

Dynamic Cheatsheet: Test-Time Learning with Adaptive Memory

Mirac Suzgun^π Mert Yuksekgonul^π Federico Bianchi^γ Dan Jurafsky^π James Zou^{π,γ}

^πStanford University ^γTogether AI

Abstract

Despite their impressive performance on complex tasks, current language models (LMs) typically operate without retaining insights from previous attempts. Here, we present Dynamic Cheatsheet (DC), a lightweight framework that endows a blackbox LM with a persistent, evolving memory. Instead of repeatedly re-discovering the same solutions and mistakes, DC enables models to store and reuse accumulated strategies, code snippets, and general problem-solving insights at inference time. This test-time learning enhances performance substantially across a range of tasks without needing explicit ground-truth labels or human feedback. Leveraging DC, Claude 3.5 Sonnet’s accuracy more than doubled on AIME once it began retaining algebraic insights across questions. GPT-4o’s success on Game of 24 jumped from 10% to 99%, and DC achieved 9% gains on knowledge-demanding tasks like GPQA and MMLU-Pro. Crucially, DC’s memory is self-curated, focusing on concise, transferable snippets rather than entire transcripts, avoiding context ballooning. Unlike fine-tuning or static retrieval methods, DC adapts LMs’ problem-solving skills on the fly, without modifying their underlying parameters, and offers a practical method for continuously refining responses and cutting routine errors. Overall, our findings present DC as a promising approach for augmenting LMs with persistent memory.¹

1 Introduction

Modern large language models (LLMs) can tackle complex reasoning tasks, answer various questions, and generate extensive texts. Yet they still suffer from one critical limitation: once deployed, these models are fixed prior to deployment and typically retain no explicit or implicit memory of past questions, successes, or mistakes during inference.

¹✉ Correspondence to: msuzgun@stanford.edu. We release all our data, results, and code at <http://github.com/suzgunmirac/dynamic-cheatsheet>.

They approach each new problem *de novo*, often re-deriving the same insights—and re-committing the same errors. In contrast, human cognition stands on a foundation of incremental learning, continuously internalizing new experiences and solutions into a persistent mental model.

In this work, we present *Dynamic Cheatsheet* (DC), a simple and intuitive framework that endows black-box LLMs with a persistent, evolving memory at inference time. Rather than fine-tuning weights (for instance, through dynamic evaluation (Krause et al., 2019) or domain adaptation (Gururangan et al., 2020)) or retrieving facts from a massive *static* corpus (as in traditional retrieval-augmented generation systems (Guu et al., 2020; Zhang et al., 2024b)), DC dynamically curates a compact library of reusable strategies, solution sketches, and code snippets. Either before or after each query, DC enables the system to decide which lessons to store, what to discard, and how to refine existing entries—thus effectively “learning” from successes and failures. It is a flexible online-learning approach that enables a black-box LLM to improve itself without needing any explicit ground truth labels or human feedback.

The overall workflow of DC is intuitive and compelling. In one version of DC (DC-Cu.), when presented with a new query, the LM first consults its external memory to see if any prior insights, strategies or relevant model solutions have been stored. It then proposes a solution by combining the retrieved insights with its own internal reasoning capabilities. Upon generating an answer, it proceeds to a curation phase that updates the memory: If the solution seems to be correct, useful, or practical, DC codifies it in its memory for future use; if an error surfaces, DC may revise or prune faulty heuristics. This all happens without gradient updates, so computational overhead remains modest, and compatibility with black-box APIs (e.g., GPT or Claude) is fully preserved. See Figure 3.



Figure 1: Comparison of different baselines and Dynamic Cheatsheet (DC) variants on challenging reasoning benchmarks, including AIME exams and GPQA-Diamond. Baseline represents standard prompting with minimal guidance, while DC- \emptyset (a stronger baseline) contains explicit structured instructions for problem solving and Python execution but lacks a memory component. Our proposed DC-Cu and DC-RS variants incorporate an evolving memory for inference-time learning. Results (accuracy, %) demonstrate substantial improvements, with Claude 3.5 Sonnet gaining 27% on AIME 2024 and 30% on AIME 2025 under DC-Cu. In Game of 24, GPT-4o leaps from 10% (baseline) to 99% under DC-RS. Similarly, Claude 3.5 Sonnet’s accuracy more than doubles in Math Equation Solver, reaching 98%.

We tested DC across multiple challenging benchmarks and observed that it increases performance and reduces repetitive mistakes. On AIME 2024, Claude 3.5 Sonnet jumped from 23% to 50% accuracy, more than doubling its baseline score, by retaining algebraic and combinatorial insights. Likewise, it gained 30% accuracy on AIME 2025. Notably, these gains hold in knowledge-intensive tasks as well. On GPQA-Diamond, which tests specialized domain questions, DC lifted Claude by over 9%. In MMLU-Pro Engineering and Physics, it provided up to an 8% boost in performance by allowing the model to maintain a “toolkit” of formulas and general problem-solving patterns.

An even more striking and compelling example is the Game of 24, a puzzle that requires the solver to combine four digits into an arithmetic expression equaling 24. GPT-4o’s baseline performance (10%) increased to 99% under DC. Early in the test sequence, the model discovered that an efficient Python brute-force solver eliminated all manual guesswork. Once this snippet was stored, GPT-4o simply retrieved it for subsequent queries, avoiding manual arithmetic entirely.

Nonetheless, DC is not a panacea. We found that smaller models, such as GPT-4o-mini, benefit from DC in limited amounts. These models generate too few correct solutions in these challenging tasks in the first place, leaving the memory populated with flawed or incomplete strategies. Worse, they struggle to refine stored content. DC can empower models that can already produce high-quality outputs, but not fix foundational gaps in reasoning.

We also note that DC differs from naive “append the entire conversation history” in-context learning methods. Under DC, memory is continuously curated, focusing on succinct and useful knowledge over raw transcripts. This prevents ballooning context lengths (Liu et al., 2024a) and helps ensure that repeated retrieval remains tractable. Indeed, part of DC’s contribution is in designating a mechanism for selective, evolving retention—storing just enough to solve the next set of tasks without drowning in an ever-growing buffer. Cf. (Karpicke and Roediger III, 2008; Roediger and Butler, 2011; Karpicke and Blunt, 2011)

2 Dynamic Cheatsheet (DC) Methodology

DC includes an external, non-parametric memory that evolves in tandem with the LLM’s inference process. Rather than fine-tuning the weights, DC tracks successes and failures of the model at test time, then selectively stores heuristics, strategies, or insights that can guide the LLM in future instances. Notably, this approach respects the black-box nature of many commercial LLM APIs: no gradient-based updates are required, and the model’s core parameters remain untouched.

2.1 DC: Building Blocks and Iterative Loop

The DC framework consists of two core modules: *generation* and *curation*. Both modules can easily operate on top of the same LM (prompted differently) or on separate LMs.

Solution Generation with Memory. Let’s consider a sequence of inputs (x_1, x_2, \dots, x_n) , where

each $x_i \sim \mathcal{D}_{\text{test}}$ indicates a new query or problem posed to the model sampled from the same distribution $\mathcal{D}_{\text{test}}$. The distribution $\mathcal{D}_{\text{test}}$ is unknown to us. At the i -th step, the model is provided with both the new query x_i and the current memory state M_i , which captures knowledge gleaned from previous successes and failures. We denote the solution generator by Gen:

$$\tilde{y}_i = \text{Gen}(x_i, M_i) \quad (1)$$

Here, \tilde{y}_i is the candidate solution produced by the model. M_i helps condition the model to reuse or adapt previously stored solutions, insights, techniques, or heuristics.

Memory Curation Step. After the generator produces its answer \tilde{y}_i to x_i , the curator, Cur, updates the current memory:

$$M_{i+1} = \text{Cur}(M_i, x_i, \tilde{y}_i) \quad (2)$$

During memory curation, Cur mainly considers: (i) *the usefulness and generalizability of the newly produced answer* (i.e., if \tilde{y}_i is correct or provides valuable and generalizable insights, it is distilled into a form suitable for later reference), (ii) *refinement or removal of existing memory entries* (i.e., if an existing memory entry was incorrect or superseded by a more efficient or versatile strategy, Cur may remove or update it), and (iii) *clarity and compactness of the entire memory* (i.e., memory entries are consolidated to retain succinct, high-impact references and heuristics).

Cur does not have access to ground-truth labels; so, it has to assess the correctness and efficiency of the solutions by itself before updating the memory. In our experiments, we instruct a single model to perform this crucial step. Yet, in practice, Cur can be implemented as a series of steps that instruct multiple tools and models, through different prompts, to verify the validity and efficiency of the solution and to transform the raw solution text into even more generalizable, reliable, and efficient strategies, insights, and code snippets.

We refer to this version of DC above as **DC-Cu** (short for DC-Cumulative). Under DC-Cu, the system first performs solution generation based on the current memory (Eqn. 1) and then updates the memory (Eqn. 2), by cumulatively expanding and refining the memory items thus far. Unlike DC-RS, which is discussed in the next part, DC-Cu, does not contain a retrieval component, however.

2.2 DC with Retrieval & Synthesis (DC-RS)

DC-Cu has two potential drawbacks. *First*, it updates the memory after processing an input query, rather than refining it before generating a response. This means the model lacks the opportunity to incorporate new insights from the current query while reasoning through its solution. *Second*, DC-Cu does not store or revisit past input-output pairs unless explicitly retained in memory. This omission prevents the model from directly retrieving and leveraging historical responses, which can be particularly valuable in benchmarks covering diverse topics or domains (e.g., GPQA-Diamond).

To address these issues, **DC-RS** modifies the sequence of memory updates and introduces a retrieval mechanism, Retr, into the curation process. Retr allows the model to retrieve the most relevant past input-output pairs from its knowledge base. By refining the memory before responding and retrieving prior cases when needed, DC-RS enhances the model’s adaptability and reasoning efficiency.

DC-RS first retrieves² top- k most similar inputs, along with their model-generated outputs, from previously seen examples, which we denote by $R_i^{(k)}$ (or simply R_i).³ It then passes these select examples, R_i , along with the most recent memory content, M_{i-1} , to the curator to update the memory, that is to get M_i . Finally, it uses the generator to produce \tilde{y}_i , given x_i and M_i .

$$R_i = \text{Retr}(x_i, \{(x_j, \tilde{y}_j)\}_{j < i}, k) \quad (3)$$

$$M_i = \text{Cur}(M_{i-1}, x_i, R_i) \quad (4)$$

$$\tilde{y}_i = \text{Gen}(x_i, M_i) \quad (5)$$

2.3 Baselines

(1) Baseline prompting (BL). This *vanilla* method, with minimal instructions, prompts a model without any memory or retrieval mechanism. It reflects traditional one-off inference.⁴

(2) DC- \emptyset (empty memory). This DC baseline always keeps the memory empty to isolate the effect of memory curation.⁵ DC- \emptyset allows us to measure

²Using OpenAI’s text-embedding-3-small

³We set k to 3 in all our experiments. (other k values did not show significant gains.)

⁴Refer to Figure 9 for the full prompt used in BL. We tested 0-shot CoT (Kojima et al., 2022) in our experiments, but it did not yield any gains (Arcuschin et al., 2025); hence, we did not include it as a baseline method in our experiments.

⁵We adopt the generator prompt template used in DC-RS, namely Figure 10, for DC- \emptyset , though we replace the memory placeholder with the text “(empty cheatsheet)”.

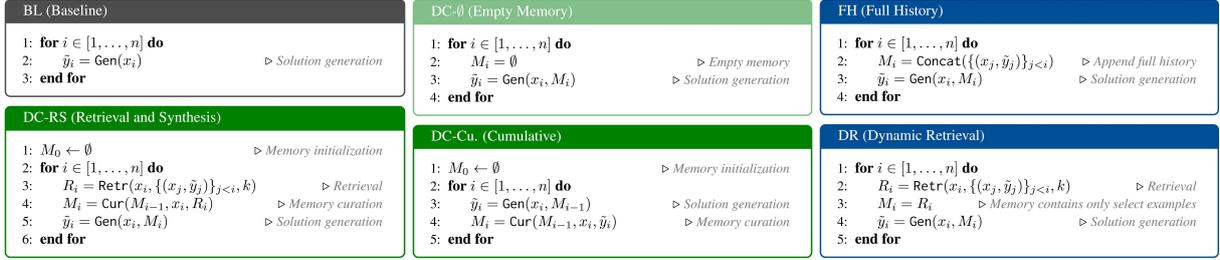


Figure 2: Dynamic Cheatsheet (DC)-based approaches. Gen represents the solution generator model, Cur the memory curator, and Retr the retriever. We use the same LLM for generation and curation with different prompts. Retrieval ranks historical inputs by cosine similarity, selecting the most relevant examples and solutions.

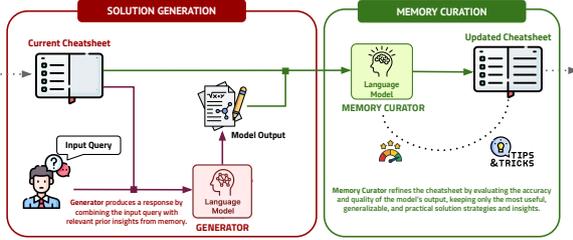


Figure 3: Illustration of *Dynamic Cheatsheet* (DC-Cu).

how much performance gain comes purely from storing and reusing knowledge over time. While there is no continuous knowledge storage or strategy reuse, this method follows the instructions in Figure 10 and is a strong baseline.

(3) **Full-History Appending (FH)**. This is a naive approach that appends the entire conversation history to the model input without any curation or truncation.⁶ FH can exceed context-window limits and include redundant or low-value information, but nonetheless, it provides a useful comparison for methods that actively curate content.⁷

(4) **Dynamic Retrieval (DR)**. This uses retrieval but no curation. Specifically, for each new query, it retrieves the most similar past interactions and directly pastes them, *verbatim*, into the prompt. DR can help the model see relevant input-output pairs but not directly codify any abstract.⁸

Figure 2 contains pseudocodes of all the primary methods and baselines considered in this paper.

3 Experimental Setup

To rigorously evaluate DC’s efficacy, we focus on challenging tasks where contemporary state-of-the-art LLMs, such as GPT-4o and Claude 3.5, still

⁶We consider and test this baseline only on AIME 2024 and AIME 2025, which are relatively small in their size (each contains 30 examples) compared to other benchmarks.

⁷We use the generator prompt template in Figure 10 again, but include the entire raw input-output pairs from the previous steps in the memory.

⁸FH is similar to DR, but we include only a select (most relevant) input-output pairs in the memory content.

face limitations. Rather than evaluating on benchmarks where performance is near saturation (e.g., BBH (Suzgun et al., 2023b), MGSM (Shi et al., 2023), GSM8K (Cobbe et al., 2021)), we prioritize tasks that demand multi-step reasoning, heuristic search, strategic adaptation, and cumulative learning—tasks where iterative memory refinement can yield tangible gains over time.⁹

Overall, the selected datasets include algorithmic, logical, and domain-specific reasoning tasks, each chosen to stress-test the model’s ability to refine its reasoning over time.

(a) **AIME 2020–2025 Exam Questions**: AIME is a prestigious high-school competition featuring complex problems across algebra, combinatorics, number theory, geometry, and probability. These questions require deep mathematical reasoning and multi-step problem-solving. We consider three subsets: AIME 2024 (30 questions), AIME 2025 (30 questions), and AIME 2020–2024 (133 questions).

(b) **GPQA-Diamond** (Rein et al., 2024): GPQA-Diamond contains 198 expert-validated questions across sciences, including biology, chemistry, and physics. These questions were correctly answered by domain experts but often missed by non-experts, making them ideal for evaluating DC’s ability to handle complex, multi-hop reasoning tasks.

(c) **Game of 24** (Yao et al., 2023; Suzgun and Kalai, 2024): A heuristic-driven arithmetic challenge where the objective is to form an expression that evaluates to 24 using four given numbers exactly once. For instance, if the input values were “7 7 8 11,” one valid answer would be “8*(7+7-11).” This task emphasizes systematic search, strategic reasoning, and pattern recognition. We use the 100 examples from (Suzgun and Kalai, 2024).

(d) **Math Equation Balancer**: Focused on elementary arithmetic reasoning, this dataset requires the model to complete equations by inserting the

⁹We release all the original input-output pairs in our codebase and in the Supplementary Materials.

appropriate operators to form valid expressions. The task emphasizes the sequential placement of operators, as illustrated by the example “1 ? 2 ? 3 = 6,” where the model must identify the correct operators to satisfy the equation (“1 + 2 + 3 = 6” or “1 * 2 * 3 = 6”). We compiled 250 arithmetic expressions for this task.

(e) **MMLU-Pro Physics & Eng.** (Wang et al., 2024b): This is a professional-level subset of the MMLU benchmark focused on physics and engineering. All questions are presented in a multiple-choice form. The original dataset contains 1,299 physics and 969 engineering questions. We sampled 250 questions from each subset.

We evaluate the efficacy of DC across a range of models. Our selection includes both recent LLMs such as GPT-4o and Claude 3.5 Sonnet and their smaller-scale counterparts (GPT-4o-mini and Claude 3.5 Haiku), as well as reasoning models such as DeepSeek R1.

4 Main Results

4.1 DC enables test-time learning and reduces repetitive errors

One of the most compelling illustrations of DC’s capabilities emerges from the Game of 24 task. As seen in Table 1, GPT-4o’s baseline accuracy on this arithmetic puzzle was just 10%. Under DC-RS, its accuracy increased to 99%, illustrating DC’s capacity for test-time learning and iterative refinement. Early in the task sequence, GPT-4o discovered a reliable, Python-based brute-force method to solve Game of 24 and later on recognized the repetitive structure of the problem. The model then encoded this approach into its memory. Once established, GPT-4o consistently retrieved and applied the more or less same Python solution for subsequent examples, leading to rapid and accurate results.

The performance under DC- \emptyset further highlights the positive impact of memory curation and retrieval. DC- \emptyset uses the same core generator but keeps the memory empty, thus lacking the mechanism to store and reuse solutions. The large gap between 19% (DC- \emptyset) and 99% (DC-RS) confirms that effective memory usage, in which past solutions are retrieved and generalized, is the main driver of GPT-4o’s transformation from ad-hoc solver to near-perfect performer in Game of 24.

In contrast, Claude 3.5 Sonnet showed marginal gain, moving from 12% to 14%. Despite DC’s scaffolding, Claude didn’t internalize a general ap-

proach but instead continued to rely on manual arithmetic solutions. This underscores that while DC provides the framework for test-time adaptation, its success hinges on the model’s capacity to identify and encode robust, reusable strategies.

4.2 DC provides substantial gain across various challenging reasoning tasks

Beyond Game of 24, DC yielded significant gains across a range of complex tasks. See Table 1. The results below illustrate how iterative solution reuse can be helpful in complex reasoning problems.

AIME Exam Problems. The AIME exams provided some of the most dramatic improvements under DC. For Claude 3.5 Sonnet, performance on AIME 2020–2024 surged from 6.7% to 40.6% under DC-RS. A similar upward trend appeared on AIME 2024 (23.3% to 50.0%) and AIME 2025 (6.7% to 36.7%) under DC-Cu. DC-Cu, where the model curates memory after processing the input and does not involve a retrieval stage, also proved potent in recent exam sets, achieving highest accuracy scores in AIME 2024 and 2025. GPT-4o also showed some noteworthy gains. Its AIME 2024 performance raised from 20.0% to 40.0% under DC-RS, while its AIME 2025 score climbed from 6.7% to 20.0%. These boosts suggest that structured test-time-produced memory can help tackle difficult math problems.

GPQA-Diamond. On GPQA-Diamond, Claude 3.5 Sonnet improved from 59.6% to 68.7% under DC-RS, a robust 9.1% gain purely from test-time adaptation. DR (63.6%) demonstrated that retrieval alone helps, but the further jump to 68.7% highlights how memory curation and synthesis can yield additional benefits. By contrast, GPT-4o experienced only a slight increase from 57.1% to 58.1% with DC-RS; our quantitative analysis of the model’s outputs and memory showed us that retrieval can, in some cases, introduce confusion, especially if suboptimal examples are recalled. This contrast between different models underscores how the success of retrieval-based adaptation partly depends on model-specific generation and curation capabilities.

Math Equation Balancer. As Table 1 shows, the baseline performance for Claude 3.5 Sonnet (44.8%) rose to 98–100% with DC-RS and DC-Cu, while GPT-4o similarly improved from 50.0% to near-perfect accuracy (99–100%). As observed in Game of 24, the models quickly learned an algorithmic or Python-based balancing routine, stored

Tasks	Claude 3.5 Sonnet					GPT-4o				
	BL	DC- \emptyset	DR	DC-Cu.	DC-RS	BL	DC- \emptyset	DR	DC-Cu.	DC-RS
AIME 2024	23.3	36.7	43.3	50.0	46.7	20.0	36.7	26.7	36.7	40.0
AIME 2025	6.7	23.3	23.3	36.7	30.0	6.7	10.0	10.0	16.7	20.0
AIME 2020–24	6.7	30.1	39.1	38.4	40.6	9.8	24.1	24.1	20.3	24.8
Game of 24	12.0	10.0	11.0	14.0	14.0	10.0	19.0	6.0	93.0	99.0
GPQA Diamond	59.6	60.1	63.6	61.1	68.7	57.1	57.1	55.1	58.1	57.1
Math Eqn. Balancer	44.8	56.4	60.4	100	97.8	50.0	88.0	100	100	99.2
MMLU Pro Eng.	61.2	57.2	65.2	66.8	67.6	53.2	51.6	48.8	44.0	51.2
MMLU Pro Physics	74.0	75.6	80.4	77.6	82.0	75.6	70.8	75.6	70.4	75.2

Table 1: Performance comparison of Dynamic Cheatsheet (DC) for Claude 3.5 Sonnet and GPT-4o across multiple benchmarks. **BL** (Baseline): standard inference without memory; **DC- \emptyset** (Empty Memory): includes structured problem-solving and explicit tool-use instructions but no memory; **DR** (Dynamic Retrieval): uses retrieval but lacks memory updates; **DC-Cu** (Cumulative Memory): iteratively accumulates model solutions but lacks retrieval; and **DC-RS** (Retrieval & Synthesis): combines retrieval with memory refinement/synthesis. These results highlight accuracy gains under DC: Claude Sonnet’s AIME 2024 accuracy jumps by 27% under DC-Cu, and GPT-4o’s Game of 24 accuracy leaps from 10% to 99% under DC-RS.

it in external memory, and repeatedly retrieved it, achieving exceptional consistency once the core method was established.

MMLU-Pro Tasks. For MMLU-Pro Eng. and Physics, Claude exhibited consistent gains, rising by up to 8.0% in Physics (from 74% to 82%). Our analysis of the memory shows that Claude temporarily stored and retrieved compact “reference guides” on engineering and physics principles, which might have proved beneficial for thematically similar questions. GPT-4o, on the other hand, observed slight decreases from the baseline on these tasks, suggesting that domain complexity and baseline knowledge gaps may attenuate DC’s benefits if curated memory is less reliable or consistent.

4.3 Memory curation fosters generalization

Whereas FH simply appends every previous dialogue turn into the prompt, DC actively filters and synthesizes high-value content. As shown in Table 2, Claude under FH reached 26.7% accuracy in 2024 questions, while DC methods hit 50.0%. Similarly, GPT-4o managed a baseline of 20.0% but fell to 6.7% using FH, in direct contrast to 40.0% with DC-RS. Excessive uncurated input-output pairs can not only overwhelm the model’s context window, dilute crucial insights and hamper retrieval efficiency, but also significantly increase inference costs over time. On the other hand, DC’s selective memory curation ensures that problem-solving tips or code snippets remain readily accessible without clutter, thus facilitating more robust and consistent improvements across consecutive queries.

4.4 DC fosters efficient tool usage

A successful behavior under DC is the LLMs’ inclination toward code generation to handle computationally intensive tasks. GPT-4o’s near-complete

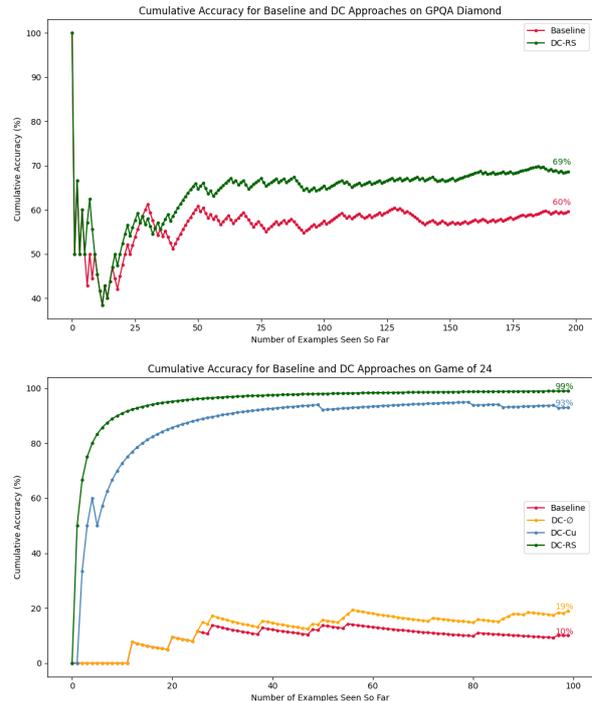


Figure 4: Cumulative performance progression under DC for GPQA-Diamond (left) and Game of 24 (right). In GPQA-Diamond, Claude 3.5 Sonnet steadily improves as it accumulates relevant knowledge snippets (first few points are noisy as y measures cumulative accuracy). Meanwhile, in Game of 24, GPT-4o rapidly gets near-perfect performance once it recognizes a Python-based solution.

reliance on Python scripts for Game of 24 exemplifies this shift. Rather than performing manual arithmetic repeatedly, GPT-4o recognized that code-based brute force is more systematic. It generated, stored, and iteratively refined a Python function that tested permutations of numbers and operations, allowing it to solve each instance of Game of 24 with high accuracy. This inclination toward automation illustrates DC’s potential to nurture efficient tool-usage: the capacity to recognize when external tools (e.g., Python, symbolic math engines, or

dedicated solvers) are more robust than internally verbalized chain-of-thought calculations. While we restricted the scope of tool usage to Python interpreter in this study, future expansions could easily explore a broader suite of tools, potentially amplifying LLM performance in specialized domains such as computational biology or legal research.

4.5 Model scale impact DC effectiveness

Our current results indicate that the effectiveness of DC is strongly tied to the model’s scale and underlying generative capacity. While Claude 3.5 Sonnet and GPT-4o showed notable gains across multiple tasks under DC, Claude 3.5 Haiku and GPT-4o-mini, showed limited and inconsistent gains. Table 4 (Appendix), for instance, shows that Claude 3.5 Haiku achieved moderate gains under DC, with its accuracy on AIME 2024 rising from 10.0% (baseline) to 36.7% under DC-Cu. Gains on AIME 2025 were weaker, reaching only 13.3% under DC- \emptyset and DC-Cu. Interestingly, GPQA-Diamond saw an improvement from 43.4% to 49.0% under DC-RS, suggesting that retrieval-based adaptation still provides utility in smaller models.

That said, GPT-4o-mini (Table 4) showed even smaller gains, with some variants leading to slight declines in performance. On AIME 2024, DC- \emptyset provided a 20.0% boost, but both DC-Cu and DC-RS performed worse than baseline. AIME 2025 showed a minor improvement, peaking at 16.7% under DC-RS. On GPQA-Diamond, GPT-4o-mini’s performance, however, remained largely stagnant or slightly declined under memory-based adaptation, suggesting that it struggled to leverage stored information effectively. These imply two drawbacks of smaller models under DC:

(a) **Generative competence.** For DC to be effective, the base model must produce correct solutions with sufficient frequency to populate the memory with high-quality, reusable strategies. Smaller models, such as GPT-4o-mini and Claude 3.5 Haiku, generate correct solutions less reliably, leading to a sparse or low-quality memory repository. So, iterative refinement stalls because the stored knowledge consists mostly of incorrect or partial attempts.

(b) **Contextual and memory curation limitations.** Smaller models struggle with long-context understanding/generation and memory retrieval, leading to inefficient or irrelevant memory usage. Unlike their larger counterparts, which can more effectively retrieve and synthesize solutions from stored heuristics, smaller models often fail to re-

trieve the most relevant past solutions or misapply retrieved knowledge to new problems. This results in inconsistent performance under DC-RS, particularly in tasks requiring complex reasoning or strategic adaptation.

4.6 Test-time task similarity and example ordering can amplify DC’s overall impact

Another central insight is that DC thrives when test examples share structural similarities. In both Game of 24 and Math Equation Balancer, once GPT-4o identified an efficient solution, it reused it consistently for subsequent tasks. Similarly, in AIME, discovering a geometry or combinatorics strategy allowed for easy transfer across questions of analogous structure. Consequently, tasks arranged to present *related* questions early may accelerate and improve the model’s test-time learning. This suggests that *curriculum-style* learning (Bengio et al., 2009), where simpler or archetypal problems are presented first to build a repository of valid heuristics, may potentially bootstrap performance. Cf. (Lopez-Paz and Ranzato, 2017; Zelikman et al., 2022; Chen et al., 2024)

5 Limitations & Additional Discussion

Reasoning and information efficiency. One key insight is that DC reduces the need to “reinvent the wheel” for each query. By encoding and reusing well-established techniques (e.g., Python-based solving for Game of 24), models can bypass repeated rediscovery of the same strategies. This significantly cuts down reasoning overhead and token usage in subsequent queries, though the initial cost of discovering a robust approach and curating it remains non-trivial.

DC performs better than majority voting (MV). To test if DC provides advantages over conventional MV at inference, we also tested Sonnet on AIME 2024 and 2025 using both approaches. MV, which selects the most common answer from three independent generations, yielded no improvements over single-shot inference. As seen in Table 3, on AIME 2024, MV performed identically to the baseline (23.3%), while on AIME 2025, it remained at 6.7%, offering no tangible gain. Even with DC- \emptyset , MV slightly underperformed (33.3% vs. 36.7%). In contrast, DC-Cu outperformed MV, reaching 50.0% on AIME 2024 and 36.7% on AIME 2025. Unlike MV, which *passively* aggregates outputs, DC *actively* refines knowledge over time, eliminating errors and improving solution

quality. This confirms that memory-based adaptation is far more effective than simple statistical voting in complex reasoning tasks.

Clustering of errors and corrections. Our experiments suggest that errors and their corrections often cluster in a latent embedding space. *See* Figure 7. Once a model acquires a high-quality heuristic for a cluster of related queries, it can apply this knowledge to tightly embedded neighbors. However, faulty heuristics that slip into memory can be equally amplified. Ensuring that the memory remains “clean” thus requires careful curation and, if necessary, pruning to avoid propagating erroneous strategies.

Transferability of memory content across models. We also observed that larger models, such as Claude 3.5 Sonnet and GPT-4o, can sometimes produce higher-quality strategies that, in principle, could benefit smaller models if the memory is transferred. However, if a smaller model lacks the generative capacity to interpret or refine those strategies correctly, its performance can stall or degrade. In our ablation experiments, we observed mixed results. This indicates that memory entries, while helpful, cannot fully compensate for inadequate base capability.

Long-context generation versus understanding. Most large LLMs excel at *processing* lengthy inputs but struggle to *generate* comparably long¹⁰ and well-organized outputs. DC’s memory curation after each query can demand precise reproduction or modification of prior knowledge. We observed instances where the model merely references or abbreviates the existing memory (e.g., “Previous content [...] preserved”) instead of explicitly rewriting it. Such truncated memory updates can reduce the quality of stored heuristics over time. Potential solutions include maintaining a structured, external database that the LM can reference without regenerating large swaths of text each time.

Retrieval bottlenecks and noise. While retrieval-based variants (e.g., DC-RS) can substantially improve accuracy, poorly filtered retrieval mechanisms can introduce confusion, particularly when presented with highly diverse or loosely related queries. For example, in our experiments, GPT-4o’s performance occasionally dipped in GPQA-Diamond due to suboptimal retrieval choices. This underscores the importance of robust retrieval methods (e.g., dense vector search,

advanced ranking algorithms) that can reliably surface higher quality exemplars or heuristics while suppressing irrelevant or contradictory texts.

Hierarchical and modular memory. As LLM deployments scale, specialized domains may benefit from subdividing or hierarchically organizing memory. For instance, a system could maintain separate curated memories for topics like biology or physics, each updated by a specialized retrieval or curator. This may reduce the load on a unified memory store and help isolate errors within their respective domains, with the goal of further improving the clarity and reliability of retrieved heuristics.

Time and token complexity. Although DC requires memory curation after each query, it optimizes efficiency over time by reducing redundant computation and token usage.¹¹ As the model retrieves and refines solutions, memory use becomes a net gain rather than a cost. However, its sequential nature still poses challenges for large parallel or batch tasks requiring independent inference.

Smaller or more specialized models and R1 experiments. Finally, we note that smaller models, such as GPT-4o-mini, show limited gains under DC, as seen in Table 4. Additional experiments with “R1” models such as DeepSeek R1 and o1 similarly showed minimal or inconsistent improvements. In these cases, these models’ generative ability appears too restricted to produce reliable strategies for storage or to interpret retrieved heuristics effectively. The solutions were far too verbose and long. Without sufficiently accurate and efficient base solutions, memory curation cannot yield substantial gains. This limitation ties back to the core premise that effective DC demands a capable foundation model to seed and refine the curated knowledge.

Overall, DC offers a useful and practical framework for continuous, test-time learning in LLMs. Our findings emphasize the synergy between model capacity and memory curation, the importance of structural task similarity and retrieval precision, and the benefits of offloading repeated computations to flexible external stores (e.g., Python scripts). At the same time, alternative mechanisms (e.g., specialized sub-memories or adaptive example ordering) and more sophisticated retrieval techniques (e.g., topological clustering) remain promising directions for further research.

¹⁰*See, e.g., (Liu et al., 2024b).*

¹¹On AIME 2024, Sonnet averaged 370 tokens under BL, 494 under DC- \emptyset , 1035 under DC-RS, and 1831 under DC-Cu.

Acknowledgments

We thank Batu El, Sabri Eyuboglu, Tayfun Gur, Adam T. Kalai, Luke Melas-Kyriazi, Emily Shen, Jake Silberg, Elana Simon, and Kyle Swanson for their helpful comments and suggestions. We also thank the members of the James Zou Lab at Stanford for their feedback in the early stages of this project. Suzgun gratefully acknowledges the support of an HAI-SAP Fellowship.

Disclaimer: Portions of this article were edited with the assistance of the GPT-4 and Claude 3.5 Sonnet models, mainly to shorten selected sections and to improve overall clarity, style, and flow. All model outputs were carefully vetted by the authors, and any revisions adopted were limited to word- or sentence-level edits. No substantive changes were made to the article.

References

- Shun-Ichi Amari. 1998. Natural gradient works efficiently in learning. *Neural computation*, 10(2):251–276.
- Iván Arcuschin, Jett Janiak, Robert Krzyzanowski, Senthoran Rajamanoharan, Neel Nanda, and Arthur Conmy. 2025. [Chain-of-thought reasoning in the wild is not always faithful](#). In *Workshop on Reasoning and Planning for Large Language Models*.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations*.
- Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and 1 others. 2024. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, and 1 others. 2022. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR.
- Léon Bottou and Yann Cun. 2003. Large scale online learning. *Advances in neural information processing systems*, 16.
- Léon Bottou and Yann Le Cun. 2005. On-line learning for very large data sets. *Applied stochastic models in business and industry*, 21(2):137–151.
- Malik Boudiaf, Romain Mueller, Ismail Ben Ayed, and Luca Bertinetto. 2022. Parameter-free online test-time adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8344–8353.
- Aydar Bulatov, Yury Kuratov, and Mikhail Burtsev. 2022. Recurrent memory transformer. *Advances in Neural Information Processing Systems*, 35:11079–11091.
- Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. 2024. Self-play fine-tuning converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Tao Feng, Pengrui Han, Guanyu Lin, Ge Liu, and Jiaxuan You. 2024. [Thought-retriever: Don’t just retrieve raw data, retrieve thoughts](#).
- Yukun Feng, Feng Li, Ziang Song, Boyuan Zheng, and Philipp Koehn. 2022. Learn to remember: Transformer with recurrent memory for document-level machine translation. *arXiv preprint arXiv:2205.01546*.
- Olga Golovneva, Sean O’Brien, Ramakanth Pasunuru, Tianlu Wang, Luke Zettlemoyer, Maryam Fazel-Zarandi, and Asli Celikyilmaz. 2023. Pathfinder: Guided search over multi-step reasoning paths. *arXiv preprint arXiv:2312.05180*.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujia Yang, Nan Duan, and Weizhu Chen. 2023. Critic: Large language models can self-correct with tool-interactive critiquing. *arXiv preprint arXiv:2305.11738*.
- Alex Graves. 2013. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*.
- Alex Graves, Greg Wayne, and Ivo Danihelka. 2014. Neural turing machines. *arXiv preprint arXiv:1410.5401*.
- Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. 2020. Don’t stop pretraining: Adapt language models to domains and tasks. *arXiv preprint arXiv:2004.10964*.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR.

- Zexue He, Leonid Karlinsky, Donghyun Kim, Julian McAuley, Dmitry Krotov, and Rogerio Feris. 2024. Camelot: Towards large language models with training-free consolidated associative memory. *arXiv preprint arXiv:2402.13449*.
- Armand Joulin and Tomas Mikolov. 2015. Inferring algorithmic patterns with stack-augmented recurrent nets. *Advances in neural information processing systems*, 28.
- Jeffrey D Karpicke and Janell R Blunt. 2011. Retrieval practice produces more learning than elaborative studying with concept mapping. *Science*, 331(6018):772–775.
- Jeffrey D Karpicke and Henry L Roediger III. 2008. The critical importance of retrieval for learning. *science*, 319(5865):966–968.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick SH Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In *EMNLP (1)*, pages 6769–6781.
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. [Generalization through memorization: Nearest neighbor language models](#). In *International Conference on Learning Representations*.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.
- Ben Krause, Emmanuel Kahembwe, Iain Murray, and Steve Renals. 2019. Dynamic evaluation of transformer language models. *arXiv preprint arXiv:1904.08378*.
- Angeliki Lazaridou, Elena Gribovskaya, Wojciech Jan Stokowiec, and Nikolai Grigorev. 2023. [Internet-augmented language models through few-shot prompting for open-domain question answering](#).
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474.
- Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024a. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173.
- Xiang Liu, Peijie Dong, Xuming Hu, and Xiaowen Chu. 2024b. Longgenbench: Long-context generation benchmark. *arXiv preprint arXiv:2410.04199*.
- Yuejiang Liu, Parth Kothari, Bastien Van Delft, Baptiste Bellot-Gurlet, Taylor Mordan, and Alexandre Alahi. 2021. Ttt++: When does self-supervised test-time training fail or thrive? *Advances in Neural Information Processing Systems*, 34:21808–21820.
- Jieyi Long. 2023. Large language model guided tree-of-thought. *arXiv preprint arXiv:2305.08291*.
- David Lopez-Paz and Marc’Aurelio Ranzato. 2017. Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30.
- Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. 2023. Chameleon: Plug-and-play compositional reasoning with large language models. *Advances in Neural Information Processing Systems*, 36:43447–43478.
- Aman Madaan, Niket Tandon, Peter Clark, and Yiming Yang. 2022. Memory-assisted prompt editing to improve gpt-3 after deployment. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2833–2861.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, and 1 others. 2023. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36:46534–46594.
- Michael McCloskey and Neal J Cohen. 1989. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier.
- Tomas Mikolov, Martin Karafiát, Lukas Burget, Jan Cernocký, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model. In *Interspeech*, volume 2, pages 1045–1048. Makuhari.
- Tsendsuren Munkhdalai, Alessandro Sordani, Tong Wang, and Adam Trischler. 2019. Metalearned neural memory. *Advances in Neural Information Processing Systems*, 32.
- Shuaicheng Niu, Jiayang Wu, Yifan Zhang, Yaofu Chen, Shijian Zheng, Peilin Zhao, and Minghui Tan. 2022. Efficient test-time model adaptation without forgetting. In *International conference on machine learning*, pages 16888–16905. PMLR.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, and 1 others. 2023. Toollm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*.
- Amal Rannen-Triki, Jorg Bornschein, Razvan Pascanu, Marcus Hutter, Andras György, Alexandre Galashov, Yee Whye Teh, and Michalis K Titsias. 2024. Revisiting dynamic evaluation: Online adaptation for large language models. *arXiv preprint arXiv:2403.01518*.

- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2024. **GPQA: A graduate-level google-proof q&a benchmark**. In *First Conference on Language Modeling*.
- Henry L Roediger and Andrew C Butler. 2011. The critical role of retrieval practice in long-term retention. *Trends in cognitive sciences*, 15(1):20–27.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551.
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023. **Hugging-GPT: Solving AI tasks with chatGPT and its friends in hugging face**. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Freda Shi, Daniel Fried, Marjan Ghazvininejad, Luke Zettlemoyer, and Sida I Wang. 2022. Natural language to code translation with execution. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3533–3546.
- Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2023. **Language models are multilingual chain-of-thought reasoners**. In *The Eleventh International Conference on Learning Representations*.
- Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2024. **REPLUG: Retrieval-augmented black-box language models**. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 8371–8384, Mexico City, Mexico. Association for Computational Linguistics.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652.
- Yu Sun, Xinhao Li, Karan Dalal, Jiarui Xu, Arjun Vikram, Genghan Zhang, Yann Dubois, Xinlei Chen, Xiaolong Wang, Sanmi Koyejo, and 1 others. 2024. Learning to (learn at test time): Rnns with expressive hidden states. *arXiv preprint arXiv:2407.04620*.
- Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. 2020. Test-time training with self-supervision for generalization under distribution shifts. In *International conference on machine learning*, pages 9229–9248. PMLR.
- Dídac Surís, Sachit Menon, and Carl Vondrick. 2023. Vipergpt: Visual inference via python execution for reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11888–11898.
- Mirac Suzgun, Sebastian Gehrmann, Yonatan Belinkov, and Stuart M Shieber. 2019. Memory-augmented recurrent neural networks can learn generalized dyck languages. *arXiv preprint arXiv:1911.03329*.
- Mirac Suzgun and Adam Tauman Kalai. 2024. **Meta-prompting: Enhancing language models with task-agnostic scaffolding**. *arXiv preprint arXiv:2401.12954*.
- Mirac Suzgun, Luke Melas-Kyriazi, and Dan Jurafsky. 2023a. Follow the wisdom of the crowd: Effective text generation via minimum bayes risk decoding. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 4265–4293.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, and 1 others. 2023b. Challenging big-bench tasks and whether chain-of-thought can solve them. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13003–13051.
- Mirac Suzgun, Stuart M Shieber, and Dan Jurafsky. 2024. string2string: A modern python library for string-to-string algorithms. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 278–285.
- Nadeem Ahmed Syed, Huan Liu, and Kah Kay Sung. 1999. Handling concept drifts in incremental learning with support vector machines. In *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 317–321.
- Sebastian Thrun and Tom M Mitchell. 1995. Lifelong robot learning. *Robotics and autonomous systems*, 15(1-2):25–46.
- Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, and 1 others. 2023. Freshllms: Refreshing large language models with search engine augmentation. *arXiv preprint arXiv:2310.03214*.
- Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. 2020. Tent: Fully test-time adaptation by entropy minimization. *arXiv preprint arXiv:2006.10726*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. **Self-consistency improves chain of thought reasoning in language models**. In *The Eleventh International Conference on Learning Representations*.

- Yu Wang, Yifan Gao, Xiusi Chen, Haoming Jiang, Shiyang Li, Jingfeng Yang, Qingyu Yin, Zheng Li, Xian Li, Bing Yin, and 1 others. 2024a. Memoryllm: Towards self-updatable large language models. *arXiv preprint arXiv:2402.04624*.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhui Chen. 2024b. **MMLU-pro: A more robust and challenging multi-task language understanding benchmark**. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Jason Weston, Sumit Chopra, and Antoine Bordes. 2014. Memory networks. *arXiv preprint arXiv:1410.3916*.
- Ling Yang, Zhaochen Yu, Tianjun Zhang, Shiyi Cao, Minkai Xu, Wentao Zhang, Joseph E Gonzalez, and Bin Cui. 2025. Buffer of thoughts: Thought-augmented reasoning with large language models. *Advances in Neural Information Processing Systems*, 37:113519–113544.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. **Tree of Thoughts: Deliberate problem solving with large language models**. *Preprint*, arXiv:2305.10601.
- Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Pan Lu, Zhi Huang, Carlos Guestrin, and James Zou. 2025. Optimizing generative ai by backpropagating language model feedback. *Nature*, 639:609–616.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. 2022. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488.
- Kai Zhang, Yangyang Kang, Fubang Zhao, and Xiaozhong Liu. 2024a. **LLM-based medical assistant personalization with short- and long-term memory coordination**. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2386–2398, Mexico City, Mexico. Association for Computational Linguistics.
- Marvin Zhang, Sergey Levine, and Chelsea Finn. 2022. Memo: Test time robustness via adaptation and augmentation. *Advances in neural information processing systems*, 35:38629–38642.
- Tianjun Zhang, Shishir G Patil, Naman Jain, Sheng Shen, Matei Zaharia, Ion Stoica, and Joseph E. Gonzalez. 2024b. **RAFT: Adapting language model to domain specific RAG**. In *First Conference on Language Modeling*.
- Zexuan Zhong, Tao Lei, and Danqi Chen. 2022. **Training language models with memory augmentation**. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5657–5673, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. 2023. **Least-to-most prompting enables complex reasoning in large language models**. In *The Eleventh International Conference on Learning Representations*.

A Background & Related Work

A.1 Test-time learning (online learning)

Test-time learning—also referred to as online or incremental learning (adaptation)—encompasses a family of methods in which a stochastic model updates its predictions by incorporating information seen during inference, without undergoing conventional, full-scale offline finetuning. Early versions of test-time adaptation focused on local or transductive learning, where a model re-fit or re-weighted its parameters with each new test instance or batch (McCloskey and Cohen, 1989; Thrun and Mitchell, 1995; Amari, 1998; Syed et al., 1999; Bottou and Cun, 2003; Bottou and Le Cun, 2005, *inter alia*). In computer vision, for example, methods like test-time training have been shown to mitigate domain shifts by optimizing a self-supervised loss on incoming data (Wang et al., 2020; Sun et al., 2020; Liu et al., 2021; Boudiaf et al., 2022; Niu et al., 2022; Zhang et al., 2022; Sun et al., 2024). In the context of natural-language generation, test-time adaptation has appeared under terms such as “dynamic evaluation” (Mikolov et al., 2010; Graves, 2013; Krause et al., 2019; Rannen-Triki et al., 2024), in which a language model is updated with gradient steps on the test-time data itself.

However, directly updating language model weights at test time can be computationally expensive and requires the capacity to modify parameters. For large-scale, black-box APIs (e.g., GPT-3 or Claude), one often lacks the ability to perform parameter updates easily, thereby making such an approach difficult, if not completely infeasible (Shi et al., 2024). To address this, a growing body of work has explored parameter-free adaptation, whereby one structurally modifies immediate model inputs (e.g., prompting) or draws from external memory to “update” the model’s effective reasoning. Our approach aligns with this direction by allowing an LM to iteratively record solutions, explanations, or heuristics in an external memory component over successive interactions, avoiding weight updates entirely.

In the broader test-time learning literature, reflexive, compositional, and iterative refinement approaches like Reflexion (Shinn et al., 2023), Self-Refine (Madaan et al., 2023), (Self-)Critic (Gou et al., 2023), Chameleon (Lu et al., 2023), Meta-Prompting (Suzgun and Kalai, 2024), and Self-RAG (Asai et al., 2023) *inter alia*, use feedback loops or verification mechanisms to correct mis-

takes in solutions. TextGrad (Yuksekgonul et al., 2025) similarly draws on the notion of “textual gradients” as an alternative to parameter-based gradients and provides a pathway for improvement based on the content of mistakes. Our proposed DC framework differs by focusing explicitly on storing generalizable heuristics, solutions, or meta-level insights that can be repeatedly retrieved and applied across tasks, not just to correct a single solution. Furthermore, DC does not require a new training loop for each batch or scenario; instead, the memory itself is updated to reflect newly found solutions, errors, or strategies without touching the model weights.

A.2 Test-time compute and reasoning

It is now widely known and accepted that contemporary LLMs such as GPT-4 can exhibit substantial improvements in reasoning and generation capability when additional compute is devoted to inference-time strategies (e.g., chain-of-thought prompting (Wei et al., 2022; Kojima et al., 2022; Zhou et al., 2023), tree-of-thought expansions (Yao et al., 2023; Long, 2023), minimum Bayes risk decoding (Suzgun et al., 2023a; Shi et al., 2022; Golovneva et al., 2023), majority-vote sampling (Wang et al., 2023)). Prompting methods such as Tree-of-Thought (Yao et al., 2023), Graph-of-Thought (Besta et al., 2024), and other non-linear compositional reasoning paradigms systematically enlarge the inference-time search space. They allow models to explore various reasoning paths and exploit consensus or iterative corrections to arrive at more accurate and reliable conclusions (Wei et al., 2022; Wang et al., 2023).

However, these expansions come at the cost of increased computational overhead per test instance (Yao et al., 2023). They are, however, typically ephemeral: once a solution is generated, subsequent tasks or input samples do not generally benefit from the heavy compute spent earlier, unless the user manually engineers advanced prompt-sharing or in-context demonstration strategies. *Cf.* (Zelikman et al., 2022). Our work, on the other hand, aims to reduce repeated overhead across multiple test instances of a similar domain by building a memory that persists from one query to the next. This memory not only reduces repetitive mistakes, but also consolidates and codifies robust solution strategies—effectively amortizing or “sharing” the

cost of initial reflection across future tasks.¹²

Another related thread involves tool usage or code execution (Schick et al., 2023; Lu et al., 2023; Shen et al., 2023; Qin et al., 2023; Surís et al., 2023; Suzgun and Kalai, 2024). These studies have explored how LLMs can call external Python interpreters, symbolic solvers, or other specialized services and APIs to offload complex computations. Our empirical findings too illustrate that once an LLM under DC recognizes a systematic way (e.g., Python-based brute force algorithm) to handle a certain class of problems (like arithmetic puzzles), it can store that approach in memory and repeatedly retrieve it. Thus, DC not only invests extra compute in a single session but spreads that computational benefit across multiple interactions, effectively learning to use tools more consistently and reliably over time.

A.3 Memory-augmented generation and reasoning

Augmenting language models with external memory has seen renewed interest in recent years (Munkhdalai et al., 2019; Guu et al., 2020; Khandelwal et al., 2020; Bulatov et al., 2022; Borgeaud et al., 2022; Zhong et al., 2022; Feng et al., 2022; He et al., 2024; Wang et al., 2024a)—see also (Graves et al., 2014; Weston et al., 2014; Joulin and Mikolov, 2015; Suzgun et al., 2019) for early studies. Modern retrieval-augmented LLM approaches generally consult an external corpus of documents (i.e., a knowledge base) to improve factuality and reduce hallucination (Lewis et al., 2020; Lazaridou et al., 2023; Vu et al., 2023; Zhang et al., 2024b), but the retrieval corpus is almost always fixed prior to inference and does not evolve over time. These methods have been especially effective for open-domain question answering (Lewis et al., 2020; Guu et al., 2020; Karpukhin et al., 2020), where the model’s own parameters may not hold all

relevant knowledge. In practice, retrieval augmentation typically involves selecting and concatenating top- k passages from a knowledge-base—while useful for factual queries, the approach, however, does not inherently solve iterative improvement or learning from mistakes in the sense of building upon prior solutions at inference time.

Another line of research more closely aligns with our vision by storing not just reference knowledge but also the reasoning processes and solution strategies of language models. Several recent works have explored this direction. Thought-Retriever (Feng et al., 2024) logs the model’s chain-of-thought from past queries and uses them for new, analogous queries. Buffer-of-Thoughts (BoT; Yang et al., 2025) takes a slightly different approach by distilling high-level “thought templates” from problem-solving processes, though it relies on predefined templates that seem to be tailored towards specific task types that were considered in their experiments. Madaan et al. (2022) have demonstrated that deployed models like GPT-3 can be improved through memory mechanisms that capture user feedback on errors, preventing similar mistakes in future interactions. Zhang et al. (2024a) have proposed a dual memory architecture combining long-term and short-term storage for medical applications, though their approach requires fine-tuning to incorporate new knowledge.

While these works reveal the many strategies for harnessing memory or feedback, DC emphasizes selectively storing the most relevant insights and heuristics. DC aims to avoid naive accumulation of full raw transcripts and ephemeral chain-of-thought expansions that can lead to memory bloat. Moreover, unlike methods that assume the model can be retrained or finetuned to incorporate memory items, we remain fully external and training-free; this aligns with “plug-and-play” usage principle, in which an off-the-shelf model is augmented by an external memory that it reads from and writes to, but does not require any gradient-based adaptation.

¹²Some lines of work—such as majority voting or sampling-based self-consistency—combine multiple inference passes for a single question but still lack a persistent knowledge base that spans different queries. DC differs in that we treat consecutive tasks in a sequence as a chance to refine a persistent, external store of learned lessons. The memory curation step selectively compiles relevant solutions, heuristics, expansions, or code blocks into a form that can be reused for upcoming queries. Thus, while the compute for the first few tasks may be higher, future tasks become simpler because the system can consult and adapt previously curated knowledge. This approach echoes the underlying motivation of test-time training—performing ongoing improvement at inference—but capitalizes on a cheap, external memory update in lieu of repeated or expensive parameter updates.

B Additional Details about Evaluation

B.1 Evaluation Protocol

To ensure standardized and reliable evaluation, all models are instructed to format their final answers in a structured, machine-readable format. All model answers are expected to be wrapped in the following XML-style tags:

```
<answer>
  (final answer)
</answer>
```

This explicit format ensures accurate and consistent parsing, eliminating errors arising from extraneous text or ambiguous outputs. Once extracted, the final answers are evaluated using their corresponding task-specific accuracy metric.

B.1.1 Accuracy Metrics

Given the diversity of the tasks, we use different accuracy metrics tailored to the specific requirements of each dataset.

Soft Match (SM) is a lenient metric that considers an answer correct if it matches the ground truth after ignoring minor formatting differences, such as punctuation or whitespace variations. We apply this metric to GPQA-Diamond, and MMLU Pro (Engineering and Physics), in which questions are presented in a multiple-choice format.

Functionally Correct (FC) is an even more flexible metric that evaluates whether the model’s output satisfies the task-specific constraints, even if the exact numeral presentation or formatting differs slightly from the reference solution. We apply this metric to the Game of 24, Math Equation Balancer, and AIME benchmarks.

C Addition Discussion of Limitations

We explicitly discuss limitations throughout the paper, particularly in Section 4 (Main Results) and Section 5 (Additional Analyses and Discussions). These sections transparently address the constraints of our approach, including generalizability concerns, potential biases in our experimental setup, computational constraints, and scenarios where our methods may not perform optimally. We believe this transparency is crucial for proper interpretation of our findings and to guide future research.

On experiment statistical significance. Due to the extensive cost and scale of conducting experiments across multiple language models, we were unable to perform multiple runs for all settings, which limits our ability to report error bars or provide a full statistical significance analysis. However, preliminary experiments involving repeated runs demonstrated consistent trends, suggesting that our results are reasonably robust despite the lack of comprehensive statistical validation.

On broader impact. While our work is primarily focused on evaluating and improving LLM capabilities on benchmark tasks, which has limited direct societal impact, we acknowledge that advances in this field could indirectly influence AI deployment in various contexts. Our research contributes to understanding LLM performance and limitations, which can help ensure more responsible implementation of these technologies. We do not foresee significant negative societal impacts specific to our technical contributions.

D Additional Figures and Tables

D.1 Performance Comparison of Baseline and DC-RS Approaches

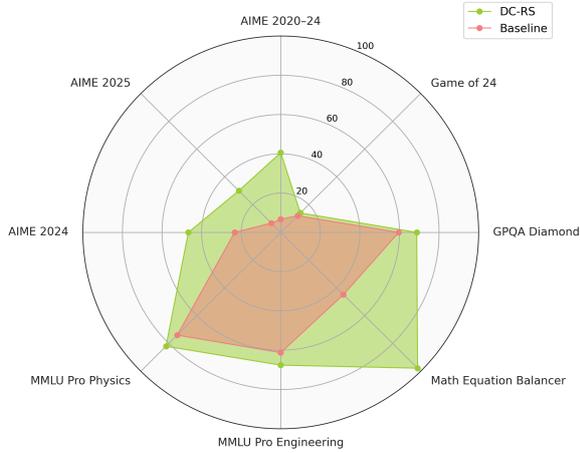


Figure 5: Overall performance of Claude 3.5 Sonnet under the baseline prompting approach with minimal instructions (Baseline) and Dynamic Cheatsheet with Retrieval & Synthesis (DC-RS).

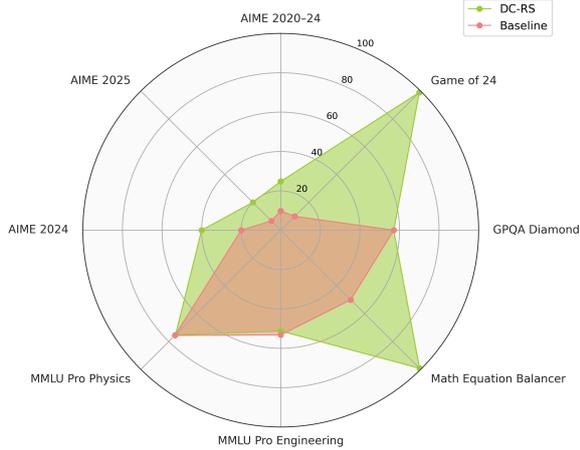


Figure 6: Overall performance of GPT-4o under the baseline prompting approach with minimal instructions (Baseline) and Dynamic Cheatsheet with Retrieval & Synthesis (DC-RS).

Tasks	Claude 3.5 Sonnet			GPT-4o		
	BL	FH	DC-Cu.	BL	FH	DC-RS
AIME 2024	23.3	26.7	50.0	20.0	13.3	40.0
AIME 2025	6.7	6.7	36.7	6.7	3.3	20.0

Table 2: Performance breakdown of BL (default baseline), **FH** (full history), **DC-Cu.** and **DC-RS** approaches under AIME 2024 and 2025. FH stores all past queries and outputs, while DC-Cu and DC-RS selectively refine stored memory. Results indicate that targeted memory curation in DC-RS leads to greater accuracy gains compared to full history retention, supporting the need for structured, self-updating knowledge mechanisms.

Tasks	Claude 3.5 Sonnet				
	BL	MV(BL)	DC- \emptyset	MV(DC- \emptyset)	DC-Cu.
AIME 2024	23.3	23.33	36.7	33.3	50.0
AIME 2025	6.7	6.7	23.3	23.3	36.7

Table 3: Comparison of majority voting (MV) with DC on AIME 2024 and 2025 examinations.

Tasks	Claude 3.5 Haiku			
	BL	DC- \emptyset	DC-Cu.	DC-RS
AIME 2024	10.0	26.7	36.7	30.0
AIME 2025	0.0	13.3	13.3	10.0
GPQA-Diamond	43.4	41.9	43.7	49.0

Tasks	GPT-4o-mini			
	BL	DC- \emptyset	DC-Cu.	DC-RS
AIME 2024	16.7	20.0	13.3	13.3
AIME 2025	10.0	13.3	13.3	16.7
GPQA-Diamond	34.3	34.3	33.8	32.3

Table 4: Performance of Claude 3.5 Haiku and GPT-4o-mini, the smaller counterparts of Claude 3.5 Sonnet and GPT-4o, across AIME and GPQA-Diamond. These smaller models struggle to fully leverage DC, suggesting that memory-based adaptation is most effective when the base LM has sufficient generative competence. Performance improvements are more muted, highlighting the dependency of DC on model-scale reasoning ability.

D.2 Clustering of Errors and Corrections

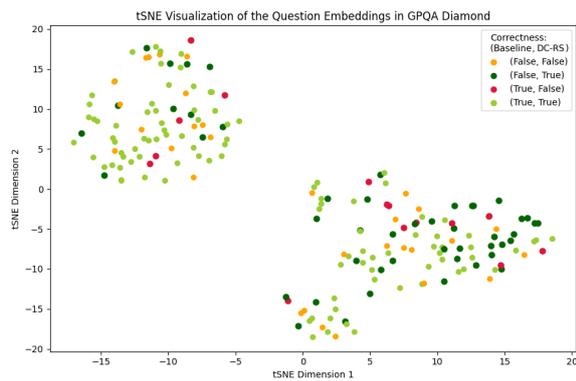


Figure 7: t-SNE visualization of the embeddings of the raw questions in GPQA-Diamond. Note that correct and incorrect answers often cluster in latent embedding space. DC can help transfer learned strategies within these clusters, but without careful curation, erroneous heuristics may also spread, thus requiring careful memory refinement and verification of solution strategies.

D.3 Evolution of Memory Content under Dynamic Cheatsheet

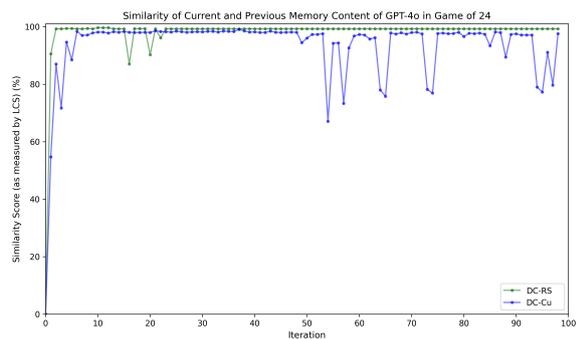


Figure 8: This figure illustrates how memory content of GPT-4o evolves over time in Game of 24, quantified using a longest-common-subsequence (LCS)-similarity metric (Suzgun et al., 2024) between consecutive states (measured at the word level). While both DC-Cu and DC-RS show high stability after the first few iterations, DC-Cu experiences slightly greater fluctuations in the second half of inference.

D.4 Solution Generator and Memory Curator Prompts

D.4.1 Prompt Used by the Generator Model in Baseline

```
Please answer the following question.

_____

## FINAL ANSWER FORMAT
ALWAYS present your final answer in the following format:
FINAL ANSWER:
<answer>
(final answer)
</answer>

N.B. Make sure that the final answer is properly wrapped inside the <answer> block.
* For multiple-choice questions: Only provide the letter choice (e.g., (A))
* For numerical answers: Only provide the final number (e.g., 42)
* For other types of answers, including free-response answers: Provide the complete final answer

Example:
Q: What is the meaning of life?
A: [...]
FINAL ANSWER:
<answer>
42
</answer>

_____

INPUT:
'''
[[QUESTION]]
'''
```

Figure 9: Prompt used in the baseline (BL) approach, where the model receives minimal instructions. The prompt simply asks the model to answer the given question without any structured guidance, additional reasoning steps, or tool-use encouragement. This setup represents a traditional one-off inference method, reflecting how LLMs typically operate by default.

D.4.2 Prompt Used by the Generator Model in DR, FH, and DC Approaches

```
# GENERATOR (PROBLEM SOLVER)
Instruction: You are an expert problem-solving assistant tasked with analyzing and solving various questions using
a combination of your expertise and provided reference materials. Each task will include:
1. A specific question or problem to solve
2. A cheatsheet containing relevant strategies, patterns, and examples from similar problems
-----

## 1. ANALYSIS & STRATEGY
- Carefully analyze both the question and cheatsheet before starting
- Search for and identify any applicable patterns, strategies, or examples within the cheatsheet
- Create a structured approach to solving the problem at hand
- Review and document any limitations in the provided reference materials
-----

## 2. SOLUTION DEVELOPMENT
- Present your solution using clear, logical steps that others can follow and review
- Explain your reasoning and methodology before presenting final conclusions
- Provide detailed explanations for each step of the process
- Check and verify all assumptions and intermediate calculations
-----

## 3. PROGRAMMING TASKS
When coding is required:
- Write clean, efficient Python code
- Follow the strict code formatting and execution protocol (always use the Python code formatting block;
furthermore, after the code block, always explicitly request execution by appending: "EXECUTE CODE!"):
```python
Your code here
```
EXECUTE CODE!
- All required imports and dependencies should be clearly declared at the top of your code
- Include clear inline comments to explain any complex programming logic
- Perform result validation after executing your code
- Apply optimization techniques from the cheatsheet when applicable
- The code should be completely self-contained without external file dependencies—it should be ready to be
executed right away
- Do not include any placeholders, system-specific paths, or hard-coded local paths
- Feel free to use standard and widely-used pip packages
- Opt for alternative methods if errors persist during execution
- Exclude local paths and engine-specific settings (e.g., avoid configurations like
chess.engine.SimpleEngine.popen_uci("/usr/bin/stockfish"))
-----

## 4. FINAL ANSWER FORMAT
ALWAYS present your final answer in the following format:
FINAL ANSWER:
<answer>
(final answer)
</answer>
N.B. Make sure that the final answer is properly wrapped inside the <answer> block.
* For multiple-choice questions: Only provide the letter choice (e.g., (A))
* For numerical answers: Only provide the final number (e.g., 42)
* For other types of answers, including free-response answers: Provide the complete final answer
Example:
Q: What is the meaning of life?
A: [...]
FINAL ANSWER:
<answer>
42
</answer>
-----

CHEATSHEET:
""
[[CHEATSHEET]]
""
-----

Now it is time to solve the following question.
CURRENT INPUT:
""
[[QUESTION]]
""
```

Figure 10: Generator prompt used in the DR, FH, and DC approaches, where the model receives structured high-level instructions on solution development, strategy selection, and tool usage. This prompt explicitly encourages Python code generation and execution for computational tasks. Notably, this same structured prompt is used in all non-BL methods, including DC- \emptyset , DR, FH, DC-Cu, and DC-RS. We also remark that during the initial phases of our experiments, we used “cheatsheet” and “memory” interchangeably to describe stored problem-solving content. However, to maintain consistency, we formally define M_i as memory throughout this paper. Since this was purely a semantic choice, we did not find it necessary to rerun our experiments to reflect this terminology shift.

D.4.3 Prompt Used by the Memory Curation Model under DC-RS

```
# CHEATSHEET CURATOR

### Purpose and Goals
You are responsible for maintaining, refining, and optimizing the Dynamic Cheatsheet, which serves as a compact yet evolving repository of problem-solving strategies, reusable code snippets, and meta-reasoning techniques. Your goal is to enhance the model's long-term performance by continuously updating the cheatsheet with high-value insights while filtering out redundant or trivial information.
- The cheatsheet should include quick, accurate, reliable, and practical solutions to a range of technical and creative challenges.
- After seeing each input, you should improve the content of the cheatsheet, synthesizing lessons, insights, tricks, and errors learned from past problems and adapting to new challenges.

### Core Responsibilities

Selective Knowledge Retention:
- Preserve only high-value strategies, code blocks, insights, and reusable patterns that significantly contribute to problem-solving.
- Discard redundant, trivial, or highly problem-specific details that do not generalize well.
- Ensure that previously effective solutions remain accessible while incorporating new, superior methods.

Continuous Refinement & Optimization:
- Improve existing strategies by incorporating more efficient, elegant, or generalizable techniques.
- Remove duplicate entries or rephrase unclear explanations for better readability.
- Introduce new meta-strategies based on recent problem-solving experiences.

Structure & Organization:
- Maintain a well-organized cheatsheet with clearly defined sections:
  - Reusable Code Snippets and Solution Strategies
  - General Problem-Solving Heuristics
  - Optimization Techniques & Edge Cases
  - Specialized Knowledge & Theorems
- Use tagging (e.g., Q14, Q22) to reference previous problems that contributed to a given strategy.

### Principles and Best Practices
For every new problem encountered:
1. Evaluate the Solution's Effectiveness
- Was the applied strategy optimal?
- Could the solution be improved, generalized, or made more efficient?
- Does the cheatsheet already contain a similar strategy, or should a new one be added?

2. Curate & Document the Most Valuable Insights
- Extract key algorithms, heuristics, and reusable code snippets that would help solve similar problems in the future.
- Identify patterns, edge cases, and problem-specific insights worth retaining.
- If a better approach than a previously recorded one is found, replace the old version.

3. Maintain Concise, Actionable Entries
- Keep explanations clear, actionable, concise, and to the point.
- Include only the most effective and widely applicable methods.
- Seek to extract useful and general solution strategies and/or Python code snippets.

4. Implement a Usage Counter
- Each entry must include a usage count: Increase the count every time a strategy is successfully used in problem-solving.
- Use the count to prioritize frequently used solutions over rarely applied ones.

### Memory Update Format
Use the following structure for each memory item:
...
<memory_item>
<description>
[Briefly describe the problem context, purpose, and key aspects of the solution.] (Reference: Q1, Q2, Q6, etc.)
</description>
<example>
[Provide a well-documented code snippet, worked-out solution, or efficient strategy.]
</example>
</memory_item>
** Count: [Number of times this strategy has been used to solve a problem.]
<memory_item>
[...]
</memory_item>
** Count: [...]
[...]
<memory_item>
[...]
</memory_item>
...

- Prioritize accuracy, efficiency & generalizability: The cheatsheet should capture insights that apply across multiple problems rather than just storing isolated solutions.
- Ensure clarity & usability: Every update should make the cheatsheet more structured, actionable, and easy to navigate.
- Maintain a balance: While adding new strategies, ensure that old but effective techniques are not lost.
- Keep it evolving: The cheatsheet should be a living document that continuously improves over time, enhancing test-time meta-learning capabilities.

N.B. Keep in mind that once the cheatsheet is updated, any previous content not directly included will be lost and cannot be retrieved. Therefore, make sure to explicitly copy any (or all) relevant information from the previous cheatsheet to the new cheatsheet! Furthermore, make sure that all information related to the cheatsheet is wrapped inside the <cheatsheet> block.
```

Figure 11: Prompt used for the memory curator under DC-RS, which is responsible for maintaining an evolving repository of problem-solving strategies, code snippets, and heuristics. The curator selectively retains high-value insights, refines existing strategies, and organizes memory efficiently. This ensures the memory (cheatsheet) remains concise, generalizable, and action-oriented, continuously improving test-time reasoning. (Once again, we note that during the initial phases of our experiments, we used “cheatsheet” and “memory” interchangeably to describe stored problem-solving content. However, to maintain consistency, we formally define M_i as memory throughout this paper. Since this was purely a semantic choice, we did not find it necessary to rerun our experiments to reflect this terminology shift.)

```

# CHEATSHEET CURATOR (cont'd)
[...]

## Cheatsheet Template
Use the following format for creating and updating the cheatsheet:

NEW CHEATSHEET:
```
<cheatsheet>
Version: [Version Number]

Reusable Code Snippets and Solution Strategies
<memory_item>
[...]
</memory_item>
[...]

General Problem-Solving Heuristics
<memory_item>
[...]
</memory_item>
[...]
</cheatsheet>
```

N.B. Make sure that all information related to the cheatsheet is wrapped inside the <cheatsheet> block.
The cheatsheet can be as long as circa 2000-2500 words.

## PREVIOUS CHEATSHEET
[[PREVIOUS_CHEATSHEET]]

## NOTES FOR CHEATSHEET
[[PREVIOUS_INPUT_OUTPUT_PAIRS]]

Make sure that the cheatsheet can aid the model tackle the next question.

## NEXT INPUT:
[[NEXT_INPUT]]

```

Figure 12: The rest of the prompt used by the memory curator under DC-RS (Figure 11).

Reusable Code Snippets and Solution Strategies

<memory_item>

<description>

Game 24 Solver Strategy: Solve the 24 Game by systematically testing combinations of four numbers with arithmetic operations (+, -, *, /) and parentheses to achieve a result of 24. Each number must be used exactly once.

</description>

<example>

Steps:

1. Understand the Problem:

- Input: Four integers.
- Goal: Combine the numbers using arithmetic operations and parentheses to evaluate to 24.
- Constraints: Each number must be used exactly once.

2. Approach:

- Use brute force or systematic trial-and-error to test all possible combinations of numbers, operations, and parentheses.
- Prioritize operations that simplify the problem (e.g., division resulting in integers, subtraction reducing values).
- Check edge cases (e.g., repeated numbers, large/small values).

3. Example Solution:

- Input: 6, 8, 8, 12
- Solution: $6 * ((8 + 8) - 12) = 24$

4. Tips:

- Use parentheses to control operation precedence.
- Division should result in valid integers or fractions.
- Test all permutations of the numbers and operations systematically.

Python Code for Automation

```
from itertools import permutations, product

def solve_24(numbers):
    ops = ['+', '-', '*', '/']
    for nums in permutations(numbers):
        for ops in product(ops, repeat=3):
            expressions = [
                f"({nums[0]} {ops[0]} {nums[1]}) {ops[1]} ({nums[2]} {ops[2]} {nums[3]})",
                f"({nums[0]} {ops[0]} {nums[1]}) {ops[1]} {nums[2]} {ops[2]} {nums[3]})",
                f"({nums[0]} {ops[0]} ({nums[1]} {ops[1]} {nums[2]}) {ops[2]} {nums[3]})",
                f"({nums[0]} {ops[0]} ({nums[1]} {ops[1]} ({nums[2]} {ops[2]} {nums[3]}))"
            ]
            for expr in expressions:
                try:
                    if abs(eval(expr) - 24) < 1e-6: # Check if result is 24
                        return expr
                except ZeroDivisionError:
                    continue
    return "No solution"

# Example Usage
print(solve_24([6, 8, 8, 12])) # Output: 6 * ((8 + 8) - 12)
```

</example>

</memory_item>

Count: 99

Figure 13: Excerpt from GPT-4o's external memory after processing 100 examples from Game of 24 under DC-RS. Early in the test sequence, the model discovered a Python-based brute-force solution, stored it, and subsequently retrieved it for subsequent puzzles. This shift to structured code reuse resulted in a dramatic performance increase from 10% to 99% accuracy, eliminating arithmetic errors and redundant problem-solving efforts.

```
GENERAL META-REASONING STRATEGIES
<memory_item>
<description>
Systematic Problem Analysis Framework (Reference: Q1-Q20)
For complex mathematical problems:
1. State problem requirements clearly
2. List key observations and theorems applicable
3. Identify patterns and relationships
4. Break into manageable sub-problems
5. Verify against examples
6. Consider computational approach when analytical solution is complex
7. For grid problems, analyze movement patterns and symmetries
8. For combinatorial problems, use appropriate counting techniques
9. Implement verification code when possible
10. Consider edge cases and constraints
11. For grid coloring problems, consider row/column patterns
</description>
<example>
Example application:
1. Requirements: list all given conditions
2. Observations: identify applicable theorems
3. Patterns: look for structural relationships
4. Sub-problems: break into steps
5. Verification: test against examples
6. Implementation: use Python for verification
</example>
</memory_item>
Count: 20
```

Figure 14: Example of Claude 3.5 Sonnet’s curated memory after processing 20 AIME 2024 questions under DC-Cu. The memory captures key solution strategies, enables the model to generalize across similar computational problems, and boosts its accuracy.

E NLP Checklist

N.B. Adopted from the NeurIPS Paper Checklist.

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction accurately reflect our paper’s contributions and scope. We present a comprehensive summary of our methodology, key findings, and results in these sections, with clear indications of what was accomplished. We also explicitly acknowledge the boundaries of our work and its limitations, ensuring that our claims are properly contextualized and not overstated.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We explicitly discuss limitations throughout the paper, particularly in Section 4 (Main Results) and Section 5 (Additional Analyses and Discussions). These sections transparently address the constraints of our approach, including generalizability concerns, potential biases in our experimental setup, computational constraints, and scenarios where our methods may not perform optimally. We believe this transparency is crucial for proper interpretation of our findings and to guide future research.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The paper does not include any theoretical results or formal mathematical proofs. Our work is primarily empirical in nature, focusing on experimental evaluation rather than theoretical proofs or derivations.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Our paper provides comprehensive information for reproducing our experimental results. Section 3 (Experimental Setup) details all tasks, models, and datasets used in our evaluation. The Appendix contains further essential details: Section B documents our complete evaluation protocol and accuracy metrics, while Section C.4 includes all prompts used by the generator and curator models across different experimental approaches. We’ve included our full codebase in the Supplementary Materials with clear documentation, and we commit to releasing all code, data, and experimental results publicly upon publication to facilitate full reproducibility.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We provide our complete codebase in the Supplementary Materials. Upon publication, we will release all code, datasets, and results through a public GitHub repository with appropriate licensing information, ensuring full transparency and accessibility for the research community.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We provide comprehensive experimental details throughout the paper. Section 3 covers the core experimental setup. The Appendix contains further information about our experiments: Section B details our full evaluation protocol and accuracy metrics with precise definitions, while Section C.4 documents all prompts used by the solution generator and memory curator models across different approaches. All implementation details necessary for reproducing our results are included, with no significant experimental parameters omitted.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: Due to the extensive cost and scale of conducting experiments across multiple language models, we were unable to perform multiple runs for all settings, which limits our ability to report error bars or provide a full statistical significance analysis. However, preliminary experiments involving repeated runs demonstrated consistent trends, suggesting that our results are reasonably robust despite the lack of comprehensive statistical validation.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [No]

Justification: We primarily used API calls to commercial LLM providers rather than running models on our own infrastructure. We spent about \$100 to run our experiments.

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: Our work fully conforms with the NeurIPS Code of Ethics.

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: While our work is primarily focused on evaluating and improving LLM capabilities on benchmark tasks, which has limited direct societal impact, we acknowledge that advances in this field could indirectly influence AI deployment in various contexts. Our research contributes to understanding LLM performance and limitations, which can help ensure more responsible implementation of these technologies. We do not foresee significant negative societal impacts specific to our technical contributions.

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our work evaluates existing LLMs on established public benchmarks without creating new models or datasets that could pose misuse risks.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We properly cite and acknowledge the original authors of all code packages and datasets used in our work, with clear references in Section 3 and throughout the paper. To the best of our knowledge, all datasets used in our experiments are licensed under CC-BY 4.0,

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The new assets introduced in our paper include experimental results, analysis data, and evaluation code. These assets are thoroughly documented in Sections 3-5 of the paper, with detailed methodological explanations. Our codebase includes comprehensive documentation with README files, requirements specifications, and usage examples.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This paper does not involve any crowdsourcing experiments or research with human subjects. All evaluations were conducted using automated metrics and existing benchmark datasets without human annotators or participants.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This research did not involve human subjects or participants in any capacity. All experiments were conducted using computational methods on existing public datasets and models. Therefore, IRB approval was not required or sought for this work.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [No]

Justification: The core method development in this research did not involve LLMs as any important, original, or non-standard components.