
HOW TO DP-FY YOUR DATA: A PRACTICAL GUIDE TO GENERATING SYNTHETIC DATA WITH DIFFERENTIAL PRIVACY

Natalia Ponomareva
Google Research
New York, New York, USA
nponomareva@google.com

Zheng Xu
Google Research
Mountain View, CA, USA
xuzhustc@gmail.com

H. Brendan McMahan
Google Research
Seattle, WA, USA
mcmahan@google.com

Peter Kairouz
Google Research
Seattle, WA, USA
kairouz@google.com

Lucas Rosenblatt
NYU*
New York, New York, USA
lurosenb@google.com

Vincent Cohen-Addad
Google Research
New York, New York, USA
cohenaddad@google.com

Cristóbal Guzmán
Institute for Mathematical
and Computational Engineering[†]
Pontificia Universidad Católica de Chile
crguzmanp@mat.uc.cl

Ryan McKenna
Google Research
Seattle, WA, USA
mckennar@google.com

Galen Andrew
Google Research
Seattle, WA, USA
galenandrew@google.com

Alex Bie
Google Research
New York, New York, USA
alexbie@google.com

Da Yu
Google Research
Mountain View, CA, USA
dayuwork@google.com

Alex Kurakin
Google DeepMind
Mountain View, CA, USA
kurakin@google.com

Morteza Zadimoghaddam
Google Research
Zurich, Switzerland
zadim@google.com

Sergei Vassilvitskii
Google Research
New York, New York, USA
sergeiv@google.com

Andreas Terzis
Google DeepMind
Mountain View, CA, USA
aterzis@google.com

December 4, 2025

ABSTRACT

High quality data is of vital importance for unlocking the full potential of AI for end users. However finding new sources of such data is getting harder: most publicly-available human generated data will soon have been used (Villalobos et al., 2024). Additionally, publicly available data often is not representative of users of a particular system — for example, a research speech dataset of contractors interacting with an AI assistant will likely be more homogeneous, well articulated and self-censored than real world commands that end users will issue. Therefore unlocking high-quality data grounded in real user interactions is of vital interest to both system creators and end users themselves. However, the direct use of user data comes with significant privacy risks, which must be addressed before the data can be used. Differential Privacy (DP) is a well established framework for reasoning about and limiting information leakage, and is a gold standard for protecting user privacy. The focus of this work, *Differentially Private Synthetic data*, refers to synthetic data that preserves the overall trends of source data (often user-generated), while providing strong privacy guarantees to individuals

*Work done at Google as part of student researcher engagement.

[†]Work done while at Google.

that contributed to the source dataset. DP synthetic data can unlock the value of datasets that have previously been inaccessible due to privacy concerns. Additionally, DP synthetic data can replace the use of sensitive datasets that previously have only had rudimentary protections like ad-hoc rule-based anonymization.

In this survey we explore the full suite of techniques surrounding DP synthetic data, the types of privacy protections different generation approaches can offer, and the state-of-the-art for various modalities including image, tabular, text and federated (decentralized) data. We outline all the components needed in a system that generates DP synthetic data, from sensitive data handling and preparation, to tracking the use of synthetic data and empirical privacy testing.

We hope that work will result in increased adoption of DP synthetic data, spur additional research in still underexplored domains, and additionally increase trust in DP synthetic data approaches.

Contents

1	Introduction	5
1.1	Preview of the later sections	6
1.2	Types of synthetic data	7
1.3	Motivation	7
1.4	Challenges in DP synthetic data generation	8
1.5	Evaluating synthetic data quality	8
2	A brief overview of Differential Privacy	9
2.1	Definition of Differential Privacy	9
2.2	Neighboring datasets and the privacy unit	10
2.3	Key Properties and Mechanisms of Differential Privacy*	12
2.4	DP synthetic data and methods for creating it	13
2.5	DP-Training	14
3	DP Synthetic Tabular data	15
3.1	Privacy unit for tabular data	16
3.2	Foundations of DP tabular synthetic data	17
3.3	Practical Workload-Adaptive Algorithms	19
3.4	Algorithms with strong worst-case guarantees*	23
3.5	Numerical feature handling	26
3.6	Modern end-to-end methods for tabular data generation	26
3.7	Comparison of the methods	28
3.8	Evaluating synthetic tabular data quality	29
3.9	Open questions and challenges	30
4	DP Synthetic Image data	31
4.1	The privacy unit for image data	31
4.2	Methods for DP Synthetic image generation	32
4.3	DP training/finetuning	33
4.4	Methods that avoid DP-Training: Private Evolution	36
4.5	Alternative methods	38
4.6	Comparison of the methods	40
4.7	Evaluating synthetic image data quality	41
4.8	Open questions and challenges	42
5	DP Synthetic Text data	43
5.1	Privacy unit for text data	43
5.2	DP training and finetuning	44
5.3	Methods that avoid DP-Training	45

5.4	Summary and comparison of methods	50
5.5	Resampling DP synthetic text	52
5.6	Evaluating synthetic text data quality	53
5.7	Open questions and challenges	54
6	DP Synthetic Data in Federated Learning	54
6.1	Federated Learning Primer	55
6.2	Differential Privacy in Federated Learning	56
6.3	DP-Training in FL	57
6.4	DP Synthetic Data in FL	58
6.5	Methods for Federated DP Synthetic Data	58
6.6	Evaluating synthetic data in FL	61
6.7	Discussion	61
7	Practical Components of DP synthetic data generation system	61
7.1	Privacy-related decisions	62
7.2	User contribution bounding	62
7.3	Pitfalls to avoid when implementing DP-training	64
7.4	Best-effort PII Removal	65
7.5	Empirical privacy auditing	66
7.6	Security and data minimization during DP synthetic data generation	71
7.7	Lineage tracking of synthetic data usage*	71
8	Conclusion	74
9	Acknowledgement	75

1 Introduction

The modern world thrives on data. Vast quantities of information, encompassing sensitive domains like healthcare records, financial transactions, personal user preferences and others are continuously generated and utilized. Leveraging this data for AI training, fine-tuning, evaluation, or inference can unlock scientific breakthroughs, business insights, and societal advancements. However, these practices carry significant privacy risks that need to be properly addressed. Traditional methods for protecting privacy like data sanitization and de-identification have been proven increasingly inadequate (Brown et al., 2022). Sophisticated attacks, including re-identification through linking anonymized datasets with external information sources (*data linkage attacks*) can successfully uncover individuals’ identities and sensitive attributes, undermining privacy protections (Narayanan and Shmatikov, 2008). The inherent difficulty in predicting all possible ways large datasets might be queried or linked makes simple anonymization a fragile defense (Narayanan and Shmatikov, 2008). Consequently, there is a critical need for more robust, mathematically grounded approaches to enable data analysis and data sharing, while rigorously protecting individual privacy (Jordon et al., 2022).

Synthetic data generation has emerged as a compelling approach to address this challenge (Jordon et al., 2022). The process involves creating artificial data that mimics the statistical patterns and properties of a real dataset, often using algorithms or generative models trained on the original data. This artificial data ideally contains no real values from the source dataset, offering a potential pathway to share insights, facilitate research, develop and test software, and train machine learning models without directly exposing sensitive individual records (Afonja et al., 2025a).

However, a common and dangerous misconception is that synthetic data is inherently private (Jordon et al., 2022). Generative models, especially powerful modern architectures, can and do inadvertently memorize and replicate specific details from their training data (Carlini et al., 2021). Synthetic records therefore might be near-copies of real ones, or the generated dataset might leak aggregate information that still allows inferences about individuals present in the original data. Consequently, simply generating synthetic data based on the original sensitive data does not automatically guarantee privacy; significant care and formal techniques are required to produce synthetic data that is both useful and demonstrably private.

Differential Privacy (DP) (Dwork and Roth, 2014) provides the necessary rigorous framework to address the privacy shortcomings of naive synthetic data generation. DP offers a rigorous quantified guarantee of privacy protection against a wide range of attacks, including those unknown at the time of deployment (Wood et al., 2018). The core promise of DP is intuitive yet powerful: the outcome of a differentially private analysis or data release should be roughly the same whether or not any single individual’s data was included in the input dataset (Dwork and Roth, 2014). This effectively masks the contribution of any single person, ensuring they remain “hidden in the crowd” and limiting what can be learned specifically about them from the output (Wood et al., 2018).

This survey provides a holistic overview of the field of DP synthetic data generation.³ This is a self-contained guide on how to create DP synthetic data for various modalities including image, tabular and text. It provides analysis and comparison of various methods that exist, highlights current state-of-the-art methods and their ideal usecases as well as shortcomings. We cover the full life cycle of DP synthetic data: from data preparation to safe data handling, best practices for testing the synthetic data and lineage tracking for usage of DP synthetic data.

Our target audience is practitioners, privacy advocates, and academic researchers alike. Practitioners will benefit from a gentle introduction to DP synthetic data field and clear definitions of what DP synthetic data guarantees and protects, comparisons of competing methods for each modality and an overview of what functionality a DP synthetic data creation system should provide. Additionally we cover important implementation details, that when overlooked could weaken privacy guarantees of the system. For example, the need to understand all the components of a generative model that is being used for creating DP synthetic data and the need to track DP synthetic data lineage. For academic researchers this work can serve as a one-stop in-depth survey of current state-of-the-art and advances in the field of DP synthetic data. We also outline open questions and challenges for each of the modalities, hoping that it will spur additional academic interest for these directions.

Finally, we note that this survey assumes a background in ML including modern auto-regressive models like large language models (LLMs). Throughout the content we mark some sections as * to indicate that they contain additional in-depth theoretical details and can be safely skipped if reading/skimming for the broad picture.

Attention boxes

We also introduce grey boxes like this one to highlight an important statement, conclusion or suggestion.

³Unintended omissions and errors are possible, especially given the breadth and challenging nature of the topic. We welcome feedback on this work from all interested parties.

1.1 Preview of the later sections

This survey is organized as follows.

1. We begin this Section with an introduction to DP synthetic data and position it w.r.t. other types of synthetic data in Section 1.2. Section 1.3 clarifies the motivation behind DP synthetic data and outlines its pros and cons. Section 1.4 introduces challenges that need to be overcome to obtain useful DP synthetic data for any modality. Finally, Section 1.5 discusses in broad strokes the evaluation of DP (and non DP) synthetic data for any modality, with later sections (Sections 3.8, 4.7, 5.6 and 6.6) covering modality specific metrics.
2. Section 2 provides background information on Differential Privacy and DP synthetic data. We start with the definition of DP itself (Section 2.1) and introduce the concept of the privacy unit, which is paramount for understanding DP guarantees (Section 2.2). We touch upon important properties and mechanisms of DP in Section 2.3*. Section 2.4 compares DP synthetic data with alternatives like direct DP training of ML models and outlines the types of methods for DP data synthesis that we will be using for categorizing methods for each modality. We also briefly introduce a technique called *DP-Training* (Section 2.5) which is a workhorse for DP synthetic data generation for all modalities and will be mentioned throughout this work extensively.
3. In Section 3 we tackle the ubiquitous case of *tabular data*. Section 3.1 introduces a discussion of appropriate privacy unit for tabular domain. We introduce foundations of *Query release* (the task of creating useful DP synthetic tabular data for answering statistical queries) and the concept of a query *workload* in Section 3.2. Section 3.3 focuses on practical algorithms and outlines the prominent Select-Measure-Estimate paradigm that these algorithms often follow. This class of algorithms still dominates the field of DP tabular data synthesis. For an interested reader, we also dive into less scalable algorithms that are nevertheless of theoretical interest in Section 3.4* and discuss data preprocessing (like imputation, domain calculation etc.) needed for both empirically and theory oriented algorithms in Section 3.5. We explore the new end-to-end approaches that employ the power of foundational models like LLMs for tabular data synthesis in Section 3.6 and offer a fair comparison of all the above methods in Section 3.7. Section 3.8 surveys metrics that can be used for evaluating the quality of tabular synthetic data w.r.t the original data, and Section 3.9 concludes with open questions and challenges for DP data synthesis methods for this domain.
4. Section 4 is devoted to the *image* modality. We start with discussing an appropriate privacy unit for visual data in Section 4.1 and introduce broad categorization of methods for DP image synthesis in 4.2. Section 4.3 is a deep dive on *DP-Training* of vision models – the workhorse method for DP synthetic image synthesis – where we explore *DP-Training* and *DP-finetuning* of GANs and Diffusion models. We additionally explore training-free method called *Private Evolution* in Section 4.4 and outline alternative methods (Section 4.5), including new promising line of work that generates images via intermediary representation (Section 4.5.2). Section 4.6 offers a comparison and outlines ideal use cases for each of the explored methods, and Section 4.7 investigates metrics for evaluating the quality of synthetic images. We conclude with open questions and challenges in Section 4.8.
5. We explore the text modality in Section 5. We outline difficulties in choosing an appropriate privacy unit for text data (Section 5.1), and explore *DP-finetuning* – the method that produces the best quality synthetic data given sufficient amount of original data and compute in Section 5.2. We outline alternative methods that don't require computationally expensive modifications to LLMs - namely DP inference (5.3.1) and Private evolution algorithms (Section 5.3.2). Proxy metrics that can be used to gauge the quality of the synthetic text are presented in Section 5.6. We draw readers' attention to open research problems pertinent to text DP synthetic data generation in Section 5.7.
6. DP synthetic data for decentralized data is covered in 6 where we first introduce *Federated Learning* (Section 6.1) and highlight difficulties of achieving DP in FL in Section 6.2. We examine the meaning of DP synthetic data in FL in Section 6.4 and outline unique challenges one encounters when attempting to generate DP synthetic data from federated datasets. We briefly survey work on FL DP synthetic data generation in Section 6.5 and compare and discuss such methods in Section 6.7. Section 6.6 introduces evaluation challenges unique to distributed nature of the sensitive data.
7. Section 7 is devoted to practical considerations that designers of a DP synthetic data end-to-end generation system need to take into account. Section 7.1 highlights the most important questions that should be answered, including which privacy unit and privacy guarantees to target (Section 7.1). We then focus on data preparation, including user-level contribution bounding (Section 7.2), and safe sensitive data handling practices (Section 7.6). We investigate a crucial component that needs to be in place before any DP synthetic data is used in downstream tasks or shared – namely *Empirical Privacy Auditing* (Section 7.5). We finally bring forward the need of lineage tracking for DP synthetic data in Section 7.7.

1.2 Types of synthetic data

Synthetic data is admittedly an overused term that can mean different things to different people. In this work, we assume that some “real” data (drawn from the “original distribution”) exists or could exist, for example from real-world physical interactions or end-user interactions with a product or service. However, the real data cannot be used directly, perhaps because it is too limited in quantity, expensive to collect, or direct use would compromise privacy. Instead, alternative *synthetic data* is desired. We assume the goal of synthetic data is to match some important properties of this real data. Such synthetic data can be obtained in many ways which can be in general categorized into two broad groups.

Methods that don’t directly utilize real data. For example, simply choosing a readily available proxy dataset (e.g., using the public corpus of Enron emails rather than real user emails) might be sufficient; alternatively synthetic data can be generating by running a physical (e.g. robotic arm manipulation experiments) or virtual simulations; by prompting a generic LLM or other GenAI model to produce data with specific desired properties (e.g. "write a scientific article about types of synthetic data, include abstract, main body and conclusion"); by asking writers or contractors to produce the data; and many more. Importantly, none of these techniques directly utilize real data. Depending on the original distribution, it may be very hard to use these techniques to generate high-quality synthetic data (data that preserves the desired properties of the real data), or in fact to even evaluate the quality of the synthetic data.

Methods that use real data. Grounding synthetic data using real data allows for much higher fidelity and utility of the synthetic data. One example of techniques in this groups is fitting parameters of a physics simulation based on the real data; classic techniques like SMOTE (Bowyer et al., 2011); prompting a human or an LLM to "Produce more documents like this text, but with different tones and styles, and on slightly different topics" or "Write chat bot interaction with a customer who is having problem with a banking interface"; or fine-tuning a generative model on the real data, and then using that fine-tuned model to sample as many examples as necessary. For any approach that utilizes the real data in producing the synthetic data, one must address the question of privacy: could private information in the real data be revealed in the synthetic data? DP, introduced more formally in Section 2.1, provides a rigorous framework for addressing these privacy concerns, and forms the basis for the DP Synthetic Data approaches described in this work which can generate high-quality synthetic data with strong privacy guarantees.

1.3 Motivation

In context of ML, Differential Privacy (Dwork and Roth, 2014) previously was applied mostly during the training stage by using specialized DP-algorithms like DP-SGD (Abadi et al., 2016) or DP-FTRL (Kairouz et al., 2021a) in lieu of standard training techniques like SGD. The state-of-the-art for using differentially-private training techniques has substantially advanced over the past decade, from early research that showed such approaches could be practical (Abadi et al., 2016) to present time, when a variety of advanced approaches are available and DP-trained models have been used in large-scale production deployments (Xu et al., 2023b; Zhang et al., 2023). Nevertheless, directly training production models on privacy-sensitive data raises numerous practical challenges — most importantly, only training infrastructure and models that support DP can be used, and because the private sensitive data cannot be inspected, there is no ability to debug with, filter, or human-label sensitive training data, and any analysis of sensitive data (e.g. like calculating counts, variance etc) should be done with DP.

The success of generative AI models has led to an alternative, 2-stage approach: we first create DP synthetic version of the private data that is highly representative of the raw private data. With suitable protections including appropriate DP parameters and audits, this DP synthetic data can be treated as fully anonymous. Thus, it can be used relatively freely in any standard ML workflow: filtering, augmentation, human-labeling or analysis, etc. The (still DP) outputs of any of these processes can then be used freely for tasks including model evaluation, fine-tuning of a production model (possibly one where its size or infrastructure would preclude DP training which is extremely computationally expensive and requires custom implementation of data processing and training), and more. In its ideal form, the formation of DP synthetic data concentrates all the challenges of anonymization in a single step, and provides an output that can be much more useful than a single DP model fine-tuned for a specific purpose.

DP synthetic data offers other advantages, on top of aforementioned ease-of-use, as well. From a privacy point of view, when fine-tuned GenAI models (e.g. LLMs) were used to create this data, the derived DP synthetic data offers a strictly smaller attack surface than the DP GenAI model from which it was derived. For example, numerous papers (Nasr et al., 2025; Carlini et al., 2022, 2021) have now shown that clever and malicious attacks can extract information from LLMs. However once DP synthetic data is generated, such attacks are structurally not possible — only DP synthetic data is released and all the artifacts for preparing this data (including GenAI models if they were employed) can be erased or remain locked down. Further, the DP synthetic data itself can be additionally audited for potentially privacy-sensitive information before it is used for any downstream purpose.

The area of DP synthetic data has been extensively explored in context of tabular data, where workload-based algorithms have been perfected over the last 2 decades (Section 3). (Non-DP) Image synthetic data, which requires larger and complex ML models like GANs, saw limited success until the emergence of powerful models like Diffusion, VAE and Autoregressive models (Section 4). The field of synthetic text benefited immensely from new powerful autoregressive models like LLMs (Section 5). Nevertheless, the area of DP synthetic data generation for complex modalities like text and images is very much still nascent.

While there are indications of increased trust in DP synthetic data for navigating the inherent trade-off between data utility and individual privacy protection (Afonja et al., 2025a), the adoption is not yet wide and DP synthetic data field remains contentious even among the experts. For example, a recent interview study with “data experts” (defined as professionals like researchers, data scientists, and practitioners in fields such as economics, medicine and public policy who directly engage with data) highlighted the challenge in practical adoption of differentially private synthetic data (Rosenblatt et al., 2024b). This study found that while experts recognize the potential for broad data access through privacy preserving mechanisms, they express many concerns about the potential for DP noise to lead to incorrect conclusions. Recommendations in the face of these concerns included (1) validating DP synthetic data utility against real world use cases and (2) establishing clear standards of use.

We hope that this self-contained work will serve to increase the trust of experts (and non experts alike) in technology that underpins DP synthetic data generation, spur wider adoption of DP synthetic data in industry, and encourage research into methods and questions that are still not resolved.

1.4 Challenges in DP synthetic data generation

Generating high-quality DP synthetic data, especially for complex modalities like images, long text, multi-modal data, still faces significant hurdles.

- **Balancing Utility and Privacy** remains the paramount challenge (Jordon et al., 2022). The noise or constraints introduced by DP inevitably distort the data (Jordon et al., 2022). The key difficulty lies in minimizing this distortion to ensure the synthetic data retains sufficient fidelity (e.g. visual quality for images, statistical accuracy, and usefulness for downstream tasks while adhering to a chosen privacy budget (ϵ, δ) (Perez et al., 2024).
- **Preventing Memorization:** Generative models that are commonly used for DP synthetic data creations, particularly large ones like LLMs, can memorize parts of their training data (Carlini et al., 2021). A DP mechanism must effectively prevent the model from simply outputting copies or near-copies of sensitive training examples. While DP provides a theoretical guarantee against exact replication influencing the output distribution, empirical validation (Section 7.5) is often needed to assess practical memorization risks, especially under weaker privacy settings (larger ϵ) which are commonly used in practice.
- **High Dimensionality of original sensitive data:** Both text and images data have high dimensionality representations and capturing the intricate distributions of such data requires complex models. Applying DP noise in high-dimensional spaces (either original space or large models gradient space during training) can easily overwhelm the signal, leading to poor quality synthetic images (Chen et al., 2022b).
- **Computational Cost:** Training large-scale generative models like GANs, diffusion models or LLMs is computationally intensive. DP training (Section 2.5, which is one of the main workhorse methods for DP synthetic data creation), often adds significant overhead due to per-example gradient computations, clipping, and noise addition, need to scale batch size and tune the hyperparameters increasing training time substantially (Ponomareva et al., 2023)

1.5 Evaluating synthetic data quality

Finally, the evaluation of synthetic data is also a complex problem, as the data must be assessed on a wide range of metrics to ensure its multi-purpose utility.

The utility of DP synthetic data is not absolute but highly dependent on the specific context and intended application. Data generated with a certain DP guarantee might exhibit good visual fidelity (e.g., low FID score for image data (Heusel et al., 2017)) but perform poorly when used to train a downstream classifier, or vice-versa (Lin et al., 2023). Similarly, data suitable for general trend analysis might yield incorrect results in specific statistical hypothesis tests due to DP-induced distortions (Perez et al., 2024). This implies that evaluating DP synthetic data requires a use-inspired approach (Ramesh et al., 2024). Simply stating data is “high quality” based on generic metrics is insufficient; its fitness for the specific purpose must be validated.

When choosing standard evaluation metrics that often work well, [Yang et al. \(2024\)](#) suggests to think about *fidelity* and *utility* angles.

Fidelity Fidelity metrics evaluate how closely the synthetic data matches statistical properties of the original data. Fidelity metrics are therefore often specific to each modality of the data we consider.

Utility For many practical applications the goal is to use synthetic data for some downstream task, for example training/finetuning an ML model. However obtaining such ultimate quality metric can be expensive during modeling and hyperparameter tuning process, which can require a large number of iterations.

A lot of works go the route of directly checking the performance of DP synthetic data on some simplified downstream tasks that server as proxy for the actual application. [Borisov et al. \(2023\)](#) referred to the metric that compares performance of various ML models on original and synthetic data as *Machine Learning efficiency (MLE)*. There is a design choice of what actual model one should train here - while academic works often train simple classifiers like Bert or XGBoost, the model should be as close in architecture and behavior to the final downstream model that will use this data.

Additionally to utility and fidelity angles, [Ramesh et al. \(2024\)](#) suggests to use the following metrics to gauge usefulness of the synthetic data.

Privacy leakage metrics Accessing potential privacy leakage from the sensitive data to the synthetic data is a very important task. There are a number of metrics that fall into this category (e.g. *Membership Inference* or *reconstruction attacks* success rates), and we cover the ways to calculate them in depth in Section 7.5

Fairness metrics In labeled original data (or data with some attributes that can be treated as labels of interest, e.g. cluster assignments) with available demographic information (e.g. race of humans in pictures, sex in tabular data, pronouns used with professions in text data etc.) fairness metrics provide an important angle. ([Liu et al., 2021a](#); [Ramesh et al., 2024](#)) suggest metrics that evaluate the sum of absolute value difference between the whole dataset and each subgroups (e.g. *false positive equality difference (FPED)* and *false negative equality difference (FNED)*) and *equalized odds* which evaluates whether false positive and true positive rates are the same for all groups.

We will cover modalities specific utility and fidelity metrics in subsequent Sections (3.8, 4.7, 5.6 and 6.6), however we want to emphasize once again that the most suitable metric for evaluating DP synthetic data quality will be always application specific. If for example downstream use of the data relies on preservation of temporal component between rows in a tabular dataset, standard fidelity and utility metrics will not be sufficient and metrics incorporating temporal correlations should be used for guiding synthetic data quality improvements.

2 A brief overview of Differential Privacy

Differential Privacy is a well established framework for reasoning about information leakage ([Dwork and Roth, 2014](#)). In this section we provide a concise overview of differential privacy (DP) concepts pertinent to DP synthetic data. First we introduce a definition of what DP guarantees (Section 2.1) and a concept of privacy unit that we will be using throughout this paper (Section 2.2). We will talk about where DP is commonly introduced in a lifecycle of machine learning and explain what DP synthetic data means (Section 2.4). We also briefly explore the concept of DP-Training, which is the workhorse algorithm that will be utilized as one of the main methods for DP synthetic data generation for all modalities in Section 2.5. For a more comprehensive introduction of DP, we refer the reader to ([Dwork and Roth, 2014](#); [Ponomareva et al., 2023](#)).

2.1 Definition of Differential Privacy

Differential privacy (DP) is not a single technique but rather a provable promise. DP provides a formal, mathematical definition of privacy centered on the principle of *indistinguishability*. The guarantee is that the outcome of a computation will remain almost exactly the same, regardless of whether any single record is included or not. This promise can be captured by the definition of (ϵ, δ) -DP ([Dwork and Roth, 2014](#)), also known as *approximate DP*.

Definition 1 ((ϵ, δ) -Differential Privacy). *A randomized algorithm \mathcal{M} satisfies (ϵ, δ) -DP if for any two neighboring datasets \mathbb{D}, \mathbb{D}' and for all $\mathcal{S} \subset \text{Range}(\mathcal{M})$:*

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

In definition 1, the randomized algorithm \mathcal{M} is the process imposing differential privacy and whose outcome is released. For machine learning (ML), \mathcal{M} is often the model training process; for synthetic data generation, \mathcal{M} can be the data generation process and the DP synthetic data is the outcome that is being released. The notion of *neighboring datasets* (datasets that are exact copies of each other and only differ by a single record, so called “privacy unit”) is discussed further in section 2.2.

The two parameters (ϵ, δ) , quantify the level of privacy protection. ϵ is the privacy budget, where a smaller non-negative value implies a stronger privacy guarantee. δ is a relaxation term, representing a probability of failure of satisfying the strict ϵ -bound. The presence of δ is why definition 1 is also called approximate DP. This relaxation is crucial for the practical application of many important DP mechanisms, such as the Gaussian mechanism common in ML. Practically speaking, δ is typically set to be small, and often smaller than the inverse the dataset size⁴, for example $1/|\mathbb{D}|^{1.1}$ ((Ponomareva et al., 2023) Section 5.3.2), so that the violation of pure ϵ -DP is rare. For ϵ , values below ten are considered a reasonable choice, while those under one offer a strong privacy guarantee. We provide a discussion on what values of (ϵ, δ) are appropriate for DP synthetic data generation in practice in Section 7.1.

While many other DP definitions exist (Ponomareva et al., 2023; Dwork and Roth, 2014), (ϵ, δ) -DP remains the popular choice for its flexibility and practical implication. We use the (ϵ, δ) -DP definition throughout this work unless otherwise specified.

DP provides a **worst-case guarantee**, protecting an individual’s data against an adversary even if an attacker has access to all other data in the dataset, possesses large computational resources, and has access to external datasets to link information. The threat model can be hypothetical and in practice the attacker often only has access to a single released outcome from $\mathcal{M}(\mathbb{D})$ - in context of this survey, this is DP synthetic data. DP guarantee is also data-independent, as it makes no assumptions about the underlying distribution of the dataset \mathbb{D} , or the one record in the adjacency definition. While the strong, future-proof guarantee simplifies the analysis and comparison of different mechanisms, the pessimistic assumptions can make achieving a favorable privacy-utility trade-off challenging. As a result, applying DP in practice, including for creating DP synthetic data, requires a large amount of computation.

2.2 Neighboring datasets and the privacy unit

The guarantee of differential privacy hinges on the definition of “neighboring datasets” (the datasets are the same and only differ in one record). This relationship is determined by two factors: the privacy unit, which specifies what constitutes one record, and adjacency, which describes how two neighboring datasets can differ.

The privacy unit is the fundamental element of data being protected and defines what constitutes “a single record” within the dataset. The choice of a privacy unit is a critical modeling decision that shapes the privacy guarantee. For example, under *instance/example-level DP*, the neighboring dataset differs on a single data entry. This is a straightforward approach, but it can be insufficient if one person contributes multiple records. In *user-level privacy*, all data points belonging to a single user are grouped and treated as a single privacy unit. This ensures the protection of an individual’s entire contribution, offering a stronger and more intuitive guarantee.

The adjacency definition determines how two datasets differ in one record: whether it is an addition or removal of one record (*add-or-remove adjacency*), *replace-one* (e.g. datasets are the same but one record in one dataset is replaced with another record to produce the second dataset) or *zero-out* (one record in one dataset is replaced with a dummy record representing a missing record) (Ponomareva et al., 2023). The *replace-one* definition is considered a stronger adjacency definition (“twice as strong” as zero-out or add-or-remove), essentially representing a combination of an addition and a removal of one record.

Pay attention to adjacency definition

When reporting privacy guarantees for DP artifacts including DP synthetic data, it is important to take note of the adjacency definition used. While *zero-out* and *add-or-remove* adjacency definitions have comparable semantics for a fixed ϵ , the guarantees for *replace-one* are considered approximately twice as strong (Ponomareva et al., 2023). Nevertheless, care should be taken when attempting to compare performance or guarantees under different adjacency definitions.

There are many important considerations in properly choosing the DP setting, and guarantees reported under different privacy units and adjacency definitions are generally not comparable. While it is possible to convert DP guarantees from example-level to user-level via group privacy theorems (Dwork and Roth, 2014; Charles et al., 2024), the resulting bound can be loose and the reverse conversion is not obvious. For a detailed discussion on selecting and reporting

⁴Dataset size is measured as a number of privacy units (Section 2.2.1) in the dataset

DP settings, we refer the reader to (Ponomareva et al., 2023). In the following sections, we will assume *add/remove adjacency*, and discuss privacy unit setting as a key consideration for various synthetic data modalities.

2.2.1 Privacy units for DP synthetic data

In subsequent section for each modality (Image, Tabular, Text, Federated Learning) we will provide a discussion of the most commonly used privacy unit and the ramifications of this choice. Below we will familiarize the reader with types of privacy units that will be encountered throughout this paper.

Choice of privacy unit

One of the most important decisions for any application of DP, including DP synthetic data generation, is the appropriate choice of privacy unit. The privacy unit determines the scope of the privacy guarantee (e.g., protecting examples versus protecting collections of examples owned by a user versus collections of examples that belong to a group of users who are part of an organization etc.). The choice of privacy unit also determines how the original data must be prepared for creating its DP synthetic version, influencing for example how training examples are formed and how data is sampled in DP fine-tuning.

Ponomareva et al. (2023) outlined several choices for the unit of protection in context of DP. Most of the time, in context of DP synthetic data generation, the choice will be among the following four:

1. *Example-level privacy unit* - in this setting, neighboring datasets from DP definition are defined as two datasets that differ in one record/instance/example only (Ponomareva et al., 2023). A definition of a record also offers some leeway based on the end use of the data. For example, for tabular data, each record is a row; however for a collection of the text documents each record can be defined as one full text document or a portion of a text document. For example, if the downstream task is to mark a paragraph as a summarizing paragraph vs a paragraph introducing new details, paragraph-level definition of record might be appropriate.
2. *User-level privacy unit*. If the original private data was generated by a large number of users who each may have contributed multiple records, this privacy unit might be more appropriate than example-level. Neighboring datasets in this setting will differ in inclusion/exclusion of *all* the data from a particular user. Consider the case of a user interacting with a banking assistant. In this setting, each user might contribute multiple interactions to the private dataset. Example-level DP protection can be insufficient to prevent the model from memorizing information common to multiple examples from the same user, like the bank account number or preferred bank location. Achieving user-level DP becomes more challenging for domains like chats or emails, where an individual example (chat or thread) might be associated with multiple individuals. For such domains, *multi-attribution user-level privacy* (Ganesh et al., 2025) is appropriate; achieving this requires more complex contribution-bounding strategies, which we discuss in Section 7.2. Figure 1 from Ganesh et al. (2025) compares various privacy units using an example of emails dataset.
3. *Larger than user-level privacy unit*. Finally, a larger privacy unit can be appropriate when private data is constructed from data from multiple groups, for example various specialized hospitals can contribute their patients data for creation of DP Synthetic data and each hospital wants to be assured that all their patients are protected as a whole group. However, this very-large privacy unit makes achieving DP much more difficult. For example, one typically needs a dataset of at least double digit (XX) thousands privacy units in order to utilize state-of-the-art techniques like DP finetuning of LLMs. While a dataset of this number of users in an industry setting might be considered “small”, a dataset of tens of thousands distinct hospitals or businesses might be very hard to obtain.
4. *Sub-user-level privacy units*. Units smaller than user level, e.g. user-month and other weaker notions of privacy can be suitable for some applications, particular when complimented by empirical privacy auditing (see Section 7.5) against real-world attacks. However in general we advocate for user-level and larger privacy units for DP synthetic data whenever possible

Discussion In order to implement DP algorithms, it must be possible to decompose the private dataset into these privacy units, or more generally, reason about how a change to any one unit might influence the final output. Thus, achieving true person-level privacy would require dataset where each example was perfectly annotated with unique and stable identifiers for each real person associated with that example. This is essentially never the case in practice. For example, user datasets might have metadata like user account ids, email address, but some people might have multiple accounts (weakening the formal privacy guarantee), while others might share an account. Thus, formally we will often only achieve “user-id-level DP” but not the true “person-level DP.” Fortunately, DP guarantees are well-behaved with respect to small changes in the privacy unit. *Group privacy* (Vadhan, 2017a) theorems can be utilized to convert guarantees for a “smaller” unit of privacy to the ones for a “larger” unit of privacy. For example, if one person had

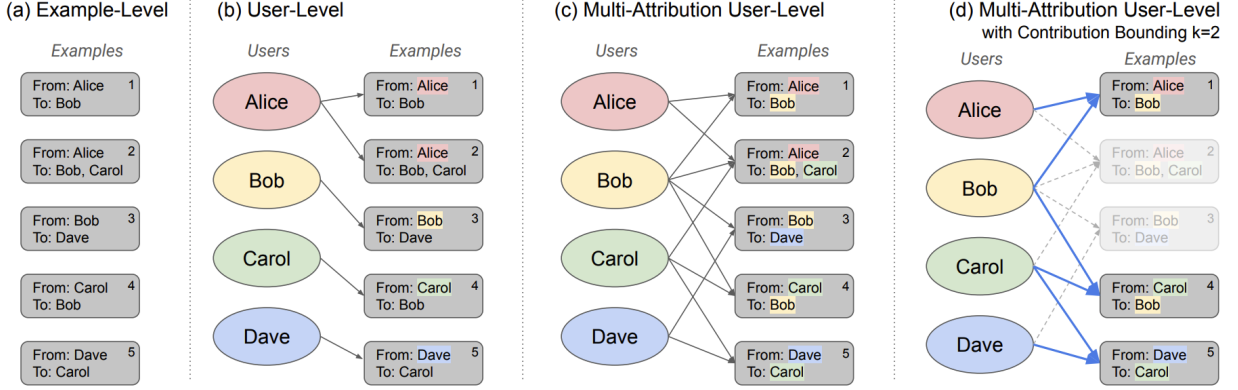


Figure 1: Visual representation of choice of privacy unit for a toy private dataset of emails, reproduced with permission from (Ganesh et al., 2025). Figure (a) treats the dataset as a flat collection of emails (*example-level privacy*). In this setting, secrets occurring in multiple emails are not protected. Figure (b) demonstrate the meaning of *user-level privacy unit* where each email is considered to belong only to the sender (*single-attribution user-level privacy unit*). User's secrets that are also present in emails the user received but didn't author might not be protected. Figure (c) attributes emails to all senders and receivers, and (d) bounds the contribution of each user to at most $k = 2$ examples from (c) by selecting a subset of emails that will be used for training.

two accounts in a system offering account-level DP, their effective DP ϵ might go from say 1.0 to 2.0 but it would not go to ∞ . A related issue is that even if the user-attribution metadata was perfect, it is possible that some private fact occurs in examples associated with multiple users — for example, many members of a family might independently discuss one family member's medical condition. Again, as long as such information only occurs in say 5 user's data (if using user-level DP), group privacy still offers protection for such information (though weaker than that offered to information only in one user's examples).

This discussion suggests that one must select the unit-of-privacy, as well as the strength of the protection (ϵ), with pragmatism. In particular, the sensitivity of the source data, the potential harm that would be caused by privacy leakage, and the likelihood and strength of potential attacks should all be considered when evaluating the privacy protections of an end-to-end DP synthetic data deployment. Empirical privacy auditing (Section 7.5) is a valuable tool in calibrating these assessments, as carefully constructed audits can both measure the effectiveness of attacks that are more realistic than the worst case assumptions of DP, and can also be used to measure risks that go beyond the chosen unit-of-privacy (for example, auditing risk when one real person contributed 10 examples, even when contribution bounding and the DP guarantee correspond to one example per username).

It is also worth noting that some methods allow one to simultaneously make DP claims for multiple units of privacy (e.g., a stronger example-level guarantee, and a somewhat weaker user-level guarantee), which can allow different people with different concerns to pick the guarantee that best protects against the threats they care about.

2.3 Key Properties and Mechanisms of Differential Privacy*

DP has several key properties that enable its practical adoption. We highlight a few key properties and discuss their relevance for synthetic data generation.

Postprocessing. One powerful property of DP is its invariance to post-processing. Any computation performed on the output of a DP algorithm (without re-accessing the original private data) incurs no additional privacy cost. This property is a natural result of the worst-case consideration of DP definition, and is exceptionally useful for synthetic data generation. Once a synthetic dataset is created using a DP algorithm, it can be analyzed, shared, and used to train any number of models without weakening the privacy guarantee. Moreover, as we will discuss in detail later, many modern synthetic data generation techniques focus on training a generative model with DP. Because of post-processing, one can safely sample from a DP generator to create an infinite number of synthetic data tailored for intended usage, and the resulting DP synthetic dataset will satisfy the same DP guarantee as the generator.

Composition* properties provide a framework for analyzing DP guarantees when multiple mechanisms are applied together. Common examples of composition include the following

1. *Parallel Composition*: If a dataset is partitioned into disjoint subsets and DP mechanisms are applied independently to each subset, the overall privacy cost is determined by the worst privacy cost of any single computation. Formally, if (ϵ_i, δ_i) -DP are satisfied on k disjoint data partitions, the combined mechanism satisfies $(\max_i \epsilon_i, \max_i \delta_i)$ -DP.
2. *Sequential Composition*: When a sequence of mechanisms is applied to the same dataset, their privacy costs accumulate. For the combined mechanism of a sequence of (ϵ_i, δ_i) -DP mechanisms, *basic composition* provides a loose bound from simple summation of the privacy parameters $(\sum_i \epsilon_i, \sum_i \delta_i)$. *Advanced composition* offers a tighter bound, which is essential for iterative algorithms. For example, repeating the (ϵ, δ) -DP mechanism for k times results in a mechanism satisfies $(\sqrt{2k \ln(1/\delta')}\epsilon + k\epsilon(e^\epsilon - 1), k\delta + \delta')$ -DP for all $\delta' \geq 0$ using Theorem 3.20 of (Dwork and Roth, 2014).

These composition principles are critical for machine learning and synthetic data generation, which often use complex algorithms that access the data multiple times and in different partitions. They are often used with **privacy amplification by sampling** (Kasiviswanathan et al., 2008), another powerful tool to make DP analyses more efficient: a key technique where applying a mechanism to a random data sample, rather than the whole dataset, significantly reduces the privacy cost. Together, these tools allow for a principled analysis (called *accounting*) of the rigorous DP guarantee for the complete learning process across all steps, which is fundamental to algorithmic design (Abadi et al., 2016).

Gaussian Mechanism The Gaussian mechanism is a foundational building block for many differentially private algorithms in machine learning. We discuss it here to explain core DP principles. We first focus on applying Gaussian mechanism once, and then discuss the usage in iterative algorithms like DP-SGD (Abadi et al., 2016) and DP-FTRL (Kairouz et al., 2021a). For a broader discussion of various DP mechanisms, we refer the reader to (Ponomareva et al., 2023; Dwork and Roth, 2014).

The Gaussian mechanism achieves DP by adding noise from a Gaussian distribution to the output of a real-valued function f . The amount of noise added depends on the function’s ℓ_2 -sensitivity, which bounds the maximum possible impact of any single privacy unit on the function’s output.

Definition 2 (ℓ_2 -sensitivity). For a function $f : \mathcal{D} \rightarrow \mathbb{R}^d$, the L_2 -sensitivity is:

$$\Delta_2(f) = \max_{\mathbb{D}, \mathbb{D}'} \|f(\mathbb{D}) - f(\mathbb{D}')\|_2$$

where the maximum is taken over all neighboring datasets \mathbb{D} and \mathbb{D}' .

A single application of Gaussian mechanism is defined as:

$$\mathcal{M}(\mathbb{D}) = f(\mathbb{D}) + \mathcal{N}(0, \sigma^2 \cdot \Delta_2(f)^2 \cdot \mathbf{I})$$

Here, noise is drawn from a zero-mean Gaussian distribution with a standard deviation of $\sigma \Delta_2(f)$, where σ is the *noise multiplier*. A larger σ provides a stronger privacy guarantee of smaller ϵ , and we can use $\sigma = \sqrt{2 \ln(1.25/\delta)}/\epsilon$ to achieve (ϵ, δ) -DP when $\epsilon < 1$, refer to Theorem A.1 of (Dwork and Roth, 2014).

2.4 DP synthetic data and methods for creating it

DP synthetic data, from the point of view of the consumer of the data, can be viewed as a privatized version of the original dataset (this corresponds to introducing DP at the input level/dp-fying the data (Ponomareva et al., 2023) Section 3.2). Once the data is dp-fied, anything done with this data is also DP due to postprocessing (Section 2.3) - any downstream ML model trained on DP synthetic data is DP w.r.t. the original training data, and any prediction of such model is also DP.

An alternative DP synthetic data is training each downstream ML model with DP via algorithms like DP-SGD (*DP-Training*) which we discuss next. This was previously prevalent in academic literature on DP ML. In this setting, by postprocessing property, once an ML model is dp-fied, any output of such ML model is DP as well (so both the model and model’s predictions are also DP w.r.t the original training data). Another alternative to DP synthetic data is introducing DP at the level of prediction of each ML model (*DP prediction or DP inference*) of an ML is the “weakest” form of DP protection - in this setting only predictions of an ML model are DP and can be shared. However in theory this type of DP entry point might require much less noise to realize, which can potentially mean much better utility of the predictions (Ponomareva et al., 2023).

DP synthetic data offers the strongest protection among the alternatives we discussed. The flip side is that creating a useful DP synthetic dataset is a task that in general is harder to achieve than for example DP-Training a particular ML

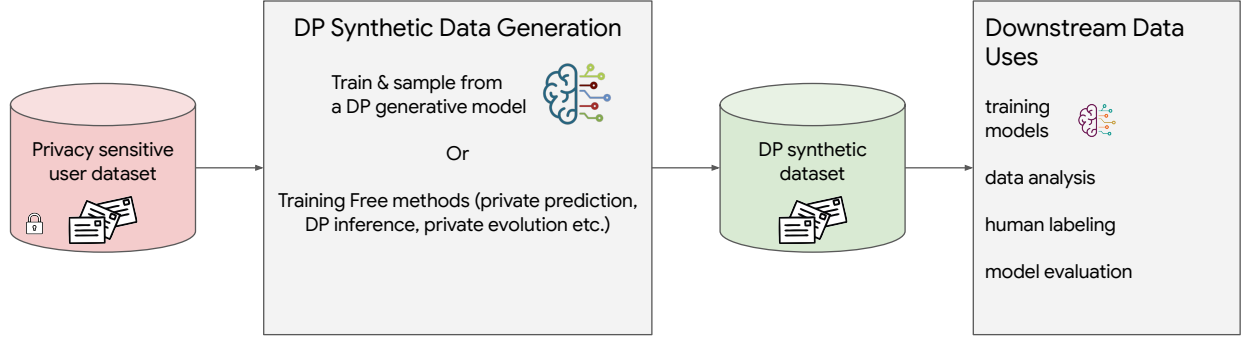


Figure 2: DP synthetic data can be used as a drop-in replacement for the original private data and can be safely used for training downstream models, data analysis, human labeling and other applications. DP synthetic data is created either with the help of DP-Training of ML models or Training free methods.

model on private data with DP, as a synthetic dataset must serve multiple, even unforeseen, purposes. The future-proof, worst-case nature of the DP guarantee makes the privacy-utility-computation trade-off difficult to navigate. However, a major advantage of DP synthetic data is its utility as a safe and shareable asset, a direct result of DP’s worst-case guarantee and invariance to post-processing, as discussed in Section 2.3. Synthetic data without explicit privacy protection are vulnerable to attacks; for example, LLMs can memorize and reconstruct training data (Carlini et al., 2021; Nasr et al., 2025).

While DP can be applied to directly train downstream models to be deployed, or during inference time for a specific task, DP synthetic dataset is more versatile. Once a DP synthetic dataset is generated, it can be distributed, analyzed, and used for unlimited analysis and modeling tasks without incurring additional privacy costs. DP synthetic data additionally can allow one to reduce computation cost by creating DP synthetic dataset once and reusing it (without DP) for training various ML models, instead of training each model with DP (a computationally expensive process).

Broadly speaking, there are two groups of methods for creating DP synthetic data: **those that involve modification of the training process of models that will be used to generate DP synthetic data**, to ensure differential privacy (DP-Training, Section 2.5) and **those that avoid training ML models on private data** and instead rely on inference-only access to trained ML models. This categorization will be present when we describe methods for modalities like Images (Section 4) and Text (Section 5). Tabular data, which is unique modality with well established history of theory backed research into dp synthesis method, is the only modality that enjoys a more diverse suite of methods that we explore in Section 3.

2.5 DP-Training

DP-Training is a workhorse algorithm for obtaining differentially private ML models that, as we will see in subsequent sections, can be successfully employed to create DP synthetic data.

*DP-Training*⁵ commonly refers to a DP variant of SGD-like algorithms (DP-SGD (Abadi et al., 2016) and DP-FTRL (Kairouz et al., 2021a)) that iteratively apply Gaussian mechanism to achieve DP protection of ML model gradients. DP-fied gradients are then used to do gradient step update, resulting in all the gradient sequences, and consequently, all the checkpoints of a trained ML model to be DP w.r.t. the original training data.

Algorithm 1 demonstrates modifications that are introduced to SGD algorithm to ensure DP. The changes include per example gradient clipping to limit ℓ_2 sensitivity and an addition of the noise via Gaussian mechanism to aggregated gradients. Similar modifications can be used to obtain DP versions of algorithms like Adafactor, Adam etc. Clipping the gradients is computationally expensive step and is the slowest part of DP-SGD, since most modern auto-differentiation based frameworks do not provide easy access to per example gradients. Additionally, in order to obtain good utility, DP-training requires significantly larger batch sizes in order to tolerate introduced noise, and potentially more iterations to converge (Ponomareva et al., 2023). Finally, other hyperparameters like learning rate and clipping norm require extensive hyperparameter tuning to achieve good utility. All of the above make DP-Training much more computationally expensive than training an ML model without DP.

⁵Throughout this work, we will use terms like *DP-Training* and *DP-finetuning*. Both terms imply an application of DP algorithm like DP-SGD. DP-training is most commonly used to describe training a model from scratch with DP, whereas DP-finetuning often is used to describe finetuning a publicly pretrained model like an LLM with DP. Some works use these terms interchangeably.

In order to provide DP guarantees of DP-SGD algorithm, tracking the cumulative privacy cost across multiple iterations of DP-SGD needs to be performed (this task is often referred to as *privacy accounting*). Combining Gaussian mechanism $\epsilon = \sqrt{2 \ln(1.25/\delta)}/\sigma$ and advanced composition $\epsilon' = \sqrt{2T \ln(1/\delta')} \epsilon + T \epsilon (e^\epsilon - 1)$ from Section 2.3 gives us a rough understanding of the relationship between DP guarantee and the noise multiplier for full-batch DP-SGD. The number of iterations T also affects the DP guarantee, and data processing pattern is another important design choice when applying advanced techniques of privacy amplification by sampling (Abadi et al., 2016; Balle et al., 2018; Chua et al., 2024b) or correlated noise mechanisms (Kairouz et al., 2021a; McMahan et al., 2024).

While basic and advanced compositions (Section 2.3) provide a theoretical understanding, they often yield loose bounds in practice. Modern privacy accounting is typically handled by specialized libraries (DP Team, 2022) that implement more sophisticated methods to achieve tighter bounds. These methods evolved from moments accountant and Rényi Differential Privacy (Abadi et al., 2016), to DP variants (Gaussian DP (Dong et al., 2022) and zero-concentrated DP (Bun and Steinke, 2016)) more suitable to Gaussian mechanism and can be cleanly converted to (ϵ, δ) -DP. The recent development of Privacy Loss Distribution (PLD) accounting tracks a more detailed representation of the privacy loss at each step, often yielding the tightest known (ϵ, δ) bounds for a given process (DP Team, 2022). These advanced accounting methods are essential for designing high-utility DP mechanisms, including training synthetic data generators, within a reasonable privacy budget.

Algorithm 1 The DP-SGD algorithm, adapted from (Ponomareva et al., 2023), based on (Abadi et al., 2016).

Require: Training data, consisting of features $X := \{x_1, x_2, \dots, x_N\}$ and labels $Y := \{y_1, y_2, \dots, y_N\}$.
 $h(x; \theta)$ is the model applied to an input x and parameterized by θ .
 $L(y, y')$ is the loss function for label y and prediction y' .
SGD hyperparameters: η learning rate, T number of iterations, B batch size.
DP hyperparameters: C clipping threshold, σ noise multiplier, δ (for privacy accounting).
Ensure: θ_T final model parameters
 $\theta_0 \leftarrow$ initialization
for $t \leftarrow 1$ to T **do**
 Randomly sample a batch B_t with sampling probability B/N for each data point.
 Data are sampled with replacement for each batch.
 for $i \in B_t$ **do**
 $g_t(x_i) \leftarrow \nabla_{\theta_t} L(y_i, h(x_i; \theta_t))$ ▷ Compute per-example gradient wrt the weights
 $\tilde{g}_t(x_i) \leftarrow g_t(x_i) / \max(1, \frac{\|g_t(x_i)\|_2}{C})$ ▷ Clip the per-example gradient
 $\bar{g}_t \leftarrow \frac{1}{B} (\sum_i g_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$ ▷ Add noise
 $\theta_{t+1} \leftarrow \theta_t - \eta \bar{g}_t$ ▷ Gradient descent step

3 DP Synthetic Tabular data

Tabular data (data that can be represented as tables of values organized into rows and columns), has long been the most prominent way to organize data in various domains like the sciences, healthcare, business, finance, retail and others. Unsurprisingly, initial attempts at generating DP synthetic data focused on tabular data, which has developed into one of the richest lines of work in the DP synthesis literature.

The shared goal among the various works in this space is to generate a synthetic data that shares the same “statistical characteristics” as the real data, such that it can be a drop-in replacement for the real-data in downstream pipelines with minimal utility degradation. Exactly what form these “statistical characteristics” take varies across approaches, typical evaluations focus on distributional similarity of low-order marginals or train-on-synthetic, test-on-real classification tasks. The space of mechanisms can be partitioned into two groups: *workload-adaptive* methods, where the mechanism tailors the synthetic data towards a specific (finite) set of downstream tasks, and *workload-agnostic* methods, where the mechanism aims for general distributional similarity, without targeting accuracy w.r.t. any specific set of tasks.

For example, some algorithms focus on ensuring the synthetic data distribution is “close” to the real one in a principled statistical distance, like Wasserstein distance. One could argue that this *implicitly* defines the workload as a broad (potentially infinite) class of smooth, Lipschitz queries. If the two distributions are close in Wasserstein distance, then we also get a bound on the error of *any* of these smooth queries, thus giving general assurances against unanticipated future analyses of this kind. However, this does not *adapt* to the workload in that this broad guarantee surely covers a greater class of queries than the user could reasonably pre-specify. It is also worth mentioning that this same rationale applies more broadly to *any integral probability metric*, e.g. Kolmogorov distance and total variation distance, for a suitable class of (infinitely-many) linear queries Müller (1997).

Marginal-based algorithms

Tabular mechanisms typically operate by taking noisy measurements and creating synthetic data that conforms to them. These measurements may or may not be related to a user’s potential workload. When an algorithm measures marginal queries, they are often referred to as **marginal-based algorithms**.

Mechanisms can also be classified as *theoretically oriented algorithms*, which aim to provide rigorous bounds on the utility of the algorithm across all possible datasets, and *empirically oriented algorithms*, which aim to provide good utility across a representative set of benchmark datasets and measurements. Empirically-oriented algorithms often outperform theoretically-oriented algorithms on real world datasets, both in terms of utility and scalability, at the cost of limited (if any) formal guarantees on the quality of the generated synthetic data. Indeed, it is not even possible to run many theoretical algorithms on real-world benchmark datasets due to challenges of scale. While the algorithmic ideas underlying both classes of approaches are largely similar, small implementation details can make *big* differences in practice.

Remark 1 (Workload vs. Measurements). *In the context of DP tabular data synthesis, it is crucial to distinguish between two concepts. The **workload** refers to the set of statistical queries a synthetic dataset is intended to be accurate on for downstream analysis. The **measurements** are the set of queries an algorithm actually executes on the private data to guide the generation process. All algorithms (apart from some of the end-to-end models from Section 3.6) take measurements; for instance, even methods that provide distributional distance guarantees can be seen as implicitly targeting certain workloads (e.g. the Chebyshev moments of a distribution (Musco et al., 2024)).*

Remark 2 (Workload-Adaptive vs. Workload-Agnostic). *The philosophies behind workload-adaptive and workload-agnostic approaches are fundamentally different. The workload-adaptive approach tailors the synthetic data to a known set of analytical needs. By leveraging this explicit guidance, these mechanisms can strategically allocate the privacy budget, significantly enhancing accuracy for the tasks that matter to the downstream analyst. In contrast, the workload-agnostic philosophy aims for a more general, “one-size-fits-all” solution by attempting to learn the most salient features of the joint probability distribution automatically. However, this ambitious goal can lead to higher-than-expected error on a specific downstream workload that was not implicitly prioritized by the algorithm’s measurement strategy. While the workload-adaptive approach risks poor performance on “off-workload” queries, this can be mitigated by providing a broad workload (like all 3-way marginals).*

The core design choice in existing algorithms for DP tabular synthetic data is not whether to take measurements, but how the measurement strategy is selected in relation to a potential workload.

Roadmap As a roadmap for this chapter, we will explore methods that fit into both of the aforementioned categorizations. We must first, however, discuss the fundamental notion of the **privacy unit** for tabular data in Section 3.1, as this underpins all subsequent algorithms. Additionally, in Section 3.2, we describe the problem of *query release* that lays the foundation of many methods for DP tabular data synthesis, as well as the *histogram* representation of tabular data that many methods use either implicitly or explicitly.

In Section 3.3, we will delve into the family of **select-measure-estimate** algorithms, which represent the most prominent and successful class of practical, empirically-oriented methods. We subsequently cover theoretically-oriented algorithms in Section 3.4, detailing the foundational algorithms that come with strong theoretical guarantees: those targeting additive error and those targeting distributional distance. We defer details on data preparation and discretizations, which are often necessary for practical implementations, to Section 3.5.

Finally, in Section 3.6, we return to empirically motivated algorithms and introduce a newer strain of work based on **end-to-end generative models** that take advantage of modern advances in deep learning, another important class of potentially practical algorithms that generally offer a different structure than traditional algorithms from the select-measure-estimate paradigm. We give a brief comparison of these different families of methods in Section 3.7, and finally, in Section 3.8, discuss methods for evaluating the quality of synthetic tabular data, which is itself an active research area.

3.1 Privacy unit for tabular data

For tabular DP synthetic data, the privacy unit is most commonly assumed to be on an *example-level* (row-level in context of tables) i.e. where each user contributes at most 1 row to the dataset. As such, DP synthetic tabular data guarantees that the resulting synthetic dataset would not be significantly different whether a particular user contributed their row to the original tabular dataset or not. Neighboring datasets are then two tabular datasets that differ in at most 1 row. It is worth noting that supporting user-level privacy with each user contributing at most x rows could be handled by creating one aggregate row for each user using x user’s datapoints, or by using a group privacy property and running

the algorithms we describe below with ϵ/x and δ/x in order to achieve x group level privacy of (ϵ, δ) (Dwork and Roth, 2014). Ensuring that each user has at most x rows can be done with appropriate user contribution bounding (Section 7.2). Developing better user-level algorithms that directly handle varying number of rows per each user remains an open research direction.

We note here an alternative, distinct privacy unit setup, which to the best of our knowledge remains relatively unexplored in the literature, where each user is assumed sole authorship of an entire table i.e. the privacy unit is *itself* a full tabular dataset. In this setup, each user can contribute one or more tabular dataset(s), over potentially varying domains (e.g. different columns headers). Thus, the neighboring definition is two tabular collections differing in at most 1 (or x) tabular datasets. The goal is then still to create synthetic data that satisfies a DP guarantee, but in this case the output is a *collection* of tabular datasets (as opposed to a *single* tabular dataset). While such a use case might be boiled down to the previous setup (e.g., if each tabular dataset has the same columns and their domains), in general this problem differs fundamentally and is significantly harder than the single table input/output model. All methods discussed in the following sections assume the former *row* level privacy unit definition, where we have a single input table on a fixed domain and output a synthetic version of it over the same domain with a DP guarantee, though we believe exploring the *table* level privacy unit is a promising avenue for future work.

3.2 Foundations of DP tabular synthetic data

DP synthetic tabular data was originally explored in context of *query release* – the task of creating useful synthetic data that gives accurate answers to a number of predefined statistical queries (**workload**) performed on this data, e.g. sums and counts. While one can answer such statistical queries by adding appropriate noise to each query result, DP synthetic tabular data allows one to both answer the predefined queries and perhaps additional queries that will be defined later (there are generally no guarantees about the utility of answers to such queries).

3.2.1 Foundation of workload-adaptive algorithms: The Query Release Problem*

More formally, the synthetic tabular data *query release* problem can be formulated as follows (refer to Table 1 for the full list of notations throughout this chapter). Let \mathcal{X} be the data universe from which rows of the tabular dataset are drawn (for now we assume it to be a finite set); for example, if our data contains d binary attributes, then $\mathcal{X} = \{0, 1\}^d$. Let \mathcal{F} be a finite set of functions that we will refer to as the *counting queries* or a (*query*) *workload*, where for each $f \in \mathcal{F}$, $f : \mathcal{X} \mapsto \{0, 1\}$ and f is a per-row predicate, and we desire an estimate of the sum $f(D) = \sum_{i=1}^n f(x_i)$ (the count of input rows that satisfy the predicate f). Note that the choice of (binary) counting queries can be easily extended to continuously-valued queries $f : \mathcal{X} \mapsto [0, 1]$, with the same error rates. This follows from the fact that these queries enjoy the same sensitivity bounds as their binary counterparts.

Recall that, as we discussed at the beginning of this section, the relationship between a DP synthetic data algorithm’s internal *measurement strategy* and a user-defined evaluation *workload* is the core of our taxonomy, a distinction that has its roots in the query release problem. As we noted previously, a *workload-adaptive* algorithm is one whose measurement strategy can be guided by the user’s workload to optimize accuracy (e.g., an algorithm where we could explicitly align workload queries with measurement queries). A *workload-agnostic* algorithm employs a measurement strategy designed to preserve general distributional properties, irrespective of any specific downstream workload (e.g., an algorithm that does not incorporate the workload queries into noisy measurements explicitly).

The goal of the workload-adaptive strategies for *query release* problem is to design an algorithm $\mathbf{A} : \mathcal{X}^n \mapsto \mathcal{X}^n$ (here, n represents the number of rows in the original dataset)⁶ that is (ϵ, δ) -DP (w.r.t. the privacy unit) and that is α -accurate with respect to \mathcal{F} for some $\alpha \geq 0$ and some norm function $\|\cdot\|$. Namely, for any input tabular dataset $D \in \mathcal{X}^n$, and its DP synthetic counterpart $\hat{D} = \mathbf{A}(D)$, it must hold that its error over the finite set of queries is bounded in **expectation** by α . With some abuse of notation, letting $\mathcal{F}(D) = (f(D))_{f \in \mathcal{F}}$ be the vector of workload query answers, this utility measure can be expressed as:

$$\text{Err}_p(\mathbf{A}, D) := \mathbb{E}_{\hat{D} \sim \mathbf{A}(D)} \left[\left\| \mathcal{F}(\hat{D}) - \mathcal{F}(D) \right\|_p \right] \leq \alpha, \quad (1)$$

Taking $p = \infty$ is a common choice in the theory literature, while $p = 1$ or $p = 2$ is common in the empirical literature. All of these options have their merits and the right choice ultimately depends on the downstream use case. We note that this objective is an *absolute* error rather than a relative one – it is typically easier to analyze and optimize mechanisms

⁶We opt here for algorithms whose input and output have the same number of rows. This is made for simplicity and can be generalized. Furthermore, it may be the case (e.g., for add/remove neighboring relations) that the number of rows is itself private; in this case, one can privatize the number of rows and produce synthetic data with that number of rows. Finally, for datasets with varying cardinality one needs to approximately preserve averages, rather than counts, as done in Equation 1.

for this setting. In the rest of the section, we will drop p from $\|\cdot\|$ with the understanding that it can in principle be anything, unless it is important to explicitly call out.

Symbol	Description
\mathcal{X}	The data universe from which tabular dataset rows are drawn; for example, if our data contains d binary features, then $\mathcal{X} = \{0, 1\}^d$.
n	Number of rows in the original tabular dataset and number of rows in the synthetic dataset.
D	Original tabular dataset, $D \in \mathcal{X}^n$
\hat{D}	DP synthetic tabular dataset, $\hat{D} \in \mathcal{X}^n$
f	Row-level predicate $f : \mathcal{X} \mapsto \{0, 1\}$
$f(D)$	$f(D) = \sum_{i=1}^n f(x_i)$ - number of rows in the dataset satisfying predicate f
\mathcal{F}	Finite set of queries f that synthetic tabular data must perform well on
\mathbf{A}	(ϵ, δ) -DP mechanism $\mathbf{A} : \mathcal{X}^n \mapsto \mathcal{X}^n$ for generating DP synthetic tabular data from the original tabular dataset. Alternatively, when talking about privatizing histograms, $\mathbf{A} : \mathcal{X}^n \mapsto \mathbb{N}_+^{ \mathcal{X} }$
α	Expected error (expected distance between synthetic and private datasets, in some norm) of the DP synthetic dataset
$h(D)$	Histogram representation of a dataset D , $h(D) \in \mathbb{N}^{ \mathcal{X} }$, where $h(D)_x := \{i \in [n] : x_i = x\} $. This histogram essentially represents, for each combination of column values, number of rows that have this combination. For example for a dataset with columns A and B (and possible values a_1, a_2 and b_1, b_2 respectively), the histogram would have bins (a_1, b_1) , (a_1, b_2) , (a_2, b_1) and (a_2, b_2) with counts of rows that fall into such combination of attribute values.
\hat{h}	DP synthetic data histogram

Table 1: Table of notations for tabular data discussion

3.2.2 Histogram representation: from data and back to data

Most DP synthetic tabular data generation approaches leverage a linear-algebraic reformulation of the problem by representing tabular datasets as *histograms*. Some methods directly instantiate a histogram while others do not; in either case it is important to understand how a dataset could be represented as a histogram and how one can “go back” from this histogram representation to form a corresponding tabular synthetic dataset (with a DP guarantee).

Let $D \in \mathcal{X}^n$ be a dataset with n records, and let $h(D) \in \mathbb{N}^{|\mathcal{X}|}$, where the histogram value for a bin x is $h(D)_x := |\{i \in [n] : x_i = x\}|$. This histogram representation counts the frequencies of rows for all possible combinations of column values, and it is useful as it linearizes the search space. Namely, using this histogram representation, query results can be expressed as inner products: $f(D) = \sum_{x \in \mathcal{X}} h_x f(x) = \langle h(D), f \rangle$.

There is a one-to-one correspondence between datasets $D \in \mathcal{X}^n$ and *integer-valued* histograms $h \in \mathbb{N}^{|\mathcal{X}|}$. However, many approaches require relaxing integrality constraints, i.e., working with real-valued histograms $\hat{h} \in \mathbb{R}_+^{|\mathcal{X}|}$, which prevents a direct construction of a dataset given a histogram. We highlight some techniques in the following bullet points to bridge this gap.

- In the context of linear programming approaches, there exist a host of *rounding techniques*, which approximate a continuously-valued vector $\hat{h} \in \mathbb{R}_+^{|\mathcal{X}|}$ by an integer-valued vector $\hat{h}' \in \mathbb{N}_+^{|\mathcal{X}|}$ (Williamson and Shmoys, 2011). The latter can be directly used to construct a dataset. Examples of this technique for DP synthetic tabular data can be found in e.g. (Barak et al., 2007; Dwork et al., 2009), as well as in the Top Down Algorithm used in 2020 US Census Disclosure Avoidance System (Abowd et al., 2022).
- A more general approach is, given a continuously-valued histogram $\hat{h} \in \mathbb{R}_+^{|\mathcal{X}|}$, its normalization $\hat{h}/\|\hat{h}\|_1$ induces a probability distribution over \mathcal{X} . Sampling n datapoints, which we can call a sample dataset $\tilde{D} = (\tilde{x}_1, \dots, \tilde{x}_n)$ i.i.d. from this distribution satisfies the following approximation bounds with high probability (this is known as Maurey’s Empirical Method or the Approximate Carathéodory Theorem (Pisier, 1981))

$$\left\| (f(\tilde{D}) - \langle \hat{h}, f \rangle)_{f \in \mathcal{F}} \right\|_\infty = O\left(\sqrt{n \log |\mathcal{F}|}\right).$$

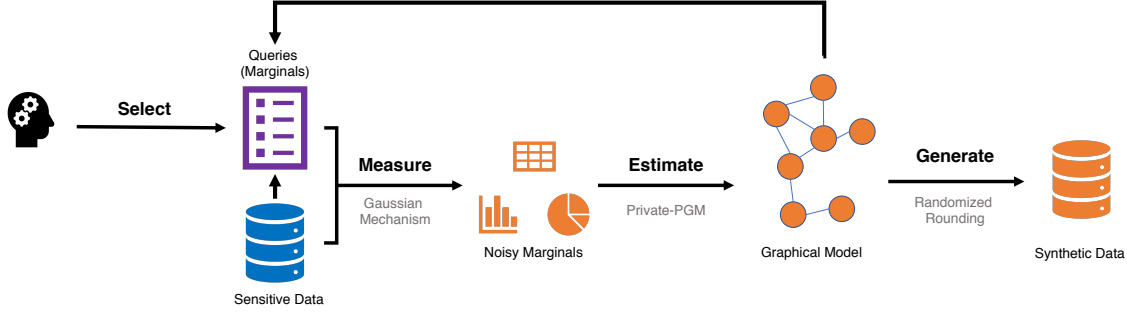


Figure 3: The select-measure-estimate paradigm. (1) A mechanism or domain expert *selects* a set of queries, possibly based on the workload and/or data. (2) The selected queries are measured by a standard noise addition mechanism like the Gaussian mechanism. (3) The data distribution is *estimated* from the noisy measurements. (4) The process optionally repeats from step (1) until the privacy budget is spent. (5) Synthetic data generated from the data distribution via a sampling or randomized rounding procedure.

This approach is advised in settings where the number of rows is relatively small compared to the number of bins or queries.

These techniques can be applied in a post-hoc fashion, and their results will be DP as long as the histogram \hat{h} was computed privately, due to the standard DP post-processing argument.

Remark 3. *Histogram representations of data provide significant mathematical convenience for designing and reasoning about mechanisms, although their linear size in $|\mathcal{X}|$ (which results in an exponential size in the number of attributes) can be prohibitive in high-dimensional settings.*

3.3 Practical Workload-Adaptive Algorithms

The following methods exemplify a practical, empirically-driven class of **workload-adaptive** algorithms that prioritize utility for a user-specified workload. Many mechanisms in this class are members of what we refer to as the *select-measure-estimate* paradigm.⁷

3.3.1 Practical Select-Measure-Estimate Algorithms

The general framework is given schematically in Figure 3. This framework encompasses a majority of practical algorithms for tabular data generation. Some additional algorithms that can’t be easily represented by this paradigm will also be discussed in Section 3.4.

This framework provides a modular approach to the privacy-preserving synthetic tabular data generation problem. It decomposes the challenge into a sequence of more tractable subproblems, most notably decoupling the task of *query selection* (i.e., deciding what statistical properties of the data to preserve) from *data estimation* (i.e., reconstructing a full data distribution from noisy answers to those queries). This separation of concerns has proven to be a powerful research catalyst, allowing for independent innovation on selection strategies (McKenna et al., 2022; Cai et al., 2021a, 2023; Fuentes et al., 2024) and estimation engines (McKenna et al., 2019, 2021b; Aydoore et al., 2021; Liu et al., 2021b,c; Mullins et al., 2024), which can then be combined in a “plug-and-play” manner. This modularity has enabled the development of general-purpose tools that free the mechanism designer to focus on the critical question of what statistics to measure to maximize the utility of the final synthetic data for downstream tasks.

The paradigm is generally understood through a sequence of five foundational steps, which together transform a sensitive dataset into a shareable, synthetic one. The process begins with **Select**, where a collection of queries is chosen

⁷This paradigm has been referred to under different names in the literature, including select-measure-generate (McKenna et al., 2021a), select-measure-reconstruct (McKenna et al., 2018), and Adaptive Measurements (Liu et al., 2021c).

Algorithm 2 AIM (Simplified)**Require:** Dataset D , workload W , privacy parameter ρ **Ensure:** Synthetic Dataset \hat{D} **Hyper-Parameters:** rounds T , $\alpha = 0.9$, MAX-SIZE (in Megabytes)1: Initialize \hat{p}_0 using an independent set procedure (see (McKenna et al., 2022) for exact algorithm)2: $\rho_{\text{used}} = 0, t = 0$ 3: Initialize schedulers for ϵ_t and σ_t 4: **while** $\rho_{\text{used}} < \rho$ **do**5: $t \leftarrow t + 1$ 6: $\rho_{\text{used}} \leftarrow \rho_{\text{used}} + \frac{1}{8}\epsilon_t^2 + \frac{1}{2\sigma_t^2}$ 7: $C_t = \{r \in W \mid \text{size}(r_1, \dots, r_t) \leq \frac{\rho_{\text{used}}}{\rho} \cdot \text{MAX-SIZE}\}$ 8: **select** $r_t \in C_t$ using exponential mechanism with ϵ_t budget:

$$q_r(D) = w_r \left(\|M_r(D) - M_r(\hat{p}_{t-1})\|_1 - \sqrt{2/\pi} \cdot \sigma_t \cdot n_r \right)$$

9: **measure** marginal on r_t :

$$\tilde{y}_t = M_{r_t}(D) + \mathcal{N}(0, \sigma_t^2 \mathbb{I})$$

10: **estimate** data distribution using PGM:

$$\hat{p}_t = \underset{p \in S}{\operatorname{argmin}} \sum_{i=1}^t \frac{1}{\sigma_i} \|M_{r_i}(p) - \tilde{y}_i\|_2^2$$

11: **anneal** parameters $\epsilon_{t+1}, \sigma_{t+1}$ using a scheduling procedure (see (McKenna et al., 2022) for exact algorithm)12: **generate** synthetic data \hat{D} from \hat{p}_t using PGM13: **return** \hat{D}

to be measured on the private data. In the context of tabular data, these queries are typically low-dimensional marginals, which are contingency tables that capture the counts of various combinations of attribute values. The choice of these marginals fundamentally determines which correlations and statistical patterns will be preserved in the final output. The second step is **Measure**, where the selected queries are executed on the sensitive dataset, and carefully calibrated noise is added to the true answers to satisfy (ϵ, δ) -differential privacy. The Laplace and Gaussian mechanisms are frequently employed for this purpose, with the Gaussian mechanism often being preferable when a larger number of marginals need to be measured (McKenna et al., 2022).

The third step is **Estimate**, which is a post-processing phase that takes the noisy, and often inconsistent, measurements from the previous step and seeks to find a coherent probability distribution over the entire data domain that best explains them. This is typically formulated as a maximum likelihood estimation (MLE) problem, where the goal is to find a distribution that minimizes the distance (e.g., ℓ_2 distance for Gaussian noise) to the noisy measurements. For some advanced strategies, this process is not a one-shot operation but is instead iterative, leading to the optional fourth step, **Repeat**. In adaptive mechanisms like PMW or MWEM (Hardt et al., 2012a) and its more sophisticated and scalable variants (McKenna et al., 2022; Cai et al., 2021a; Aydore et al., 2021; Liu et al., 2021c), information from the current estimated distribution is used to guide the selection of the next marginal query to measure, often targeting the one for which the current model exhibits the highest error. This creates a cycle of progressive refinement, focusing the privacy budget where it is most needed. The final step is **Generate**, where the fully estimated probability distribution can be converted to either a sampling process or randomized rounding process as described in Section 3.2.2.

A cornerstone of this entire paradigm is the post-processing property of differential privacy. The “Estimate” and “Generate” steps operate exclusively on the already-privatized noisy measurements and do not access the original sensitive data. Consequently, any computational effort expended during these steps – no matter how complex the optimization or how sophisticated the generative model – incurs no additional privacy cost. This crucial property has motivated research focused on building powerful, general-purpose reconstruction algorithms. By creating robust and efficient estimation engines, the paradigm allows mechanism designers to abstract away the complexities of data generation and concentrate on the distinct and equally challenging problem of query selection, which is more important for synthetic data utility (Rosenblatt et al., 2024a).

3.3.2 Selection Methods

The **Select** step is the first and arguable most critical component of the select-measure-estimate paradigm (McKenna et al., 2022; Rosenblatt et al., 2024b). The principle is simple: we cannot guarantee that information that is not *directly* measured will be *preserved*. Even a perfect, error-free estimation algorithm is powerless to reconstruct statistical patterns if they were not captured by the noisy measurements taken. This reality underscores the importance of developing intelligent strategies for selecting which queries to measure, a challenge that has motivated a significant and evolving line of research. The quality of a selection strategy can be evaluated along multiple dimensions; the work on the **AIM** mechanism (Algorithm 2) proposes a useful taxonomy of four criteria for an ideal strategy: *workload-awareness*, *data-awareness*, *budget-awareness*, and *computation-awareness* (McKenna et al., 2022).

Workload-awareness refers to a mechanism’s ability to prioritize measurements that are most important for a known, downstream analytical task, (user-specified “workload”). Mechanisms like **AIM** are highly workload-aware; they employ a quality score that explicitly weights candidate marginals based on their relevance to the user’s specified workload queries. Similarly, **MWEM+PGM** and **RAP** are designed to optimize for a given set of queries (Aydore et al., 2021). In contrast, methods like **MST** and **PrivBayes** are workload-agnostic (Zhang et al., 2017; McKenna et al., 2021a). They select marginals to measure based on general statistical properties of the data, such as attribute correlations, without considering any specific end use. This makes them suitable for general-purpose data exploration but potentially suboptimal and wasteful of the privacy budget when the analytical goal is known in advance. The trade-off is between specialized, high-utility data for a specific task and more broadly applicable, but potentially less accurate, data for general use.

Data-awareness describes the extent to which a selection strategy adapts to the underlying statistical properties of the data itself. Early strategies were often data-oblivious; a major step forward came with data-driven methods like **PrivBayes** and **MST**, which first use a portion of the privacy budget to privately learn a dependency structure from the data – a Bayesian network or a maximum spanning tree, respectively – and then use this structure to guide the selection of marginals to measure. A more advanced form of data-awareness is found in iterative and adaptive mechanisms like **AIM** and **PMW**. These algorithms greedily select the next marginal to measure based on a quality score that evaluates the error of the *current* synthetic data model. This allows them to focus the privacy budget on the aspects of the data distribution that are currently least well-approximated, leading to a more efficient learning process.

Budget-awareness concerns the intelligent allocation of the privacy budget, ϵ . Naive strategies might simply divide the total budget evenly among a predetermined number of measurements. More sophisticated mechanisms, however, recognize that the available budget dictates the achievable signal-to-noise ratio and should influence the selection strategy. For instance, with a larger budget (and thus lower noise), it becomes feasible to accurately measure larger, higher-dimensional marginals. Mechanisms like **PrivBayes** and **PrivSyn** are considered budget-aware because they adjust the number and size of selected marginals based on the total budget (Zhang et al., 2021). **AIM** introduces a particularly advanced form of budget-awareness with its “annealing” procedure. It adaptively adjusts the per-round privacy budget during its iterative selection process. If the model’s accuracy is not improving, it increases the budget for subsequent rounds, effectively “unlocking” new candidate marginals that were previously too noisy to be useful. This allows **AIM** to perform robustly across a wide range of ϵ values without manual hyperparameter tuning. A related concept, “frugal budgeting,” appears in **JAM-PGM**, which can use public data for some measurements (Fuentes et al., 2024). When a noise-free public measurement is chosen, the privacy budget allocated for that round is saved and rolled over, concentrating the budget for more difficult private measurements later on.

Computation-awareness is a crucial practical consideration, especially for mechanisms that rely on **Private-PGM** for estimation. The selection of marginals directly determines the structure of the graphical model used in the estimation step. If the selected marginals induce a graph with high treewidth, the estimation phase can become computationally intractable, failing due to excessive time or memory requirements. Computation-aware mechanisms are designed to prevent this. Methods like **MST**, **PrivMRF**, and **AIM** have this awareness built into their selection process, greedily adding marginals only if they do not violate a constraint on the complexity of the resulting graphical model (Cai et al., 2021a). **AIM** makes this safeguard explicit by maintaining a hard filter on its candidate set, discarding any marginal that would cause the estimated junction tree size to exceed a predefined memory limit (line 7 in algorithm 2). This guarantees a predictable runtime and memory footprint, a feature lacking in some other adaptive methods like **MWEM+PGM**, which can fail on certain workloads.

The design of the **AIM** (McKenna et al., 2022) algorithm represents a culmination of these ideas, demonstrating that the four “awareness” criteria are not merely an independent checklist of desirable features but are deeply interconnected aspects of a single, holistic optimization problem. The core of **AIM**’s selection process is its quality score function, where here, we let $M_r(D)$ denote the vector of counts for the marginal query r on the dataset D . The quality score is then given by,

$$q_r(D) = w_r \times (\|M_r(D) - M_r(p_{t-1})\|_1 - \sqrt{2/\pi} \times \sigma_t \times n_r), \quad (2)$$

which co-optimizes for these competing pressures in each greedy decision. The w_r term directly encodes workload-awareness by weighting candidates based on their relevance to the user’s task. The error term, $\|M_r(D) - M_r(p_{t-1})\|_1$, provides data-awareness by measuring the deficiency of the current data-dependent model. The noise penalty term, $\sqrt{2/\pi} \times \sigma_t \times n_r$, instills budget-awareness, as the noise scale σ_t is a direct function of the available privacy budget, naturally penalizing large, noisy marginals when the budget is tight. Finally, this entire selection process is constrained by an explicit check for computation-awareness, ensuring tractability at subsequent estimation steps.

3.3.3 Estimation Methods: From Structured Optimization to Generative Models

Once a set of queries has been chosen in the **Select** step and measured privately, the next challenge is to reconstruct a coherent data distribution from those noisy answers. This is the goal of the **Estimate** step, which can be formally cast as a constrained optimization problem. Given a vector of noisy measurements y , produced by applying a DP mechanism **A** to the true data’s histogram $h(D)$, the objective is to find an estimated probability distribution p that best explains the observations. This is a classic inverse problem, often framed as finding p that maximizes the likelihood of observing y , e.g. $\max_p \text{Likelihood}(p|y)$. In the common case where the mechanism **A** adds zero-mean Gaussian noise to a set of answers to linear queries (marginals) represented by a matrix Q , this simplifies to a least squares minimization problem over the probability simplex, or,

$$\min_p \|Q(p) - y\|_2^2. \quad (3)$$

While the objective function is convex, the primary challenge is the dimensionality of the search space for p , which is the size of the data universe (i.e., $|\mathcal{X}|$) and thus exponential in the number of attributes. Many works that crucially rely on a histogram representation of the dataset have been proposed to solve variants of this problem (Hay et al., 2010; Nikolov et al., 2013; Lee et al., 2015; Li et al., 2015; Abowd et al., 2019); again, we defer a discussion of these methods to Section 3.4.

An alternative and principled approach to this problem is **Private-PGM** (McKenna et al., 2019). Instead of optimizing over the full probability vector, Private-PGM leverages probabilistic graphical models (PGMs). A key insight is that when the measurements are a set of low-dimensional marginals, an optimal solution to the ℓ_2 minimization problem is guaranteed to be a distribution that can be represented by a PGM whose structure (i.e., its factors) is determined by the measured marginals. By optimizing over the parameters of this compact graphical model, Private-PGM can achieve exponential savings in computation. Critically, it solves a convex optimization problem over the *marginal polytope*, which ensures the resulting estimated distribution is *globally consistent*. Private PGM was validated by its use in winning entries of the 2018 and 2020 NIST DP competitions (McKenna et al., 2021a). The main trade-off is scalability: the computational cost is tied to the graph’s treewidth, making it intractable if the measured marginals induce a dense dependency graph.

To address these scalability limitations, a family of approximation methods has emerged. **Approx-Private-PGM** (APPGM) (McKenna et al., 2021b) relaxes the global consistency requirement, instead enforcing only a specified set of *local consistency* constraints over the simpler *local polytope*. Other methods, like the **Gradually Update Method** (GUM) (Zhang et al., 2021), are more heuristic, iteratively “massaging” a dataset to conform to the noisy marginals.

A distinct line of research bypasses the explicit structure of PGMs by using alternative, highly expressive representations. The **Relaxed Adaptive Projection** (RAP) framework (Aydore et al., 2021) uses a relaxed tabular representation amenable to gradient-based optimization. Similarly, **GEM** (Generative networks with the Exponential Mechanism) (Liu et al., 2021c) parameterizes the data distribution with a generative neural network. Both offer great flexibility but operate in a non-convex landscape, lacking the formal guarantees of Private-PGM. Other methods, such as **PrivBayes** (Zhang et al., 2017), use a Bayesian Network, though its heuristic estimation step can be improved by replacing it with a principled engine like Private-PGM (McKenna et al., 2022). Finally, some approaches leverage a **public dataset** to improve utility, e.g., by restricting the synthesizer’s domain to values seen in the public data (Liu et al., 2021b), though this carries the risk of degrading performance if the public and private distributions are dissimilar.

More recent innovations have introduced entirely new reconstruction primitives. **GREM** (Gaussian Residuals-to-Marginals) (Mullins et al., 2024) proposes reconstructing marginals not from noisy measurements of other marginals, but from noisy measurements of *residuals*. Residuals are a different basis of linear queries that are mutually orthogonal and possess a Kronecker product structure, properties that allow for extremely efficient and scalable reconstruction using pseudo-inverse operations. **GReM-LNN** further refines the output to enforce local non-negativity and consistency, which can significantly reduce error. It is important to note, however, that GREM is a mechanism for *marginal reconstruction* (i.e., query answering) and does not inherently produce a generative model. While it can produce highly accurate answers for a workload of marginals, it does not directly yield a synthetic dataset, limiting its use for analysts who need a dataset for downstream model training.

The evolution of these estimation methods highlights a fundamental design tension between rigor and scalability. Principled methods like **Private-PGM** perform an “inner approximation,” guaranteeing a valid probabilistic model at the cost of high computational complexity. Their search space is restricted to the set of valid, globally consistent probability distributions, and the resulting solutions are guaranteed to be a realizable probabilistic model from which one can directly sample. The cost of this rigor is computational complexity, as the search space (the marginal polytope) is difficult to navigate. To achieve scalability, methods like **APPGM**, **GUM**, and **GRem** perform an outer approximation. They relax the problem by optimizing over a larger space of “pseudo-marginals,” which are only required to be locally consistent. This makes the optimization problem much easier to solve but comes with a crucial caveat: the solution may not correspond to any single, valid data distribution. This introduces a new challenge, as a final (often heuristic) step is required to project the inconsistent model back into a valid dataset, a process that can introduce its own unquantified approximation errors. Therefore, the choice of an estimation method involves a deep architectural decision between a slower, principled approach that guarantees a valid model and a faster, more scalable approach that may fit the noisy data better but whose resulting model may not be mathematically sound. This tension has motivated the exploration of entirely different paradigms, such as end-to-end deep generative models, which offer an alternative path to capturing complex data distributions, as we will discuss further in Section 3.6.

3.4 Algorithms with strong worst-case guarantees*

Having detailed the paradigm that underpins many empirically-oriented algorithms, we now turn to their theoretical counterparts. Unlike methods that are optimized for performance on benchmark datasets, the algorithms in this section are designed to provide strong, provable error guarantees and give rise to precise error theorems which allow a theoretical comparison between workload-adaptive and workload-agnostic methods.

Specifically, many theoretically oriented algorithms, like their empirical counterparts discussed in Section 3.3, operate in a setting where the goal is to provide high accuracy for a given workload of queries \mathcal{F} , aiming to minimize the additive error as defined in Equation 1. Some of these foundational methods achieve their guarantees by operating on a full histogram representation of the data to produce a valid and globally consistent probability distribution. While this theoretical rigor is their primary strength, it often comes at the cost of scalability, making many of these methods impractical for large-scale problems without significant modification.

We begin by discussing a class of foundational algorithms that are data-independent. In this context, “data-independent” means that the choice of which statistics to measure is fixed in advance and does not depend on the input data D .

These approaches are often simpler to analyze and produce unbiased estimates of the workload-query answers (before post-processing, such as enforcing non-negativity constraints). One could view them as basic instances of the select-measure-estimate paradigm, but with a trivial “select” or “estimate” step. As the nuances are important and they are foundational, we believe it is appropriate to separate them out and describe them directly.

Noisy Histograms. A quintessential *workload-agnostic* approach, this method’s measurement strategy is fixed and exhaustive: it measures the count of every possible record in the universe \mathcal{X} . One directly adds discrete Laplace noise to each histogram bin: for each bin $x \in \mathcal{X}$, compute $\hat{h}_x = h(D)_x + \text{Lap}(1/\epsilon)$, and post-process it by making the noisy counts nonnegative. This approach is effective when the number of bins is much smaller than the number of queries, as its error rates scale polynomially with the former but only poly-logarithmically with the latter (e.g., Theorem 2.9 in (Vadhan, 2017c)).

Noisy Query-Answering. In direct contrast, this is a canonical *workload-adaptive* algorithm. Here, the algorithm’s measurement strategy is one and the same as the user’s evaluation workload. It adds noise (e.g., Laplace or Gaussian) directly to the answers of the queries in \mathcal{F} : $\tilde{h} = (\langle h(D), f \rangle + \text{Lap}(\Delta\mathcal{F}/\epsilon))_{f \in \mathcal{F}}$. The entire privacy budget is spent on preserving information known to be critical for the downstream task. To recover a histogram from the noisy query answers, one can solve an (integer) optimization problem over the histogram variable \hat{h} that aims to minimize the difference between its query answers and the noisy ones, as proposed by (Barak et al., 2007):

$$\min_{\hat{h} \in \mathbb{R}_+^{\mathcal{X}}: \|\hat{h}\|_1 = n} \left\| \langle \hat{h}, f \rangle - \tilde{h} \right\|_{f \in \mathcal{F}}. \quad (4)$$

While this approach is suboptimal in general, it has proved effective in settings such as synthetic graphs that are accurate for cut queries (Gupta et al., 2012). More involved randomization strategies (Steinke and Ullman, 2016; Ganesh and Zhao, 2021; Dagan and Kur, 2022; Ghazi et al., 2021), including the generalized Gaussian mechanism and the sparse vector technique, can yield optimal rates. Error rates can be made independent of the universe size (number of bins) and the number of rows, making this approach advisable for a relatively small number of queries, such as low-dimensional marginals as we discussed in Section 3.3.1.

The Matrix Mechanism. The Matrix mechanism (Li et al., 2015) interpolates between these two baselines and includes them as special cases. Its key idea is to compile the workload into a different set of queries known as the *strategy*, add noise to the strategy query answers, and then post-process them to estimate the workload answers and the underlying histogram. Setting the strategy to be the identity matrix recovers the noisy histogram approach, while setting it to the workload recovers the noisy query-answering approach. The mechanism also proposes a *strategy optimization* problem to Pareto dominate both baselines. However, it requires representing the workload and strategy in matrix form, which can be prohibitive for large domain sizes. The high-dimensional matrix mechanism (HDMM) (McKenna et al., 2018) addresses this by utilizing implicit matrix representations, enabling its use in higher-dimensional settings when the workload consists of conjunctive queries.

Data-dependent Approaches with Guarantees If we relax the requirement of unbiasedness, we can trade a small amount of bias for a large reduction in variance. This is often a worthwhile trade-off in practice, especially for smaller datasets and privacy budgets.

Reference	ℓ_∞ Upper Error bound	DP
Steinke and Ullman (2016)	$O\left(\frac{ \mathcal{F} }{\epsilon}\right)$	ϵ -DP
Dagan and Kur (2022), Ghazi et al. (2021)	$O\left(\frac{\sqrt{ \mathcal{F} \log \frac{1}{\delta}}}{\epsilon}\right)$	(ϵ, δ) -DP
Vadhan (2017c)	$O\frac{\sqrt{ \mathcal{X} \log \mathcal{F} }}{\epsilon}$	ϵ -DP
Blum et al. (2013); Hardt and Rothblum (2010); Hardt et al. (2012b)	$O n^{\frac{2}{3}} \frac{\log \mathcal{F} \log \mathcal{X} ^{\frac{1}{3}}}{\epsilon}$	ϵ -DP
Hardt and Rothblum (2010); Hardt et al. (2012b)	$O\left(n^{\frac{1}{2}} \frac{\sqrt{\log \mathcal{X} \log \mathcal{F} \log \frac{1}{\delta}}}{\epsilon}\right)$	(ϵ, δ) -DP

Table 2: Best-known upper bounds on the error for tabular synthetic data. Algorithms from the first two rows work for answering counting queries, but they can be turned into synthetic data generators by solving an integer linear program over $n \times \log |\mathcal{X}|$ boolean variables. The rest of the algorithms directly provide synthetic data. For approximate-DP, the optimal error rates are determined by the minimum of the corresponding rates in the table.

One such data-dependent approach is known as **Private Multiplicative Weights (PMW)** (Hardt and Rothblum, 2010). **PMW** works with a histogram variable initialized with uniform weights and subsequently updates the histogram to progressively make the counts more accurate (see Algorithm 3). The updates are based on the most inaccurate queries for the current histogram, which must be identified privately using methods like Report Noisy Max or the Exponential Mechanism (McSherry and Talwar, 2007). **PMW** attains error rates that are poly-logarithmic with respect to $|\mathcal{X}|, |\mathcal{F}|$, but polynomial in n , making it advisable for a relatively small number of rows. These running times are generally unavoidable (Dwork et al., 2009; Ullman and Vadhan, 2011). Considering that the data domain or the number of queries can be very large, this limits its practicality.

Following work on **PMW**, Hardt et al. (2012b) consider the case where the domain can be decomposed in product form, $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_t$, such that each query depends only on a single \mathcal{X}_i , and propose a Multiplicative Weights-Exponential Mechanism (**MWEM**) approach. **MWEM** allows for parallelization and a much smaller product-distribution representation of the histogram, though it rules out queries across attributes from different groups.

Another alternative is to solve the dual problem to (5) (Gaboardi et al., 2014). This approach, based on the Minimax Theorem (Neumann, 1928), performs PMW updates on a distribution over *queries*. Privacy is achieved by sampling from the resulting query distribution (an application of the exponential mechanism). With privacy guaranteed by sampling, the inner minimization subproblems can be solved with effective but *nonprivate* heuristics for NP-hard problems (e.g., integer linear programs). This is advantageous when the number of queries is much smaller than the number of bins.

$$\min_{\hat{h} \in \mathbb{R}_+^{\mathcal{X}}: \|\hat{h}\|_1 = n} \left\| \left(\langle \hat{h}, f \rangle - \langle h(D), f \rangle \right)_{f \in \mathcal{F}} \right\|. \quad (5)$$

Preferred Approach	Regime
Noisy histograms	$ \mathcal{X} \ll \max\{n, \mathcal{F} \}$
Noisy-query answering	$ \mathcal{F} \ll \max\{n, \mathcal{X} \}$

Private multiplicative weights

 $n \ll \max\{|\mathcal{F}|, |\mathcal{X}|\}$

Table 3: Advised choice of approach, based on the error rates from Table 2.

We emphasize that as far as the theoretical guarantees are concerned – there is no one approach (noisy histograms, noisy-query answering or PMW) that outperforms uniformly across different problem settings (Zhang et al., 2017; Rosenblatt et al., 2020). This conclusion is consistent with the fact that methods’ performance differ on different regimes of data size, data domain, and number of queries⁸. Table 3 provides rough recommendations on which approach is preferred in which regime.

Algorithm 3 Private Multiplicative Weights Update Method

Require: Private dataset: $D \in \mathcal{X}^n$, query workload: finite set $\mathcal{F} \subseteq \{0, 1\}^{\mathcal{X}}$, stepsize $\eta > 0$.

Ensure: Privatized histogram: $\hat{h} \in \mathbb{R}_+^{\mathcal{X}}$, $\|\hat{h}\|_1 = n$, representing DP synthetic tabular data

```

1:  $h^1 \leftarrow \frac{n}{|\mathcal{X}|} \mathbf{1}$ 
2: for  $t = 1$  to  $T$  do
3:   Select worst performing query  $f_t$  that approximately solves  $\max_{f \in \mathcal{F}} |\langle h^t - h(D), f \rangle|$  (in DP manner), together with  $v_t$ 
   private estimate of query error  $\langle h^t - h(D), f_t \rangle$ 
4:   if  $|v_t| \leq \alpha$  then return  $\hat{h} = h^t$ 
5:   else
6:      $h^{t+1} = n \frac{(h_x^t \cdot \exp(-\eta \cdot \text{sgn}(v_t) \cdot f_t(x)))_{x \in \mathcal{X}}}{\sum_{y \in \mathcal{X}} h_y^t \cdot \exp(-\eta \cdot \text{sgn}(v_t) \cdot f_t(y))}$ 
7: return  $\hat{h} = \frac{1}{T} \sum_{t=1}^T h^t$ 

```

3.4.1 Extensions and Variants*

Relative error guarantees Instead of focusing on achieving worst-query utility in additive terms (e.g. under the ℓ_∞ -norm in (1)), recent work has focused on assessing utility in terms of *relative error*. Here, the goal is to construct synthetic histograms with a mixed multiplicative/additive approximation guarantee: for error parameters $\zeta, \alpha > 0$, the goal is computing a synthetic histogram \hat{h} such that for all $f \in \mathcal{F}$

$$(1 - \zeta)\langle h(D), f \rangle - \alpha \leq \langle \hat{h}, f \rangle \leq (1 + \zeta)\langle h(D), f \rangle + \alpha.$$

These guarantees are particularly relevant for analytics settings, where the goal is to distinguish ‘large’ from ‘small’ counts, rather than accurately predicting each count (Cormode et al., 2012; Qardaji et al., 2013,?; Epasto et al., 2023; Ghazi et al., 2023b). Recently, it has been shown that in this setting there exist DP algorithms with constant multiplicative approximation (i.e., ζ constant) whose additive error (i.e. α) scales only poly-logarithmically with $n, |\mathcal{X}|, |\mathcal{F}|$ (Ghazi et al., 2025).

Algorithms targeting distributional guarantees. Algorithms in this group can be considered as having an infinite workload: this analogy is precise for the broad class of *integral probability metrics* (Müller, 1997). As mentioned earlier, a notable and widely studied particular case is synthetic data which are accurate in terms of the Wasserstein/Kantorovich distance, characterized by approximation of expectations over arbitrary 1-Lipschitz functions (Villani, 2008). This example is relevant, as Lipschitz queries commonly arise as the result of stable data analyses, e.g. models trained on regularized convex losses. In this context, early hierarchical histogram algorithms already approached a $\tilde{O}(1/n)$ Wasserstein-1 error bound with hierarchical histogram methods in one dimension, although they did not state the analysis as such (Hay et al., 2010). More recently, a variety of results have matched or improved these early results such that they have achieved nearly optimal Wasserstein error rates with algorithms that run in polynomial time (Boediardjo et al., 2024; He et al., 2023; Musco et al., 2024; Feldman et al., 2024); these algorithms also leverage effective ways to hierarchically discretize the continuous space, together with noise addition approaches to privately count the number of data points falling into each bucket for the discretizations. Some take the approach of constructing a random walk that diverges logarithmically with the number of steps (as opposed to with the square root of the step count) (Boediardjo et al., 2024), while others look at matching a class of distributional moments closely tied to Wasserstein distance (Musco et al., 2024), and still others look at alternative hierarchical histogram constructions (Feldman et al., 2024).

Unfortunately, the error rates in this setting suffer from the curse of dimensionality, and these rates are essentially unimprovable in general (there is a known $\Theta(1/n^{1/d})$ lower bound (Boediardjo et al., 2024)). Improvements have

⁸It may seem odd to refer here to ‘number of queries’. For clarifications on what queries may be of interest here, see Section 3.2.1.

been obtained under additional assumptions, such as (intrinsic) low dimensionality of the data (Ghazi et al., 2023a; Donhauser et al., 2023; He et al., 2025), as well as extensions to online settings (He et al., 2024). Characterizations of the sample complexity of DP synthetic data generation that are accurate for infinitely-many counting queries have also been obtained in the literature (Bousquet et al., 2020).

3.5 Numerical feature handling

Almost all of the methods we’ve described so far assume the data comes from a categorical domain with a fixed, known set of possible values. Many real-world tabular datasets have a mix of both categorical and numerical features. One simple way to bridge this gap is to discretize numerical attributes to intervals. This numerical-to-categorical transformation yields a fully categorical dataset and enables using the mechanisms described earlier. There are multiple strategies one might use to discretize numerical attributes, the simplest being a uniform discretization which maps values to fixed-width intervals. The benefit of this strategy is that it is data-independent, and only requires knowledge of the minimum and maximum possible value for a numerical attribute. This number of bins is an important hyper-parameter to consider. Larger values yield finer-grained discretizations that may better reflect the original distribution, but results in bins with smaller counts which are not preserved as well in the presence of noise, and larger categorical domains, which can increase runtime of downstream mechanisms that operate on the categorical data. When the distribution is heavily concentrated on a single value, uniform discretizations sometimes have undesirable properties. For example, the capital gain attribute of the adult dataset spans from 0 to 100,000, but 91% of the values are 0. Under a uniform discretization with 100 bins, these 0-values will be mapped to the interval $[0, 1000]$, and a lot of information is lost. Converting this interval back to a numerical value can be done in a number of ways, choosing the midpoint deterministically (500) or sampling uniformly at random from the interval are two natural options. In both cases, the numerical data is heavily skewed away from the distribution of the real data.

Adaptive discretization techniques (Zhang et al., 2016; Gillenwater et al., 2021) can overcome this problem, by identifying a partition of the domain into roughly equal-density intervals. Under this approach, smaller intervals will be chosen where data is more concentrated, while larger intervals will be chosen where it is more spread out. In the example above, the interval $[0, 0]$ would typically be chosen because it contains so much probability mass, and hence this class of approaches do not suffer from the same distribution-skew as uniform discretization.

Mapping numerical attributes to categorical ones is just one way they can be handled. However, it’s important to note that this discretization step can introduce inaccuracies and affect the utility of the resulting synthetic data, as highlighted by (Ganev et al., 2025b). The number of bins and the choice of discretization strategy can significantly impact the quality of the generated data. As shown by (Ganev et al., 2025b), optimizing both the discretizer (uniform, quantile, k-means, or PrivTree) and the number of bins can lead to substantial improvements in utility. Technically speaking, one must apply DP to the discretization process itself, to avoid introducing new privacy vulnerabilities.

Alternatively, some mechanisms, like the **RAP++** algorithm (Vietri et al., 2022) are designed to natively handle numerical attributes *without* requiring discretization. This is a key advantage, as it allows for preserving more information about the original data distribution. **RAP++** achieves this by directly optimizing for “class conditional linear threshold queries” and “mixed marginal queries,” which can capture the relationships between numerical features and class labels without relying on a binning strategy. These types of queries are defined directly on the numerical values, instead of on categories. This allows for a more accurate representation of the underlying data, and can lead to faster runtimes compared to methods that require discretization.

3.6 Modern end-to-end methods for tabular data generation

In contrast to the *workload-adaptive* algorithms that take explicit guidance on which statistics to preserve (Section 3.3), another significant class of practical algorithms is fundamentally *workload-agnostic*. These methods use modern, end-to-end deep generative models. Their “measurement strategy” is implicit in the model’s architecture and the training objective (e.g., minimizing an adversarial loss or a reconstruction loss), which aims to learn a holistic, flexible representation of the joint data distribution. This strategy is designed for general-purpose fidelity and is not structured to adapt to a user-provided workload of specific statistical queries (unless the method allows the user to specify measurement queries of interest, which is sometimes the case.) The hope is that a model that accurately captures the overall distribution will prove useful for a wide range of future, unanticipated analyses.

These approaches often rely on DP-finetuning or DP-training (from scratch) techniques, just like the winning methods in image (Section 4.3) and text (Section 5.2) data synthesis. With privacy unit being a row in table, these models conceptually learn to generate one row of synthetic data, rather than 1 full image or a piece of text.

We survey evolution of these approaches next.

Early attempts: In non-DP settings, both GANs (Generative Adversarial Networks) (Goodfellow et al., 2014) and Variational Autoencoders (VAEs) (Kingma and Welling, 2022) have been successfully applied to generative tasks, especially in image domains. GANs are adversarial networks that approximate the joint data distribution through an adversarial training process that pits a generator, which attempts to produce a synthetic sample good enough to fool a discriminator, against a discriminator that learns to distinguish real vs. synthetic samples. VAEs, on the other hand, attempt to model the joint distribution of data by mapping a sample to a distribution in the latent space (encoding step), then sampling from this distribution (decoding step).

DP-GANs either employ DP-SGD like training: DP-GAN (Xie et al., 2018), RDP-CGAN (Torfi et al., 2022), CTAB-GAN+ (Zhao et al., 2023), or attempt private inference via ensembling (PATE-GAN (Yoon et al., 2019)). However, an empirical study that looked at the performance of downstream ML models trained on DP synthetic data (Machine Learning Efficiency) by Tao et al. (2021) demonstrated that previously discussed workload-based methods outperform GAN-based models, with the latter not being able to accurately preserve even one-dimensional statistics of the tabular data. GAN-based methods may be hard to train with DP-SGD due to GANs’ adversarial nature and their inherent instability (even in a non-DP setting) (Tran and Xiong, 2024). VAE-based approaches include DPGM (Acs et al., 2018), DP-SYN (Abay et al., 2018), and P3GM (Takagi et al., 2020); DP combinations of both GANs and VAEs also exist (Yang et al., 2024).

Modern models: *Diffusion models* have shown great promise in image domains recently (Nakkiran et al., 2024), and these approaches have also been applied in the tabular domain as well. Diffusion models learn to progressively distort the input by introducing noise over a number of steps and then learning to reverse this process for generation. For example, the TableDiffusion model (Truda, 2023) outperformed DP-GAN-based models in a limited study involving two datasets (Truda, 2023). The following two groups of generative models (Autoregressive transformer-based and Pretrained LLM-based) employ LLM-like architectures and model the tabular data via textual encoding, then essentially employing the chain rule on the tabular data on a token-by-token basis, modeling next token prediction based on previous token values.

Autoregressive models: Castellon et al. (2023) introduced a probabilistic model DP-TBART, an autoregressive transformer-based model. The difference between this model and a standard LLM is the exact architecture, which is a 3-layer decoder closely resembling DistilGPT-2, and a custom tokenizer that assigns different tokens to each column’s distinct values (after discretization) so no two column encodings share the same tokens. Each tabular row is then encoded as values of columns (assuming some ordering of features), with feature names not being used in this encoding. DP-SGD is used to train the model. During sampling, to force the output to be valid, post-processing is employed, removing all disallowed tokens for each column. The authors provide a comparison with a number of methods including AIM, with AIM’s workload-based method outperforming the proposed methods (and all other considered baselines). As was noted in Section 3.3.1, higher-order marginals suffer from exponential scaling, and in practice, this leads to the use of lower-order marginals, which can fail to capture complex interactions of the features. To demonstrate this phenomenon, Castellon et al. (2023) propose two new datasets with complex joint feature interactions – Dyck-k, an artificially constructed dataset representing a worst-case dataset for workload-based methods, and WikiText-103 (a real-world dataset of Wikipedia articles). DP-TBART outperforms AIM by a significant margin on both these challenge sets.

Recently, Tiwald et al. (2025) introduced a Tabular AutoRegressive Generative Network (TabularARGN), departing from the trend of using complex models like Transformers to model joint distribution. They argued that treating tabular data as text and using LLMs is counter-intuitive. Their architecture is essentially a multi-tower model with each tower being dedicated to each feature. At the top of the model, a dedicated embedding encodes each feature; the embeddings are then combined (with some permutation of order) and forwarded to the dedicated towers. Each tower sees only the feature values on which it conditions: for the first feature, no additional embeddings are passed; the second feature is conditioned on the first and receives its embedding; the third is conditioned on the first two, etc. These towers are standard feedforward layers with ReLU and dropout. Each tower has a softmax head, and the cross-entropy loss of all heads is summed. During sampling, the column values are inferred one by one and passed as the input to the next column. To achieve DP, DP-SGD can be used. While the simplicity and performance of the model in a non-DP setting are appealing, no competing DP models were used for comparison in the DP setting, making it hard to judge the effectiveness of this method for DP synthetic tabular data generation.

Pretrained LLM based: With the rise of LLMs, research focus shifted towards attempting to harness their capabilities for tabular synthetic data generation. Using LLMs can allow researchers to forego extensive preprocessing needed for query release methods (categorization of continuous numeric attributes, data imputation and normalization, outlier removal, and data smoothing). Additionally, purely numerical data representations like histograms erase subtle connections that even a rudimentary language model can infer. For example, Borisov et al. (2023) highlighted that for the Adult dataset (Becker and Kohavi, 1996), features like age, marital status, and education exhibit clear hidden

relationships that a pretrained LLM can uncover and utilize easily: e.g., higher degrees cannot be obtained at an early age, and in many countries, it is illegal for teens to be married. While this knowledge can be discovered given sufficient data by histogram-based methods, LLMs can readily take feature names into account.

Early attempts at harnessing the emerging powers of using standard LLMs included pioneering work by [Borisov et al. \(2023\)](#) in the *non-DP setting*. Their method GReaT works by first creating a purely textual encoding of all tabular data attributes, including the label attribute (e.g., “Bachelors Education, Adult male, income <50K”), while permuting the order of the attributes in the encoding. Then, a pretrained LLM is finetuned on such textual encodings. Permutation of feature names is introduced to avoid the model’s reliance on a pseudo-ordering of attributes that the model would learn otherwise. Additionally, it unblocks conditional sampling by providing either any feature name or any combination of feature names and initial values. Sampled text data is then converted back to tabular, with samples that do not conform to the expected format (e.g., out of range, wrong names of attributes) discarded.

While this approach can be easily converted to a DP setup by swapping an optimizer with its DP version during LLM finetuning, [Tran and Xiong \(2024\)](#) argue that it will not work out of the box and the DP-finetuned model will not conform to the expected format due to the noise introduced during DP training. [Tran and Xiong \(2024\)](#) instead separate the DP finetuning into two stages: learning format compliance and actual private data modeling. For the first stage, [Tran and Xiong \(2024\)](#) treat feature names and categorical values, as well as the range of values for numerical attributes, as non-sensitive data and create random data represented in a similar format. This data does not exhibit realistic feature distributions and dependencies and serves to learn the proper output format of the data. An LLM is then finetuned non-privately on this data with LoRA. The second stage is DP LoRA finetuning on the actual private data textual encodings. These encodings are similar to ([Borisov et al., 2023](#)) (<attribute> is <value>). For this stage, the authors additionally replace the standard next token Cross Entropy loss with a combination of weighted cross-entropy (WCEL) and numerical understanding loss (NUL). WCEL up-weights the actual private tokens and down-weights formatting tokens (feature names, connection ‘is’, commas, etc.). NUL, on the other hand, attempts to solve the problem that CE-based learning does not distinguish the magnitude of error when it comes to numerical data: an error in only one token can have a different magnitude of error when the prediction is converted to a numeric format. For example, assuming each digit is a token, for a 10.0 ground truth, predictions 19.0 and 10.1 both have one token error, while the magnitude of the numerical error is widely different. NUL loss thus penalizes the squared errors between predicted numerical tokens and actual ground truth numerical values. The authors compare their method DP-LLMTGen with GAN-based methods, PATE-based methods, and some workload-based methods (RAP ([Aydore et al., 2021](#)), RAP++ ([Vietri et al., 2022](#)), and GSD ([Liu et al., 2023b](#))), with DP-LLMTGen exhibiting impressive results, outperforming competitive methods in terms of statistical fidelity in most cases, often by a large margin, with RAP++ being the only method that comes close. However, when considering the ML efficiency metric (performance of a gradient boosted tree model trained on synthetic data), DP-LLMTGen outperforms workload-based methods in only 3 out of 10 cases. While the comparison is in no way exhaustive and some modern methods like AIM ([McKenna et al., 2022](#)) are missing, this study hints at the potential of LLMs for DP tabular data synthesis.

Concurrently, [Afonja et al. \(2025b\)](#) reached similar conclusions and proposed a two-stage method similar to [Tran and Xiong \(2024\)](#). For the first stage, they considered using either public data encoding (that does not reference correct feature names, ranges, or categorical values but conforms to the same encoding scheme) or random data generated using the aforementioned knowledge. Their findings indicate that a public dataset works sufficiently well. They also argued that column shuffling makes learning in DP harder and advocate against such shuffling (doing this, however, will fix the order of feature names during sampling). For DP finetuning, they use only a weighted loss that similarly down-weights formatting tokens. For sampling, they suggest a modification of the sampling mechanism where, instead of fully discarding the sample that does not conform to the expected formatting, only the offending tokens are discarded, and sampling is re-attempted by prompting with valid tokens (imputation). Their ablations on weighted loss suggest that it is beneficial for situations when a public dataset was used in the first Stage, and less needed for when a uniformly generated dataset was used. They compare their method to modern methods like MST and AIM, and while their method demonstrates good performance across most metrics, underperforming on certain fidelity metrics, especially for large dimensional datasets, AIM ([McKenna et al., 2022](#)) still significantly outperforms it.

3.7 Comparison of the methods

When comparing these diverse approaches, it is helpful to return to our initial categorization. On one hand, practical workload-based methods have solid foundations, impressive empirical performance and, in some cases, come with strong performance guarantees. On the other hand, generative models, particularly those based on LLMs, lack this theoretical backing but show great promise and may implicitly leverage knowledge that is difficult to capture with explicit marginal measurements.

Table 4 Taxonomy over algorithms for DP tabular data synthesis. A checkmark (✓) indicates the algorithm’s measurement strategy is designed to adapt to a user’s downstream workload.

Method	Year	Paradigm	Workload Adaptive	Data Aware	Budget Aware	Computation Aware
<i>Marginal-Based Mechanisms</i>						
PMW / MWEM (Hardt and Rothblum, 2010), (Hardt et al., 2012b)	2010/12	Marginal	✓	✓		
Noisy Query-Answering (Barak et al., 2007)	2007	Marginal	✓			
Matrix Mechanism (Li et al., 2015) / HDMM (McKenna et al., 2018)	2015/18	Marginal	✓			
RAP (Aydore et al., 2021)	2021	Marginal	✓	✓		✓
GEM (Liu et al., 2021c)	2021	Marginal	✓	✓		✓
AIM (McKenna et al., 2022)	2022	Marginal	✓	✓	✓	✓
JAM-PGM (Fuentes et al., 2024)	2024	Marginal	✓	✓	✓	✓
PrivBayes (Zhang et al., 2017)	2017	Marginal		✓	✓	✓
MST (McKenna et al., 2021a)	2021	Marginal		✓		✓
PrivSyn (Zhang et al., 2021) (w/ GUM)	2021	Marginal		✓	✓	✓
PrivMRF (Cai et al., 2021a)	2021	Marginal		✓	✓	✓
<i>End-to-End Generative Mechanisms (Workload-Agnostic by design)</i>						
GAN-based (Xie et al., 2018), (Yoon et al., 2019)	2018/2019	Generative		✓		✓
VAE-based (Acs et al., 2018), (Abay et al., 2018)	2018	Generative		✓		✓
Diffusion-based (Truda, 2023)	2023	Generative		✓		✓
LLM-based (Tran and Xiong, 2024), (Afonja et al., 2025b)	2024/25	Generative		✓		✓

Non-DP tabular data generation is a rapidly developing area, and many models that are introduced and successful in the non-DP setting were either not yet attempted and adapted for DP (e.g., TabPFNGen (Ma et al., 2024)) or may require modifications to make them powerful under the noise introduced by DP. This was already demonstrated for GReaT-inspired models Borisov et al. (2023), where DP variants had to be trained in a two-stage setup. Additionally, many evaluations are not exhaustive, often offering comparisons only against other generative models of the same type and excluding well-established empirically oriented workload-based methods like AIM (McKenna et al., 2022). For example, a recent benchmark study by Chen et al. (2025) that showed AIM’s superiority specifically excluded LLM-based methods from comparison, arguing that such methods use knowledge embedded in pretrained LLMs that prevents fair comparison.

There is, however, evidence that LLM-based modeling can have an edge over workload based methods. Firstly, one interesting observation by Borisov et al. (2023) was that for tabular data generation, training generative models using meaningful feature names like Age, Education, etc., results in better synthetic data compared to the same data that had random feature names assigned. This demonstrates that LLMs utilize pretraining knowledge, which is different from GANs, VAEs, and workload-based methods. Further, pretrained (but not SFT or RLHF LLM checkpoints) work marginally but not significantly better, suggesting that pretraining is important. Castellon et al. (2023) pointed out the exponential scaling of workload-based methods to high-order marginals and demonstrated two use cases where an LLM-based method outperforms them significantly.

While finetuning-based approaches clearly exhibit potential, early attempts at using an LLM without finetuning have encountered difficulties in obtaining good quality tabular synthetic data, e.g., a private evolution-based approach (Swanberg et al., 2025) and in-context learning based (Hu et al., 2024). This highlights that tabular data is often out-of-distribution for pretrained LLMs, and thus finetuning is necessary.

To summarize, practical workload-based methods (Section 3.3) have demonstrated impressive performance, solid foundations and can come with performance guarantees. While LLM-based methods currently lack wide adoption and rich body of research, and while practical workload-based methods like AIM (McKenna et al., 2022) outperform generative-based methods currently on simpler and/or more realistic datasets, generative models, in particular LLM-based ones, show great promise.

3.8 Evaluating synthetic tabular data quality

On top of the generic metrics applicable to all data modalities mentioned in Section 1.5, below we outlined metrics specific to tabular synthetic data.

Fidelity: For workload adaptive algorithms which create synthetic data to answer a set of predefined queries accurately, the fidelity of achieved synthetic data can be evaluated and iterated upon using the Equation 1 (namely, the performance of synthetic data on workload queries). This metric also gives an idea of performance on unseen queries that were not in the query set for the algorithm, albeit no guarantees are provided.

Alternatively, for all types of tabular data, the following metrics can be employed. Aydore et al. (2021) introduced *Statistical Fidelity* which is an average of total variation distances of joint distributions (1-5 way marginals) between

synthetic and set aside real data. Similarly, *Pairwise Attribute Distribution Similarity* measures the similarity of all two way marginals by averaging histogram intersections with numerical attributes discretized into bins Afonja et al. (2025b) also suggests *pairwise correlation similarty*, which estimates how well the synthetic data preserves pairwise column correlations.

Castellon et al. (2023) suggested *Kolmogrov-Smirnov test* for numerical attributes and Chi-square test for cateogrical columns. Distributional distance metrics like *Maximum mean discrepancy (MMD)* and $\alpha - Precision$ can also be used for this purpose and beta-recall (Alaa et al., 2022)

Many additional metrics for comparing tabular datasets surfaced in non-dp setting, for example the ℓ_1 -distances of synthetic examples to the closest real datapoint (Borisov et al., 2023). However they are hard to properly dp-fy.

Utility: For many practical applications the goal is to use synthetic data for some downstream task, for example training/finetuning an ML model. None of the workload based algorithms we discussed previously directly optimize for this measure, however all the algorithms for query release can be used to create tabular synthetic data in this case, and Equation 5 can still be used as a proxy metric that is assumed to be correlated with downstream ML model performance. *This is what commonly happens in practice.*⁹

Systematic: Some papers propose utility evaluations that are systematic, in that they consider fully specified downstream use cases of differentially private synthetic data as the metric itself. For example, Hod and Canetti (2024) perform a national scale data release of sensitive medical data using DP synthesizers; as part of their real world criteria, they considered *faithfulness*. This is a finer granularity metric, requiring that some proportion of the released dataset records “look like” a proportion of the original records (as opposed to just matching moments or aggregate statistics). In particular, they operationalize their metric as (α, β) -faithfulness. Given a cost function between two records $c : X \times X \rightarrow \mathbb{R}_{\geq 0}$, a dataset $S \in X^n$ is (α, β) -faithful with respect to a dataset $R \in X^n$ if there exists a bijection $\pi : S \rightarrow R$ such that

$$\frac{1}{n} \sum_{i=1}^n \mathbb{I}[c(s_i, r_{\pi(i)}) \leq \alpha] \geq \beta, \quad (6)$$

where \mathbb{I} is the indicator function. In words, at least a β proportion of records in S are matched to records in R with a cost no greater than α . Additionally, they define maximal- β -faithfulness of a dataset $S \in X^n$ with respect to a dataset $R \in X^n$ and a cost function c as,

$$\beta_{\max}(R, S) = \max_{\pi \text{ matching}} \frac{1}{n} \sum_{i=1}^n \mathbb{I}[c(s_i, r_{\pi(i)}) \leq 1], \quad (7)$$

which is the highest possible proportion of records in S that can be matched to records in R with a cost no greater than 1.

Rosenblatt et al. (2024a) take a broader approach than a single real data release. They consider an *epistemic parity* metric, which evaluates the ability of a differentially private (DP) synthetic dataset to reproduce scientific findings by assessing whether a functionalized version of a statistical test yields the same conclusion when applied to both the original and the DP synthetic data (within some context-specific tolerance α accounting for sampling error and noise introduced by DP). The epistemic parity, or level of agreement, is the proportion of statistical tests that result in the same conclusion when using the real and DP synthetic datasets, and is a measure of how well a DP synthetic data approach works for reproducing the types of statistical analyses present in the paper in question.

Finally, we want to highlight that standard metrics outlined Section 1.5 remain the mostly commonly used metrics for utility evaluation of DP synthetic tabular data.

3.9 Open questions and challenges

Despite the significant amount of work in this area, various basic and central questions in private synthetic tabular data generation remain open. From a theoretical perspective, the main open problem is the optimal accuracy rates for pure DP algorithms in terms of worst-case error as in (1) where $p = \infty$. Currently there is a polynomial in n gap between best upper and lower bounds (Hardt, 2011; Nikolov and Ullman, 2021).

⁹Alternatively, the case of creating synthetic data to maximize downstream ML model utility can be understood as a situation with queries that preserve the excess risk of the solution of empirical risk minimization problems with respect to a class of loss functions. It turns out that the private multiplicative weights update method can be applied in greater generality to address convex minimization queries (Ullman, 2015). This algorithm however suffers the same limitations mentioned for private multiplicative weights for linear queries and is not considered practical.

Another underexplored direction relates to the computational complexity of producing such data for specific query families. While the hardness result of [Ullman and Vadhan \(2011\)](#) indicate that queries corresponding to hard constraint satisfaction problems (CSP) are necessarily hard for synthetic data generation (in fact, this argument implies that producing synthetic data that approximates 2-way marginal queries is hard), it is not known whether easy CSPs are such that producing synthetic data is easy. On the other hand, hardness results on synthetic data are based on cryptographic assumptions (more precisely, the existence of one-way functions), rather than more traditional (computational) hardness of approximation.

Many mechanisms assume that the data lives on a closed domain, where each attribute is categorical and can assume a fixed number of (reasonably small) possible values. While numerical attributes can be nominally supported via discretization, it remains relatively unexplored to handle other data types, like open set categorical attributes, set-valued attributes, and free-form text attributes.

From a practical perspective, closing the gap in performance between generative models and workload based methods is an important question. Furthermore, assessments of the data quality generated by generative models are still done on an ad-hoc basis, and more systematic methods are needed.

Additionally, we want to emphasize that the problem of generating a collection of tabular synthetic datasets (a setting where privacy unit is a full tabular dataset) is very much unexplored.

4 DP Synthetic Image data

Image data is ubiquitous in the modern age, resulting in the very large volume of available training data. However creating privacy-preserving (DP) realistic looking synthetic images remains a challenge even after years of research and several families of models (GANs, Diffusion). Generating high-quality DP synthetic images data faces several significant hurdles:

1. **High Dimensionality of Images:** The high dimensionality of image data poses a major obstacle ([Chen et al., 2022b](#)). Capturing the intricate distributions of natural images requires complex models. Applying DP noise in such high-dimensional spaces (either pixel space or gradient space during training) can easily overwhelm the signal, leading to poor quality synthetic images ([Chen et al., 2022b](#)). Techniques to counter this often involve working in lower-dimensional latent spaces ([Liu et al., 2023a](#)) or using specialized DP mechanisms ([Chen et al., 2022b](#)).
2. **Utility Evaluation Difficulties:** Assessing the quality of DP synthetic data is multifaceted and challenging. Standard metrics like FID for image quality might not correlate well with performance on specific downstream tasks ([Ramesh et al., 2024](#)). Furthermore, the noise introduced by DP can distort statistical properties, potentially leading to misleading results or false discoveries (e.g., inflated Type I errors) when analyzing the synthetic data ([Perez et al., 2024](#)).
3. **Computational Cost:** Training large-scale generative models like GANs and diffusion models is computationally intensive. DP training, particularly using DP-SGD, often adds significant overhead due to the need for larger batch sizes and per-example gradient clipping, potentially increasing training time substantially ([Liu et al., 2023a](#)). Additionally generative models like GANs and Diffusion are hard to dp-fy properly due to their training complexity, potential use of non-dp-able components (Section 7.3) and training process instability.

Nevertheless, the progress in synthetic data synthesis has been short of amazing. Not surprisingly, DP synthetic data synthetic field also demonstrated impressive results recently, which we will explore in this Section.

Roadmap First we discuss what a meaningful privacy unit for image data is (Section 4.1). Then we proceed to discuss general approaches for generating DP synthetic data (Section 4.2). We then focus on 2 prominent methodologies : DP finetuning (Section 4.3) of GANs (Section 4.3.1) and more recently Diffusion Models (Section 4.3.2) and training-free technique called Private Evolution (PE, Section 4.4). We also cover methods that don't fit neatly into the above categories in Section 4.5, including PATE-style GANs (Section 4.5.1) and DP-fying images via an intermediary re(Section 4.5.2). We then offer a comparative analysis of these approaches (Section 4.6), go over potential metrics that can be used to evaluate the quality of synthetic data (Section 4.7) and conclude with a discussion of ongoing challenges and potential future research directions (Section 4.8).

4.1 The privacy unit for image data

Determining the right privacy unit for image data is an important task. Each image often has a clear and complete semantic meaning making *sub-example-level privacy unit* not appropriate or extremely hard. For example, protecting

only people’s faces in images might not be sufficient — a photo taken in someone’s kitchen may still reveal private information about their home even if faces are blurred.

Most of the literature treats each individual image as a privacy unit (*example-level privacy unit*). This might be appropriate if each user contributed at most one image to the overall private dataset.

In cases when each user can be associated with more than 1 image, *user-level* privacy unit is a better choice. The notion of a “user” however must be chosen carefully. If a user is defined as the creator of an image, each user can be reliably associated with a set of images they contributed to the training dataset. Appropriate single-owner user contribution bounding (Section 7.2) can then be used to ensure user-level privacy. We might also want to treat someone who directly or indirectly appears in an image (e.g. a person’s face, or personal dwelling) as a user associated with the image. In this case, obtaining metadata to assign a set of users to each image may be technically challenging, raise other potential privacy concerns (e.g., the use of facial recognition to cluster images), or even be practically impossible. But when such attribution is possible, multi-owner user-contribution bounding techniques will be needed (Section 7.2)

4.2 Methods for DP Synthetic image generation

Several paradigms exist for generating synthetic image data that satisfies DP:

- **DP finetuning or training.** DP Training from scratch generative models (e.g. GANs and VAE) using private data or DP finetuning of pretrained models like LLMs or Diffusion using private data are the two most common and powerful approaches for complex data like images. This is achieved via previously discussed DP-SGD methods and its variants. Such training process ensures that the final model parameters are differentially private with respect to the training data. Samples drawn from this trained DP model constitute the synthetic dataset. We cover these methods in Section 4.3.
- **Methods that utilize API-only access of generative models.** This is another group of methods that foregoes expensive DP-Training and instead leverage pre-existing large foundation models accessible via APIs. DP is introduced during the stage of choosing the best examples generated by public models to best match private data distribution, Private evolution (PE) is one prominent method of this class. This algorithm queries “generation” and “variation” APIs iteratively, using a DP selection mechanism to guide the process towards the private data distribution (Lin et al., 2023). We cover this algorithm in Section 4.4.
- **Private prediction.** This group of methods uses ideas of Private Prediction to train GANs models without doing an end-to-end DP-Training with DP-SGD algorithms. We explore PATE-style algorithms for image generation in Section 4.5.1
- **Creating DP synthetic image data via an intermediary text representation.** These types of methods change the modality of the data from image to text via captions and then proceed to create DP synthetic captions using methods for DP synthetic text data. Synthetic images are then created using dp-fied captions. We cover this new but promising class of methods in Section 4.5.2
- **DP data release mechanisms.** This class of approaches have not been particularly successful. These approaches, instead of training or finetuning a complex generative model privately with DP-Training, focus on releasing certain statistics or intermediate representations of the private data in a DP manner. Synthetic data is then generated based on these sanitized statistics. Examples include:
 1. DPGEN (Chen et al., 2022b) avoids DP-SGD entirely and instead privatizes training images themselves by deriving DP-randomized recovery directions. This data is used to train (without DP) an energy network Dockhorn et al. (2023) demonstrated that the final trained model was not actually DP and outlined multiple sources of privacy leakage.
 2. DP-DRE (Wu et al., 2023) which utilizes a publicly pretrained encoder and ICGAN generator. Private data is embedded via an encoder and the distribution is DP-fied via DP density estimator. Samples are then drawn from this distribution and decoded via ICGAN (Liu and Jin, 2024)
 3. DP-MERF models (Harder et al., 2020) use random feature representations of kernel mean embeddings and learns synthetic data representation with Maximum Mean discrepancy objective, that pulls the embeddings of synthetic and private data together. Authors “release” mean private data embeddings’ with DP and use this representation in closed-form calculation for MMD.

We will next focus on the first four groups of approaches in the upcoming sections.

4.3 DP training/finetuning

DP training (from scratch) or finetuning (existing pre-trained models) are currently the most common and prominent approaches for synthetic image data generation. A generative model (e.g., GAN, VAE, Diffusion Model, LLM) is trained or finetuned on the private dataset using an optimization algorithm that incorporates DP, most notably Differentially Private Stochastic Gradient Descent (DP-SGD) and its variants (Section 2.5). As we have seen before, such training process ensures that the final model parameters are differentially private with respect to the training data. Samples drawn from this trained DP model constitute the synthetic dataset.

4.3.1 GAN based models

Early attempts at generating DP synthetic images used Generative Adversarial Networks (GANs). GAN training involves a generator G trying to fool a discriminator D , which tries to distinguish real images from generated ones. The most common method for creating DP-GANs involves integrating Differentially Private Stochastic Gradient Descent (DP-SGD) into the standard GAN training loop (Xie et al., 2018) (DP-Training). Often DP is primarily applied to the discriminator D as it directly processes the real private data, although variants that also apply DP to generator exist (e.g. (Chen et al., 2020a)).

DP-Training GANs however presents significant challenges: Firstly, standard GAN training is notoriously sensitive to hyperparameters and prone to instability, such as mode collapse (Thanh-Tung et al., 2018). The introduction of noise and gradient clipping via DP-SGD often exacerbates these issues, making it even harder to achieve stable convergence and produce high-quality results (Chen et al., 2022b): the already delicate balance between the generator and discriminator is easily disrupted by the DP modifications.

The inherent difficulty of stabilizing GAN training under the noisy and biased gradient conditions imposed by DP-SGD has been a major driver for research in this area. Many subsequent approaches have sought to either refine the application of DP-SGD within the GAN framework or to bypass the standard GAN/DP-SGD combination entirely, seeking more compatible pairings of generative models and DP mechanisms.

Several strategies have been explored to improve DP-GAN (Xie et al., 2018) performance and address the aforementioned challenges. Torkzadehmahani et al. (2019) introduce a framework for training Conditional Generative Adversarial Networks (CGAN) while providing DP guarantees. The authors’ primary contribution is a gradient clipping and perturbation strategy that treats the gradients from real and fake data separately within the discriminator, which they demonstrate improves the quality and utility of the generated synthetic data and labels compared to previous DP approaches.

Chen et al. (2020a) point out that since the generator is the only “released” object, DP should be applied only to gradients of the generator, reducing the total noise injected during training. Additionally, their method GS-WGAN (Chen et al., 2020a) modifies how gradients are processed before noise addition. Instead of relying on clipping arbitrary large gradients to limit their sensitivity (which is needed to scale the noise added to the gradients), an alternative loss (Wasserstein-1 loss) that generates bounded gradients with norms near 1 is used. This allowed authors to forego expensive hyperparameter search for the appropriate clipping norm. Additionally, since gradients were closer to the clipping norm, gradient distortion from DP was much less significant.

Bie et al. (2023) reexamines the DPGAN formulation of Xie et al. (2018) where only the discriminator is trained with DP-SGD, and find that taking many discriminator steps between consecutive generator steps (e.g. >50 in some cases) successfully addresses the imbalance between the generator and discriminator, leading to large improvements in quality. With this modification, Bie et al. (2023) report that a well-tuned application of DP-SGD to the GAN discriminator outperforms all alternative GAN privatization schemes.

Recognizing the difficulty of learning complex image distributions from scratch under DP-Training, some methods incorporate public data (DP-finetuning of publicly pre-trained models). *DP-GAN-MI* (Chen et al., 2022a) uses a pre-trained feature extractor but trains a standard DP-GAN directly in the feature space to model the private feature distribution. DP-DRE (Wu et al., 2023) first trains an autoencoder (specifically, an Invertible Conditional GAN or IC-GAN) non-privately on public data. Then, both private and public images are mapped to the IC-GAN’s feature space. A DP mechanism (specifically, training an auxiliary GAN as a density ratio estimator with DP-SGD) is used to learn the difference between the private and public feature distributions. To generate, features are sampled according to the learned DP ratio from the public feature distribution, and then decoded back to images using the IC-GAN’s generator. This also avoids direct DP-GAN training in the high-dimensional pixel space and leverages the structure learned from public data (Wu et al., 2023). DP-DRE was reported to perform better than DP-GAN-MI (Wu et al., 2023).

Several works avoid adversarial objective of standard GANs altogether, that is hard for DP-Training and finetuning and explore alternative losses than Wasserstein distance. *DP-Sinkhorn* (Cao et al., 2021): This method use an optimal

transport (OT) based loss, specifically the Sinkhorn divergence, measuring the distance between the real and generated data distributions. The generator is trained to minimize this divergence. The approach relies on a new semi-debiased Sinkhorn loss to better control the bias-variance trade-off inherent in estimating gradients under DP constraints. By avoiding the adversarial objective, DP-Sinkhorn aims for more stable training (Cao et al., 2021).

DP-GANs generally exhibit a pronounced privacy-utility trade-off. While progress has been made, achieving high visual fidelity, diversity, and resolution comparable to non-private state-of-the-art GANs, especially under strong privacy guarantees (low ϵ), remains a significant challenge. Methods like DP-Sinkhorn (Cao et al., 2021) reported improvements over earlier DP-SGD-based GANs. Techniques that effectively utilize public data, like DP-DRE (Wu et al., 2023) can also boost performance by reducing the learning burden on the private data under DP constraints.

Refer to Table 5 for a brief comparison between GAN-based and alternative methods.

Method	Core Model	DP Mechanism	Advantages	Challenges
DP-GAN (Xie et al., 2018)	GAN	DP-SGD on Discriminator (or both)	Foundational approach	Instability, Low Resolution, Poor Utility
GS-WGAN (Chen et al., 2020a)	GAN (WGAN)	DP-SGD on ensemble image gradients	Noise on lower-dim image gradient; ensemble stability	Still relies on DP-SGD dynamics
DP-DRE (Wu et al., 2023)	Autoenc. (IC-GAN) + Aux GAN	DP-SGD on auxiliary GAN in feature space	Leverages public data; avoids DP-GAN in pixel space; learns density ratio	Requires good autoencoder & public data
DP-Sinkhorn (Cao et al., 2021)	Generator (OT-based)	DP noise on Sinkhorn loss gradients w.r.t. samples	Replaces adversarial loss with OT; improved stability; semi-debiased loss for DP gradients	OT computation complexity

Table 5: Comparison of Selected DP-GAN and GAN alternatives for Image Synthesis . All methods ensure approximate (ϵ, δ) DP.

4.3.2 Diffusion based models

Diffusion probabilistic models have recently emerged as the state-of-the-art for high-fidelity image generation (Liu et al., 2023a). Naturally, significant effort has been directed towards applying differential privacy to these models to enable the generation of private synthetic images. Diffusion models work by gradually adding noise to data in a "forward process" and then learning a model to reverse this process, starting from pure noise and iteratively denoising it to generate a sample (Liu et al., 2023a).

Similar to DP-GANs, the predominant technique for achieving DP in diffusion models is DP-Training/Finetuning via DP-SGD (Liu et al., 2023a). DP-SGD is applied during the training phase where the model (typically a UNet-like architecture (Liu et al., 2023a)) learns to predict the noise added at each diffusion time step t , or equivalently, predict the less noisy image x_{t-1} or the original image x_0 . Commonly used loss for diffusion model training is mean squared error between the predicted noise and the actual noise added.

To potentially improve utility under DP-SGD, techniques like "noise multiplicity" (averaging gradients over predictions made with different noise samples added to the same image at time t) (Liu et al., 2023a) or "augmentation multiplicity" (averaging the loss over multiple random augmentations of an image before computing the gradient) (Park et al., 2024) have been proposed.

DP-training vs DP-finetuning The initial attempts at DP Diffusion models involved training (from scratch) a diffusion model on the private dataset from random initialization using DP-SGD throughout the process (e.g., DPDM (Liu et al., 2023a), DPLDM (Liu et al., 2023a)). While feasible, this can be computationally very expensive and often results in lower utility compared to methods using pre-training, as the model has to learn complex image features entirely under the constraints of DP noise (Liu et al., 2023a). Therefore DP-finetuning of a pretrained Diffusion models has become a dominant and highly effective paradigm (Liu et al., 2023a). *DP-Diffusion* is a diffusion model that is first pre-trained on a large, publicly available dataset (like ImageNet) without privacy constraints. This allows the model to learn general visual features and structures. Subsequently, this pre-trained model is fine-tuned on the target private dataset using

DP-SGD for a smaller number of steps (Liu et al., 2023a). Since the model starts from a strong initialization, the fine-tuning process requires less learning from the private data, thus accumulating less privacy loss and allowing for better utility for a given budget ϵ (Liu et al., 2023a). *Privimage* (Li et al., 2024) refines this idea by adding a preliminary step: it uses DP queries (based on semantic similarity using CLIP embeddings) to select a small, relevant subset from the large public dataset that semantically aligns with the private data. Pre-training is then performed only on this selected subset before DP fine-tuning. This aims to make the pre-training more targeted and efficient, potentially leading to better results with fewer parameters and less computation compared to pre-training on the entire public dataset (Li et al., 2024). *DP-RandP* (Tang et al., 2023) investigates pre-training on synthetically generated random data (e.g., images from random processes) before DP fine-tuning on the private data, offering an alternative when large relevant public datasets are unavailable. When only a small amount of in-distribution public data is available (e.g., a small public subset of a larger private dataset), work by Park et al. (2024) proposes to use a diffusion model trained on this initial small public dataset to synthesize more public-like data. This amplified public dataset is then used for more effective warm-up training before applying DP-SGD fine-tuning on the private data. This addresses the challenge of leveraging public data when large, diverse external datasets are not suitable or available (Park et al., 2024).

Efficiency Improvements via Latent Diffusion (DP-LDM): To tackle the computational cost and high dimensionality of pixel-space diffusion models, DP-LDM (Liu et al., 2023a) operates on Latent Diffusion Models (LDMs). LDMs use a pre-trained autoencoder to map images to a lower-dimensional latent space, where the diffusion process takes place (Liu et al., 2023a). DP-LDM applies DP-SGD during the fine-tuning stage only within this latent space. Furthermore, it drastically reduces the number of parameters updated with DP-SGD (by $\sim 90\%$) by fine-tuning only specific components critical for adaptation and conditional generation, namely the attention modules within the UNet backbone and the conditioning embedders (if applicable), while keeping the rest of the UNet frozen (Liu et al., 2023a). This targeted approach significantly reduces computational cost and has been shown to achieve a better privacy-utility trade-off (Liu et al., 2023a).

Leveraging Inherent Diffusion Noise (dp-promise): This method (Wang et al., 2024a) explicitly leverages the fact that the forward diffusion process itself injects Gaussian noise, which contributes to privacy (satisfying Gaussian DP). The intuition behind the method is that during forward diffusion pass, later iterations of images are very close to the noise already, whereas earlier iterations are much more sensitive and closer to the private data. Accordingly, the paper proposes to modify the training that happens during reverse diffusion propose. Authors introduce a two-phase training strategy. Phase I trains the model non-privately on the later, noisier diffusion steps (from time S to T), relying on the inherent DM noise for a partial DP guarantee. The gradients from this stage are already noised enough from the noise introduced during the forward pass, and sensitivity bound is derived based on image dimensions. Phase II trains the model using DP-SGD (clipping and injection of the noise) only on the earlier, less noisy steps (from time 1 to $S - 1$), where additional privacy protection is needed. By avoiding redundant DP-SGD noise injection in the later stages, dp-promise aims to improve the overall privacy-utility balance (Wang et al., 2024a).

Discussion The success of many state-of-the-art DP diffusion methods underscores the critical importance of leveraging external information, typically through large public datasets used for pre-training (Liu et al., 2023a). Pre-training allows the model to learn powerful, general-purpose visual representations without consuming any private budget (Liu et al., 2023a). The subsequent DP fine-tuning phase then specializes the model to the private data distribution, requiring fewer updates and thus less noise injection compared to training from scratch (Park et al., 2024). The relevance and quality of the public data can significantly impact the effectiveness of transfer learning (Li et al., 2024). This heavy reliance suggests that the vast knowledge encoded in models pre-trained on public data is currently almost essential to counteract the information loss imposed by DP constraints when learning from limited private data, especially for complex tasks like high-fidelity image generation. Achieving high utility purely from private data under strong DP remains largely impractical with current diffusion techniques.

DP diffusion models employing pre-training and fine-tuning strategies like DP-Diffusion and DP-LDM have demonstrated significant progress, achieving state-of-the-art results in DP image synthesis (Liu et al., 2023a). They often outperform earlier DP-GAN methods in terms of image fidelity (FID scores) and the utility of generated data for downstream tasks (e.g., classification accuracy) (Liu et al., 2023a). Techniques like DP-LDM also offer substantial improvements in computational efficiency by operating in latent space and fine-tuning only a fraction of the parameters (Liu et al., 2023a).

Innovations in DP diffusion are increasingly moving beyond simply applying DP-finetuning with DP-SGD as a black box. Instead, they focus on how and where to integrate DP most effectively within the diffusion model framework. This includes operating in latent spaces (DP-LDM (Liu et al., 2023a)), selectively fine-tuning parameters (DP-LDM (Liu et al., 2023a)), adapting the DP mechanism based on the diffusion time step (dp-promise (Wang et al., 2024a)), and

optimizing the use of auxiliary data (Privimage (Li et al., 2024), Park et al. (Park et al., 2024)). This trend signifies a shift towards more model-aware and efficient DP techniques tailored to the specific characteristics of diffusion models.

It is worth pointing out recent success of autoregressive models for image generation in non-DP setting (Sun et al., 2024). Visual autoregressive models (VARs) apply next-token-prediction paradigm that dominates text domain to vision domain. These models work by training special image tokenizers that generate sequence of image tokens that are then used for training transformer-based models. During generation, sequence of image tokens is generated token-by-token and subsequently decoded into the image by a trained decoder. Early attempts of adopting this paradigm for DP image synthesis (Shaikh et al., 2025) however significantly underperform DP-diffusion based models that have been extensively researched and iterated upon in DP setting.

Method Name	Core nique	Tech- nique	DP Mechanism	Advantages	Challenges
DPDM (Liu et al., 2023a)	Train from Scratch		DP-SGD	Foundational DP-Diffusion; Noise Multiplicity	High Compute Cost, Lower Utility
DP-Diffusion (Generic) (Dockhorn et al., 2023)	Pre-train (Public) / Fine-tune (Private)		DP-SGD during Fine-tuning	Leverages public pre-training for better utility/efficiency	Requires large public dataset; still costly
Privimage (Li et al., 2024)	Select Data / Pre-train / Fine-tune	Public	DP Query + DP-SGD Fine-tuning	Targeted public data selection for efficient pre-training; reduced model size	DP query overhead; relies on embedding
DP-Latent Diffusion Model (LDM) (Liu et al., 2023a)	Pre-train LDM (Public) / Fine-tune Attention (Private)		DP-SGD on Attention in Latent Space	Operates in latent space; tunes few parameters ($\sim 10\%$); high efficiency & utility	Requires LDM architecture
dp-promise (Wang et al., 2024a)	Two-Phase Training (DM Noise + DP-SGD)		Partial DP-SGD + Inherent DM Noise	Leverages DM noise for DP guarantee; reduces redundant DP-SGD noise injection	Complex training schedule; choice of S
Park et al. (Park et al., 2024)	Synthesize Public Data / Warm-up / Fine-tune		DP-SGD during Fine-tuning	Amplifies limited public data for better warm-up; addresses in-distribution public data	Extra step of public data synthesis

Table 6: Comparison of Selected DP-Diffusion Approaches for Image Synthesis . All methods ensure approximate (ϵ, δ) DP.

4.4 Methods that avoid DP-Training: Private Evolution

As we have seen in the previous section, DP training or finetuning models like GANs or Diffusion models is extremely computationally expensive - DP-SGD algorithms require much larger batches, expensive per example clipping and potentially many more iterations to converge (Ponomareva et al., 2023). Avoiding DP finetuning is thus a promising avenue that has been explored in a method called *Private evolution (PE)* introduced by Lin et al. (2023), which foregoes tuning of any models completely and relies on existing foundational models’ APIs. DP is introduced during the selection stage when best synthetic examples generated by public models are chosen to match private data distribution.

The central idea of PE is to leverage the generative capabilities of powerful, potentially proprietary or computationally expensive-to-train foundation models (like Stable Diffusion) as black boxes (Lin et al., 2023). Instead of training these models with DP, PE uses an evolutionary algorithm guided by DP principles to iteratively refine a population of synthetic samples steering them towards the distribution of the private data using only API calls for generation and variation (Lin et al., 2023). The fact that the approach is inference-only simplifies deployment and allows capitalizing on state-of-the-art models that might only be accessible via APIs (Lin et al., 2023). Additionally, private data is never directly input into the generative model - this means that the API provider can be potentially untrusted (Lin et al., 2023). This new paradigm makes advanced DP synthesis potentially more accessible, bypassing the need for direct model access, gradient computation, or extensive private training infrastructure, provided suitable APIs are available.

PE relies on two APIs, namely *random api* and *variation api*, that are often directly available for popular models (e.g. DALL-E (OpenAI, 2022), Stable Diffusion (Rombach et al., 2021)) or can be implemented with appropriate prompt engineering (e.g. GPT (Chen et al., 2020b)). The *random api* generates random samples (which can be conditioned

Algorithm 4 Private Evolution for Image (PE, Lin et al. (2023))

Require: Private data D , image embedding model Φ (e.g. CLIP, Inception), target number of synthetic samples N , evolution rounds T , population size multiplier L (number of variations for each synthetic image).

Ensure: Synthetic dataset \hat{D} .

```

1: Initialize  $\hat{D}_0$  of size  $L \cdot N$  with Random API.
2: for  $t = 0, \dots, T - 1$  do
3:    $E_t = \Phi(\hat{D}_t)$  //Embedding calculation for synthetic samples.
4:   Let each  $\Phi(z), z \in \hat{D}$  vote for the nearest embedding in  $E_t$ .
5:   Privatize the voting results with  $(\epsilon, \delta)$ -DP to get a DP histogram  $H_t$ .
6:    $\hat{H}_t = H_t / \text{sum}(H_t)$  //Histogram normalization.
7:   Get  $\hat{D}'_t$ : sample  $N$  images with replacement from  $\hat{D}_t$  proportionally to  $\hat{H}_t$ .
8:   if  $t < T - 1$  then
9:     Get  $\hat{D}_{t+1}$  of size  $L \times N$ : call Variation API to get  $L$  variants for each  $z \in \hat{D}'_t$ .
10:  else
11:    return  $\hat{D}'_t$ .
```

with text prompts for some models), while *variation api* creates variants of a given sample by generating images similar to the provided image.

The PE algorithm, presented in Algorithm 4, operates iteratively as follows (Lin et al., 2023):

Initialization A starting population of synthetic samples, \hat{D}_0 , is created, potentially using a generic generation API of the chosen foundation model. Alternatively, if available, a public dataset that is similar to the sensitive data distribution can be employed as the initial pool.

Iterative Refinement Once the initial pool of image data is created, the algorithm refines it through a repeated series of voting and variation steps.

1. *Embedding*: Both the private dataset $D = \{x_1, \dots, x_n\}$ and the current synthetic population $\hat{D}_t = \{z_1, \dots, z_m\}$ are mapped into a common embedding space using a pre-trained embedding function Φ (e.g., Inception network, CLIP embeddings). This allows for semantic comparison (Lin et al., 2023).
2. *Selection (Parent Selection via DP Nearest Neighbors Histogram)*:
 - (a) **Nearest Neighbor Identification**: For each private sample $x_i \in X$, its nearest neighbors $z_{\text{NN}(i)}$ is identified within the current synthetic population \hat{D}_t based on the ℓ_2 distance in the embedding space: $d(x_i, z_j) = \|\Phi(x_i) - \Phi(z_j)\|_2$
 - (b) **Histogram Construction**: A histogram is built over the synthetic samples in \hat{D}_t . The count for each synthetic sample z_j is the number of private samples x_i for which z_j is the nearest neighbor: $\text{Count}(z_j) = |\{i | z_j = z_{\text{NN}(i)}\}|$. To ensure differential privacy, noise is added to these counts. Specifically, i.i.d. Gaussian noise $\mathcal{N}(0, \sigma^2)$ is added to each $\text{Count}(z_j)$, where σ is a noise multiplier related to the desired privacy budget and H is a threshold parameter. The sensitivity is 1 since each private example gets only 1 vote. The noisy counts are optionally thresholded at H . Let the resulting noisy, thresholded count for z_j be $\widetilde{\text{Count}}(z_j)$. This DP histogram mechanism provides an (ϵ, δ) -DP guarantee for the release of the histogram information in each iteration
 - (c) **Parent Sampling**: A new set \hat{D}'_t of N "parent" samples is selected from \hat{D}_t by sampling with replacement, where the probability of selecting z_j is proportional to its noisy, thresholded count $\widetilde{\text{Count}}(z_j)$. Samples deemed more similar to the private data (according to the DP histogram) are more likely to be chosen as parents.
3. *Variation (Offspring Generation)*:

For each selected parent sample z_p from \hat{D}'_t , a "variation API" of the foundation model is called to generate L "offspring" samples. This API is designed to produce samples that are similar to the input z_p but introduce some diversity or modification (e.g., using image-to-image generation with slight changes, varying guidance scale in diffusion models) (Lin et al., 2023).
4. The collection of all generated offspring forms the new synthetic population \hat{D}_{t+1} for the next iteration.

Final Output: After T iterations, the final population D_{T-1} (or a subset thereof) constitutes the DP synthetic dataset. The overall privacy cost $(\epsilon_{\text{total}}, \delta_{\text{total}})$ is calculated by composing the privacy costs of the T DP histogram releases using advanced composition theorems (Lin et al., 2023).

The choice of the initial set, the distance metric and the variation API are crucial. Using ℓ_2 distance in a semantic embedding space (like CLIP (Radford et al., 2021)) allows comparison based on meaning rather than superficial pixel differences (Lin et al., 2023). The paper also proposes a "lookahead" distance metric modification, where the distance for selecting z_j is computed between $\Phi(x_i)$ and the mean embedding of k potential variations of synthetic sample z_j generated by the variation API (Lin et al., 2023):

$$d(x_i, z_j) = \left\| \Phi(x_i) - \frac{1}{k} \sum_{l=1}^k \Phi(\text{VARIATION_API}(z_j)_l) \right\|_2 \quad (8)$$

This heuristic attempts to select parents that are likely to produce good offspring in the next generation, potentially accelerating convergence by anticipating the evolutionary dynamics. The effectiveness of PE inherently depends on the quality of these components: the embedding must capture relevant features, and the variation API must allow effective exploration around promising regions of the sample space (Xie et al., 2024).

User-level privacy unit adjustments The original algorithm can realize the user-level privacy unit by normalizing the private votes from samples from the same user to sum to 1 (in ℓ_2 norm). This ensures the sensitivity w.r.t one privacy unit (user) remains one and that the noise is scaled appropriately.

Discussion The initial PE paper reported compelling results for image generation (Lin et al., 2023). On benchmarks like CIFAR10 (using ImageNet as the implicit public data via the foundation model), PE reportedly achieved FID scores comparable to or better than state-of-the-art training-based DP methods (like DP-Diffusion, DP-MEPF) but with significantly smaller privacy budgets (lower ϵ) (Lin et al., 2023). For example, an FID score below 7.9 was achieved with $\epsilon = 0.67$, whereas prior work required ϵ values around 32 for similar FID levels. Further, PE was successful at generating high-resolution (512x512) images using Stable Diffusion for inference, even when the private dataset was small (e.g., 100 cat images) (Lin et al., 2023). This is a challenging setting where training large models from scratch or fine-tuning with DP might be difficult due to data scarcity and computational cost. Additionally, PE showed limited adaptability to distribution shift between the foundation model’s implicit training data (e.g., ImageNet) and the private target data (e.g., Camelyon17 medical images). While downstream classification accuracy on Camelyon17 (79.56% at $\epsilon \approx 7.6$), as expected, lagged behind a specialized DP fine-tuning approach (91.1% at $\epsilon = 10$) in this data-rich scenario, the result was still non-trivial, indicating the framework’s potential adaptability (Lin et al., 2023). Finally, compared to DP fine-tuning methods which can require hundreds or thousands of GPU hours and access to checkpoint weights (Liu et al., 2023a), PE’s training-free nature can make it more computationally efficient, depending on the cost and latency of API calls (Liu et al., 2023a).

The downsides of PE is that its performance depends heavily on the quality of the foundation model APIs (generation, variation, embedding) (Xie et al., 2024). Since those apis can be parameterized by prompts in text-to-image models, extensive prompt engineering must be needed. Additionally, while PE can handle some distributional difference between pretraining data for the foundational models and private data to be mimicked, large distributional shifts will result in poor quality synthetic image data. Lin et al. (2025) demonstrated that instead of using foundational models (and relying on private data to be of similar distribution to their pretraining data), PE can use *synthesizers that do not rely on neural networks* like graphics-based image generation or physics-based robotics simulators. Such synthesizers can be further combined with foundational models by using simulators during early algorithm iterations to create diverse pool of samples and then using foundation model APIs to enhance quality of the final synthetic images. For situations involving large distributional shifts between pretraining and private data in domains where good image simulators exist, Sim-PE offers a significant boost (improving the accuracy x3 and reducing FID by 80% at $\epsilon = 10$) over the original PE algorithm.

4.5 Alternative methods

We want to briefly mention two additional lines of work that does not fit neatly into DP-Training or Training-Free categories we analyzed so far.

4.5.1 PATE-style models

PATE (Private Aggregation of Teacher Ensembles) (Papernot et al., 2017, 2018) is an example of introducing DP at the prediction stage, in contrast to DP-Training that alters training process of the models. The gist of PATE is to split the

private data into a number of disjoint datasets, train a number of non private models and introduce the noise to their aggregated prediction proportional to the level of agreement of the submodels.

PATE-style GAN models still train a GAN model, but avoid DP-SGD and instead privatize predictions or gradients of the discriminator using ideas from PATE.¹⁰

A PATE-GAN model (Yoon et al., 2019) was initially introduced in context of tabular data generation but subsequently was adopted and used as baselines for a number of DP image synthetic data generation papers. PATE part (DP-inference) is applied to obtaining a DP version of a discriminator: private data is partitioned into a number of disjoint subsets of data, a number of discriminator teacher models are trained on those subsets, each trained with a goal of improving their loss with respect to the generator G . To classify each new sample as either real or synthetic, all teacher models vote on the sample and their output is noisily aggregated. However, since discriminator and generator in GANs settings need to be trained together, a discriminator should be differentiable, which PATE style discriminator isn't. Instead, this collection of submodels is distilled into a student model S , which can be differentiated over. Student is trained by taking some public unlabelled data and using prediction from PATE ensemble as pseudo labels. Alternatively when public data is not available, student is trained using the data generated by the Generator and pseudo-labels from teachers. Generator G and student S are co-trained, where generator is trying to fool the student, student is trying to improve its loss w.r.t teachers and teachers are training to improve their loss wrt generator.

G-PATE (Long et al., 2019) argues that making a discriminator differentially private is unnecessary - as long the gradients of the discriminator used to update the generator are DP, the whole process is DP. Additionally, discriminator can be trained on the private data directly without relying on public data or using samples from the generator itself, as in PATE-GAN. Consequently, G-PATE does not dp-fy aggregated predictions of teacher discriminators and foregoes a student discriminator completely. A gradient aggregator gets gradients from teachers and accumulates them in PATE-style before passing this private gradient to generator for update. With such scheme G-PATE demonstrated impressive ability to generate high dimensional image data with high utility with very low privacy budgets ($\epsilon \leq 1$)

DataLens (Wang et al., 2021a) follows the similar idea and, instead of using random projection of gradients before introducing the noise, as in G-PATE (Long et al., 2019), authors employ dimensionality compression that allows to further reduce the noise needed for private aggregation. PATE-TripleGAN takes this solution one step further (Jiang et al., 2024) and introduces modifications that allow to use unlabeled private data during GANs training. Namely, they introduce a classifier to pre-classify unlabeled private data.

It is worth mentioning that PATE style algorithms introduced computational and architectural complexity that hinders their adoption in real world and academic research. Additionally, Ganev et al. (2025a) highlight that increased system complexity of many PATE-GAN style implementations leads to implementation bugs that weaken/invalidate privacy guarantees and result in more than expected privacy leakage. In terms of utility, Bie et al. (2023) report significantly better results with DP-Training of GANs compared to PATE-style training of GANs.

4.5.2 Methods that use an intermediary representation

Two concurrent works recently proposed creating DP synthetic images via first creating text captions of private images, then DP-fying these captions either via DP-finetuning (Section 5.2) (Kong and Syed, 2025) or Private Evolution (Section 5.3.2) (Wang et al., 2025), and subsequently creating images based on the DP-captions. Kong and Syed (2025) used standard DP-finetuning, but in a hierarchical manner: they first trained a model to generate privatized album descriptions, and then trained a model to create privatized photos descriptions given the DP-fied album descriptions. The hierarchical generation strategy ensured thematic and character consistency within each album, while circumventing the context window limitations that would be met if large photo albums were modeled in one go. Instead of DP fine-tuning, Wang et al. (2025) adopted the Private Evolution method. A key challenge with PE is that caption embeddings do not always reflect similarity in the resulting images. To address this, they proposed using embeddings of the generated images during the voting process, albeit at the cost of increased computation.

This promising family of approaches has proven powerful in photo generation settings. Using text as an intermediate representation leverages the main strength of a large-language model. Text is a rich representation that is however much less dimensional than high-quality images (Wang et al., 2025)

The downside of using text as intermediary is increased complexity and compute requirements stemming from reliance on two additional models - one that creates captions and the one that renders images given captions. Additionally, not all image data can be easily captionized - e.g. screenshots of apps on a phone or a monitor and images with text are notoriously hard to describe and generate based on descriptions.

¹⁰PATE-style GANs models sit in an awkward place, where one can consider them doing DP-Training with an alternative way/place of introducing the DP.

4.6 Comparison of the methods

We have previously outlined the main methods for creating DP synthetic text data: via DP finetuning of GANs (Section 4.3.1) and Diffusion models (Section 4.3.2) and DP-training free methods: PATE-style (Section 4.5.1) and private evolution (Section 4.4). Each of the aforementioned methods have their strengths and weaknesses, as well as different requirements. Table 7 compares various aspects of each method.

Choosing the method for DP synthetic text data generation

For large enough volume of sensitive data (>XXK datapoints) with enough compute the method that results in highest fidelity and utility of DP synthetic data is almost always DP-Finetuning of pretrained Diffusion models. Less computationally expensive methods like PE that don't require finetuning can provide reasonable data when sensitive data is somewhat in distribution for the pretraining data or for situations when were strict privacy guarantees (low ϵ) are needed.

- **DP finetuning** of generative models like GANs and Diffusion remains the workhorse method that delivers the best quality given sufficient amount data, compute and engineering and time investment. However, both GANs and Diffusion models are though hard to DP-finetune, both for computational reasons and due to the nature of the models. For DP-Diffusion, access to relevant large public datasets for pre-training (or access to a pre-trained checkpoint) is often highly beneficial. Diffusion models are more computationally expensive but are likely to provide better utility of synthetic data (e.g current state-of-the-art DP-LDM (Liu et al., 2023a)) than GANs based ones. DP-GANs have historically faced more significant utility degradation, although alternative frameworks like DP-Sinkhorn have shown promise in improving quality and stability (Cao et al., 2021).

Since the adoption of DP synthetic data is often *quality-bottlenecked*, the best option given unlimited computational resources and large amount of data is likely DP-finetuning of a diffusion model.

- **Private evolution** is most applicable when training or fine-tuning large models is computationally infeasible, or when access to powerful foundation models is restricted to inference APIs. Additionally, PE requires an access to a suitable embedding space (that captures aspects of the final use of the image data) and variation and random API must be available for the image modality. PE is the only method that has a chance of obtaining reasonably quality data on very small private datasets. Additionally, PE can provide more stringent privacy guarantees (low ϵ) that for finetuning methods will result in a very low utility.
- If high resolution/visual quality of generated examples is of paramount importance (as opposed to fidelity of matching the original image distribution), using methods that don't use DP-trained generators will produce the best quality images. For these scenarios, either using PE or image-via-proxy methods like (Kong and Syed, 2025; Wang et al., 2025) (Section 4.5) will like produce the best resolution at the cost of potential lower fidelity (in case of PE) and increased compute/complexity in case of proxy methods.

The field appears to be converging on the strategy of leveraging large, publicly pre-trained models as the most effective way to achieve high utility in DP image synthesis. Whether this leverage occurs via APIs (as in PE) or through DP fine-tuning (as in DP-Diffusion), the general knowledge captured in these large models seems crucial for offsetting the information loss imposed by DP when learning specifics from private data. Training high-quality, complex generative models like diffusion models purely from limited private data under strong DP guarantees remains largely impractical today.

Aspect	Method		
	DP-finetuning (GANs, Diffusion)	PATE-GANs	Private evolution
Amount of input private data			
<i>Small input quantity (<5K)</i>	Not recommended	Not recommended	Preferred
<i>Large input quantity (>10K)</i>	Preferred	Recommended	Not recommended
Reliance on Pre-trained Models / Public Data			
<i>Can benefit from additional public data</i>	Yes both for Diffusion and GANs	Distilling PATE teachers into students can utilize public data	No
<i>Needs pretrained models</i>	Yes for Diffusion, models should be pretrained to improve quality.	No	Implicit (via Foundation Model API)

Yield	Unlimited number of output examples, although with diminishing returns to downstream task performance.	Unlimited number of output examples, although with diminishing returns to downstream task performance.	In practice, suitable for outputting synthetic dataset of size \leq size of input private dataset.
Model access required	Weights	Weights of Generator and Student network and a number of teachers	Generations via API
Compute resources and engineering effort			
<i>Training of generative models is required</i>	Yes, 1	Yes, training a Generator, A number of teacher discriminators and a distilled student discriminator, in alternating fashion	No
<i>Inference cost multiple per synthetic example</i>	1 (same as regular inference).	\propto (number of PATE submodels to aggregate over).	\propto (number of iterations) \times (number of variants per sample).
<i>Prompt engineering required</i>	No	No	Possibly – if using multi-modal (text, image) models, potentially need to craft prompts for initial pure synthetic data and variate templates.
<i>Time to first example</i>	Long (Run finetuning on the entire dataset, then sample)	Extra long (train multiple models on disjoint subsets, infer on all of the models to get pseudo labels, distill into a student model, repeat the process for a number of iterations after Generator’s gradient updates)	Medium (Might require prompt engineering + running variate and embedding on the entire private dataset)
<i>Direct side by side comparison of inputs and outputs</i>	No	No	No
<i>Resilience to distribution gaps between private data and LLM</i>	High	Medium	Low
Target privacy guarantee			
<i>High ϵ (e.g. 10)/large amount of private data</i>	Preferred	Recommended	Not recommended
<i>Stringent (e.g. $\epsilon < 1$)/small amount of private data</i>	Not recommended (quality will be bad)	Not recommended (quality will be extremely bad)	Preferred
Data persistence requirements	Entire dataset required at once for training.	Entire dataset required at once for training.	Single examples can arrive in a streaming fashion, be used to cast votes, and then discarded immediately.

Table 7: High-level comparison of DP finetuning of GANs and Diffusion models, PATE-GANs and private evolution for DP text synthesis.

4.7 Evaluating synthetic image data quality

On top of the generic metrics applicable to all data modalities mentioned in Section 1.5, below we outlined metrics specific to evaluating image fidelity.

The fidelity of synthetic image data can be evaluated using several specialized techniques. Originating from the literature on Generative Adversarial Networks (GANs), so-called Inception-based scores are specifically designed to align with human judgment of image quality (Salimans et al., 2016). The core principle of these methods involves leveraging a pre-trained Inception model, or more generally Convolutional Neural Network (CNN), trained on public data to extract

visually relevant features from both real and synthetic images. By analyzing the distributions of these features, these scores can provide a quantitative measure of how realistic the generated data is.

Inception Score (IS) (Khetan and Oh, 2016; Salimans et al., 2016) The Inception Score (IS) operates on two key principles. First, it assumes that high-quality, realistic images should contain meaningful objects, leading to a confident classification by the Inception model. This translates to a low entropy for the predicted class probability distribution of each individual image. Second, for a diverse set of generated images, the overall distribution of predicted labels should be highly varied, reflecting a wide range of objects. This means the marginal distribution across all images should have high entropy. A significant drawback of the Inception Score is that it only evaluates the generated images in isolation, without directly comparing their statistical properties to those of the real data. This limitation prompted the development of the Fréchet Inception Distance.

Fréchet Inception Distance (FID) (Heusel et al., 2017) The Fréchet Inception Distance (FID) improves upon the IS by directly comparing the feature distributions of real and synthetic data. It models the activations from the Inception model’s coding layer for both datasets as multivariate Gaussian distributions. The distance between these two distributions is then calculated, providing a measure of their similarity. A lower FID score indicates that the distribution of synthetic data is closer to that of the real data, suggesting higher quality and diversity. The formal definition of the FID is:

$$\|\mu_1 - \mu_2\|_2^2 + \text{Tr}(M_1 + M_2 - 2(M_1 M_2)^{1/2}),$$

where μ_1, μ_2 are the mean of the two datasets to be compared and M_1, M_2 their respective covariance matrices of the Inception features for the real and generated datasets, respectively.

(Heusel et al., 2017) demonstrated that the FID is more robust and correlates better with human judgment of image quality than the Inception Score. Furthermore, by functioning as a true distance metric, it provides a more principled way to measure the dissimilarity between the distributions of real and synthetic data.

MAUVE is a newer metric that was originally introduced for text domain (Section 5.6, Pillutla et al. (2021)) and was later shown to be successful in vision domain as well (Pillutla et al., 2023). Mauve seeks to quantify the tradeoffs between type I (areas of synthetic data that are unlikely under real data) and type II errors (regions missing in synthetic data that are plausible for the real data). To calculate Mauve in practice, embeddings of real and synthetic data (obtained via an appropriate embedding model) are clustered into a number of clusters. Cluster assignments for real and synthetic data serve as a low-dimensional representation of real and synthetic data. A divergence curve that softly measures KL divergence of synthetic vs real and real vs synthetic quantifies Type I and II errors respectively. Area under this curve is the reported MAUVE metric. Experiments indicate that MAUVE metric induces the same ordering of sampling algorithms as FID, and just like FID, MAUVE ends up accounting for both quality and diversity of synthetic images (Pillutla et al., 2023).

Discussion Recent work by Gong et al. (2025) points out the lack of fair comparison in the DP image generation literature, due to different papers employing a range of different model architectures, privacy accounting assumptions, public pretraining datasets, hyperparameter ranges, and downstream evaluators. To address these issues, they develop *DPImageBench*, a unified evaluation codebase to compare methods on even footing. They support 12 methods from the literature, 9 evaluation datasets, consistent downstream evaluator training for utility measurement, and 6 fidelity metrics (including FID and IS). Their evaluations report that DP-trained diffusion models employing some form of curriculum (e.g. public pre-training) perform the best overall.

4.8 Open questions and challenges

Despite significant progress, several challenges and open questions remain critical for the field of DP synthetic image generation:

1. **Scalability and Fidelity:** Generating very high-resolution images ($> 512 \times 512$) with high diversity and photorealism under strong DP guarantees (e.g., $\epsilon \leq 1$) remains extremely difficult for all methods. DP mechanisms often smooth out fine-grained textures and details.
2. **Theoretical Understanding:** A deeper theoretical understanding of why certain techniques (e.g., the effectiveness of pre-training, the benefits of latent space fine-tuning, the dynamics of PE) work well is needed to guide future algorithm design. Understanding better challenging training dynamics of Diffusion or GANs model also has a potential to make DP-finetuning more successful.
3. **Private evolution:** While private evolution has shown promise in domains where foundation models, (Lin et al., 2023), simulators, or public datasets (Microsoft, 2025) are available, effectively applying PE frameworks in low-resource domains remains a challenge.

4. **Visual autoregressive models:** given the success of VARs in non-DP settings (Sun et al., 2024), figuring out effective ways of training them with DP is an important and promising area of research.

5 DP Synthetic Text data

Text data is information written using a natural language. Text is one of the most prominent data source on the internet and offline. With the rise of power of LLMs, a lot of data collected and trained upon is becoming textual in nature. However while other modalities like tabular and image DP synthetic data had quiet a bit of development prior to generative models expansion (e.g. workload based methods and GANs model respectively), DP synthetic text data really came to a spotlight recently as LLMs became more capable.

Roadmap This section examines methods for generating DP text data. We first focus our discussion on what is a meaningful privacy unit for text data (Section 5.1), then we dive into DP finetuning - the workhorse method for DP synthetic text generation that produces the best quality data given sufficient amount of private data and compute (Section 5.2). We also explore developments in methods that don’t require expensive finetuning of an LLM - namely DP inference (5.3.1) and private evolution (Section 5.3.2). This is followed by a high-level discussion and comparison between the various methods (Section 5.4). Section 5.5 discusses resampling from a corpus of DP synthetic text to improve utility. Next, we cover proxy metrics for evaluating synthetic text data (w.r.t. its fidelity and utility) in Section 5.6 and conclude with a discussion of open research problems pertinent to text DP synthetic data (Section 5.7).

5.1 Privacy unit for text data

Example-level privacy: The majority of papers on DP synthetic text data omit the discussion of the privacy unit, suggesting that example-level privacy level is employed. While for many applications (e.g. image classifications or classical models that work on rows of tabular data (Rosenblatt et al., 2025)) an example is a well-defined semantic concept (e.g. an image or a row), the question of what constitutes a meaningful semantic chunk of text is more ambiguous (Brown et al., 2022) and example-level DP privacy may not reflect human expectations of what it means for their data to be protected.

Importantly, the way text data is batched into examples for DP-SGD training to a large extent determines the unit-of-privacy, because per-example clipping is the primary mechanism for controlling sensitivity; hence, changes in example packing and context length could impact the privacy guarantee unless care is taken. For example, given a corpus of Wikipedia articles, an LLM can be trained treating each full article as an example, or by splitting the documents into chunks like paragraphs and treating each of them as a separate example. Employing DP-SGD in these two settings will provide DP guarantees for different privacy units. Smaller privacy units like paragraphs are much easier for DP to protect due to the artificially increased size of the dataset, but provide a weaker guarantee. Thus comparing reported DP synthetic data performance (both utility and fidelity) on different textual privacy units is largely meaningless.

User-level privacy: When each example is associated with a user (and this metadata is available), this is one the most straightforward and perhaps meaningful privacy units to use with DP with text data. To achieve user-level privacy, one can either subsample sensitive data to retain only 1 piece of textual information per user (e.g. keep only 1 text document per user) or allow multiple examples per user and adjust the privacy accounting accordingly. Additional complications can arise when attempting to do such subsampling as capping user contributions can become tricky, especially for multi-owner documents (refer to Section 7.2 for more in-depth discussion). It is worth also mentioning that some algorithms are more amenable to accounting for multiple examples per user (e.g. DP-SGD has a large body of work pertaining to such accounting changes, such as Charles et al. (2024)), while other more recent methods like private evolution (Section 5.3.2) and DP inference (Section 5.3.1) might be harder to adjust.

Sub-unit level of privacy: Brown et al. (2022) highlighted that DP assumes the need to preserve privacy of each user’s record. However a record may contain a mix of private and public information - in a sentence like ‘My social security number is XXX-XXX-XXXX’ all the tokens apart from the actual SSN could be considered public and common, and so revealing them would not be considered a privacy violation. Brown et al. (2022) argues that the ideal unit of privacy for DP is actually a “secret-level” privacy unit, that in context of DP synthetic data would mean that an addition or removal of any user secret does not change the synthetic data significantly. However realizing this definition is currently unrealistic as it will require an understanding of what constitutes a secret and automatically determining a secret’s boundaries. The field of *contextual integrity* (Nissenbaum, 2009) introduces a framework for understanding in which context private information can be shared and offers a direction for answering these questions.

Privacy unit discussion While user-level privacy unit is broadly used in industry (Xu et al., 2023b) and perhaps the most meaningful privacy unit to use, in several scenarios when modeling text, user-level DP still cannot prevent secret

leakage. In free-form text obtained from users interacting in the real world, secrets can be represented across multiple users. Under this view point, protections from example-level and user-level lie on a spectrum depending on how secrets are propagated. User-level is the appropriate fit in the case where secrets are strongly tied to a single user (e.g. personal search history), but not qualitatively stronger than example-level in other scenarios (e.g. from [Brown et al. \(2022\)](#): many people prompting an AI chatbot with same secret document provided as context). All of the algorithms that we describe next support example level privacy and can be also adjusted to account for group or user level privacy.

Finally, caution needs to be exercised when using methods that involve DP training/finetuning of LLM models (Section 5.2), as common training optimization techniques with LLMs like *packing* can subtly change the meaning of privacy unit - we cover these peculiarities in Section 7.3.

5.2 DP training and finetuning

DP-Training (or *DP-finetuning* of already publicly pretrained models) techniques that we explored in images synthetic section (Section 4.3.1) are also the workhorse methods for creating DP synthetic text data.

Many works on DP synthetic text generation explored LLMs ability for open-ended text generation ([Radford et al., 2018](#)). Works like ([Bommasani et al., 2019](#); [Yue et al., 2022](#); [Mattern et al., 2022](#); [Kurakin et al., 2023](#); [Wang et al., 2024b](#); [Yu et al., 2024](#); [Carranza et al., 2024](#)) employ text generative models (most often transformer-based LLMs ([Vaswani et al., 2017](#))), DP-Train or DP-finetune them (Section 2.5) on sensitive data and then sample synthetic examples from such models. While this recipe is conceptually straightforward ([Bommasani et al., 2019](#)), practical implementations which achieve high quality synthetic data quality required several modifications to standard open-ended text generation.

Data preparation for DP finetuning Most works format text data for finetuning in the following way. First a prefix p or set of prefixes $\{p_1, \dots, p_m\}$ is chosen. These prefixes could represent a class label or any other type of conditioning that can be used during generation of synthetic data. For example, in a benchmark task of generation of DP synthetic Yelp reviews, [Yue et al. \(2022\)](#) used the conditional prefix "Product type: P | Review score: R", while [Kurakin et al. \(2023\)](#) used "[yelp] [label]". The example for finetuning is then formed by concatenating prefix p with sensitive text t (e.g. Yelp review) that needs to be protected with DP.

Defining training objective The standard language modeling objective for LLM training (cross-entropy next token prediction loss ([Radford et al., 2018](#))) can be used without modification for DP-Training ([Yue et al., 2022](#); [Kurakin et al., 2023](#)). However, some works explored modifications to this objective.

For example, [Kurakin et al. \(2023\)](#) used the Prefix-LM ([Raffel et al., 2020](#)) formulation of loss that resulted in improvements in quality of synthetic data. Assume prefix p and suffix t tokens. Standard language modeling would use causal self-attention (each token depends only on previous tokens) for each of the tokens in pt and cross entropy loss on tokens in pt . Prefix-LM instead uses bi-directional attention for tokens in the prefix p , while causal self-attention is used for tokens in t . The generator is trained to minimize cross entropy loss calculated from only tokens in t . In the case of finetuning for synthetic data generation, t corresponds to private data. This objective teaches the LLM to produce synthetic data conditioned on the prefix p , which is known and always supplied at inference time.

Other attempts at improving conditional generation included works like [Putta et al. \(2023\)](#), which introduced an extra discriminator head that attempts to distinguish between the classes used for conditioning (e.g. Yelp "positive" and "negative" labels). The discriminator's loss is included into overall training objective, resulting in improved separability of classes representations. In a similar spirit, [Mattern et al. \(2022\)](#) introduced a loss term to penalize generation of sentences of incorrect class. [Kurakin et al. \(2023\)](#) however argued that proper hyperparameter tuning for DP-training with Prefix-LM loss can achieve better results and is sufficient ([Kurakin et al., 2023](#)).

DP training setup To achieve DP, model training could be done with DP-SGD or any of its variants. To reduce the negative effect of DP on utility, proper hyperparameter tuning ([Li et al., 2021](#); [Yu et al., 2021](#); [Ponomareva et al., 2023](#)), including batch size, clipping norm, learning rate and number of epochs should be done. Notably, optimal hyperparameters for DP can vary significantly from the optimal hyperparameters in the non-DP setup – without a sufficient sweep, it is easy to underestimate achievable model utility under DP ([Li et al., 2021](#)).

The choice of a checkpoint (pre-trained or instruction-tuned) has not been explored much in context of DP synthetic data generation, however works like ([Nasr et al., 2025](#)) suggest that pre-training checkpoints are more prone to memorization of the training data (and therefore perhaps are more suited for open-ended text generation). Instruction tuned checkpoints might require different prefixes for finetuning, including finding appropriate natural text prompts: "*Please generate a positive review about a restaurant*".

Parameter efficient finetuning (PET) approaches such as LoRA (Hu et al., 2022) has been shown to be beneficial compared to full finetuning with DP, resulting in significant improvements of quality of DP synthetic data (Kurakin et al., 2023; Yu et al., 2021).

DP training is typically slower and more memory hungry compared to standard non-DP training. Some of these issues could be alleviated with clever training techniques. Ghost clipping (Li et al., 2021) could be used to improve memory usage and thus potentially increase batch size per step. If desired batch size still does not fit into memory of accelerators, gradient accumulation (also known as virtual batch training) could be used to further increase the effective batch size per training update (Ponomareva et al., 2023). Gradient accumulation works by splitting desired training batch into smaller chunks which fit into accelerator memory, performing gradient computation for each individual chunk and then aggregating across all chunks to perform update of model parameters.

Sampling Sampling of DP synthetic data is achieved by choosing one of the training prefixes p_i and then generating continuation of the prefix using DP-finetuned LLM. This process is repeated until desired number of samples in the synthetic dataset is reached. Prefixes frequency can be chosen based on the distribution of labels in the private data (to increase the fidelity of DP synthetic data) or sampled uniformly.

Given a sampling prefix, sampling can be performed using many common sampling schemes such as temperature sampling, top-k sampling (Fan et al., 2018) or nucleus sampling (Holtzman et al., 2020). Kurakin et al. (2023) ablated various sampling approaches and advocated natural sampling, i.e. temperature sampling with $T = 1.0$ and no top-k or top-p restrictions, however Pillutla et al. (2021) argued that nucleus sampling results in higher quality text than natural (or ancestral) sampling.

Size of the foundational model In general, using larger and more capable models for DP training is expected to improve the fidelity and utility of DP synthetic data at the cost of increased computation, for example Yu et al. (2024) and Kurakin et al. (2023) used 7 and 8 billion parameter LLMs respectively. Tan et al. (2025) use significantly smaller LLMs (140M) by employing several stage process and publicly pretrained models. In particular, a universal topic model is first pre-trained on large scale public corpora. Given a new private dataset, a topic model produces a topic assignment for each private document, as well as topic distribution histogram (with DP) at the dataset level. Topics are then used as conditional prefixes for DP-finetuning a small LLM, with private documents as targets. During sampling, DP topic histogram is used to generate appropriate proportion of samples from each topic. This process produces DP synthetic text data comparable or slightly better than data from DP-finetuning of a much larger model (1.5B), especially under tight privacy budgets. This work also illustrates that for smaller backbone LLMs, DP finetuning with conditional prefixes, for example derived from clustering private data, makes learning the distribution and subsequent sampling of synthetic data ‘easier’.

5.3 Methods that avoid DP-Training

DP finetuning yields high-quality synthetic data but requires finetuning a pretrained LLM. This means that one requires: (1) access to model weights; (2) engineering effort to integrate an efficient DP-SGD implementation into existing (possibly distributed) training pipelines; and (3) sufficient computational resources to run DP-SGD on pretrained LLMs, which is typically significantly more expensive than non-private training.

In this section, we discuss *training-free* methods that can operate with only inference access to a pretrained LLM. The two main classes of methods we discuss are (1) *DP inference* (Section 5.3.1), which requires access to model prediction probabilities (*logits access*); and (2) *Private evolution* (Section 5.3.2), which only requires output text (*API access*).

5.3.1 DP inference

DP inference techniques for DP synthetic text generation use LLMs’ in-context learning ability and incorporate differential privacy into the LLM’s decoding process. These techniques require the ability to query for next-token probabilities from pretrained models prompted with private data, which are then aggregated and released with differential privacy. DP inference belongs to a class of techniques called *Private Prediction*, that we briefly describe first.

Private prediction. Recall that producing DP synthetic data via DP finetuning involves learning a generative model of the data and privatizing the release of the *model gradients (and therefore weights) itself*; from such privatized model, arbitrary amounts of synthetic data can be sampled. Hence for generating synthetic data, privatizing model release is *sufficient but not necessary* – instead one can aim to guarantee privacy of the sampled text by privatizing non-private model outputs directly. This is a well-studied problem in the DP literature referred to as *private prediction* (Dwork and Feldman, 2018).

Algorithm 5 DP Inference

Denote the vocabulary by \mathcal{X} . Let \mathcal{X}^* be the set of all strings over \mathcal{X} , and $\Delta(\mathcal{X})$ be the set of all probability distributions over \mathcal{X} . For $a, b \in \mathcal{X}^*$, ab denotes concatenation. $\text{DPTokenSelect}(\cdot)$ stands in for any DP mechanism used to select the next token.

Require: LLM: $\mathcal{X}^* \rightarrow \Delta(\mathcal{X})$, private prompts $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)} \in \mathcal{X}^*$, response length T .

Ensure: Response $s \in \mathcal{X}^*$.

```

1:  $s \leftarrow \emptyset$ 
2: for  $t = 1$  to  $T$  do
3:   for  $i = 1$  to  $n$  do
4:      $p_t^{(i)} \leftarrow \text{LLM}(\mathbf{x}^{(i)} s)$ 
5:    $\tilde{x}_t \leftarrow \text{DPTokenSelect}(p_t^{(1)}, \dots, p_t^{(n)})$ 
6:    $s \leftarrow s\tilde{x}_t$ 
7: return  $s$ 
```

Since it is a relaxation in the space of admissible algorithms, in theory, private prediction can improve the privacy-utility tradeoff (Ponomareva et al., 2023). The core problem in introducing DP with private prediction is to ensure that model predictions do not change significantly on neighbouring datasets. While some work obtains bounds for convex losses (Dwork and Feldman, 2018; van der Maaten and Hannun, 2020), such techniques have not been successfully extended to LLMs. Instead, to achieve private prediction for LLMs, variants of *subsample-and-aggregate* (Nissim et al., 2007) framework are used.

The *PATE (Private Aggregation of Teacher Ensembles)* line of work is one example of sample-and-aggregate framework (Papernot et al., 2017, 2018). In its classical form, PATE partitions the private data into several disjoint sets and trains separate models non-privately on each set. At inference time, all models are invoked and their outputs are aggregated and released with DP. Crucial to the analysis is the fact that each example can affect only a single model from the ensemble.

The first studies marrying private prediction and LLMs followed the PATE paradigm. Ginart et al. (2022) introduced *SubMix*, which, like PATE, non-privately finetunes several GPT-2 models on disjoint partitions of the private dataset and runs them in parallel for inference. SubMix employs a private prediction protocol that is better-suited for language modeling than PATE, and further makes use of public model predictions. *PMixED* (Flemings et al., 2024b) opts to use LoRA for per-partition models for better inference efficiency (LoRA allows all models to fit in HBM simultaneously; further gains are possible via LoRA batching (Wen and Chaudhuri, 2024)). They also introduce a much simpler prediction protocol that offers standard, unconditional DP guarantees.

DP inference. Generating text with PATE-style private prediction requires training a large number of models (>30), and having them all loaded simultaneously for inference. Also, each split should be reasonably large to ensure there is sufficient data to train a model. To address these issues, the next advance in private prediction for LLMs was to employ private data for *in-context learning*. We will refer to this class of techniques as *DP inference*, although the term *differentially private in-context learning* is also employed (Wu et al., 2024b). Wu et al. (2024b) propose to partition private examples into disjoint contexts to condition inference, rather than training sets for finetuning; PromptPATE (Duan et al., 2023) uses private predictions to label public examples, which are collected in a “student prompt” that can be deployed at test time. Wu et al. (2024b) employ Gaussian ReportNoisyMax (Zhu and Wang, 2022) for token selection. They demonstrate effectiveness on generative tasks like summarization by generating keywords with DP inference followed by using a public model to post-process results.

At its core, DP inference techniques for LLMs directly prompt a pretrained language model with examples from the private dataset, distributing them among many parallel prompts, with each one asking for similar text to the provided example. When generating each token of the response, the next-token prediction produced by each context is aggregated with DP to produce a differentially private output token. That token is then appended to each of the parallel contexts and the process repeats. Algorithm 5 gives pseudocode for the process.

Privacy cost of DP Inference

As opposed to DP finetuning, *the privacy cost of DP inference increases with the number of tokens generated*. This is because the privacy analysis typically proceeds by bounding the privacy cost of generating a single token, and then accumulating these costs via composition.

Synthetic data generation. As shown in Algorithm 5, generating synthetic data with DP inference simply requires feeding in sampled tokens in an auto-regressive manner. Tang et al. (2024) first demonstrate generating *few-shot synthetic data* with DP inference for classification and extraction tasks. They use the Gaussian mechanism with subsampling for token selection, and demonstrate generating ≈ 12 examples, which suffices for effective few-shot learning. Amin et al. (2024) introduce algorithmic improvements (clipped logits exponential mechanism for token sampling, as well as parallel composition) to scale generation to \approx thousands of examples. Increased quantity enables one to finetune downstream BERT models on generated synthetic data; quantity also improves downstream in-context learning accuracy. Gao et al. (2025) turn their attention to adaptively selecting the clipping radius for noise addition, noting that for certain tokens, predictions are highly similar and therefore less noise can be added; they demonstrate utility improvements for few-shot generation. Flemings et al. (2025) introduce a PMixED-inspired prediction protocol and show improvements for in-context few-shot generation when targeting more difficult generative tasks. Hong et al. (2024) skips synthetic data points as an intermediary for in-context learning and directly learns a classification prompt (which may include some synthetic data); they sample tokens with the LimitedDomain mechanism (Durfee and Rogers, 2019).

Synthetic data should capture facets of your source data beyond task accuracy. Cohen et al. (2023) introduce *HotPATE*, a drop-in replacement for PATE which aims to preserve the diversity of the source data during generation. PATE targeted classification tasks and hence did not consider this a design requirement. The key technical tool of Cohen et al. (2023) is *coordinated sampling*: rather than having all teachers sample tokens independently before aggregation, HotPATE uses shared public randomness when sampling to maximize agreement between teachers; this lets it better preserve the tail of the data distribution. Amin et al. (2025) compute representativeness metrics (i.e. MAUVE (Pillutla et al., 2021)) on data generated from DP inference, and find it is poor. Indeed, most prior work generated a small number of examples for in-context learning and could not compute dataset-level metrics for representativeness. Amin et al. (2025) show that clustering data and forming same-cluster batches improves representativeness.

Leveraging public information. A theme explored in several works is to use predictions made by a public model to improve the privacy/utility tradeoff of DP inference. The conceptual appeal is the following: in text, many tokens are predictable without private information, and therefore we should not “pay” privacy cost for them. The privacy cost of text generation should not scale with the *number of tokens produced*, but rather the *surplus information content* introduced by the private data.

Flemings et al. (2024b), Ginart et al. (2022), and Flemings et al. (2024a) mix predictions made from the public model with the private model. For the latter two, which employ data-dependent analyses, public predictions aligning with private predictions result in charging lower ϵ . Tang et al. (2024) filters token selection candidates to the top- k choices of a public model; this reduces the frequency of noise-induced error when selecting tokens from a large vocabulary, but introduces bias. Amin et al. (2024) and Koga et al. (2024) use AboveThreshold Dwork and Roth (2014) to query whether the public prediction aligns with the private prediction, and *only consumes privacy budget to use the private prediction if the check fails*; hence the privacy cost scales with the number of tokens where public and private predictions differ. InvisibleInk (Vinod et al., 2025) employs the exponential mechanism for token sampling and uses the public prediction to determine where to center the clipping range. This modification allows for a significant reduction of privacy budget, unlocking generation with a much smaller batch size. Amin et al. (2025) employs clusters derived from public data to batch together similar private data when running inference, which improves representativeness.

Data-dependent analysis. Papernot et al. (2017) introduced a *data-dependent* privacy analysis in which the privacy guarantee is a function of the input dataset, instead of a worst-case upper bound over all datasets. Accepting this relaxation can lead to significant gains in data volume. Roughly speaking, these analyses allow one to take advantage of *consensus* among predictions. In particular, on tokens where many teachers agree, one can argue the counterfactual impact of any user’s data is small; however this scenario cannot be assumed to hold in worst-case datasets considered in unconditional DP analyses. Duan et al. (2023) inherits PATE’s analysis and provide such a guarantee; they report $\epsilon = 0.147$. SubMix (Ginart et al., 2022) presents a data-dependent analysis. Coordinated sampling in HotPATE (Cohen et al., 2023) improves teacher agreement which results in tighter data-dependent guarantees. AdaPMixED (Flemings et al., 2024a) computes data-dependent ϵ and introduces a screening step to defer expensive queries to public model. They demonstrate $16\times$ smaller ϵ while maintaining similar utility as the unconditional DP variant (Flemings et al., 2024b). Amin et al. (2025) introduce a coordinate-wise median-based token selection algorithm designed to take advantage of prediction consensus, and prove a data-dependent and *ex-post* DP guarantee (Ligett et al., 2017).

User-level guarantees. To support user-level guarantees, the input to the base DP token selection mechanism (before applying subsampling and composition) must satisfy that *a user’s data is restricted to appearing in a single model context*. A user’s data appearing in multiple contexts necessitates the application of group privacy bounds.

Comparison of methods. Table 8 presents a comparison Private prediction, including DP inference methods, along several main axes: partitioning strategy, token selection mechanism, task, and the privacy guarantee.

- A simplified version of the method described in [Amin et al. \(2024\)](#) offers the best privacy-utility trade-offs, as it is capable of generating thousands of examples at reasonable privacy budget. The simplified version notably employs parallel composition, KV cache reuse, and exponential mechanism sampling; but does not implement the use of AboveThreshold with an additional public query. This method is recommended because it is relatively simple, is capable of generating a large quantity of data, and offers the standard unconditional DP guarantee. Results for this approach are reported as *Baseline++* in [Amin et al. \(2025\)](#).
- Extensions of this approach to adopt the public query logit clipping proposed by [Vinod et al. \(2025\)](#) results in further improvements in privacy budget and compute efficiency.
- For the above methods, the most straightforward choice of the model for synthetic data generation is a pretrained model (instead of instruction-tuned) with a generic prompt that puts sensitive examples directly into context with no additional description. This minimal approach performs well (see *Baseline++* results in [Amin et al. \(2025\)](#)) and requires no prompt engineering. Including at least 2 serial examples per context improves results ([Tang et al., 2024](#); [Amin et al., 2025](#)).

Method Paper	Partition	Mechanism	Task	Privacy	Approach
Ginart et al. (2022)	Train sets	SubMix	Next-token prediction	Data-dependent	Privacy by sampling
Wu et al. (2024b)	Contexts	Gaussian Noisy Max	Query answering, generative response	Standard	Partitions examples into contexts
Duan et al. (2023)	Contexts	PATE	Query answering	Data-dependent	Label public prompts & deploy
Cohen et al. (2023)	Contexts	HotPATE	Query answering	Data-dependent	Coordinated sampling to improve PATE diversity
Flemings et al. (2024b)	Train set	PMixED	Next-token prediction	Standard	Mix public distribution with private, then sample
Flemings et al. (2024a)	Train set	AdaPMixED	Next-token prediction	Data-dependent	Data-dependent extension of PMixED
Tang et al. (2024)	Contexts	Gaussian Noisy Max	Few shot synthetic data	Standard	Generate a few synthetic examples to be used for ICL
Hong et al. (2024)	Contexts	Limited-Domain	Prompt generation	Standard	Generate a prompt directly instead of few-shot examples
Amin et al. (2024)	Contexts	Exponential	Synthetic data	Standard	Privacy by sampling and parallel composition to generate much more data
Gao et al. (2025)	Contexts	GoodRadius + Gaussian	Few shot synthetic data	Standard	Pick per-token adaptive clipping radius
Flemings et al. (2025)	Contexts	PMixED	Few shot synthetic data	Standard	Mix public distribution with private, then sample
Amin et al. (2025)	Contexts	Median	Synthetic data	Ex-post data-dependent	Form batches from pre-clustering data to improve representativeness

Vinod et al. (2025)	Contexts	Exponential	Synthetic data	Standard	Recenter clipping range with public prediction for better privacy
---------------------	----------	-------------	----------------	----------	---

Table 8: A comparison of DP inference methods in the literature. We recommend a simplified version of the method of Amin et al. (2024) (omitting the use of public predictions) as a simple, strong baseline method for generating large quantities of synthetic data via DP inference.

5.3.2 Private evolution for Text

Algorithm 6 Private Evolution for Text (Aug-PE, Xie et al. (2024))

Require: Private data D , text embedding model Φ , target number of synthetic samples N , evolution rounds T , population size multiplier L (number of rewrites per chosen synthetic sample).

Ensure: Synthetic dataset \hat{D} .

```

1: Initialize  $\hat{D}_0$  of size  $(L + 1) \cdot N$  with Random API.
2: for  $t = 0, \dots, T - 1$  do
3:    $E_t = \Phi(\hat{D}_t)$  //Embedding calculation for synthetic samples.
4:   Let  $\Phi(p), p \in D$  vote for the nearest embedding in  $E_t$ .
5:   Privatize the voting results with  $(\epsilon, \delta)$ -DP to get a DP histogram  $H_t$ .
6:    $\hat{H}_t = H_t / \text{sum}(H_t)$  //Histogram normalization.
7:   Get  $\hat{D}'_t$ : top- $N$  samples from  $\hat{D}_t$  based on  $\hat{H}_t$ .
8:   if  $t < T - 1$  then
9:      $\hat{D}_{t+1}$  of size  $(L + 1) \cdot N$ :  $L$  variants for each  $z \in \hat{D}'_t$  via Variation API and  $\hat{D}'_t$ .
10:  else
11:    return  $\hat{D}'_t$ .
```

Private Evolution (PE) is an emerging approach for generating DP synthetic data using API-only (e.g. inference) access to foundation models. Its key advantage lies in leveraging the powerful capabilities of these models without requiring the design of DP-specific training or sampling algorithms, an often technically challenging or even infeasible task, particularly when the models are closed-source.

Originally introduced for DP image synthesis (Lin et al., 2023), PE was later extended to text by Xie et al. (2024) who proposed Augmented Private Evolution (Aug-PE). Aug-PE enhances PE with techniques specifically tailored to text generation, improving the quality and relevance of synthetic outputs. Despite domain differences, both PE and Aug-PE follow the same high-level pipeline: they start with an initial population of synthetic samples generated via prompting (this pool can also be public data of similar to sensitive data distribution), then iteratively refines these samples using a DP histogram-based evolutionary process that gradually aligns the distribution with the private dataset. We provide a pseudocode of the Aug-PE algorithm from Xie et al. (2024) in Algorithm 6.

Before detailing Aug-PE’s workflow, we highlight two key APIs required to be built using the target foundation model: (1) **Random API**: Generating an initial population of synthetic samples without access to private data. Alternatively, access to some public data of similar distribution to the sensitive data and (2) **Variation API**: Modifies/rewrites existing synthetic samples to introduce diversity. With these two functionalities in place, the Aug-PE process proceeds as follows.

Initialization. The process begins by generating an initial pool of synthetic samples, typically through prompting the foundation model. Although Aug-PE can technically start with arbitrary synthetic samples, empirical results show that using well-designed prompts informed by public priors leads to more relevant and effective outputs. These priors help guide the model toward producing samples that better reflect the target data domain. For example, in generating differentially private synthetic Yelp reviews Xie et al. (2024) used prompts with contextually appropriate phrases such as “reviews for steakhouse restaurants”—drawn exclusively from public sources. Importantly, these prompts must not incorporate any private data to ensure privacy is preserved. Given a target of N synthetic samples, Aug-PE generates an initial population of size $L \times N$, where L is a small constant (e.g. single or low double digit). This differs from the original PE for images, where the initial and target population sizes are identical.

Iterative Refinement. Once the initial population is constructed, Aug-PE refines it through repeated variation and voting steps. During variation, each sample is modified using one of two strategies: prompting the model to rephrase the text or masking random tokens and asking the model to fill in the blanks. In the voting step, an embedding is computed for each synthetic candidate sample using a text embedding model—either directly from the synthetic datapoint itself or by first generating a few variants via the Variation API and aggregating their embeddings. Each private data point then votes for its nearest candidate based on embedding space distance, producing a histogram of votes.

This histogram is privatized with differential privacy, and the samples with the top N vote counts are selected to form the next generation (an alternative strategy is to sample the synthetic datapoints with replacement according to probabilities induced by the histogram of votes, similar to what was done for images in Section 4.4, which can however result in the same synthetic samples being chosen multiple times; this isn’t desirable in this version of the algorithms, because each selected example is kept in the pool as-is). If $L > 1$, these N samples are expanded back to $(L + 1) \cdot N$ using the Variation API, rewriting each chosen synthetic sample L times. Note that this is different from PE for Images (Section 4.4), where only variants (but not the N synthetic examples themselves) were added to the pool for the next round. This is done to increase the chance of retaining high quality synthetic examples Xie et al. (2024).

This refinement process is repeated for a fixed number of iterations. At the final iteration, the top N samples are directly returned as the final synthetic dataset.

User-level privacy unit adjustments In the original algorithm, each private example uses only 1 vote (that can be either fully allocated to a closest synthetic example or distributed between a number of closest synthetic examples, depending on hyperparameter *number of nearest neighbours*). This means that l2 sensitivity is 1 and gaussian noise introduced to the histogram depends on this sensitivity. When user-level privacy is required and each user contributes multiple samples, the vote counts from each user can be normalized to sum to one (in l2 norm). This ensures that the sensitivity per user remains bounded by one.

Discussion The benefits of PE lie in its simplicity and practicality: it is an intuitive algorithm that is easy to explain and that relies solely on inference using off-the-shelf LLMs. This eliminates the need for computationally expensive DP fine-tuning of LLM checkpoints, does not require access to logits and allows to potentially use a more powerful model that can be effectively DP-finetuned.

However, the effectiveness of PE hinges on two key components: the initial pool of synthetic data and the rewrite prompt, both of which must be manually crafted by the user. In the absence of publicly available data to seed the initial pool, the authors suggest manually creating a prompt to generate data that mimics the private dataset. Yet constructing either the initial or rewrite prompts requires some public prior knowledge about the private data, which may be difficult to obtain in some settings. For instance, the prompts proposed in (Xie et al., 2024) reflect a deep understanding of the private data’s distribution, such as Yelp review categories, rating scales, or the typical format of research papers.

This raises a critical question: how can high-quality prompts be designed for sensitive datasets when no human-readable samples are available? It remains unclear whether structural insights, like categories or data schemas, can be derived in a privacy-preserving way without human inspection. Furthermore, since PE does not involve fine-tuning the LLM’s weights, it assumes that the private data’s general domain is already covered in the LLM’s training corpus. As a result, PE is unlikely to perform well in domains absent from pretraining. For instance, generating coherent legal text if legal language was not part of the original training data will be a challenging task for PE.

Recent approach for PE introduced by Hou et al. (2025) introduces updates to the foundational model via policy optimization tuning. The key idea is that voting information can be used to guide the model to prefer synthetic examples that obtain better scores than synthetic examples generated with the same prompts but receiving worse scores. With this approach the foundational model however is no longer fixed and access to model weights and tuning of the model and to compute to perform the updates is required.

5.4 Summary and comparison of methods

We have outlined 3 main methods for creating DP synthetic text data: via DP finetuning (Section 5.2); DP inference (Section 5.3.1); or private evolution algorithms (Section 5.3.2). We provide a detailed comparison of the methods in this section, highlighting similarities and differences.

5.4.1 Unified view of DP text synthesis methods

First, we present a unified view of the discussed methods. Existing methods can be viewed under the following recipe:

- Start from an existing non-private approach to generate synthetic text.

- Identify the *iterative primitive* in the approach used to *extract information from real data*.
- DP-fy that primitive.

Specific examples are presented in Table 9. The primitive used to extract information from real data must do so via aggregates (e.g. batch gradient) which allows it to be amenable to DP. The iterative nature of the primitive means that proving privacy guarantees of the primitive suffices to guarantee privacy of the entire algorithm.

	DP finetuning	DP inference	Private evolution
Comparison	Finetuning	Inference	Evolutionary prompting
Primitive	Gradient	Token sampling	Selection

Table 9: For each DP text synthesis method, we can identify a base non-private method, as well as the core iterative primitive being DP-ified.

Unified view of DP text synthesis methods

Existing approaches to DP text synthesis start from a non-private text generation procedure and identify an iterative primitive used to extract information from real data to DP-fy. The effectiveness of the method is predicted by its *original effectiveness* plus the degree of *distortion introduced by the DP substitute of the primitive*.

5.4.2 Comparison

Each of the aforementioned methods have their strengths and weaknesses and different requirements. Table 10 provides a detailed comparison of these methods along various axis like amount data, computer resources needed, yield, type of access needed etc.

Choosing the method for DP synthetic text data generation

For large enough volume of sensitive data (>XX K datapoints) with enough compute the method that results in highest fidelity and utility of DP synthetic data is almost always DP-Finetuning. Less computationally expensive methods that don't require LLM's finetuning can provide reasonable data when sensitive data is somewhat in distribution for the pretraining data.

- **DP finetuning** is the workhorse method that delivers the best quality given sufficient amount data, compute and engineering and time investment. **Since the adoption of DP synthetic data is often *quality-bottlenecked*, this is likely the option that best fits your needs.** It does require significant engineering expertise to implement and run DP-finetuning and edit access to model checkpoint weights, in turn it delivers the best quality when given access to a significant amount of data. DP-finetuning results in high fidelity synthetic data even if private data to mimic is significantly out of the distribution for LLM's pretraining data.
- **DP inference** is suitable for fast prototyping but is limited to generating short pieces of text only, due to its privacy cost accumulating with each output token.
- **Private evolution** is a strong candidate when one wants very stringent privacy guarantees ($\epsilon < 1$ and/or has access to very limited data. In fact, that it the only method that can work with very few datapoints (e.g. <1K). It is also extremely easy to explain and does not egress sensitive data to an LLM for inference (only using sensitive data's embeddings). It does require access to high quality public data similar to the sensitive data and significant work for prompt engineering.

Aspect	DP finetuning	Method DP inference	Private evolution
Amount of input private data			
<i>Small input quantity (<5K)</i>	Not recommended	Not recommended	Preferred
<i>Large input quantity (>10K)</i>	Preferred	Not recommended	Not recommended
Yield			

	Unlimited number of output examples, although with diminishing returns to downstream task performance.	$\approx 1K$ input private examples $\rightarrow \approx 25$ synthetic examples out (200 tokens per example at $\epsilon = 10$).	In practice, suitable for outputting synthetic dataset of size \leq size of input private dataset.
Model access required	Weights	Per-token prediction logits	Generations via API
Compute resources and engineering effort			
<i>Training of LLM is required</i>	Yes	No	No
<i>Inference cost multiple per synthetic example</i>	1 (same as regular inference).	\propto (number of contexts to aggregate over).	\propto (number of iterations) \times (number of rewrites).
<i>Prompt engineering required</i>	No	No	Yes – craft prompts for initial pure synthetic data and rewrite templates.
<i>Time to first example</i>	Long (Run finetuning on the entire dataset, then sample)	Short (Batch-by-batch inference)	Medium (Requires prompt engineering + running rewriting and embedding on the entire dataset)
<i>Direct side by side comparison of inputs and outputs</i>	No	Yes (input batch vs. generated output)	No
<i>Resilience to distribution gaps between private data and LLM</i>	High	Medium	Low
Target privacy guarantee			
<i>High ϵ (e.g. 10)/large amount of private data</i>	Preferred	Not recommended	Not recommended
<i>Stringent (e.g. $\epsilon < 1$)/small amount of private data</i>	Not recommended (quality will be bad)	Not recommended (short synthetic outputs and little data)	Preferred
Data persistence requirements			
	Entire dataset required at once for training.	A batch of data required to generate synthetic examples for that batch.	Single examples can arrive in a streaming fashion, be used to cast votes, and then discarded immediately.

Table 10: Comparison of DP finetuning, DP inference, and private evolution for DP text synthesis.

5.5 Resampling DP synthetic text

There are situations where it is beneficial to resample a subset from an initial differentially private (DP) synthetic dataset. This is often due to: (1) a distributional gap between the DP data and the original source data, caused by factors such as noise perturbation, clipping bias in DP algorithms, or suboptimal DP training hyperparameters; or (2) the need to focus on specific subsets of the data, such as synthetic text with a particular sentiment or length, even when the overall DP data distribution aligns well with the original training data. Additionally, one might wish to filter out noisy signals originating from the training corpus. These challenges have driven the development of resampling techniques for DP synthetic data, especially in the tabular data domain (Neunhoeffler et al., 2021; Liu et al., 2021c). For example, Wang et al. (2023) proposed a method that privatizes the correlation matrix of an initial DP dataset and then resamples synthetic data to better match the underlying distribution of the source data.

Xie et al. (2024); Yu et al. (2024) extend the idea of resampling to the text domain using private voting idea similar to PE: first it constructs a histogram to represent the target data distribution, then resamples the initial DP data according to this histogram. Xie et al. (2024) construct the histogram by assigning each bin to an initial synthetic sample, with each private sample voting for the nearest synthetic sample in the embedding space. The histogram is released with DP guarantees and subsequently used to guide the resampling process. It is noteworthy that when the number of private samples is small, or when the size of the dataset targeted for filtering is large, releasing the histogram under DP can result in a poor signal-to-noise ratio.

To address this, Yu et al. (2024) propose clustering the initial DP synthetic samples in the embedding space, so that each bin represents a cluster of synthetic samples. This aggregation unites votes from private samples, thereby improving the signal-to-noise ratio in the DP histogram. Synthetic examples are then sampled from synthetic clusters that received

highest number of private votes, and the number of synthetic samples is proportional to the number of votes a cluster received.

It is important to remember that such postprocessing uses the sensitive data and requires privacy budget allocation. When combining with core technique like DP finetuning using the same sensitive data for both to improve fidelity, it is recommended to allocate the majority budget to DP finetuning and use the rest for DP postprocessing, e.g. for a target $\epsilon = 10$, δ DP finetuning can be used to create DP synthetic pool of data with $\epsilon = 9$, $\delta/2$ and postprocessing can use remaining $\epsilon = 1$, $\delta/2$ budget to choose the synthetic samples that most closely resemble the sensitive data.

5.6 Evaluating synthetic text data quality

There are several fidelity metrics that can be used for comparing real and synthetic text data.

Statistics-based metrics. Tried-and-tested *statistics-based metrics*, which compare statistics on real and synthetic data, can be useful for detecting problems in data preprocessing or training. Relevant statistics are often selected in an ad-hoc fashion, depending on the task at hand. For example, identifying significant difference of distributions of sequence lengths on real and synthetic data can signal context length or sampling problems (Kurakin et al., 2023). Kurakin et al. (2023) use n-grams statistics like unigram/bigram histograms and compare the area under their divergence frontiers (Sajjadi et al., 2018), which is conceptually similar to MAUVE (Pillutla et al., 2021). Additionally, Kurakin et al. (2023) find that the ordering of model candidates induced by test-set perplexity is highly correlated with downstream model performance and is even comparable with correlation of leading metrics like MAUVE.

MAUVE. MAUVE (Pillutla et al., 2021) is a recent and widely adopted metric that compares two text distributions. The idea behind this metric is to quantify the tradeoffs between Type I errors (areas of probability distribution in synthetic data that is unlikely under real data) and Type II errors (lack of synthetic data in the regions that are plausible for the real data).

A practical algorithm for calculating MAUVE score first embeds real and synthetic data and clusters into k clusters (using k-means) jointly. An optional step before clustering can involve doing PCA on the embeddings. Cluster assignments for real and synthetic data are counted, resulting in two histograms of support size k , which form the representation of real and synthetic distributions. KL divergence between two representations can quantify type I and type II errors, however it is ∞ when supports of histograms differ. Therefore a divergence curve that softly (by varying mixing weight) measures KL divergence of synthetic vs real and real vs synthetic data using the aforementioned representation is built. Each point on the curve quantifies tradeoffs between Type I and II errors, and the area under the curve is the reported MAUVE metric Pillutla et al. (2021).

MAUVE exhibits behavior that is in line with expected human judgment. For example, larger models in general produce synthetic text that is higher quality and more similar to human created text, as judged by the human raters, and MAUVE values increase with model size increases. Similarly, MAUVE values are correlated with human judgments when comparing various decoding algorithms. MAUVE also correctly captures the fact that generating longer synthetic text is harder, as text becomes incoherent and inconsistent the longer the model decodes for.

Practically speaking, the higher the MAUVE value, the more similar two measured text distributions. While 1.0 is the upper limit of MAUVE, measuring MAUVE from samples from the same distribution will not necessarily result in the value of 1.0, as MAUVE scores depend on the sample size. MAUVE should be first measured on two disjoint subsets of the real data (*real vs. real*) to calculate a practical upper bound value that synthetic data will be measured against. Calculating *synthetic vs. real* MAUVE should be done with the same sample sizes. In general, papers reporting MAUVE should indicate sample sizes and MAUVE algorithm hyperparameters used (e.g. the number of clusters k , whether PCA was used on the embeddings, etc.) to allow for fair comparison. While values of MAUVE on different sample sizes are incomparable, the ranking induced by MAUVE scores is consistent. Because of randomness induced by k-means and sampling of data to calculate MAUVE, we suggest to run several trials of MAUVE calculation with different random samples of real and synthetic data (with a fixed sample size) and report mean and standard error.

Additionally, MAUVE strongly depends on the quality of the embeddings, with 'stronger' embeddings increasing the gap between *real vs. real* and *synthetic vs. real* values. Ideally, the embedding used captures the information that is essential for the downstream use of the synthetic data. This also somewhat blurs the line between MAUVE being a fidelity metric or utility metric. In practice it is rare to have access to an embedding from the downstream model at the stage of synthetic data generation, so encoders like RoBERTa (Liu et al., 2019), T5 (Raffel et al., 2020) or task-specific Gecko embeddings (Lee et al., 2024) can be used.

Finally, it is worth pointing out that MAUVE is a generic metric that can work with different modalities given an appropriate embedding model, and is now also commonly used for image data fidelity evaluation.

5.7 Open questions and challenges

Despite recent advances in DP synthetic text generation, several challenges remain in the field, including but not limited to the following.

Generation of long DP synthetic text. Most existing work focuses on generating synthetic text up to a few hundred tokens (Yue et al., 2022; Kurakin et al., 2023; Yu et al., 2024; Xie et al., 2024), which is significantly shorter than the context window length of modern LLMs. Yue et al. (2022) find that the length distribution of DP-generated text tends to be shorter than that of real data, likely due to the disparate impact of differential privacy Bagdasaryan et al. (2019). Additionally, it is known that in open-ended text generation, data becomes less coherent and consistent the longer the output is. This suggests that generating long DP synthetic text may pose additional challenges and might require different modelling techniques, for example breaking up large piece of text into overlapping chunks and finetuning (and subsequently sampling) such smaller chunks.

DP synthetic text without pre-trained LLMs. Recent advances in DP text generation heavily rely on pre-trained foundational LLMs, which are trained on broad and often uncured corpora. As a result, the DP guarantees typically apply only to the target dataset used for generating synthetic data, not to the pre-training data of the foundational model itself. Although using foundation models has become standard practice, there are concerns about potential privacy risks stemming from the pre-training data (Bommasani et al., 2021; Tramèr et al., 2022). Developing methods for DP text generation that do not depend on powerful pre-trained LLMs remains an open research challenge.

Challenges specific to main methods For DP-finetuning, the current state-of-the-art utilizes standard pretrained or post-trained checkpoint. However pretraining checkpoints are trained with teacher forcing next-token loss, encouraging the model to learn to generate the exact sequence of training data. This objective is arguably more suited for post-training tasks like summarization or translation but less so for open-ended text generation that does not require exact matching of the target text, rather a generation of a text that is similar but not exactly the same. Additionally, essentially the same loss is used for DP-finetuning. Research into best ways of pretraining (and DP-finetuning) LLMs to maximize the utility and fidelity of open ended text generation is clearly needed. Perhaps training objectives similar to RL can be used to choose better generated examples, instead of enforcing the generation of the exact training example, as it is done currently.

DP finetuning of models like Diffusion have been used for text generation (Li et al., 2022) and have demonstrated impressive results in non DP setting. Attempting such models with DP for synthetic data generation is a promising direction to explore.

For private evolution, finding a way to "variate" text that does not require extensive prompt engineering can greatly expand the adoption and improve algorithm's utility.

DP inference methods have still a long way to go before effective generation of long text is achieved, which will require significant advancing in accounting such that more tokens can be generated for the same privacy budget.

DP Synthetic data generation in online manner: Online synthetic data generation (e.g. in a streaming fashion when private data is available on instance-by-instance basis) is an underexplored area of DP synthetic text generation. Both PE and private inference methods are potentially suitable for this application, however to the best of our knowledge, such setups has not been yet been explored in the literature.

6 DP Synthetic Data in Federated Learning

So far, when discussing DP synthetic data, we focused on centralized setting where a trusted central server accesses a privacy-sensitive source dataset and produces DP synthetic version of this collection. In this section we discuss privacy-preserving synthetic data generation from sources decentralized across a distributed system. User data frequently originates on edge devices (e.g. user phones or internal company storage), where local storage offers users greater control. Additionally, in sensitive domains like finance and healthcare, privacy and business considerations may hinder data sharing among organizations. While this decentralization introduces challenges compared to centralized data storage with a trusted provider, it also creates opportunities for distributing computational costs. In this section, we specifically focus on the challenges and algorithms for generating synthetic data within the federated learning (FL) framework (McMahan et al., 2017).

Roadmap This section is organized as follows. In Section 6.1 we give a brief overview of FL in context of ML. Next, we discuss the interplay between FL and DP in Section 6.2, and DP training algorithms in Section 6.3. Section

6.4 describes what DP synthetic data means in context of FL and outlines unique challenges one encounters when attempting to generate DP synthetic data. We briefly survey methods for FL DP synthetic data generation in Section 6.5 and further compare and discuss them in Section 6.7. Section 6.6 outlines synthetic data evaluation peculiarities in context of FL.

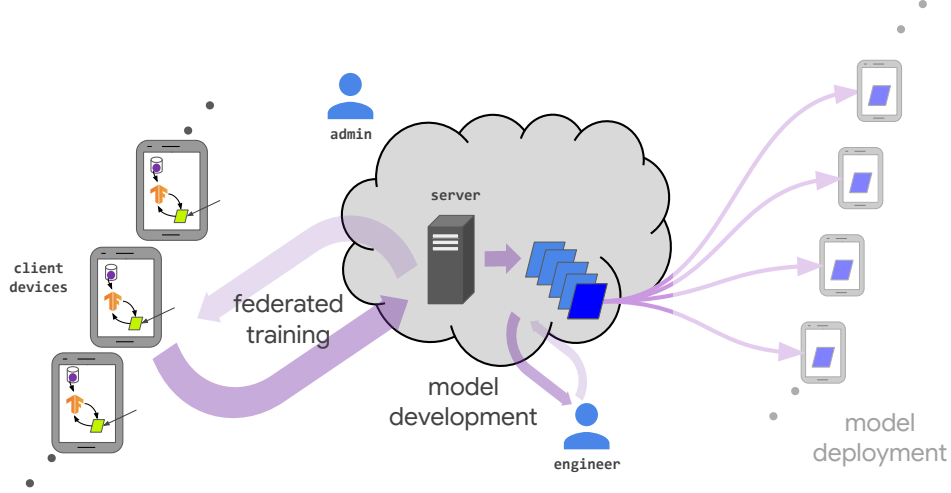


Figure 4: A typical cross-device federated learning system for mobile applications. The figure is adapted from (Kairouz et al., 2021b; Zhang et al., 2023) with permission. Private user data are decentralized on client devices, and focused updates are ephemerally aggregated in line with the data minimization principle (Bonawitz et al., 2022), an approach that can be further enhanced with techniques like encryption and differential privacy to achieve a multi-faceted model of privacy.

6.1 Federated Learning Primer

Federated learning is a machine learning paradigm for decentralized data motivated by privacy protection (McMahan et al., 2017; Kairouz et al., 2021b; Bonawitz et al., 2022). Daly et al. (2024) proposed this updated definition: **Federated learning (FL)** is a machine learning setting where multiple entities (clients) collaborate in solving a machine learning problem, under the coordination of a service provider. A complete FL system should enable clients to maintain full control over their data, the set of workloads allowed to access their data, and the anonymization properties of those workloads. FL systems should provide appropriate transparency and control to the users whose data is managed by FL clients. FL is built on the privacy principle of data minimization, as the immediate aggregation of only focused updates that improve a specific model aims to minimize the exposure of data. Clients, who own and control their data, can range from mobile devices in cross-device FL to organizations with richer resources in cross-silo FL. The server, in turn, must orchestrate the machine learning process in a trusted manner. A typical machine learning process in FL is to train a task-specific model with decentralized data, and the final trained model is deployed to the clients for application use. Figure 4 illustrates a typical federated learning system where data are decentralized on mobile devices.

While powerful, the FL paradigm presents several well-studied challenges. One of them is achieving the communication and computational efficiency in the distributed system with decentralized data. Additionally, the heterogeneous data on clients makes optimization of an ML objective on the server hard to converge, making training hard. Finally, privacy and security considerations are much more complex and nuanced in a complex system like FL.

Among the many algorithms for training ML models in FL, Federated Averaging (FedAvg) (McMahan et al., 2017) is one of the most popular ones. Algorithm 7 presents both FedAvg algorithm (without highlighted lines) and its DP version (via DP-FTRL, with highlighted lines taken into account). FedAvg algorithm has two training loops: the inner loop of locally updating local models (copy of a global model) using clients data, and the outer loop of applying the updates aggregated from clients to the global model W . FedAvg promotes data minimization and communication efficiency by distributing local updates on clients and infrequently aggregating client updates for the global model. While the aggregation step can incorporate DP and encryption methods to boost privacy in FL, the inherent data heterogeneity across clients can still pose a challenge to the FedAvg training paradigm.

Algorithm 7 A FL algorithm and its DP version (when modifications in grey are applied) used to train production language models in a cross-device FL system, as presented in (McMahan et al., 2024). The DP version adapts FedAvg (McMahan et al., 2017) with DP-FTRL (Kairouz et al., 2021a) to ensure user-level DP. Client-side optimization defaults to SGD to conserve edge device resources, while the server optimizer (ServerOpt) can use standard methods like momentum SGD. Standard DP parameters (ζ, σ) are used; details on the correlated noise function are available in (McMahan et al., 2024).

Require: clients per round m , learning rate for clients η_c and for server side updates η_s , momentum $\beta = 0.9$, total number of rounds T , clipping norm ζ , noise multiplier σ , min-separation b

Ensure: trained DP model (e.g., synthetic data generator) W^n

```

1: Initialize server side model  $W^0$  (e.g. randomly or starting from pretrained weights)
2: Initialize server optimizer state  $\mathcal{P}$ 
3: Initialize correlated noise state  $\mathcal{S}$  with  $\sigma\zeta$ 
4: for each round  $t = 1, 2, \dots, n$  do
5:    $\mathcal{Q}^t \leftarrow$  (at least  $m$  users that did not participate in the previous  $b$  rounds)
6:   for each user  $i \in \mathcal{Q}^t$  in parallel do
7:      $\Delta_i^t \leftarrow \text{ClientUpdate}(i, W^{t-1})$  ▷ Client receives a local copy of global model  $W^{t-1}$ 
8:    $\Delta^t, \mathcal{S} \leftarrow \text{AddCorrNoise}(\mathcal{S}, \sum_{i \in \mathcal{Q}^t} \Delta_i^t)$ 
9:    $W^t, \mathcal{P} \leftarrow \text{ServerOpt}(W^{t-1}, \frac{1}{m} \Delta^t, \eta_s, \beta, \mathcal{P})$ 
10: function CLIENTUPDATE( $i, W_i$ )
11:    $\mathcal{G} \leftarrow$  (batches of user  $i$ 's local data)
12:   for batch  $g \in \mathcal{G}$  do
13:      $W_i \leftarrow W_i - \eta_c \nabla \ell(W_i; g)$ 
14:    $\Delta \leftarrow W_i - W_i^{(0)}$ 
15:    $\Delta' \leftarrow \Delta \cdot \min\left(1, \frac{\zeta}{\|\Delta\|}\right)$ 
16:   return  $\Delta'$ 

```

See (Kairouz et al., 2021b) for a more detailed overview of federated learning. The clear delineation of responsibilities between clients and server in FL, where users maintain full control over their data, is a critical considerations for our subsequent discussion on differential privacy (DP) and synthetic data generation.

6.2 Differential Privacy in Federated Learning

Combining FL with DP provides robust privacy protection by leveraging both data minimization and formal anonymization guarantees (Bonawitz et al., 2022). In the distributed paradigm of FL, a key consideration of DP is whether the server is trusted. **Central DP** is defined on the global learning process and the entire collection of decentralized data. For example, training an ML model or creating a DP synthetic data using the information about all the data from the clients without explicitly storing and pooling all clients data on the server. The server is trusted and can access intermediate results (e.g., aggregated but not noised model updates) in addition to the released results with DP guarantees (e.g., all model checkpoints $\{y^t\}$). **Local DP** is defined on each client's local data and learning process (e.g. each client trains and shares their own ML model or creates and shares a DP synthetic version of their own data). The server does not need to be trusted as the client updates have effective DP guarantees before aggregation. While local DP has the advantage of not requiring a trusted server, it struggles to maintain usable model (or data) utility (Kairouz et al., 2021b). Enhancing the FL system to reduce necessary trust on server with strong privacy-utility trade-off is an active area of research, exploring among other ideas aggregation in trusted execution environments (TEEs) (Daly et al., 2024). In practice, a central DP guarantee is reported, and Figure 5 illustrates how a central DP FL algorithm is deployed in cross-device FL; with TEEs, strong DP guarantees can be achieved without fully trusting the server (Daly et al., 2024).

6.2.1 Privacy unit for FL

Apart from the choice of central or local DP, a second crucial aspect is the choice of the privacy unit, similar to the discussion in centralized setting (Section 2.2.1). **User-level privacy unit** is often a natural fit and practically achievable for cross-device FL in **central DP setting** (or more technically, device-level DP, as we use the device as a proxy user identifier), as data is already siloed by user devices, making each user an intuitive unit for privacy protection. For example Xu and Zhang (2024) launched production on-device language models for mobile virtual keyboard trained with FL.

Achieving **user-level DP for local DP** is particularly challenging and is currently not practical in cross-device FL because the corresponding noise is added on every client for each user.¹¹

Unless otherwise specified, from now on we focus our discussion on **user-level central DP** when describing methods of DP synthetic data generation in FL.

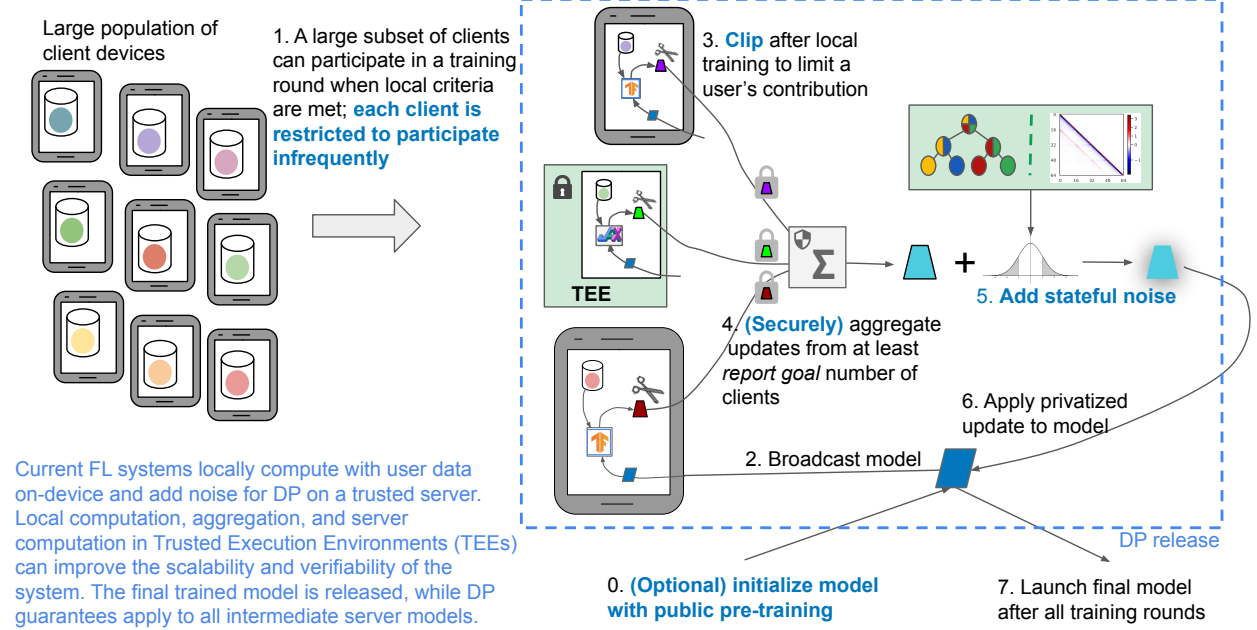


Figure 5: Federated learning with differential privacy deploying Algorithm 7. The architecture is based on current production systems with on-device model training (Xu et al., 2023b), and the integration of Trusted Execution Environments for enhanced scalability, verifiability, and privacy in future systems (Daly et al., 2024).

6.3 DP-Training in FL

Model training is an important technique in FL and it is also a fundamental component for DP synthetic data generation, echoing the discussion in the centralized setting (Sections 4 and 5). Next we provide a deeper dive of the DP FL algorithm in Algorithm 7 for model training as an example on adapting a DP algorithm to the FL setting, as well as preparing background for more discussion on synthetic data algorithms in FL in section 6.5.

Algorithm 7 (with highlighted lines) presents a DP-training version in FL setting. It is based on the generalized FedAvg (Section 6.1). Instead of using the DP-SGD (Algorithm 1) for global model training in the central DP setting, Algorithm 7 adapts the DP-FTRL algorithm (Kairouz et al., 2021a; McMahan et al., 2024) to FL. DP-FTRL uses correlated noise to achieve strong DP guarantee with high model utility without relying on data sampling for privacy amplification. In contrast, DP-SGD adds independent noise in each round, but relies on Poisson or specific form of sampling of clients to achieve strong DP guarantees. DP-FTRL overcomes a significant practical challenge in current cross-device FL: the inability to sample clients because clients are only available when restrictive criteria are met (e.g., charging, on an unmetered network, and idle) (Kairouz et al., 2021b; Wang et al., 2021b). Each client update Δ_i^t is constructed from model training on clients' local data, and then being clipped to control ℓ_2 sensitivity. The client update process before aggregation not only improves communication efficiency in FL, but also naturally fits for user-level DP. While we mostly focus on central DP setting, we want to mention that Generalized FedAvg algorithm can also be adapted to local DP setting and/or example-level privacy unit. Instead of using DP-FTRL for server-side optimization as in Algorithm 7, if the noise is added in the ClientUpdate function for each Δ_i^t , local DP can be achieved.

Local DP with user-level privacy unit As the noise is added to each Δ_i^t for user-level local DP, significantly more clients (n^2 for local DP vs n in central DP) are needed to get similar signal-to-noise ratio comparing to adding noise on aggregated $\sum_i \Delta_i^t$ in central DP. As we mentioned before, local DP on each user device with user privacy unit struggles

¹¹Local DP setting is instead can be used with example-level privacy guarantees to preserve usable utility.

to achieve usable utility and **is currently not practical**. Local DP with user-level DP in cross-silo FL, where each organization has a large number of users, is achievable.

Local DP with example-level privacy unit To get example-level local DP, a straightforward approach is to use DP-SGD (Algorithm 1) to replace SGD in the ClientUpdate function of the FedAvg algorithm.

Though not directly comparable, user-level DP is fundamentally more challenging than example-level DP (Section 2.2), and local DP is more challenging than central DP.

6.4 DP Synthetic Data in FL

The value of (DP) synthetic data is amplified in FL, as it enables privacy-preserving modeling on decentralized data. Synthetic data can be generated to mirror either the global distribution of the entire decentralized data (**global synthetic data**) or the local distributions of individual, heterogeneous clients (**client synthetic data**). In the literature, it is more common to use central DP for global synthetic data, and local DP for client synthetic data, as discussed in Section 6.5.

The applications of synthetic data in FL are numerous. It can serve as a proxy for decentralized data, which can then be safely released outside of the conventional FL cycle (Figure 5) thanks to the post-processing property of DP. For instance, early work used global synthetic data for debugging applications on private domains with unseen data (Augenstein et al., 2020). More recent work highlights the advantages of using global synthetic data as a proxy to flexibly combine public and private knowledge, potentially simplifying the system stack and reducing maintenance costs (Hou et al., 2024; Zhang et al., 2025). Furthermore, Zhang et al. (2025) highlight this proxy data can be inspected and then used in standard datacenter training pipelines to train much larger production models than are feasible in current FL systems.

Synthetic data can also be used to directly improve the FL training process itself. In cross-silo FL, client synthetic data with local DP guarantees has been used to improve federated optimization in image applications on academic benchmark datasets (Xiong et al., 2023; Abacha et al., 2024). In cross-device FL, global synthetic data has been used in pre-training to enhance the DP-FL training of on-device production language models (Wu et al., 2024a).

Generating DP synthetic data in FL is more challenging than in datacenter settings due to the inherent principle of data minimization. We highlight challenges that are mostly relevant for synthetic data generation, mirroring the more general FL challenges in (Wang et al., 2021b):

- **Communication and computation efficiency.** Generating high-quality synthetic data is a non-trivial task and often requires extensive computation resources for both fidelity and privacy. While methods like FedAvg designed for communication and computation efficiency are particularly effective and widely used (Wang et al., 2022), synthetic data generation algorithms are often far more resource-intensive. This is particularly acute in cross-device FL, where users' mobile devices have limited computational power and network bandwidth (Kairouz et al., 2021b). Current production cross-device FL systems support reliable training models of less than 100M parameters; although next-generation systems with TEEs may offer more resources, efficiency remains a primary consideration for supporting real-world applications (Daly et al., 2024).
- **Data heterogeneity and system complexity:** Heterogeneous client data that affects the model training speed and quality is extensively studied in FL (Kairouz et al., 2021b; Wang et al., 2021b). The system complexity of the distributed training paradigm on decentralized data often poses challenges, for example, the trust on server. The characteristics of the system can be challenging for achieving strong privacy-utility trade-off, and require system algorithm co-design such as the DP-FTRL algorithm. These factors, already central to FL as briefly discussed in the previous sections, add further layers of difficulty to the task of generating high-fidelity synthetic data that accurately reflects diverse client distributions while upholding strong privacy guarantees.

Overcoming these challenges presents a significant opportunity. Harnessing the vast amount of data born and siloed on user devices in a privacy-preserving manner could unblock substantial improvements in both the model capacity and user experience of next-generation LLMs, especially for mobile applications.

6.5 Methods for Federated DP Synthetic Data

Method Paper	Modality	Data Distribution	DP Setting	Key DP Method
Pentyala et al. (2024)	Tabular	Global	DP-in-MPC	select-measure-estimate

Maddock et al. (2024)	Tabular	Global	Central, example-level	select-measure-estimate
Xiong et al. (2023)	Image	Client	Local, example-level	DP dataset distillation
Abacha et al. (2024)	Image	Client	Local, example-level	private evolution
Chen et al. (2020a)	Image	Global	Local, example-level	DP model training
Augenstein et al. (2020)	Image & Text	Global	Central, user-level	DP model training
Hou et al. (2024, 2025)	Text	Global	Central, user-level	private evolution
Wu et al. (2024a); Zhang et al. (2025)	Text	Global	Central, user-level	post-processing with small DP LMs

Table 11: DP synthetic data methods in FL. This list does not include methods that are potentially feasible in FL but have not been applied in FL yet, or methods in FL do not use DP. DP-in-MPC (multi party computation) is central DP with reduced trust on the server. Dataset distillation is a learning process that trains the input instead of the model weights, using fixed model and backpropagating to the input. This input will become synthetic data. Dataset distillation can be viewed as a compression of real data into smaller amount of synthetic data.

Many methods for generating DP synthetic data can be adapted for FL, provided that the constraints on communication, computation, and the inherent privacy-utility tradeoff are carefully considered. We discuss the key algorithms of DP synthetic data for tabular, image, and text in FL and assume readers have gained familiarity with the methods from the corresponding modalities (Sections 3, 4 and 5). Table 11 provides a summary of DP synthetic data papers in FL we discuss next.

6.5.1 Methods for global and client synthetic data

We discuss methods for generating global synthetic data that can capture the combined distribution from decentralized data, as well as client synthetic data that can capture the heterogeneous client data at a finer grain. More specifically, using \mathcal{P} to denote the distribution on the population of clients, using \mathcal{D}_i to denote the data distribution of client $i \sim \mathcal{P}$, global synthetic data are sampled from the global distribution, $X \approx \{x|x \sim \mathcal{D}_i, i \sim \mathcal{P}\}$; client synthetic data are sampled separately, $\mathcal{X} = \{X_i|X_i \approx \{x \sim \mathcal{D}_i | i \sim \mathcal{P}\}\}$. The interpretation of \mathcal{P} is essential: if we view \mathcal{P} as a discrete distribution over the specific users participating in FL, then user-level DP will not be possible, as discussed below. On the other hand, if we view \mathcal{P} as a distribution over say user profiles, then we may be able to first sample DP synthetic user profiles, and then sample a client dataset for each such profile.

Methods for client synthetic data Obtaining client synthetic data can be done in local DP setting by using client’s data to finetune a local model (e.g. LLM) with DP-SGD and then subsequently sampling from it. Client synthetic data techniques are therefore often direct analogs to methods in datacenter settings that we explored in previous sections. These techniques apply DP synthetic data algorithms to client data, optionally incorporating a global model that has already learned the commonalities across clients. Only few works discuss client synthetic data, with the focus of using such data for improving FL model training for images (Xiong et al., 2023; Abacha et al., 2024).

Client synthetic data in these studies provide only example-level privacy. Beyond the general challenges with local DP discussed in section 6.3, there is a more fundamental limitation when generating synthetic data for each (real) user under user-level DP: the client synthetic data can only learn commonalities across users, and cannot reveal information of a specific user. That is, the user-level DP guarantee says that the synthetic client data for one user must be indistinguishable from the synthetic client data of any other user, defeating the point of modeling heterogeneous client data distributions. In cross-silo FL, user-level DP for a client can be possibly achieved as each client may contain data from many users. Alternatively, in cross-device FL, the recent hierarchical synthetic data generation approach can be adopted (Kong and Syed, 2025; Hu et al., 2025a): client synthetic data can be generated from synthetic user profiles. Such client synthetic data from synthetic users are compatible with user-level DP, and can satisfy both central DP and local DP. These client synthetic data are useful for improving privacy-preserving model training in federated learning.

Additionally, global DP models can be used on clients as teacher models for knowledge distillation to improve FL model training (Yao et al., 2021; Zhu et al., 2021). In this setting clients use teacher model to synthesize labels or augment features on their own data, and uses the augmented data to improve FL training without explicitly collecting synthetic data.

Methods for global synthetic data Obtaining global synthetic data can be done in many ways. One potential approach is to combine local DP client synthetic data from all clients. A better alternative is to use the central DP model

trained with Algorithm 7 and use it to sample global synthetic data. The latter approach is more commonly used and provides better utility and fidelity of the DP synthetic data.

Our primary focus from now on, unless mentioned otherwise, will be on methods for generating *global synthetic data* that reflects the overall distribution of data across the users. We will discuss tabular, image, and text data, respectively.

6.5.2 DP synthetic tabular data in FL

Synthesizing tabular data is often less resource-intensive and many methods can be straightforwardly federated. For example, many mechanisms from the select-measure-estimate paradigm (Section 3.3.1) can be distributed with federated analytics. The key insight is that these mechanisms only depend on the data through a specific set of low-order marginals, which can readily be computed by federated analytics systems. After these sufficient statistics have been collected, no more server to client communication is necessary to run the remainder of the mechanism. Additional work in this space includes Maddock et al. (2024), who demonstrated how to address data heterogeneity in FL setting. Additionally, Pentyala et al. (2024) have proposed secure multi-party computation protocols to replace the trusted aggregator in workload-based tabular synthesis methods.

6.5.3 DP synthetic image and text data in FL

The key algorithms for image and text, just as in central setting (Sections 4 and 5) can be categorized similarly into DP training/finetuning, DP inference, and Private Evolution.

DP Model Training. Early work on synthesizing image data in FL focused on Generative Adversarial Networks (GANs), carefully selecting components (e.g., discriminators) that could be federated (Augenstein et al., 2020; Chen et al., 2020a; Cai et al., 2021b).

More recently, synthetic text data came into a spotlight. The popular approach of DP fine-tuning Large Language Models (LLMs) (Yue et al., 2022; Kurakin et al., 2023) is difficult to apply directly as is in algorithm 7 because of the computation requirement of updating local LLMs on clients. While communication costs can be addressed with compression, the primary bottleneck is local computation. Parameter-efficient fine-tuning (PEFT) methods like LoRA or prompt tuning can help (Hu et al., 2022; Lester et al., 2021; Cho et al., 2024; Kurakin et al., 2023), but they often still require backpropagation through billion-parameter models, which is computationally expensive, if not impossible, in many FL environments.

DP Inference methods. DP Inference-based methods are unexplored for DP synthetic data in FL. One naive approach is to aggregate in DP manner the outputs from client models when training a central model for generating global synthetic data (Algorithm 7). However, PATE-style aggregation (section 4.5.1) would need public data, and/or a large amount of clients reliably connected at the same time. Private inference for text similar to (Amin et al., 2024, 2025) that utilize in-context learning and generates a sample token by token (Section 5) are challenging in FL because generating a single sample in central DP will require multiple communication rounds for sequential decoding and sufficient number of clients connected at the same time.

Methods that use LLM inference only Prompt-based methods that use foundation models without retraining or finetuning (Wu et al., 2024a; Zhang et al., 2025), including private evolution algorithms (Xie et al., 2024; Hou et al., 2024, 2025), are more lightweight as they only require exchanging DP statistics in each round. Wu et al. (2024a); Zhang et al. (2025) train a small DP LM (less than 10 million parameters) to reweight and filter synthetic data generated by prompting public LLMs. Hou et al. (2024, 2025) extends private evolution algorithms (Algorithms 4 and 6) by sending candidate synthetic data to clients, aggregating the voting histogram from clients, and selecting seed synthetic data for the next round. While these methods show promising early results, they often fail to capture detailed domain information due to inability to update the model that generates the synthetic data, leading to a utility gap on downstream tasks compared to DP fine-tuning on real data (Wu et al., 2024a) or to suboptimal privacy-utility tradeoffs (Xie et al., 2024; Tan et al., 2025). Hou et al. (2025) incorporates model updates into PE: authors update PE with policy optimization approach, where pairs of synthetic examples generated by the same prompt but having different voting scores used to guiding the LLM towards generating better synthetic data via a policy optimization algorithm with LoRA finetuning.

Recent research demonstrates the potential of combining DP fine-tuning with lightweight DP statistics calculated on private data (e.g. topics of the sensitive data) (Yu et al., 2024; Tan et al., 2025). Notably, Tan et al. (2025) generated high-quality DP synthetic data by fine-tuning a relatively small model (140 million parameters) that was conditioned on DP statistics of data aspects like topics. This hybrid approach represents a promising direction for balancing privacy, utility, and computational feasibility in federated text synthesis.

6.6 Evaluating synthetic data in FL

The evaluation of DP synthetic data in FL shares many considerations with the centralized setting. While metrics specific to data types like tabular, image and text data are covered in their dedicated sections (Sections 3.8, 4.7, 5.6), this section discusses evaluation methods that are commonly employed due to the decentralized nature of FL.

- *Human inspection*: The value of direct human inspection should not be underestimated. As discussed in Section 6.4, a common application of synthetic data in FL is for debugging and inspection, as the underlying decentralized data cannot be directly observed. For all applications, developer "eyeballing" of canonical synthetic samples and performing extensive data analysis are invaluable evaluation steps.
- *Downstream Model Performance*: For production applications, the most definitive evaluation often comes from deploying a downstream model trained on synthetic data and measuring its performance through live A/B testing (Wu et al., 2024a; Zhang et al., 2025). An insightful intermediate step is federated evaluation, where the model trained on synthetic data is sent to clients to be evaluated on their real, locally stored data, providing a high-fidelity measure of performance without centralizing the test data (Wu et al., 2024a).
- *Distributional Fidelity metrics*: Another approach is to measure the difference between the distribution of the synthetic data and the real user data, using scores like MAUVE (Pillutla et al., 2021) or FID (Heusel et al., 2017). While straightforward for academic benchmarks where a real evaluation dataset is available, this can be challenging in a practical FL system. It requires either sending synthetic data statistics to clients for local computation and aggregation, or first collecting differentially private statistics from clients to compute the score on the server, both of which add complexity.

6.7 Discussion

Due to computational constraints of on-device learning in FL, methods that avoid DP training of large models remain the most practical as of now (Xu and Zhang, 2025).

As a starting point, prompting LLMs with domain prior knowledge about client's data (Wu et al., 2024a; Zhang et al., 2025) is a practical and effective method, although it results in simple synthetic data, not DP synthetic data that reflects the real data distribution. In this setting, synthetic data is entirely artificial (obtained from public models). For higher-quality data suitable for debugging or developer inspection, private evolution algorithms that select and evolve instances of such public data (Xie et al., 2024; Hou et al., 2024) offer promising improvements in FL settings. For applications where maximizing the performance of downstream models is the primary objective, DP fine-tuning methods are most effective, particularly when combined with DP statistics (Yu et al., 2024; Tan et al., 2025). However, realizing the full potential of DP fine-tuning in practice will require not only algorithmic development but also critical system-level improvements to address the communication and computation constraints inherent in FL (Daly et al., 2024).

Algorithmic development for both current and future FL systems remains an important and active research direction. For research, it is an open question whether the recommendations from existing studies that considered only smaller LLMs (up to 13 billion parameters (Yu et al., 2024)) will scale effectively to much larger models. For production, a significant engineering challenge lies in advancing beyond current FL systems, which largely support only prompt-based or lightweight methods, to build future infrastructure that offers more resources while maintaining strong privacy protections. Finally, the extent to which DP synthetic data can fully represent the utility of decentralized data across all possible use cases, and how much resources are needed to generate such synthetic data remains a fundamental open question.

7 Practical Components of DP synthetic data generation system

This section outlines how one can go about building a DP synthetic data generation system end-to-end. We will touch upon important questions that should be answered, including which privacy guarantees to target (Section 7.1). We will explore data preparation including user level contribution bounding (Section 7.2), and safe sensitive data handling practices (Section 7.6). We will describe a crucial component that should be in place before any DP synthetic data is used in downstream tasks or shared – namely *empirical privacy testing* (Section 7.5). Finally, we will conclude with a discussion of lineage tracking (Section 7.7).

7.1 Privacy-related decisions

There are several important decisions that need to be made by the designer of a system that will produce and employ DP synthetic data. Namely, one must decide what privacy guarantees to target and use that to select an appropriate privacy unit.

Choice of privacy unit is one of the most important choices to make. We outlined possible privacy units in Section 2.2.1 and provided discussions on privacy units for each data modality (Sections 3.1, 4.1, 5.1 and 6.2.1). While the appropriate choice is always application specific, we emphasize that in most cases *example-level privacy unit carries limited meaning in real world applications*, and **user-level** or even larger privacy units like **group/organization privacy units** are preferred. It is also important to remember that ideal privacy unit might be not achievable for various reasons, including a technical inability to identify the same user across multiple data records in order to ensure user-level privacy. We advocate for all practical compromises to be explicitly documented and released with the DP synthetic data.

Target DP parameters (ϵ, δ) The δ parameter is usually set to be less than inverse size of the private dataset being protected (e.g. $1/n^{1.1}$, where n is the number of privacy units within the dataset). For choosing the ϵ parameter, Ponomareva et al. (2023) outlined 3 tiers of privacy guarantees in context of ML models trained with DP, which are also applicable to DP Synthetic data.

Target ϵ values for DP synthetic data (Ponomareva et al., 2023)

We assume **user privacy unit (or example-level where a single user or other appropriate group contributes at most one example)** with the add-or-remove or zero-out adjacency.

1. **Tier 1: Strong formal privacy guarantees.** Strong privacy protection achieved under $\epsilon \leq 1$ directly from the DP definition. Such low values however often result in very low quality (utility and fidelity) of synthetic data, sometimes to a point of being too low to be useful.
2. **Tier 2: Reasonable privacy guarantees.** $\epsilon \leq 10$. The values from this range are currently widely used for DP synthetic data generation in academic papers and production applications. Empirical privacy auditing (Section 7.5) and training data preprocessing (Section 7.4) are highly encouraged in this regime.
3. **Tier 3: Weak to no formal privacy guarantees.** $\epsilon > 10$. While any level of formal DP is an improvement over synthetic data with no formal DP guarantees, DP guarantee on its own for values of ϵ from this tier is vacuous. Additional measures to ensure sufficient randomization, including empirical privacy auditing and training data preprocessing are paramount.

Practitioners are encouraged to choose the lowest achievable tier from the aforementioned tiers.

Finally, it is important to remember that both the method that is chosen for DP synthetic data creation and the modality of data at hand will determine which tiers may be feasible to achieve at reasonable privacy-utility-computation tradeoffs. For example, private evolution methods for images (Section 4.4) can (with a high-quality set of seed data) potentially achieve low ϵ values including Tier 2 and even Tier 1 level guarantees, as can workload-based methods (Section 3.3) for tabular data and DP-training smaller language models (Xu and Zhang, 2024). For small language models, DP training can achieve good utility for small ϵ values from Tier 1 and 2 (Xu and Zhang, 2024). However, for DP-training based methods using large models like diffusion models and LLMs on complex modalities like text or image, such low ϵ values are currently not practically obtainable.

7.2 User contribution bounding

As we have seen in previous sections, user-level privacy often provides the most straightforward and semantically meaningful privacy guarantees. However, many algorithms, for example DP-SGD, in their vanilla forms assume example-level privacy. Further, most of the tried and tested libraries for DP training like Opacus (Yousefpour et al., 2021), Tensorflow Privacy (Abadi et al., 2015), Jax privacy (Balle et al., 2025) provide implementations for example-level DP only.

In order to support a user-level privacy unit, a practitioner has several choices. One involves reducing the problem to example-level privacy and accounting for user-level unit post factum. Straightforward application of group-level privacy (Vadhan, 2017a) (as covered in discussion in Section 2.2.1) requires a bound on group size b , which represents the maximum number of examples each user might contribute to the dataset. To obtain this bound, data can be subsampled to ensure *data level user contribution bounding*. Once the data is subsampled, any example-level DP algorithm like

DP-SGD can be used with $\epsilon/b, \delta/b$ values to account for groups of size b (Vadhan, 2017b). Further, the accounting in some cases like for DP-SGD can be significantly improved further, as we discuss next.

An alternative to reframing the problem to example-level privacy is to employ *algorithms that directly support user (or larger) privacy units* and ensure contribution bounding as part of the algorithms.

While next we cover these two types of solutions in the context of user level privacy, similar recipe can be applied to larger privacy units like groups of users, organizations etc.

User contribution bounding during preprocessing User contribution bounding during preprocessing involves selecting a maximum-size collection of examples from the input such that each user is associated with at most b of the selected examples. When $b = 1$, problems with a user-level privacy unit are essentially reframed as problems with an example-level privacy unit and all the algorithms we discussed previously can be used without additional modification. However, selecting more than 1 datapoint per user is often beneficial — this results in a larger and possibly more diverse training dataset (which in turn allows to use less noise to ensure the same level of privacy, providing better utility). Additionally, users with more data points might have higher quality data, and removing most of their datapoints via subsampling to 1 example per user can be detrimental to data quality (Amin et al., 2019).

For *single-owner data*, such data selection can be achieved independently across users as they do not share any examples (Charles et al., 2024). Each user can pick b of its own examples without interfering with the selection process of other users.

For *multi-attribution scenarios*, where one example can be associated with multiple users, bounding user contributions is more involved. Ganesh et al. (2025) propose a sequential greedy algorithm for the task of maximally selecting instances from a multi-attribution dataset subject to a user contribution bound. Ganesh et al. (2025) show that the result of their greedy algorithm is not too far from the optimum upper bound on the number of selected items achieved through a linear programming formulation. Naturally, sequential algorithms can be limited in their scalability as they must process all examples in one machine one by one. Recent work by Cohen-Addad et al. (2025) addressed this issue introducing a scalable, distributed algorithm to address data level user-level differential privacy contribution bounding problem. This algorithm can be implemented in large-scale parallel and distributed architectures including in MapReduce-like frameworks (Dean and Ghemawat, 2008; Zaharia et al., 2016). The authors modeled the data ownership as a hypergraph, where users are vertices and examples are hyperedges. This modeling allows to cast the example contribution bounding problem as maximizing the number of records while ensuring that each user’s participation in the selected records does not exceed a predefined budget b . This algorithm operates in rounds. In each round, unsaturated users propose their most preferred examples up to their remaining capacity (using an arbitrary but consistent ranking function). An example is only added to the final set if all its participating users unanimously propose it. Given that the algorithm can process each round in parallel it can scale to massive datasets beyond the capacity of a single machine. Once the data is selected, any example-level algorithm like DP-SGD, Private inference or Private Evolution can be used, however some modifications to accounting will still be required to account for at most b samples per user. Improvements over group privacy exist for DP-SGD algorithm. Charles et al. (2024) propose DP-SGD-ELS which employs modified accounting to standard DP SGD algorithm. Instead of relying on group privacy, they leverage the Mixture-of-Gaussians mechanism to derive optimal accounting for DP-SGD. This can decrease compute costs or improve accuracy at a fixed epsilon (alternatively, when holding the number of training iterations and batch size constant, increasing number of samples per user does not significantly affect the epsilon, so larger b can be used at no cost in the level of noise introduced by DP-SGD and the same ϵ guarantees). The suggested number of samples per user b is chosen as a median number of samples across the users, although it can be treated as a hyperparameter and further tuned.

The beauty of data sampling contribution bounding lies in the fact that it does not require changes in training algorithms or data ingestion, allowing the use of established frameworks like Opacus (Yousefpour et al., 2021), Tensorflow Privacy (Abadi et al., 2015), Jax privacy (Balle et al., 2025) and others.

Algorithms that directly support user-level contribution bounding Along with previously mentioned DP-SGD-ELS, Charles et al. (2024) introduces a variant DP-SGD-ULS (user-level sampling and per user gradient clipping) for handling user privacy unit directly in the single-attribution case. While this algorithm foregoes the need for contribution bounding during pre-processing, the data needs to be partitioned and sampled per user during DP-Training. Additionally, clipping of gradients happens on a per-user level as well. ULS often outperforms ELS (data subsampled) variant in fixed compute setting with the magnitude of the improvement is the largest when the data between users is highly variable, and in small ϵ and large compute budget. The downsides are significant modifications needed to the training data ingestion pipeline (propagating user ids), batch preparation (sampling per user data) and actual training algorithm (e.g changing the clipping of DP-SGD from per example to per user etc). ULS is not readily available in any of the well known DP-Training libraries so far. That being said, ULS is closely related to differential privacy in federated

learning discussed in section 6.3. For example, ULS can be constructed from algorithm 7 by using one step of gradient descent in *ClientUpdate*, and carefully adjusting the client sampling strategy to construct Q^t and noise mechanism in aggregating client updates. Previous work (McMahan et al., 2018; Xu et al., 2023a) have implemented such algorithms by combining TensorFlow Privacy (Abadi et al., 2015) and TensorFlow Federated (TensorFlow Federated Authors, 2019) libraries.

7.3 Pitfalls to avoid when implementing DP-training

As we have seen previously, for modalities like image and text, DP-training based methods remain the workhorse methods for creating DP synthetic data. These methods rely on obtaining DP models via DP-Training algorithms like DP-SGD. However, care should be taken when swapping the optimizer to a DP-SGD-based one. We will outline several known angles that need to be considered when choosing a model and its parameters, however we want to stress that in order to obtain a proper DP model, the synthetic data creator needs to have a deep understanding of the model they are using. While ideally frameworks that implement privacy should validate any privacy critical assumptions, currently this is not sufficient. Please refer to (Ponomareva et al., 2023) (Section 5.5.1) for additional discussion of model components that can invalidate privacy guarantees.

Packing For LLM models, packing is one of the most commonly used way of improving training efficiency. The training time of LLMs depends on the sequence length of the longest examples in the batch. When a batch contains examples of various lengths, short examples can be padded up to a maximum example length. However such padding wastes compute, so *packing* is often employed instead. Packing appends several shorter training example, forming one example of length up to the maximum context sequence length (Raffel et al., 2020).

In Section 5.1 we previously outlined that the choice of how to split the text into instances affects the meaning of privacy unit and obtained guarantees. Packing may similarly subtly change the interpretation of the unit of privacy if off-the-shelf algorithm is not adjusted. In the simplest case, when multiple examples are packed together, masking is commonly done to prevent tokens from different examples attending to each other (Chung et al., 2022). For implementation that don't do such masking and examples can attend to each other, an analogy of DP with microbatching can be drawn ((Ponomareva et al., 2023) Section 5.6). This means that sensitivity calculation needs to be adjusted to ensure DP guarantees apply to the original unit of privacy. Such packing also violates Poisson sampling assumptions, which can be handled with adjustment of the noise injected during the training. A microbatching analogy can be also drawn for packing algorithms that create any data dependency and choose which examples to pack together based on other examples. For example, (Staniszewski et al., 2025) introduces packing algorithms that packs together documents that are similar in the embedding space.

Some commonly used in pretraining packing algorithms like *concat-then-split* (Ritter et al., 2023) can additionally end up splitting one example into two, which can be handled with the help of *group privacy* (Dwork and Roth, 2014; Charles et al., 2024) or sensitivity adjustments.

We recommend to keep the packing off when using off-the-shelf DP finetuning methods for text synthetic data generation. If packing is used, careful sensitivity and accounting adjustments as outlined above will be needed

Inconsistencies between accounting assumptions and practice Accounting for DP-SGD-like algorithms often relies on privacy amplification via subsampling. The most commonly used amplifications theorems assume Poisson sampling of examples (where each example is selected independently with constant probability, leading to variable-sized batches). However most actual implementations shuffle the data (often incompletely) and then cycle over the fixed-size batches (Ponomareva et al., 2023). To the best of our knowledge, so far only Opacus (Yousefpour et al., 2021) implements proper Poisson sampling. This gap between infrastructure and theory has led to the common practice in academic papers of training with shuffled data, but doing privacy accounting assuming Poisson sampling. While this practice can lead to a fair comparison of algorithms in a research setting, it should be avoid in production settings including the generation of DP synthetic data from real user data.

Recent theoretical (Chua et al., 2024a,c) and empirical (Annamalai et al., 2025) works show that the gap between the actual ϵ due to shuffling and one computed assuming Poisson sampling can be significant (up to 4x), especially in low noise regimes (few finetuning steps/epochs). If proper Poisson sampling can't be used, DP introduced noise should be increased to account for shuffling.

There are other potential sources of inconsistencies that can worsen actual privacy guarantees — for example using microbatching and clipping at microbatch level, but not accounting for doubled sensitivity ((Ponomareva et al., 2023) Section 5.6).

Components that create inter-example dependencies Per-example gradient clipping (more technically, clipping the gradient associated with each row in a minibatch) is the main mechanism in DP-SGD algorithms to limit the contribution of each training example to the overall gradients and the resulting model. However, any components in modern models that create cross-example dependencies within the batch may break the ability of per-example clipping to bound the sensitivity; or alternatively, a robust implementation of per-example clipping will behave as if the gradients are all computed on batches of size one, and then these are clipped and summed, but this will essentially disable the cross-example components of the loss, likely breaking model training.

We highlight three well-known approaches that could trigger such issues:

1. *Load balancing losses* are often used in Mixture of Experts (MOE) training. *Load balancing losses* are commonly used in an attempt to nudge the gates towards distributing the tokens evenly among experts, to ensure inference time efficiency. The routing of each token when load balancing loss is enabled depends on other examples from the batches and assignment of their tokens to the experts. Such losses should be used during pre-training of MOE but not during DP-finetuning.
2. Similar to load balancing losses, *capped routing* (Lepikhin et al., 2020) which defines a maximum number of tokens that each expert can be assigned for the batch. If more tokens are routed to the expert than its maximum capacity, these extra tokens are not routed to an expert for computation and are instead handled by residual connections after MOE layers. In this setting the prediction of a token from an example depends on whether other tokens from other examples already filled experts’ capacity, again introducing an inter example dependency in a batch. Megablocks (Gale et al., 2022) is an efficient implementation of uncapped routing that voids token dropping, eliminating inter-examples dependencies and therefore should be preferred over capped routing.
3. *BatchNormalization layers*, though less common in modern LLM architectures, are still widely used (Ioffe and Szegedy, 2015). During training, BatchNormalization layers calculate the mean and standard deviation of activations in the batch and uses these statistics to renormalize layer’s input (Ponomareva et al., 2023).

Tuning the training recipe Addressing the pitfalls above requires architecture changes or special care to avoid inadvertently violating the assumptions of the DP analysis, potentially leading to an invalid DP claim. On the other hand, the common pitfall of not carefully re-tuning training hyperparameters is likely to result in a model with very poor utility. DP Training requires additional hyperparameter tuning, including typically a substantial batch size increase, finding an appropriate clipping norm and retuning learning rate and number of epochs (Ponomareva et al., 2023). On top of this, peculiarities like adversarial training (e.g. in GANs) can make DP trained models more unstable, requiring additional tuning or early stopping.

7.4 Best-effort PII Removal

While applying DP during synthetic data generation is a well established and theoretically grounded technique for protecting sensitive information, DP is not a silver bullet. Known limitations of DP include the fact that shared secrets that are repeated between privacy units are not protected with the same strength of privacy guarantees. Additionally, using DP with high values of ϵ leads to theoretical guarantees that, on their own, indicate very little privacy. Moreover, explaining the nuances of DP to non-specialists can be challenging. Finally, a production system should guard against incorrect implementations of DP algorithms (see Section 7.3) that could invalidate the protections.

Therefore, DP should be coupled with simple and straightforward techniques that promote data minimization. Techniques like best-effort PII removal, which we discuss next, are **not sufficient on their own** without DP to ensure privacy of the sensitive data. They instead provide an additional level of reassurance and for modalities like text or tabular data are easy to implement and computationally cheap to apply.

PII removal One preprocessing step that can be used to reduce the chance that DP synthetic data contains sensitive data is *Personally Identifiable Information (PII)* removal from the original sensitive data used for training. Various commercial services (e.g. Google Cloud Sensitive Data Protection (clo, 2025) and Azure PII detection service (azu, 2025)) offer text PII detection and removal capabilities. While the methods behind those services are not necessarily public, in general PII removal techniques can be categorized into rule-based and ML assisted.

Rule-based techniques work by matching the data against a large number of regex-based rules that capture PII that conforms to the expected format (e.g. SSNs, email addresses, phone numbers, employee IDs etc). Such techniques are cheap, simple and effective, however the coverage is defined by the breadth of predefined set of rules used. ML based solutions can be broadly categorized into classifier based that perform Named Entity Recognition (e.g. Name, Address, Phone number, etc) (Vakili et al., 2022) and LLM assisted models that use LLM capabilities. Classifier based models

are usually expected to have some explainability component, so simpler models like Trees or Naive Bayes can be used. LLM-based solutions include finetuning LLMs for PII removal or using prompt engineering (Labelbox, 2023).

It is also possible to combine the various methods for PII removal to maximize the true positive and minimize false negatives. It is worth noticing that after PII detection, one can undertake either PII masking or replacement. Traditional masking of PII (e.g. replace the names with placeholders like <NAME>) might hinder the resulting DP synthetic data performance. Ideally PII detection is also accompanied with replacement using random but semantically meaningful values (e.g. replacing names and addresses with randomly generated/selected names and addresses).

Finally, it is worth pointing out that some modalities like text enjoy established body of work related to PII removal, whereas for example detecting and removal PII in image data might be non-trivial. While embedded metadata like geo location, creator of the image etc. can be removed easily, identifying and removing PII from the actual visual representation of an image (Eke et al., 2021) appears challenging.

7.5 Empirical privacy auditing

Theoretical privacy guarantees, while powerful, are not always a complete picture of a system’s real-world privacy posture. Empirical Privacy Auditing (EPA) is a direct, practical assessment of a system’s privacy properties, conducted by actively mounting privacy attacks against a model or its outputs (e.g., data synthesized from the model). EPA techniques probe a trained ML model to measure how much it has “memorized” about its sensitive training data. These methods, which include membership inference and reconstruction attacks (Shokri et al., 2017; Carlini et al., 2021), essentially act as a stress test by attempting to determine if a specific person’s data was used for training, or by trying to recreate the original data points that are unique to a single person from the model’s outputs. If needed, EPA can produce a (high probability) lower bound on the DP’s ϵ ,¹² quantifying the empirical leakage under a class of realistic attacks. EPA is essential for several reasons that bridge the gap between theory and practice.

First, EPA serves as a critical end-to-end system validation. Modern machine learning pipelines, especially those implementing differentially private training, are extraordinarily complex. A correct implementation requires careful handling of per-example gradient clipping, noise generation and addition, and privacy accounting across many iterations and potentially distributed machines. A subtle bug in any of these components can silently and completely invalidate the theoretical privacy guarantee (Namatevs et al., 2025). An empirical audit, by directly attacking the final artifact (the model or its data), functions as a holistic test that can detect such implementation flaws, providing confidence that the system as a whole is behaving as expected (Nasr et al., 2023).

Second, EPA helps to quantify the conservatism of theoretical DP. For example, the provable (ϵ, δ) guarantees in the context of DP training in Section 4.3 and Section 5.2: (a) hold with respect to worst-case neighboring datasets, (b) assume an adversary who has open box (white-box) access to all intermediate model checkpoints, including the parameters of the final trained or fine-tuned model, and (c) ignore many sources of randomness (e.g. random model initialization or random shuffling and batching). Relaxing some (or all) of these assumptions leads to a more realistic threat model that better captures real-life risks. However, this leads to an intractable analysis of the provable privacy parameters. This is particularly relevant in common scenarios where models are trained to a high theoretical ϵ (e.g., $\epsilon = 10$). Such a level of privacy protection is formally almost vacuous, but may still offer substantial practical privacy, given the formidable challenges of achieving lower ϵ values in large generative models trained on vast datasets without significantly compromising model quality or incurring prohibitive computational costs. EPA provides a methodology to measure this practical level of privacy, which is invaluable for stakeholders making trade-offs between privacy and model utility (Hafner and Sun, 2024).

A deeper layer to this issue is the phenomenon of “empirical privacy variance” (Hu et al., 2025b). Recent work has shown that while the same theoretical (ϵ, δ) -DP guarantee can be achieved through multiple distinct hyperparameter configurations (e.g., different combinations of batch size, learning rate, and number of training iterations), these different configurations can lead to significantly different levels of empirical privacy leakage. Thus, the theoretical (ϵ, δ) pair alone may be insufficient to fully characterize the practical privacy risk of a model, making direct empirical measurement necessary for a complete understanding.

Finally, EPA is a powerful tool for communicating privacy to stakeholders, key opinion formers, and the end users. The formal definitions of DP are notoriously difficult for non-experts, including product managers, legal teams, and end-users, to interpret. A statement such as “the model is $(\epsilon = 8, \delta = 10^{-5})$ -differentially private” offers little intuition about the concrete risks. In contrast, an empirical finding like “we performed a suite of tests demonstrating that the model is incapable of approximately reconstructing any realistic sensitive training data” is far more accessible and provides a more tangible assurance of safety. EPA can provide the evidentiary basis for these more intuitive and

¹²See 7.5.1 for the rationale behind lower bounding ϵ .

convincing privacy claims, while also allowing privacy experts who are familiar with DP to obtain an empirical estimate of its parameters.

7.5.1 The anatomy of a standard audit

While specific techniques vary, a typical empirical privacy audit follows a well-defined, three-step workflow.

Canary insertion. The process begins by identifying or creating a set of target examples, known as “canaries,” that will be the subject of the attack. These can be a reserved subset of the actual training data (e.g., for training data regurgitation tests) or, more commonly, artificially generated records that are inserted into the training set. To simulate a worst-case scenario for a user-level privacy guarantee with a contribution bound of b , each canary is often replicated at least b times in the training data. This represents a “secret” that is present in every example contributed by the most heavily-represented user.¹³ Depending on the attack metric, a corresponding set of “held-out” canaries, drawn from the same distribution but not included in training, may also be required for comparison. The design of effective canaries is a critical and nuanced task, which is explored in detail in Section 7.5.4 below.

Attack execution. With the canaries inserted into the training data, the system is trained, producing a private model or synthetic dataset. A hypothetical adversary then mounts an attack to detect the canaries’ influence. The two most common classes of attack are *reconstruction attacks* and *membership inference attacks* (MIAs). In a reconstruction attack, the adversary attempts to use the model to regenerate each canary, for example, by providing a prefix of a canary and prompting the model to complete it. If the adversary does not have access to the model and is only given data synthesized from it, the adversary can search for (partial) matches between the canary and data synthesized from the model. In an MIA, the adversary is given a canary and must guess whether it was part of the training set (a “member”) or the held-out set (a “non-member”). Of the two, MIA is generally considered the stronger and more fundamental. Intuitively, an adversary performing MIA only needs to guess a single bit of information (whether the example was in the training data) to “win”, whereas reconstruction usually requires determining more than one bit. The existence of a successful reconstruction attack directly entails a successful MIA: guess that the canary is part of the training set if it is reconstructed from its prefix to some degree of accuracy (Shokri et al., 2017).

Risk quantification. The final step is to measure the success of the attack. For reconstruction attacks, this could be the rate of exact or partial success in reconstructing the canaries. A held-out set of i.i.d. canaries is often employed here to test if the difference between reconstruction rates on held-in vs. held-out canaries are statistically significant. For MIAs, a held-out set of i.i.d. canaries is always necessary, and success is quantified using standard binary classification metrics. While simple accuracy and even Area Under the Receiver Operating Characteristic Curve (AUROC) can be misleading,¹⁴ a more robust metric is the True Positive Rate (TPR) at a fixed, low False Positive Rate (FPR), such as TPR@1%FPR (Carlini et al., 2022). A high TPR at a low FPR indicates that the attacker can confidently identify members without making many false accusations, representing a significant privacy breach. These empirical attack results can also be used to compute a statistical lower bound on the effective privacy loss parameter ϵ , providing a direct empirical counterpart to the theoretical guarantee (Namatevs et al., 2025).

It may be counterintuitive that a *lower bound* on the privacy parameter is typically what is reported in a privacy auditing scenario, and the reason for this deserves some discussion. A concrete membership inference attack is a specification of neighboring datasets \mathbb{D} and \mathbb{D}' , and a classifier $F : \text{Range}(\mathcal{M}) \rightarrow \{0, 1\}$. Often F takes the form of a score-threshold rule: $F(o) = \mathbb{I}\{s(o) \geq \tau\}$ for some scoring function s and threshold τ , where higher values of s reflect greater confidence that the underlying dataset was \mathbb{D}' and not \mathbb{D} . The TPR of the attack is then $\Pr[s(\mathcal{M}(\mathbb{D}')) \geq \tau]$ and the FPR is $\Pr[s(\mathcal{M}(\mathbb{D})) \geq \tau]$, where the randomness is taken over the mechanism \mathcal{M} . Without loss of generality, the “true” ϵ of the mechanism is:

$$\sup_{\mathbb{D}, \mathbb{D}'} \sup_s \sup_{\tau} \log \frac{\Pr[s(\mathcal{M}(\mathbb{D}')) \geq \tau]}{\Pr[s(\mathcal{M}(\mathbb{D})) \geq \tau]}.$$

The challenge comes from estimating these probabilities from samples. In particular, for any finite set of samples, there will be some setting of τ for which $\Pr[s(\mathcal{M}(\mathbb{D})) \geq \tau] = 0$, so the estimated ϵ will be infinite. More generally, when the empirical FPR is very low, we will have extreme variance in our empirical ϵ estimate. The accepted solution is to

¹³One may also elect to use a number of canary repetitions that is a multiple of the contribution bound (say $4b$), to demonstrate protection of secrets that are shared among a small group of users. Secrets about a user that may be shared with large groups or appear in training examples that are not owned by the user are harder to protect using DP, however *group privacy* (Section 2.2.1) still offers way to quantify protection for such situations.

¹⁴Metrics like accuracy and AUROC are problematic because it is possible to have low average-case leakage even when the attacker could reconstruct the data of a small set of users with high confidence, something which would rightly be considered a serious privacy violation.

use a statistical upper bound on the FPR so that the estimate remains finite and well-behaved. This can be thought of a smoothing technique so as not to falsely declare a mechanism as non-private due only to estimation error. Thus the challenge of the auditor is to choose \mathbb{D}, \mathbb{D}' , and s , and then optimize τ to obtain the highest achievable lower bound on ε . To the extent that the bound remains low despite the auditor’s best efforts, we have greater confidence that the mechanism is truly private.

7.5.2 Auditing the synthetic data vs. auditing the model

Ultimately our goal is to demonstrate that the synthetic data is anonymous, that is, no user-specific information contained in the training data is present in the synthetic data (explicitly or implicitly). However, if the data comes from a DP finetuned model (e.g. Sections 4.3 and 5.2), it is also possible to verify the model’s privacy, which thereby guarantees the privacy of the data. This choice reflects a deeper decision about the assumed power of the adversary and has significant implications for the scope and strength of the audit’s conclusions.

The primary advantage of a data audit is that it assumes a smaller and more realistic attack surface. The end product that is released is a static dataset, not a queryable model. This aligns with a threat model where the synthetic data itself is released, but any artifacts constructed along the way are ephemeral or private. In this setting, a real attacker who wishes to reconstruct the suffix of a certain prefix they know is in the training data (perhaps “Brendan McMahan’s credit card number is ...”) cannot simply query the model with that prefix as the prompt. Of course, if that prefix happens to occur in the synthetic data, then they could still learn the private information, but careful choices of prompts or post-processing of the synthetic data might defeat such an attack even if it were possible given access to the synthetic data generating model. To reflect these limitations on a real attacker, we change the class of hypothetical attacks we might conduct as part of EPA; for example, we can no longer directly compute the loss of the model on a canary, instead we might look for features of canaries that are also present in the synthetic data.

Data auditing is also a more generally applicable methodology, since it does not depend on the mode of synthetic data generation. The same data audit methods apply equally well to datasets generated by private evolution or DP inference, whereas model-based audits only work by interacting with the DP finetuned model.

Model auditing (for synthetic data created via DP-training or finetuning methods) comes with a different set of pros and cons. First, a model audit tests against a more powerful adversary who can interact with the model directly. If a model can withstand direct interrogation—for instance, through carefully crafted prompts or inspection of its internal state—it provides a strong guarantee that no private information has been memorized in a way that could be exploited. This choice to test against a stronger threat model is a core principle of robust security and privacy engineering. The flip side is that the adversary implied by a model audit may be *too* strong: a model vulnerability may not necessarily translate to a practical data vulnerability. As a general recommendation, model auditing may be preferred when it is possible to pass a stringent model audit while still meeting synthetic data quality requirements.

Second, model auditing offers superior scalability and reusability. Once a generative model has been audited and certified as “anonymous” to a certain standard, an arbitrary amount of synthetic data can be generated from it without the need for additional audits: the privacy guarantee is a property of the generator itself. In contrast, a data audit is a property of a specific data artifact. If a new batch of synthetic data is generated, it must be re-audited from scratch, and one must account for the cumulative empirical privacy loss of multiple releases.

The field of privacy auditing and memorization attacks is rapidly evolving. Model auditing results may be impacted by additional training such as alignment or SFT, even on non-private data (Nasr et al., 2025; Borkar et al., 2025). The privacy audit attack should be chosen to be at least as strong as any conceivable real-world attack, given the full training recipe and the way the model is intended to be used. According to current best practices, synthetic data should be generated directly from a model that has passed a privacy audit, without further processing of the model.

7.5.3 A taxonomy of Privacy Attacks for Empirical Auditing

The landscape of privacy attacks is diverse and rapidly evolving. Here we will discuss some of the more influential model and data attacks.

Model-based auditing. The foundational paradigm of model-based MIAs is built on the tendency of machine learning models to overfit to their training data. Consequently, members of the training set tend to exhibit lower loss values or produce prediction vectors with lower entropy compared to non-members. A simple, though often effective MIA is the *loss attack*, which classifies a example as a member if its loss on the target model is below a certain threshold. In the case of DP finetuning, this can be significantly improved by considering the *likelihood ratio*, which looks at the *difference* in example log-likelihood between the finetuned and pretrained model (Kandpal et al., 2023).

The use of likelihood ratio partially mitigates the failure of the loss attack to account for the *intrinsic complexity* of an example. A more sophisticated mechanism to achieve this is the LiRA attack proposed by Carlini et al. (2022). LiRA operates by calibrating the confidence score of a target point x on the target model. It does this by comparing the target model’s output to the outputs from a large ensemble of “reference” or “shadow” models. Crucially, these reference models are trained on datasets that are similar to the target model’s training data, but some are trained with the target point x , and some are trained without it. By observing the distribution of confidence scores from both sets of models, the attacker can compute a likelihood ratio score for the target point. A high score suggests it is much more likely that the target model’s output was generated by a model that had seen x during training.

Several model-based reconstruction attacks have been proposed. For generative models trained without DP, very simple methods can be used to successfully extract significant amounts of training data, for example prompting an LLM with an example prefix and using greedy or beam-search decoding (Carlini et al., 2021). Even when the model is not generative, for example a text classification model, a loss-based MIA can be used to prune and rank candidate inputs for non-trivial reconstruction success (Elmahdy and Salem, 2024). Attribute inference attacks aim to deduce sensitive attributes of a data record in the training set, even if the attacker does not have the full record, for example (Fredrikson et al., 2015).

Data-based auditing. The most direct and intuitive form of data-based auditing is to search for verbatim or near-verbatim copies of sensitive data. For text, this is most commonly done using n -gram similarity attacks. Meeus et al. (2025) propose a method in which both the canary text and the generated synthetic text are decomposed into n -grams. Then the collection of n -grams from the canary is compared to the collection from the synthetic text using a set-based similarity metric such as the Jaccard similarity or the Cosine similarity. A high similarity score indicates that a significant portion of the canary’s content has been copied into the synthetic output. This method is widely used for tasks like plagiarism detection and identifying near-duplicates in databases, making it a natural fit for detecting blatant data copying in an EPA context. However, its primary limitation is its focus on syntactic overlap. It is excellent at detecting verbatim leakage of a fixed number of tokens but may fail to detect more subtle semantic leakage, where a model learns the meaning of a secret and rephrases it using different words. Data-based auditing is still very much nascent but important area of research.

7.5.4 Canary design

The success of an empirical privacy audit is not solely dependent on the attack algorithm; it is also critically reliant on the design of the “canaries” used as bait. The properties of these inserted (or identified) records determine the strength of the signal the attack attempts to detect. A weak, poorly designed canary will lead to a weak attack, which in turn produces a loose and uninformative empirical bound on privacy leakage. Therefore, the strategic design of canaries is an active and crucial area of research (Panda et al., 2025). A core challenge in canary design is navigating the trade-off between the **strength** of the attack signal and the **naturalness** of the canary (and hence the realism of the risk).

Canary strength. The key criterion is to choose canaries that will maximize the success of the attacker, giving the highest confidence that any privacy violations that could occur in practice will be detected by the audit. Often the most powerful attack signal is generated by canaries that are highly out-of-distribution relative to the rest of the training data. These examples are “surprising” to the model, forcing it to dedicate significant capacity to learn them. This results in large, distinct gradients during training that resist being “unlearned” by other data. However, canary strength also depends on the form of model finetuning, and the choice of model vs. data auditing. For example, if parameter efficient finetuning (e.g., LoRA) is used, the model may not have sufficient capacity to memorize maximally out-of-distribution canaries. Also, for a data-based audit, a canary consisting of extremely unlikely tokens may have a very low likelihood of being generated. Meeus et al. (2025) proposed to use an in-distribution prefix (to maximize the chance of being generated) with an out-of-distribution suffix (to maximize the memorization signal once generation begins).

Naturalness. There are also arguments for using “natural” canaries that are closer to the real data distribution. First, they serve an important communicative purpose. Demonstrating that a model does not leak a canary that looks like realistic user data (e.g., a plausible name and address) is more convincing to users, regulators, and other non-expert stakeholders than demonstrating protection for a string of random tokens. Second, natural canaries are less likely to degrade model utility. Injecting highly out-of-distribution examples into the training data can disrupt the learning process and harm the performance of the final model on its primary task, an undesirable side effect of the auditing process itself.¹⁵

¹⁵If strong out-of-distribution canaries are taking a toll on model quality, an alternative solution is to train two models with identical training setups but for the fact that the auditing model \mathcal{M}_A is trained with canaries added, while the generating model \mathcal{M}_G is trained on only real data. If we make the reasonable assumption that the mere presence of canaries does not significantly affect

Unfortunately, these criteria are expected often to be in conflict. In case it is not possible to satisfy both goals, it is recommended to audit with several sets of canaries to demonstrate robustness to a range of worst-case and realistic threats.

7.5.5 Concrete instantiations of empirical privacy auditing

While the three step auditing framework of canary insertion, attack execution, and risk quantification is broadly applicable, its concrete instantiation varies depending on the specific synthetic data generation mechanism and the data modality. These variations are crucial for designing effective and informative privacy audits. While a deep-dive into the specific auditing literature for various generation mechanisms and data modalities is out of scope for this manuscript, we provide an overview for the sake of completeness. The reader is encouraged to look into the references below for more details.

The influence of the data generation method. The nature of an audit is most directly influenced by the synthetic data generation process. For the widely researched area of DP-trained or DP-finetuned models, the canary-based model and data audits described previously are the standard approach. The attacks typically target the model’s memorization of specific examples introduced during the training or finetuning phase. Auditing becomes more nuanced for private prediction based methods. These approaches, such as PATE or DP inference with large language models, introduce privacy at the prediction stage rather than during training on the sensitive dataset. For PATE style algorithms, which involve an ensemble of models trained on disjoint data subsets, canary based attacks remain applicable to the teacher models. [Chadha et al. \(2024\)](#) introduce a systematic framework for auditing private prediction protocols, providing a precise privacy analysis for specific instances and demonstrating that the empirical privacy leakage can be close to the analytical bounds. For DP inference methods that leverage the in-context learning capabilities of LLMs, audits must contend with the risk of leaking information from the sensitive examples provided in the prompt. Prior work has demonstrated successful membership and reconstruction attacks against in-context learning in non-private settings, highlighting the inherent risks ([Duan et al., 2023](#); [Morris et al., 2024](#)).

Finally, methods that use foundational models for inference only like Private Evolution, use private data solely for a voting step to select or refine publicly generated candidates. These methods present a different challenge. Since the sensitive data is never used to update model parameters, the privacy leakage is more constrained and controlled by the DP guarantees of the voting mechanism. These methods are likely the most secure against direct data reconstruction. However, they can still leak private attributes if a publicly generated candidate happens to be semantically close to a private record and is selected by the private vote. More research is needed to fully quantify the empirical privacy risks for this class of methods.

The influence of the data modality. The data modality also shapes the specifics of the attack execution and canary design.

For synthetic **tabular data**, a rich body of work exists. Membership inference attacks often employ a shadow modeling technique, which is conceptually similar to canary insertion. This involves training many “shadow” models on datasets similar to the private one, some containing the target record and some not. A meta-classifier is then trained on the synthetic data from these shadow models to predict membership in the original dataset ([Stadler et al., 2020](#); [Sidorenko and Tiwald, 2025](#)). Other attacks are tailored to the generation algorithm itself; for instance, [Golob et al. \(2025\)](#) presented a successful membership attack targeting methods that rely on publishing private marginals.

For synthetic **image data**, auditing research has revealed significant vulnerabilities, particularly in modern generative models. Diffusion models, for example, have been shown to memorize and reproduce training images, including photographs and trademarked logos, at a much higher rate than older models like GANs ([Carlini et al., 2023](#)). More recent work has conducted a comprehensive privacy analysis of (non DP) Image AutoRegressive (IAR) models, a newer architecture that can achieve higher image quality and generation speed than diffusion models. [Kowalczyk et al. \(2025\)](#) developed a novel membership inference attack for IARs that achieves a remarkably high success rate (86.38% TPR@1%FPR), demonstrating that IARs are significantly more vulnerable to privacy attacks than diffusion models. Their work also presents a data extraction attack capable of reconstructing up to 698 verbatim training images from a large IAR model, underscoring a critical trade off between model performance and empirical privacy.

Finally, for synthetic **text data**, auditing has evolved from early demonstrations of memorization to sophisticated methods targeting modern LLMs. The reconstruction attacks we discussed previously, where an adversary prompts a model with a prefix of a sensitive sequence and observes if the model completes it verbatim, are often employed with text. This approach was famously used to extract hundreds of unique training sequences from (non-DP) GPT-2, including

privacy properties, and \mathcal{M}_A passes the privacy audit, we can confidently generate from \mathcal{M}_G . However, it is preferred to use a single model if possible.

personally identifiable information, even when the sequences appeared in only a single training document (Carlini et al., 2021). The success of these attacks is heavily influenced by data duplication within the training corpus; models are significantly more likely to memorize and reproduce sequences that appear multiple times, even within a single document (Kandpal et al., 2023). More recent work has demonstrated that even state-of-the-art, aligned models like ChatGPT can be induced to divulge gigabytes of their training data through a novel “divergence attack,” which causes the model to bypass its safety alignment and revert to its base language modeling objective (Nasr et al., 2025).

Beyond direct reconstruction, membership inference attacks for text aim to determine if a specific token sequence was used in training. While classic loss-based MIAs often fail on large language models (Duan et al., 2023), more advanced techniques have shown success. The Likelihood Ratio Attack (LiRA), for instance, improves upon simple loss thresholding by calibrating an example’s loss against the distribution of losses from shadow models, making it a more robust measure for membership (Carlini et al., 2022). When auditing synthetic data specifically, where an adversary only has access to the generated text and not the model, attacks must adapt. Meeus et al. (2025) demonstrate that MIAs based on n -gram frequencies in the synthetic corpus can successfully infer membership, suggesting that memorized secrets subtly alter the statistical properties of the generated text. They also find that standard out-of-distribution canaries are less effective in this data-auditing setting, as they are unlikely to be generated. To address this, they propose a hybrid canary design with an in-distribution prefix to increase the likelihood of generation and a high-perplexity suffix to ensure a strong memorization signal. The ongoing development of such attacks highlights a clear need for robust, specialized (or automated) auditing methodologies for LLM-generated text.

7.6 Security and data minimization during DP synthetic data generation

We have so far emphasized the need for anonymization and auditing algorithms to ensure that the derived synthetic data does not leak any private information. Equally important is the need for rigorous security protocols, strict access control lists, and policies to make sure that user data is protected to the highest standards, as in, if the security of the user data is compromised, even by engineers responsible for generating the synthetic data, the anonymization of the DP-fied synthetic data would not hold. This section discusses some of the security and data minimization best practices.

A cornerstone of secure data handling is strict access control and environmental isolation. All sensitive user data and associated compute workloads, including dataset preprocessing, user contribution bounding, algorithms like differentially private model training or fine-tuning or private evolution, and subsequent auditing, must operate within highly sandboxed and isolated environments. Direct human access to raw sensitive data should be eliminated wherever possible, with instead only explicitly-reviewed and allow-listed code to run on the data via an auditable process. Furthermore, intermediate sensitive data should be subject to strict Time-To-Live (TTL) policies, ensuring its automatic deletion immediately after its designated purpose is fulfilled. This minimizes the window of exposure and reduces the potential impact of a security breach.

Complementing these controls is an unwavering commitment to data minimization (Bonawitz et al., 2022). Any preprocessing or post-processing required for model training or synthetic data generation should be designed to remove unnecessary information as early as possible in the pipeline, with all unneeded artifacts being promptly discarded. Adopting these best practices creates a multi-layered defense, significantly strengthening the overall privacy posture of generative AI system

7.7 Lineage tracking of synthetic data usage*

So far when talked about DP guarantees of DP synthetic data S , these guarantees (ϵ, δ) covered a single run of an algorithm (whether it is DP-finetuning, Private Evolution or Private inference) on a private dataset D to create DP synthetic dataset. While resulting synthetic data S can be used freely (to train any other model, to perform analyses, to combine it with any datasources) due to postprocessing, all other uses of the private dataset D (or datasets that come from the same source as dataset D), including possible training runs used to tune the hyperparameters of the algorithm, or running different algorithms or training other models are not covered by this guarantee. However such isolated guarantees might be insufficient for real world applications. Suppose a company wants to make a DP guaranteed that holds across *all* GenAI/ML/Synthetic Data uses of a particular user data source. For example, a company might want to periodically create (e.g. monthly) DP synthetic datasets as private data of users interacting with a banking assistance comes in. Since such datasets can be potentially combined, providing an isolated guarantee for each dataset will not be sufficient to quantify the risks of such much broader use of user data source. Providing guarantees using *composition rules* (e.g. in the simplest form, essentially summing up ϵ s and δ s) might be too pessimistic.

Applications of DP synthetic data are still nascent and at the current stage it might be not practical yet to aim for such stronger global DP claims. However such real world applications should be considered as DP synthetic data use becomes truly widespread and we gather more knowledge of whether empirical risks compose in similar manner to

formal privacy guarantees across multiple uses of the same source of data. When this happens, data lineage will become essential for supporting such usecases.

The term *data lineage* refers to information about where the data comes from, what it is based on and how it evolves over time (Ikeda and Widom, 2009). In this section we will explore what information should be tracked for DP synthetic datasets to unlock aspirational use of private data sources. While there has been no work specifically on DP synthetic data lineage, we draw inspiration from generic lineage research surveys like (Ikeda and Widom, 2009) and works like (Heinis and Alonso, 2008).

Ikeda and Widom (2009) suggests categorizing the needs of lineage tracking along several axes. First is *where-lineage*, which can be used to understand what inputs the output depends on, and *how lineage*, which can be used for tracking how inputs were transformed to produce the output. Each type of lineage can be tracked at a coarse (dataset) or fine-grained (example) granularity. For DP synthetic data, we recommend tracking *dataset-level where lineage* and *dataset-level how lineage*. Most of the implementation work of lineage tracking would fall into the former. Indeed, the how lineage will simply describe what DP method was used for producing such data (whether it is DP-finetuning, DP inference, Private evolution or a mix of the above), what privacy unit was used and other details outlined in 5.3.3 Section of (Ponomareva et al., 2023).

Next, we will provide some motivational examples that demonstrate that careful lineage tracking is important for use of private data sources.

7.7.1 Motivating examples

The problems that can arise from untracked DP synthetic data can be broadly categorized as problems pertinent to any synthetic data (DP or not, whether it is created using real data or purely synthetic etc.) and problems **specific to DP** synthetic data. While we briefly discuss the first case, our focus will be on DP-specific needs for lineage tracking.

Potential problems with any synthetic data usage Train-eval cross-contamination is a significant risk (where the same or very similar data ends up in both datasets), potentially leading to over-fitting and over-estimation of model performance. An example of such problematic use would be when synthetic data derived from some sensitive data being included in model training dataset, while a portion of the original sensitive data is retained for evaluation of the aforementioned model. In this case DP synthetic data was influenced by the original sensitive data but can't be deduplicated against it (most likely it won't contain exact matches). Ideally the DP synthetic data should be different enough from the source data that the evaluation results would be similar whether the eval data was exactly the same data used to produce the DP synthetic data, or a separate held-out set; but empirical study is needed to test this hypothesis.

A related issue could be synthetic data and the real data from which it was derived both ending up in a the training data, leading to over-weighting this data in the training mixture and potential quality degradation.

DP-specific problems As the creation and use of DP synthetic data becomes more common-place, as we foresee, a lack of rigorous lineage tracking could lead to a gradual weakening privacy guarantees. Next we provide a non-exhaustive list of potential scenarios where this could occur.

E0: The most basic example is, given a private dataset D , running two synthetic data algorithms (e.g. DP finetuning and PE) to obtain datasets S_1 and S_2 . If both algorithms provide (ϵ, δ) guarantees, combination of S_1 and S_2 will have formal guarantees of $(2\epsilon, 2\delta)$

This example demonstrates that there is a close connection between the need for lineage tracking and the importance of carefully disclosing exactly what "private data touches" are covered by a specific DP guarantee. Quoting Ponomareva et al. (2023), *Private data can be accessed for many reasons during the process of building and deploying ML models. DP guarantees should include a description of which of these data uses are covered and which are not. E.g., does the DP guarantee apply (only) to a single training run or does it also cover hyperparameter tuning?* Expanding this to the cases considered here, as long as such precise statements are provided, the above scenario E0 would not violate either of the individual DP guarantees, but it should also be clear that if one wanted a DP guarantee that covers both synthetic datasets S_1 and S_2 , one would have to compose across the two guarantees. While this scenario might be relatively easy to detect, careful lineage tracking is necessary to surface potentially more subtle versions of the same basic issue, that we explore with the help of the next few examples.

E1: For problems that can arise with DP synthetic data built using *example-level privacy unit*, consider the following scenario. DP synthetic data S_{D_1} was produced using a model A and data D_1 with privacy guarantees ϵ , and data S_{D_1} later included into a model B (e.g. during pre-training). We now seek to produce a new version of such synthetic dataset S_{D_2} , using a more recent dataset D_2 , where intersection D_1 and D_2 is not empty. We now take a model B and DP-finetune it with data D_2 with ϵ and generate synthetic dataset S_{D_2} . This is problematic from DP point of view,

as the intersection of D_1 and D_2 is not protected with the same ϵ , as we are now essentially doing twice the number of DP-SGD steps on this portion of the data. This problem can be resolved by de-duplicating D_2 against D_1 before producing the DP synthetic data, but this requires (example-level) lineage tracking.

E2: For a user-level privacy unit, the most straightforward example of weakening of privacy guarantees is as follows. Imagine a user U contributed their data, for example their reviews of restaurants, to a sensitive dataset D_1 (x data points were sampled for each user for inclusion into the dataset). D_1 was then used to create DP synthetic version of this dataset S_{D_1} with user-level guarantees of ϵ . Such synthetic data, among other data sources, was then used to train a model A for classifying any reviews as positive or negative. This model was then subsequently finetuned on a restaurant-only reviews dataset D_2 with differential privacy (ϵ) to obtain a model B specifically for restaurants. If the user U contributed their data to both datasets D_1 and D_2 , even though no single review from the user U appears in both of the datasets (further, intersection of D_1 and D_2 is empty) the user-level privacy guarantees ϵ no longer hold since essentially a user contributed $2x$ datapoints to model B for training. This case can't be detected by duplicates checks of D_1 against D_2 . A more modern twist on this scenario would be a foundational generative model A trained on synthetic data S_{D_1} and then subsequently finetuned by a team responsible for restaurant reviews classification with DP using D_2 .

E3: Finally, consider the case of periodically updated synthetic data dumps. DP synthetic data was built using a month worth of logs of users interaction with some system (for example, all search queries of users in January) with privacy guarantees ϵ . The training dataset for January DP synthetic logs was prepared to ensure user-level privacy: at most x datapoints were sampled per user. Then the new DP synthetic dataset was prepared using February data by following the same procedure (by sampling x datapoints for each user who interacted with the system in February). This process continued for the whole year, resulting in 12 DP synthetic datasets. Clearly, each monthly DP synthetic dataset is ϵ -DP w.r.t to *user-month privacy unit*. However the combination of datasets is not ϵ DP w.r.t *user-level privacy unit*. To see it, imagine first that there was no users whose data appeared in more than 1 month dump. In this case, using parallel composition the combination of all DP synthetic datasets would have been ϵ - DP w.r.t user level privacy unit. However if there are users that had samples in more than 1 month's dump of data, then the resulting combined synthetic dataset is in worst case 12ϵ w.r.t user-level privacy unit.

Discussion An important open question is the degree to which the degradation in the formal DP guarantee in these lineage scenarios actually translates to significantly increased real-world privacy risks. The strongest possible form of DP guarantee would jointly protect *every data access of the totality of a user's privacy data* (across possibly multiple services, and even across multiple companies and organizations). This is an impossibly high bar. Thus, while conceptually the need for lineage tracking is clear, as with the discussion of the unit of privacy, we also encourage pragmatism in calibrating the need for DP guarantees that span multiple "touches" of data to the real-world risks, and we hope that further work in empirical privacy auditing will help guide such decisions.

7.7.2 Potential solutions

In this section we will outline *what* we needs to be tracked to be able to perform lineage checks before using sensitive user data with Differential Privacy.

We won't prescribe how to technically perform such lineage tracking. There are numerous decisions to be made here: what system to use for such tracking (build an in-house system for automatic tracking or manual entry tracking or using a commercial software), how the data should be stored (whether it is an entry in a database, a graph representation etc.) and retrieved efficiently (e.g. via a recursive query to a database or alternative methods like in (Heinis and Alonso, 2008) that use an optimized structure).

Information to track For each DP synthetic datasource we create, we need to keep track of

1. *Reference to the original sensitive dataset* (e.g. internal dataset name or its id)
2. *Privacy unit* - whether it is example/user level/sentence level etc.
3. *DP-related information* (method used for creating DP synthetic data, privacy guarantees details outlined in Section 5.3.3 of (Ponomareva et al., 2023))
4. *Information about any LLM checkpoints* that were used in creating DP synthetic data (and implicitly, its lineage information like which original and DP synthetic datasets it was trained on).
5. *Information necessary to run check that would ensure there is no additional loss in privacy from additional training on the source data:*
 - (a) **For user level privacy unit**
 - i. *User ids* of users who contributed the data to the dataset that was used for building DP synthetic data. These user ids will be used for de-duplication against each new related data-source. Even if the DP synthetic data is fully anonymous, this set of user ids is potentially highly privacy sensitive (consider

a DP synthetic dataset derived from a database of medical histories, for example). Thus, while use of the DP synthetic data might be relatively unencumbered, the strong security and access controls discussed in Section 7.6 might be necessary to track the lineage information, raising additional systems challenges.

- ii. Maximum number of examples per user.
- (b) **For example-level privacy unit**, if the original dataset is persisted indefinitely, saving the reference against the dataset would suffice. If the dataset is a subject to being removed (for example, due to time-to-live requirements), an alternative to running the date duplicate check is to demand that each DP synthetic dataset is generated using unique time frame (so lineage tracking will save the time-frame used for DP synthetic dataset creation).

Verifications prior to DP synthetic dataset creation Tracking the lineage should go hand-in-hand with the checks that need to be implemented before creating any new DP synthetic dataset. Namely

1. **For user-level privacy unit**
 - (a) Verify that the LLM checkpoint has not been built using related DP synthetic dataset (datasets which source is related to the current sensitive dataset that will be dp-fied) with the intersecting user ids. This check is recursive - for any DP synthetic dataset that was used for the checkpoint, the check should run using all the checkpoints in the lineage graph.
 - (b) Verify for each original sensitive dataset, if there was a related sensitive dataset that was used to create DP synthetic data. If yes, verify that user ids used are non-intersecting with the current’ sensitive dataset user ids (consider the **E3** case).
2. **For example-level privacy unit** the checks similar to the above should verify that some portion of the data was not already used in any DP synthetic data and any of the checkpoints in the lineage graph.

8 Conclusion

Recent advances in DP synthetic data have made the technology a viable option for deriving value from sensitive datasets without compromising privacy. When the private source data will only be used for one specific purpose (extracting one statistic, or training a single model), then DP technology specialized to that purposes may still be preferable in terms of utility, privacy, and computation. However, this is rarely the case in the real world. DP synthetic data shines because it can offer a drop-in replacement for the original data, which can subsequently be used for many downstream tasks: plugged into state-of-the-art ML pipelines (without any need for complex DP modifications), inspected and categorized by humans, shared broadly to researchers to further societal goals (for example, environmental or medical data), and more. We hope this work makes clear the range of sophisticated techniques already available for developing DP synthetic data.

Nevertheless, significant challenges and open problems remain. While scattered proofs-of-concept exist, generating high-quality and privacy-safe DP synthetic data requires expert knowledge and complex infrastructure, and as long as this is the case, usage will be limited. There is a clear need to more fully automate and systematize the creation of DP synthetic data in order to make the technology more widely available. A related challenge is that for many domains (e.g., large-scale text or image generation), the source dataset needs to be large (in terms of the number of privacy units it contains) in order to achieve good utility of synthetic data, and substantial computational resources are required for techniques like DP fine-tuning. Bringing down these costs is an obvious and important research direction.

Beyond the purely technical challenges, DP synthetic data deserves a broader dialogue among researchers, privacy and legal experts, regulators, and even end users. As hopefully some of our discussion around the choice of privacy unit and the role of lineage tracking show, a particular DP guarantee can easily be both much stronger than necessary (providing protections against attacks that are far more powerful than are realistically expected), but also fail to address other privacy risks that could be a concern (for example, if a particular secret is contained in many different privacy units of data).

Stepping back, we cannot escape that privacy is, at the end of the day, a humanistic value, not a mathematical construct. While DP provides a rigorous foundation for reasoning about privacy, as with any mathematical definition, it cannot fully capture the nuance, richness, and even subjectiveness of Privacy in the real world. Questions like “What is a reasonable privacy unit and choice of ϵ for this application?”, “How big are the privacy risks, and what costs are associated with them? How do we balance those risks against the value this data can offer?”, and the like require rich discussion. Our hope is that this work can provide a technical foundation for such conversations, with a goal of establishing best-practices that can provide assurances to synthetic data producers that they are meeting an appropriate standard-of-care for privacy-sensitive data.

9 Acknowledgement

The authors would like to thank Alessandro Epasto for user contribution bounding discussions, Umar Syed for reviewing the paper, Silvio Lattanzi, Andrew Tomkins, Daniel Ramage, Erik Vee, Raluca Ada Popa for helping navigate internal processes, Badih Ghazi for participating in initial discussions and brainstorming.

References

- (2025). Azure ai language personally identifiable information (pii) detection.
- (2025). Google cloud sensitive data protection.
- Abacha, F. Z., Teo, S. G., Cordeiro, L. C., and Mustafa, M. A. (2024). Synthetic data aided federated learning using foundation models. In *International Workshop on Trustworthy Federated Learning*, pages 106–118. Springer.
- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X. (2015). TensorFlow: Large-scale machine learning on heterogeneous systems. Software available from tensorflow.org.
- Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., and Zhang, L. (2016). Deep learning with differential privacy. *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*.
- Abay, N., Zhou, Y., Kantarcioglu, M., Thuraisingham, B., and Sweeney, L. (2018). Privacy preserving synthetic data release using deep learning. pages 510–526.
- Abowd, J., Ashmead, R., Simson, G., Kifer, D., Leclerc, P., Machanavajjhala, A., and Sexton, W. (2019). Census topdown: Differentially private data, incremental schemas, and consistency with public knowledge. *US Census Bureau*, pages 1–21.
- Abowd, J. M., Ashmead, R., Cumings-Menon, R., Garfinkel, S., Heineck, M., Heiss, C., Johns, R., Kifer, D., Leclerc, P., Machanavajjhala, A., Moran, B., Sexton, W., Spence, M., and Zhuravlev, P. (2022). The 2020 census disclosure avoidance system topdown algorithm.
- Acs, G., Melis, L., Castelluccia, C., and Cristofaro, E. D. (2018). Differentially private mixture of generative neural networks.
- Afonja, G., Sim, R., Lin, Z., Inan, H. A., and Yekhanin, S. (Accessed 2025a). The crossroads of innovation and privacy: Private synthetic data for generative ai. <https://www.microsoft.com/en-us/research/blog/the-crossroads-of-innovation-and-privacy-private-synthetic-data-for-generative-ai/>.
- Afonja, T., Wang, H.-P., Kerkouche, R., and Fritz, M. (2025b). Dp-2stage: Adapting language models as differentially private tabular data generators.
- Alaa, A. M., van Breugel, B., Saveliev, E., and van der Schaar, M. (2022). How faithful is your synthetic data? sample-level metrics for evaluating and auditing generative models.
- Amin, K., Avestimehr, S., Babakniya, S., Bie, A., Kong, W., Ponomareva, N., and Syed, U. (2025). Clustering and median aggregation improve differentially private inference. *CoRR*.
- Amin, K., Bie, A., Kong, W., Kurakin, A., Ponomareva, N., Syed, U., Terzis, A., and Vassilvitskii, S. (2024). Private prediction for large-scale synthetic text generation. In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 7244–7262. Association for Computational Linguistics.
- Amin, K., Kulesza, A., Munoz, A., and Vassilvitskii, S. (2019). Bounding user contributions: A bias-variance trade-off in differential privacy. In Chaudhuri, K. and Salakhutdinov, R., editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 263–271. PMLR.
- Annamalai, M. S. M. S., Balle, B., Hayes, J., and Cristofaro, E. D. (2025). To shuffle or not to shuffle: Auditing dp-sgd with shuffling.
- Augenstein, S., McMahan, H. B., Ramage, D., Ramaswamy, S., Kairouz, P., Chen, M., Mathews, R., and y Arcas, B. A. (2020). Generative models for effective ML on private, decentralized datasets. In *International Conference on Learning Representations*.

- Aydore, S., Brown, W., Kearns, M., Kenthapadi, K., Melis, L., Roth, A., and Siva, A. A. (2021). Differentially private query release through adaptive projection. In Meila, M. and Zhang, T., editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 457–467. PMLR.
- Bagdasaryan, E., Poursaeed, O., and Shmatikov, V. (2019). Differential privacy has disparate impact on model accuracy. *Advances in neural information processing systems*.
- Balle, B., Barthe, G., and Gaboardi, M. (2018). Privacy amplification by subsampling: Tight analyses via couplings and divergences. *Advances in neural information processing systems*, 31.
- Balle, B., Berrada, L., Charles, Z., Choquette-Choo, C. A., De, S., Doroshenko, V., Dvijotham, D., Galen, A., Ganesh, A., Ghalebikesabi, S., Hayes, J., Kairouz, P., McKenna, R., McMahan, B., Pappu, A., Ponomareva, N., Pratilov, M., Rush, K., Smith, S. L., and Stanforth, R. (2025). JAX-Privacy: Algorithms for privacy-preserving machine learning in jax.
- Barak, B., Chaudhuri, K., Dwork, C., Kale, S., McSherry, F., and Talwar, K. (2007). Privacy, accuracy, and consistency too: a holistic solution to contingency table release. In *PODS*, pages 273–282.
- Becker, B. and Kohavi, R. (1996). Adult. UCI Machine Learning Repository. DOI: <https://doi.org/10.24432/C5XW20>.
- Bie, A., Kamath, G., and Zhang, G. (2023). Private GANs, revisited. *Transactions on Machine Learning Research*. Survey Certification.
- Blum, A., Ligett, K., and Roth, A. (2013). A learning theory approach to noninteractive database privacy. *J. ACM*, 60(2):12:1–12:25.
- Boediardjo, M., Strohmer, T., and Vershynin, R. (2024). Private measures, random walks, and synthetic data. *Probability theory and related fields*, pages 1–43.
- Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosselut, A., Brunskill, E., et al. (2021). On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.
- Bommasani, R., Wu, S., and Schofield, X. (2019). Towards private synthetic text generation. In *NeurIPS 2019 Machine Learning with Guarantees Workshop*.
- Bonawitz, K., Kairouz, P., McMahan, B., and Ramage, D. (2022). Federated learning and privacy. *Communications of the ACM*, 65(4):90–97.
- Borisov, V., Seßler, K., Leemann, T., Pawelczyk, M., and Kasneci, G. (2023). Language models are realistic tabular data generators.
- Borkar, J., Jagielski, M., Lee, K., Mireshghallah, N., Smith, D. A., and Choquette-Choo, C. A. (2025). Privacy ripple effects from adding or removing personal information in language model training. *arXiv preprint arXiv:2502.15680*.
- Bousquet, O., Livni, R., and Moran, S. (2020). Synthetic data generators - sequential and private. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Bowyer, K. W., Chawla, N. V., Hall, L. O., and Kegelmeyer, W. P. (2011). SMOTE: synthetic minority over-sampling technique. *CoRR*, abs/1106.1813.
- Brown, H., Lee, K., Mireshghallah, F., Shokri, R., and Tramèr, F. (2022). What does it mean for a language model to preserve privacy?
- Bun, M. and Steinke, T. (2016). Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *Theory of cryptography conference*, pages 635–658. Springer.
- Cai, K., Lei, X., Wei, J., and Xiao, X. (2021a). Data synthesis via differentially private markov random fields. *Proceedings of the VLDB Endowment*, 14(11):2190–2202.
- Cai, K., Xiao, X., and Cormode, G. (2023). Privlava: synthesizing relational data with foreign keys under differential privacy. *Proceedings of the ACM on Management of Data*, 1(2):1–25.
- Cai, Z., Xiong, Z., Xu, H., Wang, P., Li, W., and Pan, Y. (2021b). Generative adversarial networks: A survey toward private and secure applications. *ACM Computing Surveys (CSUR)*, 54(6):1–38.
- Cao, T., Bie, A., Vahdat, A., Fidler, S., and Kreis, K. (2021). Don’t generate me: Training differentially private generative models with sinkhorn divergence. *Advances in Neural Information Processing Systems*, 34:12480–12492.
- Carlini, N., Chien, S., Nasr, M., Song, S., Terzis, A., and Tramer, F. (2022). Membership inference attacks from first principles. In *2022 IEEE symposium on security and privacy (SP)*, pages 1897–1914. IEEE.

- Carlini, N., Hayes, J., Nasr, M., Jagielski, M., Sehwag, V., Tramer, F., Balle, B., Ippolito, D., and Wallace, E. (2023). Extracting training data from diffusion models. In *32nd USENIX security symposium (USENIX Security 23)*, pages 5253–5270.
- Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T., Song, D., Erlingsson, U., Oprea, A., and Raffel, C. (2021). Extracting training data from large language models.
- Carranza, A., Farahani, R., Ponomareva, N., Kurakin, A., Jagielski, M., and Nasr, M. (2024). Synthetic query generation for privacy-preserving deep retrieval systems using differentially private language models. In Duh, K., Gomez, H., and Bethard, S., editors, *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3920–3930, Mexico City, Mexico. Association for Computational Linguistics.
- Castellon, R., Gopal, A., Bloniarz, B., and Rosenberg, D. (2023). Dp-tbart: A transformer-based autoregressive model for differentially private tabular data generation.
- Chadha, K., Jagielski, M., Papernot, N., Choquette-Choo, C. A., and Nasr, M. (2024). Auditing private prediction. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.
- Charles, Z., Ganesh, A., McKenna, R., McMahan, H. B., Mitchell, N., Pillutla, K., and Rush, K. (2024). Fine-tuning large language models with user-level differential privacy.
- Chen, D., Cheung, S. S., Chuah, C., and Ozonoff, S. (2022a). Differentially private generative adversarial networks with model inversion. *CoRR*, abs/2201.03139.
- Chen, D., Orekondy, T., and Fritz, M. (2020a). Gs-wgan: A gradient-sanitized approach for learning differentially private generators. *Advances in Neural Information Processing Systems*, 33:12673–12684.
- Chen, J.-W., Yu, C.-M., Kao, C.-C., Pang, T.-W., and Lu, C.-S. (2022b). Dpgen: Differentially private generative energy-guided network for natural image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8387–8396.
- Chen, K., Li, X., Gong, C., McKenna, R., and Wang, T. (2025). Benchmarking differentially private tabular data synthesis.
- Chen, M., Radford, A., Child, R., Wu, J., Jun, H., Dhariwal, P., Luan, D., and Sutskever, I. (2020b). Generative pretraining from pixels.
- Cho, Y. J., Liu, L., Xu, Z., Fahrezi, A., and Joshi, G. (2024). Heterogeneous lora for federated fine-tuning of on-device foundation models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 12903–12913.
- Chua, L., Ghazi, B., Kamath, P., Kumar, R., Manurangsi, P., Sinha, A., and Zhang, C. (2024a). How private are dp-sgd implementations?
- Chua, L., Ghazi, B., Kamath, P., Kumar, R., Manurangsi, P., Sinha, A., and Zhang, C. (2024b). Scalable dp-sgd: Shuffling vs. poisson subsampling. *Advances in Neural Information Processing Systems*, 37:70026–70047.
- Chua, L., Ghazi, B., Kamath, P., Kumar, R., Manurangsi, P., Sinha, A., and Zhang, C. (2024c). Scalable dp-sgd: Shuffling vs. poisson subsampling.
- Chung, H. W., Hou, L., Longpre, S., Zoph, B., Tay, Y., Fedus, W., Li, Y., Wang, X., Dehghani, M., Brahma, S., Webson, A., Gu, S. S., Dai, Z., Suzgun, M., Chen, X., Chowdhery, A., Castro-Ros, A., Pellat, M., Robinson, K., Valter, D., Narang, S., Mishra, G., Yu, A., Zhao, V., Huang, Y., Dai, A., Yu, H., Petrov, S., Chi, E. H., Dean, J., Devlin, J., Roberts, A., Zhou, D., Le, Q. V., and Wei, J. (2022). Scaling instruction-finetuned language models.
- Cohen, E., Lyu, X., Nelson, J., Sarlós, T., and Stemmer, U. (2023). Hot PATE: private aggregation of distributions for diverse task. *CoRR*, abs/2312.02132.
- Cohen-Addad, V., Epasto, A., Lee, J., and Zadimoghaddam, M. (2025). Scalable contribution bounding to achieve privacy. *arXiv preprint arXiv:2507.23432*.
- Cormode, G., Procopiuc, C., Srivastava, D., Shen, E., and Yu, T. (2012). Differentially private spatial decompositions. In *ICDE*, pages 20–31.
- Dagan, Y. and Kur, G. (2022). A bounded-noise mechanism for differential privacy. In *COLT*, pages 625–661.
- Daly, K., Eichner, H., Kairouz, P., McMahan, H. B., Ramage, D., and Xu, Z. (2024). Federated learning in practice: reflections and projections. In *2024 IEEE 6th International Conference on Trust, Privacy and Security in Intelligent Systems, and Applications (TPS-ISA)*, pages 148–156. IEEE.

- Dean, J. and Ghemawat, S. (2008). Mapreduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113.
- Dockhorn, T., Cao, T., Vahdat, A., and Kreis, K. (2023). Differentially private diffusion models.
- Dong, J., Roth, A., and Su, W. J. (2022). Gaussian differential privacy. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 84(1):3–37.
- Donhauser, K., Lokna, J., Sanyal, A., Boedihardjo, M., Hönig, R., and Yang, F. (2023). Sample-efficient private data release for lipschitz functions under sparsity assumptions. *CoRR*, abs/2302.09680.
- DP Team (2022). Google’s differential privacy libraries. <https://github.com/google/differential-privacy>.
- Duan, H., Dziedzic, A., Papernot, N., and Boenisch, F. (2023). Flocks of stochastic parrots: Differentially private prompt learning for large language models. *Advances in Neural Information Processing Systems*, 36:76852–76871.
- Durfee, D. and Rogers, R. M. (2019). Practical differentially private top-k selection with pay-what-you-get composition. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Dwork, C. and Feldman, V. (2018). Privacy-preserving prediction. In *Conference On Learning Theory*, pages 1693–1702. PMLR.
- Dwork, C., Naor, M., Reingold, O., Rothblum, G. N., and Vadhan, S. (2009). On the complexity of differentially private data release: efficient algorithms and hardness results. In *Proceedings of the forty-first annual ACM symposium on Theory of computing*, pages 381–390.
- Dwork, C. and Roth, A. (2014). The algorithmic foundations of differential privacy.
- Eke, D., Aasebø, I. E., Akintoye, S., Knight, W., Karakasidis, A., Mikulan, E., Ochang, P., Ogoh, G., Oostenveld, R., Pigorini, A., Stahl, B. C., White, T., and Zehl, L. (2021). Pseudonymisation of neuroimages and data protection: Increasing access to data while retaining scientific utility. *Neuroimage: Reports*, 1(4):100053.
- Elmahdy, A. and Salem, A. (2024). Deconstructing classifiers: Towards a data reconstruction attack against text classification models. In Habernal, I., Ghanavati, S., Ravichander, A., Jain, V., Thaine, P., Igamberdiev, T., Mireshghallah, N., and Feyisetan, O., editors, *Proceedings of the Fifth Workshop on Privacy in Natural Language Processing*, pages 143–158, Bangkok, Thailand. Association for Computational Linguistics.
- Epasto, A., Mao, J., Medina, A. M., Mirrokni, V., Vassilvitskii, S., and Zhong, P. (2023). Differentially private continual releases of streaming frequency moment estimations. In *ITCS*, pages 48:1–48:24.
- Fan, A., Lewis, M., and Dauphin, Y. (2018). Hierarchical neural story generation. In Gurevych, I. and Miyao, Y., editors, *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898, Melbourne, Australia. Association for Computational Linguistics.
- Feldman, V., McMillan, A., Sivakumar, S., and Talwar, K. (2024). Instance-optimal private density estimation in the wasserstein distance. *Advances in Neural Information Processing Systems*, 37:90061–90131.
- Flemings, J., Gan, H., Li, H., Razaviyayn, M., and Annavaram, M. (2025). Differentially private in-context learning via sampling few-shot mixed with zero-shot outputs.
- Flemings, J., Razaviyayn, M., and Annavaram, M. (2024a). Adaptively private next-token prediction of large language models. *arXiv preprint arXiv:2410.02016*.
- Flemings, J., Razaviyayn, M., and Annavaram, M. (2024b). Differentially private next-token prediction of large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 4390–4404, Mexico City, Mexico. Association for Computational Linguistics.
- Fredrikson, M., Jha, S., and Ristenpart, T. (2015). Model inversion attacks that exploit confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, pages 1322–1333.
- Fuentes, M., Mullins, B. C., McKenna, R., Miklau, G., and Sheldon, D. (2024). Joint selection: Adaptively incorporating public information for private synthetic data. In *International Conference on Artificial Intelligence and Statistics*, pages 2404–2412. PMLR.
- Gaboardi, M., Arias, E. J. G., Hsu, J., Roth, A., and Wu, Z. S. (2014). Dual query: Practical private query release for high dimensional data. In *ICML*, pages 1170–1178.
- Gale, T., Narayanan, D., Young, C., and Zaharia, M. (2022). Megablocks: Efficient sparse training with mixture-of-experts.
- Ganesh, A., McKenna, R., McMahan, B., Smith, A., and Wu, F. (2025). It’s my data too: Private ml for datasets with multi-user training examples.

- Ganesh, A. and Zhao, J. (2021). Privately Answering Counting Queries with Generalized Gaussian Mechanisms. In Ligett, K. and Gupta, S., editors, *2nd Symposium on Foundations of Responsible Computing (FORC 2021)*, volume 192 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 1:1–1:18, Dagstuhl, Germany. Schloss Dagstuhl – Leibniz-Zentrum für Informatik.
- Ganev, G., Annamalai, M. S. M. S., and Cristofaro, E. D. (2025a). The elusive pursuit of reproducing pate-gan: Benchmarking, auditing, debugging.
- Ganev, G., Annamalai, M. S. M. S., Mahiou, S., and De Cristofaro, E. (2025b). The importance of being discrete: Measuring the impact of discretization in end-to-end differentially private synthetic data. *arXiv preprint arXiv:2504.06923*.
- Gao, F., Zhou, R., Wang, T., Shen, C., and Yang, J. (2025). Data-adaptive differentially private prompt synthesis for in-context learning. In *The Thirteenth International Conference on Learning Representations*.
- Ghazi, B., Guzmán, C., Kamath, P., Knop, A., Kumar, R., Manurangsi, P., and Sachdeva, S. (2025). Prem: Privately answering statistical queries with relative error.
- Ghazi, B., He, J., Kohlhoff, K., Kumar, R., Manurangsi, P., Navalpakkam, V., and Valliappan, N. (2023a). Differentially private heatmaps. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(6):7696–7704.
- Ghazi, B., Kamath, P., Kumar, R., Manurangsi, P., and Wu, K. (2023b). On Differentially Private Counting on Trees. In *ICALP*, pages 66:1–66:18.
- Ghazi, B., Kumar, R., and Manurangsi, P. (2021). On avoiding the union bound when answering multiple differentially private queries. In *COLT*, pages 2133–2146.
- Gillenwater, J., Joseph, M., and Kulesza, A. (2021). Differentially private quantiles. In *International Conference on Machine Learning*, pages 3713–3722. PMLR.
- Ginart, A., van der Maaten, L., Zou, J., and Guo, C. (2022). Submix: Practical private prediction for large-scale language models. *arXiv preprint arXiv:2201.00971*.
- Golob, S., Pentyala, S., Maratkhan, A., and Cock, M. D. (2025). Privacy vulnerabilities in marginals-based synthetic data.
- Gong, C., Li, K., Lin, Z., and Wang, T. (2025). Dpimagebench: A unified benchmark for differentially private image synthesis. *arXiv preprint arXiv:2503.14681*.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial networks.
- Gupta, A., Roth, A., and Ullman, J. R. (2012). Iterative constructions and private data release. In *TCC*, pages 339–356.
- Hafner, F. and Sun, C. (2024). Empirical privacy evaluations of generative and predictive machine learning models – a review and challenges for practice.
- Harder, F., Adamczewski, K., and Park, M. (2020). Differentially private mean embeddings with random features (DP-MERF) for simple & practical synthetic data generation. *CoRR*, abs/2002.11603.
- Hardt, M. (2011). A study of privacy and fairness in sensitive data analysis. *Ph. D. Dissertation*.
- Hardt, M., Ligett, K., and McSherry, F. (2012a). A simple and practical algorithm for differentially private data release. *Advances in neural information processing systems*, 25.
- Hardt, M., Ligett, K., and McSherry, F. (2012b). A simple and practical algorithm for differentially private data release. In *NIPS*, pages 2348–2356.
- Hardt, M. and Rothblum, G. N. (2010). A multiplicative weights mechanism for privacy-preserving data analysis. In *FOCS*, pages 61–70.
- Hay, M., Rastogi, V., Miklau, G., and Suciu, D. (2010). Boosting the accuracy of differentially private histograms through consistency. *Proceedings of the VLDB Endowment*, 3(1).
- He, Y., Strohmer, T., Vershynin, R., and Zhu, Y. (2025). Differentially private low-dimensional synthetic data from high-dimensional datasets. *Information and Inference: A Journal of the IMA*, 14(1):iaae034.
- He, Y., Vershynin, R., and Zhu, Y. (2023). Algorithmically effective differentially private synthetic data. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 3941–3968. PMLR.
- He, Y., Vershynin, R., and Zhu, Y. (2024). Online differentially private synthetic data generation. *IEEE Transactions on Privacy*, 1:19–30.

- Heinis, T. and Alonso, G. (2008). Efficient lineage tracking for scientific workflows. In *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, SIGMOD '08, page 1007–1018, New York, NY, USA. Association for Computing Machinery.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. (2017). Gans trained by a two time-scale update rule converge to a local nash equilibrium. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Hod, S. and Canetti, R. (2024). Differentially private release of israel’s national registry of live births. In *2025 IEEE Symposium on Security and Privacy (SP)*, pages 101–101. IEEE Computer Society.
- Holtzman, A., Buys, J., Du, L., Forbes, M., and Choi, Y. (2020). The curious case of neural text degeneration. In *International Conference on Learning Representations*.
- Hong, J., Wang, J. T., Zhang, C., LI, Z., Li, B., and Wang, Z. (2024). DP-OPT: Make large language model your privacy-preserving prompt engineer. In *The Twelfth International Conference on Learning Representations*.
- Hou, C., Shrivastava, A., Zhan, H., Conway, R., Le, T., Sagar, A., Fanti, G., and Lazar, D. (2024). Pre-text: Training language models on private federated data in the age of llms. In *Forty-first International Conference on Machine Learning (ICML)*.
- Hou, C., Wang, M.-Y., Zhu, Y., Lazar, D., and Fanti, G. (2025). Popri:private federated learning using preference-optimized synthetic data. *arXiv preprint arXiv:2504.16438*.
- Hu, A., Guan, J., Torr, P., and Pinto, F. (2024). A cautionary tale on the evaluation of differentially private in-context learning.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. (2022). Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.
- Hu, Y., McKenna, R., Yu, D., Wu, S., Zhao, H., Xu, Z., and Kairouz, P. (2025a). Actg-arl: Differentially private conditional text generation with rl-boosted control.
- Hu, Y., Wu, F., Xian, R., Liu, Y., Zakynthinou, L., Kamath, P., Zhang, C., and Forsyth, D. (2025b). Empirical privacy variance.
- Ikeda, R. and Widom, J. (2009). Data lineage: A survey.
- Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. *CoRR*, abs/1502.03167.
- Jiang, Z., Ni, W., and Zhang, Y. (2024). Pate-triplegan: Privacy-preserving image synthesis with gaussian differential privacy.
- Jordon, J., Szpruch, L., Houssiau, F., Bottarelli, M., Cherubin, G., Maple, C., Cohen, S. N., and Weller, A. (2022). Synthetic data – what, why and how?
- Kairouz, P., McMahan, B., Song, S., Thakkar, O., Thakurta, A., and Xu, Z. (2021a). Practical and private (deep) learning without sampling or shuffling. In *International Conference on Machine Learning (ICML)*, pages 5213–5225.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., et al. (2021b). Advances and open problems in federated learning. *Foundations and Trends in Machine Learning*, 14(1–2):1–210.
- Kandpal, N., Pillutla, K., Oprea, A., Kairouz, P., Choquette-Choo, C. A., and Xu, Z. (2023). User inference attacks on large language models. *arXiv preprint arXiv:2310.09266*.
- Kasiviswanathan, S. P., Lee, H. K., Nissim, K., Raskhodnikova, S., and Smith, A. (2008). What can we learn privately? In *FOCS*, pages 531–540. IEEE Computer Society.
- Khetan, A. and Oh, S. (2016). Achieving budget-optimality with adaptive schemes in crowdsourcing. In Lee, D., Sugiyama, M., Luxburg, U., Guyon, I., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.
- Kingma, D. P. and Welling, M. (2022). Auto-encoding variational bayes.
- Koga, T., Wu, R., and Chaudhuri, K. (2024). Privacy-preserving retrieval-augmented generation with differential privacy.
- Kong, W. and Syed, U. (2025). A picture’s worth a thousand (private) words: Hierarchical generation of coherent synthetic photo albums. <https://research.google/blog/a-pictures-worth-a-thousand-private-words-hierarchical-generation-of-coherent-synthetic-photo-albums>.

- Kowalczyk, A., Dubiński, J., Boenisch, F., and Dziedzic, A. (2025). Privacy attacks on image autoregressive models. *arXiv preprint arXiv:2502.02514*.
- Kurakin, A., Ponomareva, N., Syed, U., MacDermed, L., and Terzis, A. (2023). Harnessing large-language models to generate private synthetic text. *arXiv preprint arXiv:2306.01684*.
- Labelbox (2023). How to use llms to detect and extract personal data from ai datasets.
- Lee, J., Dai, Z., Ren, X., Chen, B., Cer, D., Cole, J. R., Hui, K., Boratko, M., Kapadia, R., Ding, W., Luan, Y., Duddu, S. M. K., Abrego, G. H., Shi, W., Gupta, N., Kusupati, A., Jain, P., Jonnalagadda, S. R., Chang, M.-W., and Naim, I. (2024). Gecko: Versatile text embeddings distilled from large language models.
- Lee, J., Wang, Y., and Kifer, D. (2015). Maximum likelihood postprocessing for differential privacy under consistency constraints. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 635–644.
- Lepikhin, D., Lee, H., Xu, Y., Chen, D., Firat, O., Huang, Y., Krikun, M., Shazeer, N., and Chen, Z. (2020). Gshard: Scaling giant models with conditional computation and automatic sharding. *CoRR*, abs/2006.16668.
- Lester, B., Al-Rfou, R., and Constant, N. (2021). The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059.
- Li, C., Miklau, G., Hay, M., McGregor, A., and Rastogi, V. (2015). The matrix mechanism: optimizing linear counting queries under differential privacy. *The VLDB journal*, 24:757–781.
- Li, K., Gong, C., Li, Z., Zhao, Y., Hou, X., and Wang, T. (2024). {PrivImage}: Differentially private synthetic image generation using diffusion models with {Semantic-Aware} pretraining. In *33rd USENIX Security Symposium (USENIX Security 24)*, pages 4837–4854.
- Li, X., Tramèr, F., Liang, P., and Hashimoto, T. (2021). Large language models can be strong differentially private learners. *CoRR*, abs/2110.05679.
- Li, X. L., Thickstun, J., Gulrajani, I., Liang, P., and Hashimoto, T. (2022). Diffusion-LM improves controllable text generation. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K., editors, *Advances in Neural Information Processing Systems*.
- Ligett, K., Neel, S., Roth, A., Waggoner, B., and Wu, Z. S. (2017). Accuracy first: Selecting a differential privacy level for accuracy constrained ERM. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 2566–2576.
- Lin, Z., Baltrusaitis, T., and Yekhanin, S. (2025). Differentially private synthetic data via apis 3: Using simulators instead of foundation models.
- Lin, Z., Gopi, S., Kulkarni, J., Nori, H., and Yekhanin, S. (2023). Differentially private synthetic data via foundation model APIs 1: Images. *arXiv preprint arXiv:2305.15560*.
- Liu, H., Jin, W., Karimi, H., Liu, Z., and Tang, J. (2021a). The authors matter: Understanding and mitigating implicit bias in deep text classification. *CoRR*, abs/2105.02778.
- Liu, M. F., Lyu, S., Vinaroz, M., and Park, M. (2023a). Differentially private latent diffusion models. *arXiv preprint arXiv:2305.15759*.
- Liu, T., Tang, J., Vietri, G., and Wu, Z. S. (2023b). Generating private synthetic data with genetic algorithms.
- Liu, T., Vietri, G., Steinke, T., Ullman, J., and Wu, S. (2021b). Leveraging public data for practical private query release. In *International Conference on Machine Learning*, pages 6968–6977. PMLR.
- Liu, T., Vietri, G., and Wu, S. Z. (2021c). Iterative methods for private synthetic data: Unifying framework and new methods. *Advances in Neural Information Processing Systems*, 34:690–702.
- Liu, Y. and Jin, C. (2024). Icgan: An implicit conditioning method for interpretable feature control of neural audio synthesis.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. (2019). Roberta: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692.
- Long, Y., Lin, S., Yang, Z., Gunter, C. A., and Li, B. (2019). Scalable differentially private generative student model via PATE. *CoRR*, abs/1906.09338.
- Ma, J., Dankar, A., Stein, G., Yu, G., and Caterini, A. (2024). Tabpfgn – tabular data generation with tabpfn.
- Maddock, S., Cormode, G., and Maple, C. (2024). Flaim: Aim-based synthetic data generation in the federated setting. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 2165–2176.

- Mattern, J., Jin, Z., Weggenmann, B., Schoelkopf, B., and Sachan, M. (2022). Differentially private language models for secure data sharing. *arXiv preprint arXiv:2210.13918*.
- McKenna, R., Miklau, G., Hay, M., and Machanavajjhala, A. (2018). Optimizing error of high-dimensional statistical queries under differential privacy. *Proceedings of the VLDB Endowment*, 11(10).
- McKenna, R., Miklau, G., and Sheldon, D. (2021a). Winning the nist contest: A scalable and general approach to differentially private synthetic data. *arXiv preprint arXiv:2108.04978*.
- McKenna, R., Mullins, B., Sheldon, D., and Miklau, G. (2022). AIM: an adaptive and iterative mechanism for differentially private synthetic data. *CoRR*, abs/2201.12677.
- McKenna, R., Pradhan, S., Sheldon, D. R., and Miklau, G. (2021b). Relaxed marginal consistency for differentially private query answering. *Advances in Neural Information Processing Systems*, 34:20696–20707.
- McKenna, R., Sheldon, D., and Miklau, G. (2019). Graphical-model based estimation and inference for differential privacy. In *International Conference on Machine Learning*, pages 4435–4444. PMLR.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR.
- McMahan, H. B., Ramage, D., Talwar, K., and Zhang, L. (2018). Learning differentially private recurrent language models.
- McMahan, H. B., Xu, Z., and Zhang, Y. (2024). A hassle-free algorithm for private learning in practice: Don’t use tree aggregation, use blts. *arXiv preprint arXiv:2408.08868*.
- McSherry, F. and Talwar, K. (2007). Mechanism design via differential privacy. In *FOCS*, pages 94–103.
- Meeus, M., Wutschitz, L., Zanella-Béguelin, S., Tople, S., and Shokri, R. (2025). The canary’s echo: Auditing privacy risks of llm-generated synthetic text.
- Microsoft (Accessed 2025). DPSDA: Private evolution: Generating DP synthetic data using APIs and DP queries. <https://github.com/microsoft/DPSDA>.
- Morris, J. X., Zhao, W., Chiu, J. T., Shmatikov, V., and Rush, A. M. (2024). Language model inversion. In *The Twelfth International Conference on Learning Representations*.
- Mullins, B., Fuentes, M., Xiao, Y., Kifer, D., Musco, C., and Sheldon, D. R. (2024). Efficient and private marginal reconstruction with local non-negativity. *Advances in Neural Information Processing Systems*, 37:141356–141389.
- Musco, C., Musco, C., Rosenblatt, L., and Singh, A. V. (2024). Sharper bounds for chebyshev moment matching with applications to differential privacy and beyond. *arXiv preprint arXiv:2408.12385*.
- Müller, A. (1997). Integral probability metrics and their generating classes of functions. *Advances in Applied Probability*, 29(2):429–443.
- Nakkiran, P., Bradley, A., Zhou, H., and Advani, M. (2024). Step-by-step diffusion: An elementary tutorial.
- Namatevs, I., Sudars, K., Nikulins, A., and Ozols, K. (2025). Privacy auditing in differential private machine learning: The current trends. *Applied Sciences*, 15:647.
- Narayanan, A. and Shmatikov, V. (2008). Robust de-anonymization of large sparse datasets. In *2008 IEEE Symposium on Security and Privacy (sp 2008)*, pages 111–125.
- Nasr, M., Hayes, J., Steinke, T., Balle, B., Tramèr, F., Jagielski, M., Carlini, N., and Terzis, A. (2023). Tight auditing of differentially private machine learning. In *Proceedings of the 32nd USENIX Conference on Security Symposium, SEC ’23*, USA. USENIX Association.
- Nasr, M., Rando, J., Carlini, N., Hayase, J., Jagielski, M., Cooper, A. F., Ippolito, D., Choquette-Choo, C. A., Tramèr, F., and Lee, K. (2025). Scalable extraction of training data from aligned, production language models. In *The Thirteenth International Conference on Learning Representations*.
- Neumann, J. v. (1928). Zur theorie der gesellschaftsspiele. *Mathematische Annalen*, 100:295–320.
- Neunhoffer, M., Wu, S., and Dwork, C. (2021). Private post-{gan} boosting. In *International Conference on Learning Representations*.
- Nikolov, A., Talwar, K., and Zhang, L. (2013). The geometry of differential privacy: the sparse and approximate cases. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 351–360.
- Nikolov, A. and Ullman, J. (2021). Open problem - optimal query release for pure differential privacy. DifferentialPrivacy.org. <https://differentialprivacy.org/open-problem-optimal-query-release/>.
- Nissenbaum, H. (2009). *Privacy in Context: Technology, Policy, and the Integrity of Social Life*. Stanford University Press, USA.

- Nissim, K., Raskhodnikova, S., and Smith, A. (2007). Smooth sensitivity and sampling in private data analysis. pages 75–84.
- OpenAI (2022). DALL-E 2 [artificial intelligence system]. <https://openai.com/index/dall-e-2/>.
- Panda, A., Tang, X., Nasr, M., Choquette-Choo, C. A., and Mittal, P. (2025). Privacy auditing of large language models.
- Papernot, N., Abadi, M., Erlingsson, Ú., Goodfellow, I., and Talwar, K. (2017). Semi-supervised knowledge transfer for deep learning from private training data. In *International Conference on Learning Representations*.
- Papernot, N., Song, S., Mironov, I., Raghunathan, A., Talwar, K., and Erlingsson, U. (2018). Scalable private learning with pate. In *International Conference on Learning Representations*.
- Park, J., Choi, Y., and Lee, J. (2024). In-distribution public data synthesis with diffusion models for differentially private image classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12236–12246.
- Pentyala, S., Pereira, M., and De Cock, M. (2024). Caps: Collaborative and private synthetic data generation from distributed sources. *arXiv preprint arXiv:2402.08614*.
- Perez, I. M., Movahedi, P., Nieminen, V., Airola, A., and Pahikkala, T. (2024). Does differentially private synthetic data lead to synthetic discoveries? *Methods of Information in Medicine*, 63(01/02):035–051.
- Pillutla, K., Liu, L., Thickstun, J., Welleck, S., Swayamdipta, S., Zellers, R., Oh, S., Choi, Y., and Harchaoui, Z. (2023). Mauve scores for generative models: Theory and practice. *Journal of Machine Learning Research*, 24(356):1–92.
- Pillutla, K., Swayamdipta, S., Zellers, R., Thickstun, J., Welleck, S., Choi, Y., and Harchaoui, Z. (2021). Mauve: Measuring the gap between neural text and human text using divergence frontiers. *Advances in Neural Information Processing Systems*.
- Pisier, G. (1980-1981). Remarques sur un résultat non publié de b. maurey. *Séminaire Analyse fonctionnelle (dit "Maurey-Schwartz")*, pages 1–12.
- Ponomareva, N., Hazimeh, H., Kurakin, A., Xu, Z., Denison, C., McMahan, H. B., Vassilvitskii, S., Chien, S., and Thakurta, A. G. (2023). How to dp-fy ml: A practical guide to machine learning with differential privacy. *Journal of Artificial Intelligence Research*, 77:1113–1201.
- Putta, P., Steele, A., and Ferrara, J. W. (2023). Differentially private conditional text generation for synthetic data production.
- Qardaji, W., Yang, W., and Li, N. (2013). Differentially private grids for geospatial data. In *ICDE*, pages 757–768.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., and Sutskever, I. (2021). Learning transferable visual models from natural language supervision. *CoRR*, abs/2103.00020.
- Radford, A., Narasimhan, K., Salimans, T., and Sutskever, I. (2018). Improving language understanding by generative pre-training.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.
- Ramesh, K., Gandhi, N., Madaan, P., Bauer, L., Peris, C., and Field, A. (2024). Evaluating differentially private synthetic data generation in high-stakes domains. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, page 15254–15269. Association for Computational Linguistics.
- Ritter, M., Indyk, I., Singh, A., Audibert, A., Seelam, A., Hanes, C., Lau, E., Olesiak, J., Kang, J., and Wu, X. (2023). Grain - feeding jax models.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. (2021). High-resolution image synthesis with latent diffusion models.
- Rosenblatt, L., Herman, B., Holovenko, A., Lee, W., Loftus, J., McKinnie, E., Rumezhak, T., Stadnik, A., Howe, B., and Stoyanovich, J. (2024a). Epistemic parity: Reproducibility as an evaluation metric for differential privacy. *ACM SIGMOD Record*, 53(1):65–74.
- Rosenblatt, L., Howe, B., and Stoyanovich, J. (2024b). Are data experts buying into differentially private synthetic data? gathering community perspectives. *arXiv preprint arXiv:2412.13030*.
- Rosenblatt, L., Liu, X., Pouyanfar, S., de Leon, E., Desai, A., and Allen, J. (2020). Differentially private synthetic data: Applied evaluations and enhancements. *arXiv preprint arXiv:2011.05537*.

- Rosenblatt, L., Lut, Y., Turok, E., Avella-Medina, M., and Cummings, R. (2025). Differential privacy under class imbalance: Methods and empirical insights. *International Conference on Machine Learning*.
- Sajjadi, M. S. M., Bachem, O., Lucic, M., Bousquet, O., and Gelly, S. (2018). Assessing generative models via precision and recall.
- Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., and Chen, X. (2016). Improved techniques for training gans. *Advances in neural information processing systems*, 29.
- Shaikh, K., Kowalczyk, A., Boenisch, F., and Dziedzic, A. (2025). Implementing adaptations for vision autoregressive model.
- Shokri, R., Stronati, M., Song, C., and Shmatikov, V. (2017). Membership inference attacks against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pages 3–18. IEEE.
- Sidorenko, A. and Tiwald, P. (2025). Privacy-preserving tabular synthetic data generation using tabularargn.
- Stadler, T., Oprisanu, B., and Troncoso, C. (2020). Synthetic data - A privacy mirage. *CoRR*, abs/2011.07018.
- Staniszewski, K., Tworkowski, S., Jaszczur, S., Zhao, Y., Michalewski, H., Łukasz Kuciński, and Miłoś, P. (2025). Structured packing in llm training improves long context utilization.
- Steinke, T. and Ullman, J. R. (2016). Between pure and approximate differential privacy. *J. Priv. Confidentiality*, 7(2).
- Sun, P., Jiang, Y., Chen, S., Zhang, S., Peng, B., Luo, P., and Yuan, Z. (2024). Autoregressive model beats diffusion: Llama for scalable image generation.
- Swanberg, M., McKenna, R., Roth, E., Cheu, A., and Kairouz, P. (2025). Is api access to llms useful for generating private synthetic tabular data?
- Takagi, S., Takahashi, T., Cao, Y., and Yoshikawa, M. (2020). P3GM: private high-dimensional data release via privacy preserving phased generative model. *CoRR*, abs/2006.12101.
- Tan, B., Xu, Z., Xing, E. P., Hu, Z., and Wu, S. (2025). Synthesizing privacy-preserving text data via finetuning without finetuning billion-scale llms. In *ICML*.
- Tang, X., Panda, A., Sehwal, V., and Mittal, P. (2023). Differentially private image classification by learning priors from random processes. *Advances in neural information processing systems*, 36:35855–35877.
- Tang, X., Shin, R., Inan, H. A., Manoel, A., Miresghallah, F., Lin, Z., Gopi, S., Kulkarni, J., and Sim, R. (2024). Privacy-preserving in-context learning with differentially private few-shot generation. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*.
- Tao, Y., McKenna, R., Hay, M., Machanavajjhala, A., and Miklau, G. (2021). Benchmarking differentially private synthetic data generation algorithms. *CoRR*, abs/2112.09238.
- TensorFlow Federated Authors (2019). Tensorflow federated: Machine learning on decentralized data. <https://www.tensorflow.org/federated>.
- Thanh-Tung, H., Tran, T., and Venkatesh, S. (2018). On catastrophic forgetting and mode collapse in generative adversarial networks. *CoRR*, abs/1807.04015.
- Tiwald, P., Krchova, I., Sidorenko, A., Vieyra, M. V., Scriminaci, M., and Platzter, M. (2025). Tabularargn: A flexible and efficient auto-regressive framework for generating high-fidelity synthetic data.
- Torfi, A., Fox, E. A., and Reddy, C. K. (2022). Differentially private synthetic medical data generation using convolutional gans. *Information Sciences*, 586:485–500.
- Torkzadehmahani, R., Kairouz, P., and Paten, B. (2019). DP-CGAN: Differentially private synthetic data and label generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0.
- Tramèr, F., Kamath, G., and Carlini, N. (2022). Position: Considerations for differentially private learning with large-scale public pretraining. *arXiv preprint arXiv:2212.06470*.
- Tran, T. V. and Xiong, L. (2024). Differentially private tabular data synthesis using large language models.
- Truda, G. (2023). Generating tabular datasets under differential privacy.
- Ullman, J. R. (2015). Private multiplicative weights beyond linear queries. In Milo, T. and Calvanese, D., editors, *Proceedings of the 34th ACM Symposium on Principles of Database Systems, PODS 2015, Melbourne, Victoria, Australia, May 31 - June 4, 2015*, pages 303–312. ACM.
- Ullman, J. R. and Vadhan, S. P. (2011). Pcps and the hardness of generating private synthetic data. In *TCC*, pages 400–416.

- Vadhan, S. (2017a). *The Complexity of Differential Privacy*, pages 347–450. Springer, Yehuda Lindell, ed.
- Vadhan, S. (2017b). *The Complexity of Differential Privacy*, pages 347–450. Springer International Publishing, Cham.
- Vadhan, S. P. (2017c). The complexity of differential privacy. In *Tutorials on the Foundations of Cryptography*, pages 347–450. Springer International Publishing.
- Vakili, T., Lamproudis, A., Henriksson, A., and Dalianis, H. (2022). Downstream task performance of BERT models pre-trained using automatically de-identified clinical data. In Calzolari, N., Béchet, F., Blache, P., Choukri, K., Cieri, C., Declerck, T., Goggi, S., Isahara, H., Maegaard, B., Mariani, J., Mazo, H., Odijk, J., and Piperidis, S., editors, *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4245–4252, Marseille, France. European Language Resources Association.
- van der Maaten, L. and Hannun, A. (2020). The trade-offs of private prediction. *arXiv preprint arXiv:2007.05089*.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. (2017). Attention is all you need. In *Advances in Neural Information Processing Systems*.
- Vietri, G., Archambeau, C., Aydoore, S., Brown, W., Kearns, M., Roth, A., Siva, A., Tang, S., and Wu, S. Z. (2022). Private synthetic data for multitask learning and marginal queries. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A., editors, *Advances in Neural Information Processing Systems*, volume 35, pages 18282–18295. Curran Associates, Inc.
- Villalobos, P., Ho, A., Sevilla, J., Besiroglu, T., Heim, L., and Hobbhahn, M. (2024). Will we run out of data? limits of llm scaling based on human-generated data.
- Villani, C. (2008). *Optimal Transport: Old and New*. Grundlehren der mathematischen Wissenschaften. Springer Berlin Heidelberg.
- Vinod, V., Pillutla, K., and Thakurta, A. G. (2025). Invisibleink: High-utility and low-cost text generation with differential privacy. *CoRR*, abs/2507.02974.
- Wang, B., Wu, F., Long, Y., Rimanic, L., Zhang, C., and Li, B. (2021a). Datalens: Scalable privacy preserving training via gradient compression and aggregation. *CoRR*, abs/2103.11109.
- Wang, H., Lin, Z., Yu, D., and Zhang, H. (2025). Synthesize privacy-preserving high-resolution images via private textual intermediaries.
- Wang, H., Pang, S., Lu, Z., Rao, Y., Zhou, Y., and Xue, M. (2024a). dp-promise: Differentially private diffusion probabilistic models for image synthesis. In *33rd USENIX Security Symposium (USENIX Security 24)*, pages 1063–1080.
- Wang, H., Sudalairaj, S., Henning, J., Greenewald, K., and Srivastava, A. (2023). Post-processing private synthetic data for improving utility on selected measures. *Advances in Neural Information Processing Systems*.
- Wang, J., Charles, Z., Xu, Z., Joshi, G., McMahan, H. B., Al-Shedivat, M., Andrew, G., Avestimehr, S., Daly, K., Data, D., et al. (2021b). A field guide to federated optimization. *arXiv preprint arXiv:2107.06917*.
- Wang, J., Das, R., Joshi, G., Kale, S., Xu, Z., and Zhang, T. (2022). On the unreasonable effectiveness of federated averaging with heterogeneous data. *arXiv preprint arXiv:2206.04723*.
- Wang, W., Liang, X., Ye, R., Chai, J., Chen, S., and Wang, Y. (2024b). Knowledgesg: Privacy-preserving synthetic text generation with knowledge distillation from server.
- Wen, Y. and Chaudhuri, S. (2024). Batched low-rank adaptation of foundation models. In *The Twelfth International Conference on Learning Representations*.
- Williamson, D. P. and Shmoys, D. B. (2011). *The Design of Approximation Algorithms*. Cambridge University Press, USA, 1st edition.
- Wood, A., Altman, M., Bembenek, A., Bun, M., Gaboardi, M., Honaker, J., Nissim, K., O’Brien, D. R., Steinke, T., and Vadhan, S. (2018). Differential privacy: A primer for a non-technical audience. *Vand. J. Ent. & Tech. L.*, 21:209.
- Wu, R., Guo, C., and Chaudhuri, K. (2023). Large-scale public data improves differentially private image generation quality. *arXiv preprint arXiv:2309.00008*.
- Wu, S., Xu, Z., Zhang, Y., Zhang, Y., and Ramage, D. (2024a). Prompt public large language models to synthesize data for private on-device applications. *Conference on Language Modeling (COLM)*.
- Wu, T., Panda, A., Wang, J. T., and Mittal, P. (2024b). Privacy-preserving in-context learning for large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.

- Xie, C., Lin, Z., Backurs, A., Gopi, S., Yu, D., Inan, H. A., Nori, H., Jiang, H., Zhang, H., Lee, Y. T., et al. (2024). Differentially private synthetic data via foundation model apis 2: Text. *ICML*.
- Xie, L., Lin, K., Wang, S., Wang, F., and Zhou, J. (2018). Differentially private generative adversarial network. *CoRR*, abs/1802.06739.
- Xiong, Y., Wang, R., Cheng, M., Yu, F., and Hsieh, C.-J. (2023). Feddm: Iterative distribution matching for communication-efficient federated learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16323–16332.
- Xu, Z., Collins, M., Wang, Y., Panait, L., Oh, S., Augenstein, S., Liu, T., Schroff, F., and McMahan, H. B. (2023a). Learning to generate image embeddings with user-level differential privacy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7969–7980.
- Xu, Z. and Zhang, Y. (2024). Advances in private training for production on-device language models.
- Xu, Z. and Zhang, Y. (2025). Synthetic and federated: Privacy-preserving domain adaptation with llms for mobile applications.
- Xu, Z., Zhang, Y., Andrew, G., Choquette, C., Kairouz, P., McMahan, B., Rosenstock, J., and Zhang, Y. (2023b). Federated learning of gboard language models with differential privacy. In Sitaram, S., Beigman Klebanov, B., and Williams, J. D., editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, pages 629–639, Toronto, Canada. Association for Computational Linguistics.
- Yang, M., Chi, C.-H., Lam, K.-Y., Feng, J., Guo, T., and Ni, W. (2024). Tabular data synthesis with differential privacy: A survey.
- Yao, D., Pan, W., Dai, Y., Wan, Y., Ding, X., Jin, H., Xu, Z., and Sun, L. (2021). Local-global knowledge distillation in heterogeneous federated learning with non-iid data. *arXiv preprint arXiv:2107.00051*.
- Yoon, J., Jordon, J., and van der Schaar, M. (2019). PATE-GAN: Generating synthetic data with differential privacy guarantees. In *International Conference on Learning Representations*.
- Yousefpour, A., Shilov, I., Sablayrolles, A., Testuggine, D., Prasad, K., Malek, M., Nguyen, J., Ghosh, S., Bharadwaj, A., Zhao, J., Cormode, G., and Mironov, I. (2021). Opacus: User-friendly differential privacy library in pytorch. *CoRR*, abs/2109.12298.
- Yu, D., Kairouz, P., Oh, S., and Xu, Z. (2024). Privacy-preserving instructions for aligning large language models. *Proceedings of the 41th International Conference on Machine Learning*.
- 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. (2021). Differentially private fine-tuning of language models. *CoRR*, abs/2110.06500.
- Yue, X., Inan, H. A., Li, X., Kumar, G., McAnallen, J., Shajari, H., Sun, H., Levitan, D., and Sim, R. (2022). Synthetic text generation with differential privacy: A simple and practical recipe. *arXiv preprint arXiv:2210.14348*.
- Zaharia, M., Xin, R. S., Wendell, P., Das, T., Armbrust, M., Dave, A., Meng, X., Rosen, J., Venkataraman, S., Franklin, M. J., et al. (2016). Apache spark: a unified engine for big data processing. *Communications of the ACM*, 59(11):56–65.
- Zhang, J., Cormode, G., Procopiuc, C. M., Srivastava, D., and Xiao, X. (2017). Privbayes: Private data release via bayesian networks. *ACM Trans. Database Syst.*, 42(4).
- Zhang, J., Xiao, X., and Xie, X. (2016). Privtree: A differentially private algorithm for hierarchical decompositions. In *Proceedings of the 2016 international conference on management of data*, pages 155–170.
- Zhang, Y., Ramage, D., Xu, Z., Zhang, Y., Zhai, S., and Kairouz, P. (2023). Private federated learning in gboard.
- Zhang, Y., Xu, Z., Wu, S., Zhang, Y., and Ramage, D. (2025). Synthesizing and adapting error correction data for mobile large language model applications. *Annual Meeting of the Association for Computational Linguistics (ACL Industry Track)*.
- Zhang, Z., Wang, T., Li, N., Honorio, J., Backes, M., He, S., Chen, J., and Zhang, Y. (2021). {PrivSyn}: Differentially private data synthesis. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 929–946.
- Zhao, Z., Kunar, A., Birke, R., Van der Scheer, H., and Chen, L. (2023). Ctab-gan+: enhancing tabular data synthesis. *Frontiers in Big Data*, 6.
- Zhu, Y. and Wang, Y.-X. (2022). Adaptive private-k-selection with adaptive k and application to multi-label pate. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pages 5622–5635. PMLR.
- Zhu, Z., Hong, J., and Zhou, J. (2021). Data-free knowledge distillation for heterogeneous federated learning. In *International conference on machine learning*, pages 12878–12889. PMLR.