

# Step-by-Step Reasoning to Solve Grid Puzzles: Where do LLMs Falter?

Anonymous ACL submission

## Abstract

Solving grid puzzles involves a significant amount of logical reasoning. Hence, it is a good domain to evaluate reasoning capability of a model which can then guide us to improve the reasoning ability of models. However, most existing works evaluate only the final predicted answer of a puzzle, without delving into an in-depth analysis of the LLMs’ reasoning chains (such as where they falter) or providing any finer metrics to evaluate them. Since LLMs may rely on simple heuristics or artifacts to predict the final answer, it is crucial to evaluate the generated reasoning chain beyond overall correctness measures, for accurately evaluating the reasoning abilities of LLMs. To this end, we first develop *GridPuzzle*, an evaluation dataset comprising of 274 grid-based puzzles with different complexities. Second, we propose a new error taxonomy derived from manual analysis of reasoning chains from LLMs including GPT-4, Claude-3, Gemini, Mistral, and Llama-2. Then, we develop a LLM-based framework for large-scale subjective evaluation (i.e., identifying errors) and an objective metric, *PuzzleEval*, to evaluate the correctness of reasoning chains. Evaluating reasoning chains from LLMs leads to several interesting findings. We further show that existing prompting methods used for enhancing models’ reasoning abilities do not improve performance on *GridPuzzle*. This highlights the importance of understanding fine-grained errors, and presents a challenge for future research to enhance LLMs’ puzzle-solving abilities by developing methods that address these errors<sup>1</sup>.

## 1 Introduction

Recent advancements in LLMs such as GPT-4, Gemini, Claude-3 (Anthropic, 2024), Llama-2 (Touvron et al., 2023), and Mistral (Jiang et al., 2023) have achieved remarkable performance on

a wide range of Natural Language Understanding (NLU) tasks previously thought to be exclusive to human. Beyond NLU, exploring LLMs’ logical reasoning capabilities (Liu et al., 2021; Saparov and He, 2022; Parmar et al., 2024) on complex reasoning tasks such as puzzle-solving is under-explored. Past attempts have been made to evaluate models on logic-intensive grid-based puzzle-solving. However, they either do not focus on evaluating LLMs (Mitra and Baral, 2015; Jabrayilzade and Tekir, 2020) or do not evaluate LLMs independently, but rather use neuro-symbolic approaches (Ishay et al., 2023) that use external specialized solvers on LLM outputs. Here, we aim to evaluate puzzle-solving capabilities of LLMs by themselves, without the use of any external logic solvers.

To understand the reasoning capabilities of LLMs, it is important to evaluate reasoning chains, rather than the final predicted answer. There have been works that evaluate reasoning chains using objective metrics such as ROSCOE (Golovneva et al., 2022), CTC (Deng et al., 2021), and BARTScore (Yuan et al., 2021), however, they do not focus specifically on evaluating reasoning. Some prior works propose metrics for specific reasoning tasks, such as FOLIO (Han et al., 2022) and ProntoQA (Saparov and He, 2022). However, these methods rely on reference-based evaluation, do not focus on puzzle-solving, and do not aim to identify fine-grained errors in reasoning chains. To address these limitations, we propose a reference-free manual and automated subjective evaluation of reasoning chains to understand various fine-grained errors in reasoning chains for grid-based puzzle-solving.

Motivated by Mitra and Baral (2015), we first develop *GridPuzzle* (Figure 1), a comprehensive evaluation dataset consisting of grid-based puzzles with grid-size of  $3 \times 4$ ,  $3 \times 5$ ,  $4 \times 4$ ,  $4 \times 5$ , and  $4 \times 6$  with three levels of difficulty (easy, medium, and hard). Then, we evaluate LLMs including GPT-4, Gemini-Pro, Claude-3, Llama-2, and Mistral

<sup>1</sup>Data and source code are available at <https://anonymous.4open.science/r/GridPuzzle-00B2>

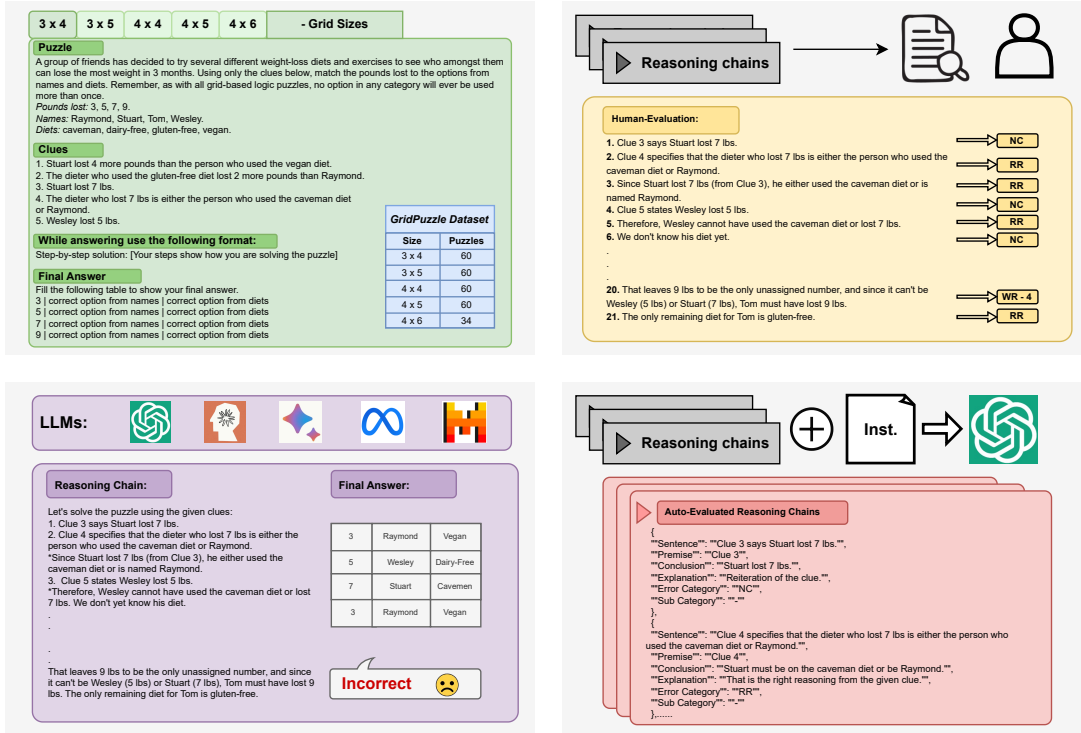


Figure 1: Schematic representation of proposed pipeline. Begins with the data collection of *GridPuzzle* dataset (top left) and evaluating various LLMs in zero-shot CoT setting (bottom left), then analyzing reasoning chains of LLMs manually to find various error types (top right) and automate this analysis process using LLM to check the correctness of reasoning chain by finding errors (bottom right).

on *GridPuzzle* in zero-shot-CoT setting (Figure 1). Experimental results show that LLMs do not fare well and achieve a maximum of 5.1% accuracy.

To investigate reasoning chains, we manually analyze them (Figure 1) to find fine-grained errors (further details in section 3.3). Based on this, we propose a new error taxonomy comprising five broad categories, and nine fine-grained sub-categories (Tables 1 and 2), providing deeper insights into the primary causes of the LLMs’ reasoning failures. However, scaling manual analysis to a larger set is time-consuming and laborious. Hence, we propose to leverage LLMs as auto-evaluators by creating prompts that utilize error taxonomy as prompts to automate the analysis of reasoning chains and help in identifying errors (Figure 1). While evaluating w.r.t. manual annotation, our auto-evaluator model achieves  $\sim 86\%$  agreement, hence providing quality error categorization.

Beyond identifying errors and the accuracy of the final answer, we propose *PuzzleEval*, a LLM-based framework to evaluate reasoning chains for grid-based puzzles. *PuzzleEval* involves a multi-stage evaluation using GPT-4o. First, we identify key logical conclusions from the reasoning chain;

second, we extract key logical concepts from these conclusions; and finally, we measure the presence of these logical concepts in the final gold answer to assess the correctness of the reasoning chain. Evaluating reasoning chains based on error categorization and *PuzzleEval* reveals interesting findings such as LLMs show lower accuracy despite having more error-free reasoning steps, open-source models lack reasoning skills compared to closed-source models, and the most dominant error categories are wrong reasoning and elimination. Additionally, we employ existing prompting methods such as Plan-and-Solve and Self-discover, demonstrating that these methods do not improve performance on *GridPuzzle*. We believe that our findings will inspire future work in the automated, in-depth evaluation of reasoning chains for broader reasoning tasks and enhance the reasoning abilities of models.

## 2 Related Work

**Puzzle-solving Task** Puzzle-solving task provides detailed insights into LLMs’ logical reasoning. Giadikiaroglou et al. (2024) categorize puzzles into (1) rule-based and (2) rule-less puzzles. Rule-less puzzles include riddles (Lin et al., 2021),

MCQs (Zhao and Anderson, 2023), programming puzzles (Schuster et al., 2021), and commonsense reasoning puzzles (Gu et al., 2023); however, in our work we focus on rule-based puzzles. In rule-based puzzles, past attempts have explored Sudoku (Never and Burdick, 2021), Rubik’s Cube, 8-puzzle, Game of 24 (Yao et al., 2024), crosswords (Yao et al., 2024), chess puzzles (Feng et al., 2024), card games (Gupta, 2023), and BoardgameQA (Kazemi et al., 2024). However, grid-based puzzle-solving is under-explored. Mitra and Baral (2015) proposed a grid-based puzzle dataset, but this work does not provide any insights into the performance of recent LLMs. Motivated by this, we propose a systematically curated grid-based puzzle dataset, *GridPuzzle*, and also provide a detailed evaluation of various LLMs in puzzle-solving.

### Automatic Evaluation of Reasoning Chains

Previous works (Dalvi et al., 2021; Saparov and He, 2022; Han et al., 2022) have focused on reference-free evaluation, which is not reliant on gold-reasoning chains. Recently, ROSCOE (Golovneva et al., 2022) proposed a suite of metrics to measure the semantic consistency, logicity, informativeness, fluency, and factuality of reasoning chains, while the ReCEval framework (Prasad et al., 2023) evaluates reasoning chains based on two key properties: correctness and informativeness. Recent evaluation methods such as *LLM evaluation* (Chiang and Lee, 2023) and G-Eval (Liu et al., 2023) leverage LLMs to measure the quality of reasoning chains. *LLM evaluation* involves presenting task instructions and a text sample to LLMs, asking them to rate the sample’s quality on a 5-point Likert scale, whereas the latter incorporates automatic chain-of-thought generated by the LLM describing the detailed evaluation steps. Additionally, Tyen et al. (2023)’s attempt to use GPT-4 as evaluator in a few-shot setting, shows that evaluating reasoning chains remains a challenge. Furthermore, AutoRace (Automatic Reasoning Chain Evaluation) (Hao et al., 2024) proposed a fully automated approach for evaluating reasoning chains that adapt to different tasks without human effort. However, these methods do not evaluate reasoning chains at the level of fine-grained error types and do not provide detailed task-specific insights. To address this, we propose LLM-based reference-free evaluation methods that identify fine-grained errors and assess the correctness of generated reasoning chains for puzzle-solving capabilities.

## 3 Evaluation of Reasoning Chains

### 3.1 GridPuzzle

To develop this dataset, we extract logic grid puzzles of various grid sizes from Puzzle Baron’s Logic Puzzles<sup>2</sup>. Specifically, we compile logic grid puzzles of size  $3 \times 4$ ,  $3 \times 5$ ,  $4 \times 4$ ,  $4 \times 5$ , and  $4 \times 6$ . Each grid size has three levels of difficulty (easy, medium, and hard) except  $4 \times 6$ . This particular grid size has only two difficulty levels (Easy and Medium). Statistics corresponding to each grid size are presented in Figure 1 (top left).

Error Category	Description
WW	Wrong Premise and Wrong Conclusion
WR	Wrong Premise and Right Conclusion
RW	Right Premise and Wrong Conclusion
RR	Right Premise and Right Conclusion
NC	No Conclusion statement or no reasoning involved

Table 1: Proposed error taxonomy for broad categories based on manual analysis. If a statement starts with “so, therefore, hence, this means, this implies, etc.” and/or is not followed by any premise, consider the previous statement’s conclusion or the previous NC as the premise.

### 3.2 Manual Evaluation

To explore where exactly these LLMs falter in performing reasoning, we conduct a detailed manual analysis of the reasoning chains generated by them while solving grid-based puzzles. Our manual analysis process consists of three steps. First, we begin by segmenting the reasoning chains into individual sentences, allowing us to categorize errors more precisely. Second, we identify the premise and conclusion for each sentence and determine their respective correctness. We refrain from subdividing sentences into multiple premises or conclusions to maintain simplicity for finding errors. At last, each sentence is categorized as either containing a single premise and conclusion or being a declarative statement without a conclusion. Afterward, we begin assessing potential issues or errors in the reasoning chains. Now, we follow an exhaustive approach to create fine-grained error categories. We begin with 30 reasoning chains (6 puzzles x 5 reasoning chains from LLMs) to manually identify potential errors. Next, we categorize these errors in a structured format. We then add another 30 reasoning chains to see if any new types of errors emerge. If new errors are identified, we refine our

<sup>2</sup><https://logic.puzzlebaron.com/>

Category	Source	Sub-Category	Description
Wrong Premise or No Conclusion	From the clues (Example: From clue 4,...)	(1) Hallucination	When information is completely out of context and not present in clues.
		(2) Incomplete Information	Lacks necessary information to make a particular conclusion.
		(3) Assumptions	Statements not derived from clues directly; might include assumed information relevant to the clue.
	Derived Conclusions using clues given in puzzle which was not inherently given in the clues.	(4) Error Propagation	Premise derived from a previous incorrect conclusion.
		(5) Incomplete Information	Lacks necessary information to make a particular conclusion.
		(6) Wrong Assumption	The derived assumption is incorrect.
Wrong Conclusion	Derived using the premise (which itself is either taken directly from the clues or derived)	(a) Wrong Reasoning	The reasoning is incorrect, regardless of the premise’s accuracy.
		(b) Error propagation	Conclusion is incorrect due to an erroneous premise.
		(c) Wrong Elimination	All premises are present, but not all conclusions are correctly derived.

Table 2: Proposed error taxonomy for sub-categories based on manual analysis. These sub-categories are defined for cases where either the conclusion or premise is incorrect (“RW” or “WR”) or both are incorrect (“WW”). For “WW”, the error sub-categories might appear in any combinations between (1-6) and (a-c) such as ‘1a’, ‘4b’, or ‘6c’.

categories accordingly. This process is repeated until we evaluate a total of 150 reasoning chains and no new types of errors are found. Based on this method, we have carefully filtered and categorized several errors made by LLMs, presenting them as five broad categories and nine sub-categories.

### 3.3 Proposed Error Taxonomy

**Broad Categories** As shown in Table 1, we present five main categories: “WW” - *Wrong Premise Wrong Conclusion*, “WR” - *Wrong Premise Right Conclusion*, “RW” - *Right Premise Wrong Conclusion*, “RR” - *Right Premise Right Conclusion*, or “NC” - *No Conclusion*. These acronyms of broad categories are self-explanatory. For instance, the category “WW” comprises sentences where a wrong premise leads to a wrong conclusion. Interestingly, we also find the “WR” category consists of instances where a wrong premise still leads to a correct conclusion. Additionally, sentences containing only information from clues or premises from previous steps fall under “NC”. We conduct further investigate why the premises and conclusions becomes incorrect.

**Sub-categories: Wrong Premise** As shown in Table 2, we identified the source of the premise to determine the origin of errors: (i) ‘*From Clues*’ – where the premise is directly borrowed from one of the clues without any further reasoning, and (ii) ‘*Derived*’ – where the premise is inferred from either the clues or the previous conclusions. From Table 2, there are six possible reasons associated with two different sources for wrong premise. When the premise originates from the source (i), we find three types of errors: **Hallucination** – When some factual information from the clues is distorted or completely made up; **Incomplete information** –

When the information is correctly borrowed from the clues but it is not sufficient to make a particular conclusion; **Assumptions** – This is a special category where the premise is not derived but also not given exactly in the clues. It is often related to one of the clues and is of the form, “Let’s assume” or “Assuming that.” When source is *derived*, we find three different errors: **Error Propagation** – This occurs when a previously incorrect conclusion becomes the basis for a flawed premise, thereby extending the error from one conclusion to the next; **Incomplete information** – When the derived premise is not sufficient to make a particular conclusion; and lastly, **Wrong Assumption** – When the LLM reasoner clearly states that a premise was an assumption but it was incorrectly derived.

**Sub-categories: Wrong Conclusion** As shown in Table 2 (source), conclusions are always logically derived from a fixed set of premises. For having a wrong conclusion in any reasoning step, we find three errors responsible: **Error Propagation** – When a conclusion is wrong strictly due to some error in the preceding premise; **Wrong Elimination** – When the conclusion is wrong because the LLM reasoner failed to eliminate all the unfit choices correctly. This case is specific to the task of the grid-based puzzle but is inherently an erroneous deduction on the LLM’s end; **Wrong reasoning** – The remaining incorrect conclusions that did not follow the patterns described in the mentioned categories are classified under this label.

### 3.4 Automated Evaluation

Manual analysis of reasoning chains provides detailed categorization of errors; however, it is tedious and, therefore, challenging to scale for the entire dataset. But analyzing the distribution of



---

**Examples of reasoning chain evaluated by GPT-4o**

---

**Sentence:** "Therefore, Zenon must be 69%, and UCLA must be 62%.",

**Premise:** "If Zenon were 55%, there would be no score 7% lower than 55% for UCLA.",

**Conclusion:** "Zenon must be 69%, and UCLA must be 62%",

**Explanation:** "The conclusion is incorrect as UCLA is already known to be 62% from clue 2.",

**Error Category:** "RW",

**Premise:** "A",

---

**Sentence:** "Since the third performer used flashlights, it must be either Lora or Carmen.",

**Premise:** "The performer who used flashlights was either Lora or Carmen.",

**Conclusion:** "The third performer must be either Lora or Carmen.",

**Explanation:** "The conclusion is based on the incorrect premise that the third performer used flashlights.",

**Error Category:** "WW",

**Sub Category:** "4B",

---

Table 3: Examples of reasoning chain evaluated by Auto-evaluator (GPT-4o).

errors from our proposed taxonomy on the whole dataset is also crucial in understanding shortcomings of LLMs’ reasoning ability. Thus we develop a LLM-based auto-evaluator to automate the process of error evaluation. To this end, we prompt GPT-4o model to identify and categorize errors in the given reasoning chain. Our prompt consists of a system prompt followed by a user prompt containing the reasoning chain to be evaluated along with the original puzzle and its gold solution. The system prompt can be further dissected into 3 key components: the instructions, the knowledge, and an exemplar. The **instruction** contains all the rules that the GPT-4o needs to follow to conduct accurate evaluation and error categorization of the reasoning chains. It incorporates similar sequential steps used during manual evaluation of reasoning chains along with the required output format. The **knowledge** has a detailed description of our error taxonomy including the broad and sub-categories. We also provide a preference order for selecting categories along with the description to minimize any ambiguity in the evaluation process. The **exemplar** consists of a puzzle, its correct solution, and a manually evaluated reasoning chain with our error categories. We termed this LLM-based evaluator as “Auto-evaluator”. The structure of the Auto-evaluator prompt is described in App B.

Using the Auto-evaluator, we evaluated a total of 1,370 reasoning chains generated by five different LLMs for solving 274 puzzles. The application of our Auto-evaluator to this large dataset allowed us

to analyze the distribution of error categories on a broader scale. To validate the accuracy of the evaluations performed by the Auto-evaluator, we randomly sampled 20 reasoning chains from the manually evaluated set. The authors then compared their error category assignments to those given by the Auto-evaluator. The agreement score for the total number of reasoning steps between the manual evaluation and the GPT-4o evaluation is  $\sim 86\%$ . Table 3 shows the example of reasoning steps evaluated by GPT-4o.

## 4 Experimental Steup

### 4.1 Experiments

We evaluate a range of closed-source LLMs including GPT-4-Turbo, Claude-3-Opus, and Gemini-Pro, and open-source models Llama-2-13B, and Mistral-7B-Instruct on *GridPuzzle* in the Zero-shot-CoT setting (Kojima et al., 2022). Our *GridPuzzle* dataset consists of a set of instances denoted as  $\mathcal{P} = \langle p_n^{i \times j}, a_n \rangle$ , where  $p_n^{i \times j}$  is  $n^{th}$  puzzle instance with grid size of  $i \times j$  and  $a_n$  as a gold answer. We prompt each LLM to generate a reasoning chain before predicting answer  $\hat{a}$ . To evaluate each model in the Zero-shot-CoT setting, we provide  $\langle I, p_n^{i \times j} \rangle$  as input to the model and predict an answer  $\hat{a}$  where  $I$  is a natural language instruction. The evaluation is conducted on the OpenAI, Google, and Anthropic model versions released in April 2024 with temperature setting 0 for deterministic predictions. NVIDIA A100 GPUs are used for conducting the inference of open-source models with a batch size of 4. The example prompts used for these experiments are provided in App. A.

### 4.2 Metrics

**Accuracy** We use accuracy to demonstrate the capability of LLMs in solving grid-based puzzles based on their ability to predict the final answer. To calculate this metric, we use the LLM-generated final answers and compare them with the available gold solution. The predicted answers and the gold solution are in the form of tables with the number of rows and columns equal to the grid size of the puzzle. We perform an Exact Match (EM) to compare the two tables and mark them as correct only when all the entries of the tables match. See the example of the final answer table in Appendix C.

**PuzzleEval** We developed this LLM-based metric to assess step-by-step reasoning chains and provide a *correctness score* for each step, as well as

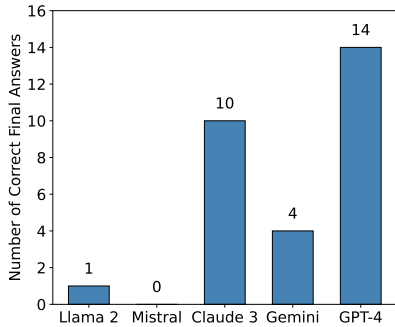


Figure 2: Performance of 5 different LLMs in terms of accuracy on the *GridPuzzle* dataset.

the Average Correctness Score (ACS) for the entire chain. *PuzzleEval* is a reference-free metric specifically designed for assessing reasoning chains generated for grid-based puzzle tasks. It evaluates the correctness of each step in the reasoning chain and reports the score using only the final answer table provided as the gold solution, without requiring any comparison to a gold reasoning chain.

As shown in Figure 13 (refer to App. F), *PuzzleEval* consists of three stage pipeline to evaluate any reasoning chain. First, we prompt GPT-4o to label all the steps sequentially to account for any discrepancies in the different formats of reasoning chains produced by various models, and to extract only the final conclusions from each step. This stage is crucial as it filters out the portion of steps where the models just reiterate clues or previous conclusions. Second, we instruct the model to extract the pair-wise relation of elements from the puzzle that have been either accepted or rejected in the extracted final conclusions. If the extracted conclusion is "Sam is assigned to the year 2015 but not 2014.", these pairs are of the form "Sam – 2015" or "Sam – not 2014". Third, we provide the gold solution table and ask the model to check if these accepted or rejected pairs match the given information. As per the validation, the pairs extracted from every step are marked as correct or incorrect. After obtaining this information for each step the *correctness score* is calculated by adding up all the correct and incorrect steps (correct pairs are marked 1 and incorrect pairs are marked 0) divided by the total number of pairs in each step. Finally, the ACS is determined by adding up all the *correctness scores* from each step and dividing by the number of steps to capture the overall quality of the reasoning chain. Hence, *PuzzleEval* provides ACS for each reasoning chain in range of 0 to 1.

Model	3 x 4	3 x 5	4 x 4	4 x 5	4 x 6	Avg
Llama	0.45	0.46	0.46	0.42	0.28	<b>0.41</b>
Mistral	0.29	0.26	0.27	0.26	0.27	<b>0.27</b>
Claude	0.60	0.56	0.52	0.55	0.46	<b>0.54</b>
Gemini	0.60	0.64	0.54	0.52	0.62	<b>0.58</b>
GPT-4	0.61	0.62	0.56	0.54	0.60	<b>0.59</b>

Table 4: The results for *PuzzleEval* on the different grid sizes available in *GridPuzzle* dataset in terms of ACS.

## 5 Results and Analysis

### 5.1 Objective Evaluation

To evaluate the performance of LLMs when solving grid-based puzzles, we assess the outputs of 5 LLMs using the accuracy and *PuzzleEval*. As shown in Figure 2, we found that all the models have low performance on the *GridPuzzle* dataset in terms of accuracy. The smaller open-source LLMs completely failed at the puzzle-solving task, with Llama-2 solving only one puzzle correctly. Close-source models with significantly larger parameter sizes also exhibited poor performance. GPT-4 had the highest accuracy at only 5.11% (14 puzzles out of 274). Despite the overall low performance of all LLMs, the closed-source models perform marginally better. We evaluate the quality of the reasoning chains using *PuzzleEval*. Table 4 provides the ACS for each grid size available in the *GridPuzzle*. Surprisingly, compared to the accuracy, the performance of the models with *PuzzleEval* was significantly better as shown in Table 4. The ACS lie in the range of 0.24 to 0.62 across all grid sizes. This higher score can be attributed to the partial correctness of reasoning chains when solving the grid-puzzle task. The disparity between metrics shows that evaluating only final answers doesn't fully capture LLMs' effectiveness in complex logical tasks like grid-puzzles.

With the increase in the sizes of the grids, the complexity of the puzzles also rises, leading to a depreciating performance by the LLMs with larger grids. Overall the performance of larger LLMs was much better than the small open-source models. Mistral-7B performed the worst in *PuzzleEval* which is in accordance with its low accuracy score. GPT-4 and Gemini models surprisingly have similar *PuzzleEval* scores (0.59 and 0.58 respectively) despite their large difference in accuracy. This difference in *PuzzleEval* could be attributed to the relatively shorter reasoning chains (fewer reasoning steps) produced by Gemini (an average of 14.91

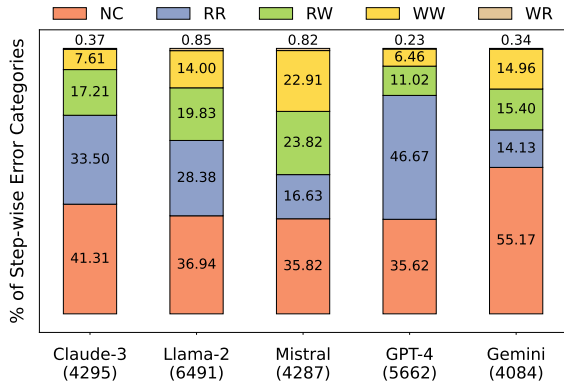


Figure 3: The percentage distribution of the broad error categories across the combined reasoning steps of all 5 LLMs. The total number of steps generated by each model is provided inside the round brackets below the model names.

steps) compared to GPT-4 (an average of 20.66 steps). Shorter reasoning chains may reduce the number of errors that occur while solving the puzzle. It is interesting to note that the smaller LLMs have consistently decreased performance with the increase in the grid size of the puzzles but the larger LLMs have mixed performance.

## 5.2 Reasoning Chain Evaluation

The relative distribution of the broad error categories over the collective reasoning steps for each model is given in Figure 3. It is important to note that, despite using the same zero-shot-CoT setting, the GPT-4 and Llama-2 used significantly more reasoning steps ( $> 5.5k$  steps) to solve the 274 puzzles compared to the other three models ( $\sim 4k$  steps). The distribution of error sub-categories for each model is presented as heatmaps in the first five sub-figures in Figure 4. Here, we present several findings based on evaluation of different error category distribution across *GridPuzzle*.

**Majority of reasoning steps are error-free.** Figure 3 shows that most reasoning steps for each model fall into the “NC” error category, indicating that many steps reiterate the facts or clues from the initial puzzle rather than focusing on reasoning. Over 55% of Gemini-Pro’s reasoning steps fall into this category, the most among all models, suggesting that Gemini spends the fewest steps on actual reasoning. The “RR” category comprises over 46% of GPT-4’s reasoning steps, highlighting its strong reasoning ability. This higher number of correct reasoning steps correlates with GPT-4’s higher *PuzzleEval* score, reflecting its overall effectiveness.

**Why is the accuracy low despite the reasoning chains being mostly error-free?** The disparity between accuracy and *PuzzleEval* arises from the relative location of errors within the reasoning chains. It has been observed that “RR” category reasoning steps mainly occur in the initial half of the chain, leading to a high overall *PuzzleEval* score. Conversely, errors in the “RW”, “WR”, and “WW” categories typically occur in the latter half of the chain, resulting in incorrect final answers and lower accuracy scores. Based on our error taxonomy, the “RW”, “WR”, and “WW” broad error categories have been further dissected into  $6 \times 3$  error sub-categories, with their distribution across reasoning steps shown in Figure 4.

**Dominant error sub-categories and their correlation with the broad categories.** The most common error sub-category across all heatmaps appears to be the “-” category which actually refers to the absence of errors. All the reasoning steps with “NC” and “RR” classifications fall in this category. To observe the actual overall trend across all 5 LLMs, the top 10 most common error sub-categories have been listed in the last sub-figure of Figure 4. The top categories ‘a’ and ‘c’ refer to the **Wrong Reasoning** and the **Wrong Elimination** sub-categories under the “RW” category. These errors arise when LLMs correctly identify the premise but fail to make accurate deductions from it. A number of the top 10 sub-error categories (‘1a’, ‘4a’, ‘4b’, ‘4c’, ‘5a’, ‘5c’, and ‘6a’) emerge from the “WW” category.

For the categories, ‘4a’, ‘4b’, and, ‘4c’ the errors in the premise are propagated from errors in previous reasoning steps showing how initially occurring errors in the chain can lead to more dependent errors. The ‘4b’ error category is the one where this behavior is maximized as here both the premise and conclusions were wrong because of previously propagated errors. The ‘5a’ and ‘5c’ errors occurred due to the incompleteness or lack of information in the premise and wrong reasoning or elimination in the corresponding conclusions. The ‘1a’ kind of error occurred when the premise consisted of hallucinated information. The only sub-category from the “WR” category making it in the top 10 is the ‘3’ category which is caused due to wrong assumptions in the premise. It can be noted here that the reasoning steps of the “WR” category do not deteriorate either of the evaluation metrics, as the conclusions ended up being correct,

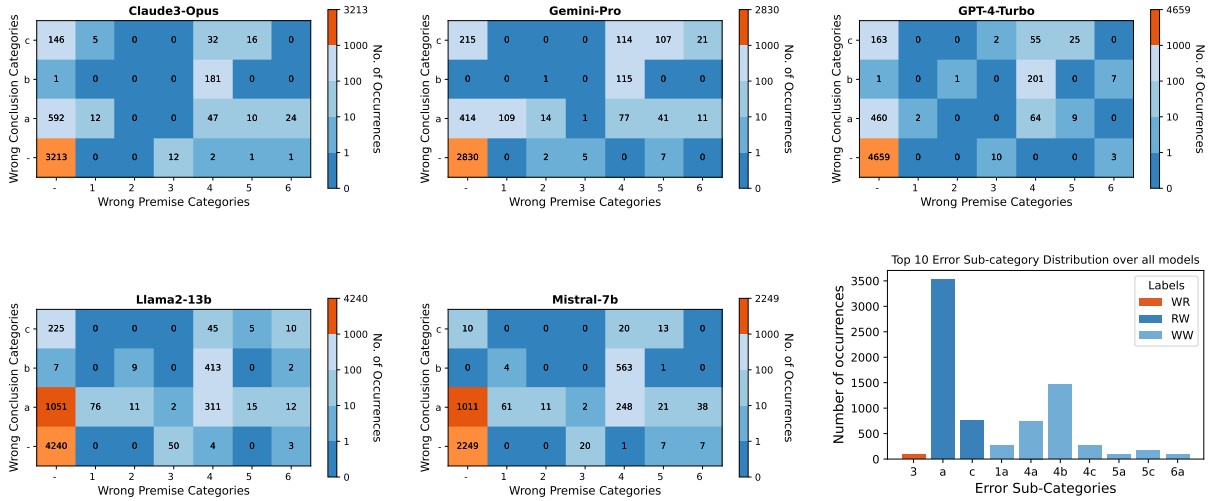


Figure 4: The first five sub-figures in the above section show the error Sub-category distribution over five LLMs. The last sub-figure denotes the top 10 error Sub category distribution across all model reasoning steps.

but rather indicate the inconsistency of the LLMs in reasoning over puzzle-solving.

Mitigation Strategy	Accuracy	<i>PuzzleEval</i>
Baseline	12	0.61
Plan-and-Solve	9	0.62
Self-correct	10	0.59
Self-discover	<b>13</b>	<b>0.65</b>
Feedback-Learning	10	0.59

Table 5: The results for accuracy and *PuzzleEval* using GPT-4-Turbo, with and without mitigation strategies.

**Performance of Open-source vs Proprietary LLMs.** From the results of our objective and subjective metrics, it is evident that the open-source models have lower performance on the grid-puzzle-solving task than the proprietary models. The Llama-2 and Mistral models have the lowest accuracy values and their low performance on the *PuzzleEval* consistently degrades with the increase in the size and complexity of the grids. The Claude-3, Gemini, and GPT-4 models have higher values of accuracy but their performance across the grid sizes in the *PuzzleEval* is inconsistent. The disparity in the performance of both kinds of models can be attributed to the difference in their parameter sizes and the low instruction following capabilities of small open-source models.

**Discussion on Mitigation Strategy** We conduct a case study on a subset of *GridPuzzle* with a 3x4 grid size using widely used prompting techniques used for improving LLMs’ reasoning. In particular, we use four strategies: (1) Plan-and-Solve

(Wang et al., 2023), (2) Self-correct (Zhang et al., 2024), (3) Self-discover (Zhou et al., 2024), and (4) Feedback-Learning. We updated the prompts corresponding to these techniques with our error taxonomy and task. App. D provides further discussion on results presented in Table 5.

## 6 Conclusion

In this work, we evaluated the logical reasoning abilities of LLMs through the lens of a grid-based puzzle-solving task. We introduced *GridPuzzle*, an evaluation dataset of 274 puzzles with various grid sizes. From a manual evaluation of reasoning chains generated by five different LLMs on *GridPuzzle*, we developed a fine-grained error taxonomy with five broad categories and nine sub-categories. We then created an Auto-evaluator to automate the identification of error categories, providing broader insights into error distributions across the dataset. Additionally, we proposed *PuzzleEval*, a reference-free metric to objectively evaluate the correctness of reasoning chains for grid-based puzzles. Our analysis of error distributions in *GridPuzzle* revealed several interesting findings and insights into the logical reasoning abilities of different LLMs. We further evaluated existing reasoning-specific prompting methods, such as self-discover and self-correct, finding that they do not improve results on *GridPuzzle*. We believe our work offers a challenging dataset, highlights where these LLMs make mistakes, and provides insights to develop better logical reasoning systems for complex tasks such as puzzle-solving.



## 589 Limitations

590 While *GridPuzzle* facilitates the evaluation of  
591 LLMs’ logical reasoning abilities, the complexity  
592 of the puzzles can be enhanced by incorporating  
593 further complex grid sizes beyond 4x6. Addition-  
594 ally, this study can be extended to different types  
595 of puzzles, such as Sudoku, Game of 24, and com-  
596 monsense puzzles. Though our study provides fine-  
597 grained error categories, it can be further refined by  
598 mapping to formal logic to identify more detailed  
599 and atomic errors, offering a deeper understand-  
600 ing of LLMs’ reasoning failures. Although we  
601 propose an effective automatic method for error  
602 identification to reduce manual analysis, explor-  
603 ing other automated methods using smaller-scale  
604 supervised learning could be a promising future  
605 research direction. We also note that this research  
606 is currently limited to the English language and can  
607 be extended to multilingual scenarios to evaluate  
608 the logical reasoning abilities of LLMs.

## 609 Ethics Statement

610 We have used AI assistants (Grammarly and  
611 ChatGPT) to address the grammatical errors and  
612 rephrase the sentences.

## 613 References

614 AI Anthropic. 2024. The claude 3 model family: Opus,  
615 sonnet, haiku. *Claude-3 Model Card*.

616 Cheng-Han Chiang and Hung-yi Lee. 2023. [Can large  
617 language models be an alternative to human evalua-  
618 tions?](#) In *Proceedings of the 61st Annual Meeting of  
619 the Association for Computational Linguistics (Vol-  
620 ume 1: Long Papers)*, pages 15607–15631, Toronto,  
621 Canada. Association for Computational Linguistics.

622 Bhavana Dalvi, Peter Jansen, Oyvind Tafjord, Zhengnan  
623 Xie, Hannah Smith, Leighanna Pipatanangkura, and  
624 Peter Clark. 2021. [Explaining answers with entail-  
625 ment trees](#). In *Proceedings of the 2021 Conference  
626 on Empirical Methods in Natural Language Process-  
627 ing*, pages 7358–7370, Online and Punta Cana, Do-  
628 minican Republic. Association for Computational  
629 Linguistics.

630 Mingkai Deng, Bowen Tan, Zhengzhong Liu, Eric Xing,  
631 and Zhiting Hu. 2021. [Compression, transduction,  
632 and creation: A unified framework for evaluating  
633 natural language generation](#). In *Proceedings of the  
634 2021 Conference on Empirical Methods in Natural  
635 Language Processing*, pages 7580–7605, Online and  
636 Punta Cana, Dominican Republic. Association for  
637 Computational Linguistics.

Xidong Feng, Yicheng Luo, Ziyang Wang, Hongrui Tang,  
Mengyue Yang, Kun Shao, David Mguni, Yali Du,  
and Jun Wang. 2024. [Chessgpt: Bridging policy  
learning and language modeling](#). *Advances in Neural  
Information Processing Systems*, 36. 638  
639  
640  
641  
642

Panagiotis Giadikiaroglou, Maria Lymperaioi, Giorgos  
Filandrianos, and Giorgos Stamou. 2024. [Puzzle  
solving using reasoning of large language models: A  
survey](#). *arXiv preprint arXiv:2402.11291*. 643  
644  
645  
646

Olga Golovneva, Moya Peng Chen, Spencer Poff, Mar-  
tin Corredor, Luke Zettlemoyer, Maryam Fazel-  
Zarandi, and Asli Celikyilmaz. 2022. [Roscoe: A  
suite of metrics for scoring step-by-step reasoning](#).  
In *The Eleventh International Conference on Learn-  
ing Representations*. 647  
648  
649  
650  
651  
652

Zhouhong Gu, Zihan Li, Lin Zhang, Zhuozhi Xiong, Si-  
hang Jiang, Xiaoxuan Zhu, Shusen Wang, Zili Wang,  
Jianchen Wang, Haoning Ye, et al. 2023. [Beyond  
the obvious: Evaluating the reasoning ability in real-  
life scenarios of language models on life scapes rea-  
soning benchmark](#) (lsr-benchmark). *arXiv preprint  
arXiv:2307.05113*. 653  
654  
655  
656  
657  
658  
659

Akshat Gupta. 2023. [Are chatgpt and gpt-4 good  
poker players?—a pre-flop analysis](#). *arXiv preprint  
arXiv:2308.12466*. 660  
661  
662

Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting  
Qi, Martin Riddell, Luke Benson, Lucy Sun, Eka-  
terina Zubova, Yujie Qiao, Matthew Burtell, et al.  
2022. [Folio: Natural language reasoning with first-  
order logic](#). *arXiv preprint arXiv:2209.00840*. 663  
664  
665  
666  
667

Shibo Hao, Yi Gu, Haotian Luo, Tianyang Liu, Xiyan  
Shao, Xinyuan Wang, Shuhua Xie, Haodi Ma,  
Adithya Samavedhi, Qiyue Gao, et al. 2024. [Llm  
reasoners: New evaluation, library, and analysis of  
step-by-step reasoning with large language models](#).  
In *ICLR 2024 Workshop on Large Language Model  
(LLM) Agents*. 668  
669  
670  
671  
672  
673  
674

Adam Ishay, Zhun Yang, and Joohyung Lee. 2023.  
[Leveraging large language models to generate an-  
swer set programs](#). In *Proceedings of the 20th In-  
ternational Conference on Principles of Knowledge  
Representation and Reasoning*, pages 374–383. 675  
676  
677  
678  
679

Elgun Jabrayilzade and Selma Tekir. 2020. [LGPSolver  
- solving logic grid puzzles automatically](#). In *Find-  
ings of the Association for Computational Linguistics:  
EMNLP 2020*, pages 1118–1123, Online. Association  
for Computational Linguistics. 680  
681  
682  
683  
684

Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-  
sch, Chris Bamford, Devendra Singh Chaplot, Diego  
de las Casas, Florian Bressand, Gianna Lengyel, Guil-  
laume Lample, Lucile Saulnier, et al. 2023. [Mistral  
7b](#). *arXiv preprint arXiv:2310.06825*. 685  
686  
687  
688  
689

Mehran Kazemi, Quan Yuan, Deepti Bhatia, Najoung  
Kim, Xin Xu, Vaiva Imbrasaite, and Deepak Ra-  
machandran. 2024. [Boardgameqa: A dataset for](#) 690  
691  
692

693	natural language reasoning with contradictory information. <i>Advances in Neural Information Processing Systems</i> , 36.	
694		
695		
696	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. <i>Advances in neural information processing systems</i> , 35:22199–22213.	
697		
698		
699		
700		
701	Bill Yuchen Lin, Ziyi Wu, Yichi Yang, Dong-Ho Lee, and Xiang Ren. 2021. <a href="#">RiddleSense: Reasoning about riddle questions featuring linguistic creativity and commonsense knowledge</a> . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 1504–1515, Online. Association for Computational Linguistics.	
702		
703		
704		
705		
706		
707		
708	Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2021. Logiqa: a challenge dataset for machine reading comprehension with logical reasoning. In <i>Proceedings of the Twenty-Ninth International Conference on Artificial Intelligence</i> , pages 3622–3628.	
709		
710		
711		
712		
713		
714		
715	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruo Chen Xu, and Chenguang Zhu. 2023. <a href="#">G-eval: NLG evaluation using gpt-4 with better human alignment</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2511–2522, Singapore. Association for Computational Linguistics.	
716		
717		
718		
719		
720		
721		
722	Arindam Mitra and Chitta Baral. 2015. <a href="#">Learning to automatically solve logic grid puzzles</a> . In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 1023–1033, Lisbon, Portugal. Association for Computational Linguistics.	
723		
724		
725		
726		
727	David A. Noever and Ryerson Burdick. 2021. <a href="#">Puzzle solving without search or human knowledge: An unnatural language approach</a> . <i>ArXiv</i> , abs/2109.02797.	
728		
729		
730	Mihir Parmar, Nisarg Patel, Neeraj Varshney, Mutsumi Nakamura, Man Luo, Santosh Mashetty, Arindam Mitra, and Chitta Baral. 2024. Towards systematic evaluation of logical reasoning ability of large language models. In <i>the proceeding of the Association for Computational Linguistics (ACL) 2024</i> .	
731		
732		
733		
734		
735		
736	Archiki Prasad, Swarnadeep Saha, Xiang Zhou, and Mohit Bansal. 2023. <a href="#">ReCEval: Evaluating reasoning chains via correctness and informativeness</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 10066–10086, Singapore. Association for Computational Linguistics.	
737		
738		
739		
740		
741		
742		
743	Abulhair Saparov and He He. 2022. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought. In <i>The Eleventh International Conference on Learning Representations</i> .	
744		
745		
746		
747	Tal Schuster, Ashwin Kalyan, Oleksandr Polozov, and Adam Tauman Kalai. 2021. Programming puzzles. <i>arXiv preprint arXiv:2106.05784</i> .	
748		
749		
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	750
		751
		752
		753
		754
		755
	Gladys Tyen, Hassan Mansoor, Peter Chen, Tony Mak, and Victor Cărbune. 2023. Llms cannot find reasoning errors, but can correct them! <i>arXiv preprint arXiv:2311.08516</i> .	756
		757
		758
		759
	Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023. <a href="#">Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models</a> . In <i>Annual Meeting of the Association for Computational Linguistics</i> .	760
		761
		762
		763
		764
		765
	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024. Tree of thoughts: Deliberate problem solving with large language models. <i>Advances in Neural Information Processing Systems</i> , 36.	766
		767
		768
		769
		770
	Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. Bartscore: Evaluating generated text as text generation. <i>Advances in Neural Information Processing Systems</i> , 34:27263–27277.	771
		772
		773
		774
	Yunxiang Zhang, Muhammad Khalifa, Lajanugen Logeswaran, Jaekyeom Kim, Moontae Lee, Honglak Lee, and Lu Wang. 2024. <a href="#">Small language models need strong verifiers to self-correct reasoning</a> .	775
		776
		777
		778
	Jingmiao Zhao and Carolyn Jane Anderson. 2023. Solving and generating npr sunday puzzles with large language models. <i>arXiv preprint arXiv:2306.12255</i> .	779
		780
		781
	Pei Zhou, Jay Pujara, Xiang Ren, Xinyun Chen, Heng-Tze Cheng, Quoc V. Le, Ed H. Chi, Denny Zhou, Swaroop Mishra, and Huaixiu Steven Zheng. 2024. <a href="#">Self-discover: Large language models self-compose reasoning structures</a> .	782
		783
		784
		785
		786

787  
788  
789  
790  
791  
792  
793

## A GridPuzzle dataset sample Puzzle

The *GridPuzzle* dataset contains 274 puzzles of various grid sizes and complexity. A sample puzzle from the dataset along with the Zero-shot-CoT prompt is described in Figure 5. All the puzzles in the dataset have a similar structure with varying numbers of clues.

**Puzzle**  
Jamie never misses an opportunity to watch the night sky for falling stars. Using only the clues below, match the shooting stars to the options from dates and locations. Remember, as with all grid-based logic puzzles, no option in any category will ever be used more than once.

**shooting stars:** 2, 9, 16, 23.  
**dates:** April 7, April 10, April 17, April 22.  
**locations:** Eastbrook, Gilmore City, Isleton, Manchester.

**Clues**

- The outing to Eastbrook logged 7 fewer shooting stars than the trip to Isleton.
- The April 17 outing was in Gilmore City.
- The outing where they saw 16 shooting stars was in Gilmore City.
- The April 10 outing logged 7 more shooting stars than the April 22 outing.

**Attached Instruction Prompt**

While answering use the following format:  
Step-by-step solution: [Your steps show how you are solving the puzzle]  
Final Answer:  
Fill the following table to show your final answer.  
2 | correct option from dates | correct option from locations  
9 | correct option from dates | correct option from locations  
16 | correct option from dates | correct option from locations  
23 | correct option from dates | correct option from locations

**Gold Answer**

2 | April 22 | Eastbrook  
9 | April 10 | Isleton  
16 | April 17 | Gilmore City  
23 | April 7 | Manchester

Figure 5: The prompt structure of a 4 x 4 grid size puzzle from *GridPuzzle* dataset. Every Zero-shot-CoT prompt from the dataset consists of Puzzle, its corresponding Clues, the Instruction for solving the puzzle, along with the Gold solution of the Puzzle.

794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804

## B GPT-4o Auto Evaluator

To expand the reasoning chain evaluation process we prompt the GPT-4o model with a detailed system prompt. The structure of this system prompt is elaborated in Figure 6. The 3 main components of this system prompt are the **Instruction** - similar to the ones given to human evaluators, the **Knowledge** - obtained from the error taxonomy, and an **Exemplar** - consisting of a Puzzle, its Gold Solution, the LLM-generated Reasoning chain, and the evaluated Reasoning Chain.

**Instruction**

**Task Definition:**  
Analyze and evaluate the given reasoning chains into specific reasoning error categories.

**Evaluation Instructions:**  
Analyze each sentence by dividing it into a premise and conclusion, evaluate errors based on the correct solution, and classify according to specific error categories and subcategories.

**Formatting Instructions:**  
Follow a JSON format for the output, detailing each sentence with its premise, conclusion, explanation, error category, and subcategory.

**Knowledge**

**Error Taxonomy Description:**  
Detailed description of each of the error categories and their subcategories.

**Main Categories:**  
Wrong Premise Wrong Conclusion (WW), Wrong Premise Right Conclusion (WR), Right Premise Wrong Conclusion (RW), Right Premise Right Conclusion (RR), and No Conclusion (NC).

**Subcategories:**  
Hallucination, Incomplete Information, Assumptions, Error Propagation, Incomplete Information, Wrong Assumption, Wrong reasoning, wrong elimination.

**Evaluation Example**

**Puzzle:**  
Nicholas and Edith are having a small dinner party...

**Categories:**  
**vintages** : 1984, 1988, 1992, 1996.  
**wines** : Annata Branco, Bianca Flaux, Ece Suss, Vendemmia.  
**types** : gewurztraminer, merlot, pinot noir, riesling.

**Clues:**  
1. The Ece Suss was bottled sometime after the Annata Branco.  
2. The Bianca Flaux was bottled 4 years before the Vendemmia.  
3. ...

**Gold Solution:**  
1984 | Annata Branco | gewurztraminer  
1988 | Ece Suss | pinot noir  
1992 | Bianca Flaux | merlot  
1996 | Vendemmia | riesling

**LLM Reasoning Chain:**  
1. From clue 6, we know that the 1984 bottle is a gewurztraminer. So,...

**Human-evaluated reasoning chain:**

```
[
  {
    "Sentence": "From clue 6, we know that the 1984 bottle is a gewurztraminer.",
    "Premise": "clue 6",
    "Conclusion": "We know that the 1984 bottle is a gewurztraminer.",
    "Explanation": "Reiteration of the clue.",
    "Error Category": "NC",
    "Sub Category": ""
  },...
```

Figure 6: The prompt structure of auto-evaluator. The prompt consists of instructions to evaluate reasoning chains, Knowledge of the error category descriptions, and a human-evaluated reasoning chain example. This system prompt is followed by a user prompt with the puzzle, its gold solution, and the reasoning chain to be evaluated.

Prompt	Model Response
<p><b>Puzzle:</b> A group of friends has decided to try several different weight-loss diets and exercises to see who amongst them can lose the most weight in 3 months. Using only the clues below, match the pounds lost to the options from names and diets. Remember, as with all grid-based logic puzzles, no option in any category will ever be used more than once.</p> <p><b>pounds lost :</b> 3, 5, 7, 9. <b>names :</b> Raymond, Stuart, Tom, Wesley. <b>diets :</b> caveman, dairy-free, gluten-free, vegan.</p> <p><b>Clues:</b> 1. Stuart lost 4 more pounds than the person who used the vegan diet. 2. The dieter who used the gluten-free diet lost 2 more pounds than Raymond. 3. Stuart lost 7 lbs. 4. The dieter who lost 7 lbs is either the person who used the caveman diet or Raymond. 5. Wesley lost 5 lbs.</p> <p>While answering use the following format: <b>Step-by-step solution:</b> [Your steps show how you are solving the puzzle] <b>Final Answer:</b> Fill the following table to show your final answer. 3   correct option from names   correct option from diets 5   correct option from names   correct option from diets 7   correct option from names   correct option from diets 9   correct option from names   correct option from diets</p> <p><b>Gold Solution:</b> 3   Raymond   vegan 5   Wesley   gluten-free 7   Stuart   caveman 9   Tom   dairy-free</p>	<p><b>Step-by-step solution:</b> 1. From clue 3, we know Stuart lost 7 lbs. 2. From clue 1, since Stuart lost 7 lbs, the person on the vegan diet lost 3 lbs (<math>7 - 4 = 3</math>). 3. From clue 4, since Stuart lost 7 lbs, he is either on the caveman diet or he is Raymond. But we know Stuart's name, so he must be on the caveman diet. 4. From clue 2, the person on the gluten-free diet lost 2 more pounds than Raymond. Since we have the pounds lost as 3, 5, 7, and 9, and Stuart already lost 7 lbs, the only way for someone to lose 2 more pounds than another (excluding the 7 lbs already accounted for) is for the gluten-free dieter to have lost 9 lbs and Raymond to have lost 7 lbs. However, since Stuart lost 7 lbs and is not Raymond, this clue helps us deduce the sequence of weight loss but needs to be revisited for exact matching once we have more information. 5. From clue 5, Wesley lost 5 lbs. 6. Now, knowing Wesley lost 5 lbs, Stuart lost 7 lbs, and the vegan diet resulted in a 3 lbs loss, we can deduce that Tom, being the only person left, lost 9 lbs. 7. Since Tom lost 9 lbs and the gluten-free diet is associated with losing 2 more pounds than Raymond, and the only diet left for Wesley (who lost 5 lbs) is dairy-free, it fits that Tom is on the gluten-free diet. 8. This leaves the vegan diet for the person who lost 3 lbs, which hasn't been directly matched to a name yet. Since Stuart is on the caveman diet and lost 7 lbs, Wesley is left with dairy-free and lost 5 lbs, and Tom is on the gluten-free diet and lost 9 lbs, Raymond must be the one who lost 3 lbs and is on the vegan diet.</p> <p><b>Final Answer:</b> 3   Raymond   vegan 5   Wesley   dairy-free 7   Stuart   caveman 9   Tom   gluten-free</p>
Manually Evaluated Reasoning Steps	Error Categories
1. From clue 3, we know Stuart lost 7 lbs.	RR
2. From clue 1, since Stuart lost 7 lbs, the person on the vegan diet lost 3 lbs ( $7 - 4 = 3$ ).	RR
From clue 4, since Stuart lost 7 lbs, he is either on the caveman diet or he is Raymond	RR
But we know Stuart's name, so he must be on the caveman diet.	RR
4. From clue 2, the person on the gluten-free diet lost 2 more pounds than Raymond.	RR
Since we have the pounds lost as 3, 5, 7, and 9, and Stuart already lost 7 lbs, the only way for someone to lose 2 more pounds than another (excluding the 7 lbs already accounted for) is for the gluten-free dieter to have lost 9 lbs and Raymond to have lost 7 lbs.	RW-C
However, since Stuart lost 7 lbs and is not Raymond, this clue helps us deduce the sequence of weight loss but needs to be revisited for exact matching once we have more information.	RR
5. From clue 5, Wesley lost 5 lbs.	RR
6. Now, knowing Wesley lost 5 lbs, Stuart lost 7 lbs, and the vegan diet resulted in a 3 lbs loss, we can deduce that Tom, being the only person left, lost 9 lbs.	WR-2
7. Since Tom lost 9 lbs and the gluten-free diet is associated with losing 2 more pounds than Raymond, and the only diet left for Wesley (who lost 5 lbs) is dairy-free, it fits that Tom is on the gluten-free diet.	RW-C
8. This leaves the vegan diet for the person who lost 3 lbs, which hasn't been directly matched to a name yet.	RR
Since Stuart is on the caveman diet and lost 7 lbs, Wesley is left with dairy-free and lost 5 lbs, and Tom is on the gluten-free diet and lost 9 lbs, Raymond must be the one who lost 3 lbs and is on the vegan diet.	WW-4B

Figure 7: The **top left** section of the figure consists of a 3x4 sample puzzle from the *GridPuzzle* dataset along with the Zero-shot-CoT prompt. Right below the prompt, we have the Gold solution for the corresponding puzzle. In the **top right** section of the figure, we have the Model-generated Reasoning chain to solve this puzzle along with the Final Answer. In this particular instance, the reasoning chain was generated by the GPT-4 model. In the **bottom half** of the figure, we have the manually evaluated reasoning chain, split into individual sentences. The corresponding error categories found in each step is given the **bottom right** column called 'Error Categories.'



Prompt	Llama2-13B Response		
<p><b>Puzzle:</b> Using only the clues below, match the days to the options from guests and rooms. Remember, as with all grid-based logic puzzles, no option in any category will ever be used more than once.</p> <p><b>days :</b> 2, 3, 4, 5. <b>guests :</b> Owens, Parker, Underwood, Zamora. <b>rooms :</b> 114, 124, 314, 320.</p> <p><b>Clues:</b> 1. Underwood is in room 114. 2. The guest in room 320 is staying 1 day fewer than Underwood. 3. Owens is staying for 3 days. 4. Parker is staying 1 day longer than the client in room 314.</p> <p>While answering use the following format: <b>Step-by-step solution:</b> Your steps showing how you are solving the puzzle <b>Final Answer:</b> Fill the following table to show your final answer. 2   correct option from guests   correct option from rooms 3   correct option from guests   correct option from rooms 4   correct option from guests   correct option from rooms 5   correct option from guests   correct option from rooms</p>	<p>Let's break down the clues and solve the puzzle step by step.</p> <p>Step 1: From clue 1, we know Underwood is in room 114. Since Underwood is staying for 2 days (as the days are 2, 3, 4, 5), we can conclude that Underwood is staying for 2 days.</p> <p>Step 2: From clue 2, we know the guest in room 320 is staying 1 day fewer than Underwood. Since Underwood is staying for 2 days, the guest in room 320 is staying for 1 day. This means the guest in room 320 is staying for day 2.</p> <p>Step 3: From clue 3, we know Owens is staying for 3 days. Since Owens is staying for 3 days, Owens must be staying for day 3.</p> <p>Step 4: From clue 4, we know Parker is staying 1 day longer than the client in room 314. Since the days are 2, 3, 4, 5, the client in room 314 must be staying for day 4. Parker is staying 1 day longer, so Parker is staying for day 5.</p> <p>Step 5: Now that we have the days for each guest, we can fill in the table: 2   Zamora   320 3   Owens   (not mentioned) 4   (not mentioned)   314 5   Parker   (not mentioned)</p>		
<p><b>Gold Solution:</b> 2   Zamora   124 3   Owens   314 4   Parker   320 5   Underwood   114</p>	<p><b>Final Answer:</b> 2   Zamora   320 3   Owens   124 4   Underwood   314 5   Parker   114</p> <p>Note: The rooms for days 3, 4, and 5 are not explicitly mentioned in the clues, but we can deduce them based on the information provided.</p>		
Auto-Evaluated Reasoning Steps		Error Categories	Sub Categories
<pre>{   "Sentence": "From clue 1, we know Underwood is in room 114.",   "Premise": "clue 1",   "Conclusion": "we know Underwood is in room 114.",   "Explanation": "Reiteration of the clue.",   "Error Category": "NC",   "Sub Category": "-" }</pre>		NC	-
<pre>{   "Sentence": "Since Underwood is staying for 2 days (as the days are 2, 3, 4, 5), we can conclude that Underwood is staying for 2 days.",   "Premise": "Underwood is in room 114.",   "Conclusion": "Underwood is staying for 2 days.",   "Explanation": "The conclusion is incorrect as it assumes Underwood is staying for 2 days without sufficient information.",   "Error Category": "RW",   "Sub Category": "a" }</pre>		RW	a
<pre>. . . .</pre>		<pre>. . . .</pre>	<pre>. . . .</pre>
<pre>{   "Sentence": "Parker is staying 1 day longer, so Parker is staying for day 5.",   "Premise": "the client in room 314 must be staying for day 4.",   "Conclusion": "Parker is staying for day 5.",   "Explanation": "The conclusion is incorrect as it is based on the wrong premise that the client in room 314 is staying for day 4.",   "Error Category": "WW",   "Sub Category": "4b" }</pre>		WW	4b

Figure 8: The **top left** section of the figure consists of a 3x4 sample puzzle from the *GridPuzzle* dataset along with the Zero-shot-CoT prompt. Right below the prompt, we have the Gold solution for the corresponding puzzle. In the **top right** section of the figure, we have the Model-generated Reasoning chain to solve this puzzle along with the Final Answer. In this instance, the reasoning chain was generated by the Llama2-13b model. In the **bottom half** of the figure, we have the GPT-4o Auto-Evaluated Reasoning chain. The auto-evaluation is done sentence-wise and the output is in a JSON-structured format consisting of 5 components: the *Sentence*, the *Premise*, the *Conclusion*, the *Error category* and the *Sub-category*. The corresponding error categories found in each sentence are given in the **bottom right** columns called ‘Error Categories’ and ‘Sub Categories.’

## C Evaluation of Reasoning Chains

In order to identify the error categories from the erroneous reasoning chains we conducted manual and auto-evaluation of the reasoning chains. The process of manual evaluation has been described in figure 7 and the process of auto-evaluation using GPT-4o has been described in figure 8.

## D Further Discussion on Mitigation Strategy

We updated the prompts corresponding to these techniques to include some of our findings from the reasoning chain evaluations as precautionary instructions. The first strategy is the Plan-and-Solve (Wang et al., 2023) strategy which prompts the model to first generate a plan to solve the given problem and then follow those steps. The second strategy is inspired by the Self-correct (Zhang et al., 2024) method which uses a combination of self-verification and self-refine to improve reasoning. Next, we used the Self-discover (Zhou et al., 2024) technique which is a 2-step structured reasoning. Lastly, we created our prompting technique called “Feedback-Learning” by providing specific feedback system instructions to the LLM based on our error taxonomy. The detailed prompt structure is described in Appendix E. The results of these strategies are provided in Table 5 (main paper). It is evident from the results that prompting-based strategies are not sufficient to significantly improve the LLM reasoning on the grid-puzzle-solving task. Compared to the rest of the strategies, Self-Discover marginally improves the performance on both accuracy and *PuzzleEval*. These results indicate the sore need of developing techniques beyond prompting by having deeper insights from LLMs’ reasoning chains.

## E Mitigation Strategy Prompts

We conducted a study on the 60, 3x4 puzzles present in *GridPuzzle* dataset to try and improve the reasoning abilities of LLMs when solving the grid-puzzle task. We used prompt-based methods, such as the Plan-and-Solve technique, which divides puzzle-solving into planning and solving steps. We also enhanced the solver with insights from our error taxonomy. The prompt structure for this technique is given in figure 9.

Next, we devised our own strategy to improve LLM reasoning by using the top error categories from our findings and teaching the LLM to rectify

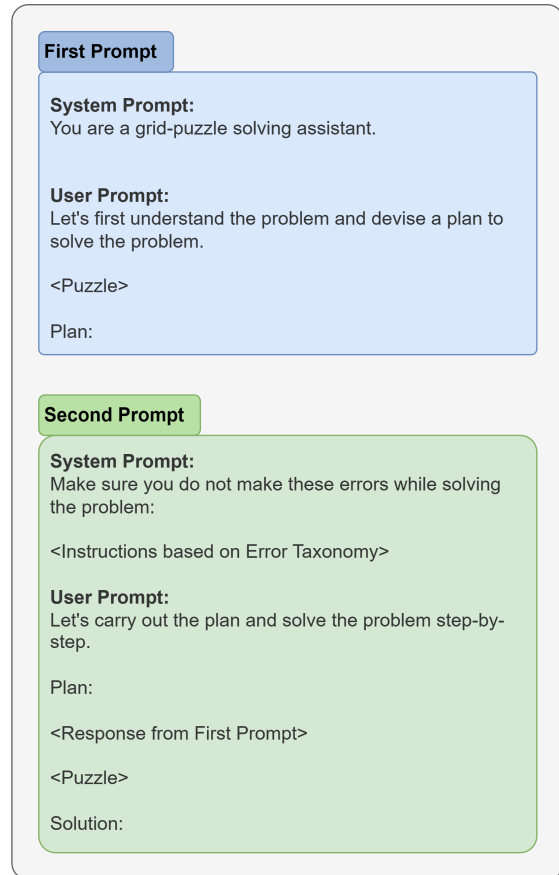


Figure 9: The prompt structure for the Plan-and-Solve strategy which is split into two prompts one for planning and the other for solving the puzzle.

those mistakes. This strategy termed as Feedback-learning makes use of a detailed system prompt that acts as a feedback-providing unit followed by a basic user prompt to solve the puzzle. The prompt structure for this strategy is shown in figure 10.

The Self-correct strategy merges Self-verify and Self-refine qualities to minimize LLM reasoning errors. It starts with solving the puzzle using a Zero-shot-CoT prompt, followed by prompting the LLM to verify and refine the solution. Finally, it integrates the model’s suggestions with insights from our error taxonomy to enhance the puzzle-solving response. The prompt structure for this strategy is shown in figure 11. The Self-Discover strategy, depicted in figure 12, proved most effective in reducing LLM reasoning errors in puzzle-solving. This approach begins by having the model analyze the problem and potential errors, follows with a list of prescribed reasoning modules, prompts the LLM to select and apply the most suitable module, and concludes by using a structured prompt to solve the puzzle.

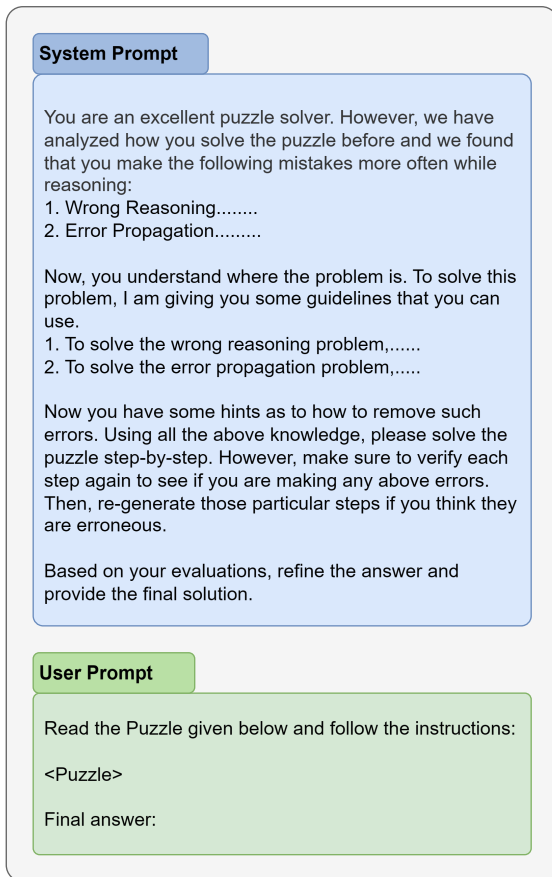


Figure 10: The prompt structure for the Feedback-learning strategy. The system prompts consist of instructions regarding the major errors as well as ways to rectify those errors.

## F Process of Calculating *PuzzleEval*

The overall flow of the process of calculating *PuzzleEval* metrics is described in figure 13. This three-stage process includes: A conclusion extraction stage, followed by a Pair-wise extraction stage and concluded by Validation of the extracted pair-wise relations.

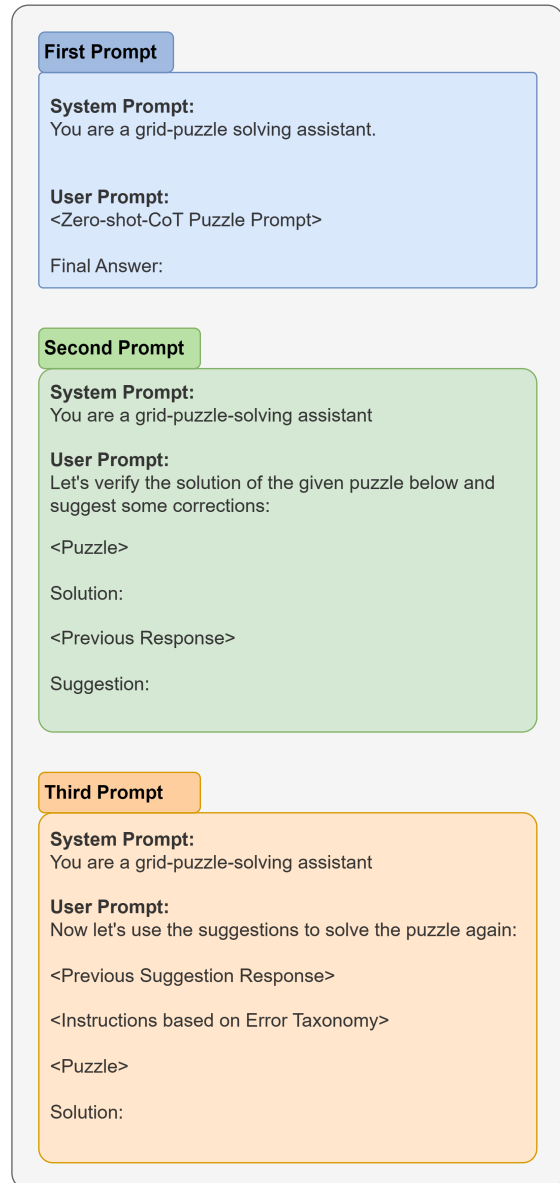


Figure 11: The prompt structure for the Self-Correct strategy is split into 3 parts. The first prompt solves the puzzle, the second prompt verifies the solution and gives suggestions to improve the solution, and the third prompt uses these suggestions along with error taxonomy-based instructions to refine the final solution.

876

877

878

879

880

881

882

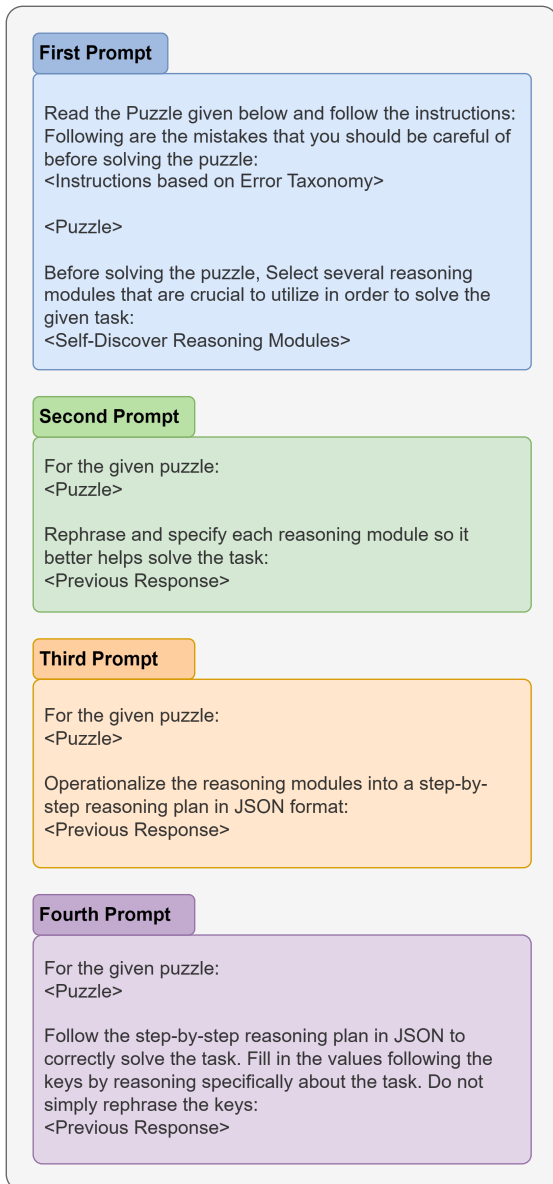


Figure 12: The prompt structure for the Self-Discover strategy. In the first part of this prompt the model is prompted to assess the problem and select the appropriate reasoning module to solve it. Then the module is modified to give a structured plan to solve the puzzle. In the second part, the model uses this structured plan along with instructions from our error taxonomy to solve the puzzle.



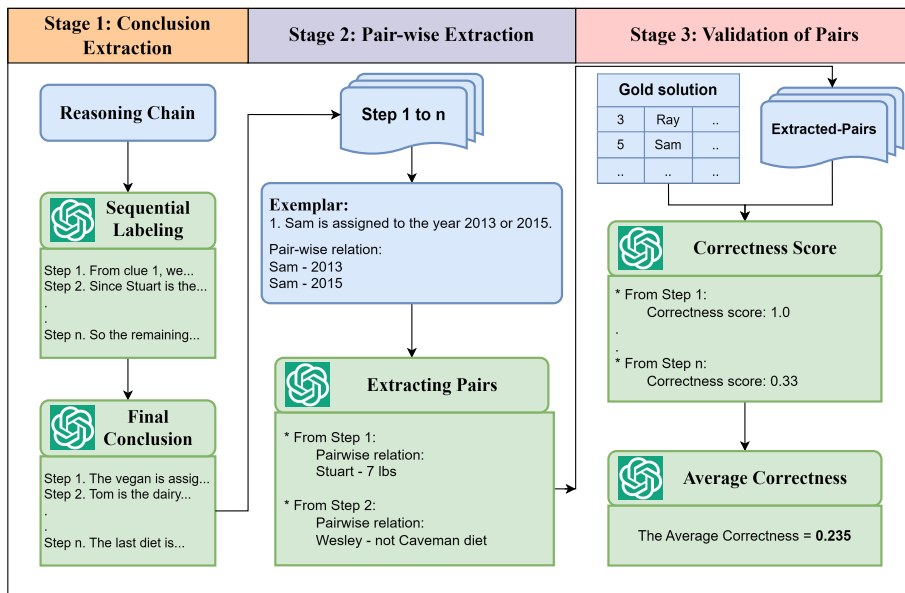


Figure 13: The process of calculating *PuzzleEval* metrics is described above. The reasoning chains are produced by our 5 LLM models and the gold solution is taken from our *GridPuzzle* dataset.