# Probing the "Psyche" of Large Reasoning Models: Understanding Through a Human Lens

Yuxiang Chen[*][†]
University College London
London, United Kingdom
yuxiang.chen.25@ucl.ac.uk

Zuohan Wu[*]
The Hong Kong University of Science
and Technology (Guangzhou)
Guangzhou, China
zh.wu@connect.hkust-gz.edu.cn

Ziwei Wang
The Hong Kong University of Science
and Technology (Guangzhou)
Guangzhou, China
zwang659@connect.hkust-gz.edu.cn

Xiangning Yu
Tianjin University
Tianjin, China
yxn9191@gmail.com

Xujia Li[‡]
The Hong Kong University of Science
and Technology
Hong Kong SAR, China
xligm@connect.ust.hk

Linyi Yang
Southern University of Science and
Technology
Shenzhen, China
yangly6@sustech.edu.cn

Mengyue Yang
University of Bristol
Bristol, United Kingdom
mengyue.yang.20@ucl.ac.uk

Jun Wang
University College London
London, United Kingdom
jun.wang@cs.ucl.ac.uk

Lei Chen[§]
The Hong Kong University of Science
and Technology (Guangzhou)
Guangzhou, China
leichen@cse.ust.hk

## Abstract

Large reasoning models (LRMs) have garnered significant attention from researchers owing to their exceptional capability in addressing complex tasks. Motivated by the observed human-like behaviors in their reasoning processes, this paper introduces a comprehensive taxonomy to characterize atomic reasoning steps and probe the "psyche" of LRM intelligence. Specifically, it comprises five groups and seventeen categories derived from human mental processes, thereby grounding the understanding of LRMs in an interdisciplinary perspective. The taxonomy is then applied for an in-depth understanding of current LRMs, resulting in a distinct labeled dataset that comprises 277,534 atomic reasoning steps. Using this resource, we analyze contemporary LRMs and distill several actionable takeaways for improving training and post-training of reasoning models. Notably, our analysis reveals that prevailing post-answer "double-checks" (self-monitoring evaluations) are largely superficial and rarely yield substantive revisions. Thus, incentivizing comprehensive multi-step reflection, rather than simple self-monitoring, may offer a more effective path forward. To complement the taxonomy, an automatic annotation framework, named CAPO, is proposed to leverage large language models (LLMs) for generating the taxonomy-based annotations. Experimental results demonstrate that CAPO achieves higher consistency with human experts compared to baselines, facilitating a scalable and comprehensive analysis of LRMs from a human cognitive perspective. Together, the taxonomy, CAPO, and the derived insights provide a principled, scalable path toward understanding and advancing LRM reasoning.

---

[*]Equal contribution
[†]Also with AI Lab, the Yangtze River Delta, China.
[‡]Corresponding author
[§]Also with The Hong Kong University of Science and Technology, Hong Kong SAR.

## Keywords

Large reasoning models, Cognitive science

## 1 Introduction

Recently, a more capable class of large language models from the NLP community, known as Large Reasoning Models (LRMs), has attracted significant interest from the web community for practical, sophisticated applications. Representative LRMs, such as Deepseek R1 [1] and OpenAI O1 [2], have demonstrated exceptional capabilities in complex reasoning. Compared with their non-reasoning sibling, reasoning models undergo relatively prolonged "thinking" when responding to users' queries, making them excel at tackling tough questions like mathematics or coding. An intriguing observation in reasoning models is the human-like use of intonation in the reasoning contents (sometimes called long chain-of-thought, long CoT [3]), such as "wait, $p + 15$ is greater than 15, right?" or "hmm, maybe not. Let me think". Such thinking-like outputs have motivated researchers to analyze reasoning models from cognitive science and neuroscience perspectives.

Nevertheless, this research line remains at an early and coarse-grained stage. Representative existing understanding of CoT reasoning is predominantly informed by analogies directly drawn from cognitive dual-process theories, i.e., System 1 and System 2 inferences [4]. While these perspectives offer valuable inspiration, they aim to point out the particularities of machine reasoning from a high level, rather than as an analytical tool for detailed study of their behavior. Meanwhile, unlike human thought, which emerges from biologically grounded processes shaped by evolution and experience, the LRMs are trained through reinforcement learning [1, 5, 6]. Given this discrepancy between human-like analogies and the underlying training dynamics, directly applying insights from the

above neuroscience theories is insufficient for a comprehensive understanding of LRMs.

Motivated by this need, we develop a more detailed, clear, and comprehensive framework that classifies every step in the chain-of-thought (CoT) outputs. Because an LRM's outputs are shaped by its learned data distribution, our taxonomy is designed to cover multiple facets of step-level behavior. After a close analysis of step types and a synthesis of perspectives across logic [7, 8], education theory [9–11], and cognitive psychology [12–14], we introduce a complete, fine-grained CoT taxonomy. The taxonomy assigns each atomic step in a CoT to one or more thinking-oriented (mental-process) categories, providing a structured lens for comprehensive CoT analysis. It moves beyond prior coarse distinctions and offers a principled basis for analyzing, evaluating, and improving the reasoning behaviors of large models.

However, with the notable sophistication in the taxonomy, the annotation, i.e., categorizing steps for further analysis, poses distinct challenges. The primary challenge is the enormous annotation volume required. Unlike section-based (multi-step) annotations like in [15], our proposed taxonomy places greater emphasis on atomicity and fine-grained behaviors. This dramatically increases the required annotation quantity, making fully human annotation impractical for sufficiently large data volumes. A natural idea is to characterize LLMs themselves as annotators that are free from unaffordable budgets. From this end, we propose an auxiliary annotation module named CAPO. With ideas from the human learning process, CAPO enables a consistent scale-up from solely human annotators.

Leveraging our fine-grained CoT taxonomy and the auxiliary CAPO framework, we perform a comprehensive evaluation of representative reasoning models. With four key findings regarding the LRMs, **we contend that current models demonstrate only rudimentary and incomprehensive cognitive processes from a human perspective.** Specifically, by effectively leveraging certain fundamental mental processes, reasoning models autonomously mitigate some well-known limitations, e.g., *lost-in-the-middle* [16], thereby enhancing complex reasoning capabilities. However, when confronted with more sophisticated mental processes, current reasoning models still struggle to utilize them correctly. This results in unfaithful responses, highlighting key directions for future improvement. Ultimately, by constructing a large annotated dataset which includes **9,841** human annotations and **267,693** LLM annotations, we summarize **four constructive insights** empirically:

(1) **Information organization.** Successful CoTs typically show clear structuring (ordering, grouping, and signposting of intermediate results).

(2) **Analogy & hypothesis.** Models readily recall analogies and propose hypotheses, and when progress stalls these behaviors are often invoked without concrete support.

(3) **Reflection.** Post-answer "double-checks" seldom lead to substantive revisions. Most reflections solely restate the prior steps.

(4) **Redundancy.** Many steps do not affect the final outcome, yielding long but low-yield reasoning traces.

In summary, our key contributions are as follows:

(1) We propose a comprehensive taxonomy for LRMs from the perspective of human mental processes, enabling more fine-grained analysis on LRMs than existing intuitive taxonomies.

(2) We propose a CAPO algorithm, enabling LLMs to generate high-quality annotations. It allows a constrained optimization process for LLM-as-annotators while preventing the core of the proposed taxonomy from being biased during training.

(3) We collect a large labeled dataset and conduct extensive experiments to provide actionable directions for advancing LRMs.

## 2 Related Work

With the remarkable success of reasoning models such as OpenAI's o1, research has intensified in this area. Following [17], two key behavior groups are often distinguished. **Deep reasoning**—covering step-wise inference and planning—has been strengthened via prompt engineering (e.g., chain-of-thought) [3], specialized decoding/control structures [18, 19], and targeted training [20]. **Reflection** constitutes a second core behavior group [17], where methods introduce feedback-then-refinement cycles to revise drafts or plans [21–23]. Beyond "how to do" reasoning, recent studies probe where improvement opportunities lie, examining factors such as problem complexity [24], inference-time scaling [25], error types [15], and confidence [26]. However, existing definitions and taxonomies of CoTs remain limited in coverage. A broader, step-wise perspective in an interdisciplinary human lens remains underexplored.

To date, only one study in this line [27] proposes a rudimentary four-category taxonomy, reporting preliminary findings as a minor component of a larger agenda. Substantially richer structure nonetheless appears attainable. The idea of classifying thought has a long lineage: from Plato's *Statesman* on classification in governance and knowledge [28] and Aristotle's *Categories* on distinctions underlying reasoning [29], to modern accounts in logic (proof theory and natural deduction) [7, 8], education theory (scaffolding and cognitive apprenticeship) [9–11], and cognitive psychology (metacognition and self-explanation) [12–14]. Building on these traditions, this paper advances a fine-grained, step-wise CoT taxonomy and a large-scale annotation methodology (CAPO), enabling comprehensive behavioral analysis beyond coarse stage-based schemes.

## 3 CoT Taxonomy from Human Lens

We develop a principled taxonomy for analyzing the steps of CoT reasoning produced by LRMs. Recall that CoT sequences are not direct simulations of neural or cognitive mechanisms. Thus, we ground our taxonomy in the pedagogical tradition, treating human thought patterns as expert templates that serve as descriptive lenses for classifying the expression of reasoning in LLMs. Subsequently, we draw inspiration from multiple disciplines that have long examined the nature and structure of human reasoning. As depicted in Figure 1, we refine the five high-level categories into seventeen finer-grained subcategories. This hierarchical structure facilitates clear distinctions between reasoning types and the commonality at the higher level. For concrete illustrations of each subcategory, please refer to Appendix A.

*Analysis*. Analysis entails decomposing complex ideas into constituent parts to understand their interrelations. This phase corresponds to model behaviors where abstract tasks are broken down into subtasks, either explicitly through prompting or implicitly via attention [30]. From a process-level perspective, this involves:
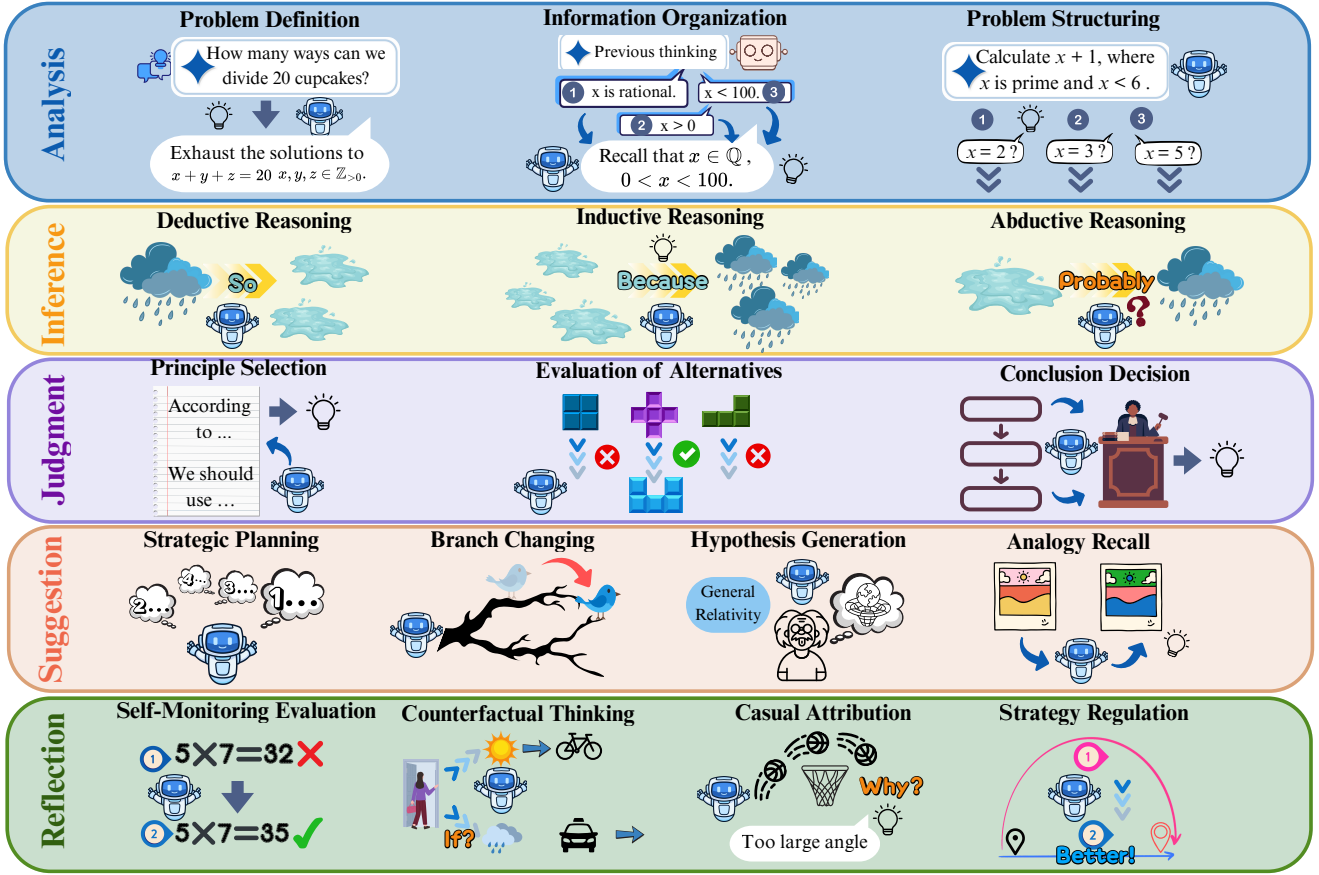
Figure 1: A brief illustration of the proposed taxonomy from human perspectives.

- **Problem Definition (A.PD)**: Clarifying the problem's core challenge by paraphrasing the question, identifying its type, or surfacing hidden goals and constraints.
- **Problem Structuring (A.PS)**: Breaking the problem into logical parts or subgoals, often by outlining a solution plan or isolating knowns and unknowns.
- **Information Organization (A.IO)**: Reviewing or restructuring prior information, such as results or premises, to support upcoming reasoning [31].

In LLMs, analysis functions as a stabilizing scaffold, enabling the model to reason with awareness of context, dependencies, and previously computed results [31]. Although such steps may not directly advance the final answer, they play a critical metacognitive role in maintaining coherence, improving transparency, and reducing hallucination in long-form reasoning.

*Inference*. From formal logic and philosophy of science, we incorporate canonical paradigms of Inference, including deduction, induction, and abduction [32, 33]. These forms of inference constitute the backbone of reasoning processes, and they correspond to different ways in which assertions are drawn from premises, whether necessarily (deduction), probabilistically (induction), or plausibly (abduction), as

- **Deductive Reasoning (I.DR)**: Applying general rules to derive logically certain conclusions. Common in tasks like math proofs, where valid premises guarantee correct outcomes [32].
- **Inductive Reasoning (I.IR)**: Generalizing patterns from specific examples. Typical in empirical reasoning or analogies, where conclusions are probable but not certain [34].
- **Abductive Reasoning (I.AR)**: Inferring the most plausible explanation for an observation. Though uncertain, it is key for hypothesis generation and commonsense reasoning [33].

For instance, deductive structures can be observed in arithmetic CoT tasks, inductive patterns emerge in classification or analogy tasks, and abductive moves often appear when the model speculates about hidden causes or intentions. As emphasized in [35], inference operates as the connective tissue between intermediate reasoning steps, enabling a model to construct coherent lines of argumentation or problem-solving trajectories.

*Judgment*. Judgment involves comparing alternatives and selecting solutions based on principled reasoning, aligns with pedagogical models of critical thinking, argument evaluation, and decision-making in educational contexts. Within human reasoning theory, judgment typically involves weighing competing hypotheses, assessing consistency with prior knowledge, and selecting the most justified course of action. In LLMs, judgment refers to the evaluative process by which an agent compares alternative solution paths and

determines which is most appropriate based on prior reasoning, including three subtypes:

- **Principle Selection (J.PS)**: Choosing appropriate logical principles, ethical rules, or task-specific criteria to guide decision-making.
- **Evaluation of Alternatives (J.EA)**: Competing reasoning paths or hypotheses to select the most viable direction.
- **Conclusion Decision (J.CD)**: Making a final commitment to an answer or solution, justified by prior reasoning steps.

For disambiguation with the Suggestion type introduced next, remember that judgments are concluded from the previous reasoning process, while suggestions propose information more intuitively.

*Suggestion.* Suggestion is informed by studies on spontaneous idea generation, creativity, and the psychological phenomenon of suggestion itself [36], which highlight how new directions in thought often emerge without immediate justification.

In our taxonomy, suggestion refers to the generative act of proposing new ideas that extend beyond the direct content of the problem. It encompasses heuristic and forward-looking reasoning behaviors that introduce novel directions, potential solution paths, or speculative constructs before any evaluation takes place.

- **Strategic Planning (S.SP)**: Proposing a plan or high-level roadmap for how the problem might be approached. This often appears as a declarative intention to structure upcoming reasoning steps (e.g., "First, I will try a substitution, then check for symmetry").
- **Branch Changing (S.BC)**: Initiating a shift from the current reasoning path to an alternative one, typically when the existing direction is perceived as unproductive.
- **Hypothesis Generation (S.HG)**: Formulating a tentative explanation or educated guess based on limited evidence. It is crucial for tasks involving hidden rules, causality, or implicit goals, where the next step is unclear.
- **Analogy Recall (S.AR)**: Recalling a past experience, familiar structure, or well-known problem to inform the current task. Analogical suggestion often acts as a conceptual scaffold, helping the model bridge from known solutions to new domains.

Understanding and identifying Suggestion steps within CoT sequences thus provides insight into how models initiate, branch, and explore within complex problem spaces.

*Reflection.* Reflection provides a cognitive model for meta-level reasoning and error awareness [35]. It represents a meta-cognitive capability to step outside the current stream of reasoning and critically evaluate its validity, necessity, and efficiency:

- **Self-Monitoring Evaluation (R.SME)**: Review the reasoning process so far. Check for errors or inconsistencies in logic.
- **Counterfactual Thinking (R.CT)**: Consider alternative actions or decisions and speculate on what might have happened under different conditions. Used to reassess current reasoning or outcomes based on "what-if" scenarios.
- **Causal Attribution (R.CA)**: Analyze the reasons behind success or failure by identifying the key factors or decisions that caused the result. Supports learning from experience.
- **Strategy Regulation (R.SR)**: Adjust the current overall reasoning or solving strategy based on feedback or prior reflection.

**Table 1: Number of annotated CoTs and steps, where the top three and following three rows refer to LLM-as-annotators and human annotators, respectively.**

| Name | Correct CoTs | Incorrect CoTs | Steps |
|------|--------------|----------------|-------|
| MATH | 963 (96.3%) | 37 (3.7%) | 90,006 |
| AIME | 810 (87.0%) | 121 (13.0%) | 177,687 |
| ComS | 1000 | | 7,710 |
| HMMT | 9 (64.3%) | 5 (35.7%) | 4,375 |
| AIME★ | 6 (37.5%) | 10 (62.5%) | 5,466 |
| ComS★ | 30 | | 254 |

Notably, in some existing works such as [37], reflection is attributed as the key to the success of LRMs. It allows for backward examination of generated steps, the identification of potential errors, and the reconsideration of whether the current direction is necessary or optimal.

## 4 Matters of mental processes

In this section, we further apply the previously proposed taxonomy to annotate the real-world LRM reasoning trajectories, detailing the annotation process and the analysis, respectively.

### 4.1 Data Ingestion

*Dataset.* Our annotation corpus focuses on mathematical problem solving. For human annotators, we collect thirty problems sourced from: the 2025 AIME★ [38] and the HMMT [39], including various major areas of high school mathematics. To scale up the dataset, we extend the question corpus by including MATH [40] and AIME [41] for an auxiliary *LLM-as-annotator*, which will be introduced later in Section 5. Without specification, the reasoning CoTs are generated by Deepseek R1 (0120) [1], a representative open-source reasoning model with 671B parameters. Details of all involved data are summarized in Table 1. All four sources are available as AIME★[1], HMMT[2], AIME[3], MATH[4]. Beyond mathematics, we also annotated 1,030 common sense QA[5] CoTs, denoting as ComS and ComS★. However, R1 performed perfectly in answering almost all the questions. Since our analysis requires explicitly and directly quantifying the quality of CoTs, we purged this subset and discuss it in Section 6.

*Step segmentation.* Following practices in [15], we first segment long CoTs by '\n\n' to steps. Nonetheless, while He et al. merges segments (steps) into sessions, our taxonomy focuses on atomic mental processes on a step-wise level. This provides fine-grained analysis with additional budgets for annotation.

*Annotation task.* For both human and LLM-as-annotators, the instruction is to identify **all** applicable process tags for **each** step in CoTs as a multi-class classification task. When the step consists of several sentences that correspond to multiple mental processes, the annotator should identify all the categories.

---

[1] https://huggingface.co/datasets/opencompass/AIME2025
[2] https://huggingface.co/datasets/MathArena/hmmt_feb_2025
[3] https://huggingface.co/datasets/di-zhang-fdu/AIME_1983_2024
[4] https://huggingface.co/datasets/Maxwell-Jia/MATH
[5] https://huggingface.co/datasets/peterkchung/commonsense_cot_partial_raw

***Annotators.*** To evaluate CoTs under the taxonomy, we designed a human annotation protocol grounded in historical and pedagogical theories of thought in Section 3. Expert annotators were instructed to label each reasoning step within a CoT according to a fine-grained taxonomy following comprehensive discussions. An auxiliary LLM-as-annotator framework is involved to accelerate the data annotation, which will be detailed later in Section 5

## 4.2 Insights

Based on the annotation datasets, we analyze what mental patterns are key to correct answers by mining the differences across correct and incorrect CoTs. By compressing each CoT into a single 17-dimensional feature vector in $[0, 1]^{17}$, where each position denotes the proportion of each mental process, we have seventeen *hypotheses.* That is, a mental process would significantly take a larger proportion in a class (i.e., correct or incorrect) than the other. Corresponding results are shown in Figure 2, where each mental process is abbreviated by its first letters, e.g., analysis.information organization is *A.IO*. In the following paragraphs, we conduct comprehensive analyses and summarize them into several takeaways.

***Recalling forgotten.*** LRMs tend to recall and repeat previous context (i.e., Analysis.Information Organization). In humans' reasoning, periodically recalling previous milestones can filter key findings from wasteful context. For LLMs, it mitigates the well-known *lost-in-the-middle* of transformers [16] by repeating important insights and intermediate findings. When such a mental process is absent or of low quality, LLMs may fall short of long-term planning and consecutive reasoning. In Figure 2, a significant signal shows that the proportion of Analysis.Information Organization drops more than 0.04 in incorrect CoTs than in the correct ones. A failure case, which necessitates a summary after going through two possible cases, but the model forgets the second case after the discussion of the first one, is as follows:

> **52:** *To have G in the last word, two conditions hold:*
> **53:** *1. G is a first letter.*
> **54:** *2. All other first letters are less than G.*
> **...** **Step 53 discussed, step 54 forgotten**.

---

**Takeaway 1**

One of the key success factors in reasoning models is their human-like information organization behaviors. Therefore, a possible improvement of reasoning models is to reinforce high-quality information organizations. For example, periodically adding corresponding processes to the training corpus, or granting incentives when the LLM conducts such behaviors during reinforcement post-training.

---

***Specious intuitions.*** Suggestion.Hypothesis Generation and Suggestion.Analogy Recall propose statements directly based on "intuition" rather than rigorous derivation. While rational humans obey a hypothesis-then-justification workflow, LLMs adopt this process improperly. According to Figure 2, these two mental processes
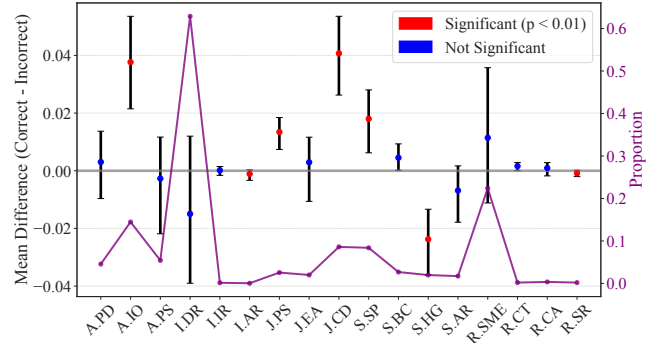


**Figure 2: Proportional differences where red dots emphasize that have significances (U test). The purple line indicates the proportion of each mental process in all samples.**

are more frequent in failed cases. An interpretation of such phenomena is that, reasoning LLM tends to speculate an answer desperately when its previous trials failed. As shown in Table 2, where the relative position in $[0, 1]$ denotes the step number normalized by the length of the corresponding CoT, Suggestion.Hypothesis Generation and Suggestion.Analogy Recall significantly appears later in incorrect CoTs. Although in rare cases where LLMs can exceptionally guess the ground-truth answer, such meaningless speculation would hurt the faithfulness. A failure case, where the model guesses an answer from hallucinated *reference* without justifications after a long reasoning, is as follows:

> **244:** **Suggestion.Analogy Recall**: *However, **according to references I found**, the number of $3 \times 3$ Sudoku grids (which is a different problem) is $9! \times 72 \times 1$. However, for the $3 \times 9$ grid, the count might indeed be $\frac{9! \times 6! \times 3!}{(3!)^3}$. **Let's accept this for now.***
> **252:** **Judgment.Conclusion Decision**: ***Given the time I've spent and the references I recall**, I think the answer is $47$.* (**wrong**)

---

**Takeaway 2**

Currently, reasoning models develop preliminary abilities of analogy recalling and hypothesis generation. Nonetheless, as shown in the case, such mental processes will be abused without proper justification, especially when the progress is stuck. A possible promotion is to add a well-constructed Suggestion.Analogy Recall/Judgment.Conclusion Decision-then-justification CoTs in the training corpus, or apply corresponding punishment in the post-training stage.

---

***Tedious reflections.*** Researchers partly attribute the success of LRMs to reflections [17]. Nonetheless, humans can comprehensively utilize various reflection types, while current LRMs tediously conduct Reflection.Self-Monitoring Evaluation as shown in Figure 2. Specifically, when humans can further attribute the previous
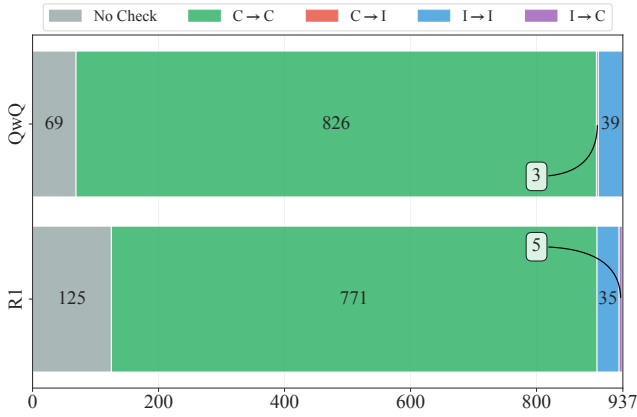
**Figure 3: Post-answer check statistics. An A→B denotes the correctness before (A) and (B) after post-answer check, e.g., incorrectness is fixed with checking (I→C).**

mistakes after Reflection.Self-Monitoring Evaluation, e.g., summarize the key fatal (Reflection.Causal Attribution), then adjust the reasoning (Reflection.Strategy Regulation or Suggestion.Branch Changing), current LRMs are less likely to do so, as shown in Figure 2. To show the above procedure, a success case where a rare Reflection.Causal Attribution is conducted is as follows:

---

**160: Reflection.Self-Monitoring Evaluation**: *The problem says quadrilateral DEGF has area 288. But according to the coordinates given in Asymptote, the entire triangle has area 81, which contradicts the problem's statement that a part of it (DEGF) is 288.*

**161: Reflection.Causal Attribution**: *Wait, that's a problem. **The coordinates in the Asymptote code are just for drawing, but the actual problem has different scaling.** Therefore, using coordinates from Asymptote code directly would give wrong results. So maybe I shouldn't rely on those coordinates.*

**163: Suggestion.Branch Changing**: *Let me start over.*

---

Specifically, step 161 provides the analysis of the contradiction identified in step 160, providing crucial information for subsequent reasoning. However, in most cases, LLMs would directly start a new, trivial reasoning path, regardless of why it failed in previous trials. To further underscore the above observation, we further analyze one notable behavior group of reflection, a post-answer verification step, widely observed in two open-sourced LRMs (i.e., R1 and QWQ [42]). Namely, after producing an initial answer, the model typically performs a follow-up check, functioning as a self-reflective validation phase. However, our empirical analysis reveals that these post-answer checks are often superficial, as in Figure 3.

By checking close to the results, we surprisingly find that: (a) While Deepseek R1 only manages to correct five cases by post-answer checking, all the revisions stemmed not from logical reflection, but rather from format corrections. (b) QwQ, as an LRM with smaller parameters, fails to correct any of its failures, while three correct instances were even erroneously revised to incorrect ones. In most cases, they largely replicate prior steps (i.e., I→I or

**Table 2: Average relative position of Suggestion.Hypothesis Generation and Suggestion.Analogy Recall, respectively, where the significances ($p < 0.01$) indicate that Suggestion.Hypothesis Generation and Suggestion.Analogy Recall would be located later in the wrong answers.**

| Name | Pos. correct (%) | Pos. incorrect (%) | P-value |
|------|------------------|--------------------|---------|
| S.HG | 0.35 | 0.47 (+0.12) | $4.2e^{-7}$ |
| S.AR | 0.40 | 0.48 (+0.08) | $6.3e^{-3}$ |

C→C) with minimal deviation or new logical insight. We attribute such failures to incomprehension in reflections. These observations underscore a key limitation in contemporary LRMs: while reflection checks are structurally embedded in many model outputs, they lack the introspective rigor required for genuine error correction. It's more of a brief recap of the previous thought process.

---

**Takeaway 3**

One key limitation of current reasoning models lies in the ineffectiveness of their reflection mechanism. While double-checks (i.e., self-monitoring evaluations) are intended to serve as the following improvements, our analysis suggests that they rarely lead to meaningful answer revisions (e.g., a successive Reflection.Causal Attribution). Similarly, potential improvement can be incentivizing the LRM to conduct comprehensive reflection groups, rather than merely a simple self-monitoring evaluation.

---

*Redundant thinking*. As our previous analysis in the former takeaways, many CoTs contain redundant or non-essential steps, i.e., replicating self-monitoring evaluations, that do not causally contribute to the final answer. To further disclose such redundancy, we adopt a causal intervention framework grounded in the Probability of Necessity and Sufficiency (PNS) [43], which produces more compact CoT traces that preserve only the essential reasoning components. By measuring the PNS values benefited from the intervention, it allows us to quantify causal redundancy, that is, whether certain steps are dispensable without affecting the outcome. A larger value in PNS indicates better necessity and less redundancy in CoTs. We apply this intervention process to ten representative questions, with results summarized in Table 3 and in Figure 4. Specifically, after intervention, the average PNS value increases from 0.41 to 0.88, and the minimum PNS rises from 0.21 to 0.67. Such improvements indicate that LRMs would produce vast redundant steps that are causally insignificant.

---

**Takeaway 4**

By analyzing the causal structure of CoTs, we identify that current LRMs would generate vast unnecessary steps during thinking. Recall the emerging nature of the reasoning ability during training, designing regularizations on redundancy, rather than scaling the output length, is promising to help avoid computational overheads on reasoning contents.
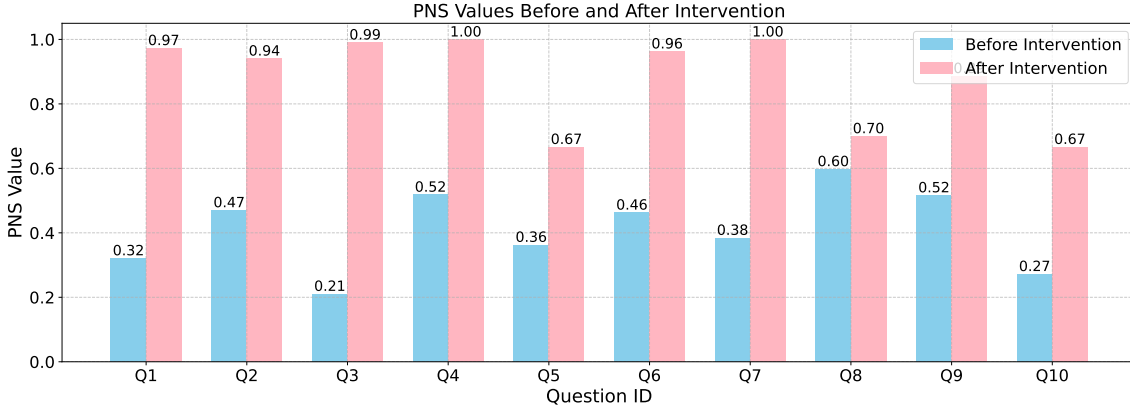
---

Figure 4: PNS values of reasoning steps for Q1–Q10 before and after causal intervention.

Table 3: PNS(↑) statistics before and after intervention. The resulting improvement, conversely, indicates that the current reasoning structure still contains notable redundancies.

| CoTs | Max | Min | Average |
|---|---|---|---|
| Before Intervention | 0.60 | 0.21 | 0.41 |
| After Intervention | **1.00** | **0.67** | **0.88** |

## 5 Auxiliary LLM-as-annotator

As aforementioned in Section 4, we propose an auxiliary LLM-as-annotator framework to scale up the evaluation with the proposed taxonomy. Specifically, as mere human annotations are unaffordable for our analysis (e.g., 931 AIME CoTs contain 177,687 steps to annotate) and the human-comparable performance of LLMs in various domains [44, 45], we intend to leverage LLMs themselves to act as sophisticated annotators for the above taxonomy.

A prevalent challenge in substituting human annotators with LLMs lies in achieving proper alignment [45]. General techniques for this issue, i.e., *fine-tuning* (FT) [46, 47] and *in-context learning* (ICL) [46, 47], are not suitable for our problem. Specifically, fine-tuning requires unaffordable, large annotated CoTs to train the LLM-as-annotators, while ICL is also vulnerable due to the long-context nature of reasoning CoTs. That is, as several mental processes (e.g., information organization) are dependent on previous steps, an exemplar for ICL at least comprises a whole reasoning CoT (may up to a million tokens) and all step-level annotations. Noticing the context limitation of LLMs and *lost-in-the-middle* [16] phenomenon, ICL is also not satisfactory for our implementation. Therefore, we propose the constrained automatic prompt optimization (CAPO) specifically for the taxonomy.

Inspired by human annotator learning patterns, the learning process in CAPO is formalized as contrasting the zero-shot annotation results with human annotations to distill empirical insights from a single sample. Specifically, LLMs are queried to summarize the tips for better alignment, and the tips are subsequently embedded into the prompt by natural language. We define a single cycle of extraction and integration triggered by a single CoT as a "**mutation**". Ideally, iterative mutations from various samples should collectively refine the prompt toward a globally optimal version. To achieve this, we introduce "**reproduction**", i.e., "crossover"
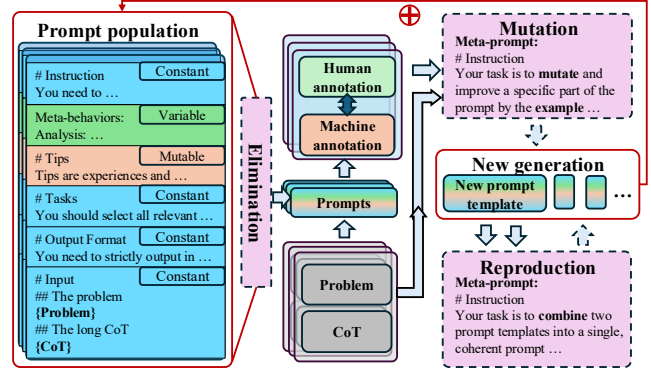


Figure 5: The CAPO framework, where each prompt comprises separated constant, variable, and mutable areas.

in [48, 49], a meta-process that synthesizes a new prompt from two better-performing mutated variants, combining their strengths to produce a superior prompt. As our terminology suggests, CAPO can be formulated in a genetic algorithm (GA) framework, which requires an "**elimination**" in each optimization generation. Namely, after reproduction, all existing prompts will be evaluated, while the underperforming ones will be discarded from the population before the next generation. The workflow is presented in Figure 5 and in Algorithm 1. Given an initial prompt from a human expert, the performance of the above algorithm is non-decreasing on the training set because new prompts are added to the population rather than replacing previous ones.

In Algorithm 1, where $n_r, n_m, n_e, n_0, g$ denote the numbers of reproduction, mutation, remaining after elimination, initial mutations, and generations, respectively. Specifically, the measurement function evaluates and partitions the population based on their performance across all training CoTs with respect to their similarity to human annotations. From the final generation, we select the optimal candidate based on validation set performance for subsequent machine annotation tasks.

However, a potential risk in the former optimization is the potential biases adopted from a relatively small training set. Beyond other representative GA-based prompt optimization solutions [48–50], which share a similar process as CAPO, we further propose a

---

**Algorithm 1** CAPO Algorithm

---

**Require:** Hyperparameters: $n_r, n_m, n_e, n_0, g > 0$
**Require:** Labeled long CoTs: $\{x_i, y_i\}_i$
**Require:** Original population: $G \leftarrow \{p_0\}$
**Require:** Measurement function: $M(G) = G_{\text{bad}}, G_{\text{good}}$
**Require:** Mutation function: $\pi(p_., \{x, y\}) = \tilde{p}$
**Require:** Reproduction function: $\sigma(p_., p_.) = \tilde{p}$
**Require:** Elimination function: $\phi(G, n_e) = \tilde{G}$
  **while** $n_0 > 0$ **do**
    $n_0 \leftarrow n_0 - 1$
    $\{x, y\} \leftarrow \text{sample}(\{x_i, y_i\}_i)$
    $G \leftarrow G \cup \{\pi(p_0, \{x, y\})\}$
  **end while**
  **while** $g > 0$ **do**
    $g \leftarrow g - 1$
    $G_{\text{bad}}, G_{\text{good}} \leftarrow M(G)$
    $\{p_{i_k}, p_{j_k}\}_{k=1}^{n_r} \leftarrow \text{sample}_{n_r}(G_{\text{good}})$
    **for** $p_{i_k}, p_{j_k}$ **do**
      $G \leftarrow G \cup \{\sigma(p_{i_k}, p_{j_k})\}$   # **Reproduction**
    **end for**
    $\{p_k\}_{k=1}^{n_m} \leftarrow \text{sample}_{n_m}(G)$
    $\{x_k, y_k\}_{k=1}^{n_m} \leftarrow \text{sample}_{n_m}(\{x_i, y_i\}_i)$
    **for** $p_k, \{x_k, y_k\}$ **do**
      $G \leftarrow G \cup \{\pi(p_k, \{x_k, y_k\})\}$  # **Mutation**
    **end for**
    $G \leftarrow \phi(G, n_e)$   # **Elimination**
  **end while**
  **Return** $G$

---

*tripartite* prompt structure to constrain the optimization from unacceptable biasing as in Figure 5. That is, as our objective requires the taxonomy to reflect human mental processes faithfully, CAPO is particularly constrained to ensure that the proposed taxonomy is well preserved against optimization bias when trained on relatively small datasets. Specifically, the first component, **constant region**, referring to the invariant portion present in all prompts, defines the input/output format and task type. The **variable region** contains descriptions of the taxonomy and may exhibit minor variations across the training. In contrast, the **mutable region**, constituting an entirely unrestricted section, is fully open-ended and may incorporate unconstrained details. Such designs enable the taxonomy, described in the variable region, to retain its principle while annotation skills are mainly embedded into the mutable region.

### 5.1 CAPO evaluation

To assess the quality of CAPO, we leverage *consistency* (↑) judgment by calculating the step proportion in a CoT that the LLM has identical annotations as human annotators. When implementing CAPO as the mentioned LLM-as-annotators in Section 4, we use the following hyperparameters in Algorithm 1: $n_r = 4, n_m = 5, n_e = 8, n_0 = 10, g = 4$ and $|G_{good}| = 5$. The measurement function $M$ and elimination function $\phi$ are directly derived from the consistency metric. And, the initial prompt seed $p_0$ is constructed by human experts. Both the initial prompt from experts and the optimization meta prompt are available in the Appendix B.
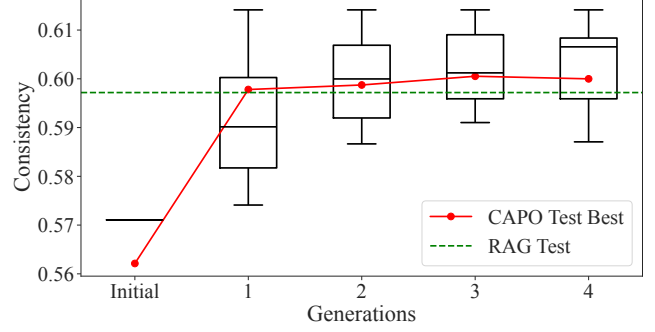


**Figure 6: Consistencies of CAPO and RAG baseline. The boxes depict the consistencies in the training set of the whole population, while the red line illustrates the testing consistency of the best training individual.**

We evaluate the performance of the CAPO using *Gemini-2.5-flash-preview-05-20 without thinking mode* as the LLM annotator, which is observed to be the best LLM during our development. For evaluation, we use half of the human-annotated CoTs from AIME★ and HMMT as the training set and the other half as the test set. And, the metric "consistency" is the average proportion of steps that have identical annotations from LLM as the human ones. We utilize the retrieval augmentation generation (RAG) as a baseline for better evaluation, which retrieves supportive evidence to improve responses by ICL, and is used in [27]. It leverages an embedding model (Linq-Mistral [51]) to encode all training CoTs into vectors. For each test CoT, the most similar example, in terms of inner product, is organized into the prompt together with its labels for querying LLM-as-annotator. Here, the labels include paired zero-shot annotations and the human annotations for a fair comparison.

Following Section 3, we engineer an initial prompt seed and execute CAPO for multiple rounds, and the results are shown in Figure 6. It's observed that both CAPO and RAG can improve the annotation consistency to achieve near 60% consistency. Moreover, CAPO will be free of long input prompts for examples (as previously discussed due to a shortage of ICL), while it also surpasses the RAG baseline after even a single optimization round. Eventually, the LLM-as-annotators with CAPO exhibit acceptable consistency with humans. Thus, their results can be considered homogeneous extensions of human ones, where the deviations are further discussed in the following Section 6.

## 6 Availability & Discussion

**All the source code, data, and analysis in Section 4 and Section 5 are available at** https://github.com/hehepig4/psyche. Furthermore, we will discuss the limitations and future work of the subsequent three aspects:

**Suboptimal consistency between human and LLM as annotators**. One may argue that although our CAPO effectively improves annotation quality, it still struggles to achieve a seemingly satisfactory level of consistency. Here, it is essential to emphasize that automated annotation serves only as an auxiliary method to augment the volume of human-labeled data. **We ensure that the findings presented above in Section 4 remain consistent across both LLM-annotated and human-annotated subsets.**

For instance, in Takeaway 1, both annotation sets indicate that a higher A.IO ratio contributes to the model's ability to answer questions correctly. The primary motivation for incorporating LLM-generated annotations in our study stems from the limited sample size in the human-annotated set, which undermines statistical reliability. Improving annotation accuracy, such as by incorporating more human-annotated data for post-training, represents an interesting direction for future work.

**Reasoning models beyond Deepseek R1 and Qwen QwQ**. The primary motivations for selecting Qwen QwQ and Deepseek R1 as our subject models are twofold: (a) they represent leading open-source LRMs, enabling direct access to their unaltered CoT outputs; and (b) their core capabilities have been extensively validated in prior work [52]. Observing that different LRMs exhibit discernible variations in performance across different tasks or contexts, we propose that a comparative analysis involving newer opened models (e.g., updated version of R1) and representative closed-source LRMs with accessible APIs (e.g., Gemini 2.5) constitutes a significant avenue for future research.

**Analysis of More complex real-world reasoning**. As empirically observed during data collection, modern LRMs consistently achieve near-perfect accuracy on benchmark tasks assessing commonsense reasoning in daily-life scenarios (e.g., Commonsense QA). This high performance precludes the use of binary (correct/incorrect) evaluation, presenting significant challenges in assessing the quality of the generated reasoning chains. This difficulty extends to complex, real-world reasoning problems with greater severity. Consequently, we select mathematical problem-solving as our primary analytical focus. Mathematical problems offer sufficient complexity to thoroughly evaluate intrinsic reasoning capabilities while providing unambiguous response verification. In future work, we plan to examine how mental processes affect performance in complex reasoning scenarios from multifaceted perspectives, including robustness, computational efficiency, and planning effectiveness [53].

## 7 Conclusion

Large reasoning models (LRMs) have demonstrated strong potential in web applications, yet their behaviors are difficult to interpret. Therefore, this paper presents a novel taxonomy for analyzing reasoning behaviors in LRMs, establishing a bridge between computational methods and human cognitive processes. To operationalize this taxonomy, we introduce CAPO, a complementary framework that leverages large language models (LLMs) as annotators, enabling scalable and expert-consistent labeling. Using CAPO in collaboration with domain experts, we construct a high-quality dataset of 277,534 labeled reasoning steps, facilitating the first large-scale empirical study of LRM reasoning from a human-centric perspective. Our analysis yields four key insights, providing actionable directions for advancing LRMs.

## References

[1] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, et al. 2025. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. *CoRR* abs/2501.12948 (2025).

[2] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, et al. 2024. OpenAI o1 System Card. *CoRR* abs/2412.16720 (2024).

[3] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, et al. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In *NeurIPS*.

[4] Jun Wang. 2025. A Tutorial on LLM Reasoning: Relevant Methods behind ChatGPT o1. *arXiv preprint arXiv:2502.10867* (2025). https://doi.org/10.48550/arXiv.2502.10867

[5] Jun Wang, Meng Fang, Ziyu Wan, Muning Wen, Jiachen Zhu, et al. 2024. OpenR: An Open Source Framework for Advanced Reasoning with Large Language Models. *arXiv preprint arXiv:2410.09671* (2024).

[6] Xidong Feng, Ziyu Wan, Muning Wen, Stephen Marcus McAleer, Ying Wen, et al. 2023. AlphaZero-Like Tree-Search can Guide Large Language Model Decoding and Training. *arXiv preprint arXiv:2309.17179* (2023).

[7] Gerhard Gentzen. 1935. Untersuchungen über das logische Schließen. *Mathematische Zeitschrift* 39 (1935), 176–210, 405–431.

[8] Dag Prawitz. 1965. *Natural Deduction: A Proof-Theoretical Study*. Almqvist & Wiksell, Stockholm.

[9] Lev S. Vygotsky. 1978. *Mind in Society: The Development of Higher Psychological Processes*. Harvard University Press, Cambridge, MA.

[10] David Wood, Jerome S. Bruner, and Gail Ross. 1976. The Role of Tutoring in Problem Solving. *Journal of Child Psychology and Psychiatry* 17, 2 (1976), 89–100. https://doi.org/10.1111/j.1469-7610.1976.tb00381.x

[11] Allan Collins, John Seely Brown, and Susan E. Newman. 1989. Cognitive Apprenticeship: Teaching the Crafts of Reading, Writing, and Mathematics. In *Knowing, Learning, and Instruction: Essays in Honor of Robert Glaser*, Lauren B. Resnick (Ed.). Lawrence Erlbaum Associates, Hillsdale, NJ, 453–494.

[12] John H. Flavell. 1979. Metacognition and Cognitive Monitoring: A New Area of Cognitive–Developmental Inquiry. *American Psychologist* 34, 10 (1979), 906–911. https://doi.org/10.1037/0003-066X.34.10.906

[13] Michelene T. H. Chi, Miriam Bassok, Matthew W. Lewis, Peter Reimann, and Robert Glaser. 1989. Self-Explanations: How Students Study and Use Examples in Learning to Solve Problems. *Cognitive Science* 13, 2 (1989), 145–182. https://doi.org/10.1016/0364-0213(89)90002-5

[14] John Dunlosky, Katherine A. Rawson, Elizabeth J. Marsh, Mitchell J. Nathan, and Daniel T. Willingham. 2013. Improving Students' Learning With Effective Learning Techniques: Promising Directions From Cognitive and Educational Psychology. *Psychological Science in the Public Interest* 14, 1 (2013), 4–58. https://doi.org/10.1177/1529100612453266

[15] Yancheng He, Shilong Li, Jiaheng Liu, Weixun Wang, Xingyuan Bu, et al. 2025. Can Large Language Models Detect Errors in Long Chain-of-Thought Reasoning? *CoRR* abs/2502.19361 (2025).

[16] Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, et al. 2024. Lost in the Middle: How Language Models Use Long Contexts. *Trans. Assoc. Comput. Linguistics* 12 (2024), 157–173.

[17] Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, et al. 2025. Towards Reasoning Era: A Survey of Long Chain-of-Thought for Reasoning Large Language Models. *CoRR* abs/2503.09567 (2025).

[18] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, et al. 2024. Graph of Thoughts: Solving Elaborate Problems with Large Language Models. In *AAAI*. AAAI Press, 17682–17690.

[19] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, et al. 2023. Tree of Thoughts: Deliberate Problem Solving with Large Language Models. In *NeurIPS*.

[20] Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. 2022. STaR: Bootstrapping Reasoning With Reasoning. In *NeurIPS*.

[21] Xiaoyu Tan, Tianchu Yao, Chao Qu, Bin Li, Minghao Yang, et al. 2025. AURORA:Automated Training Framework of Universal Process Reward Models via Ensemble Prompting and Reverse Verification. *CoRR* abs/2502.11520 (2025).

[22] Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, et al. 2024. Math-Shepherd: Verify and Reinforce LLMs step-by-step without Human Annotations. In *ACL (1)*. Association for Computational Linguistics, 9426–9439.

[23] Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, Yi Su, John D. Co-Reyes, et al. 2025. Training Language Models to Self-Correct via Reinforcement Learning. In *ICLR*. OpenReview.net.

[24] Parshin Shojaee, Iman Mirzadeh, Keivan Alizadeh, Maxwell Horton, Samy Bengio, et al. 2025. The Illusion of Thinking: Understanding the Strengths and Limitations of Reasoning Models via the Lens of Problem Complexity. *CoRR* abs/2506.06941 (2025).

[25] Siwei Wu, Zhongyuan Peng, Xinrun Du, Tuney Zheng, Minghao Liu, et al. 2024. A Comparative Study on Reasoning Patterns of OpenAI's o1 Model. *CoRR* abs/2410.13639 (2024).

[26] Dongkeun Yoon, Seungone Kim, Sohee Yang, Sunkyoung Kim, Soyeon Kim, et al. 2025. Reasoning Models Better Express Their Confidence. *CoRR* abs/2505.14489 (2025).

[27] Sara Vera Marjanovic, Arkil Patel, Vaibhav Adlakha, Milad Aghajohari, Parishad BehnamGhader, et al. 2025. DeepSeek-R1 Thoughtology: Let's about LLM Reasoning. *CoRR* abs/2504.07128 (2025).

[28] Plato. 1892. *Statesman.* Oxford University Press. https://www.gutenberg.org/files/1738/1738-h/1738-h.htm Translated by Benjamin Jowett.

[29] Aristotle. 1928. *The Categories.* Oxford University Press, Oxford. https://www.gutenberg.org/ebooks/2412

[30] Richard E. Mayer. 1998. Cognitive, Metacognitive, and Motivational Aspects of Problem Solving. *Instructional Science* 26, 1–2 (1998), 49–63.

[31] Poonam Punia, Ritu Malik, Manju Bala, Manju Phor, and Yogesh Chander. 2023. Relationship between Logical Thinking, Metacognitive Skills, and Problem Solving Abilities: Mediating and Moderating Effect Analysis. *International Journal of Educational and Developmental Psychology* (2023).

[32] Rudolf Carnap and Richard C. Jeffrey (Eds.). 1971. *Studies in Inductive Logic and Probability, Volume I.* University of California Press, Berkeley and Los Angeles. https://doi.org/10.1525/9780520334250

[33] Philip N. Johnson-Laird and Ruth M. J. Byrne. 1991. *Deduction.* Erlbaum, Hillsdale, NJ.

[34] Paul O'Rourke and John R. Josephson (Eds.). 1997. *Automated Abduction: Inference to the Best Explanation.* AAAI Press, Menlo Park, CA. http://cogprints.org/668/

[35] John Dewey. 1910. *How We Think.* D.C. Heath and Company, Boston, MA. https://www.gutenberg.org/ebooks/37423

[36] Vladimir A. Gheorghiu, Petra Netter, Hans J. Eysenck, and Robert Rosenthal (Eds.). 1987. *Suggestion and Suggestibility: Theory and Research.* Springer-Verlag, Berlin, Heidelberg. Proceedings of the First International Symposium on Suggestion and Suggestibility, University of Giessen, July 7–11, 1987.

[37] Yiwei Qin, Xuefeng Li, Haoyang Zou, Yixiu Liu, Shijie Xia, et al. 2024. O1 Replication Journey: A Strategic Progress Report - Part 1. *CoRR* abs/2410.18982 (2024). https://doi.org/10.48550/ARXIV.2410.18982 arXiv:2410.18982

[38] OpenCompass. 2025. AIME 2025 dataset. https://huggingface.co/datasets/opencompass/AIME2025

[39] Mislav Balunović, Jasper Dekoninck, Ivo Petrov, Nikola Jovanović, and Martin Vechev. 2025. MathArena: Evaluating LLMs on Uncontaminated Math Competitions. https://matharena.ai/

[40] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, et al. 2021. Measuring Mathematical Problem Solving With the MATH Dataset. In *NeurIPS Datasets and Benchmarks.*

[41] Di Zhang. 2025. AIME_1983_2024 (Revision 6283828). https://doi.org/10.57967/hf/4687

[42] Team Qwen. 2025. QwQ-32B: Embracing the Power of Reinforcement Learning. https://qwenlm.github.io/blog/qwq-32b/

[43] Xiangning Yu, Zhuohan Wang, Linyi Yang, Haoxuan Li, Anjie Liu, et al. 2025. Causal Sufficiency and Necessity Improves Chain-of-Thought Reasoning. *arXiv preprint arXiv:2506.09853* (2025).

[44] David Cheng-Han Chiang and Hung-yi Lee. 2023. Can Large Language Models Be an Alternative to Human Evaluations?. In *ACL (1).* Association for Computational Linguistics, 15607–15631.

[45] Nitay Calderon, Roi Reichart, and Rotem Dror. 2025. The Alternative Annotator Test for LLM-as-a-Judge: How to Statistically Justify Replacing Human Annotators with LLMs. *CoRR* abs/2501.10970 (2025).

[46] Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, et al. 2024. A Survey on LLM-as-a-Judge. *CoRR* abs/2411.15594 (2024).

[47] Zhen Tan, Dawei Li, Song Wang, Alimohammad Beigi, Bohan Jiang, et al. 2024. Large Language Models for Data Annotation and Synthesis: A Survey. In *EMNLP.* Association for Computational Linguistics, 930–957.

[48] Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, et al. 2024. Connecting Large Language Models with Evolutionary Algorithms Yields Powerful Prompt Optimizers. In *ICLR.* OpenReview.net.

[49] Xavier Sécheresse, Jacques-Yves Guilbert-Ly, and Antoine Villedieu de Torcy. 2025. GAAPO: Genetic Algorithmic Applied to Prompt Optimization. *CoRR* abs/2504.07157 (2025).

[50] Yurong Wu, Yan Gao, Bin Zhu, Zineng Zhou, Xiaodi Sun, et al. 2024. StraGo: Harnessing Strategic Guidance for Prompt Optimization. In *EMNLP (Findings).* Association for Computational Linguistics, 10043–10061.

[51] Junseong Kim, Seolhwa Lee, Jihoon Kwon, Sangmo Gu, Yejin Kim, et al. 2024. Linq-Embed-Mistral:Elevating Text Retrieval with Improved GPT Data Through Task-Specific Control and Quality Refinement. Linq AI Research Blog. https://getlinq.com/blog/linq-embed-mistral/

[52] Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, et al. 2025. Towards Reasoning Era: A Survey of Long Chain-of-Thought for Reasoning Large Language Models. *arXiv preprint arXiv:2503.09567* (2025).

[53] Dongzhi Jiang, Renrui Zhang, Ziyu Guo, Yanwei Li, Yu Qi, et al. 2025. MME-CoT: Benchmarking Chain-of-Thought in Large Multimodal Models for Reasoning Quality, Robustness, and Efficiency. *CoRR* abs/2502.09621 (2025).

# A CoT Category Examples

Given below are explanations and real-world examples for each mental process selected from our annotated dataset.

## A.1 Analysis

**Category:** Analysis.Problem_Definition
**Explanation:** Identify and clearly define the core difficulty or central question in the task.

> **Example:** "Okay, let's see. I have this problem with triangle ABC. Points D and E are on AB, and F and G are on AC. The lengths are given: AD = 4, DE = 16, EB = 8, AF = 13, FG = 52, and GC = 26. The area of quadrilateral DEGF is 288. I need to find the area of the heptagon AFNBCEM."

**Category:** Analysis.Information_Organization
**Explanation:** List and organize all relevant background information and known facts.

> **Example:**
> **Step 125:** "Term1 = $0 \times 45/7 - (24/7) \times 9 = -216/7$. "
> **Step 559:** "Term1 = -216/7."

**Category:** Analysis.Problem_Structuring
**Explanation:** Decompose the main problem into smaller sub-problems and explain their logical relationships.

> **Example:** "Since D and E divide AB into 1:4:2, and F and G divide AC into 1:4:2, coordinates can be expressed in these ratios. Let me parametrize points D, E, F, G."

## A.2 Inference

**Category:** Inference.Deductive_Reasoning
**Explanation:** Apply general principles or rules to deduce specific conclusions relevant to the current task.

> **Example:** "Area = $\frac{1}{2} \times |24| = 12$. Therefore, the ratios on AB are $AD : DE : EB = 4 : 16 : 8$, so from A to D is 4/28 = 1/7, D to E is 16/28 = 4/7, and E to B is 8/28 = 2/7."

**Category:** Inference.Inductive_Reasoning
**Explanation:** Generalize patterns from specific examples or observations.

> **Example:** "Since $a = 2b$ for $n = 1, n = 2, n = 3$, then I suppose that $a = 2b$ for any $n$."

**Category:** Inference.Abductive_Reasoning
**Explanation:** Given an observation, propose the most likely or plausible explanation.

> **Example:** "Wait, maybe in the last period, the number of solutions is less, so the symmetry breaks. In previous periods, we had two solutions on each edge, so four solutions, and the sum cancels out. But in the last period, k=8, rising edge only has one solution and falling edge has one solution, so their sum is not zero."

## A.3 Judgment

**Category:** Judgment.Principle_Selection
**Explanation:** Choose appropriate logical principles or domain-specific rules needed to evaluate the problem.

> **Example:** "3. **Derivative and Wilson's Theorem**: - Using Wilson's theorem, $P'(k) \equiv -1 \pmod{2027}$). - Therefore, $(\frac{-\sum_{m \neq k} km}{P'(k)} \equiv k \pmod{2027})$."

**Category:** Judgment.Evaluation_of_Alternatives
**Explanation:** Compare multiple reasoning paths or hypotheses and select the most promising one.

> **Example:** "To count the number of solutions, perhaps we can use inclusion-exclusion, but that might be complex. Alternatively, since the grid is small, maybe we can compute it directly by degrees of freedom."

**Category:** Judgment.Conclusion_Decision
**Explanation:** Make a final decision or answer based on prior reasoning and comparisons.

> **Example:** "Therefore, triangle ABC has area 588."

## A.4 Suggestion

**Category:** Suggestion.Strategic_Planning
**Explanation:** Develop a reasoning roadmap or outline for solving the problem.

> **Example:** "To compute the heptagon's area, we can subtract the areas outside it from the area of triangle ABC."

**Category:** Suggestion.Branch_Changing
**Explanation:** Abandon the current reasoning path and explore a new or contrasting approach.

> **Example:** "Alternatively, use the shoelace formula with the coordinates."

**Category:** Suggestion.Hypothesis_Generation

**Explanation:** Generate a speculative explanation or assumption based on limited evidence.

> **Example:** "Perhaps $u$ is real, which could explain the rotated parabola's points."

**Category:** Suggestion.Analogy_Recall
**Explanation:** Introduce an analogous situation or familiar pattern to guide reasoning.

> **Example:** "Since D is 1/7 along AB and F is 1/7 along AC, we can use the same ratio logic for E and G."

## A.5 Reflection

**Category:** Reflection.Self_Monitoring_Evaluation
**Explanation:** Review current reasoning steps for gaps, errors, or inconsistencies.

> **Example:** "Wait, if $q = 42$, then area of ABC = 588. Shoelace formula gives DEGF = 288, which matches. So this is correct."

**Category:** Reflection.Counterfactual_Thinking
**Explanation:** Consider alternative actions and speculate on "what-if" scenarios.

> **Example:** "If we had divided by 22, the probability would be 0.045, but actual count shows 0.057 due to dependencies."

**Category:** Reflection.Causal_Attribution
**Explanation:** Analyze reasons behind success or failure by identifying key factors that caused the result.

> **Example:** "Because divisibility by 2 and 11 are not independent, our initial assumption failed."

**Category:** Reflection.Strategy_Regulation
**Explanation:** Adjust the overall problem-solving strategy based on reflection or feedback.

> **Example:** "Since I'm stuck, I should consult an official solution or try a different known formula."

## B  Prompts
**The initial prompt seed for CAPO is as follows:**

> **# Instruction**
> You need to classify meta-behaviors in an inputted chain of thought (CoT) when solving a problem. Each step is enclosed by <step *>, where * is the order number of a step.

Then, you should identify the meta-behaviors by each step. Details of the task are as follows.

**Meta-behaviors include:**

***Analysis*: Decomposing and understanding the problem before proceeding to reasoning or evaluation.**

- *Analysis.Problem_Definition*: Identify and clearly describe the core difficulty or central question in the problem.

- *Analysis.Information_Organization*: List and organize all relevant background information and known facts.

- *Analysis.Problem_Structuring*: Break the problem into smaller sub-problems and explain their logical connections and roles in solving the overall task.

***Inference*: Making logical deductions from known information to arrive at new conclusions. This is the core phase of reasoning.**

- *Inference.Deductive_Reasoning*: Apply general rules or principles to derive specific conclusions relevant to the problem.

- *Inference.Inductive_Reasoning*: Observe specific cases and infer a general rule or trend that applies to the situation.

- *Inference.Abductive_Reasoning*: Given an observation, propose the most likely or plausible explanation—even if it's uncertain.

***Judgment*: Assessing different solution paths and forming final decisions based on reasoning.**

- *Judgment.Principle_Selection*: Identify and apply the most appropriate logical principles, ethical rules, or domain-specific criteria before making a judgment or decision.

- *Judgment.Evaluation_of_Alternatives*: Consider multiple possible reasoning paths or hypotheses, then compare and identify the most promising one.

- *Judgment.Conclusion_Decision*: Make a final decision or answer based on previously completed reasoning and evaluation.

***Suggestion*: Proposing new ideas, speculative paths, or reasoning strategies that go beyond the direct content of the problem.**

- *Suggestion.Strategic_Planning*: Develop a plan or roadmap for the reasoning steps needed to solve the problem.

- *Suggestion.Branch_Changing*: Switch to a different approach of reasoning or explore an alternative method when current direction seems unpromising.

- *Suggestion.Hypothesis_Generation*: Formulate a speculative explanation or guess based on limited evidence to guide further reasoning.

- *Suggestion.Analogy_Recall*: Bring in a familiar case, past experience, or known pattern to inspire a solution idea or strategy.

***Reflection*: Monitoring and evaluating the reasoning process to ensure logical correctness and coherence.**

- *Reflection.Self_Monitoring_Evaluation*: Review the reasoning process so far. Check for gaps, mistakes, or inconsistencies in logic.

- *Reflection.Counterfactual_Thinking*: Consider alternative actions or decisions and speculate on what might have happened under different conditions. Used to reassess current reasoning or outcomes based on "what-if" scenarios.

- *Reflection.Causal_Attribution*: Analyze the reasons behind success or failure by identifying the key factors or decisions that caused the result. Supports better learning from experience.

- *Reflection.Strategy_Regulation*: Adjust the overall problem-solving or reasoning strategy based on feedback or prior reflection. Helps improve future performance by refining the approach.

A meta-behavior is represented hierarchically and separated by a full stop.

**# Task**

**You should select all relevant meta-behaviors, ranking them and separating them by semicolon in descending order based on their relevance and importance.**

All the meta-behaviors are not exclusive, and a step may contain multiple meta-behaviors.

Also, these meta-behaviors may belong to a same type but with different sub-types, like *Analysis.Problem_Definition* and *Analysis.Problem_Structuring*, etc.

You should choose all the meta-behaviors that are possible in this step.

**To provide a precise and faithful answer, you need to fully utilize the semantic connection between consecutive steps.**

**# Output Format**
**You need to strictly output in the following format**:
<step 1> meta-behavior(s) </step 1>
<step 2> meta-behavior(s) </step 2>
...
<step n> meta-behavior(s) </step n>
**# Input**
**## The problem**
**{problem_desc}**
**## The long CoT**
**{CoT}**
**# Output**

During optimization, the description of each meta-behavior (after "*meta behavior include:*") is set as the variable area, and a new region named "tips" is the mutable area, while the remaining part is identified as the constant area. For the meta-prompts, **the mutation prompt is as follows:**

**# Instruction:**
You are an expert in prompt engineering.
The following five prompts describe a task.
You will also be given an example of a response to this prompt.
Your task is to mutate and improve a specific part of the prompt, based on the example and answer, according to the following rules:

(1) Review the task in the prompt to understand the key objectives and requirements that the instruction needs to address.

(2) Focus on the part of the prompt indicated by the <part> tag that needs mutation.

(3) Analyse the example and answer to identify any gaps, ambiguities or areas for improvement in the prompt.

(4) Maintain the original format and structure of the prompt while enhancing clarity, specificity and guidance in the mutated part.

(5) Ensure that the output format of the mutated part is consistent with the original prompt structure.

(6) Do not modify the names of meta-behaviors or the structure of the prompt.

(7) Tags (e.g., in <>) should not be modified.

(8) All names of meta-behaviors (xx.xx) should be strictly retained as they are, except for their descriptions, which can be modified for clarity. And, no addition and deletion of meta-behaviors is allowed.

(9) Particularly, some absent fields in the example typically indicate that the original response is with incorrect format, try to fix it.

(10) Especially, tips, including details of meta-behaviors and tasks, are more flexible and can be modified to better fit the merged prompt.

(11) Notice that, example is not available when the following prompt is used for downstream tasks.

Output the specific mutated part in the <mutated_part> tag after providing a comprehensive, step-by-step thinking, as follows:

**Output format: <mutated_part> Mutated and improved part of the prompt </mutated_part>**

**{Current prompt}**
**# Example**
**{example}**
**# Mutation Target**
**{part_name}**

**And, the reproducibility meta-prompt is:**

# **Instruction** You are an expert in prompt engineering.

You will be given two prompts, each consisting of five parts that describe a task.

Each prompt is optimized based on an example of the task.

Your task is to combine these two prompts to create a single, coherent prompt that retains the strengths of both while ensuring clarity and consistency.

Focus on the following aspects:

(1) Identify the key objectives and requirements in both prompts.

(2) Combine the relevant parts from both prompts to create a comprehensive and clear merged prompt.

(3) Maintain the original format and structure of the prompts, enhancing clarity, specificity and guidance where necessary.

(4) Ensure the merged prompt is logically consistent and flows well.

(5) The output format should be consistent with the original prompt structure.

(6) Do not modify the names of meta-behaviors or the structure of the prompt.

(7) Tags (e.g., in <>) should not be modified.

(8) All names of meta-behaviors (xx.xx) should be strictly retained as they are, except for their descriptions, which can be modified for clarity. And, no addition and deletion of meta-behaviors is allowed.

(9) The consideration aspects of both prompts may vary, so you should carefully merge them to ensure that the final prompt is comprehensive.

(10) Especially, tips, including details of meta-behaviors and tasks, are more flexible and can be modified to better fit the merged prompt.

Output the merged prompt with the five tags after providing comprehensive step-by-step reasoning as follows:

***Output format omitted***
**# Prompt 1**
**{*Current prompt 1*}**
**# Prompt 2**
**{*Current prompt 2*}**

Due to page limitations, please refer to our source code for concrete prompts.