

WHEN TABLES LEAK: ATTACKING DIGIT MEMORIZATION IN LLM-BASED TABULAR DATA GENERATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) have recently demonstrated remarkable performance in generating high-quality tabular synthetic data. In practice, two primary approaches have emerged for adapting LLMs to tabular data generation: (i) fine-tuning smaller models directly on tabular datasets, and (ii) prompting larger models with examples provided in context. In this work, we show that popular implementations from both regimes exhibit a tendency to compromise privacy by reproducing memorized patterns of numeric digits from their training data. To systematically analyze this risk, we introduce a simple No-box Membership Inference Attack (MIA) called LevAtt that assumes adversarial access to only the generated synthetic data and targets the string sequences of numeric digits in synthetic observations. Using this approach, our attack exposes substantial privacy leakage across a wide range of models and datasets, and in some cases, is even a perfect membership classifier on state-of-the-art models. Our findings highlight a unique privacy vulnerability of LLM-based synthetic data generation and the need for effective defenses. To this end, we propose two methods, including a novel sampling strategy that strategically perturbs digits during generation. Our evaluation demonstrates that this approach can defeat these attacks with minimal loss to the fidelity and utility of the synthetic data.¹

1 INTRODUCTION

Machine learning systems across diverse domains—from healthcare databases to financial risk assessment platforms—rely heavily on structured tabular datasets (Giuffré & Shung, 2023; Assefa et al., 2021; Flanagan et al., 2022). This widespread dependence has driven significant interest in synthetic tabular data generation, where computational models learn to produce artificial records that statistically mirror the patterns of real datasets while avoiding direct replication (Fonseca & Bação, 2023). The utility of synthetic data stems from two primary advantages: it can supplement limited training samples to improve model performance on underrepresented populations or rare events, and it facilitates private data sharing by generating records that do not directly correspond to actual individuals. These benefits are especially valuable in privacy-sensitive sectors such as medicine (Vallevik et al., 2024) and banking (Wu et al., 2023), where regulatory constraints often limit access to original data. As a result, synthetic data generation has emerged as an essential technique for expanding machine learning capabilities in data-constrained and privacy-regulated environments (Platzer & Reutterer, 2021; McKenna et al., 2022).

Large Language Models (LLMs) have recently emerged as state-of-the-art tabular synthetic data generators. Unlike traditional approaches—such as Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs), and diffusion models—that operate directly over original feature space, LLMs encode tabular rows into string representations and leverage two primary methodologies for generation. In-context learning (ICL) approaches include existing tabular rows within the LLM’s context window and prompt the LLM to generate additional rows following the observed patterns (Seedat et al., 2024; Kim et al., 2024; Hollmann et al., 2025b; Ma et al., 2024), while supervised fine-tuning (SFT) methods train LLMs on larger quantities of string-encoded tabular data to learn the underlying distribution before sampling synthetic records (Borisov et al., 2023; Wang et al., 2025). Both approaches generate synthetic data by producing new string sequences that are subsequently decoded back into tabular format. In both cases, LLM-based generators have been shown to produce synthetic data that exhibits superior statistical fidelity to original datasets and maintains high utility for downstream machine learning tasks. Moreover, initial evaluations suggest these methods may offer enhanced privacy protection relative to conventional approaches (Solatorio & Dupriez, 2023).

However, these architectural differences raise unresolved questions about privacy. Extensive research has demonstrated that LLMs exhibit tendencies to memorize training data, particularly when exposed to repeated

¹An anonymous code repository is available [here](#).

057 patterns (Carlini et al., 2021b; Kandpal et al., 2022), longer sequences (Carlini et al., 2023; Wang et al., 2024),
 058 or through supervised fine-tuning processes (Chu et al., 2025). These memorization behaviors, which can be
 059 beneficial for language modeling tasks, present unique privacy risks in the context of tabular data generation
 060 where training data often contain structured, long sequences of repeated values across observations that have
 061 little semantic meaning.

062 In order to measure privacy risk, Membership Inference Attacks (MIAs) (Shokri et al., 2017; Carlini et al.,
 063 2021a) have emerged as the linchpin of privacy auditing for tabular synthetic data generators, serving as the
 064 primary tool for evaluating privacy risks across diverse generative approaches (Chen et al., 2020; Hilprecht
 065 et al., 2019; van Breugel et al., 2023; Ward et al., 2024; 2025b). However, these methods focus exclusively
 066 on the feature space over which traditional generative models operate, potentially missing the string-space
 067 vulnerabilities introduced by LLM-based generation entirely.

068 To bridge this gap, we examine privacy risks in LLM-based tabular data generation by introducing LevAtt, an
 069 MIA that exploits Levenshtein Distance on the string representations of tabular data—the actual format LLMs
 070 generate—rather than the feature space alone. We find through extensive experimentation in both the ICL
 071 and SFT regimes that state-of-the-art LLMs often memorize and replicate numeric values from training data,
 072 exposing sensitive information digit-for-digit in synthetic outputs. Even under a conservative no-box threat
 073 model, where we assume only adversarial access to the synthetic data, we uncover that LLMs can leak private
 074 data through memorized digit patterns, revealing vulnerabilities that conventional feature-space MIAs fail to
 075 detect.

076 To address this new-found risk, we study two defenses: an ad-hoc post-processing algorithm we call Digit
 077 Modifier (DM) that flips digits to break sequential patterns and a novel Tendency-based Logit Processor (TLP)
 078 that strategically perturbs digits at sample time. We find that both strategies can defeat these attacks, however
 079 TLP can effectively control for privacy leakage with minimal effect on the statistical fidelity or downstream
 080 machine learning utility of the synthetic data.

082 2 RELATED WORK

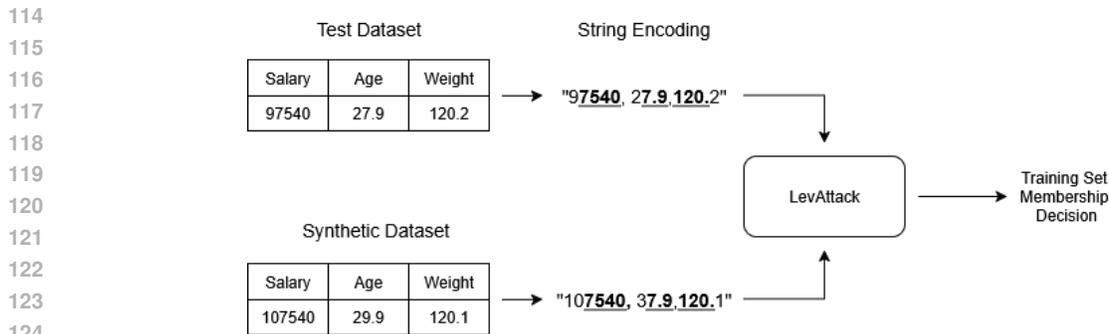
085 2.1 TABULAR DATA GENERATION

086 Tabular generative models aim to learn a generator G from training data T that approximates the true data distri-
 087 bution $p_X(X)$. We represent tabular data as a matrix $\mathbf{X} \in \mathcal{X}^{n \times d}$, where n denotes the number of samples, d the
 088 number of features, and \mathcal{X} is the feature value domain. Each row $\mathbf{x}_i \in \mathcal{X}^d$ represents a data point drawn from
 089 the underlying distribution $p_X(X)$, while columns correspond to features that may have heterogeneous data
 090 types. We denote the j -th feature value of the i -th sample as $\mathbf{x}_{i,j}$. The training dataset $T = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$
 091 comprises n independent samples from $p_X(X)$. The learned model generates synthetic samples $\tilde{\mathbf{x}} \sim G$, form-
 092 ing a synthetic dataset $S = \{\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2, \dots, \tilde{\mathbf{x}}_m\}$. In this work, we assume that features can be of a continuous,
 093 ordinal, or categorical type.

094 **Deep Learning-Based Generation.** In recent years, a variety of conventional tabular synthetic data generators
 095 have been proposed including Generative Adversarial Networks (Yoon et al., 2019; 2020; Xu et al., 2019),
 096 likelihood-based methods Ankan & Panda (2015); Durkan et al. (2019); Watson et al. (2023), and diffusion
 097 models Kotelnikov et al. (2022); Suh et al. (2023); Zhang et al. (2024) Each of these operate directly over the
 098 feature space \mathcal{X}^d , learning mappings $\mathcal{G} : \mathcal{Z} \rightarrow \mathcal{X}^d$ that model the joint distribution $p_X(X)$. Crucially, these
 099 approaches fundamentally treat each sample as an atomic unit, generating a complete feature vector $\mathbf{x} \in \mathcal{X}^d$
 100 simultaneously.

101 **LLM-Based Generation** In contrast, LLM-based approaches reframe tabular generation as autoregressive text
 102 generation. Training samples are encoded into strings, tokenized into sequences from vocabulary \mathcal{W} , and the
 103 LLM generates according to $p(t) = \prod_{k=1}^j p(w_k | w_1, \dots, w_{k-1})$. Rather than modeling the joint distribution
 104 $p_X(X)$ directly, this approach decomposes it through sequential conditioning, generating feature values $x_{i,j}$
 105 token by token based on previously generated values. This sequential generation process fundamentally breaks
 106 the atomic unit assumption, introducing new dynamics where generated tokens influence later ones through the
 107 autoregressive chain.

108 Within this paradigm, two complementary generation strategies have emerged based on data availability and
 109 computational resources. In-context-based methods leverage large foundation models by presenting tabular ex-
 110 amples directly within the context window, enabling few-shot generation without parameter updates in low-data
 111 regimes (Seedat et al., 2024; Kim et al., 2024; Hollmann et al., 2025b; Ma et al., 2024). When larger datasets are
 112 available, SFT-based methods instead fine-tune smaller language models through direct optimization on tabular
 113 generation tasks (Borisov et al., 2023; Solatorio & Dupriez, 2023; Wang et al., 2025). Both strategies main-



125 Figure 1: Diagram of Levenshtein Attack. We simply encode rows of tabular data into a string representation
 126 from which to attack. LevAtt finds signal in the highly constrained and often duplicated sequences of digits in
 127 synthetic tabular data generated by LLMs. In bold and underline: copied sequences of such patterns that are the
 128 source of LevAtt’s adversarial advantage.

129

130

131

132

tain the core autoregressive formulation while differing in how they leverage available data and computational resources.

133 2.2 MEMBERSHIP INFERENCE ATTACKS ON SYNTHETIC TABULAR DATA

134

135 Membership Inference Attacks (MIAs) aim to classify whether a specific observation was a member of the
 136 original dataset used to train a model (Shokri et al., 2017; Chen et al., 2020; Carlini et al., 2021a). Given the
 137 generative model G trained on dataset T as defined above, which generates synthetic dataset S , an adversary
 138 $\mathcal{A} : X \rightarrow \{0, 1\}$ aims to determine if a test sample x^* is an element of T . Formally, this classification or MIA
 139 can be expressed as:

$$140 \mathcal{A}(x^*) = \mathbb{I}[f(x^*) > \gamma] \quad (1)$$

141 where \mathbb{I} is the indicator function, $f(x^*)$ is a scoring function of the test observation x^* , and γ is an adjustable
 142 decision threshold. The success of the attack can be measured using traditional binary classification metrics and
 143 can be interpreted as a measure of privacy leakage from a model of the training data.

144 MIAs have become a mainstay tool to audit the privacy of tabular synthetic data as they represent a material
 145 privacy risk associated with the output of a model (Houssiau et al., 2022; Ward et al., 2025a). Indeed, a number
 146 of distance (Chen et al., 2020; Ward et al., 2024), density (Hilprecht et al., 2019; van Breugel et al., 2023), and
 147 likelihood-based (Stadler et al., 2022; Meeus et al., 2024; Ward et al., 2025b) attacks have been proposed under
 148 a wide variety of threat models. However, these existing methods all attack over the tabular feature space \mathcal{X}^d .
 149 LLM-based generators introduce a fundamentally different vulnerability: they generate in an intermediate string
 150 representation space before parsing to a tabular format. This autoregressive string generation process creates a
 151 novel attack surface that bypasses traditional feature-space-targeting attacks.

152 3 ATTACKING STRINGS OF DIGITS

153

154

155 LLM-based tabular data generators operate on records encoded as sequences of characters, often representing
 156 numeric values in fixed formats. From a privacy perspective, these numeric strings constitute a **distinct threat**
 157 **surface**: even minor changes at the character level can correspond to meaningful alterations in sensitive infor-
 158 mation, such as financial amounts, medical measurements, or personally identifiable metrics. Unlike free-form
 159 text, where approximate matches may be semantically ambiguous, we hypothesize that numeric strings are rigid
 160 and highly informative, making them particularly vulnerable to memorization by LLMs and thus adversarial
 161 signal.

162 In this section, we focus on **attacks that exploit these implied strings of digits**. We show that by measuring
 163 string similarity between synthetic outputs and potential training records, an adversary can infer membership
 164 without access to the model internals, queries, or auxiliary data. This approach exposes a realistic privacy risk
 165 inherent in current LLM-based tabular data generation pipelines, emphasizing the need to examine character-
 166 level leakage beyond traditional feature-space analyses.

167 3.1 THREAT MODEL

168

169 In this work, we adopt a conservative No-box threat model (Houssiau et al., 2022; Ward et al., 2025a), in which
 170 the adversary has access only to the synthetic data S produced by the model. We make no assumptions that the

adversary has knowledge of the model implementation and hyperparameters, query access, or even a reference dataset for constructing calibrated attacks. This threat model is motivated by two key considerations:

1. **Realism.** No-box threat models represent realistic privacy scenarios for tabular synthetic data practitioners. A common scenario involves a user training a tabular data generator on private datasets and sharing the generated data with public or private parties without disclosing implementation details. From the adversary’s perspective, this synthetic data may be maliciously acquired or discovered on open data sharing platforms, with no additional context about the generation process.
2. **Difficulty.** No-box threat models are recognized as particularly challenging for constructing effective attacks due to the severe limitations placed on the adversary [Chen et al. \(2020\)](#); [Houssiau et al. \(2022\)](#). Specifically, the adversary cannot analyze the model’s loss function, construct shadow models, or query the model directly. Successful attacks under these restrictive conditions therefore highlight fundamental vulnerabilities in these synthetic data generation methods.

3.2 LEVENSHEIN ATTACK

String similarity is a well-established concept, with Levenshtein edit distance providing a fundamental measure of character-level overlap ([Levenshtein, 1966](#); [Navarro, 2001](#); [Yujian & Bo, 2007](#)). Recent work has applied normalized variants of this distance to quantify approximate memorization in LLMs, showing that models can reproduce training examples with minor variations in natural language text ([Ippolito et al., 2023](#); [Shilov et al., 2025](#)). However, these studies stop short of embedding string similarity into an adversarial framework, as open-ended natural language provides many valid paraphrases and thus weak signals for membership inference.

In contrast, **the string representations of numeric values in tabular data are highly constrained**: numbers, delimiters, and fixed column formats mean that even small character-level deviations correspond to meaningful differences. This makes edit distance a far more reliable signal of memorization in our setting.

Motivated by this observation, we design an MIA based on the Distance of Closest Record ([Chen et al., 2020](#)), instantiated with Levenshtein distance. Formally, for each test observation and x^* and synthetic set S , we extract the ordered numeric and categorical values for each observation and encode them as strings (See Figure 1): x_{str}^* and S_{str} . The LevAtt score for Equation 1 is then:

$$f_{\text{LevAtt}}(x_{\text{str}}^*, S_{\text{str}}) = - \min_{s \in S_{\text{str}}} \ell(x_{\text{str}}^*, s), \quad (2)$$

where ℓ denotes the Levenshtein edit distance. Lower distances (i.e., higher scores) indicate closer matches and stronger evidence of membership. Because LevAtt requires only synthetic outputs, it applies directly to our No-box threat model highlighting a realistic privacy risk in LLM-based tabular data generation.

4 EXPERIMENTS

We evaluate the privacy leakage of digit memorization for in-context learning (ICL) and supervised fine-tuning (SFT) based LLM-based tabular generators. Here, we use differing experimental designs to correspond to their popular use-case settings. For clarity, we organize the section into three subsections: experimental design details for both regimes and then separate subsections for results. We include the full experimental and implementation details in Appendix: [A](#).

4.1 EXPERIMENTAL DESIGN

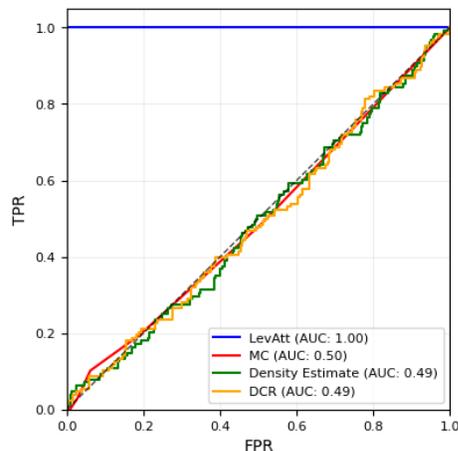
4.1.1 ICL APPROACHES

Replicating the design of [Byun et al. \(2025\)](#), we evaluate ICL approaches on the OpenML CTR23 benchmark ([Fischer et al., 2023](#)), consisting of 35 real-world classification tables with 500–100,000 rows and up to 5,000 features, including both numerical and categorical attributes. We partition each dataset into 80/20 training and test sets. To simulate the low-data regime these methods are commonly used for, we subsample 32, 64, 128 training rows without replacement. These sampled rows are provided as exemplars to the ICL models.

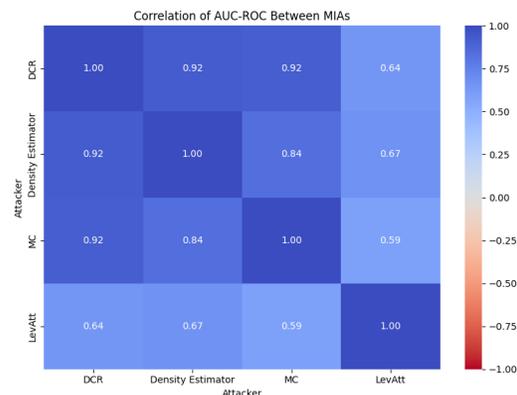
We consider three ICL models: LLaMA 3.3 70B ([Meta AI, 2024](#)) an open-source foundation model, GPT-4o-mini ([OpenAI, 2024](#)) a close-source foundation model, and TabPFN-v2 [Hollmann et al. \(2025a\)](#) a transformer-based model trained on tabular data due to their wide use and availability. For each sampled training subset, we generate an equal amount of synthetic data. Evaluation is conducted on all training rows plus an equal holdout partition of the test set. We evaluate privacy leakage using the following No-box MIAs using the

Table 1: LevAtt performance against ICL models. Mean (STD) values are reported for all datasets, training sizes, and seeds (Overall) and for the top 20 runs for each model with the highest AUC-ROC (Top 20). GPT-4o-mini shows relatively little privacy leakage, whereas LLama-3.3-70b and TabPFN-V2 show significant privacy failure- especially amongst their worst datasets.

Generator	AUC-ROC		TPR@FPR=0		TPR@FPR=0.1	
	Overall	Top 20	Overall	Top 20	Overall	Top 20
LLama-3.3-70b	0.63 (0.12)	0.91 (0.03)	0.08 (0.20)	0.64 (0.27)	0.28 (0.21)	0.81 (0.09)
TabPFN-V2	0.58 (0.05)	0.91 (0.07)	0.04 (0.01)	0.57 (0.41)	0.19 (0.06)	0.94 (0.14)
GPT-4o-mini	0.54 (0.05)	0.60 (0.09)	0.01 (0.01)	0.01 (0.03)	0.13 (0.05)	0.27 (0.09)



(a) ROC plot for various No-box MIAs against TabPFN-V2 with 128 in-context samples from the MoneyBall dataset. LevAtt (blue) is able to achieve perfect classification for all in-context samples whereas MIAs that target the feature space of tabular data fail to capture the privacy leakage.



(b) Correlation plot for No-box MIA AUC-ROC across the ICL experiment. While the feature-space targeting DCR, Density Estimate, and MC are nearly perfectly correlated, LevAtt is much less correlated. This highlights that while privacy leakage over tabular string representations and the feature space are related, LevAtt finds unique adversarial advantage.

Synth-MIA package (Ward et al., 2025a): LevAtt, Euclidean Distance to Closest Record (DCR) (Chen et al., 2020), a Monte Carlo density estimation (MC) method (Hilprecht et al., 2019), and a kernel density estimation method (Houssiau et al., 2022; van Breugel et al., 2023). We report MIA success using mean AUC-ROC and TPR at fixed FPR thresholds (Carlini et al., 2021a) to capture potential information leakage. We repeat this experimentation across 3 seeded runs.

4.1.2 SFT APPROACHES

For SFT-approaches, we benchmark the original SFT-based tabular generation method GREAT Borisov et al. (2023) and RealTabFormer (Solatorio & Dupriez, 2023), which reports improved privacy performance due to enforcing a minimum Euclidean Distance to Closest to Record distribution in its training. As both of these methods use GPT-2 (Radford et al., 2019) as a base model, we modify RealTabFormer to accept more modern, larger foundation models: LLaMA 3.2 (1B, 3B) (Grattafiori et al., 2024), Qwen2.5-3B (Qwen et al., 2025), and Mistral v0.3 7B Jiang et al. (2023). Additionally, we introduce as a control CT-GAN and TVAE Xu et al. (2019), conventional deep learning-based generators. Training follows default hyperparameters from the original GREAT and RealTabFormer implementations while CTGAN and TVAE are implemented through Synthcity Qian et al. (2023).

Experiments are conducted on five tabular datasets: CASP, Abalone, Diabetes, CA-Housing, and Faults, selected for their common use in synthetic tabular data benchmarking and containing numeric data. Similarly to ICL, we create 80/20 train-test partitions and following common synthetic tabular data benchmarking (Zhang et al., 2024), synthetic datasets equal in size to the original training sets are generated post-training. For privacy evaluation, up to 1,000 training and holdout samples are partitioned as test sets, and the same MIAs are applied. Experiments are repeated across three seeds.

Table 2: Mean (STD) LevAtt performance for each model and dataset across seeds. RealTabFormer experiences significant privacy leakage and LevAtt finds some signal for most LLM-based models. However, LevAtt identifies no leakage in conventional deep learning-based methods CT-GAN and TVAE as they do not generate strings.

Model	LevAtt Metric	CA-Housing	CASP	Abalone	Diabetes	Faults
RealTabFormer	AUC-ROC	0.70 (0.11)	0.72 (0.12)	0.60 (0.08)	0.66 (0.04)	0.61 (0.12)
	TPR@FPR=0.1	0.34 (0.23)	0.37 (0.11)	0.21 (0.05)	0.25 (0.05)	0.32 (0.10)
LLaMA 3.2-1B	AUC-ROC	0.68 (0.22)	0.52 (0.15)	0.50 (0.03)	0.56 (0.03)	0.58 (0.15)
	TPR@FPR=0.1	0.21 (0.21)	0.10 (0.06)	0.10 (0.02)	0.21 (0.02)	0.16 (0.05)
LLaMA 3.2-3B	AUC-ROC	0.63 (0.06)	0.50 (0.07)	0.62 (0.04)	0.61 (0.04)	0.58 (0.07)
	TPR@FPR=0.1	0.24 (0.09)	0.09 (0.03)	0.22 (0.11)	0.15 (0.11)	0.16 (0.03)
Qwen 2.5-3B	AUC-ROC	0.61 (0.08)	0.54 (0.02)	0.52 (0.04)	0.64 (0.04)	0.56 (0.01)
	TPR@FPR=0.1	0.25 (0.05)	0.12 (0.02)	0.15 (0.07)	0.28 (0.07)	0.11 (0.02)
Mistral-7B	AUC-ROC	0.67 (0.13)	0.51 (0.01)	0.51 (0.03)	0.53 (0.06)	0.52 (0.01)
	TPR@FPR=0.1	0.23 (0.04)	0.11 (0.02)	0.10 (0.05)	0.10 (0.03)	0.10 (0.01)
GREAT	AUC-ROC	0.66 (0.09)	0.55 (0.04)	0.52 (0.03)	0.50 (0.05)	0.54 (0.07)
	TPR@FPR=0.1	0.33 (0.10)	0.14 (0.03)	0.11 (0.02)	0.10 (0.03)	0.14 (0.02)
CT-GAN	AUC-ROC	0.48 (0.03)	0.51 (0.04)	0.48 (0.03)	0.47 (0.04)	0.50 (0.02)
	TPR@FPR=0.1	0.09 (0.01)	0.11 (0.01)	0.12 (0.01)	0.10 (0.01)	0.08 (0.01)
TVAE	AUC-ROC	0.47 (0.03)	0.49 (0.02)	0.51 (0.03)	0.53 (0.02)	0.51 (0.03)
	TPR@FPR=0.1	0.08 (0.01)	0.10 (0.01)	0.11 (0.01)	0.10 (0.01)	0.10 (0.02)

4.2 RESULTS: ICL METHODS

We summarize our key findings for our ICL experiments as follows:

LevAtt Identifies Substantial Privacy Leakage in ICL Tabular Synthetic Data Approaches. We report the Mean (STD) for LevAtt’s AUC-ROC and $\text{TPR@FPR} = \{0.0, 0.1\}$ over all datasets, training sizes, and seeds in Table 1. As there is considerable variability in the performance of LevAtt across runs and MIAs are often concerned with the worst-case performance of a model, we also report corresponding values for the top 20 runs for each model by AUC-ROC. Here, LLaMA 3.3-70b and TabPFN-V2 experience substantial privacy leakage based on memorizing approximate patterns in the training data with mean $\text{TPR@FPR}=0$ values of 0.64 and 0.57 respectively. In some cases, we even find that LevAtt is a perfect classifier of training membership (see Figure 1a).

Privacy Leakage Likely Scales with Model Size. Prior work (Carlini et al., 2021b; 2023) has shown that as the parameter count of LLMs increases, so does their capacity for memorization. While the exact sizes of GPT-4o-mini and TabPFN-V2 are not publicly disclosed, benchmarking suggests GPT-4o-mini performs similarly to models with roughly 7B parameters, and TabPFN-V2 is small enough to run on a single GPU. In contrast, the 70B-parameter LLaMA 3.3 exhibits markedly higher privacy leakage on our benchmark, consistent with the expected scaling behavior.

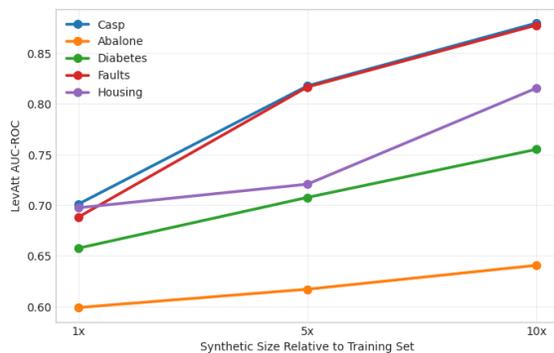
LevAtt Captures a Unique String-Similarity Signal Relative to Other MIAs. Figure 1b shows the correlation of AUC-ROC for LevAtt with other No-box threat model MIAs targeting the feature space of synthetic data—DCR, Density Estimate, and MC—across all datasets, training sizes, seeds, and models. Feature-space attacks are nearly perfectly correlated with each other, as expected. In contrast, LevAtt exhibits lower correlation, suggesting that it uncovers privacy leakage signals not captured by traditional feature-space attacks.

4.3 RESULTS: SFT METHODS

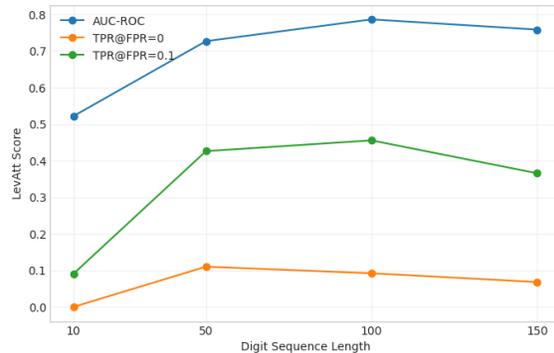
We summarize our key findings for our SFT experiments as follows:

While Less Sizeable, SFT-Approaches Similarly Exhibit Privacy Leakage. We report the Mean (STD) of LevAtt’s AUC-ROC and $\text{TPR@FPR}=0.1$ across datasets and seeds for each model Table in 2. Overall, SFT-based approaches generally exhibit lower privacy leakage than ICL counterparts. Surprisingly, RealTabFormer often shows the highest LevAtt scores, despite using privacy-aware training and sharing a GPT-2 base with GREAT, which typically leaks little. LevAtt also detects leakage in LLaMA-3.2 (1B and 3B) and Qwen-2.5-3B. In contrast, LevAtt is completely ineffective against CT-GAN and TVAE, as these models generate tabular data at the level of complete observations rather than token by token, inherently protecting against this type of attack.

LevAtt Success is Related to the Synthetic Dataset Size. To examine how synthetic data volume affects privacy leakage, we conducted an ablation study where RealTabFormer generated synthetic datasets at 1x, 5x, and 10x the size of each original training dataset. The results in Figure 2a reveal a clear monotonic relationship: LevAtt AUC-ROC scores increase consistently with larger synthetic datasets, with the Faults dataset showing



(a) LevAtt AUC-ROC for various datasets generated by RealTabFormer with increasing synthetic dataset sizes relative to the training set.



(b) LevAtt performance on RealTabFormer at various training digit sequence lengths.

a particularly striking 20% improvement in attack success. This pattern demonstrates that generating more synthetic data systematically increases the probability of reproducing memorized training examples, thereby amplifying privacy risks.

LevAtt Success is Related to the Sequence Length of Training Data. As an additional ablation, we investigate the impact of digit sequence length in the training data. Here, we generate both training and holdout datasets from the same multivariate standard Gaussian distribution each with 10,000 observations. To control for sequence length, we select an increasing number of columns for the training data, each containing exactly 10 digits. We then train RealTabFormer on datasets with progressively longer digit sequences to assess its effect on model behavior. We report LevAtt’s performance in Figure 2b where we see that increasing the count of digits or sequence length in the training data correlates with increased memorization in RealTabFormer. This further replicates findings from Carlini et al. (2021b) that found training sequence length factors into increased memorization in LLMs.

5 DEFENSES AGAINST LEVENSHTAIN ATTACK

The efficacy of LevAtt against LLM-based tabular models motivates us to explore methodologies that can defend generated synthetic data. Here, recognizing that LevAtt gains adversarial advantage from replicated patterns in strings of digits we devise two strategies based on introducing controlled “noise” into string sequences: Digit Modifier (DM) a post-processing method that operates independently of the generative model and is applied after the synthetic data have been produced, and Tendency-based Logit Processor (TLP) a method compatible with any open-source LLM that strategically perturbs digits at sample time.

5.1 DIGIT MODIFIER AND TENDENCY-BASED LOGIT PROCESSOR

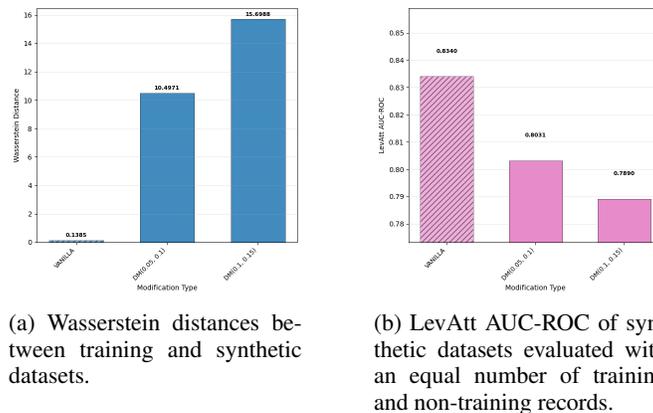
Digit Modifier

We propose a Digit Modifier (DM), a method for injecting controlled randomness into tabular data by perturbing numerical digits. Motivated by bit-flipping techniques for adding noise to binary representations of relational databases Agrawal & Kiernan (2002); Alfagi et al. (2016), DM operates by replacing tokens in a record’s tokenized sequence with alternatives sampled from a noise distribution. When tokens represent numerical values, these substitutions correspondingly alter the digits of the underlying sequences, yielding perturbed records. In this sense, DM can be viewed as a principled method for randomly replacing digits within a record.

To balance data fidelity and protection, we parameterize the mechanism as $DM(p_{min}, p_{max})$ with $0 \leq p_{min} < p_{max} \leq 1$. For a numerical column \mathbf{X}_i and entry $x_{k,i}$, a probability function $g : x_{k,i} \times \mathbf{X}_i \rightarrow [p_{min}, p_{max}]$ assigns digit-flipping probabilities such that larger-magnitude entries receive higher perturbation probabilities. Each digit of $x_{k,i}$ is then independently flipped according to $g(x_{k,i}, \mathbf{X}_i)$, while smaller-magnitude values—being more sensitive—undergo smaller changes to preserve fidelity. The design of the probability function g and full algorithm (Alg 1) are described in Appendix B.1.

Tendency-Based Logit Processor

We additionally propose Tendency-based Logit Processor (TLP), a mechanism for injecting controlled noise into synthetic data by perturbing the generator’s logits at inference time. The $TLP(t)$ selectively amplifies lower-



(a) Wasserstein distances between training and synthetic datasets.

(b) LevAtt AUC-ROC of synthetic datasets evaluated with an equal number of training and non-training records.

Figure 3: Privacy-fidelity comparison of DM on RealTabFormer synthetic data from the simulated Gaussian dataset (Section 5.2). VANILLA corresponds to plain sampling without protection. Panel (a) reports Wasserstein distances; panel (b) shows LevAtt AUC-ROC. While DM is able to induce reductions in privacy leakage, the resulting synthetic data are of low fidelity.

valued logits while suppressing higher-valued ones, making the generator more likely to select tokens that were originally less probable. The strength of this effect is controlled by the tendency parameter t : higher values of t induce stronger curvature, increasing the randomness of the generated sequence. In this way, TLP(t) acts as a principled method for introducing variability into synthetic outputs while still preserving the generator’s learned distribution.

Formally, TLP(t) first maps the generator’s logits $l = (l_1, l_2, \dots, l_k)$ into the range $[0, 1]$ using a shifted min-max scaler S_l . It then applies a monotone increasing, concave curving function $f_t : [0, 1] \rightarrow [0, 1]$, parameterized by t and satisfying $f_t(0) = 0$, to each scaled logit. Finally, the processed logits are transformed back to their original scale using the inverse scaler S_l^{-1} .

The design of f_t is central to TLP(t). Monotonicity ensures that the relative order of logits is preserved, so high-probability tokens remain more likely than lower-probability ones, allowing controlled noise injection without overwhelming the generator’s learned signal. Concavity, combined with $f_t(0) = 0$, guarantees that lower logits are curved upward, increasing their chance of being sampled. Appendix B.2 provides full details of the scalers S_l and S_l^{-1} , the conditions for valid curving, the specific f_t used in our experiments, and a full algorithm description (Alg 2).

5.2 EVALUATING DIGIT MODIFIER AND TENDENCY-BASED LOGIT PROCESSOR

To evaluate the effectiveness of DM and TLP, we use RealTabFormer, the SFT model exhibiting the highest degree of memorization, and measure how much privacy improvement each method provides against LevAtt while preserving synthetic data fidelity. Following the experimental designs of Section 4.3, we generate a 20-column training and holdout set from a Multivariate Gaussian $N(300, 5)$, and we use the CASP dataset at varying synthetic set sizes to induce privacy leakage. We then apply DM and TLP across different scaler and tendency parameter levels respectively, and assess performance in terms of privacy (LevAtt AUC-ROC), fidelity (Wasserstein Distance and Maximum Mean Discrepancy between training and synthetic datasets), and utility (RMSE of XGBoost models trained on synthetic data and evaluated on real holdout data) (Qian et al., 2023).

Overall, we find that while effective in reducing privacy leakage, DM suffers large fidelity costs (See Figure 3). DM is particularly challenged by low-bandwidth datasets or columns with restricted ranges, where naive modifications can shift values into low-density regions and result in significant fidelity loss. On the other hand, TLP, with a tuned tendency parameter t , can controllably reduce attack efficacy while effectively preserving the fidelity between the processed synthetic data and the training data. In Figure 4, in both highly unprivate Simulation and CASP datasets TLP reduces LevAtt’s AUC-ROC to 55% with effectively no penalty in the Maximum Mean Discrepancy of TLP-generated data.

We include additional experiments for TLP in Appendix B.3. There, we show that TLP is able to also control for TPR@FPR values with minimal fidelity loss (see Figure 6) as well as preserve utility in machine learning models (See Figure 7). Overall, TLP offers a competitive defense against string-based attacks.

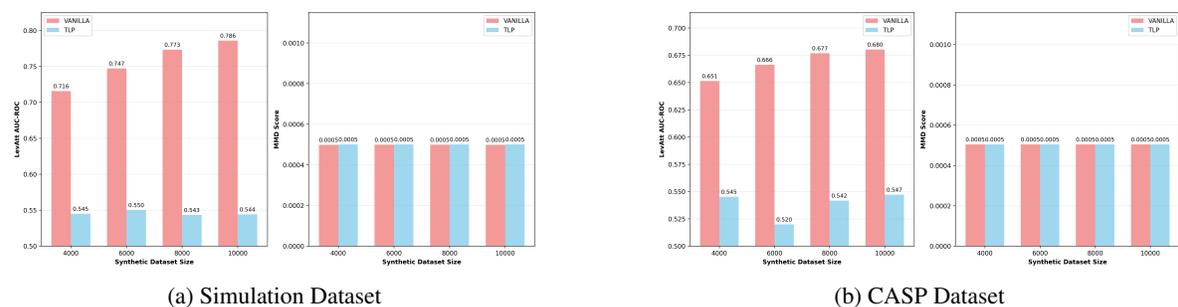


Figure 4: Privacy-fidelity trade-off of TLP in RealTabFormer synthetic data from the simulation (Section 5.2) and CASP dataset. For both simulation and CASP sets, TLP reduces LevAtt AUC-ROC to 55% across synthetic dataset sizes, improving privacy over the vanilla inference, while the Maximum Mean Discrepancy (MMD) remains nearly unchanged, demonstrating that TLP preserves fidelity while mitigating privacy leakage.

6 DISCUSSION

LevAtt reveals that LLM-based tabular data generation is uniquely unsafe relative to conventional deep learning approaches. While being a simple string-similarity attack operating under an extremely restrictive threat model, LevAtt uncovered that state-of-the-art models could catastrophically leak training membership by copying sequential patterns of digits and text from training examples. This phenomenon was further demonstrated in SFT-generators, revealing that the base implementation of RealTabFormer—a method recognized for its privacy-preserving capabilities—was susceptible to significant privacy leakage. At the same time, LevAtt was found to be substantially less correlated with feature space-oriented MIAs, and conventional tabular data generators such as CT-GAN and TVAE proved resistant to our string-based attack. These findings suggest that tabular models relying on autoregressive token-based generation expose a distinct attack surface that prior deep-learning architectures do not share.

While highly deployable and powerful, LevAtt can be defeated. In this work, we proposed an intuitive post-hoc defense called Digit Modifier (DM), which alters digits after generation. However, we found that DM failed to preserve the fidelity of the synthetic dataset. This limitation motivated the development of TLP, which can be appended to any open-source LLM through the Hugging Face API. TLP can effectively control for privacy leakage from LevAtt by smoothing the logits of digits with disproportionately high probabilities. We found that this perturbation did not substantially alter the statistical fidelity of the resulting synthetic data.

While TLP provides an effective defense against LevAtt, the need for such inference-time mitigation raises a broader question: *What are these models exactly learning?* A prevailing assumption is that LLM-based tabular generators approximate the joint distribution of the training set T through sequential string modeling, similar to conventional generative approaches. However, our findings suggest that in some cases these models behave more like perturbation mechanisms, producing outputs that closely resemble training examples with only minor modifications. This behavior naturally yields high fidelity and downstream utility but does so precisely because of its proximity to the training data—thereby exposing the system to privacy leakage. Understanding when LLMs genuinely learn tabular distributions versus when they rely on approximate memorization remains an important direction for future work.

7 CONCLUSION

In this paper we introduce LevAtt, a No-box threat model MIA that exposes substantial privacy risk for in-context learning and supervised-finetuned LLM-based tabular data generators. LevAtt shows that LLMs are vulnerable to memorization from the structured, often duplicated patterns of tabular data. By attacking the string encodings of autoregressively generated tabular data, LevAtt finds unique adversarial signal compared to existing methods. While less restrictive threat models would likely lead to a better attack, we believe No-box carries a powerful message: an attack with the minimal assumptions reasonably possible for an adversary can perfectly classify training membership in state-of-the-art generators. Lastly propose two defenses against LevAtt, showing that Tendency-Based Logit Processor can effectively defeat LevAtt with minimal loss in synthetic data fidelity. Future research directions could involve developing even more powerful attacks under less restrictive threat models, finding more efficient and provable defenses, and studying mechanistically how LLMs represent tabular distributions.

8 STATEMENT OF ETHICS

The potential for adversaries to determine whether an individual’s data was included in the original dataset presents significant privacy risks, especially in fields such as healthcare, finance, and social sciences, where sensitive personal information is commonly used. Synthetic data that fails to sufficiently mask membership information could inadvertently enable re-identification. Although this work introduces a method for assessing such risks, its primary objective is to empower researchers and practitioners to perform more rigorous privacy evaluations before deploying synthetic datasets. We emphasize that adversarial approaches are essential for advancing the development of robust privacy-preserving systems.

9 REPRODUCIBILITY STATEMENT

Our code is available through the link provided on the abstract page. The main paper explains all algorithms in detail and provides the dataset and simulation descriptions, while model hyperparameters are reported in the appendix.

REFERENCES

- Rakesh Agrawal and Jerry Kiernan. Watermarking relational databases. In *VLDB’02: Proceedings of the 28th International Conference on Very Large Databases*, pp. 155–166. Elsevier, 2002.
- Abd S Alfagi, A Abd Manaf, B Hamida, S Khan, and Ali A Elrowayati. Survey on relational database watermarking techniques. *ARPN-JEAS*, 11:422–423, 2016.
- Ankur Ankan and Abinash Panda. pgmpy: Probabilistic graphical models using python. In *Proceedings of the Python in Science Conference*, SciPy. SciPy, 2015. doi: 10.25080/majora-7b98e3ed-001. URL <http://dx.doi.org/10.25080/Majora-7b98e3ed-001>.
- Samuel A. Assefa, Danial Dervovic, Mahmoud Mahfouz, Robert E. Tillman, Prashant Reddy, and Manuela Veloso. Generating synthetic data in finance: opportunities, challenges and pitfalls. In *Proceedings of the First ACM International Conference on AI in Finance*, ICAIF ’20, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450375849. doi: 10.1145/3383455.3422554. URL <https://doi.org/10.1145/3383455.3422554>.
- Vadim Borisov, Kathrin Seßler, Tobias Leemann, Martin Pawelczyk, and Gjergji Kasneci. Language models are realistic tabular data generators, 2023. URL <https://arxiv.org/abs/2210.06280>.
- Jessup Byun, Xiaofeng Lin, Joshua Ward, and Guang Cheng. Risk in context: Benchmarking privacy leakage of foundation models in synthetic tabular data generation, 2025. URL <https://arxiv.org/abs/2507.17066>.
- Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, A. Terzis, and Florian Tramèr. Membership inference attacks from first principles, 2021a. URL <https://api.semanticscholar.org/CorpusID:244920593>.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, Alina Oprea, and Colin Raffel. Extracting training data from large language models, 2021b. URL <https://arxiv.org/abs/2012.07805>.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. Quantifying memorization across neural language models, 2023. URL <https://arxiv.org/abs/2202.07646>.
- Dingfan Chen, Ning Yu, Yang Zhang, and Mario Fritz. Gan-leaks: A taxonomy of membership inference attacks against generative models. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, CCS ’20, pp. 343–362, Virtual Event, USA, October 2020. ACM. doi: 10.1145/3372297.3417238. URL <http://dx.doi.org/10.1145/3372297.3417238>.
- Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V Le, Sergey Levine, and Yi Ma. SFT memorizes, RL generalizes: A comparative study of foundation model post-training. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=dYur3yabMj>.

- 570 Conor Durkan, Artur Bekasov, Iain Murray, and George Papamakarios. Neural spline flows. In *Advances*
571 *in Neural Information Processing Systems*, volume 32, pp. 7627–7638, Vancouver, Canada, 2019. Curran
572 Associates Inc.
- 573 Sebastian Felix Fischer, Matthias Feurer, and Bernd Bischl. OpenML-CTR23 – a curated tabular regression
574 benchmarking suite. In *AutoML Conference 2023 (Workshop)*, 2023. URL <https://openreview.net/forum?id=HebAOoMm94>.
- 575
576
577 Brendan Flanagan, Rwitajit Majumdar, and Hiroaki Ogata. Fine grain synthetic educational data: Challenges
578 and limitations of collaborative learning analytics. *IEEE Access*, 10:26230–26241, 03 2022. doi: 10.1109/
579 ACCESS.2022.3156073.
- 580
581 Joao Fonseca and Fernando Bação. Tabular and latent space synthetic data generation: a literature review.
582 *Journal of Big Data*, 10, 07 2023. doi: 10.1186/s40537-023-00792-7.
- 583
584 Mauro Giuffré and Dennis L. Shung. Harnessing the power of synthetic data in healthcare: innovation, appli-
585 cation, and privacy. *NPJ Digital Medicine*, 6, 2023. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:263802405)
586 [CorpusID:263802405](https://api.semanticscholar.org/CorpusID:263802405).
- 587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
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
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

- 627 Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia
628 Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine,
629 Delia David, Devi Parikh, Diana Liskovich, Didem Foss, DingKang Wang, Duc Le, Dustin Holland, Ed-
630 ward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan
631 Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Fil-
632 ippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez,
633 Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang,
634 Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun
635 Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj,
636 Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski,
637 James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny
638 Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill,
639 Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan
640 Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena,
641 Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender
642 A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt,
643 Madian Khabza, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias
644 Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal
645 Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike
646 Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam,
647 Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta,
648 Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar,
649 Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre
650 Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao,
651 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,
652 Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin
653 Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru
654 Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun
655 Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang,
656 Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen
657 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho,
658 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best,
659 Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked,
660 Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad
661 Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang,
662 Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv
663 Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin
664 Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef
665 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024.
666 URL <https://arxiv.org/abs/2407.21783>.
- 667 Benjamin Hilprecht, Martin Härterich, and Daniel Bernau. Monte carlo and reconstruction membership infer-
668 ence attacks against generative models. *Proceedings on Privacy Enhancing Technologies*, 2019:232 – 249,
669 2019. URL <https://api.semanticscholar.org/CorpusID:199546273>.
- 670 Noah Hollmann, Samuel Müller, Lennart Purucker, Arjun Krishnakumar, Max Körfer, Shi Bin Hoo, Robin Tibor
671 Schirrmeister, and Frank Hutter. Accurate predictions on small data with a tabular foundation model. *Nature*,
672 637(8045):319–326, 2025a.
- 673 Noah Hollmann, Samuel Müller, Lennart Purucker, Arjun Krishnakumar, Max Körfer, Shi Hoo, Robin
674 Schirrmeister, and Frank Hutter. Accurate predictions on small data with a tabular foundation model. *Nature*,
675 637:319–326, 01 2025b. doi: 10.1038/s41586-024-08328-6.
- 676 Florimond Houssiau, James Jordon, Samuel N Cohen, Owen Daniel, Andrew Elliott, James Geddes, Callum
677 Mole, Camila Rangel-Smith, and Lukasz Szpruch. Tapas: a toolbox for adversarial privacy auditing of
678 synthetic data, 2022.
- 679 Daphne Ippolito, Florian Tramer, Milad Nasr, Chiyuan Zhang, Matthew Jagielski, Katherine Lee, Christopher
680 Choquette Choo, and Nicholas Carlini. Preventing generation of verbatim memorization in language models
681 gives a false sense of privacy. In C. Maria Keet, Hung-Yi Lee, and Sina Zarrieß (eds.), *Proceedings of
682 the 16th International Natural Language Generation Conference*, pp. 28–53, Prague, Czechia, September
683 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.inlg-main.3. URL [https://
684 aclanthology.org/2023.inlg-main.3/](https://aclanthology.org/2023.inlg-main.3/).

- 684 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego
685 de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L elio Renard Lavaud,
686 Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth ee Lacroix, and
687 William El Sayed. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>.
- 688
689 Nikhil Kandpal, Eric Wallace, and Colin Raffel. Deduplicating training data mitigates privacy risks in language
690 models. *ArXiv*, abs/2202.06539, 2022. URL [https://api.semanticscholar.org/CorpusID:
691 246823128](https://api.semanticscholar.org/CorpusID:246823128).
- 692 Jinhee Kim, Taesung Kim, and Jaegul Choo. EPIC: Effective prompting for imbalanced-class data synthesis
693 in tabular data classification via large language models. In *The Thirty-eighth Annual Conference on Neural
694 Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=d5cKDHCrFJ>.
- 695 Akim Kotelnikov, Dmitry Baranchuk, Ivan Rubachev, and Artem Babenko. Tabddpm: Modelling tabular data
696 with diffusion models, 2022.
- 697
698 Vladimir I Levenshtein. Binary codes capable of correcting deletions, insertions and reversals. *Soviet Physics
699 Doklady*, 10:707, February 1966.
- 700 Junwei Ma, Apoorv Dankar, George Stein, Guangwei Yu, and Anthony Caterini. Tabpfgn – tabular data
701 generation with tabpfn, 2024. URL <https://arxiv.org/abs/2406.05216>.
- 702
703 Ryan McKenna, Brett Mullins, Daniel Sheldon, and Gerome Miklau. Aim: an adaptive and iterative mechanism
704 for differentially private synthetic data. *Proc. VLDB Endow.*, 15(11):2599–2612, July 2022. ISSN 2150-8097.
705 doi: 10.14778/3551793.3551817. URL <https://doi.org/10.14778/3551793.3551817>.
- 706 Matthieu Meeus, Florent Guepin, Ana-Maria Creţu, and Yves-Alexandre de Montjoye. *Achilles’ Heels:
707 Vulnerable Record Identification in Synthetic Data Publishing*, pp. 380–399. Springer Nature Switzer-
708 land, Cham, Switzerland, 2024. ISBN 9783031514760. doi: 10.1007/978-3-031-51476-0_19. URL
709 http://dx.doi.org/10.1007/978-3-031-51476-0_19.
- 710 Meta AI. LLaMA-3.3 70B instruct model. [https://huggingface.co/meta-llama/Llama-3.
711 3-70B-Instruct](https://huggingface.co/meta-llama/Llama-3.3-70B-Instruct), 2024. Released December 6, 2024; accessed 2025-06-13.
- 712
713 Gonzalo Navarro. A guided tour to approximate string matching. *ACM Computing Surveys*, 33(1):31–88, 2001.
- 714 OpenAI. GPT-4o Mini model in chat completions api. [https://platform.openai.com/docs/
715 models/gpt-4o-mini](https://platform.openai.com/docs/models/gpt-4o-mini), 2024. Released July 18, 2024; accessed 2025-06-13.
- 716
717 Michael Platzer and Thomas Reutterer. Holdout-based empirical assessment of mixed-type synthetic data.
718 *Frontiers in big Data*, 4:679939, 2021.
- 719
720 Zhaozhi Qian, Bogdan-Constantin Cebere, and Mihaela van der Schaar. Synthcity: facilitating innovative use
721 cases of synthetic data in different data modalities, 2023. URL [https://arxiv.org/abs/2301.
722 07573](https://arxiv.org/abs/2301.07573).
- 723 Qwen, ., An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Day-
724 iheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi
725 Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng
726 Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xu-
727 ancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and
728 Zihan Qiu. Qwen2.5 technical report, 2025. URL <https://arxiv.org/abs/2412.15115>.
- 729 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are
730 unsupervised multitask learners. 2019.
- 731
732 Nabeel Seedat, Nicolas Huynh, Boris van Breugel, and Mihaela van der Schaar. Curated LLM: Synergy of
733 LLMs and data curation for tabular augmentation in low-data regimes. In *Forty-first International Conference
734 on Machine Learning*, 2024.
- 735 Igor Shilov, Matthieu Meeus, and Yves-Alexandre de Montjoye. The mosaic memory of large language models,
736 2025. URL <https://arxiv.org/abs/2405.15523>.
- 737
738 R. Shokri, M. Stronati, C. Song, and V. Shmatikov. Membership inference attacks against machine
739 learning models. In *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 3–18, Los Alamitos,
740 CA, USA, may 2017. IEEE Computer Society. doi: 10.1109/SP.2017.41. URL [https://doi.
ieeecomputersociety.org/10.1109/SP.2017.41](https://doi.ieeecomputersociety.org/10.1109/SP.2017.41).

- 741 Aivin V Solatorio and Olivier Dupriez. Realtabformer: Generating realistic relational and tabular data using
742 transformers. *arXiv preprint arXiv:2302.02041*, 2023.
- 743
744 Theresa Stadler, Bristena Oprisanu, and Carmela Troncoso. Synthetic data – anonymisation groundhog day.
745 In *31st USENIX Security Symposium (USENIX Security 22)*, pp. 1451–1468, Boston, MA, August 2022.
746 USENIX Association. ISBN 978-1-939133-31-1. URL [https://www.usenix.org/conference/
747 usenixsecurity22/presentation/stadler](https://www.usenix.org/conference/usenixsecurity22/presentation/stadler).
- 748 Namjoon Suh, Xiaofeng Lin, Din-Yin Hsieh, Mehrdad Honarkhah, and Guang Cheng. Autodiff: combining
749 auto-encoder and diffusion model for tabular data synthesizing, 2023. URL [https://openreview.
750 net/forum?id=XhxOCXlXSh](https://openreview.net/forum?id=XhxOCXlXSh).
- 751 Vibeke Binz Vallevik, Aleksandar Babic, Serena E. Marshall, Severin Elvatun, Helga M.B. Brøgger, Sharmini
752 Alagaratnam, Bjørn Edwin, Narasimha R. Veeraragavan, Anne Kjersti Befring, and Jan F. Nygård. Can i
753 trust my fake data – a comprehensive quality assessment framework for synthetic tabular data in healthcare.
754 *International Journal of Medical Informatics*, 185:105413, 2024. ISSN 1386-5056. doi: [https://doi.org/
755 10.1016/j.ijmedinf.2024.105413](https://doi.org/10.1016/j.ijmedinf.2024.105413). URL [https://www.sciencedirect.com/science/article/
756 pii/S1386505624000765](https://www.sciencedirect.com/science/article/pii/S1386505624000765).
- 757
758 Boris van Breugel, Hao Sun, Zhaozhi Qian, and Mihaela van der Schaar. Membership inference attacks against
759 synthetic data through overfitting detection, 2023.
- 760 Yuxin Wang, Duanyu Feng, Yongfu Dai, Zhengyu Chen, Jimin Huang, Sophia Ananiadou, Qianqian Xie, and
761 Hao Wang. Harmonic: harnessing llms for tabular data synthesis and privacy protection. In *Proceedings
762 of the 38th International Conference on Neural Information Processing Systems, NIPS '24*, Red Hook, NY,
763 USA, 2025. Curran Associates Inc. ISBN 9798331314385.
- 764 Zhepeng Wang, Runxue Bao, Yawen Wu, Jackson Taylor, Cao Xiao, Feng Zheng, Weiwen Jiang, Shangqian
765 Gao, and Yanfu Zhang. Unlocking memorization in large language models with dynamic soft prompt-
766 ing. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Con-
767 ference on Empirical Methods in Natural Language Processing*, pp. 9782–9796, Miami, Florida, USA,
768 November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.546. URL
769 <https://aclanthology.org/2024.emnlp-main.546/>.
- 770 Joshua Ward, Chi-Hua Wang, and Guang Cheng. Data plagiarism index: Characterizing the privacy risk of
771 data-copying in tabular generative models, 2024. URL <https://arxiv.org/abs/2406.13012>.
- 772
773 Joshua Ward, Xiaofeng Lin, , Chi-Hua Wang, and Guang Cheng. Synth-mia: A testbed for auditing privacy
774 leakage in tabular data synthesis, 2025a. URL <https://arxiv.org/abs/2509.18014>.
- 775
776 Joshua Ward, Chi-Hua Wang, and Guang Cheng. Privacy auditing synthetic data release through local likelihood
777 attacks, 2025b. URL <https://arxiv.org/abs/2508.21146>.
- 778 David S. Watson, Kristin Blesch, Jan Kapar, and Marvin N. Wright. Adversarial random forests for density
779 estimation and generative modeling. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent (eds.),
780 *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of
781 *Proceedings of Machine Learning Research*, pp. 5357–5375, Valencia, Spain, 25–27 Apr 2023. PMLR. URL
782 <https://proceedings.mlr.press/v206/watson23a.html>.
- 783
784 Jinhong Wu, Konstantinos Plataniotis, Lucy Liu, Ehsan Amjadian, and Yuri Lawryshyn. Interpretation for
785 variational autoencoder used to generate financial synthetic tabular data. *Algorithms*, 16:121, 02 2023. doi:
786 10.3390/a16020121.
- 787
788 Lei Xu, Maria Skoularidou, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. Modeling tabular data using
789 conditional gan. In *Advances in Neural Information Processing Systems*, volume 32, pp. 7335–7345, Van-
790 couver, Canada, 2019. Curran Associates, Inc. URL [https://proceedings.neurips.cc/paper/
791 2019/hash/254ed7d2de3b23ab10936522dd547b78-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/254ed7d2de3b23ab10936522dd547b78-Abstract.html).
- 792
793 Jinsung Yoon, James Jordon, and Mihaela van der Schaar. PATE-GAN: Generating synthetic data with dif-
794 ferential privacy guarantees. In *International Conference on Learning Representations*, pp. 1–15, New
795 Orleans, LA, USA, May 2019. OpenReview.net. URL [https://openreview.net/forum?id=
796 S1zk9iRqF7](https://openreview.net/forum?id=S1zk9iRqF7).
- 797
798 Jinsung Yoon, Lydia N Drumright, and Mihaela Van Der Schaar. Anonymization through data synthesis using
799 generative adversarial networks (ads-gan). *IEEE journal of biomedical and health informatics*, 24(8):2378–
800 2388, 2020.

798 Li Yujian and Liu Bo. A normalized levenshtein distance metric. *IEEE Trans. Pattern Anal. Mach. Intell.*,
799 29(6):1091–1095, June 2007. ISSN 0162-8828. doi: 10.1109/TPAMI.2007.1078. URL <https://doi.org/10.1109/TPAMI.2007.1078>.

801 Hengrui Zhang, Jiani Zhang, Zhengyuan Shen, Balasubramaniam Srinivasan, Xiao Qin, Christos Faloutsos,
802 Huzefa Rangwala, and George Karypis. Mixed-type tabular data synthesis with score-based diffusion in
803 latent space. In *The Twelfth International Conference on Learning Representations*, pp. 4Ay23yeuz0, Vienna,
804 Austria, May 2024. OpenReview.net. URL <https://openreview.net/forum?id=4Ay23yeuz0>.

807 A SECTION 4 EXPERIMENT DETAILS

808 A.1 IN-CONTEXT LEARNING DETAILS

809 A.1.1 GENERATOR DETAILS

810 **TabPFN Generation** [Hollmann et al. \(2025a\)](#): We followed the Prior Labs tutorial (<https://priorlabs.ai/tutorials/unsupervised/>) for unsupervised TabPFN. Each training split was loaded, shuffled, and
811 divided into batches of 200 rows. Numeric features were cast to `float32`, while categorical variables were
812 label-encoded (assigning unseen categories to `-1`). Columns with zero variance were removed prior to model
813 fitting and reintroduced after sampling. For each batch, we fit TabPFN and generated synthetic rows using
814 temperature $t = 1.0$ across three random permutations. Outputs were decoded, constant columns reattached,
815 batches concatenated, and finally truncated to match the size of the original dataset.

816 **LLaMA Generation** [Meta AI \(2024\)](#): We used LLaMA 3.3 70B via the Groq API (<https://console.groq.com/docs/models>). Each training split was divided into batches of up to 32 rows, ensuring that all
817 rows were fully included in the prompt. For each batch, we computed per-column summary statistics and serial-
818 ized the data to CSV. We then queried `llama-3.3-70b-versatile` with temperature $t = 1.0$, requesting
819 N rows in JSON format. If outputs contained parse errors or incorrect row counts, we retried the generation
820 up to five times. Valid generations were concatenated, truncated, or re-prompted as necessary, and validated for
821 type and dimensional consistency.²

822 **GPT-4o-mini** [OpenAI \(2024\)](#): We applied the same prompting and inference pipeline as with LLaMA 3.3 70B.
823 However, we used OpenAI’s structured output API, defining the target format as a JSON schema with column
824 names as keys and corresponding cell values as entries.

825 A.1.2 PROMPT TEMPLATE

826 Listing 1: Prompt passed to Groq API

```
827 System role: You are a tabular synthetic data generation model.
828
829 Your goal is to produce data that mirrors the given examples in
830 causal structure and feature/label distributions,
831 while maximizing diversity.
832
833 Context: Leverage your in-context learning to generate realistic,
834 diverse samples.
835
836 Output format: JSON.
837
838 Dataset name: {dataset_name}
839
840 Column names (in order): {col_names}
841
842 Summary statistics:
843 {summary_stats}
844
845 CSV of full data:
846 {data}
```

847 ²LLaMA 3.3-70B failed on *geographical-origin-of-music*, *pumadyn32nh*, *student-performance-por*, *superconductivity*,
848 and *wave-energy* due to token limitations. TabPFN failed on *geographical-origin-of-music* due to extreme dimensionality.

```

855 Please generate {batch_size} rows of synthetic data.
856
857 Treat the rightmost column as the target. Return only a JSON object:
858 {
859   "synthetic_data": "<CSV string>"
860 }
861 Do not include any additional text.
862

```

864 A.2 SFT-GENERATION DETAILS

865 A.2.1 DATASETS

- 867 1. **Abalone** (OpenML): <https://www.openml.org/search?type=data&sort=runs&id=183&status=active>
- 868 2. **CA Housing** (OpenML): <https://www.openml.org/search?type=data&status=active&id=45578&sort=runs>
- 869 3. **CASP** (OpenML): <https://www.openml.org/search?type=data&status=active&id=42903>
- 870 4. **Diabetes** (OpenML): <https://archive.ics.uci.edu/dataset/34/diabetes3>
- 871 5. **Faults** (UCI): <https://archive.ics.uci.edu/dataset/198/steel+plates+faults>

872 A.3 MODEL DETAILS

873 We modify the original RealTabFormer implementation to use LLaMA 3.2 (1B, 3B) (Grattafiori et al., 2024),
874 Qwen2.5-3B (Qwen et al., 2025), and Mistral v0.3 7B Jiang et al. (2023). Here, we follow RealTabFormer’s
875 base training and sampling hyperparameters in Table 3. To SFT GREAT, we also use its original implementation
876 of which the base hyperparameters can be found in Table 4. For CT-GAN and TVAE, we use the default
877 hyperparameters and implementation found in Synthcity (Qian et al., 2023).

Hyperparameter	Default Value
epochs	1000
batch_size	8
train_size	1
output_max_length	512
early_stopping_patience	5
early_stopping_threshold	0
mask_rate	0
numeric_nparts	1
numeric_precision	4
numeric_max_len	10

884 Table 3: Numeric hyperparameters of REaLTabFormer.

Hyperparameter	Default Value
epochs	100
batch_size	8
float_precision	None
temperature	0.7
top-k sampling	100
max_length	100

885 Table 4: GReaT training and sampling hyperparameters.

890 A.4 SIMULATED DATA GENERATION DETAILS

891 For Figure 2b, we initialize a Multivariate Gaussian of $N(1e10, 1e9)$. This ensures that there are up to 10 digits
892 for a column as RealTabFormer can struggle to process exceptionally long decimal strings. We then sample

10,000 training and holdout rows for our experiment. At each level, we add additional columns to the training, holdout, and therefore synthetic dataset observation sequence lengths increase by 10 digits for each column.

B SECTION 5 EXPERIMENT DETAILS AND ADDITIONAL RESULTS

B.1 DETAILS OF DIGIT MODIFIER

Algorithm 1 Digit Modifier Protection Mechanism

```

1: Input: real training dataset  $D$ , generator  $G$ , parameters  $p_{\min}, p_{\max}$ , probability function  $g$ 
2: train  $G$  on  $D$ 
3: generate synthetic table  $D_{\text{syn}}$  with entries  $x_{k,i}$ 
4: let  $\mathcal{N}$  be the set of indices of numerical columns; set  $X_i = \{x_{1,i}, \dots, x_{n,i}\}$ 
5: for all  $(k, i, r)$  with  $i \in \mathcal{N}$  and  $r$  a digit place considered for  $x_{k,i}$  do
6:    $p \leftarrow g(x_{k,i}, X_i; p_{\min}, p_{\max})$ 
7:   sample  $b \sim \text{Bernoulli}(p)$ 
8:   if  $b = 1$  then
9:     replace the digit at position  $r$  in  $x_{k,i}$  with a different digit
10:  end if
11: end for

```

For the digit modifier, we define a probability function

$$g(\mathbf{x}_{k,i}, \mathbf{X}_i; p_{\min}, p_{\max})$$

that assigns a probability to each value $\mathbf{x}_{k,i}$ in column X_i , constrained to lie within the range $[p_{\min}, p_{\max}]$. In our experiment, we set

$$M(\mathbf{X}_i) = \max_{\mathbf{x}_{k,i} \in \mathbf{X}_i} |\mathbf{x}_{k,i}|,$$

which represents the largest absolute value in the numerical column. The probability function is then given by

$$g(\mathbf{x}_{k,i}, \mathbf{X}_i; p_{\min}, p_{\max}) = p_{\min} + (p_{\max} - p_{\min}) \frac{|\mathbf{x}_{k,i}|}{M(\mathbf{X}_i)}.$$

B.2 DETAILS OF TENDENCY-BASED LOGIT PROCESSOR

Algorithm 2 Tendency-Based Logit Processor Protection Mechanism

```

1: Input: real training dataset  $D$ , LLM-based generator  $G$ , parameter  $t$ , and curve-up function  $f_t$ 
2: train  $G$  on  $D$ 
3: while generator  $G$  outputs logits at inference time do
4:   scale all logits to  $[0, 1]$  using a shifted min-max scaler
5:   apply the curve-up function  $f_t$  to all scaled logits
6:   rescale the curved-up logits back to the original scale
7:   apply softmax to obtain the output token distribution
8:   sample the next token from this distribution, decode it to a digit in the original data space
9: end while

```

For the tendency-based logit processor, we need a min-max scaler S_l , a reverse scaler S_l^{-1} , and a curving up function f_t . Specifically, let $l = (l_1, l_2, \dots, l_k)$ be the k logits. Let $m_l = \min_j l_j$, $M_l = \max_j l_j$, and $\varepsilon > 0$ (small constant). Then the min-max scaler S_l is defined componentwise as

$$[S_l(l_1, l_2, \dots, l_n)]_i = \frac{l_i - m_l}{M_l - m_l + \varepsilon}, \quad i = 1, \dots, k.$$

Similarly, we can also define S_l^{-1} componentwise as

$$[S_l^{-1}(s_1, s_2, \dots, s_n)]_i = m_l + s_i (M_l - m_l + \varepsilon), \quad i = 1, \dots, n,$$

Too see why we need concavity, let's fix $0 < a < b \leq 1$. Since $a = \theta b$ for some $\theta \in (0, 1)$, concavity of f_t implies

$$f_t(a) = f_t(\theta b + (1 - \theta) \cdot 0) \geq \theta f_t(b) + (1 - \theta) f_t(0) = \theta f_t(b),$$

where we used $f(0) = 0$. Dividing both sides by $a = \theta b$ gives

$$\frac{f_t(a)}{a} \geq \frac{f_t(b)}{b}.$$

Thus, the ratio $f_t(x)/x$ is non-increasing in x . This means smaller logits receive a proportionally larger boost compared to larger logits. Therefore, f_t preserves the ordering of the logits (by monotonicity) while compressing their differences (by concavity), which corresponds to a “curving-up” transformation that favors lower logits.

In particular, the curving function we use in our experiment is $f_t(x) = x^{\frac{1}{t}}$. Figure 5 below demonstrates the graph of our f_t under different t .

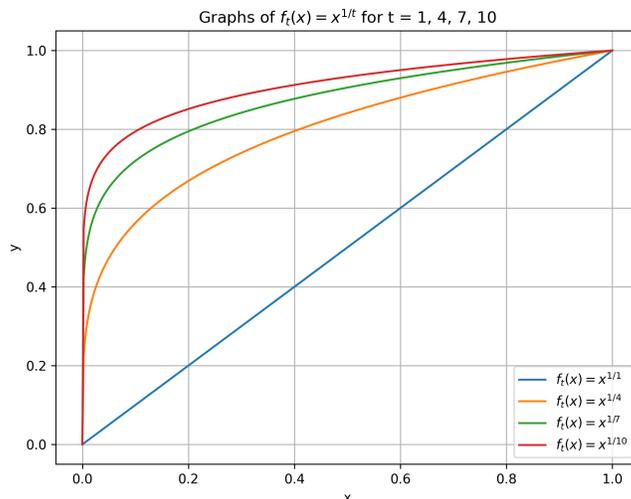


Figure 5: Visualization of the transformation function $f_t(x) = x^{1/t}$ under varying values of t . As t increases, the function becomes more concave, amplifying smaller logits proportionally more than larger ones. This behavior helps compress logit differences while preserving their ordering.

B.3 TENDENCY-BASED LOGIT PROCESSOR EFFICACY IN LOWERING LEVATT TRUE POSITIVE RATE

Please refer to Figure 6 for experimental results showing that TLP lowers the LevAtt TPR to **12.5%** at an FPR of **10%**, while maintaining good fidelity across different synthetic dataset lengths.

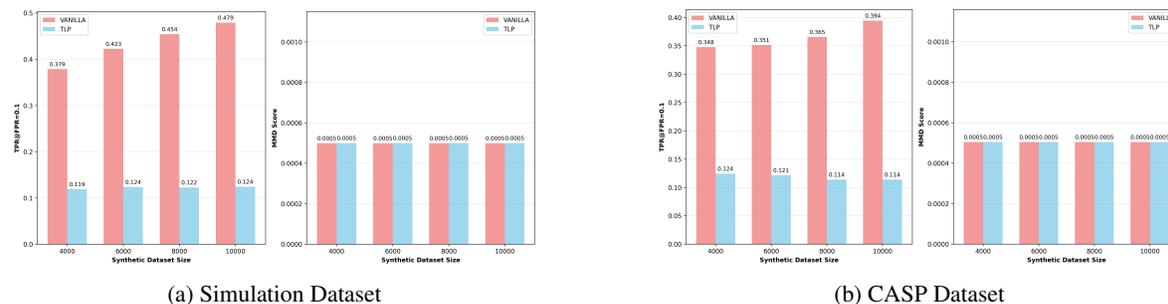
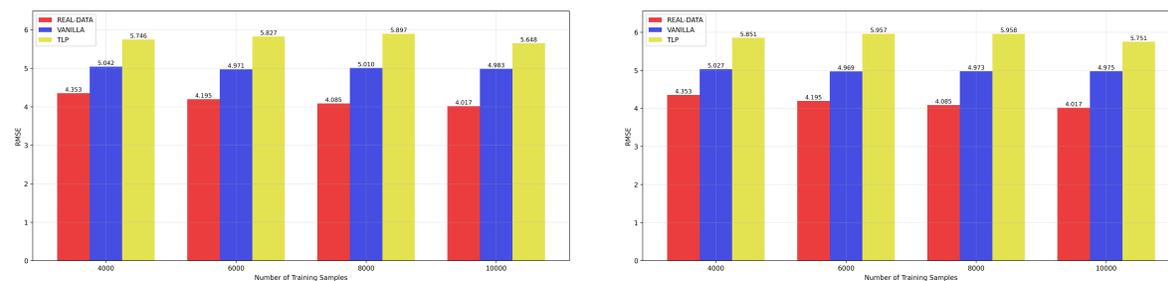


Figure 6: Privacy–utility trade-off of TLP across different datasets using RealTabFormer as the base generative model. Panel (a) reports results on the simulation dataset (5.2), while panel (b) reports results on the CASP dataset. In both cases, TLP consistently reduces the LevAtt TPR to about 12.5% at 10% FPR across all synthetic dataset sizes, indicating a substantial improvement in privacy protection compared to the vanilla inference procedure. Here, the vanilla synthetic dataset is generated using the default RealTabFormer inference without any protection mechanism. Meanwhile, the Maximum Mean Discrepancy (MMD), computed between the training data and the synthetic data generated by either vanilla or TLP, remains nearly unchanged, demonstrating that TLP preserves fidelity while mitigating privacy leakage.

B.4 TENDENCY-BASED LOGIT PROCESSOR PRIVACY VS MACHINE LEARNING EFFICACY

We refer for to Figure 7 for a full description of this experiment.



(a) CASP dataset: Utility comparison of XGBoost models trained on real data, vanilla synthetic data, and TLP-protected synthetic data, demonstrating that TLP reduces the LevAtt AUC-ROC to below 55%.

(b) CASP dataset: Utility comparison of XGBoost models trained on real data, vanilla synthetic data, and TLP-protected synthetic data, demonstrating that TLP reduces the LevAtt TPR to below 12.5% at FPR = 10%.

Figure 7: Utility comparison of XGBoost models trained on real, vanilla synthetic, and TLP-protected synthetic data. Models trained on real data achieve the lowest RMSE, while vanilla synthetic data incurs a moderate increase and TLP-protected data shows a larger degradation. Importantly, the performance gap between vanilla and TLP remains relatively stable as training size increases. Even with stronger tendency parameters required for larger synthetic datasets to meet privacy thresholds, the utility degradation appears bounded, highlighting the desirable property of TLP in controlling the privacy–utility trade-off.

C STATEMENT OF LLM USE

We used LLMs to assist in the programming of our experiments, designing data visualizations, making the codebase more user-friendly, finding clearer phrasings in our writing, and formatting Latex code. All additions by LLMs were checked by the authors.