

# GENERATING DATA-DRIVEN REASONING RUBRICS FOR DOMAIN-ADAPTIVE REWARD MODELING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

An impediment to using Large Language Models (LLMs) for reasoning output verification is that LLMs struggle to reliably identify errors in thinking traces, particularly in long outputs, domains requiring expert knowledge, and problems without verifiable rewards. We propose a data-driven approach to automatically construct highly granular reasoning error taxonomies to enhance LLM-driven error detection on unseen reasoning traces. Our findings indicate that classification approaches that leverage these error taxonomies, or “rubrics”, demonstrate strong error identification compared to baseline methods in technical domains like coding, math, and chemical engineering. These rubrics can be used to build stronger LLM-as-judge reward functions for reasoning model training via reinforcement learning. Experimental results show that these rewards have the potential to improve models’ task accuracy on difficult domains over models trained by general LLMs-as-judges by +45%, and approach performance of models trained by verifiable rewards while using as little as 20% as many gold labels. Through our approach, we extend the usage of reward rubrics from assessing qualitative model behavior to assessing quantitative model *correctness* on tasks typically learned via RLVR rewards. This extension opens the door for teaching models to solve complex technical problems without a full dataset of gold labels, which are often highly costly to procure.

## 1 INTRODUCTION

Using Large Language Models (LLMs) to dynamically self-correct thinking traces is a promising avenue for improving reasoning model performance on complex problems. LLM-as-a-judge verification demonstrates performance benefits at inference time, where it is used to perform rejection sampling and value estimation of trajectories (Lightman et al., 2023b; Gu et al., 2025), as well as at training time, where it can function as an outcome reward model whose feedback informs model fine-tuning (Hosseini et al., 2024). However, many LLMs struggle to recognize errors reliably, especially smaller architectures and those lacking domain-specific customization (??). This limitation impacts their viability in training settings and undercuts the perceived data efficiency benefits of LLM judges over other forms of supervised feedback.

Some research suggests that to be effective at most forms of verification, LLMs require either external tools like search engines or large-scale fine-tuning on feedback data (?). These interventions are computationally costly and impact the efficiency of training and inference. Furthermore, feedback signals from approaches like preference mining and self-refinement are purely discriminative, which impacts their robustness to new reasoning outputs. Specifically, these signals do not create explicit mechanisms for determining *why* one output is preferable over another, and so it is difficult to distill a true model of the rules underlying the domain to instill in an LLM judge.

We seek a generative form of external feedback that provides the same benefits as these frameworks. We hypothesize that general-domain LLMs are more adept at recognizing task-specific reasoning errors when given explicit error patterns to check against, as opposed to detecting errors with only abstract guidance as to what they may look like. We draw inspiration from inverse-constitutional AI (Bai et al., 2022), which focuses on the problem of *constructing* constitutions, or lists of desiderata, that LLMs should adhere to. We derive constitutions by directly extracting empirical errors made by a model in a given domain, resulting in a granular and representative knowledge store of

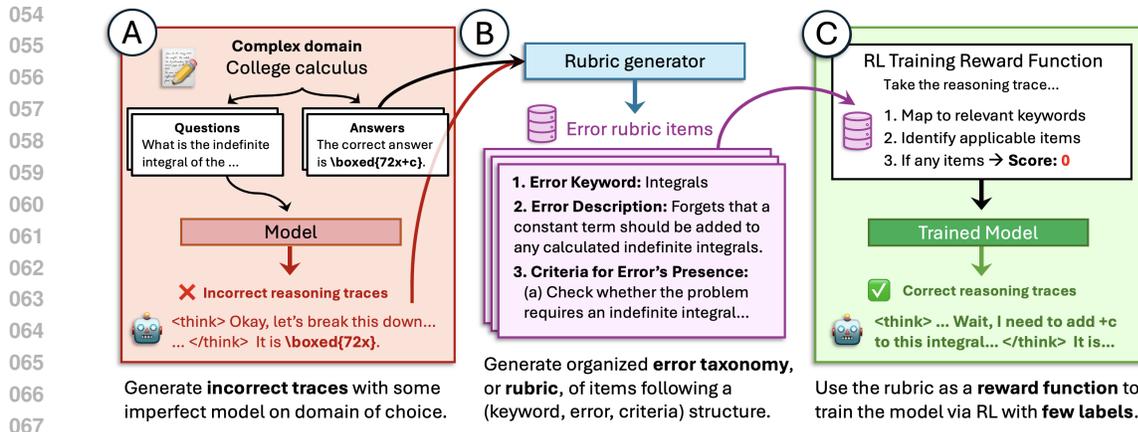


Figure 1: We propose using a knowledge store of errors from inference rollouts to enhance LLM-driven reward functions for model training. We pass incorrect reasoning traces alongside ground-truth answers (Box A) to our “rubric generator” that extracts a set of organized failure modes learned from those incorrect traces (Box B). We then pass this rubric to a LLM classifier that identifies whether a reasoning trace will result in an incorrect answer to serve as a reward function (Box C).

potential errors learned from earlier instances. We refer to this store as a “rubric”: a list of error patterns that an LLM should systematically check for in new reasoning traces. We organize our domain-specific error taxonomies as a keyword-indexed text artifact, and use them to inform LLM verifiers during training, illustrated in Figure 1.

We demonstrate that LLM judges equipped with rubrics directly constructed from individual instances of model failure can improve reasoning trace error classification accuracy by up to 11.6% in technical domains (math, coding, etc.) by improving error recall. We then establish that these enhanced judges can effectively train reasoning models in reinforcement learning settings using minimal ground truth labels, approaching downstream validation accuracy of a Qwen3-4B model fine-tuned with verifiable rewards while using less than 20% as many gold labels. Our findings suggest that rubric-enhanced LLM reward functions represent a promising direction for model training without humans in-the-loop or external feedback providing granular error signals.

In summary, our contributions are:

1. An extension of the inverse constitutional AI task to verifiable domains, in which the desired model behavior targeted by the constitution is *correct reasoning processes* leading to *correct answers*.
2. An approach to automatically generate such constitutions, or rubrics, that generalizes to arbitrary technical domains and uses minimal gold labels.
3. Empirical results demonstrating that these automatically generated rubrics improve trace correctness classification as well as downstream model task accuracy when the rubrics are used to generate RL reward signals.

## 2 RELATED WORK

### 2.1 CONSTITUTIONAL AI

Constitutional AI (CAI, Bai et al. (2022)) leverages a list of explicit governing principles (the “constitution”) to guide AI model behavior during data curation and preference tuning. This paradigm enables AI systems to proactively identify and correct undesirable outputs through self-critique. Constitutions can be used to synthesize preference data from LLMs and then train a “trait preference model” that scores model responses to encourage or discourage the explicitly stated behavior (Kundu et al., 2023). CAI can teach models safety policies (Guan et al., 2024; Mu et al., 2024) or other general principles (Fränken et al., 2024). Our work can be considered a step towards automatically inducing constitutions of undesirable *reasoning* behaviors from a dataset and then using

108 the constitution to inform a preference model. Other work has explored methods to discover traits  
109 that separate good behavior from bad behavior using a dataset, e.g. VibeCheck (Dunlap et al., 2025),  
110 which identifies user-aligned “vibe” statements, while Inverse Constitutional AI (ICAI; Findeis et al.  
111 (2025)) and C3AI (Kyrychenko et al., 2025) derive NL principles from annotated preference data.  
112 In contrast, our method uses data annotated only with end-task correctness, not preference, in mind.  
113

## 114 2.2 ERROR TAXONOMIES

115  
116 A number of papers in topics adjacent to technical reasoning have touched on this notion of “er-  
117 ror taxonomies” beyond work in constitutional AI. Many of these papers derive similar ideas from  
118 agentic pipelines and robotics, establishing frameworks that are optimized for these domains but can  
119 still inspire more “reasoning”-centric work. REFLECT (Liu et al., 2023) collects hierarchical sum-  
120 maries of past experiences to inform failure analysis in robotics. Jain et al. (2022) use SVM bound-  
121 aries to identify captions in a dataset that “summarize” failure modes of ResNets. Tong et al. (2023)  
122 scrape datasets for “erroneous agreement” among outputs of generative multimodal frameworks and  
123 generate natural language descriptions of them, Sagar et al. (2024) use scraped errors alongside hu-  
124 man feedback to adjust models, and Cornelio & Diab (2024) introduce an online, neuro-symbolic  
125 failure identification and recovery framework. Weir et al. (2024) build on the rubric-adjacent argu-  
126 ment analysis work of Jansen et al. (2021) to assess the logical consistency of reasoning arguments  
127 via rubric-equipped LLMs, and Hashemi et al. (2024) use rubrics to improve human alignment of  
128 freeform generated text evaluations. Notably, ? introduce a method for improving self-correction of  
129 complex reasoning by self-asking verification questions based on identified key conditions. Our ap-  
130 proach differs from these prior investigations in that it automatically collects and organizes granular  
131 error classes for *reasoning tasks* while targeting the complete and detailed coverage of error classes  
132 in an individual domain.

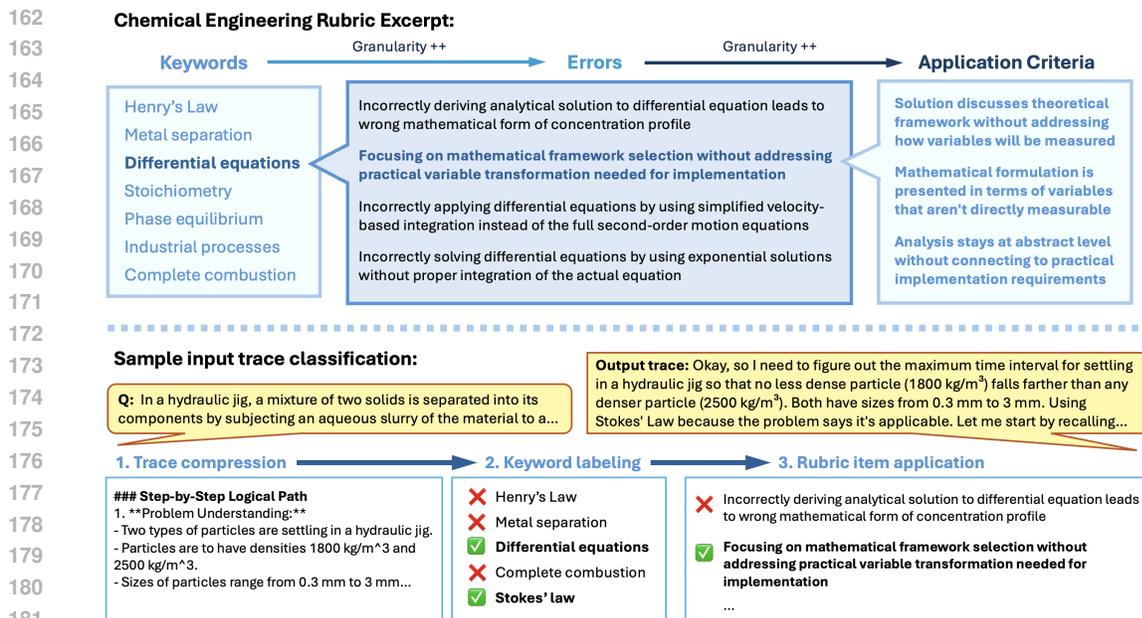
## 133 2.3 REWARD MODELING FOR REASONING MODELS

134  
135 Recent reward function approaches in reinforcement learning for LLMs and reasoning models in-  
136 clude a wide range of methods, spanning from rule-based metrics such as RLVR to learned reward  
137 models that can evaluate both process and outcomes of multi-step reasoning (Zhong et al., 2025).  
138 Rule-based rewards were critical to the development of “early” reasoning models (Guo et al., 2025),  
139 checking for accuracy as well as syntax. However, as rule-based rewards only work for verifiable do-  
140 mains, and may produce false negatives, models trained to assess the correctness of model responses  
141 have been adopted for more complex domains (Li et al., 2023; Liu et al., 2025).

142 Lightman et al. (2023a)’s Verify Step by Step paper demonstrates the benefits of using process su-  
143 pervision, or providing feedback at each reasoning step, for improving LLM alignment on domains  
144 like math. Wang et al. (2023) extend this to function without human annotations by using automati-  
145 cally constructed step labels. Luo et al. (2023), Sun et al. (2023), and Yang et al. (2024) propose  
146 methods for using LLMs-as-judges for reward modeling during training, and RewardBench was in-  
147 troduced by Lambert et al. (2024) as a benchmark for reward models themselves. Sun et al. (2024)  
148 use human-written constitutions as reward models for tasks with non-verifiable rewards, and Cui  
149 et al. (2024) produce an LLM feedback dataset, promoting the importance of “scale and diversity”.  
150 Recently, researchers have begun to adopt rubric-based rewards, but this is so far predominantly  
151 constrained to human-generated rubrics (Jia et al., 2025).

## 152 3 FORMULATING RUBRIC CONSTRUCTION

153  
154 Our goal in this paper is to develop a system that can autonomously extract reasoning errors from  
155 natural language reasoning traces produced by a reasoning model like DeepSeek-R1 (Guo et al.,  
156 2025) or Qwen (Yang et al., 2025) over a specific domain and organize them into a rubric. A rubric  
157 acts as a checklist where each item represents some catastrophic error, meaning that if a checklist  
158 item applies to a reasoning trace, the error in question is highly likely to cause the final answer  
159 to be incorrect. We center catastrophic errors to focus on, and hopefully mitigate, the measurable  
160 downstream performance impact of reasoning errors. An LLM judge, equipped with the rubric  
161 artifact for improved trace classification, takes a trace and runs through the rubric like a checklist,  
classifying traces with any rubric items “checked” as incorrect. This can then be installed as a



182 Figure 2: Top: An excerpt from a rubric constructed from chemical engineering problems in the NaturalReasoning dataset, illustrating their hierarchical organization. Bottom: An sample classification made by Claude 3.5 Sonnet using the rubric (showing trace compression, keyword labeling, and rubric item application). The depicted keyword set and rubric items are subsets of the full sets.

188 reward function that applies this checklist on traces during training for model reinforcement learning, scoring a trace as “correct” only if no rubric items apply to it. We formulate the task of rubric construction below.

192 **Input** Our inputs are training trace set  $S = \{(q_i, o_i, t_i, y_i)\}_{i=1}^N$  and validation trace set  $S' = \{(q'_i, t'_i, y'_i)\}_{i=1}^{N'}$  s.t.  $q \in Q \sim D$  and  $q' \in Q' \sim D$ , where:

- 195 •  $q \in Q$  are natural language strings that describe the parameters of a reasoning problem and serve as input to a reasoning model  $f$ .
- 196
- 197 • Distribution  $D$  is some domain of reasoning problems that share some unknown error taxonomy.
- 198
- 199 •  $o \in O$  are correct ground truth solutions to the reasoning problems.
- 200
- 201 •  $t \in T$  are reasoning traces produced by the reasoning model. These are natural language strings that are generated by the model leading up to (and including) its final answer  $\hat{o}_i = f(q_i, t_i)$ .
- 202
- 203 •  $y_i \in \{0, 1\}$  are binary correctness labels following some scoring function  $y_i = g(o_i, \hat{o}_i)$ . For example,  $g$  could be an LLM providing with a scoring prompt.
- 204
- 205

207 **Output** Our desired output is a failure taxonomy  $\Phi$ , which parameterizes the predefined trace classifier  $h_\Phi(t_i) = \hat{y}_i$  (distinct from our evaluation function  $g(o_i, \hat{o}_i) = y_i$ ).  $\Phi$  is a set of rubric items  $\phi$ , or natural language strings describing some behavior that can be observed in a hypothetical trace.

212 **Hierarchical organization** As the resulting rubrics can be large, our error items are each labeled with a general keyword that helps an LLM map the correct rubric items to the relevant traces. Then, at inference time, the LLM judge first labels each trace with relevant keywords using the full list of keywords from the taxonomy, and then these labels can be used to filter which rubric items are compared against the traces in a second forward pass. Hierarchy details are elaborated in §4.

## 216 4 METHOD

217  
218 In this section, we detail how rubrics are generated. The system can be broken down into two  
219 core components: Trace compression (§4.1) and rubric item extraction (§4.2). We describe each  
220 component in detail alongside its motivation below. We provide an illustration of a generated rubric  
221 and its application in Figure 2.

### 222 4.1 TRACE COMPRESSION

223  
224  
225 There is a distinction that needs to be made between the logical path to a solution, and the explo-  
226 ration process a model may undergo to arrive at that final solution. While in theory, the “optimal”  
227 reasoning trace is simply the former, in many models the exploration process is a core component  
228 enabling it to eventually converge to the correct answer. In some domains, unguided exploration  
229 is explicitly necessary to reach a well-formed solution. In open-ended philosophy problems, for  
230 example, multiple angles of an idea must be considered to arrive at a satisfactory answer.

231 The errors that we aim to identify are those that directly and negatively impact the final solution,  
232 and so exploration of incorrect approaches that are not incorporated into the final solution should  
233 not be included in a rubric. However, this material makes up a large portion of traces outputted  
234 by some models, especially those trained primarily via self-supervised methods. Therefore, we  
235 first “compress” traces into summaries that outline all of the logical steps taken by the model that  
236 influence the eventual final answer. We accomplish this through an LLM call. The downstream  
237 impact of this compression is explored more in §C.

### 238 4.2 EXTRACTING RUBRIC ITEMS

239  
240 As the intended downstream use case of these generated rubrics is to serve as a resource for LLMs  
241 classifying trace correctness, the rubric items must be written such that they can easily be applied  
242 to new traces by LLMs. While rubric items must be detailed enough to be applied consistently by a  
243 judge model, they must be concise enough to not result in an overly long rubric, which increases the  
244 likelihood of application error unnecessarily. To enable rubrics to be long without impacting classifi-  
245 cation performance or compute, we aim to classify traces via two distinct classification stages: (1) A  
246 high-recall, low-granularity classification stage, followed by (2) a high-precision, high-granularity  
247 classification stage. This classification approach requires rubric items to comprise three distinct  
248 fields:

- 249 1. **The main error description:** A simple, straightforward explanation (less than 25 words) de-  
250 scribing the error, and possibly a very brief overview of why it can occur.
- 251 2. **A keyword:** A word or short phrase that can be used to identify traces that might be susceptible  
252 to the error. This should be low granularity and optimized for recall: If this keyword is not  
253 relevant to a trace, then the corresponding error cannot be applicable.
- 254 3. **Verification details:** One or more descriptive explanations for how a model could tell if the  
255 error exists in that argument. If one of these descriptions matches up with the contents of a new  
256 reasoning argument, then that means the error has occurred in the trace.

257  
258 The classification approach using rubrics is illustrated in the bottom half of Figure 2. First, the com-  
259 pressed trace is compared against the full set of keywords, and appropriate keywords are tagged.  
260 Then, all rubric items associated with those keywords are presented to a model alongside the com-  
261 pressed trace, and the classifier determines which ones (if any) apply to the trace.<sup>1</sup> If a trace is  
262 tagged with any rubric items, our classifier labels it as incorrect.

263 We extract these rubric items by first passing the incorrect (compressed) traces from the training set  
264 to an LLM alongside the problem they attempted to solve and its correct solution, when available.  
265 We ask the model to identify potential issues with the compressed trace that would cause it to  
266 produce the incorrect answer. We provide an example rubric item to guide the model’s output.  
267 To further reduce rubric size, we then prompt an LLM to group related keywords, which typically  
268

269 <sup>1</sup>When classifying unseen traces, the classifier is also given the trace the error item was extracted from as  
an example of an applicable reasoning pattern.

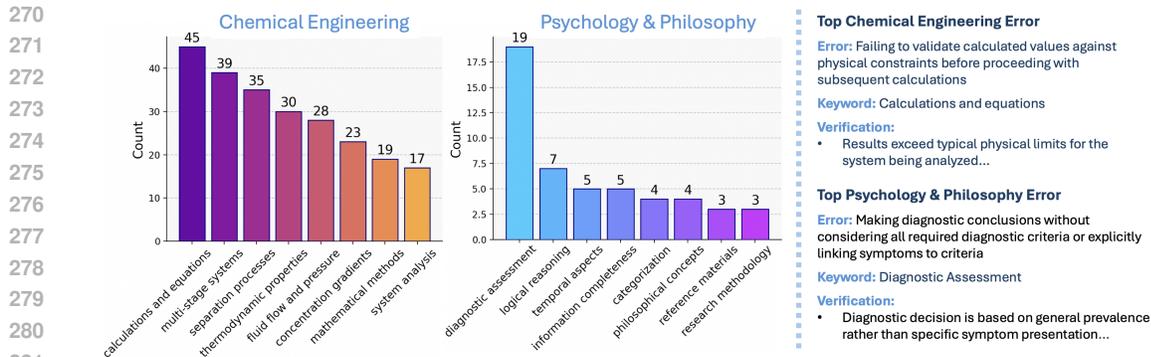


Figure 3: Distribution of most applied rubric item keywords over a technical domain and non-technical (and qualitative) domain, alongside the top most-applied rubric item for each trace set. Both domain problems taken from the Meta NaturalReasoning dataset.

lowers the number of keywords overall by approx. 50%. A comparison of common keywords for technical and non-technical domains are shown in Figure 3.

## 5 EXPERIMENTS

We aim to answer three main questions: **(R1)** Can rubric artifacts improve the specificity and overall accuracy of an LLM trace correctness classifier? **(R2)** Can rubric-augmented classifiers make stronger reward functions than a traditional LLM-as-a-judge reward function? **(R3)** How do rubric-augmented reward functions compare against verifiable reward functions (e.g. string matching) on downstream task accuracy?

We address these questions across two primary experiments: First, we assess how well rubric-augmented LLMs can classify reasoning trace correctness compared to standard LLMs **(R1)**, and then we assess how well rubric-augmented LLMs can serve as reward signals during RL training compared to standard LLMs and verifiable reward setups **(R2 and R3)**. Prompts and additional experimental details for these tasks are included in §B, §D, and §E.

For experiments assessing rubric quality as an artifact for trace classification, we consider the following metrics: **Specificity**, or  $\frac{TN}{FP+TN}$  (the percentage of incorrect traces classified as incorrect by the LLM), **balanced accuracy**, or  $\frac{TN}{2(FP+TN)} + \frac{TP}{2(TN+TP)}$  (the average of the specificity and recall), and **F0.5**,  $(1 + \beta^2) \frac{PR}{\beta^2 P + R}$  where  $\beta = 0.5$ ,  $P = \frac{TP}{TP+FP}$ , and  $R$  is recall. This is the balanced precision and recall with additional weight on precision.

### 5.1 DATASETS

Our primary data constraints are that the domains should be complex enough that they elicit sufficiently long traces, and that the answers to the problems should be verifiable in a relatively consistent and meaningful way. We explore classification behavior on more qualitative data in §B.5, but for main experiments, we select the following datasets:

**SWE-Bench** (Jimenez et al., 2024) tests our approach in the domain of coding and software development. SWE-Bench is a collection of issue-pull request pairs from GitHub repositories that evaluates Python code patches via a set of unit tests. **NuminaMath** (Li et al., 2024) tests complex mathematical reasoning. NuminaMath is a collection of high school/competition-level math problems paired with chain-of-thought style ground truth solutions collected via PDF scanning. We filter the dataset for problems labeled as calculus, geometry, and number theory. Finally, we extract a subset of chemical engineering

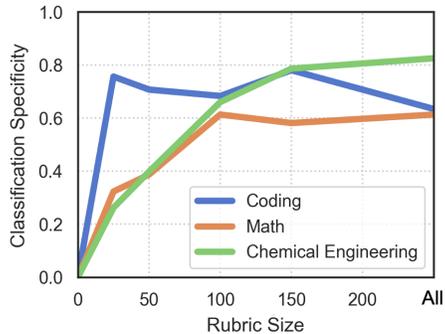


Figure 4: Experiment 1 ablation: # of rubric items vs. classification specificity.

324 problems from Meta’s **NaturalReasoning** (Yuan et al., 2025) dataset for more open-ended problem-  
 325 solution pairs. NaturalReasoning is a broad collection of problems from various disciplines ranging  
 326 from social sciences, coding, math, physical sciences, and humanities.

## 328 5.2 EXPERIMENT 1: CLASSIFYING REASONING ERRORS

329 We evaluate the performance of a rubric-augmented LLM on identifying traces that lead to incorrect  
 330 answers. This experiment targets **R1**, or whether rubric artifacts can improve the *specificity and*  
 331 *accuracy* of an LLM-based trace correctness classifier.

332  
 333 **Data** We collect long-form reasoning traces for each dataset listed in the previous section. SWE-  
 334 Bench traces are generated with a tool-calling variant of Claude 3.5 Haiku <sup>2</sup>, and NuminaMath and  
 335 NaturalReasoning traces are generated with DeepSeek-R1 (Guo et al., 2025), sourced from existing  
 336 HuggingFace resource datasets <sup>3 4</sup>. For NuminaMath and NaturalReasoning, R1-generated final  
 337 answers are scored as correct or incorrect by a Claude 3.5 Sonnet model that is provided access to  
 338 the question, R1 answer, and ground truth solution. This prompt is designed to frame the scoring  
 339 task as a homework grading scenario. For each domain, we sample 450 to 800 problems, depending  
 340 on data availability, and split 80/20 between train and validation.

341  
 342 **Methods** For each domain, a rubric is constructed using the incorrect traces from the training  
 343 question-trace pairs from each dataset with Claude 3.5 Sonnet v2, resulting in rubrics of approxi-  
 344 mately 250 items each. Keywords are illustrated in Figure 3. Keywords are not clustered for Nu-  
 345 minaMath due to the narrow range of problem types, keyword clustering influence on performance  
 346 is explored in §C. We compare a rubric-augmented LLM classifier against a standard LLM (using a  
 347 prompt with no rubric). For LLM calls, we use Claude 3.5 Sonnet v2<sup>5</sup>.

348 The baseline LLM classifier judges traces using a single prompt<sup>6</sup>, while the rubric-augmented clas-  
 349 sifier use a two-stage classification pipeline. In the first pass, the classifier identifies which keywords  
 350 generated alongside rubric items map to each validation trace, and then in the second pass, all rubric  
 351 items associated with the mapped keywords are compiled into a mini-rubric that is compared against  
 352 the reasoning trace. If any of the rubric items are tagged by the classifier as applying to the trace,  
 353 the trace is marked as incorrect. If no items are applied, the trace is marked as correct.

354  
 355 **Results** Experimental results are included in Figure 5. The results illustrate that the rubric aug-  
 356 mentation can be highly effective at improving specificity, or recognizing when an incorrect trace is  
 357 incorrect. The rubric-augmented classifier also consistently results in higher overall balanced accu-  
 358 racy when compared against the baseline classifier, although this margin is generally smaller. These  
 359 results clearly illustrate the strong tradeoff between specificity and recall, or the ability to correctly  
 360 classify correct traces. Errors in reasoning traces do not always negatively impact the downstream fi-  
 361 nal answer, as the relationships between individual model steps are often tenuous and complex Levy  
 362 et al. (2025). We believe that jointly optimizing specificity and recall further in trace classification  
 363 is an exciting direction for future research.

## 364 5.3 EXPERIMENT 2: RUBRIC REWARD FUNCTIONS

365  
 366 **Overview** Current approaches to reward modeling for reinforcement learning with language mod-  
 367 els often rely on simple evaluation metrics that focus primarily on the correctness of final answers  
 368 against ground truth solutions, requiring a large set of verifiable ground truth solutions or compute-  
 369 intensive answer verification systems. By enhancing the capacity of LLM judges to reliably evaluate  
 370 reasoning traces, we can develop more effective reward functions that optimize not just for correct  
 371 answers but for sound reasoning processes using a small set of ground truth solutions. In this sec-  
 372 tion, we explore how our rubric-augmented trace classifiers can be used to improve RL training in

373 <sup>2</sup>Link to SWE-Bench submission link

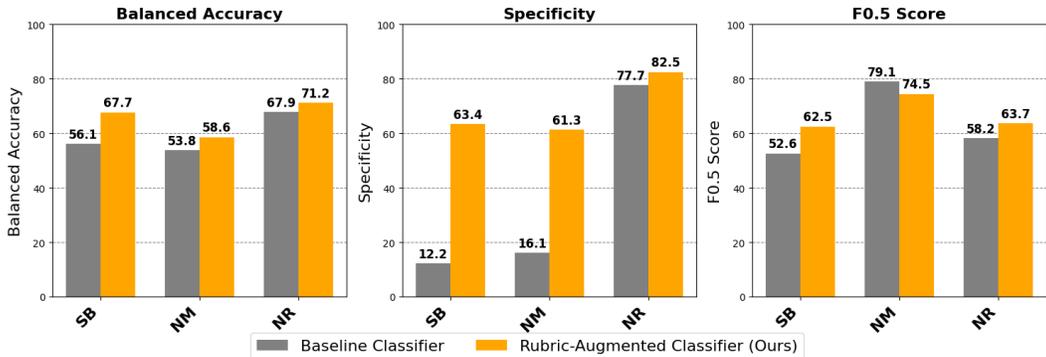
374 <sup>3</sup>Link to NuminaMath source

375 <sup>4</sup>NaturalReasoning source

376 <sup>5</sup>Link to Claude 3.5 model card

377 <sup>6</sup>Some SWE-Bench traces were longer than the model’s context size. The prefixes of these traces were cut  
 s.t. they could fit within the model’s context window.

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389



390  
391  
392  
393  
394  
395

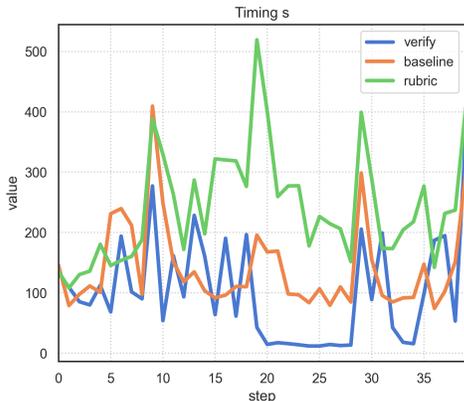
Figure 5: Rubric-augmented Claude 3.5 Sonnet outperforms standard prompted Sonnet at identifying incorrect reasoning traces across three diverse domains: SWE-Bench (SB), NuminaMath (NM), and chemical engineering problems from NaturalReasoning (NR). We present balanced accuracy (mean accuracy with correct and incorrect traces weighted equally), specificity (# of incorrect traces classified as incorrect), and F0.5 (measuring precision and recall with additional precision weight).

396  
397  
398  
399

complex domains where ground truth labels are scarce. This experiment targets **R2** and **R3**: We assess whether a rubric-augmented classifier’s ability to act as a reward function for RL training can improve downstream task accuracy over a standard LLM or verifiable rewards.

400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417

**Data** We perform this experiment in the math and coding domains, as they enable explicit verifiable rewards for a completely grounded comparison (independent of LLMs) against the new LLM-based reward approach. For each of these datasets, we first fine-tune a Qwen3-4B model (thinking mode) (Yang et al., 2025) using SFT with a subset of 1.6K samples of question-ground truth solution pairs, using a cross-entropy loss and a learning rate of 1e-6. We fine-tune for only one epoch as convergence occurs quickly for these domains due to limited drift from the original pretraining data. For NuminaMath, we partition the remaining dataset on problem type and use it to randomly sample a second subset of 1.4K training problems and 340 validation problems to serve as our (disjoint) reinforcement learning dataset. For SWE-Bench, we sample 1.4K problems from the train dataset for training and 80 problems from the test set for validation.<sup>7</sup>



418  
419  
420  
421  
422  
423

**Training Setup** We train Qwen3-4B using the DAPO RL algorithm (Yu et al., 2025). For DAPO, the model was configured with a maximum response length of 5120 tokens for math and 10000 tokens for coding. For our advantage estimation, we employed GRPO. Models are trained on 8 40GB Nvidia A100s. For our validation reward on NuminaMath we used MathVerify package<sup>8</sup>.

424  
425  
426  
427  
428

**Methods** As comparison models, we include Claude 3.7 Sonnet, Claude 3.5 Haiku, DeepSeek-R1, and Mistral-7B<sup>9</sup>. We again cut off math responses at 5k tokens and coding responses at 10k. For each domain, we use the fine-tuned Qwen3-4B model (post-SFT, pre-DAPO) to generate traces for a rubric generation dataset. For math we generate traces using a sample of our RL training problems

429  
430  
431

<sup>7</sup>Data limited to shorter inputs (less than 25000 characters for RL and 35000 characters for rubric construction for memory and compute efficiency).  
<sup>8</sup><https://github.com/huggingface/Math-Verify>  
<sup>9</sup>Link to Mistral model card. Mistral traces are cut off at 5k for both domains.

	Baseline Models				RL w/ GT	RL w/o GT		
Benchmark Model Size <sup>10</sup>	Sonnet ~100B+	R1 37B	Haiku ~20B+	Mistral 7B	VR 4B	No RL 4B	BL 4B	Rubric 4B (Ours)
NM (Acc)	21.7	15.5	16.3	3.7	16.1	5.5	10.5	<b>15.2</b>
SB (Acc)	28.8	8.8	1.3	0.0	n/a	0.0	0.0	<b>2.5</b>
SB (Comp)	77.5	30.0	2.5	15.0	n/a	0.0	2.5	<b>20.0</b>

Table 1: Validation accuracy and patch completion results post-RL training demonstrate that the rubric-augmented LLM-based reward function produces models with higher output quality than the baseline LLM reward function (using a prompt with no error rubric) on both domains, as well as the standard verifiable rewards approach on the math domain. The rubric-augmented reward function performs best when paired with a penalization term that discourages overly short reasoning traces.

set and for SWE-Bench use another subset of the original dataset’s test data (as these problems are verifiable, unlike the train set) and verify them with ground truth answers to produce a set of 200 verifiably incorrect traces. Rubrics are generated with these incorrect trace sets using Claude 3.5 Sonnet v2, following the setup detailed in the previous experiment. Due to the brevity of the produced traces post-fine-tuning, we omit the trace compression step before extracting rubric items. We compare three different training reward functions:

**Gold Match** is identical to the validation reward function. Model responses are taken as-is and compared against a ground-truth answer using MathVerify. SWE-Bench version is not included due to high resource costs of running the benchmark unit tests in an RL setting. **Baseline Classifier** replicates the baseline classifier used in Experiment 1 using Amazon Nova Lite as the LLM judge.<sup>11</sup> Provided with the question and model response, the LLM assesses whether or not the answer is likely to be correct or not, outputting binary labels. **Rubric-Augmented Classifier** replicates the rubric classifier used in the classification experiments (§5.2), also using Amazon Nova Lite with a rubric constructed by Claude 3.5 Sonnet, detailed above and also using binary [0, 1] integer labels depending on the trace correctness classification. As with the rubric construction, we omit the trace compression step at inference time.

**Results** Results are included in Table 1. Timing information is shown in Figure 6. For both datasets we report validation accuracy, and for SWE-Bench we also report completed patch rate (the number of not-empty patches that do not cause an execution error) due to task difficulty and low performance on final patch correctness of most of the tested models. We report the highest performance achieved during training (computed every ten steps over the full validation set for each domain). The higher performance of the rubric-augmented LLM reward on both domains indicates that it is a strong alternative to traditional training methods while requiring significantly fewer ground truth labels: The rubric method approaches the verifiable rewards setup for NuminaMath and improves patch completion percentage by a full 17.5% over the baseline for SWE-Bench,

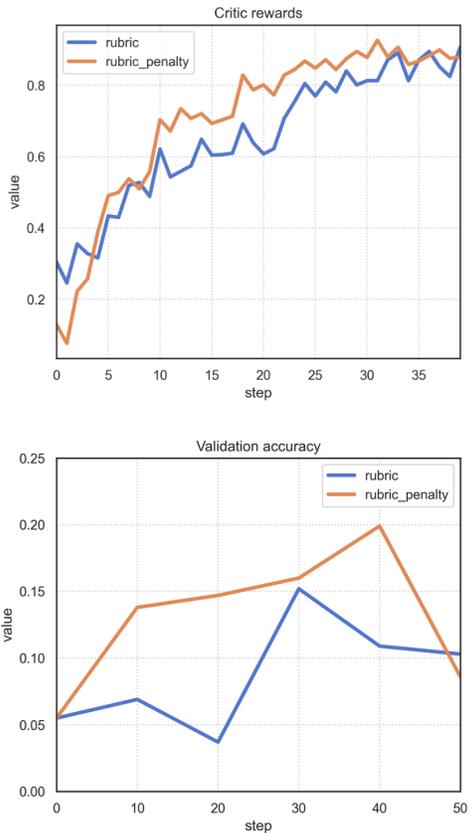


Figure 7: Train/validation for rubric rewards with and without a trace length penalty.

<sup>10</sup>Claude model parameter counts are reported as a rough, estimated lower bound, sourced here. We report the number of activated paramers at runtime, so report 37B for R1 instead of 671B.

<sup>11</sup>Link to Amazon Nova model.

notably outperforming the larger Claude 3.5 Haiku model on both metrics. We hypothesize that further investigation into leveraging rubrics as an artifact to improve reward functions can yield greater performance improvements, such as incorporating other qualitative feedback signals, indicated by the improvement achieved by adding a trace length penalty in Figure 7.

## 6 CONCLUSION

This work explore methods to automatically create domain-specific reasoning error taxonomies, significantly enhancing the capability of LLMs to identify their own errors and addressing a fundamental challenge in employing LLM-as-judge frameworks. This methodology reduces the need for extensive manual curation of gold labels while accounting for the unique reasoning patterns and potential pitfalls characteristic of each domain.

Our work enables efforts in training reasoning models in domains where ground truth data is ambiguous or expensive to obtain. By providing a structured framework for error identification requiring only a few gold labels, we address a critical bottleneck in the development of LLMs for specialized knowledge domains. The potential for these rubric-enhanced judges to support human-AI collaborative reasoning is a line of exciting future work enabled by this research.

## REFERENCES

- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- Cristina Cornelio and Mohammed Diab. Recover: A neuro-symbolic framework for failure detection and recovery, 2024. URL <https://arxiv.org/abs/2404.00756>.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie, Ruobing Xie, Yankai Lin, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with scaled ai feedback, 2024. URL <https://arxiv.org/abs/2310.01377>.
- Lisa Dunlap, Krishna Mandal, Trevor Darrell, Jacob Steinhardt, and Joseph E Gonzalez. Vibecheck: Discover and quantify qualitative differences in large language models, 2025. URL <https://arxiv.org/abs/2410.12851>.
- Arduin Findeis, Timo Kaufmann, Eyke Hüllermeier, Samuel Albanie, and Robert Mullins. Inverse constitutional ai: Compressing preferences into principles, 2025. URL <https://arxiv.org/abs/2406.06560>.
- Jan-Philipp Fränken, Eric Zelikman, Rafael Rafailov, Kanishk Gandhi, Tobias Gerstenberg, and Noah D. Goodman. Self-supervised alignment with mutual information: Learning to follow principles without preference labels, 2024. URL <https://arxiv.org/abs/2404.14313>.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Saizhuo Wang, Kun Zhang, Yuanzhuo Wang, Wen Gao, Lionel Ni, and Jian Guo. A survey on llm-as-a-judge, 2025. URL <https://arxiv.org/abs/2411.15594>.
- Melody Y Guan, Manas Joglekar, Eric Wallace, Saachi Jain, Boaz Barak, Alec Helyar, Rachel Dias, Andrea Vallone, Hongyu Ren, Jason Wei, et al. Deliberative alignment: Reasoning enables safer language models. *arXiv preprint arXiv:2412.16339*, 2024.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Helia Hashemi, Jason Eisner, Corby Rosset, Benjamin Van Durme, and Chris Kedzie. Llm-rubric: A multidimensional, calibrated approach to automated evaluation of natural language texts. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 13806–13834. Association for Computational Linguistics,

- 540 2024. doi: 10.18653/v1/2024.acl-long.745. URL [http://dx.doi.org/10.18653/v1/](http://dx.doi.org/10.18653/v1/2024.acl-long.745)  
541 [2024.acl-long.745](http://dx.doi.org/10.18653/v1/2024.acl-long.745).  
542
- 543 Arian Hosseini, Xingdi Yuan, Nikolay Malkin, Aaron Courville, Alessandro Sordoni, and Rishabh  
544 Agarwal. V-STAR: Training verifiers for self-taught reasoners. In *First Conference on Language*  
545 *Modeling*, 2024. URL <https://openreview.net/forum?id=stmqBSW2dV>.  
546
- 547 Saachi Jain, Hannah Lawrence, Ankur Moitra, and Aleksander Madry. Distilling model failures as  
548 directions in latent space, 2022. URL <https://arxiv.org/abs/2206.14754>.  
549
- 550 Peter Alexander Jansen, Kelly Smith, Dan Moreno, and Huitzililn Ortiz. On the challenges of evalu-  
551 ating compositional explanations in multi-hop inference: Relevance, completeness, and expert  
552 ratings. *ArXiv*, abs/2109.03334, 2021. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:237439283)  
553 [CorpusID:237439283](https://api.semanticscholar.org/CorpusID:237439283).  
554
- 555 Ruipeng Jia, Yunyi Yang, Yongbo Gai, Kai Luo, Shihao Huang, Jianhe Lin, Xiaoxi Jiang, and  
556 Guanjin Jiang. Writing-zero: Bridge the gap between non-verifiable tasks and verifiable rewards,  
557 2025. URL <https://arxiv.org/abs/2506.00103>.  
558
- 559 Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik  
560 Narasimhan. Swe-bench: Can language models resolve real-world github issues?, 2024. URL  
561 <https://arxiv.org/abs/2310.06770>.  
562
- 563 Sandipan Kundu, Yuntao Bai, Saurav Kadavath, Amanda Askell, Andrew Callahan, Anna Chen,  
564 Anna Goldie, Avital Balwit, Azalia Mirhoseini, Brayden McLean, Catherine Olsson, Cassie  
565 Evraets, Eli Tran-Johnson, Esin Durmus, Ethan Perez, Jackson Kernion, Jamie Kerr, Kamal  
566 Ndousse, Karina Nguyen, Nelson Elhage, Newton Cheng, Nicholas Schiefer, Nova Das-  
567 Sarma, Oliver Rausch, Robin Larson, Shannon Yang, Shauna Kravec, Timothy Telleen-Lawton,  
568 Thomas I. Liao, Tom Henighan, Tristan Hume, Zac Hatfield-Dodds, Sören Mindermann, Nicholas  
569 Joseph, Sam McCandlish, and Jared Kaplan. Specific versus general principles for constitutional  
570 ai, 2023. URL <https://arxiv.org/abs/2310.13798>.  
571
- 572 Yara Kyrzhenko, Ke Zhou, Edyta Bogucka, and Daniele Quercia. C3ai: Crafting and evaluat-  
573 ing constitutions for constitutional ai. In *Proceedings of the ACM on Web Conference 2025*,  
574 WWW '25, pp. 3204–3218, New York, NY, USA, 2025. Association for Computing Machin-  
575 ery. ISBN 9798400712746. doi: 10.1145/3696410.3714705. URL [https://doi.org/10.](https://doi.org/10.1145/3696410.3714705)  
576 [1145/3696410.3714705](https://doi.org/10.1145/3696410.3714705).  
577
- 578 Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi  
579 Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh  
580 Hajishirzi. Rewardbench: Evaluating reward models for language modeling, 2024. URL  
581 <https://arxiv.org/abs/2403.13787>.  
582
- 583 Mosh Levy, Zohar Elyoseph, and Yoav Goldberg. Humans perceive wrong narratives from ai rea-  
584 soning texts, 2025. URL <https://arxiv.org/abs/2508.16599>.  
585
- 586 Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Huang, Kashif  
587 Rasul, Longhui Yu, Albert Q Jiang, Ziju Shen, et al. Numinamath: The largest public dataset in  
588 ai4maths with 860k pairs of competition math problems and solutions. *Hugging Face repository*,  
589 13(9):9, 2024.  
590
- 591 Junlong Li, Shichao Sun, Weizhe Yuan, Run-Ze Fan, Hai Zhao, and Pengfei Liu. Generative judge  
592 for evaluating alignment, 2023. URL <https://arxiv.org/abs/2310.05470>.  
593
- 594 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan  
595 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step, 2023a. URL  
596 <https://arxiv.org/abs/2305.20050>.  
597
- 598 Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan  
599 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth*  
600 *International Conference on Learning Representations*, 2023b.

- 594 Shudong Liu, Hongwei Liu, Junnan Liu, Linchen Xiao, Songyang Gao, Chengqi Lyu, Yuzhe Gu,  
595 Wenwei Zhang, Derek F. Wong, Songyang Zhang, and Kai Chen. Compassverifier: A unified and  
596 robust verifier for llms evaluation and outcome reward, 2025. URL [https://arxiv.org/  
597 abs/2508.03686](https://arxiv.org/abs/2508.03686).
- 598 Zeyi Liu, Arpit Bahety, and Shuran Song. Reflect: Summarizing robot experiences for failure  
599 explanation and correction, 2023. URL <https://arxiv.org/abs/2306.15724>.
- 601 Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qing-  
602 wei Lin, Shifeng Chen, and Dongmei Zhang. Wizardmath: Empowering mathematical reasoning  
603 for large language models via reinforced evol-instruct. *arXiv preprint arXiv:2308.09583*, 2023.
- 604 Tong Mu, Alec Helyar, Johannes Heidecke, Joshua Achiam, Andrea Vallone, Ian Kivlichan, Molly  
605 Lin, Alex Beutel, John Schulman, and Lilian Weng. Rule based rewards for language model  
606 safety. *Advances in Neural Information Processing Systems*, 37:108877–108901, 2024.
- 607 Som Sagar, Aditya Taparia, and Ransalu Senanayake. Failures are fated, but can be faded: Char-  
608 acterizing and mitigating unwanted behaviors in large-scale vision and language models, 2024.  
609 URL <https://arxiv.org/abs/2406.07145>.
- 611 Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan,  
612 Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with  
613 factually augmented rlhf. *arXiv preprint arXiv:2309.14525*, 2023.
- 614 Zhiqing Sun, Yikang Shen, Hongxin Zhang, Qinhong Zhou, Zhenfang Chen, David Cox, Yiming  
615 Yang, and Chuang Gan. Salmon: Self-alignment with instructable reward models, 2024. URL  
616 <https://arxiv.org/abs/2310.05910>.
- 617 Shengbang Tong, Erik Jones, and Jacob Steinhardt. Mass-producing failures of multimodal systems  
618 with language models. *Advances in Neural Information Processing Systems*, 36:29292–29322,  
619 2023.
- 620 Peiyi Wang, Lei Li, Zhihong Shao, RX Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang  
621 Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. *arXiv  
622 preprint arXiv:2312.08935*, 2023.
- 623 Nathaniel Weir, Kate Sanders, Orion Weller, Shreya Sharma, Dongwei Jiang, Zhengping Jiang,  
624 Bhavana Dalvi Mishra, Oyvind Tafjord, Peter Jansen, Peter Clark, and Benjamin Van Durme.  
625 Enhancing systematic compositional natural language inference using informal logic, 2024.  
626 URL <https://arxiv.org/abs/2402.14798>.
- 627 An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jian-  
628 hong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2. 5-math technical report: Toward mathematical  
629 expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024.
- 630 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang  
631 Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,  
632 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin  
633 Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,  
634 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui  
635 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang  
636 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger  
637 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan  
638 Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- 639 Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian  
640 Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system  
641 at scale. *arXiv preprint arXiv:2503.14476*, 2025.
- 642 Weizhe Yuan, Jane Yu, Song Jiang, Karthik Padthe, Yang Li, Ilia Kulikov, Kyunghyun Cho, Dong  
643 Wang, Yuandong Tian, Jason E Weston, et al. Naturalreasoning: Reasoning in the wild with 2.8  
644 m challenging questions. *arXiv preprint arXiv:2502.13124*, 2025.

648 Jialun Zhong, Wei Shen, Yanzeng Li, Songyang Gao, Hua Lu, Yicheng Chen, Yang Zhang, Wei  
649 Zhou, Jinjie Gu, and Lei Zou. A comprehensive survey of reward models: Taxonomy, applica-  
650 tions, challenges, and future, 2025. URL <https://arxiv.org/abs/2504.12328>.  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## A ADDITIONAL BASELINE EXPERIMENTS

We evaluated a range of baseline prompts to compare against in our experiments detailed in §5.2 and §5.3. In Table 2, we report our primary classification metrics for a representative set of 6 baseline approaches. These approaches are:

1. **Baseline:** The baseline used in the main experiments, using the prompt listed in §E.11 and §E.14 alongside the full model reasoning trace. The classifier is simply asked whether, given the model output and the original question, the full reasoning trace results in a correct answer or not.
2. **Baseline Alternate 1:** The same setup as Baseline, but only the first 75% of the trace is presented to the classifier (determined by line count). The motivation behind this is that, empirically, the baseline classifier is highly likely to score a trace as correct if it appears to have converged on a final answer, regardless of answer soundness. By omitting the last decisions made by the model in its thought process, the overall error specificity is likely to increase.
3. **Baseline Alternate 2:** Uses the prompt listed in §E.12 and §E.15. Following the logic detailed for Alternate 1, we provide the full trace, but tell the classifier model that it is only a snippet of the full reasoning passage.
4. **Baseline Alternate 3:** The same setup as Alternate 2, but the last 75% of the trace is again cut off as in Alternate 1.
5. **Baseline Alternate 4:** Uses the prompt listed in §E.13 and §E.16. The prompt indicates that only a portion of the reasoning trace is provided, but instead of asking if the trace will result in a correct output, asks whether the model should continue thinking or provide an answer immediately. The full trace is provided.
6. **Baseline Alternate 5:** The same setup as Alternate 4, but the trace is cut off in the manner described in Alternate 3 and Alternate 1.

	SWE-Bench			NuminaMath			NaturalReasoning		
Domain	BA	S	F0.5	BA	S	F0.5	BA	S	F0.5
Rubric	.677	.634	.625	.586	.613	.745	.712	.825	.637
Baseline	<b>.561</b>	.122	<b>.526</b>	.538	.161	<b>.791</b>	.679	.777	.582
BL Alt. 1	.497	.463	.452	.538	.258	.777	.674	<b>.893</b>	<b>.628</b>
BL Alt. 2	<b>.561</b>	.122	<b>.526</b>	.527	.129	.789	<b>.697</b>	.777	.601
BL Alt. 3	.442	.415	.399	.533	.216	.780	.689	.787	.596
BL Alt. 4	.537	.073	.513	<b>.548</b>	.290	.780	.598	.524	.464
BL Alt. 5	.444	<b>.732</b>	.260	.511	<b>.355</b>	.736	.546	.728	.405

Table 2: Alternate baselines compared against the rubric approach and the primary baseline used in §E.12 and §E.15. All experiments are run with Claude 3.5 Sonnet. BA = Balanced Accuracy and S = Specificity.

Generally, the original baseline method achieves a strong balance between balanced accuracy and F0.5 score across the three domains compared to the other methods. As shown in the results, there is a clear tradeoff between the different primary metrics. As hypothesized, only providing the first 75% of the trace increases specificity, but generally negatively impacts balanced accuracy and F0.5. However, none of the baseline approaches outperform the rubric method on more than one of the three metrics for any domain, indicating that the rubric approach is generally superior in terms of error identification for technical content.

Additionally, we re-ran the classification experiment using an open-source model, Qwen3-235b, for both rubric generation and classification. We report these results alongside the Claude 3.5 Sonnet results in Table 3.

Domain	SWE-Bench			NuminaMath			NaturalReasoning		
	BA	S	F0.5	BA	S	F0.5	BA	S	F0.5
Claude Baseline	.561	.122	.526	.538	.161	<b>.791</b>	.679	.777	.582
Claude Rubric	<b>.677</b>	<b>.634</b>	<b>.625</b>	.586	<b>.613</b>	.745	<b>.712</b>	<b>.825</b>	<b>.637</b>
Qwen Baseline	<b>.634</b>	.268	<b>.571</b>	.526	.108	<b>.791</b>	.600	.272	.456
Qwen Rubric	.536	<b>.854</b>	.417	<b>.544</b>	<b>.231</b>	.786	<b>.674</b>	<b>.476</b>	<b>.518</b>

Table 3: Open-source classification results.

## B RUBRIC ANALYSIS

### B.1 KEYWORD VISUALIZATION

Below in Figure 8 we provide histograms of the top 20 most common keywords (in terms of # of rubric items, as opposed to in terms of the number of *applied* rubric items over the validation set of traces as in Figure 3) for the rubrics produced in §5.2.

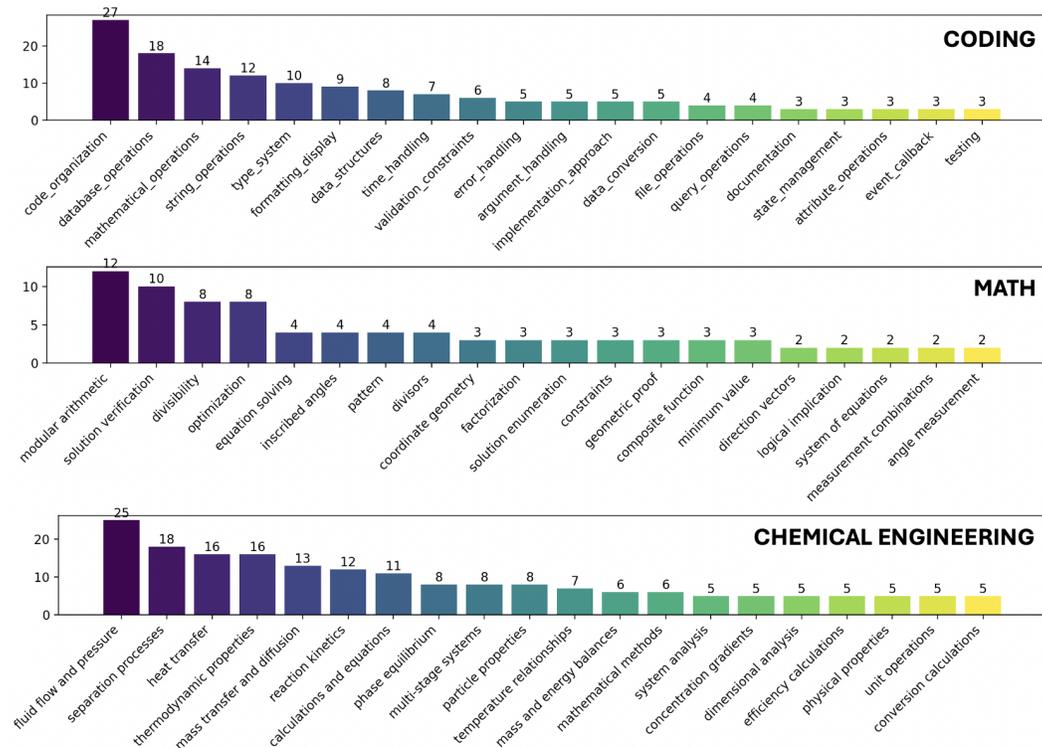


Figure 8: Frequency of top twenty keywords from generated rubrics on target technical domains, in terms of the total number of rubric items.

### B.2 RUBRIC STATISTICS

Below, we report the total number of training and test traces for the datasets described in §5.2, the corresponding rubric size, and the number of keywords (after clustering for chemical engineering and coding).

Domain	Train set	Val. set	Rubric size	# Keywords
Chemical Eng.	636	158	296	90
Math	476	124	250	169
Coding	406	73	234	89

Table 4: Dataset and corresponding rubric statistics for the §5.2 classification experiment setup.

### B.3 CLASSIFICATION CONFUSION MATRICES

In Figure 9, we show the confusion matrices for the rubric-augmented classifier’s outputs on the validation set for the experiment described in §5.2.

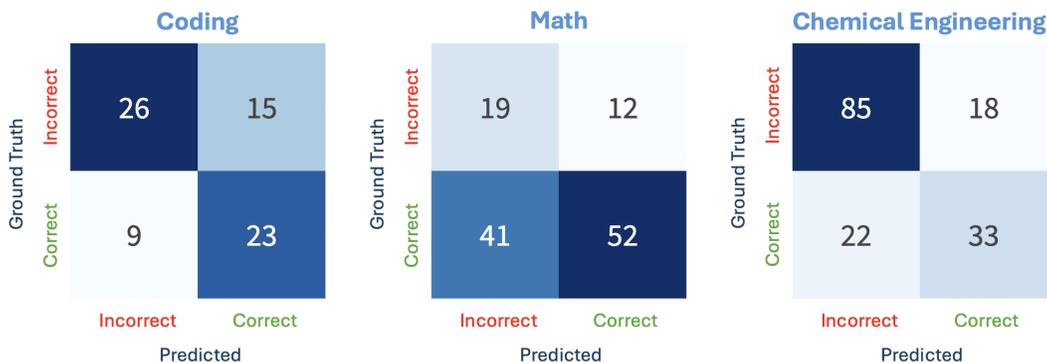


Figure 9: Confusion matrices corresponding to the experimental results illustrated in Figure 5. As shown, false negatives tend to be more common among the rubric approach than false positives, given the total number of positives and negatives per validation set.

### B.4 TRAINING SET CLASSIFICATION ACCURACY

Below in Table 5 we report the primary classification metrics over the training trace set used to construct the rubrics used for classification over the validation set in §5.2.

Domain	SWE-Bench			NuminaMath			NaturalReasoning (Chem)		
	Bal. Acc.	Spec.	F0.5	Bal. Acc.	Spec.	F0.5	Bal. Acc.	Spec.	F0.5
Baseline	.562	.123	.509	.515	.098	.524	.662	.676	.687
Rubric	.669	.438	.586	.566	.431	.545	.703	.662	.722

Table 5: Baseline Classifier vs. Rubric-Augmented Classifier results when evaluated on the training set used to construct the rubrics in §5.2.

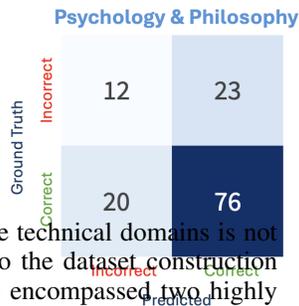
As illustrated above, the results achieved on the training set are not notably better than those achieved over the validation set as reported in Figure 5. This indicates that the approach is highly generalizable and robust to unseen traces, but simultaneously suggests that the classification approach outlined in §4.2 is not necessarily optimal for aligning erroneous traces with appropriate rubric items. All of the errors constructed for the rubric were mined specifically from incorrect traces in the training set, and so the low specificity values of  $\leq .662$  indicate that the classifier is unable to re-identify these errors without the ground truth answers being provided for guidance. This highlights an interesting area for future work that may improve the quality of the rubric approach, and its consequent downstream performance on classification and reward modeling.

### B.5 NON-TECHNICAL RESULTS

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

We select a sample of ~ 600 questions from NaturalReasoning that concern philosophy and psychology questions, set aside a test set of approx. 150 questions, and generate a rubric on the remaining trace set. We evaluate a rubric-based classifier and the baseline approach on the test set, documented in the table below alongside the corresponding confusion matrix in Table B.4.

Our hypothesis as to why the method’s performance is lower than on the technical domains is not due to inherent drift towards more qualitative failure modes, but due to the dataset construction itself: The joint collection of psychology and philosophy problems (1) encompassed two highly distinct failure mode sets, and (2) has no simple way of obtaining ground truth trace correctness scores without expert trace annotation. Both of these factors make the problem inherently more challenging, and point to tackling them as exciting areas for future work.



Method	Balanced Acc	Specificity	F0.5
Baseline	.571	.225	.821
Rubric	.567	.343	.772

## C ABLATION EXPERIMENTS

We consider different alternate rubric construction and classification setups for the three domains considered in §5.2:

- No trace compression for trace classification.
- No keyword clustering.
- A second rubric item classification filter (uses the same prompt as the first rubric item classification, shown in §E.5 and §E.10). This in effect requires the LLM trace classifier to confirm once again that the rubric items applied by the classifier in the previous step are correctly applied, in an attempt to reduce false negatives during classification.

We report the exact numbers below in Table 6. We also consider the impact of different rubric sizes on classification efficacy: 25 items, 50 items, 100 items, 150 items, and the full rubric, the results of which are shown in brief in Figure 4.

Domain	SWE-Bench			NuminaMath			NaturalReasoning		
	BA	S	F0.5	BA	S	F0.5	BA	S	F0.5
<b>Baseline</b>	.561	.122	.526	.538	.161	<b>.791</b>	.679	.777	.582
<b>Rubric</b>	<b>.677</b>	.634	<b>.625</b>	<b>.586</b>	<b>.613</b>	.745	<b>.712</b>	<b>.825</b>	<b>.637</b>
<b>No compression</b>	.491	<b>.951</b>	.114	.527	.323	.757	.601	.602	.470
<b>No clustering</b>	.581	.537	.532	<b>.586</b>	<b>.613</b>	.745	.702	.767	.601
<b>Second filter</b>	.601	.390	.551	.457	.129	.740	.657	.495	.508
Rubric Size	BA	S	F0.5	BA	S	F0.5	BA	S	F0.5
<b>25 items</b>	.628	.756	.588	.468	.323	.704	.567	.262	.436
<b>50 items</b>	.604	.707	.556	.409	.387	.608	.572	.398	.439
<b>100 items</b>	.560	.683	.500	.554	<b>.613</b>	.708	.657	.660	.531
<b>150 items</b>	.609	<b>.780</b>	.565	.495	.581	.640	.657	.786	.560
<b>All items</b>	<b>.677</b>	.634	<b>.625</b>	<b>.586</b>	<b>.613</b>	<b>.745</b>	<b>.712</b>	<b>.825</b>	<b>.637</b>

Table 6: Ablation experiments corresponding to the experiment setup in §5.2.

The performance drop from removing the trace compression at test time indicates that the unavoidable information loss of this step is less critical than the importance of distilling the trace to a smaller size. In implementation, the resulting prompts were sometimes longer than the context window of Claude 3.5 Sonnet, and so the prompts had to be cut off to obtain classifications. For the domains with more failure modes (not the NuminaMath subset), removing the keyword clustering similarly hurt overall performance.

Somewhat surprisingly, the second filter did not improve F0.5 scores over the normal rubric approach, despite the filtering method being designed to improve recall and precision of the classifier. This result may be related to that of §B.4, in that the classifier itself is likely leveraging the rubric artifact sub-optimally which may be introducing noise.

The math and chemical engineering rubric size trends are not particularly surprising: There is a general positive correlation between increased rubric size and metric performance across the board, although the results demonstrate slightly higher balanced accuracy and F0.5 scores for math at n=25 over respectively larger (but still small, overall) rubrics. The coding results are an outlier, showing higher scores for all three metrics for the smallest rubric size over most other rubrics except for n=150 and the full-sized rubric. This indicates that there are a small number of rubric items that may account for a large portion of the failure modes in the training and validation sets for this domain. Exploring the relative importance of different rubric items, as well as determining the optimal rubric size for a given domain, depending on the domain’s breadth and potential for a variety of failure types, is interesting future work. Learning to leverage this information could produce more intelligent and effective classifiers overall.

## D ADDITIONAL RL TRAINING FIGURES

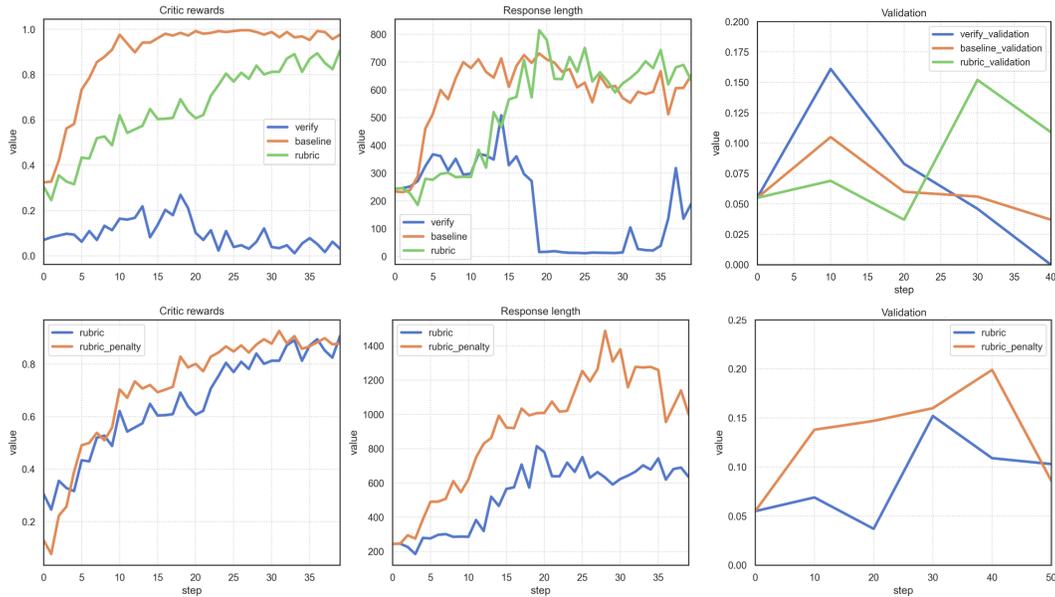


Figure 10: Mean critic reward, mean response length, and mean validation for DAPO training on NuminaMath. Top row shows the training run using the verifiable rewards reward function (Math-Verify), the baseline classifier, and the rubric-augmented classifier, and the bottom row shows the rubric-augmented classifier with and without a length penalty applied to the model output.

As illustrated above in Figure 10, the model is able to learn the LLM-as-judge reward functions much more easily than the verifiable rewards feedback signal, despite selecting the training hyperparameters based on performance using the latter reward function. This suboptimal behavior is possibly partially explained by the response length, which drops dramatically around when the model achieves its best performance on the training data. The model learns the baseline reward function most quickly, and then seems to proceed to overfit for that signal, resulting in a poor translation to downstream validation set performance.

The second row (partially reported in Figure 7) demonstrates the potential of “hybrid rewards” that incorporate the rubric-augmented classifier into a more complex reward function. Adding a length penalty,  $\max(0, \frac{1}{200}(200 - \xi))$ , to the reward results in a more stable validation curve that better resembles the critic reward, and improves the best validation accuracy achieved by 4.7% and outperforms the verifiable rewards setup by 3.8%.

Seconds-per-step timing information is included in Figure 6, which shows that the rubric-augmented classifier does not incur a significant increase in runtime compared to the other two reward function approaches.

1026 E PROMPTS

1027

1028

1029

E.1 INITIAL TRACE COMPRESSION (GENERAL)

1030

1031

1032

1033

You are given a reasoning trace output by a model trying to answer a question. The reasoning trace involves the model’s exploration process to find the logical path to the correct answer. Your job is to distill out the entire logical path that it ultimately converged to and describe it to me in detailed, enumerated steps.

1034

1035

1036

1037

Potential errors should not be omitted from your logical path; provide the final deductive steps as faithfully to the model’s trace as possible. The final trace should resemble the deductive steps that a college student would write on a exam to ”show their work” to their professor (ideally in complete sentences).

1038

1039

1040

Write your distillation of the model’s logical reasoning path, and nothing else, i.e. don’t include ”here is the distilled path” or anything like that.

1041

1042

1043

1044

1045

1046

TRACE:

““

{t}

““

1047

1048

FINAL DISTILLED LOGICAL PATH:

1049

1050

E.2 RUBRIC EXTRACTION (GENERAL)

1051

1052

You are a reasoning analysis system that (1) identifies errors in students’ logical arguments and (2) explains how others can identify these same errors in other arguments.

1053

1054

1055

1056

1057

1058

1059

1060

1061

You will be given a question, correct answer, and the reasoning behind a student’s incorrect answer. Identify the error made by the student that caused them to arrive at an incorrect answer. An error description should be concise (max 25 words), but descriptive enough that it could be easily identified in other students’ responses in the future. Do not mention the student in an error: Your description should describe a reasoning mistake in the abstract that can be used by other students when solving similar problems, without calling out this specific student. For similar reasons, do not make the error overly specific to the problem, as the error should be sufficiently generalizable to other similar problems that may have different details. Anyone who makes the error on future problems will likely lead to a wrong answer.

1062

1063

1064

You will write instructions in a specific format that will describe how someone could identify the error if it appeared in a different logical argument for a different question. Your instructions should be organized into two parts: A ”key concept” and ”verification details”.

1065

1066

1067

1068

1069

\* KEY CONCEPT: This is a short keyword or concept that someone can quickly identify as being present or not present in a reasoning argument to determine if the error could possibly occur in that argument. This should be the lowest level of granularity: If this keyword or concept is present, then the error cannot possibly be present in the argument. Key concepts should be general: You want to aim for recall, not precision.

1070

1071

1072

1073

1074

\* VERIFICATION DETAILS: Provided that one or more of the key concepts are present in an argument, write one or more descriptive explanations for how someone could tell if the error exists in that argument. If one of these descriptions matches up with the contents of a new reasoning argument, then that means the error has occurred there.

1075

1076

1077

Write your answers in JSON format with your error as a key and your instructions in another dictionary as the corresponding value, i.e. `{{(error_description): {{(key_concept): (scanning_keyword), (verification): (verification_details)}}}}`. Here is an example output:

1078

1079

`{{(Failing to verify that atomic/molecular balances are maintained in chemical equations results in a stoichiometric balance violation): {{(key_concept): (chemical equations), (verification): [(Number of atoms of each element isn’t equal on both sides of reactions), (Molar ratios in calculations`

1080 don't match stoichiometric coefficients", "Product quantities exceed theoretical maximum based on  
 1081 reactant stoichiometry", "Component balances don't account for stoichiometric relationships"]}}}

1082  
 1083 Now, here is the data you should use to formulate your output:

1084 Question: "{tq}"

1085 Correct answer: "{ra}"

1086 Student reasoning: "{tt}"

1087 Output the error JSON and \*\*\*nothing else\*\*\*:  
 1088  
 1089

### 1090 E.3 KEYWORD AGGREGATOR (GENERAL) 1091

1092 A model is trying to match pre-written feedback to different samples of students' work. Here are  
 1093 some examples of the pre-written feedback in the corpus:

1094 {samples}  
 1095

1096 However, there are more possible feedback pieces than can fit in a single prompt. So, the model first  
 1097 labels each student's work with one or more "key concepts", which it then uses to figure out which  
 1098 feedback might be applicable and should be included in the final prompt.

1099 The following candidate "key concept" labels have been generated. Your job is to consolidate highly  
 1100 similar labels that describe the same concepts, to produce a shorter final list of key concepts.

1101 {keywords}  
 1102

1103 Group the concepts by the underlying sort of material they describe, and then for each group write  
 1104 a new key concept in the same style that can replace the labels in that group. Organize your output  
 1105 in JSON format with your new labels as keys and the lists of original label indices as the values,  
 1106 e.g., {"description 1": [1, 15, 62], "description 2": [4, 5, 22, 23, 45], ...}. Output the JSON and  
 1107 nothing else.

### 1108 E.4 RUBRIC KEYWORD APPLICATOR (GENERAL) 1109

1110 You are a reasoning analysis system that identifies core features of reasoning traces produced by  
 1111 other models. You are an expert at tagging these reasoning traces with labels to accurately cover  
 1112 all relevant aspects of the trace so the trace can be searched for in a database later. Given a reason-  
 1113 ing trace and list of concept labels, identify each label in the list that describes a concept or idea  
 1114 present in the reasoning trace. You are aiming for complete coverage - every relevant label should  
 1115 be included.

1116 Go down the list item-by-item to see if each one applies to the trace. Write the name of each label  
 1117 that relates to the reasoning trace in a Python list, e.g., ["label1", "label2", "label3", ...]. Output your  
 1118 Python list and nothing else.

1119 All concept labels:

1120 ""  
 1121 ""  
 1122 {r}

1123 ""  
 1124 ""  
 1125 Question:

1126 ""  
 1127 ""  
 1128 {q}

1129 ""  
 1130 Reasoning trace:

1131 ""  
 1132 ""  
 1133 {t}

1134 ""  
1135  
1136 Matching keywords and concepts (in Python list form):  
1137  
1138 E.5 RUBRIC ITEM APPLICATOR (GENERAL)  
1139  
1140 You are a teacher's assistant. You grade student responses to questions based on whether or not  
1141 they will arrive at the correct answer and why. For this class, the quality of a student's reasoning  
1142 generally doesn't matter, and they are only graded on final responses. Due to a server failure, one  
1143 student's final answer was deleted from the database, leaving only their work leading up to their  
1144 answer. Your job is to use the student's work to best identify whether the student's final answer was  
1145 correct or incorrect.  
1146  
1147 You have been given the teacher's grading rubric that identifies specific reasons why a student's  
1148 answer may be incorrect based on their reasoning process, and provides specific ways to tell if one  
1149 of these errors has occurred in a student's work. Use this rubric to predict whether the student's  
1150 answer should be marked "correct" or "incorrect".  
1151  
1152 First, read through the student's work. Then, go through the teacher's grading rubric, line-by-line,  
1153 and see if any of the entries and their descriptions describe behavior in the student's work that  
1154 will cause them to arrive at the wrong answer. To provide transparent documentation of your own  
1155 reasoning, write an explanation for each rubric item describing why or why not the item explains  
1156 why the student got the answer correct or not.  
1157  
1158 NOTE: The teacher's rubric are suggestions to check for, but just because the student exhibits the  
1159 erroneous behavior described in a given rubric item DOESN'T mean that they will definitely arrive  
1160 at the wrong answer. ONLY report back rubric items if you are confident that the erroneous behavior  
1161 is a highly likely cause for why the student's reasoning arrived at the wrong answer.  
1162  
1163 Do not skip any items - provide an explanation for each one. Explain one error type per line, and  
1164 do not abridge the process by saying things like [Continuing through all list items...]. I want all of  
1165 the items listed, you do not have to ask me again. I do not care how long the output is. Write your  
1166 reasons in JSON format, e.g., {1: "explanation for error 1", 2: "explanation for error 2", ...}.  
1167  
1168 Then, at the end after your JSON of explanations, you will enter a final grading database entry  
1169 using <ANS> brackets. Write a final response of "N/A" if you are confident the student ended up  
1170 at a correct answer. Write the error number(s) of rubric reasons if the student's work ends up at an  
1171 incorrect answer FOR THOSE SPECIFIC REASONS. Write your answer in <ANS> brackets, e.g.  
1172 <ANS>3, 41</ANS>, <ANS>22</ANS>, <ANS>N/A</ANS>, etc.  
1173  
1174 Error rubric: {rubric}  
1175  
1176 Question: {question}  
1177  
1178 Student's work: {trace}  
1179  
1180 Responses:  
1181  
1182  
1183  
1184  
1185  
1186  
1187

## 1188 E.6 INITIAL TRACE COMPRESSION (CODING)

1189

1190 You are given a reasoning trace output by a model trying to solve a PR request. The reasoning trace  
 1191 involves the model’s exploration process to find the correct solution. A critic model will be used  
 1192 to determine whether any mistakes were made in the process, but the existing trace is too long to  
 1193 be passed in directly. Derive the entire approach and summarize it for the critic model, describing  
 1194 decisions made and any possible mistakes in detailed, enumerated steps.

1195 ““ {t} ““

1196

## 1197 E.7 RUBRIC EXTRACTION (CODING)

1198

1199 You are a reasoning analysis system that (1) identifies errors in models’ solutions to GitHub issues  
 1200 and (2) explains how a critic model can identify these same errors in other failed solutions.

1201 You will be given a GitHub issue and the summary of a model’s incorrect solution to the GitHub  
 1202 issue. Identify the specific error made by the model that caused it to fail. An error description should  
 1203 be concise (max 25 words), but descriptive enough that it could be easily identified in other models’  
 1204 responses in the future. Do not mention the model in an error: Your description should describe a  
 1205 reasoning mistake in the abstract that can be used by a critic model when evaluating other solutions.  
 1206 For similar reasons, do not make the error overly specific to the problem, as the error should be  
 1207 sufficiently generalizable to other similar issues that may have different details.

1208 You will write instructions in a specific format that will describe how someone could identify the  
 1209 error if it appeared in a different solution for a different issue. Your instructions should be organized  
 1210 into two parts: A ”key concept” and ”verification details”.

1211 \* KEY CONCEPT: This is a short keyword or concept that someone can quickly identify as being  
 1212 present or not present in a solution to determine if the error could possibly occur in that issue. This  
 1213 should be the lowest level of granularity: If this keyword or concept is present, then the error cannot  
 1214 possibly be present in the solution. Key concepts should be general: You want to aim for recall, not  
 1215 precision.

1216 \* VERIFICATION DETAILS: Provided that one or more of the key concepts are present in a solu-  
 1217 tion, write one or more descriptive explanations for how someone could tell if the error exists in that  
 1218 solution. If one of these descriptions matches up with the contents of a solution, then that means the  
 1219 error has occurred there.

1220 Write your answers in JSON format with your error as a key and your instructions in another dictio-  
 1221 nary as the corresponding value, i.e. {{{error\_description}: {{"key\_concept": <scanning\_keyword>,  
 1222 "verification": <verification\_details>}}}}.

1223 Now, here is the data you should use to formulate your output:

1224 GitHub issue: ”{tq}”

1225 Summary of the incorrect solution: ”{tt}”

1226 Output the error JSON and \*\*\*nothing else\*\*\*:

1227

## 1228 E.8 KEYWORD AGGREGATOR (CODING)

1229

1230 A model is trying to match pre-written feedback to different incorrect solutions to GitHub issues.  
 1231 Here are some examples of the pre-written feedback in the corpus:

1232 {samples}

1233 However, there are more possible feedback pieces than can fit in a single prompt. So, the model first  
 1234 labels each issue solution with one or more ”key concepts”, which it then uses to figure out which  
 1235 feedback might be applicable and should be included in the final prompt.

1236 The following candidate ”key concept” labels have been generated. Your job is to consolidate highly  
 1237 similar labels that describe the same concepts, to produce a shorter final list of key concepts.

1238 {keywords}

1239

1242 Group the concepts by the underlying sort of material they describe, and then for each group write  
 1243 a new key concept in the same style that can replace the labels in that group. Organize your output  
 1244 in JSON format with your new labels as keys and the lists of original label indices as the values,  
 1245 e.g., `{ "description 1": [1, 15, 62], "description 2": [4, 5, 22, 23, 45], ... }`. Output the JSON  
 1246 and nothing else. Make sure your JSON output has correct Python syntax and all lists have correct  
 1247 closing brackets: Lists should always have an opening bracket `"["` as well as a closing bracket `"]"`,  
 1248 e.g., `"[4, 55, 120]"`.

#### 1249 1250 E.9 RUBRIC KEYWORD APPLICATOR (CODING)

1251 You are a reasoning analysis system that is an expert at tagging summaries of solutions to GitHub  
 1252 issues with labels to accurately cover all relevant aspects of the solution so that it can be searched for  
 1253 in a database later. Given a solution summary and list of keyword labels, identify each label in the  
 1254 list that describes a concept or idea present in the solution. You are aiming for complete coverage -  
 1255 every relevant label should be included.

1256 Go down the list item-by-item to see if each one applies to the solution. Write the name of each  
 1257 label that relates to the solution in a Python list, e.g., `["label1", "label2", "label3", ...]`. Output your  
 1258 Python list and nothing else.

1259 All concept labels: `"{r}"`

1260 GitHub issue: `"{q}"`

1261 Solution: `"{t}"`

1262 Matching keywords and concepts (in Python list form):

#### 1263 1264 1265 1266 E.10 RUBRIC ITEM APPLICATOR (CODING)

1267 You are a critic model that takes summaries of AI solutions to GitHub issues and determines whether  
 1268 or not they will produce a correct solution and why. The quality of the path taken to the final pull  
 1269 request doesn't matter, as long as the final solution submitted is correct and solves the issue.

1270 You have been given an error rubric that identifies specific reasons why a solution may be incorrect  
 1271 based on the model's reasoning process, and provides specific ways to tell if one of these errors has  
 1272 occurred in a solution. Use this rubric to predict whether the solution should be marked "correct" or  
 1273 "incorrect".

1274 First, read through the solution. Then, go through the error rubric, line-by-line, and see if any of  
 1275 the entries and their descriptions describe behavior in the solution that will cause it to fail. To  
 1276 provide transparent documentation of your own reasoning, write an explanation for each rubric item  
 1277 describing why or why not the item explains why the solution will fail or not. Do not skip any items  
 1278 - provide an explanation for each one. Explain one error type per line, and do not abridge the process  
 1279 by saying things like [Continuing through all list items...]. I want all of the items listed, you do not  
 1280 have to ask me again. I do not care how long the output is. Write your reasons in JSON format, e.g.,  
 1281 `{1: "explanation for error 1", 2: "explanation for error 2", ...}`.

1282 Then, at the end after your JSON of explanations, you will enter a final database entry using `<ANS>`  
 1283 brackets. Write a final response of "N/A" if you are confident the GitHub issue solution is cor-  
 1284 rect. Write the error number(s) if the solution will fail, for one or more of the reasons in the  
 1285 error rubric. Write your answer in `<ANS>` brackets, e.g. `<ANS>3, 41</ANS>`, `<ANS>22</ANS>`,  
 1286 `<ANS>N/A</ANS>`, etc.

1287 Error rubric: `{rubric}`

1288 GitHub issue: `{question}`

1289 Solution summary: `{trace}`

1290 Responses:

1291  
1292  
1293  
1294  
1295

1296 E.11 BASELINE CLASSIFIER (GENERAL)  
1297

1298 A reasoning system is given a question to answer, and it outputs a reasoning trace as an intermediate  
1299 step towards outputting a final answer. Your job is to determine whether the reasoning trace is likely  
1300 to result in a correct final answer. Respond with "correct" or "incorrect" in  $\langle \text{ANS} \rangle \langle / \text{ANS} \rangle$  brackets,  
1301 and nothing else.

1302 Question: {q}

1303 Reasoning trace: {t}  
13041305 E.12 BASELINE, ALT. 2& 3 (GENERAL)  
1306

1307 A reasoning system is given a question to answer, and it outputs a reasoning trace as an intermediate  
1308 step towards outputting a final answer. Given an excerpt from the outputted reasoning trace, your  
1309 job is to determine whether the system is likely to result in a correct final answer. Respond with  
1310 "correct" or "incorrect" in  $\langle \text{ANS} \rangle \langle / \text{ANS} \rangle$  brackets, and nothing else.

1311 Question: {q}

1312 Reasoning trace: {t}  
1313  
13141315 E.13 BASELINE, ALT. 4 & 5 (GENERAL)  
1316

1317 A reasoning system is given a question to answer, and it outputs a reasoning trace as an intermediate  
1318 step towards outputting a final answer. Given an excerpt from the outputted reasoning trace, your  
1319 job is to determine whether the system should output an answer now or keep thinking and extend  
1320 the reasoning trace. Respond with "correct" or "incorrect" in  $\langle \text{ANS} \rangle \langle / \text{ANS} \rangle$  brackets, and nothing  
1321 else.

1322 Question: {q}

1323 Reasoning trace: {t}  
13241325 E.14 BASELINE CLASSIFIER (CODING)  
1326

1327 A model is asked to solve a PR request from GitHub, and it outputs a reasoning trace documenting  
1328 its path towards completing the PR request. Your job is to determine whether the reasoning trace is  
1329 likely to result in a correct completion of the PR request. Respond with "correct" or "incorrect" in  
1330  $\langle \text{ANS} \rangle \langle / \text{ANS} \rangle$  brackets, and nothing else.

1331 Question: {q}

1332 Reasoning trace: {t}  
1333  
13341335 E.15 BASELINE, ALT. 2& 3 (CODING)  
1336

1337 A model is asked to solve a PR request from GitHub, and it outputs a reasoning trace documenting  
1338 its path towards completing the PR request. Given an excerpt from the outputted reasoning trace,  
1339 your job is to determine whether the system is likely to result in a correct completion of the PR  
1340 request. Respond with "correct" or "incorrect" in  $\langle \text{ANS} \rangle \langle / \text{ANS} \rangle$  brackets, and nothing else.

1341 Question: {q}

1342 Reasoning trace: {t}  
13431344 E.16 BASELINE, ALT. 4 & 5 (CODING)  
1345

1346 A model is asked to solve a PR request from GitHub, and it outputs a reasoning trace documenting  
1347 its path towards completing the PR request. Given an excerpt from the outputted reasoning trace,  
1348 your job is to determine whether the system should output an answer now or keep thinking and  
1349 extend the reasoning trace. Respond with "correct" or "incorrect" in  $\langle \text{ANS} \rangle \langle / \text{ANS} \rangle$  brackets, and  
nothing else.

1350 Question: {q}  
1351 Reasoning trace: {t}  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

## 1404 F EXAMPLE DATA

1405

## 1406 F.1 CHEMICAL ENGINEERING QUESTION

1407

1408 In a mixture to be fed to a continuous distillation column, the mole fraction of phenol is 0.35, o-  
 1409 cresol is 0.15, m-cresol is 0.30 and xylenols is 0.20. A product is required with a mole fraction of  
 1410 phenol of 0.952, o-cresol 0.0474 and m-cresol 0.0006. If the volatility to o-cresol of phenol is 1.26  
 1411 and of m-cresol is 0.70, estimate how many theoretical plates would be required at total reflux.

1412

## 1413 F.2 COMPRESSED CHEMICAL ENGINEERING TRACE

1414

1415 1. Given a multicomponent distillation column operating at total reflux with:

- 1416 - Feed composition: phenol (0.35), o-cresol (0.15), m-cresol (0.30), xylenols (0.20)
- 1417 - Desired distillate composition: phenol (0.952), o-cresol (0.0474), m-cresol (0.0006)
- 1418 - Relative volatilities to o-cresol: phenol (1.26), m-cresol (0.70)

1419 2. Set up material balances for each component:

- 1420 - Phenol:  $0.35F = 0.952D + B_{ph}$
- 1421 - O-cresol:  $0.15F = 0.0474D + B_o$
- 1422 - M-cresol:  $0.30F = 0.0006D + B_m$

1423 Where F is feed rate, D is distillate rate, B is bottoms rate

1424 3. Define  $q = D/F$  (distillate fraction) and establish two Fenske equations:

- 1425 - For phenol/o-cresol:  $20.0844 * (0.15 - 0.0474q) / (0.35 - 0.952q) = 1.26^N$
- 1426 - For phenol/m-cresol:  $79.0 * [(0.30 - 0.0006q) / (0.15 - 0.0474q)] = 1.42857^N$

1427 4. Through iterative calculations, found  $q \approx 0.265$  where both equations give  $N \approx 14.45$ :

- 1428 - With  $q = 0.265$ :
- 1429 \*  $B_{ph} = 0.0977F$
- 1430 \*  $B_o = 0.13744F$
- 1431 \*  $B_m = 0.29984F$

1432 5. Verified result by checking separation ratios:

- 1433 - For phenol/o-cresol:  $(D_{ph}/D_o) / (B_{ph}/B_o) \approx 28.26 = 1.26^{14.45}$
- 1434 - For phenol/m-cresol:  $(D_{ph}/D_m) / (B_{ph}/B_m) \approx 4865.5 \approx 1.8^{14.45}$

1435 6. Since number of plates must be an integer and we need to ensure separation, round up to 15  
 1436 theoretical plates.

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

1458 F.3 MATH QUESTION  
14591460 4. Let in the  $xOy$  plane, the area enclosed by  $0 < y \leq x^2, 0 \leq x \leq 1$  be  $\frac{1}{3}$ . Then the intersection of  
1461 the sets

1462

1463 
$$M = \{(x, y) \mid |y| - |x| \leq 1\},$$

1464

1465 
$$N = \{(x, y) \mid |y| \geq x^2 + 1\}$$

1466 representing the area of  $M \cap N$  is ( ). (A)  $\frac{1}{3}$  (B)  $\frac{2}{3}$  (C) 1 (D)  $\frac{4}{3}$   
14671468 F.4 COMPRESSED MATH TRACE  
14691470 1. First identify that set  $M$  is defined by  $|y| \leq |x| + 1$ , which forms a region bounded by  $y = |x| + 1$   
1471 and  $y = -(|x| + 1)$ 

1472

1473 2. Set  $N$  is defined by  $|y| \geq x^2 + 1$ , which forms two regions:

1474

- Above the parabola  $y = x^2 + 1$ 

1475

- Below the parabola  $y = -x^2 - 1$ 

1476

3. The intersection  $M \cap N$  requires points satisfying  $x^2 + 1 \leq |y| \leq |x| + 1$ 

1477

4. In the lower half-plane ( $y \leq 0$ ), the intersection contains only discrete points and contributes no  
1478 area

1479

5. In the upper half-plane ( $y \geq 0$ ), we need points where:

1480

-  $x^2 + 1 \leq y \leq |x| + 1$ 

1481

- This is only possible for  $|x| \leq 1$  (since  $x^2 + 1 > |x| + 1$  when  $|x| > 1$ )

1482

6. Due to symmetry about the  $y$ -axis, calculate area in first quadrant and multiply by 2:

1483

- For  $0 \leq x \leq 1$ , area between  $y = x + 1$  and  $y = x^2 + 1$ 

1484

- Vertical distance between curves is  $(x + 1) - (x^2 + 1) = x - x^2$ 

1485

7. Calculate area in first quadrant:

1486

1487 
$$\int_0^1 (x - x^2) dx = \left[ \frac{1}{2}x^2 - \frac{1}{3}x^3 \right]_0^1 = \frac{1}{2} - \frac{1}{3} = \frac{1}{6}$$

1488

8. Total area =  $2 \times \frac{1}{6} = \frac{1}{3}$ 

1489

Therefore, the area of  $M \cap N$  is  $\frac{1}{3}$ .

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

## 1512 F.5 CODING QUESTION

1513

1514 **caplog.get\_records** and **caplog.clear** conflict

1515

1516 **Description**

1517 `caplog.get_records()` gets decoupled from actual `caplog` records when `caplog.clear()`  
 1518 is called. As a result, after `caplog.clear()` is called, `caplog.get_records()` is frozen: it  
 1519 does not get cleared, nor does it get new records.

1520 During test set up it is set to the same list as `caplog.records`, but the latter gets replaced rather  
 1521 than cleared in `caplog.clear()`, which diverges the two objects.

1522

1523 **Reproductive example**

1524

```
1524 import logging
```

1525

```
1526 def test(caplog) -> None:
1527     def verify_consistency() -> None:
1528         assert caplog.get_records("call") == caplog.records
1529
1530         verify_consistency()
1531         logging.warning("test")
1532         verify_consistency()
1533         caplog.clear()
1534         verify_consistency() # fails: assert [<LogRecord: ..., 8, "test">] == []
```

1535 **Environment details** Arch Linux, Python 3.9.10:

1536

1537 Package	Version
1538 -----	-----
1539 attrs	21.4.0
1540 iniconfig	1.1.1
1541 packaging	21.3
1542 pip	22.0.4
1543 pluggy	1.0.0
1544 py	1.11.0
1545 pyparsing	3.0.8
1546 pytest	7.1.1
1547 setuptools	60.10.0
1548 tomli	2.0.1
1549 wheel	0.37.1

## 1550 F.6 COMPRESSED CODING TRACE

1551

## 1552 1. Problem Identification:

- 1553 - The issue was in the `LogCaptureHandler.reset()` method
- 1554 - When `caplog.clear()` was called, it created a new list instead of clearing the existing one
- 1555 - This broke the reference between `caplog.get_records()` and `caplog.records`

1556

## 1557 2. Investigation Steps:

- 1558 - Created a reproduction script to simulate the issue
- 1559 - Verified that creating a new list breaks references
- 1560 - Tested alternative approaches to clearing the list

1561

## 1562 3. Solution:

- 1563 - Changed `self.records = []` to `del self.records[:]` in `LogCaptureHandler.reset()`
- 1564 - This clears the list in-place while maintaining the original reference
- 1565 - Verified the fix with a test script

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

4. Key Decisions:

- Used `del list[:]` instead of alternatives like `list.clear()` or `list[:] = []`
- Kept the change minimal and focused on the core issue
- Maintained backward compatibility
- Did not require changes to test files

5. Potential Issues:

- None identified - the change is straightforward and maintains existing behavior
- The fix is well-tested with both reproduction and verification scripts

6. Validation:

- Created test scripts to verify both the bug and the fix
- Ensured the fix maintains list reference consistency
- Confirmed the solution works with `pytest`'s fixture behavior

The approach was methodical and focused on the root cause, resulting in a minimal but effective fix.

1620 G AI USAGE STATEMENT

1621

1622

1623

1624

1625

1626

1627

1628

1629

1630

1631

1632

1633

1634

1635

1636

1637

1638

1639

1640

1641

1642

1643

1644

1645

1646

1647

1648

1649

1650

1651

1652

1653

1654

1655

1656

1657

1658

1659

1660

1661

1662

1663

1664

1665

1666

1667

1668

1669

1670

1671

1672

1673

LLMs were used to format small sections of LaTeX in the paper.