

# Valid $\neq$ Necessary: Diagnosing Latent Inefficiency in Chain-of-Thought

Anonymous ACL submission

## Abstract

Chain-of-Thought (CoT) prompting has significantly advanced the reasoning capabilities of Large Language Models (LLMs), yet it often incurs substantial computational costs due to “over-reasoning”—the generation of redundant, verbose, or irrelevant steps. While existing reasoning step evaluators effectively detect logical fallacies and factual errors, our analysis reveals a critical blind spot: they fail to penalize “valid but inefficient” reasoning steps that inflate token usage without contributing to the solution. To systematically diagnose this limitation, we introduce **RIV-GSM8K**, a diagnostic benchmark injected with five distinct types of inefficiencies, including circular reasoning and excessive decomposition. Diagnostic experiments reveal that state-of-the-art evaluators struggle to distinguish these inefficiencies from necessary reasoning. To address this, we propose **CAID** (Context-Aware Information Density), a training-free metric grounded in information theory that effectively identifies low-utility steps. To validate the metric’s practical utility, we apply it within **PACE**, a post-hoc compression strategy. Empirical results on GSM8K, StrategyQA, and ARC-Challenge demonstrate that PACE reduces token consumption by 31–53% while maintaining accuracy, confirming that CAID successfully distills informational “froth” from reasoning chains without compromising deductive validity.

## 1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities in complex reasoning tasks, largely driven by the Chain-of-Thought (CoT) prompting strategy (Wei et al., 2022). By decomposing complex problems into intermediate steps, LLMs can bridge the gap between question and answer. However, this performance gain often comes at the cost of inference efficiency. Recent studies indicate that LLMs exhibit a tendency towards “over-reasoning”—generating verbose ex-

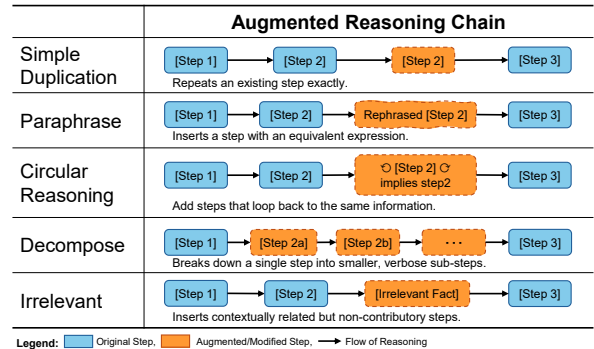


Figure 1: Taxonomy of reasoning inefficiencies in RIV-GSM8K. The diagram illustrates how five distinct types of redundant steps are synthetically injected into the reasoning chain to simulate valid but dispensable “froth.”

planations, repetitive statements, or contextually irrelevant details that inflate computational costs without adding deductive value (Turpin et al., 2023; Wang et al., 2023; Chiang and Lee, 2024).

To assess and improve reasoning quality, various Process Reward Models (PRMs) and reasoning step evaluators have been proposed, such as ReasonEval and Math-Shepherd (Xia et al., 2025; Wang et al., 2024). These methods have primarily focused on **correctness** and **logical validity**, aiming to penalize factual errors or hallucinations. While effective for validity, our analysis reveals a critical blind spot in these state-of-the-art evaluators: they struggle to distinguish *inefficiency* from *reasoning*. Specifically, they often assign high scores to “valid but redundant” steps—such as excessive decomposition or circular logic—merely because these steps remain factually true and linguistically coherent. Consequently, the fundamental question in current reasoning evaluation is limited to “Is this step true?”, neglecting the equally important dimension: **“Is this step truly necessary?”**

In this paper, we aim to shift the evaluation paradigm from verifying correctness to optimizing **information density**. To systematically diagnose the limitations of current evaluators, we first

introduce **RIV-GSM8K**, a diagnostic benchmark derived from GSM8K (Cobbe et al., 2021). As illustrated in Figure 1, RIV-GSM8K is synthetically injected with five distinct types of inefficiencies ranging from simple duplication to subtle circular reasoning. Using this benchmark, we empirically demonstrate that existing validity-focused PRMs are largely insensitive to explicit redundancy.

To address this gap, we propose **CAID** (Context-Aware Information Density), a novel unsupervised metric grounded in information theory that evaluates reasoning steps based on local novelty, global goal alignment, and information density. Unlike previous metrics, CAID effectively identifies informational “froth” within reasoning chains. To **empirically validate the diagnostic precision** of this metric, we present **PACE** (Pruning And Compression for Efficiency), a post-hoc compression strategy. Notably, PACE goes beyond simple deletion; it identifies *latent inefficiency* in reasoning chains and predominantly **compresses** verbose steps (Merge) while pruning irrelevant ones, ensuring that the process remains logically sound but significantly more compact.

Our main contributions are summarized as follows:

- We uncover the “efficiency blind spot” of current reasoning evaluators through **RIV-GSM8K**, a stress-test benchmark designed to diagnose specific types of reasoning inefficiencies.
- We propose **CAID**, an interpretable, training-free metric that quantifies the informational utility of reasoning steps, distinguishing essential logic from redundant “froth.”
- We validate our approach via **PACE**, which reduces token consumption by 31–53% across arithmetic, commonsense, and scientific reasoning tasks without compromising accuracy. This **serves as empirical evidence** for the existence of significant *latent inefficiency* in standard CoT reasoning, demonstrating that high-quality reasoning data can be far more compact than previously assumed.

## 2 Related Work

### 2.1 Over-reasoning and Inference Efficiency

While Chain-of-Thought (CoT) prompting (Wei et al., 2022) has revolutionized LLM reasoning, it

has also introduced the challenge of inference inefficiency. Recent studies highlight that LLMs suffer from “**over-reasoning**”—a tendency to generate verbose explanations, repetitive loops, or contextually irrelevant details (Turpin et al., 2023; Chen et al., 2024; Chiang and Lee, 2024). Jiang et al. (2023) demonstrated that such redundant context not only inflates computational costs but can also degrade performance by distracting the model.

To mitigate computational costs, token pruning methods such as H2O (Zhang et al., 2023) and Learned Token Pruning (Kim et al., 2022) have been proposed. These approaches reduce sequence length by discarding tokens with low attention scores during inference. However, they operate at the **token level**, focusing primarily on latency reduction (KV cache optimization) rather than the informational quality of the content. In contrast, our strategy, PACE, addresses inefficiency at the **semantic step level**. Rather than competing with runtime token pruners, we focus on identifying and filtering out *latent inefficiency* embedded in the reasoning chain itself. This semantic compression serves a complementary role, potentially enhancing the quality of training data for future models.

### 2.2 Reasoning Step Evaluation Methods

Moving beyond outcome-based evaluation, step-wise evaluation methods have emerged to provide granular supervision. Process Reward Models (PRMs) like PRM800K (Lightman et al., 2024) and Math-Shepherd (Wang et al., 2024) train verifiers to distinguish between correct and incorrect reasoning paths. More recently, **ReasonEval** (Xia et al., 2025) has attempted to assess reasoning quality beyond mere validity, incorporating scores for redundancy and clarity.

However, existing methods face two critical limitations. First, standard PRMs primarily focus on **correctness verification**, penalizing factual errors but often rewarding valid yet inefficient steps (e.g., stating the obvious). Second, while ReasonEval attempts to detect redundancy, it relies on supervised training with human-annotated data. **Crucially, human annotators often prioritize validity over conciseness, leading to sparse and ambiguous labels regarding what constitutes “unnecessary.”** In contrast, our approach utilizes **RIV-GSM8K**, where inefficiencies are synthetically injected based on a distinct taxonomy. This allows for a deterministic evaluation of detection capability.

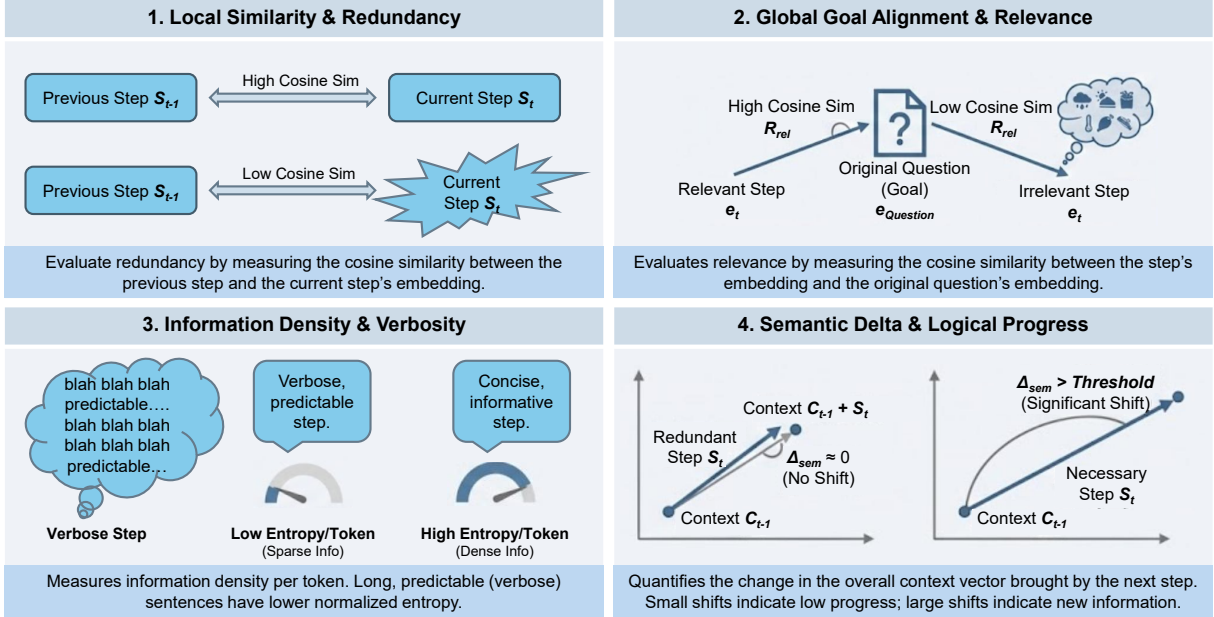


Figure 2: Conceptual overview of the four information-theoretic indicators in CAID. (1) Local Similarity detects surface-level redundancy via adjacent step comparison. (2) Global Goal Alignment filters irrelevant steps by measuring drift from the original question. (3) Information Density identifies verbose, low-entropy steps using length-normalized perplexity. (4) Semantic Delta quantifies the vector shift in the context, ensuring substantial logical progress.

### 2.3 Redundancy Detection Metrics

A few studies have addressed the evaluation of text redundancy. ROSCOE (Golovneva et al., 2023) proposes a suite of metrics for text generation, using semantic similarity (e.g., SimCSE) to detect repetitions. Similarly, LI et al. (2023) introduced Selective Context, which uses self-information (perplexity) to prune input prompts.

These approaches, however, are limited in the context of multi-step reasoning. ROSCOE relies heavily on **surface-level similarity** and often requires a **reference text** (golden truth) for accurate evaluation. **However, relying on references is problematic because, as our findings suggest, even human-written golden truths often contain significant inefficiency.** Furthermore, metrics like Selective Context measure static information content but do not account for the logical flow. CAID advances these by introducing **Semantic Delta** and **Goal Alignment**, which measure not just the static information of a step, but its dynamic contribution to the logical *progress* towards the solution without requiring a reference chain.

## 3 Methodology

We propose a comprehensive framework to diagnose, measure, and mitigate reasoning inefficiency, shifting the evaluation paradigm from validity-

centric to **efficiency-aware**. Our approach operates in three stages: (1) **Diagnosis**: We introduce **RIV-GSM8K** to expose the “efficiency blind spots” of current evaluators. (2) **Measurement**: We propose **CAID**, an information-theoretic metric designed to detect valid but inefficient steps without human annotations. (3) **Validation**: We apply **CAID** within the **PACE** strategy to empirically validate its diagnostic precision through post-hoc chain compression.

### 3.1 RIV-GSM8K: Diagnosing Inefficiency

A critical challenge in evaluating reasoning efficiency is the lack of an absolute ground truth for “necessary” steps. Human annotations are often subjective and sparse. To address this, we construct **RIV-GSM8K**, a diagnostic benchmark derived from GSM8K (Cobbe et al., 2021) based on a “**relative inefficiency**” paradigm. By synthetically injecting controlled perturbations into baseline reasoning chains (Golden CoT), we create steps that are *strictly less efficient* than the original ones while remaining factually correct. This allows for a deterministic evaluation of whether an evaluator can distinguish “necessary” logic from “redundant” froth.

**Construction Process.** The overall construction procedure is detailed in Appendix A. We employ a hybrid generation approach to balance diversity

and control. **Simple Duplication** is generated via rule-based repetition to serve as a baseline. For complex types (**Paraphrase, Decompose, Circular, Irrelevant**), we utilize **GPT-4o** as the generator  $\mathcal{G}$ . Crucially, we enforce strict constraints during generation to ensure that the injected steps are *valid but inefficient*:

- **No New Progress:** The generated step must not advance the reasoning state beyond the target step  $S_t$ . It must remain logically stationary or redundant.
- **Contextual Coherence:** Particularly for the *Irrelevant* type, the generator is prompted to produce sentences that are mathematically true and linguistically blended with the current context, but strictly “non-contributory” to the solution logic.

The detailed prompts used for each perturbation type are provided in Appendix B. In total, we constructed 7,473 samples containing over 20,000 augmented steps (see Appendix A.2 for detailed statistics).

As illustrated in Figure 1, we define five distinct taxonomies of reasoning inefficiency designed to mimic the “over-reasoning” behaviors of LLMs (visual examples in Appendix C). We first introduce **Simple Duplication** and **Paraphrase** to serve as baselines for detecting lexical and semantic redundancy, respectively. To address informational dilution, we define **Decompose**, which models the excessive fragmentation of a single logical step into multiple *low-density* micro-steps. Finally, we simulate stalled logical progress and deviations from the problem objective through **Circular Reasoning** (self-referential verification) and **Irrelevant** steps (non-contributory facts), thereby testing the model’s sensitivity to *global goal alignment*.

**Quality Verification via Human Evaluation.** To ensure the validity of the synthetically generated inefficiencies, we conducted a human evaluation on the four GPT-4o-generated types (Paraphrase, Decompose, Circular, Irrelevant), excluding the rule-based Simple Duplication. We randomly sampled 30 instances per type (120 total) and verified adherence to the generation constraints, specifically checking for *Logical Equivalence* and the *No New Progress* rule. The evaluation revealed a high overall success rate. The *Irrelevant* and *Circular Reasoning* types achieved near-perfect validity (29/30

and 30/30, respectively). The *Decompose* type showed a slightly higher failure rate (4 failures), primarily due to “atomic” steps (e.g., simple equations like  $3 \times x = 30$ ) that are inherently indivisible, leading the model to either hallucinate details or mistakenly include future reasoning steps. Despite these edge cases, the vast majority of generated steps correctly introduced the intended inefficiency without altering the ground truth logic.

### 3.2 CAID: Context-Aware Information Density

Current evaluators primarily focus on factual validity, often overlooking the inefficiencies modeled in RIV-GSM8K. To bridge this gap, we propose **CAID**, a reference-free, unsupervised metric designed to quantify the *informational utility* of a reasoning step. Rather than relying on heuristics tailored to specific error types, CAID is grounded in **information theory**, evaluating steps based on their contribution to the reasoning process relative to their length and context. CAID integrates four core indicators, as conceptualized in Figure 2:

**1. Local Similarity (Redundancy).** To capture redundancy at the surface and semantic levels, we measure the cosine similarity between the current step  $S_t$  and its immediate predecessor  $S_{t-1}$  using a lightweight encoder  $\mathcal{E}$  (e.g., MiniLM) as  $\mathcal{M}_{sim}(S_t) = \text{CosSim}(\mathcal{E}(S_t), \mathcal{E}(S_{t-1}))$ . A high similarity score indicates that the step provides negligible new information compared to the immediate history, signaling potential repetition or inefficient paraphrasing.

**2. Global Goal Alignment (Relevance).** Reasoning steps must remain relevant to the problem objective. We measure the alignment between the step and the original question  $Q$  by computing  $\mathcal{M}_{rel}(S_t) = \text{CosSim}(\mathcal{E}(S_t), \mathcal{E}(Q))$ . A significantly low alignment score suggests that the step has drifted from the core problem goal, identifying content that may be linguistically coherent but contextually irrelevant.

**3. Information Density (Verbosity).** Efficient reasoning should convey maximum information with minimum tokens. We define density as the length-normalized perplexity using a causal language model  $\mathcal{M}$  (e.g., GPT-2):

$$\mathcal{M}_{density}(S_t) = \frac{\log(\text{PPL}_{\mathcal{M}}(S_t))}{\text{Length}(S_t)} \quad (1)$$

319	Steps with exceptionally low density imply a lack	• <b>PRUNE:</b> Removes steps flagged as high re-	367
320	of information content relative to their verbosity.	dundancy ( $\mathcal{M}_{sim}$ ) or low relevance ( $\mathcal{M}_{rel}$ ).	368
321	This metric effectively penalizes steps that ex-	• <b>MERGE:</b> Compresses steps exhibiting <i>latent</i>	369
322	cessively decompose simple logic into multiple,	<i>inefficiency</i> (valid but verbose/fragmented) us-	370
323	low-entropy micro-steps, detecting “informational	ing an LLM re-writer. To prevent semantic	371
324	froth.”	drift or information overload, we enforce two	372
325	<b>4. Semantic Delta (Logical Progress).</b> Cru-	safety constraints before merging step $S_t$ into	373
326	cially, a valid reasoning step must advance the log-	the accumulated step $S'_{last}$ :	374
327	ical state toward the solution. We define Semantic	1. <b>Semantic Consistency:</b>	375
328	Delta as the shift in the context vector induced by	$\text{CosSim}(\mathcal{E}(S'_{last}), \mathcal{E}(S_t)) \geq \tau_{merge}$ .	376
329	adding $S_t$ :	We ensure the new content is logically	377
	$\mathcal{M}_{delta}(S_t) = 1 - \text{CosSim}(\mathcal{E}(C_{t-1}), \mathcal{E}(C_{t-1} \oplus S_t))$	compatible with the current context.	378
330	(2)	2. <b>Information Saturation:</b> $\mathcal{I}(S'_{last}) \leq$	379
331	A near-zero delta implies that $S_t$ fails to update	$\tau_{max}$ . We prevent merging if the cur-	380
332	the semantic state of the context, characterizing	rent step is already information-dense,	381
333	tautologies or circular reasoning. As reasoning	avoiding readability loss.	382
334	naturally converges toward the final answer, the	If constraints are violated, the merge is halted,	383
335	marginal information gain of each subsequent step	and a new step is initiated.	384
336	often diminishes. To address this, we employ an	• <b>KEEP:</b> Retains steps essential for logical	385
337	<b>adaptive decaying threshold</b> $\tau_{\delta}(t) = \tau_{base} \cdot \lambda^t$ ,	progress.	386
338	which dynamically adjusts sensitivity based on the		
339	step position.	<b>Addressing the “Trivial Accuracy” Concern.</b>	387
340	<b>Decision Logic.</b> CAID aggregates these indi-	One might assume that maintaining accuracy is triv-	388
341	cators to classify steps into an action set $\mathcal{A} =$	ial if the final conclusion step is preserved. How-	389
342	$\{\text{PRUNE}, \text{MERGE}, \text{KEEP}\}$ via a hierarchical deci-	ever, in our evaluation, we construct the prompt	390
343	sion process. Specifically, steps exhibiting high	using the compressed chain $C'$ <b>excluding the final</b>	391
344	redundancy ( $\mathcal{M}_{sim}$ ) or low relevance ( $\mathcal{M}_{rel}$ ) are	<b>numerical answer</b> (i.e., the “### Result” token).	392
345	considered strictly unnecessary and flagged for	The model is required to <i>regenerate</i> the final answer	393
346	<b>PRUNE</b> . In contrast, steps that are relevant but	solely based on the logic provided in $C'$ . Since rea-	394
347	exhibit low information density ( $\mathcal{M}_{density}$ ) or	soning chains are causal, removing or altering a	395
348	marginal logical progress ( $\mathcal{M}_{delta}$ ) are targeted	necessary intermediate step would break the log-	396
349	for <b>MERGE</b> , an action that preserves the underly-	ical dependency required to derive the correct so-	397
350	ing logic while condensing verbose or fragmented	lution. Therefore, the fact that PACE maintains	398
351	expressions. This multi-view approach ensures ro-	accuracy with significantly fewer tokens serves as	399
352	burst optimization, selectively removing distractions	robust evidence that the removed <b>or compressed</b>	400
353	while refining the density of valid reasoning.	steps were indeed functionally redundant and that	401
354	<b>3.3 Application: Validating CAID via PACE</b>	CAID correctly identified the core reasoning path.	402
355	To empirically validate the diagnostic precision of	<b>4 Experiments</b>	403
356	CAID, we introduce <b>PACE</b> (Pruning And Com-	<b>4.1 Experimental Setup</b>	404
357	pression for Efficiency) as a <b>post-hoc compression</b>	<b>Datasets.</b> We employ diverse benchmarks to con-	405
358	<b>strategy</b> . Our primary objective here is not	duct a two-stage evaluation, assessing both diag-	406
359	to accelerate real-time inference, but to serve as	nostic sensitivity and practical compression utility.	407
360	a diagnostic probe: demonstrating that the steps	• <b>RIV-GSM8K:</b> A controlled diagnostic set	408
361	flagged by CAID are indeed dispensable “froth”	used to measure the sensitivity of metrics to	409
362	that can be distilled without breaking the deductive	explicit, synthetically injected inefficiencies	410
363	chain.	(Section 3.1).	411
364	PACE operates in a <i>Generate-then-Refine</i>		
365	pipeline. Based on the classification from CAID,		
366	we apply three actions:		

Model	simple_duplication		paraphrase		decompose		circular_reasoning		irrelevant	
	Aug PR (↓)	Gold PR	Aug PR (↓)	Gold PR	Aug PR (↓)	Gold PR	Aug PR (↓)	Gold PR	Aug PR (↓)	Gold PR
ReasonEval 7B	0.6555	<b>0.9897</b>	0.7338	<b>0.9867</b>	0.7917	<b>0.9853</b>	0.4809	<b>0.9746</b>	<u>0.1509</u>	<b>0.9571</b>
ReasonEval 34B	0.7779	0.9558	0.7251	<u>0.9533</u>	0.7604	<u>0.9571</u>	<u>0.3762</u>	0.9345	<b>0.0331</b>	0.9118
ThinkPRM 1.5B	<u>0.6264</u>	0.7342	<u>0.6821</u>	0.7486	<u>0.7187</u>	0.7663	0.7179	0.7810	0.6667	0.7353
ThinkPRM 7B	0.6849	0.8003	0.7834	0.8149	0.7326	0.8046	0.7619	0.8881	0.7999	0.8682
ThinkPRM 14B	0.8582	0.9142	0.8859	0.9140	0.8150	0.8706	0.9118	0.9274	0.8474	0.9188
Qwen2.5-Math-PRM-7B	0.9679	<u>0.9606</u>	0.9512	0.9521	0.9350	0.9433	0.8614	<u>0.9562</u>	0.9746	<u>0.9512</u>
Qwen2.5-Math-PRM-72B	0.8368	0.9517	0.8896	0.9457	0.9108	0.9502	0.8703	0.9554	0.8961	0.9504
CAID (Ours)	<b>0.0000</b>	0.5752	<b>0.0174</b>	0.5596	<b>0.2006</b>	0.5142	<b>0.0190</b>	0.4799	0.1598	0.5155

Table 1: Step Preservation Rate (SPR) comparison by augmentation type. Aug PR: Augmented Step Preservation Rate (lower is better), Gold PR: Gold Step Preservation Rate (higher indicates retention). **Bold** indicates the best performance, and underlined indicates the second-best performance.

- **Standard Benchmarks:** To validate PACE on real-world reasoning, we use **GSM8K** (Cobbe et al., 2021), **StrategyQA** (Geva et al., 2021), and **ARC-Challenge** (Clark et al., 2018). These datasets cover arithmetic, commonsense, and scientific reasoning, respectively, allowing us to test whether CAID generalizes across different domains without task-specific tuning.

**Baselines.** For the diagnostic comparison on RIV-GSM8K, we evaluate state-of-the-art PRMs including **ReasonEval-34B** (Xia et al., 2025), **ThinkPRM-7B** (Khalifa et al., 2025), and **Qwen-Math-PRM-72B** (Yang et al., 2024). For the compression validation via PACE, we use **Llama-3.1-8B-Instruct** as the backbone model and compare the compressed chains against the standard **Zero-shot CoT** baseline to measure the trade-off between token reduction and accuracy.

**Implementation of CAID.** We utilize lightweight off-the-shelf models for efficiency: **all-MiniLM-L6-v2** (22M) for semantic encoding and **GPT-2 Small** (124M) for density estimation. To ensure robustness, we use a fixed set of hyperparameters across all datasets without task-specific tuning. Detailed model configurations, threshold values, and sensitivity analysis are provided in Appendix D.

## 4.2 Results 1: Diagnostic Capability on RIV-GSM8K

We evaluate how well different evaluators handle the inefficiencies injected into RIV-GSM8K. Instead of binary accuracy, we report the **Step Preservation Rate (SPR)**, defined as the ratio of steps retained after evaluation.

- **Augmented PR (Aug PR):** Measures recall on inefficient steps. **Lower is better**, indicat-

ing the model successfully removed or flagged the inefficient step.

- **Gold PR:** Measures retention of original human-written steps. **Higher typically indicates safety**, assuming human steps are perfectly efficient. However, as discussed below, we challenge this assumption.

**Blind Spot of Validity-Focused Evaluators.** As shown in Table 1, existing methods exhibit surprisingly high Aug PR scores. A critical finding is the performance of **ReasonEval**. Despite being a specialized reasoning step evaluator equipped with an explicit *redundancy score*, it fails to effectively penalize redundant steps, retaining approximately 70% of *Simple Duplication*, *Paraphrase*, and *Decompose* types. Similarly, even the 72B-parameter Qwen-Math preserves 83.68% of *Simple Duplications*. This confirms that validity-focused models, regardless of their size or specific scoring sub-metrics, remain essentially blind to inefficiency as long as the statement is factually correct.

**Effectiveness and Efficiency of CAID.** In contrast, CAID achieves near-perfect detection on redundancy, with an Aug PR of **0.0000** for Duplication and **0.0174** for Paraphrasing. It also effectively identifies complex inefficiencies like Circular Reasoning (0.0190) and Decomposition (0.2006), where baselines struggle significantly. Regarding *Irrelevant* steps, while the large-scale supervised model **ReasonEval-34B** achieves the best performance (0.0331), CAID (0.1598) demonstrates competitive capability, performing comparably to **ReasonEval-7B** (0.1509). Crucially, CAID achieves this with a total of only **146M parameters** (124M GPT-2 + 22M MiniLM), whereas ReasonEval requires 7B to 34B parameters. This demonstrates that CAID delivers robust diagnostic

precision with orders of magnitude greater computational efficiency than large-scale supervised evaluators.

### Redefining “Gold”: Deletion vs. Compression.

A distinct characteristic of CAID is its lower Gold PR ( $\approx 0.55$ ) compared to baselines ( $> 0.90$ ). While this might initially appear as over-penalization, a granular analysis of the action distribution reveals that CAID is not “wrong,” but rather stricter regarding information density. Out of 11,899 Gold steps not fully preserved by CAID:

- Only **1.5% (184 steps)** were flagged for removal (PRUNE), primarily due to high redundancy (169 steps) or irrelevance (15 steps).
- The remaining **98.5% ( $\approx 11,700$  steps)** were flagged for MERGE.

As qualitatively analyzed in Appendix E (Table 5), these flagged steps are factually valid but functionally inefficient. For instance, steps that merely restate a calculated value (e.g., “Child = 4”  $\rightarrow$  “Ticket is \$4”) trigger the **Low Semantic Delta** criteria due to a lack of logical progress. Similarly, steps that verbally describe an operation before executing it (e.g., “Then multiply the number..”) are flagged for **Low Information Density**. This empirical evidence suggests that standard datasets contain significant **latent inefficiency**, validating our approach of *compression* over blind retention.

### 4.3 Results 2: Efficiency via PACE

To empirically validate the diagnostic precision of CAID, we applied PACE as a compression probe to reasoning chains generated by Llama-3.1-8B. Table 2 summarizes the trade-off between token reduction and reasoning accuracy.

Dataset	Method	Performance		Efficiency	
		Acc (%)	$\Delta$	Tokens	Red (%)
GSM8K	Baseline	82.03	-	214.6	-
	PACE	81.12	-0.91	<b>148.0</b>	<b>-31.0%</b>
StrategyQA	Baseline	70.31	-	327.7	-
	PACE	69.93	-0.37	<b>154.4</b>	<b>-52.9%</b>
ARC-C	Baseline	83.70	-	277.8	-
	PACE	<b>84.64</b>	<b>+0.94</b>	<b>156.7</b>	<b>-43.6%</b>

Table 2: Comparison of accuracy and token usage between baseline CoT and PACE. PACE significantly reduces tokens while maintaining or improving accuracy.

### Validating Latent Inefficiency via Compression.

PACE successfully reduces token consumption by 31.0% to 52.9% across all datasets with negligible

impact on accuracy ( $< 1\%$ ). This empirical finding serves as strong evidence for the prevalence of **latent inefficiency** in standard CoT reasoning. It confirms that a significant portion of generated tokens acts as informational “froth”—valid but functionally unnecessary—and that CAID correctly identifies this redundancy without disrupting the deductive chain.

### Enhancing Clarity in Information-Heavy Domains.

Interestingly, on ARC-Challenge, PACE actually improves accuracy by +0.94%p while reducing tokens by 43.6%. Unlike arithmetic tasks where steps are strictly sequential, scientific reasoning often suffers from hallucinated tangents or excessive context. By pruning these distractions (via Goal Alignment) and condensing verbose explanations (via Density), PACE effectively increases the **signal-to-noise ratio** of the reasoning context, helping the model focus on critical scientific principles.

### 4.4 Ablation Study

To investigate the contribution of each information-theoretic component, we conducted a cumulative ablation study on GSM8K. Figure 3 visualizes the performance gain as we integrate each metric.

**Necessity of Semantic Delta.** As shown in the figure, reliance on *Local Similarity* alone effectively detects lexical repetitions but fails completely on *Decomposition* and *Irrelevance* (near 0% detection). Adding *Semantic Delta* triggers a sharp performance boost, confirming that measuring the “vector shift” or logical velocity is essential for identifying stalled reasoning.

**Role of Information Density.** *Information Density* proves critical for detecting *Circular Reasoning*. Tautological statements often exhibit low perplexity (high probability), which CAID successfully flags as low-density steps. This confirms our hypothesis that verbosity is not just about raw length, but about information content relative to length (entropy rate).

**Impact of Compression Strategy.** We further evaluated the effectiveness of different compression actions and safety constraints. Our results show that simply deleting steps with low logical progress causes a sharp accuracy drop (-5.38%), implying that these steps serve as essential “connective tissue.” In contrast, PACE’s merging strategy, reinforced by safety constraints (Consistency and

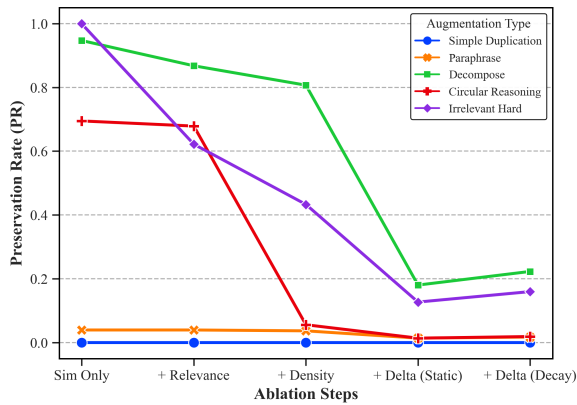


Figure 3: Impact of cumulative CAID components on Preservation Rate (PR). Lower PR indicates better detection. The integration of **Information Density** and **Semantic Delta** proves critical for detecting complex inefficiencies (e.g., *Circular Reasoning*, *Decomposition*) that bypass surface-level similarity.

Saturation), successfully recovers accuracy while maintaining high efficiency. Detailed ablation results and analysis are provided in Appendix F.

## 5 Discussion

**The Phenomenon of Latent Inefficiency.** A significant revelation in our experiments is that CAID identifies approximately 50% of human-written (Golden) steps as **compressible** rather than strictly necessary. As detailed in Section 4.2, the vast majority (98.5%) of these flagged steps were categorized for MERGE due to **low information density** or **minimal logical progress**, with only 1.5% identified as redundant or irrelevant enough for removal. This points to a prevalent phenomenon of *Latent Inefficiency* in current CoT datasets: valid reasoning is often diluted by excessive verbosity or fragmented micro-steps. It challenges the prevailing assumption that human-written reasoning is the “gold standard” for efficiency, highlighting the need for metrics like CAID to prioritize **information density** alongside correctness.

**Beyond Post-hoc: Implications for Data-Centric AI.** While PACE is presented here as a diagnostic compression strategy, its capability to distill informational “froth” offers substantial value in two key data-centric applications:

- **Efficiency-Aware Instruction Tuning:** PACE can serve as a high-quality filtration pipeline for Data-Centric AI. By constructing **efficiency-aware training datasets** from dense reasoning chains, we can potentially fine-tune smaller models (Students) to reason

efficiently from scratch. This effectively transfers the information density inherent in CAID to the model weights, eliminating the need for runtime post-processing.

- **RAG Context Compression:** In Retrieval-Augmented Generation (RAG) systems, retrieved reasoning paths often consume significant memory. PACE can pre-compress these chains offline, allowing systems to retrieve and process more relevant context within the same token budget, thereby enhancing multi-turn reasoning efficiency via ‘Context Reuse’.

## 6 Conclusion

In this work, we investigated the critical inefficiency of “over-reasoning” in Large Language Models, diagnosing that current models often generate valid but redundant, verbose, or irrelevant steps that inflate computational costs without contributing to the solution. While existing reasoning step evaluators excel at verifying correctness, we demonstrated their inability to detect these subtle forms of inefficiency due to a lack of information-theoretic criteria.

To address this, we proposed a comprehensive framework comprising diagnosis, measurement, and validation. First, we introduced **RIV-GSM8K**, a diagnostic benchmark that exposed the “efficiency blind spots” of state-of-the-art evaluators. Second, we developed **CAID**, a reference-free metric that quantifies the informational utility of reasoning steps by integrating local novelty, information density, and global goal alignment. Finally, we validated the diagnostic precision of this metric through **PACE**, a post-hoc compression strategy.

Our experiments confirmed that PACE achieves substantial token reductions of 31–53% across arithmetic, commonsense, and scientific domains while maintaining reasoning accuracy. This empirically validates the existence of significant “latent inefficiency” in standard CoT reasoning, proving that a large portion of tokens represents informational “froth” that can be safely compressed without breaking the deductive chain. By shifting the paradigm of reasoning evaluation from “Correctness-only” to “**Correctness-and-Efficiency**,” this work establishes a foundation for building more sustainable reasoning systems and constructing efficiency-aware training datasets.

652	<b>Limitations</b>		
653	While our framework effectively diagnoses and	ensure that valid dialectal variations are not unfairly	699
654	mitigates reasoning inefficiency, we acknowledge	classified as redundant “froth.”	700
655	several limitations.		
656	<b>Computational Overhead of Post-hoc Process.</b>	<b>Use of AI Assistants.</b> We utilized GPT-4o as a	701
657	PACE operates as a generate-then-refine strategy,	component of our data augmentation pipeline to	702
658	meaning it does not reduce the latency of the initial	generate synthetic reasoning inefficiencies for the	703
659	inference pass. Consequently, it is best suited	RIV-GSM8K benchmark, as detailed in Section	704
660	for offline applications—such as compressing re-	3.1 and Appendix A. Additionally, AI assistants	705
661	trieved context for RAG or synthesizing efficiency-	were used for preliminary code implementation	706
662	aware training data—rather than real-time latency	and linguistic polishing of the manuscript. The	707
663	reduction.	authors have proofread all AI-generated content	708
664		and remain responsible for the final output.	709
665	<b>Dependency on Writer Capability.</b> The MERGE		
666	action relies on the semantic capability of the under-	<b>References</b>	710
667	lying LLM to rewrite verbose steps without infor-	Lingjiao Chen, Matei Zaharia, and James Zou. 2024.	711
668	mation loss. While our safety constraints ( $\tau_{merge}$ ,	<a href="#">FrugalGPT: How to use large language models while</a>	712
669	$\tau_{sat}$ ) mitigate semantic drift, there remains a resid-	<a href="#">reducing cost and improving performance</a> . <i>Trans-</i>	713
670	ual risk that smaller models may over-simplify com-	<i>actions on Machine Learning Research</i> . Featured	714
671	plex logic during compression.	Certification.	715
672		Cheng-Han Chiang and Hung-yi Lee. 2024. <a href="#">Over-</a>	716
673	<b>Domain Generalizability.</b> Our experiments fo-	<a href="#">reasoning and redundant calculation of large lan-</a>	717
674	cused on arithmetic (GSM8K), commonsense	<a href="#">guage models</a> . In <i>Proceedings of the 18th Confer-</i>	718
675	(StrategyQA), and scientific reasoning (ARC-C).	<i>ence of the European Chapter of the Association for</i>	719
676	The definition of “efficiency” may differ in creative	<i>Computational Linguistics (Volume 2: Short Papers)</i> ,	720
677	or open-ended domains (e.g., storytelling) where	pages 161–169, St. Julian’s, Malta. Association for	721
678	verbosity serves a stylistic purpose. Future work is	Computational Linguistics.	722
679	required to adapt information-theoretic criteria for	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,	723
680	such subjective tasks.	Ashish Sabharwal, Carissa Schoenick, and Oyvind	724
681		Tafjord. 2018. <a href="#">Think you have solved question an-</a>	725
682	<b>Hyperparameter Sensitivity.</b> Although CAID	<a href="#">swering? try arc, the AI2 reasoning challenge</a> . <i>CoRR</i> ,	726
683	demonstrated robustness across our benchmarks	abs/1803.05457.	727
684	with a fixed set of thresholds, applying the metric	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,	728
685	to domains with vastly different linguistic densities	Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias	729
686	(e.g., code generation or legal texts) may require	Plappert, Jerry Tworek, Jacob Hilton, Reiichiro	730
687	task-specific hyperparameter tuning.	Nakano, Christopher Hesse, and John Schulman.	731
688		2021. <a href="#">Training verifiers to solve math word prob-</a>	732
689	<b>Ethics Statement</b>	<a href="#">lems</a> . <i>Preprint</i> , arXiv:2110.14168.	733
690	This work aligns with the goals of <b>Green AI</b> by	Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot,	734
691	promoting sustainable computation through the re-	Dan Roth, and Jonathan Berant. 2021. Did Aristotle	735
692	duction of token consumption in LLMs.	Use a Laptop? A Question Answering Bench-	736
693		mark with Implicit Reasoning Strategies. <i>Trans-</i>	737
694	<b>Data Usage and Privacy.</b> Our research utilizes	<i>actions of the Association for Computational Linguis-</i>	738
695	publicly available datasets and does not involve	<i>tics (ACL)</i> .	739
696	private or personally identifiable information. The	Olga Golovneva, Moya Peng Chen, Spencer Poff, Mar-	740
697	synthetic perturbations in RIV-GSM8K were gen-	tin Corredor, Luke Zettlemoyer, Maryam Fazel-	741
698	erated using standard LLMs (GPT-4o) and do not	Zarandi, and Asli Celikyilmaz. 2023. <a href="#">ROSCOE: A</a>	742
	contain harmful content.	<a href="#">suite of metrics for scoring step-by-step reasoning</a> . In	743
		<i>The Eleventh International Conference on Learning</i>	744
		<i>Representations</i> .	745
		Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing	746
		Yang, and Lili Qiu. 2023. <a href="#">LLMLingua: Compressing</a>	747
		<a href="#">prompts for accelerated inference of large language</a>	748
		<a href="#">models</a> . In <i>The 2023 Conference on Empirical Meth-</i>	749
		<i>ods in Natural Language Processing</i> .	750



### Task 1: Paraphrase

#### System Prompt:

You are an expert paraphrasing assistant. Your task is to rephrase the given text while preserving its original meaning. Ensure that the paraphrased text is clear, coherent, and maintains the same intent as the original. Avoid using overly complex language or altering the tone significantly. Crucially, wrap all calculations in «calculation=result» format.

#### User Input Template:

"Please paraphrase the following text: [INPUT\_TEXT]"

### Task 2: Decompose

#### System Prompt:

You are an expert at breaking down complex tasks into smaller, manageable, verbose subtasks. Your goal is to analyze the given task description and decompose it into a list of sequential steps that can be followed to accomplish the task effectively. Ensure that decomposed steps maintain same intent as the original. Avoid using overly complex language or altering the tone significantly. Crucially, wrap all calculations in «calculation=result» format.

#### User Input Template:

"Please decompose the following task into smaller subtasks: [INPUT\_TASK]"

Figure 4: Configuration and prompt details for basic augmentation strategies (*Paraphrase* and *Decompose*). The top block shows shared hyperparameters.

tency, we employ a **temperature of 1.0** and a **maximum token limit of 5,000**. Additionally, we enforce a **strict JSON Schema** for all outputs to facilitate robust parsing.

Below, we detail the specific system prompts and input templates used for our pipeline. Figure 4 outlines the basic augmentation types, while Figure 5 details the context-aware strategies.

## C Qualitative Examples of Reasoning Augmentations

To better understand the nature of the perturbations introduced by our pipeline, we provide concrete examples of generated reasoning steps in Table 4. These examples are derived from a single original step in the GSM8K dataset: "*Adults = 10 \* 8 = 80*".

As shown in Table 4, our augmentation strategies cover a spectrum of “over-reasoning” behaviors observed in Large Language Models:

- **Surface-level Redundancy:** *Simple Duplication* and *Paraphrase* retain the exact logic of the original step but introduce lexical variations or repetitions, testing the model’s robustness to verbose phrasing.
- **Granularity Expansion:** *Decompose* breaks down a single atomic operation into a verbose

chain of micro-steps (e.g., identifying variables, stating the operation, calculating, and restating the result), significantly inflating the token count without adding deductive value.

- **Logical Loops:** *Circular Reasoning* mimics a model’s tendency to “double-check” itself unnecessarily. It uses inverse operations (division/multiplication) to verify an already established fact, creating a closed logical loop that adds computational cost.
- **Contextual Noise:** *Irrelevant* introduces distractors that differ from hallucinations; they are mathematically true and contextually plausible (e.g., discussing ticket prices or family groups) but contribute nothing to the solution path.

These qualitative samples illustrate the diverse challenges our dataset poses to reasoning evaluators.

## D Implementation Details

### D.1 Model Configuration

CAID is designed to be computationally efficient and widely applicable. We employ the following off-the-shelf models:

- **Semantic Encoder ( $\mathcal{E}$ ):** We use all-MiniLM-L6-v2 (22M parameters)

### Task 3: Circular Reasoning

#### System Prompt:

You are an expert at inserting circular reasoning into mathematical solutions. Your task is to generate a sequence of steps that redundantly verifies a previously established fact or calculated number using inverse operations or self-referential logic.

You will be provided with:

1. The Question
2. Previous Reasoning Steps
3. The Current Target Step

Generate a reasoning section that:

- Takes a number or fact already established in the 'Previous Reasoning Steps'.
- Performs a set of operations that eventually lead back to the original number (e.g., "Since X is 5, multiplying by 2 gives 10, and dividing by 2 returns 5, confirming X is indeed 5.").
- Is mathematically true but strictly unnecessary for solving the problem.
- Does not alter the final answer or the logical path required for the solution.

Crucially, wrap all calculations in «calculation=result» format.

#### User Input Template:

"Based on the context below, generate circular reasoning sentences that could be inserted after the Current Step:"

### Task 4: Irrelevant (Hard)

#### System Prompt:

You are an expert at generating context-aware distractions. Your task is to generate mathematically correct but irrelevant sentences that sound like they belong to the solution flow but do not advance the solution logic or provide any new information needed for the answer.

You will be provided with:

1. The Question
2. Previous Reasoning Steps
3. The Current Target Step
4. Next Reasoning Steps

Generate a reasoning section that:

- naturally fits between the previous reasoning, the current target step, and the next reasoning steps,
- maintains the same tone, context, and mathematical domain,
- uses a smooth transitional phrase to connect the surrounding steps,
- is mathematically true but does not contribute to solving the problem,
- does not alter any variables, numbers, or assumptions in the reasoning,
- and does not suggest new solution paths or constraints.

Crucially, wrap all calculations in «calculation=result» format if any numbers appear.

#### User Input Template:

"Based on the context below, generate irrelevant sentences that could be inserted after the Current Step:"

Figure 5: Prompt configurations for context-aware augmentation types (*Circular Reasoning* and *Irrelevant*). Note that these tasks require full context inputs (question, previous/next reasoning steps).

884 to compute cosine similarity for Local Sim-  
885 ilarity ( $\mathcal{M}_{sim}$ ) and Global Goal Alignment  
886 ( $\mathcal{M}_{rel}$ ). This model was selected for its high  
887 speed and strong performance on semantic  
888 textual similarity tasks.

• **Density Estimator ( $\mathcal{M}$ ):** We use GPT-2 889  
Small (124M parameters) to calculate the per- 890  
plexity for Information Density ( $\mathcal{M}_{density}$ ). 891  
The total parameter count for CAID is approxi- 892  
mately 146M, which is significantly smaller than 893  
the baseline PRMs (e.g., ReasonEval-34B). 894

Augmentation Type	Generated Reasoning Step
Original Input	Adults = 10 * 8 = \$«10*8=80»80
1. Simple Duplication	Adults = 10 * 8 = \$«10*8=80»80
2. Paraphrase	Adults: 10 multiplied by 8 equals \$«10*8=80»80
3. Decompose	Identify the number of adults involved, which is given as 10. → Identify the multiplication factor for the adults, which is 8 in this case. → Multiply the number of adults (10) by the factor (8): «calculation=80». → State the result of the multiplication, which is \$80.
4. Circular Reasoning	Having calculated \$80 for the adults, we can verify this by noting that \$80 divided by the 10 adults gives us \$«80/10=8»8 per adult. → Re-multiplying this \$8 per adult by 10 adults returns \$«8*10=80»80, ensuring our number is correct. → Alternatively, if we consider twice the accounted \$80, we get \$«80*2=160»160, and half of this sum is \$«160/2=80»80, which confirms our original total remains consistent at \$80.
5. Irrelevant	Since the adults paid a total of \$80, this also means they chose to spend their money on tickets rather than concessions at the game, an interesting financial decision given the usual pricing. → It’s notable that the sum of \$124 offers no remainder when considering integer ticket prices, suggesting an elegance in the problem structure. → The relationship between the number of adults and children indicates a typical family group could have attended this game together.

Table 4: Examples of augmented reasoning steps generated from a single original step. The generated text aims to mimic specific reasoning flaws or stylistic variations.

## D.2 Hyperparameters

We utilize a fixed set of thresholds across all experiments (GSM8K, StrategyQA, ARC-Challenge) to demonstrate the generalizability of our metric. The specific values are:

- **Removal Thresholds (PRUNE):**
  - High Redundancy:  $\tau_{sim} = 0.85$
  - Low Relevance:  $\tau_{rel} = 0.25$
- **Compression Candidates (MERGE):**
  - Low Information Density:  $\tau_{density} = 0.1$
  - Low Semantic Delta (Base):  $\tau_{delta} = 0.03$
- **Adaptive Decay:**
  - Decay Factor:  $\lambda = 0.95$  (Applied as  $\tau_{\delta}(t) = \tau_{delta} \cdot \lambda^t$ )

**Sensitivity Analysis.** We observed that the performance of CAID is relatively stable around these threshold values. For instance, varying  $\tau_{sim}$  between 0.80 and 0.90 or  $\tau_{rel}$  between 0.20 and 0.30 resulted in minimal fluctuations in the Step Preservation Rate (SPR) on the RIV-GSM8K validation set. This suggests that the chosen hyperparameters are robust and not overfitted to a specific dataset distribution.

## E Qualitative Analysis of Latent Inefficiency

To better understand the nature of “Latent Inefficiency” in human-written Gold data, we provide a detailed qualitative analysis of steps flagged for MERGE by CAID. Table 5 presents concrete examples from the GSM8K dataset.

Previous Step ( $S_{t-1}$ )	Target Step ( $S_t$ ) [Gold]	Reason
Child = 44/11 = \$«44/11=4»4	Each child’s ticket is \$«4=4»4.	<b>Low Delta</b>
→ <i>Diagnosis:</i> The target step merely repeats the value ‘4’ established in the previous step, contributing no new deductive information (No Progress).		
...when self-checkout is broken... 160 * 1.2 = 192 complaints/day	Then multiply the number of complaints per day by the number of days...: 192 * 3 = 576...	<b>Low Density</b>
→ <i>Diagnosis:</i> The step explicitly describes the operation before performing it, inflating token usage (Verbose).		

Table 5: Qualitative examples of Gold steps flagged for MERGE. By utilizing the full width for diagnosis, we clarify why valid steps are identified as inefficient (e.g., lack of progress or verbosity).

## F Detailed Ablation on Compression Strategy

In this section, we provide the extended ablation study on the compression strategies employed in

ID	Method Description	Performance		Efficiency		Ratio (Tok)
		Acc (%)	$\Delta$	Tok Red (%)	Step Red (%)	
0	Baseline (Original CoT)	82.03	0.00	0.00	0.00	1.00
1	+ Similarity (Remove)	84.46	+2.43	7.36	15.47	1.08
2	+ Relevance (Remove)	84.38	+2.35	8.91	17.06	1.10
3	+ Density (Remove)	<b>84.53</b>	<b>+2.50</b>	20.12	21.69	1.25
4	+ Delta (Remove)	76.65	-5.38	11.18	<b>72.52</b>	1.13
5	+ Merge (No Safety)	76.95	-5.08	<b>45.09</b>	70.41	<b>1.82</b>
6	<b>PACE (Full Method)</b>	81.12	-0.91	31.05	53.53	1.45

Table 6: Ablation study of PACE components. We analyze the impact of each module on accuracy and compression efficiency. **Modes 1–4** use removal-only logic, while **Modes 5–6** introduce the merging mechanism. **PACE (Mode 6)** achieves the best balance between accuracy recovery and token reduction.

931 PACE, validating the necessity of the MERGE ac-  
932 tion and safety constraints.

933 **Deletion vs. Compression.** As shown in Ta-  
934 ble 6, while removing redundant steps (Modes 1–  
935 3) yields slight accuracy gains, simply deleting  
936 steps with low logical progress (Mode 4, Delta)  
937 causes a sharp accuracy drop (-5.38%). This im-  
938 plies that even repetitive or slow-progressing steps  
939 serve as essential “connective tissue” in the reason-  
940 ing chain, carrying implicit dependencies. They  
941 cannot be blindly removed (Prune) but must be  
942 **merged** to preserve logical continuity while reduc-  
943 ing verbosity.

944 **Necessity of Safety Constraints.** Mode 5  
945 (Merge without constraints) achieves high token  
946 reduction (45%) but suffers significant accuracy  
947 degradation (-5.08%) due to semantic drift and  
948 information overload. By enforcing our safety  
949 constraints (Consistency and Saturation), **PACE**  
950 **(Mode 6)** successfully recovers the accuracy (Acc  
951 81.12%) while still delivering substantial efficiency  
952 (31% token reduction), demonstrating that our  
953 density-aware merging strategy achieves the op-  
954 timal trade-off between compression and validity.