003

008 009 010

# Do You Really Need Public Data? Surrogate Public Data for Differential Privacy on Tabular Data

Anonymous Authors<sup>1</sup>

#### Abstract

011 Differentially private (DP) machine learning often 012 relies on the availability of public data for tasks like privacy-utility trade-off estimation, hyperparameter tuning, and pretraining. While public data 015 assumptions may be reasonable in text and image data, they are less likely to hold for tabular data. This work introduces the notion of "surrogate" 018 public data - datasets generated independently of 019 sensitive data, which consume no privacy budget 020 and are constructed solely from publicly available metadata. We automate the process of generating surrogate public data with large language models (LLMs): in particular, we propose two methods: direct record generation as CSV files, and 025 automated structural causal model (SCM) construction for sampling records. Through extensive experiments, we demonstrate that surrogate pub-028 lic tabular data can effectively replace traditional 029 public data when pretraining differentially private 030 tabular classifiers. To a lesser extent, surrogate public data are also useful for hyperparameter tuning of DP synthetic data generators, and for estimating the privacy-utility tradeoff.

### 1. Introduction

034

035

049

050

051

052

053

054

Differential privacy (DP) is a mathematical framework for protecting individuals' privacy in statistical analysis and 038 machine learning (Dwork et al., 2016), and was deployed 039 in multiple recent high-stakes releases and systems (Abowd et al., 2022; Hod & Canetti, 2025; Miklau, 2022; Burman 041 et al., 2019; Wilson et al., 2020; Fitzpatrick & DeSalvo, 2020) (see (Desfontaines, 2021) for a more complete list). 043 It is common in the design of differentially private algorithms to assume access to a relevant public dataset that can 045 guide hyperparameter tuning, pretraining, or performance 046 improving mechanisms (Bassily et al., 2019; 2020b; Liu 047 et al., 2021a; Zhou et al., 2021). Executing these tasks

with *sensitive* data would require an additional allocation of the privacy budget, resulting in weaker overall privacy guarantees or reduced utility. However, using this assumed *public* data in a private mechanism avoids additional privacy budget consumption. This leads to the following informal definition of *public* data in our work:

#### Public Data (informal)

A dataset is considered *public* if a computation taking it as input does not consume privacy loss budget with respect to any fixed private, sensitive dataset.

For text and image domains, assuming public data availability is often reasonable: public image collections or large-scale textual corpora are readily available, and it has been shown that even out-of-distribution data can serve as a valuable prior in these contexts, whether through pretraining or foundation models (Nasr et al., 2023; Ganesh et al., 2023). However, this assumption does not often hold in a tabular data setting. Tabular data is heterogeneous, highdimensional, subject to strict privacy or legal restrictions, and has few universal priors (Müller et al., 2022). In many real-world domains like healthcare, finance, and government administration, tabular data encodes sensitive information that drives high-stakes decisions. It is thus rare to find truly public, non-sensitive samples with sufficient alignment to a private distribution to be used for private hyperparameter tuning or pretraining.

Nevertheless, recent theoretical insights confirm that if a public dataset is "close enough" to a sensitive data distribution, then private learning can still achieve strong utility, even when the public and private datasets are not perfectly matched (Bassily et al., 2019). In practice, however, identifying or constructing such a surrogate is often far from trivial. Real-world deployments of differentially private methods face numerous hurdles related to data availability (Cummings & Sarathy, 2023; Cummings et al., 2024a). As an example, a recent release of Israel's Live Birth Registry (Hod & Canetti, 2025) underscores the challenges of obtaining an end-to-end differential privacy guarantee.

Public data served two purposes for Hod & Canetti (2025): it helped constrain the hyperparameter space within a computationally locked-down enclave environment, and it enabled

<sup>&</sup>lt;sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

the estimation of the privacy-utility trade-off when allocating privacy budget. Yet, in general, sensitive datasets (e.g.,
birth records) are not readily available as public data. Hod &
Canetti (2025) reported finding only one open-access birth
dataset worldwide (in the U.S.); without it, estimating the
necessary parameter settings for their release would have
been significantly more challenging.

062

063

064

065

066

067

068

098

099

100

104

105

106

109

A recent practical guide for differentially private machine learning recommends that "the simplest approach, when possible, is to do all model architecture search and hyperparameter tuning on a proxy public dataset (with a distribution similar to the private data), and only use the private training dataset to train the final DP model" (Ponomareva et al., 2023).



Figure 1: (**Top**) The premise of this work: Can we utilize LLMs to generate surrogate public data to solve DP auxiliary tasks? (**Bottom**) The answer is **yes**; for example, pretraining on *surrogate* public data generated through various LLM based methods (green, orange) nearly matches the performance of pretraining on regular public data ( magenta), and outperforms no pretraining (red), or pretraining on baselines (gray, blue). Results on the EDAD dataset.

These two examples highlight a fundamental challenge; many differentially private algorithms require informed decisions *a priori* that, ideally, do not consume extra privacy budget. This leads us to consider a class of *DP auxiliary tasks*, which we define informally as:

#### **Differential Privacy Auxiliary Task**

A *differential privacy auxiliary task*, with respect to a differentially private mechanism for conducting an analysis of interest, is a required decision or procedure for execution. The auxiliary task may or may not incur privacy loss. Examples include hyperparameter tuning, setting  $\varepsilon$ , mechanism initialization, model selection, etc.

Motivated by the example of Hod & Canetti (2025) and the recommendation of Ponomareva et al. (2023), we imagine a world where we could convene a panel of domain experts, and ask them to manually encode an approximate data-generating process. In the birth registry example, epidemiologists and bio-statisticians could approximate highlevel relationships among the birth-related variables (e.g., premature birth correlated with infant weight), yielding a sufficiently similar distribution. From this data generating process, one could then generate "public" samples for differential privacy auxiliary tasks . Indeed, for many tabular settings that must accommodate strict privacy and legal constraints, we hypothesize that such an expert-driven approach could offer a practical surrogate to *traditional* public data (Hasani et al., 2024).

#### Surrogate Public Data

We consider a dataset generated independently of a sensitive dataset, consuming no privacy loss budget, and based only on publicly available schema or metadata to be a *surrogate* public dataset.

Surrogate public data is positioned in contrast to "*traditional*" *public data*, which shares a similar generation process (often *naturally occurring* then *collected*) as the private data. Then, the main question of this paper is: how useful is *surrogate* public (1) relative to *traditional* public data, or (2) relative to the *lack* of any public data? Is, for example, *automating* the process of expert panel data generation with large language models (LLMs) a suitable surrogate?

To investigate these questions, we evaluate automated data generation approaches that leverage LLMs (Borisov et al., 2023; Zhao et al., 2023; Kim et al., 2024). LLMs are trained on enormous and diverse datasets, including vast amounts of tabular data (Borisov et al., 2023; Hegselmann et al., 2023) as well as scientific literature (Phan et al., 2025) that captures rich structural and contextual knowledge of relationships between variables. This allows for the *direct* generation of coherent, causally informed relationships among variables that can lead to the generation of reasonable tabular data.

0 This leads us to our main contributions.

(1.) Methods for generating *surrogate* public data (Section 2). We introduced an agent-based strategy with a blackbox LLM access assumption to automatically construct a plausible structural causal model for surrogate public data generation. We also introduce a number of simpler baselines methods for comparison.

118 (2.) Benchmark of DP auxiliary tasks with surrogate 119 public data (Section 3). Auxiliary DP tasks are part of a 120 wider private pipeline. Consequently, evaluating the use-121 fulness of surrogate public data requires a careful design 122 across the DP downstream task, datasets, baselines, com-123 parison conditions, and aggregated metrics. In this work, 124 we propose such a benchmark framework and provide an 125 extensible, method-agnostic implementation. 126

127 (3.) In-depth experimental results evaluating surrogate 128 public data on some DP auxiliary tasks (Section 4). We 129 find that pretraining with LLM-generated surrogate public 130 data can substantially improve differentially private clas-131 sification performance; this holds true in the low dataset 132 size regime in particular. Additionally, we show that LLM-133 generated surrogate public data can be useful for hyperparameter tuning of private data synthesizers. We further 134 135 present a complicated story on using surrogate public data 136 for privacy-utility tradeoff estimation . 137

# 1381202. Producing Public Data Surrogates

139 We evaluate multiple methods for generating surrogates to 140 public data, categorizing them into baseline and LLM-based 141 approaches. For these methods, we assume that the private 142 data's metadata - consisting of the dataset schema and a 143 brief description of its topic (e.g., demographics, epidemiol-144 ogy) – is publicly available. All methods we introduce rely 145 solely on this metadata.<sup>1</sup> A schema provides a description 146 of the dataset domain and structure, specifying for each vari-147 able: (1) its name, (2) a very brief description, (3) the data 148 type (e.g., integer, string), and (4) either allowed values and 149 their meanings for categorical columns or value ranges for 150 continuous columns. (See Figure 2 for an example.) Each 151 LLM-based method is applied to the three models presented 152 in Table 2. Below we briefly summarize each approach; see 153 Appendix E for full details, and Figure 1 for an overview. 154

Baselines. We evaluate three baseline methods – Uniform,
Univariate, and Arbitrary – that rely solely on the
public schema. The Uniform method samples records i.i.d.
from each variable's full domain, while the Univariate
approach samples columns independently using empirical

1-way marginal distributions from the private data (this violates privacy, but serves as a competitive baseline). The Arbitrary method constructs a random Bayesian network over the high-dimensional domain by sequentially building a DAG (maximum in-degree of 5) and parameterizes random conditional probability tables via Dirichlet sampling Algorithm 1).

**CSV (direct generation).** The CSV method prompts LLMs to generate CSV records that strictly adhere to the schema, including exact header rows, data types, allowed values, and realistic inter-field relationships (see Figure 3). The prompt specifies rules to ensure statistical plausibility and the inclusion of realistic edge cases, while the generation is executed in batches with schema-based validation of each record. This process relies *solely* on the LLM's pretrained knowledge without any direct access to the private data.

Agent (state machine) approach. The Agent method (implemented as a state machine, see Figure 4) is a multistep process to construct a structural causal model (SCM) from the schema metadata. It begins by describing the full set of variables and domain-specific constraints, then sequentially constructs a causal DAG – first identifying root nodes and then establishing edges (ensuring acyclicity deterministically) before mapping variables to structural equations. The final output is an integrated Python program based on the Pyro package that enforces variable ranges and constraints. The Agent method has two variants: we experiment with generating multiple expert models whose records are re-sampled (using uniform or facility location-based sampling (Wang & Zhou, 2020)).

## 3. Evaluation Framework

Our evaluation framework assesses the viability of the surrogate public data in three DP auxiliary tasks: (1) classifier pretraining, (2) hyperparameter optimization, and (3) privacy-utility trade-off estimation. Each task is assessed using three datasets, and corresponding DP mechanisms. Our strategy in evaluating each task is guided by a high level question: how useful is each surrogate public data method relative to traditional public data and relative to the lack of any public data?

**Task 1: Model Pretraining for Classification.** We assess the benefit of surrogate public data as a pretraining step for binary classification on tabular data using an FTTransformer model (Gorishniy et al., 2021; Rosenblatt et al., 2024b). Public and private datasets are split into train, validation, and test sets (72:8:20), and performance is measured via AUC along with an AUC Advantage metric comparing models with and without public pretraining, aggregated over multiple experiment seeds. See Appendix F.1.1 for the detailed specification of this task.

<sup>&</sup>lt;sup>1</sup>With one exception: the *univariate* baseline, which samples
directly from the sensitive data *without* correlation between variables. This method is introduced purely for comparison, and is *not* a valid public data surrogate under our working defition.

165 Task 2: Hyperparameter Tuning for Synthetic Data Generation. We evaluate whether surrogate public datasets 167 can effectively guide hyperparameter selection for DP syn-168 thetic data generators across privacy budgets  $\varepsilon$ . For each 169 synthesizer and hyperparameter configuration, we generate synthetic datasets matching the size of the private data 170 171 and measure their performance on the private data using 172 multiple metrics (marginals, correlations, and classification-173 based; Table 3). We quantify performance degradation by 174 comparing results obtained when tuning hyperparameters on 175 surrogate public data versus directly on private data (optimal 176 choice of hyperparameters), aggregating outcomes across 177 multiple experimental seeds via Pareto frontier analysis. See 178 Appendix F.1.2 for the detailed specification of this task.

179 Task 3: Privacy-Utility Rstimation for Synthetic Data 180 Generation. Finally, we assess how accurately surrogate 181 public datasets can approximate the privacy-utility trade-off 182 curve of a DP synthetic data mechanism fitted to private 183 data. For each combination of synthesizer, dataset, and met-184 ric group, we select the best-performing hyperparameters 185 based on surrogate public datasets, and then evaluate the 186 DP mechanism across a range of  $\varepsilon$  values. This produces 187 paired performance curves: one based on the private data 188 (true curve) and others based on surrogate public datasets. 189 We quantify the dissimilarity between these curves using 190 the  $\ell_1$  and  $\ell_2$  distances. See Appendix F.1.3 for the detailed 191 specification of this task.

193 Datasets. We run the experiments on three datasets (ACS, EDAD, and WE; high-level details presented in Table 1). 195 Each dataset has a private, sensitive split; additionally, we 196 pair each dataset with a reasonable public analogue. These 197 public datasets have inherent distribution shift between them: for ACS this is a geographical variation or, for EDAD 199 and WE, temporal differences. The private split serves as 200 ground truth to benchmark the contribution of the "traditional" approach of using a public split compared to our 202 surrogate generation methods. To mitigate the risk of data memorization in LLMs, we specifically selected the private 204 splits for EDAD and WE to be *recently* published, i.e., after the training data cutoff of some of the LLMs we evaluate. 206 To this end, we include a memorization analysis, based on the methodology of Bordt et al. (2024). For the complete 208 details for each dataset and an in-depth discussion of LLM 209 memorization, refer to Appendix F.2. 210

### 4. Results

211

212

Task 1: Model Pretraining for Classification. Our experiments provide strong evidence that LLM-based methods – both CSV and Agent generation (notably with Claude 3.5
Sonnet) – offer a competitive alternative to traditional public data in the small dataset regime (fewer than 10K records).
In many settings, the surrogate public data matched (or

even occasionally exceeded) the performance of regular public data as a starting point for DP fine-tuning. Figure 7 presents our experimental results on the EDAD and WE datasets, demonstrating how pretraining on the surrogate public data can vastly improve the starting point of model performance. For the ACS dataset, we do not observe a benefit from pretraining, either with regular or surrogate public data (see Figure 13). However, when the ACS dataset is sub-sampled to a smaller dataset (e.g., 5% of the records), we observe a similar pattern regarding the usefulness the traditional and surrogate public data as with the EDAD and WE datasets. See Appendix G.2 for a further analysis of the role of dataset size. See Appendix G.1 for an in-depth treatment of all pretraining results.

Task 2: Hyperparameter Tuning for Synthetic Data Generation. No single surrogate public-data strategy dominates across all evaluation criteria (Figure 7; Tables 1-3). Each generation method achieves superior performance on a subset of metrics - whether predictive accuracy, preservation of pairwise correlations, or marginal distribution fidelity. Pareto-frontier analysis suggests that the critical determinant is the extent to which a surrogate captures higher-order dependence structure; even the Arbitrary baseline appears on the frontier even though it does not encode similar statistical relationship of the private data (for evidence, see Figure 7; Table 3). Nonetheless, the LLM-generation methods (CSV and Agent) exhibit the most favorable overall trade-offs, underscoring their utility for hyperparameter tuning of DP synthetic data generators. See Appendix G.3 for an in-depth treatment of all hyperparameter tuning results.

**Task 3: Privacy-utility Estimation for Synthetic Data Generation.** The story for this task is less clear-cut. While surrogate public data generally provides a reasonable approximation for the privacy-utility tradeoff curves, the differences between various generation methods were not pronounced. We observed that regular public data provided the best or second-best estimation of the privacy-utility tradeoff curve in the vast majority of cases. This observation suggests that data similarity may be an important contributing factor. However, a subsequent analysis (Section G.5) examining the role of similarity did not reveal a clear pattern that explains this result. See Appendix G.4 for an in-depth treatment of all privacy-utility estimation results.

### 5. Conclusion

In this work, we asked whether LLMs can be used to generate effective surrogate public data for solving DP auxiliary tasks in settings where traditional public tabular data is limited or unavailable. Overall, our results provide an affirmative answer: for the DP auxiliary tasks we considered, generating surrogate public data with LLMs *can* overcome tabular public data scarcity.

#### 220 References

Abadi, M., Chu, A., Goodfellow, I. J., McMahan, H. B., Mironov, I., Talwar, K., and Zhang, L. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, Vienna, Austria, October 24-28, 2016*, pp. 308– 318. ACM, 2016.

Abdulaal, A., Hadjivasiliou, A., Brown, N. M., He, T., Ijishakin, A., Drobnjak, I., Castro, D. C., and Alexander, D. C. Causal modelling agents: Causal graph discovery through synergising metadata- and data-driven reasoning. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11*, 2024. OpenReview.net, 2024.

Abowd, J. M., Ashmead, R., Cumings-Menon, R., Garfinkel,
S. L., Heineck, M., Heiss, C., Johns, R., Kifer, D., Leclerc,
P., Machanavajjhala, A., Moran, B., Sexton, W., Spence,
M., and Zhuravlev, P. The 2020 census disclosure avoidance system topdown algorithm. *CoRR*, abs/2204.08986,
2022.

Abowd, J. M., Adams, T., Ashmead, R., Darais, D., Dey,
S., Garfinkel, S. L., Goldschlag, N., Kifer, D., Leclerc, P.,
Lew, E., Moore, S., Rodr'iguez, R. A., Tadros, R. N., and
Vilhuber, L. The 2010 Census confidentiality protections
failed, here's how and why. Technical report, National
Bureau of Economic Research, 2023.

- Almeida, D. R. Synthetic data generation (part 1). https: //cookbook.openai.com/examples/sdg1, 2024. OpenAI Cookbook.
- Amid, E., Ganesh, A., Mathews, R., Ramaswamy, S., Song,
  S., Steinke, T., Suriyakumar, V. M., Thakkar, O., and
  Thakurta, A. Public data-assisted mirror descent for private model training. In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore,*Maryland, USA, volume 162 of Proceedings of Machine
  Learning Research, pp. 517–535. PMLR, 2022.
- Amin, K., Bie, A., Kong, W., Kurakin, A., Ponomareva,
  N., Syed, U., Terzis, A., and Vassilvitskii, S. Private
  prediction for large-scale synthetic text generation. In *Findings of the Association for Computational Linguis- tics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pp. 7244–7262. Association for Computational Linguistics, 2024.
- Anthropic. Claude API Documentation. https://docs. anthropic.com/claude/reference/, 2025.
- Aydöre, S., Brown, W., Kearns, M., Kenthapadi, K., Melis,
  L., Roth, A., and Siva, A. A. Differentially private query
  release through adaptive projection. In Meila, M. and

Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event,* volume 139 of *Proceedings of Machine Learning Research,* pp. 457–467. PMLR, 2021.

- Bassily, R., Moran, S., and Alon, N. Limits of private learning with access to public data. In Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pp. 10342–10352, 2019.
- Bassily, R., Cheu, A., Moran, S., Nikolov, A., Ullman, J. R., and Wu, Z. S. Private query release assisted by public data. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 695–703. PMLR, 2020a.
- Bassily, R., Moran, S., and Nandi, A. Learning from mixtures of private and public populations. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020b.
- Ben-David, S., Bie, A., Canonne, C. L., Kamath, G., and Singhal, V. Private distribution learning with public data: The view from sample compression. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 -16, 2023, 2023.
- Bie, A., Kamath, G., and Singhal, V. Private estimation with public data. 2022.
- Bingham, E., Chen, J. P., Jankowiak, M., Obermeyer, F., Pradhan, N., Karaletsos, T., Singh, R., Szerlip, P. A., Horsfall, P., and Goodman, N. D. Pyro: Deep universal probabilistic programming. *J. Mach. Learn. Res.*, 20: 28:1–28:6, 2019.
- Block, A., Bun, M., Desai, R., Shetty, A., and Wu, S. Oracleefficient differentially private learning with public data. *CoRR*, abs/2402.09483, 2024.
- Bordt, S., Nori, H., Rodrigues, V., Nushi, B., and Caruana, R. Elephants never forget: Memorization and learning of tabular data in large language models. In *Conference on Languge Modeling (COLM)*, 2024.
- 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, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net, 2023.

- Bu, Z., Wang, Y., Zha, S., and Karypis, G. Differentially
  private bias-term fine-tuning of foundation models. In *Forty-first International Conference on Machine Learn- ing, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024.
- Burman, L. E., Engler, A., Khitatrakun, S., Nunns, J. R., Armstrong, S., Iselin, J., MacDonald, G., and Stallworth, P. Safely expanding research access to administrative tax data: creating a synthetic public use file and a validation server. *Technical report US, Internal Revenue Service*, 2019.
- Bynum, L. E. J. and Cho, K. Language models as causal
  effect generators. *CoRR*, abs/2411.08019, 2024.
- Cai, K., Lei, X., Wei, J., and Xiao, X. Data synthesis via differentially private markov random field. *Proc. VLDB Endow.*, 14(11):2190–2202, 2021.
- Carlini, N., Tramèr, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T. B., Song, D., Erlingsson, Ú., Oprea, A., and Raffel, C. Extracting training data from large language models. In *30th USENIX Security Symposium, USENIX Security 2021, August 11-13, 2021, 2021.*
- Carlini, N., Ippolito, D., Jagielski, M., Lee, K., Tramèr, F.,
  and Zhang, C. Quantifying memorization across neural
  language models. In *The Eleventh International Confer*-*ence on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023, 2023.*
- Cattan, Y., Choquette-Choo, C. A., Papernot, N., and Thakurta, A. Fine-tuning with differential privacy necessitates an additional hyperparameter search. *CoRR*, abs/2210.02156, 2022.
- Chen, S., Peng, B., Chen, M., Wang, R., Xu, M., Zeng, X.,
  Zhao, R., Zhao, S., Qiao, Y., and Lu, C. Causal evaluation
  of language models. *CoRR*, abs/2405.00622, 2024.
- Chen, T. and Guestrin, C. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, pp. 785–794. ACM, 2016.
- Cheng, J., Marone, M., Weller, O., Lawrie, D. J., Khashabi,
  D., and Durme, B. V. Dated data: Tracing knowledge cutoffs in large language models. In *Conference on Languge Modeling (COLM)*, 2024.
- Cummings, R. and Sarathy, J. Centering policy and practice: Research gaps around usable differential privacy. In 5th IEEE International Conference on Trust, Privacy and Security in Intelligent Systems and Applications, TPS-ISA 2023, Atlanta, GA, USA, November 1-4, 2023, pp. 122–135. IEEE, 2023.

- Cummings, R., Desfontaines, D., Evans, D., Geambasu, R., Huang, Y., Jagielski, M., Kairouz, P., Kamath, G., Oh, S., Ohrimenko, O., Papernot, N., Rogers, R., Shen, M., Song, S., Su, W., Terzis, A., Thakurta, A., Vassilvitskii, S., Wang, Y.-X., Xiong, L., Yekhanin, S., Yu, D., Zhang, H., and Zhang, W. Advancing Differential Privacy: Where We Are Now and Future Directions for Real-World Deployment. *Harvard Data Science Review*, 6(1), 2024a.
- Cummings, R., Hod, S., Sarathy, J., and Swanberg, M. AT-TAXONOMY: unpacking differential privacy guarantees against practical adversaries. *CoRR*, abs/2405.01716, 2024b.
- Darvariu, V., Hailes, S., and Musolesi, M. Large language models are effective priors for causal graph discovery. *CoRR*, abs/2405.13551, 2024.
- DeepSeek-AI, Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., Zhang, X., Yu, X., Wu, Y., Wu, Z. F., Gou, Z., Shao, Z., Li, Z., Gao, Z., Liu, A., Xue, B., Wang, B., Wu, B., Feng, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C., Dai, D., Chen, D., Ji, D., Li, E., Lin, F., Dai, F., Luo, F., Hao, G., Chen, G., Li, G., Zhang, H., Bao, H., Xu, H., Wang, H., Ding, H., Xin, H., Gao, H., Qu, H., Li, H., Guo, J., Li, J., Wang, J., Chen, J., Yuan, J., Qiu, J., Li, J., Cai, J. L., Ni, J., Liang, J., Chen, J., Dong, K., Hu, K., Gao, K., Guan, K., Huang, K., Yu, K., Wang, L., Zhang, L., Zhao, L., Wang, L., Zhang, L., Xu, L., Xia, L., Zhang, M., Zhang, M., Tang, M., Li, M., Wang, M., Li, M., Tian, N., Huang, P., Zhang, P., Wang, Q., Chen, Q., Du, Q., Ge, R., Zhang, R., Pan, R., Wang, R., Chen, R. J., Jin, R. L., Chen, R., Lu, S., Zhou, S., Chen, S., Ye, S., Wang, S., Yu, S., Zhou, S., Pan, S., and Li, S. S. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. CoRR, abs/2501.12948, 2025.
- Desfontaines, D. A list of real-world uses of differential privacy - Ted is writing things — desfontain.es. https://desfontain.es/privacy/ real-world-differential-privacy.html, 2021.
- Devlin, J., Chang, M., Lee, K., and Toutanova, K. BERT: pre-training of deep bidirectional transformers for language understanding. pp. 4171–4186, 2019.
- Dong, Y., Jiang, X., Liu, H., Jin, Z., Gu, B., Yang, M., and Li, G. Generalization or memorization: Data contamination and trustworthy evaluation for large language models. In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024.* Association for Computational Linguistics, 2024.

- Dwork, C., Naor, M., Reingold, O., Rothblum, G. N., and
  Vadhan, S. P. On the complexity of differentially private
  data release: efficient algorithms and hardness results. In
  Mitzenmacher, M. (ed.), *Proceedings of the 41st Annual*
- ACM Symposium on Theory of Computing, STOC 2009,
   Bethesda, MD, USA, May 31 June 2, 2009, pp. 381–390.
   ACM, 2009.
- Dwork, C., McSherry, F., Nissim, K., and Smith, A. D.
  Calibrating noise to sensitivity in private data analysis.
  volume 7, pp. 17–51, 2016.
- Ehrgott, M. *Multicriteria optimization*, volume 491.
  Springer Science & Business Media, 2005.
- Fitzpatrick, J. and DeSalvo, K. Helping public health officials combat covid-19. https: //blog.google/technology/health/ covid-19-community-mobility-reports/, 2020.
- Gallegos, I. O., Rossi, R. A., Barrow, J., Tanjim, M. M.,
  Kim, S., Dernoncourt, F., Yu, T., Zhang, R., and Ahmed,
  N. K. Bias and fairness in large language models: A
  survey. *Computational Linguistics*, 50(3):1097–1179,
  2024.
- 355 Ganesh, A., Haghifam, M., Nasr, M., Oh, S., Steinke, T., 356 Thakkar, O., Thakurta, A. G., and Wang, L. Why is 357 public pretraining necessary for private model training? 358 In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., 359 Sabato, S., and Scarlett, J. (eds.), International Confer-360 ence on Machine Learning, ICML 2023, 23-29 July 2023, 361 Honolulu, Hawaii, USA, volume 202 of Proceedings of 362 Machine Learning Research, pp. 10611–10627. PMLR, 363 2023. 364
- Ginart, A., van der Maaten, L., Zou, J., and Guo, C. Submix: Practical private prediction for large-scale language
  models. *CoRR*, abs/2201.00971, 2022.
- Golatkar, A., Achille, A., Wang, Y., Roth, A., Kearns, M.,
  and Soatto, S. Mixed differential privacy in computer
  vision. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, pp. 8366–8376. IEEE, 2022.*
- Golchin, S. and Surdeanu, M. Time travel in llms: Tracing data contamination in large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024, 2024.*
- Gorishniy, Y., Rubachev, I., Khrulkov, V., and Babenko,
  A. Revisiting deep learning models for tabular data. In
  Advances in Neural Information Processing Systems 34:
  Annual Conference on Neural Information Processing

Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 18932–18943, 2021.

- Gorishniy, Y., Rubachev, I., and Babenko, A. On embeddings for numerical features in tabular deep learning. 2022.
- Gu, X., Kamath, G., and Wu, Z. S. Choosing public datasets for private machine learning via gradient subspace distance. *CoRR*, abs/2303.01256, 2023.
- Gulati, M. and Roysdon, P. F. Tabmt: Generating tabular data with masked transformers. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023, 2023.
- Hagberg, A., Swart, P. J., and Schult, D. A. Exploring network structure, dynamics, and function using networkx. Technical report, Los Alamos National Laboratory (LANL), Los Alamos, NM (United States), 2008.
- Hardt, M., Ligett, K., and McSherry, F. A simple and practical algorithm for differentially private data release. In Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States, pp. 2348–2356, 2012.
- Hasani, W. S. R., Musa, K. I., Chen, X. W., and Cheng, K. Y. Constructing causal pathways for premature cardiovascular disease mortality using directed acyclic graphs with integrating evidence synthesis and expert knowledge. *Scientific Reports*, 14(1):28849, 2024.
- He, J., Li, X., Yu, D., Zhang, H., Kulkarni, J., Lee, Y. T., Backurs, A., Yu, N., and Bian, J. Exploring the limits of differentially private deep learning with group-wise clipping. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net, 2023.
- Hegselmann, S., Buendia, A., Lang, H., Agrawal, M., Jiang, X., and Sontag, D. A. Tabllm: Few-shot classification of tabular data with large language models. In Ruiz, F. J. R., Dy, J. G., and van de Meent, J. (eds.), *International Conference on Artificial Intelligence and Statistics*, 25-27 April 2023, Palau de Congressos, Valencia, Spain, volume 206 of Proceedings of Machine Learning Research, pp. 5549–5581. PMLR, 2023.
- Hod, S. and Canetti, R. Differentially private release of Israel's national registry of live births. In 46th IEEE Symposium on Security and Privacy, SP 2025, San Francisco, CA, USA, May 12-15, 2025. IEEE, 2025.

- Instituto Nacional de Estadística. Disabilities survey - results - microdata. https://www.ine. es/dyngs/INEbase/en/operacion.htm? c=Estadistica\_C&cid=1254736176782& menu=resultados&idp=1254735573175# \_tabs-1254736195313, 2024.
- Iyengar, R., Near, J. P., Song, D., Thakkar, O., Thakurta,
  A., and Wang, L. Towards practical differentially private
  convex optimization. In 2019 IEEE Symposium on Security and Privacy, SP 2019, San Francisco, CA, USA, May
  19-23, 2019, pp. 299–316. IEEE, 2019.
- 397 Jaech, A., Kalai, A., Lerer, A., Richardson, A., El-Kishky, 398 A., Low, A., Helyar, A., Madry, A., Beutel, A., Carney, 399 A., Iftimie, A., Karpenko, A., Passos, A. T., Neitz, A., 400 Prokofiev, A., Wei, A., Tam, A., Bennett, A., Kumar, 401 A., Saraiva, A., Vallone, A., Duberstein, A., Kondrich, 402 A., Mishchenko, A., Applebaum, A., Jiang, A., Nair, 403 A., Zoph, B., Ghorbani, B., Rossen, B., Sokolowsky, B., 404 Barak, B., McGrew, B., Minaiev, B., Hao, B., Baker, B., 405 Houghton, B., McKinzie, B., Eastman, B., Lugaresi, C., 406 Bassin, C., Hudson, C., Li, C. M., de Bourcy, C., Voss, 407 C., Shen, C., Zhang, C., Koch, C., Orsinger, C., Hesse, C., 408 Fischer, C., Chan, C., Roberts, D., Kappler, D., Levy, D., 409 Selsam, D., Dohan, D., Farhi, D., Mely, D., Robinson, D., 410 Tsipras, D., Li, D., Oprica, D., Freeman, E., Zhang, E., 411 Wong, E., Proehl, E., Cheung, E., Mitchell, E., Wallace, 412 E., Ritter, E., Mays, E., Wang, F., Such, F. P., Raso, 413 F., Leoni, F., Tsimpourlas, F., Song, F., von Lohmann, 414 F., Sulit, F., Salmon, G., Parascandolo, G., Chabot, G., 415 Zhao, G., Brockman, G., Leclerc, G., Salman, H., Bao, 416 H., Sheng, H., Andrin, H., Bagherinezhad, H., Ren, H., 417 Lightman, H., Chung, H. W., Kivlichan, I., O'Connell, 418 I., Osband, I., Gilaberte, I. C., and Akkaya, I. Openai o1 419 system card. CoRR, abs/2412.16720, 2024. 420
- Jaynes, E. T. Information theory and statistical mechanics. *Physical review*, 106(4):620, 1957.
- Kairouz, P., Diaz, M. R., Rush, K., and Thakurta, A.
  (nearly) dimension independent private ERM with adagrad ratesvia publicly estimated subspaces. In Belkin,
  M. and Kpotufe, S. (eds.), *Conference on Learning Theory, COLT 2021, 15-19 August 2021, Boulder, Colorado, USA*, volume 134 of *Proceedings of Machine Learning Research*, pp. 2717–2746. PMLR, 2021.
- Kandpal, N., Wallace, E., and Raffel, C. Deduplicating training data mitigates privacy risks in language models. In *International Conference on Machine Learning, ICML* 2022, 17-23 July 2022, Baltimore, Maryland, USA, 2022.
- Ke, S., Hou, C., Fanti, G., and Oh, S. On the convergence of differentially-private fine-tuning: To linearly probe or to fully fine-tune? *CoRR*, abs/2402.18905, 2024.

- Khan, S. H., Naseer, M., Hayat, M., Zamir, S. W., Khan, F. S., and Shah, M. Transformers in vision: A survey. *ACM Comput. Surv.*, 54(10s):200:1–200:41, 2022.
- Kiciman, E., Ness, R., Sharma, A., and Tan, C. Causal reasoning and large language models: Opening a new frontier for causality. *CoRR*, abs/2305.00050, 2023.
- Kim, J., Kim, T., and Choo, J. Group-wise prompting for synthetic tabular data generation using large language models. *CoRR*, abs/2404.12404, 2024.
- Kurakin, A., Ponomareva, N., Syed, U., MacDermed, L., and Terzis, A. Harnessing large-language models to generate private synthetic text. *CoRR*, abs/2306.01684, 2023.
- Le, H. D., Xia, X., and Chen, Z. Multi-agent causal discovery using large language models. *CoRR*, abs/2407.15073, 2024.
- Lemieux, C., Taylor, S., Stone, A., Wooden, P., and Chauhan, C. Workplace equity survey 2023, 2024. https://doi.org/10.3886/E202701V1.
- Lewis, P. S. H., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W., Rocktäschel, T., Riedel, S., and Kiela, D. Retrievalaugmented generation for knowledge-intensive NLP tasks. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.
- Liu, T., Vietri, G., Steinke, T., Ullman, J. R., and Wu, Z. S. Leveraging public data for practical private query release. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 6968–6977. PMLR, 2021a.
- Liu, T., Vietri, G., and Wu, S. Iterative methods for private synthetic data: Unifying framework and new methods. In Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 690–702, 2021b.
- Liu, T., Vietri, G., and Wu, S. Z. Iterative methods for private synthetic data: Unifying framework and new methods. Advances in Neural Information Processing Systems, 34:690–702, 2021c.
- Long, S., Piché, A., Zantedeschi, V., Schuster, T., and Drouin, A. Causal discovery with language models as imperfect experts. *CoRR*, abs/2307.02390, 2023a.
- Long, S., Schuster, T., and Piché, A. Can large language models build causal graphs? *CoRR*, abs/2303.05279, 2023b.

- Lowy, A., Li, Z., Huang, T., and Razaviyayn, M. Optimal differentially private model training with public
  data. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024.*OpenReview.net, 2024.
- Ma, J., Dankar, A., Stein, G., Yu, G., and Caterini, A. L.
  Tabpfgen tabular data generation with tabpfn. *CoRR*, abs/2406.05216, 2024.
- Magar, I. and Schwartz, R. Data contamination: From memorization to exploitation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022, 2022.*
- McKenna, R., Miklau, G., and Sheldon, D. Winning the
  NIST contest: A scalable and general approach to differentially private synthetic data. *J. Priv. Confidentiality*, 11
  (3), 2021.
- 460 McKenna, R., Mullins, B., Sheldon, D., and Miklau, G.
  461 AIM: an adaptive and iterative mechanism for differen462 tially private synthetic data. *Proc. VLDB Endow.*, 15(11):
  463 2599–2612, 2022.
- McKenna, R., Miklau, G., Hay, M., and Machanavajjhala, A.
  Optimizing error of high-dimensional statistical queries under differential privacy. J. Priv. Confidentiality, 13(1), 2023.
- McSherry, F. and Talwar, K. Mechanism design via differential privacy. In 48th Annual IEEE Symposium on Foundations of Computer Science (FOCS 2007), October 20-23, 2007, Providence, RI, USA, Proceedings, pp. 94–103. IEEE Computer Society, 2007.
- 475
  476
  476
  477
  477
  478
  478
  478
  475
  478
  475
  478
  475
  478
  478
  475
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
  478
- 479 Müller, S., Hollmann, N., Pineda-Arango, S., Grabocka, J.,
  480 and Hutter, F. Transformers can do bayesian inference.
  481 2022.
- 482 Nasr, M., Mahloujifar, S., Tang, X., Mittal, P., and 483 Houmansadr, A. Effectively using public data in privacy 484 preserving machine learning. In Krause, A., Brunskill, E., 485 Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), 486 International Conference on Machine Learning, ICML 487 2023, 23-29 July 2023, Honolulu, Hawaii, USA, volume 488 202 of Proceedings of Machine Learning Research, pp. 489 25718-25732. PMLR, 2023. 490
- Olatunji, I. E., Funke, T., and Khosla, M. Releasing graph neural networks with differential privacy guarantees. *Trans. Mach. Learn. Res.*, 2023, 2023.

- **OpenAI. OpenAI API Documentation.** https://platform.openai.com/docs, 2025.
- Phan, L., Gatti, A., Han, Z., Li, N., Hu, J., Zhang, H., Shi, S., Choi, M., Agrawal, A., Chopra, A., Khoja, A., Kim, R., Hausenloy, J., Zhang, O., Mazeika, M., Anderson, D., Nguyen, T., Mahmood, M., Feng, F., Feng, S. Y., Zhao, H., Yu, M., Gangal, V., Zou, C., Wang, Z., Wang, J. P., Kumar, P., Pokutnyi, O., Gerbicz, R., Popov, S., Levin, J., Kazakov, M., Schmitt, J., Galgon, G., Sanchez, A., Lee, Y., Yeadon, W., Sauers, S., Roth, M., Agu, C., Riis, S., Giska, F., Utpala, S., Giboney, Z., Goshu, G. M., of Arc Xavier, J., Crowson, S., Naiya, M. M., Burns, N., Finke, L., Cheng, Z., Park, H., Fournier-Facio, F., Wydallis, J., Nandor, M., Singh, A., Gehrunger, T., Cai, J., McCarty, B., Duclosel, D., Nam, J., Zampese, J., Hoerr, R. G., Bacho, A., Loume, G. A., Galal, A., Cao, H., Garretson, A. C., Sileo, D., Ren, Q., Cojoc, D., Arkhipov, P., Qazi, U., Li, L., Motwani, S., de Witt, C. S., Taylor, E., Veith, J., Singer, E., Hartman, T. D., Rissone, P., Jin, J., Shi, J. W. L., Willcocks, C. G., Robinson, J., Mikov, A., Prabhu, A., Tang, L., Alapont, X., Uro, J. L., Zhou, K., de Oliveira Santos, E., Maksimov, A. P., Vendrow, E., Zenitani, K., Guillod, J., Li, Y., Vendrow, J., Kuchkin, V., and Ze-An, N. Humanity's last exam. CoRR, abs/2501.14249, 2025.
- Ponomareva, N., Hazimeh, H., Kurakin, A., Xu, Z., Denison, C., McMahan, H. B., Vassilvitskii, S., Chien, S., and Thakurta, A. G. How to dp-fy ML: A practical guide to machine learning with differential privacy. *J. Artif. Intell. Res.*, 77:1113–1201, 2023.
- Roberts, M., Thakur, H., Herlihy, C., White, C., and Dooley, S. To the cutoff... and beyond? A longitudinal perspective on LLM data contamination. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net, 2024.
- Rosenblatt, L., Liu, X., Pouyanfar, S., de Leon, E., Desai, A., and Allen, J. Differentially private synthetic data: Applied evaluations and enhancements. *CoRR*, abs/2011.05537, 2020.
- Rosenblatt, L., Herman, B., Holovenko, A., Lee, W., Loftus, J. R., McKinnie, E., Rumezhak, T., Stadnik, A., Howe, B., and Stoyanovich, J. Epistemic parity: Reproducibility as an evaluation metric for differential privacy. *Proc. VLDB Endow.*, 16(11):3178–3191, 2023.
- Rosenblatt, L., Howe, B., and Stoyanovich, J. Are data experts buying into differentially private synthetic data? gathering community perspectives. *CoRR*, abs/2412.13030, 2024a.
- Rosenblatt, L., Lut, Y., Turok, E., Avella-Medina, M., and Cummings, R. Differential privacy under class imbalance:

- 495 Methods and empirical insights. *CoRR*, abs/2411.05733,496 2024b.
- 497
  498
  498
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  499
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
  400
- Solatorio, A. V. and Dupriez, O. Realtabformer: Generating
  realistic relational and tabular data using transformers. *CoRR*, abs/2302.02041, 2023.
- Spilka, S., Taylor, S., and Wachter, J. Workplace equity survey, 2020. https://doi.org/10.3886/
  E116922V2.
- 511
  512
  513
  514
  513
  514
  514
  514
  515
  515
  516
  516
  517
  518
  518
  511
  511
  512
  512
  513
  514
  515
  514
  515
  516
  517
  518
  518
  511
  511
  512
  512
  513
  514
  514
  515
  515
  516
  517
  518
  518
  511
  511
  512
  512
  512
  514
  515
  514
  515
  516
  517
  518
  518
  518
  518
  511
  511
  512
  512
  512
  512
  512
  512
  512
  512
  512
  512
  514
  515
  515
  516
  517
  518
  518
  518
  518
  518
  518
  518
  517
  518
  518
  518
  518
  518
  518
  518
  518
  514
  515
  516
  517
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
  518
- Swanberg, M., McKenna, R., Roth, E., Cheu, A., and
  Kairouz, P. Is API access to llms useful for generating
  private synthetic tabular data? *CoRR*, abs/2502.06555,
  2025.
- Tao, Y., McKenna, R., Hay, M., Machanavajjhala, A., and
  Miklau, G. Benchmarking differentially private synthetic
  data generation algorithms. *CoRR*, abs/2112.09238, 2021.
- Task, C., Bhagat, K., Sen, A., Streat, D., Simpson, A.,
  and Howarth, G. The NIST data excerpt benchmarks.
  https://github.com/usnistgov/SDNist/
  blob/main/BenchmarkData/README.md, 2023.
  NIST CRC.
- Thaker, P., Setlur, A., Wu, S. Z., and Smith, V. On the benefits of public representations for private transfer learning under distribution shift. In Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024, 2024.
- Together AI. Together AI LLaMA API Documentation.
   https://docs.together.ai/reference/
   chat-completions-1, 2025.
- Tramèr, F. and Boneh, D. Differentially private learning needs better features (or much more data). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.

- Tramèr, F., Kamath, G., and Carlini, N. Position: Considerations for differentially private learning with large-scale public pretraining. In *Forty-first International Conference* on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024. OpenReview.net, 2024.
- Vashishtha, A., Reddy, A. G., Kumar, A., Bachu, S., Balasubramanian, V. N., and Sharma, A. Causal inference using llm-guided discovery. *CoRR*, abs/2310.15117, 2023.
- Vietri, G., Tian, G., Bun, M., Steinke, T., and Wu, Z. S. New oracle-efficient algorithms for private synthetic data release. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event,* volume 119 of *Proceedings of Machine Learning Research,* pp. 9765–9774. PMLR, 2020.
- Wang, D., Hu, L., Zhang, H., Gaboardi, M., and Xu, J. Generalized linear models in non-interactive local differential privacy with public data. *J. Mach. Learn. Res.*, 24: 132:1–132:57, 2023a.
- Wang, J. and Zhou, Z. Differentially private learning with small public data. In *The Thirty-Fourth AAAI Conference* on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020, pp. 6219–6226. AAAI Press, 2020.
- Wang, X., Wei, J., Schuurmans, D., Le, Q. V., Chi, E. H., Narang, S., Chowdhery, A., and Zhou, D. Selfconsistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023, 2023b.*
- Wilson, R. J., Zhang, C. Y., Lam, W., Desfontaines, D., Simmons-Marengo, D., and Gipson, B. Differentially private SQL with bounded user contribution. *Proc. Priv. Enhancing Technol.*, 2020(2):230–250, 2020.
- Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., et al. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pp. 38–45, 2020.
- Wu, S., Xu, Z., Zhang, Y., Zhang, Y., and Ramage, D. Prompt public large language models to synthesize data for private on-device applications. *CoRR*, abs/2404.04360, 2024.
- Xu, C., Guan, S., Greene, D., and Kechadi, M. T. Benchmark data contamination of large language models: A survey. *CoRR*, abs/2406.04244, 2024.

- Xu, L., Skoularidou, M., Cuesta-Infante, A., and Veeramachaneni, K. Modeling tabular data using conditional
  GAN. pp. 7333–7343, 2019.
- Yu, D., Naik, S., Backurs, A., Gopi, S., Inan, H. A., Kamath,
  G., Kulkarni, J., Lee, Y. T., Manoel, A., Wutschitz, L.,
  Yekhanin, S., and Zhang, H. Differentially private finetuning of language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net, 2022.
- Yu, D., Backurs, A., Gopi, S., Inan, H., Kulkarni, J., Lin, Z.,
  Xie, C., Zhang, H., and Zhang, W. Training private and
  efficient language models with synthetic data from llms.
  In Socially Responsible Language Modelling Research,
  2023.
- 566 Zhang, J., Cormode, G., Procopiuc, C. M., Srivastava, D.,
  567 and Xiao, X. Privbayes: private data release via bayesian
  568 networks. In Dyreson, C. E., Li, F., and Özsu, M. T.
  569 (eds.), *International Conference on Management of Data,*570 *SIGMOD 2014, Snowbird, UT, USA, June 22-27, 2014,*571 pp. 1423–1434. ACM, 2014.
  - Zhang, Y., Zhang, Y., Gan, Y., Yao, L., and Wang, C. Causal graph discovery with retrieval-augmented generation based large language models. *CoRR*, abs/2402.15301, 2024.
  - Zhao, Z., Birke, R., and Chen, L. Y. Tabula: Harnessing language models for tabular data synthesis. *CoRR*, abs/2310.12746, 2023.

Zhou, Y., Wu, S., and Banerjee, A. Bypassing the ambient
dimension: Private SGD with gradient subspace identification. In 9th International Conference on Learning *Representations, ICLR 2021, Virtual Event, Austria, May*3-7, 2021. OpenReview.net, 2021.

## 605 A. Related Work

611

Public data in differential privacy Empirical evidence demonstrates that public data can improve the performance
of differentially private machine learning models through a two-stage approach: pretraining on public data followed by
differentially private fine-tuning on sensitive data. This approach has been extensively studied across NLP and vision tasks
(Tramèr & Boneh, 2021; Amid et al., 2022; Yu et al., 2022; Golatkar et al., 2022; Ginart et al., 2022; He et al., 2023; Bu
et al., 2024).

Ganesh et al. (2023) identify two phases in neural network optimization within non-convex loss landscapes. The first locates an optimal basin, where public data suffices and using the privacy budget is unnecessary. The second performs local optimization within that basin; here, if the public and target distributions differ – as they often do – consuming privacy budget to update weights with sensitive data is beneficial. Supporting this, Thaker et al. (2024) show that public pretraining outperforms fully private training in vision tasks, even under significant distribution shifts. This advantage holds even when private fine-tuning is limited to the final layer, as in Ke et al. (2024).

618 Another research direction incorporates public data directly into differentially private computations, rather than treating it as 619 a separate preprocessing step. This approach spans private estimation (Bie et al., 2022), statistical queries (Bassily et al., 620 2020a; Liu et al., 2021a), and learning and optimization (Bassily et al., 2019; Wang & Zhou, 2020; Bassily et al., 2020b; 621 Kairouz et al., 2021; Zhou et al., 2021; Nasr et al., 2023; Ben-David et al., 2023; Gu et al., 2023; Olatunji et al., 2023; 622 Wang et al., 2023a; Block et al., 2024; Lowy et al., 2024). An emerging line of research finetunes pretrained, open-source 623 LLMs on private, sensitive data with DP-SGD (Abadi et al., 2016) to generate training data for downstream models, such as 624 classifiers or other LLMs (Kurakin et al., 2023; Yu et al., 2023; Amin et al., 2024; Wu et al., 2024). For a broader survey on 625 recent advances in privacy research, see (Cummings et al., 2024a). 626

627 Finally, contemporary work by Swanberg et al. (2025) is closely related to ours, but with three key differences. First, 628 while they evaluate LLM-generated public data in a single experimental setting (for public pretraining of private synthetic 629 data mechanisms), we assess its utility across several DP auxiliary tasks — including hyperparameter tuning for synthetic 630 data generation, privacy/utility tradeoff estimation, and private classifier pretraining. Second, our evaluation is broader in 631 scope, incorporating multiple datasets (with different data-origins), diverse metrics and additional baselines / methods for 632 leveraging an LLM to produce surrogate public data. Third, we designed our experiments to mitigate the risk of positive 633 results due to memorization, including an explicit test based on Bordt et al. (2024), and provide results and analysis to assess 634 the impact of data leakage on the performance of our methods. 635

636 Generating tabular data with LLMs Transformer-based models can be used to generate synthetic samples from tabular 637 data. The fundamental approach involves treating each record as a "sentence" for the transformer architecture to process. 638 Two overall strategies exist: training transformers from scratch specifically for tabular data and adapting pretrained LLMs 639 for tabular generation tasks. For the first strategy, one variant trains a transformer on an individual dataset or distribution 640 to produce synthetic records (Solatorio & Dupriez, 2023; Zhao et al., 2023; Gulati & Roysdon, 2023; Zhao et al., 2023); 641 another variant pretrains a general tabular foundation model on multiple datasets and then adapts this model to novel unseen 642 datasets through in-context learning (Ma et al., 2024). The second strategy uses existing pretrained LLMs, adapting them for 643 tabular data generation through either fine-tuning (Borisov et al., 2023) or in-context learning (Seedat et al., 2024; Kim et al., 644 2024). These methods *cannot* be directly applied to our setting, as we only consider DP auxiliary tasks (e.g., pretraining, 645 hyperparameter tuning) that do not consume privacy budget. Both approaches condition on sensitive data and thus require 646 accounting for privacy loss. 647

648 Recent work has explored the potential of LLMs for causal modeling tasks, including pairwise causal discovery, causal 649 model generation, and counterfactual reasoning (Kiciman et al., 2023; Chen et al., 2024). While causality itself is not the 650 primary focus of our project, the ability to produce plausible causal models is highly relevant since causal models are also 651 generative, capable of producing realistic records. LLM-based causal model discovery methods can operate either with 652 metadata alone (using only dataset descriptions and schema) (Vashishtha et al., 2023; Long et al., 2023b;; Zhang et al., 653 2024; Darvariu et al., 2024; Bynum & Cho, 2024) or with additional observations (Abdulaal et al., 2024; Le et al., 2024) 654 - with the metadata-only approach being particularly relevant to our project as we have no access to observations. Most 655 literature in this area that operates without observations focuses solely on discovering causal graphs – descriptions of causal 656 dependencies without specifying conditional distributions. However, (Bynum & Cho, 2024) extends this approach by adding 657 a second step that prompts LLMs with the topological order over variables, embedded in a prompt structure, to generate 658 records directly. 659

## 660 **B. Future Work**

661 The strong performance of the Arbitrary method in hyperparameter tuning is intriguing, as it suggests that finding 662 good-enough hyperparameter configuration might depend on the record domain and the synthetic data generators, and not 663 necessarily on the private data. This raises questions about potential theoretical justifications for this observation.

The fact that traditional public data often performs best for privacy/utility trade-off estimation would lead us to believe that dataset similarity plays an important role for this task. We hypothesize that the two similarity metrics used in this work, while being natural candidates, may not adequately capture dataset characteristics relevant to estimating the behavior of data synthesizers across privacy budget settings. Identifying metrics that better predict which surrogate data provides accurate trade-off estimations would be beneficial. Such a metric could enable, e.g., the exponential mechanism (McSherry & Talwar, 2007) to select similar datasets (or combinations of datasets), if such a metric had low sensitivity with respect to the private dataset.

<sup>672</sup> Several additional DP auxiliary tasks remain unexplored in our study, such as using public data for seeding synthetic data <sup>673</sup> generation (Swanberg et al., 2025) and assessing the success rate of privacy attacks as a function of  $\varepsilon$  (Cummings et al., <sup>674</sup> 2024b). We leave these avenues for future research.

We propose three approaches to improve the quality of surrogate data produced by Agent-based methods, making it more closely resemble private data. First, subject matter experts can review and refine the generated SCM to better encode experts' domain knowledge. Second, Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) could be beneficial to surface specific knowledge from scientific literature, enabling the model to incorporate both accurate causal relationships and their quantitative parameters as established in peer-reviewed research. Third, recent advancements in reasoning LLMs (Sun et al., 2023; Jaech et al., 2024; DeepSeek-AI et al., 2025) may enhance LLMs' ability to consider causal relationships.

682 Finally, some recent work on Sequence Driven Structural Causal Models (SD-SCMs) shows how to simulate counterfactual 683 outcomes and treatment scenarios that are often inaccessible in sensitive datasets (by allowing an LLM to specify structural 684 equations implicitly, given a topological order and a specific prompting structure) (Bynum & Cho, 2024). Similarly to the 685 surrogate public data approaches explored in this paper, the SD-SCM approach does not require access to a downstream 686 private dataset of interest; instead, it only requires a schema over the data to be generated, and a user to specify the prompting 687 structure and topological order over variables (which could be generated e.g., by the first few steps of the Agent procedure 688 given in Figure 4). There may be many potential uses for SD-SCM generated surrogate public data for private causal 689 algorithms; for example, we believe that future work could explore how it can be used to improve the performance of 690 hyperparameter tuning for private causal effect estimators. 691

## 693 C. Limitations

692

The overarching goal of this work is to address a significant barrier to the real-world adoption of differential privacy (Dwork et al., 2016). Enabling more widespread deployment of differential privacy, when appropriate, can promote the protection of individuals' privacy (Cummings et al., 2024a). However, our work does have several risks and limitations. First, there is a risk that LLM memorization may lead to overly optimistic performance estimates (we attempted to mitigate this risk by carefully checking for evidence of memorization, see Appendix F.2.4).

Second, the normative implications of employing LLMs to generate surrogate public data should be carefully analyzed in 700 this context. Recent work by Tramèr et al. (2024) cautions against treating web-scraped LLM training data as "public" or 701 non-sensitive. Traditionally, differentially private algorithms have assumed data is either fully private (restricted) or fully public (freely available and safe to reuse). However, Tramèr et al. (2024) emphasize a messier reality; social media and other sources of personally identifiable information, for example, may be both accessible to language models for training data and 704 contain sensitive information specific to individuals. When an LLM is trained on such data, it may memorize fragments of it; 705 regurgitating these private fragments could be interpreted as a privacy violation. Indeed, even if a final model is fine-tuned 706 under DP constraints, privacy violations may originate from the pretrained model (e.g., a base model memorized private details during pretraining, and a subsequent DP fine-tuning step does not noise those probabilities sufficiently to obfuscate). This undermines trust, as an individual may be told that the entire pipeline is "privacy preserving," yet see their personal 709 data re-emerge in the final model's outputs. 710

We carefully position our work under the paradigm shift identified by Tramèr et al. (2024). Using LLMs to emulate *expert-driven* data-generating processes risks inadvertently exposing sensitive information that is publicly available, as mediated by the LLM. Thus, we propose that best practice is to *report empirical measurements of memorization levels*. We 715 do this by leveraging work by Bordt et al. (2024) on verbatim memorization of tabular data by LLMs; see Appendix F.2.4.

716 Additionally, we report on datasets (see Table 1) that post-date LLM training (for the models we evaluate; see Table 2). 717 Choosing tasks where the LLM's prior knowledge is outdated or non-existent demonstrates performance on truly unseen

718 data (Cheng et al., 2024). We stress the importance of communicating these nuances, and of reporting, to the best of one's

719 knowledge/ability, the empirical level of memorization and the potential LLM data regurgitation risks when presenting these

720 methods.

726 727

728

729 730

731

732

733

734

735

736

737 738 739

741

745

747

749

750

751

752

758

759

761

768 769

721 Finally, given substantial evidence that LLMs encode biases (Gallegos et al., 2024), these biases could be reflected in the 722 generated data – either implicitly in CSV generation or explicitly via the causal relationships in the Agent-based approach. 723 For instance, a stereotypical correlation could persist through pretraining and DP fine-tuning, ultimately resulting in an 724 unfair classifier. We leave a detailed investigation of these issues for future work. 725

## **D.** Preliminaries

We now provide relevant background on differential privacy, large language models, and statical distance metrics.

## **D.1. Differential Privacy**

Differential privacy (DP) ensures that the presence or absence of a single individual's data has only a limited influence on an output statistic; in other words, it restricts how much any single record can affect the outcome of an analysis. To define this, we consider two datasets  $D, D' \in \mathcal{X}^n$ , which are *neighboring* if they differ in at most one data entry. Let  $\mathcal{X}$  denote the universe of records.

**Definition D.1** (Differential Privacy (Dwork et al., 2016)). An algorithm  $\mathcal{M} : \mathcal{X}^n \to \mathbb{R}$  satisfies  $(\varepsilon, \delta)$ -differentially private if, for every pair of neighboring datasets  $D, D' \in \mathcal{X}^n$ , and for every subset of possible outputs  $\mathcal{S} \subseteq \mathbb{R}$ ,

$$\Pr[\mathcal{M}(D) \in \mathcal{S}] \le e^{\varepsilon} \Pr[\mathcal{M}(D') \in \mathcal{S}] + \delta$$
.

740 The  $\varepsilon$  parameter is considered the leading privacy parameter. ( $\varepsilon, \delta$ )-DP is also referred to as *approximate differential privacy*. When  $\delta = 0$ , i.e.,  $(\varepsilon, 0)$ -DP, we refer to it as *pure differential privacy* and denote it with  $\varepsilon$ -DP. 742

#### 743 The following definition of public data is inspired by Ben-David et al. (2023). 744

**Definition D.2** (Public Data). A dataset  $\hat{D} \in \mathcal{X}^m$  is *public* if incorporating it into any computation does not incur additional privacy loss. That is, for any sensitive dataset  $D \in \mathcal{X}^n$  and for every  $(\varepsilon, \delta)$ -differentially private mechanism  $\mathcal{M}$ , the privacy 746 guarantee is identical whether  $\hat{D}$  is used or not, i.e.,  $\mathcal{M}(D, \cdot)$  and  $\mathcal{M}(D, \hat{D})$  both satisfy identical  $(\varepsilon, \delta)$ -differential privacy guarantees. 748

### **D.2.** Large Language Models

Large Language Models (LLMs) are trained to generate sequences of tokens by modeling the probability of the next token given its preceding context (Devlin et al., 2019). Let V be the vocabulary of tokens. Formally, given a sequence of tokens  $\mathbf{x} = (x_1, x_2, \dots, x_n) \in V^n$ , a generative language model estimates,

$$\Pr[x_1, x_2, \dots, x_n] = \prod_{i=1}^n \Pr[x_i \mid x_1, \dots, x_{i-1}].$$

In other words, a language model is a function  $f: V^n \to \mathcal{P}(V)$ ;  $f(\mathbf{x})$  maps a sequence to a probability distribution over the vocabulary V of possible tokens for the next token, where  $\mathcal{P}(V)$  is the space of probability distributions over V. This autoregressive formulation enables, for example, the generation of new samples from a tabular distribution when prompted 760 with known samples (Borisov et al., 2023; Zhao et al., 2023).

#### **D.3. Statistical Distance Metrics** 763

764 We now introduce metrics for comparing probability distributions and datasets used throughout this paper.

765 **Definition D.3** (Total Variation Distance). For discrete probability distributions P and Q over  $\mathcal{X}$ , the Total Variation 766 Distance (TVD) is defined as: 767

$$\mathsf{TVD}(P,Q) = \frac{1}{2} \sum_{x \in \mathcal{X}} |P(x) - Q(x)|$$

The Total Variation Similarity (TVS) is simply 1 - TVD(P, Q), representing the similarity rather than the distance between distributions. Both TVD and TVS can be naturally extended to datasets by considering the empirical probability distributions induced by the datasets over the universe  $\mathcal{X}$ .

Now we turn to a more specific measurement of disparity between two datasets based on the results of statistical queries.

**Definition D.4** (Linear Query). Given a predicate  $\phi : \mathcal{X} \to \{0, 1\}$  that maps database records to binary values, a linear query  $q_{\phi} : \mathcal{X}^n \to \mathbb{N}_0^+$  is a function that, for a dataset  $D \in \mathcal{X}^n$ , computes:

$$q_{\phi}(D) = \sum_{r \in D} \phi(r)$$

<sup>780</sup> In other words, a linear query counts the number of records in dataset D that satisfy the predicate  $\phi$ .

**Definition D.5** (Workload Error). Given a workload  $W = \{q_1, \ldots, q_k\}$  of linear queries, and a pair of datasets  $D, D' \in \mathcal{X}^n$ , 783 the workload error is defined as:

WError
$$(D, D') = \sum_{q \in W} |q_i(D) - q_i(D')|$$

The *average k-way marginal error* can be defined as a special case of the workload error where the workload W consists of all possible *k*-way marginal queries. For instance, the 3-way marginal error uses all possible triplet combinations of attributes as queries. Assuming datasets of equal size, the average *k*-way marginal error is normalized by both the number of queries in the workload |W| and the size of the datasets |D|:

$$\operatorname{AvgError}_{k\operatorname{-way}}(D,D') = \frac{1}{|W|\cdot |D|}\sum_{q\in W} |q(D)-q(D')|$$

5 where W is the set of all k-way marginal queries, and  $|W| = {d \choose k}$  for a dataset with d attributes.

#### E. Details of Surrogate Public Data Generation 825

#### 826 E.1. Baselines

834

835

836

837

838

839 840

827 Before discussing the LLM-based approach, we present a series of baseline generation processes to systematically evaluate 828 which aspects of public data characteristics are useful for differential privacy tasks: pretraining, hyperparameter tuning, and 829 estimating the privacy-utility trade-off. The baselines differ in statistical structure and in the information available about the 830 private data. 831

#### 832 E.1.1. UNIFORM DISTRIBUTION OVER THE DOMAIN 833

The dimensionality of the data plays a critical role in differentially private algorithms (McKenna et al., 2021; Rosenblatt et al., 2023), as it could affect, for example, the magnitude of noise introduced to satisfy DP or the ratio between signal and that noise (e.g., when tuning data synthesizers like PrivBayes or AIM). This Uniform distribution baseline captures the scenario where we have no prior knowledge about the underlying data distribution beyond the schema itself by using the maximum entropy probability distribution (Jaynes, 1957): for each record, Uniform samples i.i.d. from either the set of possible values (for categorical columns) or the specified range (for continuous columns), both given in the schema.

## **E.1.2.** UNIVARIATE DISTRIBUTION

841 Beyond knowledge of the record domains, organizations and researchers might have access to prior information about 842 the *univariate* distributions of individual columns, either precisely or approximately. This prior knowledge is available 843 in cases where organizations may have released various statistical measures of private data, such as histograms, means, 844 medians, and standard deviations, with or without differential privacy (Rosenblatt et al., 2024a; Hasani et al., 2024). As 845 a facsimile for data generated with knowledge of the distributions along individual columns, the Univariate baseline 846 samples independently from each column according to the empirical univariate distribution drawn directly from the private 847 data. To make this baseline more realistic – assuming only an approximate PDF (e.g., the distribution's "shape") is known – 848 we round the probabilities to two decimal places, normalize to 1, and rescale during sampling. 849

#### 850 E.1.3. ARBITRARY DISTRIBUTION 851

The previous two baselines are limited by *column independence* in their sampling, preventing them from capturing complex 852 statistical structures needed for higher-order analysis and predictive tasks (Rosenblatt et al., 2023). To isolate the role 853 of structural dependencies in our DP auxiliary tasks with surrogate public data, we consider whether only capturing the 854 existence of relationships between columns could make surrogate public data a useful prior. To test this, we generate an 855 arbitrary dataset from a random but structured distribution that adheres to the schema. 856

857 Algorithm 1 details the full Arbitrary baseline procedure; here, we provide a high-level overview of the two-step 858 generation process. First, we construct a random Directed Acyclic Graph (DAG) representing a Bayesian network over 859 the column variables. The DAG is built sequentially, with each new node potentially connecting to any previously added 860 nodes, subject to a maximum in-degree (here we used 5). This ensures a structured yet arbitrary dependency pattern between 861 variables. Second, we parameterize the network by sampling conditional probability tables for each node. For a given 862 node, we use a Dirichlet distribution with concentration parameter  $\alpha = 1$  to generate probability distributions for each 863 configuration of its parent variables. Specifically, for each parent value combination, we sample a categorical distribution 864 from the k-simplex, where k is the cardinality of the node's domain. This yields a distribution with meaningful dependencies 865 (e.g., correlations) while remaining *entirely independent* of the true empirical distribution of the private data. 866

#### 867 **E.2. CSV Direct Generation** 868

We evaluate a direct approach to data generation using LLMs. The generation process involves prompting the LLM to create 869 CSVs – tabular records that adhere to the schema while following specific guidelines (Almeida, 2024). These guidelines 870 instruct the model to ensure realistic value distributions and relationships between fields, maintain real-world patterns and 871 constraints, and incorporate edge cases at frequencies that mirror their natural occurrence. Similarly to the other surrogate 872 public data methods we evaluate, this approach operates without access to the private dataset, relying solely on the LLM's 873 pretrained knowledge. 874

875 To ensure data quality, each generated record is validated against the schema, and only valid records are retained. Due to 876 context window limitations and API constraints, the generation process is executed in multiple batches until the desired 877 number of records is obtained (OpenAI, 2025; Anthropic, 2025; Together AI, 2025). Note that, due to the autoregressive 878 nature of LLMs (see e.g., Appendix D), records within the same generation batch are not sampled independently, in contrast 879

to the baseline methods.

881

## 882 E.3. Agent (State Machine) Approach

As a final approach, we employ a multi-step, Agent-based process to elicit a structural causal model (SCM) from an LLM given only text-based access through prompts and responses. Our goal is to arrive at a coherent directed acyclic graph (DAG) that captures the inter-dependencies among variables in the schema, along with associated structural equations (e.g., the actual distributional parameters, probabilities, etc.). Each step concludes with an automated validation of the LLM's output; so, if any contradictions or omissions are detected, the Agent (implemented as a state machine, see Figure 4) automatically refines our prompt and re-queries the LLM.

889 First, we prompt the LLM to • list out all variables (keys) from the provided schema, ensuring the response exactly matches 890 the schema's variable set. Next, we ask it to 2 propose realistic consistency constraints among these variables; these 891 constraints should capture domain knowledge such as permissible value ranges (e.g., "age must be at least 0") or logical 892 relationships (e.g., "an individual who is 10 years old must have fewer than 10 years of education"). We then instruct 893 the LLM to S identify a subset of variables that can serve as the "root nodes" in a causal graph, typically those deemed 894 exogenous or less likely to be influenced by other variables in the schema. From there, the LLM proposes parent-child 895 relationships **④** from root nodes to non-root nodes, and then **⑤** among all remaining variables, **⑥** ensuring no cycles are 896 introduced so that the final structure is a DAG (which we validate with a graph library to confirm it contains all variables 897 exactly once and remains acyclic (Hagberg et al., 2008)). 898

899 Having obtained a DAG, we prompt the LLM to **1** map each variable to a structural equation that references its parents. 900 For instance, if a node depends on two parents, the LLM might generate a formula specifying a probabilistic distribution 901 conditional on parent values. These structural equations encode marginal distributions for root variables and conditional 902 distributions for their descendants. Sometimes the structural equations are not fully specified (e.g., the probability parameter 903 in Bernoulli distribution is parameterized), so we instruct the LLM to ③ assign values to all parameters. Then, ④ we combine 904 the DAG and structural equations automatically into a single code snippet (we use the Pyro library (Bingham et al., 2019)), 905 which lets us generate synthetic data automatically. Finally, we ask the LLM to amend the Python code to **(**) enforce the 906 range or valid values for each column, and **1** include the constraints elicited at the beginning of the interaction. This entire 907 interaction is a stateful, automatic, closed loop, allowing the LLM to act on its own as an "expert" to design a plausible 908 causal model solely from schema level information (containing short descriptions of each variable), without a need to inspect 909 any real-world sensitive records. 910

To extend this approach from a "single expert" to a "panel of experts," we execute the complete generation workflow 911 multiple times to produce a *collection* of datasets, inspired by prior work on "self-consistency" prompting methods (Wang 912 et al., 2023b). These datasets are then combined to yield a single mixed dataset using two approaches. The first approach, 913 Unif., involves uniform sampling of records across all generated datasets. The second approach, Max Cov., solves the 914 Facility Location submodular problem (Wang & Zhou, 2020) by finding a subset of datasets that maximizes the sum of 915 pairwise Total Variation similarities. This optimization selects a subset of datasets that aims to represent the space of all 916 generated datasets (Wang & Zhou, 2020). Then, similarly to the Unif. approach, we sample records uniformly from the 917 selected datasets. 918

One important advantage of agent-generated SCMs is that domain experts can modify the causal structure, structural
 equations, and constraints based on their expertise, scientific literature, and common sense. We leave this for future work.

- 922
- 923
- 924 925
- 926
- 928
- 929
- 930
- 931
- 932
- 933 934

943

### 944 Algorithm 1 Random Bayesian network generation for the arbitrary dataset.

945 1: **procedure** GENERATERANDOMBN( $S, d_{\max}, \alpha$ ) 946 Input: 947  $\mathcal{S} = \{(v_1, \mathcal{D}_1), \dots, (v_n, \mathcal{D}_n)\}$ : Schema where  $v_i$  is a variable and  $\mathcal{D}_i$  is its domain of possible values 948  $d_{\max}$ : Maximum parent degree 949  $\alpha$ : Dirichlet concentration parameter 950 Output: 951 Bayesian network  $\mathcal{B} = (\mathcal{G}, \Theta)$  where: 952  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ : Directed acyclic graph with nodes  $\mathcal{V}$  and edges  $\mathcal{E}$ 953  $\Theta = \{\theta_{v|\Pi_v} : v \in \mathcal{V}\}$ : Set of conditional probability distributions, where  $\theta_{v|\Pi_v}$  represents 954 the distribution of v given its parent set  $\Pi_v$ 955 Initialization: 956 Extract variables  $\mathcal{V} = \{v_1, \ldots, v_n\}$  from schema  $\mathcal{S}$ 2: 957 Define indexing function  $\phi_v : \mathcal{D}_v \to \{1, \dots, |\mathcal{D}_v|\}$  for each  $v \in \mathcal{V}$ 3: 958 Network Structure Generation: 959 4: Randomly permute the ordering of variables in  $\mathcal{V}$ 960 Initialize edge set  $\mathcal{E} \leftarrow \emptyset$ 5: 961 Initialize parameter set  $\Theta \leftarrow \emptyset$ 6: 962 7: for i = 1 to n do 963 8: Define candidate parent set  $C_i = \{v_1, \ldots, v_{i-1}\}$ 964 Select  $\Pi_i \subseteq C_i$  randomly with  $|\Pi_i| \leq \min(d_{\max}, i-1)$ 9: 965 Add edges  $\{(u, v_i) : u \in \Pi_i\}$  to  $\mathcal{E}$ 10: 966 Parameter Generation: 967 11: Let  $\Omega_{\Pi_i}$  be the set of all configurations of  $\Pi_i$  where each configuration  $\pi \in \Omega_{\Pi_i}$  is a tuple of values 968 Let  $k_i = |\mathcal{D}_{v_i}|$  be the cardinality of variable  $v_i$ 's domain 12: 969 if  $\Pi_i = \emptyset$  then 13: 970  $\theta_{v_i} \sim \operatorname{Dir}(\alpha \cdot \mathbf{1}_{k_i})$ 14:  $\triangleright$  Sample from Dirichlet with symmetric  $\alpha$  parameter 971 15: else 972 16: for all  $\pi \in \Omega_{\Pi_i}$  do 973 17:  $\theta_{v_i|\pi} \sim \operatorname{Dir}(\alpha \cdot \mathbf{1}_{k_i})$  $\triangleright$  Conditional probability distribution of  $v_i$  given parent configuration  $\pi$ 974 end for 18: 975 19: end if 976  $\Theta \leftarrow \Theta \cup \{\theta_{v_i \mid \Pi_i}\}$ 20: 977 21: end for 978 22: return  $\mathcal{B} = ((\mathcal{V}, \mathcal{E}), \Theta)$ 979 23: end procedure 980 981 982 983 984 985

- 987
- 988 989

990 991 992

```
996
        {
997
          . . .
998
          "RELACT": {
999
            "description": "Main labour market activity status",
1000
            "dtype": "int64",
1001
            "values": {
1002
              "1": "Employed"
1003
              "2": "Unemployed",
1004
              "3": "Retired",
1005
              "4": "Student",
1006
              "5": "Unable to work",
1007
              "6": "Doing unpaid social work or charitable activities",
1008
              "7": "Other inactive person"
1009
            }
          },
1011
          "CERTIG": {
1012
            "description": "Degree of disability",
1013
            "dtype": "int64",
1014
            "values": {
1015
              "1": "0-32%",
1016
              "2": "33-44%",
              "3": "45-64%".
1018
              "4": "65-74%",
1019
              "5": "75% or more",
              "6": "Not known"
            }
          },
1023
          "AUDI_7_1": {
1024
            "description": "Has significant difficulty hearing a conversation with
1025
                \hookrightarrow several people without a hearing aid",
1026
            "dtype": "int64",
1027
            "values": {
1028
              "1": "Yes",
1029
               "2": "No"
            }
          },
          . . .
1033
        }
1034
```

1035 1036

Figure 2: Excerpt from the schema of the EDAD dataset (Spanish disability, autonomy, and dependency survey) (InstitutoNacional de Estadística, 2024).

1039

1040 1041

10/1

1043

1045	
1046	
1047	
1048	
1040	
1049	
1050	
1051	
1052	
1053	
1054	
1055	
1056	
1057	
1057	
1058	
1059	
1060	
1061	System: You are an expert in (domain) who generates synthetic data that
1062	System: fou are an expert in {domain} who generates synthetic data that
1063	$\hookrightarrow$ closely mirrors real-world (domain) data. Your goal is to create data
1064	$\hookrightarrow$ that would be indistinguishable from real {domain} records.
1064	
1065	Follow exactly these rules:
1066	1. Only output the CSV data with no additional text or explanations
1067	2 Always include a header row matching the schema exactly
1068	2. Always include a header low matching the schema exactly
1069	5. Strictly addete to the provided schema's data types and possible values
1070	$\hookrightarrow$ for all fields
1070	4. Use comma as the separator
1071	5. Ensure all values and relationships between fields are realistic and
1072	$\hookrightarrow$ statistically plausible
1073	6. Generate diverse data while maintaining real-world patterns and
1074	$\leftrightarrow$ constraints
1075	7 Include occasional adma cases at realistic frequencies
1076	7. Include occasional cuge cases at realistic inequencies
1077	
1078	User: Generate {num_rows} rows of data with these fields:
1070	
1079	{schema}
1080	
1081	
1082	
1083	Figure 3: The prompt template used for CSV generation with an LLM.
1084	
1001	
1005	
1080	
1087	
1088	
1089	
1090	
1091	
1092	
1002	
1023	
1094	
1095	
1096	
1097	
1098	
1099	
+ ~ / /	



Figure 4: State machine for the SCM Agent showing state transitions. Each state can transition to itself upon failure or advance to the next state upon success, following a zigzag pattern.

Dataset	Topic	Features	imesDims		Private Sp	plit	Pu	ıblic Split	
				Name	Size	Published	Name	Size	Publishe
ACS	Census	7	116.640	National	23.006	Sep 2020	Massachusetts	23,006	Sep 202
DAD	Disability	11	2,188,800	2023	1,469	Apr 2024	2020	1,469	Apr 202
VE	Workplace	12	1,924,560	2023	1,400	Apr 2024	2018	837	Dec 201
			Table 2. Lar	ge Language	Models (1	I Ms) used in	n this work		
			Table 2. Lai	ge Language	, widdels (1	LLIVIS) used ii	n uns work		
Name			Provider		Vancia	n		Cu	toff Date
			TTOVILLET		versio	11			ton Dute
GPT-40			OpenAI		gpt-4	u lo-2024-08	3-06	Oc	tober 2023
GPT-40 Claude 3 Llama 3	3.5 Sonnet 3 70B Instru	ict-Turbo	OpenAI Anthropio Meta via	c Together AI	gpt-4 claud	lo-2024-08 le-3-5-sor	3-06 nnet-2024102	Oc 2 Ap rbo De	tober 2023 ril 2024 cember 202
GPT-40 Claude 3 Llama 3.	3.5 Sonnet 3 70B Instru	ıct-Turbo	OpenAI Anthropio Meta via	c TogetherAI	gpt-4 claud Llama	lo-2024-08 le-3-5-sor l-3.3-70B-	3-06 nnet-2024102 -Instruct-Tu	Oc 2 Ap rbo De	tober 2023 ril 2024 cember 202
GPT-40 Claude 3 Llama 3. <b>F. Deta</b> Traditi Surrog Arbi	3.5 Sonnet 3 70B Instru ills of Eval Ional Public ATE Public Dor itrary	Uation I Data	Framewor Sample	c TogetherAI	tility Est.	(1) Parameter 1) P 10 Parameter 10 P 10 P 1	B-06 Inet-2024102 Instruct-Tu Tuning Tuning CLASS Ination	Oc 2 Ap rbo De odel Pretrainin, SIFICATION: ablic Data so rd. or Alt.)	tober 2023 ril 2024 cember 202 g SETUP $\xrightarrow{D} M_{pre}$
GPT-40 Claude 3 Llama 3. F. Deta Traditi Surrog. Surrog.	3.5 Sonnet 3 70B Instru ils of Eval Ional Public ATE Public Do itrary Do SV GPT-40	ULLLMS	Framewor Sample Inivariate	c TogetherAI	tility Est.	( ) Parameter $( ) Parameter$	Tuning	Oc 2 Ap rbo De odel Pretrainin SIFICATION: ablic Data <u>so</u> rd. or Alt.)	tober 2023 ril 2024 cember 202 SETUP $\xrightarrow{D} M_{pre}$



 $\frac{1207}{1208}$  privacy-utility estimation – and one to classification model pretraining.

Category	Metric	Description
Marginals	Total Variation Distance	Distance between the joint distributions of the original and synthetic datasets.
	Max 3-Way Marginal Error	Maximum absolute difference error for 3-way marginals be tween original and synthetic datasets, normalized by dataset size.
	Avg. 3-Way Marginal Error	Average absolute difference error for 3-way marginals between original and synthetic datasets, normalized by dataset size and query count.
	Max Binarized Marginal Error	Maximum absolute difference error for 3-way marginals after thresholding continuous variables to binary values, normalized by dataset size.
	Avg. Binarized Marginal Error	Average absolute difference error for 3-way marginals after thresholding continuous variables to binary values, normalized by dataset size and query count.
Correlations	Max Pearson Correlation Diff	Maximum absolute difference between Pearson correlation coefficients of original and synthetic datasets.
	Avg. Pearson Correlation Diff	Average absolute difference between Pearson correlation coef- ficients of original and synthetic datasets.
	Max Cramer's V Diff	Maximum absolute difference between Cramer's V correlation coefficients of original and synthetic datasets.
	Avg. Cramer's V Diff	Average absolute difference between Cramer's V correlation coefficients of original and synthetic datasets.
Classification	Error Rate Diff	Difference in classification error rates between models trained on original vs. synthetic data and evaluated on the same test set.
	AUC Diff	Difference in Area Under the ROC Curve (AUC) between models trained on original vs. synthetic data and evaluated on the same test set.

## Do You Really Need Public Data? Surrogate Public Data for Differential Privacy on Tabular Data

Table 3: Overview of quality evaluation metrics for a synthetic dataset against the original dataset. All metrics range from 0to 1, with lower values indicating better synthetic data quality.

1249 1250

## 1251 **F.1. Tasks**

### 1252 F.1.1. TASK 1: MODEL PRETRAINING FOR CLASSIFICATION

1253 A common practice in machine learning with DP is to first *pretrain* a model on public data (incurring no privacy loss) before 1254 *fine-tuning* it privately on sensitive data (using e.g., DP-SGD, incurring fixed  $(\varepsilon, \delta)$ -privacy loss). We apply this method to 1255 evaluate surrogate public data for binary classification tasks on tabular data (recall that this is a less common setting than 1256 public pretraining with image data (Ganesh et al., 2023; Thaker et al., 2024), due to a general lack of publicly available 1257 priors for tabular datasets).

We divided public and private datasets into train, validation, and test subsets using a 72 : 8 : 20 ratio, and used an FTTransformer deep neural attention based classification model architecture (Gorishniy et al. (2021); see Appendix F.3.1 for more details). Our classification evaluation framework follows three steps: (1) standard pretraining, updating model weights with gradients calculated from (surrogate) public data; (2) DP fine-tuning on private training data; and (3) performance assessment on the private test data. For comparison, we include a control condition that omits the pretraining phase. We

measure classification performance using AUC metric and ensure balanced datasets by downsampling the majority class to

match the minority class size. We also consider an AUC Advantage metric, which we define as the difference in AUC
between models *with* public data pretraining and a model *without* pretraining, which directly quantifies the incremental
benefit provided by pretraining before private finetuning.

To account for the multiple hyperparameters in both pretraining and fine-tuning stages, we conduct a comprehensive grid search, further running each configuration 10 times to mitigate variations inherent to differential privacy training and model initialization. We analyze results using two complementary approaches: averaging performance across all hyperparameter combinations, and simulating a real-world scenario by selecting the optimal pretraining hyperparameters based on public validation data before averaging results across fine-tuning hyperparameters. Refer to Appendix F.4 for the complete hyperparameter space details.

1275

### 1276 F.1.2. TASK 2: HYPERPARAMETER TUNING FOR SYNTHETIC DATA

Hyperparameters play an important role in training machine learning models, especially when differential privacy is involved
(Ponomareva et al., 2023). While selecting the best performing hyperparameters in the non-private setting can be done with
many model training runs using a validation split or cross-validation, this is not feasible in a straightforward manner with
differential privacy due to the privacy loss incurred on each run. Public data may be helpful in this case, allowing researchers
to run multiple experiments without consuming the privacy loss budget (Iyengar et al., 2019; Cattan et al., 2022).

To assess the usefulness of surrogate public data for this DP auxiliary task, we run a large-scale DP synthetic data evaluation across multiple dimensions: (1) datasets (including private and public splits, and various public data surrogates); (2) privacy loss budget  $\varepsilon$ ; (3) different DP synthetic data generators (see Section F.3.2; GEM (Liu et al., 2021b), AIM (McKenna et al., 2022), PrivBayes (Zhang et al., 2014)); and (4) their associated hyperparameter spaces. For each configuration, we fit a synthetic data generator and produce a synthetic dataset of the same size as the original, private data. We then evaluate across a variety of metrics, which fall into three general categories: marginal-based metrics, correlational metrics, and classification-based metrics, as shown in Table 3.

We conduct our analysis (1) *per synthetic data generator*, because each has a different hyperparameter space and different sensitivity to changes in hyperparameter configuration; (2) *per metric*, because the best-performing hyperparameter is defined with respect to a specific metric; and (3) *per privacy loss budget*  $\varepsilon$ . We quantify the degradation in performance when using the synthetic generator *on the private data* by comparing the best hyperparameter setting that we would have chosen *with the private data* (i.e., the optimal case) relative to the hyperparameters we would have chosen *with each of the* (*potentially surrogate*) *public datasets*.

1296 To aggregate the usefulness of public data in choosing hyperparameter configurations across different evaluation metrics, we computed a *relative* performance degradation metric for each configuration. Concretely, for every private synthetic 1298 data generator, privacy level  $\varepsilon$  and dataset (ACS, EDAD, and WE), we first identified the hyperparameter configuration 1299 that yielded the best performance on the private reference dataset (i.e. the real data). We then determined, for each 1300 candidate surrogate public dataset (and the regular public data), the hyperparameter configuration that would have been 1301 chosen based solely on its corresponding performance. Our benchmark quantifies degradation as the relative difference 1302 between the performance achieved by the surrogate-chosen hyperparameters on the private reference and the optimal 1303 performance on the reference dataset (measured as either absolute error or percent degradation, depending on the metric). 1304 We conducted this process independently for each metric – across classification, correlation, and marginal-based metrics. We 1305 averaged across multiple experimental seeds to obtain aggregate performance with standard error; we then conducted a Pareto 1306 frontier analysis (Ehrgott, 2005) across the frontier defined by aggregating into the three metric categories: classification, 1307 correlation and marginal-based metrics.

1308 1309

## 1310 F.1.3. TASK 3: PRIVACY-UTILITY ESTIMATION FOR SYNTHETIC DATA

Understanding the privacy-utility trade-off of a mechanism for a specific private dataset is *extremely* useful for producing a
differentially private release in the real world (Rosenblatt et al., 2024a). For example, it may provide guidance on setting the
privacy loss budget by exposing its impact on the fidelity of private synthetic data (e.g., (Abowd et al., 2023; Hod & Canetti,
2025)).

<sup>1315</sup> In this task, we evaluate how well a public dataset – either traditional or surrogate – can estimate the privacy-utility *curve* for  $\frac{1316}{1316}$  and still the maximum tip in a surrogate dataset – either traditional or surrogate – can estimate the privacy-utility *curve* for  $\frac{1316}{1316}$  and still the maximum tip in a surrogate – can estimate the privacy-utility *curve* for  $\frac{1316}{1316}$  and  $\frac{1316}{1316}$ 

each utility metric. This experiment is, in a sense, the "dual" of the hyperparameter tuning task described in the previous

- section: here, we compare the privacy-utility curve computed on the public data with the curve obtained on the private
- data. To mimic real-world usage, we run the DP mechanism with the best-performing hyperparameters determined from the

public data (Table 3), selecting the optimal configuration independently at each tested  $\varepsilon$  value.

For each dataset, synthetic data generator, and evaluation metric, we created both public-based and private-based curves over

a range of privacy loss budgets  $\varepsilon$ . To aggregate the results across different evaluation metrics, we first compute, for each

metric group (classification, correlation, and marginals) and each synthesizer (PrivBayes, AIM, and GEM), an aggregated

- performance value that is the average "chosen value" across all metrics in that group. For a given synthesizer and for each  $\varepsilon$ , we group the results by dataset and reference dataset and then pivot these averages so that each row corresponds to a dataset
- 1326 and each column to an  $\varepsilon$  level. This representation enables us to generate line plots to visually assess the similarity between
- and each column to an enever. This representation enables us to generate line plots to visually a performance curves (see, e.g., Figure 28 for an example with the PrivBayes synthesizer).
- 1329 Since the line plots alone are insufficient to quantify aggregate closeness, we compute both  $\ell_1$  and  $\ell_2$  distances between
- 1330 each pair of curves. The  $\ell_1$  distance is more interpretable being in the same units as the evaluation metric while the  $\ell_2$
- 1331 distance is less sensitive to outliers. We average the  $\ell_1$  and  $\ell_2$  distances across the different metric categories (weighting
- each category equally). To reduce variability, each configuration is run 10 times. Finally, we perform a Pareto frontier
- 1333 analysis across both  $\ell_1$  and  $\ell_2$  distances for each dataset (Ehrgott, 2005).
- 1334
- 1335 **F.2. Datasets**
- 1336 F.2.1. ACS
- The ACS data excerpt was released by the US Census Bureau in September 2020 and provided by the NIST CRC to assess synthetic data generation methods. We designated the "National" dataset (27,254 records) as the private split and the "Massachusetts" dataset (7,634 records) as the public split. Since the differential privacy synthetic data generators assessed in this project are primarily designed for categorical data, we used the "demographic" subset containing 7 categorical
- <sup>1341</sup> features provided by NIST CRC. After removing records with missing values, we retained 23,006 and 6,514 records for the <sup>1342</sup> private and public splits, respectively. The public split was up-sampled to match the size of the private split. For a complete
- 1342 private and public splits, respectively. The public split was up-sampled to match the size of the private split. For a complete 1343 description of the dataset and its curation refer to its documentation (Task et al. 2023)
- $^{1343}$  description of the dataset and its curation, refer to its documentation (Task et al., 2023).

## <sup>1345</sup> F.2.2. EDAD

1346 The EDAD (Survey on Disability, Personal Autonomy and Dependency Situations) datasets were released by the Spanish 1347 National Statistics Institute (INE) in April 2022 and April 2024, containing responses from their 2020 (164,254 records) 1348 and 2023 (12,518 records) surveys respectively. We designated the 2023 survey responses as the private split and the 2020 1349 survey responses as the public split. Since our synthetic data generators are primarily designed for categorical data, we used 1350 a subset of 11 categorical features from both surveys. After removing records with missing values, we retained 8,922 and 1351 1,469 records for the private and public splits, respectively. The private split was down-sampled to match the size of the 1352 public split. For a complete description of the datasets and their curation, refer to the documentation given by (Instituto 1353 Nacional de Estadística, 2024). 1354

# <sup>1355</sup> F.2.3. WE

1356 The Workplace Equity Survey datasets (WE) consist of responses from two global surveys conducted in 2018 (released 1357 December 2019) and 2023 (released April 2024) by the Coalition for Diversity and Inclusion in Scholarly Communications 1358 C4DISC). We designated the 2023 survey responses (1,755 records) as the private split and the 2018 survey responses (1,182 1359 records) as the public split. Since our synthetic data generators are primarily designed for categorical data, we used a subset 1360 of 12 categorical features from both surveys. In this dataset, we kept the missing values as another category. We retained 1361 837 and 1,400 records for the public and private splits, respectively, and no upsampling or downsampling was done. The 1362 slight reduction in records is due to filtering response with high levels of missingness and only using respondents from the 1363 top 10 most common country affiliations in the survey (to reduce dimensionality). For a complete description of the datasets 1364 and their curation, refer to their documentation (Spilka et al., 2020; Lemieux et al., 2024). 1365

# 1366 F.2.4. DATASET MEMORIZATION BY THE LLMS

Recent research has highlighted growing concerns that, because LLMs are exposed to benchmark data from the internet during training, their performance those and other benchmarks may be inflated when assessing performance post-training (Magar & Schwartz, 2022; Golchin & Surdeanu, 2024; Roberts et al., 2024; Xu et al., 2024; Dong et al., 2024). For example, it is well known that LLMs have a large capacity for training data memorization (Carlini et al., 2021; Kandpal et al., 2022;

- <sup>1371</sup> Carlini et al., 2023); this is one mechanism by which they could "hack" existing benchmarks, by simply memorizing the
- examples and their answers. This memorization consideration is particularly relevant for our experimental setup, where we
- utilize LLMs to generate records both directly and indirectly. Thus, any prior exposure to our evaluation datasets (ACS,

EDAD, and WE) could significantly impact model performance in our evaluations (of particular concern is exposure to the split of these datasets that we consider *private* in our evaluations, e.g., the national version of the ACS dataset). We address this memorization concern through two mitigation strategies.

First, we considered the temporal relationship between dataset releases and *model knowledge cutoff dates* when selecting

two of our datasets for evaluation. Namely, the private splits of EDAD and WE were released in April 2024, which is

later than the knowledge cutoff dates of most models used in our study (Table 2): GPT-40 (October 2023), Llama 3.3 70B
 (December 2023), and Claude 3.5 Sonnet (April 2024). While there is a one-month overlap with Claude, the analysis of

(December 2023), and Claude 3.5 Sonnet (April 2024). While there is a one-month overlap with Claude, the analysi
 Cheng et al. (2024) suggests that the effective knowledge cutoff dates of LLMs typically *precede* their reported dates.

Second, we executed the LLM memorization assessment methodology proposed by (Bordt et al., 2024); they provide an extensive package & benchmark for LLM memorization detection *specific to tabular data*. We ran their assessment across all private and public splits. In the data generation tests from (Bordt et al., 2024) – the most relevant to our setting – both header tests (generating the first few rows) and row completion tests (generating random-location rows) indicated *no evidence* of record-level memorization by any of the three LLMs across all datasets. Refer to Figure 6 for an example of the header test results for the ACS dataset with Claude 3.5 Sonnet.

- 1390
- Additional tests examining an LLM's metadata knowledge of tabular datasets, rather than record generation capabilities,
   revealed varying levels of dataset familiarity. The models unsurprisingly demonstrated strong familiarity with ACS, but
   limited knowledge of EDAD and minimal recognition of WE. This pattern aligns with the relative public visibility of these
- datasets: ACS is a core and official product of the US Census, EDAD is an official product of the Spanish National Statistics Institute, and WE is a small-scale survey conducted by a coalition of professional and trade organizations.
- When provided with header columns and the first few rows, all models successfully identified the name of the ACS dataset, and sometimes could identify the EDAD dataset name (where the 2020 public split consists of multiple raw files). However, for the WE dataset, even when given headers and first rows, no model generated the correct dataset name – instead, they provided thematically related names such as "work-life-and-career-survey" and "publishing-industry-diversity-survey." We hypothesize that this pattern emerges from the survey questions themselves serving as column names, which inherently reveal the overall topic of the survey (e.g., "How long have you worked in publishing and/or related industries?").
- We observed similar patterns regarding column name completion. When given the dataset name and the first few features, all
  models failed to generate the correct column names for both EDAD and WE datasets. For ACS, the models could generate
  some of the column names, but not in the correct order. We hypothesize that this is due to the fact that the ACS datasets we
  used were sub-sampled, modified, and adopted from the US Census release by NIST.

# 1408 **F.3. Private Mechanisms**

#### 1409 F.3.1. CLASSIFICATION

Differentially private pretraining is usually conducted in domains where strong, publicly available priors with matching data-dimensionality are available (e.g., text or image data). In these fields, neural transformer models dominate (Wolf et al., 2020; Khan et al., 2022).

1414 For an adequate analog to this space in the tabular setting, we consider an FTTransformer model (Gorishniy et al., 2021), 1415 which is a transformer based architecture for tabular data classification. FTTransformer has demonstrated strong performance 1416 against established powerful gradient boosting approaches such as XGBoost (Chen & Guestrin, 2016). Its effectiveness 1417 stems from specialized data transformations that mitigate information loss in transformer-based attention layers (Gorishniy 1418 et al., 2022). Prior work shows how simple it can be to adapt FTTransformer to the private setting (Rosenblatt et al., 2024b) 1419 by making minor modifications to its architecture to support DP-SGD (Abadi et al., 2016). Importantly, it can also be easily 1420 pre-trained with public data through standard gradient updates before private training. The differentially private variant of 1421 FTTransformer is  $(\varepsilon, \delta)$ -DP, for which we set  $\delta = 10^{-5}$ . 1422

# 1423 F.3.2. DATA SYNTHESIS

We considered three representative state-of-the-art private data release methods: PrivBayes (Zhang et al., 2014), GEM (Liu et al., 2021c) and AIM (McKenna et al., 2022). Each of these synthesizers follows the "Select-Measure-Project" paradigm, in that they *privately* select statistical queries (marginals or correlations) to run on a sensitive distribution, *privately* measure these queries, and then as *post-processing* project these measurements onto a synthetic distribution (from which we can draw arbitrary samples) that approximates the original, sensitive distribution.

1430	
1421	PUMA,AGEP,SEX,MSP,HISP,RAC1P,NOC,NPF,HOUSING_TYPE,OWN_RENT,DENSITY,INDP,
1431	INDE CAT EDIT PINCE PINCE DECILE POVELD DVET DREM DEHY DEVE DEAR PWGTP
1432	
1433	WGIP
1424	01-01301,18,2,6,0,9,N,N,3,0,2731.2,N,N,7,0.0,0,N,N,2,2,2,2,79,0
1434	01-01301,27,1,6,0,1,N,N,3,0,2731.2,3291,4,7,15400.0,4,116,N,2,2,2,2,5,0
1435	01-01301,74,2,3,0,2,N,N,2,0,2731.2,N,N,9,12900.0,3,N,N,2,1,2,2,19,0
1436	01-01301.22.1.6.0.1.N.N.3.0.2731 2.N.N.7.0 0.0.N.N.2.2.2.2.10.0
1437	
1438	01-01301,18,2,6,0,1,N,N,3,0,2/31.2,N,N,/,0.0,0,N,N,2,2,2,15,0
1420	01-01301, 52, 2, 1, 0, 1, N, N, 1, 1, 2731.2, 7860, 8, 10, 52000.0, 8, 433, N, 2, 2, 2, 2, 25, 0
1439	01-01301,54,1,1,0,1,N,N,1,1,2731.2,7860,8,10,55000.0,8,458,N,2,2,2,2,25,0
1440	01-01301, 20, 2, 6, 0, 1, N, N, 3, 0, 2731.2, N, N, 7, 35400.0, 0, N, N, 2, 2, 2, 2, 12, 0
1441	01-01301.48.2.1.0.1.N.N.1.1.2731.2.8680.9.10.45000.0.7.375.N.2.2.2.2.20.0
1442	
1//3	01-01301,49,1,1,0,1,N,N,1,1,2/31.2,7860,8,9,48000.0,7,400,N,2,2,2,2,20,0
1443	01-01301, 15, 1, 6, 0, 1, N, N, 3, 0, 2731.2, N, N, 6, 9300.0, 0, N, N, 2, 2, 2, 2, 18, 0
1444	01-01301, <b>45</b> , <b>2</b> , <b>1</b> , 0, <b>1</b> , N, N, <b>1</b> , <b>1</b> , 2731. 2, N, N, <b>5</b> , <b>27</b> 860. 0, <b>8</b> , N, N, <b>2</b> , <b>1</b> , <b>0</b> , <b>2</b> , <b>142</b> 0
1445	01 01 01 27 1 250 1 N 2 2 2 27 1 8 0
1446	01-01.01,27,1,500,1,10,2,2,2,27,1.0,0

Figure 6: The header test output on the ACS dataset on Claude 3.5 Sonnet. The LLM is prompted with the column names as well as a few *first rows* of the dataset (black), and its completion is presented. The output is colored according to its Levenshtein string distance compared to the original records: correct, incorrect, and missing. We observe that the LLM failed to reproduce the header, as many errors occur within columns with variability.

PrivBayes builds a Bayesian network (BN) and adds noise to all k-way correlations to ensure differential privacy. Despite having been published in 2017, PrivBayes is still considered state-of-the-art and was chosen to produce the differentially private release of the Israel National Live Birth Registry (Hod & Canetti, 2025). GEM parameterizes a neural model to represent a synthetic distribution that approximates the true distribution by minimizing a linear query error based loss (with linear queries implemented as k-way marginals, where by default k = 3). AIM relies on the Private-PGM graphical model (McKenna et al., 2021) to parameterize the underlying distribution, and utilizes an iterative process to take advantage of higher values of  $\varepsilon$ . Both AIM and GEM are considered the state-of-the-art approaches to generating private synthetic data (Tao et al., 2021; Rosenblatt et al., 2023). Outside of these methods, we acknowledge that many other methods exist for generating DP data (Dwork et al., 2009; Hardt et al., 2012; Vietri et al., 2020; McKenna et al., 2023; Xu et al., 2019; Rosenblatt et al., 2020; Avdöre et al., 2021; Cai et al., 2021), but we believe that PrivBayes, GEM and AIM are a representative set of what can be currently considered state-of-the-art. 

PrivBayes and GEM are  $\varepsilon$ -DP, whereas AIM is ( $\varepsilon$ ,  $\delta$ )-DP, for which we set  $\delta = 10^{-9}$ . All three methods come with hyperparameters that need to be tuned. Detailed lists of hyperparameters per-synthetic data generator, and their associated values, are given in Appendix F.4.

### **F.4. Hyperparameter Spaces**

Table 4: Hyperparameters fo	FTTransformer Classifier
-----------------------------	--------------------------

Hyperparameter	Description	Values	
pre_num_epochs pre_batch_size pre_lr	Number of epochs for pre-training Batch size for pre-training Learning rate for pre-training	$ \begin{array}{c} \{1,9\} \\ \{32,128\} \\ \{3 \times 10^{-4}, 3 \times 10^{-5}\} \end{array} $	
dp_num_epochs dp_batch_size dp_lr	Number of epochs for differential private fine-tuning Batch size for differential private fine-tuning Learning rate for differential private fine-tuning	$20 \\ 128 \\ \{3 \times 10^{-3}, 3 \times 10^{-4}\}$	

#### Table 5: Hyperparameters for GEM

Hyperparameter	Description	Values	
k	Maximum degree of measured marginals	$\{2,3\}$	
Т	Number of iterations	$\{50, 100\}$	
alpha	Learning rate	$\{0.1, 0.5\}$	
ema_weights_beta	EMA weights coefficient	$\{0.1, 0.9\}$	

#### Table 6: Hyperparameters for AIM

Hyperparameter	Description	Values
degree	Maximum degree of measured marginals	$\{2,3\}$
rounds	Number of iterations	$\{20, 40\}$

#### Table 7: Hyperparameters for PrivBayes

Hyperparameter	Description	Values
theta	SNR heuristic to set max node degree	$\{2, 8, 32, 64\}$
epsilon_split	Prop. of privacy budget allocated to structure learning	$\{0.1, 0.5, 0.75\}$



Figure 7: Task 1 – Pretraining: (a, b) Comparing the mean AUC on the *test* subset *private* split for the **Pretraining** model vs. 1561 the Fine-tuned model, grouped by generation method (mean calculated across DP finetuning parameter space when the best 1562 configuration is chosen with the validation subset of the public split for the pretraining step, across 10 runs). Note how the 1563 starting point of model AUC differs, while the improvement from private finetuning (i.e. the increase in AUC) is relatively 1564 stable. Task 2 – Hyperparameter tuning: (Tables 1, 2 and 3) show some Pareto frontiers for the *performance degradation* 1565 metric when hyperparameter tuning using (surrogate) public data methods for tuning relative to tuning on private data. Note 1566 how CSV and Agent methods are competitive with tuning on the regular public data. See Section 4 and Appendix G for 1567 complete results. Note that we adopt the Olympic medal convention in each table in our paper: gold, silver and bronze 1568 cells signify first, second and third best performance, respectively. 1569

# 1570 G. Additional Result Discussion

In this section, we present the results of our evaluation framework (Section 3) for the following DP auxiliary tasks:
 pretraining (Section G.1), hyperparameter tuning (Section G.3), and estimating the privacy-utility trade-off (Section G.4).
 Appendix H provides additional results details.

1575 All of the experiments were done with  $\varepsilon \in \{1, 2, 4, 8, 16\}$ , and each hyperparameter configuration (Appendix F.4) was run 1576 10 times.

## <sup>1578</sup> G.1. Results for Task 1: Pretraining for DP Classification

In our analysis, the best pretraining hyperparameter configuration was selected based on the public validation subset (see
 Figure 15 and Table 22a for full hyperparameter averaging results, which show similar trends).

**EDAD and WE.** Overall, we find strong evidence that LLM-based methods – both CSV and Agent surrogate public data generation (particularly with Claude 3.5 Sonnet) – offer a competitive alternative to traditional public data. Figure 8 presents our experimental results on the WE and EDAD ( $\varepsilon = 1$ ), demonstrating how pretraining on the surrogate public data can vastly improve the starting point of model performance.

Figure 12 shows a diminishing pretraining advantage when increasing  $\varepsilon$  for both EDAD and WE. This is an expected behavior: high epsilon allows for the extraction of more signal from the private dataset, and may reduce the usefulness of public data, regular or surrogate (Thaker et al., 2024).

<sup>1590</sup> Under a more granular analysis, the EDAD dataset benefits substantially from pretraining, with average AUC advantages per method ranging from 0.09 to 0.19. Here, the traditional public dataset delivers the highest improvement across  $\varepsilon$ values. When aggregated by generation method, CSV-based methods perform slightly worse than the regular public dataset,

followed by the Agent-based method. A more careful examination of surrogate approaches in Table 8b reveals that the



Figure 8: Mean AUC on the test subset of the private dataset split for the pretraining model and the fine-tuned model, grouped by generation method. The mean is calculated across the DP finetuning hyperparameter space when best pretraining hyperparameter configuration is chosen for the pretraining step, with 10 runs per hyperparameter configuration.

1626 CSV (Claude) (AUC advantages ranging from 0.07-0.17) and CSV (Llama) (ranging from 0.08-0.17) perform on par 1627 with or slightly worse than the regular public data (ranging from 0.09-0.19). For example, at  $\varepsilon = 1$ , the AUC advantages of 1628 traditional public data, CSV (Claude), CSV (Llama) are 0.19 and 0.17, respectively. As expected, pretraining with 1629 baselines (Uniform and Univariate) and Arbitrary yields almost no benefit, because they contain essentially no 1630 signal about the relationship between the target variable and the features in the classification task.

The WE dataset exhibits trends similar to EDAD. Although the traditional public dataset achieves the best advantage at  $\varepsilon = 2$ , its performance is not consistently top-ranked across all privacy levels. In fact, for  $\varepsilon = 1, 4, 16$ , it is not in the top three. Notably, the two Claude Agent-based variants have the best performance across most  $\varepsilon$  values, with AUC improvements ranging from 0.07 to 0.21.

#### 1636 1637 G.2. ACS and the Role of Dataset Size

1625

For the ACS dataset, we do not observe any benefit from pretraining, either with traditional or surrogate public data (e.g., as Figure 13b shows for  $\varepsilon = 1$ ). However, **a follow-up analysis reveals that this is due to the relatively large size of the dataset.** Dataset size is a key factor in differentially private mechanisms, as it directly influences the noise level added to achieve a specific level of privacy protection (Dwork et al., 2016). The relatively large size of the ACS dataset partly explains why the benefit of regular public data pretraining appears marginal in, e.g., Table 8a; as privacy sensitivity scales inversely with dataset size, when the private dataset is sufficiently large, the magnitude of noise necessary for a DP guarantee decreases.

To investigate this effect, we repeated the full pretraining experiment on four ACS subsets obtained by subsampling at 5%, 10% - 20% and 50%. In these superiments are formed at the AUC advertises of a subset of multiplets is more than the base of a subset of multiplets in the superimentation.

10%, 20%, and 50%. In these experiments we focus on the AUC advantage at  $\varepsilon = 1$ , where the benefit of public data is most pronounced. Figure 13b (in Appendix G.1) shows that with 5% subsampling, the ACS dataset exhibits a similar pattern of

performance to the one we found with the EDAD and WE datasets.



Figure 12: Mean AUC Advantage of the DP model after pretraining, grouped by generation method. The mean is calculated
across the DP finetuning hyperparameter space when best pretraining hyperparameter configuration is chosen for the
pretraining step, with 10 runs per hyperparameter configuration.

In fact, the LLM-based methods (using Claude Sonnet 3.5) outperformed the traditional public dataset. Figure 14 presents the relationship between the (subsampled) dataset size and the AUC advantage per generation method category to  $\varepsilon = 1$ . As we examine smaller datasets, the differences we observe align with results on the EDAD and WE datasets. Both the CSV and Agent surrogate datasets perform on average similarly to traditional public data. This observation may also help

1669 CSV and Agent surrogate datasets perform on average similarly to traditional public data. This observation may also help 1670 explain the negative findings reported by Swanberg et al. (2025) regarding the use of LLM-generated public data for DP

 $_{1671}$  synthetic data generation on the Adult dataset, which is substantially larger than some of the datasets considered here. We

1672 hypothesize that with smaller datasets, LLM-generated public data surrogates could provide some benefit in pretraining

1673 differentially private data synthesizers, but leave a closer examination of that DP auxiliary task to future work.



1751

1749 1750

Figure 14: Mean AUC Advantage of the DP classification model with  $\varepsilon = 1$  after pretraining for each subsampled dataset, 1753 grouped by generation method category. The mean is calculated across the DP finetuning hyperparameter space when best 1754 pretraining hyperparameter configuration is chosen for the pretraining step, with 10 runs per hyperparameter configuration. 1755

Subsampling

0.0

-0.05.05.1

Agent Without pretraining

- 1756 1757
- 1758
- 1759

Table 8: Mean AUC Advantage (AUC in parentheses) of the DP model after pretraining, grouped by generation method. The
 mean is calculated across the DP finetuning hyperparameter space *when the best pretraining hyperparameter configuration is chosen* for the pretraining step, with 10 runs per hyperparameter configuration.

1764		(a) ACS				
1765	Method	$\varepsilon = 1$	$\varepsilon = 2$	$\varepsilon = 4$	$\varepsilon = 8$	$\varepsilon = 16$
1766	Without pretraining	.00 (.74)	.00 (.74)	.00 (.74)	.00 (.75)	.00 (.75)
1767	Public	.01 (.75)	.01 (.76)	.01 (.76)	.00 (.75)	.00 (.75)
768	Baseline (Domain)	= 03(71)	= 03(71)	- 03 (71)	- 03 (72)	- 05 ( 70)
769	Baseline (Domain) Baseline (Univariate)	01 (.73)	.00 (.74)	03 (.71)	02 (.73)	.00 (.75)
770	Arbitrary	.00 (.74)	.01 (.75)	.00 (.74)	.00 (.75)	.00 (.75)
.771	CSV (Claude 3.5 Sonnet)	.01 (.74)	.01 (.75)	.01 (.75)	.00 (.75)	.01 (.76)
.772	CSV (GPT-4o)	.01 (.74)	.01 (.75)	.01 (.76)	.00 (.75)	.01 (.76)
.773	CSV (Llama 3.3 70B)	.01 (.75)	.01 (.75)	.01 (.75)	.00 (.75)	.01 (.76)
774	Agent (Claude 3.5 Sonnet, Unif.)	.01 (.74)	.01 (.75)	.01 (.75)	.00 (.75)	.00 (.75)
775	Agent (Claude 3.5 Sonnet, Max Cov.)	.01 (.74)	.00 (.75)	.01 (.75)	.00 (.75)	.00 (.75)
776	Agent (GPT-40, Unif.)	.00 (.74)	.00 (.75)	.01 (.75)	.00(.75)	.00 (.75)
777	Agent (Upri-40, Max Cov.)	.00(.74)	.00 (.74)	.00(.75)	.00 (.75)	.00 (.75)
778	Agent (Llama 3.3.70B, Max Cov.)	01(.73)	01(75)	01(.75)	00(.75)	01 (76)
770	Agent (All, Unif.)	.01 (.75)	.01 (.75)	.01 (.75)	.01 (.76)	.01 (.76)
780	Agent (All, Max Cov.)	.01 (.74)	.01 (.75)	.01 (.75)	.00 (.75)	.00 (.75)
781						
782		(b) ED.	AD			
783	Method	$\varepsilon = 1$	$\varepsilon = 2$	$\varepsilon = 4$	$\varepsilon = 8$	$\varepsilon = 16$
784	Without pretraining	.00 (.65)	.00 (.69)	.00 (.71)	.00 (.74)	.00 (.76)
785	Public	.19 (.84)	.12 (.81)	.13 (.85)	.08 (.82)	.09 (.85)
786	Baseline (Domain)	.00 (.65)	.00 (.69)	02 (.70)	04 (.70)	.02 (.78)
/8/	Baseline (Univariate)	.04 (.69)	06 (.63)	.05 (.76)	07 (.67)	.03 (.79)
700	Arbitrary	.04 (.69)	.03 (.72)	.03 (.75)	01 (.74)	.04 (.80)
769	CSV (Claude 3.5 Sonnet)	.17 (.82)	.12 (.80)	.10 (.81)	.07 (.82)	.07 (.83)
790	CSV (GPT-40)	.15 (.81)	.09 (.77)	.11 (.83)	.07 (.81)	.07 (.83)
791	CSV (Llama 3.3 70B)	.17 (.82)	.12 (.80)	.12 (.83)	.08 (.82)	.08 (.84)
792	Agent (Claude 3.5 Sonnet, Unif.)	.14 (.79)	.10 (.79)	.10 (.81)	.06 (.81)	.07 (.82)
793	Agent (Claude 3.5 Sonnet, Max Cov.)	.16 (.81)	.10 (.79)	.09 (.81)	.06 (.80)	.08 (.84)
794	Agent (GPT-40, Unif.)	.15 (.80)	.05 (.74)	.10 (.81)	.06 (.81)	.07 (.83)
795	Agent (GPI-40, Max Cov.)	.14 (.80)	.08 (.77)	.07 (.78)	.04 (.79)	.07 (.83)
796	Agent (All, Max Cov.)	.15 (.78)	.09 (.78)	.08 (.79)	.07 (.81)	.07 (.83)
797		.10 (.01)	.07 (.70)	.12 (.04)	.07 (.01)	.07 (.05)
798		(c) W	E			
799	Method	$\varepsilon = 1$	$\varepsilon = 2$	$\varepsilon = 4$	$\varepsilon = 8$	$\varepsilon = 16$
800	Without pretraining	.00 (.53)	.00 (.55)	.00 (.58)	.00 (.63)	.00 (.66)
801	Public	.11 (.64)	.18 (.73)	.13 (.71)	.06 (.69)	.06 (.72)
802	Baseline (Domain)	- 01 (53)	- 01 (54)	- 06 ( 52)	- 06 ( 57)	- 07 ( 59)
804	Baseline (Univariate)	.00 (.53)	.02 (.58)	.00 (.52)	.01 (.64)	05 (.61)
805	Arbitrary	.01 (.55)	.01 (.56)	.01 (.59)	02 (.61)	05 (.61)
806	CSV (Claude 3.5 Sonnet)	.12 (.65)	.09 (.65)	.09 (.67)	.02 (.65)	.05 (.70)
807	CSV (GPT-4o)	.05 (.58)	.08 (.64)	.06 (.64)	.05 (.69)	.03 (.69)
808	CSV (Llama 3.3 70B)	01 (.52)	.01 (.57)	.04 (.61)	.00 (.63)	04 (.62)
000	Agent (Claude 3.5 Sonnet, Unif.)	.21 (.74)	.15 (.70)	.17 (.75)	.07 (.70)	.11 (.77)
809	Agent (Claude 3.5 Sonnet, Max Cov.)	.20 (.73)	.17 (.72)	.15 (.73)	.06 (.69)	.07 (.73)
810	Agent (GPT-40, Unif.)	01 (.52)	.02 (.58)	01 (.57)	04 (.59)	06 (.60)
811	Agent (GPT-40, Max Cov.)	05 (.48)	05 (.50)	07 (.51)	05 (.58)	02 (.63)
812	Agent (Llama 3.3 70B, Unif.)	.00 (.54)	.03 (.59)	.04 (.62)	.02 (.66)	.02 (.68)
813	Agent (Llama 3.3 70B, Max Cov.)	01 (.52)	01 (.54)	.02 (.60)	.02 (.65)	.02 (.68)
<u>81</u> /	Agent (All, Unif.)	.09 (.62)	.15 (.71)	.12 (.70)	.05 (.68)	.07 (.73)
014	Agent (All, Max Cov.)	.08 (.61)	.14 (.69)	.13 (.71)	.07 (.71)	.06 (.72)

Table 9: Dataset similarity assessment against the private data for ACS, EDAD and WE. The datasets are evaluated based on two distance metrics (Section D.3): (1) Total Variation Distance (TVD); and (2) Average error on 3-Way Marginals (3WM). Both metrics are in range [0, 1], inverted to represent similarity (1 - x), and scaled by 100. Zero values (rounded) are omitted for readability. 

ACS		EDAD		WE	
1-TVD	1-3WM	1-TVD	1-3WM	1-TVD	1-3WM
48.5	50.4	4.9	26.1	6.7	34.1
4.3		0.1		0.2	
44.6	63.8	7.1	66.7	15.4	78.5
2.8		0.1			
14.4	15.0				10.9
25.7	30.2		11.5		14.2
16.6	10.0				2.4
41.5	48.3		5.5		11.7
40.1	40.0		6.8		8.0
27.3	23.3		7.2		
27.4	20.4		6.9		
13.8					
10.3					
30.5	26.6				
24.6	15.7				
	A           1-TVD           48.5           4.3           44.6           2.8           14.4           25.7           16.6           41.5           40.1           27.3           27.4           13.8           10.3           30.5           24.6	ACS           1-TVD         1-3WM           48.5         50.4           4.3         63.8           2.8         2.8           14.4         15.0           25.7         30.2           16.6         10.0           41.5         48.3           40.1         40.0           27.3         23.3           27.4         20.4           13.8         10.3           30.5         26.6           24.6         15.7	$\begin{array}{ c c c c } ACS & ED \\ \hline 1-TVD & 1-3WM & 1-TVD \\ \hline 48.5 & 50.4 & 4.9 \\ \hline 4.3 & 0.1 \\ 44.6 & 63.8 & 7.1 \\ \hline 2.8 & 0.1 \\ \hline 14.4 & 15.0 \\ 25.7 & 30.2 \\ 16.6 & 10.0 \\ \hline 41.5 & 48.3 \\ 40.1 & 40.0 \\ 27.3 & 23.3 \\ 27.4 & 20.4 \\ 13.8 \\ 10.3 \\ 30.5 & 26.6 \\ 24.6 & 15.7 \\ \hline \end{array}$	$\begin{array}{ c c c c } \hline ACS & EDAD \\ \hline 1-TVD & 1-3WM & 1-TVD & 1-3WM \\ \hline 48.5 & 50.4 & 4.9 & 26.1 \\ \hline 48.5 & 50.4 & 4.9 & 26.1 \\ \hline 44.6 & 63.8 & 7.1 & 66.7 \\ \hline 2.8 & 0.1 & & & \\ 14.4 & 15.0 & & & \\ 25.7 & 30.2 & 0.1 & & \\ 14.4 & 15.0 & & & \\ 25.7 & 30.2 & 11.5 \\ \hline 16.6 & 10.0 & & & \\ 5.5 & & & \\ 40.1 & 40.0 & & & \\ 5.5 & & & \\ 40.1 & 40.0 & & & \\ 5.5 & & & \\ 40.1 & 40.0 & & & \\ 27.3 & 23.3 & & & \\ 7.2 & & & \\ 27.4 & 20.4 & & & \\ 27.4 & 20.4 & & & \\ 13.8 & & & & \\ 10.3 & & & & \\ 10.3 & & & & \\ 30.5 & 26.6 & & \\ 24.6 & 15.7 & & \\ \end{array}$	$\begin{array}{ c c c c } \hline ACS & EDAD & W \\ \hline \hline 1-TVD & 1-3WM & 1-TVD & 1-3WM & 1-TVD \\ \hline 48.5 & 50.4 & 4.9 & 26.1 & 6.7 \\ \hline 4.3 & 0.1 & 0.2 \\ \hline 44.6 & 63.8 & 7.1 & 66.7 & 15.4 \\ \hline 2.8 & 0.1 & & & \\ \hline 14.4 & 15.0 & & & \\ 25.7 & 30.2 & 11.5 & & \\ \hline 16.6 & 10.0 & & & \\ \hline 41.5 & 48.3 & 5.5 & & \\ \hline 40.1 & 40.0 & 6.8 & & \\ 27.3 & 23.3 & 7.2 & & \\ 27.4 & 20.4 & 6.9 & & \\ \hline 13.8 & & & \\ 10.3 & & & \\ 30.5 & 26.6 & & & \\ 24.6 & 15.7 & & & \\ \hline \end{array}$

Do You Really Need Public Data? Surrogate Public Data for Differential Privacy on Tabular Data

1870	Method	Classification	Correlation	Marginals
1871	CSV (Claude)	0.002	0.033	0.121
1872	CSV (GPT)	0.001	0.149	0.096
1873	CSV (Llama)	0.003	0.052	0.041
1874	Agent (Llama, Unif.)	0.002	0.061	0.086

Table 10: Pareto Efficient Methods (Task 2: Hyperparameter tuning for private synthetic data) for PrivBayes on ACS.

1070				
1879	Method	Classification	Correlation	Marginals
1000	Arbitrary (Baseline)	0.043	0.052	0.056
1001	Agent (All, Max Cov.)	0.016	0.096	0.172
1882	Agent (Claude, Unif.)	0.019	0.040	0.070
1883				

Table 12: Pareto Efficient Methods (Task 2: Hyperparameter tuning for private synthetic data) for PrivBayes on WE.

Method	Classification	Correlation	Marginals
Public	0.008	0.047	0.097
CSV (Claude)	0.033	0.046	0.227
Agent (Claude, Unif.)	0.004	0.134	0.225

Table 11: Pareto Efficient Methods (Task 2: Hyperparameter tuning for private synthetic data) for PrivBayes on EDAD.

Method	Classification	Correlation	Marginals	
CSV (Claude)	0.013	0.002	0.045	
Agent (Claude, Max Cov.)	0.010	0.003	0.125	
Agent (Claude, Unif.)	0.004	0.003	0.024	

Table 13: Pareto Efficient Methods (Task 2: Hyperparameter tuning for private synthetic data) for AIM on ACS.

## 1888 G.3. Results for Task 2: Hyperparameter Tuning for DP Synthetic Data Generation

ACS. On ACS, where LLMs are likely to possess well-calibrated priors due to extensive training on U.S. Census data (see Appendix F.2.4), the AIM synthesizer (Table 13) shows that Agent (Claude, Unif.) is best for both classification 1890 (0.004) and marginal consistency (0.024), while CSV (Claude) has the best correlation metric (0.002) (although the 1891 Agent based Claude methods here are close behind). For the PrivBayes synthesizer (Table 10), the CSV-based approaches 1892 1893 are impressive: CSV (GPT) achieves the best in terms of classification metrics (0.001), CSV (Llama) is best in marginal metrics (0.041), and CSV (Claude) is best for correlation metrics (0.033). For GEM on ACS, the Agent (Claude, 1894 Max Cov.) approach is dominant along with the Arbitrary baseline. Recall that the Arbitrary baseline directly 1895 encodes relationships into the data (via the Bayesian approach described in Appendix E.1.3). In the case of GEM, whether 1896 relationships between variables are accurate to *true* relationships in the private data is less important when tuning its 1897 hyperparameters. 1898

1899

1878

1887

**EDAD.** We now turn to the EDAD dataset (a Spanish disability survey); EDAD was published after many LLMs' training cutoffs, so we expect the LLMs to have less, if any, prior exposure to tabular data in the same domain as the schema we present. For the AIM synthesizer (Table 14), *several* agent-based methods (e.g., Agent (All, Unif.)) are similarly strong for classification metrics (0.001). Although the correlation metrics are tightly grouped (ranging from 0.014 to 0.019), the overall Pareto frontier is defined by a mix of the CSV and Agent approaches. For the PrivBayes synthesizer (Table 11), the agent-based method Agent (Claude, Unif.) again leads on classification (0.004) while CSV (Claude) remains on the Pareto frontier for correlation (0.046); meanwhile, the real Public data yields the best marginal consistency (0.097).

WE. For WE - the Workplace Equity survey dataset, also from a period after many LLM training cutoffs - for the AIM synthesizer (Table 15) the best-performing methods are exclusively Agent-based methods. Here, Agent (Claude, Unif.) leads in classification metrics (0.016), Agent (GPT, Unif.) attains the best correlation metrics (0.007), and Agent (GPT, Max Cov.) provides the strongest marginal consistency (0.025). In contrast, for the PrivBayes synthesizer (Table 12), although the Arbitrary baseline dominates on marginal consistency (0.056) and is competitive on correlation (0.052), agent-based methods (both All, Max Cov. and Claude, Unif.) yield a substantial improvement in classification performance (0.016 - 0.019).

- 1915
- 1916
- 1918
- 1910
- 1920
- 1921
- 1922
- 1923
- 1924

1925				
1926				
1927				
1928				
1920				
1930	Method	Classification	Correlation	Marginals
1931	CSV (Claude)	0.004	0.014	0.037
1032	CSV (Llama)	0.003	0.018	0.010
1932	Agent (All, Max Cov.)	0.001	0.019	0.013
1933	Agent (All, Unif.)	0.001	0.018	0.040
1934	Agent (Claude, Max Cov.)	0.004	0.014	0.010
1935	Agent (Claude, Unif.)	0.003	0.014	0.025
1036	Agent (GPT, Max Cov.)	0.003	0.017	0.012
1930	Agent (GPT, Unif.)	0.003	0.015	0.011
1937			T 1 2 11	
1930	Table 14: Pareto Efficie	ent Methods (	Task 2: Hy	perparame
1939	ter tuning for private sy	inthetic data)	for AIM of	1 EDAD.
1940				
1042				
1942				
1944				
1945				
1946				
1947				
1948				
1949				
1950				
1951	Method	Classification	Correlation	Marginals
1952	Arbitrary (Baseline)	0.002	0.023	0.043
1953	Agent (Claude, Max Cov.)	0.002	0.039	0.072
1954				
1955	Table 16: Pareto Efficie	ent Methods (	Task 2: Hy	perparame
1956	ter tuning for private sy	inthetic data)	for GEM o	n ACS.
1957				
1958				
1959				
1960				
1961				
1962				
1963				
1964				
1965				
1966				

Method	Classification	Correlation	Marginals	
Agent (Claude, Unif.)	0.016	0.016	0.198	
Agent (GPT, Max Cov.)	0.033	0.013	0.025	
Agent (GPT, Unif.)	0.020	0.007	0.030	
Agent (Llama, Unif.)	0.017	0.010	0.047	

Table 15: Pareto Efficient Methods (Task 2: Hyperparameter tuning for private synthetic data) for AIM on WE.

Method	Classification	Correlation	Marginals
Public	0.008	0.222	0.146
CSV (GPT)	0.004	0.172	0.166
Agent (Claude, Max Cov.)	0.007	0.104	0.147

Table 17: Pareto Efficient Methods (Task 2: Hyperparameter tuning for private synthetic data) for GEM on EDAD.

Method	Classification	Correlation	Marginals
CSV (LLaMA)	0.025	0.057	0.521
Agent (All, Unif.)	0.007	0.071	0.059
Agent (GPT, Unif.)	0.028	0.056	0.058

Table 18: Pareto Efficient Methods (Task 2: Hyperparameter tuning for private synthetic data) for GEM on WE.

#### 1980 G.4. Results for Task 3: Privacy-Utility Trade-off Estimation for DP Synthetic Data Generation

As shown in Table 19, Table 20, and Table 21, the distances between the performance vectors – measured in both  $\ell_1$  and  $\ell_2$  norms – vary considerably across datasets and synthesizers. For example, in the AIM synthesizer (Table 19), methods such as Agent (All, Max Cov.) achieve an ACS  $\ell_1$  of 0.039 and an ACS  $\ell_2$  of 0.023, while CSV (Llama) attains similar values (ACS  $\ell_1$ : 0.044, ACS  $\ell_2$ : 0.023). In the GEM setting (Table 20), a similar trend is observed. Here, the Arbitrary baseline exhibits impressively low EDAD  $\ell_1$  (0.028) and EDAD  $\ell_2$  (0.013) distances, while other methods, such as Agent (All, Unif.) and CSV (GPT), also display competitive performance on certain metrics. For the PrivBayes synthesizer (Table 21), CSV (GPT) achieves an ACS  $\ell_1$  of 0.091 and an ACS  $\ell_2$  of 0.042 – values that are generally lower than those produced by several agent-based approaches on other metrics.

l	9	8	9	
	_	_	_	

2006

2007 2008

1990	Method	ACS $\ell_1$	EDAD $\ell_1$	WE $\ell_1$	ACS $\ell_2$	EDAD $\ell_2$	WE $\ell_2$
1001	Arbitrary (Baseline)	0.353	0.510	0.367	0.184	0.231	0.166
1//1	Agent (All, Max Cov.)	0.364	0.258	0.330	0.186	0.116	0.148
1992	Agent (All, Unif.)	0.519	0.126	0.373	0.274	0.057	0.168
1002	Agent (Claude, Unif.)	0.705	0.257	0.254	0.355	0.115	0.124
1993	Agent (Llama, Max Cov.)	0.543	0.559	0.260	0.288	0.251	0.119
1994	Agent (Llama, Unif.)	0.337	0.696	0.295	0.176	0.312	0.133

Method	ACS $\ell_1$	EDAD $\ell_1$	WE $\ell_1$	ACS $\ell_2$	EDAD $\ell_2$	WE $\ell_2$
Univariate (Baseline)	0.321	0.028	0.294	0.144	0.013	0.133
CSV (GPT)	0.091	0.155	0.402	0.042	0.070	0.180
Agent (All, Max Cov.)	0.094	0.071	0.318	0.043	0.033	0.144
Agent (All, Unif.)	0.112	0.051	0.280	0.051	0.024	0.126
Agent (GPT, Max Cov.)	0.127	0.061	0.232	0.058	0.027	0.105

Table 19: Priv/Util Pareto Efficient Methods (Task 3: Privacy/utility tradeoff estimation) for AIM.

Table 20: Priv/Util Pareto Efficient Methods (Task 3: Privacy/utility tradeoff estimation) for GEM.

Method	ACS $\ell_1$	EDAD $\ell_1$	WE $\ell_1$	ACS $\ell_2$	EDAD $\ell_2$	WE $\ell_2$
CSV (Llama)	0.044	0.387	0.376	0.023	0.188	0.171
Agent (All, Max Cov.)	0.039	0.100	0.191	0.023	0.051	0.092
Agent (All, Unif.)	0.082	0.091	0.167	0.041	0.056	0.081
Agent (Claude, Max Cov.)	0.068	0.152	0.111	0.033	0.085	0.063
Agent (Claude, Unif.)	0.070	0.151	0.114	0.035	0.082	0.059
Agent (GPT, Max Cov.)	0.065	0.164	0.194	0.034	0.099	0.092
Agent (Llama, Max Cov.)	0.048	0.332	0.158	0.024	0.180	0.073
Agent (Llama, Unif.)	0.042	0.442	0.177	0.023	0.216	0.082

Table 21: Priv/Util Pareto Efficient Methods (Task 3: Privacy/utility tradeoff estimation) for PrivBayes.

# 2009 G.5. Dataset Similarity May Be Less Important Than You'd Think

By using public data in DP auxiliary tasks, we implicitly assume statistical similarity to the private, sensitive data. Our results generally back this up; in pretraining (Task 1) and privacy-utility trade-off estimation (Task 3), we observe consistently better traditional public data performance. To explore whether the traditional public data dominance (and the relative performance rankings of the surrogate public data) could be explained by dataset similarity, we measure the similarity between all datasets using two common metrics from the DP literature (see Appendix D.3): Total Variation Distance (TVD) and average error across all 3-way marginal queries (3WM) (Liu et al., 2021c; McKenna et al., 2022).

**Comparing Private vs. Public.** The first dataset similarity question we ask is: how similar is a public data variant to 2017 the true, private data? Both traditional and surrogate private vs. public data similarity results are shown in Table 9. In 2018 general, our metrics suggest that the traditional public dataset and the Univariate baseline (recall, this baseline samples 2019 independently with a little noise from the true private distribution) are most similar to the private data. For EDAD and WE datasets, we can explain the lower overall similarity scores due to their higher dimensionality (defined as the Cartesian product of possible unique variable values; see the "× Dims" column in Table 1); the dataset distance is exacerbated 2022 by sparsity (particularly for TVD). However, even accounting for the limitations of these metrics, we did not observe a 2023 clear relationship between the similarity rankings of public datasets and their usefulness rankings in the pretraining and 2024 privacy-utility tasks. Our explanatory hypothesis: common similarity metrics, like TVD and 3WM, may not adequately 2025 capture dataset characteristics relevant to the DP auxiliary tasks we frame. We leave further exploration of suitable metrics 2026 to future research. 2027

Comparing Among (Traditional or Surrogate) Public Data. The second dataset similarity question we ask is: how similar are public data variants to each other, and does this partially explain their relative performance rankings? To this end, we compared similarity metrics among traditional and surrogate public data, with heatmap plots provided in Appendix I. The most consistent pattern observed across datasets and metrics is the strong similarity between Agent pairs using the same LLM but differing only in mixing methods (Unif. vs. Max Cov.) (this is barring TVD for EDAD and WE, where most entries are zero due to the aforementioned dimensionality constraints). This pattern extends to similarities

2035	between the overall mixing datasets and individual Agent datasets (with the exception of the Llama datasets on EDAD).
2030	This is expected since these pairs share the same underlying source of sampled records. Interestingly, we did not find stable
2037	similarities across generated data between different LLMs within either Agent or CSV methods, or between the same LLM
2038	across these two methods. Again, this could be an artifact of the metrics we use, but we leave a deeper exploration of this for
2039	future work.
2040	
2041	
2042	
2042	
2045	
2044	
2045	
2046	
2047	
2048	
2049	
2050	
2051	
2052	
2053	
2054	
2055	
2055	
2057	
2057	
2058	
2059	
2060	
2061	
2062	
2063	
2064	
2065	
2066	
2067	
2068	
2060	
2009	
2070	
2071	
2072	
2073	
2074	
2075	
2076	
2077	
2078	
2079	
2080	
2081	
2082	
2082	
2005	
2004	
2003	
2086	
2087	
2088	
2089	

### 2090 H. Additional Detailed Results

<sup>2091</sup> In this section, we present additional detailed results of our evaluation framework (Section 3 and Appendix F) for the <sup>2092</sup> following DP auxiliary tasks: pretraining, hyperparameter tuning, and estimating the privacy-utility trade-off.

#### 2093 2094 H.1. Results for Task 1: Private Pretraining for Classification 2095 0.20 2096 2097 0.152098 Public AUC Advantage 2099 Baseline 0.102100 Arbitrary 2101 $\operatorname{CSV}$ 0.052102 Agent 2103 Without pretraining 0.00 2104 2105 -0.052106 4 Epsilon 8 16 2107 2108 (a) ACS 2109 2110 0.20 2111 2112 0.152113 AUC Advantage Public 2114 Baseline 0.10Arbitrary 2115 $\operatorname{CSV}$ 2116 0.05 Agent 2117 Without pretraining 0.00 2118 2119 -0.052120 2121 $\dot{2}$ 4 8 16 1 Epsilon 2122 2123 (b) EDAD 2124 0.20 2125 2126 0.152127 Public AUC Advantage 2128 Baseline 0.10 2129 Arbitrary 2130 $\operatorname{CSV}$ 0.05 2131 Agent 2132 Without pretraining 0.00 2133 2134 -0.052135 $\dot{2}$ 4 Epsilon 16 8 1 2136 2137 (c) WE 2138

Figure 15: Mean AUC Advantage of the DP model after pretraining, grouped by generation method. The mean is calculated across the hyperparameter space, with 10 runs per hyperparameter configuration.

2141

2143



Do You Really Need Public Data? Surrogate Public Data for Differential Privacy on Tabular Data

Figure 16: Mean AUC Advantage of the DP model after pretraining, grouped by generation method for the sub-sampled
 ACS dataset. The mean is calculated across the DP finetuning hyperparameter space when best pretraining hyperparameter
 configuration is chosen for the pretraining step, with 10 runs per hyperparameter configuration.

Table 22: Mean AUC Advantage (AUC in parentheses) of the DP model after pretraining, grouped by generation method.
The mean is calculated across the hyperparameter space, with 10 runs per hyperparameter configuration.

2202

2229

2230

2231

2232 2233

2234

2235

2236

2237

2238 2239

Arbitrary

CSV (GPT-4o)

CSV (Claude 3.5 Sonnet)

Agent (Claude 3.5 Sonnet, Unif.)

Agent (Claude 3.5 Sonnet, Max Cov.)

CSV (Llama 3.3 70B)

Agent (GPT-40, Unif.)

Agent (All, Max Cov.)

Agent (All, Unif.)

Agent (GPT-40, Max Cov.)

	(a) AC	CS			
Method	$\varepsilon = 1$	$\varepsilon = 2$	$\varepsilon = 4$	$\varepsilon = 8$	$\varepsilon = 16$
Without pretraining	.00 (.74)	.00 (.74)	.00 (.74)	.00 (.75)	.00 (.75)
Public	.01 (.75)	.01 (.75)	.01 (.75)	.00 (.75)	.00 (.75)
Baseline (Domain)	03 (.71)	02 (.72)	01 (.73)	01 (.74)	01 (.74)
Baseline (Univariate)	02 (.72)	02 (.73)	01 (.73)	01 (.74)	01 (.74)
Arbitrary	.00 (.74)	.00 (.74)	.00 (.74)	.00 (.75)	.00 (.75)
CSV (Claude 3.5 Sonnet)	.00 (.74)	.00 (.74)	.01 (.75)	.00 (.75)	.00 (.75)
CSV (GPT-40)	.00 (.74)	.00 (.74)	.01 (.75)	.00 (.75)	.00 (.75)
CSV (Llama 3.3 70B)	.01 (.74)	.01 (.75)	.01 (.75)	.00 (.75)	.00 (.75)
Agent (Claude 3.5 Sonnet, Unif.)	.01 (.74)	.00 (.75)	.01 (.75)	.00 (.75)	.00 (.75)
Agent (Claude 3.5 Sonnet, Max Cov.)	.01 (.74)	.00 (.75)	.01 (.75)	.00 (.75)	.00 (.75)
Agent (GPT-40, Unif.)	.00 (.74)	.00 (.74)	.01 (.75)	.00 (.75)	.00 (.75)
Agent (GPT-40, Max Cov.)	.00 (.74)	.00 (.74)	.00 (.75)	.00 (.75)	.00 (.75)
Agent (Llama 3.3 70B, Unif.)	.00 (.74)	.00 (.74)	.00 (.75)	.00 (.75)	.00 (.75)
Agent (Llama 3.3 70B, Max Cov.)	.00 (.74)	.00 (.74)	.01 (.75)	.00 (.75)	.00 (.75)
Agent (Allm Unif.)	.01 (.74)	.01 (.75)	.01 (.75)	.00 (.75)	.00 (.75)
Agent (All, Max Cov.)	.00 (.74)	.00 (.75)	.01 (.75)	.00 (.75)	.00 (.75)
	(b) EDA	AD			
Method	$\varepsilon = 1$	$\varepsilon = 2$	$\varepsilon = 4$	$\varepsilon = 8$	$\varepsilon = 16$
Without pretraining	.00 (.65)	.00 (.69)	.00 (.71)	.00 (.74)	.00 (.76)
Public	.11 (.76)	.09 (.78)	.07 (.79)	.06 (.80)	.06 (.82)
Baseline (Domain)	02 (.63)	01 (.67)	02 (.69)	02 (.73)	01 (.75)
Baseline (Univariate)	03 (.62)	01 (.68)	02 (.70)	03 (.71)	02 (.74

.01 (.66)

.11 (.76)

.09 (.74)

.11 (.76)

.08 (.73)

.09 (.74)

.07 (.72)

.07 (.72)

.08 (.73)

.09 (.74)

.01 (.69)

.09 (.78)

.08 (.77)

.09 (.78)

.07 (.76)

.07 (.76)

.06 (.74)

.06 (.75)

.07 (.75)

.07 (.76)

.00 (.71)

.08 (.79)

.06 (.78)

.08 (.79)

.06 (.77)

.06 (.78)

.05 (.77)

.04 (.76)

.05 (.77)

.06 (.78)

-.01 (.74)

.07 (.81)

.06 (.80)

.07 (.81)

.05 (.80)

.04 (.79)

.04 (.78)

.04 (.79)

.05 (.79)

.05 (.79)

.01 (.77)

.06 (.82)

.05 (.81)

.05 (.81)

.04 (.81)

.05 (.81)

.04 (.80)

.04 (.80)

.04 (.81)

.05 (.81)

(c) WE

Method	$\varepsilon = 1$	$\varepsilon = 2$	$\varepsilon = 4$	$\varepsilon = 8$	$\varepsilon = 16$
Without pretraining	.00 (.53)	.00 (.55)	.00 (.58)	.00 (.63)	.00 (.66)
Public	.06 (.59)	.06 (.61)	.05 (.63)	.03 (.66)	.03 (.69)
Baseline (Domain)	02 (.51)	03 (.52)	02 (.56)	04 (.59)	04 (.62)
Baseline (Univariate)	.00 (.53)	01 (.55)	.01 (.58)	02 (.61)	03 (.63)
Arbitrary	.00 (.53)	.01 (.56)	.00 (.58)	.00 (.63)	01 (.65)
CSV (Claude 3.5 Sonnet)	.06 (.59)	.06 (.61)	.06 (.64)	.03 (.66)	.02 (.68)
CSV (GPT-40)	.04 (.58)	.05 (.60)	.04 (.62)	.03 (.66)	.02 (.67)
CSV (Llama 3.3 70B)	.00 (.53)	.00 (.56)	.02 (.59)	01 (.62)	02 (.64)
Agent (Claude 3.5 Sonnet, Unif.)	.10 (.64)	.10 (.65)	.10 (.68)	.06 (.69)	.05 (.71)
Agent (Claude 3.5 Sonnet, Max Cov.)	.11 (.64)	.11 (.66)	.09 (.67)	.07 (.70)	.05 (.71)
Agent (GPT-40, Unif.)	01 (.53)	.00 (.55)	.01 (.59)	01 (.62)	02 (.64)
Agent (GPT-40, Max Cov.)	04 (.49)	03 (.52)	03 (.55)	03 (.60)	03 (.62)
Agent (Llama 3.3 70B, Unif.)	.00 (.53)	.01 (.56)	.02 (.60)	.00 (.63)	01 (.65)
Agent (Llama 3.3 70B, Max Cov.)	01 (.52)	01 (.55)	.01 (.59)	01 (.62)	02 (.64)
Agent (All, Unif.)	.06 (.54)]	.06 (.61)	.06 (.64)	.03 (.66)	.02 (.68)
Agent (All, Max Cov.)	.03 (.56)	.03 (.59)	.05 (.63)	.01 (.64)	.01 (.67)



Figure 17: Granular hyperparameter tuning results for ACS on PrivBayes. Note the poor relative performances of the Baselines relative to the other methods; encoding relationships between variables is clearly very important to tuning hyperparameters on the PrivBayes Classifier.



Figure 18: Granular hyperparameter tuning results for EDAD on PrivBayes. Note that the agent-based method Agent (Claude, Unif.) leads in classification (0.004) while CSV (Claude) dominates the correlation metric (0.046); meanwhile, real public data yields the best marginal consistency (0.097).



Figure 19: Granular hyperparameter tuning results for WE on PrivBayes. Observe that although the Arbitrary baseline excels in marginal consistency (0.056) and is competitive on correlation (0.052), Agent-based approaches (e.g., All, Max Cov. and Claude, Unif.) offer improvement in classification performance (0.016–0.019).



Figure 20: Granular hyperparameter tuning results for ACS on the AIM synthesizer. Here, Agent (Claude, Unif.) outperforms on both classification (0.004) and marginal consistency (0.024), while CSV (Claude) is best on correlation (0.002).



Figure 21: Granular hyperparameter tuning results for EDAD on the AIM synthesizer. Several agent-based methods, such as Agent (All, Unif.), deliver strong classification performance (0.001), with the Pareto frontier defined by a mix of CSV and agent-based approaches (correlation metrics ranging from 0.014 to 0.019).

- 2362
- 2363



Figure 22: Granular hyperparameter tuning results for WE on the AIM synthesizer. Exclusively agent-based methods dominate, with Agent (Claude, Unif.) leading in classification (0.016), Agent (GPT, Unif.) achieving the best correlation (0.007), and Agent (GPT, Max Cov.) strong marginal consistency (0.025).



Figure 23: Granular hyperparameter tuning results for ACS on the GEM synthesizer. The Agent (Claude, Max Cov.) method, alongside the Arbitrary baseline that directly encodes variable relationships, is dominant – reinforcing that structure in the data is beneficial.



Figure 24: Granular hyperparameter tuning results for EDAD on the GEM synthesizer. As in ACS, both the agent-based approach and the Arbitrary baseline perform competitively.

- 2418
- 2419



Figure 25: Granular hyperparameter tuning results for WE on the GEM synthesizer. The trends mirror those in ACS, with the Arbitrary baseline maintaining strong performance and Agent-based methods showing similar competitiveness. 

10 20 Avg. % Degradation

20 40 Avg. % Degradation

CSV (L)

-

Avg. Abs. Degradation



Figure 26: Privacy/utility tradeoff estimation results in terms of  $\ell_1$  distance from the true sensitive data tradeoff. Note the relatively consistent performance across synthesizers for each dataset between some methods (e.g., poor privacy/utility tradeoff estimation for CSV on WE), while other methods have higher variance (e.g., Agent (Claude 3.5 Sonnet, Max Cov. on ACS, between GEM and AIM).

2527

2528



Figure 27: Privacy/utility tradeoff estimation results in terms of  $\ell_2$  distance from the true sensitive data tradeoff. These results largely mirror the  $\ell_1$  distance results, although the increased sensitivity to outliers leads to some interchanges of ranking (e.g., Agent (Claude 3.5 Sonnet, Unif.) and CSV (GPT-40) interchange places on ACS PrivBayes).



Figure 28: To provide intuition for *exactly* what the  $\ell_1$  and  $\ell_2$  scores in Figures 26 and 27 attempt to capture, we plot the average performance across epsilon that constitutes each vector, relative to the true performance of the sensitive data (which, in these plots, is the black dotted line). Ideally, for privacy/utility estimation, any public data (surrogate or otherwise) would *match* the performance of the private data across privacy loss budget parameters. This would allow a practitioner to, say, choose the correct  $\epsilon$  based on a performance threshold in absolute terms. Clearly, given the noisiness of the lines (which generally cluster around, but inconsistently track, the black dotted line for private data performance), this is a difficult estimation problem.

- 2633 2634
- 2635
- 2636
- 2637
- 2638
- 2639



Figure 29: Heatmap of similarity metrics based on the Total Variation Distance (TVD) between the datasets based on the ACS data. The metric is in range [0, 1], inverted to represent similarity (1 - x), and scaled by 100, and rounded to a single digit.

2695																		
2696																		
2697																		
2698																		
2699																		
2700																		
2701																		
2702																		
2703													1.)					
2704											.ş	Q .	<sup>2</sup> 00			ංර	4.)	
2705								_			UN	10	,	<i>(</i> )	ili.	197 C		
2706					~		6	Jet)	à		They an	They 'S	.) 4	Cor &	N 1) - S	Ar.		
2707				in	riate		13 502		27000	1.3 No	1.3 Se	Unit	Mor	2310	2310		4.)	
2708			00	nia Thi	7.0	,de	5. 5.	()	5.0 and	e and	e or	to ot	XO '	1 3 . MR	,	С Ф	5	
2709		• 61	ine	The Co	rary .	Clark	Gr Y	Libit.	Oro,	Oro,	Gr '	œ,	Dia.	Ula ,	On.	Aio		
2710		Public	Basel Base	L' Aibit	" CSA	654	_ c54	Agen	r Agel	Ager	A veen	, Agen	r Ager	re Neet	Asen	, ,		
2711		Ŷ	v v	۰ ۲				×	×	· · · · ·	*	*	,	,	*			
2712	Private	50.4	<b>0</b> 63.8	0	15	30.2	10	48.3	40	23.3	20.4	0	0	26.6	15.7			
2713																		
2714	Public		0 36	0	18.9	16.6	0	33.9	27.5	9.5	8	0	0	16.6	6.8			
2715	Baseline (Domain)	. –	0	9.5	0	0	5.2	0	0	79	0	0	0	0	0			100
2716	Dasenne (Domann)		0	5.0	0	0	0.2	0	0	1.5	0	0	0		0			100
2717	Baseline (Univariate)			0	0	14.3	1.9	52.7	42.3	29.3	25.3	0	0	34	22			
2718																		
2719	Arbitrary				0	0	0	0	0	0	0	0	0	0	0		-	80
2720	CSV (Claude 3.5 Sonnet)					40.9	0	4.4	6	0	0	0	0	0	0		-	
2721						10.00												
2722	CSV (GPT-4o)	-					15.8	16.8	11.4	3.5	1.6	0	0	3.8	0		-	
2723									ō	10.0	10.1		ō		_		-	60
2724	CSV (Llama 3.3 70B)	-						12.1	0	18.3	13.1	0	0	15				
2725	Agent (Claude 3.5 Sonnet, Unif.)	_							70.1	47.5	45.1	3.7	0	55.7	37.9		-	
2726																	-	
2727	Agent (Claude 3.5 Sonnet, Max Cov.)	-								25.3	26.3	0	0	34.1	20.5		-	40
2728											70.0	10.0	0.1	00.4	50.0			
2729	Agent (GP 1-40, Unif.)										79.6	12.2	2.1	03.4	50.9			
2730	Agent (GPT-40, Max Cov.)											14.3	1	61.6	51.4			20
2731																		20
2132	Agent (Llama 3.3 70B, Unif.)												56.8	40.6	52.3			
2133	Amont (Ilema 9.9.70D M. C. )													07.4	11-0-			
2725	Agent (Liama 3.3 (0B, Max Cov.)	-												-21.4	44.8			0
2133	Agent (Unif.)	-													71.9			
2130																		
2131																		

Figure 30: Heatmap of similarity metrics based on the Average Error on 3-Way Marginals (3WM) between the datasets based on the ACS data. The metric is in range [0, 1], inverted to represent similarity (1 - x), and scaled by 100, and rounded to a single digit.

2754																	
2755																	
2756																	
2757																	
2758																	
2759												6	్రం			1 a.)	
2760											UN UN	No	μ.,		it.	at CO	
2761					~		S	er	à		They, of	They, it.		<sup>م</sup> ر مر	NY N	Ar.	
2762				a;	riate		12 Sor		200	్సిం	1,00	Unit	Mar	3701	370	4.	1
2763			.D05	no. This	<u>8</u>	de	3." × /	5 (S	5.5 all	e nid	e or	x0.	10, 118	,	·	4.00°	
2764			sine	The Co	253 1	Clair	Gr v	Lian,	Olo	Or	Car )	Gr	(Lio	(Dia	On 1	Nia	
2765		Public	38501 5850	r sibili	r cest	کچی `	_ CSA	, ser	r veer	y Peet	r veer	, Agen	, Asen	, reer	r veene		
2766		Ŷ	• •	,	Ũ	Ū	e	*	*	7	*	*	,	\$	,		
2767	Private	4.9 0	).1 7.1	0.1	0	0	0	0	0	0	0	0	0	0	0		
2768																	
2769	Public	- 0	).1 1.8	0	0	0	0	0	0	0	0	0.1	0	0	0.1		
2770	Pagaling (Domain)		0.1	0.1	0	0	0	0	0	0	0.1	0	0.1	0	0		100
2771	Baseline (Domain)	-	0.1	0.1	0	0	0	0	0	0	0.1	U	0.1	0	0		- 100
2772	Baseline (Univariate)	-		0.1	0	0	0	0	0	0	0	0	0	0	0		
2773																	
2774	Arbitrary	-			0	0	0	0	0	0	0.1	0.1	0	0	0.1		- 80
2775	CSV (Clauda 2.5 Seppet)			I		<u>م</u> ا	05	0	0	0	0	0	0	0	0		
2776	CSV (Claude 3.5 Sonnet)					2.4	0.0	0	0	0	0	0	0	0	0		
2777	CSV (GPT-40)	-					0.7	0	0	0	0	0	0	0	0		
2778	· · · · ·																- 60
2779	CSV (Llama 3.3 70B)	-						0	0	0.1	0.1	0	0	0	0		
2780	Ament (Claude 2 5 Connet Unif)						1		0	0	0	0	0	16	0		-
2781	Agent (Claude 5.5 Sonnet, Unit.)	-							0	. 0	0	U	U	4.0	0		
2782	Agent (Claude 3.5 Sonnet, Max Cov.)	-								0.1	0	0	0	0	0.1		- 40
2783																	
2784	Agent (GPT-40, Unif.)	-									0.9	0.1	0.5	0	0.3		-
2785	A rest (CDT 4- Mars Cars)											0.1	0 5	0.1	06		-
2786	Agent (GP1-40, Max Cov.)	-										0.1	0.5	0.1	0.6		- 20
2787	Agent (Llama 3.3 70B, Unif.)												9.3	0.7	4.4		
2788																	
2789	Agent (Llama 3.3 70B, Max Cov.)	-												0.8	13.1		
2790	A /TT . C \														0.0	_	
2791	Agent (Unif.)	-													0.6		
2792																	

Figure 31: Heatmap of similarity metrics based on the Total Variation Distance (TVD) between the datasets based on the EDAD data. The metric is in range [0, 1], inverted to represent similarity (1 - x), and scaled by 100, and rounded to a single digit.



Figure 32: Heatmap of similarity metrics based on the Average Error on 3-Way Marginals (3WM) between the datasets based on the EDAD data. The metric is in range [0, 1], inverted to represent similarity (1 - x), and scaled by 100, and rounded to a single digit.

2860																			
2861																			
2862																			
2863																			
2864																			
2865																			
2866																			
2867																			
2868														<i>.</i>					
2869												.9	;	_00~			. 6	. k	
2870												UN UN	10 10	yik"		if.	. A C	,	
2871						~		5	Jet)	2		met.	met		00 <sup>0</sup>	n a	Pro		
2872					a.	riatel		5 SOL		1015	1,00	్లి	Unit	Mar	3702	3702			
2873				OOT	las vii	Ser	Se	3.º ~ W	S	3.3 . Jd	2° . W	e at	10. A	10. 10?	.5. 	.». 	, <u></u>	o'	
2874			<i>.</i>	se V	De C	ary .	Claur	Er .	Diante.	Olor,	Olo.	CX.	CX .	(1)an	(1)an	Our	Oya.		
2875		oud	10 2326	1 23 <sup>201</sup>	r vibit	1.64	654	لهجم	, ser	i seri	> ~set	r veen	N Ser	P Ref	r Ser	» Ser	2		
2876		Х,	v	Ŷ	۶ <b>۰</b>	0	U	U	ç.	ç.	ç <b>ı</b>	ç.	ç.	ç.	ζ <b>τ</b> .	ç.			
2877	Private	6.7	0.2	15.4	0	0	0	0	0	0	0	0	0	0	0	0			
2878																			
2879	Public		0	5.4	0	0	0	0	0	0	0	0	0	0	0	0			
2880	Basalina (Domain)			0.1	0.1	0	0	0	0	0	0	0	0	0	0	0.1			100
2881	Dasenne (Domain)			0.1	0.1	U	0	0	0	0	0	0	0	0	0	0.1			100
2882	Baseline (Univariate)	-			0	0	0	0	0	0	0	0	0	0	0	0			
2883																			
2884	Arbitrary	-				0	0	0	0	0	0	0.1	0	0	0	0		_	80
2885	CSV (Claude 3.5 Seppet)						54	1.6	0	0	0	0	0	0	0	0			
2886	CSV (Claude 5.5 Sonnet)						0.4	1.0	0	0	0	0	0	U	0	0			
2887	CSV (GPT-4o)	-						1	0	0	0	0	0	0	0	0			
2888																		-	60
2889	CSV (Llama 3.3 70B)	- -							0	0	0	0	0	0	0	0			
2890	Agent (Claude 3.5 Sennet Unif)									4.5	0	0	0	0	5.1	0			
2891	Agent (Claude 5.5 Sonnet, Chin.)									4.0	0	0	0	0	0.1	0			
2892	Agent (Claude 3.5 Sonnet, Max Cov.)	_									0	0	0.1	0	1.1	0		-	40
2893																		-	
2894	Agent (GPT-40, Unif.)											8.2	0.1	0.1	0.6	0.3			
2895	Agent (CDT 4c, May Cov.)												0.1	0.9	0.4	0.1		-	
2896	Agent (GF 1-40, Max COV.)	-											0.1	0.2	0.4	0.1			20
2897	Agent (Llama 3.3 70B, Unif.)	-												21.1	0.6	3.1			
2898		-																	
2899	Agent (Llama 3.3 $70\mathrm{B},\mathrm{Max}$ Cov.)	-													0.5	4.1			0
2900																	i		U

Figure 33: Heatmap of similarity metrics based on the Total Variation Distance (TVD) between the datasets based on the WE data. The metric is in range [0, 1], inverted to represent similarity (1 - x), and scaled by 100, and rounded to a single digit.

0.4

Agent (Unif.)



Figure 34: Heatmap of similarity metrics based on the Average Error on 3-Way Marginals (3WM) between the datasets based on the WE data. The metric is in range [0, 1], inverted to represent similarity (1 - x), and scaled by 100, and rounded to a single digit.

## 2970 J. Compute and Resources

Benchmarking DP synthesizers and training models for differentially private tasks is computationally intensive (Rosenblatt et al., 2023). We executed our experiments on a combination of high-performance GPU and CPU clusters hosted on AWS EC2. Specifically, we utilized three q4dn.12xlarge instances - each equipped with NVIDIA T4 GPUs - for approximately 17.3 days of continuous up-time per instance, amounting to roughly 52 GPU-days in total (although it is hard to assess the true GPU utilization). In addition to local compute, we used LLM APIs provided by OpenAI, Anthropic, and TogetherAI (for the Llama 3 model) for both our direct CSV generation and multi-step Agent-based approaches. We conducted substantial inference for our experiments; as an example, during January, our queries to Claude alone amounted to a total of 38,092,225 input tokens and produced 7,099,403 output tokens, in February, we recorded 11,922,046 input tokens and 226,998 output tokens, and in March, 9,027,827 input tokens and 124,484 output tokens were consumed (imbalance between input output due to re-inputting previously generated tokens as context on each call in the state machine for the Agent). These resources allowed for extensive hyperparameter searches, multiple runs per privacy setting, and a comprehensive evaluation across DP auxiliary tasks.