LIMITED REFERENCE, RELIABLE GENERATION: RULE-GUIDED TABULAR DATA GENERATION WITH DUAL-GRANULARITY FILTERING

Anonymous authorsPaper under double-blind review

000

001

002

004

006

008 009 010

011 012 013

014

015

016

017

018

019

021

023

025

026

027

028

031 032 033

034

037

040

041

042

043

044

045

046

047

048

051

052

ABSTRACT

Synthetic tabular data generation is increasingly essential in machine learning, supporting downstream applications when real-world and high-quality tabular data is insufficient. Existing tabular generation approaches, such as generative adversarial networks (GANs), diffusion models, and fine-tuned Large Language Models (LLMs), typically require sufficient reference data, limiting their effectiveness in domain-specific datasets with scarce records. While prompt-based LLMs offer flexibility without parameter tuning, they often generate distributionally drifted data with localized redundancy, leading to degradation in downstream task performance. To overcome these issues, we propose **ReFine**, a framework that (i) derives symbolic *if*—then rules from interpretable models and embeds them into prompts to explicitly guide the generation process toward the domainspecific distribution, and (ii) applies a dual-granularity filtering that suppresses over-sampling patterns and selectively refines rare but informative samples to reduce localized redundancy. Extensive experiments on various regression and classification benchmarks demonstrate that ReFine consistently outperforms state-ofthe-art methods, achieving up to 0.36 absolute improvement in R^2 for regression and 7.50% relative improvement in F_1 for classification tasks.

1 Introduction

Tabular data serves as a foundational modality in machine learning, underpinning critical applications in domains such as healthcare, finance, and scientific research (Benjelloun et al., 2020; Ghosh, 2012; Yang et al., 2024; Shankar et al., 2024). However, the collection of large-scale, high-quality tabular datasets is often hindered by strict privacy regulations and the prohibitive costs of expert-driven annotation (Voigt & Von dem Bussche, 2017; Miceli et al., 2020). Such limitations severely restrict effective model training in many tabular applications, thereby motivating the use of synthetic data generation (Hernandez et al., 2022; Shankar et al., 2024).

Previous mainstream methods for tabular data generation are based on non-LLM generative models such as VAEs (Xu et al., 2019), GANs (Zhao et al., 2021) and diffusion models (Kotelnikov et al., 2023). Within the LLM-based paradigm, fine-tuning approaches have demonstrated notable performance (Borisov et al., 2023). Both approaches share a fundamental prerequisite: access to sufficient reference data, which used as the foundation for learning underlying distributions. In practice, this prerequisite is often violated in high-stakes domains, where extremely strict privacy regulations and the rarity of critical events make large-scale, high-quality tabular datasets difficult to obtain (Kovalerchuk & Vityaev, 2005; Ji et al., 2014). For example, in rare disease diagnosis, available datasets may contain only a few dozen records, a scenario commonly referred to as a low-data regime (Seedat et al., 2023), which poses severe challenges for data-driven modeling (Raghavan & El Gayar, 2019; Li et al., 2023; Wang et al., 2024a). Consequently, in low-data regimes where only a handful of samples are available, existing generative methods fail to capture underlying distributions and thus struggle to produce high-quality synthetic data (Bommareddy et al., 2022; Fang et al., 2024a; Zhang et al., 2024a) In contrast, prompt-based methods exploit in-context learning to synthesize data without training, offering an alternative for low-data regimes (Seedat et al., 2023; Kim et al., 2024). However, prompt-based methods face two challenges in low-data regimes:

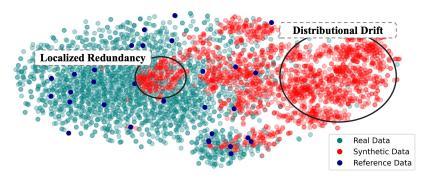


Figure 1: Two key challenges in prompt-based LLM tabular data generation in low-data regimes: (i) **Distributional Drift**: LLMs often generate biased samples by relying on spurious patterns from pretraining rather than real data distribution. (ii) **Localized Redundancy**: Synthetic samples tend to concentrate in limited regions of the feature space due to repeated use of identical prompts.

- (i) **Distributional Drift**: Prompt-based methods often over-rely on LLMs' pretrained knowledge, which poorly captures dataset-specific distribution (Zhong et al., 2024). Consequently, generated samples tend to follow spurious patterns inherited from pretraining corpora rather than the true patterns of the target dataset (Sahoo et al., 2024; Stoian et al., 2024a). This mismatch leads to synthetic distributions that drifted from the *real data*. Figure 1 illustrates this issue: when conditioned on a small reference set (blue), the LLM produces *synthetic data* (red) that drift away from the manifold of *real data* (cyan). Many samples populate low-density or unsupported regions, highlighting the LLM's inability to reconstruct the target joint distribution under data scarcity.
- (ii) **Localized Redundancy**: Prompt-based generation tends to overproduce high-frequency attribute combinations present in the reference data, while rare but informative patterns are rarely synthesized (Amatriain, 2024; Liu et al., 2023; Zevallos et al., 2023). This overconcentration, which we term localized redundancy, is further amplified when identical prompt templates are reused to generate data in multiple batches. As a result, synthetic samples cluster around a few dominant modes (Kim et al., 2024; Seedat et al., 2023), forming high-density regions (red) that contrast sharply with the broader and more balanced distribution of real data (cyan), as shown in Figure 1.

To systematically address the two key challenges of prompt-based tabular generation in low-data regimes, we propose *ReFine* (Rule-Guided Generation and Dual-Granularity Filtering), a framework comprising two components. To mitigate the distribution of synthetic data often drifted by LLMs in low-data regimes (Challenge i), we introduce **Rules-Guided Generation**, which extracts symbolic *if—then* formulas from interpretable tree-based models. These association rules are embedded into prompts to guide the LLM toward the distribution of the real data. To mitigate localized redundancy that persists despite prompt-based generation (Challenge ii), we propose **Dual-Granularity Filtering**. This component informs a two-level filtering process: chunk-level pruning of dominant high-density modes, and instance-level refinement to retain low-density but informative samples. Our key contributions can be summarized as follows:

- We identify two key challenges of LLM prompt-based methods in tabular data generation in low-data regimes: (i) distributional drift of the synthetic data; and (ii) localized redundancy in the synthetic data.
- To address the two challenges, we propose *ReFine*¹, a framework that constructs *association rules* to guide LLM for tabular data generation, and applies proxy-based distribution estimation with *dual-granularity* filtering to reduce localized redundancy.
- Experimental results demonstrate that ReFine consistently outperforms strong baselines, achieving up to 0.36 absolute gain in R^2 for regression and 7.5% relative improvement in F_1 for classification. Comprehensive ablations further highlight the respective contributions of Rules-Guided Generation and Dual-Granularity Filtering components.

Anonymous code: https://anonymous.4open.science/r/ReFine-5328

2 RELATED WORK

2.1 Non-LLM Tabular Generation Method

Generative Model-based generation. Many work on tabular data synthesis relied on GANs, diffusion models, and score-based models (Hernandez et al., 2022). Among classical methods, CT-GAN (Xu et al., 2019) extends GANs to handle mixed-type variables but suffers from mode collapse and requires heavy preprocessing. TabDDPM (Kotelnikov et al., 2023) applies diffusion for continuous attributes, yet its iterative denoising is computationally costly and unstable with scarce samples. TABSYN (Zhang et al., 2024b) integrates diffusion with a VAE backbone to better support mixed-type data, but still depends on sufficient training density to avoid spurious correlations. Overall, while effective in abundant-data settings, these models degrade sharply in low-data regimes.

Constraint-based tabular generation. Another line of work enforces domain validity through hard logical constraints. Early methods embed linear constraints into generative models via *Constraint Layers* Stoian et al. (2024b). More recently, *Disjunctive Refinement Layers (DRL)* extend this idea to quantifier-free real linear arithmetic (QFLRA), enabling non-convex and disjunctive feasible regions Stoian & Giunchiglia (2025). While effective for ensuring semantic validity, such methods face two key drawbacks: they rely on exhaustively specified domain rules, which is rarely feasible in practice, and hard constraints restrict the generation space, thereby limiting diversity.

2.2 LLM-based Tabular Generation Method

LLMs have recently gained attention for tabular data generation, which exploit pretrained knowledge to make them well-suited for structured data tasks. Existing approaches fall into two categories: *fine-tuning methods* and *prompt-based methods*.

Fine-Tuning Methods. Fine-tuning methods adapt LLM parameters to tabular formats and domain constraints. For instance, GReaT (Borisov et al., 2023) fine-tunes GPT-2.5 on tabular corpora, while HARMONIC (Wang et al., 2024b) introduces instruction signals derived from nearest-neighbor relationships. Although effective with abundant reference data, these methods risk severe overfitting when reference data is small, as parameter updates dominate the limited supervision.

Prompt-based Methods. Prompt-based methods leverage the in-context learning ability of LLMs, enabling them to generate tabular data by conditioning on a few labeled examples embedded directly in the prompt. Without modifying model parameters, these methods use prompt design to guide the generation process. EPIC improves representation balance across classes by formatting grouped data and crafting class-aware prompts (Kim et al., 2024). CLLM enhances data quality in low-resource scenarios by combining prompt design with a curation step that filters samples based on model confidence and uncertainty estimates (Seedat et al., 2023). However, prompt-based methods struggle to capture logical dependencies and often suffer from distributional imbalance due to repeated use of identical prompts, limiting their effectiveness in low-data regimes.

3 METHODOLOGY

Prompt-based tabular generation in low-data regimes suffers from two key issues: *Distributional Drift* and *Localized Redundancy*. To tackle these challenges, we unify our intuition and methodological pipeline in Figure 2. Here, the *Real Distribution* (cyan) denotes the target data manifold, while the *Original Small Dataset* (orange) provides limited supervision. To mitigate *Distributional Drift*, we propose (1) Rules-Guided Generation, which extracts *Association Rules* (red) from tree-based models to explicitly capture key feature dependencies. These rules are embedded into prompts, guiding the LLM toward generation that better align with the real distribution. However, static prompting still induces *Localized Redundancy*. To address this, we introduce (2) Dual-Granularity Filtering, which prunes dominant high-density chunks while retaining informative but rare samples at the instance level. Together, these two components expand the effective generation space and enhance the fidelity and diversity of the augmented dataset for downstream learning. In what follows, we first formally define the tabular generation task under low-data regimes. Then, we describe our two components for addressing the two primary challenges.

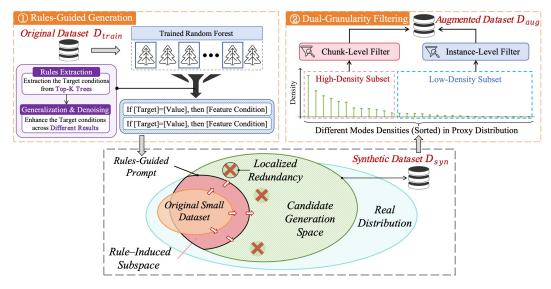


Figure 2: Overall framework of *ReFine* (*upper panel*), which consists of two components: (1) **Rules-Guided Generation**, which leverage rules to guide LLM generation toward the underlying *Target Distribution*; and (2) **Dual-Granularity Filtering**, which suppresses overrepresented patterns and preserves informative, low-density instances. The geometric view (*lower panel*) illustrates our insight—how rule-guided prompting anchors generation in a faithful *Generation Space*.

3.1 PROBLEM SETUP

Definition 3.1 Tabular Data Generation in Low-data Regimes. Let $X \subseteq \mathbb{R}^d$ be the d-dimensional feature space and Y the label space (discrete for classification, real-valued for regression). Let $D_{train} = \{(x_i, y_i)\}_{i=1}^N$ be a small labeled dataset independently and identically drawn from the unknown real distribution $p_R(X, Y)$, where $x_i \in X$ and $y_i \in Y$, with N limited to only a few examples. The task defines a generation function G that maps D_{train} to a synthetic dataset:

$$D_{syn} = \mathcal{G}(D_{train}) = \{(\tilde{x}_j, \tilde{y}_j)\}_{j=1}^M.$$

where $(\tilde{x}_j, \tilde{y}_j)$ is a synthetic sample, and M is the number of synthetic samples, typically satisfying $M \gg N$. The task goal is to generate D_{syn} such that a model trained on it achieves strong predictive performance when evaluated on a held-out real test set $D_{\text{test}} \sim p_R$.

3.2 COMPONENT I: RULES-GUIDED GENERATION

We mitigate Distributional Drift with Rules-Guided Generation, which extracts association rules from limited reference data $D_{\rm train}$ as structural priors to guide generation. Unlike other modalities, tabular data often lack inherent structure and contain many irrelevant features, making it difficult for an LLM to capture feature dependencies from limited data (Fang et al., 2024b). To provide informative guidance, we define association rules as if—then formulas extracted from decision paths in tree-based models, capturing supervised dependencies between features and labels. Unlike classical association rule mining based on co-occurrence statistics, our rules reflect predictive patterns learned via training. We apply a two-stage LLM-driven procedure to extract reliable association rules.

- 1) Rules Extraction From Top-Performing Trees. We train a Random Forest (RF) on $D_{\rm train}$ and rank trees based on in-sample accuracy, selecting the top-k trees (e.g., k=3). Each selected tree provides a set of if-else decision rules (Kulkarni & Sinha, 2012; Azad et al., 2025; Khan et al., 2024), yielding diverse local patterns under low-data regimes.
- 2) Rule Generalization and Denoising. To construct reliable association rules, we apply:
- (a) Rule Generalization. We prune and consolidate decision paths across top-performing trees, retaining only their core (i.e., high-support, low-depth) branches that reflect stable feature—label dependencies. This produces *if*—then association rules that generalize beyond individual training

instances. Treated as conditional templates with the label as premise, these rules enable inverse reasoning and steer generation toward broader yet distributionally consistent samples.

- (b) Rule Denoising. To reduce variance introduced by symbolic extraction and decoding noise (He et al., 2024), we apply self-consistency techniques from reasoning tasks (Wang et al., 2023; Lewkowycz et al., 2022): we perform multiple generations with different seeds and retain only the most frequently occurring rules. This ensures the resulting rule set is stable and reliable.
- 3) Tabular Data Generation via Rules-Guided Prompt. Association rules obtained in Component I are converted into structured prompts encoding their if—then formulas. These prompts constrain the LLM's decoding space to enforce meaningful distributional patterns, yielding the synthetic dataset D_{syn} that offers distribution-consistent diversity.

3.3 COMPONENT II: DUAL-GRANULARITY FILTERING

We introduce **Dual-Granularity Filtering** (Component II) to mitigate *localized redundancy* in the D_{syn} . The procedure operates at two levels: *chunk-level filtering*, which prunes over-represented modes in high-density regions, and *instance-level filtering*, which filters unreliable samples in low-density regions. Both stages are guided by a *reference model* \mathcal{M} trained exclusively on D_{train} , ensuring that filtering remains consistent with the real data distribution.

3.3.1 PROXY-BASED DENSITY ESTIMATION

We quantify local redundancy by assigning each synthetic sample to its nearest real anchor and examining the concentration of synthetic mass around anchors. Concretely, let $\{s_j\}_{j=1}^N$ denote the $N = |\mathcal{D}_{\text{train}}|$ and let DCR be the mixed-type distance from (Borisov et al., 2023). This induces a discrete *proxy density* over anchors,

$$p_{j} = \frac{1}{|\mathcal{D}_{\text{syn}}|} \cdot \left| \left\{ x_{i} \in \mathcal{D}_{\text{syn}} \mid j = \arg \min_{1 \le k \le n} \text{DCR}(x_{i}, s_{k}) \right\} \right|, \quad j = 1, \dots, N.$$
 (1)

where $x_i \in \mathcal{D}_{\mathrm{syn}}, \, s_j \in \mathcal{D}_{\mathrm{train}}$ and p_j is the fraction of synthetic samples whose nearest neighbor in $\mathcal{D}_{\mathrm{train}}$ is s_j (thus $\sum_j p_j = 1$). We quantify overall concentration with the Gini coefficient G(p); larger G(p) indicates more localized redundancy. For targeted filtering, we split $\mathcal{D}_{\mathrm{syn}}$ into high-density ($\mathcal{D}_{\mathrm{high}}$) and low-density ($\mathcal{D}_{\mathrm{low}}$) sets by G(p).

3.3.2 Chunk-Level Filtering

Samples in \mathcal{D}_{high} are grouped into chunks of size S, each denoted as \mathcal{C}_r . Each chunk receives a score based on the average prediction correctness across its members under \mathcal{M} :

$$Score(\mathcal{C}_r) = \frac{1}{|\mathcal{C}_r|} \sum_{(x_i, y_i) \in \mathcal{C}_r} \frac{1}{T} \sum_{t=1}^T \mathbb{1} \left(\mathbb{P}_{\mathcal{M}_t}(y_i \mid x_i) > 0.5 \right), \tag{2}$$

Chunks are then ranked by score, and only a top fraction is retained, with the *pruning ratio* adaptively scaled by the *redundancy*:

$$ratio_{prune} = A \ln(G(p)) + B, \tag{3}$$

where A and B are fixed constants. This adaptive schedule increases pruning strength with greater localized redundancy, while maintaining flexibility when the proxy distribution is balanced.

3.3.3 Instance-Level Filtering

In contrast, \mathcal{D}_{low} reflects low-density regions with potentially informative but noisy samples. We filter samples using *confidence* and *uncertainty* scores derived from the same reference model \mathcal{M} . Thresholds are modulated by the G(p):

$$Conf_{thresh} = \mu_{conf} - G(p) \cdot \sigma_{conf},
Uncert_{thresh} = \mu_{uncert} + G(p) \cdot \sigma_{uncert}.$$
(4)

Only samples satisfying both $Conf(x) \ge Conf_{thresh}$ and $Uncert(x) \le Uncert_{thresh}$ are retained, ensuring filtering that adapts to dataset-specific sparsity.

3.3.4 Joint Tuning via Surprisal Minimization

The only tunable hyperparameter is chunk size S. We determine the optimal S^* by minimizing model surprisal on \mathcal{D} train:

$$S^* = \arg\min_{S} \left[-\frac{1}{|\mathcal{D}\text{train}|} \sum_{i} \log \mathbb{P}_{\mathcal{M}}(y_i \mid x_i) \right]. \tag{5}$$

This process produces a filtered dataset that enables the model to align with the target distribution.

4 EXPERIMENT

In this section, we evaluate *ReFine* on downstream tasks and analyze how its two components address the key challenges:

- 1. **Mitigating Distributional Drift**: *How do rules enhance data quality?* Section 4.3 shows that rule-guided generation achieves rule-compliance rates consistent with real data, thereby aligning synthetic samples more faithfully with the target distribution.
- 2. **Reducing Localized Redundancy**: *How does filtering balance the distribution?* Section 4.4 shows that a reliable redundancy metric together with dual-granularity filtering prevents overconcentration, leading to more useful synthetic datasets.

4.1 Experiment Settings

Datasets: To avoid potential *data contamination*—where strong performance on popular benchmarks such as *Adult*, *Heart*, and *Housing* may arise from LLM memorization rather than genuine generalization (Xu et al., 2024a; Ronval et al., 2025), we use tabmemcheck (Bordt et al., 2024), a recently proposed tool for detecting memorization in tabular data. Specifically, we apply two of its tests—the *Feature Names* and *Header Test*—to eight candidate datasets. Based on these tests, only *adult* and *heart* are classified as seen datasets; the remaining six show no evidence of memorization and are classified as unseen; full results are in Appendix B.

Baselines: (1) Non-LLM baselines: (i)CTGAN (Xu et al., 2019), a GAN-based model designed for tabular data generation; and (ii) TABSYN (Zhang et al., 2024b), a score-based generative model that achieves strong performance in sufficient-data settings. (2) LLM-based baselines: (i) GREAT (Borisov et al., 2023), which fine-tunes a pretrained GPT 2.5 for tabular data generation; (ii) EPIC (Kim et al., 2024), a prompting-based method that automates dataset construction through instruction-driven generation; and (iii) CLLM (Seedat et al., 2023), which enhances LLM-generated data via instance-level curation. (3) Constraint-based baseline: DRL (Stoian & Giunchiglia, 2025), a constraint-driven generator that synthesizes data by solving logical constraints.

Experimental Setup: We evaluate all models under low-data regimes with $N \in \{30, 60, 90, 120\}$. For each dataset and value of N, we randomly sample 10 training splits, while the remaining data serve as test sets. Evaluation follows the Machine Learning Efficiency (MLE) setting: each method generates 1,000 synthetic samples per split, on which we train an XGBoost (Chen & Guestrin, 2016). Performance is reported as average **F1 score** (classification) or \mathbb{R}^2 (regression) over all splits. Further implementation details are provided in Appendix C.

4.2 MAIN RESULTS

Table 1 reports the downstream performance of **ReFine** and representative baselines under four low-data regimes ($N \in \{30, 60, 90, 120\}$). **ReFine** achieves the strong results across all regimes, which improves R^2 by as much as **0.36** and F_1 by up to **7.5%** relative to the strongest prior method (CLLM) on unseen datasets, demonstrating robust generalization. The full **ReFine** framework (I+II) outperforms its individual components, achieving higher average ranks across benchmarks. This suggests the integration of both components contributes positively to overall effectiveness. As the amount of training data (N) increases, distribution-modeling capacity of non-LLM baselines, especially TabSyn, becomes more evident—they narrow the gap and occasionally approach LLM methods. This is consistent with their need for denser reference data. At the same time, DRL shows

Table 1: Main results on benchmark datasets. We report F_1 score for classification tasks and \mathbb{R}^2 for regression tasks. Unseen datasets (i.e., not included in the LLM's memory for LLM-based Generator) are highlighted in blue, and seen datasets are highlighted in red. The best result in each row is shown in **bold**, and the second-best result is underlined.

Original Da	Original Data			hods			LLM-Base	d Methods		
Datasets	Real	CTGAN	DRL	TABSYN	GREAT	EPIC	CLLM	I \ II	II \ I	I+II
Disease (N=30)	93.68	44.58	39.49	54.79	46.54	32.01	61.89	59.61	64.44	70.22
Game (N=30)	86.0	32.03	33.54	44.13	53.65	13.93	54.12	56.44	45.0	59.13
Apple (N=30)	86.60	47.21	51.33	32.32	58.91	56.10	58.18	57.80	58.66	59.43
GPA (N=30)	47.29	15.76	14.46	19.22	31.57	32.03	40.17	41.21	35.88	43.14
Student (N=30)	0.67	-0.91	-0.01	0.14	0.21	0.35	-0.11	0.37	0.08	0.38
Farm (N=30)	-0.04	-0.51	-0.07	-0.38	-1.03	-0.68	-0.29	-0.44	-0.23	-0.30
Adult (N=30)	76.92	50.47	46.59	62.83	67.45	61.94	73.11	73.04	73.15	73.91
Heart (N=30)	86.71	48.16	38.66	79.40	80.74	81.20	80.47	75.17	80.00	80.14
Disease (N=60)	93.68	30.03	38.25	65.34	52.80	44.05	66.86	75.65	62.58	72.41
Game (N=60)	86.0	31.83	30.10	61.10	54.76	13.16	<u>67.54</u>	61.39	59.81	70.87
Apple (N=60)	86.60	40.40	49.23	27.07	60.33	67.92	68.60	58.71	69.92	66.70
GPA (N=60)	47.29	18.11	15.64	28.48	35.60	32.19	33.17	44.34	21.57	44.58
Student (N=60)	0.67	-0.07	-0.04	-0.14	0.02	-0.63	-0.48	0.27	-0.15	0.34
Farm (N=60)	-0.04	-0.22	-0.16	-0.17	-0.16	-0.21	-0.30	-0.21	-0.12	<u>-0.15</u>
Adult (N=60)	76.92	47.55	49.02	63.58	69.70	62.78	73.94	73.28	72.87	71.48
Heart (N=60)	86.71	49.77	38.14	81.35	80.64	77.55	80.33	81.37	78.15	80.36
Disease (N=90)	93.68	47.18	39.04	69.90	65.32	30.03	74.04	69.74	65.47	76.45
Game (N=90)	86.0	33.67	27.81	65.93	61.17	13.16	59.97	59.44	41.77	62.17
Apple (N=90)	86.60	43.08	38.67	20.69	64.23	73.67	73.43	65.23	72.78	74.88
GPA (N=90)	47.29	14.91	12.24	34.51	36.19	16.28	41.79	47.21	31.56	48.89
Student (N=90)	0.67	-0.02	-0.03	-0.12	0.22	0.32	-0.69	0.34	-0.27	0.47
Farm (N=90)	-0.04	-0.34	-0.05	-0.15	-0.98	-0.16	-0.16	-0.15	-0.14	<u>-0.13</u>
Adult (N=90)	76.92	41.67	47.47	69.77	71.52	66.73	74.11	78.45	73.75	74.10
Heart (N=90)	86.71	42.04	39.03	81.43	81.05	77.65	81.23	82.44	78.74	80.00
Disease (N=120)	93.68	45.30	73.92	62.37	55.70	55.37	78.30	78.97	61.82	81.96
Game (N=120)	86.0	26.95	31.04	63.50	66.32	48.99	61.32	45.48	54.76	61.67
Apple (N=120)	86.60	39.83	34.29	80.54	60.91	<u>79.16</u>	70.66	70.62	68.53	74.96
GPA (N=120)	47.29	16.58	10.63	38.57	51.29	46.61	45.20	40.42	46.95	<u>47.77</u>
Student (N=120)	0.67	0	0.01	0.41	0.23	0.29	0.22	0.21	0.17	0.29
Farm (N=120)	-0.04	-0.05	-0.01	-0.22	-0.18	-0.30	-0.23	-0.10	-0.29	-0.07
Adult (N=120)	76.92	48.12	39.22	68.87	72.58	68.35	71.73	73.36	73.19	72.93
Heart (N=120)	86.71	46.29	38.29	83.80	84.17	71.33	75.24	84.97	75.02	75.93

nearly flat performance across N, since it enforces rules as hard constraints during generation. While this guarantees strict rule adherence, it severely limits the diversity of generated samples and prevents DRL from exploiting richer evidence when more real data become available. By contrast, performance varies significantly among LLM-based methods. EPIC and GREAT lack explicit generation guidance, which limits their ability to enforce meaningful structure in synthetic data. The relative advantage of $\it ReFine$ becomes smaller; in some cases (e.g., GPA), CLLM achieves a slightly higher $\it R^2$. This suggests that rules, while highly beneficial under data scarcity, may add mild constraints on numeric variability once sufficient reference data is available. Some LLM-based methods perform well on seen datasets but degrade notably on unseen datasets. This discrepancy aligns with the data contamination risk highlighted in our datasets selection—performance gains may in part reflect memorization from pretraining rather than true generalization. In contrast, $\it ReFine$ maintains stable performance across both, validating its capability for genuine generalization.

4.3 MITIGATING DISTRIBUTIONAL DRIFT

How do rules enhance data quality? To assess the extent to which generated samples adhere to the extracted association rules, we define the Rule Compliance Rate (RCR). Given a dataset S and a set of association rules \mathcal{R} , RCR is the percentage of samples in S that satisfy all rules in \mathcal{R} . We compare rule-guided and prompt-only generation across three representative datasets and four low-data regimes (N=30/60/90/120), evaluating both rule adherence (RCR) and downstream utility. Table 2 shows that under extreme low-data regimes (N=30), rule-guided generation significantly compensates for the lack of distributional evidence: compared to prompt-only generation, it produces samples with both higher rule adherence and stronger downstream utility (MLE). As N increases (N=60/120), the difference becomes most evident in distributional alignment. For in-

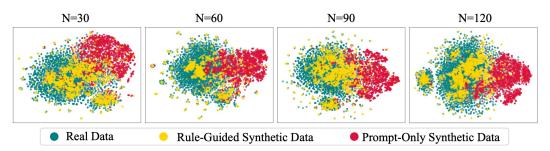


Figure 3: t-SNE visualization of *Real Data*, *Rule-Guided* synthetic data, and *Prompt-only* synthetic data on the *Disease* dataset across varying training sizes (N=30/60/90/120).

stance, in the Disease dataset, prompt-only data attains superficially higher RCR but deviates from the real distribution and yields much lower utility. By contrast, rule-guided data achieves RCR values closer to real data and substantially improves utility. This indicates that high but distorted

RCR is not reliable, and that the true advantage of rules lies in correcting such bias and aligning the synthetic distribution with the target one. The t-SNE visualizations (Figure 3) further support this finding: rule-guided samples form distributions that remain close to the real manifold, whereas prompt-only samples drift into unsupported regions. Overall, rules effectively mitigate distributional drift and enhance the utility of synthetic data across regimes.

We conduct further studies on *Component I* to assess both robustness and design

Table 2: Comparison of *rule-guided* data vs. *prompt-only* data across three datasets and varying training sizes (N=30/60/90/120). The better results are in bold.

	D_{test}	D_{train}	Rule Guid	led D_{syn}	Prompt-o	nly D_{syn}
	RCR(%)	RCR(%)	RCR(%)	MLE	RCR(%)	MLE
Disease (N=30)	21.4	33.3	16.3	59.61	12.2	48.70
GPA (N=30)	64.5	70.0	56.7	41.21	33.8	26.85
Student (N=30)	33.9	36.7	46.4	0.37	15.9	-0.11
Disease (N=60)	25.5	31.7	30.9	75.65	71.1	54.32
GPA (N=60)	51.5	61.7	49.1	44.34	22.5	15.44
Student (N=60)	80.2	13.3	16.1	0.27	15.6	-0.49
Disease (N=90)	31.3	44.4	41.6	69.74	66.8	68.85
GPA (N=90)	61.1	51.1	52.7	47.21	26.4	27.07
Student (N=90)	53.2	46.7	55.4	0.34	26.8	-0.70
Disease (N=120)	28.7	38.3	30.7	76.56	52.8	55.52
GPA (N=120)	52.6	53.3	40.8	40.42	31.6	44.13
Student (N=120)	57.7	54.2	67.7	0.21	17.2	-0.46

effectiveness (Appendix D). Results show that while performance remains stable across different top-k values, structured rule formats (i.e., "if-then") and self-consistency denoising yield clear improvements, validating the effectiveness of our design.

4.4 REDUCING LOCALIZED REDUNDANCY

How does filtering balance the distribution? To answer this, we compare different redundancy metrics and granularity settings and evaluate their impact on downstream utility (MLE) across datasets and data regimes. Table 3 reveals two key insights. (i) Not all redundancy metrics are equally effective. Although both Gini and entropy cab measure redundancy, Gini consistently yields higher MLE across datasets and training sizes (N), particularly under data-scarce conditions. This indicates that redundancy in prompt-generated data is primarily driven by a small number of highly repeated patterns, rather than widespread noise. Gini better captures this structure by emphasizing inequality, whereas entropy averages over all regions, giving undue weight to rare, possibly noisy instances and thereby underestimating redundancy. These results emphasize dominant concentrations (e.g. Gini) are better suited to detect and suppress it. (ii) Redundancy in LLM-generated data manifests at multiple scales—both within individual samples and across batch-level concentrations. As shown in Table 3, dual-granularity filtering consistently outperforms instance-only and chunk-only strategies across datasets and training sizes, confirming the necessity of addressing both levels. Instance-level filtering effectively removes noisy outliers but fails to capture global distributional imbalances.

Table 3: Evaluation of different redundancy metrics (Gini vs. Entropy) and filtering granularities (Instance-only, Chunk-only, Dual) within *Component II*.

	Redundancy Metric				Different Granularity			
	G	ini	Entropy		Instance	Chunk	Dual	
	Value	MLE	Value	MLE	Only	Only	Granularity	
Disease (n=30)	0.40	70.22	0.13	68.87	69.29	63.88	70.22	
GPA (n=30)	0.23	43.14	0.06	39.88	39.39	41.02	43.14	
Student (n=30)	0.58	0.38	0.31	0.36	0.39	0.35	0.38	
Disease (n=60)	0.68	72.41	0.30	72.03	71.77	68.81	72.41	
GPA (n=60)	0.23	44.58	0.03	38.53	43.25	45.04	44.58	
Student (n=60)	0.61	0.34	0.20	0.31	0.32	0.19	0.34	
Disease (n=90)	0.40	72.86	0.21	76.45	71.65	71.97	76.45	
GPA (n=90)	0.25	48.89	0.02	46.91	47.62	46.55	48.89	
Student (n=90)	0.58	0.47	0.20	0.43	0.45	0.38	0.47	
Disease (n=120)	0.35	81.96	0.14	78.97	80.48	69.74	81.96	
GPA (n=120)	0.29	47.77	0.08	45.12	42.99	40.77	47.77	
Student (n=120)	0.31	0.29	0.09	0.33	0.05	0.45	0.29	

In contrast, chunk-level filtering suppresses dominant high-density modes but may retain *localized redundancy*. By integrating both scales signals, dual-granularity filtering achieves more balanced coverage and higher downstream utility. This effect is visually illustrated in Figure 4: chunk-level filtering flattens overrepresented modes, while instance-level filtering enriches the long-tail regions—together restoring a more uniform and informative synthetic distribution.

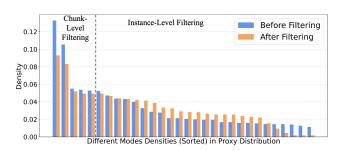


Figure 4: Proxy modes distribution before and after *dual-granularity filtering*.

We further validate the robustness of *Component II* in Appendix E. The Gini-driven pruning function peaks at intermediate retention, striking a balance between diversity and redundancy removal. Notably, Gini values stabilize with as few as 1,000 samples, confirming its suitability as a stable redundancy signal across generations.

5 CONCLUSION

 We present **ReFine**, a two-component framework that addresses fundamental limitations of prompt-based LLM tabular generation in low-data regimes. Our approach integrates symbolic *if-then* rules derived from interpretable models to enforce domain-specific feature distribution, while employing dual-granularity filtering to mitigate localized redundancy inherent in batch generation processes. The framework demonstrates consistent improvements across diverse benchmarks, achieving up to 0.36 absolute gain in R^2 and 7.5% relative improvement in F_1 over existing methods. Componentwise analysis reveals that rule-guided generation effectively captures dataset-specific dependencies often overlooked by pre-trained models, while dual-granularity filtering successfully rebalances synthetic distributions through coordinated chunk-level and instance-level selection strategies. Nevertheless, the current implementation of chunk-level retention utilizes empirically derived logarithmic scaling, potentially limiting generalizability under extreme distributional scenarios. Future research directions include exploring adaptive retention functions and establishing theoretical foundations for improved cross-domain robustness and generalization. We leave this as an avenue for future work.

REFERENCES

- Xavier Amatriain. Prompt design and engineering: Introduction and advanced methods. *arXiv* preprint arXiv:2401.14423, 2024.
- Mohammad Azad, Tasnemul Hasan Nehal, and Mikhail Moshkov. A novel ensemble learning method using majority based voting of multiple selective decision trees. *Computing*, 107(1): 42, 2025.
- Omar Benjelloun, Shiyu Chen, and Natasha Noy. Google dataset search by the numbers. In *International semantic web conference*, pp. 667–682. Springer, 2020.
 - Sahithi Bommareddy, Javed Ahmad Khan, Rohit Anand, et al. A review on healthcare data privacy and security. *Networking Technologies in Smart Healthcare*, pp. 165–187, 2022.
 - Sebastian Bordt, Harsha Nori, Vanessa Rodrigues, Besmira Nushi, and Rich Caruana. Elephants never forget: Memorization and learning of tabular data in large language models. In *Conference on Language Modeling (COLM)*, 2024.
 - Vadim Borisov, Kathrin Seßler, Tobias Leemann, Martin Pawelczyk, and Gjergji Kasneci. Language models are realistic tabular data generators. In *ICLR*, 2023.
 - Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pp. 785–794, 2016.
 - Xi Fang, Weijie Xu, Fiona Anting Tan, Jiani Zhang, Ziqing Hu, Yanjun Qi, Scott Nickleach, Diego Socolinsky, Srinivasan Sengamedu, and Christos Faloutsos. Large language models (llms) on tabular data: Prediction, generation, and understanding—a survey. *arXiv preprint arXiv:2402.17944*, 2024a.
 - Xi Fang, Weijie Xu, Fiona Anting Tan, Jiani Zhang, Ziqing Hu, Yanjun (Jane) Qi, Scott Nickleach, Diego Socolinsky, "SHS" Srinivasan Sengamedu, and Christos Faloutsos. Large language models (llms) on tabular data: Prediction, generation, and understanding a survey. *Transactions on Machine Learning Research*, 2024b.
 - Sakti P Ghosh. Statistical relational tables for statistical database management. *IEEE Transactions on Software Engineering*, (12):1106–1116, 2012.
 - Mingqian He, Yongliang Shen, Wenqi Zhang, Zeqi Tan, and Weiming Lu. Advancing process verification for large language models via tree-based preference learning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 2086–2099, 2024.
 - Mikel Hernandez, Gorka Epelde, Ane Alberdi, Rodrigo Cilla, and Debbie Rankin. Synthetic data generation for tabular health records: A systematic review. *Neurocomputing*, 493:28–45, 2022.
 - Zhanglong Ji, Zachary C Lipton, and Charles Elkan. Differential privacy and machine learning: a survey and review. *arXiv preprint arXiv:1412.7584*, 2014.
 - Azal Ahmad Khan, Omkar Chaudhari, and Rohitash Chandra. A review of ensemble learning and data augmentation models for class imbalanced problems: Combination, implementation and evaluation. *Expert Systems with Applications*, 244:122778, 2024.
- Jinhee Kim, Taesung Kim, and Jaegul Choo. Epic: Effective prompting for imbalanced-class data synthesis in tabular data classification via large language models. *Advances in Neural Information Processing Systems*, 37:31504–31542, 2024.
- Akim Kotelnikov, Dmitry Baranchuk, Ivan Rubachev, and Artem Babenko. Tabddpm: Modelling tabular data with diffusion models. In *International Conference on Machine Learning*, pp. 17564–17579. PMLR, 2023.
 - Boris Kovalerchuk and Evgenii Vityaev. *Data mining in finance: advances in relational and hybrid methods*, volume 547. Springer Science & Business Media, 2005.

Vrushali Y Kulkarni and Pradeep K Sinha. Pruning of random forest classifiers: A survey and future directions. In 2012 International Conference on Data Science & Engineering (ICDSE), pp. 64–68. IEEE, 2012.

- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in neural information processing systems*, 35:3843–3857, 2022.
- Zhenyu Li, Xiuxing Li, Zhichao Duan, Bowen Dong, Ning Liu, and Jianyong Wang. Toward a unified framework for unsupervised complex tabular reasoning. In 2023 IEEE 39th International Conference on Data Engineering (ICDE), pp. 1691–1704. IEEE, 2023.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM computing surveys*, 55(9):1–35, 2023.
- Milagros Miceli, Martin Schuessler, and Tianling Yang. Between subjectivity and imposition: Power dynamics in data annotation for computer vision. *Proceedings of the ACM on Human-Computer Interaction*, 4(CSCW2):1–25, 2020.
- Pradheepan Raghavan and Neamat El Gayar. Fraud detection using machine learning and deep learning. In 2019 international conference on computational intelligence and knowledge economy (ICCIKE), pp. 334–339. IEEE, 2019.
- Benoît Ronval, Pierre Dupont, and Siegfried Nijssen. Detection of large language model contamination with tabular data. In *International Symposium on Intelligent Data Analysis*, pp. 234–245. Springer, 2025.
- Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. A systematic survey of prompt engineering in large language models: Techniques and applications. *arXiv preprint arXiv:2402.07927*, 2024.
- Nabeel Seedat, Nicolas Huynh, Boris Van Breugel, and Mihaela Van Der Schaar. Curated llm: Synergy of llms and data curation for tabular augmentation in low-data regimes. *arXiv preprint arXiv:2312.12112*, 2023.
- Aditya Shankar, Hans Brouwer, Rihan Hai, and Lydia Chen. S i 1 o f use: Cross-silo synthetic data generation with latent tabular diffusion models. In 2024 IEEE 40th International Conference on Data Engineering (ICDE), pp. 110–123. IEEE, 2024.
- Mihaela C Stoian and Eleonora Giunchiglia. Beyond the convexity assumption: Realistic tabular data generation under quantifier-free real linear constraints. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Mihaela C Stoian, Salijona Dyrmishi, Maxime Cordy, Thomas Lukasiewicz, and Eleonora Giunchiglia. How realistic is your synthetic data? constraining deep generative models for tabular data. In *The Twelfth International Conference on Learning Representations*, 2024a.
- Mihaela Cătălina Stoian, Salijona Dyrmishi, Maxime Cordy, Thomas Lukasiewicz, and Eleonora Giunchiglia. How realistic is your synthetic data? constraining deep generative models for tabular data. *arXiv preprint arXiv:2402.04823*, 2024b.
- Paul Voigt and Axel Von dem Bussche. The eu general data protection regulation (gdpr). *A practical guide, 1st ed., Cham: Springer International Publishing*, 10(3152676):10–5555, 2017.
- Alex X Wang, Stefanka S Chukova, Colin R Simpson, and Binh P Nguyen. Challenges and opportunities of generative models on tabular data. *Applied Soft Computing*, pp. 112223, 2024a.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2023.

Yuxin Wang, Duanyu Feng, Yongfu Dai, Zhengyu Chen, Jimin Huang, Sophia Ananiadou, Qianqian Xie, and Hao Wang. Harmonic: Harnessing llms for tabular data synthesis and privacy protection. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024b.

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837, 2022.
- Cheng Xu, Shuhao Guan, Derek Greene, M Kechadi, et al. Benchmark data contamination of large language models: A survey. *arXiv preprint arXiv:2406.04244*, 2024a.
- Jundong Xu, Hao Fei, Liangming Pan, Qian Liu, Mong-Li Lee, and Wynne Hsu. Faithful logical reasoning via symbolic chain-of-thought. In *The 62nd Annual Meeting of the Association for Computational Linguistics*, 2024b. URL https://arxiv.org/abs/2405.18357.
- Lei Xu, Maria Skoularidou, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. Modeling tabular data using conditional gan. *Advances in neural information processing systems*, 32, 2019.
- Xu Yang, Meihui Zhang, Ju Fan, Zeyu Luo, and Yuxin Yang. A multi-task learning framework for reading comprehension of scientific tabular data. In 2024 IEEE 40th International Conference on Data Engineering (ICDE), pp. 3710–3724. IEEE, 2024.
- Rodolfo Zevallos, Mireia Farrús, and Núria Bel. Frequency balanced datasets lead to better language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 7859–7872, 2023.
- Baoquan Zhang, Chuyao Luo, Demin Yu, Xutao Li, Huiwei Lin, Yunming Ye, and Bowen Zhang. Metadiff: Meta-learning with conditional diffusion for few-shot learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pp. 16687–16695, 2024a.
- Hengrui Zhang, Jiani Zhang, Balasubramaniam Srinivasan, Zhengyuan Shen, Xiao Qin, Christos Faloutsos, Huzefa Rangwala, and George Karypis. Mixed-type tabular data synthesis with score-based diffusion in latent space. In *The twelfth International Conference on Learning Representations*, 2024b.
- Zilong Zhao, Aditya Kunar, Robert Birke, and Lydia Y Chen. Ctab-gan: Effective table data synthesizing. In *Asian conference on machine learning*, pp. 97–112. PMLR, 2021.
- Tianyang Zhong, Zhenyuan Yang, Zhengliang Liu, Ruidong Zhang, Yiheng Liu, Haiyang Sun, Yi Pan, Yiwei Li, Yifan Zhou, Hanqi Jiang, et al. Opportunities and challenges of large language models for low-resource languages in humanities research. *arXiv preprint arXiv:2412.04497*, 2024.

A USE OF LARGE LANGUAGE MODELS (LLMS)

In preparing this manuscript, we employed a large language model (LLM) to assist with polishing the writing. Specifically, the LLM was used to improve clarity, grammar, and readability of the text. All substantive ideas, analyses, and conclusions are solely those of the authors.

B DATASET CONTAMINATION TEST

B.1 TABMEMCHECK

To rigorously evaluate the generalization capabilities of prompt-based tabular data generation, we assess potential training set contamination in popular benchmark datasets. Following concerns raised in recent studies (Xu et al., 2024a; Ronval et al., 2025), we employ tabmemcheck (Bordt et al., 2024), a diagnostic tool for detecting dataset memorization in LLMs. Specifically, we use two tests:

- Feature Names Test: Given a few sample rows, the LLM is prompted to infer corresponding Column Names.
- Header Test: The LLM is prompted to reconstruct the header (column names) and the first *few data rows* of the dataset in CSV format.

B.2 EXAMPLE: HEART DATASET (GPT-40)

As shown in Figure 5, GPT-40 accurately recovers the feature names and generates realistic rows for the heart dataset, despite limited prompting. Fields such as ST_Slope and RestingECG are generated verbatim, indicating memorization.

This demonstrates that **strong LLMs may replicate datasets**, compromising fair evaluation.

```
Dataset: heart.csv
Feature Names Test: Age, Sex, ChestPainType, RestingBP, Cholesterol, FastingBS, RestingECG, MaxHR, ExerciseAngina, Oldpeak, ST_Slope, HeartDisease Legend: Prompt Correct Incorrect Missing
Dataset: heart.csv
Header Test: Age, Sex, ChestPainType, RestingBP, Cholesterol, FastingBS, RestingECG, MaxHR, ExerciseAngina, Oldpeak, ST_Slope, HeartDisease
40, M, ATA, 140, 289, 0, Normal, 172, N, 8, Up, 0
40, F, NAP, 160, 180, 0, Normal, 172, N, 8, Up, 0
40, F, NAP, 160, 180, 0, Normal, 150, N, 1, Flat, 1
54, M, ATA, 150, 195, 0, Normal, 108, Y, 1.5, Flat, 1
54, M, ATA, 150, 195, 0, Normal, 170, N, 0, Up, 0
45, F, ATA, 130, 237, 0, Normal, 170, N, 0, Up, 0
45, F, ATA, 130, 237, 0, Normal, 170, N, 0, Up, 0
50, F, NAP, 120, 219, 0, Normal, 170, N, 0, Up, 0
50, F, NAP, 1220, 219, 0, Normal, 155, N, 1.0, FlUp, 0
```

Figure 5: **GPT-40 reproduces the heart dataset.** Field names and values are copied precisely, indicating contamination.

B.3 EXAMPLE: ADULT DATASET (QWEN2.5-32B)

In contrast, Qwen2.5-32b fails to regenerate exact feature names but instead **identifies the dataset by name and describes its contents and use cases** (as shown in Figure 6). This demonstrates semantic-level exposure rather than strict memorization.

Such dataset-wide familiarity still violates the assumption of data independence and motivates filtering benchmarks during evaluation.

B.4 IMPLICATIONS AND DATASET FILTERING

These examples reveal that both closed-source (GPT-4o) and open-source (Qwen2.5) models can memorize benchmarks—either by reproducing format or semantically recognizing dataset identity. Table 4 summarizes model accuracy across these tests. High scores on these tasks suggest that the LLM has memorized structural or content-level information about the dataset, potentially inflating downstream performance.

```
Info: Removed the few-inord dataset adult-train.civ because it is similar to the dataset being tested.

Dataset: adult.civ

Feature Names Test: app, workclass, finley, It seems like the list for the "adult" dataset as cut off. The Adult dataset is commonly used in machine learning tasks and typically includes information, related to income prescribed as a new complete list of feature names for the Adult dataset, which usually consists of 15 features (excluding the target variable 'income'):

Info: Removed the few-shot dataset adult train.cav because it is similar to the dataset being tested.

Dataset: adult.civ

Reader Test: app, workclass, finleyt, education, education-may, marital-storus, occupation, relationship, race, sex, capital-gain, capital-loss, hours-per-week, country, salary

Sy, State-gov, 1754, Beachers, 23, Never-maried, Adm-clerical, Not-in-family, White, Rel, 2,174, Opt, United-States, <-6000

Sy, Private, 133069, Master, 34, Maride-cit-v-spouse, Exec-managerial, Notament, No
```

Figure 6: **Qwen2.5-32b recognizes the adult dataset.** While exact headers are missing, the model describes the dataset and its usage.

Table 4: **Feature Names and Header Test** Results. **Feature Names Test**: The LLM infers column names from sample rows. **Header Test**: The LLM completes the header and first few rows of the CSV. Percentages indicate accuracy for each dataset and LLM.

	GPT-40-0806		GPT-3.5-turbo-1106		Qwen2.5-32b-Instruct		Qwen2.5-14b-Instruct	
	Feature Names	Header	Feature Names	Header	Feature Names	Header	Feature Names	Header
Other Datasets	0%	0%	0%	0%	0%	0%	0%	0%
Adult	86.67%	75.0%	86.67%	100.0%	0%*	25.0%	0%*	80.0%
Heart	100.0%	55.56%	25%	22.22%	0%	11.11%	0%	12.50%

^{*}Although no correct column names were produced, the LLM identified "the Adult dataset from the UCI Machine Learning Repository"

C EXPERIMENTAL SETUP DETAILS

Model Configuration. We use GPT-3.5-Turbo-1106 as the backend model for both association rule extraction and data generation. Each baseline is configured to generate roughly 2,000 synthetic samples per dataset.

Rule-Guided Generation. We set k=3 for selecting top-performing trees in Random Forests, and apply self-consistency by aggregating 5 independently sampled generations for rule extraction.

Dual-Granularity Curation. We use XGBoost as the reference model \mathcal{M} for evaluating chunk-level informativeness, and search chunk sizes $S \in \{20, 25, \dots, 60\}$ during tuning.

Evaluation. For each run, we train XGBoost on 1,000 synthetic samples under 10 random seeds and evaluate on the corresponding real dataset. Classification tasks are measured with F1 score; regression tasks with R^2 .

D FURTHER STUDY ON COMPONENT I

We conduct a series of supplementary analyses to unpack this question. In Section D.1, we test the stability of extracted rules by varying the number of top-k trees and find that rule quality remains robust across settings. In Section D.2, we compare rule formats and show that explicit if—then clauses provide clearer guidance than natural-language paraphrases. In Section D.3, we examine rule denoising strategies and demonstrate that self-consistency aggregation yields more reliable rules than single-pass or CoT prompting. In Section D.4, we validate transferability across different LLM backbones, confirming that rule guidance consistently improves synthetic data quality regardless of model scale. In Section D.5, a case study illustrates how noisy tree paths are distilled into compact and interpretable rules through merging and denoising, highlighting both robustness and interpretability. Together, these studies confirm that rules enhance data quality by providing stable, precise, and broadly applicable generation guidance.

D.1 EFFECT OF TOP-K TREES IN RULE EXTRACTION.

To examine the robustness of the extracted rules, we vary the number of top-k trees used for rule extraction (k=1, 2, 3, 5, 10). Across all settings, the resulting rule-compliance rates (RCR) remain stable, suggesting that the induced rules capture consistent feature-label dependencies that are not sensitive to the choice of k. While k=1 or 2 already achieve similar RCR, we adopt k=3 as a default since it provides broader rule coverage than very small k while avoiding the additional overhead

Table 5: Rule Compliance Rate (RCR) Results for different k values (D_{test} and D_{train}).

	k =	= 1	k =	= 2	k =	= 3	k =	= 5	k =	= 10
	D_{test}	D_{train}								
Disease (N=30)	63.44	83.33	36.38	53.33	29.79	40.00	44.61	50.00	72.99	63.33
GPA (N=30)	62.45	70.00	62.45	70.00	58.17	66.67	62.57	70.00	63.76	60.00
Student (N=30)	46.12	56.67	28.83	50.00	26.53	50.00	50.66	73.33	36.45	30.00
Disease (N=60)	60.67	76.67	34.92	35.00	43.04	53.33	26.41	28.33	50.93	56.67
GPA (N=60)	41.90	36.67	42.02	56.67	58.02	70.00	49.01	65.00	49.70	58.33
Student (N=60)	17.72	13.33	54.77	48.33	34.46	36.67	43.07	35.00	57.95	61.67
Disease (N=90)	50.56	70.00	34.61	46.67	69.18	80.00	13.67	17.78	61.10	73.33
GPA (N=90)	50.61	57.78	50.74	57.78	74.85	70.00	66.12	62.22	54.00	48.89
Student (N=90)	31.58	26.67	26.67	22.22	31.13	27.78	55.65	50.00	49.38	51.11
Disease (N=120)	51.21	63.33	27.16	41.67	30.26	39.17	28.63	38.33	28.78	38.33
GPA (N=120)	40.98	44.17	48.50	48.33	60.74	41.67	54.52	39.11	64.73	54.62
Student (N=120)	29.57	31.67	65.82	63.33	42.58	43.33	32.36	34.17	28.97	34.17

Natraul Language Rule Form

- There is a noticeable trend where patients diagnosed with Alzheimer's often have lower MMSE scores. This trend is particularly apparent among older patients, suggesting a strong association between advanced age, decreased cognitive function, and an Alzheimer's diagnosis. The lower functional assessment scores in these individuals further support this link, indicating that cognitive decline is a key factor in distinguishing those with Alzheimer's.

"If-Then" Rule Form

- If Diagnosis = 1:
- Then MMSE ≤ 16.85 and Functional Assessment < 6.90
- If Diagnosis = 0:
- Then MMSE > 18.85 and Functional Assessment > 6.90

Figure 7: Illustrative "if-then" Form and its Natural-Language paraphrase derived from the *Disease* dataset (N=30).

of larger k values. This balance ensures that the extracted rules are both robust and efficient for downstream prompting.

D.2 DIFFERENT RULE FORMAT COMPARISON.

To gauge whether the *explicit* "if—then" representation is essential to the success of Rule-Guided Generation, we recast every Random-Forest path into two formats: (i) its original "if—then" clause and (ii) a concise natural-language paraphrase, which automatically produced by prompting the LLM to restate each path in natural language. These two formats reflect two distinct rule-extraction schemes. A side-by-side example of the two rule forms is shown in Figure 7. As shown in Table 6, *if—then* rules outperform both natural-language rules and the No-Rule baseline, demonstrating the benefit of symbolic structure. While both rule-based approaches surpass the baseline, the symbolic form delivers more reliable performance. This advantage stems from two key factors: (i) Random Forests extract concise and faithful feature—label dependencies even in low-data settings, and (ii) the *if—then* format retains explicit numeric boundaries and dependent constraints, unlike natural language, which tends to weaken precision (Xu et al., 2024b). By guiding generation through explicit symbolic rules, the LLM is directed toward semantically coherent subspaces, resulting in higher-quality samples.

D.3 ROBUSTNESS OF RULE DENOISING STRATEGY.

We evaluate how different Rule Denoising strategies affect data quality. We compare our proposed *Self-Consistency Rule Denoising* (described in Section 3.2) against two alternatives: (i) **Single-Pass**, which denoise rules in one step without verification, and (ii) **Chain-of-Thought** (**CoT**) prompting, which guides the LLM to reason step-by-step during rule denoising (Wei et al., 2022). Results in Table 7 reveal that both **Single-Pass** and **CoT** aggregation yield less consistent performance

Table 6: Performance under different rule representations (n=30).

	No Rule	Natural Language	"if-then" Form
Disease	49.13 ± 1.2	57.54 ± 1.7	59.61 ± 1.5
GPA	24.80 ± 4.0	$37.78 \pm .60$	$\textbf{41.21} \pm \textbf{.97}$
Student	-0.11 ± 12	-0.79 ± 53	$0.37 \pm .02$

Table 7: Performance under different *rule de*noising strategies (n = 30).

	Single-Pass	СОТ	Self-Consistency
Disease	54.26 ± 3.2	$58.46 \pm .98$	$\textbf{59.61} \pm \textbf{1.5}$
GPA	24.80 ± 4.0	$38.48 \pm .82$	$\textbf{41.21} \pm \textbf{.97}$
Student	$0.35 \pm .06$	$0.35 \pm .02$	$\textbf{0.37} \pm \textbf{.02}$

Table 8: Main results on benchmark datasets (n=30). We report **F1 score** for classification tasks and \mathbf{R}^2 for regression tasks. The best result in each row is shown in **bold**.

Origi	inal Data	Rules Generator	Data Generator					
Datasets	Real Data		Qwen2.5-14b-Instruct	Qwen2.5-32b-Instruct	GPT-4o-0806			
Disease	93.68 ± 1.4	No Rules Qwen2.5-32b-Instruct GPT-40-0806			$ \begin{vmatrix} 62.26 \pm 3.2 \\ 67.38 \pm .85 \\ 64.04 \pm 2.7 \end{vmatrix} $			
GPA	$47.29 \pm .90$	No Rules Qwen2.5-32b-Instruct GPT-40-0806	$ \begin{array}{c} 20.67 \pm .87 \\ 40.62 \pm 1.4 \\ 42.25 \pm 1.4 \end{array} $	$ \begin{vmatrix} 30.23 \pm 1.4 \\ 39.13 \pm 3.4 \\ 31.99 \pm 2.9 \end{vmatrix} $	$ \begin{vmatrix} 27.58 \pm 2.3 \\ \textbf{47.07} \pm \textbf{0.5} \\ 44.34 \pm 0.6 \end{vmatrix} $			
Student	$0.67 \pm .07$	No Rules Qwen2.5-32b-Instruct GPT-40-0806	$ \begin{vmatrix} 0.25 \pm .04 \\ 0.30 \pm .03 \\ 0.18 \pm .05 \end{vmatrix} $	$\begin{array}{c} 0.32 \pm .03 \\ \textbf{0.33} \pm .03 \\ -0.38 \pm .16 \end{array}$	$ \begin{vmatrix} -1.29 \pm .27 \\ 0.25 \pm .01 \\ -0.01 \pm .14 \end{vmatrix}$			

across datasets. Their reliance on one-shot reasoning makes them vulnerable to local inconsistencies and stochastic behavior in LLM outputs. In contrast, self-consistency aggregation enforces cross-run agreement and filters out unstable logic fragments, leading to more robust rule sets and better downstream fidelity.

D.4 RULE-BASED METHODS GENERALIZE ACROSS LLMS

To evaluate the robustness of *Rule-Guided Generation* in different LLM backbones, we compare performance using the MLE under varying combinations of rules generators and data generators. The result as shown in Table 8, across all tested LLMs which ranging from mid-sized (Qwen2.5-14b) to larger models (Qwen2.5-32b and GPT-4o), the introduction of association rules consistently improves the quality of synthetic data. Notably, even when rules are extracted by weaker models (e.g., Qwen2.5-32b), stronger models like GPT-4o still benefit from them—highlighting that rule-based guidance contributes independently of LLMs capacity. This demonstrates the generality and transferability of our design.

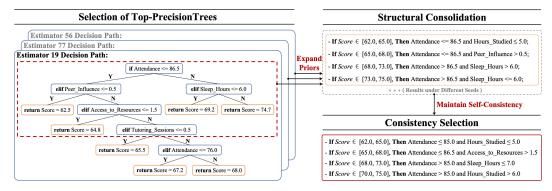


Figure 8: Case study for Component I on the *Student* dataset. Noisy tree paths are denoised into a self-consistent symbolic formulas set and later guides data generation.

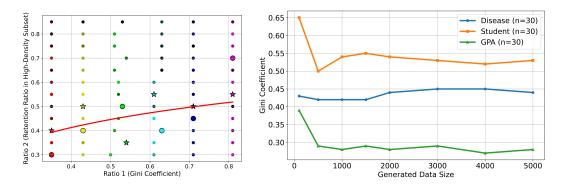


Figure 9: (Left) Scatter plot of G(p) (Gini coefficient) versus $Ratio_{prune}$. Point color indicates downstream model performance after filtering, with lighter colors representing higher performance and darker colors indicating lower performance. In each group, the star (\star) marks the best-performing point and the circle (\circ) marks the second-best. (Right) Gini coefficient under different synthetic data sizes.

D.5 ILLUSTRATIVE EXAMPLE OF RULE GENERALIZATION & DENOISING.

To illustrate the practical functioning of **Component I**, we present a case study demonstrating how symbolic rules are distilled into interpretable if—then forms through rule merging and aggregation. Intermediate outcomes are visualized in Figure 8. In this example, decision trees from the trained random forest exhibit conflicting paths—for example, assigning different $Exam_Score$ values to overlapping input regions such as $Attendance \le 86.5$ and $Attendance \le 76.0$. The merging phase addresses these inconsistencies by consolidating noisy rule fragments into coherent patterns, such as those involving $Sleep_Hours$. To improve robustness, the training and generalization process is repeated under multiple random seeds. The denoising step then retains only those patterns that appear consistently across runs, thereby filtering out unstable conditions and enforcing self-consistency. This case highlights two key strengths of Component I: (1) it distills noisy and fragmented tree logic into compact rules that capture meaningful relationships in low-data regimes; (2) it produces one concise association rule per segment, enhancing both interpretability and generation quality.

E FURTHER STUDY ON COMPONENT II

To further examine *how filtering balances the distribution*, we study the behavior of **Component II** under different control settings. Specifically, we analyze (i) the effect of the log-scaled retention function on pruning schedules, and (ii) the stability of the Gini coefficient under varying generation sizes. These studies show that Gini-based filtering prunes aggressively only when redundancy is severe while maintaining diversity in balanced regimes, and that Gini values quickly stabilize as generation size grows. Together, these results confirm that our filtering strategy provides a robust and scale-invariant mechanism for mitigating localized redundancy and restoring distributional balance. In following experiments, we fix the backbone settings (GPT-40-0806 as Rule Generator and GPT-3.5-turbo-1106 as Data Generator) and evaluate results via MLE score. We report **F1 score** for *Disease* and *GPA*, and **R**² for *Student*.

E.1 IMPACT ON THE LOG-SCALED RETENTION FUNCTION.

Using the log-scaled mapping in (3), we vary G(p) across its empirical range and record the resulting ratio_{prune} as well as downstream F1. Figure 9 (Left) shows that performance peaks at intermediate retention, validating that the Gini-driven schedule prunes aggressively only when redundancy is severe, while preserving diversity in more balanced settings. We observe that the Gini values across datasets tend to lie in a relatively narrow and stable range, which contributes to the robustness of the fitted retention schedule. As a result, the coefficients A=0.15 and B=0.55, derived from cross-dataset regression, generalize well without tuning.

E.2 STABILITY UNDER VARYING GENERATION SIZES

We next test whether Gini is sensitive to the number of generated samples, since an unstable control signal would undermine filtering. Figure 9 (Right) shows that the Gini coefficient stabilizes after roughly 1,000 samples across tasks, with only minor fluctuations thereafter. This indicates that dominant distributional modes emerge early in generation, and that Gini provides a consistent, scale-invariant signal for downstream filtering decision.