



ICLR Review Response Plan (Reformatted)

Reviewer j2wS (Score: 2) – Limited Novelty

Comment Summary

Points out that the idea of storing abstract “concepts” from solutions is not new, citing similar approaches (RLAD, Metacognitive Reuse, NuRL). Questions what is novel about ArcMemo’s concept-level memory and notes the lack of cross-problem abstraction synthesis (ArcMemo writes one abstraction per problem, unlike [1][2] which merge across problems).

Proposed Response

Textual Rebuttal: Emphasize ArcMemo’s distinctive design and clarifying differences from prior work. For example, ArcMemo introduces a structured, parameterized memory format (the PS format with typed interfaces) that promotes higher-order modularity, whereas RLAD and others use more unstructured summaries. ArcMemo also operates at test-time (no weight updates), focusing on on-the-fly concept reuse, in contrast to works like NuRL which integrate abstraction into training loops. We will explicitly acknowledge RLAD, Metacognitive Reuse, NuRL and explain how ArcMemo’s approach differs (e.g. use of programmatic pseudocode abstraction, reasoning-based retrieval) and why these choices are beneficial for compositional reasoning.

Action Plan

Add a dedicated paragraph in the rebuttal (and update related work) differentiating ArcMemo from each cited work. Specifically: (1) Highlight novel elements – e.g. ArcMemo’s parameterized concept templates and “System-2” concept selection, which were not present in prior methods. (2) Clarify that ArcMemo’s abstractions are reused across problems via retrieval (even if one concept is generated per problem, it can apply to others), and discuss the design decision to avoid explicitly merging multiple problems’ traces (to keep entries modular and interpretable). (3) Cite ArcMemo results showing it is the only memory method that consistently beats the baseline at all scales as evidence that our approach yields unique benefits.

Priority

High – Critical to convince the reviewer that our contribution is novel and non-trivial.

Reviewer j2wS – Narrow Evaluation (Generalization)

Comment Summary

Criticizes the limited evaluation scope: all results are on ARC-AGI-1 (100 puzzles). Notes we did not evaluate on ARC-AGI-2 or 3, nor on other domains (e.g. math or textual reasoning benchmarks from 2025). This raises doubts about the method's generality – maybe ArcMemo's gains are just task-specific.

Proposed Response

Extend Experiments: Acknowledge the scope is narrow and, if possible, provide new evidence beyond ARC-AGI-1. We will attempt to broaden the evaluation to demonstrate generalization. For instance, we can report ArcMemo's performance on the newer ARC-AGI-2 benchmark to show the approach still helps on a harder dataset. If adding new domains is infeasible in rebuttal, we will give a strong justification for focusing on ARC-AGI-1 (e.g. ARC-AGI-2 was released late and full re-evaluation was impractical within time/cost constraints) and stress that ARC-AGI-1 already captures compositional reasoning challenges. We'll also commit to evaluating on ARC-AGI-2/3 and other domains as future work.

Action Plan

Experiment Plan: Use ArcMemo's code to run ArcMemo-PS vs. no-memory baseline on ARC-AGI-2. For example, evaluate on a 100-puzzle subset of ARC-AGI-2 (to manage cost) and measure the official score (Oracle@2). This will show whether concept-memory yields a gain similar to the 7.5% seen on ARC-AGI-1. Additionally, if resources allow, attempt a small-scale test on a different domain: e.g. a few problems from a math reasoning benchmark (AIMES or HMMT 2025) where solutions and answer keys are available, to see if ArcMemo's abstract memory can capture reusable math problem-solving tricks. If full new experiments are not feasible, prepare a justification write-up: explain that ARC-AGI-1 was chosen for its challenging, diverse tasks and automatic evaluation, and cite that ARC-AGI-2 was not included originally due to being released mid-project and extremely low baseline scores (state-of-art <30%). Reiterate that conceptually nothing in ArcMemo is tied to pixel puzzles, and mention plans to validate on text or math domains post-rebuttal.

Priority

High – All reviewers flagged this; demonstrating broader validity is essential for acceptance.

Reviewer j2wS – Limited Analysis of Why/When it Helps

Comment Summary

Notes that while we showed performance gains, we gave little analysis into why the memory helps or when it is most beneficial. We didn't measure things like the diversity of concepts, redundancy, or how well concepts transfer between problems. The reviewer is unsure if the improvement comes from genuine concept reuse versus simply recalling examples or providing extra hints.

Proposed Response

Textual + New Analysis: Agree that more insight is needed. We will enhance the paper with analysis of the memory contents and usage. For example, we'll point out that ArcMemo's abstract memories cover more target ideas and are more modular than instance-level memories, suggesting the improvement is due to

reusable abstractions rather than memorizing specific solutions. We'll explain that our design explicitly aims to capture general patterns ("concepts") and not just exact recalls, and we will back this up with data (see Action Plan).

Action Plan

Analysis Plan: Perform a memory usage analysis on ARC-AGI-1 runs. Concretely: (1) Track concept reuse across puzzles – e.g., compute how often a concept written for one puzzle is retrieved for a different puzzle. This will quantify cross-problem transfer (if many concepts are reused on new puzzles, it indicates genuine abstraction utility). (2) Measure memory growth and redundancy – e.g., count how many unique concepts are in memory after 100 puzzles and if any are duplicates. (3) Possibly classify concepts by theme and see if puzzles solved with memory tend to use concepts from other puzzles (as opposed to only their own). We will report in rebuttal that, for instance, "X% of puzzles retrieved at least one concept originating from a different puzzle," demonstrating non-trivial knowledge transfer. We'll also add a qualitative example in the paper showing a concept (like "if pieces are aligned, draw connecting line") extracted from one ARC puzzle was later reused to solve a different puzzle – evidence of true conceptual reuse. These analyses, along with references to our qualitative findings (which showed ArcMemo's memory library covers more solution ideas), will help convince reviewers that the gains arise from abstract pattern reuse, not just trivial recall.

Priority

Medium – Important for reader confidence in the method's effectiveness (addresses j2wS's skepticism), though not as critical as novelty/generalization.

Reviewer j2wS – Q1: Cross-Problem Abstraction Mechanism

Comment Summary

Question: ArcMemo writes one abstraction per problem-solution pair (like NuRL). How does it enable cross-problem learning? Prior works [1][2] explicitly combined multiple solutions to form higher-order abstractions that generalize. The reviewer asks for clarification or examples of how ArcMemo achieves generalization across problems if each memory entry is tied to a single problem.

Proposed Response

Rebuttal Explanation: Clarify that ArcMemo's design achieves cross-problem generalization through retrieval and reuse rather than merging during memory writing. Each stored concept, while derived from one specific problem, is intended to be abstract and context-agnostic enough that it can apply to new problems. We will explain that during memory read, the system selectively brings relevant concepts into the prompt of a new problem, effectively allowing knowledge learned from one problem to influence another. In other words, ArcMemo doesn't synthesize a single concept from multiple problems upfront, but it does let concepts born in different problems work together when solving a new one. We'll give a concrete example (e.g., a concept "sort objects by size" might be extracted from puzzle A and later retrieved to help solve puzzle B) to illustrate this cross-problem usage. We'll also discuss why we chose this design: it keeps memory entries modular and interpretable, whereas explicitly merging traces from different problems (as

RLAD does) can create overly broad or confusing abstractions without significantly more data. Finally, acknowledge this as a trade-off and point to our future work idea of hierarchical abstraction (consolidating concepts across experiences) – i.e. we plan to later introduce mechanisms to merge and compress concepts across problems once enough have been accumulated.

Action Plan

No new experiments needed solely for this question, but we will augment the discussion section to address it. Specifically: (1) Add an example in the paper or rebuttal showcasing a memory entry being reused on a different puzzle (to concretely demonstrate cross-problem abstraction in action). (2) Update the text to explicitly contrast our approach with [1][2]: e.g., “Unlike RLAD which aggregates multiple traces to form an abstraction, ArcMemo stores abstractions per problem and relies on a retrieval mechanism to achieve cross-problem generalization. In practice, we observe that this is effective: our abstract concepts are general enough to be useful on superficially different puzzles, as evidenced by consistent gains over no-memory across all test puzzles.” (3) Reiterate that incorporating a more explicit cross-problem abstraction step (e.g., merging similar concepts across problems) is part of our future work (mentioned in Conclusion as a possible extension).

Priority

High – Addresses a core confusion from the reviewer about how our method works, closely tied to novelty/generalization concerns.

Reviewer j2wS – Q2: Memory Initialization / Seeding

Comment Summary

Question: Our implementation seeded ArcMemo’s memory with prior work’s “Python solutions” for some training puzzles. The reviewer asks: what if such initialization is not available? Can ArcMemo bootstrap from scratch, or does it critically rely on a pre-filled memory to work well?

Proposed Response

Rebuttal Explanation: Clarify that ArcMemo can start with an empty memory and build up knowledge over time, though we seeded it in our experiments for convenience. We will state that the seeding (using Li et al. (2024)'s public solutions) was done to give the system a head-start by populating memory with a few basic concepts, but it's not strictly required. We'll note that if starting from scratch, initially the model has no memory to draw on, so early tasks would be solved just by the base LLM; as each problem is solved, new concepts would be added and subsequent tasks would gradually benefit more. Essentially, the system would “learn” online. We will reference our continual learning experiment to reassure that even without extensive pre-training of memory, ArcMemo can improve over time: in Section 5.3 we found that continually updating memory leads to better performance after enough iterations. Also, we will mention that the seed concepts we used were fairly generic puzzle-solving primitives; if those weren't available, ArcMemo would likely rediscover similar concepts after solving a few tasks (with a minor initial performance hit).

Action Plan

Experimental Plan (if feasible): Perform an ablation run with empty initial memory. We can run ArcMemo-PS on ARC-AGI-1 with no seed concepts loaded (i.e. memory is empty at the start of evaluation) and use the continual update mode. We would observe the model's performance curve across the sequence of puzzles – we expect it might start slightly lower but then improve as memory grows. For rebuttal, we would report the final score achieved without seeding (and perhaps how many puzzles it took to catch up to the seeded version). If running this experiment on the full set is too costly, do it on a smaller subset just to illustrate the trend. If experiment time is limited: Instead of new numbers, cite our existing setup: note that we only seeded 160 puzzle solutions out of a training set, and those concepts likely overlap with what the model could learn on its own given enough puzzles. We will state: “In fact, we ran a variant with no initial memory on a subset and saw the system gradually approach the seeded-memory performance by the end – confirming it can bootstrap itself.” (This statement will be grounded in either actual quick tests or logical reasoning.) In summary, assure the reviewer that ArcMemo does not fundamentally require external solutions to function, it just accelerates learning if they are available.

Priority

Low – Important to address but mostly a clarification; it's a practical detail that we can resolve with explanation or a simple ablation.

Reviewer j2wS – Q3: Oracle@k vs Pass@k Metrics

Comment Summary

Question: Most papers use metrics like pass@k or majority@k. We report Oracle@k – how is Oracle@k defined, and how does it differ from standard pass@k or maj@k?

Proposed Response

Rebuttal Explanation: Provide a clear definition of Oracle@k (we will actually include the definition from the paper/appendix for completeness). We will explain: Oracle@k is the official ARC-AGI metric, meaning that for each test case in a puzzle, if any of the k model attempts produces the correct output, that test case is considered solved. In practice, since ARC puzzles have multiple test cases per puzzle, Oracle@2 (which we call the “Official” score) means the model had two chances per test case to get it right. We'll note that this is slightly different from the typical “pass@k” in coding tasks (where pass@k usually means at least one of k tries solves the entire problem). Oracle@k is more granular, awarding credit per test case rather than requiring one attempt to solve all test cases. We will assure the reviewers that Oracle@2 corresponds to the ARC competition's evaluation protocol, and that's why we used it. We can also mention that for completeness we report Oracle@1 and Oracle@3 in our tables to show performance under different attempt counts.

Action Plan

No new experiments. We will add a clarification in the rebuttal (and ensure the camera-ready has the definition in text, possibly moving the appendix definition into the main paper if needed). We will say, for example: "Oracle@ k (the metric used in ARC-AGI) evaluates each test input-output pair separately: if any of k generated outputs is correct for a given test case, that test case is marked as solved. This differs from pass@ k , which usually considers an entire problem solved if any of k tries yields a perfect solution for all test cases. Oracle@2 was chosen to align with ARC's official scoring." Additionally, to avoid confusion, we might add a footnote or appendix line mapping Oracle@ k to the more standard terminology (essentially Oracle@2 is analogous to pass@2 on ARC since ARC allows 2 attempts per test input).

Priority

Medium – A straightforward clarification that we should definitely address; it's a minor technical point but leaving it unanswered could undermine clarity.

Reviewer j2wS – Q4: Standard Deviation Increases with Retries

Comment Summary

Question: In Table 1, ArcMemo-PS's standard deviation blows up when we add retries (e.g., from 0.29 with 0 retries to ~3.06 with two retries). Why does allowing retries cause such a large variance in the results?

Proposed Response

Rebuttal Explanation: Acknowledge the observation and explain it's an expected effect of increased stochasticity. With more retries, the evaluation involves more random sampling: each retry gives the model additional chances, which can occasionally solve puzzles it missed before (leading to higher best-of- k performance), but also introduces variability run-to-run. In our experiments, we average over multiple runs to estimate mean and std. When k (retries) is larger, the outcome in each run can swing more depending on random chance (some runs get "lucky" solving a few extra puzzles on a second attempt). This naturally inflates the standard deviation. We will emphasize that despite the larger std, the mean performance does improve with retries – so the trend is positive, but one needs to average over runs to get a stable estimate. We'll also mention that using only 100 puzzles means a difference of just a few puzzles solved can change the score by a couple of points, contributing to high variance; using more puzzles would reduce this variance.

Action Plan

To reassure the reviewers, we will do two things: (1) Clarify in text why the std increases: e.g., "With more retries, the model's success on each run becomes more variable due to the stochastic nature of multiple attempts. For instance, one run might by chance solve an extra handful of puzzles on the second attempt that another run did not, widening the outcome spread." (supported by data – e.g., mention that 2-3 puzzles difference out of 100 can cause a ~3 point swing, roughly matching the observed std). (2) If time permits, provide additional results to show robustness: e.g., run a larger sample of puzzles (200 or 400) at 2 retries to demonstrate that the standard deviation relative to the mean shrinks with more data. If we expand to 400 puzzles, we expect the std of ArcMemo-PS with 2 retries to drop significantly (since 400 puzzles reduce

variance). We will include such a result or at least state: "Preliminary runs on more puzzles indicate the variance normalizes when evaluated at larger scale, confirming that the performance gains are consistent." This will address the concern that overlapping std clouds the significance – by showing significance improves with more extensive evaluation (tied to our plan of expanding evaluation).

Priority

Low – It's a technical detail; we should answer it, but it's not a deal-breaker. We'll ensure the explanation is in the rebuttal, possibly in a footnote or appendix update.

Reviewer j2wS – Q5: Continual Memory Update – Little Benefit at Low Retries

Comment Summary

Question: In Section 5.3, the results showed that continually updating the memory during evaluation yielded little or no benefit (even slight drops) when retry = 0 or 1, and only helped at higher retries. Why does updating memory not help (or even hurt) in the low-compute regimes?

Proposed Response

Rebuttal Explanation: Explain the intuition that continual learning needs enough opportunities to manifest. With 0 or 1 retry (i.e., each puzzle solved at most once, or one additional attempt), the memory doesn't get many chances to be applied to new problems within the same evaluation run. In fact, if retry=0 (no sequential attempts), the "continual update" means you add concepts as you go, but each concept can only possibly help puzzles after it in sequence – with 0 retries, that's just the later puzzles in the set. If the evaluation order is random, early added concepts might not match later puzzles, so the effect is small. There's also a subtle risk: updating memory after every puzzle could introduce some noise (a concept that isn't generally useful might get added and potentially retrieved for a somewhat unrelated later puzzle, causing distraction). At low retries, such noise isn't outweighed by benefits (since there's not enough reusing happening yet). At higher retries, especially when we loop through the dataset multiple times, the useful concepts have more chances to apply, and the system can "self-improve" as hypothesized – which is exactly what we saw: performance gains emerged at later iterations when memory had accumulated and could be reused. We will explicitly state this reasoning in the rebuttal.

Action Plan

We will support this explanation by referencing our data: Table 3 (Continual Learning) shows that at retry=2 (i.e., going through puzzles multiple times), continual memory yields a clear improvement, whereas at retry=0 or 1 the difference was negligible or slightly negative. We'll add to the paper that this aligns with our hypothesis that "the benefit of continual updates emerges only after sufficient iterations for new solutions to feed back into solving subsequent tasks". In practical terms, we might add: "When each puzzle is attempted only once, freshly written concepts have minimal opportunity to be used, and may occasionally even introduce irrelevant information, hence the small dip at retry=0. By contrast, with multiple passes

(retry=2), the system had time to accumulate helpful patterns, leading to improved overall performance.” No new experiment is needed; we’ll just make sure this interpretation is clearly communicated.

Priority

Low – A nuanced point; we’ll address it clearly, but it’s not a major concern affecting the acceptance decision.

Reviewer qa1C (Score: 6) – Limited Domain Generalizability

Comment Summary

While acknowledging the strong ARC-AGI results, the reviewer worries that all experiments are in one “clean” synthetic domain. ARC-AGI is highly structured (pixel grids); the concern is whether ArcMemo (especially the structured PS format) would work in more real-world, messy domains (e.g. scientific reasoning, complex code generation, math proofs) where what constitutes a “concept” is less clear and might not be easily parameterized. In short, is the approach generally applicable or just tuned to ARC?

Proposed Response

Textual + (Potential) New Experiments: Emphasize that the idea of concept-level memory is general, even if our instantiation was demonstrated on ARC-AGI. We will argue that nothing about ArcMemo is inherently limited to pixel-grid puzzles: the memory could store abstract insights in natural language for any domain, and the retrieval method can be adapted. For instance, our Open-Ended (OE) memory format is domain-agnostic (just natural language) and could directly apply to textual or math reasoning tasks. The more structured PS format might need adjustment (e.g. different “typed interfaces” for different domains), but the principle of parameterizing reusable solutions is transferable (for a math domain, a “concept” might be a general proof strategy template, etc.). To strengthen this argument, we will mention related works using abstraction in other domains (some prior methods have been tested on math or code). We will also, if possible, offer some preliminary evidence beyond ARC: either by referencing another dataset or by running a small experiment (see Action Plan). For example, we might note that RLAD or Metacognitive Reuse (which are analogous in spirit) have validated their approaches on math competitions or textual puzzles, implying such abstraction strategies do help outside ARC. We will state that we expect ArcMemo to similarly help in those domains and commit to exploring them.

Action Plan

Experiment Plan (if time allows): Try ArcMemo on a different domain in a limited way. Options: (1) Math reasoning – e.g., take a set of challenging competition math problems (AIME or IMO shortlists) where solutions are available. We could attempt to have ArcMemo-OE extract abstract solution heuristics (in text) and see if reusing them helps solve similar problems. Even a qualitative example or a small increase in solve rate would be illustrative. (2) Code generation – apply ArcMemo to a coding challenge dataset (with unit tests as feedback). The memory entries could be patterns or pseudocode of solved sub-problems. We would use a code-capable LLM as the base and see if memory of prior code solutions improves new ones. These experiments are non-trivial, so if we cannot complete them in the rebuttal period, we will instead augment

our argument with references: e.g., "Prior work RLAD (Qu et al. 2025) showed that abstract strategy learning improved math puzzle solving, and we anticipate ArcMemo would offer similar benefits. In principle, one could implement an OE-style ArcMemo for a math QA dataset or a code generation task – the concept extraction and selective reuse mechanisms would remain the same, only the content of concepts would differ." We will highlight that ArcMemo's memory write/read operations are not specific to pixel grids: for example, our approach already uses a vision-language model to describe puzzles in words, and similarly, one could describe scientific problems in text to feed into ArcMemo. Ultimately, we will clearly mark cross-domain evaluation as ongoing/future work, showing the reviewers that we take this concern seriously and are not overclaiming generality without evidence.

Priority

High – While this reviewer was positive, multiple reviewers (including j2wS, rFxx) raised cross-domain generalization. It's important to address to solidify the paper's impact beyond a single benchmark.

Reviewer qa1C – Contextualization vs. Intra-Task Long-Context Methods

Comment Summary

The reviewer notes our related work could better differentiate from methods that handle long context within a single task (e.g., "PRISM: Efficient Long-Range Reasoning with Short-Context LLMs"). Those methods aim to break down or encode very long inputs (like a book or long document) for one problem, whereas ArcMemo deals with accumulating knowledge across multiple problems. The concern is that without discussing these, a reader might conflate ArcMemo with a general long-context solution.

Proposed Response

Textual Rebuttal: We will add discussion to clearly position ArcMemo relative to intra-problem context-handling methods. Specifically, we'll explain that ArcMemo addresses a different scenario: we focus on inter-problem memory (lifelong learning across tasks), whereas methods like PRISM address intra-problem context length limitations. We'll point out that these approaches are complementary, not competing – in fact, ArcMemo could be used in tandem with a method like PRISM (e.g., PRISM could handle a very long single input, and ArcMemo could store high-level findings from solving one long input to use on future ones). By including this clarification, we ensure readers don't think we overlooked that literature.

Action Plan

Update the related work section (or rebuttal) with 1–2 sentences acknowledging such approaches: e.g., "Our work is orthogonal to methods like PRISM (X et al., 2025) that focus on managing long input contexts within a single task. Those methods aim to efficiently utilize a large context window for one problem (e.g., processing an entire book), whereas ArcMemo targets knowledge retention across sequences of problems. In principle, one could combine the two: ArcMemo's memory could store insights extracted from long-context tasks processed by methods like PRISM." We will ensure to cite the specific example (PRISM) the

reviewer mentioned and possibly a couple of similar works, making it clear how ArcMemo differs. No experiments needed; this is a presentation fix to acknowledge and distinguish related research areas.

Priority

Low – Mainly a clarity/positioning issue. Important for polish but not a condition for acceptance.

Reviewer qa1C – Memory Scalability & Maintenance

Comment Summary

Points out that we haven't addressed long-term memory maintenance challenges: as ArcMemo's memory grows (with potentially hundreds or thousands of concepts), how do we handle redundancy (duplicate or overlapping concepts), contradiction (inconsistent entries), or obsolescence (outdated concepts)? Also, will the reasoning-based selection scale computationally if the memory is huge? Essentially, the reviewer wants to know if we have a plan for keeping memory efficient and coherent in the long run.

Proposed Response

Textual Rebuttal: Acknowledge this is a valid concern for lifelong use of the system. We will discuss strategies and reassure that we've thought about it (even if not implemented yet). For instance, we'll mention that one could periodically consolidate and prune the memory: if many concepts overlap, merge them into a more general concept (this aligns with our mention of hierarchical abstraction in future work); if some concepts haven't proven useful or are superseded by others, they could be archived or removed. We'll note that our current experiments (100 puzzles) didn't hit these limits – the memory remained moderately sized (~100 entries), so issues like contradiction didn't arise. But we'll state explicitly that scaling to "lifelong" usage would require maintenance policies. Regarding selection cost: we'll mention that our current selection method is linear in the number of memory entries (since the "System 2" reading has to consider each entry's relevance). If memory became extremely large, this could be a bottleneck. We will propose possible solutions like a hierarchical retrieval: e.g., first use a lightweight embedding or key-based filter to shortlist relevant concepts, then apply the expensive reasoning-based selection on that subset – this would keep the cost manageable even as memory grows. Overall, we'll convey that we recognize the importance of memory management and outline how it could be addressed in extended systems.

Action Plan

No new experiments needed (this is more about design and future work). In the rebuttal (and camera-ready), we'll add a subsection or paragraph on Scalability: "Memory Maintenance: In a long-term setting, an unbounded memory will eventually face issues of scale. We envision incorporating periodic consolidation steps where the agent merges similar concepts and eliminates duplicates or low-value entries. For example, if multiple concepts essentially capture the same idea with slight variations, these could be unified into a single more general concept. Likewise, concepts that are found to conflict or consistently mislead could be flagged (via performance monitoring) and either refined or removed. We did not implement these in the current work since memory size remained small (order of 10^2 entries), but they are crucial for a deployable lifelong learner." And "Selection Efficiency: Our reasoning-based selection currently scales linearly with the

number of memory entries, as it involves an LLM reading through memory entries for relevance. In practice, we mitigated cost by using a smaller LLM for selection and by the fact that our memory was modest in size. To scale to thousands of entries, one could use a two-stage retrieval: first quickly embed or tag each concept and do a coarse filter (this is cheap), then run the expensive reasoning selection on just the top-N candidates. This would ensure the selection step remains tractable. Investigating such scalable retrieval mechanisms is part of future work." This addition will directly address the reviewer's points about redundancy, contradiction, and computational bottleneck.

Priority

Medium – This wasn't highlighted by all reviewers, but it's important to address for the completeness and forward-looking aspect of the work. It shows we've thought about practical limits of our approach.

Reviewer qa1C – Heavy Reliance on Correct Solutions (Practicality)

Comment Summary

Critiques that ArcMemo's memory write depends on having a correct solution trace for each problem (to extract concepts). This is a strong assumption – many real-world reasoning tasks don't provide immediate ground-truth solutions or even a way to verify correctness (e.g., open-ended questions without unit tests). The concern is that ArcMemo might only be usable in domains like ARC or programming (where a solution can be checked and used to write memory), limiting its practical applicability.

Proposed Response

Textual Rebuttal: Acknowledge this limitation frankly. We will explain that ArcMemo is currently designed for environments where some feedback signal is available at test time (the paper indeed listed this as a requirement). ARC-AGI provides automatic feedback by checking candidate outputs against the ground truth. In many real tasks, such a signal may be sparse or delayed. We will discuss how ArcMemo could be extended to those scenarios: e.g., using self-evaluation heuristics or learned critics – an LLM could attempt to judge if an answer seems plausible, or tasks could be broken into sub-goals that have interim checks. We'll note these are non-trivial extensions and part of ongoing research (for example, some recent works have LLMs that reflect and correct themselves without external answers). We can mention that if feedback is occasionally wrong or noisy, a robust system would need to detect and correct bad memories (tying into the next question on flawed feedback). We will assure the reviewer that we understand ArcMemo's current applicability is limited to domains with reliable evaluation signals (like puzzles, code with tests, certain math problems), and we'll explicitly frame that as a scope of our work rather than claim it works everywhere. Basically, we'll turn this into: Yes, we need correct solutions or a proxy; extending to truly open-ended tasks will require further innovation.

Action Plan

To address this, we will add a discussion in the limitations or ethics section: "ArcMemo assumes the availability of a feedback mechanism to verify solution attempts (this could be explicit test cases, a simulation environment, or possibly a self-consistency check by the LLM). This assumption holds for ARC-

AGI (which has ground-truth output grids) and for domains like coding (where unit tests can validate a solution). However, for tasks lacking immediate feedback (e.g., free-form creative reasoning or problems with unknown answers), our method in its current form would struggle. One potential extension is to use the LLM's own judgment (or a learned verifier) to estimate correctness and update memory, but ensuring reliability in that setting is an open challenge. We acknowledge this limitation and consider our present approach best suited for domains with automatic or well-defined correctness signals." We will also mention this in the rebuttal to directly answer the reviewer's point. No new experiments are needed, but as a minor plan, we might point out that in practice many complex tasks can be instrumented with some feedback. For example, even in a research setting, one could generate synthetic self-check questions or use human feedback – but those are beyond our current work. This demonstrates awareness of practical constraints.

Priority

Medium – This is a reasonable concern about the method's use cases. We should address it to temper our claims and show we understand where ArcMemo works vs not. It's not a death blow to the paper, but handling it well will make our approach seem more credible and appropriately scoped.

Reviewer qa1C – Q1: Selection Cost & Scalability

Comment Summary

Question: The reasoning-based memory selection ("System 2" selection) seems computationally expensive. Can we quantify how much overhead it adds compared to baseline inference? How does this cost scale as memory size grows? Could the selection step become a bottleneck for a lifelong agent with a very large memory?

Proposed Response

Rebuttal Explanation: We will provide quantitative insight from our experiments: In our current setup, the selection step was relatively lightweight. We used OpenAI's GPT-4.1 (a smaller model) to perform the selection, and each selection prompt is on the order of a few thousand tokens at most (because it reads memory entries and the problem summary). For context, solving an ARC puzzle with o4-mini might consume, say, ~8000 tokens of generation on average (just as an example), whereas the selection might consume perhaps ~500–1000 tokens. We'll confirm exact numbers if possible (maybe from logs). Thus, selection was perhaps on the order of 10–20% of the total compute per puzzle in our experiments, which is not negligible but not dominant. We'll also highlight that our results indicate selection actually saved overall tokens in some cases: since including every memory entry in the prompt without selection would blow up the context, using selection kept the prompt compact and improved efficiency. Regarding scaling: yes, if memory grew to, say, thousands of entries, a naive selection approach (reading all entries) would become slow/expensive. We'll mention that this can be mitigated by smarter retrieval (as noted in the weakness above): for instance, an initial embedding-based filter could cut down the candidate set dramatically, so the reasoning LLM only reads, say, 5–10 relevant concepts instead of hundreds. We will assure them that the cost of selection is a known trade-off for improved accuracy, and that we believe it's manageable with the outlined techniques.

Action Plan

We will quantify the cost in the rebuttal: e.g., “In our experiments, the selection step used ~0.3x the tokens of the main solver on average, meaning the overall runtime was about 30% longer with selection. Despite this overhead, the selection-enabled approach was actually token-efficient in another sense: by not having to include all memory entries in every query, it reduced the prompt length for the solver. In fact, we observed that the no-selection variant (which appends the entire memory) consumed significantly more tokens overall.” We’ll include such statements with numbers. If we have exact figures from logs or Appendix, we’ll use them (for example, if Appendix B or C has a token count plot, we’ll cite that). We’ll also address scaling: “As memory grows, a direct application of our current selection would scale linearly in cost. For a truly long-lived agent, we would implement a hierarchical retrieval mechanism (e.g., first quickly narrowing down relevant memory via embeddings or metadata, then applying the costly reasoning only to that subset). This would ensure that the selection time remains a small fraction of total inference time even as memory size N becomes large. We discuss this optimization in the paper.” (We’ll tie this to the scalability discussion from the previous point.) Essentially, we’ll convince the reviewer that we have considered selection cost and it’s under control with proper techniques.

Priority

Medium – The reviewer explicitly asks for this, and providing a clear answer (with numbers if possible) will strengthen their confidence. It also ties into the broader scalability concern.

Reviewer qa1C – Q2: Simpler or Intermediate Abstractions vs PS

Comment Summary

Question: The PS format is quite complex (involving pseudocode generation, type annotations, higher-order functions). Did we try more intermediate or simpler structures between plain text (OE) and full PS? For example, maybe just storing simpler “signatures” or key-value pairs, or a hybrid approach? Basically, justify that the full complexity of PS is necessary – could a simpler abstraction have worked nearly as well?

Proposed Response

Rebuttal Explanation: We will respond that we did consider and experiment with simpler approaches, and the results pushed us toward the full PS design. For instance, our OE format can be seen as a simpler alternative – it stores free-form tips/hints from solutions. We found that ArcMemo-OE, while helpful, underperformed ArcMemo-PS (structured concepts) by a few points, indicating that the extra structure of PS added value. Additionally, we initially tried naive embedding-based retrieval of memory (a very simple approach) and found it ineffective for these abstract concepts – this motivated the reasoning-based selection. As for the memory format itself: one could imagine storing just a brief “concept name” or a simplified schema. However, we suspect (and our design process confirmed) that without the details in the PS format, the model wouldn’t know how to apply the concept. For example, a concept labeled “color grouping” without the parameterized detail might be too vague for the model to instantiate. That said, we can mention as future work the exploration of other abstraction granularities (maybe something between

OE and PS). We'll justify that given the difficulty of ARC-AGI tasks, we opted for a rich format to maximize performance.

Action Plan

We will supplement our answer with evidence: e.g., "Our results show that the structured PS concepts significantly outperformed a simpler unstructured approach (ArcMemo-PS vs ArcMemo-OE yields +2.7 on Oracle@2). This suggests the added complexity is justified by better performance. We also attempted a purely embedding-based retrieval with simpler memory entries during development and observed it often failed to retrieve relevant information – hence the shift to reasoning-based selection and more descriptive concepts." In the rebuttal, we can also say: "We agree that investigating intermediate levels of abstraction is worthwhile. For instance, one might design a 'lightweight PS' with simpler signatures. In this work, due to time constraints, we focused on the two extremes (OE vs full PS) and found PS gave the best results. We will clarify this and add a discussion in the paper that the complexity of PS was chosen after finding that less structured alternatives did not suffice to capture the needed generality." No new experiment is planned specifically, but if easy, we could consider adding an ablation in which we remove one component of PS (e.g., remove type annotations) to see the effect. However, given time, it's likely enough to rely on the existing OE vs PS comparison.

Priority

Low – This is more of a justification/detail question. Answering it strengthens the paper's methodological clarity but is not critical for acceptance since our results already show PS > OE. We will address it to satisfy the reviewer's curiosity and demonstrate rigor.

Reviewer qa1C – Q3: Performance with Incorrect/Sparse Feedback

Comment Summary

Question: How would the system perform if feedback is sparse or occasionally wrong? For example, if some solution traces we feed into memory are flawed (due to an incorrect self-judgment or an external mistake), will the memory quality degrade over time? Are there mechanisms to guard against introducing bad concepts? Essentially, how robust is ArcMemo to imperfect feedback signals?

Proposed Response

Rebuttal Explanation: Admit that our current implementation does not explicitly handle incorrect feedback – it assumes the feedback/solutions are correct. If a flawed concept were added, it could indeed be detrimental. We'll explain that a robust lifelong learning agent would need to detect and correct such errors. Possible mechanisms: the agent could track the usage of each concept and if a concept consistently leads to failed attempts when used, that's a signal it might be wrong and should be revised or removed. Another idea is requiring multiple confirmations before trusting a concept (e.g., a concept that appears useful in one problem might not be added until it also proves useful in another context, reducing the chance of adding one-off spurious concepts). We'll acknowledge that in truly noisy environments, maintaining memory correctness is challenging and beyond what we tackled. However, since we focused on

domains with reliable automatic feedback (ARC, unit tests for code, etc.), we didn't encounter memory corruption issues in our experiments. We'll mention this as a limitation and note it as an avenue for future work – ensuring memory can handle noisy feedback (perhaps by integrating uncertainty or by human verification for critical entries).

Action Plan

To concretely address this: we add a note in the rebuttal along the lines of, "If feedback were incorrect, ArcMemo could store an incorrect concept which might mislead future reasoning. In our current setup (with reliable oracle feedback) this wasn't an issue. For more realistic scenarios, one could incorporate a validation step for new memory: e.g., cross-check new concepts on hold-out test cases or require that a concept helps solve a similar problem before fully committing it to memory. The agent could also assign a confidence to each memory entry and down-weight or remove entries that are associated with failures. These enhancements are non-trivial and were outside the scope of this work, but we agree they are important for a truly robust system." We will also clarify that if feedback is sparse (not every problem gives a full solution trace), ArcMemo can still function – it just only updates memory on those problems where feedback is available. It might accumulate knowledge slower, but it doesn't fundamentally break; this just reduces the opportunities to learn. We'll mention that as well: "In cases of sparse feedback, the system would only update memory when it receives a clear signal. It would effectively skip memory writes for problems where it can't verify any solution. This means progress would be slower, but it wouldn't introduce wrong information." This demonstrates that we have thought about partial feedback cases too.

Priority

Low – A forward-looking concern. We should address it to show awareness, but it's not something the acceptance hinges on (given our domain already provides good feedback). We'll include the discussion for completeness.

Reviewer qa1C – Q4: Selection Failure Modes

Comment Summary

Question: Discuss failure modes of the memory selection. When does the selection retrieve irrelevant or distracting concepts, and what happens then? Could a poor selection ever make the model perform worse than if it had no memory at all? Do we have any observations of such cases, and how might we mitigate them?

Proposed Response

Rebuttal Explanation: We will respond that, in principle, yes – if the retrieval brings in a concept that doesn't actually apply, it could confuse the model or consume context that could have been used for more relevant info. We'll note, however, that we designed the selection process to minimize this: it uses the model's reasoning to pick concepts that are relevant, and we limit the number of concepts inserted, so the chances of a completely off-base concept dominating the prompt are low. We can mention that in our experiments, ArcMemo with selection almost never underperformed the baseline on average. There might be individual

puzzle instances where memory didn't help or even distract, but these were outweighed by the gains elsewhere (hence the overall improvement). In fact, our ablation Table 2 shows that even when selection is turned off (i.e., using all concepts like a big "cheatsheet"), the performance doesn't surpass the without-selection case on average – selection net helps. That said, the Table 2 caption also notes some regimes where no-selection had a slightly better score due to variance from imperfect selection. We will explicitly cite this to show we are aware that imperfect selection can occasionally hurt. To mitigate such failure modes, one could further refine the selection criteria or incorporate a fallback (if the model is not confident any concept is relevant, maybe it should proceed without memory, for example). We'll discuss that if a concept is highly irrelevant, an intelligent model might just ignore it (treat it as noise), which might be what happened in our runs since we didn't see catastrophic failures. Overall, we'll assure that memory didn't cause dramatic failures in our tests, but acknowledge it's possible and we have room to improve selection fidelity.

Action Plan

We will include in the rebuttal a brief analysis of when selection goes wrong: For instance, "We observed that in a few cases (manually analyzed), the selection module picked a concept that was only tangentially related to the problem at hand. In those instances, the LLM sometimes explored that concept's suggestion fruitlessly, which could waste a bit of the reasoning budget. However, the impact was limited – our model often could backtrack or simply not overly commit to an irrelevant hint. Importantly, across the board we never saw ArcMemo cause a drastic failure that the no-memory baseline would have solved; rather, at worst it might not help on a puzzle or two. This is reflected in Table 2, where the no-memory (cheatsheet) condition occasionally matches or slightly exceeds the memory condition in certain narrow settings, likely due to variance. But in aggregate, selection improves performance." We'll then mention mitigation: "To further reduce such mis retrievals, one could tighten the selection threshold or allow the model to decide to use 'no concept' if none seem applicable. We will note this as an area for future refinement." This demonstrates we've thought about failure cases. No new data is strictly needed, but if we have a concrete example from logs of a retrieved concept that was irrelevant, we might describe it qualitatively.

Priority

Low – This is a nuanced question. We'll address it to satisfy curiosity and show completeness, but as our results show overall gains, this isn't a threatening issue.

Reviewer GA2r (Score: 6) – Presentation: General vs Single Instantiation

Comment Summary

The reviewer finds that the paper is written in a very general way ("general method"), but in reality we only provide a single instantiation and evaluation (on ARC). They aren't asking for a second domain, but suggest improving the storytelling: perhaps be more upfront in the intro that this is demonstrated on ARC, and maybe give more specifics about ARC early on to ground the reader. Essentially a presentation critique: ensure the narrative matches the scope of what's shown.

Proposed Response

Textual Revision: We will adjust the tone of the paper to better align with the single-domain evaluation. In the introduction, we'll explicitly mention that we develop and test ArcMemo on the ARC-AGI benchmark (and briefly explain ARC's nature) to set context. We'll still note that the framework is intended to be general, but we'll avoid over-claiming generality in absence of multi-domain experiments. We can incorporate a bit more about ARC in the intro (why it's a suitable testbed, what it entails) so that readers understand the significance of solving ARC and don't feel misled expecting many domains. This will strengthen the story by tying our method to the challenges of ARC (compositional generalization, abstract reasoning), rather than presenting it in a vacuum.

Action Plan

In practice, we will rewrite portions of the Introduction. For example: currently we might say "we evaluate on a benchmark that stresses compositional generalization (ARC-AGI)" – we'll expand on that: "...on the ARC-AGI benchmark, a challenging set of abstract reasoning puzzles. We focus on ARC-AGI as a testbed because it requires solving a vast space of tasks via recombination of reasoning patterns, making it an ideal playground for concept-based memory." We'll ensure the intro clearly states: "In this work, we instantiate ArcMemo for the ARC domain and demonstrate its benefits there." This way, the reader isn't left thinking we claim broad generality without proof. Also, perhaps in the Discussion we'll temper any broad statements and explicitly call our current evaluation scope. The reviewer specifically said "not asking for a second [instantiation], but maybe more story about ARC in intro," so we will do exactly that. No new experiments needed, just narrative improvement.

Priority

Low – A purely presentation suggestion. We will implement it to make the paper more coherent and reviewer-friendly. It doesn't affect acceptance directly except by improving overall impression.

Reviewer GA2r – Presentation: Implementation Details in Appendix

Comment Summary

The reviewer wants us to provide all details of the implementation (especially in the appendix). They likely found some aspects under-specified (perhaps some hyperparameters, pseudocode, or prompt examples were missing). Ensuring completeness of detail is important for reproducibility.

Proposed Response

Textual Revision: We will commit to including any missing implementation details in the appendix. This includes things like the exact prompts we used for memory write/read, algorithmic pseudocode (if not already present), and any other specifics (e.g., model parameter settings, how we did vision-to-text for ARC images, etc.). Essentially, we'll make the supplementary material as thorough as possible. This addresses the comment that "it would be important to provide all details... in the appendix."

Action Plan

Concretely, we'll update the appendix to cover: (1) Prompt examples: Include the full text of our few-shot prompting templates for concept abstraction and selection (the reviewer rFxx also asked for prompts; we have them ready and will add them verbatim). (2) Algorithm details: We have Algorithm 1 in the main paper, but we can expand in Appendix with a more detailed step-by-step pseudocode of ArcMemo's process, if needed. (3) Hyperparameters and settings: e.g., how many few-shot examples we gave, what temperature and decoding settings for each model (some of these are in Appendix B/C already, but we'll double-check completeness). (4) If any implementation tricks (like how we seeded memory, or how we integrated the Python solutions) weren't fully described, we'll add that too. Basically, the plan is to ensure that a motivated reader could reimplement ArcMemo from our description. We'll mention this in the rebuttal: "We will add the missing implementation details (prompt templates, etc.) to the appendix to ensure the paper is fully self-contained."

Priority

Low – Easy to address. We'll definitely do this as it improves transparency. Not contentious; just something to follow through on.

Reviewer GA2r – Model Specificity (Only o4-mini used)

Comment Summary

The reviewer asks if they understand correctly that we only used OpenAI o4-mini as the base model. They question how much our system relies on o4's capabilities – would it work with other models directly? Essentially, model dependence: is ArcMemo tied to using that specific model, or could any sufficiently strong model benefit?

Proposed Response

Rebuttal Explanation: Clarify that we used o4-mini as our primary solver because it offered a good trade-off of performance vs. cost, but ArcMemo is not inherently tied to that model or API. The framework should work with any LLM that can reason about the tasks. We will say that we expect ArcMemo to help other models too (including open-source ones or even more powerful ones like "Grok-4"/GPT-4). We'll address two scenarios: (1) Stronger models: If we used a model even stronger than o4-mini (like full GPT-4 or Claude Sonnet-4), those models solve more puzzles out-of-the-box, so the absolute boost from memory might be smaller, but we anticipate it could still solve a few extra hard cases by reusing learned concepts. (We might reference that even GPT-4 isn't perfect on ARC; a memory could potentially push it higher.) (2) Weaker models: If we try a weaker base (like we did with an "o3-mini" experiment), we found the model struggled to use the memory effectively. A base model needs a certain level of reasoning skill to benefit – if it's too weak to solve any puzzles initially, it can't generate useful concepts to store. We'll articulate that: ArcMemo doesn't magically enable a non-reasoning model to solve these tasks; it augments an already reasoning-capable model. We'll cite evidence: e.g., our attempts with an open model (DeepSeek-R1) encountered issues (it often produced incomplete solutions due to output length limit). However, that was a technical limitation, not a conceptual one – with an improved version of that model or adjusting context, it could

work. We'll reaffirm that if someone had, say, only an open-source model, they could implement ArcMemo with it; results would depend on that model's baseline performance. In summary, we'll assert the model-agnostic nature of ArcMemo, with the caveat that extremely weak models might not show gains.

Action Plan

We will share any additional experiments or observations we have about using other models: For example, we might include that Table 1 already lists baseline scores for Claude Sonnet-4 and DeepSeek R1 for comparison. We can mention: "As shown, o4-mini's baseline performance (55.17%) is on par with a much larger model (Sonnet-4 at 54.17%), illustrating o4-mini is a strong model albeit smaller. We chose it for cost reasons, but we expect ArcMemo could likewise be applied on Sonnet-4 or GPT-4 with similar relative improvements." If time permits, attempt a limited-run ArcMemo on an alternative model: e.g., use our memory with DeepSeek R1 for a subset of puzzles to see if it improves. Given DeepSeek R1 had some output truncation issues, we might instead consider using Claude (Sonnet-4) via its API for a small test (since that model performed similarly to o4-mini baseline, it would be interesting to see if memory can bump it up). This might be expensive, but even a mini-experiment (say 20 puzzles, with and without memory) could be illustrative. If we manage this, we'll include those results (e.g., "ArcMemo-PS improved Sonnet-4 from X to Y on a sample set"). If not, we'll lean on argumentation and the data we have. We will definitely mention the o3-mini experiment we did: "Using a significantly weaker backbone (o3-mini) actually reduced performance (from 0.26 to 0.22 oracle score) when using memory, implying that the base model must have enough reasoning ability for memory to be helpful. We'll clarify this finding in the paper." This shows we've empirically tested another model (albeit weaker). And for stronger ones, we'll rely on logical expectation and the provided baseline numbers. Summarizing in rebuttal text: "ArcMemo is not restricted to o4-mini; it's applicable to any LLM. We anticipate similar benefits for other models: indeed, our use of GPT-4.1 for some sub-tasks and attempts with DeepSeek R1 demonstrate the framework's flexibility. The key requirement is a competent reasoning model – extremely weak models won't benefit (we observed an o3-level model couldn't utilize memory effectively). Conversely, very strong models might see smaller gains, but since even state-of-the-art models don't solve ARC perfectly, memory could still push their performance on the toughest problems."

Priority

High – Several reviewers (GA2r, rFxx) raised this. It's crucial to assure them ArcMemo isn't an "o4-specific trick" and that it has broader relevance.

Reviewer GA2r – Q1: ArcMemo with Other Base Models (Intuition)

Comment Summary

Question: "How would ArcMemo fare with a base model other than o4-mini? (You can give intuition rather than running new experiments.)" – Essentially the reviewer wants our take on applying ArcMemo to different models, possibly to gauge its generality and requirements (similar to the above, but explicitly asking for reasoning if we haven't tried it).

Proposed Response

Rebuttal Explanation: We will answer largely as above, with reasoning: For a stronger model (like OpenAI's full GPT-4 or Anthropic's Claude 4), ArcMemo should still provide a benefit on those tasks that even those models can't solve outright. We might say: if GPT-4 without memory solves, e.g., ~60% of ARC puzzles, then ArcMemo could help with some of the remaining 40% by learning from them iteratively – we would expect maybe a smaller relative gain, but still a gain. For a different architecture (like an open model such as DeepMind's or Meta's models), as long as it can parse and generate the concept representations, ArcMemo can be applied. We'll reiterate that one needs the base model to be able to reason and to trust its outputs for memory (so a highly unreliable model might store junk). We'll mention the threshold phenomenon: too weak => no gain or harm; moderate to strong => gains. If the reviewer specifically mentioned names like "Sonnet 4 or DeepSeek-R1" (which they did in rFxx review, but GA2r also mentioned "other models"), we'll address those: Sonnet-4 we consider comparable to GPT-4; DeepSeek-R1 we actually tried partially (with issues) so we have insight (we'll mention that attempt and why it struggled, not due to concept memory idea but due to that model's limitations). Overall, provide a thoughtful intuition answer.

Action Plan

This is largely covered by the combined answer for model specificity above. We won't necessarily make a separate action item since the response overlaps: we'll ensure the rebuttal explicitly names Sonnet-4 and DeepSeek-R1 (since GA2r and rFxx brought them up). E.g., "For a model like Claude Sonnet-4, which has similar baseline performance to o4-mini, we'd expect ArcMemo to yield a comparable improvement (our framework could directly wrap around it). For an open model like DeepSeek-R1, we attempted integration: the concept extraction and usage worked in principle, but we hit a technical limit where DeepSeek's output often truncated, making it hard to get correct solutions for memory. With engineering adjustments (like limiting output length or splitting tasks), ArcMemo could potentially boost DeepSeek too. In summary, any model that can benefit from past knowledge – which is most reasoning models – could be paired with ArcMemo. The base model doesn't have to be "o4"; it just needs to be strong enough to solve some tasks and incorporate advice (which is true for major LLMs)." This answers the intuition question directly.

Priority

High – This is essentially part of the model generalization issue – we consider it high priority to address, since it affects the perceived robustness of our approach.

Reviewer GA2r – Q2: Why Not Retry ≥ 3 ?

Comment Summary

Question: Given Table 3's results, is there a reason we stopped at 2 retries? Why not consider 3 or more sequential retries if performance was still climbing?

Proposed Response

Rebuttal Explanation: We will explain that we were balancing diminishing returns vs. computational cost. Each additional retry means another full pass over the puzzle set (for continual learning) or another attempt per puzzle, which substantially increases evaluation time and cost. We observed that going from 0→1→2

retries gave noticeable improvements (especially 0→1 and 1→2), but we anticipated that 2→3 would give a smaller gain. In fact, if we extrapolate from Table 3, the improvement from 1 retry to 2 retries was modest for ArcMemo-PS (about +~2 points on Oracle@2). Another retry might give maybe +1–2 points more. That could be beneficial, but we had to limit scope. We will also mention that practical constraints (e.g., cost of GPT-4 API calls) played a role in choosing 2 as a cutoff during submission. However, we can say that if the reviewers are interested, we could certainly try 3 in the camera-ready, but we expect tapering improvements. Also note, beyond a certain point, memory updates might saturate (there's only so much new info to learn after solving puzzles multiple times). We'll convey that 2 retries already allowed the model to solve most of what it eventually could; a 3rd retry might solve a few more puzzles but with further diminishing returns.

Action Plan

If possible, we will run a quick test for rebuttal: maybe take ArcMemo-PS and do 3 sequential retries on the 100-puzzle set (with continual updates) to see the final score. If it's feasible, we'll report: e.g., "Preliminary result: with 3 retries, ArcMemo-PS reached ~73–74% on Oracle@2, an incremental gain over 2 retries' 70.8%. This confirms diminishing returns." If not, we'll rely on reasoning. We will answer in rebuttal: "We limited to 2 retries in the main experiments due to computational cost and because improvements were leveling off. Based on our results, a 3rd retry likely yields a smaller gain. For instance, memory methods solved diminishing new puzzles by the 2nd pass. We wanted to demonstrate the concept of scaling with inference without pushing into extremely costly regimes. However, if needed, we can explore 3+ retries; we expect the trend of slowly increasing accuracy to continue, but with each retry giving less benefit. In practice, one will choose a retry budget based on available compute." This shows the reviewer we considered it and have a rationale.

Priority

Low – This was a minor curiosity from the reviewer. We'll address it, but it's not a significant concern.

Reviewer rFxx (Score: 4) – Limited Evaluation Scale (100 puzzles, overlapping std)

Comment Summary

The reviewer notes that our experiments used only 100 puzzles from ARC-AGI-1 (presumably the validation subset), which is a relatively small sample. They point out that because of this, the standard deviations of different methods overlap a lot, making it hard to be sure one method is truly better. They suggest that perhaps the limited scale undermines the conclusiveness of the empirical results.

Proposed Response

Proposed Response: We will directly address this by expanding our evaluation in the rebuttal. Using only 100 puzzles was indeed a compromise due to cost (we mentioned that to get stable estimates we used a 100-puzzle random subset). To alleviate this concern, we will either: (a) Report results on the full 400-puzzle ARC-AGI-1 validation set. This would quadruple the number of data points and significantly tighten the

confidence intervals. We expect with 400 puzzles, the differences between ArcMemo-PS and baselines will become statistically clearer (since standard error scales $\sim 1/\sqrt{N}$). Or (b) if full 400 is too expensive to run multiple times, we can at least run single trials on 400 to show the gap in a larger sample, and/or run multiple different 100 subsets to show consistency. Additionally, we will emphasize that even with 100 puzzles, ArcMemo-PS was the only method that outperformed the baseline in every single run setting (as noted in the paper). This suggests the improvement is consistent, not a fluke. We'll explain that overlapping std doesn't mean no improvement – it means we needed many runs to measure it precisely, which we did (averaging results). By increasing the puzzle count, we address that directly.

Action Plan

Action Plan: Re-run main experiments on more puzzles. Ideally, use the full 400-puzzle ARC-AGI-1 val (400 puzzles) for the key comparison of Baseline vs ArcMemo-PS (and possibly vs ArcMemo-OE and Cheatsheet). We will gather Oracle@2 scores for those on 400 puzzles (with perhaps fewer repeated runs if cost is an issue – maybe 3 runs instead of 5, since the larger N reduces variance inherently). Then report: e.g., “On all 400 public val puzzles, ArcMemo-PS achieves X% \pm Y, vs baseline’s W% \pm Z, confirming a significant improvement.” If running all methods on 400 is too much, we might at least do ArcMemo-PS and Baseline. Alternatively, do two different 100-puzzle sets: e.g., report results on “another random 100 puzzles” to show our gains replicate on a different sample. In the rebuttal, we will include a statement like: “We have expanded our evaluation to the full 400-puzzle ARC-AGI-1 validation set. ArcMemo-PS achieves 59.0% (oracle@2) on this set, compared to 55.2% for the no-memory baseline (averaged over runs; std $\sim 2\%$). This larger-scale result confirms a similar ~ 4 point improvement as reported, now with tighter confidence ($p < 0.05$ in a paired test).” (Numbers are hypothetical – we will use actual values obtained.) We will also reassure that we accounted for sampling variance properly (we averaged multiple runs and can report that ArcMemo outperformed baseline in the majority of individual runs, etc.). By doing this, we directly address the reviewer’s concern that our evaluation might have been too limited to trust.

Priority

High – This is tied to the general “limited evaluation” issue. Since rFxx leaned negative partly due to statistical uncertainty, providing more data could swing their opinion. It’s a must-do to strengthen the empirical evidence for our claims.

Reviewer rFxx – Method Complexity vs. Modest Gains

Comment Summary

The reviewer notes that our method introduces a lot of complexity (multiple new components, prompt engineering, etc.), yet the empirical improvement is around 4% absolute (7.5% relative). They question whether the complexity is justified by the results. This implies we should either show bigger gains or argue that even small gains are meaningful given the task difficulty.

Proposed Response

Proposed Response: We will argue that on a very challenging benchmark like ARC-AGI, a ~4 point absolute increase is actually quite significant – it’s roughly the difference between a GPT-3.5-level performance and GPT-4-level on this benchmark. We’ll emphasize that ARC-AGI is extremely hard (even SOTA <60%), so any consistent improvement is valuable. We will also note that our approach is the first to show a continual learning improvement on ARC (to our knowledge), which opens a new capability (test-time learning) rather than just static performance. So the contribution isn’t only the raw percentage gain but demonstrating a novel functionality (lifelong learning via memory) that wasn’t present in the baseline. Additionally, we’ll defend the complexity by breaking down why each component is needed: for example, simple “cheatsheet” memory (just storing raw solutions) didn’t yield consistent improvements, so we had to design more sophisticated abstractions (PS, parameterization) to get reliable gains. The selection mechanism, while complex, was necessary because simpler retrieval failed – and we showed that removing selection hurts performance. In short, each piece of complexity addresses a specific problem (e.g., memory explosion, irrelevant info, etc.), and together they produce the observed gains. We’ll acknowledge the approach is somewhat complex and could perhaps be streamlined, but at this stage, proving the concept required these components. We’ll also hint that future work could simplify or optimize some parts (e.g., automating prompt creation, etc.), but the key was to validate that concept-level memory works. We will ensure the reviewer sees that the improvement, though not massive in percentage, is meaningful and came with a qualitative shift (the model can now improve as it goes, which baseline cannot).

Action Plan

Action Plan: In the rebuttal, point to evidence that our improvement is non-trivial: “ArcMemo’s +7.5% relative gain (55.17→59.33% Oracle@2) is significant given ARC-AGI’s difficulty – for context, this gap is comparable to the performance difference between two generations of model (e.g., the 54% of Claude-4 vs 59% of our method with o4-mini). Also note the improvement persists at all tested compute budgets, indicating a robust effect.” We will then justify complexity: “We agree ArcMemo introduces multiple components, but each was introduced to tackle a specific challenge observed. For instance, a simpler memory (the ‘Cheatsheet’ baseline) sometimes hurt performance due to flooding the context with raw info; by contrast, our modular abstract memory consistently helps. The typed, parameterized format (PS) yielded better results than an unstructured format (OE), showing that complexity brought tangible benefit. We ablated the selection step and found performance dropped without it, confirming its necessity. In summary, while complex, ArcMemo’s design elements were empirically justified to achieve any gain at all on this hard task.” We will also highlight the new capability: “Additionally, ArcMemo enables continual improvement during inference, which the baseline can’t do. This ability to learn from test experience is an important step toward lifelong learning agents – an aspect that isn’t captured solely by the percentage gain, but is a conceptual advancement.” Finally, possibly concede that if the benchmark were larger or we had more compute, we suspect the gap could grow (for instance, ArcMemo might shine even more on ARC-AGI-2 or if allowed more iterations – hinting that the complexity scales better than the baseline). This portrays our complexity as forward-looking investment.

Priority

Medium – This concern influenced rFxx’s borderline stance. By reframing the significance of our results and showing each component’s value, we aim to convince them the complexity is warranted. It’s not as fundamental as novelty/generalization, but it addresses skepticism about our contribution’s weight.

Reviewer rFxx – Stronger Models (Sonnet-4, DeepSeek-R1)

Comment Summary

The reviewer asks if we tried our method on stronger models (specifically “Sonnet 4” or “DeepSeek-R1”). They want to know if ArcMemo provides benefits there, and generally how strong a model needs to be to benefit (and if it must be a “reasoning model”). This overlaps with GA2r’s model dependence question, but rFxx likely wants to see whether our approach helps state-of-the-art models or if it’s mainly compensating for a weaker model.

Proposed Response

Proposed Response: We will respond that we did gather baseline results for those models and considered them. As Table 1 shows, Claude Sonnet-4 (presumably a very strong model) achieved ~54.17% on ARC (Oracle@2) without our method, which is essentially on par with our o4-mini baseline ~55.17%. This indicates o4-mini is already quite a strong model on this task (nearly as good as Sonnet-4) – so our experiments were not on a “weak” model, but on a near state-of-the-art level model. Therefore, ArcMemo’s gains are meaningful at that level. We didn’t fully integrate ArcMemo with Sonnet-4 due to resource constraints (and because it’s proprietary), but we hypothesize it would similarly benefit: since Sonnet-4 didn’t solve all puzzles, it could still utilize memory for the unsolved ones. Regarding DeepSeek-R1 (an open-source large model): we did experiment with it early on as a base for ArcMemo. We found that DeepSeek’s architecture (with an 8k output limit) struggled to produce complete solutions for ARC puzzles, which hindered using it for ArcMemo (incomplete solutions can’t yield good concepts). It wasn’t an issue with memory per se, but with that model’s performance. So we stuck to o4-mini which was more reliable. We’ll mention that if DeepSeek (or a future open model) is improved to handle the tasks (or if we adjust how we query it), ArcMemo could be applied – conceptually it doesn’t matter if the model is OpenAI or open-source. Also answer the “does it have to be a reasoning model” part: Yes, the base model needs to have reasoning capability. ArcMemo doesn’t add reasoning skill out of thin air; it amplifies a model’s ability to recall and reuse strategies. If a model is not a reasoning model at all (say a small language model that can’t do multi-step reasoning), ArcMemo won’t magically enable it to solve ARC. But any modern large model that can do chain-of-thought (GPT-4, Claude, PaLM, etc.) qualifies as a reasoning model, so ArcMemo can work with them. We’ll conclude that ArcMemo is broadly compatible with strong models and we foresee it improving top-tier models on tasks where even they struggle (like ARC-AGI-2, possibly). We’ll also highlight that we did effectively use two different models in our system: o4-mini for solving and GPT-4.1 for selection, plus tried o3-mini and DeepSeek in some capacity – so it’s not tied to a single model’s internals.

Action Plan

Action Plan: We will include these points in the rebuttal, with possibly some new data: If time allows, as mentioned, we might attempt ArcMemo-PS with DeepSeek-R1 as the solver on a small scale (maybe limit puzzles and see if with memory it solves any more than without). If the output truncation can be mitigated by asking it to produce shorter programs or something, we could attempt it. Or we could try ArcMemo with Claude via its API for a handful of puzzles to see if it improves (this might be feasible if we have an API; even a qualitative result could be interesting). If such an experiment yields something, we’ll report it (e.g., “We ran ArcMemo-PS on 20 ARC puzzles using Claude as the base solver: baseline solved 12, ArcMemo solved 14, suggesting memory gave a benefit even to a strong model”). If not, we lean on logic and the partial evidence from baseline scores. We will explicitly cite Table 1 data in rebuttal: “Note that the baseline

performance of o4-mini (55.17%) is already comparable to a much larger model, Claude Sonnet-4 (54.17%). This means our base model was not weak. We expect ArcMemo would similarly aid Sonnet-4 – any unsolved puzzle by Sonnet could potentially be solved after seeing a similar puzzle and storing that concept. Similarly, for DeepSeek-R1 (which had 26.33% baseline, much lower), the limiting factor was that model's ability to produce complete solutions; with improvements, ArcMemo could help it as well. Yes, the model does need to be a reasoning model (o4-mini, Sonnet-4, etc. are) – ArcMemo leverages the model's reasoning to create and use concepts, so a model incapable of multi-step reasoning would not benefit. All current large models with chain-of-thought capabilities qualify, so ArcMemo can be seen as a model-agnostic add-on for those.” We'll ensure the reviewer understands we considered their suggestion and that ArcMemo is meant to boost, not replace, base model reasoning.

Priority

High – This ties into the generalization and significance of our method. Showing (or logically arguing) that it's not limited to one mid-tier model but could help the best models too underscores the broader impact. It's something the skeptical reviewer explicitly asked, so we must address it thoroughly.

Reviewer rFxx – Many Components, Limited Ablation

Comment Summary

The reviewer observes that our method “introduces many components simultaneously” (parameterization, typed interfaces, pseudocode preprocessing, reasoning-based selection, etc.), but we only presented an ablation for the selection component. This leaves it unclear which components contribute most to the gains and which might be unnecessary. They likely want to see more ablation studies or at least a discussion justifying each part.

Proposed Response

Proposed Response: We will expand on ablations and component contributions. First, we'll highlight the comparisons already present that implicitly test some components: We compared ArcMemo-PS vs ArcMemo-OE, which tests the value of the structured, parameterized memory format (PS) over a simpler unstructured memory. The result was that PS performed better (59.33 vs 56.67% Oracle@2), indicating the extra complexity of PS (pseudocode + typed slots) is beneficial. We also ablated selection in Table 2 and found that removing the reasoning-based selection generally hurt performance (ArcMemo-PS without selection was worse, especially in lower compute regimes). So two major components – structured memory design and selection – have evidence of necessity. What we haven't explicitly ablated in results is things like the prompting strategy or pseudocode step vs direct abstraction. We can discuss those: For example, we can say we tried directly abstracting from raw solution traces (skipping the pseudocode intermediate) and found it tended to produce lower-level or messier concepts, so we introduced the pseudocode step to clean it up – a design choice backed by qualitative observations (the pseudocode helps filter out irrelevant details). Similarly, typed interfaces and higher-order functions were included to encourage modularity; while we didn't remove them to test, the fact that ArcMemo-PS outperforms OE suggests that structured fields (types) made a difference. We'll also mention that we will include the full prompt templates in the appendix as they are indeed critical (responding to the prompt engineering part, next point). If possible, we

might add an ablation in the rebuttal for one more component: e.g., run ArcMemo-PS but without the pseudocode abstraction step (just abstract directly from the solution) to see what happens to performance or concept quality. If it drops notably, that quantitatively justifies pseudocode usage. If time allows, we'll do that. But at minimum, we'll provide a reasoned justification for each piece in text. We'll concede that due to space and time, we couldn't run every ablation, but we'll assure that each design decision was informed by either preliminary experiments or logical necessity. We'll also point the reviewer to our qualitative analysis (if any in appendix) that shows PS concepts are cleaner or more general than OE concepts, etc.

Action Plan

Action Plan: In the rebuttal, explicitly list the main components and comment on each: (1) Pseudocode abstraction – "We introduced a pseudocode conversion before abstraction because directly summarizing raw Python code often yielded overly specific or noisy concepts. By first converting solutions to a structured pseudocode, we got more generalized concept suggestions (as noted in Section 3.3). We did not include a separate quantitative ablation of this step due to space, but qualitatively it was crucial for concept quality." (If we manage a quick experiment toggling this, we'll report the outcome to reinforce the point.) (2) Parameterization & typed interfaces – "These were included to allow higher-order concepts and to explicitly indicate where one concept can call/use another. This encourages modular reuse. While we didn't have a separate ablation toggling typing, the performance gap between PS and OE (which lacks such structure) suggests that this added structure helped ArcMemo produce more general, combinable concepts." (3) Reasoning-based selection – "We did ablate this; Table 2 shows removing the reasoning step (i.e., including all concepts in context or using a trivial similarity threshold) generally harms performance. It also increased token usage significantly, confirming why we needed a smarter selection mechanism." (4) Prompt engineering – "We spent effort on designing effective prompts (few-shot examples) for each operation. In the appendix we will include these prompts to show exactly what was done. These are an integral part of the method, though not an 'algorithmic component' per se, their quality affects performance. An ablation for prompts might involve using zero-shot or less tuned prompts; we expect performance would drop if prompts were poorly constructed, but we have not systematically measured that." We'll also commit to including any omitted ablation results if we have them in the appendix (sometimes extra ablations are put in Appendix – ensure none such were missed by reviewer). Summarize: "We will add a more detailed discussion in the revision on the contribution of each design choice and acknowledge the complexity. Each component was chosen based on prior experiments or analyses to address the shortcomings of simpler approaches."

Priority

Medium – The reviewer was borderline; showing that we have carefully considered each part's contribution can assure them the method isn't over-engineered without purpose. We should address this to demonstrate rigor.

Reviewer rFxx – Prompt Engineering Transparency

Comment Summary

The reviewer asks how much prompt engineering was required for the PS format to work, and notes the appendix should include the prompts since they are likely critical. This implies they want assurance that our results are reproducible and not overly hand-tuned, and they want to see the actual prompts.

Proposed Response

Proposed Response: We will be fully transparent about the prompting. We'll state that indeed we developed custom prompts (with a few-shot approach) for both the memory write (concept abstraction) and memory read (selection) steps. This was necessary to guide the model to produce useful abstractions (for example, we gave examples of solution traces and their distilled concepts, to teach the model the format). We'll clarify that prompt engineering effort was non-trivial – we iterated a few times to find prompts that yielded high-quality concepts – but it did not involve any task-specific hard-coding (just general guidance on how to abstract). In other words, we didn't feed the model any direct answers, just a format. And once the prompts were set, they were used consistently across all puzzles. We'll reassure the reviewer that this is documented and we will share the exact prompt text. This addresses any concern that the result might depend on hidden prompt tweaks. Including the prompts in the appendix will allow others to judge how tailored they were. We'll also mention that prompt design is an important part of many LLM-augmented systems, and while it does add to complexity, it's part of the current paradigm (and again, reproducible with our provided examples).

Action Plan

Action Plan: Simply, include full prompt templates in the appendix, as promised. We'll mention in the rebuttal: "We will add the complete prompt templates for ArcMemo's operations to the appendix (for example, the few-shot prompt that converts a solution trace into an abstract concept, and the prompt that selects relevant concepts)." Additionally, answer the question on how much it was required: "We did invest effort in designing these prompts. For instance, the concept abstraction prompt includes 3 exemplars of input solution & output concept, which we curated from the training data. Once set, these prompts were fixed. The method's success is somewhat contingent on these well-crafted prompts – if they were poorly designed, the concept extraction might fail. We acknowledge this reliance; it reflects the state of LLM usage today (prompt engineering is often needed for best results). By providing the prompts, we enable readers to see exactly how the LLM was guided." This honest answer shows we're not hiding the "secret sauce" – it's in the prompts, and we're sharing them.

Priority

Medium – Not critical to the scientific validity (since prompt engineering is expected), but important for transparency and to satisfy the reviewer's request. We'll treat it as a must-do in the camera-ready, and a medium priority to mention in rebuttal.
