MALLM-GAN: MULTI-AGENT LARGE LANGUAGE MODEL AS GENERATIVE ADVERSARIAL NETWORK FOR SYNTHESIZING TABULAR DATA

Anonymous authors

Paper under double-blind review

ABSTRACT

In the era of big data, access to abundant data is crucial for driving research forward. However, such data is often inaccessible due to privacy concerns or high costs, particularly in healthcare domain. Generating synthetic (tabular) data can address this, but existing models typically require substantial amounts of data to train effectively, contradicting our objective to solve data scarcity. To address this challenge, we propose a novel framework to generate synthetic tabular data, powered by large language models (LLMs) that emulates the architecture of a Generative Adversarial Network (GAN). By incorporating data generation process as contextual information and utilizing LLM as the optimizer, our approach significantly enhance the quality of synthetic data generation in common scenarios with small sample sizes. Our experimental results on public and private datasets demonstrate that our model outperforms several state-of-art models regarding generating higher quality synthetic data for downstream tasks while keeping privacy of the real data.

025 026 027

006

008 009 010

011

013

014

015

016

017

018

019

021

1 INTRODUCTION

Tabular data is the most common data format in high-stakes sectors like healthcare. There are many fundamental problems in dealing with tabular data, such as data scarcity, missing values, and irregularity. Among them, the data scarcity problem has been the main roadblock. Many datasets in healthcare, such as clinical trial data, have small data sizes due to data collection costs and privacy risks, and consequently, these data cannot afford modern machine learning (e.g., deep learning), which generally has thousands of parameters, at minimum.

Recent advancements in generative models, particularly in text and image,(7; 28) have shown the benefits of technology for generating synthetic data that resembles real data. Despite this potential, generating tabular data has not been fully tapped into; it has evolved through traditional statistical approaches, like Bayesian networks (38), to deep learning techniques, including autoencoders and Generative Adversarial Networks (GANs) (42). However, these methods require large amounts of data for their training, contradicting our objective of solving data scarcity. Also, the sparsity and heterogeneity in tabular data make GAN or other deep learning a suboptimal choice, as evidenced by the fact that the tree-based method (e.g., XGBoost) works better than deep learning model (11).

Recently, advancements in large language models (LLMs) have also enabled researchers to use their
general intelligence to synthesize tabular data.(6; 14) The premise is that prior knowledge encoded in
the parameters of LLMs can provide contextual knowledge for coherent semantics that is required
to learn the underlying data generation process. Several studies transformed tabular data to natural
language via serialization, and used pre-trained LLMs to generate text containing the synthetic tabular
data (6; 14; 21). While fine-tuning LLMs has led to the creation of more nuanced synthetic data,
this process again requires a significant sample size, contradicting the objective of addressing data
scarcity.

In contrast, in-context learning presents a promising alternative by allowing LLMs to customize
 without compromising their general reasoning abilities (24). Particularly, few-shot learning in in context learning is to provide a few "examples" of data to let LLM to learn the patterns and mimic the
 examples (17). Our study aims to utilize this few-shot capability for synthetic tabular data generation.
 Currently, tabular generative model by in-context learning, such as (31) faces a critical shortcoming;

it only accommodate too few "examples" (not an entire dataset available), thus discarding remaining
data that cannot manage to fit into limited context length (input token). For example, in our real-world
trial data ATACH2 with 37 variables, just ten samples consume 4,232 input tokens, whereas the
common input token size of LLMs, such as GPT3, is 2,048. This failure to utilize all available scarce
data can make LLMs perceive the underlying data-generating process merely based on "educated
guess" with its prior knowledge instead of the data itself. This, consequently, causes distribution
discrepancy between real data and synthetic data.

Therefore, we aim to bridge this critical gap. Our key idea is to make the *data generation process* explicit; the objective of our in-context learning is to generate a better data generation process, as
well as to generate individual data instances. Here, the data generation process is a prompt text that consists of the context of data and any simple "model" that describes the relationship between data variables. We chose to use a Bayesian network or causal structure due to its interpretability and simplicity.

However, another challenge is to identify the ground-truth data generation process. Motivated by
GAN's adversarial training, we optimize the data generation process ("generator") in adversarial
training with "discriminator" (Table 1). The discriminator's role is to discriminate real data from the
generated data, and we use the accuracy of the discriminator as a "loss function" to be minimized to
optimize the generator. Unlike GAN, our generator is a text format, which doesn't have derivatives.
We address it by Optimization by Prompting, which leverages an independent LLM as an optimizer
(43). After optimizing the data generation process, the LLM as generator uses it to finally generate
synthetic data.

Table 1: Comparison of Generative Adversarial Network (GAN) and Our Model

	GAN	Our Model
Generator	Neural network	Frozen LLM and prompt
Discriminator	Neural network	Tabular data classifier
Optimizer	Gradient descent	Frozen LLM and prompt

Contribution of this paper can be summarized as below

- *Novelty*: We propose a novel concept of optimizing data generation process using in-context learning of LLM. This leverages both data-driven supervised model (discriminator) and knowledge-driven in-context learning of LLM (generator, optimizer).
- *Few-shot synthetic data generation*: Our model works when there is too little data to train a parametric model. It mitigates data scarcity problem in many small-size tabular data in healthcare.
- *Conditional sampling*: Our generator is based on LLM, which enable conditional sampling seamlessly by prompting.
- *Explainability*: Our LLM-based generator explicitly reveals data generation process and its reasoning, which are explainable by design. This enables transparency of our model and facilitates human feedback, such as refining the knowledge.

2 RELATED WORK

075

081 082 083

084

085

090

092

093

095

096

098

This section provides a brief overview of the most relevant prior work in the field of tabular data generation, LLM-supported synthetic data generation, and the different roles that LLM plays in real-world applications.

Synthetic tabular data generation. Many studies have been proposed to generate high-quality
 tabular data for privacy-preserving data sharing and augment training data size for machine learning
 models. Traditional method including Bayesian network (44; 38), approximation Bayesian computa tion (5), and SMOTE (8). Particularly, a Bayesian network can represent a pairwise causal structure
 via directed acyclic graph (DAG) (26), in which a directed edge between variable A and B exists if A
 causes B. Causal structure is a compact representation of underlying variable relationship, but does
 not fully capture the nonlinear numerical relationship between the different types of variables. Deep



Figure 1: Overview. In each optimization step, the LLM as Optimizer generates a data generation process θ in (Generator prompt) based on the pairs of previous θ and its score in **Optimizer prompt**. Then the LLM as Generator use the current θ and a few examples to generate data. We evaluate the θ using the accuracy (score) of Discriminator. The more the data generation process θ is optimized, the lower the discriminator's accuracy. This adversarial optimization finishes when data generation process is no more improved.

117

118

119

120

121

122

generative models have also been widely utilized. This includes variational autoencoders such as TVAE (42) and others (2; 40); GAN such as CTGAN (42), and diffusion model such as TabDDPM (19). The limitation of these models is that they require a sufficient sample size to train the generative model, which contradicts our aim to solve data scarcity.

Synthetic data generation using LLMs. Recently, large language models emerged and have 130 demonstrated their powerful performance in generating natural text. It has also shown a great 131 potential in tabular data (10), such as for predicting (14; 12; 45; 21) and generating tabular data 132 (6; 32; 47; 12). Among them, GReaT (6), the first model in this line, transformed the tabular data 133 into text and fine-tuning LLM (GPT-2), including a feature order permutation step for added realism. 134 However, these prior methods require fine-tuning LLMs, which requires large data size and computing 135 resources for fine-tuning. This limitation of existing models motivates us to develop *few-shot* tabular 136 data generative model. 137

Multiple roles of LLM in real-world applications. In addition to typical natural text generation tasks it was originally trained for, LLMs have been utilized in various tasks. LLM has been used as an optimizer for the data type that we can not calculate derivatives directly (43). This LLM as optimizer was used to optimize prompt (43) and optimize heuristic algorithms written in codes (29). LLMs also have been utilized for communicating with other LLMs (i.e., multi-agent LLM); multiple LLMs play different roles to accomplish a task collaboratively, such as coding,(13) question answering,(41) and online decision making.(35; 16)

LLM and causal structure discovery Causal structure discovery involves either data-driven or 145 expert knowledge-driven approaches. Data-driven approaches discover causal structure from data 146 by conditional independence tests (33; 34), score-based heuristics (36), or relaxing the discrete 147 constraints into continuous optimization (48; 46). Despite these advances, identifying the ground-148 truth causal structure from data remains a significant challenge, particularly in complicated domains 149 like healthcare or when data is scarce. Expert-driven approaches can be an alternative option, but 150 these methods are time-consuming and require significant expert involvement. Recently, many studies 151 (18; 23) suggest that LLMs, which encode prior knowledge in their parameters, can support causal 152 discovery by complementing expert knowledge. In this paper, we leverage multiple LLMs with 153 different roles to mimic adversarial training in GAN and use the heuristic causal structure discovery 154 to guide the data generation process.

155 156

3 Methodology

157 158 3.1 PROBLEM FORMULATION

Given a small labeled tabular dataset with n instances and d features, denoted as $D_{real} = (\mathbf{x}, y)$ where x represents a d-dimensional vector of features and y indicates label. The features are described by natural-language strings like "age" or "gender". For synthetic data generation, we train a generator on a training subset D_{train} of D_{real} , generating synthetic dataset D_{syn} .

162 3.2 MULTI-AGENT LLM AS GAN

163 Overview. We propose to develop a multi-agent LLM as GAN (MALLM-GAN) that generates tabular 164 data by mimicing adversarial optimization (Fig. 1). The objective is to optimize the data generation 165 process θ , which is a natural language description of i) the problem description and ii) the simple data 166 generation process or causal structures representing relationships between variables. In MALLM-GAN, 167 for each iteration i, an LLM agent **Generator** generates data D_{syn} with θ_i and a batch in D_{train} ; 168 a supervised model **Discriminator** is accordingly optimized using $[D_{train}, D_{syn}]$ and evaluates θ_i 169 using D_{test} ; and another LLM agent **Optimizer** improves θ_i to decrease the discriminator's accuracy 170 (Algorithm 1). We repeat the iterations until the discriminator's accuracy converges or the iteration reaches the maximum epoch. 171

172 3.2.1 GENERATOR

185

199

173 **Data generation process.** The data generation process θ is described in natural language and prompts 174 the generator LLM to create synthetic data. It includes: i) context of data collection, ii) data schema, 175 iii) causal structure describing relationships between variables, and iv) task instruction. The context 176 provides external knowledge on data collection (e.g., "this dataset includes subject's socioeconomic 177 *factors...*"). The data schema contains the meta-information of variables (e.g., name, description, 178 type, and categorical values). These elements remain constant during optimization. The causal structure, represented as a DAG and converted into text format (x_1, x_2) , indicates x_1 causes x_2 . 179 Various serialization techniques were tested, but the original structured format proved most effective. 180 The initial causal structure is heuristically determined (e.g., Hill climbing (37)). The task instruction 181 guides the goal, such as "produce accurate and convincing synthetic data". Through adversarial 182 optimization, the causal structure and instructions are refined to reduce the discriminator's accuracy. 183 Thus, for each iteration i, θ_i is: 184

$$\theta_i = [\text{context}][\text{schema}], [\text{causal structure}]_i[\text{task instruction}]_i.$$
 (1)

Note that subscription for iteration i will be omitted for simplicity without loss of generalizability. Also, note that we used causal structure as a means to convey the relationship between variables within the prompt; thus, obtaining ground-truth causal structure is not our ultimate goal.

Few shot examples. The data generation process θ is supplemented with n examples to leverage 190 in-context few-shot learning. Structured data (x, y) is serialized into JSON format, e.g., *[age: 53,* 191 work class: self-emp, ...] (Supplement listing 2 Lines 25-28). Various natural language serializations 192 were tested but had minimal impact on performance. The number n of examples is crucial; a large 193 n allows learning from diverse examples but is constrained by context length, while a small n194 avoids overflow but underutilizes data. Our solution, "batches in a batch," splits a batch into smaller 195 pieces that fit the input token size, generates a set of synthetic data, and collates them into D_{syn} (see 196 Algorithm 1 Line 6). This approach balances the trade-offs in in-context few-shot learning. The final 197 input to the generator LLM is:

$$\theta_i, [JSON((\mathbf{x}, y)_1)], [JSON((\mathbf{x}, y)_2)], \dots, [JSON((\mathbf{x}, y)_n)]$$
(2)

²⁰⁰ for each optimization iteration. See Supplement listing 2 for a full example.

LLM as generator. With the prompt in Eq. 2, the pre-trained, frozen LLM (e.g., GPT-3.5) generates synthetic data. The goal is to create similar but not identical text to the *n* samples, with the temperature parameter controlling variability. The temperature is set low enough to maintain the original data distribution but high enough to avoid copying. The generator LLM runs multiple times with smaller examples in a batch, and the generated data is collated into $D_{syn.i}$ denotes the synthetic data generated at iteration *i*. See Supplement listing 3 for an example.

Conditional sampling. As MALLM-GAN generates synthetic data using LLM, it seamlessly inherits the benefits of LLM, such as conditional sampling. LLM predicts the next tokens given a user-provided context, even when the specified condition is rare or given as a range. To conditionally sample the synthetic data, we modify the task instruction to contain specific conditions (Supplement Listing 4).

- 213 3.2.2 DISCRIMINATOR
- Based on the generated data, we evaluate and score the quality of θ by assessing how easy it is to distinguish generated synthetic data from real data. Naturally, this is a supervised learning rather than a reasoning task with LLMs. We build a discriminator f such that $f : \mathcal{X} \to c$ where $\mathbf{x} \in \mathcal{X}$ and $f(\mathbf{x})$

is the predicted label c, which is 1 if $\mathbf{x} \in D_{train}$ and 0 if $\mathbf{x} \in D_{syn}$. Specifically, at each iteration i, a new set of synthetic data $D_{syn,i}$ is generated. We form the combined dataset $D_{train} \cup D_{syn,i}$. We assign labels to the combined dataset by

$$D_i = \{ (\mathbf{x}, c) \mid \mathbf{x} \in D_{train}, c = 1 \} \cup \{ (\mathbf{x}, c) \mid x \in D_{syn,i}, c = 0 \}.$$

We update the discriminator f_i incrementally based on f_{i-1} . We evaluate the accuracy of the discriminator with D_{test} and pass a pair of $(\theta_i, L(f_i))$ to the optimizer where L(f) denotes the discriminatory power of f (e.g., accuracy, likelihood). We prefer to use accuracy (rather than likelihood) because this is a direct measurement we aim to increase and because our optimizer does not require numerical derivatives.

The discriminator obtains better discriminatory accuracy to distinguish real or synthetic data as the discriminator accumulates the discriminatory power of past iterations 0, ..., i - 1 and is updated with newly generated, more realistic synthetic data from the current iteration *i*. However, on the other hand, as the D_{syn} becomes more realistic over the iterations, it gets easier to fool the discriminator, and the discriminator's accuracy decreases. Therefore, our discriminator obtains better discriminatory power during this adversarial optimization.

232 3.2.3 OPTIMIZER

The next task is to optimize θ_i based on its score $L(f_i)$. Our parameter to optimize is θ , a text, which doesn't have derivatives. So we use Optimization by Prompting, which leverages LLM as an optimizer (43). To make LLM acts as a optimizer, we provide a meta-prompt, which consists of the causal structure and the optimization task descriptions such as "Your task is to optimize prompts for generating high-quality synthetic data. Aim ..." (see Example in Supplement listing 5 Line 3-6).

To leverage LLM's in-context few-shot learning in the optimizer (43), we provide a few "examples" of possible solutions $(\theta, L(f))$. Note that the example here is different from data (\mathbf{x}, y) . We keep the top k solution pairs over the past iteration as the optimization trajectory to guide the optimization. We sort the score, so that the more desirable θ goes to the end of prompt. This will allow the LLM to recognize patterns among the data generation process with better score. See example in Supplement listing 5 Line 9-31.

A potential pitfall is that the L(f) of past iteration 0, ..., i-1 is not comparable to the L(f) of current iteration *i*. The past discriminators $f_0, ..., f_{i-1}$ have much lower performance in discriminating real and fake, thus the score $L(f_0), ..., L(f_{i-1})$ are not reliable to compare θ of past iterations with θ of current iteration. Thus we adjust the score L(f) of past iterations with the latest discriminator f_i , so that all the scores are directly comparable to select best θ .

In all, the optimizer LLM takes as input the meta prompt and a series of data generation process θ and adjusted scores $L(f_i)$). The optimizer outputs the revised data generation process, particularly focusing on causal structure and task instruction. We repeat the iterative optimization and generation until reaching to the maximum iteration.

254	1	def optimize_theta(theta):
255	2	theta_score_pairs = []
256	3	<pre>for _ in range(max_epoch):</pre>
257	4	for batch in D_train:
231	5	# 1. Run generator
258	6	D_syn = [LLM_generator(theta + example) for example in
259		batch]
260	7	# 2. Run discriminator
261	8	labels_syn, labels_train = [0] * len (D_syn), [1] * len (
262		D_train)
202	9	(train, test), (train_label, test_label) =
263		<pre>train_test_split(concat(D_train, D_syn), concat(</pre>
264		labels_train, labels_syn))
265	10	discriminator.update(train, train_label)
266	11	<pre>score = get_accuracy(discriminator.predict(test),</pre>
267		test_label)
207	12	# 3. Run optimizer
268	13	theta_score_pairs.append((theta, score))
269	14	<pre>theta = LLM_optimizer(instruction + str(theta_score_pairs))</pre>

return theta

270 271

15

16

17

18

- 272
- 273 274
- 275 276

277

3.3 COMPARISON TO GAN AND CONVERGENCE

def generate_synthetic_data(theta):

278 MALLM-GAN's adversarial training is motivated by GAN, but it differs fundamentally from traditional GANs in that it operates in a natural language optimization space, using LLMs for prompt-based 279 generation, which lacks formal mathematical guarantees. Unlike gradient-based optimization in 280 GANs, MALLM-GAN's optimization relies on empirical refinement of prompts through adversarial 281 optimization with a discriminator. Theoretical convergence analysis is challenging due to the absence 282 of numerical gradients. However, empirical convergence, demonstrated in our experiments in Section 283 4 and prior work (43), shows stable convergence, where discriminator accuracy declines as prompts 284 are refined. This practical convergence criterion serves as a reliable alternative to formal guarantees 285 in real-world tasks. 286

return [LLM_generator(theta + example) **for** example **in** D_train]

Listing 1: Python style pseudocode for MALLM-GAN's optimization and generation

²⁸⁷ 4 EXPERIMENTS

We present the evaluation results of MALLM-GAN. Our extensive experiments demonstrated that
 MALLM-GAN outperforms baselines in generating high-quality synthetic data while preserving
 data privacy, thanks to the adversarial optimization of the data generation process. Additionally,
 MALLM-GAN' provides explainable data generation through natural textual representation, effectively
 generating high-quality synthetic data based on user-provided conditions, even for rare categorical
 values or numeric ranges.

295 4.1 SETTING

LLM. We used HIPPA-compliant Azure OpenAI GPT-3.5(7) as our generator and GPT-4 (25)(gpt-4-32k-0613) as our optimizer. Due to the extensive workload of the generator, we opted for the lighter, faster gpt-35-turbo-0125 model with a 16k context length. For the optimizer, requiring combinatorial search and high-level reasoning, we used the up-to-date gpt-4-32k-0613 model. As the optimizer requires more "creativity" than the generator, the generator's temperature was set to 0.5, and the optimizer's to 1.0 after multiple trials.

302 Discriminator. Strong discriminators do not always contribute to a better generator (3). We tested
 303 Logistic regression, XGBoost, and neural network; we used the logistic regression model because
 304 it showed the highest performance while ensuring tractability during incremental updates over the
 305 iterations (Supplementary 6).

Data. Our benchmarks include several datasets from various domains: three public datasets (Adult(4), Medical Insurance(1), Asia(30)), and two private medical datasets (ATACH2, ERICH) (22). To ensure fair comparison without memorization concerns of LLM (e.g., public datasets are in the training corpus of LLM), private datasets were included. Details are in Supplement Table 5.

Baselines. We compared MALLM-GAN with multiple state-of-the-art tabular generative models such as: traditional over-sampling techniques, SMOTE (9), the variational auto-encoder, TVAE (42), the generative adversarial network, CTGAN (42), LLM-based synthetic data generation model, Be-GReaT(6), and a diffusion model, TabDDPM (19). Similar to MALLM-GAN, a prior work (31) uses in-context few-shot learning of pre-trained LLMs but incorporates post-hoc data selection, which is beyond our scope. A comparison without post-hoc selection is available in Table 3.

- Other hyperparameters. Various serialization techniques have been proposed to transform tabular data into natural language text (14). We tested several serializations, such as Manual Template (14) ("*Age is 53, Work class type is self-employed,* ..."), which proved ineffective for moderate feature sizes; this serialization made the input prompt lengthy and talkative, only worked when feature size |x| is very small. Feature order permutation (6) also had negligible impact on performance. Specific hyperparameters and computing resources are available in Supplement Section 2.
- **Training data size vs. quality of synthetic data**. We evaluated the impact of training data size $N = |D_{\text{train}}|$ on synthetic data quality by sampling subsets of different sizes (N = 100, 200, 400, 800).

324 We particularly aimed to compare performances in low and moderate data size. For fair comparison 325 between real and synthetic data, synthetic data was generated to match the size of real data ($|D_{train}|$ = 326 $|D_{syn}|$). We held out 200 samples as the test set D_{test} before sampling and replicated experiments 327 for each sub sample five times to estimate the standard error of the evaluation metrics. The batch size 328 was set to be 50, with maximum iterations set to be 5, 4, 3, 2 for data sizes N = 100, 200, 400, 800,respectively.

4.2 PERFORMANCE EVALUATION 331

330

345

332 We evaluate the performance of synthetic data generation models from two perspectives: Privacy leakage by Distance to Closest Records (DCR) and Machine Learning Efficiency (MLE) (10; 42). 333

334 **MLE.** To assess the utility of our synthetic data, we use it to train supervised model and test prediction 335 accuracy on real data (D_{test}). The Adult data is used for a classification task, while the other three 336 datasets are used for regression. For classification, we fit logistic regression model, random forest, 337 and Support Vector Machine model, XGBoost Classifier, calculating F1 score. For regression, we 338 fit linear regression, random forest, XGBoost Regressor and calculate R^2 . For each model, we 339 report the average of the best scores for each random seed. We also fit models using real data D_{train} 340 as a gold standard of the MLE for comparison. As a result, MALLM-GAN generated high-quality 341 synthetic tabular data across multiple datasets and training data size, outperforming baselines (Table 2), specially with small training sizes (N = 100). This indicates MALLM-GAN's robustness to smaller 342 sample sizes, unlike baselines that require more data. MALLM-GAN also outperformed baselines on 343 both public and private datasets, suggesting it does not rely on the pre-trained LLM's memorization. 344

Table 2: Benchmark MLE results over 5 datasets. Baseline results were obtained from training the 346 supervised models directly on the real data. SMOTE* interpolates data within the training set, thus it 347 gets higher accuracy by copying training data and compromising DCR. 348

349				Public dataset		Private	lataset
350			$\operatorname{Adult}(F1)$	Asia (F1)	Insurance (R^2)	$ATACH(R^2)$	$\operatorname{ERICH}(R^2)$
351		Real data	0.86	0.83	0.82	0.26	-0.04
352	N=100	SMOTE*	0.78 ± 0.01	0.83 ± 0.00	0.80 ± 0.01	0.27 ± 0.03	-0.15 ± 0.13
353	11-100	TabDDPM CTGAN	0.75 ± 0.01 0.66 ± 0.06	-0.63 ± 0.19	-5.26 ± 0.42 -0.09 ± 0.11	-0.99 ± 0.33 -0.40 ± 0.21	-0.19 ± 0.05 -0.33 ± 0.11
354		TVAE	0.67 ± 0.05	0.83 ± 0.01	0.39 ± 0.15	-0.01 ± 0.07	-0.11 ± 0.12
355		Be-GReaT MALLM-GAN	$\begin{array}{c} 0.71 \pm 0.03 \\ 0.79 \pm 0.02 \end{array}$	0.83 ± 0.00 0.83 ± 0.00	0.54 ± 0.10 0.72 ± 0.00	-0.25 ± 0.23 0.27 ± 0.07	$-0.38 \pm 0.12 \\ -0.03 \pm 0.07$
356		Real data	0.85	0.83	0.83	0.27	0.16
357	N=200	SMOTE*	0.78 ± 0.04	0.83 ± 0.00	0.79 ± 0.02	0.31 ± 0.04	0.05 ± 0.06
358	11-200	TabDDPM CTGAN	$\begin{array}{c} 0.60 \pm 0.15 \\ 0.61 \pm 0.02 \end{array}$	$-$ 0.71 \pm 0.10	$0.56 \pm 0.14 \\ -0.12 \pm 0.08$	$-0.55 \pm 0.33 \\ -0.27 \pm 0.05$	$-0.30 \pm 0.06 \\ -0.19 \pm 0.10$
359		TVAE	0.67 ± 0.05	0.82 ± 0.01	0.62 ± 0.05	0.08 ± 0.06	-0.08 ± 0.07
360		MALLM-GAN	0.69 ± 0.05 0.77 ± 0.03	0.82 ± 0.00 0.83 ± 0.01	0.72 ± 0.03 0.69 ± 0.04	0.16 ± 0.06 0.28 ± 0.07	-0.18 ± 0.16 0.02 ± 0.02
361		Real data	0.83	0.84	0.85	0.31	0.18
362	N=400	SMOTE*	0.85 ± 0.03	0.84 ± 0.00	0.83 ± 0.00	0.32 ± 0.02	0.07 ± 0.05
363		TabDDPM CTGAN	0.82 ± 0.03 0.63 ± 0.02	- 0.59 \pm 0.17	0.79 ± 0.03 -0.18 \pm 0.10	$\begin{array}{c} 0.36 \pm 0.02 \\ -0.08 \pm 0.07 \end{array}$	$0.09 \pm 0.04 \\ -0.24 \pm 0.10$
364		TVAE Be GReeT	0.71 ± 0.07 0.79 ± 0.04	0.71 ± 0.07 0.79 ± 0.00	0.62 ± 0.05 0.72 + 0.03	0.16 ± 0.08 0.20 ± 0.06	-0.19 ± 0.06 -0.13 ± 0.07
365		MALLM-GAN	0.79 ± 0.04 0.79 ± 0.02	0.83 ± 0.00	0.72 ± 0.03 0.71 ± 0.03	0.20 ± 0.00 0.27 ± 0.04	0.02 ± 0.03
366		Real data	0.71	0.84	0.85	0.40	0.21
367	N=800	SMOTE*	0.71 ± 0.03	0.84 ± 0.00	0.83 ± 0.00	0.37 ± 0.03	0.10 ± 0.05
368		TabDDPM CTGAN	0.70 ± 0.03 0.64 ± 0.05	-0.48 ± 0.06	0.83 ± 0.01 -0.41 ± 0.06	$-0.53 \pm 0.45 \\ -0.05 \pm 0.06$	$0.12 \pm 0.04 = -0.04 \pm 0.02$
369		TVAE Be-GReaT	$0.77 \pm 0.02 \\ 0.75 \pm 0.07$	$0.82 \pm 0.01 \\ 0.82 \pm 0.00$	0.68 ± 0.01 0.53 ± 0.21	0.12 ± 0.07 0.00 ± 0.07	-0.05 ± 0.03 -0.04 ± 0.05
370		MALLM-GAN	0.80 ± 0.02	0.84 ± 0.00	0.72 ± 0.01	0.36 ± 0.02	0.02 ± 0.02

37 371

372 **DCR distributions.** The DCR metric assesses the realism and diversity of synthetic data. It 373 determines whether synthetic data points are too similar to the real data points (potential pri-374 vacy leakage) or too dissimilar (hurting the utility of the synthetic data). The DCR is defined 375 as $d(\mathbf{x}_{syn}, D_{real}) = \min_{\mathbf{x}_{real} \in D_{real}} l_1$ -norm $(\mathbf{x}_{syn}, \mathbf{x}_{real})$. Low DCR indicates that synthetic data are very close to real data points, implying a privacy leakage, as synthetic data too closely mimic the real 376 data. We chose to use l_1 -norm distance to measure the distance between two data points (6). For 377 the categorical variables, the distance is 1 if two categories are different; otherwise, 0. As a result,

MALLM-GAN achieved similar or higher DCR levels compared to baseline models (Fig. 2), imply ing effective privacy protection without compromising MLE. Overall, MALLM-GAN demonstrated
 superior performance in generating synthetic data with small data as balancing privacy and utility.



Figure 2: DCR between the synthetic data and the real data. DCR were calculated based on training data and held-out test data for each model. A good model should have similar distributions between the DCR to training and the DCR to held-out dataset.

4.3 ABLATION STUDY

Number n of example in in-context few shot learning. Due to the LLM's limited context length, we implemented a "batches in a batch" method to leverage all training data within these constraints (Section 3.2.1). We varied the number n of examples and found n = 1 to be optimal, achieving high DCR without compromising MLE (Supplement Section 2.5).

395 **Causal structure and Optimization**. To assess the impact of each component on overall performance. we examined the contribution of the causal structure in the data generation process θ and the LLM as 397 an optimizer. We compared the full model, which includes both components, to a version without 398 them, similar to CLLM (31) without post-processing data selection (Table 3). The ablation study 399 showed that incorporating the causal structure alone did not significantly improve the MLE compared 400 to a model with only in-context few-shot learning. However, the LLM optimizer improved θ using 401 prior knowledge encoded in LLM and finally achieved the highest MLE. Incorporating external 402 knowledge into LLMs has been shown to significantly improve the quality of generated text, similar to retrieval-augmented generation (RAG) (20). Our approach shares this concept by incorporating a 403 "knowledge" graph but optimizes the knowledge itself through adversarial optimization. 404

Table 3: MLE of ablated models to evaluate the effects of causal structure in data generation process and optimization via LLM. Causal: Causal structure in data generation process, Opt: Optimization by LLM.

	Few-shot	Few-shot+Causal	Few-shot+Causal+Opt (ours)
Adult (F1)	0.7550 ± 0.0454	0.7503 ± 0.0393	0.7892 ± 0.0358
Asia $(F1)$	0.2335 ± 0.0000	0.2756 ± 0.2842	0.8282 ± 0.0041
Insurance (R^2)	0.6821 ± 0.0193	0.6718 ± 0.0916	0.7152 ± 0.0447
ATACH (R^2)	0.1581 ± 0.0850	0.1326 ± 0.0637	$\bf 0.2726 \pm 0.0707$
ERICH (R^2)	-0.0647 ± 0.0701	0.0281 ± 0.0424	-0.0253 ± 0.0671

414 415 416

417

405

406

407

386

387

389 390

391

392

394

4.4 Optimization trajectory of data generation process

A key advantage of MALLM-GAN is its transparent data generation process, described in natural 418 text, which allows us to observe the evolution trajectory of the data generation mechanism during 419 adversarial optimization. We present examples of optimization trajectories. We showed how the 420 causal structure evolves to ground truth (Fig. 3) over iteration. We used the Asia dataset because it has 421 known ground-truth causal structures and reported graph edit distance (GED) between ground truth 422 and identified causal structures. In this example, the heuristically initialized causal structure gradually 423 converges to ground truth, thanks to the knowledge obtained from the pre-trained LLM. Different 424 convergence patterns were observed with different initialization strategies (Table 9), supporting the 425 benefit of our heuristic initialization. We also investigated how the task instruction in the generator 426 prompt gets sophisticated and how the discriminator's accuracy changes over iterations. In Table 4, 427 the task instructions evolved to include specific details, and the discriminator's accuracy decreased, 428 implying that synthetic data gets indiscriminative to real data.

429 430 4.5 CONDITIONAL SAMPLING

431 We leverage the generator's conditional generative capability to create synthetic data with userprovided conditions, focusing on categorical values and numerical ranges. We compare MALLM-GAN

Epoch 0 (GED=5)

Ground Truth

Figure 3: An example of trajectory of causal structure in data generation process over adversarial optimization using Asia dataset. T: Tuberculosis, V: Visit to Asia, S: Smoke, LC: Lung cancer, T/L: Tuberculosis or Lung cancer, CX: Chest X-ray, D: Dyspnea, B: Bronchitis

Epoch 1 (GED=3)

Epoch 2 (GED=1)

Epoch 3, 4 (GED=0)

Iteration	Task instruction	Score
Epoch 1	"The ultimate goal is to produce accurate and convincing synthetic data that	100.0%
	dutifully represents these causal relationships given the user provided samples."	
Epoch 2	"The ultimate goal is to create a detailed and convincing dataset that accurately mirrors these causal pathways. While synthesizing your data, keep in mind the fol- lowing key relationships: a 'visit to Asia' increases the likelihood of 'tuberculosis', 'smoking' can lead to 'lung cancer' and 'bronchitis', and both 'tuberculosis' and 'lung cancer' can contribute to 'either tuberculosis or lung cancer', which in turn can lead to 'Dyspnea'. Also, take note of how both 'tuberculosis' and 'lung cancer' are associated with 'chest X-ray' results. Your data should reflect these intricate relationships while remaining consistent and realistic."	76.19%
Epoch 4	"You are tasked with generating a synthetic dataset that faithfully demonstrates the given causal connections. Make sure the dataset illustrates how a 'visit to Asia' can cause 'tuberculosis', how 'smoking' can lead to 'lung cancer' and 'bronchitis', and how either 'tuberculosis' or 'lung cancer' can eventually incite 'Dyspnea'. Also, the dataset should reasonably reveal how a 'chest X-ray' ties in with 'tuberculosis' and 'lung cancer'. Ensure the synthetic data reflects realistic scenarios where these factors interact, affecting each other exactly as per these defined causal relationships."	66.67%

Table 4: Trajectory of task instruction in data generation process over adversarial optimization. Lower score is the better.

and baseline models by visualizing lower-dimensional projections (UMAP). For categorical conditions, we selected three rare conditions in the ERICH dataset: i) *hematoma location = right putaminal*, ii) *GCS score = 13*, and iii) *prior history in vascular disease*. The conditions were met by 187, 83, and 29 patients, respectively. All three baselines failed to generate synthetic data due to insufficient training data. In contrast, MALLM-GAN successfully generated data with distributions similar to the real data (Fig. 4). For numeric range conditions, we selected 'age' > 65 in the ERICH dataset, met by 534 patients. The baselines were unable to incorporate numeric range conditions by design. However, MALLM-GAN successfully generated data satisfying the condition (Fig. 4), demonstrating its ability to understand and flexibly apply conditions in natural text format.





4.6 LIMITATIONS

The proposed framework has several shortcomings. Firstly, due to the limited context length of
 the LLM, our model struggles with high-dimensional datasets with too many categorical variables,
 which make the context information lengthy and reduce the success rate of data generation. Another
 limitation of synthetic data generation introduced by LLM is that the LLM struggles with random

number generation as pointed out in (15), which cast negative effects on our framework's potential
when dealing with datasets of many continuous variables. Additionally, mimicking the traditional
GAN framework, it suffers from a theoretical convergence guarantee. While our model performs well
with small sample sizes, showing better results than other baselines, the improvement diminishes
with larger datasets. Moreover, the training and generation process is costly when dealing with large
data volumes.

5 CONCLUSION

We propose a novel framework to generate synthetic tabular data by leveraging multi-agent LLMs to address the limited sample size issues that are prevalent in healthcare. Compared with other LLM-based methods, we propose an in-context learning approach that does not require fine-tuning on LLM but still leverages the whole data. We use causal structure to guide the data generation process and mimic a GAN architecture to optimize the process. We demonstrate that our model can generate high-quality synthetic data while preserving the privacy of real data. Moreover, compared with other black box models, our proposed work enables transparent data generation that allows domain experts to control the process.

540 REFERENCES

543

544

545

546

547

548

549 550

551

552 553

554

555 556

558

559

561

562

563

564

565 566

567

568

569

570

571

572 573

574

575

576

577

578 579

580

581

582 583

584

585

586

588

589

590

- 542 [1] Medical cost personal datasets, May 2018.
 - [2] AN, S., AND JEON, J.-J. Distributional learning of variational autoencoder: Application to synthetic data generation. In *Advances in Neural Information Processing Systems* (2023), A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, Eds., vol. 36, Curran Associates, Inc., pp. 57825–57851.
 - [3] ARJOVSKY, M., AND BOTTOU, L. Towards principled methods for training generative adversarial networks, 2017.
 - [4] BECKER, B., AND KOHAVI, R. Adult. UCI Machine Learning Repository, 1996. DOI: https://doi.org/10.24432/C5XW20.
 - [5] BERNTON, E., JACOB, P. E., GERBER, M., AND ROBERT, C. P. Approximate Bayesian Computation with the Wasserstein Distance. *Journal of the Royal Statistical Society Series B: Statistical Methodology* 81, 2 (02 2019), 235–269.
 - [6] BORISOV, V., SESSLER, K., LEEMANN, T., PAWELCZYK, M., AND KASNECI, G. Language models are realistic tabular data generators. In *The Eleventh International Conference on Learning Representations* (2023).
 - [7] BROWN, T., MANN, B., RYDER, N., SUBBIAH, M., KAPLAN, J. D., DHARIWAL, P., NEELAKANTAN, A., SHYAM, P., SASTRY, G., ASKELL, A., ET AL. Language models are few-shot learners. Advances in neural information processing systems 33 (2020), 1877–1901.
 - [8] CHAWLA, N., BOWYER, K., HALL, L., AND KEGELMEYER, W. Smote: Synthetic minority over-sampling technique. J. Artif. Intell. Res. (JAIR) 16 (06 2002), 321–357.
 - [9] CHAWLA, N. V., BOWYER, K. W., HALL, L. O., AND KEGELMEYER, W. P. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research 16* (2002), 321– 357.
 - [10] FANG, X., XU, W., TAN, F. A., ZHANG, J., HU, Z., QI, Y., NICKLEACH, S., SOCOLINSKY, D., SENGAMEDU, S., AND FALOUTSOS, C. Large language models(llms) on tabular data: Prediction, generation, and understanding – a survey, 2024.
 - [11] GRINSZTAJN, L., OYALLON, E., AND VAROQUAUX, G. Why do tree-based models still outperform deep learning on typical tabular data? In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track* (2022).
 - [12] GULATI, M. S., AND ROYSDON, P. F. TabMT: Generating tabular data with masked transformers. In *Thirty-seventh Conference on Neural Information Processing Systems* (2023).
 - [13] GUO, T., CHEN, X., WANG, Y., CHANG, R., PEI, S., CHAWLA, N. V., WIEST, O., AND ZHANG, X. Large language model based multi-agents: A survey of progress and challenges. arXiv preprint arXiv:2402.01680 (2024).
 - [14] HEGSELMANN, S., BUENDIA, A., LANG, H., AGRAWAL, M., JIANG, X., AND SONTAG, D. Tabllm: Few-shot classification of tabular data with large language models. In *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics* (25–27 Apr 2023), F. Ruiz, J. Dy, and J.-W. van de Meent, Eds., vol. 206 of *Proceedings of Machine Learning Research*, PMLR, pp. 5549–5581.
 - [15] HOPKINS, A. K., RENDA, A., AND CARBIN, M. Can llms generate random numbers? evaluating llm sampling in controlled domains. In *ICML 2023 Workshop: Sampling and Optimization in Discrete Space* (2023).
 - [16] HUANG, J.-T., LI, E. J., LAM, M. H., LIANG, T., WANG, W., YUAN, Y., JIAO, W., WANG, X., TU, Z., AND LYU, M. R. How far are we on the decision-making of llms? evaluating llms' gaming ability in multi-agent environments. arXiv preprint arXiv:2403.11807 (2024).

- 594 [17] KAPLAN, J., MCCANDLISH, S., HENIGHAN, T., BROWN, T. B., CHESS, B., CHILD, R., 595 GRAY, S., RADFORD, A., WU, J., AND AMODEI, D. Scaling laws for neural language models, 596 2020. 597 [18] KICIMAN, E., NESS, R., SHARMA, A., AND TAN, C. Causal reasoning and large language 598 models: Opening a new frontier for causality. arXiv preprint arXiv:2305.00050 (2023). 600 [19] KOTELNIKOV, A., BARANCHUK, D., RUBACHEV, I., AND BABENKO, A. Tabddpm: Mod-601 elling tabular data with diffusion models. In International Conference on Machine Learning 602 (2023), PMLR, pp. 17564–17579. 603 [20] LEWIS, P., PEREZ, E., PIKTUS, A., PETRONI, F., KARPUKHIN, V., GOYAL, N., KÜTTLER, 604 H., LEWIS, M., TAU YIH, W., ROCKTÄSCHEL, T., RIEDEL, S., AND KIELA, D. Retrieval-605 augmented generation for knowledge-intensive nlp tasks, 2021. 606 607 [21] LI, T., SHETTY, S., KAMATH, A., JAISWAL, A., JIANG, X., DING, Y., AND KIM, Y. 608 Cancergpt for few shot drug pair synergy prediction using large pretrained language models. 609 npj Digital Medicine 7, 1 (Feb 2024), 40. 610 [22] LING, Y., TARIQ, M. B., TANG, K., ARONOWSKI, J., FANN, Y., SAVITZ, S. I., JIANG, 611 X., AND KIM, Y. An interpretable framework to identify responsive subgroups from clinical 612 trials regarding treatment effects: Application to treatment of intracerebral hemorrhage. PLOS 613 Digital Health 3, 5 (05 2024), 1–17. 614 615 [23] LONG, S., PICHÉ, A., ZANTEDESCHI, V., SCHUSTER, T., AND DROUIN, A. Causal discovery 616 with language models as imperfect experts, 07 2023. 617 [24] NORI, H., LEE, Y. T., ZHANG, S., CARIGNAN, D., EDGAR, R., FUSI, N., KING, N., 618 LARSON, J., LI, Y., LIU, W., LUO, R., MCKINNEY, S. M., NESS, R. O., POON, H., QIN, T., 619 USUYAMA, N., WHITE, C., AND HORVITZ, E. Can generalist foundation models outcompete 620 special-purpose tuning? case study in medicine. November 2023. 621 622 [25] OPENAI. Gpt-4 technical report, 2023. 623 [26] PEARL, J. Causality, 2 ed. Cambridge University Press, Cambridge, UK, 2009. 624 625 [27] QURESHI, A. I., PALESCH, Y. Y., BARSAN, W. G., HANLEY, D. F., HSU, C. Y., MARTIN, 626 R. L., MOY, C. S., SILBERGLEIT, R., STEINER, T., SUAREZ, J. I., TOYODA, K., WANG, Y., 627 YAMAMOTO, H., AND YOON, B.-W. Intensive blood-pressure lowering in patients with acute 628 cerebral hemorrhage. New England Journal of Medicine 375, 11 (2016), 1033-1043. 629 [28] RAMESH, A., PAVLOV, M., GOH, G., GRAY, S., VOSS, C., RADFORD, A., CHEN, M., AND 630 SUTSKEVER, I. Zero-shot text-to-image generation, 2021. 631 632 [29] ROMERA-PAREDES, B., BAREKATAIN, M., NOVIKOV, A., BALOG, M., KUMAR, M., 633 DUPONT, E., RUIZ, F., ELLENBERG, J., WANG, P., FAWZI, O., KOHLI, P., AND FAWZI, A. 634 Mathematical discoveries from program search with large language models. *Nature* 625 (12) 2023). 635 636 [30] SCUTARI, M. Learning bayesian networks with the bnlearn r package. arXiv preprint 637 arXiv:0908.3817 (2009). 638 639 [31] SEEDAT, N., HUYNH, N., VAN BREUGEL, B., AND VAN DER SCHAAR, M. Curated llm: 640 Synergy of llms and data curation for tabular augmentation in ultra low-data regimes, 2024. 641 [32] SOLATORIO, A. V., AND DUPRIEZ, O. Realtabformer: Generating realistic relational and 642 tabular data using transformers, 2023. 643 644 [33] SPIRTES, P., GLYMOUR, C., AND SCHEINES, R. Causation, prediction, and search. MIT 645 press, 2001. 646 [34] SPIRTES, P., MEEK, C., AND RICHARDSON, T. An algorithm for causal inference in the 647 presence of latent variables and selection bias (vol. 1), 1999.
 - 12

- [35] TALEBIRAD, Y., AND NADIRI, A. Multi-agent collaboration: Harnessing the power of intelligent Ilm agents. *arXiv preprint arXiv:2306.03314* (2023).
 - [36] TSAMARDINOS, I., BROWN, L. E., AND ALIFERIS, C. F. The max-min hill-climbing bayesian network structure learning algorithm. *Machine learning* 65 (2006), 31–78.
 - [37] TSAMARDINOS, I., BROWN, L. E., AND ALIFERIS, C. F. The max-min hill-climbing bayesian network structure learning algorithm. *Machine Learning* 65, 1 (Oct 2006), 31–78.
 - [38] UPADHYAYA, P., ZHANG, K., LI, C., JIANG, X., AND KIM, Y. Scalable causal structure learning: Scoping review of traditional and deep learning algorithms and new opportunities in biomedicine. *JMIR Med Inform 11* (Jan 2023), e38266.
 - [39] WOO, D., ROSAND, J., KIDWELL, C., MCCAULEY, J. L., OSBORNE, J., BROWN, M. W., WEST, S. E., RADEMACHER, E. W., WADDY, S., ROBERTS, J. N., ET AL. The ethnic/racial variations of intracerebral hemorrhage (erich) study protocol. *Stroke* 44, 10 (2013), e120–e125.
 - [40] WU, J., PLATANIOTIS, K., LIU, L., AMJADIAN, E., AND LAWRYSHYN, Y. Interpretation for variational autoencoder used to generate financial synthetic tabular data. *Algorithms* 16, 2 (2023).
 - [41] WU, Q., BANSAL, G., ZHANG, J., WU, Y., ZHANG, S., ZHU, E., LI, B., JIANG, L., ZHANG, X., AND WANG, C. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. arXiv preprint arXiv:2308.08155 (2023).
 - [42] XU, L., SKOULARIDOU, M., CUESTA-INFANTE, A., AND VEERAMACHANENI, K. *Modeling tabular data using conditional GAN*. Curran Associates Inc., Red Hook, NY, USA, 2019.
 - [43] YANG, C., WANG, X., LU, Y., LIU, H., LE, Q. V., ZHOU, D., AND CHEN, X. Large language models as optimizers.
 - [44] YOUNG, J., GRAHAM, P., AND PENNY, R. Using bayesian networks to create synthetic data, 2009.
 - [45] YU, B., FU, C., YU, H., HUANG, F., AND LI, Y. Unified language representation for question answering over text, tables, and images. In *Findings of the Association for Computational Linguistics: ACL 2023* (Toronto, Canada, July 2023), A. Rogers, J. Boyd-Graber, and N. Okazaki, Eds., Association for Computational Linguistics, pp. 4756–4765.
 - [46] YU, Y., CHEN, J., GAO, T., AND YU, M. Dag-gnn: Dag structure learning with graph neural networks, 2019.
 - [47] ZHANG, T., WANG, S., YAN, S., JIAN, L., AND LIU, Q. Generative table pre-training empowers models for tabular prediction. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing* (Singapore, Dec. 2023), H. Bouamor, J. Pino, and K. Bali, Eds., Association for Computational Linguistics, pp. 14836–14854.
 - [48] ZHENG, X., ARAGAM, B., RAVIKUMAR, P., AND XING, E. P. Dags with no tears: Continuous optimization for structure learning, 2018.

702 SUPPLEMENTARY

704

705 706

707

708

1 EXAMPLE OF PROMPTS AND OUTPUT

Here, we provided example of generator prompt and optimizer prompt. Note that generator prompt evolves over the iteration.

709 System role: 710 2 % Specify role **and** task 3 You are a data generation model. Your task is to understand the 711 instruction below $\ensuremath{\text{and}}$ generate tabular data. 712 4 713 % Context of data 5 714 <context>The dataset include subject's social economic factors and 6 715 demographics with the label that indicates whether their income is higher than 50k. </context> 716 7 717 8 % Data schema 718 <schema> age (numerical), workclass (categorical), education (9 719 categorical), education-num (numerical), marital-status (categorical), occupation (categorical), relationship (720 categorical), race (categorical), sex (categorical), capital-721 gain (numerical), capital-loss (numerical), hours-per-week (722 numerical), native-country (categorical), Income (categorical) 723 </schema> 724 10 725 %Categorical variables and their available categories <categorical variables> workclass: {'Private', 'Local-gov', ' 726 Without-pay', 'Self-emp-not-inc', 'State-gov', 'Federal-gov', 727 'Self-emp-inc'}, education: {'Some-college', 'Masters', '11th', 728 '1st-4th', '7th-8th', 'Bachelors', 'Doctorate', '12th', '5th -6th', 'Prof-school', 'Assoc-voc', 'Assoc-acdm', '10th', '9th', 729 'HS-grad'}, marital-status: {'Divorced', 'Married-spouse-730 absent', 'Married-civ-spouse', 'Never-married', 'Widowed', ' Separated'}, occupation: {'Handlers-cleaners', 'Transport-731 732 moving', 'Sales', 'Prof-specialty', 'Farming-fishing', ' 733 Machine-op-inspct', 'Adm-clerical', 'Other-service', 'Craft-734 repair', 'Protective-serv', 'Exec-managerial', 'Tech-support', 'Priv-house-serv'}, relationship: {'Wife', 'Not-in-family', 735 Other-relative', 'Unmarried', 'Own-child', 'Husband'}, race: {
 Black', 'Amer-Indian-Eskimo', 'Other', 'Asian-Pac-Islander', 736 737 'White'}, sex: {'Male', 'Female'}, native-country: {'Vietnam', 'Mexico', 'Hong', 'Taiwan', 'Italy', 'Portugal', 'Ireland', 739 Guatemala', 'El-Salvador', 'United-States'}, Income: {'>50K', 740 ' <=50K' } </categorical variables> 741 14 742 %causal structure 15 743 <causal structure> Consider this optimized causal graph of the 16 744 data, where a pair (A, B) is used to represent a scenario where A affects B: [('age', 'workclass'), ('education', ' 745 education-num'), ('education-num', 'Income'), ('marital-status
', 'relationship'), ('occupation', 'Income'), ('hours-per-week 746 747 ', 'Income'), ('workclass', 'Income')] 748 749 This adjusted graph introduces 'education-num', which is a key 18 determinant of 'Income'. Be sure to reflect 'age' impact on ' 750 workclass' and 'marital-status' effect on 'relationship'. When 751 creating the 'Income' data, pay careful attention to the 752 roles of 'education', 'education-num', 'occupation', and ' 753 hours-per-week' as stated in the causal graph. 754 19 </causal structure> 755 20 %Task 21

less than 70.0%. </task>

<example> Here are examples **from** real data:

User role:

% Example

28	<pre>', 'education-num': 6.0, 'marital-status': 'Married-civ-spouse ', 'occupation': 'Farming-fishing', 'relationship': 'Husband', 'race': 'White', 'sex': 'Male', 'capital-gain': 0.0, 'capital- loss': 0.0, 'hours-per-week': 60.0, 'native-country': 'United- States', 'Income': '<=50K'}, {'age': 23.0, 'workclass': ' Private', 'education': 'HS-grad', 'education-num': 9.0, ' marital-status': 'Never-married', 'occupation': 'Adm-clerical', 'relationship': 'Own-child', 'race': 'White', 'sex': 'Female', 'capital-gain': 0.0, 'capital-loss': 0.0, 'hours-per-week': 40.0, 'native-country': 'United-States', 'Income': '<=50K'}] </pre>
29	
30	<instruction></instruction>
31	Generate two synthetic samples mimic the provided samples. DO NOT COPY the samples and try to make the generated samples diverse. The response should be formatted strictly as a list in JSON format , suitable for direct use in data processing scripts such as conversion to a DataFrame in Python. No additional text or numbers should precede the JSON data.
	Listing 2: Example of generator prompt
1	ison
2	<pre>[{"treatment": 0, "age": 68.2, "ICH volume": 4.1, "ICH Location": "L Lobar", "IVH volume": 0.2, "GCS score": 14.3, "NIHSS score": 11.7, "Systolic blood pressure": 195.0, "Diastolic Blood Pressure": 83.0, "Hypertension": 1, "Hyperlipidemia": 1, "Type I Diabetes": 0, "Type II Diabetes": 0, "Congestive heart failure": 0, "Atrial Fibrillation": 0, "PTCA": 0, "Peripheral Vascular Disease": 0, "Myocardial fraction": 0, "Anti-diabetic ": 0, "Antihypertensives": 1, "White blood count": 4.3, " Hemoglobin": 12.5, "Hematocrit": 37.7, "Platelet count": 129.0, "APTT": 35.3, "INR": 1.1, "Glucose": 148.0, "Sodium": 145.0, "Potassium": 4.1, "Chloride": 106.0, "CD": 30.1, "Blood urea nitrogen": 18.0, "Creatinine": 1.2, "race": "White", "sex": "</pre>

<task> The ultimate goal is to produce accurate and convincing

relationships. As such, strive **for** a quality score that **is**

[{'age': 53.0, 'workclass': 'Self-emp-not-inc', 'education': '10th

synthetic data that dutifully represents these causal

Listing 3: Example of generator output. We presented an example with sufficiently high DCR (39.7) to protect patient data privacy

Female", "ethnicity": "Hispanic", "mRS score after 30 days":

<Instruction> Generate {number of samples in real data meeting the conditions} synthetic samples with {user-provided conditions}. Response should be formatted strictly as a list in JSON format, suitable for direct use in data processing scripts such as conversion to a DataFrame in Python. No additional text or numbers should precede the JSON data. </Instruction>

Listing 4: Modified instruction in generator prompt for conditional sampling

System role:

 $2.7\}]$

% Specify role **and** task

```
810
        3 Your task is to optimize prompts for generating high-quality
811
              synthetic data. Aim to lower the scores associated with each
812
              casual structure and prompt, where a lower score reflects
813
              better quality. Here are the steps:
          1. Examine the existing prompt-score pairs.
814
          2. Adjust the causal graph to better represent the underlying
815
              relationships by adding or removing connections, and consider
816
              incorporating new features from the list {self.cols}.
817
          3. Modify the prompt guidance to align with the revised causal
       6
818
              graph, ensuring it aids in reducing the score.
819
          User role:
       8
820
          <pair>
       9
821
          Reflecting the adjusted causal graph of the data, where each tuple
       10
822
               (A, B) indicates that A impacts B:
823
          [('age', 'workclass'), ('marital-status', 'relationship'), ('
       11
              marital-status', 'Income'), ('relationship', 'sex'), ('
824
              education', 'Income'), ('occupation', 'Income'), ('workclass',
825
               'Income'), ('hours-per-week', 'Income')]
826
       12
827
          Use this causal graph as a guide to generate synthetic data that
       13
828
              closely mirrors the real-world dataset. Remember to factor in
              the influence of 'age' on 'workclass', and 'marital-status' on
829
               'relationship' and 'Income'. The 'relationship' should guide
830
              the generation of the 'sex' attribute. Further, take into
831
              consideration the effects of 'education', 'occupation', and '
832
              hours-per-week' on 'Income' when synthesizing your data. The
833
              goal is to produce synthetic data that convincingly mimic
              these causal relationships.
834
          Set your aim to achieve a score below 75.0%.
       14
835
       15
          Score: 80.0%
836
          </pair>
       16
837
       17
838
       18
          <pair>
          Consider the revised and detailed causal graph of the data, which
839
       19
              includes ('age', 'workclass'), ('marital-status', '
840
              relationship'), ('relationship', 'sex'), ('education', 'Income
841
              '), ('occupation', 'Income'), ('workclass', 'Income'), ('hours-
              per-week', 'Income'):
842
843
       20
          In light of the causal graph, generate synthetic samples that
       21
844
              mimic the structure in the provided dataset. Values such as '
845
              age' should reflect on 'workclass'; 'marital-status' and '
846
              relationship' should collaborate to inform 'sex', while '
847
              education', 'occupation', 'workclass', and 'hours-per-week'
              should exhibit their influence on 'Income'. Also consider '
848
              marital-status' influence on 'Income'. Your aim is to generate
849
              synthetic data that fully embody the interconnections within
850
              this causal graph.
851
          Aim to achieve a score lower than 75%
       22
852
       23
          Score: 80.95%
853
       24
          </pair>
       25
854
       26
          <pair>
855
          Here is the causal graph of the data, where a tuple (A, B)
856
              indicates A causes B:
857
          [('marital-status', 'relationship'), ('marital-status', 'Income'),
       28
               ('relationship', 'sex')]
858
          Given the description of the data, generate synthetic samples that
       29
859
               mimic the provided samples.
860
          Score: 85.71%
       30
861
       31
          </pair>
862
       32
863
```

869

870 871 872

873

874

875

876

877

878

879

880

881

882

883 884

885

886

887

888

889

890

891

892

893 894

895 896 897

898 899

900

901 902

903

34

35

864

865

Your updated prompt should explicitly include **any** modifications to the causal graph **and** guidance. The aim **is** to create a prompt that leads to the lowest possible score. The updated prompt:

Listing 5: Example of optimizer prompt

<Causal structure> The optimized causal network, suggesting the influence of variable A on variable B, includes the following relationships: [('Age', 'Hyperlipidemia'), ('Hyperlipidemia', 'Type II Diabetes'), ('Type II Diabetes', 'Blood urea nitrogen '), ('Blood urea nitrogen', 'Creatinine'), ('Hypertension', Congestive heart failure'), ('Congestive heart failure', ' Atrial Fibrillation'), ('Atrial Fibrillation', 'GCS score'), ('GCS score', 'mRS score after 30 days'), ('Anti-diabetic', Type I Diabetes'), ('Type I Diabetes', 'Antihypertensives'), ('Antihypertensives', 'Potassium'), ('Potassium', 'Sodium'), (' PTCA', 'Peripheral Vascular Disease'), ('Peripheral Vascular Disease', 'Myocardial fraction'), ('Myocardial fraction', ' Hemoglobin'), ('Hemoglobin', 'Hematocrit'), ('race', ' ethnicity'), ('Sex', 'Hyperlipidemia')]</Causal structure> <Task> Your task is to create realistic synthetic patient data, keeping the altered causal relationships as your guiding principle. Ensure the data reflects a diverse set of potential patient scenarios, evidencing the variety of health conditions one might find in a clinical setting. Remember that the engineered data should present unique, individual patient scenarios, each portraying a different, **complex** clinical situation. The synthetic data needs to be representative of different demographics ('Sex', 'race', 'ethnicity') and should also take into consideration different health conditions $\boldsymbol{\mathsf{and}}$ treatment plans.</Task>"

Listing 6: Example of optimizer output

2 EXPERIMENT DETAILS

2.1 BENCHMARK DATASETS DESCRIPTIONS

We provide detailed description on the benchmark data in Table 5

Table 5: Datasets description

	# samples	# features	Description	Source
Adult	32,561	14	The dataset include people's social economic factors and	(4)
			demographics with the label that indicates whether their	
			income is higher than 50k.	
Medical Insurance	2,772	7	This is a dataset used to describe the paitents' demograph-	(1)
			ics with their health insurance bills.	
Asia	10000	8	This is the dataset used to illustrate the utility of Baysian	
			network to do causal structure discovery. The dataset is	
			available in the R-package(30).	
ATACH2	1,000	37	This is an RCT data that investigate in treatment for In-	(27)
			tracerebral hemorrhage patients.	
ERICH	1,521	29	The data is from a case-control study of Intracerebral	(39)
			Hemorrhage study which aims to investigate in the Eth-	
			nic/Racial variations.	

918	2.2	Hyperparameters
919	Spec	ific hyperparameters for each model are provided below.
921		• CTGAN: Default parameters
922		• TVAF: Default parameters
923		• I VAL. Default parameters
924		• BeGReaT:
925		- Base LLM: Distiled-GPT2
926		Detab sizes 40
927		- Batch size: 40
928		- Epochs: Depend on the feature numbers and the total sample size. (200-400)
929 930		• MALLM-GAN:
931		– Temperature for generator: 0.5
932		Tomporature for antimizer 1.0
933		- Temperature for optimizer. 1.0
934		- Batch size: 50
935		- Discriminator: XGBoost (max denth: 3 eta: 0.3 objective: binary:logistic)
936		- Diserminator. AODoost (max depuit. 5, eta. 0.5, objective. binary.iogistic)
937		• TabDDPM: Default parameters
938	22	COMPARISON AMONG DIFFERENT KINDS OF DISCRIMINATORS
939	2.5	COMPARISON AMONG DIFFERENT KINDS OF DISCRIMINATORS

Table 6: Comparison of different discriminators effects on the quality of the synthetic data. An experiments on sub-sample of Adult data.

		N = 100	N = 200	N = 400	N = 800
Adult (F1 score)	XGBoost Logistic regression Neural Network	$0.78 \pm 0.03 \\ 0.79 \pm 0.02 \\ 0.80 \pm 0.02$	$\begin{array}{c} 0.73 \pm 0.01 \\ \textbf{0.77} \pm \textbf{0.02} \\ 0.57 \pm 0.12 \end{array}$	$\begin{array}{c} 0.76 \pm 0.06 \\ \textbf{0.79} \pm \textbf{0.03} \\ 0.78 \pm 0.06 \end{array}$	$0.72 \pm 0.00 \\ 0.80 \pm 0.02 \\ 0.67 \pm 0.12$

2.4 COMPUTING RESOURCE DETAILS

940

949

953

The model proposed in this study does not require extensive computing resource for fine-tuning.
However, this model require access to Azure service. For other baseline models, they are implemented
on one NVIDIA H100 80GB HBM3 GPU.

2.5 NUMBER OF EXAMPLES IN IN-CONTEXT FEW SHOT LEARNING

954 Given the limited context length that the LLM can understand, we proposed "batches in a batch" 955 method to leverage all training data in limited context length in the generator LLM (Section 3.2.1). 956 We varied the number n of few-shot examples by n = 1, ...5 and measured the MLE (Fig. 5) and 957 DCR distribution (Table 7, 8) to find the optimal number n. As a result, the increasing number n958 of examples did not always increase the MLE of synthetic data (Fig. 5) but decreased the DCR 959 (Table 7), thus increasing privacy concerns. Instead, n = 1 achieved sufficiently high DCR without 960 compromising MLE. The MLE did not increase with more examples because the more examples will 961 increase the context length and the generator LLM overlook some key context information. On the other hand, the DCR decreased with more examples because the generator LLM is more likely to 962 stick to copy the provided examples. Interesting, the increasing number of examples does not affect 963 the DCR of synthetic data generated from public dataset (Adult, Insurance). 964

965 2.6 Optimization trajectory: example on insurance dataset

As seen in Table 9, the causal structure does not converge to the ground truth after 5 epochs. When initialized with the ground truth, the causal structure maintains slight fluctuations. Another example of causal structure discovery on insurance dataset were presented in Supplementary 2.6.

Here is the example for the Insurance dataset. Initially, the causal structure derived from heuristics had
 no edges (Fig. 6). Over iterative optimization, a stable causal structure emerged. The task instructions
 evolved to include specific details (Table 10). Our objective was to reduce the discriminator's





_		1	2	3	4	5
_	Adult Insurance Asia ATACH2 ERICH	4, 7, 10 30, 115, 337 0, 0, 0 84, 100, 120 70, 87, 110	5, 7, 11 34, 91, 405 0, 0, 0 82, 99, 122 66, 82, 111	5, 6, 1036, 76, 2450, 0, 081, 97, 12551, 82, 104	$\begin{array}{c} 4, 7, 10 \\ 24, 64, 170 \\ 0, 0, 0 \\ 79, 98, 124 \\ 62, 80, 108 \end{array}$	$\begin{array}{r} 4, 7, 11 \\ 27, 70, 150 \\ 0, 0, 0 \\ 82, 103, 128 \\ 62, 80, 117 \end{array}$
	Table 9: Gi	raph Edit Distar	nce (GED) betw	ween causal st	ructure in θ an	d ground truth
		Epoch No Ini	tialization Heuristic	Initialization Grou	nd Truth Initialization	-
		0	16 4	5 1	0	-
		2 3	4 3	1 0	0 0	
		4	4	0	2	-
Iterat	ion Task	instruction	to much sin data	al at accurate l		Sumal nota - 9
Iterat Epocl	ion Task h 1 "The tionsl	instruction task is to genera hips. The data sho	te synthetic data puld include vari	ı that accurately jables such as 'a	mirrors these c ge', 'sex', 'bmi',	rausal rela- 8 'children',
Iterat Epocl	ion Task h 1 "The tionsl 'smok	instruction task is to genera hips. The data sho er', 'region', and	te synthetic data ould include vari l 'charges'. Ead	that accurately ables such as 'a ch variable shou	mirrors these c ge', 'sex', 'bmi', uld influence the	vausal rela- 8 'children', e others as
Iterat Epocl	ion Task h 1 "The tionsl 'smok per th scena	instruction task is to genera hips. The data sho er', 'region', and he causal structur trios''	te synthetic data ould include vari d 'charges'. Eac re, creating a rea	i that accurately iables such as 'a ch variable shou ilistic representa	y mirrors these c ge', 'sex', 'bmi', uld influence the ttion of possible	rausal rela- 8 'children', e others as real-world
Iterat Epocl	ion Task h 1 "The tionsl 'smok per th scena h 2 "You	instruction task is to genera hips. The data sho er', 'region', and he causal structur urios." r task is to genera	te synthetic data ould include vari d 'charges'. Eac e, creating a rea rate synthetic da	n that accurately iables such as 'a ch variable shou ulistic representa ata that faithful	mirrors these c ge', 'sex', 'bmi', uld influence the ution of possible ly represents th	rausal rela- 8 'children', e others as real-world ese causal 7
Iterat Epocl	ion Task h 1 "The tionsl 'smok per th scena h 2 "You relati	instruction task is to genera hips. The data she er', 'region', and he causal structur trios." r task is to gener conships. The dat	te synthetic data ould include vari l'charges'. Eau re, creating a rea rate synthetic da a should encom	a that accurately iables such as 'a ch variable shou ilistic representa ata that faithful ipass variables	mirrors these c ge', 'sex', 'bmi', uld influence the tion of possible ly represents th such as 'age', 'l	sausal rela- 'children', e others as real-world ese causal bmi', 'chil-
Iterat Epocl	ion Task h 1 "The tionsl 'smok per th scena h 2 "You relati dren'	instruction task is to genera hips. The data sho er', 'region', and he causal structur urios." r task is to gener ionships. The data , 'smoker', 'regio	te synthetic data ould include vari d'charges'. Eac re, creating a rea rate synthetic da ta should encom n', and 'charges useal structure.	a that accurately iables such as 'a ch variable show ilistic representa ata that faithful pass variables '. Each variable	y mirrors these c ge', 'sex', 'bmi', uld influence the ttion of possible ly represents th such as 'age', 'h should affect th ible simulation	sausal rela- 'children', e others as real-world ese causal for cotential
Iterat Epocl	ion Task h 1 "The tionsl 'smok per th scena h 2 "You relati dren' accor real-v	instruction task is to genera hips. The data sho er', 'region', and the causal structur trios." task is to generationships. The data , 'smoker', 'regio dance with the ca world scenarios"	te synthetic data ould include vari d'charges'. Eac e, creating a rea rate synthetic da a should encom n', and 'charges ausal structure, p	a that accurately iables such as 'a ch variable shou ilistic representa ata that faithful pass variables '. Each variables providing a cred	mirrors these c ge', 'sex', 'bmi', uld influence the ttion of possible ly represents th such as 'age', 't should affect th ible simulation o	rausal rela- 8 'children', e others as real-world ese causal 7 bmi', 'chil- te others in of potential
Iterat Epocl Epocl	ion Task h 1 "The tionsl 'smok per th scena h 2 "You relati dren' accor real-v h 4 "The	instruction task is to genera hips. The data she er', 'region', and te causal structur trios." r task is to gene ionships. The dat , 'smoker', 'regio dance with the co vorld scenarios" ultimate goal is to	te synthetic data ould include vari d'charges'. Eau re, creating a rea rate synthetic da ta should encom n', and 'charges ausal structure, p o generate synthe	that accurately iables such as 'a ch variable shou ulistic representa ata that faithful upass variables '. Each variables providing a cred	y mirrors these c ge', 'sex', 'bmi', uld influence the ttion of possible ly represents th such as 'age', 'h should affect th ible simulation c urately reflects th	e ausal rela- 'children', e others as real-world ese causal bmi', 'chil- the others in of potential these causal 3.
Iterat Epocl Epocl	ion Task h 1 "The tionsl 'smok per th scena h 2 "You relati dren' accon real-v h 4 "The relati 'bmi'	instruction task is to genera hips. The data sho er', 'region', and he causal structur vrios." r task is to gener onships. The data , 'smoker', 'regio dance with the ca vorld scenarios" ultimate goal is to onships. The syn 'children' 'smo	te synthetic data ould include vari d'charges'. Ead rate synthetic da ta should encom n', and 'charges ausal structure, p o generate synthe thetic data shou wher' 'region'	a that accurately ables such as 'a ch variable shou ilistic representa ata that faithful pass variables '. Each variable providing a cred etic data that acc ild incorporate f and 'charges' a	y mirrors these of ge', 'sex', 'bmi', uld influence the tion of possible ly represents th such as 'age', 'h should affect th ible simulation of urately reflects th factors such as ' nd their influen	sausal rela- 'children', e others as real-world ese causal bmi', 'chil- te others in of potential hese causal age', 'sex', ce on each
Iterat Epocl	ion Task h 1 "The tionsl 'smok per th scena h 2 "You relati dren' accor real-v h 4 "The relati 'bmi', other	instruction task is to genera hips. The data sho er', 'region', and he causal structur trios." r task is to generation ordence with the ca vorld scenarios" ultimate goal is to conships. The syn 'children', 'smo as indicated in th	te synthetic data ould include vari d'charges'. Ead rate synthetic da a should encom n', and 'charges ausal structure, p o generate synthe thetic data shou ker', 'region', a ne causal structu	a that accurately iables such as 'a ch variable shou distic representa ata that faithful pass variables '. Each variables providing a cred etic data that acc di incorporate f und 'charges', a re. The synthetic	o mirrors these c ge', 'sex', 'bmi', uld influence the ttion of possible 'ly represents th such as 'age', 'h should affect th ible simulation c urately reflects th factors such as ' nd their influence data should be	sausal rela- 'children', e others as real-world ese causal of potential hese causal age', 'sex', ce on each convincing