FACTOR: FACTORING COMPLEXITY AND CONTEXT LENGTH IN LONG-CONTEXT MODEL EVALUATION

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027 028 029

030

Paper under double-blind review

ABSTRACT

Large language models (LLMs) with extended context windows are gaining attention. However, whether these resource-intensive LLMs indeed surpass simpler Retrieval Augmented Generation (RAG) techniques remains debatable. We precisely delineate differences between long-context LLMs and RAG methods, emphasizing the long-context reasoning abilities of LLMs. Existing benchmarks for long-context models often focus on information retrieval, hindering the assessment of reasoning over extended contexts. We introduce the FACTOR benchmark (Factoring Analysis of Complexity and Textual Context in Reasoning). FACTOR consists of two suites of tasks, covering both the symbolic and real-world facets of reasoning evaluation. Both suites are carefully created to delineate task complexity and context length when evaluating LLMs. We present detailed evaluations of popular LLMs on FACTOR. Besides accuracy scores, we also model the relationship between accuracy and task complexity. A simple but consistent log-linear relationship works surprisingly well across various models. From the log-linear relationship, two explainable parameters, the slope or Complexity Decay Factor (CDF) and the y-intercept or Contextual Decay Offset (CDO), are shown to offer separate and insightful measures of the models' complex reasoning and long context innate ability. Our findings highlight distinct failure modes linked to task complexity and context length, underscoring the unique reasoning capabilities of long-context LLMs unattainable by RAG methods.¹

1 INTRODUCTION

Recently, large language models (LLMs) with extended context windows have gained attention in real-world applications (Achiam et al., 2023; Team et al., 2024). With the advent of next-generation models (Dubey et al., 2024), context lengths of up to 128K tokens are becoming the new norm. Despite these advancements, a persistent debate exists (Li et al., 2024; Yu et al., 2024) regarding whether sophisticated and resource-intensive long-context LLMs genuinely offer advantages over more straightforward and cost-effective Retrieval Augmented Generation (RAG) techniques. This paper aims to precisely delineate the differences between what LLMs can accomplish with extended context capabilities and what is attainable through RAG methods. We contend that long-context LLMs possess unique long-context reasoning abilities that are inherently challenging for RAG-based methods to replicate.

While existing benchmarks, such as RULER (Hsieh et al., 2024) and ∞ Bench (Zhang et al., 2024), cover 040 a range of tasks and are becoming the new paradigm for evaluating the long-context models, they often 041 fail to capture the fundamental distinctions between the reasoning capabilities of long-context LLMs and 042 the retrieval strengths of RAG methods. Specifically, existing benchmarks exhibit two key limitations that 043 inadvertently favor RAG-based approaches: (1) they heavily present tasks focused on tasks that assess the 044 model's retrieval ability for long context (key characteristics being the complexity independent of context length) and (2) even though some tasks indeed see complexity increases with context length, e.g. Variable Tracking (Hsieh et al., 2024), they are often still too simple, and doesn't require model's reasoning ability 046 to solve. We empirically showed that simple RAG techniques easily get perfect scores on the Variable Tracing 047 tasks. Using the MPnetv2 Song et al. (2020) as the encoder and Llama3.1 8B Instruct Dubey et al. (2024) 048 model as the generator, the system achieves 100% and 98% accuracy on 131k and 1M context length with only 1024 actual context during the retrieval. Experiments show that solving VT is viable through simple backtracking only the occurrence of the query variable without grasping the entire problem setting and states 051 of other variables. We contend that benchmarks hide the unique but impressive reasoning advantages that 052

¹Revision during Rebuttal are marked in color blue.



Figure 1: (a) presents representative performance of Llama-3.1-70B-Instruct on FACTOR. From the trend, we can see that as N or the complexity increases, the accuracy decreases. Also, as the context length of the question prompt increases, the curve shifts downward. (b) shows the ranking of mainstream long-context LLMs on the average accuracy taken from all 39 different complexity settings N. (c) presents our attempt to model accuracy versus complexity for various models generally. Surprisingly, we observe that the logarithm of FACTOR accuracy linearly correlates nicely with the complexity N.

074

079

082

083

084

085

087

066

067

068

069

long-context LLMs can provide over RAG techniques. Therefore, the long-context language model needs a benchmark that can evaluate the LLM ability beyond retrieval. More detailed evaluation in Section 3.

075 Therefore, in this paper, we introduce the FACTOR (Factoring Analysis of Complexity and Textual COntext 076 in Reasoning) benchmark, designed to systematically evaluate language models' long context complex 077 reasoning ability through carefully curated synthetic tasks that require models' grasp of the entire problem to 078 succeed. Similar to previous benchmarks, FACTOR disentangles task complexity and context length, each can be independently varied. Specifically, the FACTOR benchmark relies on a suite of synthetic task generators by adjusting the following two key knobs for independently controlling task complexity and context length. 081

- Number of Variables (Task Complexity): Defines the complexity of the reasoning task via the number of interdependent variables. For most evaluations, the number of variables is limited to less than 40.
- Length of Filler Text (Context Length): To independently control the task complexity, we insert the question prompt with text irrelevant to the necessary portion of logic arguments, referring to them as Filler Tex. Filler text lengths are selected from predefined lengths: 0, 4K, 8K, 16K, 32K, 64K, and 128K tokens.

FACTOR consists of three subsets: Easy, Medium, and Hard. Easy subset consists of symbolic reasoning tasks, where the model evaluated is asked to deduce from variables in "Vx" for integer x. For Medium and Hard, two of context in natural language and contains rich hidden operations from semantics. In all 090 three suites, the model is presented with long chains of variables, and they can only provide correct output 091 when they correctly capture the relationship of all variables that appeared. The rest of the context length 092 is filler text. Computation operators are limited to grade school level similar to GSM8K (Cobbe et al., 2021).

- The general trend of FACTOR is shown in Figure 1 (a). We comprehensively evaluate state-of-the-art 094 pre-trained LLMs on FACTOR, where Figure 5 presents model names enumeration. Besides, the 095 aforementioned rag technique only achieves 4.5% accuracy on the realistic tasks of FACTOR for op=5 and 096 4K, far worse than the full attention LLM counterpart (33%). A snippet of our evaluation is shown in Figure 1 (b). Besides, we also perform explainable mathematical modeling of the performance of various models 098 on FACTOR. To our surprise, we observe that the logarithm of accuracy generally correlates with the task 099 complexity linearly across all LLMs evaluated, as shown in Figure 1 (c). Through modeling, we obtain 100 more insightful comparisons between different models.
- 101 Moreover, we found that two parameters (slope and y-intercept) used in the linear regression possess 102 explainable meanings and be used as quantitative metrics for describing models' abilities and behaviors. 103 For a given context length, the slope, referred to as **Contextual Decay Factor** (CDF), indicates the rate 104 of degradation of the model when solving increasingly longer context. The y-intercept, or Contextual Decay 105 **Offset** (CDO) captures the model's baseline performance at the given context length. We can separately 106 conclude both the model's reasoning ability and the long-context tracking ability from CDF and CDO. 107 Specifically, benefiting from FACTOR design to isolate complexity and context length, we found that most



- **Reproducing Failure Modes Through Fine-Tuning Strategies**: In Section **??**, our experiments show that the observed failure modes can be reproduced using different fine-tuning methods—*course learning* (as in Llama models) and *mixed sequence length training* (as in GPT-40-mini). This highlights the impact of training methodologies on models' abilities to handle complex reasoning over long contexts.
- Unveiling Limitations via Repeated Sampling: Further in Section ??, we investigate inference-time strategies like repeated sampling and find that, although they have potential to improve overall performance, inherent biases limit models' abilities to indefinitely extend their reasoning capabilities. The longer the context, the more challenging it becomes to recover performance levels seen with clean context training.
- **2 D**E

2 RELATED WORK

2.1 LONG-CONTEXT LANGUAGE MODELS

Various works related to the Long-context Language Model have been proposed. Flash attention(Dao et al., 2022), Flash attention2(Dao, 2024), Ring attention(Liu et al., 2023a), and Tree attention(Shyam et al., 2024) significantly reduced the memory footprint and communication overhead for processing long context in engineering level across multiple nodes. Architectural level innovations such as sparse attentions represented by sliding window attention(Beltagy et al., 2020), are also widely used to reduce the overhead caused by the increasing sequence length. New training strategies, such as gradually extending the training context length in the final stages of pretraining have been applied to support a long context window(Dubey et al., 2024).

162 2.2 LONG CONTEXT BENCHMARKS AND TASKS 163

There have been a quite a few works benchmarking long-context language models. Existing comprehensive 164 benchmarks like ∞ bench(Zhang et al., 2024) cover realistic tasks including document QA, summary, and 165 synthetic tasks including information retrieval, expression calculation, extending the context length in the 166 benchmark to over 200k tokens. ∞ bench(Zhang et al., 2024) does have mathematical reasoning tasks, 167 however the most relevant math.calc part seems to be too difficult for SOTA models to work out. Synthetic 168 tasks often offer more control and are less affected by parametric knowledge in comparison with realistic 169 tasks. One comprehensive synthetic benchmark is RULER(Hsieh et al., 2024), a synthetic benchmark with 170 tasks including retrieval, variable tracking and so on, offering some controls over context length and task 171 complexity. Experiments with various complexities were done, but it does not provide a quantitative analysis of complexity and context length on the correctness of the task, let alone isolate two separate patterns of 172 performance decay. Other benchmarks usually focus on simple retrieval(Github, 2023; Liu et al., 2023b), 173 fact reasoning(Kuratov et al., 2024), the impact of long context on natural language reasoning(Levy et al., 174 2024) and other real-world knowledge involved tasks. 175

176 177

2.3 SYNTHESIZED DATASETS FOR LONG-CONTEXT REASONING

178 Synthesized tasks are simple to build and absolutely deterministic, data contamination safe, but highly 179 effective to evaluate certain aspects of LLM performance. Its use in long-context benchmarks is profound. Needle-in-the-haystack Kamradt (2023), a pioneering long-context synthesized task, now becomes the go-to 180 task for evaluating LLM long-context retrieval ability. On the other hand, LLM reasoning benchmarks also 181 see recent efforts in synthesized tasks. Mirzadeh et al. (2024) recently proposes to use build synthesized 182 dataset upon GSM8K Cobbe et al. (2021) to study the robustness of LLM reasoning. Part of our work draws 183 a strong inspiration from a series of works (Ye et al. (2024a), Ye et al. (2024b)) which systematically 184 studies the intricacies of decoder transformers in solving grade-school level problems. Following their 185 footsteps, we carefully redesign the process of generating the problems so current LLMs can solve without 186 training, and together with thoughtful steps in noise addition, we effectively construct effective reasoning 187 benchmarks for the long-context community.

188 189

RETRIEVAL-AUGMENTED GENERATION VS. LONG CONTEXT BENCHMARKS 3

190 In this section, we first show RAG's strong performance on current standard long context LLM benchmarks, 191 and we then show RAG's limitation to grasp the logic sequences and flow on our newly proposed tasks 192 that specifically aim to test the LLM's complex reasoning ability. 193

We build a standard RAG system for reference. We use MPnet-base-V2 Song et al. (2020) as the 194 sentence retriever and cosine similarity as the metrics for retrieval. For the RAG system decoder, we use 195 Llama-3.1-8B-Instruct. Then, we run the RAG system on standard long context benchmarks in Table 1 and 196 Table 2 with a budget of 2K next to Llama-3.1-8B-Instruct full context. Apart from very minimal tasks, RAG 197 performance on conventional long-context benchmarks is solid and is a close rival to the long-context LLMs.

200	Dataset	LLM (0-8k)	LLM + RAG 2K (0-8k)	LLM (>8k)	LLM + RAG 2K (>8k)
201	qasper	45.86	41.40	43.65	37.97
202	multifieldqa_en	52.95	51.82	48.43	43.41
203	hotpotqa	56.19	56.05	44.24	47.49
204	2wikimqa	48.05	50.41	33.40	36.01
205	multi_news	26.72	26.01	24.95	21.02
206	triviaqa	90.85	78.60	94.36	79.60
207	samsum	41.88	37.68	44.49	41.95

Table 1: Performance comparison of LLM and LLM + RAG 2K across different datasets.

208 209

199

210 Then, we run the RAG on our tasks. We select the FACTOR medium subset. The RAG cannot solve the 211 task completely. The results are presented in Table 3. We still use 8K context with budget of 2k. 212

213 To make the case even more convincing, we now enable iterative prefilling of the RAG system. In other words, we allow the RAG to reselect the context for every token generated. The augmented model then produce 100% 214 accuracy on VT of RULER, which is even higher than 8B decoder in full context. Even though it is better than 215 normal RAG, it still achieves almost 0 accuracy for op > 3, as shown in Table 3. RAG's poor performance

225 226 227

235

237

238 239 240

241

253 254

255 256

257

258

259

260

261

262

263 264

216	Models	n_s_1	n_s_2	n_s_3	n_mk_1	n_mk_2	n_mk_3	n_mv
217 218 219	Baseline Score RAG Score	100.0 100.0	100.0 100.0	100.0 100.0	100.0 100.0	100.0 100.0	100.0 98.0	98.5 100.0
220 221 222	Models Baseline Score RAG Score	n_mq 100.0 100.0	vt 99.2 86.4	cwe 94.8 34.8	fwe 86.67 80.67	qa_1 84.0 74.0	qa_2 62.0 60.0	

Table 2: RAG on RULER. "n" stands for needle-in-the-haystack, "s" for single, "mk" for multikey, "mv" for multivalue, "mq" for multi-query, "vt" for Variable Tracking.

ор	RAG (with 8B)	Llama 3.1 8B	RAG (with 70B)	Llama 3.1 70B
2	0.13	0.8024	0.175	0.975
3	0.11	0.5322	0.07	0.835
4	0.055	0.6794	0.145	0.92
5	0.095	0.2641	0.1	0.87
6	0.08	0.2338	0.12	0.705
7	0.065	0.2802	0.085	0.655
8	0.055	0.2459	0.05	0.59
9	0.05	0.1693	0.07	0.38
10	0.045	0.1370	0.045	0.39

Table 3: Performance comparison for RAG and Llama models at different operations.

on FACTOR is due to the high-quality noise we designed to increment context length from essential hundreds of tokens for solving the problem to 8k. We will give more details of the design in Section 4.4.

Dataset	Iterative Prefill RAG	Llama-3.1-8B-Instruct
VT from RULER	100%	99.2%
2	0.52	0.8024
4	0.22	0.6794
5	0.07	0.2641
6	0.02	0.2338
8	0.06	0.2459

Table 4: Comparison of VT from RULER and Noise Type performance between 8K Iterative Prefill RAG and 8K (Llama-3.1-8B-Instruct).

4 THE FACTOR BENCHMARK

We introduce the **FACTOR** (Factoring Analysis of Complexity and Textual COntext in Reasoning) benchmark. In the design, we intentionally disentangle the effects of task complexity and context length on language model performance, providing a systematic framework for evaluating reasoning over long contexts. **FACTOR** is subdivided into three different subsets Easy, Medium, and Hard. The task division is classified based on the hierarchical depth of the semantics. The Easy subset consists of assignments and operations acting on symbols, like "Vi" for integer i. Therefore, every operation in this subset is explicit, meaning the model directly follows the instruction to perform operations. More details in 4.1. The medium and hard are questions that have the maximum hierarchical depth of 2 and 3 respectively. Detailed definition and construction in 4.2.

265 4.1 SYMBOLIC PORTION

Tasks are generated by first creating a set of variables $\{v_0, v_1, ..., v_N\}$, where N represents the *task complexity*. Mathematical relationships among these variables (e.g., $v_i = v_j \pm 1$) are then established to form a dependency graph. Consistent values that satisfy all relationships are assigned. These relationships are embedded within filler text, the randomly generated text irrelavant to the logic components, to create contexts of varying lengths, representing different *context lengths*. To distinguish the variable relationships from the filler text, they are enclosed within triple angle brackets <<< and >>>, and further enclosed within @ symbols to separate them from the surrounding content. This ensures that the relationships are clearly identifiable and not affected by the filler text.

To avoid falling back to Variable Tracking-like tasks that use backtracking for answer generation, we choose not to allow the model to calculate only the value of certain query variables. Instead, as shown in the example, we ask the model to output all the variables of a query value, which none can be a valid answer. In practice, our approach works well in differentiating between strong and weak reasoning abilities across models. No partial answer is allowed, so to consistently do the question correctly, the model is asked to traverse through the entire computation graph.

A Symbolic task of FACTOR benchmark

This is the beginning of the text: $@\langle\langle\langle \operatorname{assign} v_1 = v_4 - 1 \rangle\rangle\rangle @@\langle\langle\langle \operatorname{assign} v_0 = v_4 - 1 \rangle\rangle\rangle @@\langle\langle\langle \operatorname{assign} v_3 = v_4 + 1 \rangle\rangle\rangle @@\langle\langle\langle \operatorname{assign} v_2 = 1 \rangle\rangle\rangle @@\langle\langle\langle \operatorname{assign} v_4 = v_2 \rangle\rangle\rangle @$ This the end of the text. The text contains relationships between variables enclosed by ' $\langle\langle\langle ' \operatorname{and} \rangle\rangle\rangle$ '. These relationships are not sequential assignments in a programming language. They are independent mathematical equations that are all true simultaneously. Using only these relationships, determine what variable(s), if any, are equal to 2. Show your step-by-step reasoning and calculations, and then conclude your final answer in a sentence.

Answer: v3

279 280 281

282 283

284

287

288

289

291 292 293

4.2 COMMONSENSE PORTION

Then, we attempt to map the computation graph into the real-world context. Similar to the symbolic datasets, we strictly control the difficulty level of each problem by the number of binary operations (two variables each) needed to get to the final answer. Our method draws strong inspiration from a previous study Ye et al. (2024a).

298 **Defining Hierarchical Depth** - The key to designing the commonsense reasoning benchmark similar to 299 Cobbe et al. (2021) is to craft the hidden operations. The difficulty of hidden operations is controlled by 300 the hierarchical structure of the context. For example, the context of "Animals in the Location" contains two classes of objects ("Animals" and "Location") where Animals is the possession of Location. It is a 301 structure with hierarchical depth two. Given a context of two layers of objects, we have the most fundamental 302 hidden operation: addition. For example, "The number of Beverly Forest's Fox is 2. The number of Beverly 303 Forest's Wolves is 3. Assuming there isn't any other type of animal in Beverly Forest." The total number 304 of animals in Beverly Forest equals the sum of the first two objects, even though there isn't any description 305 of addition in the text. Context of higher hierarchical depth usually contains all possible relationships of 306 the previous plus the newly emerged relationship pattern, making them harder to deduce. 307

On the other hand, adding "number of children" to the previous context increments the hierarchical depth 308 to three. Continuing the previous sentence, "The average number of children of Beverly Forest's Wolves 309 is 3, while the average number of children of Beverly Forest's Fox is 2. What is the total number of Beverly 310 Forest's animal children?" To compute that we need to multiply each animal count with their average count 311 of children and sum up all the animals. Similarly, the chain of operations isn't narrated explicitly in the 312 text, but it is intuitive for humans to apply. As an example, for depth of three, it contains the direct operations 313 in the Easy subset as well as the addition hidden operation, plus the sum of multiply. We naturally designate 314 problems with a depth of two to be in the Medium subset, and a depth of three to be in the Hard subset. 315 Below we present an example with FACTOR HARD where blue signifies abstract variables that appear 316 starting from Medium (depth 2), whereas cyan signifies abstract variables that only appear starting from Hard (depth 3). Commonsense reasoning problems mostly contain depth 2 and depth 3 relationships (Ye 317 et al. (2024a)), where deeper hierarchy is rare even in normal natural language. 318

319 Making the solution - Similar to the Easy subset, we start from the randomly generated computational 320 graph, which is a DAG, and then attach numbers, forming abstract parameters to the graph. Out from the 321 computational graph, we randomly select a variable as the query variable, then perform a topological sort of the 322 computation graph. For the nodes on the topological sort list that aren't pointed to by other variables, we initial-323 ize their initial values and continue going through the sort list until the query variable. The entire chain is the solution, while the value of the query variable is the answer. Below we present a full example of FACTOR Hard.

3	2	4
3	2	5
3	2	6

328

330

331

332

333

334

335 336 337

350

351

A typical problem of FACTOR Hard

Problem: The number of adult deer in Oakridge Riverside equals 2 times the total number of newborn animal children in Cedar Valley. The number of adult deer in Cedar Valley equals 2. The average number of newborn children per adult deer in Oakridge Riverside equals 4. The average number of newborn children per adult deer in Cedar Valley equals the number of adult deer in Cedar Valley. Question: What is the total number of adult animals in Oakridge Riverside? **Solution**: Define adult deer in Cedar Valley as s; so s = 2. Define average number of newborn children per adult deer in Cedar Valley as i; so i = s = 2. Define total number of newborn animal children in Cedar Valley as M; so M = s * i = 2 * 2 = 4. Define adult deer in Oakridge Riverside as h; z = M = 4; so h = 2 * z = 2 * 4 = 8. Define total number of adult animals in Oakridge Riverside as Z; so Z = h = 8. Answer: 8.

4.3 **REVERSE REASONING PATH**

338 Our problem construction mandates the use of a topological sort list, so all the hidden operations we added are 339 constructive, always either addition or multiplication. However, to enable hidden operations of subtraction and 340 division naturally through our generation pipeline, we also designate half of the problem in the benchmark 341 to be of the reverse reasoning path. For the reverse reasoning path, we first still construct the topological 342 list first. Then we initialize the previous query variable at the end of the sort list. Then, we assign the other 343 side of the sorted list, variables that no other nodes point to as the query variable. Therefore, since we go in the complete opposite direction of the normal forwarding logic consisting of addition and multiplication, 344 we now have natural subtraction and division. 345

Overview of Commonsense subsets - For both the medium and hard subsets, we have half of the problem
in the normal forward, and half in the reverse path. Also, we deploy three different contextual templates
to diversify the problem text. These templates are "location-animal-children", "city-school-teacher", and
"festival-movie-nomination".

4.4 NOISE ADDITION

352 Adding noise is essential to enriching the length of the problem in the long context regime. For FACTOR, we mainly explore two different directions of making the noise. For FACTOR Easy, we follow conventional 353 long context benchmarks and generate random noise to increase the context length. On the other hand, 354 for FACTOR Medium and Hard, Benefiting from our commonsense reasoning problem generator, we can 355 generate noise statements that are in the same format and semantics as the essential logic statements. In 356 practice, we also encourage the noise statements to be noise variables pointed by essential variables to enlarge 357 the connection between the core graph and the noise graph. This tight connection is the key to why RAG 358 cannot solve the problem effectively. On the other hand, long context LLM isn't sensitive to close noise. 359



Figure 3: (a) Illustration of the two-phase accuracy behavior as a function of task complexity N. (b) presents the two different accuracy patterns between GPT-4o-mini and Gemini-1.5-Pro. GPT-4o-mini outperforms Gemini-1.5-Pro in low complexity tasks, while Gemini-1.5-Pro is more capable of dealing with more complex tasks. (c) shows the amazing performance of o1-mini: its has a higher $N_{\rm eff}$ than any other models tested and a decent CDF. This illustrates the effectiveness of the inference time strategy on reasoning.

374 375

376

5 EVALUATION ON PRETRAINED MODELS

We evaluate a range of pre-trained language models using the FACTOR benchmark to understand how they handle increasing task complexity and long context lengths, identifying their failure modes. The evaluation

is structured into two parts: (1) Benchmarking with Zero Filler Context: Assessing models' abilities to
 handle task complexity independently of context length. (2) Benchmarking with Long Contexts: Analyzing
 how models that perform well without filler context degrade when exposed to long contexts. Here due
 to page limit, we mainly focus on the FACTOR Easy task, we put the evaluation on the Medium and
 Hard tasks later in the Appendix.

384 5.1 EVALUATION METRICS

Besides just getting the accuracy, we also care about modeling the relationships between the accuracy and the complexity, thanks to the FACTOR design to isolate complexity and context length.

Two-Phase Accuracy Behavior - Models generally exhibit a characteristic two-phase behavior (Figure 3 (a)) in accuracy as the number of variables N increases (see Figure 3). **Phase 1**: For small values of N, models maintain near-perfect accuracy, effectively handling tasks with low complexity. **Phase 2**: Beyond a critical complexity threshold (N_{eff}), accuracy declines exponentially with increasing N, indicating rapid degradation in performance for more complex tasks. This pattern suggests a limit to the task complexity that models can handle before performance significantly deteriorates to close zero, where it plateaus.

It is intuitive to fit the accuracy versus operations using logistic regression that can fit all two phases elegantly. 394 Also, another candidate is a stepwise function with a constant for phase 1 following a exponential decay function. The only difference is at the transition between the two phases, where the logistic regression predicts 396 a smooth transition, while exponential decay is a sharper one. Later, these two functions are mathematically 397 similar, even close to identical when x gets better. We compare these two methods rigorously on O1-mini's 398 hehavior, where we found that the step-wise gives 0.006 MSE score from ops 0-150, compared with 0.013 399 from logistic regression, suggesting that the true transition from top LLM is sharp. On the other hand, due to 400 lack of reasoning ability, most LLMs barely have the phase 1 displayed. Therefore, the focus of our analysis 401 now is on the decay stage, where the two candidates are very similar. We proceed with exponential decay function as the modeling tool for later analysis because of its surprisingly strong vicinity to our collected 402 data. Extensive studies and fitting to justify using exponential function is presented in Appendix B. 403

Evaluation Metrics Definition - In Phase 2, accuracy A decreases exponentially with increasing N. By taking the natural logarithm of accuracy, we linearize this decay:

407

$$\log(A) = \text{CDF} \times N + \text{CDO}$$

Where **Complexity Decay Factor** (**CDF**) is the negative slope of the line (CDF < 0), representing the rate at which accuracy decays with increasing task complexity; **Contextual Decay Offset** (**CDO**) is the intercept of the line, capturing the baseline performance level influenced by context length. Our two-phase model explains well the pattern of decreasing model accuracy with task complexity, and this model does not predict accuracy outside the 0-1 range.

From these parameters, we define the **Effective Complexity** N_{eff} , indicating the maximum task complexity the model can handle before significant performance degradation. It is calculated as the value of N when the extrapolated $\log(A) = 0$ (i.e., when accuracy A = 1): $N_{\text{eff}} = -\frac{\text{CDO}}{\text{CDF}}$. However, if the CDO is negative, the extrapolated N_{eff} becomes negative, which is not meaningful since we cannot have a negative number of variables. In such cases, the two-phase behavior is not observed—the model's accuracy declines from the outset without an initial phase of high accuracy. Therefore, for negative CDO values, we focus on the exponential decay characterized by the CDF and CDO. A more detailed description of data processing can be found in the Appendix D.

420

421 5.2 Symbolic Tasks: Benchmarking with Zero Filler Context

We evaluated 15 models, comparing their FACTOR benchmark metrics with ELO scores from the LMSYS
Chatbot Arena(Chiang et al., 2024) on hard prompts with style control (see Table 5). We use the area under the
accuracy curve (AUC) to rank the models. In particular, AUC40 denotes the AUC calculated for N less than 40.

The evaluation reveals key insights into model performance concerning task complexity: Models with similar LMSYS ELO scores exhibit different behaviors as task complexity N increases(see Figure ??): models like **GPT-40-mini-2024-07-18** (CDO = 0.4303, CDF = -0.0401) excel on simple tasks (high CDO) but degrade rapidly with increasing complexity (more negative CDF); Models like **Gemini-1.5-Pro-002**(Team et al., 2024) (CDO = -0.0696, CDF = -0.0081) struggle with simple tasks (low CDO) but handle complex tasks better as they degrade more slowly (less negative CDF). **o1-mini** achieves both high CDO (0.6303) and a less negative CDF (-0.0117), resulting in a high N_{eff} (53.87). This indicates that inference-time strategies can significantly enhance performance across task complexities.(see Figure 3)

Iı	ndex	Model	CDF	CDO	$N_{\rm eff}$	AUC	AUC40	ELO
1		o1-mini	-0.0117	0.6303	53.87	139.34	38.25	1294
2		Gemini-1.5-Pro-002	-0.0081	-0.0696	-8.59	114.87	31.42	_
3		GPT-40-2024-05-13	-0.0220	0.2298	10.45	55.90	31.31	1251
4		GPT-40-2024-08-06	-0.0205	0.0899	4.39	53.17	29.45	1237
5		Claude-3-5-Sonnet-20240620	-0.0187	-0.1447	-7.74	46.27	25.52	1268
6		Qwen2.5-72B-Instruct	-0.0265	0.2056	7.76	45.50	28.26	1223
7		Mistral-Large-Instruct-2407	-0.0279	0.2100	7.53	43.37	28.12	1231
8		Gemini-1.5-Flash-002	-0.0244	-0.0180	-0.74	40.25	25.04	
9		GPT-4o-mini-2024-07-18	-0.0401	0.4303	10.73	35.67	27.19	1219
1	0	GPT-4-Turbo-2024-04-09	-0.0378	0.2514	6.65	33.10	25.06	1226
1	1	Llama-3.1-70B-Instruct	-0.0302	-0.0481	-1.59	31.56	21.6	1187
1	2	Qwen2-72B-Instruct	-0.0467	0.0123	0.26	21.67	17.91	1178
1	3	Claude-3-Haiku-20240307	-0.0471	-0.0848	-1.80	19.50	16.87	1173
1	4	Mistral-Nemo-Instruct-2407	-0.0608	0.0735	1.21	17.66	15.85	
1	5	Llama-3.1-8B-Instruct	-0.0694	-0.4615	-6.65	9.08	8.10	1132

Table 5: Metrics with Zero Filler Context are listed in this table. We observe a strong correlation between model capability and these metrics, together with different accuracy patterns of different models.

Models may perform similarly on general benchmarks but differ on tasks requiring complex reasoning. Selecting models for applications should consider CDO and CDF to match the complexity needs. The FACTOR benchmark highlights the necessity to evaluate models on complexity handling rather than solely on overall scores. By focusing on CDO and CDF, we are able to model distinct failure modes among models as task complexity increases. Understanding these metrics aids in selecting appropriate models for specific tasks and emphasizes the importance of specialized benchmarks like FACTOR.

5.3 SYMBOLIC TASKS: BENCHMARKING WITH LONG CONTEXTS

We analyze how models degrade when exposed to long contexts by examining the CDF and CDO metrics across different context lengths. Models were tested with varying amounts of filler context, extending up to their maximum context lengths (4K to 128K tokens). (see Table 6 and 7)

Table 6: Complexity Decay Factor (CDF) for Models at Different Context Lengths. (Values scaled by 10^2)

Model	4 K	8K	16K	32K	64K	128K
Mistral-Large-Instruct-2407	-3.26	-3.88	-4.04	-4.41	-15.19	
Qwen2.5-72B-Instruct	-2.82	-3.16	-3.27	-3.08	—	—
GPT-40-mini-2024-07-18	-4.93	-4.90	-4.95	-5.13	-5.01	-6.02
Llama-3.1-70B-Instruct	-3.88	-4.27	-4.82	-5.55	—	—
Qwen2-72B-Instruct	-4.66	-5.22	-5.56	-6.10	—	—
Mistral-Nemo-Instruct-2407	-5.15	-6.58	-29.43	-4.35	-3.58	-3.64
Llama-3.1-8B-Instruct	-7.24	-7.12	-8.48	-9.98	-10.19	

The evaluation reveals distinct failure modes when models are exposed to longer contexts.(all CDF values are scaled by 10^2 , see Figure 4) (1) Stable CDO, Degrading CDF (e.g., Llama Series): Models like Llama-3.1-**70B-Instruct** maintain relatively stable CDO across context lengths (from -0.014 at 4K to -0.104 at 32K), indicating consistent baseline performance. However, their CDF becomes more negative as context length increases (from -3.88 to -5.55), signifying increased difficulty with task complexity in longer contexts. 2. Stable CDF, Degrading CDO (e.g., GPT-4o-mini-2024-07-18): models like GPT-4o-mini-2024-07-18 maintains a consistent CDF across context lengths (approximately -4.93 to -5.13 up to 64K tokens), suggesting stable handling of task complexity. However, its CDO decreases steadily (from -0.020 at 4K to -0.711 at 64K), indicating declining baseline performance with longer contexts. **3. Degrading CDO** and CDF: models like Mistral-Large-Instruct-2407 show degradation in both CDO and CDF as context length increases, facing compounded difficulties with baseline performance and task complexity.

Model	4K	8K	16K	32K	64K	128K
Mistral-Large-Instruct-2407	0.085	-0.027	-0.153	-0.497	-0.510	_
Qwen2.5-72B-Instruct	0.055	-0.021	-0.088	-0.281	_	_
GPT-40-mini-2024-07-18	-0.020	-0.212	-0.312	-0.494	-0.711	-0.700
Llama-3.1-70B-Instruct	-0.014	-0.065	-0.054	-0.104		_
Qwen2-72B-Instruct	-0.172	-0.210	-0.197	-0.296	_	
Mistral-Nemo-Instruct-2407	-0.459	-0.609	-0.482	-1.154	-1.287	-1.287
Llama-3.1-8B-Instruct	-0.591	-0.679	-0.686	-0.608	-0.621	_

Table 7: Contextual Decay Offset (CDO) for Models at Different Context Lengths.



Figure 4: CDF/CDO of three models at different context length. Llama-3.1-70B-Instruct maintains a relatively stable CDO across context lengths, while its CDF decreases over context length. GPT-40-mini, on the contrary, maintains a relatively stable CDF while its CDO decreases over context length. Mistral-Large-Instruct is observed to have both metrics decreasing over context length.

Models exhibit different failure modes with long contexts: (1) Stable CDO, Degrading CDF: Models
 maintain baseline performance but increasingly struggle with task complexity as context lengthens; (2) Stable
 CDF, Degrading CDO: Models handle task complexity consistently but suffer declining overall performance
 with longer contexts; (3) Degrading CDO and CDF: Models face difficulties with both baseline performance
 and task complexity. Understanding these patterns is crucial for selecting appropriate models based on task
 requirements and highlights the importance of enhancing model robustness in processing long contexts
 and complex tasks.

6 CONCLUSION

In this paper, we introduced the FACTOR benchmark, a novel framework designed to systematically evaluate the complex reasoning abilities of large language models (LLMs) over long contexts. A key innovation of our work is the modeling of performance over task complexity, moving beyond traditional scalar evaluation metrics to capture the dynamic two-phase behavior in accuracy as complexity increases. Specifically, we observed that LLMs maintain high accuracy up to a certain complexity threshold, after which performance declines exponentially. By characterizing this behavior through the **Complexity Decay Factor (CDF)** and the **Contextual Decay Offset (CDO)**, we provided a nuanced understanding of how task complexity and context length independently affect model performance.

Our analysis revealed that these metrics not only quantify the degradation of logical reasoning ability and
 baseline accuracy but also highlight the limitations of current LLMs in handling complex reasoning tasks
 over extended contexts. Furthermore, we demonstrated that different fine-tuning strategies can reproduce
 these failure modes, emphasizing the significant impact of training methodologies on model capabilities.
 By modeling the two-phase behavior in accuracy rather than relying on a single performance score, the
 FACTOR benchmark offers a more detailed and insightful evaluation framework. This approach allows
 researchers to identify specific areas for improvement in LLMs and guides future developments in creating
 more robust language models capable of complex reasoning over long contexts.

540 REFERENCES

547

551

552

553

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo
 Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer, 2020.
 URL https://arxiv.org/abs/2004.05150.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li,
 Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. Chatbot arena: An open
 platform for evaluating llms by human preference, 2024.
 - Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168, 2021.
- ⁵⁵⁵ Tri Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. In *International Conference on Learning Representations (ICLR)*, 2024.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, 561 Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, 564 Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe 565 Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, 566 Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego 567 Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily 568 Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, 569 Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, 570 Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der 571 Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie 572 Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, 573 Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin 574 Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, 575 Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo 576 Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat 577 Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, 578 Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, 579 Nikolay Bogovchev, Niladri Chatterji, Olivier Duchenne, Onur Celebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, 581 Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, 582 Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana 583 Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan 584 Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, 585 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor 588 Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, 592 Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples,

594 Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita 595 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, 596 Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi 597 Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, 598 Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, 600 Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward 601 Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, 602 Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco 603 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, 604 Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, 605 Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun 606 Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, 607 Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, 608 Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan 609 Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay 610 Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun 611 Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro 612 Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian 613 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew 614 Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. 615 Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, 616 Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, 617 Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas 618 Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, 619 Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip 620 Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond 621 Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh 622 Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh 623 Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, 624 Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, 625 Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan 626 Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney 627 Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun 628 Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, 629 Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu 630 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, 631 Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen 632 Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and 633 Zhiwei Zhao. The llama 3 herd of models, 2024. URL https://arxiv.org/abs/2407.21783. 634

- Github. Needle in a haystack pressure testing llms, 2023. URL https://github.com/gkamradt/
 LLMTest_NeedleInAHaystack/tree/main.
- 637
 638
 639
 640
 640
 641
 642
 643
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
 644
- Gregory Kamradt. Needle in a haystack pressure testing llms, 2023. URL https:
 //github.com/gkamradt/LLMTestNeedleInAHaystack/tree/main.

Yuri Kuratov, Aydar Bulatov, Petr Anokhin, Ivan Rodkin, Dmitry Sorokin, Artyom Sorokin, and Mikhail Burtsev. Babilong: Testing the limits of llms with long context reasoning-in-a-haystack, 2024. URL https://arxiv.org/abs/2406.10149.

647 Mosh Levy, Alon Jacoby, and Yoav Goldberg. Same task, more tokens: the impact of input length on the reasoning performance of large language models, 2024. URL https://arxiv.org/abs/2402.14848.

657

658

659

660 661

662

663 664

- ⁶⁴⁸
 ⁶⁴⁹
 ⁶⁴⁹ Thuowan Li, Cheng Li, Mingyang Zhang, Qiaozhu Mei, and Michael Bendersky. Retrieval augmented generation or long-context llms? a comprehensive study and hybrid approach, 2024. URL https://arxiv.org/abs/2407.16833.
- Hao Liu, Matei Zaharia, and Pieter Abbeel. Ring attention with blockwise transformers for near-infinite context, 2023a. URL https://arxiv.org/abs/2310.01889.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni,
 and Percy Liang. Lost in the middle: How language models use long contexts, 2023b. URL
 https://arxiv.org/abs/2307.03172.
 - Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. arXiv preprint arXiv:2410.05229, 2024.
 - Vasudev Shyam, Jonathan Pilault, Emily Shepperd, Quentin Anthony, and Beren Millidge. Tree attention: Topology-aware decoding for long-context attention on gpu clusters, 2024. URL https://arxiv.org/abs/2408.04093.
 - Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. Mpnet: Masked and permuted pre-training for language understanding, 2020. URL https://arxiv.org/abs/2004.09297.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien 667 Vincent, Zhufeng Pan, Shibo Wang, Soroosh Mariooryad, Yifan Ding, Xinyang Geng, Fred Alcober, Roy 668 Frostig, Mark Omernick, Lexi Walker, Cosmin Paduraru, Christina Sorokin, Andrea Tacchetti, Colin 669 Gaffney, Samira Daruki, Olcan Sercinoglu, Zach Gleicher, Juliette Love, Paul Voigtlaender, Rohan Jain, 670 Gabriela Surita, Kareem Mohamed, Rory Blevins, Junwhan Ahn, Tao Zhu, Kornraphop Kawintiranon, 671 Orhan Firat, Yiming Gu, Yujing Zhang, Matthew Rahtz, Manaal Faruqui, Natalie Clay, Justin Gilmer, 672 JD Co-Reyes, Ivo Penchev, Rui Zhu, Nobuyuki Morioka, Kevin Hui, Krishna Haridasan, Victor Campos, 673 Mahdis Mahdieh, Mandy Guo, Samer Hassan, Kevin Kilgour, Arpi Vezer, Heng-Tze Cheng, Raoul 674 de Liedekerke, Siddharth Goyal, Paul Barham, DJ Strouse, Seb Noury, Jonas Adler, Mukund Sundararajan, 675 Sharad Vikram, Dmitry Lepikhin, Michela Paganini, Xavier Garcia, Fan Yang, Dasha Valter, Maja Trebacz, 676 Kiran Vodrahalli, Chulayuth Asawaroengchai, Roman Ring, Norbert Kalb, Livio Baldini Soares, Siddhartha Brahma, David Steiner, Tianhe Yu, Fabian Mentzer, Antoine He, Lucas Gonzalez, Bibo Xu, Raphael Lopez 677 Kaufman, Laurent El Shafey, Junhyuk Oh, Tom Hennigan, George van den Driessche, Seth Odoom, 678 Mario Lucic, Becca Roelofs, Sid Lall, Amit Marathe, Betty Chan, Santiago Ontanon, Luheng He, Denis 679 Teplyashin, Jonathan Lai, Phil Crone, Bogdan Damoc, Lewis Ho, Sebastian Riedel, Karel Lenc, Chih-Kuan 680 Yeh, Aakanksha Chowdhery, Yang Xu, Mehran Kazemi, Ehsan Amid, Anastasia Petrushkina, Kevin Swersky, Ali Khodaei, Gowoon Chen, Chris Larkin, Mario Pinto, Geng Yan, Adria Puigdomenech Badia, 682 Piyush Patil, Steven Hansen, Dave Orr, Sebastien M. R. Arnold, Jordan Grimstad, Andrew Dai, Sholto Douglas, Rishika Sinha, Vikas Yadav, Xi Chen, Elena Gribovskaya, Jacob Austin, Jeffrey Zhao, Kaushal 684 Patel, Paul Komarek, Sophia Austin, Sebastian Borgeaud, Linda Friso, Abhimanyu Goyal, Ben Caine, 685 Kris Cao, Da-Woon Chung, Matthew Lamm, Gabe Barth-Maron, Thais Kagohara, Kate Olszewska, Mia 686 Chen, Kaushik Shivakumar, Rishabh Agarwal, Harshal Godhia, Ravi Rajwar, Javier Snaider, Xerxes 687 Dotiwalla, Yuan Liu, Aditya Barua, Victor Ungureanu, Yuan Zhang, Bat-Orgil Batsaikhan, Mateo Wirth, James Qin, Ivo Danihelka, Tulsee Doshi, Martin Chadwick, Jilin Chen, Sanil Jain, Quoc Le, Arjun Kar, 688 Madhu Gurumurthy, Cheng Li, Ruoxin Sang, Fangyu Liu, Lampros Lamprou, Rich Munoz, Nathan 689 Lintz, Harsh Mehta, Heidi Howard, Malcolm Reynolds, Lora Aroyo, Quan Wang, Lorenzo Blanco, Albin 690 Cassirer, Jordan Griffith, Dipanjan Das, Stephan Lee, Jakub Sygnowski, Zach Fisher, James Besley, 691 Richard Powell, Zafarali Ahmed, Dominik Paulus, David Reitter, Zalan Borsos, Rishabh Joshi, Aedan 692 Pope, Steven Hand, Vittorio Selo, Vihan Jain, Nikhil Sethi, Megha Goel, Takaki Makino, Rhys May, 693 Zhen Yang, Johan Schalkwyk, Christina Butterfield, Anja Hauth, Alex Goldin, Will Hawkins, Evan Senter, Sergey Brin, Oliver Woodman, Marvin Ritter, Eric Noland, Minh Giang, Vijay Bolina, Lisa Lee, Tim Blyth, Ian Mackinnon, Machel Reid, Obaid Sarvana, David Silver, Alexander Chen, Lily Wang, Loren Maggiore, Oscar Chang, Nithya Attaluri, Gregory Thornton, Chung-Cheng Chiu, Oskar Bunyan, 697 Nir Levine, Timothy Chung, Evgenii Eltyshev, Xiance Si, Timothy Lillicrap, Demetra Brady, Vaibhav Aggarwal, Boxi Wu, Yuanzhong Xu, Ross McIlroy, Kartikeya Badola, Paramjit Sandhu, Erica Moreira, Wojciech Stokowiec, Ross Hemsley, Dong Li, Alex Tudor, Pranav Shyam, Elahe Rahimtoroghi, Salem 699 Haykal, Pablo Sprechmann, Xiang Zhou, Diana Mincu, Yujia Li, Ravi Addanki, Kalpesh Krishna, Xiao Wu, Alexandre Frechette, Matan Eyal, Allan Dafoe, Dave Lacey, Jay Whang, Thi Avrahami, Ye Zhang, Emanuel Taropa, Hanzhao Lin, Daniel Toyama, Eliza Rutherford, Motoki Sano, HyunJeong Choe, Alex

702 Tomala, Chalence Safranek-Shrader, Nora Kassner, Mantas Pajarskas, Matt Harvey, Sean Sechrist, Meire 703 Fortunato, Christina Lyu, Gamaleldin Elsayed, Chenkai Kuang, James Lottes, Eric Chu, Chao Jia, Chih-Wei 704 Chen, Peter Humphreys, Kate Baumli, Connie Tao, Rajkumar Samuel, Cicero Nogueira dos Santos, 705 Anders Andreassen, Nemanja Rakićević, Dominik Grewe, Aviral Kumar, Stephanie Winkler, Jonathan 706 Caton, Andrew Brock, Sid Dalmia, Hannah Sheahan, Iain Barr, Yingjie Miao, Paul Natsev, Jacob Devlin, Feryal Behbahani, Flavien Prost, Yanhua Sun, Artiom Myaskovsky, Thanumalayan Sankaranarayana Pillai, Dan Hurt, Angeliki Lazaridou, Xi Xiong, Ce Zheng, Fabio Pardo, Xiaowei Li, Dan Horgan, Joe 708 Stanton, Moran Ambar, Fei Xia, Alejandro Lince, Mingqiu Wang, Basil Mustafa, Albert Webson, Hyo 709 Lee, Rohan Anil, Martin Wicke, Timothy Dozat, Abhishek Sinha, Enrique Piqueras, Elahe Dabir, Shyam 710 Upadhyay, Anudhyan Boral, Lisa Anne Hendricks, Corey Fry, Josip Djolonga, Yi Su, Jake Walker, Jane 711 Labanowski, Ronny Huang, Vedant Misra, Jeremy Chen, RJ Skerry-Ryan, Avi Singh, Shruti Rijhwani, Dian 712 Yu, Alex Castro-Ros, Beer Changpinyo, Romina Datta, Sumit Bagri, Arnar Mar Hrafnkelsson, Marcello 713 Maggioni, Daniel Zheng, Yury Sulsky, Shaobo Hou, Tom Le Paine, Antoine Yang, Jason Riesa, Dominika 714 Rogozinska, Dror Marcus, Dalia El Badawy, Qiao Zhang, Luyu Wang, Helen Miller, Jeremy Greer, 715 Lars Lowe Sjos, Azade Nova, Heiga Zen, Rahma Chaabouni, Mihaela Rosca, Jiepu Jiang, Charlie Chen, 716 Ruibo Liu, Tara Sainath, Maxim Krikun, Alex Polozov, Jean-Baptiste Lespiau, Josh Newlan, Zeyncep 717 Cankara, Soo Kwak, Yunhan Xu, Phil Chen, Andy Coenen, Clemens Meyer, Katerina Tsihlas, Ada Ma, Juraj Gottweis, Jinwei Xing, Chenjie Gu, Jin Miao, Christian Frank, Zeynep Cankara, Sanjay Ganapathy, 718 Ishita Dasgupta, Steph Hughes-Fitt, Heng Chen, David Reid, Keran Rong, Hongmin Fan, Joost van 719 Amersfoort, Vincent Zhuang, Aaron Cohen, Shixiang Shane Gu, Anhad Mohananey, Anastasija Ilic, Taylor 720 Tobin, John Wieting, Anna Bortsova, Phoebe Thacker, Emma Wang, Emily Caveness, Justin Chiu, Eren 721 Sezener, Alex Kaskasoli, Steven Baker, Katie Millican, Mohamed Elhawaty, Kostas Aisopos, Carl Lebsack, 722 Nathan Byrd, Hanjun Dai, Wenhao Jia, Matthew Wiethoff, Elnaz Davoodi, Albert Weston, Lakshman 723 Yagati, Arun Ahuja, Isabel Gao, Golan Pundak, Susan Zhang, Michael Azzam, Khe Chai Sim, Sergi 724 Caelles, James Keeling, Abhanshu Sharma, Andy Swing, YaGuang Li, Chenxi Liu, Carrie Grimes Bostock, 725 Yamini Bansal, Zachary Nado, Ankesh Anand, Josh Lipschultz, Abhijit Karmarkar, Lev Proleev, Abe 726 Ittycheriah, Soheil Hassas Yeganeh, George Polovets, Aleksandra Faust, Jiao Sun, Alban Rrustemi, Pen 727 Li, Rakesh Shivanna, Jeremiah Liu, Chris Welty, Federico Lebron, Anirudh Baddepudi, Sebastian Krause, Emilio Parisotto, Radu Soricut, Zheng Xu, Dawn Bloxwich, Melvin Johnson, Behnam Neyshabur, Justin 728 Mao-Jones, Renshen Wang, Vinay Ramasesh, Zaheer Abbas, Arthur Guez, Constant Segal, Duc Dung 729 Nguyen, James Svensson, Le Hou, Sarah York, Kieran Milan, Sophie Bridgers, Wiktor Gworek, Marco 730 Tagliasacchi, James Lee-Thorp, Michael Chang, Alexey Guseynov, Ale Jakse Hartman, Michael Kwong, 731 Ruizhe Zhao, Sheleem Kashem, Elizabeth Cole, Antoine Miech, Richard Tanburn, Mary Phuong, Filip 732 Pavetic, Sebastien Cevey, Ramona Comanescu, Richard Ives, Sherry Yang, Cosmo Du, Bo Li, Zizhao 733 Zhang, Mariko Iinuma, Clara Huiyi Hu, Aurko Roy, Shaan Bijwadia, Zhenkai Zhu, Danilo Martins, 734 Rachel Saputro, Anita Gergely, Steven Zheng, Dawei Jia, Ioannis Antonoglou, Adam Sadovsky, Shane Gu, 735 Yingying Bi, Alek Andreev, Sina Samangooei, Mina Khan, Tomas Kocisky, Angelos Filos, Chintu Kumar, 736 Colton Bishop, Adams Yu, Sarah Hodkinson, Sid Mittal, Premal Shah, Alexandre Moufarek, Yong Cheng, Adam Bloniarz, Jaehoon Lee, Pedram Pejman, Paul Michel, Stephen Spencer, Vladimir Feinberg, Xuehan Xiong, Nikolay Savinov, Charlotte Smith, Siamak Shakeri, Dustin Tran, Mary Chesus, Bernd Bohnet, George Tucker, Tamara von Glehn, Carrie Muir, Yiran Mao, Hideto Kazawa, Ambrose Slone, Kedar 739 Soparkar, Disha Shrivastava, James Cobon-Kerr, Michael Sharman, Jay Pavagadhi, Carlos Araya, Karolis 740 Misiunas, Nimesh Ghelani, Michael Laskin, David Barker, Qiujia Li, Anton Briukhov, Neil Houlsby, Mia 741 Glaese, Balaji Lakshminarayanan, Nathan Schucher, Yunhao Tang, Eli Collins, Hyeontaek Lim, Fangxiaoyu 742 Feng, Adria Recasens, Guangda Lai, Alberto Magni, Nicola De Cao, Aditya Siddhant, Zoe Ashwood, Jordi 743 Orbay, Mostafa Dehghani, Jenny Brennan, Yifan He, Kelvin Xu, Yang Gao, Carl Saroufim, James Molloy, 744 Xinyi Wu, Seb Arnold, Solomon Chang, Julian Schrittwieser, Elena Buchatskaya, Soroush Radpour, 745 Martin Polacek, Skye Giordano, Ankur Bapna, Simon Tokumine, Vincent Hellendoorn, Thibault Sottiaux, 746 Sarah Cogan, Aliaksei Severyn, Mohammad Saleh, Shantanu Thakoor, Laurent Shefey, Siyuan Qiao, 747 Meenu Gaba, Shuo yiin Chang, Craig Swanson, Biao Zhang, Benjamin Lee, Paul Kishan Rubenstein, Gan 748 Song, Tom Kwiatkowski, Anna Koop, Ajay Kannan, David Kao, Parker Schuh, Axel Stjerngren, Golnaz 749 Ghiasi, Gena Gibson, Luke Vilnis, Ye Yuan, Felipe Tiengo Ferreira, Aishwarya Kamath, Ted Klimenko, 750 Ken Franko, Kefan Xiao, Indro Bhattacharya, Miteyan Patel, Rui Wang, Alex Morris, Robin Strudel, Vivek Sharma, Peter Choy, Sayed Hadi Hashemi, Jessica Landon, Mara Finkelstein, Priya Jhakra, Justin 751 Frye, Megan Barnes, Matthew Mauger, Dennis Daun, Khuslen Baatarsukh, Matthew Tung, Wael Farhan, 752 Henryk Michalewski, Fabio Viola, Felix de Chaumont Quitry, Charline Le Lan, Tom Hudson, Qingze 753 Wang, Felix Fischer, Ivy Zheng, Elspeth White, Anca Dragan, Jean baptiste Alayrac, Eric Ni, Alexander 754 Pritzel, Adam Iwanicki, Michael Isard, Anna Bulanova, Lukas Zilka, Ethan Dyer, Devendra Sachan, 755 Srivatsan Srinivasan, Hannah Muckenhirn, Honglong Cai, Amol Mandhane, Mukarram Tariq, Jack W. Rae,

756 Gary Wang, Kareem Ayoub, Nicholas FitzGerald, Yao Zhao, Woohyun Han, Chris Alberti, Dan Garrette, Kashyap Krishnakumar, Mai Gimenez, Anselm Levskaya, Daniel Sohn, Josip Matak, Inaki Iturrate, 758 Michael B. Chang, Jackie Xiang, Yuan Cao, Nishant Ranka, Geoff Brown, Adrian Hutter, Vahab Mirrokni, 759 Nanxin Chen, Kaisheng Yao, Zoltan Egyed, Francois Galilee, Tyler Liechty, Praveen Kallakuri, Evan 760 Palmer, Sanjay Ghemawat, Jasmine Liu, David Tao, Chloe Thornton, Tim Green, Mimi Jasarevic, Sharon Lin, Victor Cotruta, Yi-Xuan Tan, Noah Fiedel, Hongkun Yu, Ed Chi, Alexander Neitz, Jens Heitkaemper, 761 Anu Sinha, Denny Zhou, Yi Sun, Charbel Kaed, Brice Hulse, Swaroop Mishra, Maria Georgaki, Sneha 762 Kudugunta, Clement Farabet, Izhak Shafran, Daniel Vlasic, Anton Tsitsulin, Rajagopal Ananthanarayanan, Alen Carin, Guolong Su, Pei Sun, Shashank V, Gabriel Carvajal, Josef Broder, Iulia Comsa, Alena Repina, 764 William Wong, Warren Weilun Chen, Peter Hawkins, Egor Filonov, Lucia Loher, Christoph Hirnschall, 765 Weiyi Wang, Jingchen Ye, Andrea Burns, Hardie Cate, Diana Gage Wright, Federico Piccinini, Lei Zhang, 766 Chu-Cheng Lin, Ionel Gog, Yana Kulizhskaya, Ashwin Sreevatsa, Shuang Song, Luis C. Cobo, Anand Iyer, 767 Chetan Tekur, Guillermo Garrido, Zhuyun Xiao, Rupert Kemp, Huaixiu Steven Zheng, Hui Li, Ananth 768 Agarwal, Christel Ngani, Kati Goshvadi, Rebeca Santamaria-Fernandez, Wojciech Fica, Xinyun Chen, 769 Chris Gorgolewski, Sean Sun, Roopal Garg, Xinyu Ye, S. M. Ali Eslami, Nan Hua, Jon Simon, Pratik 770 Joshi, Yelin Kim, Ian Tenney, Sahitya Potluri, Lam Nguyen Thiet, Quan Yuan, Florian Luisier, Alexandra Chronopoulou, Salvatore Scellato, Praveen Srinivasan, Minmin Chen, Vinod Koverkathu, Valentin Dalibard, 771 Yaming Xu, Brennan Saeta, Keith Anderson, Thibault Sellam, Nick Fernando, Fantine Huot, Junehyuk Jung, Mani Varadarajan, Michael Quinn, Amit Raul, Maigo Le, Ruslan Habalov, Jon Clark, Komal Jalan, Kalesha Bullard, Achintya Singhal, Thang Luong, Boyu Wang, Sujeevan Rajayogam, Julian Eisenschlos, 774 Johnson Jia, Daniel Finchelstein, Alex Yakubovich, Daniel Balle, Michael Fink, Sameer Agarwal, Jing 775 Li, Dj Dvijotham, Shalini Pal, Kai Kang, Jaclyn Konzelmann, Jennifer Beattie, Olivier Dousse, Diane 776 Wu, Remi Crocker, Chen Elkind, Siddhartha Reddy Jonnalagadda, Jong Lee, Dan Holtmann-Rice, Krystal 777 Kallarackal, Rosanne Liu, Denis Vnukov, Neera Vats, Luca Invernizzi, Mohsen Jafari, Huanjie Zhou, Lilly 778 Taylor, Jennifer Prendki, Marcus Wu, Tom Eccles, Tianqi Liu, Kavya Kopparapu, Francoise Beaufays, 779 Christof Angermueller, Andreea Marzoca, Shourya Sarcar, Hilal Dib, Jeff Stanway, Frank Perbet, Nejc 780 Trdin, Rachel Sterneck, Andrey Khorlin, Dinghua Li, Xihui Wu, Sonam Goenka, David Madras, Sasha 781 Goldshtein, Willi Gierke, Tong Zhou, Yaxin Liu, Yannie Liang, Anais White, Yunjie Li, Shreya Singh, Sanaz Bahargam, Mark Epstein, Sujoy Basu, Li Lao, Adnan Ozturel, Carl Crous, Alex Zhai, Han Lu, 782 Zora Tung, Neeraj Gaur, Alanna Walton, Lucas Dixon, Ming Zhang, Amir Globerson, Grant Uy, Andrew 783 Bolt, Olivia Wiles, Milad Nasr, Ilia Shumailov, Marco Selvi, Francesco Piccinno, Ricardo Aguilar, Sara 784 McCarthy, Misha Khalman, Mrinal Shukla, Vlado Galic, John Carpenter, Kevin Villela, Haibin Zhang, 785 Harry Richardson, James Martens, Matko Bosnjak, Shreyas Rammohan Belle, Jeff Seibert, Mahmoud 786 Alnahlawi, Brian McWilliams, Sankalp Singh, Annie Louis, Wen Ding, Dan Popovici, Lenin Simicich, 787 Laura Knight, Pulkit Mehta, Nishesh Gupta, Chongyang Shi, Saaber Fatehi, Jovana Mitrovic, Alex Grills, 788 Joseph Pagadora, Dessie Petrova, Danielle Eisenbud, Zhishuai Zhang, Damion Yates, Bhavishya Mittal, 789 Nilesh Tripuraneni, Yannis Assael, Thomas Brovelli, Prateek Jain, Mihajlo Velimirovic, Canfer Akbulut, 790 Jiaqi Mu, Wolfgang Macherey, Ravin Kumar, Jun Xu, Haroon Oureshi, Gheorghe Comanici, Jeremy Wiesner, Zhitao Gong, Anton Ruddock, Matthias Bauer, Nick Felt, Anirudh GP, Anurag Arnab, Dustin Zelle, Jonas Rothfuss, Bill Rosgen, Ashish Shenoy, Bryan Seybold, Xinjian Li, Jayaram Mudigonda, Goker Erdogan, Jiawei Xia, Jiri Simsa, Andrea Michi, Yi Yao, Christopher Yew, Steven Kan, Isaac 793 Caswell, Carey Radebaugh, Andre Elisseeff, Pedro Valenzuela, Kay McKinney, Kim Paterson, Albert 794 Cui, Eri Latorre-Chimoto, Solomon Kim, William Zeng, Ken Durden, Priya Ponnapalli, Tiberiu Sosea, Christopher A. Choquette-Choo, James Manyika, Brona Robenek, Harsha Vashisht, Sebastien Pereira, 796 Hoi Lam, Marko Velic, Denese Owusu-Afriyie, Katherine Lee, Tolga Bolukbasi, Alicia Parrish, Shawn Lu, Jane Park, Balaji Venkatraman, Alice Talbert, Lambert Rosique, Yuchung Cheng, Andrei Sozanschi, Adam 798 Paszke, Praveen Kumar, Jessica Austin, Lu Li, Khalid Salama, Wooyeol Kim, Nandita Dukkipati, Anthony 799 Baryshnikov, Christos Kaplanis, XiangHai Sheng, Yuri Chervonyi, Caglar Unlu, Diego de Las Casas, 800 Harry Askham, Kathryn Tunyasuvunakool, Felix Gimeno, Siim Poder, Chester Kwak, Matt Miecnikowski, 801 Vahab Mirrokni, Alek Dimitriev, Aaron Parisi, Dangyi Liu, Tomy Tsai, Toby Shevlane, Christina Kouridi, 802 Drew Garmon, Adrian Goedeckemeyer, Adam R. Brown, Anitha Vijayakumar, Ali Elgursh, Sadegh 803 Jazayeri, Jin Huang, Sara Mc Carthy, Jay Hoover, Lucy Kim, Sandeep Kumar, Wei Chen, Courtney Biles, Garrett Bingham, Evan Rosen, Lisa Wang, Qijun Tan, David Engel, Francesco Pongetti, Dario de Cesare, 804 Dongseong Hwang, Lily Yu, Jennifer Pullman, Srini Narayanan, Kyle Levin, Siddharth Gopal, Megan Li, Asaf Aharoni, Trieu Trinh, Jessica Lo, Norman Casagrande, Roopali Vij, Loic Matthey, Bramandia Ramadhana, Austin Matthews, CJ Carey, Matthew Johnson, Kremena Goranova, Rohin Shah, Shereen 807 Ashraf, Kingshuk Dasgupta, Rasmus Larsen, Yicheng Wang, Manish Reddy Vuyyuru, Chong Jiang, Joana 808 Ijazi, Kazuki Osawa, Celine Smith, Ramya Sree Boppana, Taylan Bilal, Yuma Koizumi, Ying Xu, Yasemin Altun, Nir Shabat, Ben Bariach, Alex Korchemniy, Kiam Choo, Olaf Ronneberger, Chimezie Iwuanyanwu,

810	Shubin Zhao, David Soergel, Cho-Jui Hsieh, Irene Cai, Shariq Iqbal, Martin Sundermeyer, Zhe Chen, Elie
811	Bursztein, Chaitanya Malaviya, Fadi Biadsy, Prakash Shroff, Inderjit Dhillon, Tejasi Latkar, Chris Dyer,
812	Hannah Forbes, Massimo Nicosia, Vitaly Nikolaev, Somer Greene, Marin Georgiev, Pidong Wang, Nina
813	Martin, Hanie Sedghi, John Zhang, Praseem Banzal, Doug Fritz, Vikram Rao, Xuezhi Wang, Jiageng Zhang,
814	Viorica Patraucean, Dayou Du, Igor Mordatch, Ivan Jurin, Lewis Liu, Ayush Dubey, Abhi Mohan, Janek
815	Nowakowski, Vlad-Doru Ion, Nan Wei, Reiko Tojo, Maria Abi Raad, Drew A. Hudson, Vaishakh Keshava,
816	Shubham Agrawal, Kevin Ramirez, Zhichun Wu, Hoang Nguyen, Ji Liu, Madhavi Sewak, Bryce Petrini,
817	DongHyun Choi, Ivan Philips, Ziyue Wang, Ioana Bica, Ankush Garg, Jarek Wilkiewicz, Priyanka Agrawal,
818	Xiaowei Li, Danhao Guo, Emily Xue, Naseer Shaik, Andrew Leach, Sadh MNM Khan, Julia Wiesinger,
819	Sammy Jerome, Abhishek Chakladar, Alek Wenjiao Wang, Tina Ornduff, Folake Abu, Alireza Ghaffarkhah,
820	Marcus Wainwright, Mario Cortes, Frederick Liu, Joshua Maynez, Andreas Terzis, Pouya Samangouei,
020	Riham Mansour, Tomasz Kepa, François-Xavier Aubet, Anton Algymr, Dan Banica, Agoston Weisz,
021	Andras Orban, Alexandre Senges, Ewa Andrejczuk, Mark Geller, Niccolo Dal Santo, Valentin Anklin,
822	Majd Al Merey, Martin Baeuml, Trevor Strohman, Junwen Bai, Slav Petrov, Yonghui Wu, Demis Hassabis,
823	Koray Kavukcuoglu, Jeffrey Dean, and Oriol Vinyals. Gemini 1.5: Unlocking multimodal understanding
824	across millions of tokens of context, 2024. URL https://arxiv.org/abs/2403.05530.
825	Tion Vo. Zichong Vu. Wuonzhi Li and Zaman Allan Zhu. Dhusias of language models, Dart 2.1 grade school
000	Han te. Zicheng Au, tuanzhi Li, and Zevuan Ahen-Zhu. Physics of language models: Part 2.1. grade-school

- Tian Ye, Zicheng Xu, Yuanzhi Li, and Zeyuan Allen-Zhu. Physics of language models: Part 2.1, grade-school math and the hidden reasoning process, 2024a. URL https://arxiv.org/abs/2407.20311.
- Tian Ye, Zicheng Xu, Yuanzhi Li, and Zeyuan Allen-Zhu. Physics of language models: Part 2.2, how to learn from mistakes on grade-school math problems, 2024b. URL https://arxiv.org/abs/2408.16293.
- Tan Yu, Anbang Xu, and Rama Akkiraju. In defense of rag in the era of long-context language models, 2024. URL https://arxiv.org/abs/2409.01666.
- Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han, Zhen Leng
 Thai, Shuo Wang, Zhiyuan Liu, and Maosong Sun. ∞bench: Extending long context evaluation beyond
 100k tokens, 2024. URL https://arxiv.org/abs/2402.13718.



Figure 5: Left shows the difference between hierarchical depth in LLM evaluation under zero-context; Middle shows the comparison between different templates, showing they are the same; Right shows the performance comparison between forward logic and reverse logic

A MEDIUM AND HARD SUBSETS

In this section, we present more results on Medium and Hard subsets of FACTOR. We first compare all three different subsets in A.1. We then present more results on Llama-3.1-8B-Instruct and Llama-3.1-70B-Instruct on these two subsets in A.2. We compare forward reasoning and reverse reasoning in A.3. Comparison with different templates in A.4.

A.1 EASY VS. MEDIUM VS. HARD

First, we compare the Easy, Medium, and Hard subsets in the FACTOR dataset. The results are summarized in Table 8. You can see that from easy to medium to hard, the accuracy curve decays progressively quickly. The difference between Easy, Medium, and Hard is essentially the difference in hierarchical depth. To study the effect of hierarchical depth, we also make the depth-1 version of the commonsense problem generator. We plot them together in Figure 5 (left).

ор	Easy	Medium	Hard
5	0.808	0.86	0.7
8	0.832	0.875	0.512
10	0.708	0.612	0.41
12	0.628	0.56	0.36
15	0.556	0.43	0.19
20	0.516	0.2	0.05

Table 8: Performance comparison across difficulty levels (Easy, Medium, Hard) for different operations.

A.2 FULL RESULTS

Llama-3.1-70B-Instruct

1666					
400.0	302.2	310.0	223.5	201.3	152.2
1000.7	736.5	769.7	699.5	345.9	355.2
k medium	16k hard	32k medium	32k hard		
	1000.7 5k medium	1000.7 736.5 k medium 16k hard	1000.7 736.5 769.7 item it	1000.7 736.5 769.7 699.5 5k medium 16k hard 32k medium 32k hard	1000.7 736.5 769.7 699.5 345.9 5k medium 16k hard 32k medium 32k hard

215.1

169.6

Table 9: Comparison of models AUC across different context lengths. For zero-context and 4k, the AUC is computed between 2 and 20 inclusively, while the rest uses 2 to 10, because after then, they are very close to zero

234.2

253.3

915 Due to the main body's space limitations, we place the results of the commonsense reasoning Medium and 916 Hard subsets of FACTOR here. Currently, Llama-3.1-8B-Instruct and Llama-3.1-70B-Instruct are thoroughly 917 evaluated. We found that the trend we observed on FACTOR Easy also holds for FACTOR Medium and 918 Lord As a grade study, we plat the FACTOR medium in Figure 6 (a) and (b). The decaying trend is also

Hard. As a case study, we plot the FACTOR medium in Figure 6 (a) and (b). The decaying trend is also



Figure 6: Realistic Portion Results. (a) shows that the log-linear relationship also holds true for more complex datasets as well across three different LLMs. (b) shows that the performance generally degrades as the context length increases. (c) shows that for two different noise addition methods, which shows that the two noise aren't much different to LLMs across both 4K and 8K context length.

exponential, and if we take the accuracy logarithm, we see that the linear curve can also fit it neatly. We also summarize the AUC score for both 8B and 70B. For both medium and hard subsets and each ops split, we have 200 problems of forward logic, and 200 problems of reverse logic, as described in 4.3. We can compare the 8B and 70B models on FACTOR Medium zero context. We can see that the 8B model plateau way faster than the 70B model. Therefore, FACTOR potentially offers a much more comprehensive evaluation of LLMs performance in a easy to scale up difficulty manner.

A.3 FORWARD VS. REVERSE

Image shown in Figure 5 (right). The forward and backward generations are run using the 8B model. We can see that overall two curves have a similar trend. The reverse (yellow) is consistently less than or on par (not surpass).

A.4 TEMPLATE

Figure 5 (right) shows the comparison between the three contextual template. The dataset is collected for
Llama-3.1-8B-Instruct under zero-context. The difference might be subtle and small, which supports the
decision to use multiple templaThe forward reasoning is slightly better for LLM to reverse logic. Adding
multiple templates is leading huge improvement in text diversity.

B VISUALIZATION OF MODELING EXPONENTIAL DECAYINGJ CURVES

	Loglinear Fit	Logistic Regression Fit
O1-mini on FACTOR easy	0.01349	0.006215
Llama-3.1-8B-Instruct on FACTOR easy	0.003196	0.026159
Llama-3.1-70B-Instruct on FACTOR-easy	0.002697	0.041390
Llama-3.1-8B-Instruct on FACTOR medium	0.001932	0.003372
Llama-3.1-70B-Instruct on FACTOR medium	0.004230	0.016529

Table 10: Comparison of Loglinear Fit and Logistic Regression Fit across different models and datasets.

The experiment results show that loglinear is a strong modeling method, better than logistic regression.

We also show a study with much higher precision in Figure 8 on Gemini-1.5-Flash. We only focus on the decaying and then plateauing at near-zero performance. We fit logistic regression fit in red, but exponential decay in blue. We report MSE in the legends, we can see that exponential fit does provide a better fit, because of its sharp decrease.



Figure 7: Comparing Between Logistic Regression and LogLinear fit, Loglinear offers much closer fit to the plotted scatter points



Figure 8: High precision fitting study on the decaying phase and near zero plateauing phases of Gemini-1.5 Flash. Logistic regression fit and Exponential decaying curve are both presented in red and blue respectively.
 The exponential decaying curve fits the high-precision scatter plots better according to the MSE score

C TASK GENERATION AND DISTRIBUTION

In this section, we provide a detailed explanation of the process used to generate the synthetic tasks in the
 FACTOR benchmark. These tasks are meticulously designed to evaluate the reasoning abilities of language
 models over varying levels of task complexity and context length while ensuring that these two factors are
 independently controlled.

1026 C.1 TASK GENERATION PROCESS

The generation of each task involves several steps: creating variables and relationships, generating the payload, preparing the context, inserting the payload into the context, and forming the question prompt. This process is carefully constructed to introduce variability and prevent models from relying on simple heuristics or memorization.

1032 C.1.1 VARIABLE AND RELATIONSHIP CREATION

To control task complexity, we vary the number of variables N, where each variable is denoted as v_i for i=0,1,...,N-1. We establish interdependent relationships among these variables by generating a directed forest—a collection of trees representing dependencies.

1037 We start by randomly shuffling the variables to introduce variability. For each variable v_i where $i \ge k$ (with 1038 k being a randomly selected parameter between 1 and N), we randomly select a parent variable v_p from the 1039 set $\{v_0, v_1, ..., v_{i-1}\}$. This process creates directed edges from v_p to v_i , establishing dependency relationships.

To define the mathematical relationships between the variables, we assign simple operations to each edge. These operations are randomly chosen from {no operation, +1, -1}. The use of these simple operations ensures that the calculations required to solve the tasks remain within the realm of basic arithmetic, preventing models from facing overly complex computations.

Variable Value Assignment After establishing the relationships, we assign integer values to the variables
 while satisfying the dependencies. For root variables (those without parents), we assign random integer
 values between 0 and 10. This range is chosen to keep the numerical values small, again to prevent the
 need for complicated calculations.

1049 For each child variable v_i , its value is computed based on its parent's value and the assigned operation:

1050

1057

1044

1033

1051
1052
1053
1054value(v_i) =value(v_p)+1, if the operation is +1
value(v_p)-1, if the operation is -1
value(v_p), if there is no operation

This approach maintains simplicity in the calculations required, ensuring that the tasks assess the models' reasoning abilities rather than computational prowess.

1058 C.1.2 PAYLOAD GENERATION

The payload consists of the variable relationships formatted as textual statements. Each relationship is expressed as an assignment statement, enclosed within triple angle brackets <<< and >>> to clearly distinguish them from the filler text. These statements are further enclosed within single @ symbols when inserted into the context.

- An example of a payload statement is:
- 106
- 1066

1072

<<assign $v_i = v_p$ [operation]>>>

where [operation] is either +1, -1, or left empty if there is no operation. The payload statements are shuffled to present the relationships in a non-sequential order, adding to task complexity by preventing models from relying on the order of presentation.

1071 C.1.3 CONTEXT PREPARATION

1073 To control the context length independently, we generate filler text of specified lengths. The filler text is 1074 irrelevant to the variable relationships and serves to increase the context length, simulating scenarios where 1075 critical information is embedded within large amounts of unrelated data.

The filler text is created by randomly selecting words from a predefined list related to computational topics,
such as "algorithm," "data," "performance," and so on. These words are concatenated to form sentences,
with occasional punctuation added to simulate natural language text. The filler text may also include random
sentences or phrases inserted between the payload statements to further obscure the relationships and mimic
real-world text where key information is interleaved with irrelevant content.

1080 C.1.4 PAYLOAD INSERTION

The payload of variable relationships is inserted into the filler text at random positions. The filler text is
first tokenized into sentences using the NLTK library's sentence tokenizer. Insertion points are determined
based on specified intervals, ensuring that the payloads are dispersed throughout the context rather than
clustered together.

Each payload statement, enclosed within @ symbols and <<<>>> brackets, is inserted at the selected points. The maximum group size parameter controls how many payload statements can be inserted together at a single insertion point, adding another layer of variability. Additionally, filler content may be placed between payload statements, further increasing the challenge by requiring the model to discriminate between relevant and irrelevant information.

1091

C.1.5 QUESTION PROMPT FORMATION

Finally, we generate the question prompt that instructs the model on the task to perform. The prompt includes:

Delimiters indicating the beginning and end of the text. - An explanation that the relationships enclosed by <<< and >>> are independent mathematical equations that are all true simultaneously, and not sequential assignments in a programming language. - A task instruction, asking the model to determine which variable(s), if any, are equal to a randomly selected target value. The target value is chosen from among the variable values or from values slightly outside the range of assigned values, which may result in no variables matching the target, thereby introducing unpredictability. - A request for the model to show step-by-step reasoning and conclude with the final answer in a sentence.

This structure ensures that the model must process and reason over the entire context, filtering out irrelevant information and correctly interpreting the relationships to arrive at the answer.

1105

1106 C.2 TASK DISTRIBUTION

We generate tasks across a wide range of complexities and context lengths to create a comprehensive evaluation suite.

1110

1111 C.2.1 COMPLEXITY LEVELS

The number of variables N varies from 1 to 39, covering a broad spectrum of task complexities. For each value of N, we generate multiple tasks with different configurations to ensure diversity. The parameter k, controlling the number of trees in the forest, is randomly selected for each task and ranges from 1 to N.

1116

1117 C.2.2 CONTEXT LENGTHS

Context lengths are set to predefined values: 0 (no filler text), 1K, 2K, 4K, 8K, 16K, 32K, 64K, and 128K tokens. Our benchmark in the main context only include subsets of 0 and 4K or more. For each context length, filler text is generated accordingly, and the payloads are inserted as described. The inclusion of different context lengths allows us to evaluate how models handle tasks when key information is embedded within varying amounts of irrelevant text.

- 1124
- 1125 1126 C.2.3 SAMPLE SIZES

For each combination of N and context length, we generate 50 distinct task instances. This results in a total of 1,950 tasks for each context length setting (39 values of N times 50 tasks). The large sample size ensures that our evaluation is statistically robust and that any observed trends are not due to random chance.

- 1131
- 1132 C.3 EXAMPLE TASK

An example of a generated task is as follows:

34 Sar	nple Task
Co Thi $@ < v_0$ pet. [fThiTherelamaUsi $stepAnv3$	ntext: s is the beginning of the text: $<>>@ algorithm data optimization. @<<= v_4 - 1>>>@ performance analysis. @<<>>@ code snip @<<>>@ best practice. @<<>>@ iller text continues] is is the end of the text. = text contains relationships between variables enclosed by '<<<' and '>>>'. Theseationships are not sequential assignments in a programming language; they are independentthematical equations that are all true simultaneously.ing only these relationships, determine what variable(s), if any, are equal to 2. Show yourp-by-step reasoning and calculations, and then conclude your final answer in a sentence.swer:$
In this e Filler co the pay	example, the variables v_0 to v_4 have interdependent relationships defined by the equations provided. ontent, such as "algorithm data optimization" and "performance analysis," is interspersed between load statements, increasing the context length and complexity. The task requires the model to:
 Extrac Detern in a col 	t the relevant variable relationships from the payloads Interpret and solve the system of equations. nine which variable(s) equal the target value (in this case, 2) Present the reasoning and final answer herent manner.
This ex extende	ample illustrates how the task assesses the model's ability to perform multi-step reasoning over ad contexts that include irrelevant information.
C.4 I	Ensuring Diversity and Avoiding Data Leakage
To prev	ent models from exploiting patterns or memorizing specific instances, we employ several strategies:
• Variab order an options chosen meet th	le names are randomly shuffled for each task instance Relationships are presented in a random nd interleaved with filler content Operations assigned to edges are randomly selected from simple to maintain calculation simplicity while adding variability Target values for the query are randomly and may not correspond to any variable value in the task, introducing the possibility that no variables the condition.
These n Models pattern	neasures create a wide variety of task instances, reducing the likelihood of data leakage or overfitting. must genuinely understand and reason through each task rather than relying on memorization or recognition.
C.5 S	SUMMARY
The tas task con number of speci and cor reasoni	k generation process in the FACTOR benchmark is carefully designed to independently control mplexity and context length while preventing models from relying on shortcuts. By varying the of variables and their interdependencies, we manipulate task complexity. By inserting filler text fied lengths and interleaving filler content between payload statements, we manipulate context length mplexity. The use of simple operations and small integer values ensures that the focus remains on ng rather than computation.
This sys abilities perform in reaso	stematic approach allows us to create a diverse set of tasks that robustly evaluate models' reasoning s over long contexts. By disentangling the effects of task complexity and context length on model ance, the FACTOR benchmark facilitates a deeper understanding of models' strengths and limitations oning over extended textual inputs, guiding future improvements in language model development.
D D	DATA ANALYSIS AND LINEAR REGRESSION METHODOLOGY
This se	ction details the methodology used for the linear regression analyses in our study including data

This section details the methodology used for the linear regression analyses in our study, including data selection criteria, regression techniques, handling of low accuracy values, and confidence interval calculations. The approach aligns with the code used in our interactive analysis.

1188 D.1 TWO-PHASE ACCURACY BEHAVIOR AND DATA SELECTION 1189 Our experiments reveal a characteristic two-phase accuracy behavior as task complexity N increases: 1190 (1) Phase 1: High accuracy plateau where models perform well on simpler tasks, typically with accuracies 1191 close to 100%. 1192 1193 (2) **Phase 2**: Exponential decay of accuracy as task complexity exceeds a certain threshold $N_{\rm eff}$. 1194 To model the rate of accuracy decline in Phase 2, we perform linear regression on the natural logarithm 1195 of accuracy versus task complexity N. It is essential to focus on data from Phase 2 only, to avoid distortion 1196 from the Phase 1 plateau. 1197 1198 D.2 DATA SELECTION CRITERIA 1199 Selecting appropriate data ranges is crucial for accurate regression. We use different accuracy ranges for 1200 different models and experimental stages to focus on Phase 2 data: 1201 (1) Pretrained Models: Accuracy between [0.1,0.9]. (2) Train-from-Scratch Models (Pretraining Stage): 1202 Accuracy between [0.1,0.8]. (3) Fine-tuned Models: Accuracy between [0.2,0.8]. (4)Repeated Sampling 1203 **Experiments**: Accuracy between [0.02, 0.8]. 1204 1205 These ranges exclude Phase 1 data (high accuracy plateau) and avoid extremely low accuracies that can 1206 lead to numerical instability. 1207 D.3 LINEAR REGRESSION METHODOLOGY 1208 1209 For each model and experimental condition, we perform linear regression to fit: 1210 $\log(A) = aN + b$ 1211 where: 1212 1213 - A is the accuracy for task complexity N. 1214 - a is the Complexity Decay Factor (CDF), indicating the rate of exponential decay. 1215 - b is the Contextual Decay Offset (CDO), representing baseline performance at the onset of Phase 2. 1216 1217 D.4 THEORETICAL ERROR ANALYSIS FOR LINEAR RREGRESSION 1218 1219 We verify that the error in estimating the regression coefficients a and b in the model 1220 $\log(p_t) = at + b + \epsilon_t$ 1221 scales as 1222 $\operatorname{Error} \approx \operatorname{Error}_0 \times \left(\sum_{t=1}^k N_t\right)^{-1}$ 1223 1224 1225 where N_t is the number of samples for x = t, and Error₀ is a constant encapsulating design matrix properties and variance factors. 1226 1227 D.4.1 PROBLEM SETUP 1228 1229 For each $t \in \{1, 2, ..., k\}$, we sample N_t times. Each sample yields $y_{ti} \in \{0, 1\}$, and the average accuracy 1230 is computed as: 1231 $\hat{p}_t = \frac{1}{N_t} \sum_{i=1}^{N_t} y_{ti}.$ 1232 1233 1234 Assuming large N_t , \hat{p}_t follows approximately: $\hat{p}_t \sim \mathcal{N}\left(p_t, \frac{p_t(1-p_t)}{N_t}\right),$ 1235 1236 1237 where p_t is the true probability. The regression model assumes that $\log(p_t)$ is linearly related to t: 1239 $\log(p_t) = at + b$, 1240 and after transformation, the response variable $Y_t = \log(\hat{p}_t)$ satisfies: 1241 $Y_t = \log(p_t) + \epsilon_t,$



1285 D.4.4 STABILITY OF FITTED CURVE WITH RESPECT TO MORE EXAMPLES EVALUATION

Here we show what happens to the accuracy with the number of operations relationships if we increase from 50 examples per data point to 2000 examples. We select Qwen-2.5-7B-Instruct as the model of interest. The result is presented in Figure 9. We can see that the overall curve hasn't changed much. In fact, if we calculate the area under the curve for (a), it only changes from 1470 (with 50 examples) to 1490 (with 2000 examples): 1.3% difference. The result showcases that the pattern is stable very quickly.

1292 D.4.5 THEORETICAL ERROR ANALYSIS

1286

1291

1293

1294 The variances of the OLS estimates depend on the covariance structure of the errors. In linear regression, the covariance matrix of the estimates is:

$$\operatorname{Cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

1296 where X is the design matrix and σ^2 is the error variance. For simple linear regression, the variances are: 1297 $\operatorname{Var}(\hat{a}) = \sigma^2 \frac{1}{S_{tt} - \frac{S_t^2}{L}}, \quad \operatorname{Var}(\hat{b}) = \sigma^2 \left(\frac{1}{k} + \frac{\overline{t}^2}{S_{tt} - \frac{S_t^2}{L}}\right),$ 1298 1299 1300 where $\bar{t} = \frac{S_t}{h}$. 1301 D.4.6 VARIANCE OF ϵ_t 1302 1303 The variance of ϵ_t is $\sigma_t^2 = \frac{1-p_t}{n_t N_t}$. To simplify the analysis, assume approximate homoscedasticity: 1304 1305 $\sigma_t^2 \approx \frac{1 - \bar{p}}{\bar{p}\bar{N}}, \quad \text{where } \bar{p} = \frac{1}{k} \sum_{t=1}^k p_t, \bar{N} = \frac{\sum_{t=1}^k N_t}{k}.$ 1306 1307 Define the total sample size as $\mathcal{N} = \sum_{t=1}^{k} N_t$. Since $\bar{N} = \mathcal{N}/k$, the effective error variance scales as: 1308 1309 $\sigma^2 \propto \frac{1}{M}$ 1310 1311 1312 D.4.7 Scaling of Errors with \mathcal{N} 1313 $\operatorname{Var}(\hat{a}) \propto \frac{1}{\mathcal{N}}, \quad \operatorname{Var}(\hat{b}) \propto \frac{1}{\mathcal{N}}.$ Taking square roots to obtain standard errors: Using the approximation $\sigma^2 \propto \frac{1}{N}$, the variances of the OLS estimates become: 1314 1315 1316 1317 $\operatorname{SE}(\hat{a}) \propto \frac{1}{\sqrt{\mathcal{N}}}, \quad \operatorname{SE}(\hat{b}) \propto \frac{1}{\sqrt{\mathcal{N}}}.$ 1318 1319 1320 D.4.8 FINAL ERROR EXPRESSION 1321 1322 The final expression for the regression error is: 1323 $\operatorname{Error} \approx \operatorname{Error}_0 \times \left(\sum_{i=1}^k N_t\right)^{-1/2},$ 1324 1325 1326 where Error_0 encapsulates variability due to p_t and the properties of the design matrix (e.g., $S_{tt} - \frac{S_t^2}{L}$). This 1327 confirms the scaling of estimation error with the total sample size. 1328 D.4.9 CONFIDENCE INTERVAL CALCULATION 1330 1331 We calculate confidence intervals for the regression coefficients using their standard errors and the inverse 1332 error function (erf^{-1}) : 1333 1334 $CI = Coefficient \pm \left(SE \times \sqrt{2} \times erf^{-1}\left(\frac{C}{100}\right)\right)$ 1335 1336 1337 where C is the confidence level (e.g., 95%). 1338 1339 D.5 SUMMARY 1340 Our data analysis and regression methodology accurately capture the relationship between task complexity 1341 and model accuracy during the exponential decay phase. By carefully selecting data ranges, we ensure reliable 1342 regression results that provide valuable insights into models' abilities to handle complex reasoning tasks over extended contexts. 1344

1345 E FACTORS AFFECTING THE METRICS

E.1 PRETRAINING STRATEGIES

We investigate how different pretraining strategies influence the Complexity Decay Factor (CDF), Contextual Decay Offset (CDO), and especially the Effective Complexity (N_{eff}) in the FACTOR benchmark. By training models from scratch using various methodologies, we aim to understand their effects on the models' abilities to handle task complexity and context length. See Appendix F for training settings and Appendix G for training data distribution. Models were trained using the following strategies:

(1) Baseline: Trained directly on the FACTOR benchmark training set (regenerated to avoid data leakage), mirroring the benchmark's organization. The max context length is limited to 1000 tokens. (2) Naive Packed Training (*Packed*): Sequences in a mini-batch are concatenated to form longer training inputs. (3)
Question-Masked Training (*Masked*): The question portion (payload) is masked during loss computation.
(4) Clean Context Training (*Clean*): The model is exposed only to the payloads, without filler text. (5)
Packed Training with Diagonal Attention Mask (*Diagonal*): Similar to Packed, but with diagonal attention masks to prevent cross-sequence attention; position IDs are retained.

We evaluated the models on the FACTOR benchmark with zero filler context to assess their baseline performance. The results are presented in Table 11.

1361 1362 1363

1364 1365

1367

1369

Table 11: Effects of Pretraining Strategies on Model Performance (Zero Filler Context).

Model	CDF	CDO	$N_{\rm eff}$
Baseline	-0.3292	8.4685	25.72
Packed	-0.1339	2.6381	19.70
Masked	-0.3278	8.5196	25.99
Clean	-0.2672	7.0388	26.34
Diagonal	-0.3780	10.0887	26.69

1370 1371

The Effective Complexity (N_{eff}) indicates the maximum task complexity the model can handle before performance significantly degrades. Higher N_{eff} values reflect better performance over a broader range of complexities.

From Table 11, we observe: (1) The **Baseline**, **Masked**, **Clean**, and **Diagonal** models achieve similar N_{eff} 1376 values around 26, indicating they handle task complexities up to $N \approx 26$ effectively. (2) The **Packed** model 1377 has a lower $N_{\rm eff}$ of 19.70, suggesting worse performance over the range of task complexities. Naive packed 1378 training leads to poorer handling of complex tasks compared to other strategies. Despite a more negative 1379 CDF (-0.3780), the **Diagonal** model achieves the highest N_{eff} (26.69) and a higher CDO (10.0887). This 1380 means the model starts with higher baseline accuracy (due to higher CDO) and maintains decent performance 1381 for N between 25 and 33, extending the effective complexity range. **Packed Training** reduces $N_{\rm eff}$, leading 1382 to worse performance across task complexities. **Diagonal Training**, although it has a more negative CDF, extends $N_{\rm eff}$, improving performance for some higher N values. Therefore, while naive packed training 1384 negatively impacts the model's ability to handle complex tasks, the Diagonal strategy enhances performance at higher complexities, evidencing its benefit in extending the range over which the model maintains accuracy. 1385

1386

1387 E.2 POST-TRAINING STRATEGIES

1388 We examine how different fine-tuning strategies affect model performance on the FACTOR benchmark, 1389 particularly the CDF and CDO metrics, with a maximum context length of 16K tokens. Notably, although 1390 the maximum training length is 16K, the majority of sequences in the training distribution are less than 1391 8K. The models are trained based on the checkpoing of the Baseline model in the last section. We evaluate 1392 three approaches: (1) Gradual Sequence Length Increase (Course): Starting with shorter sequences and progressively increasing the length during fine-tuning. (2) Mixed-Length Sequence Training (Mixed): 1393 Training on sequences of varied lengths simultaneously. (3) Direct Long Sequence Training (Full): 1394 Fine-tuning directly on the maximum sequence length without gradual adaptation. 1395

Table 12 and 13 present the CDF and CDO metrics at different context lengths.

The fine-tuning strategies exhibit distinct effects on CDF and CDO:

Despite training up to 16K contexts, the *Course* and *Mixed* strategies significantly outperform the *Full* strategy, likely due to the training distribution containing mostly shorter sequences. Sudden exposure to full-length sequences in the *Full* strategy causes a distribution shift, leading to degraded performance.

1401 **Course** strategy results in models with stable CDO and degrading CDF, resembling Llama series models.

This indicates that the models maintain baseline performance but struggle more with task complexity as context length increases. **Mixed** strategy yields models with stable CDF and degrading CDO, similar to

Table 12: CDF at Different Context Lengths. (Values scaled by 10^{-2})

Model	0	1K	2K	4 K	8K	16K
Course	-7.91	-6.44	-8.28	-8.31	-10.29	-14.56
Mixed	-8.43	-6.66	-7.79	-6.86	-8.11	-11.25
Full	-12.48	-9.34	-17.28	-12.42	-5.76	-18.82

Table 13: CDO at Different Context Lengths.

Model	0	1K	2K	4 K	8K	16K
Course Mixed	$0.4775 \\ 0.3227$	$0.3539 \\ 0.1351$	$0.4916 \\ 0.2076$	0.3383 0.0080	$0.3392 \\ -0.0922$	$0.1038 \\ -0.3469$
Full	-0.3222	-0.4789	-0.2658	-0.3260	-0.7427	-0.2435

some pretrained models that handle task complexity consistently but whose baseline performance declines
with longer contexts, resembling models like gpt-40-mini. Full strategy leads to poor performance in both
metrics, due to abrupt changes in data distribution. Gradual adaptation to longer sequences during fine-tuning
helps models cope better with increasing context lengths and task complexities, mitigating failure modes
observed in pretrained models.

1427 E.3 REPEATED SAMPLING STRATEGIES

We analyze how increasing computational efforts during inference affects model performance, particularly focusing on the accuracy and the effective complexity N_{eff} . To investigate the impact of increased computational effort during inference, we employ **repeated sampling**. This approach measures the probability of correctly solving at least one instance out of t samples. We conduct experiments on models trained with the Gradual Increase strategy at different filler context lengths, as well as on the model pre-trained with a clean context (zero filler context).

1434Table 14 summarizes the linear regression results for different filler context lengths. The regression equation1435is given by:

$$\log(N_{\text{eff}}^{\text{clean}} - N_{\text{eff}}^{\text{model}}) = k \log(t) + k$$

Table 14: Linear Regression Results for Different Filler Context Length.

k (Slope)	h (Tratamaant)	
(S10pt)	0 (Intercept)	R-squared
-0.358287	2.865435	0.992356
-0.333536	2.670190	0.987057
-0.376972	2.865891	0.982553
-0.261895	2.999738	0.995347
-0.245732	3.242366	0.992648
-0.150702	3.524461	0.998234
	-0.358287 -0.333536 -0.376972 -0.261895 -0.245732 -0.150702	-0.3582872.865435-0.3335362.670190-0.3769722.865891-0.2618952.999738-0.2457323.242366-0.1507023.524461

Figure 10 illustrates the relationship between $N_{\rm eff}$ and the number of tries t.

Key observations from the results are: (1) For models fine-tuned with gradually increasing context lengths, increasing the number of samples t leads to a slow improvement in $N_{\rm eff}$, but it remains below the $N_{\rm eff}$ of the clean context model. (2) **Effect of Context Length**: As the filler context length increases, the value of kbecomes less negative, indicating a decrease in the rate of change. Simultaneously, the intercept b increases. (3) **Interpretation:** While increasing computational efforts during inference provides marginal benefits, it does not fully overcome the challenges posed by high task complexity and long contexts. There exists spontaneous bias within the model. For higher complexity tasks, the benefit of repeated sampling is diminishing, revealing that for certain questions, achieving the correct answer though repeated sampling is nearly impossible. This raises concerns about the model's generalization performance with respect to out-of-training-distribution complexity.



Figure 10: (a) exhibits the strong linear correlation between $log(N_{eff}^{clean} - N_{eff}^{model})$ and log(t). This implies a slow improvement in N_{eff} with the increase of number of samples, and it does not fully overcome the challenges posed by high task complexity and long contexts. (b) is the average accuracy as a function of the number of variables (N) for different sample sizes. The plot reveals that For higher complexity tasks, the benefit of repeated sampling is diminishing.

¹⁴⁷⁶ F TRAINING CONFIGURATIONS

In this section, we provide a concise overview of the training configurations used for developing our language
models evaluated on the FACTOR benchmark. Key configurations are summarized in Table 15, and additional
details are described to clarify the training setup.

Table 1	5:	Summary	of	Training	C	Configurations	5.
---------	----	---------	----	----------	---	----------------	----

1484	Configuration Flement	Setting
1485	Comguration Element	Setting
1486	Model Architecture	LlamaForCausalLM
1487	Number of Layers	12
1400	Number of Attention Heads	12
1400	Hidden Size	768
1489	Intermediate Size	3,072
1490	Vocabulary Size	424
1491	Maximum Position Embeddings	32,768
1492	RoPE Theta	500,000.0
1493	Data Type	bfloat16
1494	Peak Learning Rate	4e-4
1495	Batch Size	192
1496	Learning Rate Scheduler	Cosine decay with warmup ratio 0.01
1497	Number of Training Epochs	1
1498	Optimizer	AdamW

1499 1500

1475

1481 1482 1483

1501 F.1 MODEL ARCHITECTURE

We employ a Llama3 architecture configured to handle long sequences and complex reasoning tasks. The model consists of 12 Transformer layers, each with 12 attention heads, resulting in a hidden size of 768 (calculated as num_heads \times 64). The intermediate size is set to 3,072 (four times the hidden size), following common practice to enhance model capacity.

To effectively handle long contexts, we use Rotary Position Embeddings with a RoPE theta value of 500,000.0.
 This allows the model to capture positional information over extended sequences. We utilize FlashAttention-2 for efficient attention computation on long sequences, which improves training speed and reduces memory consumption.

1511 All computations are performed using bfloat16 precision to optimize memory usage and computational efficiency without significantly affecting model accuracy.

1512 F.2 TOKENIZER

A custom tokenizer is used, tailored to the specific needs of the FACTOR benchmark tasks. The tokenizer has a vocabulary size of 424 tokens.

Optimizer and Scheduler The AdamW optimizer is utilized with default parameters except for the learning rate. The cosine learning rate scheduler with warmup helps in gradually adapting the learning rate, reducing the risk of training divergence at the start.

1519 1520

1534 1535

1539

1544 1545

1557

1561

G TRAINING DATA GENERATION

This section details the generation of the training data used in both the pretraining and posttraining phases.
Our approach carefully controls task complexity and context length to effectively train the models for evaluating their reasoning abilities on the FACTOR benchmark.

1525 G.1 PRETRAINING DATASET

The pretraining dataset comprises 3 million synthetic Question-Solution pairs. It is designed to teach the model fundamental reasoning skills over a range of task complexities and shorter context lengths.

1529 1530 G.1.1 TASK COMPLEXITY DISTRIBUTION

To control task complexity, we vary the number of variables N in each synthetic example. The value of N is determined by:

$$N = \min(N_1, N_2), N_1, N_2 \sim \text{Uniform}(1, 29)$$

where Uniform(1,29) denotes a discrete uniform distribution over integers from 1 to 29 (excluding 30). By taking the minimum of two independently sampled values, we bias the distribution toward smaller values of N, emphasizing simpler tasks while still including examples with higher complexities up to N = 29.

1540 G.1.2 FILLER CONTEXT LENGTH DISTRIBUTION

The length of the filler context L_{pretrain} is sampled from a log-normal distribution to introduce variability in context lengths while maintaining manageable sequence sizes for pretraining. Specifically, we use:

$$\ln L \sim \mathcal{N}(\mu_{\ln L}, \sigma_{\ln L}^2), \quad L_{\text{pretrain}} = \exp(\ln L) - L_{\min}$$

1546 where: $\mu_{\ln L} = 5$, $\sigma_{\ln L} = 2.5$, $L_{\min} = 100$, $L_{\max} = 1000$

After sampling $\ln L$, we compute L, round it to the nearest integer, and clip it to the range $[L_{\min}, L_{\max}]$. We then subtract L_{\min} to adjust the lengths to start from zero. This results in a distribution of filler context lengths biased toward shorter lengths but including a range up to 900 tokens.

1551 1552 G.2 POST-TRAINING DATASET

The post-training dataset consists of 60,000 synthetic long-context Question-Solution pairs. It is used to fine-tune the pretrained model, enhancing its ability to handle extended contexts and complex reasoning tasks.

1555 We employ three different training strategies to investigate their effects on model performance:

- Gradually Increasing Context Length: The dataset is divided into four equal parts, each containing 15,000 examples. Each part corresponds to filler context lengths that are progressively increased. Specifically, the filler context length distribution in each stage is scaled by factors of 2, 4, 8, and 16 times the pretraining filler context length L_{pretrain}, respectively. By gradually increasing the context length, the model is incrementally adapted to handle longer sequences, reaching the maximum context length in the final stage.
- Mixture of Different Context Lengths: In this strategy, the examples from all four stages of the gradually increasing context length dataset are combined and shuffled. Training on this mixed dataset exposes the model to different context lengths simultaneously, which may encourage better generalization across various sequence lengths.

3. Training Directly with Full Context Length: The model is trained starting from a filler context length distribution scaled by $16 \times L_{\text{pretrain}}$. This approach tests the model's ability to handle long contexts without prior adaptation through shorter sequences.

G.2.1 TASK COMPLEXITY DISTRIBUTION

For all posttraining strategies, the number of variables N is sampled using the same method as in pretraining:

 $N = \min(N_1, N_2), N_1, N_2 \sim \text{Uniform}(1, 29)$

This ensures consistency in task complexity across both pretraining and posttraining datasets.

G.2.2 FILLER CONTEXT LENGTH GENERATION

For each stage in the gradually increasing context length strategy, the filler context lengths are generated by scaling the pretraining filler context lengths L_{pretrain} :

 $L_{\text{stage}} = s \times L_{\text{pretrain}}$

where s is the scaling factor for each stage (s=2,4,8,16).

The filler context lengths are generated using the same log-normal distribution as in pretraining but adjusted for the scaling factor. This adjustment ensures that the filler context lengths are appropriately scaled for each stage.

G.2.3 SYNTHETIC DATA GENERATION PROCESS

For each example in both the pretraining and posttraining datasets, we follow the same synthetic data generation steps as described in Section C, adjusting the filler context lengths according to the stage and strategy.

G.3 SUMMARY

By carefully generating training data with controlled task complexity and systematically varied context lengths, we aim to investigate the impact of different training strategies on the models' abilities to handle complex reasoning tasks over long contexts. The datasets are designed to provide comprehensive exposure to the challenges posed by increased sequence lengths and to facilitate a detailed analysis of model performance under various training conditions.