

Quantifying Tolerance to Errors in Synthetic Data: An Atomic-level Operand vs. Operator Perturbation Study

Anonymous ACL submission

Abstract

Synthetic data generation has become a cornerstone for advancing large language models. However, the absence of a quantitative analysis for error tolerance remains a critical bottleneck. Consequently, current filtering strategies fluctuate between two extremes: they are either overly aggressive, risking the exclusion of potentially valuable samples, or overly permissive, failing to eliminate erroneous samples effectively. To bridge this gap, we introduce **Atomic Tree Operation Modeling (ATOM)**, a framework that decomposes data into functional units ($y = f(x)$) to precisely differentiate between *benign Operand perturbations* and *fatal Operator perturbations*. Our experiments reveal a **double dissociation**: models are robust to operand noise but collapse under operator disruption. By prioritizing operator over strict operand precision, our ATOM-synthesized data significantly outperforms rigorous baselines (e.g., +3.3% gain over LIMA), validating that structural diversity is the decisive factor for synthetic data.

1 Introduction

The ability of Large Language Models (LLMs) heavily relies on training data (Grattafiori et al., 2024; OpenAI et al., 2024; Qwen-Team, 2024; DeepSeek-AI et al., 2025; Yang et al., 2025a; Gemini-Team et al., 2025). As high-quality human-generated data becomes increasingly scarce, the field has largely adopted synthetic data generation (Taori et al., 2023; Zhou et al., 2023; Wang et al., 2023b; Xu et al., 2024, 2025; Maosong et al., 2025) to distill capabilities from stronger models. However, a major bottleneck remains: teacher models inevitably hallucinate, injecting factual errors into the training signals.

Studies such as LIMA (Zhou et al., 2023) underscored that data quality is a decisive factor in supervised fine-tuning (SFT). However, this conclusion lacks a theoretical framework to quantita-

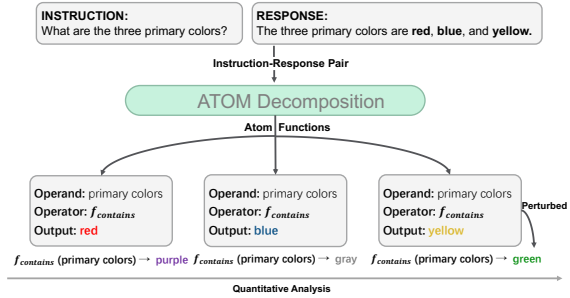


Figure 1: The illustration of extract atomic functions.

tively define this metric. Consequently, existing data filtering methodologies fluctuate between two extremes: they are either *overly aggressive* (Zhou et al., 2023; Xu et al., 2024; Maosong et al., 2025), risking the exclusion of potentially valuable samples, or *overly permissive* (Li et al., 2024; Zhao et al., 2024; Li et al., 2025), failing to eliminate erroneous samples effectively. This dilemma stems from the inability of current methods to strictly quantify the impact of data defects. Analytically, existing works are limited to either coarse, qualitative assessments that neglect error quantification (Cho, 2024a; Raghavendra et al., 2025), or linguistic surface taxonomies that overlook essential logical properties (Alajrami et al., 2025). **Ultimately, these research gaps leave the boundary of error tolerance undefined.**

To rigorously identify the boundaries of error tolerance, we propose the **Atomic Tree Operation Modeling (ATOM)** framework, grounded in classic Problem Space Theory (Newell, 1972). As illustrated in Figure 1, ATOM first decomposes instruction-response pairs into atomic operations, denoted as $\mathcal{A} = \langle x, f, y \rangle$, such as $f_{\text{contains}}(\text{primary colors}) \rightarrow \text{red/blue/yellow}$. We further classify synthetic data errors into two distinct types based on the preservation of the logical operator f : **Operand Errors**, which perturb x and update y to maintain the validity of f ; and **Operator Errors**, which modify y given a fixed x to explicitly

073	violate f (e.g., red \rightarrow purple). By isolating errors	124
074	at the atomic level, ATOM enables precise control	125
075	over error type and ratio, allowing us to identify the	126
076	effects of different error mechanisms and pinpoint	127
077	the specific noise ratio at which performance degrada-	128
078	tion becomes significant. Specifically, guided by	129
079	this framework, we construct the diverse Atomic-	130
080	QA dataset to systematically re-evaluate the costs	
081	of data filtering and establish a quantification of	
082	error tolerance with generalizable Atomic-X.	
083	To validate the quantify the cost of <i>aggressive fil-</i>	
084	<i>tering</i> , we conduct a comprehensive evaluation by	
085	fine-tuning four distinct models and assessing their	
086	performance across 10 benchmark tasks. When	
087	benchmarking against standard baselines such as	
088	Alpaca (Taori et al., 2023), LIMA (Zhou et al.,	
089	2023), and WizardLM (Xu et al., 2025), Atomic-	
090	QA demonstrates superior performance in all set-	
091	tings. It implies that previous aggressive filtering	
092	strategies likely discard training samples critical	
093	for achieving optimal performance and further	
094	indicates the diversity of Atomic-QA.	
095	In addition, we systematically inject the two de-	
096	defined errors into Atomic-QA to construct corrupted	
097	datasets for model fine-tuning. By evaluating the	
098	models trained on these distinct variations, we ob-	
099	serve a distinct asymmetry: the model exhibits	
100	strong tolerance to operand errors. Even when	
101	operands are replaced with random strings, perfor-	
102	mance on downstream tasks remains stable with	
103	relative decrease of only 2.0% , and mechanistic	
104	analysis with Centered Kernel Alignment (Korn-	
105	blith et al., 2019) confirms that internal represen-	
106	tations remain consistent with the original Atomic-	
107	QA. In contrast, the model is highly sensitive to	
108	operator errors leading to a relative decrease of	
109	15.7% . <i>Leveraging the flexibility of ATOM, we extend</i>	
110	<i>this verification to general datasets Atomic-X.</i>	
111	By modifying Alpaca (Taori et al., 2023) and Meta-	
112	Math (Yu et al., 2023), we also find that increasing	
113	operator errors leads to severe model collapse. Fur-	
114	thermore, we find that existing <i>permissive filtering</i>	
115	<i>methods</i> have limitations in detecting these errors.	
116	Generally, our contributions are as follows:	
117	• A Quantifiable Atomic Framework: We intro-	
118	duce ATOM to decompose data into functional	
119	triplets ($\langle x, f, y \rangle$). This formalism provides the	
120	first precise boundary between tolerable <i>Operand</i>	
121	<i>Errors</i> and intolerable <i>Operator Errors</i> .	
122	• Controllable Training Dataset: Using the	
123	Atomic Tree constructed with ATOM, we pro-	
	pose two controllable training datasets for re-	124
	search purposes: diverse Atomic-QA and gener-	125
	alizable Atomic-X.	126
	• Asymmetric Error Tolerance: We discover a	127
	fundamental dichotomy: models maintain stable	128
	internal representations under operand perturba-	129
	tions but collapse under operator violations.	130
	2 Related Work	131
	Construction and Filtering of Synthetic Data	132
	Prior research establishes that SFT prioritizes data	133
	quality over sheer scale, with works like LIMA	134
	(Zhou et al., 2023) and Textbooks Are All You	135
	Need (Gunasekar et al., 2023) proving that effective	136
	alignment can emerge from small, high-quality cor-	137
	pora. Consequently, pipelines such as Self-Instruct	138
	(Wang et al., 2023b) and WizardLM (Xu et al.,	139
	2025) employ generate-then-filter strategies to syn-	140
	thesize data, while other approaches focus on se-	141
	lection via heuristic-based strategies (Zhao et al.,	142
	2024; Li et al., 2024, 2025) or contribution esti-	143
	mation (Dai et al., 2025; Chen et al., 2025b; Jiang	144
	et al., 2025). However, current methods lack a	145
	quantitative framework for defining data quality,	146
	leading to a trade-off between overly aggressive	147
	filtering and overly permissive retention.	148
	Learning Mechanisms in Noisy Data Prior	149
	work on noisy SFT investigates various settings.	150
	FTNI (Alajrami et al., 2025) found that perturbed	151
	instructions can surprisingly improve performance.	152
	TInt (Havrilla and Iyer, 2024) and FACO (Cho,	153
	2024b) explored model robustness to erroneous rea-	154
	soning chains. Recent evidence challenges the su-	155
	perficial alignment hypothesis (Zhou et al., 2023),	156
	showing that SFT goes beyond style adaptation	157
	to genuinely enhance reasoning and follow pre-	158
	training scaling laws (Chen et al., 2025a; Raghaven-	159
	dra et al., 2025). However, current research lacks	160
	fine-grained analysis, failing to determine noise	161
	thresholds. Details are listed in Appendix A.	162
	3 Methodology	163
	This section details the Atomic Tree Operation	164
	Modeling (ATOM) framework. As depicted in	165
	Figure 2, the proposed framework operates in two	166
	distinct phases: 1) Atomic Tree Construction (§3.1)	167
	and 2) Atomic Data Synthesis (§3.2).	168
	3.1 Atomic Tree Construction	169
	To formalize the construction process, we first in-	170
	troduce the concept of the Atomic Function .	171

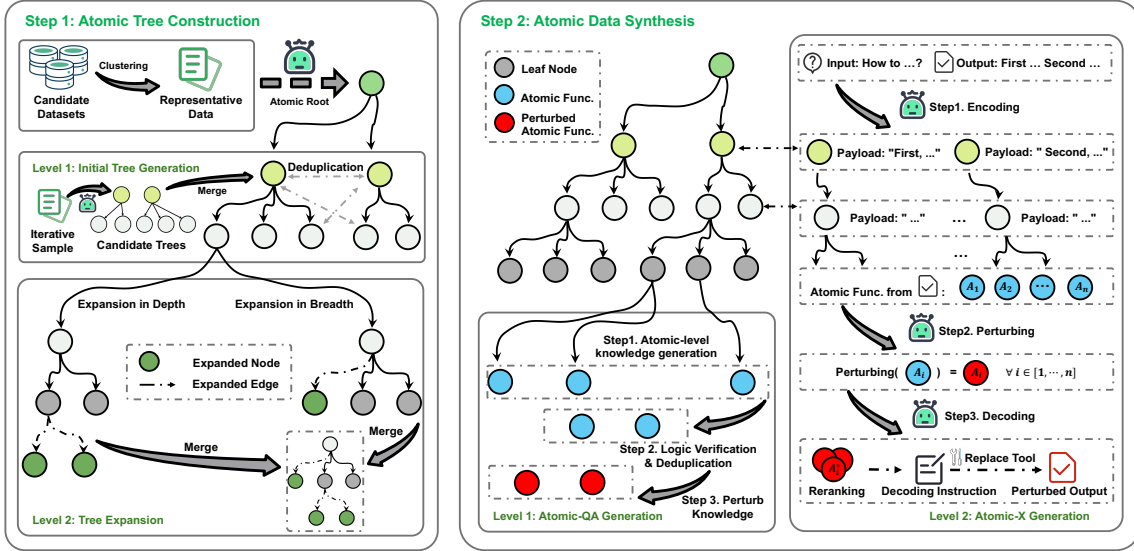


Figure 2: The pipeline of our Atomic Tree Operation Modeling (ATOM) framework. ATOM mainly consists of two steps: 1) Atomic Tree Construction and 2) Atomic Data Synthesis.

3.1.1 Formulation of Atomic Function

Formulation 1 (Atomic Function). An Atomic Function is defined as the minimal, logically indivisible unit of execution, denoted as a triplet $\mathcal{A} = \langle x, f, y \rangle$.

Formally, it represents a deterministic mapping $y = f(x)$, where:

- x (**The Operand**): Represents the irreducible logic or rule, such as arithmetic operations, relation schemas, or causal transitions.
- f (**The Operator**): Represents the irreducible structural operation logic or rule, such as an arithmetic operation, a relation extraction schema, or a causal transition.
- y (**The Output**): Represents the execution result derived from applying f to x .

A function is characterized as Atomic if any further decomposition of the operator f would result in a loss of semantic intent or invalidate the executable logic within the task context. This definition holds across various domains, though its specific manifestation may vary. For example, in the mathematical domain, an Atomic Function acts as a minimal indivisible calculation, such as simplifying $\sqrt{(-2)^2 + (-1)^2}$ to $\sqrt{5}$. In the knowledge domain, it represents a fundamental relational mapping, exemplified by $f_{\text{CapitalOf}}(\text{Paris}) \mapsto \text{France}$, which retrieves the target entity based on a given subject and relation. **To further clarify the atomic function across different scenarios, we provide additional examples in Appendix G.**

3.1.2 Pipeline of Atomic Tree Construction

To precisely extract or synthesize atomic functions, an appropriate data structure is essential. We adopt the tree structure as the canonical representation, given its inherent alignment with recursive functional composition. Since complex tasks can be modeled as nested compositions of functions, the tree structure explicitly maps this hierarchical logic. Through this framework, a complex objective function is recursively decomposed until the nodes become semantically irreducible terminal nodes—identified as atomic functions. Consequently, our goal is to construct a *functional hierarchy* where leaf nodes represent these operators.

As illustrated in Step 1 of Figure 2, the phase commences with **Representative Data Selection**. To mitigate computational overhead while ensuring coverage, we encode candidate datasets using bge-m3 (Chen et al., 2024). We then apply KCenterGreedy clustering to select samples that maximize the coverage of the semantic space, thereby ensuring the extracted subset embodies diversity.

The process then progresses to **Level 1: Initial Tree Generation**, which establishes the skeletal structure of the functional hierarchy. We employ an iterative strategy where the o3-mini model iteratively processes the selected samples to explicitly generate candidate sub-trees. These sub-trees serve as the primary functional units extracted directly from the raw data. As illustrated in the left-middle panel of Figure 2, the construction process is dynamic. Generated candidate trees are systematic-

cally integrated via a merge operation to construct the global functional tree. To guarantee the quality of this integration, we employ an LLM-based semantic assessment module. This module acts as a strict filter with to identify and resolve semantic redundancies among nodes, ensuring that the tree remains structurally compact and functionally distinct. **The initial tree are ultimately checked by four human experts' voting results.**

Following initialization, the pipeline advances to **Level 2: Tree Expansion**, aimed at satisfying the strict "atomic" constraint defined in our formulation. This phase refines the hierarchy through a dual-path expansion strategy inspired by EpiCoder (Wang et al., 2025). We employ Expansion in Depth to recursively refine nodes until irreducible judged by LLM or capped, alongside Expansion in Breadth to supplement lateral operators for full coverage. As shown in the left-bottom panel of Figure 2, the resulting expanded nodes are consolidated back into the main hierarchy through a final Merge operation, ensuring the evolving tree remains rigorous, logically consistent, and complete. **The detailed examples and prompts are listed in Appendix B.1.**

3.2 Atomic Data Synthesis

After obtaining the atomic tree, to accommodate different experimental requirements, we support two Atomic Data Synthesis strategies: **Atomic-QA Generation** and **Atomic-X Generation**, which respectively support simple data synthesis and complex data analysis.

3.2.1 Atomic-QA Generation

To enable large-scale synthetic data generation for robustness analysis, we propose the **Atomic-QA Generation** module. This process transforms leaf nodes into synthetic QA pairs through a three-step pipeline, where the second and third steps respectively generate Atomic-QA and its variants:

Step 1 (Generation): We instantiate operands (x) from leaf nodes and prompt the model to synthesize corresponding atomic triplets $\mathcal{A} = \langle x, f, y \rangle$.

Step 2 (Verification): To ensure data quality and diversity, we deduplicate overlapping operands x across neighboring nodes and rigorously verify the logical validity of the generated triplets.

Step 3 (Perturbation): Finally, we inject controlled noise to create *Perturbed Atomic Functions*. This involves two strategies: i) replacing the input x with random strings and adjusting y accordingly

to preserve the operator nature of f as much as possible, or ii) altering only the output y to break the deterministic mapping $f(x)$, thereby creating logically invalid samples.

3.2.2 Atomic-X Generation

To extend our framework to existing complex datasets such as Alpaca (Taori et al., 2023) and MetaMath (Yu et al., 2023), we propose the **Atomic-X Generation** pipeline. The right panel of Figure 2 shows this process operates through an Encoder-Perturber-Decoder architecture:

Step 1 (Encoding): The process initiates by taking an instruction and its corresponding output from existing datasets as inputs. The encoder performs a top-down decomposition of the output text into a hierarchical tree structure. As illustrated in Figure 2, the model recursively parses the content into distinct logical branches based on the identified child nodes, where each node preserves a specific text segment as its semantic payload. This decomposition process repeats until the logic is no longer divisible, corresponding to the final leaf nodes in atomic tree, yielding an ordered sequence denoted as $\mathbb{A} = \{\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_n\}$.

Step 2 (Perturbing): In this phase, a perturbation function is applied to the extracted atomic functions. As shown in the formula $\text{Perturbing}(\mathcal{A}_i) = \mathcal{A}'_i$, we inject noise or logical alterations (e.g., modifying operations from $1 + 1 = 2$ to $1 + 1 = 3$) to generate a set of corrupted atomic functions.

Step 3 (Decoding): The decoder reconstructs the final output through a multi-stage process. First, the perturbed atomic functions are reranked based on the magnitude of their logical impact as $\mathbb{A}^* = \{\mathcal{A}_1^*, \mathcal{A}_2^*, \dots, \mathcal{A}_n^*\}$. Next, the model generates specific decoding instructions for modification. Finally, a deterministic replace Tool utilizes these instructions to edit the original text, producing the final *Perturbed Output*. **The detailed pipeline and prompts is listed in Appendix B.2.**

4 Experimental Setup

In the following sections, we want to answer the following research questions (RQs):

- **RQ1:** Does aggressive filtering come at the cost of structural diversity, thereby hindering model performance? (§5)
- **RQ2:** Which plays the dominant role in training: the **Operand** (x) or the **Operator** (f)? (§6)
- **RQ3:** What is the quantitative threshold for noise

Table 1: Overall performance across all training datasets on LLaMA3.1-8B and Qwen3-8B. Color-coded cells (green/red) are used to indicate performance above or below the mean across methods (All Average), with darker colors representing better or worse performance. **The results for Qwen2.5-7B and Qwen2.5-14B are in Table 3.**

Model Family		ARC-Challenge	ARC-Easy	BoolQ	Hotpot-QA	Open-BookQA	PIQA	RACE	SIQA	SQuADv2	Wino-Grande	Avg
LLaMA3.1-8B	Base	0.549	0.847	0.820	0.106	0.366	0.807	0.433	0.529	0.375	0.765	0.560
	All Average	0.564	0.850	0.855	0.128	0.387	0.812	0.464	0.548	0.381	0.768	0.576
	Instruct	0.558	0.854	0.870	0.209	0.362	0.811	0.464	0.548	0.420	0.764	0.586
	LIMA	0.552	0.848	0.828	0.134	0.368	0.808	0.432	0.529	0.377	0.773	0.565
	Condor	0.563	0.854	0.855	0.117	0.382	0.811	0.470	0.546	0.378	0.770	0.575
	Magpie	0.575	0.845	0.859	0.122	0.380	0.808	0.464	0.563	0.374	0.761	0.575
	WizardLM	0.580	0.862	0.858	0.107	0.396	0.823	0.478	0.573	0.385	0.766	0.583
	Alpaca	0.562	0.846	0.867	0.115	0.406	0.819	0.469	0.576	0.393	0.776	0.583
	Self-Instruct	0.469	0.811	0.839	0.097	0.352	0.803	0.430	0.511	0.321	0.739	0.537
	Wiki	0.562	0.856	0.845	0.132	0.370	0.810	0.453	0.530	0.386	0.778	0.572
	Wiki-Rewrite	0.560	0.855	0.847	0.100	0.364	0.811	0.456	0.528	0.388	0.777	0.569
	Atomic-LIFD	0.607	0.858	0.875	0.133	0.440	0.812	0.494	0.568	0.364	0.770	0.592
Atomic-QA	0.613	0.864	0.865	0.145	0.442	0.816	0.500	0.562	0.399	0.774	0.598	
Qwen3-8B	Base	0.631	0.871	0.871	0.233	0.370	0.799	0.473	0.566	0.444	0.743	0.600
	All Average	0.631	0.870	0.875	0.183	0.375	0.798	0.489	0.562	0.420	0.746	0.595
	Instruct	0.629	0.873	0.871	0.114	0.378	0.782	0.491	0.567	0.430	0.707	0.584
	LIMA	0.632	0.870	0.873	0.237	0.374	0.803	0.478	0.558	0.445	0.755	0.602
	Condor	0.637	0.875	0.875	0.230	0.374	0.800	0.486	0.567	0.435	0.751	0.603
	Magpie	0.629	0.869	0.875	0.125	0.352	0.798	0.482	0.553	0.436	0.742	0.586
	WizardLM	0.639	0.870	0.877	0.053	0.360	0.799	0.491	0.572	0.428	0.756	0.585
	Alpaca	0.634	0.871	0.880	0.042	0.362	0.803	0.495	0.580	0.445	0.753	0.586
	Self-Instruct	0.515	0.820	0.866	0.103	0.338	0.787	0.439	0.513	0.341	0.721	0.544
	Wiki	0.638	0.876	0.869	0.191	0.366	0.801	0.489	0.570	0.424	0.759	0.598
	Wiki-Rewrite	0.641	0.877	0.871	0.212	0.370	0.800	0.497	0.568	0.417	0.758	0.601
	Atomic-LIFD	0.671	0.886	0.884	0.364	0.420	0.800	0.523	0.568	0.398	0.751	0.626
Atomic-QA	0.672	0.886	0.883	0.338	0.426	0.807	0.510	0.562	0.421	0.750	0.625	

tolerance in real synthesis data setting? (§6)

Dataset: To address RQ1 and RQ2, we construct an Atomic Tree following Section 3.2.1 using FB15k (Schlichtkrull et al., 2018), ZsRE (Tran et al., 2022), and wiki_recent (Wang et al., 2023a) as seeds. This yields the Atomic-QA dataset with 1.37 million instances. To answer question RQ3, we modify real-world complex data. Specifically, we apply the methods from Section 3.2.2 to Alpaca and MetaMath. See Appendix C.1 for details.

Model: To ensure generalizability, we evaluate LLaMA3.1 (8B) (Grattafiori et al., 2024), Qwen2.5 (7B & 14B) (Qwen-Team, 2024), and Qwen3 (8B) (Yang et al., 2025a). Training is conducted via Llama-Factory (Zheng et al., 2024b) and evaluation via the LM Evaluation Harness (Gao et al., 2024). Configuration details are in Appendix C.2.

5 Quantifying the Cost of Strict Filtering

In this section, we quantify the cost of strict filtering in existing synthetic data pipelines from two perspectives to answer the RQ1: the degradation in model performance and exclusion of potentially valuable samples. To assess the performance impact, we conduct a comparative evaluation across 10 benchmarks using multiple backbones. Subsequently, to analyze the exclusion of valuable data, we employ multidimensional diversity metrics to visualize the semantic contraction in traditional methods.

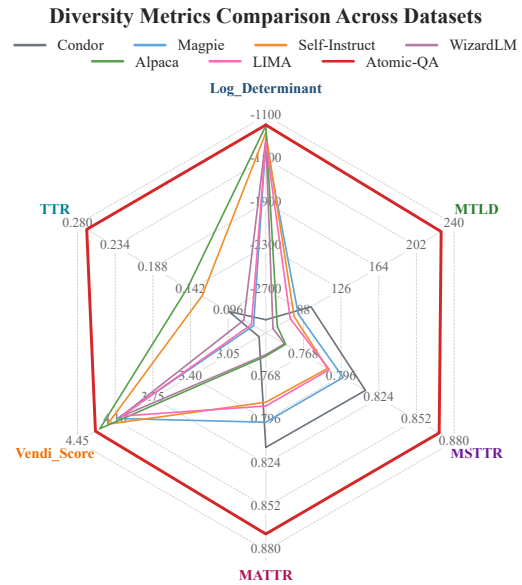


Figure 3: Radar plot of corpus diversity.

5.1 Degradation in Model Performance

Table 1 and Table 3 presents the evaluation results across 10 benchmarks using four backbones. The empirical data explicitly quantifies the performance penalty incurred by traditional pipelines.

The Cost of Aggressive Filtering: A direct comparison with LIMA reveals a significant performance gap. On LLaMA3.1-8B, LIMA achieves an average score of 56.5%, whereas our **Atomic-QA** reaches 59.8%—a substantial absolute gain of **+3.3%**. This trend persists across model scales, with Atomic-QA outperforming LIMA by 1.8%

LLaMA3.1-8B Performance: Atomic-QA vs Replaced Atomic-QA vs Rewrite-to-Hallucination

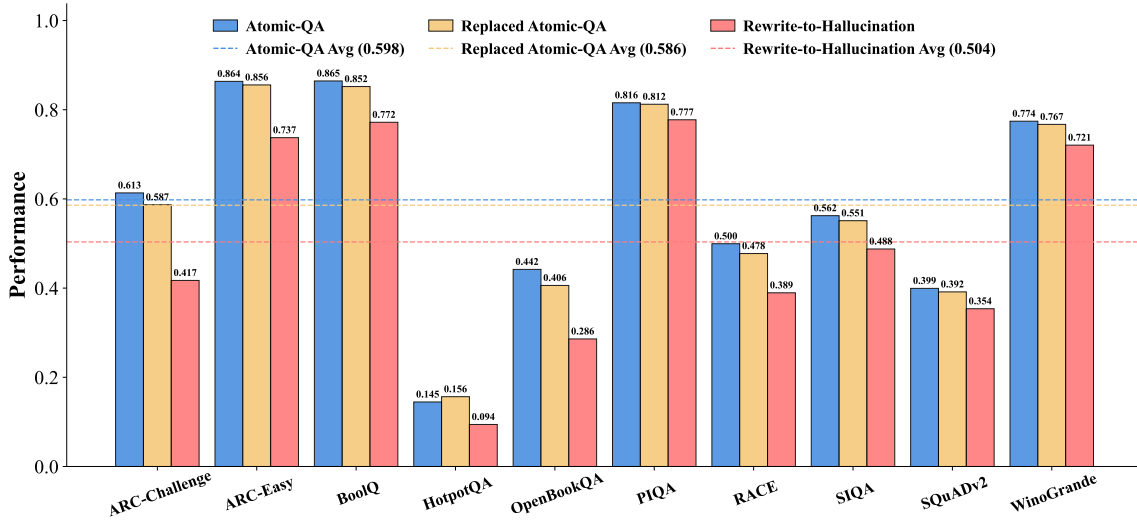


Figure 4: Performance comparison of LLaMA3.1 across the Atomic-QA, Replaced Atomic-QA, and Rewrite-to-Hallucination settings. Each bar shows the model’s accuracy on individual tasks, while the dashed lines denote the average performance for each setting.

on Qwen2.5-7B and 3.8% on Qwen3-8B. This gap underscores that aggressively filtering methods discard valuable structural signals, capping the model’s potential. **Despite using a minimal number of training tokens (Appendix D.1), Atomic-LIFD achieves competitive results. This strongly demonstrates that diversity, rather than data scale, is the primary driver of performance.**

Superiority over Standard Synthetic Baselines: Compared to widely used synthetic datasets such as Alpaca and WizardLM, Atomic-QA demonstrates consistent superiority. For instance, on the Qwen2.5-7B backbone, Atomic-QA (61.2%) surpasses Alpaca (56.5%) and WizardLM (56.8%) by margins exceeding 4%. Notably, the Atomic-QA excels in reasoning-intensive tasks like OpenBookQA and ARC-Challenge, suggesting that preserving structural diversity is more critical.

5.2 Exclusion of Potentially Valuable Samples

Figure 3 visualizes the multidimensional diversity of training corpora using metrics such as Vendi Score (Friedman and Dieng, 2023), Log Determinant (Yang et al., 2025b), TTR-variants (Covington and McFall, 2010) and MTLD (McCarthy and Jarvis, 2010). The radar plot reveals a stark trade-off: datasets characterized by aggressive filtering exhibit a significantly contracted diversity.

In contrast, Atomic-QA (represented by the outermost red boundary) consistently achieves the highest scores across all six metrics, completely encompassing the baselines. This empirical evi-

dence confirms that previous aggressive methods inadvertently discard potentially valuable samples, thereby stripping away the rich structural and semantic variations valuable for training. By relaxing atomic operand verification, Atomic-QA successfully preserves this critical structural breadth. By shifting the focus to operator consistency, Atomic-QA effectively retains these potentially valuable samples, thereby recovering the performance capabilities that are typically lost during aggressive data pruning.

6 Measuring Operator–Operand Dominance

To evaluate operator–operand dominance during training, we modify Atomic-QA in the following ways. First, we replace the operand string in each Atomic-QA question–answer pair with a random string of the same character length, while keeping the operator unchanged, resulting in *Replaced Atomic-QA*. This ensures that the operand in Replaced Atomic-QA is incorrect while maximizing the correctness of the operator. Following Xie et al. (2024), we then use an LLM to rewrite the answers and verify hallucinations to modify the operator accordingly to construct *Rewrite-to-Hallucination*.

6.1 Performance Comparison

Figure 4 presents the comparative results on LLaMA3.1-8B, revealing a striking double dissociation between structural alignment and factual

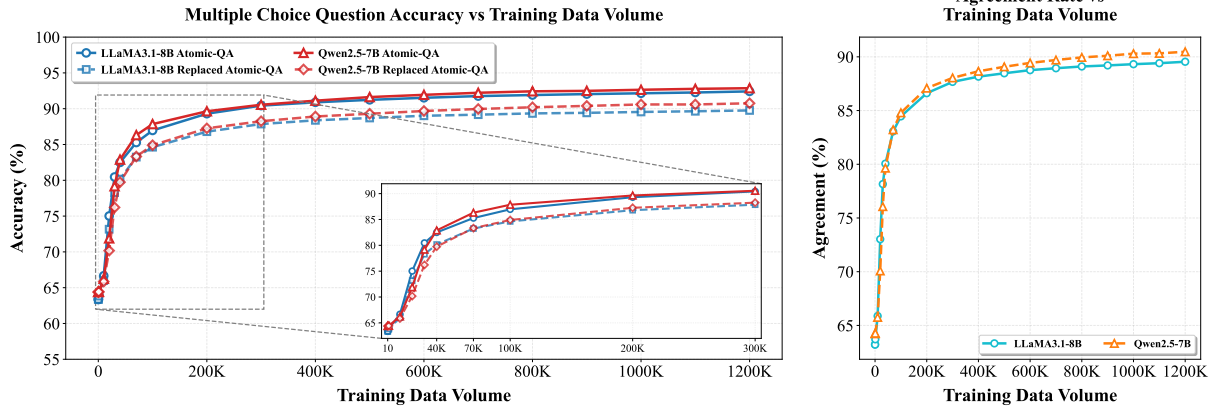


Figure 5: MCQ accuracy as a function of training data scale (main plot), together with model agreement under the replaced vs. original entity settings (inset).

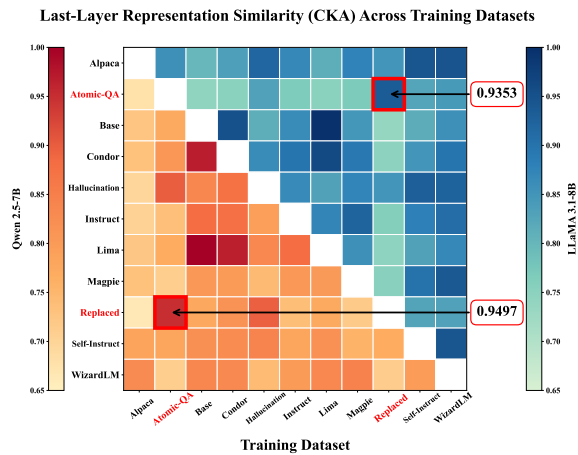


Figure 6: CKA similarity of the output hidden states between different trained models. The lower part represents the CKA similarity of Qwen 2.5-7B, while the upper part represents LLaMA 3.1-8B.

correctness to answer the RQ2. Results of other models are listed in Figure 14, 15 and 16.

Robustness to Operand Noise: The resilience to operand noise is not limited to simple tasks but generalizes consistently across downstream evaluations of varying complexity. As shown in Figure 4, the performance of Replaced Atomic-QA closely tracks the baseline Atomic-QA across all datasets. On reasoning-heavy tasks like ARC-Challenge, the performance gap is minimal (0.589 vs 0.613), and on the complex multi-hop reasoning task HotpotQA, the model trained on random strings even slightly outperforms the baseline (0.158 vs 0.145). This demonstrates that even in complex scenarios, the model’s performance relies primarily on acquiring the operator logic f rather than memorizing the specific atomic operands x . **The results shown in Figures 14, 15, and 16 indicate that our findings generalize across different model structures.**

Sensitivity to Operator Noise: In stark contrast, when trained on Rewrite-to-Hallucination, where the structural logic is disrupted (corrupting the operator f), performance plummets to **50.4%**. This represents a substantial absolute decline of **9.4%** (a relative drop of $\sim 15.7\%$), highlighting the model’s acute sensitivity to structural integrity.

6.2 Mechanism Validation

A potential critique of the strong performance of Replaced Atomic-QA is that the model might be "hacking" the evaluation metric or learning superficial heuristics rather than genuine task structures. To refute this and demonstrate that the model is indeed learning the underlying operator logic (f), we analyze the training dynamics and latent representations.

Behavioral Alignment: We constructed a held-out diagnostic test set by sampling subsets from Atomic-QA and rewriting them into multiple-choice questions (MCQs). We then monitored the validation metrics throughout the training process. As shown in Figure 5 (Left), the accuracy curves for models trained on both Atomic-QA and Replaced Atomic-QA rise in unison, converging to nearly identical performance levels. Complementing this, Figure 5 (Right) illustrates the Agreement Rate—the proportion of test instances where both models predict correctly. The steady increase in agreement demonstrates that as training progresses, the models do not just achieve similar scores; they converge on the same behavioral patterns. This synchrony indicates that the optimization landscape is dominated by the structural operator f , rendering atomic operands (x) irrelevant to learning.

Representational Alignment: Do these models merely behave similarly, or do they actually learn

the same internal representations? We utilized Centered Kernel Alignment (CKA) (Kornblith et al., 2019) to measure the similarity of the last-layer hidden states between different models. Figure 6 presents the CKA heatmap. Strikingly, the model trained on Replaced Atomic-QA exhibits an exceptionally high similarity score (e.g., **0.9497** on Qwen2.5-7B) with the model trained on the standard *Atomic-QA*. This alignment suggests that the models have learned structurally equivalent internal representations, despite their differing training data. We also present the training dynamics of PPL and CKA in Figure 12 and 13.



Figure 7: Accuracy degradation under increasing simulated error rates for LLaMA 3.1-8B and Qwen 2.5-7B on Alpaca. The graph shows model performance as a function of training error rate. The results for MetaMath is listed in Figure 17.

6.3 Operator Noise Tolerance Limits

To simulate the error patterns of teacher models with different capability levels under real-world conditions (*RQ3*), we adopt Beta-distribution-based sampling as described in Appendix E and inject operator errors with varying levels of atomic-operation noise into Alpaca and MetaMath. We analyze the impact of increasing atomic error rates (from 0% to 100%) on model performance, visualized in Figure 7. The results indicates that model exhibits a distinct two-stage non-linear degradation pattern in response to operation noise.

Phase I: Gradual Logic Degradation: As shown in Figure 7, injecting noise into atomic operations results in a slow but observable performance decay during the initial stage (0% \rightarrow 40% noise). Unlike operand noise, operator errors disrupt structural learning, causing proportional capability loss.

Phase II: Structural Collapse: As the error rate exceeds a critical threshold (roughly $> 60\%$ in Figure 7), the degradation pattern shifts distinctly, characterized by a significantly steepened negative slope. Visually, this is represented by the rapid plunge in accuracy towards the tail end of the curve (0.6 \rightarrow 1.0). This accelerated decline indicates a structural collapse, suggesting that the accumulation of logical inconsistencies has overwhelmed the model’s reasoning schema (\mathcal{T}).

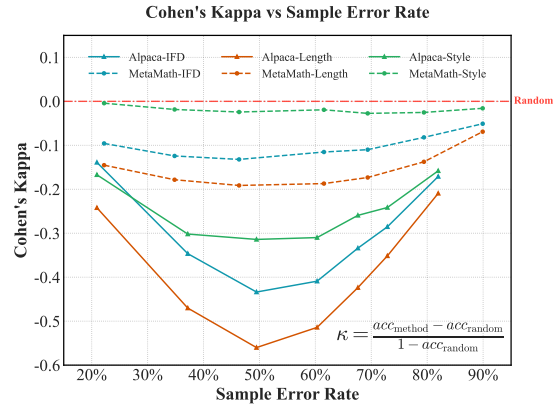


Figure 8: Cohen’s Kappa versus Sample Error Rate across filtering methods.

6.4 Upper Bound vs. Filtering Reality

To assess filtering potential, we established an upper bound by retraining on error-filtered subsets with sample error rates less than 50% (dashed lines in Figure 7). These models matched the noise-free baseline, confirming that degradation drives from operator errors. However, as shown in Figure 8, current filtering strategies based on metrics such as IFD (Li et al., 2024), sequence length (Zhao et al., 2024), and style-consistency (Li et al., 2025) fail to identify these atomic faults. The resulting negative Cohen’s Kappa coefficients indicate that these algorithms are unable to differentiate between atomic logic faults and correct samples

7 Conclusion

This study establishes error tolerance in synthesis data, moving beyond the binary high-low quality definitions. We demonstrate that SFT is fundamentally a logic-learning process in atomic view where models learn functional atomic operators while specific operand-values are changeable variables. We show that models learn atomic logic even from randomized operands. Thus, synthetic data should favor operators over operands for generalizability.

8 Limitation

Due to computational resource constraints, we do not conduct experiments on models with 32B parameters or larger. Furthermore, while our ATOM framework demonstrates robustness across general instruction following and mathematical reasoning tasks, we have not yet extended the atomic decomposition methodology to code generation or other domains.

References

- Ahmed Alajrami, Xingwei Tan, and Nikolaos Aletras. 2025. [Fine-tuning on noisy instructions: Effects on generalization and performance](#). *Preprint*, arXiv:2510.03528.
- Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. [Mathqa: Towards interpretable math word problem solving with operation-based formalisms](#). *Preprint*, arXiv:1905.13319.
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024. [Longbench: A bilingual, multitask benchmark for long context understanding](#). *Preprint*, arXiv:2308.14508.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2019. [Piqa: Reasoning about physical commonsense in natural language](#). *Preprint*, arXiv:1911.11641.
- Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. [Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation](#). *Preprint*, arXiv:2402.03216.
- Runjin Chen, Gabriel Jacob Perin, Xuxi Chen, Xilun Chen, Yan Han, Nina S. T. Hirata, Junyuan Hong, and Bhavya Kailkhura. 2025a. [Extracting and understanding the superficial knowledge in alignment](#). *Preprint*, arXiv:2502.04602.
- Yicheng Chen, Yining Li, Kai Hu, Zerun Ma, Haochen Ye, and Kai Chen. 2025b. [Mig: Automatic data selection for instruction tuning by maximizing information gain in semantic space](#). *Preprint*, arXiv:2504.13835.
- Hyunsoo Cho. 2024a. [Unveiling imitation learning: Exploring the impact of data falsity to large language model](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 62–73, Bangkok, Thailand. Association for Computational Linguistics.
- Hyunsoo Cho. 2024b. [Unveiling imitation learning: Exploring the impact of data falsity to large language model](#). *Preprint*, arXiv:2404.09717.

- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. [Boolq: Exploring the surprising difficulty of natural yes/no questions](#). *Preprint*, arXiv:1905.10044.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. [Think you have solved question answering? try arc, the ai2 reasoning challenge](#). *Preprint*, arXiv:1803.05457.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *Preprint*, arXiv:2110.14168.
- Michael A. Covington and Joe D. McFall. 2010. [Cutting the gordian knot: The moving-average type–token ratio \(mattr\)](#). *Journal of Quantitative Linguistics*, 17:100 – 94.
- Qirun Dai, Dylan Zhang, Jiaqi W. Ma, and Hao Peng. 2025. [Improving influence-based instruction tuning data selection for balanced learning of diverse capabilities](#). *Preprint*, arXiv:2501.12147.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *Preprint*, arXiv:2501.12948.
- Dan Friedman and Adji Bousso Dieng. 2023. [The vendi score: A diversity evaluation metric for machine learning](#). *Transactions on Machine Learning Research*.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, and 5 others. 2024. [The language model evaluation harness](#).
- Gemini-Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, and 1332 others. 2025. [Gemini: A family of highly capable multimodal models](#). *Preprint*, arXiv:2312.11805.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh

883 proposed methods that estimate sample contribu-
884 tions to downstream performance, such as BIDS
885 (Dai et al., 2025), MIG (Chen et al., 2025b), and
886 MIWV (Jiang et al., 2025).

887 However, most data synthesis, refinement, and
888 selection methods lack a guiding framework for
889 quantitatively defining data quality. Consequently,
890 they oscillate between two extremes: overly ag-
891 gressive filtering that removes valuable supervisory
892 signals, and overly permissive filtering that retains
893 training instances with potential errors.

894 **Training Data Noise Effects and Learning Mech-**
895 **anisms in SFT** Prior work has investigated noisy
896 supervision in SFT under a variety of controlled
897 settings. FTNI (Alajrami et al., 2025) observes that
898 training on perturbed instructions (such as remov-
899 ing stop words or shuffling words) can improve
900 downstream performance in some cases. For rea-
901 soning supervision, TInt (Havrilla and Iyer, 2024)
902 studies robustness to noise in algorithmic chains of
903 thought by injecting errors during reasoning, and
904 FACO (Cho, 2024b) studies how different propor-
905 tions of erroneous reasoning-chain samples in the
906 training data affect model performance.

907 From a learning-mechanism perspective, the su-
908 perflicial alignment hypothesis characterizes SFT as
909 primarily focusing on adopting the language style
910 of responsible AI assistants and relying largely on
911 the knowledge already acquired by base LLMs
912 (Zhou et al., 2023; Lin et al., 2023). Beyond
913 this view, subsequent evidence shows that SFT
914 enhances models’ reasoning and contextual under-
915 standing compared to their base counterparts (Chen
916 et al., 2025a). Raghavendra et al. (2025) further ob-
917 serve that, similar to the pre-training scaling laws,
918 post-training task performance scales as a power
919 law against the number of finetuning examples, re-
920 flecting that language models are not necessarily
921 confined to using only the knowledge learned dur-
922 ing pretraining.

923 Current analyses of noise in training data and
924 SFT learning mechanisms remain confined to spe-
925 cific types of datasets or lack fine-grained, sample-
926 level structural analysis. As a result, they are un-
927 able to determine which errors are tolerable to the
928 model and which genuinely disrupt learning, leav-
929 ing the critical boundary between benign and detri-
930 mental noise unclear and the quantification of error
931 tolerance in synthetic data unexplored.

B Pipeline for ATOM 932

This section presents the detailed process of ATOM.
We elaborate from the perspectives of the Pipeline
for Atomic Tree Construction and Atomic Data
Synthesis, respectively. 933
934
935
936

B.1 Pipeline for Atomic Tree Construction 937

To facilitate clarity, we first detail the prompts em-
ployed in our pipeline, followed by a formal de-
scription using pseudo-code. 938
939
940

ATOM initiates the process by extracting atomic
functions from representative data, guided by pre-
defined definitions via Prompt 1. Given that
Atomic-QA is grounded in factual datasets, we
leverage entities as cues to ensure the precision of
constructed functions. Subsequently, we organize
these functions into an Initial Atomic Tree using
Prompt 2. To guarantee structural completeness,
we generate trees for various subsets of atomic
functions under diverse sampling schemes, while
simultaneously pruning duplicate nodes at the same
hierarchy level via Prompt 3. After the above steps,
we obtain the tree shown in Figure 9, 10 and 11. Fi-
nally, building upon this initial structure and its enu-
merated first-layer nodes, we expand the tree across
both depth and breadth dimensions using Prompt
4. The comprehensive procedure for Atomic Tree
Construction is outlined in Algorithm 1. 941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958

B.2 Pipeline for Atomic Data Synthesis 959

This section introduces Atomic-QA Generation and
Atomic-X Generation, respectively. 960
961

B.2.1 Atomic-QA Generation 962

In the Atomic-QA generation phase, we first em-
ploy an LLM to generate operands (i.e., entities)
belonging to a specific category within the given
atomic tree, guided by Prompt 5. Subsequently, the
model generates potential operations and their cor-
responding outputs for each operand. Finally, we
utilize Prompt 6 to validate the logical correctness
of the generated atomic operations. 963
964
965
966
967
968
969
970

B.2.2 Atomic-X Generation 971

The Atomic-X generation framework comprises
three core modules: the Encoder, the Perturber, and
the Decoder. Guided by Prompt 7, the Encoder
initiates the process at the root node of the atomic
tree. It recursively identifies operation-bearing pay-
loads and decomposes the structure until reaching
the leaf nodes, which are subsequently resolved
into atomic operations. Next, the Perturber applies 972
973
974
975
976
977
978
979

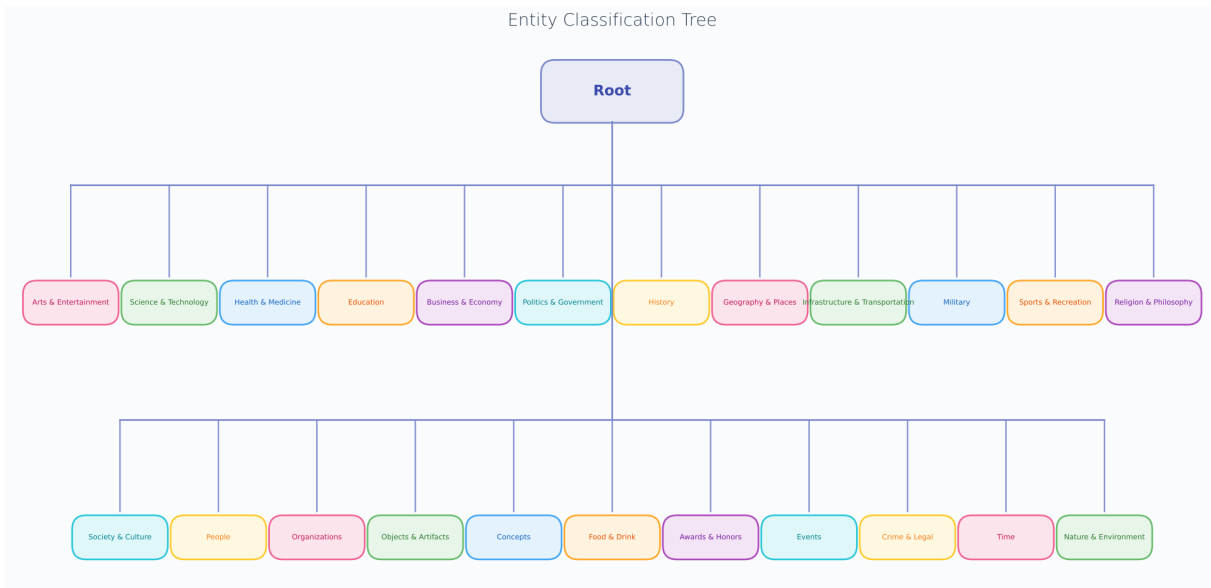


Figure 9: The initial atomic tree of Atomic-QA.

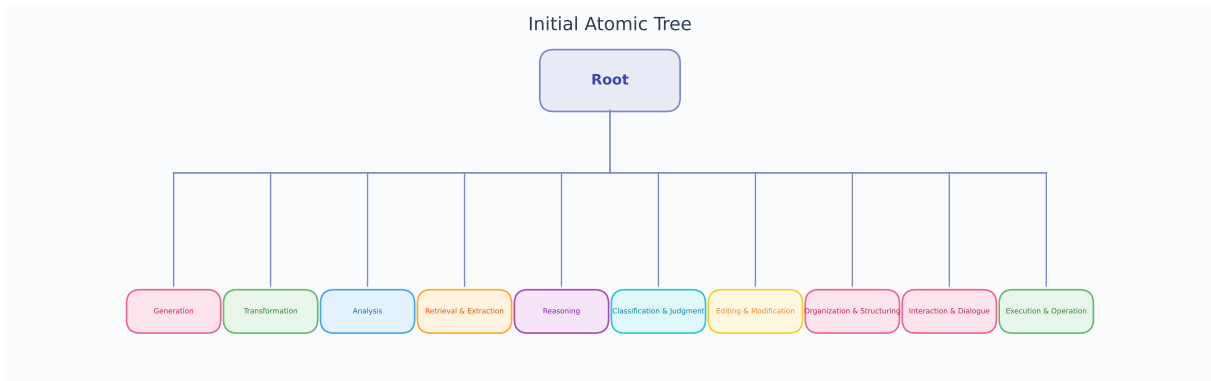


Figure 10: The initial atomic tree of Alpaca.

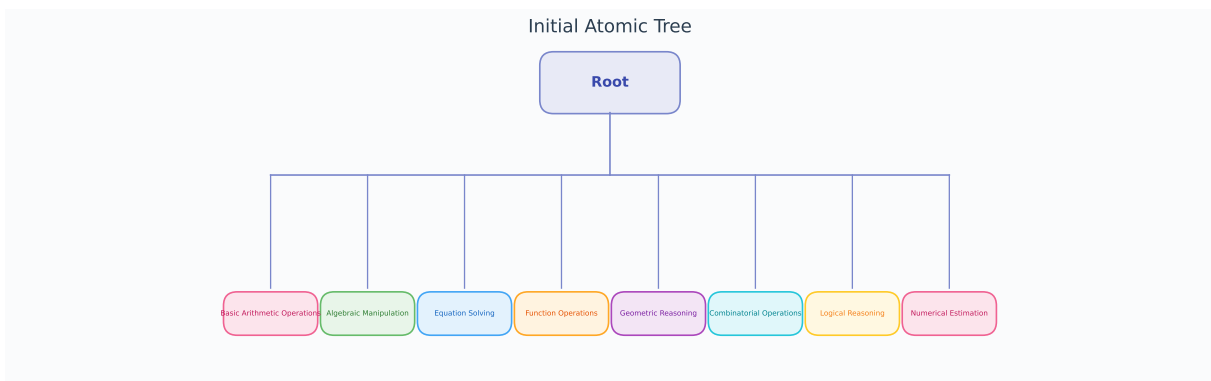


Figure 11: The initial atomic tree of MetaMath.

980 perturbations to each atomic operation utilizing the
 981 Prompt 8. Finally, the Decoder modifies the original
 982 data by employing the SEARCH-REPLACE
 983 and TERMINATE tools defined in Prompt 9.

984 C Details about Experimental Setup

985 In this section we will introduce our experimental
 986 details. First, we introduce the details of data
 987 generation used in the experiments, followed by
 988 a description of the model training and evaluation
 989 procedures.

Algorithm 1: Atomic Tree Construction Pipeline

Input : Candidate Datasets \mathcal{D} , Embedding Model \mathcal{M}_{emb} , Generation Model \mathcal{M}_{gen} , Max Interaction \mathcal{I}

Output : Atomic Tree \mathcal{T}

// Phase 1: Representative Data Selection

```
1  $V \leftarrow \text{Encode}(\mathcal{D}, \mathcal{M}_{bge})$ ,  
    $S_{rep} \leftarrow \text{KCenterGreedy}(V)$ ,  
    $\mathcal{T} \leftarrow \{root\}$   
2  $\mathbb{A}_{rep} \leftarrow \text{Extract}(S_{rep})$   
   // Phase 2: Level 1 - Initial Tree Generation  
3 for  $batch \mathbb{A}_{sub} \in \text{sample}(\mathbb{A}_{rep})$  do  
4    $\mathcal{T}_{sub} \leftarrow \text{GenSubTree}(\mathbb{A}_{sub}, \mathcal{M}_{gen})$   
    $\mathcal{T} \leftarrow \text{Merge}(\mathcal{T}, \mathcal{T}_{sub})$   
   // Phase 3: Level 2 - Tree Expansion  
5 for  $iter i \in \mathcal{I}$  do  
6    $\mathcal{T}_i \leftarrow \emptyset$   
7   for  $node n \in \text{LeafNodes}(\mathcal{T})$  do  
8      $\mathcal{N}_d \leftarrow \text{DepthExpansion}(n, \mathcal{M}_{gen})$   
      $\mathcal{N}_b \leftarrow \text{BreadthExpansion}(n, \mathcal{M}_{gen})$   
9     if  $\mathcal{N}_d \neq \emptyset \vee \mathcal{N}_b \neq \emptyset$  then  
10       $\mathcal{T}_i \leftarrow \text{Merge}(\mathcal{T}_i, \mathcal{N}_d \cup \mathcal{N}_b)$   
11   if  $\mathcal{T}_i = \emptyset$  then  
12     break  
13    $\mathcal{T} \leftarrow \text{Merge}(\mathcal{T}, \mathcal{T}_i)$   
14 return  $\mathcal{T}$ 
```

C.1 Details about Dataset Generation

For the Atomic-QA generation process, we commence by employing the KCenterGreedy algorithm to identify 5,000 cluster centers. The initial tree is constructed through an iterative procedure spanning 500 rounds, with 50 subjects randomly sampled in each round. During the subsequent evolution stage, the tree evolves for 1,715 iterations. Notably, in each iteration, a single leaf node undergoes five depth expansions and five breadth expansions. This results in a tree with 6,931 leaf nodes (max depth 4), yielding the Atomic-QA dataset with approximate 1.37 million instances.

Regarding the Atomic-X generation process, we initiate construction by randomly sampling 10,000 instances, generating 100 cluster centers at each step. The initial tree construction involves 200 evo-

lution rounds, wherein 50 data points are selected per round. In the final stage, the evolution iterations differ by model: Alpaca evolves for 1,203 iterations, while MetaMath evolves for 1,347 iterations. This results in a tree with 276 and 398 leaves respectively for Alpaca (max depth 5) and MetaMath (max depth 6). From the remaining samples, we randomly select 10,000 instances and construct the corresponding Atomic-Alpaca and Atomic-MetaMath datasets using the method described in Section 3.2.2.

C.2 Details about Model Training and Evaluation

We conduct experiments on four representative open-source pre-trained LLMs that span different architectures and scales, including LLaMA3.1-8B, Qwen3-8B, Qwen2.5-7B, and Qwen2.5-14B. All models are fine-tuned with LLaMA-Factory (Zheng et al., 2024a) and evaluated with the LM Evaluation Harness (Gao et al., 2024).

Datasets. We compare models fine-tuned separately under different training dataset configurations, as summarized in Tables 1 and 3. The comparison set includes Atomic-QA and its variants (Wiki, Wiki-Rewrite, and Atomic-LIFD), the official *Instruct* versions, and six representative open-source datasets: LIMA (Zhou et al., 2023), Condor (Maosong et al., 2025), Magpie (Xu et al., 2024), WizardLM (Xu et al., 2025), Alpaca (Taori et al., 2023), and Self-Instruct (Wang et al., 2023b). Specifically, *Instruct* refers to the official tuned versions released for each pretrained model. *Wiki* consists of synthetic instructions paired with Wikipedia entity summaries as responses, where the entities are drawn from the Atomic-QA entity set with available Wikipedia pages. *Wiki-Rewrite* reformulates these instructions to increase diversity, serving as a variant derived from Wiki. In addition, *Atomic-LIFD* is a compact distilled subset of Atomic-QA obtained via perplexity-based selection, with construction details provided in Appendix D.1.

Training Configuration. All models are fine-tuned using Low-Rank Adaptation (LoRA) with rank $r = 16$ and a weight decay of 0.01 for 2 epochs. We adopt a learning rate of 1×10^{-4} by default. For experiments involving scaled training data, specifically Atomic-QA MCQ, perplexity analysis, and CKA variance analysis, we use a learning rate of 1×10^{-5} , as shown respectively in Figure 5, Figure 12, and Figure 13.

Table 2: Evaluation benchmarks and configurations.

Category	Benchmark	Shots	Metric
Text Comprehension and Information Integration	BoolQ	3	Accuracy
	LongBench-HotpotQA	0	F1
	RACE	3	Accuracy
	SQuADv2	3	F1
Knowledge-Driven Reasoning	ARC-Challenge	3	Accuracy
	ARC-Easy	3	Accuracy
	OpenBookQA	3	Accuracy
Commonsense Reasoning	PIQA	3	Accuracy
	SIQA	3	Accuracy
	WinoGrande	3	Accuracy
Mathematical Reasoning	GSM8K	3	Exact Match
	Hendrycks Math500	3	Exact Match
	MathQA	3	Accuracy

Evaluation Configuration. We evaluate Atomic-QA alongside other dataset baselines across a benchmark suite spanning multiple capability dimensions. In the domain of text comprehension and information integration, assessment is conducted using BoolQ (Clark et al., 2019), LongBench-HotpotQA (Bai et al., 2024), RACE (Lai et al., 2017), and SQuADv2 (Rajpurkar et al., 2018). With regard to knowledge-driven reasoning, capabilities are measured through ARC-Easy/Challenge (Clark et al., 2018) and OpenBookQA (Mihaylov et al., 2018). Additionally, for commonsense reasoning, evaluation is performed via PIQA (Bisk et al., 2019), SIQA (Sap et al., 2019), and WinoGrande (Sakaguchi et al., 2019).

To analyze the impact of increasing atomic error rates (from 0% to 100%) on model performance, we evaluate datasets generated via the Atomic-X pipeline. Specifically, Atomic-Alpaca is tested on the identical ten benchmarks introduced above to observe how different error levels influence general capability dimensions. In parallel, Atomic-MetaMath is evaluated on specialized mathematical reasoning tasks, including GSM8K (Cobbe et al., 2021), Math500 (Hendrycks et al., 2021), and MathQA (Amini et al., 2019). Regarding the evaluation protocol, the majority of benchmarks are conducted under a 3-shot setting, while LongBench-HotpotQA and Atomic-QA MCQ follow a zero-shot paradigm. Results on Atomic-QA MCQ are reported in Figure 5, and detailed evaluation configurations for all benchmarks are summarized in Table 2.

D Additional Results

In this section, we present additional experimental results, structured to mirror the organization of the main text. We begin with supplementary find-

ings on Quantifying the Cost of Strict Filtering, followed by further details on Measuring Operator-Oper and Dominance.

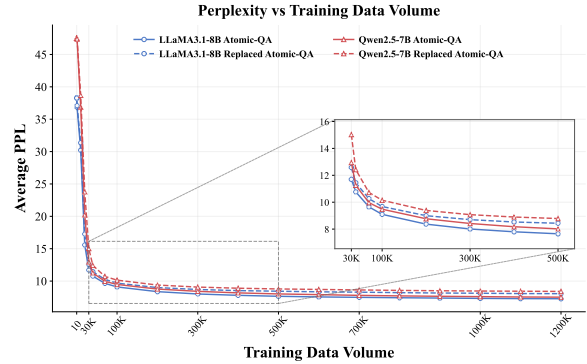


Figure 12: Perplexity trends of LLaMA and Qwen models as a function of training data volume. The curve displays average PPL for both the original and entity-replaced variants, with an inset showing a zoomed region for the 30K–500K data range.

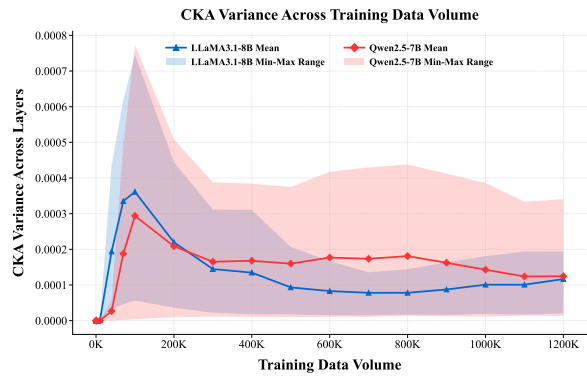


Figure 13: CKA variance across layers versus training data volume for LLaMA3.1-8B and Qwen2.5-7B, showing mean values (lines) and min–max ranges (shaded areas).

D.1 Training Tokens Comparison

Table 4 reports the total number of training tokens for each dataset. The datasets vary substantially in scale, ranging from fewer than one million tokens (e.g., LIMA) to over two hundred million tokens (e.g., Magpie), with several datasets clustered in the range of a few million to several tens of millions of tokens. Given this substantial variation, it is important to determine whether dataset scale or structural diversity serves as the primary determinant of performance. To investigate this question, we construct Atomic-LIFD, a compact dataset of approximately 2.7M tokens that is larger than LIMA but smaller than other considered datasets. Specifically, this subset is derived by applying perplexity-based

Table 3: Overall performance across all training datasets on Qwen2.5-7B and Qwen2.5-14B. Color-coded cells (green/red) are used to indicate performance above or below the mean across methods (All Average), with darker colors representing better or worse performance.

Model Family		ARC-Challenge	ARC-Easy	BoolQ	Hotpot-QA	Open-BookQA	PIQA	RACE	SIQA	SQuADv2	Wino-Grande	Avg
Qwen2.5-7B	Base	0.590	0.860	0.872	0.216	0.390	0.803	0.459	0.576	0.421	0.754	0.594
	All Average	0.586	0.859	0.875	0.190	0.399	0.800	0.474	0.572	0.364	0.747	0.587
	Instruct	0.616	0.873	0.865	0.143	0.404	0.786	0.512	0.566	0.324	0.735	0.582
	LIMA	0.586	0.861	0.871	0.207	0.396	0.801	0.459	0.574	0.424	0.757	0.594
	Condor	0.597	0.864	0.874	0.185	0.392	0.801	0.471	0.580	0.396	0.755	0.591
	Magpie	0.582	0.857	0.876	0.241	0.364	0.795	0.455	0.572	0.412	0.740	0.589
	WizardLM	0.606	0.863	0.879	0.050	0.412	0.811	0.464	0.584	0.259	0.756	0.568
	Alpaca	0.586	0.862	0.883	0.037	0.392	0.808	0.482	0.587	0.265	0.747	0.565
	Self-Instruct	0.490	0.816	0.859	0.062	0.376	0.792	0.430	0.524	0.304	0.723	0.537
	Wiki	0.583	0.868	0.876	0.248	0.398	0.803	0.470	0.576	0.426	0.763	0.601
	Wiki-Rewrite	0.584	0.870	0.876	0.204	0.404	0.802	0.470	0.576	0.422	0.751	0.596
	Atomic-LIFD	0.606	0.859	0.882	0.383	0.432	0.797	0.500	0.576	0.392	0.743	0.617
Atomic-QA	0.609	0.859	0.882	0.328	0.414	0.801	0.507	0.581	0.386	0.751	0.612	
Qwen2.5-14B	Base	0.637	0.878	0.888	0.147	0.404	0.822	0.487	0.579	0.392	0.793	0.603
	All Average	0.630	0.879	0.884	0.167	0.409	0.816	0.508	0.586	0.393	0.791	0.606
	Instruct	0.699	0.902	0.883	0.160	0.430	0.821	0.577	0.611	0.411	0.785	0.628
	LIMA	0.634	0.880	0.886	0.139	0.404	0.820	0.484	0.582	0.395	0.800	0.602
	Condor	0.631	0.879	0.885	0.133	0.396	0.818	0.490	0.588	0.379	0.800	0.600
	Magpie	0.630	0.882	0.893	0.037	0.368	0.811	0.496	0.580	0.415	0.776	0.589
	WizardLM	0.640	0.888	0.893	0.076	0.400	0.824	0.520	0.590	0.351	0.789	0.597
	Alpaca	0.609	0.880	0.898	0.143	0.394	0.822	0.500	0.597	0.427	0.796	0.607
	Self-Instruct	0.518	0.826	0.850	0.263	0.368	0.800	0.463	0.534	0.319	0.763	0.570
	Wiki	0.629	0.880	0.874	0.068	0.416	0.817	0.508	0.597	0.417	0.807	0.601
	Wiki-Rewrite	0.626	0.880	0.882	0.063	0.408	0.820	0.515	0.592	0.416	0.811	0.601
	Atomic-LIFD	0.653	0.878	0.894	0.414	0.466	0.808	0.507	0.582	0.392	0.770	0.636
Atomic-QA	0.664	0.890	0.882	0.337	0.450	0.816	0.534	0.595	0.398	0.807	0.637	

Table 4: The number of training tokens in different datasets.

Dataset	Training Tokens
LIMA	627,993
Condor	21,131,550
Magpie	203,256,865
WizardLM	27,673,025
Alpaca	3,766,455
Self-Instruct	4,536,225
Wiki	11,082,377
Wiki-Rewrite	11,356,362
Atomic-QA	34,901,525
Atomic-LIFD	2,701,510

selection (IFD) to Atomic-QA, a logically verified and structurally diverse source.

Remarkably, Atomic-LIFD achieves performance comparable to that of the full Atomic-QA dataset, as shown in Table 1 and Table 3. Given that data scale is the sole variable distinguishing these two settings, this minimal performance gap strongly suggests that data quantity is not the decisive factor behind Atomic-QA’s superiority over other datasets.

D.2 Quantifying the Cost of Strict Filtering

Table 3 extends our evaluation to Qwen2.5-7B and Qwen3-8B to verify the robustness of our approach. Consistent with the findings in the main text, Atomic-QA demonstrates strong generalizability, achieving the highest average performance

across all baselines on both backbones (61.2% on Qwen2.5-7B and 62.5% on Qwen3-8B).

The persistent performance gap over LIMA (+1.8% to +2.3%) and other synthetic methods (Alpaca, WizardLM) confirms that our conclusion is model-agnostic: prioritizing structural scale over strict filtering consistently yields better instruction-following and reasoning capabilities, regardless of the underlying model architecture or size.

D.3 Measuring Operator-Operand and Dominance

We further validated our findings on the Qwen model family, as illustrated in Figures 14, 15, and 16. Across varying model sizes (7B, 8B, and 14B), Replaced Atomic-QA consistently aligns with standard Atomic-QA performance with negligible degradation (e.g., an average drop of only 0.002 for Qwen3-8B). This stability demonstrates that robustness to operand errors is a generalizable phenomenon, independent of specific model architectures.

To substantiate our main findings, we visualize the training dynamics in Figures 12 and 13. Figure 12 shows that the perplexity curves for the Replaced Atomic-QA and the original dataset are nearly indistinguishable, demonstrating that operand errors impose almost no additional optimization burden. This indicates the model learns the structural logic with equal efficiency. Furthermore, Figure 13 illus-

Qwen2.5-7B Performance: Atomic-QA vs Replaced Atomic-QA vs Rewrite-to-Hallucination

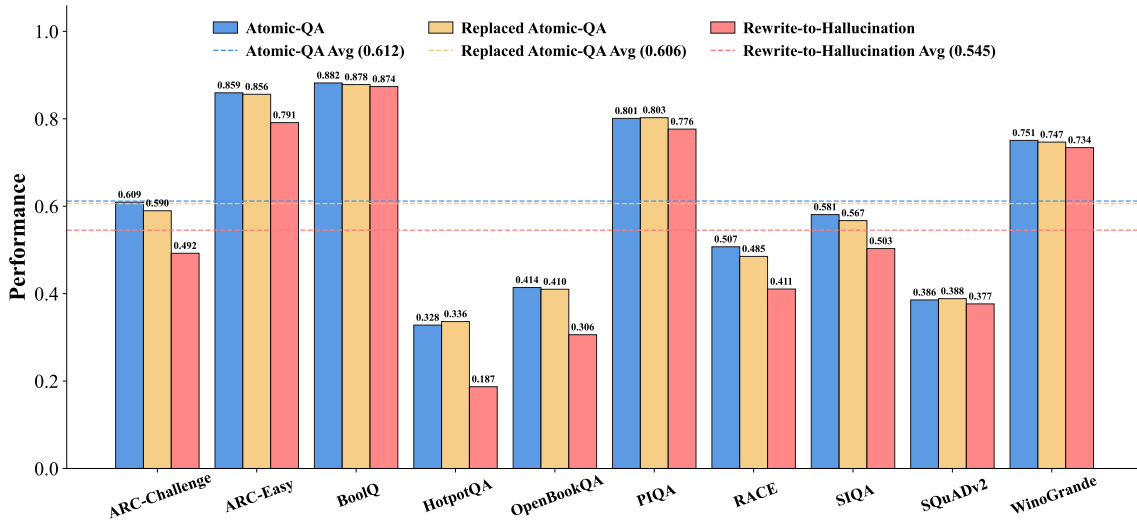


Figure 14: Performance comparison of Qwen2.5-7B across the Atomic-QA, Atomic-QA Replaced, and Rewrite-to-Hallucination settings. Each bar shows the model’s accuracy on individual tasks, while the dashed lines denote the average performance for each setting.

Qwen2.5-14B Performance: Atomic-QA vs Replaced Atomic-QA vs Rewrite-to-Hallucination

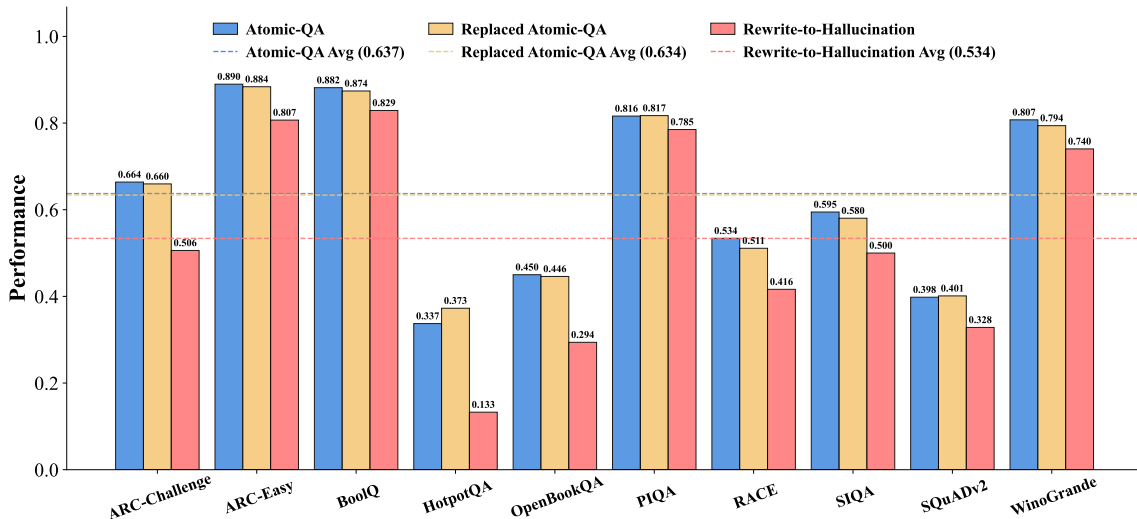


Figure 15: Performance comparison of Qwen2.5-14B across the Atomic-QA, Atomic-QA Replaced, and Rewrite-to-Hallucination settings. Each bar shows the model’s accuracy on individual tasks, while the dashed lines denote the average performance for each setting.

1157 rates that the variance in CKA similarity initially
 1158 spikes but subsequently stabilizes and decreases as
 1159 data volume increases. This convergence confirms
 1160 that the high representational alignment discussed
 1161 in the main text is a robust and stable outcome of
 1162 training, rather than a transient state.

1163 Finally, to verify the generalizability of the
 1164 degradation patterns in operation errors, we ex-
 1165 tended our analysis to mathematical reasoning us-
 1166 ing the MetaMath dataset. As illustrated in Figure
 1167 17, consistent with the collapse observed in AI-

1168 paca, both LLaMA 3.1-8B and Qwen 2.5-7B even-
 1169 tually exhibit a sharp performance decline. How-
 1170 ever, a key distinction lies in the critical threshold.
 1171 While models trained on Alpaca begin to collapse
 1172 around an error rate of 0.6, those trained on Meta-
 1173 Math demonstrate higher robustness, maintaining
 1174 stability until the error rate approaches saturation
 1175 (> 0.9). This confirms that while the non-linear
 1176 degradation dynamics are universal across tasks,
 1177 the tipping point for structural collapse depends
 1178 on domain complexity. Overall, the models exhibit

Qwen3-8B Performance: Atomic-QA vs Replaced Atomic-QA vs Rewrite-to-Hallucination

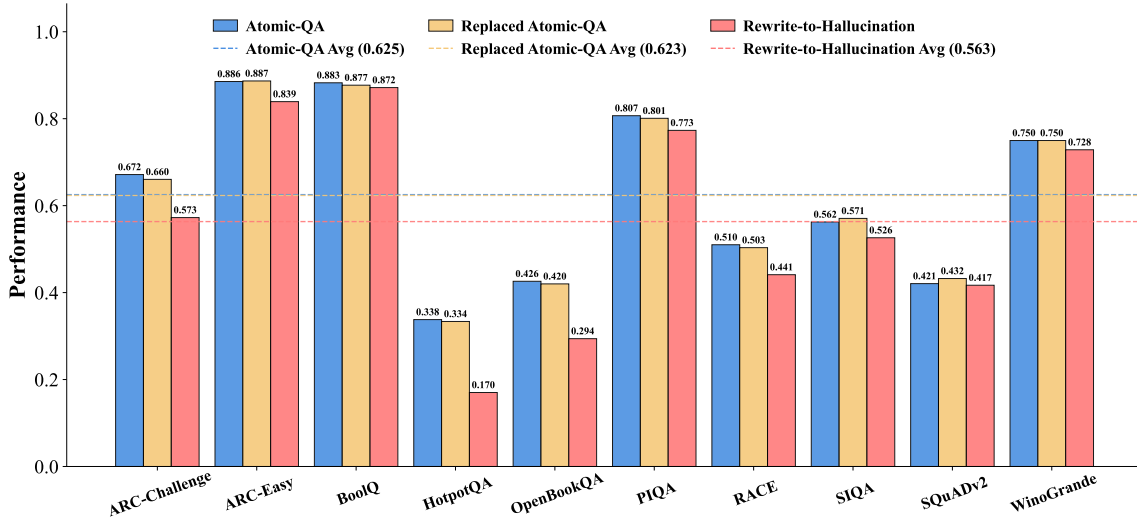


Figure 16: Performance comparison of Qwen3-8B across the Atomic-QA, Atomic-QA Replaced, and Rewrite-to-Hallucination settings. Each bar shows the model’s accuracy on individual tasks, while the dashed lines denote the average performance for each setting.

exceptional resilience to operand perturbations. Remarkably, even when the operand error rate reaches almost 100%, performance remains largely unaffected. In stark contrast, a 100% error rate in operators precipitates a significant performance decline across both the Alpaca and MetaMath datasets.

E Beta-Distribution Based Error Assignment

In this section, we present a Beta-Distribution based approach for assigning errors to samples while maintaining a target global error rate. We first formulate the problem, then introduce our stochastic assignment method that leverages the Beta distribution to create realistic error diversity across samples.

E.1 Problem Formulation

We consider a dataset of M samples, where sample i comprises n_i atomic functions. Our objective is to introduce errors while matching a prescribed global error rate $r \in [0, 1]$ at the atomic-function level. The total number of atomic functions is $N = \sum_{i=1}^M n_i$.

For each sample i , we pre-generate $n_i + 1$ candidate versions indexed by $k \in \{0, 1, \dots, n_i\}$, where the k -th version contains exactly k corrupted functions. Each version is constructed by introducing one additional corruption relative to the $(k - 1)$ -th version. Exactly one version must be selected for

Model Performance: Accuracy vs Sample Error Rate (MetaMath)

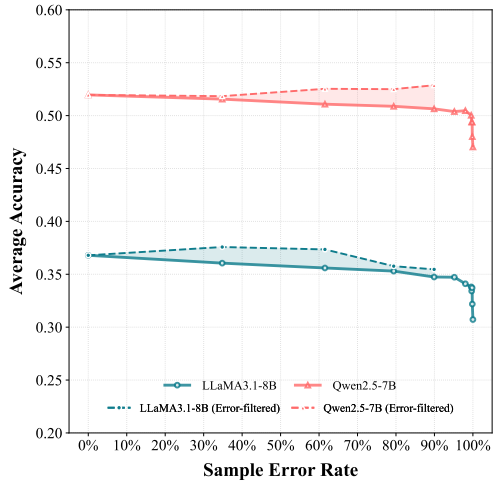


Figure 17: Accuracy degradation under increasing simulated error rates for LLaMA 3.1-8B and Qwen 2.5-7B on MetaMath. The x-axis denotes the proportion of erroneous data in the training set, while the y-axis represents the performance of the trained model.

each sample to satisfy the following condition:

$$\sum_{i=1}^M e_i = \text{round}(r \cdot N),$$

where e_i is the number of errors in the selected version of sample i , and $\text{round}(\cdot)$ rounds to the nearest integer.

E.2 Beta Distribution Based Assignment

Motivation. A naive uniform approach assigns errors deterministically by setting $e_i = \lfloor r \cdot n_i \rfloor$ for

each sample. While this guarantees the global error rate, it produces homogeneous error distributions where every sample exhibits approximately the same error rate. Such uniformity fails to capture the natural heterogeneity observed in real-world data, where some samples are nearly perfect while others contain multiple errors. To address this limitation, we introduce controlled stochasticity through the Beta distribution.

Beta Distribution Properties. The Beta distribution is a continuous probability distribution defined on the interval $[0, 1]$, making it particularly well suited for modeling proportions and probabilities. Its flexible shape is controlled by two parameters α and β . Under the parameterization $\alpha = r\kappa$ and $\beta = (1 - r)\kappa$, the distribution has mean

$$\mathbb{E} = \frac{\alpha}{\alpha + \beta} = r$$

and the concentration parameter κ controls the shape of the distribution.

Stochastic Sampling Phase. We model each sample’s error probability using a Beta distribution parameterized by the target error rate r and concentration parameter κ . For each sample i , we independently draw an error probability

$$q_i \sim \text{Beta}(r\kappa, (1 - r)\kappa).$$

The parameterization ensures that the expected per-sample error probability satisfies $\mathbb{E}[q_i] = r$, so that the expected global error rate converges to the target value r . This stochastic sampling introduces natural heterogeneity across samples, reflecting realistic scenarios in which data quality varies at the instance level. In our experiments, we set $\kappa = 10$, which provides moderate variance and produces the relationship between atomic-level error rates and the proportion of corrupted samples shown in Figure 19.

Normalization, Adjustment, and Quantization. To ensure that the expected total number of errors matches the target, we first compute the aggregate statistic

$$S = \sum_{j=1}^M n_j q_j,$$

and derive the normalization constant as

$$c = \frac{r \cdot N}{S}.$$

The adjusted error probability for sample i is then given by :

$$p_i = \text{clip}(c \cdot q_i, 0, 1).$$

This correction enforces $\sum_{i=1}^M n_i p_i \approx r \cdot N$ while maintaining the sample-level heterogeneity introduced by the Beta sampling. The clipping operation ensures all probabilities remain within $[0, 1]$.

Since error counts must be integers, we set the error count for each sample to

$$e_i = \text{round}(n_i \cdot p_i),$$

which rounds $n_i \cdot p_i$ to the nearest integer, and compute the residual $\Delta = \text{round}(r \cdot N) - \sum_{i=1}^M e_i$. To correct this discretization error, we iteratively adjust error assignments to ensure that the final assignment satisfies

$$\sum_{i=1}^M e_i = \text{round}(r \cdot N).$$

The complete procedure is detailed in Algorithm 2, where $\text{sign}(x)$ denotes the sign function that returns $+1$ if $x > 0$, -1 if $x < 0$, and 0 if $x = 0$.

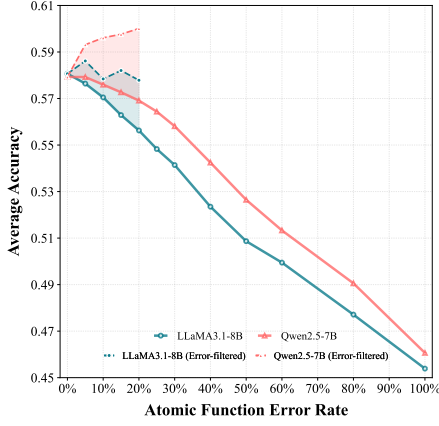
F Prompts

This section enumerates all prompts utilized in our framework. Prompts 1, 2, 3, and 4 detail the construction process of the Atomic Tree. Subsequently, Prompts 5 and 6 outline the creation of Atomic-QA, while Prompts 7, 8, and 9 correspond to the construction of Atomic-X.

G Examples of Atomic Function

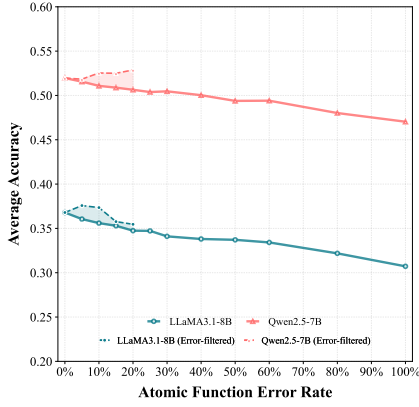
In this section, we present examples of atomic functions. Figures 20, 21, and 22 illustrate instances from the Atomic-QA, Alpaca, and Meta-math datasets, respectively. In each figure, the leftmost panel displays the original input prompt and its corresponding response. The central columns list the extracted atomic functions followed by the perturbed atomic functions. Finally, the rightmost panel shows the output response after perturbation.

Model Performance: Accuracy vs Atomic Function Error Rate (Alpaca)



(a) Alpaca accuracy degradation under the perspective of erroneous atomic functions ratio.

Model Performance: Accuracy vs Atomic Function Error Rate (MetaMath)



(b) MetaMath accuracy degradation under the perspective of erroneous atomic functions ratio.

Figure 18: Accuracy degradation under increasing atomic function error rates for (a) LLaMA and (b) Qwen. The x-axis denotes the proportion of erroneous atomic functions in the training set, while the y-axis represents the performance of the trained model.

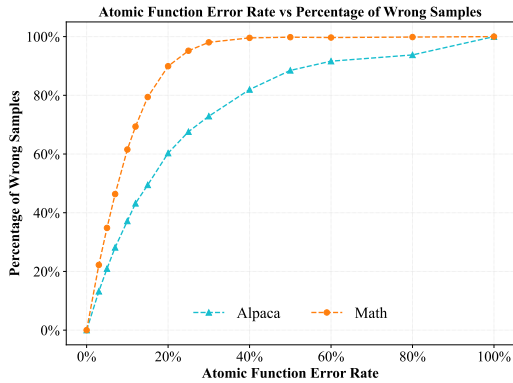


Figure 19: Relationship between the atomic function error rate and the resulting percentage of wrong samples.

Algorithm 2: Beta Distribution Error Assignment

Input : $\mathbf{n} = \{n_1, n_2, \dots, n_M\} \in \mathbb{N}^M$,
atomic function counts for M
samples

Input : $r \in [0, 1]$, target global error rate

Param : $\kappa \in \mathbb{R}_+$, concentration parameter
(default: $\kappa = 10.0$)

Output : $\mathbf{e} = \{e_1, e_2, \dots, e_M\} \in \mathbb{N}^M$,
error counts per sample, where
 $\sum_i e_i = \text{round}(r \cdot N)$

- 1 $N \leftarrow \sum_{i=1}^M n_i$
 - 2 $T \leftarrow \text{round}(r \cdot N)$
 - 3 **forall** $i \in \{1, \dots, M\}$ **do**
 - 4 $q_i \sim \text{Beta}(r \cdot \kappa, (1 - r) \cdot \kappa)$
 - 5 $S \leftarrow \sum_{j=1}^M n_j \cdot q_j$
 - 6 $c \leftarrow \frac{r \cdot N}{S}$
 - 7 **forall** $i \in \{1, \dots, M\}$ **do**
 - 8 $p_i \leftarrow \text{clip}(c \cdot q_i, 0, 1)$
 - 9 $e_i \leftarrow \text{round}(n_i \cdot p_i)$
 - 10 $\Delta \leftarrow T - \sum_{i=1}^M e_i$
 - 11 **if** $\Delta \neq 0$ **then**
 - 12 **if** $\Delta > 0$ **then**
 - 13 $\sigma \leftarrow \text{argsort}(-\mathbf{n})$
 - 14 **else**
 - 15 $\sigma \leftarrow \text{argsort}(\mathbf{n})$
 - 16 **for** $i \in \sigma$ **do**
 - 17 **if** $\Delta = 0$ **then**
 - 18 **break**
 - 19 **if** $0 \leq e_i + \text{sign}(\Delta) \leq n_i$ **then**
 - 20 $e_i \leftarrow e_i + \text{sign}(\Delta)$
 - 21 $\Delta \leftarrow \Delta - \text{sign}(\Delta)$
 - 22 **return** $\mathbf{e} = \{e_1, e_2, \dots, e_M\}$
-

Prompt 1: Atomic Function Extraction

You are the ENCODER agent. Your job is to extract atomic functions from the Sample Output and represent each input as (operand, operator, output) triples.

```
{atomic function definition}
```

Return **exactly one JSON array** of atomic functions. If no atomic functions can be extracted, return [].

```
## Atomic Function Definition  
{atomic function definition}
```

```
## Rules
```

- 1) Minimality: one predicate per triple. Do not merge multiple assertions. If the current sample is already the smallest unit satisfying the definition, no splitting is required.
- 2) Order: emit triples in the order their evidence appears in the Sample Output.
- 3) Numbers/dates: keep literal values and units as written (e.g., "18", "2019", "3.14", "USD")
- 4) Determinism: no randomness; if uncertain, omit the item.
- 5) Strict output: JSON array only, no comments, no trailing commas, no extra keys.

```
## Sample Input
```

This section provides the original input of the data sample. Such as a question, instruction, or code comment.
{sample input}

```
## Sample Output
```

This section provides the original answer of the data sample. It may contain multiple sentences, equations, or code snippets.
{sample output}

```
## Output Format
```

Your answer should be wrapped in a proper JSON code block using triple backticks like this:

```
```json  
{example}
```
```

Prompt 2: Initial Tree Generation

The following are anatomic functions, followed by their definitions.

Input anatomic functions:
{anatomic functions}
The definitions of the anatomic functions:
{definitions of anatomic functions}

Requirements:

1. Format the output as a JSON tree structure.
2. Consider both coarse-grained and fine-grained classification of anatomic functions' categories, building parent-child relationships.
3. Each node is a dict with category names as keys.
4. The root node should be {root name}.
5. Leaf nodes should contain empty dicts.
6. Wrap the JSON with <begin> and <end>.

Think step by step:

1. Identify the categories of the anatomic functions
2. Organize the categories (not anatomic functions) into a hierarchical tree structure.

Output format:

```
<think>
[your reasoning process]
</think>
<begin>
{example}
<end>
```

Prompt 3: Tree Duplication

Please perform strict deduplication on the following subcategories under the parent category "{parent category}".

Requirements:

1. Remove any semantic duplicates (even if names are different)
2. Ensure all items strictly belong to "{parent category}", remove any unrelated items
3. Eliminate any hierarchical relationships between subcategories, remove subcategories that are included in other subcategories
4. Keep the most appropriate names
5. Output the cleaned list

Current subcategories:

```
{subcategories}
```

Output format:

```
```json
{{"unique_categories": ["example1", "example2"]}}
```
```

Prompt 4: Tree Expansion

As a taxonomy expert, suggest 3-5 possible {width/depth} expansion candidates for category: [{category}]

Related Nodes:

- Parent: {parent}
- Siblings: {siblings}

Definitions:

- Current Node: {definition}
- Parent: {parent_definition}
- Siblings: {siblings_definitions}

Guidelines:

1. Ensure logical consistency with parent/sibling definitions
2. Avoid overlaps with existing nodes
3. Use clear and specific naming

Output JSON format:

```
```json
{"new_nodes": ["name1", "name2"]}
```
```

Prompt 5: Operand Generation

You are given the following information:

- Category
- Existing Operand

Your task is to generate new operands that fit the category based on the context provided.

Your generation should consider the following aspects:

1. The generated operands should be specific and fall under the given category.
2. The generated operands should not be similar to the existing operands.
3. Each generated operand should have a brief description that help understand the operands.

Category

This section contains the category that you need to generate operands for.

- Category Name: {category}
- Category Definition: {definition}

Existing operands

This section contains the existing operands that are already classified under the category.

You should avoid generating operands that are similar to the existing ones.

- Existing operands:

{existing operands}

Output format

You should output the generated operands in the JSON format, wrapped in a proper Json code block using triple backticks like this:

```
```json
{{
 "operand_name1": "Description of operand 1",
 "operand_name2": "Description of operands 2",
 ...
}}
```

### Prompt 6: Verification

You are given the following information:

- Function: {function}
- Operand: {operand}
- Output: {output}

Your task is to verify whether the function is logically consistent, that is, whether **Function(Operand)** equals the **Output**.

## Output format ##

You should output the generated operands in the JSON format, wrapped in a proper Json code block using triple backticks like this:

```
```json
{{
  "judgement": "pass"/"fail",
  "reason": "reason for your judgement"
}}
```

Prompt 7: Encoder

You are the ENCODER agent. Your job is to extract atomic functions from the Sample Output and represent each atomic function as a (subject, relation, object) triple.

Categories

This chapter enumerates the possible candidate categories for decomposing atomic functions. You need to determine whether the **## Payload** can be further decomposed within some category.

{categories}

Sample Input

This section provides the original input of the data sample. Such as a question, instruction, or code comment.

{sample input}

Sample Output

This section provides the original answer of the data sample. It may contain multiple sentences, equations, or code snippets.

{sample output}

Payload

This section provides the original answer of the data sample. It may contain multiple sentences, equations, or code snippets.

{sample output}

Output Format

Your answer should wrapped in a proper Json code block using triple backticks like this:

```
```json
{example}
```
```

Prompt 8: Perturber

You are the PERTURBER agent. Your job is to apply controlled edits to specific atomic functions in the Sample Output, guided by the provided (subject, relation, object) triple and payload location.

A *controlled edit* is a minimal, deterministic change that makes the function incorrect while preserving overall coherence and fluency.

Return **exactly one JSON object** describing the edit.

Use the Sample Input only to preserve coherence (e.g., pronouns); all edits **must** be grounded in the Sample Output and the given function.

Sample Input

This section provides the original input of the data sample (e.g., question/instruction/code comment).

{input_sample}

Sample Output

This section provides the original answer of the data sample (sentences, equations, or code).

{output_sample}

Atomic function to be Perturbed (from ENCODER)

This section provides the function to be perturbed, along with its unique fid and the corresponding payload location in the Sample Output.

function: {triple}

Payload: {payload}

Output Format

Your answer should be wrapped in a proper Json code block using triple backticks like this:

Output Format

Your answer should be wrapped in a proper Json code block using triple backticks like this:

```
```json
{example}
```
```

Prompt 9: Decoder

You are the DECODER agent. Your job is to reconstruct an output sample from the perturbed atomic function that is provided.

You should use the SEARCH-REPLACE strategy to ensure the perturbed function is correctly integrated into the Sample Output while maintaining overall coherence and fluency.

You are allowed to use the following tools, and you must use one tool at a time:

- SEARCH-REPLACE: Locate the content in the Sample Output and replace it with the given content.

For each change, follow this format exactly:

```
```[TOOL_CALL]
```

```
>>> SEARCH
```

```
<original snippet copied from the Sample Output>
```

```
=====
```

```
>>> REPLACE
```

```
<modified snippet you want to replace the original snippet with>
```

```
=====
```

```
```
```

- TERMINATE: Indicate the end of the editing process and output the final revised Sample Output.

To finish the editing process, follow this format exactly:

```
```[TOOL_CALL]
```

```
Terminate(reason="The reason for termination, e.g., no more changes needed")
```

```
```
```

Sample Input

This section provides the original input of the data sample (e.g., question/instruction/code comment).

```
{input_sample}
```

Sample Output

This section provides the original answer of the data sample (sentences, equations, or code).

```
{output_sample}
```

Perturbed Atomic function (from PERTURBER)

This section provides the perturbed function along with the type of edit applied. You should use change the corresponding content that matches the Original function in the Sample Output to the Perturbed function.

Original function: {pre_triple}

Perturbed function: {post_triple}

Output Format

IMPORTANT: Only use one tool call per response. Do not use multiple tool calls in a single response or code block.

IMPORTANT: You ONLY need to edit the specific function that is perturbed. Do NOT make any other changes to the Sample Output. Do NOT consider the effect of edited function on other parts of the Sample Output.

Your tool calls should be wrapped in a proper [TOOL_CALL] code block using triple backticks like this:

```
```[TOOL_CALL]
```

```
<your tool call here>
```

```
```
```



Figure 20: The example of Atomic-QA.

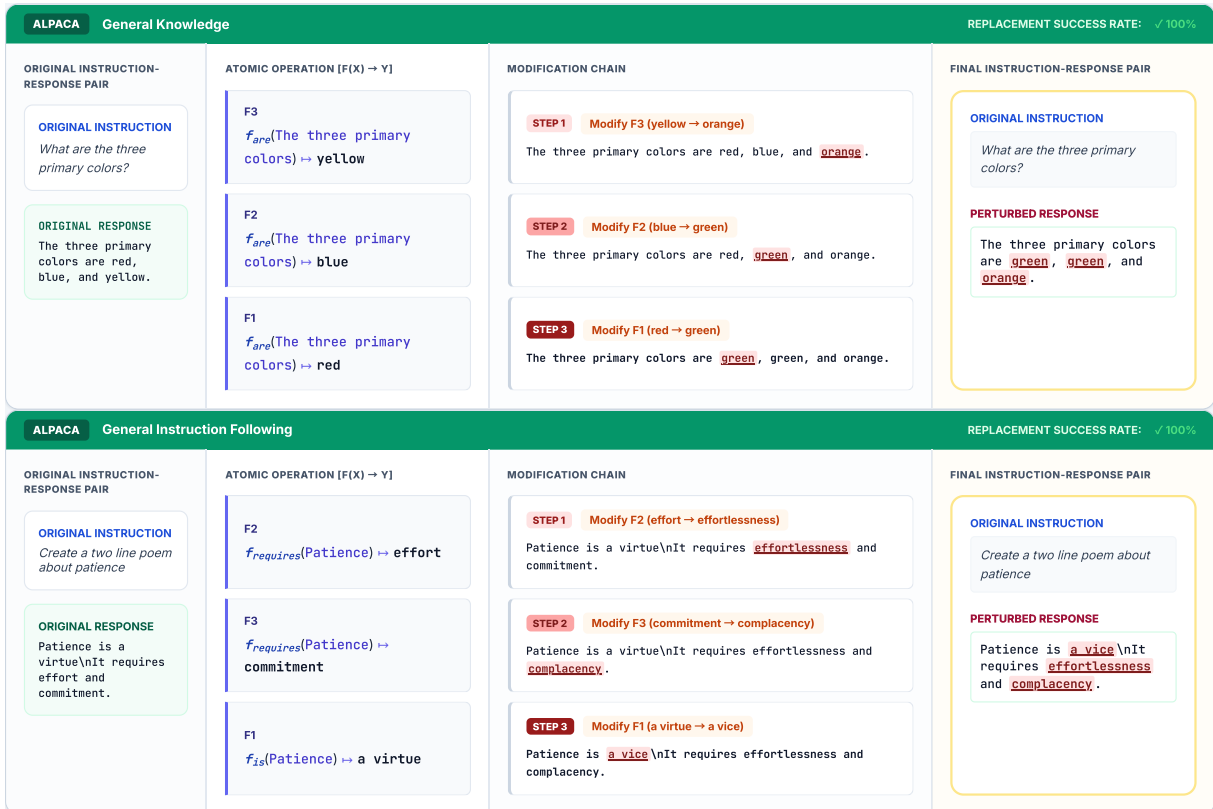


Figure 21: The example of Alpaca.

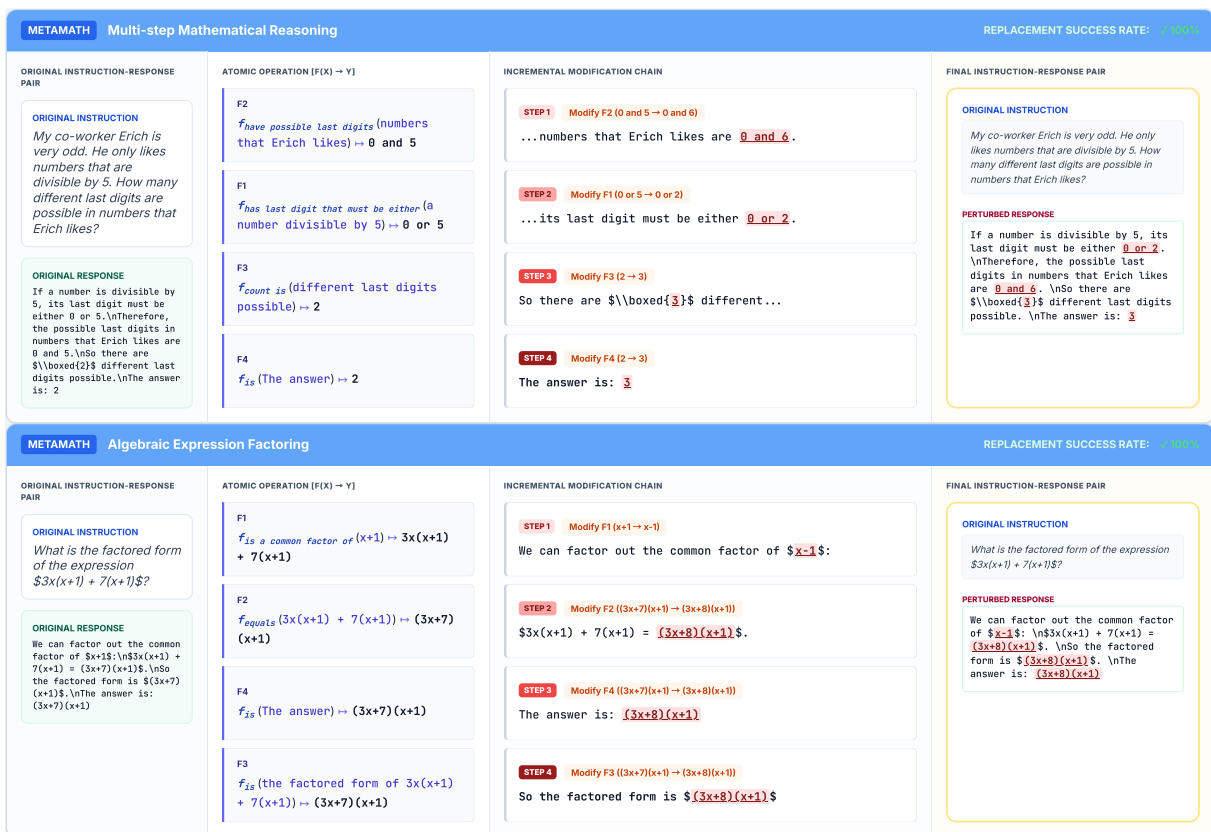


Figure 22: The example of MetaMath.