        | 6          | $W \leftarrow WIILE(Z_1,, Z_E)$   |
| 339        | 7          | $\tilde{\gamma} \leftarrow \text{queryUpdate}(\chi, \alpha, \tau)$                                      |
| 340        | '          | $Z_r$ (queryoptate $(Z_q, a, r)$ )<br>// Map to latent prior to performing temporal encoding            |
| 341        | 8          | $i^* \leftarrow \arg\min_{i \in \{1, \dots, p\}} \ \hat{z}_i - \hat{z}_r\ _2$                           |
| 342        | 9          | <b>return</b> decode $(z_{i*}, W_M, P_c)$ // generate the answer with the                               |
| 343        | -          | decoder, $W_M$ is a learnable parameter which interfaces the  |
| 344        |            | $z_{i^*}$ with the decoder  |
| 345        | 10         |   |
| 340        | 11         |   |
| 347        | 12 F       | unction <code>temporalEncoding</code> ( $\{z_1,,z_E\}$ , method):                                       |
| 348        |            | // temporally encode the sequence $\{z_1,,z_E\}$  |
| 349        | 13         | if method = $PE$ then   |
| 350        | 14         | $  \operatorname{return} \{z_i + PE(i)   \forall i \in \{1,, E\} \}$                                    |
| 351        | 15         | else il metnod = $GRU$ (nen   |
| 352        | 16         | return $O(\{z_1,, z_E\})$   |
| 353        | 17<br>10 F | unction guery Undate $(x, \alpha, \tau)$ .  |
| 354        | 10 1       | // given the guery encoding $\alpha_{\tau}$ and threshold $\tau_{\tau}$ perform                         |
| 300        |            | iterative reading and update guery  |
| 330        | 19         | $\tilde{z}_r \leftarrow \text{read}(W_a z_a, M)$ // $W_a$ is learned parameter                          |
| 357        | 20         | $z_a = z_a + \alpha \tilde{z_r}$ // query update  |
| 358        | 21         | $\tilde{z}_{r,\text{next}} \leftarrow \text{read}(W_q z_q, M)$  |
| 359        | 22         | do  |
| 360        | 23         | $\tilde{z}_r \leftarrow \tilde{z}_{r,\text{next}}$  |
| 361        | 24         | $z_q = z_q + \gamma \tilde{z}_r$  |
| 362        | 25         | $ z_{r,\text{next}} \leftarrow \text{read}(W_q z_q, M)$   |
| 363        | 26         | while $  \tilde{z}_{r,next} - \tilde{z}_r  _2 > \tau$   |
| 364        | 27         | return $z_r$  |
| 365        | 28         |   |

data. The fourth term is a loss for encouraging the index of the most similar entry (by 12 distance) to the ordered readout at each iteration to match the index of the supporting fact through computing the cross entropy.

### 4 EXPERIMENTAL DETAILS AND RESULTS

## 4.1 DATASETS AND DATA PRE-PROCESSING

In the main paper, we utilize tasks 1 and 2 from the synthetic bAbi benchmark as our testbed. We
also report results on Variable Tracking task from the RULER benchmark (Hsieh et al., 2024) in
appendix. The bAbi datasets were prepared by synthesizing relations among characters and objects
across various locations, each represented as a fact, such as "Mary traveled to the garden". Task 1

378 requires performing a single hop to find answer, whereas task 2 requires gathering two supporting 379 facts in the right order (see Fig 1). These single to multi-hop QA tasks from BABILong benchmark 380 together provide a controlled setting for evaluating LLMs' ability to reason over long context, where 381 the difficulty of the task can be varied by changing the length of irrelevant text. The nature of this 382 benchmark, where the synthetic sentences corresponding to the actual reasoning task are hidden inside irrelevant but lengthy naturally occurring text, keeps it at a low risk of data contamination to 383 training sets of todays' LLMs. And finally, BABILong leaderboard shows tasks 1 and 2, while being 384 simple enough, are challenging enough for off-the-shelf LLMs to solve. 385

We finetune MemReasoner separately on original bAbi task 1 and task 2 training split, each consisting of 10k samples (Weston et al., 2015). We then evaluate on the test set of the corresponding task from bAbi as well as from BABILong (Kuratov et al., 2024), in which the core reasoning facts from bAbi is distributed over arbitrarily long documents. Here we benchmark MemReasoner on BABILong test samples of up to 128k tokens.

For preprocessing bAbi data, we treat each training sample comprised of multiple facts as a single context episode, and individual sentence within that context as an instance within that episode. Each fact within an episode contains up to 64 tokens.For BABILong and for Wikipedia, if sentences are longer than 64 tokens, we split the sentences at multiples of 64 tokens.

We initiate MemReasoner finetuning from Larimar checkpoint pretrained on Wikitext (obtained 396 by following the training protocol described in (Das et al., 2024)), which uses a Bert-large as the 397 encoder and a GPT2-large as the decoder (For extension to MemReasoner with GPTJ-6B, see ap-398 pendix). The number of parameters in MemReasoner is 1.4B. The slot size in the memory is 512. 399 During finetuning, we randomly sample a batch of pretraining data (Wikipedia) of the same size as 400 the batch of finetuning data (bAbi) for computing the autoencoding loss on the pretrain dataset of 401 2M samples. We generate the answer to the question by passing a prompt to the decoder (i.e. in the 402 case of bAbi Task1-2, the prompt has the from " $\langle BOS \rangle X$  is in the" where X denotes subject of 403 the query).

404 We train MemReasoner models for 200 epochs using Adam optimizer with learning rate 5e-6. We 405 set batch size to be 10. Additionally, we set query update parameter  $\alpha = 1$ . The maximum episode 406 length varies from 14 (bAbi Task 1) to 72 (bAbi Task 2). Which means that MemReasoner has been 407 exposed to a maximum of 90 and 573 tokens during finetuning on task 1 and task 2, respectively, 408 whereas at test-time the model is exposed to contexts that are up to 128k tokens long. Since bAbi 409 Task 1 is a single hop task, we do not perform query update during either training or inference. When 410 fine-tuning on bAbi Task 2, we perform a fix number of 2 hop (equivalent to 1 query update) during 411 the training. With bAbi Task 2 fine-tuned MemReasoner, we re-use the "2 hop" setting at inference 412 on all tasks, including bAbi Task 2 and BABILong Task1/2. We consistently use query update parameter  $\alpha$ =8 throughout our experiments and include an ablation study on  $\alpha$  in the appendix. 413 Due to the page limit, we also defer ablation studies on the episodic memory, temporal encoding 414 schemes, level of supervision, and the number of training epochs to the appendix. 415

416 4.2 BASELINE METHODS

Off-the-shelf Baselines. We show published results from (Yang et al., 2023) obtained using GPT-3 (175 B parameters) as an off-the-shelf baseline, with few-shot and chain-of-thought prompting, for comparison with MemReasoner on original bAbi test set. We also report performances of a recurrent memory transformer-0.77B and of a Mamba-1.4B model, which we fine-tune on bAbi samples, on bAbi test set. For BABILong benchmarking, we include the following models from BABILong leaderboard: (1) Meta-Llama-3-8B-Instruct with an 8K context window size, (2) Phi3-mini-128k-instruct – a long-context LLM consisting of 3.8B parameters and a 128k token long context window, and (3) Llama3-ChatQA-1.5-8B with a nvidia/dragon-multiturn-query-encoder – a RAG framework.

Fine-tuned Baselines. We add RMT-137M and Mamba-130M performances from BABILong leaderboard, which has been finetuned on a set of samples that belong to the same distribution as BABILong (with PG19 padding) but is not included in BABILong benchmarking test set. These models were finetuned by using a curriculum schedule that progressively increases sequence lengths: 1, 2, 4, 6, 8, 16 and 32 segments (Kuratov et al., 2024).

431 We further benchmark RMT and Mamba models finetuned on bAbi on BABILong test samples. The goal is to figure out if those alternative recurrent models perform well on BABILong leaderboard due

432 to their true learning ability of the underlying task or due to their exposure to BABILong samples 433 during finetuning. We fine-tune off-the-shelf RMT(0.14b/0.77b) and Mamba (0.13b/1.4b) models 434 using the next token prediction loss on final answer reconstruction on bAbi Task 1 and 2 separately 435 till the testing accuracy on the task is sufficiently high (near 100%). In practice, we use 5 epochs 436 to reach above 99% accuracy on RMT and 20 epochs for the accuracy to plateau on Mamba, all using Adam optimizer with learning rate 1e - 5. RMT training was done with multiple segments 437 using a curriculum learning procedure. In order to train with more segments while exposing the 438 model to only bAbi data, we reduce the segment size to 64 for task 1 and 128 for task 2. This leads 439 to 2 segments in training for task 1 and 2-4 segments in training for task 2. In order to mimic the 440 curriculum learning process, we filter the data so that we train with inputs with token length up to 441 the segment size for 10 epochs, up to 2 times segment size for another 10 epochs, and so on. 442

We also add a Larimar-1.3B baseline, which is finetuned on bAbi and Wikipedia samples with first 443 and fifth terms from eqn. 1. The purpose of comparing MemReasoner with respect to Larimar is 444 to disambiguate the benefits of temporal feature learning and iterative query and read updates on 445 top of the episodic memory. Larimar fine-tuning shares the same training setups as MemReasoner. 446 We further add experiments with Qwen2.5-0.5B (https://huggingface.co/Qwen/Qwen2.5-0.5B) and 447 Qwen2.5-1.5B (https://huggingface.co/Qwen/Qwen2.5-1.5B) models (both of which support long 448 context windows up to 128k tokens), as well as a memory network (Sukhbaatar et al., 2015) that is 449 not coupled to transformer-based LLMs. See appendix for results. 450

It should be mentioned that all baselines used in this study are trained with supervision on final answer, whereas MemReasoner uses both supporting fact and final answer supervision. As mentioned
earlier, the goal is to check if this additional supervision, when tied to the operations around the
latent memory, enables better reasoning generalization. In that sense, MemReasoner offers a principled, model-agnostic approach for augmenting memory-based LLMs with robust reasoning, which
can be complimentary or used together with continual training. Throughout the paper, we report
task accuracy as the performance metric, so higher the better.

458 4.3 RESULTS

## 459 4.3.1 PERFORMANCE ON BABI TEST SET

460 Table 1 reports the performance of MemReasoner, which is independently finetuned on original 461 bAbi task 1 and task 2, along with the baselines on the corresponding bAbi test set of 1k sam-462 ples. Results show that, while prompting techniques such as few-shot learning and chain-of-thought 463 prompting (Yang et al., 2023) work well on task 1 which requires a single hop to find the entity location, those baselines perform much poorly on task 2 that requires learning temporal dependence 464 and performing multiple hops across facts to generate the final answer of object location. MemRea-465 soner, as well as RMT, Mamba and Larimar, all finetuned on bAbi achieves near-perfect accuracy 466 on both tasks. Importantly, Larimar baseline falls behind MemReasoner on both tasks, while the gap 467 being bigger on more complicated task 2, implying that read/write to episodic memory alone is not 468 sufficient. 469

470 4.3.2 PERFORMANCE ON BABILONG TEST SET

471 Table 2 and Table 3 report accuracy of MemReasoner, together with baseline methods, on BABI-Long task 1 and task 2 samples, respectively. '-' means unavailable due to out of memory errors 472 or maximal input length constraints. For task 1, the following observations can be made: (i) at 473 half of model's context window, the accuracy of Llama-3-8B-Instruct drops to 80% and Phi-3-mini-474 128k drops to 63% of the corresponding model's performance at 0k samples, indicating LLMs are 475 not good at utilizing their full context window. With RAG, the performance stays at a flat  $\approx 60\%$ 476 all throughout. Interestingly, while RMT and Mamba, when finetuned on BABILong samples of 477 up to 16k tokens, are the best models reported on BABILong leaderboard, they perform poorly on 478 BABILong samples beyond 0k as we finetune them on bAbi samples. This suggests exposure to 479 BABILong during training helps RMT and Mamba, as the models have seen facts embedded in-480 side the background distractor text from PG19. Larimar finetuned on bAbi, while performing much 481 poorly on bAbi test set and BABILong 0k set to begin with, the accuracy on longer BABILong 482 samples is higher than bAbi-tuned RMT and Mamba baselines. In contrast, MemReasoner trained 483 on bAbi with supervision on supporting fact(s) and final answer generalizes well on BABILong for task 1, providing an average accuracy of 84.6% and 68.5% on  $\leq$  8k and  $\geq$  16k BABILong samples, 484 respectively. These results suggest, that models can benefit on long-context reasoning from having 485 access to longer similar sequences or to reasoning processes during training.

| 486 | Model type              | Task 1 | Task 2 |
|-----|-------------------------|--------|--------|
| 487 | CoT - GPT-3             | 97.3   | 72.2   |
| 488 | Few-shot - GPT-3        | 98.4   | 60.8   |
| 489 | RMT77B (bAbi)           | 97.7   | 97.5   |
| 490 | Mamba-1.4B (bAbi)       | 100    | 95     |
| 491 | Larimar-1.3B (bAbi)     | 60.6   | 44.9   |
| 492 | MemReasoner-1.4B (bAbi) | 100    | 100    |

Table 1: Performance on bAbi tasks. Best model is highlighted in bold. GPT-3 (=text-davinci-003) baselines are from (Yang et al., 2023). Finetuning data, if any, seen by a model is specified within parentheses.

|                                      | Avg.      | Avg.       |     |     |     |     |     |     |     |     |      |
|--------------------------------------|-----------|------------|-----|-----|-----|-----|-----|-----|-----|-----|------|
| Model type                           | $\leq 8k$ | $\geq 16k$ | 0k  | 1k  | 2k  | 4k  | 8k  | 16k | 32k | 64k | 128k |
| RMT14B (BABILong)*                   | 100       | 97         | 100 | 100 | 100 | 100 | 100 | 100 | 99  | 96  | 94   |
| Mamba13B (BABILong)*                 | 100       | 100        | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100  |
| Few-shot - Meta-Llama-3-8B-Instruct* | 84.4      | -          | 98  | 93  | 90  | 79  | 62  | -   | -   | -   | -    |
| Few-shot - Phi-3-mini-128k-instruct* | 78.4      | 38         | 97  | 84  | 72  | 69  | 70  | 60  | 53  | 38  | 1    |
| RAG - Llama3-ChatQA-1.5-8B*          | 59.6      | 60         | 60  | 62  | 60  | 58  | 58  | 60  | 60  | 56  | 64   |
| RMT14B (bAbi)                        | 32.8      | 15.5       | 96  | 4   | 26  | 19  | 19  | 12  | 22  | 12  | 16   |
| RMT77B (bAbi)                        | 37.2      | 16.7       | 99  | 27  | 21  | 25  | 14  | 14  | 19  | 16  | 18   |
| Mamba13B (bAbi)                      | 20.4      | -          | 85  | 11  | 5   | 0   | 1   | 0   | 0   | 0   | -    |
| Mamba-1.4B (bAbi)                    | 44.2      | -          | 100 | 60  | 42  | 19  | 0   | 0   | 0   | 0   | -    |
| Larimar-1.3B (bAbi)                  | 44.8      | 14.3       | 63  | 59  | 55  | 28  | 19  | 14  | 16  | 13  | 14   |
| MemReasoner-1.4B (bAbi)              | 84.6      | 68.5       | 99  | 91  | 83  | 76  | 74  | 71  | 68  | 70  | 65   |

Table 2: BABILong Task 1 Results. Baseline results marked with "\*" are cited from (Kuratov et al.,2024). The finetuning data, if any, seen by each model is specified within parentheses.

|                                      | Avg.      | Avg.       |     |     |    |    |    |     |     |     |      |
|--------------------------------------|-----------|------------|-----|-----|----|----|----|-----|-----|-----|------|
| Model type                           | $\leq 8k$ | $\geq 16k$ | 0k  | 1k  | 2k | 4k | 8k | 16k | 32k | 64k | 128k |
| RMT14B (BABILong)*                   | 98.8      | 68.5       | 100 | 100 | 99 | 98 | 97 | 94  | 82  | 59  | 39   |
| Mamba13B (BABILong)*                 | 98.0      | 94.5       | 98  | 98  | 98 | 98 | 98 | 98  | 98  | 95  | 87   |
| Few-shot - Meta-Llama-3-8B-Instruct* | 40.2      | -          | 47  | 46  | 49 | 39 | 20 | -   | -   | -   | -    |
| Few-shot - Phi-3-mini-128k-instruct* | 40.6      | 15.5       | 57  | 38  | 38 | 36 | 34 | 23  | 22  | 15  | 2    |
| RAG - Llama3-ChatQA-1.5-8B*          | 21.6      | 8.75       | 28  | 25  | 22 | 19 | 14 | 13  | 9   | 7   | 6    |
| RMT14B (bAbi)                        | 36.6      | 12         | 97  | 31  | 19 | 16 | 20 | 12  | 12  | 14  | 10   |
| RMT77B (bAbi)                        | 41.2      | 17.5       | 100 | 36  | 21 | 27 | 22 | 18  | 23  | 13  | 16   |
| Mamba13B (bAbi)                      | 16.2      | -          | 64  | 10  | 3  | 3  | 1  | 0   | 0   | 0   | -    |
| Mamba-1.4B (bAbi)                    | 31.6      | -          | 94  | 44  | 15 | 5  | 0  | 0   | 0   | 0   | -    |
| Larimar-1.3B (bAbi)                  | 31        | 20.3       | 42  | 41  | 29 | 22 | 21 | 19  | 16  | 22  | 24   |
| MemReasoner-1.4B (bAbi)              | 60.6      | 18.5       | 100 | 73  | 61 | 46 | 23 | 20  | 19  | 17  | 20   |

Table 3: BABILong Task 2 Results. Baseline results marked with "\*" are cited from (Kuratov et al., 2024).

For more complicated task 2, which requires learning temporal dependence between the facts and finding and using two supporting facts in correct order for generation, both few-shot prompting and RAG with different base LLM show poor performance to begin with, and sharply degrade with context length increase of test samples. Again, RMT and Mamba, when fitted to BABILong, perform well on test samples, both struggle to generalize from bAbi to BABILong. For example, the accuracy drops from near 100% at 0k to 18% for RMT and to 36% for Mamba at 1k. Poor results at short context length for Larimar also indicates model's failure to learn the task. MemReasoner, in comparison, provides an accuracy of 100% at 0k, 73% at 1k, and 46% at 4k, while performance degrades to  $\approx 18.5\%$  beyond 16k. The modest ( $\approx 18.5\%$ ) performance of bAbi-tuned MemReasoner at 16k or longer context suggests that there remains significant room for MemReasoner to improve, which will be investigated in future. One possible direction is to train MemReasoner on longer sequences and/or with different levels of supervision.

#### 4.3.3 GENERALIZATION TO OUT-OF-DISTRIBUTION TEST SETS

To test if the models have indeed learned to solve the tasks, we create a new testbed where the construct of the tasks remains the same, but the answer changes from training to test set. Specifically, we change the location information present in the answer set of bAbi training  $\rightarrow$  test as follows:

| Model type             | Task 1 | Task 2 |
|------------------------|--------|--------|
| RMT77B (bAbi)          | 44.7   | 0.6    |
| Mamba-1.4B (bAbi)      | 67     | 44     |
| Larimar-1.3B (bAbi)    | 24.9   | 7.8    |
| MemReasoner-1.4B (bAbi | ) 87.2 | 52.7   |
|                        |        |        |

Table 4: Robustness to location changes in bAbi test set.

| Model type              | 0k  | 1k | 2k | 4k |
|-------------------------|-----|----|----|----|
| RMT77B (bAbi)           | 100 | 19 | 20 | 12 |
| Mamba-1.4B (bAbi)       | 81  | 8  | 0  | 0  |
| Larimar-1.3B (bAbi)     | 45  | 19 | 20 | 11 |
| MemReasoner-1.4B (bAbi) | 83  | 58 | 50 | 45 |

Table 5: Performance on bAbi task  $2 \rightarrow$  BABILong task 1 generalization.

556office  $\rightarrow$  library, garden  $\rightarrow$  garage, kitchen  $\rightarrow$  cafe, bathroom  $\rightarrow$  attic, bedroom  $\rightarrow$  basement, hallway557 $\rightarrow$  gym. This now becomes a more stringent test, to which we subject all alternative architectures558including MemReasoner. As shown in Table 4, RMT struggles in this setting across both tasks. On559task 1, MemReasoner shows  $\approx 20\%$  higher accuracy than Mamba, whereas on task 2 MemReasoner560wins by  $\approx 8\%$ .

Finally, we also check if the models trained on 2-hop bAbi task 2 can solve the simpler 1-hop task 1
but on the corresponding BABILong samples. Results are shown in Table 5, indicating that the best
performing model on 0k BABILong task 1 samples is RMT, while MemReasoner being a second.
However, both RMT and Mamba perform very poorly on longer (1-4k tokens) BABILong samples,
whereas MemReasoner's accuracy remains strong.

#### 567 4.4 CONCLUSION

554 555

566

568 In this work, we introduce a new memory-augmented LLM architecture that comes with two essen-569 tial abilities required to perform robust multi-step reasoning, *i.e.*, learning temporal relations and to hop meaningfully between facts within a context. Our formulation and implementation of the multi-570 step reasoning mechanisms around the episodic memory, textcolortogether with supervised training 571 using reasoning steps and final answer, is generic and in principle model-agnostic, and therefore 572 can be leveraged to enhance other memory-augmented LLMs, including the ones used in this study 573 as baselines. We examine MemReasoner on BABILong, a benchmark purposed to test models' 574 reasoning ability when relevant facts are distributed in background of very large textual corpora. 575 This deceptively lengthy nature of BABILong samples, along with the presence of distracting text 576 that is naturally occurring, makes the underlying reasoning task more challenging on which even 577 bigger LLMs that have seen samples with long context during training fails. We show here that, 578 MemReasoner trained on bAbi samples provides strong performance on BABILong, compared to 579 the off-the-shelf powerful LLM baselines and alternative recurrent architectures that are also finetuned on bAbi data, though only with final answer supervision. We further show that MemReasoner 580 generalizes better in the setting where answers in training set differs from those in the test within the same task. MemReasoner also shows good adaptation from two-hop to single-hop QA task, 582 whereas the test samples are much longer and mixed with natural irrelevant text. Additional ex-583 periments show generality of MemReasoner approach across decoder scale (GPT2-1 to GPTJ-6B) 584 and across different multi-hop tasks. Results indicate, supervision on both supporting facts and final 585 outcome, together with multi-hop search over the context in latent memory space, enables more 586 robust reasoning generalization of LLMs. Taken together, designing alternative architectures with new loss objectives that encourage the model to learn the underlying reasoning skills is a potential 588 path toward more robust reasoners.

# 590 REFERENCES

589

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

#### 11

| 594<br>595<br>596        | Aydar Bulatov, Yury Kuratov, and Mikhail Burtsev. Recurrent memory transformer. Advances in Neural Information Processing Systems, 35:11079–11091, 2022.   |
|--------------------------|--|
| 597<br>598<br>599<br>600 | Payel Das, Subhajit Chaudhury, Elliot Nelson, Igor Melnyk, Sarath Swaminathan, Sihui Dai, Aurélie<br>Lozano, Georgios Kollias, Vijil Chenthamarakshan, Jiří, Navrátil, Soham Dan, and Pin-Yu Chen.<br>Larimar: Large language models with episodic memory control, 2024. URL https://arxiv.<br>org/abs/2403.11901. |
| 601<br>602<br>603        | Yiran Ding, Li Lyna Zhang, Chengruidong Zhang, Yuanyuan Xu, Ning Shang, Jiahang Xu, Fan<br>Yang, and Mao Yang. Longrope: Extending llm context window beyond 2 million tokens, 2024.<br>URL https://arxiv.org/abs/2402.13753.  |
| 604<br>605<br>606<br>607 | Yao Fu, Rameswar Panda, Xinyao Niu, Xiang Yue, Hannaneh Hajishirzi, Yoon Kim, and Hao Peng.<br>Data engineering for scaling language models to 128k context, 2024. URL https://arxiv.<br>org/abs/2402.10171.   |
| 608                      | Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces, 2023.   |
| 609<br>610<br>611<br>612 | Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, Yang Zhang, and Boris Ginsburg. Ruler: What's the real context size of your long-context language models?, 2024. URL https://arxiv.org/abs/2404.06654.  |
| 613<br>614<br>615        | Yichen Jiang and Mohit Bansal. Avoiding reasoning shortcuts: Adversarial evaluation, training, and model development for multi-hop qa, 2019. URL https://arxiv.org/abs/1906.07132.   |
| 617<br>618<br>619<br>620 | Daniel Kahneman. <i>Thinking, fast and slow.</i> Farrar, Straus and Giroux, New York, 2011. ISBN 9780374275631 0374275637. URL https://www.amazon.de/<br>Thinking-Fast-Slow-Daniel-Kahneman/dp/0374275637/ref=wl_it_dp_o_pdT1_nS_nC?ie=UTF8&colid=151193SNGKJT9&coliid=I30CESLZCVDFL7.                             |
| 621                      | Pentti Kanerva. Sparse distributed memory. MIT press, 1988.  |
| 622<br>623<br>624<br>625 | Muhammad Khalifa, Lajanugen Logeswaran, Moontae Lee, Honglak Lee, and Lu Wang. Grace:<br>Discriminator-guided chain-of-thought reasoning, 2023. URL https://arxiv.org/abs/<br>2305.14934.  |
| 626<br>627               | Georgios Kollias, Payel Das, and Subhajit Chaudhury. Generation constraint scaling can mitigate hallucination, 2024. URL https://arxiv.org/abs/2407.16908.   |
| 628<br>629<br>630<br>631 | Yuri Kuratov, Aydar Bulatov, Petr Anokhin, Ivan Rodkin, Dmitry Sorokin, Artyom Sorokin, and Mikhail Burtsev. Babilong: Testing the limits of llms with long context reasoning-in-a-haystack. <i>arXiv preprint arXiv:2406.10149</i> , 2024.  |
| 632<br>633<br>634        | Mosh Levy, Alon Jacoby, and Yoav Goldberg. Same task, more tokens: the impact of input length on the reasoning performance of large language models, 2024. URL https://arxiv.org/abs/2402.14848.   |
| 636<br>637<br>638        | Hanmeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji Zhou, and Yue Zhang. Evaluating the log-<br>ical reasoning ability of chatgpt and gpt-4, 2023a. URL https://arxiv.org/abs/2304.<br>03439.  |
| 639<br>640               | Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts, 2023b.   |
| 041<br>642<br>643<br>644 | Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke<br>Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work?,<br>2022. URL https://arxiv.org/abs/2202.12837.   |
| 645<br>646               | Melanie Mitchell. How do we know how smart ai systems are?, 2023.  |

647 Tsendsuren Munkhdalai, Manaal Faruqui, and Siddharth Gopal. Leave no context behind: Efficient infinite context transformers with infini-attention. *arXiv preprint arXiv:2404.07143*, 2024.

| 648<br>649<br>650        | Marianna Nezhurina, Lucia Cipolina-Kun, Mehdi Cherti, and Jenia Jitsev. Alice in wonderland:<br>Simple tasks showing complete reasoning breakdown in state-of-the-art large language models,<br>2024. URL https://arxiv.org/abs/2406.02061.  |
|--------------------------|--|
| 651<br>652<br>653<br>654 | Richard Yuanzhe Pang, Weizhe Yuan, Kyunghyun Cho, He He, Sainbayar Sukhbaatar, and Jason Weston. Iterative reasoning preference optimization, 2024. URL https://arxiv.org/abs/2404.19733.  |
| 655<br>656<br>657        | Bhargavi Paranjape, Scott Lundberg, Sameer Singh, Hannaneh Hajishirzi, Luke Zettlemoyer, and Marco Tulio Ribeiro. Art: Automatic multi-step reasoning and tool-use for large language models, 2023. URL https://arxiv.org/abs/2303.09014.  |
| 658<br>659<br>660        | Kha Pham, Hung Le, Man Ngo, Truyen Tran, Bao Ho, and Svetha Venkatesh. Generative pseudo-<br>inverse memory. In <i>International Conference on Learning Representations</i> , 2022.  |
| 661<br>662<br>663        | Ofir Press, Noah Smith, and Mike Lewis. Train short, test long: Attention with linear biases enables input length extrapolation. In <i>International Conference on Learning Representations</i> , 2022. URL https://openreview.net/forum?id=R8sQPpGCv0.  |
| 664<br>665               | 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.  |
| 666<br>667<br>668<br>669 | Nir Ratner, Yoav Levine, Yonatan Belinkov, Ori Ram, Inbal Magar, Omri Abend, Ehud Karpas, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. Parallel context windows for large language models, 2023. URL https://arxiv.org/abs/2212.10947.  |
| 670                      | Sebastian Ruder. Challenges and opportunities in nlp benchmarking, 2021.   |
| 671<br>672<br>673        | Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer,<br>Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to<br>use tools, 2023. URL https://arxiv.org/abs/2302.04761.   |
| 675<br>676               | Kaya Stechly, Karthik Valmeekam, and Subbarao Kambhampati. Chain of thoughtlessness? an analysis of cot in planning, 2024. URL https://arxiv.org/abs/2405.04776.   |
| 677<br>678<br>679        | Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2023. URL https://arxiv.org/abs/2104.09864.  |
| 680<br>681<br>682        | Sainbayar Sukhbaatar, Arthur Szlam, Jason Weston, and Rob Fergus. End-to-end memory networks, 2015. URL https://arxiv.org/abs/1503.08895.  |
| 683<br>684<br>685<br>686 | Karthik Valmeekam, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. Large lan-<br>guage models still can't plan (a benchmark for LLMs on planning and reasoning about change).<br>In <i>NeurIPS 2022 Foundation Models for Decision Making Workshop</i> , 2022. URL https:<br>//openreview.net/forum?id=wUU-7XTL5XO. |
| 687<br>688<br>689        | 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.   |
| 691<br>692<br>693        | Chaojie Wang, Yanchen Deng, Zhiyi Lyu, Liang Zeng, Jujie He, Shuicheng Yan, and Bo An. Q*:<br>Improving multi-step reasoning for llms with deliberative planning, 2024. URL https://<br>arxiv.org/abs/2406.14283.  |
| 694<br>695<br>696        | Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837, 2022.   |
| 697<br>698<br>699        | Jason Weston, Sumit Chopra, and Antoine Bordes. Memory networks. arXiv preprint arXiv:1410.3916, 2014.   |
| 700<br>701               | Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M Rush, Bart Van Merriënboer, Armand Joulin, and Tomas Mikolov. Towards ai-complete question answering: A set of prerequisite toy tasks. <i>arXiv preprint arXiv:1502.05698</i> , 2015.  |

| 702<br>703<br>704               | Jian Wu, Linyi Yang, Zhen Wang, Manabu Okumura, and Yue Zhang. Cofca: A step-wise counter-<br>factual multi-hop qa benchmark, 2024. URL https://arxiv.org/abs/2402.11924.   |
|---------------------------------|---|
| 705<br>706                      | Yan Wu, Greg Wayne, Alex Graves, and Timothy Lillicrap. The kanerva machine: A generative distributed memory, 2018.   |
| 707<br>708<br>709<br>710<br>711 | Wenhan Xiong, Jingyu Liu, Igor Molybog, Hejia Zhang, Prajjwal Bhargava, Rui Hou, Louis Martin,<br>Rashi Rungta, Karthik Abinav Sankararaman, Barlas Oguz, Madian Khabsa, Han Fang, Yashar<br>Mehdad, Sharan Narang, Kshitiz Malik, Angela Fan, Shruti Bhosale, Sergey Edunov, Mike Lewis,<br>Sinong Wang, and Hao Ma. Effective long-context scaling of foundation models, 2023. URL<br>https://arxiv.org/abs/2309.16039. |
| 712<br>713<br>714<br>715        | Zhun Yang, Adam Ishay, and Joohyung Lee. Coupling large language models with logic program-<br>ming for robust and general reasoning from text. In <i>The 61st Annual Meeting Of The Association</i><br><i>For Computational Linguistics</i> , 2023.  |
| 716<br>717<br>718<br>719        | Tao Yuan, Xuefei Ning, Dong Zhou, Zhijie Yang, Shiyao Li, Minghui Zhuang, Zheyue Tan, Zhuyu Yao, Dahua Lin, Boxun Li, Guohao Dai, Shengen Yan, and Yu Wang. Lv-eval: A balanced long-context benchmark with 5 length levels up to 256k, 2024. URL https://arxiv.org/abs/2402.05136.   |
| 720<br>721<br>722               | Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. Star: Bootstrapping reasoning with reasoning, 2022. URL https://arxiv.org/abs/2203.14465.  |
| 723                             |   |
| 724                             |   |
| 725                             |   |
| 726                             |   |
| 727                             |   |
| 728                             |   |
| 729                             |   |
| 730                             |   |
| 731                             |   |
| 732                             |   |
| 733                             |   |
| 734                             |   |
| 735                             |   |
| 736                             |   |
| 737                             |   |
| 738                             |   |
| 739                             |   |
| 740                             |   |
| 741                             |   |
| 742                             |   |
| 743                             |   |
| 744                             |   |
| 745                             |   |
| 746                             |   |
| 747                             |   |
| 748                             |   |
| 749                             |   |
| 750                             |   |
| 751                             |   |
| 752                             |   |

# 756 A APPENDIX

## 758 A.1 WHY NOT FINETUNE ON LONGER SEQUENCES?

759 The goal of this work is to propose an LLM-based architecture that learns to reason over unseen 760 long-context in an efficient, robust, and generalizable manner. As such, the evaluation framework 761 corresponds to a set-up where the core reasoning facts are diluted in the presence of irrelevant natural text distractors distributed over the context. This setup allows one to test how consistently 762 language model can solve the same reasoning task across different input lengths. This is inspired by the recent research showing that current LLMs' reasoning performance degrade at much shorter 764 input lengths than their technical maximum (Levy et al., 2024). At the same time, finetuning on 765 longer sequences presents several practical challenges: (i) The longer sequences with proper (human 766 or machine) annotation should be available during training – which is typically expensive and is 767 difficult to scale in real-world. (ii) Expanding the context window usually incurs a quadratic increase 768 in computational and memory cost for transformer-based LLMs. For example, the training setup 769 used in (Fu et al., 2024) shows that extending the Llama-2 7B model's context window from 4k 770 to 80k requires 8 A100 GPUs (80G each) for five days. The costs of resources and time further 771 increase significantly for larger models, for longer context length, and for more extended training 772 period. (iii) The test distribution is still expected to match to longer sequences seen during training 773 (e.g. mixture of bAbi with text from PG-19 in the case of BABILong) – which may not be always possible. As a result, straightforward continual pre-training or fine-tuning on longer sequences may 774 still not fully solve the fundamental problem of learning to reason over (long) context in a robust 775 and generalizable manner, and such approaches can benefit from using additional supervision (when 776 available) and training-time search in the latent memory space. 777

While MemReasoner is not trained on longer samples that are similar to the test distribution (whereas
RMT-BABILong and MAMBA-BABILong models in Tables 2 and 3 are), we train the model with
additional supervision on supporting facts. In that sense, our MemReasoner approach with reasoning
process and outcome supervision is complementary to the continual pre-training only with outcome
supervision.

783

# A.2 ADDITIONAL DATASET PREPROCESSING DETAILS

In the unprocessed bAbi data, a single data instance consists of a sequence of lines representing facts to reason over with questions interspersed throughout the facts. We preprocess the bAbi data such that after pre-processing, a single training sample consists of a single question with facts for reasoning being the lines before it, with previous questions replaced by an empty line. On average, this leads to about 2 empty lines per training sample. For batches containing training samples with different lengths of context episodes, we pad shorter episodes with rows of the encoder padding token at the beginning.

791 792

### A.3 COMPARISON OF INFERENCE-TIME COMPLEXITY

793 Let  $H_1$ ,  $H_2$  and  $d_1$ ,  $d_2$  be the number of transformer layers and hidden state dimension in the 794 encoder and decoder, respectively. Let E denote the number of context lines in a sample, L be 795 the max context length,  $L_1$  be the max query length, D be the latent space dimension, and m be 796 the memory size. The inference-time computational complexity for MemReasonr can be estimated by the encoder complexity  $\mathcal{O}(H_1((EL^2 + L_1^2)d_1 + (EL + L_1)d_1^2)))$ , temporal encoding complexity 797  $\mathcal{O}(Ed^2)$ , memory operation complexity  $\mathcal{O}(Edm^2)$ , decoding complexity  $\mathcal{O}(H_2(|P_a|^2d_2 + |P_a|d_2^2))$ , 798 and broadcasting complexity  $\mathcal{O}(d_1 dE)$  and  $\mathcal{O}(d_2 dH_2)$ . For a typical GPT decoding, the inference-799 time computational complexity is  $\mathcal{O}(H_2((EL+L_1)^2d_2+(EL+L_1)d_2^2)))$ . 800

To provide a more direct comparison, we give in Table 6 the inference cost measured in seconds per input for evaluating with BABILong in comparison to the base decoder (gpt2-large). We note that gpt2-large does not support context lengths longer than 1024 tokens. Overall, we observe that the increase in inference time for MemReasoner is very small for 0k and MemReasoner is more efficient for 1k context length. This is because of utilizing the latent encodings of context, performing oneshot write to the memory, and executing multiple hops over that memory in latent space.

806 807

## A.4 COMPARISON TO DECODER-ONLY LMS THAT SUPPORT LONG CONTEXTS

809 We experiment with Qwen2.5-0.5B and Qwen2.5-1.5B models both of which are decoder-only LMs that support long context windows (up to 128k tokens). The performance of Qwen models on bAbi

1k

Model type 0k

| 811 | gpt2-large              | 0.28      | 1.13      | -          | -      |      | -      | -      |      | -       | -      |        | -         |
|-----|-------------------------|-----------|-----------|------------|--------|------|--------|--------|------|---------|--------|--------|-----------|
| 812 | MemReasoner             | 0.30      | 0.33      | 0.40       | 0.61   | . (  | ).98   | 1.9    | 94   | 3.26    | 11.    | 25     | 13.77     |
| 813 |                         |           |           |            |        |      |        |        |      |         |        |        |           |
| 814 | Table 6: Th             | e infere  | ence co   | st meas    | ured i | n se | cond   | s per  | inp  | ut on l | BABI   | Long   |           |
| 815 |                         |           |           |            |        |      |        |        |      |         |        |        |           |
| 816 |                         |           | Moo       | lel type   |        |      | Tas    | k 1    | Tas  | sk 2    |        |        |           |
| 817 |                         | Q         | wen2.5    | -0.5B (t   | oAbi)  |      | 10     | 0      | 9    | 6       |        |        |           |
| 818 |                         | Q         | wen2.5    | -1.5B (t   | oAbi)  |      | 99     | .9     | 98   | 3.9     |        |        |           |
| 010 |                         | Mem       | Reasor    | her-1.4E   | 8 (bAt | oi)  | 10     | 0      | 1    | 00      |        |        |           |
| 015 |                         |           |           |            |        |      |        |        |      |         |        |        |           |
| 020 | Table 7: Performance of | n bAbi    | tasks. I  | Best mo    | del is | higl | hligh  | ted in | n bo | ld. GF  | PT-3 ( | =text  | -davine   |
| 821 | baselines are from (Yar | ng et al. | , 2023)   | . Finetu   | ining  | data | , if a | ny, s  | een  | by a n  | nodel  | is spo | ecified v |
| 822 | parentheses.            |           |           |            |        |      |        |        |      |         |        |        |           |
| 823 | -                       |           |           |            |        |      |        |        |      |         |        |        |           |
| 824 |                         |           | Avg.      | Avg.       |        |      |        |        |      |         |        |        |           |
| 825 | Model type              |           | $\leq 8k$ | $\geq 16k$ | Ok     | 1k   | 2k     | 4k     | 8k   | 16k     | 32k    | 64k    | 128k      |
| 025 | Qwen2.5-0.5B (1         | oAbi)     | 45.4      | -          | 94     | 66   | 34     | 23     | 10   | 3       | 1      | -      | -         |
| 826 | Qwen2.5-1.5B (          | oAbi)     | 61.6      | -          | 100    | 81   | 57     | 42     | 28   | 32      | 18     | -      | -         |
| 827 | MemReasoner-1.4         | B (bAbi)  | 84.6      | 68.5       | 99     | 91   | 83     | 76     | 74   | 71      | 68     | 70     | 65        |
| 828 |                         |           |           |            |        |      |        |        |      |         |        |        |           |

2k

4k

8k

16k

32k

64k

128k

Table 8: BABILong Task 1 Results - Qwen family models.

| Model type              | $Avg. \le 8k$ | Avg. $\geq 16k$ | 0k  | 1k | 2k | 4k | 8k | 16k | 32k | 64k | 128k |
|-------------------------|---------------|-----------------|-----|----|----|----|----|-----|-----|-----|------|
| Qwen2.5-0.5B (bAbi)     | 57.8          | -               | 96  | 76 | 59 | 39 | 19 | 11  | 3   | -   | -    |
| Qwen2.5-1.5B (bAbi)     | 46.6          | -               | 99  | 67 | 32 | 25 | 10 | 6   | 2   | -   | -    |
| MemReasoner-1.4B (bAbi) | 60.6          | 18.5            | 100 | 73 | 61 | 46 | 23 | 20  | 19  | 17  | 20   |

834 835 836

837

810

Table 9: BABILong Task 2 Results - Qwen family models.

838Task 1 and Task 2 is similar to the best in MemReasoner (Table 7).Overall, we find that MemRea-839840840840840841841841841842842842843843

844 A.5 EXTENSION TO GPTJ-6B

845 MemReasoner is a model-agnostic way to augment current decoder-only LLMs with dynamically 846 updatable memory. Via end-to-end training, the architecture learns to write the latent encodings in 847 a fixed-size memory, order them in their order of appearance in the context, and perform multiple hop over that context and update the latent query accordingly. The decoder learns a differentiated 848 attention mechanism to the readout from the memory, to accurately generate the final answer and 849 supporting facts (intermediate hops). Below, we provide the results when we train a GPTJ-6B de-850 coder with MemReasoner training protocol, suggesting more or less similar performance compared 851 to MemReasoner-1.3B. 852

# 853 A.6 BEYOND BABI DATASET

854 In this section, we explore the generalization of MemReasoner on another dataset, variable tracking 855 (VT) from RULER (Hsieh et al., 2024). In the VT task, the model is given context with lines with 856 information about variable value assignment such as "VAR AAAAA = 16438" or "VAR BBBBB = 857 AAAAA'' and the model is prompted to obtain all variables with a specific value. Variable names 858 have the format of 5 repeating letters randomly sampled from the alphabet. We train and evaluate 859 with chains of length 2, 4, 6, 8, and 10 and return the average accuracy over all chain lengths for 860 the 1 hop and 2 hop VT tasks. In order to pad the context for lengths 1k, 4k, and 16k, we follow the 861 approach taken from RULER of padding with the sentence "The grass is green. The sky is blue. The sun is yellow. Here we go. There and back again.\n" until the context reaches the desired length. 862 This noise is not present during training and the 0k data follow the same distribution as the training 863 data.

| Model type              | Task 1 | Task 2 |
|-------------------------|--------|--------|
| Qwen2.5-0.5B (bAbi)     | 44.2   | 14.5   |
| Qwen2.5-1.5B (bAbi)     | 75.2   | 63.5   |
| MemReasoner-1.4B (bAbi) | 87.2   | 52.7   |

Table 10: Robustness to location changes in bAbi test set.

| Model type              | 0k  | 1k | 2k | 4k |
|-------------------------|-----|----|----|----|
| Qwen2.5-0.5B (bAbi)     | 97  | 71 | 47 | 32 |
| Qwen2.5-1.5B (bAbi)     | 100 | 58 | 36 | 18 |
| MemReasoner-1.4B (bAbi) | 83  | 58 | 50 | 45 |

Table 11: Performance on bAbi task  $2 \rightarrow$  BABILong task 1 generalization.

Since VT asks for all variables with a specific value, for MemReasoner, we take all unordered
readouts of the model and pass them individually to the decoder to get the variables from each
reasoning hop, and then concatenate these variables in order to obtain the final answer. For RMT,
we train with 2 segments, with segment size set to the median length on the train dataset. From
Table 13 and Table 14, we observed that it is difficult to train RMT with 2 segments for the 2-hop
VT task, RMT can easily learn a shortcut and have high accuracy on the training set, but does not
generalize well to the test set at 0k length and performance degrades further at longer context length.
Larimar also learned short cuts on 2-hop VT tasks and could not perform well on test sets.

#### 886 A.7 COMPARISON WITH TRADITIONAL NON-LLM MEMORY NETWORKS

We have performed additional experiments to evaluate the performance of MemN2N (Sukhbaatar et al., 2015), which we trained on bAbi task 1 and task 2 data with final answer supervision and achieved 100% test accuracy on both. The results are summarized in the following tables and demonstrate the lack of generalization ability of MemN2N compared to MemReasoner.

- 891 892 A.8 ABLATION STUDIES
- 893 A.8.1 MEMORY

In Table 16, we conduct the ablation study on the episodic memory module in MemReasoner on bAbi and BABILong, task 1 and 2. Specifically, MemReasoner w/o memory module uses the same architecture of encoder and decoder (BERT-Large and GPT2-Large respectively) but does not use the memory module for encoding the context. Instead, the MemReasoner w/o memory uses the encoder to encode only the question and this is passed in to the decoder as kv-cache. Additionally, the context and question are passed to the decoder as part of the prompt with the format:

- 900 Context:
- 901 {context}
- 902 Question:
- 903 {question}
- 904 Answer:

where {context} and {question} represent the context and the question for the datapoint. We train
the model with reconstruction loss to ensure that the model is able to fill in the answer given this
prompt and with autoencoding loss on the pretraining dataset (see last term of Equation 1) in order
to reduce overfitting on bAbi data. We train MemReasoner w/o memory module for 5 epochs.

MemReasoner w/o memory module trained on bAbi task 1 obtains almost perfect accuracy on bAbi task 1 and BABILong task 1 0k. However, its generalization ability to long context (BABILong 1k and 2k) is much inferior to MemReasoner (MemReasoner\memory 0% vs. MemReasoner 91% on BABILong 1k). Similar trends can also be seen from bAbi task 2 trained MemReasoner\memory, implying the significance of the episodic memory module and the operations around it in MemReasoner.

#### 915 916 A.8.2 TEMPORAL ENCODING

917 In Table 17, we experiment with different temporal encoding schemes, including non-parametric method (Positional Encoding) and parametric method (GRU). In the table, we show MemRea-

<sup>874</sup> 875 876 877

| 918 |   | Model type                       | 0k      | 1k                | 2k                    | 4k     | 8k               | 16k      | 32k              | 64k       | 128k      |                  |  |  |
|-----|---|----------------------------------|---------|-------------------|-----------------------|--------|------------------|----------|------------------|-----------|-----------|------------------|--|--|
| 919 |   | Task 1                           | 98      | 82                | 77                    | 65     | 60               | 68       | 70               | 65        | 67        | -                |  |  |
| 920 |   | Task 2                           | 98      | 65                | 50                    | 34     | 35               | 32       | 22               | 27        | 30        |                  |  |  |
| 921 |   |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 922 | Table 12: Performance of MemReasoner with a GPTJ-6B decoder on BABILong.  |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 923 |   |                                  |         |                   |                       | 1      | 01               | 11       | 41               | 1.0       |           |                  |  |  |
| 924 |   |                                  | Vlode   | I type            |                       |        | 0K               | 1K       | 4K               | 16        | <u> </u>  |                  |  |  |
| 925 | KMI//B(VI)  |                                  |         |                   |                       |        |                  | 5.7      | 5.0              | 4.2       | )<br>6    |                  |  |  |
| 926 | Larimar-1.5B (VI)<br>MemPersoner 1 $AP_{i}(VT)$   |                                  |         |                   |                       |        |                  | 92.5     | 94.0             | 95.<br>00 | 0<br>0    |                  |  |  |
| 927 |   | wichikk                          | 20011   | -1-1.7            | <b>D</b> ( <b>v</b> ) |        | <i>)).)</i>      | 100.0    | ,,,,             | "         | /         |                  |  |  |
| 928 |   | Та                               | ble 13  | 3: Sin            | gle h                 | op va  | ariable          | trackin  | g resul          | ts.       |           |                  |  |  |
| 929 |   |                                  |         |                   | C                     | 1      |                  |          | C                |           |           |                  |  |  |
| 930 |   |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 931 | soner's accurac   | cy on BABILo                     | ng Ta   | sk 1.             | It can                | be s   | seen that        | at GRU   | encod            | ing ha    | s signifi | icant advantage  |  |  |
| 932 | over Positional Encoding, with much slower decay in the accuracy as the context length increases.   |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 933 | Additionally, though showing higher accuracy compared with Positional Encoding, uni-directional   |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 934 | GRU's accuracy decreases faster than bi-directional GRUs. Since 1-layer bi-directional GRU has  |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 935 | similar performance with 2-layer bi-directional GRU, we choose the lighter model and use 1-layer bi-directional CPU throughout the experiments in this paper.   |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 936 | bi-directional GRU throughout the experiments in this paper.  |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 937 | A.8.3 QUER  | by Update $\alpha$               |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 938 | In Table 18, v  | ve exploit test                  | -time   | infer             | ence                  | hype   | er-parai         | meter a  | $\alpha$ and i   | ts effe   | ect in re | easoning tasks'  |  |  |
| 939 | performance.  | We draw inspir                   | ation   | from              | (Kol                  | lias e | et al., 2        | 2024), v | where a          | uthors    | investi   | gated the effect |  |  |
| 940 | of scaling readout vectors to improve generation quality. In Line 20 of Algorithm 1, when using an  |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 941 | $\alpha > 1$ , we equ   | ivalently scale                  | up th   | e read            | lout v                | rector | rs whic          | ch great | ly help          | our g     | eneraliz  | zation to Task 1 |  |  |
| 942 | BABILong acc  | cording to Tabl                  | le 18 ( | (e.g. f           | from 1                | 4%     | to 45%           | on 4k    | contex           | t toker   | ı task).  |                  |  |  |
| 943 | A 9 EFFECT  | OF ARBITRAI                      | Y NI    | IMRE              | ROF                   | норя   | S WITH           | IWFAK    | FR SU            | PFRVI     | SION      |                  |  |  |
| 944 | Table 10 share  |                                  | -ENA    | D                 | K OI                  |        |                  | :        | 41               | 1         |           |                  |  |  |
| 945 | Table 19 shows performance of MemReasoner that is trained with weaker supervision on bAbi task<br>2 and is tasted on BABII and task 2 tast set. In this case, during training on arbitrary number (5) |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 946 | of hops was us  | sed together w                   | ig tasi | nerv <sup>1</sup> | st set.               | only   | on the           | final s  | ing trai         | no fac    | t and th  | ne final answer  |  |  |
| 947 | While perform   | ance on longe                    | r sam   | ples of           | drops                 | com    | pared            | to the r | nodel t          | rained    | with fi   | ull supervision. |  |  |
| 948 | the model generalizes well on 1k tokens long BARII ong samples compared to other baselines (see   |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 949 | Table 3. This c   | lirection will b                 | e furt  | her ex            | xplore                | d in   | future           | work.    |                  | 1         |           |                  |  |  |
| 950 |   |                                  |         |                   | •                     |        |                  |          |                  |           |           |                  |  |  |
| 951 | A.9.1 IRAI  | NING EPOCHS                      | P       |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 952 | In Table 20, w  | e evaluate Me                    | mRea    | isone             | r's pe                | rtorn  | nance v          | when fi  | ne-tune          | ed on     | bAbi ta   | sk 2 as a func-  |  |  |
| 953 | tion of the num   | nber of training                 | g epoc  | ens. S            | pecin                 |        | , with           | iewer e  | pocns,           | Memi      | Keasone   | er demonstrates  |  |  |
| 954 | stronger robus  | iness to location and 50% as the | on cha  | inge,             | ntinu                 | ing a  | $\frac{11}{100}$ | DOth o   | 1 /9% ;<br>noch) | at the    | ooth ep   | side MemPee      |  |  |
| 955 | soper's accura  | cy on shorter                    | conte   | nig Co<br>vt tas  | ks in                 | RAR    | RII ong          | Task 1   | $\mid$ and 2     | (ie       | (0.4k) i  | moroves as the   |  |  |
| 956 | training contin   | ues.                             | conte   | at tub            | K5 III                | DITL   | JILONG           | Tusk     | und 2            | (1.0.     | 0 11() 1  | inproves us the  |  |  |
| 957 | a anna g comm   |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 958 | A.10 LIMIT.   | ATIONS AND H                     | TUTU    | re W              | ORK                   |        |                  |          |                  |           |           |                  |  |  |
| 959 | The current w   | ork is limited                   | to tes  | ting t            | the M                 | emR    | leasone          | er fram  | ework            | on syr    | thetic 1  | reasoning tasks  |  |  |
| 960 | only. Future  | work will exte                   | end th  | ie fra            | mewc                  | ork to | o evalu          | ating 1  | easoni           | ng ger    | neraliza  | tion on natural  |  |  |
| 961 | language datas  | sets. Another p                  | otenti  | al dir            | rection               | ı is e | xtendi           | ng Mer   | nReaso           | ner to    | scenari   | os with weaker   |  |  |
| 962 | and noisy supe  | ervision.                        |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 963 |   |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 964 |   |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |
| 965 |   |                                  |         |                   |                       |        |                  |          |                  |           |           |                  |  |  |

| 972  |               | Μ                  | [oda]             | tuno                     |                 | I.              | 012              | 11-                           | /         | Ŀ                | 16b             |       |            |                |            |
|------|---------------|--------------------|-------------------|--------------------------|-----------------|-----------------|------------------|-------------------------------|-----------|------------------|-----------------|-------|------------|----------------|------------|
| 973  |               |                    |                   |                          | <b>F</b> )      |                 | 74.6             | 57                            | 1         | 2                | 0.2             | -     |            |                |            |
| 974  |               | 1//1<br>19r_1 (    | 3 R (1            | 1)<br>/T)                |                 | 0.1             | 0                | 0                             | .5        | 0.2              |                 |       |            |                |            |
| 975  |               | sonei              | 3D (∖<br>-14F     | 7 1 )<br>R (V1           | <u>г</u> )      | 0.1<br>98 4     | 97.6             | 5 97                          | .1<br>7 0 | 98.0             |                 |       |            |                |            |
| 976  |               | , ,                | .0                | 70.0                     |                 |                 |                  |                               |           |                  |                 |       |            |                |            |
| 977  |               | Tab                | le 14             | : Two                    | o hor           | o vari          | able t           | rackir                        | ig res    | ults.            |                 |       |            |                |            |
| 978  |               |                    |                   |                          | 1               |                 |                  |                               | 0         |                  |                 |       |            |                |            |
| 979  |               |                    |                   |                          |                 |                 |                  |                               |           |                  |                 |       |            |                |            |
| 980  |               |                    | N                 | /lodel                   | type            | 2               |                  | 0                             | k         | lk i             | $\frac{2k}{15}$ |       |            |                |            |
| 981  |               |                    | BAB               | SILON                    | ig Tas          | sk I            |                  | 1(                            | )0 .      | 36<br>7 4        | 15              |       |            |                |            |
| 982  |               | h / h; t;          | BAB               | SILON                    | g la            | sk 2            | to alt 1         |                               | )U :      | 04 I             | 21<br>10        |       |            |                |            |
| 983  |               | DADI ta            | ask Z             | $\rightarrow \mathbf{D}$ | ADII            | Long            | task             | 1   3                         | 5         | 20               | 19              |       |            |                |            |
| 984  |               |                    | Table             | 15:                      | Perfo           | orma            | nce of           | Mem                           | N2N       |                  |                 |       |            |                |            |
| 985  |               |                    | ruon              | . 10.                    | 1 0110          | JIIIa           |                  |                               |           | •                |                 |       |            |                |            |
| 986  |               | Nr. 1.1.           | 1                 | <b>T</b> 1               | 1               | 01              | 11               | 01                            |           | 1.0              | 01              | 11    | ~          | 1              |            |
| 987  |               | Model type         |                   | 10                       | <u> </u>        | UK              | 1K               | 2K                            |           | $\frac{5K 2}{2}$ | UK              | 1K    |            | K              |            |
| 988  | Me            | mReasoner\mem      | ory               | 10                       | U               | 100             | 0                | - 02                          | 99        | 9.3<br>00        | 100             | 29    |            | -              |            |
| 989  |               | Memkeasoner        |                   | 10                       | U               | 99              | 91               | 03                            | 1 1       | UU               | 100             | 13    |            | )1             |            |
| 990  |               | Table              | 16 <sup>.</sup> A | blati                    | on st           | udv a           | on the           | eniso                         | dic n     | nemor            | v               |       |            |                |            |
| 991  |               | Tuble              | 10.1              | loiuti                   | 011 50          | aay             |                  | epibe                         | are n     | lennor           | 5               |       |            |                |            |
| 992  |               |                    |                   |                          |                 |                 |                  |                               |           |                  |                 |       |            |                |            |
| 993  |               |                    | Enco              | oding                    | ; sche          | eme             |                  | 0k                            | 1k        | 2k               |                 |       |            |                |            |
| 994  |               |                    | Positi            | onal                     | Enco            | oding           | DII              | 100                           | 27        | 20               |                 |       |            |                |            |
| 995  |               | 2-la               | yer b             | 1-dire                   | ectior          | hal G           | RU               | 100                           | 90        | 80               |                 |       |            |                |            |
| 996  |               | 2-la               | yer ur            | 11-d1r                   | ectio           | nal C           |                  | 94                            | /3        | 61<br>93         |                 |       |            |                |            |
| 997  |               | 1-1a               | yer b             | 1-dire                   | ectior          | ial G           | KU               | 99                            | 91        | 83               |                 |       |            |                |            |
| 998  |               | Table 17. A        | blatic            | on stu                   | idv o           | n the           | temp             | oral e                        | ncodi     | ng sc            | heme            | s     |            |                |            |
| 999  |               | 14010 17711        |                   | 511 500                  | uj e            |                 | to mp            | 01010                         |           |                  |                 |       |            |                |            |
| 1000 | <u> </u>      |                    |                   |                          |                 |                 |                  | БИ                            |           |                  |                 |       |            |                |            |
| 1001 | Query upda    | location change    | Ok                | 1k                       | 2k              | 1as<br>4k       | K 2 BA<br>8k     | BILON<br>16k                  | g<br>32k  | 64k              | 128k            | 1 as  | SKIE<br>1k | 5AB11<br>2k    | Long<br>4k |
| 1002 | 1             | 52.6               | 100               | 46                       | 25              | 18              | 18               | 13                            | 16        | 12               | 13              | 78    | 21         | 17             | 14         |
| 1003 | 4             | 54.2               | 100               | 73                       | 61              | 46              | 26               | 22                            | 19        | 19               | 27              | 83    | 47         | 44             | 40         |
| 1004 | 8             | 52.7               | 100               | 13                       | 01              | 40              | 23               | 20                            | 19        | 17               | 20              | 03    | 29         | 50             | 45         |
| 1005 |               | Table 18: A        | Ablati            | ion st                   | udy (           | on th           | e que            | ry upc                        | late p    | aram             | eter $\alpha$   |       |            |                |            |
| 1006 |               |                    |                   |                          | 5               |                 | 1                | 5 1                           | 1         |                  |                 |       |            |                |            |
| 1007 |               |                    | Ν.                | 1.1.4                    |                 | 1               | 01               | 11                            | 21        | 41               |                 |       |            |                |            |
| 1008 |               | <b>D</b>           |                   | $\frac{101}{14R}$        | pe<br>hAhi      | )               | 07               | 21<br>21                      | 2K        | 4K               | _               |       |            |                |            |
| 1009 |               | RI<br>RI           | MT_ 7             | 17B (                    | hAhi            | 8               | 97<br>100        | 36                            | 21        | 27               |                 |       |            |                |            |
| 1010 |               | Ma                 | mba-              | .13R                     | (hAł            | bi)             | 64               | 10                            | 3         | 3                |                 |       |            |                |            |
| 1011 |               | Ma                 | mba-              | 1.4B                     | (bAł            | bi)             | 94               | 44                            | 15        | 5                |                 |       |            |                |            |
| 1012 |               | Lar                | imar-             | 1.3B                     | (bAl            | bi)             | 42               | 41                            | 29        | 22               |                 |       |            |                |            |
| 1013 |               | Me                 | mRea              | asone                    | r (fu           | 11)             | 100              | 73                            | 61        | 46               | _               |       |            |                |            |
| 1014 |               | Mer                | nRea              | soner                    | (we             | ak)             | 100              | 58                            | 31        | 22               |                 |       |            |                |            |
| 1015 |               |                    |                   |                          |                 | 1               |                  |                               |           |                  | _               |       |            |                |            |
| 1016 | Table 19: Con | mparison of Mem    | Rease             | oner t                   | raine           | ed wi           | th full          | l supe                        | rvisic    | on wit           | h Me            | mRea  | ason       | er(w           | eak) on    |
| 1017 | BABILong ta   | isk 2 samples, wh  | here t            | he we                    | eak s           | uper            | vision           | cons                          | 1ders     | an ar            | bitrar          | y fiv | e ho       | ps ai          | nd only    |
| 1018 | supervision o | n iinai supporting | ract              | and fi                   | inal a          | inswe           | er.              |                               |           |                  |                 |       |            |                |            |
| 1019 |               |                    |                   |                          |                 |                 |                  |                               |           |                  |                 |       |            |                |            |
| 1020 |               | Task 2 bAbi        |                   |                          | 1               | Task 2          | BABI             | Long                          |           |                  |                 | Task  | 1 BA       | ABIL           | ong        |
| 1021 | #epochs       | location change    | 0k<br>aa          | 1k 1                     | $\frac{2k}{54}$ | $\frac{4k}{30}$ | 8k 10            | $\frac{5k}{3}$ $\frac{32}{1}$ | 2k 6      | 4k 1<br>18       | 128k            | 0k    | 1k<br>51   | 2k             | 4k<br>37   |
| 1022 | 100           | 47.3 1             | 100               | 70                       | 57 3            | 38 2            | $\frac{2}{28}$ 3 | 1 2                           | .5        | 12               | 19              | 82    | <b>58</b>  | <del>5</del> 0 | <b>46</b>  |
| 1023 | 200           | 52.7               | 100               | 73                       | 61 4            | 46 2            | 23 2             | 0 1                           | 9         | 17               | 20              | 83    | 58         | 50             | 45         |
| 1024 |               | T-11- 00           | A h 1 - 4         |                          |                 | on 41.          |                  | hare                          | f         |                  | -1-             |       |            |                |            |
| 1020 |               | iable 20: 7        | hdiat             | ion st                   | .udv (          | on th           | e num            | iper o                        | i tran    | iing e           | DOCU            | Ś     |            |                |            |

