

From Generation to Selection: Findings of Converting Analogical Problem-Solving into Multiple-Choice Questions

Anonymous ACL submission

Abstract

As the reasoning abilities of artificial intelligence gain more attention, generating reliable benchmarks to evaluate reasoning capabilities is becoming increasingly important. The Abstract and Reasoning Corpus (ARC) is one of the introduced reasoning benchmarks, providing challenging problems that artificial intelligence has yet to solve. While ARC has been recognized for assessing reasoning abilities, it has a limitation in that its evaluation method through generation fails to consider other aspects of assessment. Bloom’s taxonomy, widely known in education, argues that good evaluation methods should evaluate the six stages of *Remember*, *Understand*, *Apply*, *Analyze*, *Evaluate*, and *Create* in a step-by-step manner. To make ARC, which primarily evaluates the *Create* stage, suitable for assessing stages like *Understand* and *Apply*, we developed MC-LARC. This new multiple-choice format fits well on evaluating large language models (LLMs) across different cognitive stages. We evaluated the analogical reasoning abilities of ChatGPT4V with MC-LARC, confirming that 1) a multiple-choice format can support the language model’s reasoning capabilities and 2) facilitate evidence analysis. However, we noticed LLMs relying on shortcuts when tackling MC-LARC. By analyzing this, we identified areas to consider in multiple-choice synthesis and specified criteria for what constitutes good choices based on these findings.

1 Introduction

Research on artificial intelligence with reasoning capabilities is attracting attention, leading to the proposal of benchmarks to measure such abilities. The Abstraction and Reasoning Corpus (ARC) is one such benchmark designed to evaluate reasoning abilities. Each ARC task consists of 2–5 examples where both input and output are provided, along with one task where only the input is given. The

goal is to infer the rule from the examples and deduce the answer to the task. The input and output grids in ARC can range from a minimum 1×1 grid to a maximum 30×30 grid, with each grid filled with up to 10 different colors. Unlike existing reasoning benchmarks, ARC’s strength lies in its specialization in evaluating reasoning abilities alone by reducing the amount of prior knowledge and data required to solve the tasks.

However, ARC has limitations in that it is an overly difficult benchmark requiring multiple stages of reasoning to solve. According to Bloom’s Taxonomy (Anderson et al., 2001), proposed in traditional educational theory, evaluation consists of the following six stages: *Remember*, *Understand*, *Apply*, *Analyze*, *Evaluate*, and *Create*. In this taxonomy, ARC assesses creation, which encompasses all prior levels of cognitive processes, making it difficult to pinpoint which specific stage may be problematic when a solution is not reached. Even if the logical reasoning process is correct, the entire response is marked wrong if there is a slight error in the generated grid. This issue is also found in derived datasets with reduced difficulty, such as Mini-ARC (Kim et al., 2022) and 1D-ARC (Xu et al., 2023). Although these datasets changed grid sizes or reduced 2D arrays to 1D arrays, it remains difficult to identify which part of the model’s reasoning process is flawed when the task is not solved due to the evaluation format that includes creation. Therefore, a new evaluation method is needed to identify which step of reasoning is problematic in solving ARC.

Therefore, this paper proposes a modified benchmark called MC-LARC to provide an intermediate step in solving ARC tasks. MC-LARC aims to convert the evaluation format from generation to selection, assessing the areas corresponding to *Understand* and *Apply* in Bloom’s Taxonomy. It converts the dataset into a multiple-choice language format by using Large Language Models (LLMs)

to generate four alternative options based on the correct answer to ARC tasks. We conducted experiments to investigate the impact of the transformation into multiple-choice form and found the following two points: 1) The accuracy of LLMs on ARC tasks increased from about 10% to 75%. This indicates that the options in MC-LARC have served a supportive role in the inference of LLMs, which are more aligned with language generation and comprehension than image processing. 2) Evaluating the extent of the inferential abilities of LLMs becomes more clearly feasible. However, it was observed that LLMs used shortcuts while solving MC-LARC, finding the correct answer by considering the form or internal context of the choices to eliminate inappropriate options, rather than utilizing reasoning abilities. Based on this analysis, it was confirmed that when synthesizing data into a multiple-choice format using LLMs, sufficient and accurate context information should be provided to avoid unnecessary additional information. Additionally, this analysis established criteria for what constitutes good multiple-choice options.

2 Related Works

2.1 Evaluation Methods for LLM Abilities Based on Bloom’s Taxonomy

Bloom’s Taxonomy (Anderson et al., 2001) provides a hierarchical classification of cognitive skills that educators can use to structure learning objectives, assessments, and activities. The taxonomy categorizes cognitive skills into six levels as illustrated in Figure 1, each representing a different level of complexity and depth of understanding, from the most basic (*Remembering*) to the most advanced (*Creating*).

By utilizing Bloom’s Taxonomy, educators and researchers can more effectively design, evaluate, and enhance learning experiences and assessments, ensuring that they address all levels of cognitive skills, from basic recall of information to the creation of new and original work.

(Shojaee-Mend et al., 2024) employed Bloom’s Taxonomy to assess the cognitive levels of neurophysiology questions answered by large language models, revealing strengths in basic knowledge recall and weaknesses in higher-order reasoning and knowledge integration. Similarly, (Joshi et al., 2024) used this taxonomy to analyze the cognitive depth of recommendations made by ChatGPT and Bard for teaching Parallel Coordinate Plots.

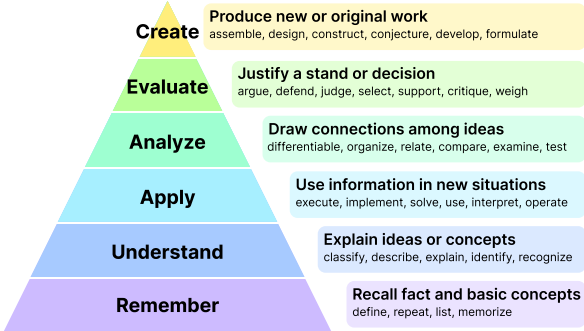


Figure 1: The six cognitive skills in Bloom’s Taxonomy. These skills begin with basic tasks like recalling facts and remembering concepts at the pyramid’s base, progressing to creating original work based on a comprehensive understanding of a specific concept at the top. Image credits: Center for Teaching, Vanderbilt University (Armstrong, 2010).

Human-expert evaluations showed that ChatGPT’s suggestions were generally more appropriate and effective across various cognitive stages, while Bard’s recommendations were often less reliable. Additionally, the BloomGPT project (Spanos et al., 2024) structured a ChatGPT-powered web application around Bloom’s Taxonomy, enhancing students’ cognitive and metacognitive learning in an undergraduate history course. Expert evaluations indicated that the application effectively supported diverse cognitive processes.

2.2 Benchmark for Analogy Abstraction Tasks

Abstraction and Reasoning Corpus (ARC) The Abstraction and Reasoning Corpus (ARC) benchmark (Chollet, 2019) was created for the purpose of measuring intelligence in computer systems. This benchmark requires inference based on complex prior knowledge such as arithmetic abilities, geometric understanding, and topological understanding. The goal is to derive common rules from examples and apply them to infer the appropriate output image for a given test input image. Each task provides 2–5 pairs of example input and output images. The original ARC benchmark consists of 400 training set, 400 evaluation set, and 200 test set. Moreover, the ARC benchmark is represented as 2D matrices.

Language-complete ARC (LARC) The LARC (Acquaviva et al., 2022) dataset consists of 400 ARC training data, each accompanied by 1) a description of the input image and 2) a natural language description of the rules between

167 the input and output images. Both the input
 168 description and the output description must be
 169 language-complete. Language-complete refers to
 170 having sufficient relevant information even when
 171 neither input nor output images are provided. In
 172 other words, humans should be able to understand
 173 the task sufficiently based solely on the description
 174 of LARC without the presence of images. A
 175 language-complete LARC is shown in the Refined
 176 LARC in Figure 2.

177 **Modified Benchmark with Transformed Evalu-**
 178 **ation Format** Abstract and reasoning tasks often
 179 face problems in setting task objectives due to their
 180 attempt to measure unclearly defined reasoning
 181 abilities. Therefore, there have been previous stud-
 182 ies that tried to perform new tasks by modifying
 183 or expanding existing tasks. Bongard-LOGO (Nie
 184 et al., 2020) is an example of simplifying a complex
 185 task. Bongard (Bongard, 1967), one of the Visual
 186 Reasoning benchmarks, is a task that expresses
 187 the difference between two given abstract image
 188 groups as a natural language description. It has
 189 long been a notable task as it requires high abstrac-
 190 tion and reasoning ability to solve the problem, but
 191 it had limitations in analyzing the cause when a spe-
 192 cific model could not solve it, as it is a description
 193 task requiring natural language processing abilities.
 194 To address this, Bongard-LOGO transformed the
 195 type of Bongard problem from a description task
 196 to a classification task. On the other hand, there
 197 are also cases where simple tasks were changed
 198 into complex tasks. VQA (Antol et al., 2015) is a
 199 task that evaluates how well one can answer when
 200 given an image and a question. However, VQA
 201 only assesses whether the given image and natural
 202 language problem are well understood, making it
 203 unsuitable for evaluating reasoning abilities. To
 204 overcome this limitation, a modified benchmark,
 205 TGIF-QA (Jang et al., 2017), which added ques-
 206 tions requiring reasoning about visual images, was
 207 proposed. Thus, especially in the field of Visual
 208 Reasoning, attempts are being made to establish
 209 intermediary results through task transformation.

210 3 Methodology

211 We created MC-LARC through the following two
 212 steps: 1) manually refining the existing LARC, and
 213 2) utilizing ChatGPT4 to generate wrong options
 214 based on LARC.

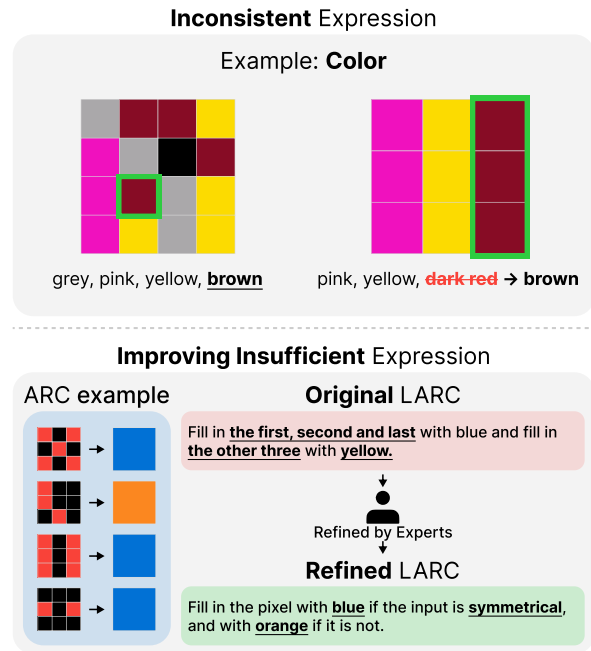


Figure 2: Two main issues of LARC. (Upper part) There are instances where different expressions are used for the same concept within LARC. For example, some LARC expressions describe brown as “dark red”. (Lower part) This task involves identifying the symmetry of the input grid to predict the output image result. However, some original LARC expressions provide insufficient information necessary for ARC problem-solving. These have been revised to contain sufficient and accurate information by experts.

215 **Refining process** The original LARC exhibited
 216 significant quality issues, as evidenced by Figure 2.
 217 These issues appeared primarily in 1) inconsisten-
 218 cies across expressions for the same concept and 2)
 219 a lack of information in the provided explanations.
 220 For instance, the upper part of Figure 2 illustrates
 221 different representations for the same concepts,
 222 leading to confusion. Additionally, the explana-
 223 tions accompanying the tasks often omitted crucial
 224 information necessary for their successful comple-
 225 tion. These issues emerged as a consequence of
 226 the dataset’s compilation by numerous non-experts
 227 using Amazon Mechanical Turk.

228 In addition to the issues highlighted in Figure 2,
 229 there were further cases of inconsistency through-
 230 out the dataset. These inconsistencies involved not
 231 only color but also shape representations and grid
 232 manipulation operations. The presence of these
 233 multiple issues complicates the process of generat-
 234 ing new datasets based on LARC, emphasizing the
 235 challenges of relying on flawed data sources.

236 To address these issues, we conducted a refining

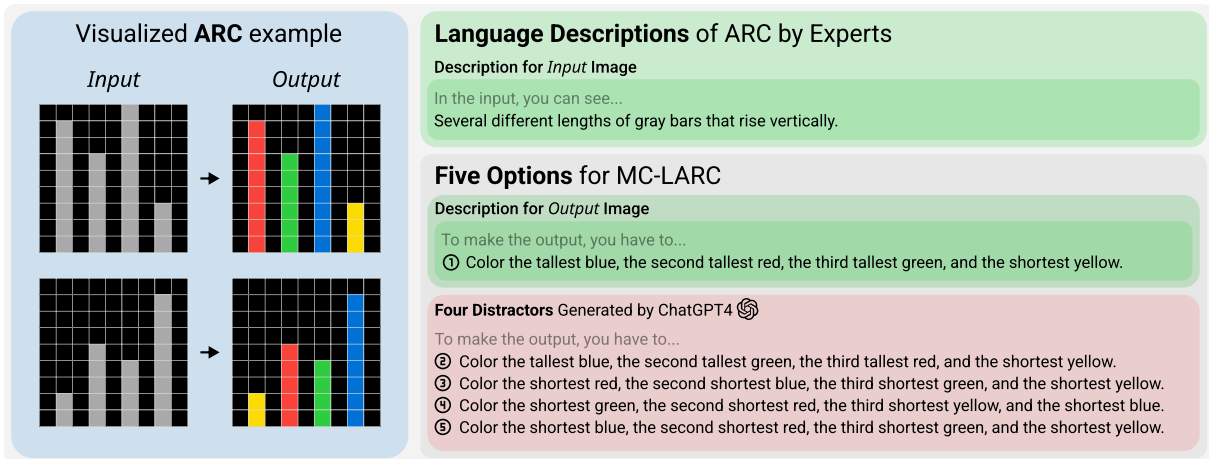


Figure 3: The composition of MC-LARC. It consists of a visualized ARC example and five multiple-choice options. The five multiple-choice options consist of the correct solution and four distractors. (Blue part) It visualizes ARC represented as a 2D matrix. (Green part) It is LARC refined manually by experts. (Red part) Using ChatGPT4, four distractors were generated from the output description (Red boundary) of the refined LARC. To solve MC-LARC, the solver must identify common rules from the visualized ARC example and select the one option from the *Five Options for MC-LARC* that best describes those rules.

process to enhance quality. This process prioritized ensuring consistency in expressions and rectifying erroneous representations. Figure 2 provides an overview of this refining process.

Generating wrong options with ChatGPT4

Based on the given output description of LARC, we generated four distractors through ChatGPT4, as illustrated in Figure 3. However, allowing unrestricted generation of distractors led to issues such as creating out-of-context choices unrelated to the task. To address this problem, we improved by adding constraints during the prompt level. The constraints added to the prompt are as follows:

- **In context vocabulary:** To generate plausible distractors, it was necessary to limit the expressions within the context that aligns with the ARC domain. To achieve this, two contextual constraints were imposed. One involved adding descriptions about the ARC environment, while the other entailed mentioning specific words that should not be used.
- **Length of options:** When generating distractors for lengthy options, there were cases where LLM produced relatively short options, leading to easily solvable problems. Therefore, we restricted the LLM to generate incorrect options of similar lengths to the correct options.
- **Format:** When creating distractors, we en-

sured that the opening phrases of the sentences exactly matched the correct answer option’s ‘*To make the output, you have to...*’. If the opening phrases of the incorrect options vary, it could lead to selecting the correct answer based on the format rather than the meaning of the sentence.

As shown in Figure 7, before constraints were added, the model generated options that were either completely irrelevant to the ARC problem context or altered parts that were not core concepts. These were classified as either bad or moderate. However, after the constraints were applied, the model did not produce any bad options, and the options were classified only as best or moderate. Despite this improvement, the model still faces the challenge of not being able to produce best options for all tasks.

4 Experiments

To verify that the augmented multiple-choice options generated by the LLM did not inadvertently reveal more information than intended, we conducted a control test, as illustrated in Figure 4, where the LLM was presented with only the options, devoid of any accompanying images. If the options were crafted appropriately and free from informational bias, the LLM’s expected accuracy rate would approximate 20%. Additionally, this image-free experiment required the LLM to justify its choice for each option.

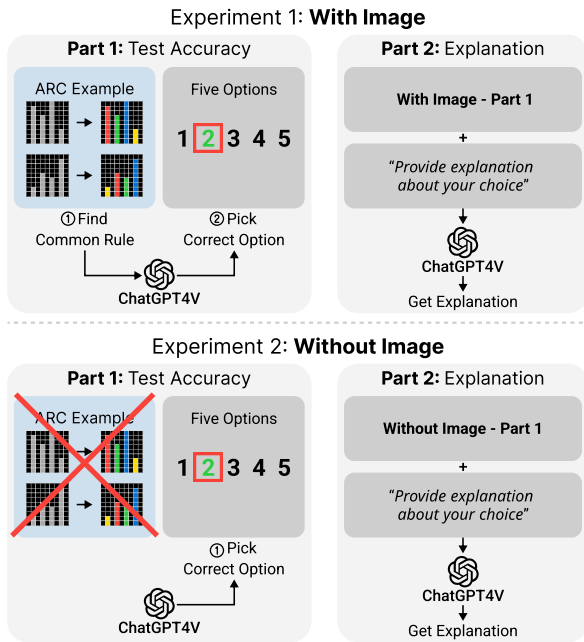


Figure 4: Overview of the conducted experiments. The upper part illustrates the first experiment, which includes visualized ARC example images, while the lower part depicts the second experiment, which does not include these images. Each experiment is divided into two parts. In Part 1, ChatGPT4 is tasked with solving the MC-LARC to measure accuracy. In Part 2, it is requested to provide explanations for its choices, in addition to completing the tasks from Part 1.

4.1 Influence of Multiple Choices

Table 1: A table summarizing the results of experiments where ChatGPT4V solved MC-LARC five times. It shows statistics on the accuracy and Krippendorff’s Alpha score. The statistics show the mean, standard deviation, and 95% confidence interval for the accuracy. Krippendorff’s Alpha score evaluates whether ChatGPT4V’s responses are reliable across the five repeated experiments.

Category	Mean (%)	Std.	95% CI (%)	Alpha
With images	75.81	1.11	74.93 - 76.70	0.8329
Without images	64.61	1.75	63.08 - 66.14	0.7995

For the MC-LARC, we asked the ChatGPT4V model 5 times per problem, and as shown in Table 1, the accuracy of correctly answered tasks out of the total 400 tasks was about 75%. Considering that the accuracy of LLMs on ARC tasks is around 10% (Qiu et al., 2024), this is certainly a high score. Additionally, Krippendorff’s Alpha score of approximately 0.83 confirmed that the LLM was consistently reasoning the answers.

To further evaluate the reasoning process of the

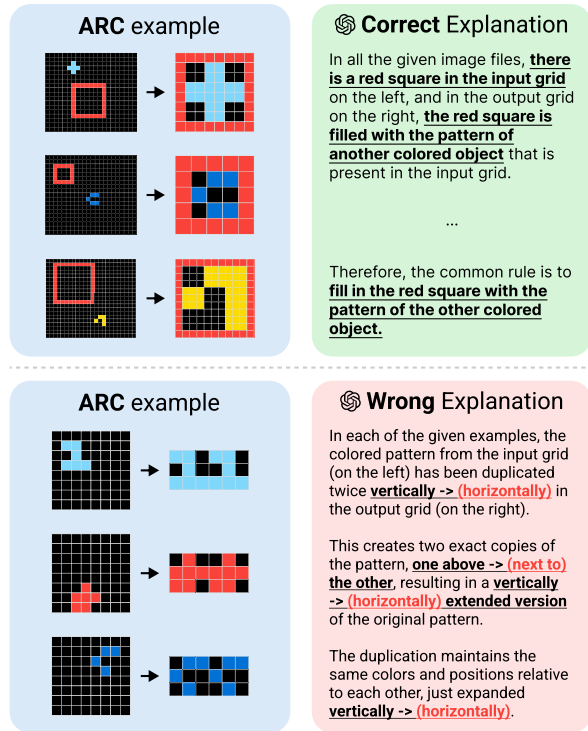


Figure 5: A result of requesting an explanation of the experiments with provided images. (Upper part) It shows an example where the answer to MC-LARC is correctly chosen. (Lower part) It demonstrates the incorrect answers due to failure to infer the correct solution.

LLM, we additionally asked for the reason behind selecting each option. As a result, there were cases where both the answer and the reasoning process were correct or both were incorrect, but there were almost no cases where the answer was correct but the explanation was wrong, or where the answer was wrong but the explanation was correct. This indicates a decrease in the errors of generating correct answers through incorrect reasoning processes or giving inconsistent answers, which tend to occur when LLMs directly solve ARC tasks (Lee et al., 2024). Therefore, even when multiple-choice options, including incorrect options along with the answer description, were provided, we could confirm that the LLM’s reasoning ability was partially improved.

4.2 Problems on Augmentation

However, there were indications that the LLM found a shortcut when solving MC-LARC. MC-LARC should be solved by inferring the rule from the given images and choosing the correct option, but the LLM achieved an accuracy of 65% even when the task was provided without images. The Krippendorff’s Alpha score was also 0.79, not

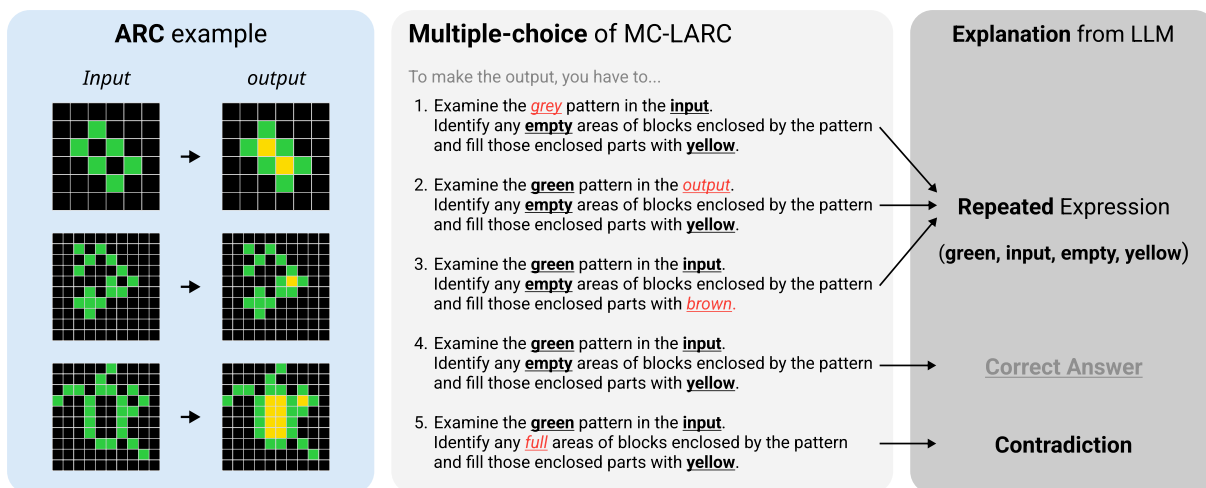


Figure 6: Example of an experiment without an image. When given five options, the LLM solves the problem by analyzing them in the following manner. By examining the options, the LLM identifies repeated expressions and excludes the options that use different vocabulary from the others. Additionally, it excludes options that cannot be represented in the ARC grid by identifying semantic contradiction within the sentences themselves.

331 much lower than the experiment with images provided. This can be understood as evidence that
 332 the LLM found a consistent logic for getting the
 333 correct answers.
 334

335 To analyze how the LLM solved MC-LARC
 336 without the problem images, we additionally asked
 337 the LLM to explain the reasoning behind its
 338 answers. As shown in Figure 6, we found that the
 339 LLM inferred the correct option by 1) choosing the
 340 option with the most repeated expressions and 2)
 341 eliminating options that were self-contradictory.
 342

343 We point out two problems in the generation
 344 process: First, generating four different incorrect
 345 options from one correct option became problem-
 346 atic, as the correct option naturally included more
 347 repeated words than the incorrect options. Second,
 348 not providing image and context information for op-
 349 tion generation led to contradictory or incompatible
 350 expressions in some options. Therefore, from this
 351 experiment, we can conclude that to fairly evaluate
 352 reasoning ability, the process of generating choices
 353 should be improved to avoid providing additional
 information that could serve as a shortcut.

354 4.3 Good Option and Bad Option

355 From the two experiments above, we confirmed
 356 that converting to a multiple-choice format has ad-
 357 vantages as an inference problem in two aspects:
 358 1) providing additional information to solve the
 359 reasoning problem, and 2) allowing for a more
 360 transparent evaluation of the reasoning process.
 361 However, we also found cases where unintended
 362 shortcuts were discovered, and to address this is-

sue, the process of augmenting choices needs to be
 363 improved. But before improving the choice gen-
 364 eration process, this question must be answered
 365 first: What distinguishes a good choice from a bad
 366 choice?
 367

368 As we examined the augmented choice examples
 369 generated by the LLM, we were able to categorize
 370 the choices into three levels of quality, as shown in
 371 Figure 7. The best choices modified the core part of
 372 the problem that fits the context. In ARC, the core
 373 is the part where a change occurs between images,
 374 so in the given examples, completing a square by
 375 filling in orange pixels is the core. Thus, choices
 376 questioning the change to orange can be considered
 377 the best type of choice. Next, choices that were
 378 possible to predict from the input image but did
 379 not capture the core of the problem were of mod-
 380 erate quality. Examples include using colors not
 381 present in the input image or specifying grid sizes
 382 that were not present. Finally, choices that included
 383 cases that cannot occur in the ARC domain at all
 384 were the worst. Commands like “Write an essay”
 385 are irrelevant to ARC and do not require any rea-
 386 soning process to solve the problem, making them
 387 poor choices. Therefore, good text descriptions
 388 and choices should 1) include the core of the prob-
 389 lem in the choices, and 2) be consistent within the
 390 context of the problem. Identifying the criteria in
 391 form and content needed to generate good choices
 392 during the augmentation process is the contribution
 393 of this study.



Figure 7: Three examples of multi-choice options augmented differently by the LLM. The given problem is to fill in an object with holes with the color orange to make a 3×3 square, where the size of the square and the color are the core aspects of the problem. The good example demonstrates an understanding of the core of the problem and provides consistent variations, while the poorer examples increasingly include choices that are unrelated to the problem and inconsistent.

5 Discussion

5.1 Limitations in the Multi-Choice Generation Method

While the experimental results confirmed that the multiple-choice problem format provided sufficient additional information to adequately assess *Understand* and *Apply* aspects, the issue of finding shortcuts during the solving process was raised. This problem is not unique to LLM evaluation. The issue of imbalance among options in multiple-choice questions has already been raised in classical test theory (Alagumalai and Curtis, 2005). The following are suggestions for improving the options in MC-LARC:

- **Option Quality Improvement:** The multiple-choice evaluation method has been criticized for the existence of shortcuts such as *Logical cues*, *Long correct answer*, *Word repeats*, and *Convergence strategy*, even in the case of humans (Case and Swanson, 1998). It has also been pointed out that when there is a lack of discrimination power, the quality of the options decreases. The most intuitive way to address this issue is for humans to consider constraints when creating options.
- **Modification on the Benchmark Format:** Not only the content of the options but also the format of the options can affect the benchmark. Currently, MC-LARC follows a format

where one correct answer option is chosen among five options. On the other hand, another study reported that the selection ratio between options remained similar when there were four or three options compared to five options (Vyas and Supe, 2008). It is also noteworthy that problems with multiple correct answers tend to be more difficult than those with a single correct answer (Case and Swanson, 1998). However, it is not yet known how these various multiple-choice formats differ for LLMs, and therefore, they need to be considered as hyperparameters in the future.

- **Changing the Evaluation Objective:** Modifying the content of the multiple-choice options to measure various areas of reasoning such as application and creation is another possible improvement. Currently, the options in MC-LARC are focused on finding the correct way to solve the ARC task, which is aimed at assessing the understanding of the task. To extend the assessment to other reasoning abilities, the application and creation stages of the task need to be evaluated. Converting the problem into a multiple-choice format where images are selected instead of answer texts, similar to MARVEL (Jiang et al., 2024), could be one possible way to shift the problem format to the creation stage. To transition to the application stage, instead of using an entire problem description, it may be necessary to

454	consider separating the steps required to solve	et al., 2006). Similarly, 1) three or more peo-	503
455	the problem and have the option to select steps	ple can evaluate whether there are errors in the	504
456	that are not necessary for solving the given	options, and 2) the quality of the options can	505
457	ARC task.	be compared with human-created questions.	506
458	5.2 Limitations in the Evaluation	6 Conclusion	507
459	Methodology		
460	One of the current limitations of MC-LARC is	To overcome the limitations of the existing ARC in	508
461	the lack of sufficient evaluation metrics for the	measuring inferential reasoning ability, we created	509
462	proposed benchmark. Therefore, it is difficult to assess	a new multiple-choice dataset called MC-LARC.	510
463	how much the addition of multiple-choice has con-	As a result, the multiple-choice format allowed for	511
464	tributed to securing intermediate reasoning stages	a clearer analysis of logical flow during problem-	512
465	leading up to ARC, and how well the options are	solving and provided supplementary support for	513
466	constructed. The following describes existing meth-	the solver’s reasoning abilities. However, in an	514
467	ods for evaluating options:	additional control experiment without images, we	515
468		found that the LLM solved problems by finding	516
469	• Using Scoring Models: (Ding and Beichner,	shortcuts instead of using reasoning abilities. This	517
470	2009) has proposed statistical and numerical	highlights the regulation needed when using LLMs	518
471	methods for evaluating the quality of multiple-	to synthesize multiple-choice questions. Based on	519
472	choice questions (MCQs). They propose three	these findings, we propose specific conditions for	520
473	methods for individual item evaluation (<i>Item</i>	designing multiple-choice questions that effectively	521
474	<i>Difficulty Level</i> , <i>Item Discrimination Index</i> ,	evaluate the required reasoning abilities without	522
475	<i>Point Biserial Coefficient</i>) and two methods	enabling shortcuts.	523
476	for overall test evaluation (<i>Kuder-Richardson</i>	These findings have several important implica-	524
477	<i>Reliability Index</i> , <i>Ferguson’s Delta</i>). <i>Item</i>	tions. Firstly, they offer valuable insights into	525
478	<i>Difficulty Level</i> and <i>Item Discrimination In-</i>	the appropriate methods for evaluating inferential	526
479	<i>Index</i> measure item difficulty and discrimina-	reasoning, demonstrating the potential of using	527
480	tive power, while <i>Point Biserial Coefficient</i>	multiple-choice questions for this purpose. Sec-	528
481	assesses each item’s appropriateness by compar-	ondly, by identifying the constraints to consider	529
482	ing item scores with the total test score. The	when using LLMs to synthesize multiple-choice	530
483	<i>Kuder-Richardson Reliability Index</i> deter-	questions, this research paves the way for the de-	531
484	mines whether the test is suitable for indi-	velopment of more sophisticated and automated	532
485	vidual or group assessments, and <i>Ferguson’s</i>	high-quality question generators.	533
486	<i>Delta</i> measures the test’s ability to distinguish	7 Limitation	534
487	between varying levels of proficiency. Addi-		
488	tionally, they introduce clustering analysis for	Our study has two main limitations. First, the	535
489	analyzing respondent patterns and model use-	generated options lack quality, allowing LLMs to	536
490	age. Therefore, using metrics to measure the	find shortcuts. Second, there is a lack of metrics	537
491	quality of MCQs is one method for improving	to measure the quality of the options. We have	538
492	MC-LARC.	found issues such as repeated words and contradic-	539
493		tory content in the current multiple-choice options.	540
494	• Comparison with Human-Created Ques-	However, these issues are inherent limitations of	541
495	tions: One issue with the current MC-LARC	multiple-choice questions (Alagumalai and Curtis,	542
496	is that both question generation and evalua-	2005), and therefore, do not undermine the funda-	543
497	tion are done through a single model, Chat-	mental purpose of MC-LARC to assess cognitive	544
498	GPT4V. This evaluation approach does not	features of LLMs such as understanding and appli-	545
499	reveal whether MC-LARC can be properly	cation, which are difficult to confirm solely through	546
500	evaluated on other models, including other	solving ARC problems.	547
501	LLMs. In existing test theory, to compare with	Secondly, our current analysis is limited to the	548
502	human-created options, a large number of peo-	accuracy of LLMs. In existing test theory, met-	549
	ple directly participated in the evaluation to	rics such as discrimination are used to evaluate the	550
	minimize errors as much as possible (Palmer	quality of options. This requires the use of various	551

552	LLMs and analysis of human cases. Nonetheless,	Subin Kim, Prin Phunyaphibarn, Donghyun Ahn, and	602
553	this study lays the foundation for identifying cog-	Sundong Kim. 2022. Playgrounds for Abstraction	603
554	gnitive features that cannot be confirmed through	and Reasoning. In <i>NeurIPS Workshop on nCSI</i> .	604
555	ARC alone, with significant potential for future		
556	expansion.		
557	References		
558	Samuel Acquaviva, Yewen Pu, Marta Kryven,	Seungpil Lee, Woochang Sim, Donghyeon Shin, Sanha	605
559	Theodoros Sechopoulos, Catherine Wong, Gabrielle	Hwang, Wongyu Seo, Jiwon Park, Seokki Lee,	606
560	Ecanow, Maxwell Nye, Michael Tessler, and	Sejin Kim, and Sundong Kim. 2024. Reason-	607
561	Joshua B. Tenenbaum. 2022. Communicating Natu-	ing Abilities of Large Language Models: In-Depth	608
562	ral Programs to Humans and Machines. In <i>NeurIPS</i> .	Analysis on the Abstraction and Reasoning Corpus.	609
563		<i>arXiv:2403.11793</i> .	610
564	Sivakumar Alagumalai and David D. Curtis. 2005. <i>Classical Test Theory</i> . Springer.	Weili Nie, Zhiding Yu, Lei Mao, Ankit B Patel, Yuke	611
565		Zhu, and Anima Anandkumar. 2020. Bongard-Logo:	612
566	Lorin W. Anderson, David R. Krathwohl, Peter W.	A New Benchmark for Human-Level Concept Learn-	613
567	Airasian, Kathleen A. Cruikshank, Richard E. Mayer,	ing and Reasoning. In <i>NeurIPS</i> .	614
568	Paul R. Pintrich, James Raths, and Merlin C. Wit-		
569	trock. 2001. <i>A Revision of Bloom's Taxonomy of Educational Objectives</i> . Pearson.	Edward Palmer, Peter Devitt, et al. 2006. Constructing	615
570		Multiple Choice Questions as a Method for Learning.	616
571	Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Mar-	<i>Annals-Academy of Medicine Singapore</i> , 35(9):604.	617
572	garet Mitchell, Dhruv Batra, C Lawrence Zitnick,		
573	and Devi Parikh. 2015. VQA: Visual Question An-	Linlu Qiu, Liwei Jiang, Ximing Lu, Melanie Sclar,	618
574	swering. In <i>ICCV</i> , pages 2425–2433.	Valentina Pyatkin, Chandra Bhagavatula, Bailin	619
575		Wang, Yoon Kim, Yejin Choi, Nouha Dziri, and Xi-	620
576	Patricia Armstrong. 2010. Bloom's taxonomy. <i>Vander-</i>	ang Ren. 2024. Phenomenal Yet Puzzling: Testing	621
577	<i>bilt University Center for Teaching</i> , pages 1–3.	Inductive Reasoning Capabilities of Language Mod-	622
578		els with Hypothesis Refinement. In <i>ICLR</i> .	623
579	Mikhail Moiseevich Bongard. 1967. The Recognition	Hassan Shojaee-Mend, Reza Mohebbati, Mostafa Amiri,	624
580	Problem. <i>M.: Nauka</i> .	and Alireza Atarodi. 2024. Evaluating the Strengths	625
581		and Weaknesses of Large Language Models in An-	626
582	Susan M. Case and David B. Swanson. 1998. <i>Con-</i>	swering Neurophysiology Questions. <i>Scientific Re-</i>	627
583	structing Written Test Questions for the Basic and	<i>ports</i> , 14(1):10785.	628
584	<i>Clinical Sciences</i> . National Board of Medical Exam-		
585	iners Philadelphia.	Apostolos Spanos et al. 2024. Bloomgpt: Using chatgpt	629
586		as learning assistant in relation to bloom's taxonomy	630
587	François Chollet. 2019. On the Measure of Intelligence.	of educational objectives. In <i>Conference Proceed-</i>	631
588	<i>arXiv:1911.01547</i> .	<i>ings. The Future of Education 2024</i> .	632
589			
590	Lin Ding and Robert Beichner. 2009. Approaches to	Rashmi Vyas and Avinash Supe. 2008. Multiple Choice	633
591	Data Analysis of Multiple-Choice Questions. <i>Phys-</i>	Questions: a Literature Review on the Optimal Num-	634
592	<i>ical Review Special Topics-Physics Education Re-</i>	ber of Options. <i>Natl Med J India</i> , 21(3):130–3.	635
593	<i>search</i> , 5(2):020103.		
594		Yudong Xu, Wenhao Li, Pashootan Vaezipoor, Scott	636
595	Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim,	Sanner, and Elias B Khalil. 2023. LLMs and the Ab-	637
596	and Gunhee Kim. 2017. TGIF-QA: Toward Spatio-	straction and Reasoning Corpus: Successes, Failures,	638
597	Temporal Reasoning in Visual Question Answering.	and the Importance of Object-based Representations.	639
598	In <i>CVPR</i> , pages 2758–2766.	<i>Transactions on Machine Learning Research</i> .	640
599			
600	Yifan Jiang, Jiarui Zhang, Kexuan Sun, Zhivar Sourati,		
601	Kian Ahrabian, Kaixin Ma, Filip Ilievski, and Jay		
602	Pujara. 2024. MARVEL: Multidimensional Abstrac-		
603	tion and Reasoning through Visual Evaluation and		
604	Learning. <i>arXiv:2404.13591</i> .		
605			
606	Alark Joshi, Chandana Srinivas, Elif E. Firat, and		
607	Robert S. Laramée. 2024. Evaluating the recom-		
608	mendations of llms to teach a visualization technique		
609	using bloom's taxonomy. <i>Electronic Imaging</i> , pages		
610	1–8.		
611			