

Retrieval-Augmented World Models Enhanced with Reinforcement Learning

Anonymous ACL submission

Abstract

World models achieve remarkable success in predicting future states and planning in complex environments and Large Language Models (LLMs) serve as promising foundation to build general world models. However, their performances are usually constrained by the limited knowledge to the specific environments. Existing research attempts to enhance LLM-based world models through prompting or fine-tuning approaches, which are either requiring human knowledge or computational extensive. Therefore, we introduce **Retrieval-Augmented World Models (RAWM)**, a novel framework that leverages retrieval-augmented generation to improve the LLM-based world models. Our main contributions are threefold: (i) We introduce a memory system and design an embedding model to retrieve and incorporate relevant experiences, significantly improving the world model’s predictive accuracy. (ii) We develop a reinforcement learning (RL) training pipeline that fine-tunes a small MLP head on the pre-trained embedding model using Proximal Policy Optimization (PPO), further enhancing prediction performance. (iii) We conduct extensive experiments across three diverse RL environments, i.e., Game24, BlocksWorld, and BabyAI, demonstrating that RAWM consistently outperforms baseline models and exhibits strong generalizability. By leveraging external memory and retrieval techniques and training embedding with RL pipeline, RAWM dynamically utilizes relevant historical experiences and equips LLMs with environment-specific knowledge without retraining, enabling more accurate and generalizable predictions.

1 Introduction

The world model (Ha and Schmidhuber, 2018) has demonstrated to be an important module in decision making due to the celebrating success of MuZero (Schrittwieser et al., 2020) and Dreamer (Hafner et al., 2019, 2021, 2023). As

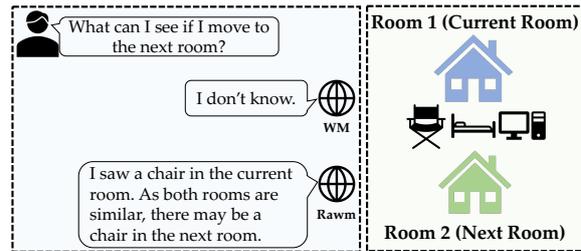


Figure 1: Why retrieval is needed?

learned accurate simulators, world models encode rich representations of the complex dynamics of the environment to predict the future states and the rewards. World models are critical for several key capabilities, such as generalization to novel tasks (Byravan et al., 2020; Robey et al., 2021; Young et al., 2023), efficient planning (Sekar et al., 2020; Hamrick et al., 2021; Schrittwieser et al., 2020), and offline learning (Schrittwieser et al., 2021; Yu et al., 2020, 2021). Beyond decision making, recent works such as Genie (Bruce et al., 2024) and Vista (Gao et al., 2024) demonstrate that world models can serve as general-purpose world simulators and users can directly interact with these models for playing and planning.

The past five years witness the remarkable success of large language models (LLMs) in enormous text generation and understanding tasks (Brown et al., 2020; OpenAI, 2023). LLMs serve as the world model explicitly in Reasoning via Planning (RAP) (Hao et al., 2023) and Reason for Future, Act for Now (RAFA) (Liu et al., 2023), where the LLMs predict the next states based on the actions executed at current states, e.g., the states of blocks in the BlocksWorld (Valmeekam et al., 2023), which is used to assist the planning methods. LLMs serve as the world model implicitly in the widely-used Tree of Thoughts (ToT) (Yao et al., 2023), as well as Graph of Thoughts (GoT) (Besta et al., 2024), where the LLMs need to predict the states and evaluate the thoughts to help the selection of the thoughts to advance the reasoning. The

075 generalizability of LLMs presents them as promis- 125
076 ing foundations for the world models. 126

077 However, the pre-trained LLMs may lack the 127
078 knowledge of the specific environments, which pro- 128
079 hibits them to be accurate world models. For the 129
080 example displayed in Figure 1, the LLM cannot 130
081 provide the accurate predictions whether there is a 131
082 chair in the next room if the next room has never 132
083 been visited. To address this issue, we can carefully 133
084 design the prompts to add the specific knowledge 134
085 to help the LLMs in making predictions, e.g., the 135
086 rules for objects and actions (Wang et al., 2024; 136
087 Gu et al., 2024). However, these knowledge is 137
088 even usually not available for users. Alternatively, 138
089 we can fine-tune the LLMs on the specific envi- 139
090 ronments (Xiang et al., 2023; Chae et al., 2025). 140
091 However, the training of LLMs brings additional 141
092 complexities for building the world models with 142
093 LLMs and may also hurt the generalizability of 143
094 LLMs across different tasks. 144

095 To tackle these challenges, we propose retrieval- 145
096 augmented world models (RAWM). Specifically, 146
097 our contributions are threefold. First, inspired by 147
098 the retrieval-augmented generation (RAG) (Lewis 148
099 et al., 2020), we introduce the memory, which 149
100 stores the pre-collected experiences from the en- 150
101 vironments, and the embedding model, which is 151
102 used for querying relevant experience to assist the 152
103 world model to make predictions. Second, we in- 153
104 troduce the reinforcement learning (RL) training 154
105 pipeline, which adds a small MLP head to the pre- 155
106 trained embedding model and trains the MLP layer 156
107 with proximal policy optimization (PPO) (Schul- 157
108 man et al., 2017). Third, we collect the data from 158
109 Game24, BlocksWorld and BabyAI, and extensive 159
110 experiments demonstrate RAWM can significantly 160
111 outperform the world model without retrieved ex- 161
112 periences and the pre-trained embedding models 162
113 and demonstrate the generalizability. RAWM is an 163
114 efficient way for LLMs to obtain the environment- 164
115 specific knowledge to build the better world mod- 165
116 els without training LLMs, and our RL training 166
117 pipeline can further improve the accuracy of the 167
118 predictions of LLM-based world models efficiently.

119 2 Related Work 167

120 **World Models and LLMs.** MuZero (Schrit- 168
121 twieser et al., 2020) and Dreamer (Hafner et al., 169
122 2019) are the two prominent examples of the world 170
123 model for complex decision making tasks. Tra- 171
124 jectory transformer (Janner et al., 2021) leverages 172
173

125 transformer to model the decision making as a 126
127 sequence modeling problem. The world models 128
129 trained in these methods are environment specific 130
131 and cannot generalize to other environments. Re- 132
133 cently, researchers leverage the generalizability of 134
135 LLMs to build general world models for reasoning 136
137 and decision making (Hao et al., 2023; Wang et al., 138
139 2024; Yang et al., 2024b; Lin et al., 2024). Specifi- 140
141 cally, RAP (Hao et al., 2023) and RAFA (Liu et al., 142
143 2023) use LLMs to predict the next states explicitly 144
145 and planning methods for decision making. While 146
147 ToT (Yao et al., 2023) and GoT (Besta et al., 2024) 148
149 use LLMs as the world model implicitly to advance 150
151 and evaluate the different thoughts. 152

153 **Retrieval-Augmented Generation.** RAG is an 154
155 efficient way for LLMs to incorporate the ex- 156
157 ternal knowledge for generation and understand- 157
158 ing (Lewis et al., 2020; Gao et al., 2023). Specifi- 158
159 cally, RAG leverages the retrieval model to query 159
160 the relevant experiences from the memory, which 160
161 are further provided to the LLMs as the in-context 161
162 examples. Different from simple prompting, where 162
163 the external knowledge is provided by human- 163
164 written prompts (Wang et al., 2024), and simple 164
165 in-context learning, where the in-context examples 165
166 are randomly picked (Hao et al., 2023), RAG can 166
167 provide better examples for accurate predictions. 167
168 Compared with fine-tuning (Xiang et al., 2023), 168
169 RAG can provide a more efficient way of trans- 169
170 forming LLMs into world models. 170

171 **RL for LLM Optimization.** RL is a powerful 171
172 method to train the model with trial and error (Sut- 172
173 ton and Barto, 2018). In addition to the applica- 173
174 tions of RL in games and robotics (Silver et al., 174
175 2017) to optimize the LLMs, such as optimizing 175
176 the prompts (Deng et al., 2022) and the decoding 176
177 process (Wan et al., 2024), recent works also lever- 177
178 age RL to improve the reasoning capabilities of 178
179 LLMs (Lambert et al., 2024; Guo et al., 2025). In 179
180 this work, we leverage the RL method to train the 180
181 retriever to find the better examples to boost the 181
182 prediction of the world model. 182

183 3 Preliminaries 177

184 In this section, we present the preliminaries of 178
185 RAWM, including the formulation of the decision 179
186 making, the LLMs, and the world models. 180

181 **Markov Decision Process (MDP).** A decision 171
182 making problem is usually represented as a Markov 172
183 decision process (MDP) (Sutton and Barto, 2018), 173

which is defined by the tuple $\mathcal{M} = (S, A, T, R, \gamma)$, where S is the state space, A is the action space, $T : S \times A \rightarrow S$ is the transition dynamics, which specifies the next state s' given the current state s and action a , $R : S \times A \rightarrow \mathbb{R}$ is the reward function, which specifies the agent's reward given the current state s and action a , and γ is the discount factor. The agent's policy is defined by $\pi_\theta : S \times A \rightarrow [0, 1]$, parameterized by θ , which takes the state s as the input and outputs the action a to be executed. The objective of the agent is to learn an optimal policy to maximize the expected return, i.e., $\pi^* := \arg \max_\pi \mathbb{E}_\pi [\sum_{t=0}^{\infty} \gamma^t r_t | s_0]$ with initial state s_0 .

Large Language Models (LLMs). Large Language models (LLMs) learn from text data using unsupervised learning. LLMs optimize the joint probabilities of variable-length symbol sequences as the product of conditional probabilities by $P(x) = \prod_{i=1}^n P(s_i | s_1, \dots, s_{i-1})$, where (s_1, s_2, \dots, s_n) is the variable-length sequence of symbols. With the billions of parameters and extensive training data, the vast amounts of common knowledge encoded in LLMs lead to the remarkable generalization across various NLP tasks with simple prompting and in-context learning, and without task-specific fine-tuning (Touvron et al., 2023; OpenAI, 2023). Among them, RAG (Lewis et al., 2020) is viewed as a powerful method to incorporate external knowledge to LLMs for generation. Given the generalizability, LLMs present promising foundations for general world models.

World Models. The world model Ω is introduced to predict the dynamics of the environment, thus supporting the decision making process. Specifically, the world model is trained or prompted to predict the next state s' , the reward r , and the terminal function d , given the current state s and action a . The world model can be one or multiple neural networks specially trained on the environments for the three prediction tasks (Hafner et al., 2019; Schrittwieser et al., 2020), which cannot generalize across different environments. Recent works leverage LLMs to build the general world models, where the prompting (Xie et al., 2024), in-context learning (Wang et al., 2024), and even fine-tuning methods (Xiang et al., 2023; Lin et al., 2024) are used to build the LLM-based world models. In this work, we primarily focus on the prediction of the next state, which is the most important feature, as both the reward and terminal are usually derived from the next state visited.

4 Retrieval-Augmented World Models

In this section, we introduce **Retrieval-Augmented World Models (RAWM)**. We will first introduce the architecture of RAWM and then introduce the RL training pipeline for the retrieval process.

4.1 Architecture

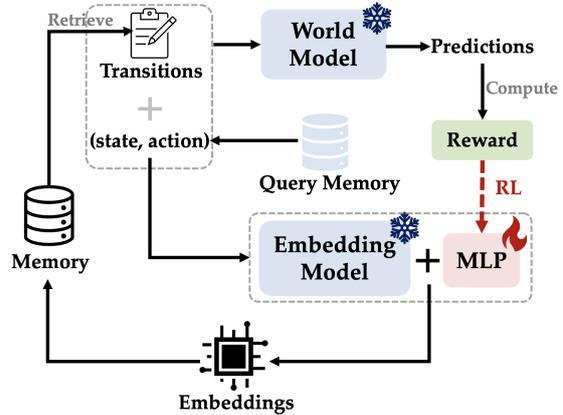


Figure 2: The overview of RAWM.

The architecture of RAWM is displayed in Figure 2. We introduce a memory Ξ , which stores the pre-collected experiences, an embedding model, which is used to rank and retrieve the relevant experiences. Specifically, given the query $q = (s, a) \in Q$, where Q is the query dataset, we will use the embedding model to query top K relevant experiences $c = \langle c_k \rangle$, where $c_k = (s_k, a_k, s'_k)$, $k = 1, \dots, K$. The retrieved experiences c will be concatenated with the query q to form the input to the world model Ω . We note that for the environments where the states are not texts, e.g., BabyAI (Chevalier-Boisvert et al., 2019a), we need to first transform them into the text representation.

Prompt Design. For the prompt design, any information related to the environments will not be provided to the world model, including the tasks, the object and action rules. We expect that all the environment knowledge is provided by the in-context examples retrieved from the memory. The prompt template is displayed as follows:

```

Prompt Template
System prompt:
"After being given a current state and an action, directly
give the next state after performing the action."
Content prompt:
Current state: <text of the current state>
Action: <text of the selected action>
Next state: <text of the next state> or <for prediction>

```

The system prompt provides a general description of the prediction tasks, and the content prompt

includes the query and the context examples. For the context examples, the next state is provided, while for the query, the next state is predicted by the world model Ω . Similarly, this content template is also used to get the embeddings of both query dataset and the memory for the retrieval process.

Trainable Embedding of Transitions. We use the pre-trained embedding model ϕ to encode the transitions into the M -dimensional vector representation. Specifically, for the query dataset, we only encode the state and the action, and for the memory, we encode the state, the action and the next state. However, the embedding model is trained over general corpus, which would be not suitable to the specific environment, so adapting the embedding model is needed. There are several methods to adapt the embedding model to the specific environment: i) fine-tuning all parameters in ϕ , which is not training efficient, ii) low-rank adaption (LoRA) (Hu et al., 2022), which introduces trainable low-rank decomposition matrices for each layer to reduce the parameters to be trained. Though the number of trained parameters is reduced, LoRA still requires to leverage the full embedding model to inference. Besides, both full-parameter fine-tuning and LoRA requires that the access of the parameters of the pre-trained embedding model and cannot be applied to close-source models, e.g., text-embedding-3. Therefore, inspired by the linear probe (Radford et al., 2021), we introduce a trainable MLP module above the pre-trained embedding model, which is denoted as ψ . Therefore, the embedding process for both query data and the memory can be represented as:

$$e_q = \psi(\phi(s, a)), \forall (s, a) \in \mathcal{Q}, \quad (1)$$

$$e_c = \psi(\phi(s, a, s')), \forall (s, a, s') \in \Xi. \quad (2)$$

We will introduce the RL training pipeline of ψ in the next section and the parameters in ϕ are frozen. Compared with the full parameter fine-tuning and the LoRA, this method only requires the pre-trained embedding to encode the data in the query dataset and the memory once, and the number of trainable parameters is even significantly less than LoRA.

Retrieval-Augmented Predictions. To query the relevant experiences, a similarity measure, e.g., cosine similarity, is used to rank the examples in the memory, which is denoted as $\text{sim}(\cdot)$. Therefore,

$$\mathbf{c} = \{c_k | k \in \text{topK}(\text{sim}(e_q, e_c)), \forall c \in \Xi\}, \quad (3)$$

where $\text{topK}(\cdot)$ is selecting the indices with the top- K maximum values. The K retrieved examples \mathbf{c} will be formed the in-context examples and append before the query for the prediction. We concatenate the in-context examples with the query in a **reverse order**, i.e., the examples with larger similarities will be the later examples, and the query is the last one. We found that this reverse order is important for the generalization of the embedding model in different K values, as the reverse order can ensure the last several examples be the same, (e.g., for $K \in \{1, 2\}$, the top-1 example is the same, which is the last example before the query in the prompt), thus leading to a more stable generalization performance of the world model.

Evaluation Measure. The evaluation measure is important for the RL training. We follow RAP (Hao et al., 2023) to design the reward: given the output o from the world model, which may include a set of the conditions, e.g., the predicted state of blocks, and s' is the target, we will calculate the accuracy of the prediction, denoted as $v(o, s')$. Alternatively, we can calculate the log likelihood of the target s' , which is used in the original RAG (Lewis et al., 2020). However, this may require the access of the logits of the LLMs and cannot be applied to the closed-source models, e.g., GPT-4o.

4.2 Training

In this section, we introduce the efficient RL pipeline to train the embedding models, i.e., training of the MLP head ψ specifically. Typically, the retriever in RAG is trained with supervised learning (Lewis et al., 2020). However, in RAWM, the world models are not trained and we cannot compute the gradient of the embedding directly. Besides, as the retriever needs to explore to choose the examples for the better prediction with the world model, RL is one of the straightforward methods to optimize the embedding model.

One-step Decision Making. To apply RL methods to optimize the embedding model, we need to build the MDP \mathcal{M}^ψ for the embedding ψ^1 :

- State space $S^\psi : \{\phi(s, a), \forall (s, a) \in \mathcal{Q}\} \cup \{\phi(s, a, s'), \forall (s, a, s') \in \Xi\}$, i.e., the embeddings of all data from query dataset and the memory generated by the pre-trained model ϕ .
- Action space $\mathcal{A}^\psi \in \mathbb{R}^{M'}$, where M' is the output dimension of ψ , i.e., ψ will transform the

¹Please distinguish \mathcal{M}^ψ with the one used for the environment \mathcal{M} , where \mathcal{M}^ψ is introduced only for the training.

embeddings by ϕ to M' -dimensional vectors.

- Reward $r = v(\Omega(q, c), s')$, where $\Omega(q, c)$ is the output of the world model Ω with the input (q, c) . We note that \mathcal{M}^ψ is a one-step decision making problem, i.e., \mathcal{M}^ψ always ends after the first time step, so the transition function and the discount factor are not necessary for the RL training.

Design of ψ . Before diving into the RL training, we first discuss about the design of ψ . A simple setting for ψ is a randomly initialized MLP, which means this initialization will start with the random embedding for the training and ignore the embeddings generated by the pre-trained model ϕ . On the other hand, we can initialize the MLP with an identify matrix, i.e., $\psi = I$.² Both methods have their own advantages and disadvantages: for the random initialization, we can arbitrarily choose the output dimension and the activation function of ψ , but the training will start with a relatively worse performance, while for the identify initialization, the output dimension of ψ must be the same with ϕ , i.e., $M' = M$, and the training will start with the performance of the pre-trained embedding model.

RL Training. RL methods rely on the trial-and-error process to explore the solution space for better policies. The primary RL method is Q-learning (Watkins and Dayan, 1992; Mnih et al., 2015), which can only be used on the problems with discrete actions, and the policy gradient methods are proposed for the problems with both discrete and continuous actions (Sutton et al., 1999; Mnih et al., 2016; Haarnoja et al., 2018). PPO (Schulman et al., 2017) is an on-policy policy gradient method, which is a simplified, but more data efficient and reliable, variant of Trust Region Policy Optimization (TRPO) (Schulman et al., 2015), which leverages the “trust region” to bound the update of the policy to avoid training collapse. Compared with TRPO, PPO is more data efficient and with more reliable performances than TRPO, while only using the first-order optimization for computational efficiency. Specifically, PPO is maximizing the objective

$$J(\psi) = \mathbb{E} [\min(\rho_\psi \cdot r, \text{clip}(\rho_\psi, 1 - \epsilon, 1 + \epsilon) \cdot r)], \quad (4)$$

where ρ_ψ is the importance sampling ratio conditional on ψ , r is the reward, and ϵ is the hyperparameter which controls the boundary of the trust

²With a slight abuse of notations, we use ψ to represent both the MLP and the parameters.

region. We note that the advantages in the general PPO implementation is replaced with the reward. We only provide a short introduction of PPO in this section, as we take PPO as a blackbox for optimizing ψ . The full training procedure is displayed in Algorithm 1. Other RL methods, e.g., soft actor critic (SAC) (Haarnoja et al., 2018), can also be used and for more details of RL, we refer readers to the book (Sutton and Barto, 2018).

Algorithm 1 Training of RAWM

- 1: **Input:** World model Ω , pre-trained embedding model ϕ , memory \mathcal{M} , Query dataset \mathcal{Q} , number of retrieval candidates K
 - 2: Initialize the MLP ψ .
 - 3: Computing the embeddings with ϕ , i.e., $\mathcal{Q}_\phi = \{\phi(s, a), \forall (s, a) \in \mathcal{Q}\}$ and $\mathcal{M}_\phi = \{\phi(s, a, s'), \forall (s, a, s') \in \mathcal{M}\}$.
 - 4: **for** iter $\in \{1, 2, \dots\}$ **do**
 - 5: Update the memory embedding $\mathcal{M}_\psi = \{\psi(\phi(s, a, s')), \forall (s, a, s') \in \mathcal{M}\}$.
 - 6: **for** (s, a) in \mathcal{Q} **do**
 - 7: Compute query embedding $\psi(\phi(s, a))$.
 - 8: Select top- K relevant transitions c from \mathcal{M} with the embedding in \mathcal{M}_ψ .
 - 9: Generate the prediction o and compute the reward $v(o, s')$.
 - 10: **end for**
 - 11: Train ψ with PPO, i.e., Eq. (4).
 - 12: **end for**
-

5 Experiments

In this section, we present the extensive experiments to evaluate the effectiveness of RAWM. We will first introduce the experiment setup and then the experiment results and analysis.

5.1 Setup

Environments. The environments considered in this work include (as shown in Figure 3)

- **Game24:** a mathematical puzzle game where four numbers are given (e.g., 10, 3, 6, and 4) and the player can only use the basic arithmetic operations, i.e., $(+, -, \times, \div)$, to obtain 24 (e.g., $10 \times (6 \div 3) + 4$). This puzzle is widely used to benchmark the LLMs’ reasoning capabilities (Yao et al., 2023) and the LLMs need to generate a sequence of operations to obtain 24. In this game, the world model needs to correctly generate the remaining number when an opera-

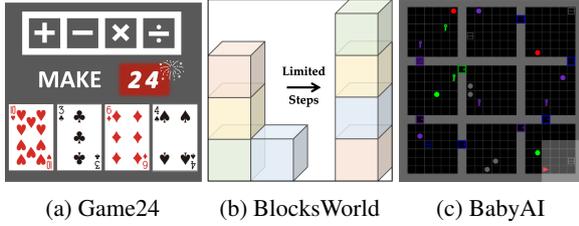


Figure 3: Environments

tion is executed, i.e., 10, 2, 4 are the remaining numbers when $6 \div 3$ is executed.

- **BlocksWorld**: a simple world of blocks where a set of blocks is placed on the plat and the player needs to perform the basic actions, i.e., pick up, put down, stack, and unstack, to transform the blocks to a target configuration (Valmeekam et al., 2023; Hao et al., 2023). In this game, the world model needs to predict the states for all blocks (e.g., the blue block is on top of the red block) after an action is executed (e.g., stack blue block on the red block).
- **BabyAI**: a grid world with objects where the agent needs to complete the tasks defined with language instructions (Chevalier-Boisvert et al., 2019a) with the actions, i.e., turn left, turn right, move forward and pick up. We use the text description of the states in (Carta et al., 2023) for the environments. In this environment, the world model needs to predict the locations of the objects after performing the action.

Datasets. Given the environments, we need to collect the datasets for the memory, query and test datasets, respectively. We use the query dataset to train the embedding with RL and use the test dataset to validate the performance of the trained models. For Game24 and BlocksWorld, the number of all possible transitions are less than 10K, therefore, we use the Depth-First Search (DFS) to enumerate all transitions to form the full datasets. While for BabyAI, we cannot enumerate all transitions due to the complexity of the environments. Therefore, we utilize the bot provided in (Chevalier-Boisvert et al., 2019b) to collect the data, where we enumerate all valid actions to gather the transitions along the action sequences generated by the bot. After the collection, we choose the separate subsets to form the three datasets without any overlapping to avoid any data leakage. We provide the detailed introduction of the environments and the protocol for data collection in Appendix C.

Model Selection. We use the embedding model Alibaba-NLP/gte-Qwen2-1.5B-instruct as the

pre-trained ϕ , which is the leading open-source text embedding model on MTEB (Li et al., 2023). For the world model, we choose the Qwen-2.5 instruct model series with the model sizes as {1.5B, 3B, 7B} (Yang et al., 2024a)³. The AWQ quantized models are chosen for efficient inference. For the configuration of ψ , we consider a three layer MLP with Tanh() activation function for the random initialization and a single layer without any activation function for the identity initialization.⁴ We provide the details about the model selection and the initialization in Appendix E.

RL Training. For the efficiency, we consider several implementation tricks. i) Compared with the training of the MLP ψ , the inference of the world model is much more time-consuming. Therefore, we enlarge the number of batch sizes and for each batch, we sample multiple times, which can stabilize the training. ii) We also consider fixing the embeddings in the memory, i.e., only the embeddings of the query datasets are trained, and do not observe the advantages. Therefore, we update the embedding of both datasets. iii) The output dimension of the random initialization is much smaller than the output dimension of the identity initialization, which enjoys the training stabilities with larger learning rates and smaller memory usages when retrieval. The hyperparameters for the RL training of ψ is provided in Appendix F.

Methods Evaluated. The methods evaluated in the experiments are: i) zero-shot: the world models give the prediction without any in-context examples (Wang et al., 2024), ii) random: the world models give the prediction with randomly selected in-context examples from \mathcal{M} (Hao et al., 2023), iii) RAWM $_{\psi,rand}$: RAWM with the randomly initialization of ψ , which differs from the previous method, iv) RAWM $_{\psi,eye}$: RAWM with the identity initialization of ψ , equivalent to the pre-trained embedding model ϕ , v) RAWM $_{\psi,rand}^{RL}$: RAWM with randomly initialized ψ and RL training, and vi) RAWM $_{\psi,eye}^{RL}$: RAWM with identity initialized ψ and RL training.

5.2 Evaluation

We present the extensive evaluations of RAWM in this section. There are three main research questions (RQs) investigated:

³<https://huggingface.co/spaces/Qwen/Qwen2.5>

⁴We would note that RAWM can work for both close-source and open-source embedding and world models. We choose open-source models for efficient training and inference.

Model	Method	Game24				BlocksWorld				BabyAI			
		$K = 1$		$K = 2$		$K = 1$		$K = 2$		$K = 1$		$K = 2$	
		train	test	train	test	train	test	train	test	train	test	train	test
1.5B	zero-shot	0.5224	0.5455	0.5224	0.5455	0.3804	0.3849	0.3804	0.3849	0.3786	0.3772	0.3786	0.3772
	random	0.5586	0.5664	0.5714	0.5959	0.4848	0.4822	0.4975	0.4991	0.3851	0.3856	0.3973	0.4030
	RAWM $_{\psi,rand}$	0.5156	0.5219	0.5322	0.5534	0.5386	0.5402	0.5597	0.5589	0.3415	0.3479	0.3527	0.3484
	RAWM $_{\psi,eye}$	0.5352	0.5474	0.5510	0.5600	0.5659	0.5697	0.5878	0.5888	0.4427	0.4446	0.4710	0.4671
3B	zero-shot	0.4888	0.4971	0.4888	0.4971	0.3644	0.3661	0.3644	0.3661	0.3303	0.3330	0.3303	0.3330
	random	0.6703	0.6719	0.6984	0.7010	0.4717	0.4706	0.5089	0.5083	0.3912	0.3908	0.4073	0.4052
	RAWM $_{\psi,rand}$	0.7041	0.7043	0.7269	0.7292	0.5729	0.5739	0.6005	0.6019	0.3855	0.3892	0.3985	0.3991
	RAWM $_{\psi,eye}$	0.7022	0.7179	0.7313	0.7463	0.6127	0.6102	0.6440	0.6397	0.4355	0.4297	0.4646	0.4633
7B	zero-shot	0.5957	0.6121	0.5957	0.6121	0.5215	0.5207	0.5215	0.5207	0.4201	0.4254	0.4201	0.4254
	random	0.8241	0.8267	0.8712	0.8667	0.5897	0.5838	0.6021	0.6072	0.4084	0.4181	0.4178	0.4221
	RAWM $_{\psi,rand}$	0.8362	0.8375	0.8724	0.8703	0.6274	0.6240	0.6332	0.6314	0.4301	0.4322	0.4403	0.4355
	RAWM $_{\psi,eye}$	0.8511	0.8527	0.8781	0.8734	0.6472	0.6452	0.6556	0.6541	0.4484	0.4501	0.4633	0.4693

Table 1: Performance of RAWM with the retrieval mechanism over three environments.

- **RQ1:** Can the retrieved methods in RAWM improve the performance of world model?
- **RQ2:** Can the RL training pipeline in RAWM improve the performance of the world model, compared with pre-trained models?
- **RQ3:** Can the learned model generalize across different settings, e.g., different values of K ?

5.2.1 Analysis of RQ1

To investigate the RQ1, we conduct the experiments of RAWM on the different sizes of the world model, i.e., 1.5B, 3B and 7B, over the three environments. We consider the values of K as $\{1, 2\}$. The experiment results are displayed in Table 1. From the results, we observe that the performances over the train and test yield the same trend, which avoids the over-fitting to the specific dataset.

With more in-context examples selected, the performance of the world model is significantly improved, which is consistent with other research (Agarwal et al., 2024). Another interesting observation is that increasing the model sizes of LLMs does not necessarily improve the performance of the world models. For example, the 3B world model performs worse than the 1.5B world model in BabyAI. This is because the LLMs do not have the specific knowledge of the environments and increasing the model size cannot solve this.

We also observe that given the same number of the in-context examples, the pre-trained model (i.e., RAWM $_{\psi,eye}$) can retrieve more relevant examples for the world models across different sizes in BlocksWorld and BabyAI. While for the 1.5B world model of Game24, the pre-trained models perform worse than the random examples. There-

fore, optimizing for a better embedding model can potentially further improve the performance.

5.2.2 Analysis of RQ2

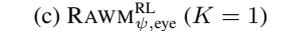
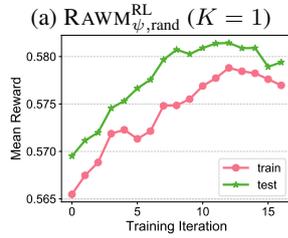
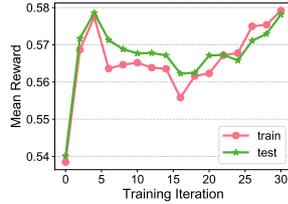


Figure 4: Training curves on BlocksWorld.

We then present the results of the RL training pipeline of RAWM. Due to the limitation of the resource, we only conduct the training on the world models with 1.5B LLMs. The results of different configurations of ψ across different environments are displayed in Figure 5.

From the results, we observe that the RL training can improve upon the initialization, which indicates the capability of RL to optimize the embedding model through exploration. We observe that both initialization can outperform the pre-trained embedding model, i.e., RAWM $_{\psi,eye}$, in Game24 and BlocksWorld, while the random initialization fails to find a better embedding than the pre-trained one in BabyAI. The training curves are displayed in

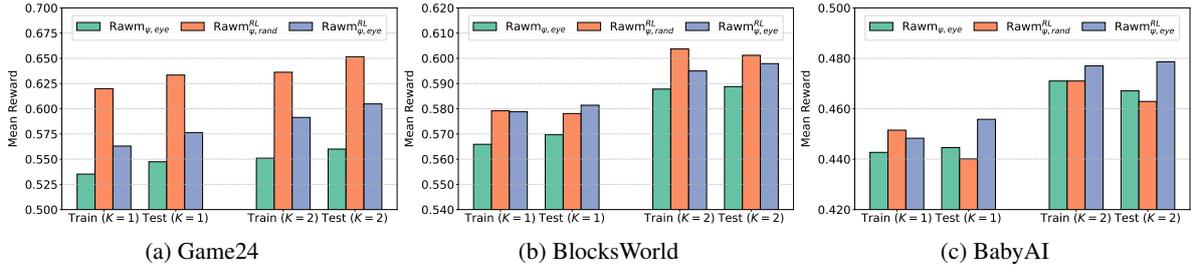


Figure 5: Performance of the RL training pipeline in RAWM over three environments.

Method	Game24				BlocksWorld				BabyAI			
	$K = 3$		$K = 5$		$K = 3$		$K = 5$		$K = 3$		$K = 5$	
	train	test	train	test	train	test	train	test	train	test	train	test
random	0.5745	0.5862	0.5669	0.5866	0.5096	0.5125	0.5261	0.5228	0.4044	0.4071	0.4220	0.4165
RAWM $_{\psi,rand}$	0.5431	0.5443	0.5538	0.5635	0.5702	0.5711	0.5730	0.5738	0.3522	0.3551	0.3696	0.3636
RAWM $_{\psi,eye}$	0.5533	0.5528	0.5660	0.5765	0.6016	0.5994	0.6178	0.6149	0.4838	0.4753	0.4816	0.4860
RAWM $_{\psi,rand}^{RL} (K = 1)$	0.5766	0.5974	0.6002	0.6106	0.6001	0.6022	0.6199	0.6200	0.4716	0.4624	0.4745	0.4702
RAWM $_{\psi,eye}^{RL} (K = 1)$	0.5893	0.5950	0.5976	0.6053	0.6038	0.6042	0.6222	0.6220	0.4878	0.4877	0.4982	0.4872
RAWM $_{\psi,rand}^{RL} (K = 2)$	0.6097	0.6344	0.5901	0.5999	0.6100	0.6129	0.6198	0.6202	0.4732	0.4711	0.4738	0.4741
RAWM $_{\psi,eye}^{RL} (K = 2)$	0.5912	0.6020	0.5981	0.6067	0.6049	0.6052	0.6215	0.6205	0.4864	0.4852	0.4976	0.4888

Table 2: Shot generalization of the 1.5B world model trained with RL.

Figure 4. Typically, the random initialization will let the model train from a relative low performance and we observe a drop of the performance due to the exploration for better embedding model (i.e., Figure 4b). And for the identity initialization, the training is more stable with smaller learning rates (i.e., Figures 4c and 4d). These results indicate the effectiveness of our RL training pipeline.⁵

Our RL training pipeline can also be used to diagnose the failure of the retrieval-augmented generation systems. If the RL pipeline cannot find a better embedding to improve the world model’s performance, then the user would replace the LLMs for the world models and the datasets.⁶

5.2.3 Analysis of RQ3

The results of shot generation are displayed in Table 2, where the embedding models trained with random and identity initializations of $K \in \{1, 2\}$ are evaluated over the $K \in \{3, 5\}$, i.e., the generalization over shots. From the results, we observe that with larger values of K , the performance of the world model will be further improved. The embedding models trained with RL pipeline demonstrate

⁵We note that the improvement that RL training can bring will largely be influenced by the LLMs’ capabilities.

⁶Different from the factual QA (Gao et al., 2023) where we can manually check whether the retrieved examples are correct or not, RAWM relies on the LLM’s inherent understanding capabilities for the prediction and human cannot manually check the correctness of the retrieval. Therefore, a systematic method, e.g., RL, is needed for diagnosing the system.

to be more capable for the generalization over shots, compared with the pre-trained embedding model.

Method	Game24		BlocksWorld		BabyAI	
	Train	Test	Train	Test	Train	Test
RAWM $_{\psi,rand}$	0.8724	0.8703	0.6332	0.6314	0.4403	0.4355
RAWM $_{\psi,eye}$	0.8781	0.8734	0.6556	0.6541	0.4633	0.4693
RAWM $_{\psi,rand}^{RL}$	0.8799	0.8829	0.6631	0.6630	0.4518	0.4484
RAWM $_{\psi,eye}^{RL}$	0.8852	0.8812	0.6597	0.6560	0.4700	0.4721

Table 3: Model generalization of 1.5B \rightarrow 7B ($K = 2$).

We also consider the generalization over different LLMs, which is more difficult than the shot generalization. Table 3 displays the results of generalizing the RL trained embedding model from 1.5B to 7B. We observe that the RL trained embedding admits better generalization performance than the pre-trained embedding model.

6 Conclusions

LLMs present promising foundations to build general world models. However, LLMs usually lack the specific knowledge of environments. Therefore, we introduce **Retrieval-Augmented World Models** (RAWM), which leverages the retrieval-augmented generation to improve LLM-based world models. We then introduce an efficient RL training pipeline to further improve the performance. Extensive experiments demonstrate the effectiveness and the generalizability of RAWM. RAWM is an efficient method to build the highly capable LLM-based world models without fine-tuning the LLMs.

612 Limitations

613 There are several limitations of current work.

- 614 • Current RAWM focuses on prediction and the
615 prediction can be used for decision making. We
616 will extend current RAWM to support the better
617 decision making in future work.
- 618 • Current RAWM is based on the pre-collected data,
619 which may require a large number of data to
620 achieve good performance. We will consider to
621 let the model to proactively collect the data and
622 improve the performance automatically.
- 623 • Current RAWM is based on LLM and the envi-
624 ronments are represented by texts. RAWM can
625 also handle the multi-modal environments, e.g.,
626 text and image, where both embedding models
627 and world models will be multi-modal models.
628 We will explore this direction in future work.

629 We expect that RAWM can be a general framework
630 to build highly capable multi-modal world model
631 with automatically data collection and training, to
632 finally support decision making in complex tasks.

633 Ethics Statement

634 We confirm that we have fully complied with the
635 ACL Ethics Policy in this study. All the environ-
636 ments are publicly available and have been exten-
637 sively used in the research.

638 References

639 Rishabh Agarwal, Avi Singh, Lei M Zhang, Bernd
640 Bohnet, Luis Rosias, Stephanie Chan, Biao Zhang,
641 Ankesh Anand, Zaheer Abbas, Azade Nova, et al.
642 2024. Many-shot in-context learning. *arXiv preprint*
643 *arXiv:2404.11018*.

644 Maciej Besta, Nils Blach, Ales Kubicek, Robert Ger-
645 stenberger, Michal Podstawski, Lukas Gianinazzi,
646 Joanna Gajda, Tomasz Lehmann, Hubert Niewiadom-
647 ski, Piotr Nyczyk, et al. 2024. Graph of thoughts:
648 Solving elaborate problems with large language mod-
649 els. In *AAAI*, pages 17682–17690.

650 Tom B Brown, Benjamin Mann, Nick Ryder, Melanie
651 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
652 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
653 Askeel, et al. 2020. Language models are few-shot
654 learners. In *NeurIPS*, pages 1877–1901.

655 Jake Bruce, Michael D Dennis, Ashley Edwards, Jack
656 Parker-Holder, Yuge Shi, Edward Hughes, Matthew
657 Lai, Aditi Mavalankar, Richie Steigerwald, Chris
658 Apps, et al. 2024. Genie: Generative interactive
659 environments. In *ICML*.

660 Arunkumar Byravan, Jost Tobias Springenberg, Ab-
661 bas Abdolmaleki, Roland Hafner, Michael Neunert,

662 Thomas Lampe, Noah Siegel, Nicolas Heess, and
663 Martin Riedmiller. 2020. Imagined value gradients:
664 Model-based policy optimization with transferable la-
665 tent dynamics models. In *CoRL*, pages 566–589.

666 Thomas Carta, Clément Romac, Thomas Wolf, Sylvain
667 Lamprier, Olivier Sigaud, and Pierre-Yves Oudeyer.
668 2023. Grounding large language models in interac-
669 tive environments with online reinforcement learning.
670 *arXiv preprint arXiv:2302.02662*.

671 Hyungjoo Chae, Namyoung Kim, Kai Tzu iunn Ong,
672 Minju Gwak, Gwanwoo Song, Jihoon Kim, Sungh-
673 wan Kim, Dongha Lee, and Jinyoung Yeo. 2025.
674 Web agents with world models: Learning and lever-
675 aging environment dynamics in web navigation. In
676 *ICLR*.

677 Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem
678 Lahlou, Lucas Willems, Chitwan Saharia, Thien Huu
679 Nguyen, and Yoshua Bengio. 2019a. BabyAI: First
680 steps towards grounded language learning with a hu-
681 man in the loop. In *ICLR*.

682 Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem
683 Lahlou, Lucas Willems, Chitwan Saharia, Thien Huu
684 Nguyen, and Yoshua Bengio. 2019b. BabyAI: First
685 steps towards grounded language learning with a hu-
686 man in the loop. In *ICLR*.

687 Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan
688 Wang, Han Guo, Tianmin Shu, Meng Song, Eric
689 Xing, and Zhiting Hu. 2022. Rlprompt: Optimizing
690 discrete text prompts with reinforcement learning. In
691 *EMNLP*.

692 Shenyuan Gao, Jiazhi Yang, Li Chen, Kashyap Chitta,
693 Yihang Qiu, Andreas Geiger, Jun Zhang, and
694 Hongyang Li. 2024. Vista: A generalizable driving
695 world model with high fidelity and versatile control-
696 lability. *arXiv preprint arXiv:2405.17398*.

697 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia,
698 Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen
699 Wang. 2023. Retrieval-augmented generation for
700 large language models: A survey. *arXiv preprint*
701 *arXiv:2312.10997*.

702 Yu Gu, Boyuan Zheng, Boyu Gou, Kai Zhang, Cheng
703 Chang, Sanjari Srivastava, Yanan Xie, Peng Qi, Huan
704 Sun, and Yu Su. 2024. Is your llm secretly a world
705 model of the internet? model-based planning for web
706 agents. *arXiv preprint arXiv:2411.06559*.

707 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,
708 Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma,
709 Peiyi Wang, Xiao Bi, et al. 2025. DeepSeek-R1: In-
710 centivizing reasoning capability in llms via reinforc-
711 e learning. *arXiv preprint arXiv:2501.12948*.

712 David Ha and Jürgen Schmidhuber. 2018. World mod-
713 els. *arXiv preprint arXiv:1803.10122*.

714 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and
715 Sergey Levine. 2018. Soft actor-critic: Off-policy
716 maximum entropy deep reinforcement learning with
717 a stochastic actor. In *ICML*, pages 1861–1870.

718	Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. 2019. Dream to control: Learning behaviors by latent imagination. In <i>ICLR</i> .	774
719		775
720		776
721	Danijar Hafner, Timothy P Lillicrap, Mohammad Norouzi, and Jimmy Ba. 2021. Mastering Atari with discrete world models. In <i>ICLR</i> .	777
722		778
723		
724	Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. 2023. Mastering diverse domains through world models. <i>arXiv preprint arXiv:2301.04104</i> .	
725		
726		
727		
728	Jessica B Hamrick, Abram L. Friesen, Feryal Behbahani, Arthur Guez, Fabio Viola, Sims Wither- spoon, Thomas Anthony, Lars Holger Buesing, Petar Veličković, and Theophane Weber. 2021. On the role of planning in model-based deep reinforcement learning. In <i>ICLR</i> .	
729		
730		
731		
732		
733		
734	Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. 2023. Reasoning with language model is planning with world model. <i>arXiv preprint arXiv:2305.14992</i> .	
735		
736		
737		
738	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In <i>ICLR</i> .	
739		
740		
741		
742	Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. 2019. When to trust your model: Model-based policy optimization. <i>arXiv preprint arXiv:1906.08253</i> .	
743		
744		
745		
746	Michael Janner, Qiyang Li, and Sergey Levine. 2021. Offline reinforcement learning as one big sequence modeling problem. In <i>NeurIPS</i> .	
747		
748		
749	Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. 2024. Tulu 3: Pushing frontiers in open language model post-training. <i>arXiv preprint arXiv:2411.15124</i> .	
750		
751		
752		
753		
754		
755	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. In <i>NeurIPS</i> , pages 9459–9474.	
756		
757		
758		
759		
760		
761	Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. 2023. Towards general text embeddings with multi-stage contrastive learning. <i>arXiv preprint arXiv:2308.03281</i> .	
762		
763		
764		
765	Jessy Lin, Yuqing Du, Olivia Watkins, Danijar Hafner, Pieter Abbeel, Dan Klein, and Anca Dragan. 2024. Learning to model the world with language. In <i>ICML</i> .	
766		
767		
768		
769	Zhihan Liu, Hao Hu, Shenao Zhang, Hongyi Guo, Shuqi Ke, Boyi Liu, and Zhaoran Wang. 2023. Reason for future, act for now: A principled framework for autonomous llm agents with provable sample efficiency. <i>arXiv preprint arXiv:2309.17382</i> .	
770		
771		
772		
773		
	Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. 2016. Asynchronous methods for deep reinforcement learning. In <i>ICML</i> , pages 1928–1937.	779
		780
		781
		782
		783
		784
	Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidje- land, Georg Ostrovski, et al. 2015. Human-level control through deep reinforcement learning. <i>Nature</i> , 518(7540):529–533.	
	OpenAI. 2023. GPT-4 technical report. <i>arXiv preprint arXiv:2303.08774</i> .	785
		786
	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas- try, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In <i>ICML</i> , pages 8748–8763.	787
		788
		789
		790
		791
		792
	Alexander Robey, George J Pappas, and Hamed Has- sani. 2021. Model-based domain generalization. In <i>NeurIPS</i> , pages 20210–20229.	793
		794
		795
	Julian Schrittwieser, Ioannis Antonoglou, Thomas Hu- bert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. 2020. Mastering Atari, Go, chess and shogi by planning with a learned model. <i>Nature</i> , 588(7839):604–609.	796
		797
		798
		799
		800
		801
	Julian Schrittwieser, Thomas Hubert, Amol Mandhane, Mohammadamin Barekatin, Ioannis Antonoglou, and David Silver. 2021. Online and offline reinforce- ment learning by planning with a learned model. In <i>NeurIPS</i> , pages 27580–27591.	802
		803
		804
		805
		806
	John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. 2015. Trust region policy optimization. In <i>ICML</i> , pages 1889–1897.	807
		808
		809
	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proxi- mal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> .	810
		811
		812
		813
	Ramanan Sekar, Oleh Rybkin, Kostas Daniilidis, Pieter Abbeel, Danijar Hafner, and Deepak Pathak. 2020. Planning to explore via self-supervised world models. In <i>ICML</i> , pages 8583–8592.	814
		815
		816
		817
	David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. 2017. Mastering the game of Go with- out human knowledge. <i>Nature</i> , 550(7676):354–359.	818
		819
		820
		821
		822
	Richard S Sutton and Andrew G Barto. 2018. <i>Reinforce- ment Learning: An Introduction</i> . MIT press.	823
		824
	Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. 1999. Policy gradient methods for reinforcement learning with function approxima- tion. In <i>NeurIPS</i> , pages 1057–1063.	825
		826
		827
		828

829	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	A Frequently Asked Questions (FAQs)	884
830	Martinet, Marie-Anne Lachaux, Timothée Lacroix,		
831	Baptiste Rozière, Naman Goyal, Eric Hambro,	A.1 Advantages of RAWM	885
832	Faisal Azhar, et al. 2023. LLaMA: Open and ef-	There are several advantages of RAWM, compared	886
833	ficient foundation language models. <i>arXiv preprint</i>	with other methods for LLM-based world models:	887
834	<i>arXiv:2302.13971</i> .	• RAWM does not require the fine-tuning of LLMs,	888
835	Karthik Valmeekam, Matthew Marquez, Alberto Olmo,	where the fine-tuning of LLMs is usually time	889
836	Sarath Sreedharan, and Subbarao Kambhampati.	and computation extensive. Besides, the fine-	890
837	2023. PlanBench: An extensible benchmark for eval-	tuning may also hurt the capabilities of LLMs	891
838	uating large language models on planning and reason-	on other tasks. RAWM can be viewed as a plug-	892
839	ing about change. In <i>NeurIPS</i> , pages 38975–38987.	and-play framework to transform the LLMs into	893
840	Ziyu Wan, Xidong Feng, Muning Wen, Stephen Marcus	world models.	894
841	McAleer, Ying Wen, Weinan Zhang, and Jun Wang.	• RAWM does not require the manually design of	895
842	2024. Alphazero-like tree-search can guide large	the prompts, i.e., instructions and in-context ex-	896
843	language model decoding and training. In <i>ICML</i> .	amples, for LLMs, which is usually labor inten-	897
844	Ruoyao Wang, Graham Todd, Ziang Xiao, Xingdi Yuan,	sive to optimize the prompts. RAWM automati-	898
845	Marc-Alexandre Côté, Peter Clark, and Peter Jansen.	cally retrieve the in-context examples from mem-	899
846	2024. Can language models serve as text-based	ory to assist the world models for predictions.	900
847	world simulators? <i>arXiv preprint arXiv:2406.06485</i> .	• RAWM introduces the efficient RL training to	901
848	Christopher JCH Watkins and Peter Dayan. 1992. Q-	further improve the world models with retrieval-	902
849	learning. <i>Machine learning</i> , 8(3):279–292.	augmented generation. We note that with the RL	903
850	Jiannan Xiang, Tianhua Tao, Yi Gu, Tianmin Shu, Zirui	training pipeline, RAWM can find the capability	904
851	Wang, Zichao Yang, and Zhiting Hu. 2023. Lan-	limit of the memory and the world model, thus	905
852	guage models meet world models: Embodied experi-	can be used to diagnose the systems.	906
853	ences enhance language models. <i>arXiv preprint</i>		
854	<i>arXiv:2305.10626</i> .	A.2 Why Focusing on Next State Prediction?	907
855	Kaige Xie, Ian Yang, John Gunerli, and Mark Riedl.	Next state prediction is the most important feature	908
856	2024. Making large language models into world mod-	for the world model (Wang et al., 2024). The re-	909
857	els with precondition and effect knowledge. <i>arXiv</i>	ward and the terminal can usually derived from	910
858	<i>preprint arXiv:2409.12278</i> .	the next state. For example, for Game24 and	911
859	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui,	BlocksWorld, we can derive the reward to check	912
860	Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu,	whether the remaining number is 24 and whether	913
861	Fei Huang, Haoran Wei, et al. 2024a. Qwen 2.5	the next state is the same as the goal state, respec-	914
862	technical report. <i>arXiv preprint arXiv:2412.15115</i> .	tively. Therefore, we focus on next state prediction.	915
863	Chang Yang, Xinrun Wang, Junzhe Jiang, Qinggang	A.3 Why Not Larger LLMs?	916
864	Zhang, and Xiao Huang. 2024b. Evaluating world	We note that Qwen/Qwen2.5-1.5B-Instruct is a	917
865	models with llm for decision making. <i>arXiv preprint</i>	highly capable LLM, which achieves 60.9% ac-	918
866	<i>arXiv:2411.08794</i> .	curacy on the MMLU benchmark. Therefore, we	919
867	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,	choose this small LLM as the base model for the	920
868	Thomas L Griffiths, Yuan Cao, and Karthik	RL training for the efficiency.	921
869	Narasimhan. 2023. Tree of thoughts: deliberate prob-	We also consider the models with sizes 3B and	922
870	lem solving with large language models. In <i>NeurIPS</i> ,	7B for inference, which achieve 65.6% and 72.4%	923
871	pages 11809–11822.	accuracy on MMLU benchmark, respectively.	924
872	Kenny Young, Aditya Ramesh, Louis Kirsch, and Jür-	A.4 What If RL Training Cannot Improve?	925
873	gen Schmidhuber. 2023. The benefits of model-based	RL training is a powerful framework. However,	926
874	generalization in reinforcement learning. In <i>ICML</i> ,	due to the trial-and-error process, RL training may	927
875	pages 40254–40276.	be more complicated than the supervised learning.	928
876	Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Ra-	Here we provide some guidance for the training:	929
877	jeswaran, Sergey Levine, and Chelsea Finn. 2021.	• Smaller learning rate with the identity initial-	930
878	COMBO: conservative offline model-based policy op-	ization would be safer for the better perfor-	931
879	timization. In <i>NeurIPS</i> , pages 28954–28967.		
880	Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon,		
881	James Zou, Sergey Levine, Chelsea Finn, and Tengyu		
882	Ma. 2020. MOPO: Model-based offline policy opti-		
883	mization. In <i>NeurIPS</i> , pages 14129–14142.		

mance than pre-trained models. While random initialization can potentially find better embedding models with longer training.

- We would also note that the improvement of RL training may also depend on the data in the memory and the LLMs for the world model. Therefore, if no good hyperparameters for the improvement, please consider larger LLMs and larger memory.

A.5 Code and Dataset Availability

We will release all the code and datasets upon the paper acceptance. The anonymous code can be access at: <https://anonymous.4open.science/r/rawm-acl>.

B Related Work

World Models in Decision Making. World models are actively explored by researchers to further improve the agent’s performance and the sample efficiency (Ha and Schmidhuber, 2018; Janner et al., 2019; Hafner et al., 2019; Schrittwieser et al., 2020). Dreamer (Hafner et al., 2019) is a practical model-based reinforcement learning algorithm that introduces the belief over states as a part of the input to the model-free DRL algorithm used. Trajectory Transformer (Janner et al., 2021) trains the transformer to predict the next state and action as a sequence modeling problem for continuous robot control. MuZero (Schrittwieser et al., 2020) is a remarkable success of model-based RL, which learns the world model and conducts the planning in the latent space. The world model with LLM in (Xiang et al., 2023) is trained to gain the environment knowledge, while maintaining other capabilities of the LLMs. Dynalang (Lin et al., 2024) proposes the multi-modal world model, which unifies videos and texts for the future prediction in decision making.

LLMs as World Simulators. World simulators are developed to model the dynamics of the world (Bruce et al., 2024). LLMs serve as the world simulators due to their generalizability across tasks. Specifically, The LLMs (i.e., GPT-3.5 and GPT-4) are evaluated to predict the state transitions, the game progress and scores with the given object, action, and score rules, where these rules are demonstrated to be crucial to the world model predictions (Wang et al., 2024). The world models with LLMs in (Xie et al., 2024) need to additionally identify the valid actions.

World Models in LLMs. The concept of world model also be explored in the deliberation reasoning of LLMs. Specifically, Reasoning via Planning (RAP) (Hao et al., 2023) leverages the planning methods (e.g., Monte Carlo Tree Search (MCTS)) with the world model with LLMs for plan generation and math reasoning, where LLMs need to predict the next state and the reward to guide the search. Tree of Thought (ToT) (Yao et al., 2023) implicitly leverages the LLMs as the world model to predict the next state and the reward for the search over different thoughts. Reason for future, act for now (RAFA) (Liu et al., 2023) combine the planning and reflection with the world model for complex reasoning tasks.

C Environments and Data Collection

C.1 Game24

Game24 is an interesting puzzle game, where four integer numbers $\{1, 2, 3, \dots, 13\}$ are given, the player needs to use the basic arithmetic operators, i.e., $+$, $-$, \times and \div , and use each number exactly at once to form 24. This puzzle game is used in (Yao et al., 2023) and (Liu et al., 2023) to benchmark the LLM’s reasoning capabilities.



Figure 6: Game24

The instances of Game24 used in this work can be accessed at <https://github.com/princeton-nlp/tree-of-thought-llm/blob/master/src/tot/data/24/24.csv>. The state of Game24 is the remaining numbers and the action is applying the operator between two remaining numbers. Here is an example of the transition:

```
{
  "state": (1.0, 1.0, 5.0, 8.0),
  "action": "1.0 + 1.0",
  "next_state": (2.0, 5.0, 8.0),
  "reward": False,
}
```

We provide the python-style code to transform the transitions to natural language examples in Algorithm 2.

Algorithm 2 Transitions to in-context examples for Game24

```

# transition is the dict with "state", "
  action", "next_state" and "reward"
def transition2example_game24(
    transition, is_query=False,
    is_next_state_prediction=True
):
    example = ""

    example += "current state: {}\n".
format(transition["state"])
    example += "action: {}\n".format(
    transition["action"])

    if not is_query:
        if is_next_state_prediction:
            example += "next state: {}\n"
.format(transition["next_state"])
        else:
            example += "reward: {}\n".
format(transition["reward"])

    return example

```

1020 **C.2 BlocksWorld**

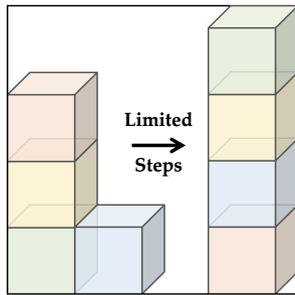


Figure 7: BlocksWorld

1021 BlocksWorld is a widely used benchmark
1022 to evaluate the planning capabilities of
1023 LLMs (Valmeekam et al., 2023; Hao et al.,
1024 2023). All the instances of the BlocksWorld can be
1025 accessed at [https://github.com/karthikv792/](https://github.com/karthikv792/LLMs-Planning/tree/main/plan-bench/instances/blocksworld)
1026 [LLMs-Planning/tree/main/plan-bench/](https://github.com/karthikv792/LLMs-Planning/tree/main/plan-bench/instances/blocksworld)
1027 [instances/blocksworld](https://github.com/karthikv792/LLMs-Planning/tree/main/plan-bench/instances/blocksworld). We build the envi-
1028 ronment by transforming the instances to MDPs,
1029 which can provide the transitions. Here is an
1030 example of the transition:

```

{
  "state": "the red block is clear,
the hand is empty, the orange block
is on top of the yellow block, the
red block is on top of the orange
block, the yellow block is on top of
the blue block, and the blue block
is on the table.",
  "action": "unstack the red block
from on top of the orange block",

```

```

"next_state": "the orange block is
clear, the red block is in the hand,
the hand is holding the red block,
the orange block is on top of the
yellow block, the yellow block is on
top of the blue block, and the blue
block is on the table.",
"reward": False,
"info": {
  "goal": "the red block is on top
of the blue block, the blue block
is on top of the yellow block and
the yellow block is on top of the
orange block"
},
}

```

We provide the python-style code to transform the transitions to natural language examples in Algorithm 3.

Algorithm 3 Transitions to in-context examples for BlocksWorld

```

# transition is the dict with "state", "
  action", "next_state" and "reward"
def transition2example_bw(transition,
    is_query=False,
    is_next_state_prediction=True):
    example = ""

    example += "goal state: {}\n".format(
    transition["info"]["goal"])
    example += "current state: {}\n".
format(transition["state"])
    example += "action: {}\n".format(
    transition["action"])

    if not is_query:
        if is_next_state_prediction:
            example += "next state: {}\n"
.format(transition["next_state"])
        else:
            example += "reward: {}\n".
format(transition["reward"])

    return example

```

1060 **C.3 BabyAI**

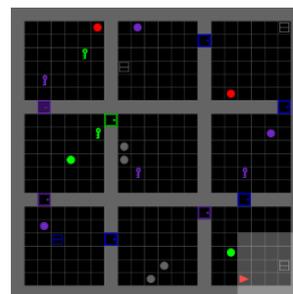


Figure 8: BabyAI

```

{

```

1061

```

1062 "mission": "go to a red box after
1063 you pick up the purple key",
1064 "state": [
1065     "You carry a purple key",
1066     "You see a wall 2 steps left",
1067     "You see a blue ball 2 steps
1068 forward",
1069     "You see a yellow ball 1 step
1070 right and 1 step forward",
1071     "You see a purple ball 2 steps
1072 right and 2 steps forward",
1073     "You see a red box 2 steps right
1074 and 1 step forward",
1075 ],
1076 "action": "turn right",
1077 "reward": 0,
1078 "done": False,
1079 "next_state": [
1080     "You carry a purple key",
1081     "You see a purple ball 2 steps
1082 left and 2 steps forward",
1083     "You see a blue ball 2 steps
1084 left",
1085     "You see a red box 1 step left
1086 and 2 steps forward",
1087     "You see a yellow ball 1 step
1088 left and 1 step forward",
1089     "You see a green key 4 steps
1090 forward",
1091     "You see a green key 1 step
1092 right",
1093     "You see a red box 2 steps right
1094 and 1 step forward",
1095     "You see a yellow key 3 steps
1096 right and 3 steps forward",
1097     "You see a red ball 3 steps
1098 right",
1099 ],
1100 }

```

1101 C.4 Statistics of Datasets

1102 Table 4 provides the statistics of the datasets used for the RL training and testing.

	Memory	Query	Test
Game24	2882	2882	5764
BlocksWorld	2416	2416	4833
BabyAI	3124	1562	3124

Table 4: Statistics of the datasets

1103

1104 D Prompts

1105 **Design of Prompts.** To make the world model
1106 as general as possible, we do not specifically de-
1107 sign the prompts. The system prompt of the world
1108 model is "After being given a current state
1109 and an action, directly give the next
1110 state after performing the action." We do
1111 not provide the description of the task, such as "I
1112 am playing with a set of blocks where I

Algorithm 4 Transitions to in-context examples for BabyAI

```

# transition is the dict with "state", "
# action", "next_state" and "reward"
def transition2example_babyai(
    transition, is_query=False,
    is_next_state_prediction=True
):
    def state_to_string(state):
        state_string = ""
        for idx, sta in enumerate(state)
        :
            state_string += sta
            if idx == len(state) - 1:
                continue
            else:
                state_string += ", "
        return state_string

    example = ""

    example += "mission: {}\n".format(
        transition["mission"])

    example += "current state: {}\n".
    format(state_to_string(transition["
    state"]))
    example += "action: {}\n".format(
        transition["action"])

    if not is_query:
        if is_next_state_prediction:
            example += "next state: {}\n"
            .format(
                state_to_string(
                    transition["next_state"])
                )
            else:
                example += "reward: {}\n".
                format(transition["reward"])

    return example

```

need to arrange the blocks into stacks.", 1113
 which is game specific and it needs human to write 1114
 the specific prompts. 1115

Content Prompt for LLMs. We present the tem- 1116
 plate for building the full prompt, i.e., the in- 1117
 context examples and the query, for the LLMs in 1118
 Algorithm 6. 1119

1120 E Model Selection

1121 E.1 World Models

1122 We expect to transform the LLMs into world mod-
 1123 els without any manually prompt engineering or
 1124 fine-tuning of LLMs. Therefore, the world models
 1125 are the general LLMs. The most capable open-
 1126 source LLM models are the Qwen-2.5-instruct se-
 1127 ries models (Yang et al., 2024a). Due to the limited

Algorithm 5 Prompt template

```

system_prompt = (
    "After being given a current state and
    an action, "
    "directly give the next state after
    performing the action."
)
message = [
    {
        "role": "system",
        "content": system_prompt,
    },
    {"role": "user", "content": prompt},
]

```

resources, we only consider the models with sizes in {1.5B, 3B, 7B} for inference and the 1.5B model for RL training. We note that RAWM can work for both open-source and close-source models.

For the embedding model, we choose the General Text Embedding (gte) family (Li et al., 2023). We choose Alibaba-NLP/gte-Qwen2-1.5B-instruct as the embedding model, which is the leading open-source model on MTEB.

Emb. Model ϕ	Alibaba-NLP/gte-Qwen2-1.5B-instruct Qwen/Qwen2.5-1.5B-Instruct-AWQ
World Model Ω	Qwen/Qwen2.5-3B-Instruct-AWQ Qwen/Qwen2.5-7B-Instruct-AWQ

Table 5: LLMs for Embedding and World Models

E.2 Architectures of MLP Head

Algorithm 7 presents the python implementation of the two types of initialization of the MLP. Table 6 displays the comparison of the two initializations.

	Random	Identity
Output dimension	Arbitrary	Same to ϕ
Initial performance	Low	High
Training instabilities	Low	High

Table 6: Comparison between two initialization

F Hyperparameters of RL Training

The hyperparameters of RL training are displayed in Table 7 and Table 8.

G Additional Experiment Results

The training curves for Game24 and BabyAI are shown in Figure 9 and Figure 10 respectively.

Algorithm 6 Generating prompts for LLMs

```

def get_query_examples_prompts(
    query_transitions,
    memory_transitions=None,
    exp_name=None,
):
    query_prompts = []
    for idx in range(len(
        query_transitions)):
        query_prompt =
            transition2example(
                query_transitions[idx],
                is_query=True, exp_name=exp_name
            )
        memory_prompt = ""
        if memory_transitions is not
            None:
            for memory_transition in
                reversed(memory_transitions[idx]):
                memory_prompt +=
                    transition2example(
                        memory_transition,
                        exp_name=exp_name
                    )
            query_memory_prompt =
                memory_prompt + query_prompt + "next
                state:"

        query_prompts.append(
            query_memory_prompt)

    return query_prompts

```

Hyperparameter	Value
norm_adv	True
clip_coef	0.2
entropy_coef	0.2
max_grad_norm	0.2
eps	1e-5

Table 7: Fixed Hyperparameters

Env	Method	Hyperparameter	Value
Game24	RAWM _{ψ,rand} ^{RL}	learning_rate	1e-4
		update_epochs	10
BlocksWorld	RAWM _{ψ,eye} ^{RL}	learning_rate	1e-5
		update_epochs	5
BabyAI	RAWM _{ψ,rand} ^{RL}	learning_rate	1e-4
		update_epochs	20
BabyAI	RAWM _{ψ,eye} ^{RL}	learning_rate	1e-5
		update_epochs	10
BabyAI	RAWM _{ψ,rand} ^{RL}	learning_rate	3e-6
		update_epochs	10
BabyAI	RAWM _{ψ,eye} ^{RL}	learning_rate	5e-5
		update_epochs	10

Table 8: Modified Hyperparameters

Algorithm 7 MLP initializations

```

# base_emb_dim: dimension of the pre-
#   trained embedding model, i.e., 1536
# final_emb_dim: dimension of the MLP,
#   36 for rand and 1536 for eye
def layer_init(layer, std=np.sqrt(2),
  bias_const=0.0, with_diag=False):
  if with_diag:
    torch.nn.init.eye_(layer.weight)
    torch.nn.init.constant_(layer.
  bias, 0.0)
  else:
    torch.nn.init.orthogonal_(layer.
  weight, std)
    torch.nn.init.constant_(layer.
  bias, bias_const)
  return layer

mlp_eye = nn.Sequential(
  layer_init(
    nn.Linear(
      base_emb_dim, final_emb_dim),
      with_diag=True
    ),
)
mlp_rand = nn.Sequential(
  layer_init(nn.Linear(
    base_emb_dim, 64)),
    nn.Tanh(),
    layer_init(nn.Linear(64,
    64)),
    nn.Tanh(),
    layer_init(nn.Linear(64,
    final_emb_dim), std=0.01),
)

```

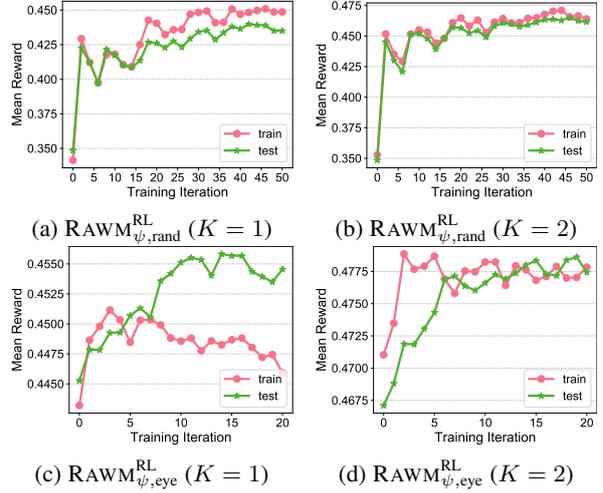


Figure 10: Training curves on BabyAI.

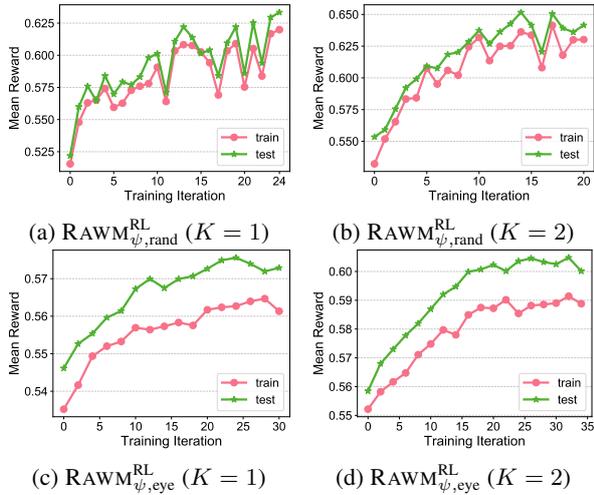


Figure 9: Training curves on Game24.