

---

# ReasonOps: Operator Segmentation for LLM Reasoning Traces

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Chain-of-thought traces from large reasoning models can span tens of thousands  
2 of tokens, yet we lack a vocabulary for describing their internal structure. Pre-  
3 vious methods developed to analyze chain-of-thought traces are either too rigid  
4 or not expressive enough, failing to capture features across domains and models.  
5 To remedy this, we develop ReasonOps, an unsupervised, expressive method for  
6 annotating chain-of-thought traces, providing succinct universal operators. Us-  
7 ing ReasonOps, we analyze 44,662 traces from 12 thinking LLMs spanning 6  
8 families across 8 reasoning benchmarks and discover that they share a common  
9 compositional structure: 7 recurring reasoning operators—discourse-level moves  
10 such as BACKTRACKING, INFERRING, and HYPOTHESIZING—that emerge from  
11 unsupervised clustering of sentence-initial 3-token pivots. These operators appear  
12 across every model family and benchmark domain, confirmed by three independ-  
13 ent LLM judges who classify held-out samples at 70–76% accuracy. We analyze  
14 the structure of operators on easy vs. hard problems, revealing that reflective  
15 operators are more helpful on hard problems and harm performance on easy prob-  
16 lems. Operator sequences are highly model-identifying: a classifier trained on  
17 operator distributions alone recovers the source model with macro-AUC = 0.987,  
18 revealing that each model family has a distinctive reasoning fingerprint. Structural  
19 operator features predict within-problem answer correctness well above baselines.  
20 Classifiers built on these operators reach WP-AUC = 0.701 and 0.801 on AIME  
21 specifically. ReasonOps further enables early quality estimation well before the  
22 trace completes: we predict at WP-AUC = 0.664 for only 50% of the trace. The  
23 ReasonOps pipeline is unsupervised and annotation-free, enabling deep insights  
24 into LLM reasoning traces as well as strong downstream results on model identi-  
25 fication and correctness prediction.

## 26 1 Introduction

27 Reasoning-capable large language models (LLMs) increasingly solve difficult problems by produc-  
28 ing long intermediate traces before emitting a final answer. LLMs have become synonymous with  
29 "large reasoning models" (LRMs) as the majority of frontier models today are post-trained to elicit  
30 extended chain-of-thought reasoning before producing a final answer [1, 2, 3, 4]. Prompting meth-  
31 ods, multi-path decoding, process supervision, and reinforcement-learning-based post-training have  
32 all improved performance on mathematics, science, and coding benchmarks, but they have also  
33 made reasoning traces longer, more expensive, and harder to analyze [5, 6, 3]. This comes at an  
34 increasing demand for monitoring and oversight of LLM decision-making processes [7, 8, 9].

35 Since the advent of chain-of-thought (CoT) prompting [5], reasoning traces have been hypothesized  
36 as a window into the problem-solving capabilities of LLMs. Reasoning traces, or chain-of-thought  
37 traces, are valuable artifacts for black-box understanding of LLMs as they can be obtained without

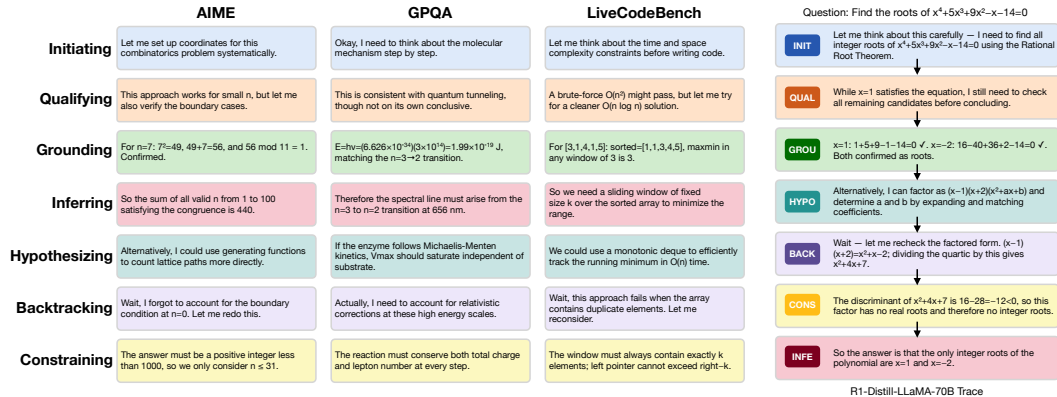


Figure 1: (Left) representative operator units extracted from AIME, GPQA, and LiveCodeBench. (Right) representative reasoning trace extracted from R1-Distill-LLaMA-70B on MATH dataset.

38 access to model weights. These traces have been shown to contain rich information [10, 11], yet  
 39 we lack a uniform vocabulary under which to characterize reasoning traces across models, domains,  
 40 and datasets. While some methods have attempted to provide systematic annotations for reasoning  
 41 traces, these methods often rely on *a priori* vocabularies and syntactic structures that are overfit  
 42 to current models or domains [12, 13, 14, 15, 16]. In this paper, we ask: can we compress diverse traces  
 43 into a small vocabulary of reasoning operators shared across model families, tasks, and domains  
 44 while preserving predictive information about success and failure?

45 **Our contributions.** We introduce ReasonOps, an unsupervised framework for inducing a compact  
 46 vocabulary of reasoning operators from visible chain-of-thought traces. We use sentence-initial  
 47 3-token pivots, frequency filtering, and semantic clustering with e5-small [17]. We show that the  
 48 resulting 7 operators are quantitatively meaningful: they generalize across 12 thinking LLMs from  
 49 6 families and 8 benchmarks, confirmed by three independent LLM judges (70–76% classifica-  
 50 tion accuracy; chance: 14%). Operator sequences identify the source model with macro-AUC  
 51 = 0.987, revealing model-family reasoning fingerprints. We show that operator features predict  
 52 within-problem correctness above all span-free baselines and the LLM self-judge (the OST reaches  
 53 0.701 cross-dataset, matching the content-augmented Op-XGB classifier while reading only opera-  
 54 tor labels), and that the OST—trained once on full sequences—enables correctness estimation from  
 55 partial traces at any depth, surpassing an Op-XGB upper bound retrained per depth. We open-source  
 56 our codebase as a resource for the community.<sup>1</sup>

## 57 2 Related Work

58 **Annotating and analyzing reasoning traces.** Deepseek-R1 creators noted that "Wait" phrases were  
 59 emergent behavior associated with reasoning capability increases [3]. This and the interpretability  
 60 of chain-of-thought traces inspired work into annotating and describing reasoning traces. However,  
 61 previous works primarily rely on *a priori* vocabularies of reasoning [12, 13, 18], operate on arbitrary  
 62 syntactic features such as sentences [14, 15], or rely on domain-specific vocabularies [16].  
 63 ReasonOps is an unsupervised method that does not assume *a priori* syntactic features or a fixed  
 64 vocabulary of annotations. Other works have included *ad hoc* methods to analyze reasoning traces  
 65 for particular hypotheses [10, 11].

66 **Monitoring and scalable oversight.** The scalable oversight community has extensively explored  
 67 ways to characterize LLM behavior [9, 19, 20], and chain-of-thought monitoring has been explored  
 68 as a prominent direction [8, 7]. This work is tangential to ReasonOps in terms of reasoning trace  
 69 characterization. Some work has also shown that chain-of-thought can be an unfaithful representa-  
 70 tion of model behavior [21], but recent work has attempted to improve this in frontier models [7, 22].  
 71 Our work occupies a complementary point: we treat visible traces as behaviorally meaningful arti-  
 72 facts and ask whether they admit a stable, annotation-free meso-scale abstraction.

73 **Correctness prediction.** Some methods have been proposed for early correctness prediction, in-  
 74 cluding LLM-only methods [23, 24], and mechanistic interpretability techniques [25]. Similar to

<sup>1</sup>Anonymized repository: <https://anonymous.4open.science/r/ReasonOps>

75 this line of work is that of process reward models [6] and outcome reward models [26], including  
76 work on verifiers in test-time scaling [27, 28]. This work has been particularly prolific in mathemat-  
77 ics problems [26, 6, 29]. Other work has explored test time scaling techniques for self-verification,  
78 including SelfCheck [23] and others [30, 31, 32].

### 79 3 Methods

80 We now describe how we build an unsupervised mechanism to infer operators from reasoning traces  
81 across 12 LLMs. Our pipeline is driven by linguistic principles and requires no *a priori* definition  
82 of an operator vocabulary.

83 **Data collection.** We collect reasoning traces from 12 thinking language models spanning 6 families  
84 (Table 4 in Appendix A) on 8 benchmarks: MATH-500 [33], GPQA Diamond [34], AIME 2024,  
85 LiveCodeBench [35], HumanEval [36], MMLU-Pro [37], ARC-Challenge [38], and BIG-Bench  
86 Hard [39]. All models are queried with a 65,536-token context budget. Raw chain-of-thought tokens  
87 are collected via the OpenRouter API for all models except Claude, which is queried directly via  
88 the Anthropic API with extended thinking enabled; per-model API details are given in Appendix A.  
89 After removing truncated and malformed traces, we retain **44,662 traces** ( $\approx 3,720$  per model).

90 **Pivot-based span segmentation.** Each trace is segmented into **spans**. A span begins at a *pivot*  
91 *sentence*: any sentence whose first three alpha tokens (the *sentence start*) appear frequently enough  
92 across the corpus to be recognized as a discourse signal. Formally,  $\text{pivot}(s) = (w_1, w_2, w_3)$  where  
93  $w_i$  are the ordered lowercase alphabetic tokens at the start of sentence  $s$ . A discourse pivot  $p$  is  
94 included if (1) it appears in  $\geq 100$  distinct traces (*frequency filter*), (2) it appears in traces from  $\geq 3$   
95 distinct datasets (*domain-diversity filter*), and (3) all three tokens  $w_i$  belong to the top-2,000 most  
96 frequent corpus tokens (*vocabulary filter*).

97 These thresholds are set to balance specificity and generality. The frequency threshold of 100  
98 traces is chosen to exclude idiosyncratic phrases that appear in only a handful of traces while re-  
99 taining common discourse moves; lowering it to 50 adds many domain-specific phrases that are  
100 not discourse-functional. The domain-diversity filter of 3 datasets ensures that any accepted pivot  
101 spans multiple task types, preventing dataset-specific phrases from entering the vocabulary. The  
102 top-2,000 vocabulary filter ensures pivots are composed of common English function words rather  
103 than domain-specific content tokens (e.g., chemistry formulae, programming keywords); we verified  
104 that varying this between 1,500 and 3,000 does not qualitatively change the resulting clusters.

105 **Semantic pivot embedding and clustering.** We embed all discourse pivots using the  
106 `intfloat/e5-small-v2` sentence encoder [17]. Sentence-level embeddings capture full-phrase  
107 meaning, correctly separating pivots such as “let me think” and “let me verify” that share a common  
108 prefix. We cluster pivot embeddings with  $K$ -means ( $K \in \{6, \dots, 11\}$ , 30 restarts), selecting  $K$  by  
109 maximizing Cohen’s  $\kappa$  against an independent LLM judge on held-out spans. The optimal  $K=7$   
110 yields  $\kappa = 0.693\text{--}0.720$  across three judges (Table 6, Appendix C).

111 **Independence from correctness.** The operator discovery pipeline uses no correctness labels: pivots  
112 are filtered by frequency and domain diversity, embeddings are fit on pivot text alone, and cluster  
113 assignments are by nearest-centroid lookup. Correctness labels enter only in the downstream predic-  
114 tion step (§6), where Op-XGB and the OST are trained in a standard supervised 5-fold CV protocol.  
115 The operators themselves are thus a fully unsupervised intermediate representation.

116 **Computational cost.** The discovery pipeline is lightweight. Pivot extraction (frequency counting  
117 over all sentence starts) takes under 3 minutes on a single CPU core for 44,662 traces. Embedding  
118 the 5,464 accepted pivots with `e5-small-v2` takes 6 seconds;  $K$ -means on 5,464 points takes 12  
119 seconds. Annotating a new trace requires only a dictionary lookup per span (no embedding at  
120 inference time), adding  $<1$  ms per trace. End-to-end operator discovery over the full corpus takes  
121 under 5 minutes, and annotating the entire 44K-trace corpus takes under 2 minutes.

### 122 4 Reasoning operators

123 Seven operators emerge from unsupervised clustering (Table 1; Figure 1). They span the space  
124 of discourse moves in extended reasoning: goal-setting, fact-anchoring, inference-drawing, explo-  
125 ration, hedging, self-correction, and constraint-specification. As a more coarse classification, opera-  
126 tors can be broken into *committal operators* (INITIATING, INFERRING, CONSTRAINING, GROUND-  
127 ING) and *reflective operators* (QUALIFYING, BACKTRACKING, HYPOTHESIZING). These operators  
128 recur across all 12 models and all 8 benchmarks, suggesting that extended reasoning—regardless of

Table 1: The 7 discovered reasoning operators. Representative pivots are the most frequent 3-tuples in each cluster.

Operator	Representative Pivots	Description
INITIATING	let me think, let me check, let me verify	Explicitly launches a new cognitive operation
QUALIFYING	hmm, but the problem, but let me	Introduces caveats or complications
GROUNDING	i need to, the question is, looking at	Anchors reasoning in facts or given information
INFERRING	so the answer, thus final answer, yes	Draws a conclusion from prior steps
HYPOTHESIZING	alternatively perhaps, for example if, so perhaps	Entertains a tentative or conditional scenario
BACKTRACKING	wait, wait no, wait let me	Signals potential error and prepares to restart
CONSTRAINING	we need to, now we need, so we need	Identifies necessary conditions, narrows solution space

129 training data, architecture, or problem domain—organizes itself into a common compositional vocabulary. Figure 1 shows representative span excerpts per operator across three of the eight benchmarks. An illustration of a full operator sequence in Figure 1 and Appendix J.

132 **Evaluation with an LLM judge.** We first test whether the discovered operators are semantically coherent from the perspective of independent judges. Each LLM judge receives only cluster names, descriptions, and the top 12 representative pivot phrases per cluster—no training spans. It labels 50 held-out spans per cluster ( $n = 350$  total). Table 2 reports classification accuracy for three judges. All reach 70–76% on the 7-way classification task, 4.9–5.3 $\times$  above the 14.3% chance baseline. Adding 3 example spans per cluster to the prompt raises accuracy from 69.8% to **74.6%**: few-shot examples aid pattern-matching from a closed 7-label set.

139 **Cluster naming stability.** To assess the semantic stability of the clusters, we re-ran the LLM naming step for the clusters by sampling 30 seeds of exemplars. For each name and resulting description, we computed pairwise cosine similarity of description embeddings across the seeds. The resulting descriptions were highly consistent within clusters with a median pairwise cosine similarity of 0.660, substantially above both an across-cluster baseline of 0.467 ( $r_{rb} = +0.72$ ) and a matched random-group control of 0.565 ( $r_{rb} = +0.39$ ;  $p < 10^{-149}$ ). All seven operators showed significant within-cluster stability, indicating that the naming procedure recovers coherent and reproducible operator semantics; by contrast, randomly grouped spans did not elicit comparably coherent names. See Appendix K.1 for more details.

Table 2: LLM-judge validation ( $n=350$  spans, 50 per cluster). Judges receive cluster names, descriptions, and top-12 representative pivots. **Exemplars:** prompt includes 3 example spans per cluster. Full  $\kappa$  values are reported in Appendix C.

Judge	Prompt variant	Accuracy
Random baseline	—	14.3%
Claude Haiku 4.5	baseline	76.0%
GPT-5.4-Mini	baseline	71.7%
Claude Sonnet 4.6	baseline	69.7%
Claude Sonnet 4.6	exemplars	<b>74.6%</b>

## 156 5 Analysis of operator distributions

157 Operators show distinct patterns when stratified by dataset and model (Figure 2). Mathematical reasoning benchmarks (AIME, MATH) show higher GROUNDING and CONSTRAINING (step-by-step derivation), while commonsense benchmarks (ARC-C, BBH) use more GROUNDING and elevated HYPOTHESIZING. Challenging multiple choice datasets like GPQA and MMLU-Pro exhibit higher use of HYPOTHESIZING and QUALIFYING than any other dataset. Coding benchmarks (LiveCodeBench, HumanEval) exhibit elevated CONSTRAINING, consistent with specification-first programming. Models further show operator usage fingerprints that roughly cluster with model families, a phenomenon that’s supported by the results in §6.1. R1-distill model profiles are both very similar, and Qwen models show almost identical operator proportions. Kimi/Moonshot models show high CONSTRAINING usage. Claude Haiku, Claude Sonnet, and Grok-3-Mini use excessive GROUNDING. Despite these differences, all 7 operators are represented in every model.

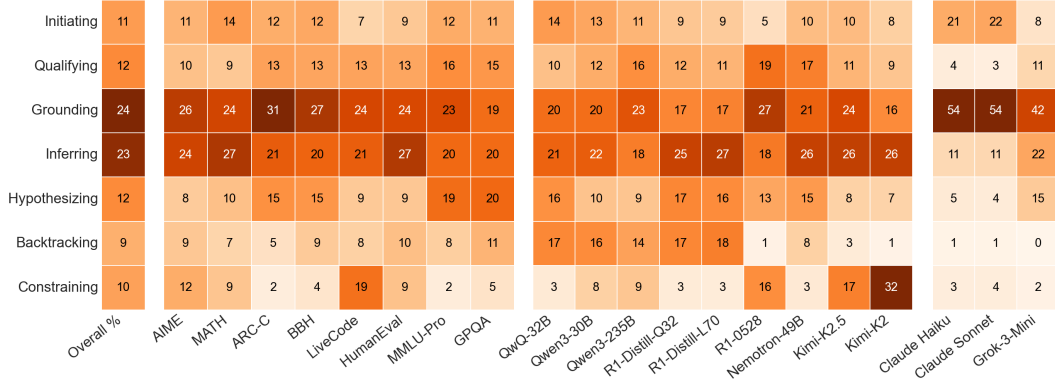


Figure 2: Distribution of operators across datasets (top) and models (bottom). Metrics are percent usage = average appearance of that operator for that given model/dataset vs. all other operators.

168 **Transition structure.** The operator transition matrix reveals Markov structure in reasoning:  
 169 **GROUNDING** strongly self-transitions (sustained fact-anchoring), while **BACKTRACKING** most of  
 170 ten leads into **INITIATING** (correction followed by a fresh attempt). **HYPOTHESIZING** frequently  
 171 precedes **INFERRING**, consistent with a generate-and-test schema. Run-length analysis reinforces  
 172 this: **GROUNDING** and **CONSTRAINING** accumulate long self-runs, while **BACKTRACKING** is al-  
 173 most always a single span—a brief interruption rather than a sustained state (Figure 11, Appendix L).  
 174 We expand on this **BACKTRACKING** phenomenon in §5.2.

## 175 5.1 Operator usage on correct vs. incorrect traces

176 We see a few trends emerge in terms of opera-  
 177 tor usage on correct vs. incorrect traces (Figure  
 178 3). First, committal operators seem to be over-  
 179 represented in correct traces than in incorrect  
 180 traces, with **GROUNDING**, **INFERRING**, **INITI-**  
 181 **ATING**, and **CONSTRAINING** all being positive  
 182 in correct - incorrect usage. Conversely, reflect-  
 183 ive operators are globally seen more often in  
 184 incorrect traces than correct traces. **HYPOTH-**  
 185 **ESIZING** and **QUALIFYING** show particularly  
 186 higher usage in incorrect traces than correct  
 187 traces. This may be indicative of correct traces  
 188 showing more confident, forward-stepping op-  
 189 erations rather than reflecting. This then led us  
 190 to investigate these dynamics on easy vs. hard  
 191 problems.

192 **Easy vs. hard problems.** We next define a  
 193 problem as “hard” if fewer than 1/3 of models solve it, and “easy” if more than 2/3 models pass. We  
 194 first filter to only **BBH**, **LiveCodeBench**, **MMLU-Pro**, and **GPQA**, which have sufficient samples  
 195 of hard problems. For easy problems, correct traces exhibit significantly higher usage of committal  
 196 operators than reflective operators, indicating confident forward-stepping through the problem. The  
 197 usage of committal - reflective operators can be defined as the *committal-reflective gap*. On easy  
 198 problems, this gap is much higher in correct traces (+44.2%) than in incorrect traces (+7.5%)  
 199 and is strongly predictive of correctness (AUC = 0.76 [0.74, 0.78] 95% CI). On average across all  
 200 problems, committal usage peaks early and late in traces while reflective operators appear mid-trace  
 201 (Figure 4a).

202 Hard problems show less obvious differences between correct and incorrect traces. The committal-  
 203 reflective gap flips, showing +30.4% on correct traces and +34.9% on incorrect traces, a -4.5%  
 204 shift which is not statistically-significant (Mann-Whitney p=0.15) (Figure 4a). The predictive signal  
 205 disappears as the AUC for correct vs. incorrect via committal-reflective gap alone drops to AUC =  
 206 0.48 [0.44, 0.51]. However, examining **HYPOTHESIZING** vs. **INFERRING**, instead of all committal

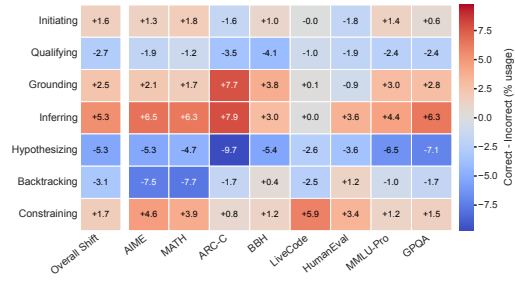


Figure 3: Shift of operator usage by model in correct and incorrect traces. % usage difference = take percent usage (Figure 2) and subtract usage in correct and incorrect traces.

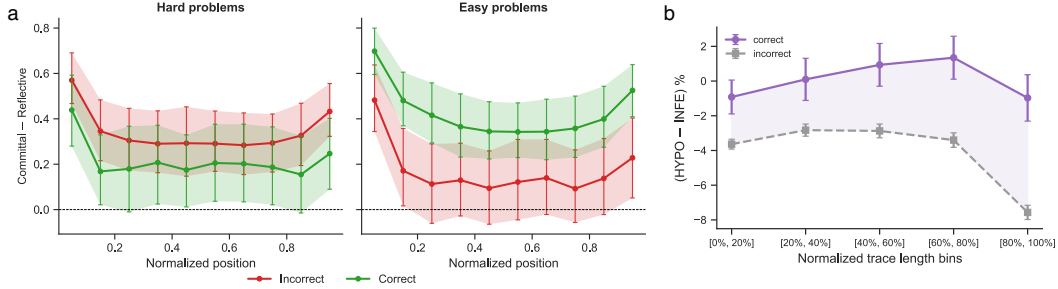


Figure 4: (a) Temporal distribution of committal - reflective operators for correct and incorrect problems. 95% CI shown at each normalized position. (b) HYPOTHESIZING- INFERRING gap over only hard problems for correct and incorrect traces. Shaded regions show the gap between correct and incorrect curves.

207 and reflective operators, yields some effects. When examining the distribution of these operators  
 208 over time, the HYPOTHESIZING- INFERRING gap is particularly representative for correct vs. incor-  
 209 rect hard problems. This gap increases monotonically over the length of the trace, reaching 6.6%  
 210 at the tail of the trace, with incorrect traces showing much higher usage of INFERRING over HYPOTH-  
 211 ESIZING (Figure 4b). This shows that reflective operators, specifically HYPOTHESIZING, might  
 212 provide some benefit on hard problems while harming performance on easy problems.

## 213 5.2 Backtracking is a local operator

214 We centered the next analysis around BACKTRACKING and HYPOTHESIZING. Both operators se-  
 215 mantically appear to interrupt a model’s current line of reasoning, but we asked whether these op-  
 216 erators actually reflected global strategic shifts in the model’s problem-solving. We used Claude  
 217 Sonnet 4.6 as an LLM judge to classify a stratified sample of GPQA events as LOCAL (re-checks  
 218 a single calculation; method unchanged), SUB-PROBLEM (re-opens a specific case), or GLOBAL  
 219 (abandons the current strategy), prompting the judge with the full reasoning trace and the event  
 220 marked from the operator span to its next same-operator span. Local backtracking is explicitly de-  
 221 fined as “The model is re-checking a single calculation, fact lookup, or specific claim...”, and global  
 222 event is defined as “model abandons or proposes to replace its current solution...”. **BACKTRACK-**  
 223 **ING and HYPOTHESIZING are overwhelmingly local revisions rather than strategy changes,**  
 224 with statistically indistinguishable scope distributions: BACKTRACKING is 85.4% LOCAL / 12.8%  
 225 SUB-PROBLEM / 1.6% GLOBAL ( $n = 1,294$ ), and HYPOTHESIZING is 80.2% / 18.2% / 1.6%  
 226 ( $n = 192$ ). This is supported by literature that shows that backtracking is often superficial [40] and  
 227 rarely changes a model’s answer [41]. More details are in Appendix H.

## 228 6 Downstream Applications

229 We now demonstrate how operators discovered by ReasonOps can be used to predict downstream  
 230 properties of reasoning traces, namely model identity and correctness of the trace.

### 231 6.1 Model Identification

232 If operators capture stable behavioral structure, traces from the same model family should cluster  
 233 together in operator space. We train **Op-XGB**, an XGBoost classifier [42] (400 trees, max  
 234 depth 6) on a curated operator feature set: a 117-dim handcrafted operator feature vector—global  
 235 and quartile-localized operator frequencies, bigram transition matrix, run-length statistics, first/last  
 236 operator one-hots, and entropy/length scalars—concatenated with an 8,000-feature anchor-phrase  
 237 TF-IDF representation; the TF-IDF is re-fit inside each training fold (full feature definitions in Ap-  
 238 pendix E). Op-XGB is trained under problem-level 5-fold CV with model identity as the target. It  
 239 achieves overall accuracy = 79.9% (chance = 8.3%) and macro-AUC = **0.987**. Per-model one-vs-  
 240 rest AUCs (Figure 5) span 0.967 for R1-Distill-Llama-70B to 0.997 for Kimi-K2, Kimi-K2.5, and  
 241 Grok-3-Mini, with the lowest AUCs concentrated in within-family pairs where intra-family confu-  
 242 sion is hardest. The full 12-class confusion matrix (Figure 12, Appendix L.2) makes this explicit:  
 243 R1-Distill-Llama-70B traces are predicted as R1-Distill-Qwen-32B at 0.41 and the reverse at 0.29  
 244 (the largest off-diagonal mass), consistent with their shared distillation lineage from DeepSeek-R1;

Table 3: Within-problem AUC (WP-AUC) for correctness prediction. **CD** = 5-fold problem-level CV trained on all datasets pooled; **ID** = 5-fold CV trained and evaluated within each dataset. Length, Backtrack, and Wait use logistic regression on a single monotonic feature; AUC is rank-invariant to the training fold (CD=ID). Op-XGB: 117-dim operator features (frequencies, quartile localization, bigram transitions, run lengths) + 8K anchor-phrase TF-IDF, concatenated; XGBoost (400 trees, depth 6). See Appendix E for full feature details. OST at 100% trace depth; see §6.3 for partial-trace results. **Bold** = best. Underline = 2nd best. **Our methods**.

Dataset	Op-XGB		OST		Length		Backtrack		Wait		SelfCheck	
	CD	ID	CD	ID	CD	ID	CD	ID	CD	ID	CD	ID
AIME	0.779	<b>0.838</b>	0.801	0.797	0.500	0.500	0.621	0.621	0.679	0.679	0.501	0.501
ARC-Challenge	<u>0.618</u>	0.559	<b>0.645</b>	0.556	0.559	0.559	0.520	0.520	0.501	0.499	0.578	0.578
GPQA	0.682	<b>0.703</b>	<u>0.691</u>	<u>0.691</u>	0.591	0.591	0.591	0.591	0.602	0.602	0.526	0.526
LiveCodeBench	0.754	<b>0.795</b>	0.727	0.732	0.426	0.574	0.616	0.616	0.591	0.591	0.496	0.496
MATH	0.597	0.570	<u>0.612</u>	<b>0.662</b>	0.407	0.407	0.464	0.464	0.483	0.517	0.504	0.504
MMLU-Pro	0.635	<b>0.639</b>	0.635	<u>0.638</u>	0.587	0.587	0.628	0.628	0.633	0.633	0.520	0.520
<b>Global</b>	0.701	<b>0.723</b>	0.701	<u>0.703</u>	0.507	0.551	0.594	0.594	0.600	0.603	0.512	0.512

245 the Claude-Haiku-4.5 / Claude-Sonnet-4.5 pair shows mutual confusion in the 0.21–0.26 range; and  
 246 Kimi-K2.5 is predicted as Kimi-K2 at 0.10 (with the reverse at 0.05).

## 247 6.2 Correctness Prediction

248 We next evaluate whether operator composition  
 249 can be used to predict answer correctness.  
 250 We focus on the 6 benchmarks where operator  
 251 structure is most informative: AIME, ARC-  
 252 Challenge, GPQA, LiveCodeBench, MATH-  
 253 500, and MMLU-Pro. We exclude BIG-Bench  
 254 Hard and HumanEval: BBH uses answer-  
 255 matching formats where operator structure has  
 256 near-chance predictive signal, and HumanEval  
 257 correctness is determined by code execution  
 258 rather than reasoning quality. We verify in  
 259 §6 that this exclusion is empirically motivated.  
 260 Our primary metric, *WP-AUC* (within-problem  
 261 AUC; adapted from [43]), computes AUC sep-  
 262 arately for each problem’s trace group and  
 263 averages—isolating the ability to discriminate  
 264 better from worse attempts on the *same* prob-  
 265 lem, with difficulty factored out. We use 5-fold  
 266 problem-level CV. All traces for a given prob-  
 267 lem are assigned to the same fold, preventing  
 268 leakage from multi-model evaluation.

269 We compare six methods: (1) **trace length** (log  
 270 response character count), (2) **backtrack count**  
 271 (fraction of spans labeled BACKTRACKING), (3) **wait count** (raw count of “wait”-prefixed spans),  
 272 (4) **SelfCheck** [23]: the same model that generated the trace reads the problem and its full reasoning  
 273 and predicts correctness without the gold answer—no training, (5) **Op-XGB**, and (6) **OST** (Opera-  
 274 tor Sequence Transformer; described in §6.3). Methods (1)–(3) use logistic regression on a single  
 275 scalar feature with the same 5-fold problem-level CV protocol; since a 1-D logistic regression is a  
 276 monotone transformation, AUC is fold-invariant (CD=ID). Method (4) requires no training. Meth-  
 277 ods (5)–(6) are evaluated under two protocols: **CD** (cross-dataset), trained on all datasets pooled,  
 278 and **ID** (in-dataset), trained and evaluated within each dataset separately.

279 Table 3 reports WP-AUC per dataset for all methods.

280 Trace length is near chance globally (WP-AUC = 0.551 within-dataset). Backtrack count (0.594)  
 281 and wait count (0.603) outperform length, confirming that reflective operator frequency carries sig-

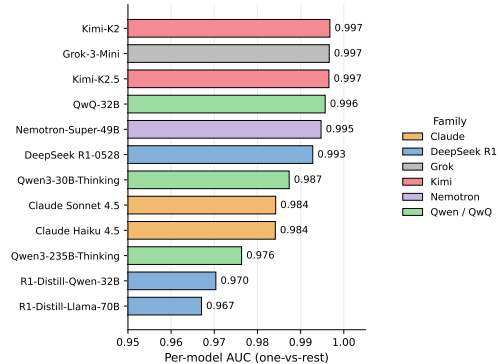


Figure 5: Per-model one-vs-rest AUC for the Op-XGB model identification classifier, colored by family. AUCs span 0.967 (R1-Distill-Llama-70B) to 0.997 (Kimi-K2, Kimi-K2.5, and Grok-3-Mini, tied); within-family pairs (R1-Distill, Claude, Kimi) sit at the lower end where intra-family confusion is concentrated.

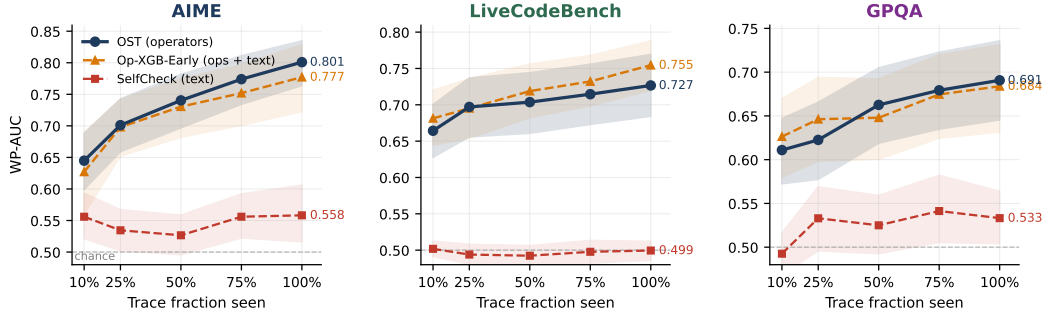


Figure 6: Correctness prediction WP-AUC at increasing trace depths for the three best-performing datasets (AIME, LiveCodeBench, GPQA). Navy solid (**OST**): operator sequence transformer trained once on full sequences, reading the discrete operator label sequence at any prefix length without retraining. Amber dashed (**Op-XGB-Early**): full Op-XGB recipe (117-dim operator features + 8,000-dim anchor TF-IDF; Appendix E), retrained from scratch at each depth as an upper-bound comparison. Red dashed (**SelfCheck**): the same model that generated the trace reads its own raw reasoning up to depth  $p\%$  and predicts correctness with no training. Shaded bands: 95% CIs.

282 nal even without full sequence modeling. SelfCheck [23]—the same model that generated the trace  
 283 reading its own reasoning and predicting correctness—reaches WP-AUC = 0.512 globally (near  
 284 chance), a striking null result: reasoning models are overconfident and predict “correct” for both  
 285 correct and incorrect traces at similar rates. Op-XGB reaches 0.701 cross-dataset and 0.723 within-  
 286 dataset globally, with strong within-dataset performance on AIME (0.838 ID) and LiveCodeBench  
 287 (0.795 ID). The OST reaches 0.701 cross-dataset and 0.703 within-dataset at full trace depth—  
 288 tied with Op-XGB cross-dataset despite reading only operator labels, with no access to anchor phrase  
 289 text. OST surpasses Op-XGB on ARC-Challenge (0.645 vs. 0.618 CD) and MATH (0.662 vs. 0.570  
 290 ID), where temporal operator dynamics capture structure that aggregate feature counts miss.

### 291 6.3 Early Correctness Prediction

292 A practical question is whether correctness can be estimated from a partial trace, before the model  
 293 finishes reasoning. We introduce the **Operator Sequence Transformer (OST)**: a lightweight Trans-  
 294 former encoder ( $\sim 800\text{K}$  parameters; 4 layers,  $d=128$ , 4 heads, pre-LayerNorm) trained end-to-end  
 295 directly on operator sequences with no pretraining. Each operator token is represented as the sum  
 296 of a unigram embedding, a *bigram* (transition) embedding encoding the previous  $\rightarrow$  current operator  
 297 pair, and a continuous sinusoidal position encoding proportional to the token’s normalized position  
 298 within the visible prefix—giving the model a fine-grained sense of *where* in the trace each operator  
 299 appears. The model is trained with a pure contrastive loss that directly encourages correct traces  
 300 to score higher than incorrect traces on the same problem, under the same 5-fold problem-level  
 301 CV protocol as the other methods. Crucially, the OST is *trained once on full operator sequences*;  
 302 at inference time it accepts any operator-sequence prefix without retraining, making it a practical  
 303 single-model early-prediction system deployable at any trace depth. More details are in Appendix D.

304 Figure 6 shows OST WP-AUC vs. trace depth for the three best-performing datasets. On AIME,  
 305 WP-AUC rises from 0.645 at 10% depth to 0.740 at 50% and 0.801 at 100%—at half the trace it  
 306 already surpasses all span-free baselines. LiveCodeBench and GPQA follow similar trajectories,  
 307 reaching 0.727 and 0.691 at full depth. Globally, the OST rises from WP-AUC = 0.605 at 10%  
 308 depth to 0.664 at 50% and 0.701 at 100%, monotonically improving with trace length.

309 We compare against **Op-XGB-Early**, which applies the same Op-XGB feature recipe at each depth  
 310  $d$ : we truncate every trace to its first  $d\%$  of spans, recompute the operator-sequence statistics and the  
 311 TF-IDF on the truncated span set, and fit a fresh XGBoost from scratch. Per-depth retraining gives  
 312 Op-XGB-Early an explicit upper-bound advantage—it is re-optimized for each truncation level—  
 313 whereas the OST is trained only once on full sequences. Globally, Op-XGB-Early reaches WP-AUC  
 314 0.620 at 10%, 0.647 at 25%, 0.661 at 50%, 0.682 at 75%, and 0.699 at 100%. The OST and Op-  
 315 XGB-Early track each other within a few thousandths across all depths: OST–Op-XGB-Early gaps  
 316 are  $-0.015$ ,  $-0.016$ ,  $+0.003$ ,  $+0.006$ , and  $+0.003$  at depths  $\{10, 25, 50, 75, 100\}\%$ . The single  
 317 OST model thus matches a per-depth-retrained Op-XGB—which has access to both the operator

318 features *and* the raw anchor text—for partial-trace prediction at 50% depth and beyond, despite  
319 reading only the discrete operator label sequence and never being retrained for partial inputs.

320 By contrast, the partial-trace SelfCheck baseline—the same model reading its own first  $p\%$  of raw  
321 reasoning text—is near chance at all depths (0.50–0.56 across datasets), confirming that the OST’s  
322 signal comes from the structural operator representation rather than from text content alone.

## 323 7 Discussion

324 ReasonOps is an unsupervised method for annotating reasoning traces, built from careful observa-  
325 tions of their semantic structure. Three-token pivots strike a useful balance for discourse analysis:  
326 they disambiguate the polysemy of single-token markers like "but" or "so" (e.g., distinguishing  
327 "wait," "wait no," and "wait let me" as backtracking, rejection, and correction respectively) while  
328 remaining short enough to recur across tens of thousands of traces. The top-2,000 vocabulary fil-  
329 ter is essential to this approach, as it excludes high-frequency but domain-specific phrases (e.g.,  
330 "the double bond" in chemistry traces) that would otherwise contaminate the pivot set with non-  
331 discourse-functional content. ReasonOps is a robust method to provide domain-agnostic, model-  
332 agnostic annotations of reasoning traces that enable universal analysis and downstream tasks such  
333 as model identification and correctness prediction.

334 **Future work and limitations.** We see several avenues for future work on annotating reasoning  
335 traces. An extension could be made to agentic traces, i.e., traces of actions taken by an LLM when  
336 enabled with a harness. The ReasonOps approach, which is unsupervised and does not assume *a*  
337 *priori* sets of operators, might be more useful for agentic traces where actions may change drasti-  
338 cally based on the harness. The operators define a structured vocabulary for step-level correctness  
339 prediction without human annotation. OST’s early correctness signal can trigger intervention (e.g.,  
340 beam expansion, early stopping) at any trace position, enabling compute-adaptive generation. There  
341 are also several limitations of ReasonOps. First, the pipeline relies on sentence-initial pivots; spans  
342 bounded at a sub-sentence level cannot be labeled. Second, the operators are validated by LLM  
343 judgment, which introduces its own biases. Third, cross-model transfer is limited: operator fre-  
344 quency features partially capture model-specific patterns that do not fully generalize. Finally, the  
345 corpus covers 8 benchmarks in English; generalization to other languages is not verified.

346 **Broader impacts.** This work paves the way to a common language for describing reasoning traces  
347 of LLMs. This could be used for further analysis of LLMs as well as downstream tasks beyond  
348 correctness prediction.

349 **Summary.** Seven reasoning operators emerge from unsupervised analysis of 44,662 chain-of-  
350 thought traces across 12 thinking LLMs and 8 benchmarks. Three independent LLM judges con-  
351 firm the taxonomy (70–76% classification accuracy; chance: 14%). Operator sequences identify the  
352 source model with macro-AUC = 0.987. Structural features predict within-problem correctness well  
353 above all span-free baselines and above the LLM self-judge (0.512, near chance). The OST reaches  
354 WP-AUC = 0.701 cross-dataset—tied with Op-XGB (0.701 CD) while reading only operator la-  
355 bels, not the additional text features Op-XGB receives. The OST predicts correctness monotonically  
356 from partial traces, reaching 0.664 at 50% trace depth and surpassing an Op-XGB baseline retrained  
357 at each depth, with no per-depth retraining required. The convergence of architecturally diverse  
358 models on the same 7 discourse-level moves suggests that extended reasoning organizes itself into a  
359 common compositional vocabulary.

## 360 References

- 361 [1] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,  
362 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*  
363 *arXiv:2505.09388*, 2025.
- 364 [2] Kimi Team, Tongtong Bai, Yifan Bai, Yiping Bao, SH Cai, Yuan Cao, Y Charles, HS Che,  
365 Cheng Chen, Guanduo Chen, et al. Kimi k2.5: Visual agentic intelligence. *arXiv preprint*  
366 *arXiv:2602.02276*, 2026.
- 367 [3] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu,  
368 Ruoyu Zhang, Shirong Ma, Xiao Bi, et al. Deepseek-r1 incentivizes reasoning in llms through  
369 reinforcement learning. *Nature*, 645(8081):633–638, 2025.

- 370 [4] Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus,  
371 Rahul K Arora, Yu Bai, Bowen Baker, Haiming Bao, et al. gpt-oss-120b & gpt-oss-20b model  
372 card. *arXiv preprint arXiv:2508.10925*, 2025.
- 373 [5] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,  
374 Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models.  
375 *Advances in neural information processing systems*, 35:24824–24837, 2022.
- 376 [6] Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee,  
377 Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The*  
378 *twelfth international conference on learning representations*, 2023.
- 379 [7] Melody Y Guan, Miles Wang, Micah Carroll, Zehao Dou, Annie Y Wei, Marcus Williams,  
380 Benjamin Arnav, Joost Huizinga, Ian Kivlichan, Mia Glaese, et al. Monitoring monitorability.  
381 *arXiv preprint arXiv:2512.18311*, 2025.
- 382 [8] Tomek Korbak, Mikita Balesni, Elizabeth Barnes, Yoshua Bengio, Joe Benton, Joseph Bloom,  
383 Mark Chen, Alan Cooney, Allan Dafoe, Anca Dragan, et al. Chain of thought monitorability: A  
384 new and fragile opportunity for ai safety. *arXiv preprint arXiv:2507.11473*, 2025.
- 385 [9] Samuel R Bowman, Jeeyoon Hyun, Ethan Perez, Edwin Chen, Craig Pettit, Scott Heiner, Kamilė  
386 Lukošūūtė, Amanda Askell, Andy Jones, Anna Chen, et al. Measuring progress on scalable  
387 oversight for large language models. *arXiv preprint arXiv:2211.03540*, 2022.
- 388 [10] Parshin Shojaee, Iman Mirzadeh, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and  
389 Mehrdad Farajtabar. The illusion of thinking: Understanding the strengths and limitations of  
390 reasoning models via the lens of problem complexity. *arXiv preprint arXiv:2506.06941*, 2025.
- 391 [11] Junsol Kim, Shiyang Lai, Nino Scherrer, James Evans, et al. Reasoning models generate  
392 societies of thought. *arXiv preprint arXiv:2601.10825*, 2026.
- 393 [12] Shiming Yang, Yuxuan Tong, Xinyao Niu, Graham Neubig, and Xiang Yue. Demystifying long  
394 chain-of-thought reasoning. In *Forty-second International Conference on Machine Learning*,  
395 2025.
- 396 [13] Constantin Venhoeff, Iván Arcuschin, Philip Torr, Arthur Conmy, and Neel Nanda. Understand-  
397 ing reasoning in thinking language models via steering vectors. In *Workshop on Reasoning and*  
398 *Planning for Large Language Models*, 2025.
- 399 [14] Jinu Lee, Sagnik Mukherjee, Dilek Hakkani-Tur, and Julia Hockenmaier. Reasoningflow:  
400 Semantic structure of complex reasoning traces. *arXiv preprint arXiv:2506.02532*, 2025.
- 401 [15] Paul C Bogdan, Uzay Macar, Neel Nanda, and Arthur Conmy. Thought anchors: Which llm  
402 reasoning steps matter? *arXiv preprint arXiv:2506.19143*, 2025.
- 403 [16] Ming Li, Nan Zhang, Chenrui Fan, Hong Jiao, Yanbin Fu, Sydney Peters, Qingshu Xu, Robert  
404 Lissitz, and Tianyi Zhou. Understanding the thinking process of reasoning models: A perspec-  
405 tive from schoenfeld’s episode theory. In *Proceedings of the 2025 Conference on Empirical*  
406 *Methods in Natural Language Processing*, pages 18278–18299, 2025.
- 407 [17] Liang Wang, Nan Yang, Xinlong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Ma-  
408 jumder, and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training. *arXiv*  
409 *preprint arXiv:2212.03533*, 2022.
- 410 [18] Priyanka Kargupta, Shuyue Stella Li, Haocheng Wang, Jinu Lee, Shan Chen, Oreaoghene  
411 Ahia, Dean Light, Thomas L Griffiths, Max Kleiman-Weiner, Jiawei Han, et al. Cognitive  
412 foundations for reasoning and their manifestation in llms. *arXiv preprint arXiv:2511.16660*,  
413 2025.
- 414 [19] Zachary Kenton, Noah Y Siegel, János Kramár, Jonah Brown-Cohen, Samuel Albanie, Jannis  
415 Bulian, Rishabh Agarwal, David Lindner, Yunhao Tang, Noah D Goodman, et al. On scalable  
416 oversight with weak llms judging strong llms. *Advances in Neural Information Processing*  
417 *Systems*, 37:75229–75276, 2024.

- 418 [20] Joshua Engels, David D. Baek, Subhash Kantamneni, and Max Tegmark. Scaling laws for  
419 scalable oversight. In *The Thirty-ninth Annual Conference on Neural Information Processing*  
420 *Systems*, 2026.
- 421 [21] Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny  
422 Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, et al. Measuring faith-  
423 fulness in chain-of-thought reasoning. *arXiv preprint arXiv:2307.13702*, 2023.
- 424 [22] Debjit Paul, Robert West, Antoine Bosselut, and Boi Faltings. Making reasoning matter: Mea-  
425 suring and improving faithfulness of chain-of-thought reasoning. In *Findings of the Association*  
426 *for Computational Linguistics: EMNLP 2024*, pages 15012–15032, 2024.
- 427 [23] Ning Miao, Yee Whye Teh, and Tom Rainforth. Selfcheck: Using llms to zero-shot check their  
428 own step-by-step reasoning. In *The Twelfth International Conference on Learning Representa-*  
429 *tions (ICLR)*, 2024.
- 430 [24] Yang Xiang, Yixin Ji, Ruotao Xu, Dan Qiao, Zheming Yang, Juntao Li, and Min Zhang. When  
431 is thinking enough? early exit via sufficiency assessment for efficient reasoning. *arXiv preprint*  
432 *arXiv:2604.06787*, 2026.
- 433 [25] Zheng Zhao, Yeskendir Koishckenov, Xianjun Yang, Naila Murray, and Nicola Cancedda. Ver-  
434 ifying chain-of-thought reasoning via its computational graph. In *The Fourteenth International*  
435 *Conference on Learning Representations*, 2026.
- 436 [26] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
437 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to  
438 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 439 [27] Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh  
440 Agarwal. Generative verifiers: Reward modeling as next-token prediction. In *The Thirteenth*  
441 *International Conference on Learning Representations*, 2025.
- 442 [28] Arian Hosseini, Xingdi Yuan, Nikolay Malkin, Aaron Courville, Alessandro Sordoni, and  
443 Rishabh Agarwal. V-STar: Training verifiers for self-taught reasoners. In *First Conference*  
444 *on Language Modeling*, 2024.
- 445 [29] Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and  
446 Zhifang Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations.  
447 In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*  
448 *(Volume 1: Long Papers)*, pages 9426–9439, 2024.
- 449 [30] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang,  
450 Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reason-  
451 ing in language models. In *The Eleventh International Conference on Learning Representations*,  
452 2023.
- 453 [31] Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré,  
454 and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated  
455 sampling. *arXiv preprint arXiv:2407.21787*, 2024.
- 456 [32] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute op-  
457 timally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*,  
458 2024.
- 459 [33] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn  
460 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset.  
461 *NeurIPS*, 2021.
- 462 [34] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien  
463 Dirani, Julian Michael, and Samuel R Bowman. GPQA: A graduate-level google-proof Q&A  
464 benchmark. *arXiv preprint arXiv:2311.12022*, 2024.

- 465 [35] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Ar-  
466 mando Solar-Lezama, Koushik Sen, and Ion Stoica. LiveCodeBench: Holistic and contami-  
467 nation free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*,  
468 2024.
- 469 [36] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto,  
470 Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul  
471 Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke  
472 Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad  
473 Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias  
474 Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex  
475 Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain,  
476 William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra,  
477 Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer,  
478 Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech  
479 Zaremba. Evaluating large language models trained on code, 2021.
- 480 [37] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo,  
481 Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and  
482 challenging multi-task language understanding benchmark. *Advances in Neural Information*  
483 *Processing Systems*, 37:95266–95290, 2024.
- 484 [38] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick,  
485 and Oyvind Tafjord. Think you have solved question answering? Try ARC, the AI2 reasoning  
486 challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- 487 [39] Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won  
488 Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, et al. Challenging big-bench  
489 tasks and whether chain-of-thought can solve them. In *Findings of the Association for Compu-*  
490 *tational Linguistics: ACL 2023*, pages 13003–13051, 2023.
- 491 [40] Hanbin Wang, Jingwei Song, Jinpeng Li, Qi Zhu, Fei Mi, Ganqu Cui, Yasheng Wang,  
492 and Lifeng Shang. Teaching large reasoning models effective reflection. *arXiv preprint*  
493 *arXiv:2601.12720*, 2026.
- 494 [41] Liwei Kang, Yue Deng, Yao Xiao, Zhanfeng Mo, Wee Sun Lee, and Lidong Bing. First try  
495 matters: Revisiting the role of reflection in reasoning models. *arXiv preprint arXiv:2510.08308*,  
496 2025.
- 497 [42] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of*  
498 *the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages  
499 785–794, 2016.
- 500 [43] Guorui Zhou, Xiaoqiang Zhu, Chengru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan,  
501 Junqi Jin, Han Li, and Kun Gai. Deep interest network for click-through rate prediction. In  
502 *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery &*  
503 *Data Mining*, pages 1059–1068, 2018.

504 **A Model Details**

505 Table 4 lists the 12 thinking models used for trace collection. All expose raw chain-of-thought tokens  
 506 accessible via API. All models except Claude are queried via the OpenRouter API; reasoning tokens  
 507 are returned in the `reasoning` field of the OpenRouter response message. The `reasoning` request  
 508 parameter is set to `{effort: high}` for Grok-3-Mini, `{max_tokens: N}` for Qwen3 thinking  
 509 variants, and omitted for DeepSeek, QwQ-32B, Kimi, and Nemotron models. Claude models are ac-  
 510 cessed directly via the Anthropic API with `thinking: {type: enabled, budget_tokens:`  
 511 `N}` enabled; raw thinking content is extracted from `type: thinking` response blocks.

Table 4: Models used for trace collection. All 12 expose raw chain-of-thought tokens.

Model	Family	Size	Notes
QwQ-32B	Qwen	32B	S[1]
Qwen3-30B-Thinking	Qwen	30B (3B act.)	MoE S[1]
Qwen3-235B-Thinking	Qwen	235B (22B act.)	MoE S[1]
DeepSeek-R1-0528	DeepSeek	671B MoE	S[3]
R1-Distill-LLaMA-70B	DeepSeek	70B	Distilled S[3]
R1-Distill-Qwen-32B	DeepSeek	32B	Distilled S[3]
Claude-Sonnet-4.5	Anthropic	—	Extended thinking S[4]
Claude-Haiku-4.5	Anthropic	—	Extended thinking S[4]
Grok-3-Mini	xAI	—	S[5]
Nemotron-Super-49B	NVIDIA	49B	S[6]
Kimi-K2-Thinking	Moonshot	MoE	S[7]
Kimi-K2.5	Moonshot	MoE	S[2]

512 **B Corpus Statistics**

513 The full corpus spans 8 benchmarks and contains 44,662 traces across 12 models. Table 5 reports  
 514 per-dataset statistics for the 6-dataset correctness prediction subset (33,209 traces), which excludes  
 515 BBH, HumanEval, and ARC-AGI-2 due to grading ambiguity or near-ceiling accuracy. Each prob-  
 516 lem was attempted by all 12 models with up to 5 independent samples per model; the AIME trace  
 517 count is lower because AIME 2024 has only 30 problems.

Table 5: Per-dataset trace counts in the 6-dataset correctness prediction corpus. All 12 models contribute to every dataset. Accuracy = fraction of traces judged correct by the official grader.

Dataset	Traces	Problems	Models	Accuracy
AIME 2024	4,604	90	12	86.4%
ARC-Challenge	5,808	97	12	94.9%
GPQA Diamond	5,615	100	12	74.0%
LiveCodeBench	5,516	99	12	33.5%
MATH-500	5,794	100	12	94.6%
MMLU-Pro	5,872	100	12	76.3%
<b>Total</b>	<b>33,209</b>	<b>586</b>		

518 The full operator-discovery corpus additionally includes BIG-Bench Hard and HumanEval (44,662  
 519 traces total across all 8 benchmarks). These two datasets are excluded from correctness prediction:  
 520 BBH uses answer-matching formats where operator structure shows near-chance predictive signal,  
 521 and HumanEval correctness is determined by code execution rather than reasoning quality.

522 **C K Selection**

523 Table 6 reports LLM-judge  $\kappa$  (Claude Sonnet 4.6) for each  $K$  in the sweep.  $K=7$  achieves the  
 524 highest  $\kappa$ .

Table 6:  $K$  selection via LLM-judge  $\kappa$  (Claude Sonnet 4.6).

$K$	$\kappa$
6	0.604
<b>7</b>	<b>0.693</b>
8	0.611
9	0.648
10	0.573
11	0.558

## 525 D Operator Sequence Transformer (OST) Architecture

526 The OST is a lightweight Transformer encoder trained end-to-end on operator token sequences. Its  
 527 design is optimized for long sequences of discrete operator labels (typical length: 50–500 tokens)  
 528 while remaining parameter-efficient enough to train in 5-fold CV.

529 **Tokenization.** Each span in a trace is assigned one of 7 operator labels (integer indices 0–6). The  
 530 OST processes the sequence of these labels for a trace prefix of specified depth.

531 **Input embeddings.** Each position  $t$  in the operator sequence receives a 128-dimensional input  
 532 vector formed as the sum of three components:

- 533 • **Unigram embedding:** a learned embedding  $E_u[o_t] \in \mathbb{R}^{128}$  for the current operator  $o_t$ ;
- 534 • **Bigram (transition) embedding:** a learned embedding  $E_b[o_{t-1}, o_t] \in \mathbb{R}^{128}$  for the  
 535 previous→current operator pair (special “start” token for  $t = 0$ );
- 536 • **Continuous sinusoidal position encoding:** standard sinusoidal encoding scaled by the  
 537 *normalized* position  $t/T$  within the visible prefix, giving the model a fine-grained sense of  
 538 where in the trace each operator appears.

539 **Transformer body.** The model uses a standard pre-LayerNorm Transformer encoder with 4 lay-  
 540 ers,  $d_{\text{model}} = 128$ , 4 attention heads, feed-forward width = 512, and 5% token dropout during  
 541 training. Attention is full (not causal), since the model reads a complete prefix.

542 **Pooling and scoring.** Sequence representations are aggregated by attention pooling (learned query  
 543 vector) into a single 128-dimensional vector, followed by a two-layer MLP ( $128 \rightarrow 64 \rightarrow 1$ )  
 544 producing a real-valued correctness score. Total parameter count:  $\approx 800,000$ .

545 **Training.** The OST is trained with a pure contrastive loss: for every problem with at least one  
 546 correct and one incorrect trace in the training fold, all pairs  $(i^+, i^-)$  of correct and incorrect traces for  
 547 that problem are formed, and the loss pushes the correct trace to score higher. This directly optimizes  
 548 the within-problem ranking that WP-AUC measures. Training uses AdamW ( $\eta = 3 \times 10^{-4}$ , weight  
 549 decay = 0.01), batch size 64, and early stopping on fold validation WP-AUC.

550 **Early prediction.** At inference time the OST is given only the first  $p\%$  of operator tokens (by span  
 551 count), allowing evaluation at any trace depth. The sinusoidal position encoding is re-normalized to  
 552 the partial prefix length.

## 553 E Op-XGB Feature Details

554 Op-XGB combines a 117-dimensional handcrafted operator-sequence feature vector with an 8,000-  
 555 feature anchor-phrase TF-IDF representation. Both feature sets are derived from the operator labels  
 556 and sentence-initial pivot phrases—no substantive math, code, or reasoning content of the trace is  
 557 used.

558 **Operator feature vector (117 dim).** For each trace’s operator sequence  $\sigma = (\sigma_1, \dots, \sigma_n) \in$   
 559  $\{0, \dots, 6\}^n$ :

- 560 • **Global frequencies** (7): the share of each operator across the full sequence,  $\text{freq}_k =$   
561  $\frac{1}{n} \sum_t \mathbf{1}[\sigma_t = k]$ .
- 562 • **Quartile frequencies** ( $7 \times 4 = 28$ ): the share of each operator within each of four equal-  
563 sized contiguous chunks of  $\sigma$ , capturing temporal localization.
- 564 • **Scalar summaries** (5): Shannon entropy of the global frequency vector, the maximum of  
565  $\text{freq}$ , the number of distinct operators that appear,  $\log(1 + n)$ , and the immediate-repeat  
566 rate  $\frac{1}{n-1} \sum_{t \geq 2} \mathbf{1}[\sigma_t = \sigma_{t-1}]$ .
- 567 • **First/last one-hot** ( $7 + 7 = 14$ ): one-hot encodings of  $\sigma_1$  and  $\sigma_n$ .
- 568 • **Bigram transition matrix** (49): the  $7 \times 7$  matrix of bigram frequencies  
569  $\frac{1}{n-1} \sum_{t \geq 2} \mathbf{1}[\sigma_{t-1}=a, \sigma_t=b]$ , flattened.
- 570 • **Run-length statistics** ( $7 + 7 = 14$ ): for each operator  $k$ , the mean and max length of  
571 consecutive runs of  $k$ , each normalized by  $n$ .

572 This vector is computed once per trace and depends only on the discrete operator label sequence; it  
573 carries no anchor text.

574 **Anchor-phrase TF-IDF (8,000 dim).** We additionally extract each span’s anchor (the sentence-  
575 initial pivot phrase, truncated to 80 characters), concatenate all anchors of a trace with single spaces,  
576 and fit `sklearn.feature_extraction.text.TfidfVectorizer` with `max_features=8000`  
577 and `sublinear_tf=True` on the training-fold anchor strings (default unigram tokenization, sub-  
578 linear TF  $1 + \log(\text{tf})$ , IDF reweighting, vocabulary re-fit inside each fold to prevent leakage). The  
579 anchor strings consist only of discourse-level pivot phrases (e.g., “let’s denote”, “so we have”, “wait,  
580 that means”)—they are exactly the inputs that the unsupervised operator-discovery pipeline embeds  
581 and clusters into the 7 operators.

582 **Classifier.** The 117-dim operator feature vector and the 8,000-dim TF-IDF vector are concatenated  
583 to form an 8,117-dim per-trace representation, which is passed to XGBoost (400 trees, max depth  
584 6, learning rate 0.05, subsample 0.8, `colsample_bytree` 0.8, logloss objective). The classifier is  
585 trained with problem-level 5-fold CV (cross-dataset) or within each dataset separately (in-dataset).  
586 XGBoost predicted probabilities are used directly as scores for WP-AUC computation. Op-XGB  
587 therefore reads operator information at *both* levels of abstraction (discrete labels and the raw pivot  
588 phrases that produce them), in contrast to the OST, which reads only the discrete labels.

589 **Op-XGB-Early variant.** For the early-prediction comparison (§6.3, Figure 6), at each depth  $d \in$   
590  $\{10, 25, 50, 75, 100\}\%$  we truncate each trace’s spans to the first  $\lceil N \cdot d/100 \rceil$  spans, recompute  
591 the 117-dim operator feature vector on the truncated operator subsequence, recompute the anchor  
592 TF-IDF on the truncated anchor text, and re-fit a fresh XGBoost on the concatenated representation  
593 under the same 5-fold problem-level CV. This per-depth retraining provides an upper bound on what  
594 the full Op-XGB feature recipe can achieve when the model is re-optimized for each truncation  
595 level—in contrast to the OST, which is trained once on full sequences and evaluated on partial  
596 prefixes without retraining.

597 **Op-XGB for model identification.** For 12-class source-model identification (Figure 12), we  
598 use the identical Op-XGB feature recipe described above (117-dim operator features + 8,000-  
599 dim TF-IDF, concatenated to an 8,117-dim vector) and train XGBoost as a multi-class classifier  
600 (`multi:softprob`, same hyperparameters as the binary version) under the same problem-level 5-  
601 fold CV protocol. Reusing the full Op-XGB recipe ensures that the same model class supports both  
602 correctness prediction and model identification.

## 603 F Operator Usage by Correctness per Benchmark

604 Figure 7 shows operator frequency separately for correct and incorrect traces. INFERRING is con-  
605 sistentally higher in correct traces; HYPOTHESIZING and BACKTRACKING are higher in incorrect  
606 traces. The gap is largest on commonsense benchmarks (ARC-C, MMLU-Pro:  $\sim 24\%$  incorrect  
607 vs.  $\sim 14\%$  correct for HYPOTHESIZING). On LiveCodeBench, CONSTRAINING is 10 pp higher in  
608 correct traces.

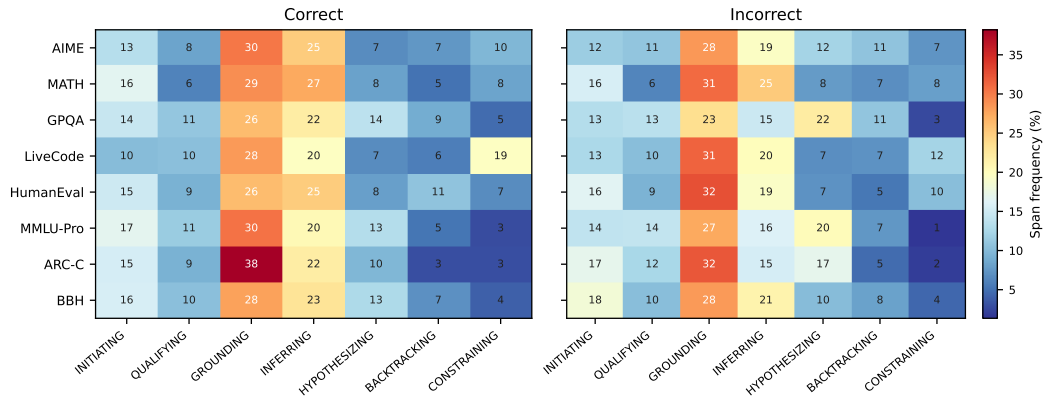


Figure 7: Operator frequency (%) per benchmark, split by correctness. Left: correct traces. Right: incorrect traces.

## 609 G Temporal Operator Profiles

610 Figure 8 shows operator density as a function of normalized trace position. GROUNDING dominates  
 611 the early trace; INFERRING concentrates near the end. BACKTRACKING shows a midtrace spike in  
 612 incorrect traces—a signature of failed exploration not followed by recovery.

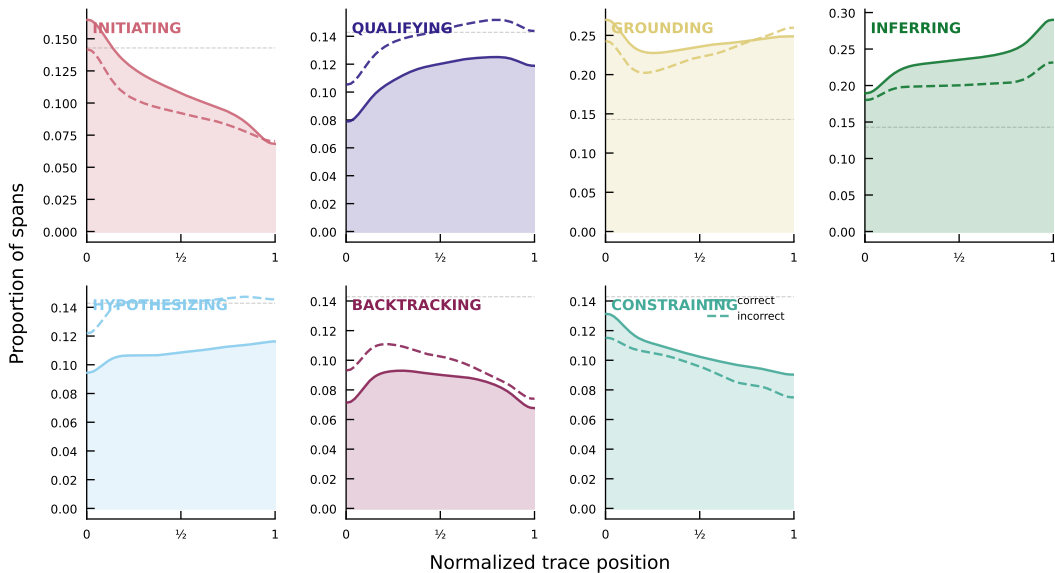


Figure 8: Operator density as a function of normalized trace position (Gaussian KDE). Dashed reference = uniform (1/7). Correct: solid; incorrect: dashed.

613 **H Backtracking analysis details**

614 We classify each BACKTRACKING and HYPOTHESIZING event on GPQA as LOCAL (re-checks  
 615 a single calculation; method unchanged), SUB-PROBLEM (re-opens a specific case or branch), or  
 616 GLOBAL (abandons the current solution strategy), using Claude Sonnet 4.6 as a judge.

617 **Sample.** Events are drawn from GPQA traces with stratification over (model, correctness), yield-  
 618 ing  $n=1,294$  BACKTRACKING and  $n=192$  HYPOTHESIZING events.

619 **Prompt.** The judge sees the GPQA problem and the full reasoning trace with the target event  
 620 marked: a contiguous region from the target operator span to the next span of the same operator class  
 621 (or trace end). The system rubric is reproduced below; the judge returns a JSON object {"scope":  
 622 ..., "rationale": ...}. See prompt below for more information.

Operator	Local	Sub-Problem	Global	$n$
BACKTRACKING	85.4%	12.8%	1.6%	1,294
HYPOTHESIZING	80.2%	18.2%	1.6%	192

Table 7: Scope classification for BACKTRACKING and HYPOTHESIZING events on GPQA.

**Prompt: Backtracking analysis judge prompt**

You are analyzing reasoning traces from large language models solving GPQA (graduate-level science) questions. Your task is to classify the SCOPE of a single REVISION-ADJACENT region marked in the trace—a contiguous section where the model re-checks, considers an alternative, hypothesizes, qualifies, or hesitates.

The marked region (between >>>MARKER<<< and <<<END>>> markers) starts at the target operator event and ends just before the next event of the SAME operator type (or at the end of the trace). Read the surrounding context to decide what the model is actually re-thinking inside the marked region.

**CLASSIFY THE SCOPE:**

**LOCAL:** The model is re-checking a single calculation, fact lookup, or specific claim (e.g., “wait,  $3 \times 7 = 21$  not 18”; “actually that should be oxygen”; “let me recompute that integral”; “I think this is the methyl group”). The OVERALL APPROACH is unchanged—same method, same plan, just verifying or fixing one step. Includes self-confirmation moves and tentative single-fact hypotheses.

**SUB\_PROBLEM:** The model is re-opening or alternating among specific cases, branches, or sub-questions within the larger problem (e.g., “let me reconsider case 2”; “what about the limit where  $x \rightarrow 0$ ?”; “could this be the (R) instead of (S) configuration?”). Strategy preserved but a discrete piece is redone or swapped.

**GLOBAL:** The model abandons or proposes to replace its current solution strategy/method with a fundamentally different approach (e.g., “this Lagrangian approach isn’t working, let me try energy conservation”; “I shouldn’t use the ideal-gas assumption—let me redo with van der Waals”; “I had the wrong mechanism entirely”; “let me use Gauss’s law instead of Coulomb”). The high-level plan changes.

**Output format:** a single JSON object on one line:

`{"scope": "LOCAL"|"SUB_PROBLEM"|"GLOBAL", "rationale": "<one short sentence>"}`

Be decisive. Compare what the model is doing in the spans IMMEDIATELY before vs. after the marked region—that local context tells you what is actually changing. Judge each case on its merits without a default fallback category.

-----  
 Prompt 1: System prompt used to classify the scope of BACKTRACKING and HYPOTHESIZING spans in GPQA reasoning traces in backtracking analysis.

623

624 **I Per-Dataset Baseline Comparison**

625 Figure 9 shows WP-AUC (cross-dataset protocol) across all six baselines for each of the 6 bench-  
 626 marks. BIG-Bench Hard and HumanEval are excluded (see §3). Op-XGB and OST dominate on  
 627 math-heavy and code benchmarks; the LLM self-judge is near chance on all datasets.

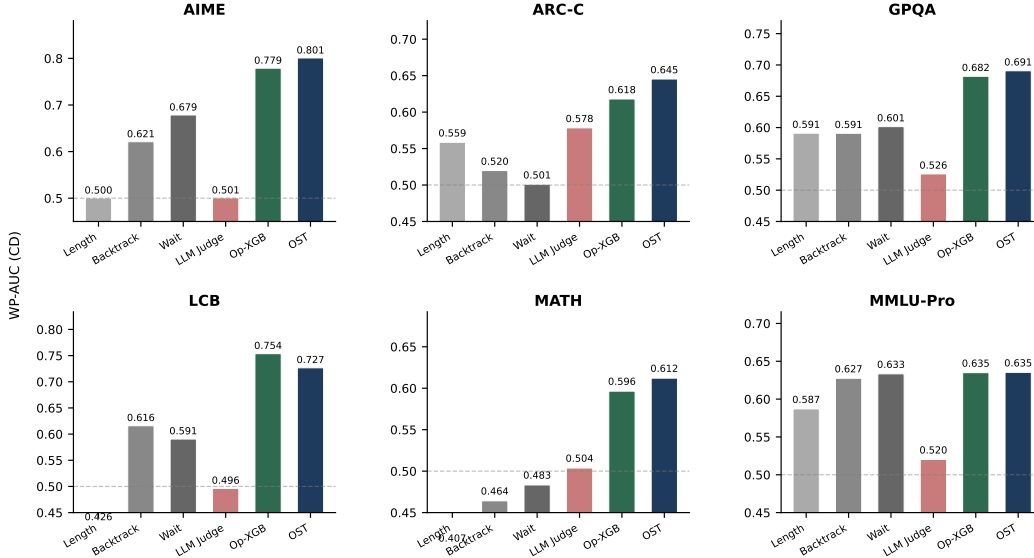


Figure 9: Per-benchmark WP-AUC (cross-dataset CV) for all six methods. Gray bars: span-free baselines (Length, Backtrack, Wait). Red bar: LLM self-judge. Green bar: Op-XGB (operator n-gram counts + anchor-phrase TF-IDF; XGBoost). Blue bar: OST (Operator Sequence Transformer).

628 **J Labeled Trace Example**

629 Table 8 shows a representative labeled trace from R1-Distill-LLaMA-70B on a MATH intermediate  
 630 algebra problem (incorrect). Each span begins at a pivot sentence (anchor, upright font) followed by  
 631 continuation sentences in *italics*.

632 **K Evaluation**

633 **K.1 Stability Analysis**

634 **Methodology.** To assess whether the induced operator labels reflect stable semantic structure  
 635 rather than arbitrary exemplar choice, we held the  $K = 7$  cluster assignments fixed and re-ran  
 636 the LLM naming step (claude-sonnet-4-6,  $T = 0$ ; same prompt as the discovery pipeline—  
 637 top-12 anchor  $n$ -grams plus 5 example spans per cluster) across 30 exemplar seeds. The top  $n$ -  
 638 grams are deterministic per cluster; only the 5 exemplars vary by seed, isolating naming stability  
 639 from the upstream clustering. Each of the  $30 \times 7 = 210$  resulting descriptions was embedded  
 640 once with text-embedding-3-small (1536-d,  $L_2$ -normalized), and we computed pairwise cosine  
 641 similarity as the inner product. We then formed three pair distributions: a *within-real* distribu-  
 642 tion (all  $\binom{30}{2} = 435$  cross-seed pairs per cluster, 3,045 pairs pooled across the seven operators);  
 643 an *across-real baseline* of 2,000 random pairs  $(i, j)$  with  $\text{cluster}(i) \neq \text{cluster}(j)$  from the same  
 644 description pool; and a *random-group control* generated by the identical naming protocol applied  
 645 to seven disjoint “fake clusters” of 2,000 spans each, drawn uniformly from the full type-A span  
 646 pool—these mix every real operator in corpus-rate proportions, so their top  $n$ -grams concentrate on  
 647 high-frequency meta-phrases (“now we need”, “let me think”) with no operator-specific signal.

648 **Statistical comparisons.** We test whether the within-real distribution stochastically dominates  
 649 each of the two reference distributions using a one-sided Mann-Whitney  $U$  test (alternative =

Table 8: Labeled trace with full span content. Model: R1-Distill-LLaMA-70B. Dataset: MATH. Correct: False.

<i>Problem:</i> Find all integer roots of $x^4+5x^3+9x^2-x-14=0$ .	
Operator	Span content
INFERRING	Okay, so I need to solve this equation and find all the integer roots. <i>And the task is to find all the integer roots.</i>
HYPOTHESIZING	Hmm, I remember that for polynomials the Rational Root Theorem can be helpful.
INITIATING	Let me try to recall how that works. <i>The Rational Root Theorem states that any rational root <math>p/q</math> has <math>p \mid a_0</math> and <math>q \mid a_n</math>.</i>
HYPOTHESIZING	Therefore, the possible integer roots are the factors of $-14$ .
INITIATING	Let me start with $x=1$ and $x=-1$ first.
GROUNDING	Substituting: $1+5+9-1-14=0$ . <i>So <math>x=1</math> is a root.</i>
INFERRING	So $x=-2$ is another integer root.
BACKTRACKING	Wait, but let me double-check my calculations because sometimes I might make an arithmetic mistake.
GROUNDING	Starting with $x=1$ : $1+5+9-1-14=0$ . Correct.
BACKTRACKING	But wait, the equation is a quartic, so it should have four roots in total.
QUALIFYING	But according to our tests, the other candidates didn't work, so perhaps the other roots are non-integer.
HYPOTHESIZING	Alternatively, maybe I missed something, so let me try another approach.
CONSTRAINING	Now, we have factored the polynomial as $(x-1)(x^3+6x^2+15x+14)$ .
QUALIFYING	Therefore, the only integer roots are 1 and $-2$ .

650 greater). For interpretability we report the rank-biserial effect size  $r_{\text{rb}} = 2U/(n_1n_2) - 1$ , which  
651 rescales  $U$  to  $[-1, +1]$  and equals  $2P(A > B) - 1$ : with sample sizes in the thousands the  $U$ -test  
652 rejects for trivial gaps, so  $r_{\text{rb}}$  is the load-bearing summary of *how much* the distributions are sepa-  
653 rated. The 435 within-cluster pair cosines per cluster are not statistically independent—they share  
654 the 30 underlying seed embeddings—which makes the  $U$ -test  $p$ -values anticonservative;  $r_{\text{rb}}$  does  
655 not depend on this independence assumption and is therefore reported as the primary result.

656 **Results and interpretation.** Figure 10 shows the three cosine-similarity distributions. The  
657 within-real pooled median cosine is 0.660 versus 0.467 for the across-real baseline and 0.565  
658 for the random-group control. Real-cluster names dominate the across-cluster baseline with  
659  $r_{\text{rb}} = +0.72$  and the random-group control with  $r_{\text{rb}} = +0.39$  (one-sided  $U$ -test  $p \approx 0$   
660 and  $p < 10^{-149}$ , respectively): roughly 86% of within-real pairs are more similar than a ran-  
661 dom across-cluster pair, and 69% more similar than a random within-fake-group pair. Every  
662 individual operator clears the across-cluster baseline (per-cluster  $U$ -test  $p < 10^{-12}$ ;  $r_{\text{rb}} \in$   
663  $[+0.22$  (GROUNDING, weakest),  $+1.00$  (BACKTRACKING, strongest)]. The random-group con-  
664 trol falls midway between within-real and across-cluster ( $r_{\text{rb}} = +0.44$  of random vs. across; median  
665 gap = 0.098): LLM descriptions of mixed-operator span pools are somewhat self-consistent, but far  
666 below the within-real ceiling, confirming that the coherence of real operator names is content-driven  
667 rather than a generic tendency of the naming prompt to produce similar-sounding sentences.

## 668 L Additional Analyses

### 669 L.1 Transition Structure and Run-Length Statistics

670 Figure 11 shows the operator transition probability matrix and run-length statistics. GROUNDING  
671 has the longest self-runs, reflecting sustained fact-anchoring that persists for multiple consecutive  
672 spans. BACKTRACKING is nearly always a single span—it interrupts but does not persist—which is  
673 consistent with its function as a brief error-signal rather than a sustained strategy change. The HY-  
674 POTHEISIZING  $\rightarrow$  INFERRING pathway has elevated probability relative to most other off-diagonal

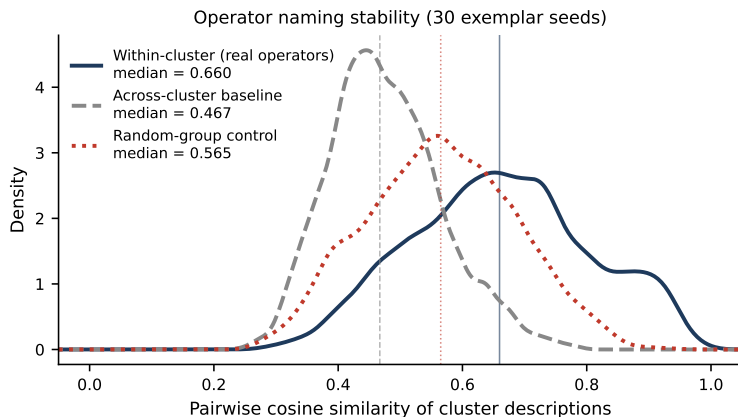


Figure 10: Naming stability: pairwise cosine similarity distributions between cluster-description embeddings generated under 30 exemplar seeds (`text-embedding-3-small`,  $L_2$ -normalized). The *within-cluster* distribution (median = 0.660) is well separated from both the *across-cluster* baseline (median = 0.467) and the *random-group control* (median = 0.565), which applies the same naming protocol to seven disjoint random span groupings. The gap confirms that operator names are stable at the level of meaning and are not an artifact of the LLM prompt.

675 entries, consistent with a generate-and-test schema in which tentative scenarios are followed by  
 676 conclusions.

## 677 L.2 Model Identification Confusion Matrix

678 Figure 12 shows the confusion matrix for the 12-way model identification task under problem-level  
 679 5-fold CV using the full Op-XGB feature recipe (Appendix E). The classifier reaches accuracy  
 680 = 79.9% and macro-AUC = 0.987. Errors concentrate within model families. The R1-Distill-  
 681 Llama-70B / R1-Distill-Qwen-32B pair contributes the largest off-diagonal mass: R1-Distill-Llama-  
 682 70B traces are predicted as R1-Distill-Qwen-32B at 0.41 and the reverse at 0.29, consistent with  
 683 their shared distillation lineage from DeepSeek-R1. The Claude-Haiku-4.5 / Claude-Sonnet-4.5 pair  
 684 shows mutual confusion in the 0.21–0.26 range, and Kimi-K2.5 is predicted as Kimi-K2 at 0.10  
 685 (with the reverse at 0.05). Off-diagonal mass outside these family clusters is small. Per-model  
 686 AUCs span 0.967 (R1-Distill-Llama-70B) to 0.997 (Kimi-K2, Kimi-K2.5, Grok-3-Mini, all tied at  
 687 this rounding).

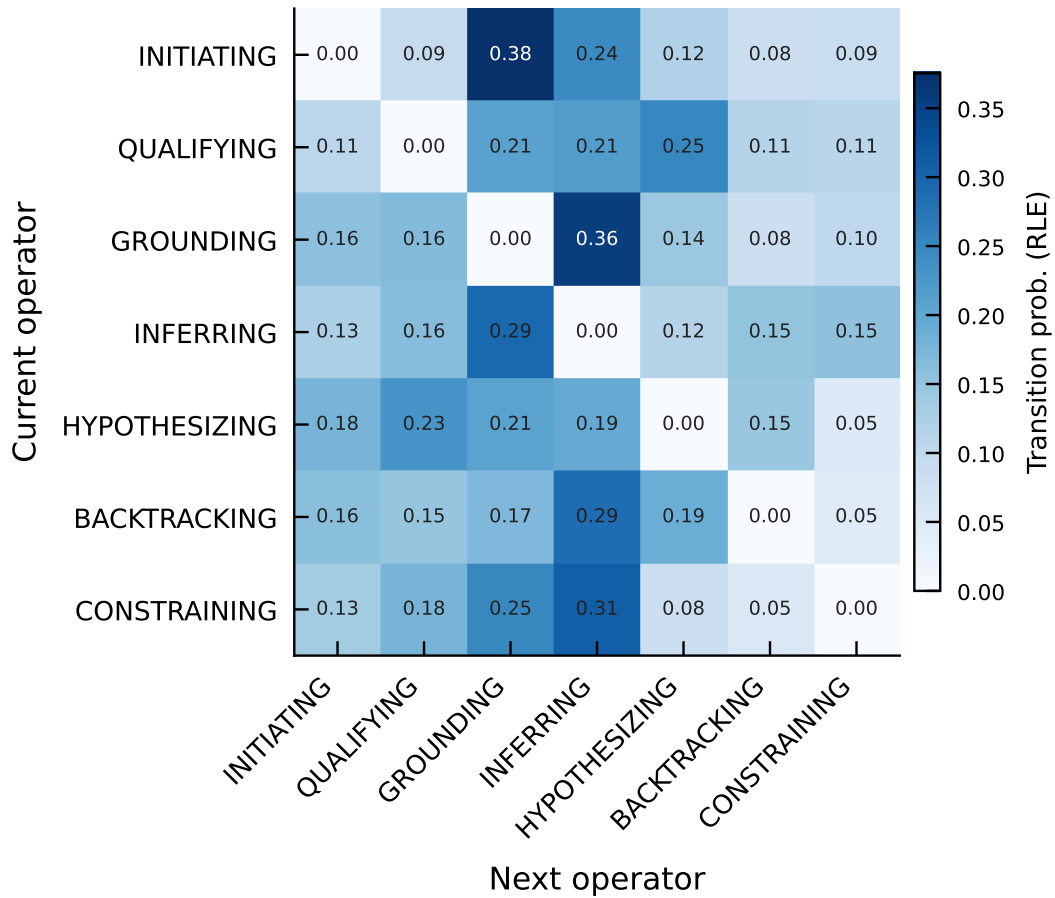


Figure 11: Left: operator transition probability matrix. Entry  $(i, j)$  = probability that operator  $j$  follows  $i$ . Right: run-length statistics per operator—mean and distribution of consecutive same-operator span counts. GROUNDING has the longest self-runs (sustained fact-anchoring); BACKTRACKING is nearly always a single span.

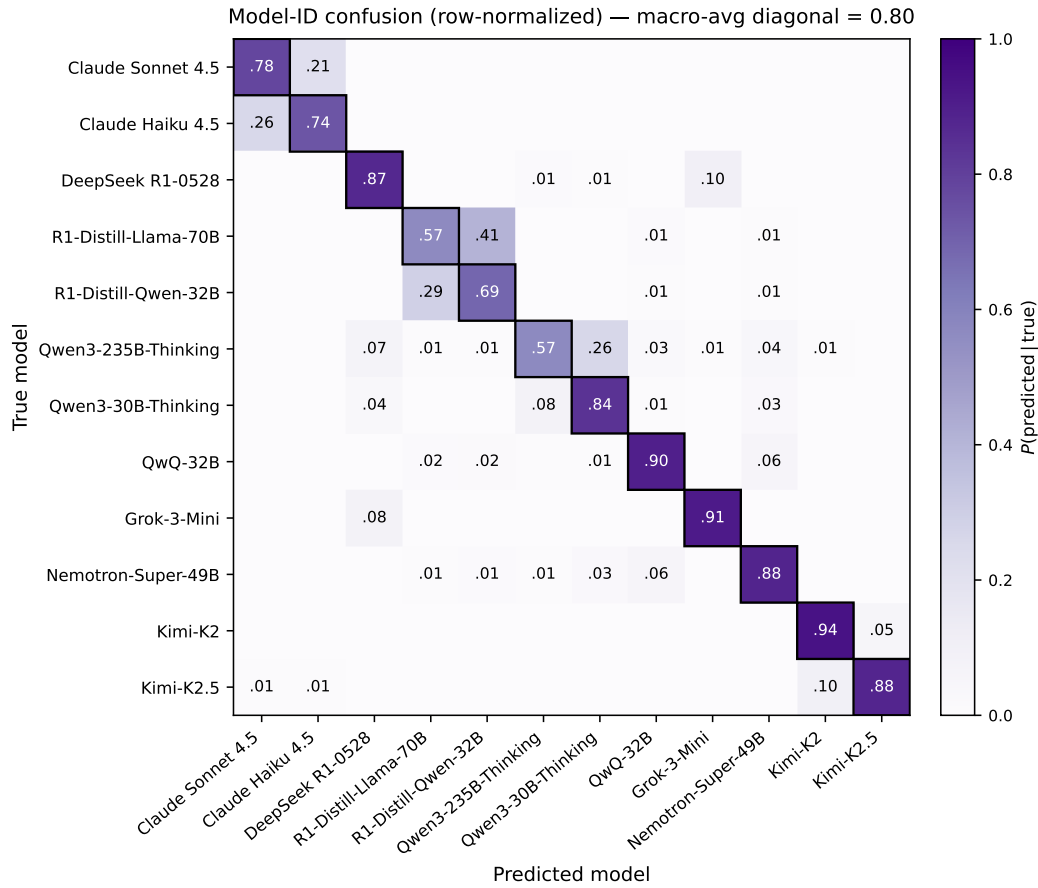


Figure 12: Model identification confusion matrix (row-normalized) under problem-level 5-fold CV with the full Op-XGB feature recipe (Appendix E). Accuracy = 0.799, macro-AUC = 0.987. Errors concentrate within model families (R1-Distill pair, Claude pair, Kimi pair).

688 **Appendix References**

- 689 [S1] Qwen Team. QwQ-32B: Embracing the power of reinforcement learning. [https://qwenlm.](https://qwenlm.github.io/blog/qwq-32b/)  
690 [github.io/blog/qwq-32b/](https://qwenlm.github.io/blog/qwq-32b/), mar 2025. Blog post, March 6, 2025.
- 691 [S2] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,  
692 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*  
693 *arXiv:2505.09388*, 2025.
- 694 [S3] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao  
695 Zhu, et al. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement  
696 learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 697 [S4] Anthropic. Introducing Claude 4. <https://www.anthropic.com/news/claude-4>, May  
698 2025. Blog post, May 22, 2025.
- 699 [S5] xAI. Grok 3 beta — the age of reasoning agents. <https://x.ai/news/grok-3>, February  
700 2025. Blog post, February 19, 2025.
- 701 [S6] Akhiad Bercovich, Itay Levy, Izik Golan, Mohammad Dabbah, Ran El-Yaniv, Omri Puny,  
702 et al. Llama-nemotron: Efficient reasoning models. *arXiv preprint arXiv:2505.00949*, 2025.
- 703 [S7] Kimi Team. Kimi k2: Open agentic intelligence. *arXiv preprint arXiv:2507.20534*, 2025.
- 704 [S8] Kimi Team, Tongtong Bai, Yifan Bai, Yiping Bao, SH Cai, Yuan Cao, Y Charles, HS Che,  
705 Cheng Chen, Guanduo Chen, et al. Kimi k2.5: Visual agentic intelligence. *arXiv preprint*  
706 *arXiv:2602.02276*, 2026.