OPTIMIZING INFERENCE-TIME REASONING IN LLMS VIA RETRIEVAL-AUGMENTED REFLECTION

Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

023

025 026 027

028

Paper under double-blind review

ABSTRACT

Empowering LLMs to improve their performance through increased inferencetime computation is a crucial step in developing self-improving agents capable of operating in open-ended natural language contexts. In this paper, we explore how iteratively revising a chain of thoughts guided by information retrieval significantly improves large language models' reasoning ability in challenging tasks, while hugely mitigating hallucination. In particular, the proposed method — retrievalaugmented reflection (RaR) — revises the generation tokens step by step, leveraging multiple pieces of retrieved information relevant to the intermediate reasoning steps and the instruction. Applying RaR during inference-time to a various set of language models substantially improves their performances on various reasoning tasks; on relatively increasing scores by up to +16.4% on code generation, +11.6%on mathematical reasoning, and 29.1% on embodied task planning. Moreover, we find that with more inference-time computation given to the LLM for multi-times retrieval-augmented reflection, the LLM can continuously improve on various reasoning benchmarks. A small LM can surpass the performance of the LM with more than 10 times parameters, when giving more computation cost.

1 INTRODUCTION

029 Large Language Models (LLMs) have achieved fruitful progress on various natural language reasoning 031 tasks (Wei et al., 2022; Yao et al., 2022; Wang et al., 2023a; Zhou et al., 2023; Brown et al., 2020), espe-033 cially when combining large-scale models (Team, 034 2022; OpenAI, 2023) with sophisticated prompting strategies, notably chain-of-thought (CoT) prompting (Wei et al., 2022; Kojima et al., 2022). However, there have been increasing concerns about the factual 037 correctness of LLMs reasoning, citing the possible hallucinations in model responses (Rawte et al., 2023) or the intermediate reasoning paths, *i.e.* CoTs (Dhu-040 liawala et al., 2023). This issue becomes more sig-041 nificant when it comes to zero-shot CoT prompting, 042 aka. "let's think step-by-step" (Kojima et al., 2022) 043 and long-horizon generation tasks that require multi-044 step and context-aware reasoning, including code generation, task planning, mathematical reasoning, etc. Factually valid intermediate thoughts could be 046 critical to the successful completion of these tasks. 047

Several prompting techniques have been proposed to mitigate this issue, one promising direction, Retrieval Augmented Generation (RAG) (Lewis et al., 2020b)
seeks insights from human reasoning (Holyoak & Morrison, 2012), and utilizes retrieved information to facilitate more factually grounded reasoning. In this paper, we explore how to synergize RAG with



Figure 1: Performance vs. Inference Cost for various models on the ClassEval benchmark (Du et al., 2023). The size of each circle reflects the model parameters. The chart highlights the trade-off between inference-time computation and performance, demonstrating the effectiveness of different model architectures, with RaR-enhanced models generally achieving better performance for the same model scale. We also list the latest OpenAI O1 model (OpenAI, 2024), and our methods achieve better performance with less inference-time computation.

054 sophisticated long-horizon reasoning. Our intuition is that the hallucination within the intermediate 055 reasoning process could be alleviated through the help of outside knowledge. The resulting prompting 056 strategy, retrieval-augmented reflection (RaR), comprises two key ideas. Firstly, the initial zero-shot 057 CoT produced by LLMs along with the original task prompt will be used as queries to retrieve 058 relevant information that could help revise the possibly flawed CoT. Secondly, instead of retrieving and revising with the full CoT and producing the final response at once, we devise a progressive approach, where LLMs produce the response step-by-step following the CoT (a series of subtasks), 060 and only the current thought step will be revised based on the information retrieved with task prompt, 061 the current and the past CoTs. This strategy can be an analogy to the human reasoning process: 062 we utilize outside knowledge to adjust our step-by-step thinking during complex long-horizon 063 problem-solving (Holyoak & Morrison, 2012). 064

- We observe that hallucinations in model outputs often originate from errors in earlier tokens, which 065 can propagate and lead to incorrect final results. To tackle this issue, we extend retrieval-augmented 066 generation to retrieval-augmented reflection. In RaR, the language model retrieves external infor-067 mation to verify and revise the original generation. This modification allows the model to correct 068 errors in previously generated tokens—something traditional generation pipelines cannot achieve. 069 By adopting retrieval-augmented reflection, we enable the model to iteratively improve its reasoning process in real-time, making it highly scalable with increased inference-time computation. This 071 approach resolves the limitation of not being able to revise erroneous tokens generated earlier in the 072 process, thus providing more reliable outputs without requiring model parameter adjustments (Ope-073 nAI, 2024). Different from methods required reinforcement learning to perform better reasoning 074 during inference-time computation (OpenAI, 2024), our method does not require any modifications 075 to the original model parameters; it only needs to provide more inference-time computation tokens, allowing the LM to automatically verify and revise the generation. As shown in Figure 1, RaR 076 demonstrates better scalability performance. 077
- 078 We evaluate Retrieval-augmented Reflection (RaR) on a wide collection of challenging long-horizon 079 tasks, including code generation, mathematical reasoning, embodied task planning, and creative writing. We employ several LLMs of varied scales: GPT-3.5 (Brown et al., 2020), GPT-4 (OpenAI, 081 2023), Deepseek-Coder (Zhu et al., 2024), Llama-3 (AI@Meta, 2024) and Gemma (Team et al., 2024). The results indicate that combing RaR with these LLMs elicits strong advantages over vanilla CoT prompting and RAG approaches. In particular, we observe new state-of-the-art level 083 of performances across our selection of tasks when given the same maximum inference-time token 084 limitation: 1) code generation: ClassEval (+16.4%), HumanEval (+4.7%), HumanEval+ (+2.4%), 085 MBPP (+2.5%), MBPP+ (+4.2%); 2) mathematical reasoning problems: GSM8K (+11.6%), and GSMHard (+3.0%); 3) Minecraft task planning (+2.2% on accuracy); 4) QA (+16.44% on accuracy). 087 Our scaling experiments show that RaR can be scalable with more inference-time computation and 880 model parameters. By allowing for more inference-time tokens, RaR can achieve up to +33.3%, 089 +12.2%, and +16.2% relative improvements in code generation, math reasoning, and task planning benchmarks, respectively. Our additional ablation studies further confirm the crucial roles played 091 by the two key ingredients of RaR: revising CoT using RAG and progressive revision & generation. 092 This work reveals how can LLMs revise their reasoning process in a zero-shot fashion with the help of outside knowledge, just as what humans do. 093
- 094

2 RELATED WORKS

096 097

Retrieval-augmented Generation (RAG). Recently, RAG has gained popularity for boosting the performance of LLMs by guiding their generation process using the retrieved knowledge (Zhao et al., 2023). Without updating model parameters that may be expensive (Lewis et al., 2020a) or unstable (Ke et al., 2022b;a), RAG is a cost-effective way for LLMs to interact with the external world (Gu et al., 2018; Lewis et al., 2020a). RAG is widely applied to downstream tasks, such as code generation (Zhou et al., 2022b; Lu et al., 2022; Nashid et al., 2023), question answering (Baek et al., 2023; Siriwardhana et al., 2023), and creative writing (Wen et al., 2023; Asai et al., 2023).

Reasoning-enhanced RAG. Some recent works also leverage reasoning to enhance the performance of RAG (Li et al., 2023b). For example, IRCoT (Trivedi et al., 2022) exploits CoT to generate better queries for retrieval, IRGR (Ribeiro et al., 2022) performs iteratively retrieval to search for suitable premises for multi-hop QA, GEEK (Liu et al., 2023a) can choose to query external knowledge or

108 perform a single logical reasoning step in long-horizon generation tasks, and ITRG (Feng et al., 109 2023a) performs retrieval based on the last-step generation. Active RAG (Jiang et al., 2023) also 110 utilizes reasoning to enhance the quality of language model retrieval for better completion of QA 111 tasks. These previous RAG methods simply adopt a single query to retrieve the knowledge for 112 question-answering tasks (Gao et al., 2023; Feng et al., 2023b; Jiang et al., 2023; Yu et al., 2023). Our proposed RaR focuses on retrieval to refine LLM reasoning outputs and ensure consistent results. 113 The reasoning and retrieval processes are interlinked to improve both aspects, as shown in Figure 114 2. Our approach is evaluated across long-horizon content generation tasks such as code generation, 115 math reasoning, embodied planning, and creative writing. 116

117 Language Model for Reasoning. The advancement of reasoning in language models has seen 118 notable methodologies emerge since CoT was proposed by Wei et al. (2022), which showcased LMs' ability to generate self-derived problem-solving strategies. This foundational work spurred further 119 innovations such as the least-to-most prompting (Zhou et al., 2022a), zero-shot CoT (Kojima et al., 120 2022), self-consistency (Wang et al., 2022), zero-shot CoT without prompting (Wang & Zhou, 2024). 121 Moving beyond basic prompting, Creswell et al. (2022) introduced the Selection-Inference framework, 122 while Zelikman et al. (2022) developed STaR to refine reasoning through model finetuning. Creswell 123 & Shanahan (2022) proposed a faithful reasoning model, segmenting reasoning into dedicated steps, 124 similar to Scratchpad's approach by Nye et al. (2021) for enhancing multi-step computation. Tree-125 of-Thought (Yao et al., 2023) and Graph-of-Thought (Besta et al., 2023) also expand the reasoning 126 paths into a complex structure instead of linear CoT. These methods usually aim to improve the 127 reasoning ability of LLM by designing prompts or providing feedback from the environment to assist 128 in better planning and decision-making (Wang et al., 2023c; Yao et al., 2022; Shinn et al., 2023; Li 129 et al., 2023a; Zhang et al., 2023). However, RaR takes a different approach by using RAG to access external knowledge that can help LLM with its reasoning process. 130

131 132

133 134

3 RETRIEVAL AUGMENTED REFLECTION

Our goal is to support reasoning and generation while mitigating hallucination when using LLMs. 135 To have satisfying performance on long-horizon tasks, two ingredients are indispensable. Firstly, 136 access to factual information can be facilitated by retrieval. Secondly, appropriate intermediate steps 137 that outline a scratchpad to finish complex tasks, can be facilitated by CoT. Yet, a naive combination 138 of the two would not necessarily yield improvements. Three questions still persist: (1) what is 139 relevant information to retrieve; (2) how to effectively correct reasoning steps with relevant factual 140 information; (3) can this combination be scalable with more inference-time computation. To better 141 appreciate our method and why our method can address these questions, we first provide a brief 142 introduction of basic RaR and interactive RaR.

143 144

145

3.1 RETRIEVAL AUGMENTED REFLECTION

We first describe the pipeline of the Retrieval-augmented Reflection (RaR) method. The key steps of the RaR process are as follows:

First, given an input question or instruction x, we first generate an initial coarse response $y^{\text{raw}} \sim p_{\text{LM}}(\cdot \mid x)$ using the language model p_{LM} , where p_{LM} denotes the probabilistic output of the language model. To enhance the quality of initial response y_{raw} , we construct a retrieval query $q \sim p_{\text{LM}}(\cdot \mid x)$ by leveraging both the input question x and the initial coarse response y^{raw} . We then retrieve the relevant document V^k from a set of K candidate documents $V = \{V^1, V^2, \dots, V^K\}$. Generally, the retrieval process involves selecting the document that maximizes a similarity function $\sin(q, V^k)$, defined over the query embedding $\operatorname{emb}(q) \in \mathbb{R}^d$ and the documents embeddings $\operatorname{emb}(V^k) \in \mathbb{R}^d$:

155 156

157 158

$$V^{k} = \arg \max_{k \in \{1, \dots, K\}} \sin(q, V^{k}) = \arg \max_{k \in \{1, \dots, K\}} \frac{\operatorname{emb}(q) \cdot \operatorname{emb}(V^{k})}{\|\operatorname{emb}(q)\| \|\operatorname{emb}(V^{k})\|},$$
(1)

where d is the dimension of text embeddings (Reimers & Gurevych, 2019). Finally, the retrieved document V^k is then used to prompt the language model to reflect the initial response y^{raw} and output reflection $y^{\text{reflection}} \sim p_{\text{LM}}(\cdot \mid x, y^{\text{raw}}, V^k)$. The reflection process identifies potential errors or hallucinations in y^{raw} . Based on the reflection output $y^{\text{reflection}}$, the language model then refine the

162 Algorithm 1 Iterative Retrieval-augmented Reflection (RaR) 163 **Input:** Task Prompt x, Autoregressive Large Language Model p_{θ} , Number of Reasoning Steps n 164 1: $y^{\text{raw}}, \{y_i^{\text{thought}}\}_{i=1}^{i=n} \leftarrow p_{\theta}(\cdot|I) \triangleright$ Generate initial response y^{raw} and step-by-step reasoning y^{thought} 165 2: $y_0^{\text{RaR}} \leftarrow \text{None}$ ▷ Initialize the intermediate RaR response 166 3: for i = 1 to n do 167 $q_i^{\text{inter}} \leftarrow p_{\text{query}}(x, y_{i-1}^{\text{RaR}}, y_i^{\text{thought}})$ \triangleright Generate query for intermediate reasoning step *i* 4: 168 $V_i^{\text{inter}} \leftarrow \text{RetrieveFromCorpus}(q_i^{\text{inter}})$ 5: \triangleright Retrieve related document for step *i* 169 $y_i^{\text{RaR}} \leftarrow p_{\theta}(\cdot | I, T_{\leq i-1}^{\text{thought}}, K_i^{\text{inter}})$ 6: \triangleright Reflect and refine step *i* 170 7: end for 171 8: $y_{n+1}^{\text{RaR}} \leftarrow p_{\theta}(\cdot|I, y^{\text{RaR}}, T^{\text{raw}}, K_{n+1}^{\text{iter}})$ ▷ Initialize the overall RaR response 172 9: $j \leftarrow n+1$ Start Overall Response RaR 173 10: repeat $\begin{array}{l} q_{j}^{\text{iter}} \leftarrow p_{\text{query}}(x, y_{j}^{\text{RaR}}) \\ V_{j}^{\text{iter}} \leftarrow \text{RetrieveFromCorpus}(q_{j}^{\text{iter}}) \\ y_{j}^{\text{RaR}} \leftarrow p_{\theta}(\cdot | x, y_{j-1}^{\text{RaR}}, V_{j}^{\text{iter}}) \\ j \leftarrow j+1 \end{array}$ 174 11: \triangleright Generate query for full reflection step j175 12: Retrieve additional related documents 176 13: ▷ Iteratively refine full response 177 14: 178 15: until Convergence (e.g., m identical generations) or maximum token limit reached 179 \triangleright Output refined response y^{RaR} as the final generation 16: **return** *T**

initial raw coarse response, incorporating the corrective information from V^k :

$$\mathbf{RAG}: \quad y^{\mathbf{RAG}} \sim p_{\mathbf{LM}}(\cdot \mid x, V^k), \quad V^k = \arg \max_{k \in \{1, \dots, K\}} \sin(x, V^k),$$

$$\mathbf{RaR}: \quad y^{\mathbf{RaR}} \sim p_{\mathbf{LM}}(\cdot \mid x, y^{\mathrm{raw}}, V^k, y^{\mathrm{reflection}}).$$
(2)

This pipeline synergizes the Retrieval-augmented Generation and self-reflection based on the same 187 Language Model, which is different from RAG methods (Lewis et al., 2020b). To improve the 188 performance, we further increase the number of RaR iterations through iterative integration of 189 retrieval-augmented generation and self-reflective reasoning. 190

3.2 ITERACTIVE RETRIEVAL AUGMENTED REFLECTION

1

Z 1

193 To further scale up reasoning accuracy and overall response quality, we extend the basic RaR 194 framework to an iterative version, referred to as **Iterative RaR**. This approach enables the model to 195 repeatedly retrieve and reflect on intermediate reasoning steps and the overall response in multiple 196 iterations, progressively improving its performance. The workflow of Iterative RaR comprises two 197 key phases:

Reflection on Intermediate Reasoning Steps. In this phase, RaR focuses on improving the step-by-199 step reasoning process by reflecting on intermediate steps. We first enhance the instruction x with a 200 CoT prompt to form an augmented instruction x^* . The language model then generates a reasoning 201 process and the raw response, represented as: 202

$$y^{\text{thought}}, y^{\text{raw}}) \sim p_{\text{LM}}(\cdot \mid x^*),$$
(3)

203 where $y^{\text{thought}} = \{y_1^{\text{thought}}, \dots, y_J^{\text{thought}}\}$ denotes the step-by-step reasoning process, divided into J 204 sections corresponding to J reasoning steps. For each iteration i, the retrieval query q^i is constructed 205 based on the causal history of reasoning steps up to step i, while masking future steps j > i. The 206 query generation is expressed as: 207

$$q^i \sim p_{\text{LM}}(\cdot \mid x^*, \{y_j^{\text{thought}}\}_{j=1}^{j < =i}), i = 1, \dots, J.$$
 (4)

209 To improve efficiency, this retrieval based on intermediate reasoning steps is parallelized through 210 causal mask (Vaswani et al., 2017), enabling the generation of queries for different reasoning steps 211 simultaneously. And the retrieved documents $\{V_k^j\}_{j=1}^J$ are used to review and correct any potential 212 errors in the reasoning steps.

thought **r** - k reflection

213

181

183

185 186

191

192

208

2

213
214
215

$$y_{i}^{\text{RaR}} \sim \begin{cases} p_{\text{LM}}(\cdot \mid x, y_{i}^{\text{thought}}, V_{i}^{k}, y_{i}^{\text{reflection}}), & \text{if } i = 1, \\ p_{\text{LM}}(\cdot \mid x, y_{i-1}^{\text{RaR}}, y_{i}^{\text{thought}}, V_{i}^{k}, y_{i}^{\text{reflection}}), & \text{if } 1 < i < J, \\ p_{\text{LM}}(\cdot \mid x, y_{i-1}^{\text{RaR}}, y_{i}^{\text{thought}}, y_{i}^{\text{raw}}, V_{i}^{k}, y_{i}^{\text{reflection}}), & \text{if } i = J. \end{cases}$$
(5)

216 This allows for correcting errors in the original thoughts y^{thought} by continually consulting different 217 reference texts and ensuring that the most accurate and relevant information informs each step of 218 reasoning. Previous methods have demonstrated that in tasks involving long-term planning and 219 rigorous reasoning, like mathematical reasoning (Lightman et al., 2023) and embodied planning (Yao 220 et al., 2022; Shinn et al., 2023), supervision of intermediate processes is necessary to ensure the accuracy of model outputs. However, these approaches typically rely on feedback from humans or the 221 environment, which can be costly in situations where exploration and annotation expenses are high, 222 such as in safe decision-making scenarios (Gu et al., 2022). In contrast, RaR can automatically access 223 relevant information from external sources to validate and revise the content of model outputs through 224 a retrieval process. This allows RaR to autonomously verify each step without requiring human 225 labels (Lightman et al., 2023), which explains its significant success in mathematical reasoning. 226

Refinement of Overall Response. After refining the intermediate reasoning steps, RaR observes the full response y_J^{RaR} with refined intermediate steps, obtained from the corrected reasoning process. The full response is then used to produce query $q_i \sim p_{\text{LM}}(\cdot \mid x^*, y_J^{\text{RaR}})$ about overall structure to identify potential errors or inconsistencies in the overall structure, where i > J. The retrieved documents guide the refinement of the complete response, ensuring consistency and correctness across the entire output. The final response will be formulated as:

234

$$y_i^{\text{RaR}} \sim p_{\text{LM}}(\cdot \mid x, y_{i-1}^{\text{RaR}}, V_k^i, y_i^{\text{reflection}}), i > J.$$
(6)

We can repeat the RaR more times for better refinement of final response with more inference-time computation. The iterative process continues until the response reaches the maximum token limitation or the reflection process produces m consecutive identical generations, indicating convergence. In practice, m is set as 3.

239 Our hypothesis why our method can address the two problems mentioned at the beginning of this 240 section is as follows. Firstly, the most straightforward way to know what information will be used 241 in complex reasoning is to "see" the reasoning steps. Our approach leverages all the generated 242 thoughts along with the task prompt to provide more clues for more effective retrieval. Secondly, some information cannot be directly retrieved, especially information related to the final answer to a 243 hard complex question. Instead, retrieval of information relevant to intermediate questions, which are 244 assumed to be easier, is more accessible. Thanks to the compositional nature of many reasoning tasks, 245 an iterative retrieval process could also be more effective. Thirdly, correcting potential hallucinations 246 needs to be targeted. Revising a complete CoT with RAG could introduce errors at otherwise already-247 correct steps. Revising every step one by one could be more reliable. The first two points address 248 question (1) and the last point addresses question (2). Quantitative evidence can be found in our 249 ablation studies in Section 4.4. 250

251 252

3.3 Optimizing RAR by Scaling Inference-time Computation

In this section, we explore how to improve the reasoning performance of LLMs by scaling inference-time computation. While the native approach to scaling inference-time performance typically
 involves increasing the parameters of the LLM itself, we present an alternative method using Iterative
 Retrieval-augmented Reflection (RaR) that scales computation by augmenting input and output tokens
 iteratively.

Native Scaling via Model Size. The conventional approach to improving LLM performance, referred to as DIRECT Scaling, relies on increasing the number of model parameters, as suggested by the LLM scaling laws (Kaplan et al., 2020). Larger models generally exhibit better performance with higher computational costs during inference. However, this method comes with significant drawbacks, such as the increased burden of training large-scale models and the accompanying demand for vast amounts of data (Kaplan et al., 2020).

Inference-time Scaling without Increasing Model Parameters. Recent research shows that
 inference-time performance can be improved for LLMs with fixed parameters by increasing the
 computational cost through additional *input tokens* or *output tokens* during inference (Snell et al.,
 2024). Scaling input tokens involves extending the user prompt with more tokens. Few-shot CoT
 adds more demonstrations (instruction-response pair) to the prompt, e.g., increasing from 1-shot
 to 5-shot CoT, enables better contextual reasoning for complex tasks (Wei et al., 2022). RAG
 retrieves and adds more related documents to the user prompt. For instance, increasing the top-k

documents from 1-shot RAG to 5-shot RAG can provide richer contextual information for reasoning.
 Scaling output tokens involves generating more detailed outputs or sampling more generations. For
 example, Self-Consistency will generate multiple outputs in parallel and merge them into a final
 response (Wang et al., 2023a).

Iterative RaR simultaneously scales **both input and output tokens** by performing multiple rounds of retrieval and reflective reasoning. Specifically, for n iterations, RaR adds $n \cdot top-k$ related documents to the input tokens and generates more than n-times the output tokens during reasoning. This iterative process enables a gradual refinement of the reasoning process and response quality. As shown in our experiments, the performance of RaR improves consistently with increasing computational cost.

Scaling input and output tokens at inference time often encounters the long-context problem, wherein the fixed context length of LLMs restricts the number of tokens that can be processed (Brown et al., 2020). When the token count approaches the model's maximum limit, performance degradation can occur (Li et al., 2024). RaR mitigates this issue by iteratively performing retrieval and reflection, thereby avoiding a single, overly lengthy context and maintaining high performance.

It is worth noting that some approaches employ agent systems that utilize multiple prompts within a single LLM to create complex pipelines for answering questions (Yao et al., 2022; Gravitas, 2024).
However, such methods are compositions of foundational techniques and are not directly comparable to standalone reasoning and generation methods like RaR. For this reason, we exclude agent-based approaches from our discussions and experiments.

290 291

292

4 EXPERIMENTS

We evaluate our proposed method RaR on a variety of benchmarks that emphasize LLM generation and reasoning. Previous methods have struggled with these benchmarks, often producing "hallucinated" steps in LLM outputs that do not align with the original query or are clearly incorrect. For a detailed discussion, please refer to subsection 4.3 (case analysis). Due to space limitations, we do not present each benchmark setting or discuss our results extensively for each benchmark. Instead, this section offers a thorough showcase of our method's performance, shedding light on the preliminary empirical analysis of when our method succeeds and when it falls short.

299 300 301

4.1 EXPERIMENTAL SETUPS

302 Benchmarks. We adopt four groups of benchmarks including class-level code generation, math 303 reasoning, task planning, and question answering.¹ Code Generation benchmarks include ClassE-304 val (Du et al., 2023), HumanEval (Chen et al., 2021), HumanEval+ (Liu et al., 2023b), MBPP (Austin 305 et al., 2021), and MBPP+ (Liu et al., 2023b). These benchmarks encompass a wide range of pro-306 gramming problems, from simple function implementations to more complex class-level generation 307 challenges, providing a robust testbed for assessing generative and reasoning capabilities. Mathe-308 matical Reasoning evaluation is conducted on GSM8K and GSM-HARD dataset, which comprises thousands of multi-step mathematical problems (Cobbe et al., 2021; Gao et al., 2022). Task Planning 309 are evaluated on open-ended environments Minecraft. A set of 100+ tasks ranging from simple 310 objectives to challenging diamond objectives are evaluated through MC-TextWorld (Lin et al., 2023). 311 Question-answering (QA) tasks are assessed using TriviaQA (Joshi et al., 2017) benchmark, where 312 systems are required to answer various questions related to factual knowledge. 313

Evaluation Metrics. For fair evaluation, we report the best results achieved by all methods under the given maximum token limitation in Table 1, and the performance of different methods under different inference computation costs is reported in Figure 2. For code generation benchmarks, we use classical pass rate pass@5 as the evaluation metrics (Chen et al., 2021; Liu et al., 2023b). In ClassEval, we tested the generation pass rates for class-level and method-level (function) separately. We compute accuracy to evaluate every question in mathematical reasoning tasks, aligning with the established

 ¹We used bigcode-evaluation as the tool library for code evaluation. The pass@1 result of DIRECT in the table is slightly different from the result in the bigcode leaderboard, because we tested our pass@1 five times in our original setup and calculated the average value. We used the same settings as DIRECT in all experiments and reported on the relative improvement of RaR compared to baselines to promise fair evaluation and comparison.

Token Limitation	Method	ClassEval Function	ClassEval Class	HumanEval	HumanEval+	MBPP	MBPP+	GSM8K	GSMHard	Planning
	DIRECT (gpt-3.5)	14.4	23.8	72.5	70.5	72.9	64.0	65.9	51.3	19.3
	Input Scaling									
	RAG (n-shot)	24.3	53.1	76.2	70.5	70.5	68.0	61.8	56.8	33.0
	IRCoT	29.2	57.6	77.4	-	-	-	-	60.3	57.3
	Active RAG	25.7	55.3	75.8	-	-	-	-	61.3	59.4
	Output Scaling									
4K	Few-shot CoT	16.0	36.2	75.8	74.8	65.4	62.9	63.8	44.7	49.3
	Self-Refine	21.4	46.4	75.8	74.2	69.4	65.6	65.8	55.4	50.5
	Self-Consistency	19.4	40.4	75.8	72.9	73.8	69.7	65.0	52.8	50.2
	Both Scaling									
	RAG+CoT	24.7	55.7	76.8	69.3	-	-	67.5	67.3	50.4
	RaR (ours)	34.0	66.9	80.4	76.0	74.7	72.6	75.3	69.3	60.5
0V	Self-Refine	22.8	45.4	73.2	72.6	67.3	63.5	62.1	50.3	55.5
oK	RaR (ours)	37.3	66.9	81.3	79.4	76.9	74.9	75.8	69.3	76.7
16K	RaR (ours)	38.9	67.2	82.5	79.6	-	-	-	-	76.7

Table 1: The evaluation results cover various benchmarks such as code generation, math reasoning, and task
 planning. All evaluations were conducted using the GPT-3.5-turbo model from the OpenAI API, which has a
 maximum token length of 4096.

341 342

metric for the GSM8K (Cobbe et al., 2021). For embodied planning tasks, we compute the plan
execution success rate in MC-TextWorld as accuracy (Lin et al., 2023). For the QA benchmarks, we
use the official accuracy metric (Asai et al., 2023) to evaluate all methods. These indicators are better
the higher they are.

347 **Baselines.** To establish a comprehensive and equitable comparison landscape, we incorporate a suite of baseline methods. Our baselines, in addition to directly using LLM for sampling (DIRECT), 348 are divided into three groups: Scaling input tokens, which includes RAG (Lewis et al., 2020b), 349 IRCoT (Trivedi et al., 2022), Active-RAG (Jiang et al., 2023); scaling output tokens, which includes 350 CoT (Wei et al., 2022), Self-consistency (Wang et al., 2023a), Self-refine (Madaan et al., 2024); and 351 scaling both, including RAG+CoT and RaR. For specific scaling methods of each approach, please 352 refer to Section 3.3. For each method, we will scale them to stay within the permitted maximum 353 token limit, e.g. by adjusting the sample size n for self-consistency and the number of documents k in 354 RAG. For the QA benchmark, we also list the result from Self-RAG (Asai et al., 2023) in Figure 2 (a). 355 For different methods, the same language model is used as the base model. All methods in the Table 1 356 are evaluated with gpt-3.5-turbo. To ensure a fair comparison, none of the methods used examples 357 from the benchmark as demonstrations for in-context learning.

358 Model and RAG Settings. RaR leverages the capabilities of Retrieval-Augmented Generation 359 methods, which enhance the performance of language models by integrating external knowledge 360 sources. Specifically, we employed the codeparrot/github-jupyter dataset as our primary 361 search vector library for code generation and mathematical reasoning tasks. For embodied planning 362 tasks in Minecraft, we utilized the Minecraft Wiki² and DigMinecraft³ websites as the information 363 sources accessible to the LLMs. For OA benchmarks, we use the wiki pages as the retrieval library, which is consistent with Asai et al. (2023). We utilized OpenAI's text-embedding-ada-002 364 API service for all embedding calculations across different methods and base models. 365

Acknowledging the risk of benchmark contamination (an issue where the code library may contain solutions to the exact problems being evaluated), we adopted a rigorous pre-processing methodology as described by Guo et al. (2024). The potential implications of benchmark contamination, along with the effectiveness of our pre-processing strategy, are discussed in detail in Appendix D.

370 371 372

4.2 RESULTS

The code generation, math reasoning and task planning results are presented in Table 1 and QA results are presented in Figure 2 (a), which demonstrate the comprehensive evaluation of the RaR across multiple benchmarks. RaR consistently outperforms the other methods across the majority of

³⁷⁶ 377

²https://minecraft.wiki/

³https://www.digminecraft.com/



Figure 2: Evaluation results on different methods with (a) scaling model parameters, (b) scaling inference tokens, and (c) average computation cost. The language model in (b) is Deepseek-Coder 6.7B (Zhu et al., 2024). The language models in (c) are all OpenAI models including gpt-3.5 (Brown et al., 2020), gpt-4 (OpenAI, 2023), and openai o1 (OpenAI, 2024). The computation cost in (c) is computed with OpenAI Pricing.

394 the benchmarks and metrics, showcasing its superior ability in reasoning. For the most challenging 395 class-level generation tasks in ClassEval, RaR based on ChatGPT outperforms the base model with 396 more than 1.8 times improvements. For the method-level generation tasks in the HumanEval and 397 HumanEval+ benchmarks of code generation, RaR achieves remarkable improvements in pass@5 398 rates, indicating a significant enhancement in first-attempt accuracy and within the top five attempts. For example, on the HumanEval benchmark, RaR improves pass@5 by up to 25.68% relative to 399 the base models' performances. This trend is observed across different underlying base models, 400 highlighting RaR's effectiveness regardless of the initial model's capabilities. For mathematical 401 reasoning tasks, RaR demonstrates a significant relative improvement, with a 14.44% increase in 402 accuracy on GSM8K and a remarkable 35.27% relative increase on GSMHard, culminating in an 403 overall average improvement of 23% when deployed on the GPT-3.5 model. RaR significantly 404 outperforms all other methods on task planning tasks in Minecraft, achieving the highest scores with 405 76.67±8.02% for accuracy and 88.73% for partial accuracy, demonstrating its superior ability to 406 generate feasible and contextually appropriate plans in the complex open-world environment. RaR 407 has achieved state-of-the-art results on the QA benchmark. Specifically, the RaR method, based on 408 the 8B model, outperforms the larger 27B model in terms of direct output results despite having 409 fewer parameters as shown in Figure 2 (a). While other methods like self-RAG (Asai et al., 2023) and Active-RAG (Jiang et al., 2023) also incorporate reasoning methods and RAG during inference, 410 they fall short in performance compared to RaR. This highlights the effectiveness of the RaR method 411 in leveraging retrieved content for reflective reasoning. 412

The tasks are extremely diverse, while RaR can have consistent improvements over all baselines.
These results underline the advantages of RaR's approach, which leverages iterative refinement of
retrieval queries based on evolving reasoning thoughts. This strategy not only enhances the relevance
and quality of the information retrieved but also significantly improves the accuracy and efficiency of
the generated context.

418

389

390

391

392 393

419 4.3 SCALING EXPERIMENTS

420

We explored the performance scaling of the RaR technique from three key perspectives: (1) increasing
the base model parameters, (2) increasing the inference-time tokens, and (3) increasing the inferencetime computation cost (influenced by both token number and model parameters) which correspond to
more iterations of Iterative Retrieval-Augmented Reflection during inference. These experiments
were conducted on the Trivial QA Benchmark (Joshi et al., 2017) and the ClassEval Code Generation
Benchmark (Du et al., 2023). The results are demonstrated in Figure 2.

427 Scaling Model Parameters. In the QA benchmark, we utilized language models with varying parameter sizes—Gemma-2-2B, Llama-3-8B, Gemma-2-27B, and Llama-3-70B—as base models (Team et al., 2024; AI@Meta, 2024). Our findings demonstrate that RaR consistently enhances performance as the model scale increases. This trend was observed across all base model sizes, with RaR maintaining a significant performance advantage over both the Gemma and Llama models in the QA task. Under the same model parameters, RaR always performs the best. Importantly, the



Figure 3: Case analysis on long-horizon task planning and question answering. RaR improves upon CoT's initial answers by continuously refining thoughts with intermediate step retrieval and overall retrieval, aligning closely with task progression and relevant item knowledge.

more reasoning-intensive the task, the more pronounced the performance gains achieved by RaR,
 underscoring its ability to improve reasoning capabilities, which are typically more pronounced in
 larger models.

460 Scaling Inference Tokens. Additionally, we investigated the impact of stricter inference-time tokens 461 on the performance of different methods. We utilize the DeepSeek-Coder-6.7B (Zhu et al., 2024) on 462 ClassEval method-level pass rates (Du et al., 2023) to analyze all methods and present the performance 463 results along with the corresponding tokens in Figure 2 (b). Almost all methods show growth with the increase in inference-time tokens used. RaR performs worse than methods like IRCoT when 464 the number of tokens used is less than 2k, because RaR has not yet completed modifications to 465 all intermediate steps. However, when given more tokens up to 4k, RaR shows significant growth 466 compared to other methods. Additionally, we found that methods like self-consistency and RAG 467 experienced a decline in performance during the later stages of token growth. We speculate that this 468 is due to these methods requiring a large number of tokens at once, leading to long text issues that 469 cause performance degradation. In contrast, RaR, by using an iterative updating approach, does not 470 encounter long text problems, thus showing continuous performance growth when given more tokens. 471 This demonstrates that RaR has good scaling potential. 472

Scaling Computation Cost. Finally, we will also combine the model scale and the number of tokens 473 used to examine the performance of different methods under computation cost scaling. The horizontal 474 axis in Figure 2 (c) represents the average API price used, where all models utilize OpenAI models, 475 including gpt-3.5-turbo, gpt-4, and openai-o1. All pricing is based on OpenAI's official pricing⁴, 476 calculated according to the actual consumption of prompt tokens and generation tokens. We found 477 that RaR is the most economical choice, as it is nearly the best under the same computation cost. 478 Additionally, we discovered that with less computation cost, the RaR based on gpt-4 demonstrated 479 performance exceeding that of openai o1. 480

481

482

485

453

454

455 456

4.4 ABLATION STUDY

Ablation on retrieval in RaR. In this ablation study, we investigate the influence of various retrieval strategies on the efficacy of RaR, focusing on the optimization of content retrieval for

⁴https://openai.com/api/pricing/

Table 2: Comparative Impact of Retrieval Strategies onRaR Performance.

Table 3:	Ablation	Study	on	Causal	vs.	Non-Causal
Reasonin	g in RaR.					

Method	HumanEval Huma		inEval+ Method		HumanEval		HumanEval+		
moulou	pass@1(Δ) \uparrow	pass@5(Δ) \uparrow	pass@1(Δ) \uparrow	pass@5(Δ) \uparrow	method	$pass@1(\Delta)\uparrow$	pass@5(Δ) \uparrow	$pass@1(\Delta)\uparrow$	pass@5(Δ) \uparrow
Baseline	50.6%	76.2%	48.2%	70.5%	Baseline	47.3%	75.8%	41.7%	74.8%
CoT+RAG	53.9(+3.3)%	76.8(<mark>+0.6</mark>)%	51.3(+3.1)%	69.3(-1.2)%	Non-Causal	57.3(+10.0)%	78.0(+2.1)%	54.9(+13.2)%	74.8(+0.0)%
RaR	59.2(+8.7)%	80.4(+7.9)%	56.3(+8.2)%	76.0(+5.5)%	Causal	59.2(+11.9)%	80.4(<mark>+4.6</mark>)%	56.3(+14.6)%	76.0(<mark>+1.2</mark>)%

493 improving generative outputs. The experimental results, detailed in Table 2, highlight the significant 494 advancements achieved through the iterative refinement of retrieval queries in RaR compared to 495 baseline methods. The baseline denoted as RAG-1, employs a direct approach by using the question 496 itself as the retrieval query. In contrast, CoT+RAG enhances this process by utilizing the entirety of the reasoning thoughts output by the language model as the query, aiming for a broader contextual 497 understanding. However, RaR introduces a more dynamic method by employing continuously 498 modified parts of reasoning thoughts as queries, which allows for a more focused and relevant 499 information retrieval process. The comparative analysis shows that RaR surpasses both the baseline 500 and the CoT+RAG method in terms of pass@1 and pass@5 metrics across the HumanEval and 501 HumanEval+ benchmarks. Specifically, RaR demonstrates an 8.7 percentage point increase in pass@1 502 and a 7.9 percentage point increase in pass@5 over the baseline in the HumanEval benchmark, and 503 similarly impressive gains in the HumanEval+ benchmark. These improvements underscore the 504 effectiveness of RaR's retrieval strategy, which by iteratively refining next queries based on evolving 505 reasoning thoughts and previous queries, ensures the retrieval of highly pertinent information. This 506 process not only enhances the relevance of the information retrieved but also significantly improves 507 the quality and accuracy of the final generated outputs. The results firmly establish the superiority of RaR's dynamic retrieval method in leveraging contextual nuances to drive more precise and effective 508 generative processes. 509

510 Ablation on causal reasoning in RaR. In this ablation study, we systematically examine the impact 511 of causal and non-causal reasoning approaches on the performance of the RaR system, with the 512 Chain of Thought (CoT) serving as our baseline. Our findings, as summarized in Table 3, reveal 513 significant enhancements in generation capabilities when incorporating causal reasoning techniques. Specifically, the causal approach, which iteratively performs reasoning and retrieval, leads to notable 514 improvements in both pass@1 and pass@5 metrics across HumanEval and HumanEval+ benchmarks. 515 For instance, the causal method outperforms the baseline (CoT) by 11.9 percentage points in pass@1 516 and by 4.6 percentage points in pass@5 on the HumanEval dataset. This approach contrasts with 517 the non-causal method, which, although also surpassing the baseline, leverages the initial reasoning 518 thought to directly retrieve all necessary steps and generate the final answer. The causal method's 519 superior performance underscores the value of sequential reasoning and information retrieval in 520 enhancing the accuracy and reliability of generated outputs. This iterative process likely aids in 521 refining the search and reasoning steps based on continuously updated context, allowing for more 522 precise and relevant information retrieval, which in turn supports more accurate final answers. These 523 results firmly establish the efficacy of causal reasoning in long-horizon problem-solving tasks.

524 525

5 DISCUSSION ON LIMITATIONS AND CONCLUSION

526 527

> 528 One limitation of this work is that the performance of RaR relies on the chain-of-thought reasoning 529 and in-context learning (or RAG) capability of the base LLM. Since this work does not involve 530 any model training, the capability of base LLM will not change when applying RaR. Despite RaR 531 achieves significant improvement on powerful LLMs such as GPT-3.5 and GPT-4, the effect on 532 smaller and weaker LLMs is questionable. Another limitation of this work is that the performance of 533 RaR also relies on the quality of the retrieved knowledge. Another limitation of this work is that the 534 performance of RaR also relies on the quality of the retrieved knowledge.

> We have presented Retrieval Augmented Reflection (RAR), a simple yet effective prompting strategy that synergies chain of thought (CoT) prompting and retrieval augmented generation (RAG) to address the challenging long-horizon reasoning and generation tasks. Our key ideas involve revising the zero-shot chain of thoughts produced by LLMs through RAG with the thoughts as queries, and causally revising the thoughts & generating the response progressively. RaR, a **zero-shot** prompting approach, has demonstrated significant advantages over vanilla CoT prompting, RAG, and other

baselines on challenging code generation, mathematics reasoning, embodied task planning, and
 creative writing tasks.

543 544 REFERENCES

547

570

576

- AI@Meta. Llama 3 model card, 2024. URL https://github.com/meta-llama/llama3/
 blob/main/MODEL_CARD.md.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*, 2023.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,
 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language
 models. *arXiv preprint arXiv:2108.07732*, 2021.
- Jinheon Baek, Alham Fikri Aji, and Amir Saffari. Knowledge-augmented language model prompting for zero-shot knowledge graph question answering. *arXiv preprint arXiv:2306.04136*, 2023.
- Bowen Baker, Ilge Akkaya, Peter Zhokhov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon
 Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (vpt): Learning to act by watching
 unlabeled online videos. *arXiv preprint arXiv:2206.11795*, 2022.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. *arXiv preprint arXiv:2308.09687*, 2023.
- 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.
- Shaofei Cai, Zihao Wang, Xiaojian Ma, Anji Liu, and Yitao Liang. Open-world multi-task control through goal-aware representation learning and adaptive horizon prediction. 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 13734–13744, 2023a.
- Shaofei Cai, Bowei Zhang, Zihao Wang, Xiaojian Ma, Anji Liu, and Yitao Liang. Groot: Learning to follow instructions by watching gameplay videos. *arXiv preprint arXiv:2310.08235*, 2023b.
- 573 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared
 574 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
 575 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Antonia Creswell and Murray Shanahan. Faithful reasoning using large language models. *arXiv preprint arXiv:2208.14271*, 2022.
- Antonia Creswell, Murray Shanahan, and Irina Higgins. Selection-inference: Exploiting large language models for interpretable logical reasoning. *arXiv preprint arXiv:2205.09712*, 2022.
- Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and
 Jason Weston. Chain-of-verification reduces hallucination in large language models. *arXiv preprint arXiv:* 2309.11495, 2023.
- Xueying Du, Mingwei Liu, Kaixin Wang, Hanlin Wang, Junwei Liu, Yixuan Chen, Jiayi Feng, Chaofeng Sha, Xin Peng, and Yiling Lou. Classeval: A manually-crafted benchmark for evaluating llms on class-level code generation. *arXiv preprint arXiv:2308.01861*, 2023.
- ⁵⁹³ Zhangyin Feng, Xiaocheng Feng, Dezhi Zhao, Maojin Yang, and Bing Qin. Retrieval-generation synergy augmented large language models. *ArXiv*, abs/2310.05149, 2023a.

594 Zhangyin Feng, Xiaocheng Feng, Dezhi Zhao, Maojin Yang, and Bing Qin. Retrieval-generation 595 synergy augmented large language models. arXiv preprint arXiv:2310.05149, 2023b. 596 Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and 597 Graham Neubig. Pal: Program-aided language models. arXiv preprint arXiv:2211.10435, 2022. 598 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and 600 Haofen Wang. Retrieval-augmented generation for large language models: A survey. arXiv 601 preprint arXiv:2312.10997, 2023. 602 Significant Gravitas. Autogpt. https://github.com/Significant-Gravitas/ 603 AutoGPT, 2024. GitHub repository. 604 605 Jiatao Gu, Yong Wang, Kyunghyun Cho, and Victor OK Li. Search engine guided neural machine 606 translation. In Proceedings of the AAAI Conference on Artificial Intelligence, 2018. 607 608 Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. arXiv 609 preprint arXiv:2205.10330, 2022. 610 611 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao 612 Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. Deepseek-coder: When 613 the large language model meets programming – the rise of code intelligence. arXiv preprint 614 arXiv:2401.14196, 2024. 615 Ralf Herbrich, Tom Minka, and Thore Graepel. TrueskillTM: a bayesian skill rating system. Advances 616 in neural information processing systems, 19, 2006. 617 618 Keith J Holyoak and Robert G Morrison. The Oxford handbook of thinking and reasoning. Oxford 619 University Press, 2012. 620 Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot 621 planners: Extracting actionable knowledge for embodied agents. ICML, 2022. 622 623 Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, 624 Jamie Callan, and Graham Neubig. Active retrieval augmented generation. Conference on 625 Empirical Methods in Natural Language Processing, 2023. doi: 10.48550/arXiv.2305.06983. 626 Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly 627 supervised challenge dataset for reading comprehension. arXiv preprint arXiv:1705.03551, 2017. 628 629 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott 630 Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. 631 *arXiv preprint arXiv:2001.08361*, 2020. 632 Zixuan Ke, Haowei Lin, Yijia Shao, Hu Xu, Lei Shu, and Bing Liu. Continual training of language 633 models for few-shot learning. arXiv preprint arXiv:2210.05549, 2022a. 634 635 Zixuan Ke, Yijia Shao, Haowei Lin, Hu Xu, Lei Shu, and Bing Liu. Adapting a language model 636 while preserving its general knowledge. In Proceedings of the 2022 Conference on Empirical 637 Methods in Natural Language Processing, pp. 10177–10188, 2022b. 638 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large 639 language models are zero-shot reasoners. Advances in neural information processing systems, 35: 640 22199-22213, 2022. 641 642 Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast inference from transformers via speculative 643 decoding. arXiv preprint arXiv: 2211.17192, 2022. 644 645 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented genera-646 tion for knowledge-intensive nlp tasks. Advances in Neural Information Processing Systems, 33: 647 9459-9474, 2020a.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in Neural Information Processing Systems</i> , 33: 9459–9474, 2020b.
Chengshu Li, Jacky Liang, Andy Zeng, Xinyun Chen, Karol Hausman, Dorsa Sadigh, Sergey Levine, Li Fei-Fei, Fei Xia, and Brian Ichter. Chain of code: Reasoning with a language model-augmented code emulator, 2023a.
Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhu Chen. Long-context llms struggle with long in-context learning. <i>arXiv preprint arXiv:2404.02060</i> , 2024.
Xingxuan Li, Ruochen Zhao, Yew Ken Chia, Bosheng Ding, Shafiq Joty, Soujanya Poria, and Lidong Bing. Chain-of-knowledge: Grounding large language models via dynamic knowledge adapting over heterogeneous sources. In <i>The Twelfth International Conference on Learning Representations</i> , 2023b.
Shalev Lifshitz, Keiran Paster, Harris Chan, Jimmy Ba, and Sheila McIlraith. Steve-1: A generative model for text-to-behavior in minecraft. <i>arXiv preprint arXiv:2306.00937</i> , 2023.
Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. <i>PREPRINT</i> , 2023.
Haowei Lin, Zihao Wang, Jianzhu Ma, and Yitao Liang. Mcu: A task-centric framework for open-ended agent evaluation in minecraft. <i>arXiv preprint arXiv:2310.08367</i> , 2023.
Chang Liu, Xiaoguang Li, Lifeng Shang, Xin Jiang, Qun Liu, Edmund Y. Lam, and Ngai Wong. Gradually excavating external knowledge for implicit complex question answering. In <i>Conference</i> on Empirical Methods in Natural Language Processing, 2023a.
Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatGPT really correct? rigorous evaluation of large language models for code generation. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023b.
Shuai Lu, Nan Duan, Hojae Han, Daya Guo, Seung-won Hwang, and Alexey Svyatkovskiy. Reacc: A retrieval-augmented code completion framework. In <i>Proceedings of the 60th Annual Meeting of</i> <i>the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pp. 6227–6240, 2022.
Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
Noor Nashid, Mifta Sintaha, and Ali Mesbah. Retrieval-based prompt selection for code-related few-shot learning. In <i>Proceedings of the 45th International Conference on Software Engineering (ICSE'23)</i> , 2023.
Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. Show your work: Scratchpads for intermediate computation with language models. <i>arXiv preprint arXiv:2112.00114</i> , 2021.
OpenAI. Gpt-4 technical report, 2023.
OpenAI. Introducing openai o1 preview. https://openai.com/index/ introducing-openai-o1-preview/, 2024. Accessed: 2024-09-25.
Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. <i>arXiv preprint arXiv:2203.02155</i> , 2022.

701 Vipula Rawte, Amit Sheth, and Amitava Das. A survey of hallucination in large foundation models. *arXiv preprint arXiv:2309.05922*, 2023.

702 703	Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. arXiv preprint arXiv:1908.10084, 2019.
704 705	Danilo Ribeiro, Shen Wang, Xiaofei Ma, Rui Dong, Xiaokai Wei, Henry Zhu, Xinchi Chen, Zhiheng
706	Huang, Peng Xu, Andrew Arnold, et al. Entailment tree explanations via iterative retrieval- generation reasoner arXiv preprint arXiv:2205.09224, 2022
707	generation reasoner. arxiv preprint arxiv:2205.09221, 2022.
700	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Tan, Yossi Adi,
709	Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, I. Evtimov, Joanna Bitton, Manish P
711	Bhatt, Cristian Cantón Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre D'efossez, Jade Copet Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel
712	Synnaeve. Code llama: Open foundation models for code. <i>ArXiv</i> , abs/2308.12950, 2023.
713	Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: an autonomous agent with dynamic memory and self-reflection arXiv preprint arXiv:2303 11366 2023
715	
716 717	Shamane Siriwardhana, Rivindu Weerasekera, Elliott Wen, Tharindu Kaluarachchi, Rajib Rana, and Suranga Nanayakkara. Improving the domain adaptation of retrieval augmented generation (rag)
718	models for open domain question answering. Transactions of the Association for Computational
719	<i>Linguistics</i> , 11:1–17, 2023.
720	Charlie Shall Jackson Lee, Kelvin Vu, and Aviral Kumar, Seeling Ilm test time compute entimely
721 722	can be more effective than scaling model parameters. <i>arXiv preprint arXiv:2408.03314</i> , 2024.
723	Gemma Team Morgane Riviere, Shreya Pathak, Pier Giusenne Sessa, Cassidy Hardin, Surva
724	Bhupatiraiu, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al.
725	Gemma 2: Improving open language models at a practical size. <i>arXiv preprint arXiv:2408.00118</i> ,
726	2024.
727	
728 729	Google PaLM Team. Palm: Scaling language modeling with pathways. <i>arXiv preprint arXiv:</i> 2204.02311, 2022.
730	H Trivadi Nironian Balacubramanian Tushar Khat and Ashish Sabharwal Interlagging ratriaval with
731 732	chain-of-thought reasoning for knowledge-intensive multi-step questions. ArXiv, abs/2212.10509,
733	2022.
734 735	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. <i>arXiv preprint arXiv: 1706.03762</i> , 2017.
736	$\mathbf{Y} = 1^{T} \mathbf{W}$
737 738	<i>arXiv:2402.10200</i> , 2024.
739	Xuezhi Wang Jason Wei Dale Schuurmans Quoc Le Ed Chi Sharan Narang Aakanksha Chowdh-
740	ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. arXiv preprint arXiv:2203.11171, 2022.
749	· · · · · · · · · · · · · · · · · · ·
742	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha
743	Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
745	models. In The Eleventh International Conference on Learning Representations, ICLR 2023,
746	2023a.
747	Zihao Wang, Shaofei Cai, Anii Liu, Yonggang Jin, Jinbing Hou, Bowei Zhang, Haowei Lin, Zhaofeng
748	He, Zilong Zheng, Yaodong Yang, Xiaojian Ma, and Yitao Liang. Jarvis-1: Open-world multi-task
749	agents with memory-augmented multimodal language models. ArXiv, abs/2311.05997, 2023b.
750	Zihao Wang Chaofai Cai Anii Liu Vicciian Ma and Vitao Liang Describe surlain alar a data ta
751	Linao wang, Shaolei Cai, Anji Liu, Alaojian Ma, and Yitao Liang. Describe, explain, plan and select:
752	niciacity pranning with large language models enables open-world multi-task agents. <i>arXiv</i>
753	preprint urxiv.2302.01300, 2023C.
754 755	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. <i>36th Conference on Neural</i> <i>Information Processing Systems (NeurIPS 2022)</i> , 2022.

756	Zhihua Wan, Zhiliang Tian, Wei Wu, Yuyin Yang, Yangi Shi, Zhan Huang, and Dongshang Li				
757 758	Grove: a retrieval-augmented complex story generation framework with a forest of evidence. <i>arXiv</i> preprint arXiv:2310.05388, 2023.				
759					
760	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.				
761	React: Synergizing reasoning and acting in language models. <i>arXiv preprint arXiv:2210.03629</i> ,				
762	2022.				
763	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik				
764	Narasimhan. Tree of thoughts: Deliberate problem solving with large language models, 2023.				
765	$\mathbf{W} = \mathbf{I} + \mathbf{V} + \mathbf{Z} \mathbf{I} + \mathbf{Z} $				
766 767	models via plug-and-play retrieval feedback. <i>arXiv preprint arXiv:2305.14002</i> , 2023.				
768	Haogi Yuan Chi Zhang Hongcheng Wang Feiyang Xie Penglin Cai Hao Dong and Zongging				
769	Lu. Plan4mc: Skill reinforcement learning and planning for open-world minecraft tasks.				
770	preprint arXiv:2303.16563, 2023.				
771					
772	Haoqi Yuan, Zhancun Mu, Feiyang Xie, and Zongqing Lu. Pre-training goal-based models for				
773	sample-efficient reinforcement learning. In The Twelfth International Conference on Learning				
774	<i>Representations</i> , 2024.				
775	Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with				
776	reasoning. Advances in Neural Information Processing Systems, 35:15476–15488, 2022.				
777					
778	Ceyao Zhang, Kaijie Yang, Siyi Hu, Zihao Wang, Guanghe Li, Yihang Sun, Cheng Zhang, Zhaowei Zhang, Anii Liu, Sang Chun Zhu, et al. Programti Puilding magazing accompanying a with large				
779	Zhang, Anji Liu, Song-Chun Zhu, et al. Proagent: Building proactive cooperative ai with language models. arXiv preprint arXiv:2308.11330, 2023				
780	ianguage mouers. <i>arxiv preprini arxiv.250</i> 0.11539, 2025.				
781	Ruochen Zhao, Hailin Chen, Weishi Wang, Fangkai Jiao, Xuan Long Do, Chengwei Qin, Bosheng				
782	Ding, Xiaobao Guo, Minzhi Li, Xingxuan Li, and Shafiq R. Joty. Retrieving multimodal infor				
783	tion for augmented generation: A survey. ArXiv, abs/2303.10868, 2023.				
784	Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans,				
785	Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-most prompting enables complex reasoning				
786	in large language models. arXiv preprint arXiv:2205.10625, 2022a.				
787	Danny Zhau Nathangal Sahörli La Hau Jasan Wai Nathan Saalas Yuazhi Wang Dala Sahuurmana				
788	Claire Cui Olivier Bousquet Ouoc V Le and Ed H Chi Least-to-most prompting enables				
789	complex reasoning in large language models. In <i>The Eleventh International Conference on</i>				
790	Learning Representations, ICLR 2023, 2023.				
791					
792	Shuyan Zhou, Uri Alon, Frank F Xu, Zhengbao Jiang, and Graham Neubig. Docprompting: Gen-				
793	erating code by retrieving the docs. In <i>The Eleventh International Conference on Learning</i>				
794	Kepresentations, 20220.				
706	Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li,				
790	Huazuo Gao, Shirong Ma, et al. Deepseek-coder-v2: Breaking the barrier of closed-source models				
797	in code intelligence. arXiv preprint arXiv:2406.11931, 2024.				
790					
800					
801					
802					
803					
804					
805					
806					
807					

A TASK DETAILS

818

819

820

821

822

823

824

825

827

828

829

A.1 CODE GENERATION 813

Benchmarks. We select HumanEval (Chen et al., 2021), HumanEval+ (Liu et al., 2023b),
MBPP (Austin et al., 2021), and MBPP+ (Liu et al., 2023b) as the code generation evaluation
benchmark. These benchmarks are commonly used to test the performance of code generation
models, which are briefly introduced below:

- **HumanEval** consists of 164 Python programming problems, each with a function signature, docstring, body, and multiple unit tests (Chen et al., 2021).
- **HumanEval+** includes the same programming problems as HumanEval, but with an additional 80 times more unit tests for each of the 164 problems (Liu et al., 2023b).
- **MBPP** is a collection of approximately 1,000 Python programming problems that are intended to be solvable by beginner programmers. Each problem includes an English task description, a code solution, and three automated test cases. We assess the sample test set from index 11 to 175 (Austin et al., 2021).
- MBPP+ consists of 399 tasks (Liu et al., 2023b), which are a subset of the original MBPP dataset. Additionally, MBPP+ includes extra unit tests for each of the 399 problems (35 times more than the original MBPP). We utilized the first 164 questions as our test set.

These benchmarks encompass a wide range of programming problems, from simple function implementations to more complex algorithmic challenges, providing a robust testbed for assessing the generative capabilities of various models.

834 Metrics. We adopt the pass@k metric for evaluating the efficacy of various code generation 835 algorithms, following the methodology proposed by Chen et al. (2021) and extended by Liu et al. (2023b). This metric quantifies the rate at which generated code snippets successfully execute and 836 pass all test cases, where k represents the number of attempts or samples generated by the model 837 for each problem. This approach allows us to rigorously assess the precision and reliability of 838 code generation models in producing functionally correct code across a diverse set of programming 839 challenges.

Baselines. To establish a comprehensive and equitable comparison landscape, we incorporate a suite 841 of baseline methods and diverse code generation models. Our baselines include the original code 842 generation language models, referred to as DIRECT, and the Retrieval-Augmented Generation (RAG) 843 methodology with n retrieved examples, instantiated in both single-shot (1 shot) and multi-shot (5 844 shots) configurations, as documented by Lewis et al. (2020b). Additionally, we examine the zero-shot 845 CoT (CoT) approach, as conceptualized by Kojima et al. (2022), which simulates a step-by-step 846 reasoning process to facilitate complex problem-solving tasks under zero demonstration. To ensure 847 a fair comparison, none of the methods used examples from the benchmark as demonstrations for 848 in-context learning.

849 The diversity of our evaluation is further enriched by testing across various language mod-850 els with differing capacities, including CodeLlama-7b (Rozière et al., 2023), along with Chat-851 GPT(gpt-3.5-turbo) (Ouyang et al., 2022), and the more advanced GPT-4(gpt-4) model (Ope-852 nAI, 2023). Recognizing the potential format discrepancies in code outputs, especially considering 853 that models like qpt-3.5-turbo and qpt-4 may produce code in markdown format which is not 854 immediately executable, we implement post-processing steps to convert the original language model 855 outputs into a form that can be executed within a sandbox environment. This normalization ensures that all models are evaluated under uniform execution conditions, thereby facilitating a fair and direct 856 comparison of their code generation capabilities. Through this methodological framework, we aim to 857 provide a detailed and nuanced understanding of the performance landscape across a spectrum of 858 LLM-driven code generation approaches. 859

RAG Settings. RaR leverages the capabilities of Retrieval-Augmented Generation methods, which
 enhance the performance of language models by integrating external knowledge sources. Specifically,
 we employed the codeparrot/github-jupyter dataset as our primary search vector library.
 This dataset is a comprehensive compilation of 452k markdown and code pairs, meticulously extracted
 from Jupyter notebooks hosted on GitHub BigQuery, representing a rich repository of programming

knowledge and examples. We utilized OpenAI's text-embedding-ada-002 API service for
 all embedding calculations across different methods and base models.

A.2 MATHEMATICAL REASONING

Benchmarks. Our evaluation framework for assessing mathematical reasoning capabilities leverages 870 two primary benchmarks: the GSM8K dataset, which comprises over 8,000 multi-step mathematical 871 problems (Cobbe et al., 2021), and the GSM-HARD dataset, an adaptation of GSM8K where 872 numbers in the questions are replaced with larger values to increase problem complexity (Gao et al., 2022). This study employs the PAL methodology to scrutinize the mathematical reasoning results, 873 involving the utilization of Large Language Models (LLMs) to parse natural language problems, 874 generate intermediary programmatic solutions, and subsequently execute these solutions via a Python 875 interpreter. The test set for each benchmark consists of samples ranging from index 1 to 200. Uniquely, 876 our approach does not use any examples for in-context learning, differing from the original PAL 877 methods. 878

Metrics and Baselines. Accuracy serves as our principal metric for evaluation, aligning with the
established metric for the GSM8K benchmark. Each question undergoes three execution attempts,
with the average score recorded as the final result. The baselines, including DIRECT, CoT, RAG
(1 shot), and RAG (5 shots), are consistent with those outlined in code generation, facilitating a
comprehensive and comparative analysis across different code generation benchmarks. The RAG
settings are consistent with the code generation tasks.

A.3 EMBODIED PLANNING

We further conduct experiments on embodied planning benchmarks on open-ended environments Minecraft (Lin et al., 2023).

Benchmarks. The complexity and vast item interconnectivity within the open-world Minecraft 890 present an ideal testbed for evaluating the LLM's capability to generate long-horizon plans (Yuan 891 et al., 2023; Wang et al., 2023; b). With thousands of items and intricate relationships between them, 892 obtaining a specific item in survival mode from scratch may involve dozens of intermediate items and 893 their quantitative relationships, such as crafting 1 crafting table from 4 planks. This setting rigorously 894 tests the planning abilities of LLMs instead of low-level control policies (Cai et al., 2023b; Baker 895 et al., 2022; Cai et al., 2023a; Lifshitz et al., 2023; Yuan et al., 2024). Moreover, Wang et al. (2023b) 896 have identified instances of hallucinations about Minecraft knowledge in OpenAI's ChatGPT and a 897 general scarcity of Minecraft-related knowledge in open-source language models, making this task a 898 suitable benchmark for assessing the RAG algorithm's effectiveness.

The planning prompts are aligned with those used in DEPS (Wang et al., 2023c), structured as Python templates and evaluated using MC-TextWorld as detailed by Lin et al. (2023). A set of 100 tasks were randomly selected for the test set, ranging from simple objectives like obtaining a crafting table to more complex goals such as crafting an iron helmet and even challenging making an enchanting table. The task instruction is formulated as:

904 905

906

907 908

909

910

867

868

885

• Give you nothing in the inventory, generate a step-by-step plan for the task of obtaining a {placeholder:acacia_boat} in Minecraft survival mode, and describe the object Minecraft item and its number at every step. For every step, start with 'STEP' as start.

- Give you nothing in the inventory, generate a step-by-step plan for the task of obtaining a {placeholder:diamond_pickaxe} boat in Minecraft survival mode, and describe the object Minecraft item and its number at every step. For every step, start with 'STEP' as start.
- 911 912 There are over 100 tasks involving different Minecraft items.

RAG Settings. For the retrieval component of the RAG algorithm, we utilized the Minecraft Wiki⁵ and
 DigMinecraft⁶ websites as the information sources accessible to the LLMs. Data from these websites
 was cleaned and formatted into markdown text, then segmented into trunks not exceeding 2000 tokens

⁹¹⁶ 917

⁵https://minecraft.wiki/

⁶https://www.digminecraft.com/

918 each, with embedding calculations performed using OpenAI's text-embedding-ada-002 API 919 service. 920

Evaluation Metrics. Based on the methodology of Huang et al. (2022), our evaluation of open-ended, 921 long-horizon planning in Minecraft focuses on both executability and plausibility. Executability 922 primarily examines whether a plan can be carried out, including the accuracy of each step's precon-923 ditions and effects. The executability is automatically calculated using MC-TextWorld (Lin et al., 924 2023). However, executability only evaluates if an objective-level plan can be executed, without 925 considering the specific details involved in executing individual objectives. For instance, crafting a 926 wooden pickaxe requires placing a crafting table and arranging three planks and two sticks in a partic-927 ular pattern, which are important details for human execution but not assessed by MC-TextWorld. 928 Therefore, we complement our evaluation with human ratings to assess the plausibility of plans.

929 930

931

935

936

937 938

939 940

941

942 943

944

945

956 957

958 959

960

A.4 CREATIVE WRITING

932 To further understand the potential of Retrieval-Augmented Generation (RAG) models in enhancing 933 the creativity and relevance of generated content, we extend our investigation to open-ended text generation tasks within the realm of creative writing. 934

Benchmarks. The versatility of RaR was tested through a series of creative writing tasks, each chosen to highlight different aspects of open-ended text generation. These tasks include:

- Write a survey paper to summarize the placeholder:Retrieval-augmented Generation methods for Large Language Models.
- Describe of placeholder: Jin-Yong's life.
- Summarize the placeholder: American Civil War according to the timeline.

For each task, three variants for placeholder were created to ensure a comprehensive evaluation of the model's performance across different contexts and requirements.

RAG Settings. Differing from previous tasks, creative writing is categorized as an open-ended 946 generation task, demanding a broader scope of information retrieval to aid content generation. To 947 accommodate this, Google was utilized as the search engine, with the top-k web pages converted into 948 markdown text to assist the LLM in generating outputs. This approach allowed LLM to leverage a 949 wide array of information sources. 950

Baselines and Evaluations. To benchmark RaR's performance, we compared it against DIRECT, 951 RAG-1 shot, and RAG-5 shot methods, all based on the gpt-3.5-turbo model. The evaluation 952 was conducted by human experts, employing the TrueSkill rating system (Herbrich et al., 2006) to 953 calculate scores for each method. This evaluation framework enabled a comprehensive assessment of 954 each model's creative output quality, accuracy, relevance, and innovativeness. 955

PROMPT DETAILS В

Our prompts consist of three parts: prompt for generating initial answer, prompt for generating search query, and prompt for revising answers according to retrieved context.

961	
962	Prompt B.1: Prompt for generating initial answers in creative writing tasks
963 964	{user}
965	##Question: {auestion}
966	##Instruction:
967	Try to answer this question/instruction with step-by-step thoughts and make the answer more structural.
968	Use /n/n to split the answer into several paragraphs.
969	Just respond to the instruction directly. DO NOT add additional explanations or introducement in the answer
970	unless you are asked to.
971	{assistant}

The process of query generation is omitted in code generation tasks. Instead, we use the generated code draft as a query and compute the embedding of it based on OpenAI Embedding services. For embodied planning and creative writing tasks, we will generate an additional query.

-	
7	Prompt B 2: Prompt for generating open-search query in creative writing tasks
8	Trompt D.2. Trompt for generating open search query in creative writing tasks
9	##Question:
0	{question}
	##Content:
	{answer}
	##Instruction.
	Please summarize the content with the corresponding question.
	This summarization will be used as a query to search with Bing search engine.
	The query should be short but need to be specific to promise Bing can find related knowledge or pages.
	You can also use search syntax to make the query short and clear enough for the search engine to find relevant
	language data.
	ity to make the query as relevant as possible to the last rew sentences in the content.
	Just output the query directly DO NOT add additional explanations or introducement in the answer unless
	you are asked to.
	{assistant}
	Dramat D 2: Dramat for ravising anowar according to rational materials in greative writing tasks
	Frompt D.S. Frompt for revising answer according to retrieved materials in creative writing tasks
	{user}
	##Existing Text in Wiki Web:
	{content}
	##Question:
	{question}
	##Answer:
	{answer} ##Instruction:
	I want to revise the answer according to retrieved related text of the question in WIKI pages
	You need to check whether the answer is correct.
	If you find some errors in the answer, revise the answer to make it better.
P	If you find some necessary details are ignored, add it to make the answer more plausible according to the

related text.

{assistant}

IMPORTANT

C TRUESKILL EVALUATION FRAMEWORK

Part of the tasks in "Embodied planning" and "creative writing" involve using humans for labeling.
Human labelers have 4 choices: "A is better", "B is better", "Tie" or "Both are bad". In this case,
"Tie" and "Both are bad" will be counted as a tie. For each task group, we have selected more than
professional annotators to provide labels. We use the Python "trueskill" package to calculate
the win rate and score. The default score for every method is set as 25. In order to facilitate user
understanding and selection, we also provide prompts when entering the system.

If you find the answer is right and do not need to add more details, just output the original answer directly.

Try to keep the structure (multiple paragraphs with its subtitles) in the revised answer and make it more

structural for understanding. Split the paragraphs with /n/n characters. Just output the revised answer directly. DO NOT add additional explanations or annoucement in the revised answer unless you are asked to.

aniectaic Planning	0			
struction				
ive you nothing in the inventory, generate a step-by-step plan for the task of obtaining a aca vvery step, start with 'STEP' as start.	ia planks in Minecraft survival mode, and describe the object Minecraft item and its number at every step. F			
idel A Response:	Model B Response:			
STEP 1:	Step 1:			
tem: Empty Inventory শan: Look around your surroundings for Acacia trees. They have a unique grey-brown trunk খ	ith Task: Your initial task requires you to locate an Acacia tree which are usually found in the Sat			
. flat canopy. You can usually find these trees in the Savanna biome.	biome, recognized by their unique gray-brown trunk and flat canopy. After finding the tree, o it down using your hand or, if available, an axe. You will acquire Acacia logs in return which a			
;TEP 2: tem: Empty hands	essential for the creation of Acacia planks.			
lan: Approach the Acacia tree. You need to start breaking the wooden part of the tree, also known as the trunk or log.	Step 2: Item: Empty Hands			
STEP 3:	Task: Move closer to the Acacia tree and start breaking the wood, specifically the trunk or log			
tem: Acacia Log	your hotbar. Make sure to store these logs in your inventory for future use.			
'an: Arter you ve collected at least one Acacia log, open your inventory and place the Acacia n the crafting section. Each log will create 4 wooden planks.	Step 3:			
STEP 4:	Item: Acacia Log Task: After collecting at least one Acacia log, you need to access your crafting table which co			
tem: Acacia Planks Plan: Once you've placed the Acacia Log in the crafting section of your inventorv. vou'll have	a 3x3 crafting grid. Position your acquired log in the first box of the first row of the grid, whic create 4 Acacia planks instantly. Move these planks to your inventory. Remember, for each A			
cacia Planks as the result in the box next to it. Drag them into your inventory.	log you put in the crafting grid, you will yield 4 wooden planks.			
By following these steps, you've successfully obtained Acacia Planks in survival mode on	Step 4:			
linecraft. You started with nothing and using the resources around you were able to create a	/ ITEM: ACACIA PLANKS			
del A From:	Model A From:			
manj				
→ A is better				
	So Tie So the Both are bad			
Usevia A ■ skip Usevia A ure C.1: The human evaluation UI. We will of ruction on the page simultaneously. The source sling, [MASK] will be replaced with the specif	Built with Gradio ● Iisplay responses from two different methods for the of the response will be marked as [MASK], and after the method name.			
Skip Usevia A ure C.1: The human evaluation UI. We will a ruction on the page simultaneously. The source eling, [MASK] will be replaced with the specif Chatbot Arena : Benchmarking LLMs in the V fRules Refresh to obtain the question and its correspon	Built with Gradio ● Iisplay responses from two different methods for the of the response will be marked as [MASK], and after be nethod name. //ild ding answers from two anonymous models.			
	Both are bad Both are bad New Round New R			
Use via A ure C.1: The human evaluation UI. We will of ruction on the page simultaneously. The source eling, [MASK] will be replaced with the specific Chatbot Arena : Benchmarking LLMs in the V fRules Refresh to obtain the question and its corresponder /ote for the better answer. And then click "New f both answers are bad, vote for "Both are bad'	Both are bad Both are bad New Round I Solut with Gradio Solution Solution of the response will be marked as [MASK], and after line of the response will be marked as [MASK], and after line to method name. Vild ding answers from two anonymous models. P Round" to get a new question.			
Skip Usevia A ure C.1: The human evaluation UI. We will o ruction on the page simultaneously. The source eling, [MASK] will be replaced with the specif Chatbot Arena : Benchmarking LLMs in the V fRules Refresh to obtain the question and its correspon /ote for the better answer. And then click "New f both answers are bad, vote for "Both are bad" f you want to skip, click "Skip".	Built with Gradie Rew Round South are bad Rew Round South are bad Rew Round South are bad			
Skip Usevia A ure C.1: The human evaluation UI. We will o ruction on the page simultaneously. The source eling, [MASK] will be replaced with the specif Chatbot Arena : Benchmarking LLMs in the V fRules Refresh to obtain the question and its correspon /ote for the better answer. And then click "New f both answers are bad, vote for "Both are bad" f you want to skip, click "Skip".	Image: Point are bad Image: Point are ba			
Skip Usevia A ure C.1: The human evaluation UI. We will o ruction on the page simultaneously. The source eling, [MASK] will be replaced with the specif Chatbot Arena : Benchmarking LLMs in the V Rules Refresh to obtain the question and its correspon /ote for the better answer. And then click "New f both answers are bad, vote for "Both are bad" f you want to skip, click "Skip". Principle	Image: The Control of the responses from two different methods for the of the response will be marked as [MASK], and after the control of the response will be marked as [MASK], and after the control of the response will be marked as [MASK], and after the control of the response will be marked as [MASK]. //ild //ild //ind			
Skip Usevia A ure C.1: The human evaluation UI. We will o ruction on the page simultaneously. The source eling, [MASK] will be replaced with the specif Chatbot Arena : Benchmarking LLMs in the V Rules Refresh to obtain the question and its correspon /ote for the better answer. And then click "New f both answers are bad, vote for "Both are bad" f you want to skip, click "Skip". Principle ou can evaluate the performance of the model f **Relevance**: Does it answer the question	Image: The Control of the responses from two different methods for the of the response will be marked as [MASK], and after the of the response will be marked as [MASK], and after the control of the response will be marked as [MASK], and after the control of the response will be marked as [MASK]. //ild //ild //ind ///			
Skip Usevia A ure C.1: The human evaluation UI. We will of ruction on the page simultaneously. The source eling, [MASK] will be replaced with the specified Chatbot Arena : Benchmarking LLMs in the V Rules Refresh to obtain the question and its correspon /ote for the better answer. And then click "New f both answers are bad, vote for "Both are bad" f you want to skip, click "Skip". Principle ou can evaluate the performance of the model f **Relevance**: Does it answer the question a **Accuracy**: Is it accurate? For example a	Image: The Control of the responses from two different methods for the of the response will be marked as [MASK], and after the of the response will be marked as [MASK], and after the control of the response will be marked as [MASK], and after the control of the response will be marked as [MASK], and after the control of the response will be marked as [MASK], and after the control of the response will be marked as [MASK], and after the control of the response will be marked as [MASK], and after the control of the response will be marked as [MASK], and after the control of the response will be marked as [MASK]. //ild //ild //ind //i			
Skip Usevia A ure C.1: The human evaluation UI. We will of ruction on the page simultaneously. The source eling, [MASK] will be replaced with the specified Chatbot Arena : Benchmarking LLMs in the V Rules Refresh to obtain the question and its correspon /ote for the better answer. And then click "New f both answers are bad, vote for "Both are bad" f you want to skip, click "Skip". Principle bu can evaluate the performance of the model f **Relevance**: Does it answer the question a **Accuracy**: Is it accurate? For example, a gs; a diamond axe requires 3 diamonds and 2 s	Image: The Control of the responses from two different methods for the of the response will be marked as [MASK], and after here the of the response will be marked as [MASK], and after here the of the response will be marked as [MASK], and after here the of the response will be marked as [MASK], and after here the of the response will be marked as [MASK], and after here the of the response of the marked as [MASK], and after here the of the response of the res			
Skip Usevia A ure C.1: The human evaluation UI. We will of ruction on the page simultaneously. The source eling, [MASK] will be replaced with the specified Chatbot Arena : Benchmarking LLMs in the V fRules Refresh to obtain the question and its correspon //ote for the better answer. And then click "New f both answers are bad, vote for "Both are bad" f you want to skip, click "Skip".	Image: The Control of the responses from two different methods for the of the response will be marked as [MASK], and after her method name. //ild //ild //ild //ind			
Skip Usevia A ure C.1: The human evaluation UI. We will of ruction on the page simultaneously. The source eling, [MASK] will be replaced with the specific Chatbot Arena : Benchmarking LLMs in the V Rules Refresh to obtain the question and its correspon /ote for the better answer. And then click "New /ote for the better answer. And then click "New /ote for the better answer. And then click "New /ote for the better answer. And then click "New /ote for the better answer. And then click "New /ote for the better answer. And then click "New /ote for the better answer. And then click "New /ote for the better answer. And then click "New /ote for the better answer. So the model f /ou want to skip, click "Skip". // Principle ucan evaluate the performance of the model f **Relevance**: Does it answer the question a **Accuracy**: Is it accurate? For example, a gs; a diamond axe requires 3 diamonds and 2 s **Completeness**: Is it complete? For st crafting wooden planks and then crafting	Image: The Control of the responses from two different methods for the of the response will be marked as [MASK], and after the of the response will be marked as [MASK], and after the or method name. //ild //ild //ind			
Skip Usevia A ure C.1: The human evaluation UI. We will of ruction on the page simultaneously. The source eling, [MASK] will be replaced with the specified Chatbot Arena : Benchmarking LLMs in the V Rules Refresh to obtain the question and its correspon /ote for the better answer. And then click "New f both answers are bad, vote for "Both are bad" f you want to skip, click "Skip". Principle ou can evaluate the performance of the model f **Relevance**: Does it answer the question a **Accuracy**: Is it accurate? For example, a gs; a diamond axe requires 3 diamonds and 2 s **Completeness**: Is it complete? For st crafting wooden planks and then crafting termediate steps cannot be ignored.	Image: The Control of the responses from two different methods for the of the response will be marked as [MASK], and after the of the response will be marked as [MASK], and after the or method name. //ild //ild //ind /ind //ind			
	Image: The Control of the responses from two different methods for the of the response will be marked as [MASK], and after here the of the response will be marked as [MASK], and after here the method name. //ild //ild //ind			
Skip Usevia A ure C.1: The human evaluation UI. We will of ruction on the page simultaneously. The source eling, [MASK] will be replaced with the specif Chatbot Arena : Benchmarking LLMs in the V Rules Refresh to obtain the question and its correspon /ote for the better answer. And then click "New f both answers are bad, vote for "Both are bad" f you want to skip, click "Skip". Principle ou can evaluate the performance of the model f **Relevance**: Does it answer the question a **Accuracy**: Is it accurate? For example, a gs; a diamond axe requires 3 diamonds and 2 s **Completeness**: Is it complete? For st crafting wooden planks and then crafting termediate steps cannot be ignored. **Readability**: Is it coherent? **Executability**: Considering the character	Image: Source and Source			
Skip Usevia A ure C.1: The human evaluation UI. We will of ruction on the page simultaneously. The source eling, [MASK] will be replaced with the specified Chatbot Arena : Benchmarking LLMs in the V Rules Refresh to obtain the question and its correspon /ote for the better answer. And then click "New f both answers are bad, vote for "Both are bad" f you want to skip, click "Skip". Principle but can evaluate the performance of the model f **Relevance**: Does it answer the question a **Accuracy**: Is it accurate? For example, a gs; a diamond axe requires 3 diamonds and 2 s **Completeness**: Is it complete? For example, a st crafting wooden planks and then crafting termediate steps cannot be ignored. **Readability**: Is it coherent? **Executability**: Considering the character ** Vote now!	Image: The Council and the second			

To avoid the code library containing solutions to the exact problems being evaluated) in code generation benchmarks, we adopted a rigorous pre-processing methodology as described by Guo et al. (2024). This process was meticulously designed to remove any direct matches or overly similar

code snippets from our search vector library, thereby ensuring that our evaluation remains fair and uncontaminated by pre-existing solutions. This examination aims to underscore the importance of maintaining the integrity of the evaluation process while utilizing external knowledge sources to augment the capabilities of language models in code-generation tasks.

Table 4: Benchmark contamination study based on CodeLlama-7b-Python models. We find using the code corpus finetuning the CodeLlama models directly will damage model performance on the code generation benchmark, while using RaR for in-context generation yields better performance.

Method	Huma	ınEval	HumanEval+		
method	pass@1	pass@5	pass@1	pass@5	
DIRECT	40.85%	53.65%	37.43%	48.78%	
FINETUNE	29.02%	40.24%	26.34%	35.98%	
RaR	45.73%	59.75%	43.29%	53.66%	

To further explore the potential benchmark contamination, we also conducted additional finetuning on CodeLLaMA-7B-Python using the code corpus in Table 4.

1098 E MORE RESULTS

1084

1093 1094

1095

1096 1097

1099

1100 E.1 EMBODDIED PLANNING

Embodied planning involves multiple steps, each of which relies on specific world knowledge and causal knowledge (i.e., preceding steps are usually preconditions for subsequent steps), including recipes for items in Minecraft, tools for performing related actions, and quantity restrictions. Therefore, although the plan generated by ChatGPT may appear complete and correct, there are often errors within the steps that can affect the proper execution of the plan.

1107 We mark the errors found in the generated plan in red.

Γ	(induction)
L	(Instruction) Give you nothing in the inventory generate a step by step plan for the task of obtaining a golden an
	Minecraft survival mode, and describe the object Minecraft item and its number at every sten. For ever
	start with 'STEP' as start
	{response from Original ChatGPT}
	STEP 1: Start by punching trees to gather at least 4 logs Minecraft item: 4x Logs
	STEP 2: Craft the logs into a crafting table Minecraft item: 1x Crafting Table
	STEP 3: Place the crafting table on the ground and right-click on it to open the crafting interface.
	STEP 4: Use the crafting interface to turn 4 of the oak planks into a wooden pickaxe.
	- Minecraft item: Wooden Pickaxe
	STEP 5: Find a cave or dig into the ground to gather some stone. Mine at least 8 blocks of stone using
l	wooden pickaxe Minecraft item: 8x Cobblestone
	STEP 6: Craft the cobblestone into a furnace using the crafting table Minecraft item: 1x Furnace
L	STEP 7: Place the furnace on the ground and right-click on it to open the furnace interface.
	STEP 8: Use the furnace interface to smell some iron ore. Place at least 3 iron ore in the top slot
	Turnace and use wooden planks as fuel in the bottom slot Minecraft item: 3x fron ingots
	STEP 9: Crait the fron ingots into a bucket using the craiting table Minecrait field in the bucket
	ster 10. Find a water source (such as a lake of a five) and fight-click on it while holding the bucket
	STEP 11: Look for an apple tree in a forest biome. Break some leaves until you find an apple - Min
	item. 1x Apple
	STEP 12: Combine the water bucket and the apple in the crafting table to craft a golden apple - Min
	item: 1x Golden Apple
	Congratulations! You have successfully obtained a golden apple in Minecraft survival mode

within it. These include recipe errors in STEP 2, where the crafting table requires planks instead of wood; missing raw materials in STEP 4, as the wooden pickaxe needs 2 sticks but lacks the relevant

1137

step in the plan; absence of instructions to use a stone pickaxe to mine iron ore in STEP 8; and an incorrect recipe for golden apple in STEP 12, which should include gold ingots and an apple rather than a water bucket.

1138 Response E.2: Response from RAR(based on ChatGPT-3.5) in embodied planning tasks 1139 {instruction} 1140 Give you nothing in the inventory, generate a step-by-step plan for the task of obtaining a golden apple in 1141 Minecraft survival mode, and describe the object Minecraft item and its number at every step. For every step, 1142 start with 'STEP' as start. 1143 1144 {response from Original ChatGPT} 1145 STEP 1: Find a tree and chop it down to get wood. You need at least 4 Oak Logs. 1146 - Minecraft items: 4x Oak Logs 1147 STEP 2: Craft the oak log into oak planks. 1148 - Minecraft items: 16x Oak Planks 1149 STEP 3: Convert some oak planks into sticks. - Minecraft items: 4x Stick 1150 STEP 4: Use the crafting table to create a wooden pickaxe. 1151 - Minecraft items: 1x Wooden Pickaxe 1152 STEP 5: Start digging down to find stone and gather at least three pieces of cobblestone. 1153 - Minecraft items: 3x Cobblestone 1154 STEP 6: Use the crafting table to create a stone pickaxe. - Minecraft items: 1x Stone Pickaxe 1155 STEP 7: Explore caves and mineshafts to find iron ore. Collect 3 iron ore for crafting an iron pickaxe. 1156 - Minecraft items: 3x Iron Ore 1157 STEP 8: Smelt the iron ore using a furnace to obtain iron ingot. 1158 - Minecraft items: 3x Iron Ingot 1159 STEP 9: Use the crafting table to create an iron pickaxe. - Minecraft items: 1x iron pickaxe 1160 STEP 10: Explore caves and mineshafts to find gold ore. Use an iron pickaxe or higher to mine the gold ore 1161 block. Collect as much gold ore as needed for crafting a golden apple. 1162 - Minecraft items: 8x Gold Ore 1163 STEP 11: Smelt the raw gold using a furnace to obtain a gold ingot. 1164 - Minecraft items: 8x Gold Ingot STEP 12: Find an oak tree and chop it down until an apple falls. 1165 - Minecraft items: 1x Apple 1166 STEP 13: Open the crafting menu and place the items in the correct pattern to craft a Golden Apple. In the 1167 3x3 crafting grid, place 8 Gold Ingots in the first and third rows, leaving the middle square empty. In the 1168 second row, place 1 Gold Ingot in the first and third squares, and place an Apple in the middle square. 1169 - Minecraft items: 1x Golden Apple Congratulations! You have successfully obtained a golden apple in Minecraft survival mode. 1170 1171 There are still errors in the plan generated by RAR, such as the lack of tools before "STEP 8: smelt 1172 iron ore into iron ingots", which should be "Mine 8 cobblestone to craft 1 furnace". However, 1173 compared to the errors in ChatGPT, the error rate in the plan has been significantly reduced. 1174 1175 We have also listed the links of the retrieved pages involved in different steps in Table 5 and Table 6. 1176 We can see that the text sources retrieved in each step generated by RaR are usually highly related to 1177 the synthesized item of that step. Traditional RAG (with 5 retrieval documents) uses instructions for 1178 retrieval and can only find the final step and other unrelated items, which may even harm the model's output. Compared to standard RAG, RaR extracts more pertinent knowledge from the database and 1179 provides more accurate answers. Furthermore, RaR does not directly generate answers based on 1180 retrieved content but evaluates previously generated content using both external information and 1181 internal knowledge. If inconsistencies or inaccuracies are detected in previous responses, revisions 1182 are made accordingly. By leveraging LLM's reasoning capabilities, RaR can decrease its reliance on 1183 retrieved data. 1184 1185

1186 E.2 CREATIVE WRITING

1187

Discussions on Computational Efficiency.

Table 5: Page link of retrieved text by RaR in embodied Minecraft planning tasks.

1100				
1109	Step	Item	Recipe	Link
1190	1	4x Oak Log	-	https://minecraft.fandom.com/wiki/Log
1191	2	16x Oak Planks	4x Oak Log	https://www.digminecraft.com/basic recipes/make oak wood plank.php
1192	3	4x Stick	2x Oak Planks	https://www.digminecraft.com/basic_recipes/make_stick.php
1102	4	1x Wooden Pickaxe	3x Oak Planks, 2 Stick	https://www.digminecraft.com/tool_recipes/make_wooden_pickaxe.php
1193	5	3x Cobblestone	Wooden Pickaxe	https://minecraft.fandom.com/wiki/Cobblestone
1194	6	1x Stone Pickaxe	3x Cobblestone, 2 Stick	https://www.digminecraft.com/tool_recipes/make_stone_pickaxe.php
1105	7	3x Iron Ore	Stone Pickaxe	https://minecraft.fandom.com/wiki/Iron_Ore
1195	8	3x Iron Ingot	3x Iron Ore	https://www.digminecraft.com/basic_recipes/make_iron_ingot.php
1196	9	1 Iron Pickaxe	3x Iron Ingot, 2x Stick	https://www.digminecraft.com/tool_recipes/make_iron_pickaxe.php
1197	10	8x Gold Ore	Iron Pickaxe	https://minecraft.fandom.com/wiki/Gold_Ore
1100	11	8x Gold Ingot	8x Gold Ore	https://www.digminecraft.com/basic_recipes/make_gold_ingot.php
1198	12	1x Apple	-	https://minecraft.fandom.com/wiki/Apple
1199	13	1x Golden Apple	8x Gold Ingot, 1x Apple	https://www.digminecraft.com/food_recipes/make_golden_apple.php

Table 6: Page link of retrieved text by conventional RAG methods in embodied Minecraft planning tasks.

Step	Item	Recipe	Link
1	4x Oak Log	-	-
2	16x Oak Planks	4x Oak Log	-
3	4x Stick	2x Oak Planks	-
4	1x Wooden Pickaxe	3x Oak Planks, 2 Stick	-
5	3x Cobblestone	Wooden Pickaxe	-
6	1x Stone Pickaxe	3x Cobblestone, 2 Stick	-
7	3x Iron Ore	Stone Pickaxe	-
8	3x Iron Ingot	3x Iron Ore	-
9	1 Iron Pickaxe	3x Iron Ingot, 2x Stick	-
10	8x Gold Ore	Iron Pickaxe	-
11	8x Gold Ingot	8x Gold Ore	-
12	1x Apple	-	https://minecraft.fandom.com/wiki/Apple
13	1x Golden Apple	8x Gold Ingot, 1x Apple	https://minecraft.fandom.com/wiki/Golden_Apple https://www.digminecraft.com/food_recipes/make_golden_apple.php https://minecraft.fandom.com/wiki/Enchanted_Golden_Apple
14	Other unrelated Items	-	https://www.digminecraft.com/game_commands/enchant_command.php https://www.digminecraft.com/tool_recipes/make_golden_pickaxe.php

In our experiment, we used RAG-5 as the baseline, which retrieves 5 content sources (commonly
 the default choice for most current RAG algorithms). We also set an upper limit on RaR retrieval
 times, ensuring that RaR does not retrieve more than 5 contents. Therefore, RaR does not exceed the
 standard RAG algorithm in terms of content retrieval.

Compared to including all 5 retrieved contents at once in context for generation in RAG-5, RaR
 incorporates one external content at a time to reduce prompt length. This approach prevents exceeding
 the language model's context length limitation that may occur when inputting all knowledge into the
 prompt with RAG-5. Using only one external content per revision in RaR avoids issues related to
 excessive content length and allows more token space for generating longer content.

Additionally, we utilized the Google search engine directly for RaR searches and compared time consumption between RaR and RAG-5 when creating creative writing prompts. Our findings indicate that while increased search time plays a role, more time is consumed during crawling and generation processes. And RaR does not have a significant efficiency gap compared to RAG-5. The detailed time consumption of every method is listed in Table 7. However, since RaR must be processed causally, it is worth further exploring how to improve the efficiency of RaR by referring to speculative decoding and other methods (Leviathan et al., 2022).

1242

1243

		001	KAG-J	Nan		
	Search Time(s)	_	0.12	0.47		
	Crawl Content Time(s)	-	5.17	4.26		
	Generation Time(s)	3.59	10.72	12.75		
	Total Time(s)	3.59	17.37	18.45		
		_				_
Response E.3: F	esponse from ORIGINAL Chat	GPT-3.	5 in creati	ve writing ta	asks	
		_	_	_	_	_
Instruction }	ry of the United States {response t	from Ori	iginal Chat(GPT 1		
Describe the filsto	Ty of the Officed States. Tesponse		igiliai Cilaw	Grij		
The history of the	United States dates back thousands	s of year	s, with indi	genous peopl	les such as the	e Native
Americans inhabi	ting the land long before the arrival	of Éuro	pean settle	rs. However,	the official his	story of
the United States I	begins in 1776 when the Founding F	athers de	eclared inde	ependence fro	m Great Brita	in. This
narked the start of	the American Revolution, a conflic	t that ult	imately led	to the creation	n of the Unite	d States
s an independent	nauon.	rough a	neriod of -	anid expansio	on in the 10th	century
The country acquir	red vast amounts of territory includi	ng the L	ouisiana Pu	rchase and the	e annexation c	of Texas
This expansion wa	as accompanied by conflicts, such as	the Me	xican-Amei	rican War and	the Civil Wa	r, which
ested the unity of	the nation.					
n the 20th centur	y, the United States emerged as a glo	obal sup	erpower. Th	he country pla	yed a signific	ant role
n both World War	I and World War II, eventually beco	oming th	e leader of	the free world	during the C	old War
the Civil Rights M	lovement and the rise of consumerior	u sociai	and cultura	u shirts during	g uns une, m	cluding
the civil Rights iv	tovenient and the rise of consumers	5111.				
Today, the United	States remains a powerful nation.	with a d	iverse popu	ulation and a	complex hist	orv that
Today, the United continues to shape	States remains a powerful nation, e its identity.	with a d	iverse popu	ulation and a	complex hist	ory that
Today, the United continues to shape	States remains a powerful nation, e its identity.	with a d	liverse popu	ulation and a	complex hist	ory that
Today, the United continues to shape	States remains a powerful nation, e its identity.	with a d	liverse popu	ulation and a	complex hist	ory that
Today, the United continues to shape	States remains a powerful nation, e its identity.	with a d	liverse popu	ulation and a	complex hist	ory that
Today, the United continues to shape Response E.4: F	States remains a powerful nation, e its identity. Response from Zero-shot CoT in	with a d	liverse popu	ulation and a	complex hist	ory that
Today, the United continues to shape Response E.4: F	States remains a powerful nation, e its identity.	with a d	liverse popu	ulation and a	complex histo	ory that
Response E.4: F	States remains a powerful nation, e its identity. Response from Zero-shot CoT in	with a d	liverse popu	ulation and a	complex histo	ory that
Foday, the United continues to shape Response E.4: F instruction } Describe the histo response from 74	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT}	with a d	liverse popu	ulation and a	complex hist	ory that
Today, the United continues to shape Response E.4: F instruction} Describe the histo response from Ze **Pre-Columbian	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era**	with a d	liverse popu	ulation and a	complex hist	ory that
Response E.4: F instruction } Describe the histo response from Zu *Pre-Columbian The history of the	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** 9 United States begins long before	with a d creativ the arriv	e writing t al of Europ	ulation and a tasks bean explorer	s in the 15th	century.
Response E.4: F instruction } Describe the histo response from Za *Pre-Columbian 'he history of the 'he land was init	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** United States begins long before ially inhabited by diverse Native A	with a d creativ the arriv merican	e writing t al of Europ t tribes, suc	ulation and a tasks pean explorer ch as the Apa	s in the 15th che, Sioux, In	century. roquois,
Response E.4: F instruction } Describe the histo response from Zo **Pre-Columbian The history of the The land was init and Cherokee. Th	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** 9 United States begins long before to ially inhabited by diverse Native A ese societies had rich cultures, with	with a d creativ the arriv merican their ov	e writing t al of Europ tribes, suc yn language	ulation and a tasks pean explorer th as the Apa es, traditions,	s in the 15th che, Sioux, Ir and social str	century. roquois, uctures.
A continues to shape and the second s	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** 9 United States begins long before to ially inhabited by diverse Native A ese societies had rich cultures, with land, flourishing through hunting,	with a d creativ the arriv merican their ov fishing,	e writing t al of Europ a tribes, suc vn language farming, an	ulation and a tasks pean explorer ch as the Apa es, traditions, nd trade. How	rs in the 15th che, Sioux, Ir and social str vever, the indi	century. roquois, uctures, igenous
Response E.4: F instruction } Describe the histo response from Zi *Pre-Columbian 'he history of the 'he land was init nd Cherokee. Th 'hey lived off the opulations faced	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** e United States begins long before ially inhabited by diverse Native A ese societies had rich cultures, with land, flourishing through hunting, significant challenges, including of	with a d creativ the arriv merican their ov fishing, conflicts	e writing t al of Europ tribes, suc n language farming, ar displacen	ulation and a tasks pean explorer ch as the Apa es, traditions, nd trade. How nent, and disc	rs in the 15th iche, Sioux, Ir and social str vever, the indi eases brought	century. roquois, ructures. igenous t by the
Action of the second se	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** e United States begins long before ially inhabited by diverse Native A ese societies had rich cultures, with land, flourishing through hunting, significant challenges, including	with a d creativ the arriv merican their ov fishing, conflicts	e writing t e writing t tribes, suc n language farming, ar displacen	tasks pean explorer ch as the Apa es, traditions, nd trade. How nent, and disc	s in the 15th che, Sioux, In and social str vever, the indi eases brought	century. roquois, uctures. igenous t by the
Response E.4: F instruction } Describe the histo response from Ze **Pre-Columbian The history of the The land was init and Cherokee. Th They lived off the populations faced European settlers. **Colonial Period in the 16th centu	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** e United States begins long before fially inhabited by diverse Native A ese societies had rich cultures, with land, flourishing through hunting, significant challenges, including a **	with a d creativ the arriv merican their ov fishing, conflicts m territe	e writing (al of Europ tribes, suc vn language farming, ar displacen	tasks bean explorer ch as the Apa es, traditions, nd trade. How nent, and disc	s in the 15th che, Sioux, Ir and social str vever, the indi eases brought	century. roquois, ructures. igenous t by the es. The
Response E.4: R instruction } Describe the histo response from Ze **Pre-Columbian The history of the The land was init and Cherokee. Th They lived off the populations faced European settlers. **Colonial Period in the 16th centu Spanish were the	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** 9 United States begins long before to ially inhabited by diverse Native A ese societies had rich cultures, with land, flourishing through hunting, significant challenges, including of ** ry, European powers began to clai first to establish colonies, initially for	with a d creativ the arriv merican their ov fishing, conflicts m territo ocusing o	e writing t al of Europ tribes, suc yn language farming, ar displacen ories in wh	tasks bean explorer ch as the Apa es, traditions, nd trade. How nent, and disc nat is now the and South Am	s in the 15th che, Sioux, If and social str vever, the indi eases brought e United State herica. Later, J	century. roquois, uctures. igenous t by the es. The English,
Response E.4: F instruction } Describe the histo response from Zi **Pre-Columbian The history of the The land was init and Cherokee. Th They lived off the opulations faced European settlers. **Colonial Period n the 16th centu Spanish were the is French, and Dutch	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** e United States begins long before ially inhabited by diverse Native A ese societies had rich cultures, with land, flourishing through hunting, significant challenges, including of ** ry, European powers began to clai first to establish colonies, initially for explorers arrived, establishing their	with a d creativ the arriv merican their ov fishing, conflicts m territo ocusing o	e writing t al of Europ a tribes, suc vn language farming, ar a, displacen ories in wh on Central a lonies along	tasks bean explorer ch as the Apa es, traditions, nd trade. How nent, and disc nat is now the and South Am g the Atlantic	e United State coast. These of	century. roquois, uctures. igenous t by the es. The English, colonies
Response E.4: F instruction } Describe the histo response from Za **Pre-Columbian The history of the The land was init and Cherokee. Th They lived off the populations faced European settlers. **Colonial Period In the 16th centu Spanish were the French, and Dutch varied in purpose	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** e United States begins long before to ially inhabited by diverse Native A ese societies had rich cultures, with land, flourishing through hunting, significant challenges, including of ** ry, European powers began to clai first to establish colonies, initially for explorers arrived, establishing their and organization, with the English colonies	with a d creativ the arriv merican their ov fishing, conflicts m territo ocusing o own col- olonies	al of Europ e writing t a tribes, suc vn language farming, ar a, displacen ories in wh on Central a lonies along primarily se	tasks tasks tasks tasks tasks tasks traditions, nd trade. How nent, and disc traditions traditions, nd trade. How nent, and disc tast is now the and South Am g the Atlantic eeking econor	e United State coast. These of mic opportuni	century. roquois, uctures. igenous t by the es. The English, colonies ties and
Response E.4: F instruction } Describe the histo response from Zd **Pre-Columbian The history of the The land was init and Cherokee. Th They lived off the populations faced European settlers. **Colonial Period In the 16th centu Spanish were the f French, and Dutch varied in purpose religious freedom	States remains a powerful nation, its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** United States begins long before ially inhabited by diverse Native A ese societies had rich cultures, with land, flourishing through hunting, significant challenges, including ** ry, European powers began to clai first to establish colonies, initially for explorers arrived, establishing their and organization, with the English c	with a d creativ the arriv merican their ov fishing, conflicts m territo ocusing o own col olonies jies emer	al of Europ e writing (a tribes, suc vn language farming, ar , displacen ories in wh on Central a lonies along primarily se ged, forgin	tasks bean explorer ch as the Apa es, traditions, nd trade. How nent, and disc nat is now the and South Am g the Atlantic eeking econor ig their own is	e United State coast. These c mic opportuni dentities, dev	century. roquois, uctures. igenous t by the es. The English, colonies ties and eloping
Response E.4: F (instruction) Describe the histo (response from Zd **Pre-Columbian The history of the The land was init and Cherokee. Th They lived off the populations faced European settlers. **Colonial Period In the 16th centu Spanish were the f French, and Dutch varied in purpose religious freedom local governments	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** United States begins long before ially inhabited by diverse Native A ese societies had rich cultures, with land, flourishing through hunting, significant challenges, including a ** ry, European powers began to clai first to establish colonies, initially for explorers arrived, establishing their and organization, with the English c . Over time, thirteen British colonies, and experiencing an increase in per	with a d creativ the arriv merican their ov fishing, conflicts m territo cousing of own col- olonies j ies emer- opulation	al of Europ e writing (a tribes, suc yn language farming, ar s, displacen ories in wh on Central a lonies along primarily se ged, forgin n due to im	tasks bean explorer ch as the Apa es, traditions, nd trade. How nent, and disc traditions is nent, and disc traditions, nd trade. How nent, and disc traditions, nd trade. How nent, and disc traditions, net is now the and South Arr g the Atlantic ecking econor ig their own is migration.	e United State e United State perica. Later, 1 coast. These c mic opportuni dentities, dev	century. roquois, uctures. igenous t by the es. The English, colonies ties and eloping
Response E.4: F (instruction) Describe the histo (response from Ze **Pre-Columbian The history of the The land was init and Cherokee. Th They lived off the populations faced European settlers. **Colonial Period In the 16th centu Spanish were the in French, and Dutch varied in purpose religious freedom local governments **American Revo	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** e United States begins long before to ially inhabited by diverse Native A ese societies had rich cultures, with land, flourishing through hunting, significant challenges, including of ** ry, European powers began to clais first to establish colonies, initially for explorers arrived, establishing their and organization, with the English colonies, and experiencing an increase in po- lution and Beyond** the British colonies and their colories	with a d creativ the arriv merican their ov fishing, conflicts m territe cousing of own col olonies ies emer opulation	al of Europ e writing (a tribes, suc yn languag farming, ar , displacen ories in wh on Central a lonies along primarily se ged, forgin n due to im	tasks bean explorer ch as the Apa es, traditions, nd trade. How nent, and disc nat is now the and South Arr g the Atlantic eaking econor ing their own is migration.	e United State e United State perican Later, I coast. These c mic opportuni dentities, dev	century. roquois, uctures. igenous t by the es. The English, colonies ties and eloping
Response E.4: F (instruction) Describe the histo (response from Zet **Pre-Columbian The history of the The land was init and Cherokee. Th They lived off the populations faced European settlers. **Colonial Period In the 16th centu Spanish were the f French, and Dutch varied in purpose religious freedom local governments **American Revo Tensions between 1775. The colonial	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** 9 United States begins long before to ially inhabited by diverse Native A ese societies had rich cultures, with land, flourishing through hunting, disignificant challenges, including of ** ry, European powers began to clais first to establish colonies, initially for explorers arrived, establishing their and organization, with the English colonies and experiencing an increase in po- lution and Beyond** the British colonies and their colonies s, unified under the Continental Cor-	with a d creativ the arriv merican their ov fishing, conflicts m territo cousing of oonnies ies emer opulation al master	e writing (al of Europ tribes, suc yn language farming, ar , displacen ories in wh on Central a lonies along primarily se ged, forgin n due to im rs grew, leae ought for in	tasks bean explorer thas the Apa es, traditions, nd trade. How nent, and disc nat is now the and South Arr g the Atlantic eeking econor ing their own is migration.	rs in the 15th che, Sioux, Ir and social str vever, the indi eases brought e United State herica. Later, I coast. These of mic opportuni dentities, dev nerican Revol rom British r	century. roquois, roquois, uctures. igenous t by the es. The English, colonies ties and eloping ution in ule The
Response E.4: F (instruction) Describe the histor (response from Ze **Pre-Columbian The history of the The land was init and Cherokee. Th They lived off the populations faced European settlers. **Colonial Period In the 16th centu Spanish were the f French, and Dutch varied in purpose religious freedom local governments **American Revo Tensions between 1775. The colonie war lasted several	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** e United States begins long before ially inhabited by diverse Native A ese societies had rich cultures, with land, flourishing through hunting, isignificant challenges, including of ** ry, European powers began to clai first to establish colonies, initially for explorers arrived, establishing their and organization, with the English colonies, of experiencing an increase in po- lution and Beyond** the British colonies and their colonies, s, unified under the Continental Cor- years, until the signing of the Trea-	with a d creativ the arriv merican their ov fishing, conflicts m territo occusing o olonies ies emer opulation al master agress, fo ty of Pa	e writing (al of Europ a tribes, suc vn language farming, ar c, displacen ories in wh on Central a lonies along grimarily sa ged, forgin n due to im- rs grew, leas ought for in ris in 1783	tasks bean explorer chas the Apa es, traditions, nd trade. How nent, and disc nat is now the and South Am g the Atlantic eeking econor ig their own is migration. ding to the Ar dependence f , which grant	s in the 15th che, Sioux, In and social str vever, the indi eases brought e United State herica. Later, I coast. These of mic opportuni dentities, dev nerican Revol from British ru	century. roquois, ructures. igenous t by the es. The English, colonies ties and eloping ution in ile. The formed
Response E.4: F instruction } Describe the histor response from Za **Pre-Columbian The history of the The land was init and Cheroke. Th They lived off the populations faced European settlers. **Colonial Period In the 16th centu Spanish were the f French, and Dutch varied in purpose religious freedom local governments **American Revo Tensions between 1775. The colonie war lasted several United States of A	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** e United States begins long before ially inhabited by diverse Native A ese societies had rich cultures, with land, flourishing through hunting, isignificant challenges, including ** ry, European powers began to clai first to establish colonies, initially for explorers arrived, establishing their and organization, with the English colonies, and experiencing an increase in po- lution and Beyond** the British colonies and their colonies, s, unified under the Continental Cor- years, until the signing of the Trea America independence. After achie	with a d creativ the arriv merican their ov fishing, conflicts m territo cousing of olonies ies emer opulation al master ngress, for ty of Pa ving ind	e writing (al of Europ a tribes, suc vn language farming, ar c, displacen ories in wh on Central a lonies along primarily sa ged, forgin n due to im- rs grew, leas ought for in ris in 1783 ependence.	tasks bean explorer ch as the Apa es, traditions, nd trade. How nent, and disc tat is now the and South Am g the Atlantic eeking econor ig their own is migration. ding to the Ar idependence f , which grant , the founding	s in the 15th che, Sioux, In and social str vever, the indi eases brought e United State perica. Later, I coast. These c mic opportuni dentities, dev nerican Revol from British ru ted the newly g fathers conv	century. roquois, ructures. igenous t by the es. The English, colonies ties and eloping ution in ale. The formed tened to
Response E.4: F (instruction) Describe the histor (response from Zu **Pre-Columbian The history of the The land was init and Cherokee. Th They lived off the populations faced European settlers. **Colonial Period In the 16th centu Spanish were the 1 French, and Dutch varied in purpose religious freedom local governments **American Revo Tensions between 1775. The colonie war lasted several United States of A draft the Constitut	States remains a powerful nation, e its identity. Response from Zero-shot CoT in ry of the United States. ero-shot CoT} Era** e United States begins long before ially inhabited by diverse Native A ese societies had rich cultures, with land, flourishing through hunting, significant challenges, including ** ry, European powers began to clai first to establish colonies, initially fo explorers arrived, establishing their and organization, with the English c . Over time, thirteen British colonies, and experiencing an increase in po- lution and Beyond** the British colonies and their colonies, unified under the Continental Cor years, until the signing of the Trea America independence. After achie ion, establishing the framework for	with a d creativ the arriv merican their ov fishing, conflicts m territo occusing o olonies ies emer opulation al master agress, for ty of Pa ving ind the new	e writing (al of Europ tribes, suc n language farming, ar displacen ories in wh on Central a lonies along primarily sa ged, forgin n due to im- rs grew, leas ought for in ris in 1783 ependence, nation's go	tasks bean explorer ch as the Apa es, traditions, nd trade. How nent, and disc tat is now the and South Am g the Atlantic eeking econor ig their own is migration. ding to the Ar idependence f , which grant , the founding povernment.	s in the 15th iche, Sioux, Ir and social str vever, the indi eases brought e United State herica. Later, I coast. These of mic opportuni dentities, dev nerican Revol from British ru ted the newly g fathers conv	century. roquois, ructures. igenous t by the es. The English, colonies ties and eloping lution in ale. The formed rened to

Table 7: Comparison of time consumed by different methods in creative writing tasks.

This is just a brief overview, and the history of the United States is vast and complex. However, these steps provide a simplified understanding of how the nation evolved from Native American civilizations to a group of British colonies, to an independent country.

297	Response E.5: Response from RAG(1-shot) in creative writing tasks
298	{instruction}
299	Describe the history of the United States.
300	{ response from KAG-1} The history of the United States is a complex and rich tanestry that spans over four centuries. The story
301	begins with the arrival of European explorers and settlers in the 16th century. In 1492, Christopher Columbus
302	reached the Americas, which led to subsequent expeditions and colonization by the Spanish, French, Dutch,
303	and English. The first successful English settlement was established at Jamestown, Virginia in 1607.
304	Over the next few centuries, the colonies grew and developed their own unique cultures, economies, and systems of government. Tensions rose between the American colonies and the British Crown leading to the
305	American Revolutionary War in 1775. The colonies declared their independence in 1776 and formed the
207	United States of America, which became a republic based on democratic ideals and principles.
307	The young nation faced numerous challenges in its early years, including establishing a central government,
300	expanding westward, and dealing with the issue of slavery. The 19th century saw rapid westward expansion,
310	and powerful nation, with the abolition of slavery and the preservation of the Union. In the 20th century, the
311	United States became a global superpower, playing a crucial role in both World Wars and the Cold War. It
312	witnessed significant social movements for civil rights and equality, as well as advancements in technology
313	and innovation.
314	and a complex political landscape. Its history continues to shape its present and future, as the nation graphies
315	with ongoing social, economic, and political challenges.
316	
817	
318	
819	
320	
821	
322	
323	
324	
325	
26	
27	
28	
9	Response E.6: Response from RAG(5-shot) in creative writing tasks
80	(instruction)
31	Describe the history of the United States.
32	{response from RAG-5}
33	The history of the United States dates back to the colonization of the Americas by European explorers in the
34	16th century. The first permanent English settlement was established in Jamestown, Virginia in 1607. This
35 00	i marked the degrinning of british colonization in North America.
30	Over the next century, more British colonies were established along the eastern seaboard including Mas-
~ /	Over the next century, more British colonies were established along the eastern seaboard, including Mas- sachusetts, New York, Pennsylvania, and Georgia. These colonies grew and developed their own distinct
21	Over the next century, more British colonies were established along the eastern seaboard, including Mas- sachusetts, New York, Pennsylvania, and Georgia. These colonies grew and developed their own distinct cultures and economies.
38	Over the next century, more British colonies were established along the eastern seaboard, including Mas- sachusetts, New York, Pennsylvania, and Georgia. These colonies grew and developed their own distinct cultures and economies. Tensions between the colonies and the British government grew in the mid-18th century, leading to the
38 39	Over the next century, more British colonies were established along the eastern seaboard, including Mas- sachusetts, New York, Pennsylvania, and Georgia. These colonies grew and developed their own distinct cultures and economies. Tensions between the colonies and the British government grew in the mid-18th century, leading to the American Revolutionary War in 1775. The colonists, led by figures such as George Washington, Thomas Laffærson and Baniamin Franklin fought for their independence from Privice rule.
38 39 40	Over the next century, more British colonies were established along the eastern seaboard, including Mas- sachusetts, New York, Pennsylvania, and Georgia. These colonies grew and developed their own distinct cultures and economies. Tensions between the colonies and the British government grew in the mid-18th century, leading to the American Revolutionary War in 1775. The colonists, led by figures such as George Washington, Thomas Jefferson, and Benjamin Franklin, fought for their independence from British rule. In 1776, the Second Continental Congress adopted the Declaration of Independence, which declared the
38 39 40 41	Over the next century, more British colonies were established along the eastern seaboard, including Mas- sachusetts, New York, Pennsylvania, and Georgia. These colonies grew and developed their own distinct cultures and economies. Tensions between the colonies and the British government grew in the mid-18th century, leading to the American Revolutionary War in 1775. The colonists, led by figures such as George Washington, Thomas Jefferson, and Benjamin Franklin, fought for their independence from British rule. In 1776, the Second Continental Congress adopted the Declaration of Independence, which declared the thirteen American colonies to be independent states. This marked the birth of the United States of America.
38 39 40 41 42	Over the next century, more British colonies were established along the eastern seaboard, including Mas- sachusetts, New York, Pennsylvania, and Georgia. These colonies grew and developed their own distinct cultures and economies. Tensions between the colonies and the British government grew in the mid-18th century, leading to the American Revolutionary War in 1775. The colonists, led by figures such as George Washington, Thomas Jefferson, and Benjamin Franklin, fought for their independence from British rule. In 1776, the Second Continental Congress adopted the Declaration of Independence, which declared the thirteen American colonies to be independent states. This marked the birth of the United States of America. Following the Revolutionary War, the Founding Fathers drafted the Constitution of the United States, which
38 39 40 41 42 43	Over the next century, more British colonies were established along the eastern seaboard, including Mas- sachusetts, New York, Pennsylvania, and Georgia. These colonies grew and developed their own distinct cultures and economies. Tensions between the colonies and the British government grew in the mid-18th century, leading to the American Revolutionary War in 1775. The colonists, led by figures such as George Washington, Thomas Jefferson, and Benjamin Franklin, fought for their independence from British rule. In 1776, the Second Continental Congress adopted the Declaration of Independence, which declared the thirteen American colonies to be independent states. This marked the birth of the United States of America. Following the Revolutionary War, the Founding Fathers drafted the Constitution of the United States, which established a federal government and outlined the rights and freedoms of American citizens. The Constitution
37 38 39 40 41 42 43 44	Over the next century, more British colonies were established along the eastern seaboard, including Mas- sachusetts, New York, Pennsylvania, and Georgia. These colonies grew and developed their own distinct cultures and economies. Tensions between the colonies and the British government grew in the mid-18th century, leading to the American Revolutionary War in 1775. The colonists, led by figures such as George Washington, Thomas Jefferson, and Benjamin Franklin, fought for their independence from British rule. In 1776, the Second Continental Congress adopted the Declaration of Independence, which declared the thirteen American colonies to be independent states. This marked the birth of the United States of America. Following the Revolutionary War, the Founding Fathers drafted the Constitution of the United States, which established a federal government and outlined the rights and freedoms of American citizens. The Constitution was ratified in 1788 and is still the supreme law of the land.
37 38 39 40 41 42 43 44 45 45	Over the next century, more British colonies were established along the eastern seaboard, including Mas- sachusetts, New York, Pennsylvania, and Georgia. These colonies grew and developed their own distinct cultures and economies. Tensions between the colonies and the British government grew in the mid-18th century, leading to the American Revolutionary War in 1775. The colonists, led by figures such as George Washington, Thomas Jefferson, and Benjamin Franklin, fought for their independence from British rule. In 1776, the Second Continental Congress adopted the Declaration of Independence, which declared the thirteen American colonies to be independent states. This marked the birth of the United States of America. Following the Revolutionary War, the Founding Fathers drafted the Constitution of the United States, which established a federal government and outlined the rights and freedoms of American citizens. The Constitution was ratified in 1788 and is still the supreme law of the land. Since its founding, the United States has experienced significant growth and development. It expanded westward through the Louisiana Purchase and the annexation of Texas, and it eventually stretched from coast
37 38 39 40 41 42 43 44 45 46 47	Over the next century, more British colonies were established along the eastern seaboard, including Mas- sachusetts, New York, Pennsylvania, and Georgia. These colonies grew and developed their own distinct cultures and economies. Tensions between the colonies and the British government grew in the mid-18th century, leading to the American Revolutionary War in 1775. The colonists, led by figures such as George Washington, Thomas Jefferson, and Benjamin Franklin, fought for their independence from British rule. In 1776, the Second Continental Congress adopted the Declaration of Independence, which declared the thirteen American colonies to be independent states. This marked the birth of the United States of America. Following the Revolutionary War, the Founding Fathers drafted the Constitution of the United States, which established a federal government and outlined the rights and freedoms of American citizens. The Constitution was ratified in 1788 and is still the supreme law of the land. Since its founding, the United States has experienced significant growth and development. It expanded westward through the Louisiana Purchase and the annexation of Texas, and it eventually stretched from coast to coast. The country also faced significant challenges, such as the Civil War, the Great Depression, and the
37 38 39 40 41 42 43 44 45 46 47 48	Over the next century, more British colonies were established along the eastern seaboard, including Mas- sachusetts, New York, Pennsylvania, and Georgia. These colonies grew and developed their own distinct cultures and economies. Tensions between the colonies and the British government grew in the mid-18th century, leading to the American Revolutionary War in 1775. The colonists, led by figures such as George Washington, Thomas Jefferson, and Benjamin Franklin, fought for their independence from British rule. In 1776, the Second Continental Congress adopted the Declaration of Independence, which declared the thirteen American colonies to be independent states. This marked the birth of the United States of America. Following the Revolutionary War, the Founding Fathers drafted the Constitution of the United States, which established a federal government and outlined the rights and freedoms of American citizens. The Constitution was ratified in 1788 and is still the supreme law of the land. Since its founding, the United States has experienced significant growth and development. It expanded westward through the Louisiana Purchase and the annexation of Texas, and it eventually stretched from coast to coast. The country also faced significant challenges, such as the Civil War, the Great Depression, and the Civil Rights Movement.

1350 1351	Response E.7: Response from RaR in creative writing tasks
1352	{instruction}
1353	Describe the history of the United States.
1354	{response from RAR}
1355	The history of the United States is characterized by several significant phases. The pre-colonial period
1356	witnessed diverse Native American populations with unique cultures innabiling the land. This era transitioned into colonialism through European exploration and settlement by powers such as the Spanish Erench Dutch
1357	and British, leading to the establishment of colonies along the East Coast with distinct economic and political
1358	frameworks.
1359	A pivotal moment in American history was the American Revolution, beginning in 1775, as the colonies
1360 1361	sought independence from British rule. The Declaration of Independence in 1776 marked the nation's founding on principles of liberty and democracy, establishing autonomy from Great Britain. This revolutionary
1362	period not only birthed a republic based on the will of the people but also forged American national identity,
1363	the establishment of republican governments elsewhere.
1364	The United States stands out as the first modern nation to achieve independence through a national war of
1365	liberation, articulating its motives and aspirations in a declaration of independence. This event had global reverberations, spreading fundamental values like liberty, equality, natural and civil rights, and civic duty. The
1366	formative years were shaped by the drafting of the Constitution in 1787, defining the federal government's
1367	framework. The subsequent ratification of the Bill of Rights further cemented individual freedoms and
1368	outlined governance principles that have endured across the nation's history, ensuring a system of checks and
1369	balances that protect cutters rights. The 19th century in the United States witnessed significant developments including westward expansion
1370	industrialization, and conflicts like the Civil War (1861-1865), which decisively addressed issues of slavery
1371	and reinforced federal authority over the states. The late 19th and early 20th centuries were marked by rapid
1072	industrial growth, extensive immigration, urbanization, and overseas expansion, including involvement in
1373	as a global superpower after World War II. engaging in the Cold War with the Soviet Union and participating
1375	in conflicts like the Korean and Vietnam Wars, significantly impacting the nation's foreign policy and global
1376	standing.
1377	Recent history includes social movements for civil rights, gender equality, and environmental protection,
1378	international relations.
1379	
1380	
1381	
1382	
1383	
1384	
1385	
1386	
1387	
1388	
1389	
1390	
1391	
1392	
1393	
1394	
1393	
1207	
1302	
1300	
1400	
1401	
1402	
1403	