

Individualized PATE: Differentially Private Machine Learning with Individual Privacy Guarantees*

Franziska Boenisch^{†‡}
 franziska.boenisch@vectorinstitute.ai
 Vector Institute
 Toronto, Canada

Christopher Mühl[‡]
 christopher.muehl@fu-berlin.de
 Free University Berlin
 Berlin, Germany

Roy Rinberg
 roy.rinberg@columbia.edu
 Columbia University
 New York, USA

Jannis Ihrig
 jannis.ihrig@fu-berlin.de
 Free University Berlin
 Berlin, Germany

Adam Dziedzic
 adam.dziedzic@utoronto.ca
 University of Toronto, Vector Institute
 Toronto, Canada

ABSTRACT

Applying machine learning (ML) to sensitive domains requires privacy protection of the underlying training data through formal privacy frameworks, such as differential privacy (DP). Yet, usually, the privacy of the training data comes at the cost of the resulting ML models’ utility. One reason for this is that DP uses one uniform privacy budget ϵ for all training data points, which has to align with the strictest privacy requirement encountered among all data holders. In practice, different data holders have different privacy requirements and data points of data holders with lower requirements can contribute more information to the training process of the ML models. To account for this need, we propose two novel methods based on the Private Aggregation of Teacher Ensembles (PATE) framework to support the training of ML models with individualized privacy guarantees. We formally describe the methods, provide a theoretical analysis of their privacy bounds, and experimentally evaluate their effect on the final model’s utility using the MNIST, SVHN, and Adult income datasets. Our empirical results show that the individualized privacy methods yield ML models of higher accuracy than the non-individualized baseline. Thereby, we improve the privacy-utility trade-off in scenarios in which different data holders consent to contribute their sensitive data at different individual privacy levels.

KEYWORDS

Differential Privacy, Machine Learning, Private Aggregation of Teacher Ensembles (PATE), Individualized Privacy

1 INTRODUCTION

Machine learning (ML) is increasingly applied in settings where training data is sensitive. At the same time, training data leakage is ubiquitous [17, 18, 35], motivating approaches that integrate *differential privacy* (DP) [12]. When properly applied, DP guarantees that the amount of sensitive information that the trained ML models can potentially leak at inference time is bounded by the privacy budget ϵ . However, there exist trade-offs between the degree of privacy introduced by DP and the model’s utility (measured as, for example, its accuracy). Furthermore, we observe an important characteristic

Table 1: Number of generated labels by Standard vs our Individualized PATE on the MNIST and SVHN datasets. *D*: distribution of privacy groups (percentage wise), ϵ : privacy budget for a given group, *N*: number of generated labels, and *A*: accuracy of all votes generated by the respective method. The percentages of the three privacy groups are chosen according to [3] (first setup/row) and [31] (second setup/row). We use the standard train-test split for MNIST and train the teacher models using the first 50K train samples while keeping the remaining 10K as the public dataset, and evaluating on the 10K standard test samples. Similarly for SVHN, we split the train set into a public set of 10K samples and use the remaining 63257 train samples as the standard train set. Then, the test set is used for evaluation. We run the experiments three times and report the standard deviation.

		SETUP		PATE	UPSAMPLE	WEIGHT
MNIST	<i>D</i>	34%-43%-23%	N	365±2	1333±1	1312±8
	ϵ	1.0-2.0-3.0	A	97.40	97.20	97.24
MNIST	<i>D</i>	54%-37%-9%	N	361±3	949±18	1894±10
	ϵ	1.0-2.0-3.0	A	97.32	97.18	97.33
SVHN	<i>D</i>	34%-43%-23%	N	90±1	394±5	409±3
	ϵ	1.0-2.0-3.0	A	64.55	64.33	66.40
SVHN	<i>D</i>	54%-37%-9%	N	96±2	284±1	558±1
	ϵ	1.0-2.0-3.0	A	62.90	62.85	65.65

of current ML applications with DP: ϵ is a single parameter that controls the protection level for the entire dataset, even if some data points in it are not sensitive at all. This coarse level of privacy parameterization seems extremely wasteful: intuitively, if large portions of the data need little protection whereas other parts are highly sensitive, then choosing an ϵ tuned to sufficiently protect the sensitive data, and using it to protect *all* data, might unnecessarily penalize the model’s utility.

In addition to different data being inherently more or less sensitive, it is also known that in society, individuals have different attitudes towards privacy protection, and therefore, require their data to be protected at different levels [5, 20]. Since current ML applications under DP only allow for setting a uniform privacy budget ϵ , even when the data holders have different privacy requirements, the privacy budget would always have to be chosen according to

*Accepted at 23rd Privacy Enhancing Technologies Symposium (PETS 2023).

[†]The work was done while the author was at Fraunhofer AISEC.

[‡]Both authors contributed equally.

the individuals with the highest requirements. However, given the privacy-utility trade-off mentioned above, it would be desirable not to always implement the highest privacy protections for all data points. Instead, allowing several individual privacy budgets according to the data holders’ respective preferences can help to better leverage the training data, and increase the utility of the resulting ML model.

While approaches for supporting the specification of individual privacy preferences exist for statistical data analyses with DP [3, 22], to the best of our knowledge, no such frameworks exist in the context of ML. Yet, there is a multitude of applications that already benefit from individualized DP for data analysis, such as smart home [47], smart grid [8], and object localization [11, 44]—underlining the relevance of the topic and the need to extend individualized DP methods to ML. To this end, in this work, we introduce two novel methods (*upsampling* and *weighting*) that extend the Private Aggregation of Teacher Ensembles (PATE) [33] algorithm—one of the standard frameworks to implement DP in ML applications—and support individualized assignment of privacy budgets among the sensitive training data. We first theoretically introduce both our methods and provide a detailed privacy analysis. Then, we experimentally evaluate our methods’ implications on the resulting model utility on the example of the MNIST [25], SVHN [30], and Adult income [23] datasets. In particular, we study how different distributions of individual privacy preferences and respective privacy budgets influence the gained utility. Our experiments highlight that in comparison to the standard PATE approach where the uniform privacy budget is determined by the data point with the highest privacy requirements, our individualized PATE variants generate significantly more labels, and thereby increase utility of the student model. The significant increase of generated labels by our upsampling and weighting method in comparison to standard PATE is visualized in Table 1. The fraction of data points assigned to each of the three privacy groups is specified according to individuals’ preferences observed within society by [5, 20].

In summary, we make the following contributions:

- Introduction of two novel individualized PATE variants;
- Theoretical analysis of the respective privacy bounds;
- Experimental evaluation of utility improvements for the MNIST, SVHN, and Adult income dataset;
- Quantification of the effects of different privacy budget distributions on the gained utility;

Ethical Implications. In general, deciding on an adequate DP budget ϵ in ML applications dealing with sensitive data is a challenging task. This results from real-world implications of concrete values for ϵ being poorly understood. Additionally, even the calculated privacy budgets ϵ for the same application and data might decrease over time, when tighter bounds for their calculation are pushed forward [1]. These inherent difficulties of choosing an adequate ϵ are also faced when assigning individual privacy budgets to data points. In particular, one needs to make sure that no entity training an ML model with individual DP guarantees abuses their power and assigns poor levels of privacy to data that actually requires privacy protection. We, therefore, suggest the use of our new individualized PATE variants in settings that contain a process for obtaining informed consent of the data holders to process their data

at a given privacy level, such as [37]. This process should consist of (1) the identification of the individual privacy preferences [24, 39], (2) the communication of the associated privacy risks and limitations (e.g. [42]), and (3) enabling meaningful decision-making processes by providing information about DP concerning sensitive data disclosure (e.g. [45]). In particular (2) must be implemented in a way that the risks are communicated clearly, such that nudging individuals into giving up their privacy will be prevented. Moreover, we argue that, due to its difficult interpretability, individual data holders should not be in charge of choosing their numeric privacy budget ϵ , but, based on the information on potential risks and benefits decide on an abstract privacy level, such as "high", "average", or "low" [5, 19, 20]. Concrete numeric values ϵ can then be fixed by the regulator or ethics committee in charge depending on the sensitive data itself and the application [7]. We argue that these values should be chosen such that even the lowest privacy budget still offers protection in practice [29]. This approach can be considered as a form of soft-paternalism [2] to protect privacy of the sensitive data held by individuals who are not concerned about the topic.

2 NOTATION & BACKGROUND

We call \mathcal{D} and \mathcal{R} the sets of all possible data points, and all possible processing results that can be produced on them, respectively. Furthermore, two concrete datasets $D, D' \subseteq \mathcal{D}$ are called neighboring (written $D \sim D'$) if D and D' differentiate exactly in one data point. More specifically, they are called neighboring on d (written $D \stackrel{d}{\sim} D'$) if they differ by any but exactly one data point $d \in \mathcal{D}$.

To refer to ϵ , we will use the term *privacy budget* when expressing the privacy preference specified for a (group of) data point(s), and the term *privacy costs* when referring to the proportion of budget being already consumed in a DP-based mechanism.

All *log* values in this work are based on the natural logarithm. Furthermore, $\mathbb{P}[\cdot]$ denotes the probability of an event according to an adequate probability measure, and $\mathbb{E}[\cdot]$ outputs the expected value of a given random variable.

2.1 Differential Privacy

DP formalizes the idea of limiting the influence of individual data points on the results of analyses conducted on a whole dataset. One relaxation of the standard definition of DP is called (ϵ, δ) -DP.

Definition 1 (cf. [13], Def. 2.4). Let $D, D' \subseteq \mathcal{D}$ be two neighboring datasets. Let $M: \mathcal{D}^* \rightarrow \mathcal{R}$ be a mechanism that processes arbitrarily many data points. M satisfies (ϵ, δ) -DP with $\epsilon \in \mathbb{R}_+$ and $\delta \in [0, 1]$ if for all datasets $D \sim D'$, and for all result events $R \subseteq \mathcal{R}$

$$\mathbb{P}[M(D) \in R] \leq e^\epsilon \cdot \mathbb{P}[M(D') \in R] + \delta. \quad (1)$$

Thereby, it expresses the guarantee that a single data point cannot alter the probability of any processing result by a factor larger than $\exp(\epsilon)$. The second parameter δ specifies a small density of probability on which the upper bound does not have to hold.

In ML, data is usually processed multiple times to train a model, e.g. by conducting several training epochs. This process can be considered as a *composition* of mechanisms that each have privacy costs. The following composition theorem states how DP behaves under composition as follows.

Proposition 1 (cf. [14], Thm. 3.16). *Let $\mathcal{R}_1, \mathcal{R}_2$ be two arbitrary result spaces. Let further $M_1: \mathcal{D}^* \rightarrow \mathcal{R}_1, M_2: \mathcal{D}^* \rightarrow \mathcal{R}_2$ be mechanisms that satisfy (ϵ_1, δ_1) - and (ϵ_2, δ_2) -DP, respectively. Then, the composition $M_3(D) \mapsto (M_1(D), M_2(D))$ satisfies $(\epsilon_1 + \epsilon_2, \delta_1 + \delta_2)$ -DP.*

The proof can be found in Appendix B of [14].

2.2 Rényi Differential Privacy

Proposition 1 shows that under composition, (ϵ, δ) -DP quickly leads to a combinatorial explosion of parameters. A smoother composition of privacy bounds can be achieved by using Rényi Differential Privacy (RDP) [27] which is based on the Rényi divergence (see Definition 9 in Appendix A).

Definition 2 (cf. [27], Def. 4). *A mechanism $M: \mathcal{D}^* \rightarrow \mathcal{R}$ satisfies (α, ϵ) -RDP with $\alpha \in \mathbb{R}_+ \setminus \{1\}$ and $\epsilon \in \mathbb{R}_+$ if for all datasets $D \sim D'$ and for all result events $R \subseteq \mathcal{R}$*

$$\mathbb{D}_\alpha [f_{M(D)} \parallel f_{M(D')}] \leq \epsilon. \quad (2)$$

Here, $f_{M(D)}$ and $f_{M(D')}$ are the probability distributions of the results of M on D and D' , respectively.

In Lemma 3, and Lemma 4 in Appendix A, we show the composition and transformation from RDP to DP guarantees, respectively.

2.3 Individualized Differential Privacy

Individualized DP, similar to [3, 15, 22, 43], allows accounting for privacy for data points individually.

Definition 3 (cf. [22], Def. 6). *For any data point $d \in \mathcal{D}$, M satisfies (ϵ_d, δ_d) -DP with $\epsilon_d \in \mathbb{R}_+$ and $\delta_d \in [0, 1]$ if for all datasets $D \stackrel{d}{\sim} D'$, and for all result events $R \subseteq \mathcal{R}$*

$$\mathbb{P}[M(D) \in R] \leq e^{\epsilon_d} \cdot \mathbb{P}[M(D') \in R] + \delta_d. \quad (3)$$

Accounting privacy per data point can also be applied to different DP variants, such as RDP. Properties like composition and transformation apply to RDP analogously to the original concepts.

2.4 PATE

The PATE framework [33] can be used to perform supervised ML with DP guarantees. Therefore, the set of private labeled training data is split among a pre-defined number of so-called *teacher* models and each teacher is trained on their partition of the data. Afterward, the knowledge gained by the teachers from the private training data is transferred to a public so-called *student* model. To do so, the teachers label a public and unlabeled dataset as training data for the student. Privacy protection for the teachers' sensitive training data is obtained by adding DP noise during the labeling process, and by the fact that the student does not get to interact with the sensitive data, but instead uses the public dataset for training. See Figure 4a in the Appendix for an overview of the approach.

The DP noise addition in the labeling process determines the privacy level of PATE. To obtain a label for a public data point, each teacher issues a vote for a specific class. These votes are aggregated with the Gaussian NoisyMax Aggregator as follows:

Definition 4 (cf. [33], Sec. 2.1). *Let \mathcal{X}, \mathcal{Y} be the feature space, and the set of classes corresponding to data distribution, respectively. Further, let $t_i: \mathcal{X} \rightarrow \mathcal{Y}$ be the i -th teacher of a teacher ensemble of*

size $k \in \mathbb{N}$. The vote count $n: \mathcal{Y} \times \mathcal{X} \rightarrow \mathbb{N}$ of any class $j \in \mathcal{Y}$ for any data point $x \in \mathcal{X}$ is:

$$n_j(x) := \sum_{i=1}^k \mathbb{1}(t_i(x) = j). \quad (4)$$

The characteristic function $\mathbb{1}: \{\perp, \top\} \rightarrow \{0, 1\}$ maps 'true' to 1 and 'false' to 0. Note that the vote count depends on the teachers and, therefore, also on their training data.

Definition 5 (cf. [34], Sec. 4.1). *Let n_j be the vote count as defined in Definition 4 for each class $j \in \mathcal{Y}$. Then, the Gaussian NoisyMax (GNMax) aggregation method with parameter $\sigma \in \mathbb{R}_+$ on any data point $x \in \mathcal{X}$ is given by:*

$$\text{GNMax}_\sigma(x) := \arg \max_{j \in \mathcal{Y}} \left\{ n_j(x) + \mathcal{N}(0, \sigma^2) \right\}. \quad (5)$$

The Gaussian noise is sampled from a normal distribution $\mathcal{N}(\mu, \sigma^2)$ with mean $\mu = 0$ and variance σ^2 .

As an extension of the original GNMax Aggregator, Papernot *et al.* [34] proposed the *Confident-GNMax Aggregator*:

$$\max_{j \in \mathcal{Y}} \{ n_j(x) \} + \mathcal{N}(0, \sigma_T^2) > T \quad (6)$$

that only labels data points for which the consensus of the teachers exceeds a pre-defined threshold T . See [34] for a formalization of this idea.

3 RELATED WORK

We present related work on individualized privacy, DP, and privacy-preserving ML techniques.

3.1 Individualizing Privacy

According to [5, 20], there exist at least three different groups of individuals, demanding high, average, and low privacy protection for their data, respectively. These groups are sometimes referred to as *privacy fundamentalists*, *privacy pragmatists* and *privacy unconcerned* [38]. In general and ML-based applications with DP that deal with data of individuals from all three groups, the privacy budget would always have to be chosen according to the privacy fundamentalists' requirements. This can lead to unfavorable privacy-utility trade-offs in the respective application. Hence, it would be desirable to use individualized privacy budgets to improve model utility while complying with each individual's personal privacy requirements.

3.2 Techniques for Individualized DP

Several techniques for implementing individualized privacy guarantees with DP for data analysis outside of the scope of ML have been proposed.

One of the first techniques for individualized DP was proposed by Alagga *et al.* [3]. Their *stretching mechanism* scales data points individually before perturbing them with statistical noise. As a consequence, the DP noise affects each point with individual intensity. Jorgensen *et al.* proposed two additional methods [22]. Their first *sample mechanism* excludes particular data points from being included in the respective data analysis with a probability according to their privacy preferences. Their second *personalized exponential*

mechanism assigns probabilities to processing results according to individual data points’ privacy requirements. These probabilities are then used to randomly select the final processing result for a dataset. Two *partitioning algorithms* were introduced by Li *et al.* [26]. These separate process groups of sensitive data, each with an individual privacy preference. In a similar vein, Niu *et al.* [32] described a utility-aware sub-sampling mechanism to implement individualized DP guarantees. Ebadi *et al.* [15] put forward a *personalized DP* mechanism that relies on excluding data points from the analysis once their respective privacy budgets are exceeded. Their algorithm is designed for live databases in mind where individual data points might not only require individual privacy protection but can also be added to the data analysis at different points in time. As a consequence, each data point also needs individual privacy budget accounting.

Since all the proposed individualized DP mechanisms are designed for privacy-preserving data analysis on datasets and databases, rather than on ML models, they are not directly applicable to our setting. Note, however, that our weighting mechanism is inspired by the stretching mechanism.

3.3 DP Mechanisms for ML

PATE is not the only approach that can be used to apply DP in ML workflows. Another commonly used approach is the *Differentially Private Stochastic Gradient Descent (DP-SGD)* [36]. In DP-SGD, privacy is achieved by first limiting the changes to an ML model that each individual data point can cause. This is done by clipping model gradients on a per-example basis during training. Then, to achieve DP guarantees, noise is added to the gradients before the model parameters are updated with them. Privacy costs of DP-SGD are accounted for through the *moments accountant* [1]. In this approach, multiple moments of the privacy loss random variable are calculated to obtain a DP bound by using the standard Markov inequality.

In the scope of DP-SGD, Feldman and Zrnic [16] proposed an individual per-data point privacy accounting using *RDP filters*. Similarly, Jordon *et al.* [21] personalized the moments’ accountant by dividing it into an *upwards* and a *downwards moments accountant* which are composed to a *personalized moments accountant* to provide data-dependent DP bounds individually per data point. In a similar vein, Yu *et al.* [46] proposed individualized privacy accounting for DP-SGD based on the gradient norms of the individual data points. While both our and these three works aim at improving the privacy-utility trade-offs in ML with DP, their work differs from ours in the problem setting. Our work sets out to address the problem of supporting data holders in *specifying and implementing* their individual privacy preferences, whereas their work aims at *accounting* for per-data point loss incurred during training of the ML model. Therefore, they assign a uniform privacy budget over the whole training dataset and then provide a tighter per-data point analysis of privacy loss. Based on this tighter analysis, data points can be excluded from training once their individual privacy budget is exhausted, while other data points can still be used for further training. So, rather than asking the question that our work is concerned with, namely *What impact does assign individual privacy*

budgets to the training data have on the resulting ML model utility?, they address the question *What privacy loss is incurred to each individual data point by the given algorithm on the given dataset?* As a consequence, while utility gain in their method is solely due to leveraging each data point based on its individual privacy loss, our method can offer an additional utility gain due to supporting individual privacy budgets per data point.

Note that, due to the different structures of the approaches, the individualized privacy accounting of DP-SGD from [16] or [21] cannot be directly applied to PATE. Therefore, our methods extend PATE’s inherent privacy accounting to individualized accounting.

4 INDIVIDUALIZED EXTENSIONS FOR PATE

To implement individual privacy requirements of sensitive training data points, we propose two novel individualized variants of PATE, namely *upsampling* and *weighting*. Each variant modifies the original PATE algorithm in some aspects to provide individualized privacy.

Each of our individualized variants overcomes the limitation of non-individualized PATE where the uniform privacy budget ϵ has to be chosen according to the highest privacy requirement encountered in the sensitive training data. Thereby, our variants allow us to generate more labels than PATE, and to train a student model with higher utility. In the case when all sensitive data points require the same privacy, our individualized PATE variants are equivalent to non-individualized standard PATE.

In this section, we first introduce the ideas behind our variants and then perform an evaluation of their privacy levels. Therefore, we rely on the privacy analysis of the original PATE algorithm [33], and extend it to our individual variants by analyzing the sensitivity of the vote counts. For both our variants, we also propose concrete algorithms illustrating how they can be implemented. Note, however, that these algorithms only represent possible instantiations of the implementations. In general, what the algorithms should ensure is that our variants yield setups in which data points with different privacy budgets exceed their respective budget at approximately the same number of generated labels. This is because label generation in individualized PATE stops once any data point exceeds their privacy budget. In practice, when training with individualized PATE, model owners can simply observe the privacy budget consumption in the labeling process. By identifying the best parameters for each variant, such that the points’ budget is exceeded at approximately the same number of generated labels, the model owner can then make sure that all privacy budgets are fully leveraged, and the highest number of labels is generated. This in turn, leads to the best student model utility.

4.1 Upsampling Mechanism

Our *upsampling* mechanism relies on duplicating sensitive data such that overlapping data-subsets can be allocated to different teachers. Thereby, data with higher privacy budgets is learned by a higher number of teachers. The upsampling mechanism stands in contrast to the original PATE algorithm where *disjoint* data partitions are passed to the teachers. Since data duplicates extend the amount of training data, they allow for two possible modifications of PATE: (1) keeping the number of teachers constant and allocating

more training data to each teacher, or (2) keeping the number of training data points per teacher constant and increasing the number of teachers. Our experimental evaluation indicates that (2) yields a higher utility gain of upsampling PATE. Intuitively, the teachers perform already reasonably well with the initial amount of training data, and allocating more data to them yields only marginal performance gains. In contrast, having more teachers participate in the voting results in more accurate vote counts with less variance due to statistical randomness. As a consequence, we implement upsampling according to (2) with a constant number of training data points per teacher as specified in Algorithm 1. The algorithm ensures that points are duplicated by an integer according to the privacy budget ratios, since only entire data points (and not fractions of a point) can be assigned to a teacher model. See Figure 4b in the Appendix for a visualization of the approach.

Algorithm 1: Prepare training data for teacher models in the **upsampling** method.

Input: Privacy budgets $\{\varepsilon_d\}$ for each data point d , precision $p \in \mathbb{N}$.

Result: Upsampling factor u_d for each data point d .

```

1  $\{\varepsilon_1, \dots, \varepsilon_j\} \leftarrow \text{unique}(\{\varepsilon_d\});$  /*Get unique budgets*/
2 for Each  $\varepsilon_j$  do
3    $\bar{\varepsilon}_j \leftarrow \varepsilon_j \cdot 10^p;$  /*Upscale budgets*/
4 end
5  $D \leftarrow \text{Greatest Common Divisor}(\bar{\varepsilon}_1, \dots, \bar{\varepsilon}_G);$ 
6 for Each  $\bar{\varepsilon}_j$  do
7    $u_d \leftarrow \frac{\bar{\varepsilon}_j}{D};$ 
8 end
```

We call the PATE aggregator for our upsampling approach *upsampling GNMax* (*uGNMax*). It applies the *upsampling vote count* which is defined as follows:

Definition 6 (Upsampling Vote Count). Let $t_i: \mathcal{X} \rightarrow \mathcal{Y}$ be the i -th out of $k \in \mathbb{N}$ teachers. Let further $N \in \mathbb{N}$ be the number of sensitive data points and $m_i \in \{0, 1\}^N$ a mapping that describes which points are learned by t_i . The upsampling vote count $\tilde{n}: \mathcal{Y} \times \mathcal{X} \rightarrow \mathbb{N}$ of any class $j \in \mathcal{Y}$ for any data point $x \in \mathcal{X}$ is

$$\tilde{n}_j(x) := \sum_{i=1}^k \mathbb{1}(t_i(x) = j) . \quad (7)$$

Although the definition for the upsampling vote count looks the same as the non-individualized vote count (Definition 4), their sensitivities differ due to data points to be learned by several teachers (see Proposition 2 in Section 4.3.2).

4.2 Weighting Mechanism

Our *weighting* mechanism modifies the aggregation of teacher votes. It does so by weighting individual teachers' votes higher or lower depending on their training data points' privacy requirements. Therefore, sensitive data points that have the same privacy budget ε_j , which we call a privacy group g_j , have to be allocated to the same teacher(s). In Algorithm 2, we present how weights w_i

can be assigned to the teachers. A visualization of the weighting mechanism is provided in Figure 4d in the Appendix.

Algorithm 2: Assign weights to teacher models in the **weighting** method.

Input: Privacy budget ε_j and number of teachers n_j for each privacy group g_j , $j \in 1, \dots, G$, and total number of teachers k .

Result: Weight w_i for each teacher t_i .

```

1  $\mathcal{E} \leftarrow \sum_{j=1}^G \varepsilon_j;$ 
2 for Each privacy group  $g_j$  do
3    $\bar{\varepsilon}_j \leftarrow \frac{\varepsilon_j}{\mathcal{E}};$  /*Relative privacy budget*/
4    $\bar{n}_j \leftarrow \frac{n_j}{k};$  /*Relative group size*/
5    $\bar{w}_j \leftarrow \bar{\varepsilon}_j \cdot \bar{n}_j;$ 
6  $\mathcal{W} \leftarrow \sum_{j=1}^G \bar{w}_j;$ 
7 for Each privacy group  $g_j$  do
8    $w_j \leftarrow \frac{\bar{w}_j}{\mathcal{W}} \cdot k;$  /*Make sum of weights match  $k$ */
9   for Each teacher  $t_i$  with data from  $g_j$  do
10     $w_i \leftarrow w_j;$ 
```

We call the aggregation method of this PATE variant *weighting GNMax* (*wGNMax*). Its vote count mechanism is defined as follows:

Definition 7 (Weighting Vote Count). Let $t_i: \mathcal{X} \rightarrow \mathcal{Y}$ be the i -th out of $k \in \mathbb{N}$ teachers. Let further $N \in \mathbb{N}$ be the number of sensitive data points and $m_i \in \{0, 1\}^N$ a mapping that describes which points are learned by t_i . Moreover, let $w_i \in \mathbb{R}_+$ be the weight of t_i for all $i \in \{1, \dots, k\}$. The weighting vote count $\tilde{n}: \mathcal{Y} \times \mathcal{X} \rightarrow \mathbb{N}$ of any class $j \in \mathcal{Y}$ for any unlabeled public data point $x \in \mathcal{X}$ is

$$\tilde{n}_j(x) := \sum_{i=1}^k w_i \cdot \mathbb{1}(t_i(x) = j) . \quad (8)$$

As a particular variant of the weighting-mechanism, we also evaluate cases where some teachers have a zero-weight during some votings. We call this variant the **Vanishing Mechanism**. Intuitively, individualized privacy guarantees in the vanishing-result from teachers contributing their information to more or less voting processes, depending on their data points' lower or higher privacy requirements, respectively. See Appendix B for details on the vanishing mechanism and its privacy assessment. However, our experimental evaluation highlights that this approach, in general, yields low utility. We suspect that this is due to the resulting reduced size of the teacher ensemble.

4.3 Privacy Evaluation

The privacy calculation of individualized PATE differs from the standard (non-individualized) PATE in that it is done for particular data points or groups of data points separately, rather than for the whole dataset. We first introduce the general elements of the privacy analysis for the standard PATE which is shared by our two novel variants. Then, we evaluate the individualized privacy guarantees of each variant depending on its vote count and aggregation mechanism. Table 2 summarizes our two methods, their differences, and their respective privacy guarantees.

Variant	Manipulation	Distributed	Privacy-budget	Sensitivity	Parameter changes	RDP privacy bound
Upsampling	dataset	no	per data-point d	u_d (how often d is upsampled)	k, σ, σ_T, T scaled according to u_d	$(\alpha, (u_d)^2 \cdot \alpha/\sigma^2)$
Weighting	teacher aggregation	yes	per teacher i	w_i (weight of teacher i)	N/A	$(\alpha, (w_i)^2 \cdot \alpha/\sigma^2)$

Table 2: Summary of our individualized PATE variants. The table shows the properties of and privacy guarantees achieved by our mechanisms. Manipulation: what part of standard PATE is adapted; Distributed: mechanism suitable when data is distributed over different parties; Privacy-budget: how fine-grained can individual privacy budget be assigned; Sensitivity: sensitivity for teacher voting; Parameter changes: what parameters of standard PATE need to be adapted; RDP privacy bound: loose bound for privacy calculation. Calculation of both variants’ tight bound is shown in Corollary 1.

4.3.1 *Privacy Evaluation of Standard PATE.* A key element of privacy calculation in PATE is the aggregation mechanism. PATE’s GNMax Aggregator is a function of a Gaussian mechanism.

Definition 8 (cf. [14], Sec. 3.5.3). Let $f: \mathcal{D}^* \rightarrow \mathbb{R}^z$ with $z \in \mathbb{N}$ be any real-valued function and let $\sigma \in \mathbb{R}_+$ be any positive real. Then, the Gaussian mechanism of f with standard deviation σ is

$$M_{f,\sigma}(x) := f(x) + \mathcal{N}(0, \sigma^2). \quad (9)$$

Note: the same random noise is added to $f(x)$ in each dimension.

Gaussian mechanisms have RDP costs depending on σ .

Lemma 1 (cf. [27], Prop. 7). Let $\sigma \in \mathbb{R}_+$ and let $f: \mathcal{D}^* \rightarrow \mathbb{R}$ be a real-valued function with sensitivity $\Delta_f := \max_{D \sim D'} \|f(D) - f(D')\|_2$.

Then, the Gaussian mechanism $M_{f,\sigma}$ satisfies $(\alpha, \Delta_f^2 \cdot \alpha/2\sigma^2)$ -RDP for all $\alpha \in \mathbb{R}_+ \setminus \{1\}$.

Lemma 1 is proven in [27]. The resulting RDP costs can be transformed into (ϵ, δ) -DP costs using Lemma 4.

The data-independent *loose bound* privacy costs that arise in PATE are given by:

Lemma 2 (cf. [34], Prop. 8). The GNMax aggregator satisfies $(\alpha, \alpha/\sigma^2)$ -RDP for all $\alpha \in \mathbb{R}_+ \setminus \{1\}$.

The intuition behind it is that in PATE, each data point of the training dataset is learned by exactly one teacher and is potentially able to change this teacher’s vote. Since DP guarantees are expressed for neighboring datasets that differ in exactly one data point d , in the worst case, d changes the vote count for two classes (reduce one class count by one, and increase another class count by one). Thus, a teacher voting can be considered as the composition of two Gaussian mechanisms each with sensitivity $\Delta_f = 1$ and parameter σ equal to the standard deviation of the Gaussian noise. Putting the standard deviation of one into Lemma 1, and applying composition of two Gaussian mechanisms, this yields the term specified in the loose bound.

In addition, it is also possible to obtain a tighter data-dependent bound for privacy estimation in PATE as defined in [34]. See Lemma 5 in Appendix A for a definition of this *tight bound*.

4.3.2 *Privacy Evaluation of Individualized PATE.* All our new aggregation mechanisms apply individualized vote counts $\bar{n}: \mathcal{Y} \times \mathcal{X} \rightarrow \mathbb{N}$ whose sensitivities are no longer $\Delta_f = 1$, but are determined individually for particular data points (or groups of data points). Therefore,

in the privacy analysis, we need to calculate their privacy bounds based on the mechanisms’ individual sensitivities and the general privacy bounds of PATE.

The individual sensitivity of any function $f: \mathcal{D}^* \rightarrow \mathbb{R}^z$ with $z \in \mathbb{N}$ regarding any data point $d \in \mathcal{D}$ can be defined as $\Delta_{f,d} := \max_{D \sim D'} \|f(D) - f(D')\|_2$. The following propositions formalize the individual sensitivity of the vote counts in our individualized PATE mechanisms.

Proposition 2 (Upsampling Sensitivity). Let $d \in \mathcal{D}$ be any sensitive data point. Let $u_d \in \mathbb{N}$ be the number of duplicates of d (incl. the original d). Then, the individual sensitivity of the vote count, regarding d , in upsampling PATE is

$$\Delta_{\text{upsampling},d} = u_d. \quad (10)$$

PROOF. In upsampling PATE, every teacher that is trained on data point $d \in \mathcal{D}$ can have a different vote for neighboring datasets that differ in d . For each duplicate of d , this results in an increase of one vote count and a decrease of another one. Let $t_{(d)}$ be the set of teachers trained on d . Assume that all u_d votes of $t_{(d)}$ would have changed if d were different. From the perspective of d , the voting can then be considered as a composition of $2 \cdot |\mathcal{Y}|$ Gaussian mechanisms (some might have a sensitivity of zero *s.t.* they have no privacy costs). For each class $j \in \mathcal{Y}$ there are two Gaussian mechanisms, one with sensitivity equal to the number of votes of $t_{(d)}$ for j if d were changed, the other if d were not changed. Applying Lemma 2 and Lemma 3 yields a sum of RDP values, each dependent on its specific sensitivity. Since the sensitivity has a quadratic impact on the RDP costs of a Gaussian mechanism, votes for the same class are more expensive than votes for different classes (see Lemma 1). Therefore, the RDP costs are the highest if all u_d teachers trained on d would consent on a class j when trained on d and would consent on class $j' \neq j$ if d would be different. \square

To perform the privacy analysis in the framework of PATE, let N and N' be the numbers of sensitive data points and the number of the upsampled data points, respectively. Then we can define the relative upsampling of training data as $u := N'/N$. Since we keep the number of data points per teacher constant, the number of teachers k has to be scaled by u . The remaining PATE hyperparameters: σ (for GNMax from Equation (5)), σ_T , and T (for Confident GNMax from Equation (6)) are scaled by u as well to achieve a comparable

voting accuracy and privacy efficiency as for the standard (non-individualized) PATE.

Proposition 3 (Weighting Sensitivity). *Let $d^{(i)} \in \mathcal{D}$ be a sensitive data point learned by teacher $t_i \in \{t_1, \dots, t_k\}$. Let w_i be the weight to determine the influence of t_i to votings. Then, the individual sensitivity of the weighting vote count, regarding $d^{(i)}$, is:*

$$\Delta_{\text{weighting},d}^{(i)} = w_i. \quad (11)$$

PROOF. In weighting PATE, every data point only influences one teacher. Therefore, on neighboring datasets, every vote count might change by the corresponding teacher’s weight w_i . \square

Note that the weighting approach does not change PATE hyper-parameters (σ , σ_T , and T). Nonetheless, sensitive data has to be grouped budget-wise before being provided to the teachers. The teachers are then given weights according to the budgets *s.t.* all weights sum up to the number of teachers k .

4.3.3 Privacy Bounds. Based on the mechanisms’ sensitivity, we can formulate the loose bound of our individualized aggregation mechanisms as follows:

Theorem 1 (Individual Loose Bound). *Let M be an individualized GNMax aggregator with noise scale $\sigma \in \mathbb{R}_+$. Let further $d \in \mathcal{D}$ be any data point, and $\Delta_{M,d}$ be the individual sensitivity of M ’s individualized vote count regarding d . Then, M satisfies an individual $(\alpha, (\Delta_{M,d})^2 \cdot \alpha/\sigma^2)$ -RDP regarding d for all $\alpha \in \mathbb{R}_+ \setminus \{1\}$.*

PROOF. Individualized GNMax aggregators can be considered as the composition of all classes’ vote counts regarding each data point. Only two of them can be changed at the same time on neighboring datasets. Thus, the two Gaussian mechanisms with an individual sensitivity per data point are composed. Therefore, the claimed RDP guarantee is achieved by using Lemma 1 on privacy guarantees of Gaussian mechanisms, and Lemma 3 from Appendix A on composition. Note that in the upsampling variant, more than two vote counts can be changed. The worst case occurs if all teachers affected by data point d change the same vote counts. This is because votes for that same class are more expensive than votes for different classes (see the proof of Proposition 2). \square

We can also compute the data-dependent tight bound (Lemma 5 in Appendix A) for our individualized PATE variants. PATE’s calculation of the tight bound builds on the loose bound, and is calibrated for a sensitivity of 1 for the specified noise scale σ . However, the sensitivity and the noise scale applied by PATE are related. Therefore, when providing a different sensitivity than 1 to the tight bound calculation, it suffices to re-scale σ according to that sensitivity. Our sensitivity values directly correspond to the parameters of our variants of PATE (upsampling duplication factors, participation frequencies in vanishing, and teachers’ weights in the weighting method).

Corollary 1 (Scaling Invariance of the Individual Loose Bound). *Let $c \in \mathbb{R}_+$ be any positive scalar. Let M be an individualized GNMax aggregator with noise scale $\sigma \in \mathbb{R}_+$ and an individual sensitivity $\Delta_{M,d} \in \mathbb{R}_+$ for some data point $d \in \mathcal{D}$. Furthermore, let \tilde{M} be another individualized GNMax aggregator with noise scale $\tilde{\sigma} = c \cdot \sigma$ and*

individual sensitivity $\Delta_{\tilde{M},d} = c \cdot \Delta_{M,d}$ regarding d . Then, M and \tilde{M} have the same individual loose bound regarding d for any $\alpha \in \mathbb{R}_+ \setminus \{1\}$.

PROOF. Fix $\alpha \in \mathbb{R}_+ \setminus \{1\}$. M, \tilde{M} satisfy individual (α, ε) - and $(\alpha, \tilde{\varepsilon})$ -RDP, respectively, regarding d . The equality of ε and $\tilde{\varepsilon}$ is verified by direct computation as follows:

$$\begin{aligned} \tilde{\varepsilon} &:= \left(\Delta_{\tilde{M},d}\right)^2 \cdot \alpha/\tilde{\sigma}^2 \\ &= (c \cdot \Delta_{M,d})^2 \cdot \alpha/(c \cdot \sigma)^2 \\ &= c^2 \cdot (\Delta_{M,d})^2 \cdot \alpha/c^2 \cdot \sigma^2 \\ &= (\Delta_{M,d})^2 \cdot \alpha/\sigma^2 \\ &=: \varepsilon \end{aligned} \quad (12)$$

\square

Note that all data points from the same privacy group share the same sensitivity, and, thereby, also have the same tight bound.

5 EXPERIMENTAL SETUP

In this section, we describe the setup for the empirical evaluation of our individualized PATE variants. Over all experiments, we use the Confident-GNMax algorithm from [34], where the privacy protection is ensured by Gaussian noise within PATE, and labels are only produced if a consensus among the teachers is reached. To isolate the performance-gain of our individualized PATE variants, we do not perform additional methods to improve utility of the student model from previous PATE papers, such as virtual adversarial training [28] or MixMatch [6]. Foregoing these methods allows us for a direct and more precise comparison between standard PATE and our new variants of the framework. However, as a consequence, our reported student accuracies cannot be compared to the accuracies reported in [34]. Therefore, as a baseline to compare our individualized variants, we implement standard PATE within our framework following [34]. Our framework includes Gaussian PATE (GNMax, Confident-GNMax, Interactive-GNMax), our proposed individualized variants, and the support for experimentation is implemented using Python (version 3.8) [41]. Our code can be accessed online.¹

5.1 Datasets and Models

We conduct the experiments presented in this section on the MNIST [25] and the Adult income dataset [23]. MNIST consists of 70,000 (28×28)-pixel gray-scale images depicting handwritten digits for classification. We scale the pixel values of all images to range $[0, 1]$. The Adult income dataset contains 48,842 tabular data points from the US census of the year 1994. The corresponding classification task is to predict if the yearly income of a person represented in the data is greater than \$50k. As a pre-processing of the data, we remove 3,620 damaged data points from the dataset and transformed categorical features into numerical values. Furthermore, we normalize these numerical values to the range of zero to one.

To train the teacher and student models on MNIST, we use a simple convolutional neural network (CNN) architecture taken from [9] (see Table 3). All weights in the output layer are initialized by values randomly sampled from the Glorot uniform distribution,

¹ <https://github.com/fraboeni/individualized-pate>

whereas all other weights are sampled from the He uniform distribution. Optimization is performed using the Adam optimizer and categorical cross-entropy loss. All other parameters are set according to the default values from TensorFlow (version 2.4.1).

For the Adult income dataset, the teacher and student models are implemented as random forest models from the scikit-learn library. Each random forest consists of 100 decision trees. Otherwise, the default parameters of the library are applied.

layer	type of layer	parameters	activation
1	convolutional	32 (3, 3)-kernels	ReLU
2	batch normalization	-	-
3	max pooling	size (2, 2)	-
4	flatten	-	-
5	fully connected	100 nodes	ReLU
6	batch normalization	-	-
7	fully connected	10 nodes	softmax

Table 3: CNN-architecture for MNIST.

Since, as done in standard PATE [33], every teacher is provided with only 240 data points for training on MNIST, we apply a custom data augmentation within each individual teacher and the student to improve model performances. Therefore, each data point within one model’s training data is randomly rotated by up to $\pm 7.5^\circ$ and randomly shifted by up to 7% both, in horizontal and vertical directions to make a larger training dataset for that model. This data augmentation does not influence the privacy costs since we augment the data points only within their respective model’s training dataset and not over different datasets. As a consequence, augmented data points are solely used to train the same model as their original data point. Since PATE is already based on the assumption that each single data point can completely determine the behavior of a corresponding teacher model (*cf.* Lemma 2 and Lemma 5), no additional DP costs are incurred by augmenting datapoints. For the experiments based on the Adult income data, no data augmentation is applied since it does not yield any performance benefits.

5.2 Evaluation Metrics

To measure the utility of the different PATE variants and privacy budget distributions, we mainly track three metrics. First, we count the *number of produced labels* until any of the specified privacy budgets is exhausted. Second, we measure the *accuracy of the student model* trained on that resulting labeled data. Additionally, we also analyze the *accuracy of the generated labels* (“voting accuracy”). As baselines to compare our individualized methods to, we conduct experiments with standard non-individualized PATE and Confident-GNMax using as the dataset-wide ϵ the *minimum* privacy budget encountered in the sensitive data. This has to be done in order not to violate any training data point’s privacy requirements.

5.3 PATE Experiments

To experimentally evaluate our two novel individualized PATE variants, we carry out the empirical analysis in four steps. (1) At first, the complete dataset is randomly divided into private, public, and test partitions. Note that for increased randomization over the experiments, we do not rely on the standard train-test split in MNIST, but instead combine all 70,000 data points and then

partition the dataset. The sizes of these partitions as well as general parameters for Confident-GNMax and its individualized variants on both datasets are described in Table 4, where we follow the setup from PATE [34] in terms of the number of data points per set. The parameters are adopted in the upsampling mechanism so that teacher accuracies and voting accuracies align with those of weighting and non-individualized experiments. (2) Privacy budgets are randomly assigned to the private data according to a given privacy budget distribution. Afterward, the data is allocated to the corresponding teacher models for training. (3) The trained teachers are used to produce labels in the voting process. Aggregation of the teacher votes is conducted according to the PATE variant under evaluation. As a baseline to evaluate our two variants, we use the standard non-individualized Confident-GNMax. After every voting, the current accumulated RDP costs of data points are computed and stored group-wise. For *upsampling*, all data points that share the same number of duplicates have the same privacy costs whereas in *weighting*, all data points that are learned by teachers of the same weight exhibit the same privacy costs. We consider all-natural RDP α values from 2 to 50. These RDP costs are transformed into standard DP costs by taking the best α at that point of the voting. After 2,000 produced labels, the voting process is terminated since we observe that all experiments could exhaust their privacy budget within that number. Tracking privacy costs above the actual budget exhaustion up to the fixed number of 2,000 generated labels is done to compare the privacy costs are spent over many votings. (4) The student model is trained on the labeled data that the respective teacher ensemble produced until any private data point’s privacy budget is exceeded. To get more reliable results, we average our measurements in all following experiments over multiple runs for the same parameters with the different random initialization, and for data shuffling and noise invoked.

5.3.1 Uniform Assignment of Privacy Budgets. We conduct our first set of experiments on both datasets with various privacy budget distributions. We use two privacy groups in the experiments reported in this section as a micro-analysis to clearly show the differences between our variants, and to avoid a combinatorial explosion of privacy budgets and distributions being depicted. We assign one of two different budgets (a higher and a lower budget) to every data point in the private dataset at random. We vary the ratio of data points having the higher budget among 25%, 50%, and 75%. The lower budget is set to $\log 2 \approx 0.69$ over all experiments while we assign the higher budget from $\log 4$, $\log 8$, and $\log 16$. Using logarithmic values provides a more intuitive comparison among the privacy budgets since the formulation of DP (Definition 1) uses $\exp(\epsilon)$. Hence, an $\epsilon = \log s$ for any real $s \geq 1$ is half of a privacy budget $\epsilon' = \log 2s$. For example with our chosen budgets, a budget of $\epsilon = \log 8$ is four times as high as a budget of $\epsilon' = \log 2$. The data and resulting labels that are produced until any data point’s privacy budget is exhausted are used to train the student models.

5.3.2 Non-Uniform Assignment of Privacy Budgets. In the previous experiment, we assign data points randomly to the given privacy groups and their respective privacy budgets. However, this might not necessarily reflect real-world use-cases where individuals’ privacy requirements can correlate with their characteristics, for example which class they belong to. In individualized PATE, a data

dataset	# teachers	# data	private	public	test	σ_T	σ	T	δ
MNIST	250		60,000	9,000	1,000	150	40	200	10^{-5}
Adult	250		37,222	7,000	1,000	200	40	300	10^{-5}

Table 4: PATE Parameters. Parameters used in the experiments for the Confident-GNMax on the MNIST and Adult income datasets. σ is the standard deviation of the noise-induced to the label aggregation. σ_T specifies the standard deviation of noise used to check if the teachers have a consensus given the threshold T .

point’s privacy budget determines how much information that data point can contribute to the voting. As a consequence, when individuals with specific characteristics, or individuals from a specific class have much higher or much lower privacy requirements than other individuals, this can introduce biases to the generated labels, and thereby, also to the student model.

In this experiment, we, therefore, evaluate how the performance of our individualized PATE variants is influenced when the privacy budget distributions vary significantly between different classes. To do so, we use the Adults income dataset, which has an unbalanced class distribution (incomes lower than \$50k make 75.2% of the dataset). We assign the higher privacy budget solely to the underrepresented high-income class to determine to what extent this shifts the trained student model’s predictions. The higher privacy budgets are again set to log 4, log 8, and log 16, the lower budget to log 2. We vary the ratio of data in the underrepresented class that receives the higher privacy budget among 25%, 50%, 75%, and 100%. To produce more reliable results for each budget and ratio combination, we train ten teacher ensembles, use each of them for five voting processes, and report the average. The data and corresponding labels that are produced until any data point’s privacy budget is exceeded represent the student model’s training data.

We restrict ourselves to the upsampling mechanism for our evaluation, as in weighting PATE, teachers are trained on data points with the same privacy budget. Since, in this experiment, we assign privacy budgets according to the classes, most teachers would be trained on data from solely one class. This would result in poor teacher performance.

6 EMPIRICAL RESULTS

We present the quantitative effects of individualization in PATE based on the number of produced labels and the accuracy of the generated student models (as our evaluation metrics described in Section 5.2).

6.1 Advantage of Individualization

To better understand how much privacy the generation of labels consumes on both data groups (lower and higher privacy), we track the privacy costs over the course of generating 2,000 labels on both datasets for our novel PATE mechanisms. Figure 1 showcases the continuous privacy costs of generating labels for the MNIST dataset using the upsampling mechanism over the respective data groups. Therefore, 50% of the data points (randomly chosen) are assigned the lower privacy budget of log 2. The remaining 50% are assigned log 8. As a baseline, we plot the continuous privacy costs for standard PATE. To evaluate the number of labels that can actually be generated for the given privacy budget distribution, we have to count how many labels are returned before any data

point’s privacy budget is exceeded. In Figure 1, this corresponds to the moment when either the lower costs reach the lower budget or the higher costs reach the higher budget, whatever happens first. In the setup depicted in Figure 1, our individualized PATE is able to generate more than three times the number of labels generated by standard PATE (890 vs 257).

For more extensive results on the privacy budget consumption for label generation on MNIST and Adult, see Figure 5 and Figure 6 in Appendix E, respectively. The resulting numbers of produced labels over the experiments are shown in Table 5 for the MNIST dataset and in Table 7 in the Appendix for Adult income.

In the results, we observe several different trends. First, when analyzing the lines corresponding to the privacy costs in Figure 5 and Figure 6, we find that both lines differ more, the more the individual budgets differ. This effect also increases when the proportion of sensitive data with a higher budget decreases. Second, with an increasing ratio of the higher budget, both costs grow slower, resulting in more generated labels, see Table 5. Thereby, the utility advantage of our individualized PATE over the non-individualized standard variant becomes visible. For half of the sensitive MNIST data having a budget of log 4, log 8, or log 16 and the other half having log 2, 492, 890, and 1239 labels can be produced by our weighting mechanism, respectively, instead of 257 in the case of non-individualized PATE. This leads to a student accuracies of 93.08%, 94.68%, and 96.32% while non-individualized PATE only achieves 88.7%. Analogously, on Adult income, in the same privacy budget configuration, 203, 349, and 530 labels can be produced, leading to student accuracies of 81.76%, 82.60%, and 82.84%, respectively for weighting. Standard non-individualized PATE, instead, produces 88 labels so that the student only achieves an accuracy of 79.85%. For a detailed overview on the final students’ accuracies for MNIST and Adult with the different individualized variants and privacy budget distributions, see Table 8 and Table 9 in Appendix E.

Our results are not directly comparable to those [34] since, in contrast to their work, we do not apply virtual adversarial training but only use the public data. Their final models’ accuracy is 98.5% on MNIST with $\epsilon = 1.97$ while our student model never surpassed 97% even for a privacy budget of $\epsilon = \log 16 \approx 2.77$ on 75% of the data. We decided not to integrate the adversarial training method in order to study the pure effect of our individualization and exclude any other effects on the resulting model utility.

However, our individualization still outperforms [34] when it comes to the voting accuracy, *i.e.* the proportion of correctly generated labels: Our generated labels are more accurate ($\approx 97.7\%$) than theirs (93.18%) when evaluating them against ground truth, which is partly due to the better accuracy of our teachers. Our teachers achieve an average test accuracy of 90.2% (81.7% on Adult) on average while theirs are at 83.86% (83.18% on Adult).

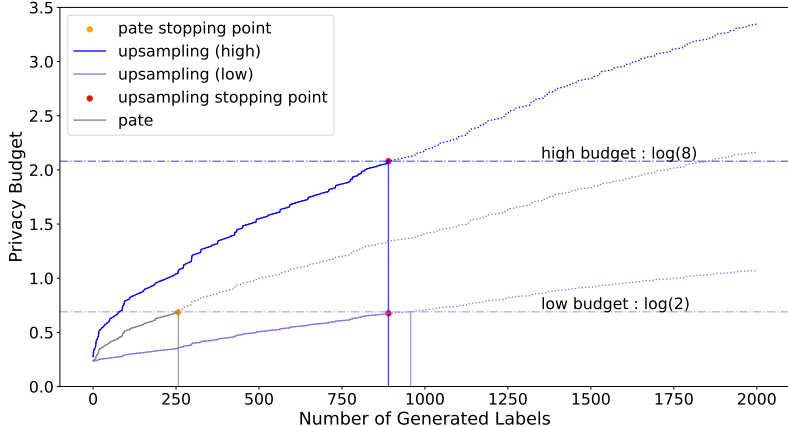


Figure 1: Individualized upsampling outperforms standard PATE. We compare the number of generated labels for given privacy budgets: $\log 2$ and $\log 8$. The lines for upsampling (high) and (low) and PATE (denoted as pate) represent the privacy costs during the generation of the first 2,000 labels for the MNIST dataset. Standard PATE generates only 257 labels before the low privacy budget of $\log 2$ is reached (pate stopping point). The upsampling method exhausts its privacy budgets much later than PATE since it takes advantage of different privacy budgets. The upsampling (high) exhausts its budget (of $\log 8$) at 890 labels (upsampling stopping point). The upsampling (low) exhausts its budget (of $\log 2$) at 957 labels. Overall, the upsampling method exhausts its budget (upsampling stopping point) at the minimum of the two privacy groups ($\log 2$ and $\log 8$), therefore returning 890 generated labels, which is more than 3 times of the number of labels returned by standard PATE. For the upsampling method, we have to adjust its parameters (the number of times the data points from different privacy groups are upsampled) so that the different privacy budgets are exhausted at approximately the same number of generated labels.

Note that the budget combinations $\log 4$ with 75% and $\log 8$ with 25% yield the same average privacy budget over the entire dataset. Nevertheless, the experiment on distribution $\log 4$ with 75% yields more labels and higher accuracy than that on $\log 8$ with 25%, see Table 5. This might indicate that having a smaller gap between the lower and the higher privacy budget leads to increased performance and that it might be better to have more data points with slightly higher privacy budgets than a few data points with very high privacy budgets.

higher budget in ϵ	25% ratio		50% ratio		75% ratio	
	U	W	U	W	U	W
$\log 4$	158	433	237	492	326	564
$\log 8$	231	474	414	890	636	1163
$\log 16$	308	648	623	1239	1038	1787
baseline	257					

Table 5: Number of labels returned by individualized PATE (Upsampling and Weighting). Non-individualized GMax using the minimum budget of $\log 2$ serves as baseline. The voting accuracies for all methods are $\approx 97.7\%$.

6.2 Generated Labels as a Function of Privacy and Relative Group Size

We observe that there are two main factors that allow the individualized PATE algorithm to increase the number of labels that are generated: the number of individuals that have a larger privacy

budget, and the actual size of the non-minimum privacy budget. Either increasing the number of individuals that have a larger privacy budget or increasing the larger privacy budgets, allows our individualized PATE to incur a smaller privacy cost on the most privacy-conscious group.

We run an experiment measuring the number of generated labels for a series of budget combinations $((1., 2.), (1., 3.)\dots)$, and for the group distributions $((25\%, 75\%), (50\%, 50\%), (75\%, 25\%))$. We find that the relationship between the contributions of these two is in fact linear. Scaling up the number of individuals that have a large privacy budget, while equivalently scaling down the privacy budget of that group, keeps the number of generated labels roughly equivalent; and vice-versa. We find this effect to be significant for both our upsampling and weighting mechanism. A more detailed analysis of how to select a scaling, given a privacy-ratio, when group size is fixed is given in Figure 3.

Note that we implement our individualization through the algorithms from Section 4. Using a different approach to implement our variants would change the curve. Hence, Figure 2 allows us to assess the selection of hyperparameters (upsampling factors and teachers' weights) for our variants of PATE. The lower the curve-of-best-fit for an algorithm is, the better it is at utilizing differences in privacy budgets, to generate more labels.

6.3 Non-Uniform Privacy Budgets

The experiment in this section serves to evaluate the influence of assigning a higher privacy budget only to (parts of) the underrepresented class in the Adult income dataset. We analyse the effects

higher budget in ϵ	25% ratio		50% ratio		75% ratio		100% ratio	
	low	high	low	high	low	high	low	high
log 4	95.93	36.77	93.11	45.93	90.13	54.74	86.24	63.39
log 8	93.25	45.91	86.82	62.25	80.57	72.61	77.91	77.36
log 16	90.42	54.00	80.74	72.68	76.68	79.42	73.82	83.59
baseline	(98.01, 24.78)							

Table 6: Per-group student accuracy (in %) for unbalanced individualization (Adult). Results reported for low-income (majority) and high-income (underrepresented) class as an average over 50 different students trained through five voting processes by ten teacher ensembles using upsampling. The ratios specify what proportion of the underrepresented class was assigned the higher privacy budget. The remaining data obtained a privacy budget of log 2. Non-individualized experiments with all data points having a privacy budget of log 2 serve as the baseline. As the proportion of data points from the underrepresented class that obtain a higher privacy budget (or their respective budget) increases, we observe an increase in student accuracy on this class. At the same time, the student accuracy on the majority class decreases.

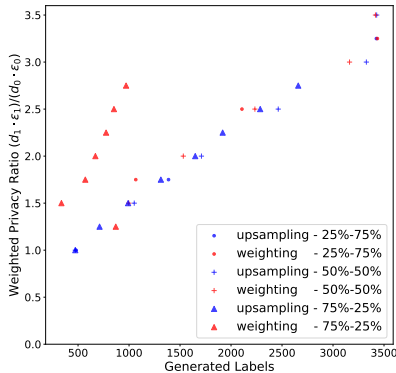


Figure 2: Weighted ratio of generated labels. The y-axis is the weighted privacy ratio, where the contribution of the privacy groups is scaled by the group size, for two different privacy groups this is $(\epsilon_1 \cdot d_1)/(\epsilon_0 \cdot d_0)$. There is a nearly-linear relationship between the number of generated labels, and the weighted privacy ratio.

on both the teachers and the student model and on the generated labels. Table 6 highlights how assigning higher privacy budgets to different proportions of the underrepresented class causes changes in the resulting student models’ predictions. We observe that with larger proportions of data from the underrepresented class that receive a high privacy budget, the resulting student model’s accuracy on this underrepresented class increases. At the same time, the student model’s accuracy on the majority class decreases significantly. The same holds when increasing the privacy budget on the underrepresented class. Table 10a and Table 10b in Appendix E show similar trends for the average teacher and voting accuracies, respectively. These observations indicate that the higher the privacy budget of the underrepresented class (*i.e.*, the lower their privacy requirement), or the more data from the underrepresented class requires lower privacy protection, the more frequently that class gets predicted. This highlights that an individualized privacy budget assignment is able influence the model predictions, and thereby, to enforce and to mitigate biases in the resulting ML models. When it

comes to label generation for the non-uniform privacy-budget assignment (see Table 11 in the Appendix E), we observe an increase in the number of generated labels depending on the fraction of underrepresented data that is assigned a higher privacy budget and the respective budgets. We also compare the the number of labels generated in this setup with the number of labels generated in the random privacy budget assignment (see Table 7 in Appendix E). The rightmost column of Table 11 (100% of the underrepresented class obtain the higher privacy budget) can be directly compared with the leftmost column in Table 7 where 25% of the overall data obtain a higher budget. This is because the underrepresented class represents roughly 25% of the data. We observe that the non-uniform assignment of privacy budgets yields fewer labels. Hence, we can conclude that even when the performance of PATE on the class that receives a higher privacy budget increases, the overall performance decreases.

7 DISCUSSION AND FUTURE WORK

This section discusses results and implications of this work, and provides an outlook on possible future research directions.

7.1 Improving Utility with Individualization

Particularly in sensitive domains such as health care, applying high utility ML models is crucial. This is because incorrect model decisions can have catastrophic consequences. Since introducing DP into training often yields decreased ML model utility, many parties still entirely forego its adaptation within their sensitive ML applications, or they assign a very high privacy budget for all data points, which results in low privacy protection.

The introduction of our individualized PATE yields multiple benefits in these scenarios. First of all, our methods allow integration of PATE in a system where individual data holders can choose what privacy level they want their data to be treated with. That option alone might make individuals more willing to share their data, which would result in the availability of more training data for the ML models. This, in turn, is known to have positive effects on the model utility when training with DP [40]. Additionally, our experiments highlight that our individualized PATE variants yield more generated labels and higher student model utility than

standard PATE which has to comply with the most strict privacy requirements encountered in the training dataset.

7.2 Comparison of our Variants

In the practical comparison of our two PATE variants, we see that the upsampling and weighting variants constantly outperform standard PATE. Additionally, both variants have different benefits and use-cases: The upsampling approach offers high flexibility in terms of individual privacy budget preferences. In theory, each data point could require a different privacy budget and would just have to be upscaled accordingly. In practice, upsampling can increase the computational costs of PATE significantly since more training data is available and more teacher models need to be trained. Moreover, the upsampling method cannot be used for distributed scenarios where the sensitive training data belong to different parties, such as hospitals that jointly want to train a student model based on their respective patients’ data. This is because, for upsampling, the sensitive data would have to be shared among the different parties which is usually restricted by privacy regulations. Weighting is well suited for such distributed scenarios because each party can train their own teacher model and assign the weight according to the privacy requirements of their sensitive data. However, to fully leverage the benefits of weighting, teachers must be trained on data points with the same privacy budgets, which reduces flexibility. It is possible to group together data points with different privacy budgets to one teacher within the weighting approach, but then this teacher’s weight and corresponding privacy level must be set to comply with the strictest requirement among all its training data points. Such an assignment results in a waste of privacy budgets among all data points with higher privacy budgets within this teacher.

One great advantage of our individualized PATE variants is that their implementation can be configured such that all data points are able to exhaust their privacy budgets at roughly the same time. This allows to fully make use of each individual data point’s privacy budget, and thereby to fully leverage the sensitive training data in order to produce higher-utility ML models.

7.3 Individualized Privacy and Biases

Our experiments on assigning higher privacy budgets to data points from one particular class highlight that individualized DP guarantees can enforce or mitigate biases in the resulting ML models. More concretely, data with higher privacy budgets has a direct influence on what classes the student model predicts. The higher a data point’s privacy budget, the higher its influence on the model’s prediction. Therefore, whenever assigning individualized privacy budgets, a thorough evaluation of the resulting ML models concerning biases and model fairness needs to be conducted. Once such negative effects are detected, the privacy budgets of the respective (groups of) data points can be scaled down to reduce their influence.

7.4 Outlook and Future Directions

Independent of individualized privacy guarantees, further theoretical research on improving the tight bound analysis in PATE would be helpful to obtain more realistic estimates of the privacy costs during the voting process. The current analysis assumes that each data point can fully change its teacher model’s prediction. In most

scenarios, this assumption is, however, too strong. With a tighter estimate of a data point’s influence, each vote consumes less privacy budget, and as a consequence, more labels can be produced.

Moreover, it would be of interest to study how applications of distributed PATE and similar frameworks, e.g. CaPC [10], can benefit from our individualized aggregation mechanisms. Our mechanisms can be applied there to implement both individual data point privacy requirements, but also different “per-party” requirements. What is more is that in particular for the weighting mechanism, these per-party privacy requirements could be extended to different weighting schemes taking into account, for example, the amount of training data a party holds, how diverse this data is, and how accurate their trained model predicts—assuming these properties can be determined without undermining the privacy guarantees of the system. Such extensions can then support more meaningful cooperative ML model training and yield models of higher utility.

Finally, in this work, we focus on individualized extensions of PATE. Due to its structure, PATE is naturally suited to support different privacy budgets among its training data. Also, data that does not require any privacy protection can directly be leveraged by the framework as public training data for the student model. However, in the future, it would also be of interest to develop extensions of DP-SGD to support individualized privacy guarantees within this framework. Such extensions could, for example, be implemented by sub-sampling the model’s training data points with non-uniform probabilities according to their privacy budgets, or adding DP noise with different magnitudes to different data points’ gradients.

8 CONCLUSION

Preserving privacy for the training data in ML is a crucial topic. Often, this privacy is achieved at the cost of the final model’s utility. To improve the privacy-utility trade-off and to cater to the requirement encountered among all data holders, we propose two novel variants for the PATE algorithm that allow for the use of individualized privacy budgets among the data points. We formally define our variants, conduct theoretical analyses of their privacy bounds, and experimentally evaluate their effect on PATE’s utility for different datasets and different privacy budget distributions within them. Our results show that through individualized PATE, we are able to generate significantly more labels in comparison to standard PATE which has to comply with the highest privacy requirements encountered in its training dataset. The increased amount of labels also translates into significant improvements in the student model’s accuracy. Our individualized PATE variants are, therefore, able to reduce the loss of utility that is usually introduced by DP.

ACKNOWLEDGMENTS

This work is supported by the German Federal Ministry of Education and Research (grant 16SV8463: WerteRadar).

REFERENCES

- [1] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. 308–318.
- [2] Alessandro Acquisti. 2009. Nudging privacy: The behavioral economics of personal information. *IEEE security & privacy* 7, 6 (2009), 82–85.
- [3] Mohammad Alaggan, Sébastien Gambs, and Anne-Marie Kermarrec. 2015. Heterogeneous differential privacy. *arXiv preprint arXiv:1504.06998* (2015).
- [4] Borja Balle, Gilles Barthe, and Marco Gaboardi. 2018. Privacy Amplification by Subsampling: Tight Analyses via Couplings and Divergences. In *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2018/file/3b5020bb891119b9f5130f1fea9bd773-Paper.pdf>
- [5] Bettina Berendt, Oliver Günther, and Sarah Spiekermann. 2005. Privacy in e-commerce: Stated preferences vs. actual behavior. *Commun. ACM* 48, 4 (2005), 101–106.
- [6] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin Raffel. 2019. MixMatch: A Holistic Approach to Semi-Supervised Learning. [arXiv:1905.02249](https://arxiv.org/abs/1905.02249) [cs.LG].
- [7] Karthik S Bhat and Neha Kumar. 2020. Sociocultural Dimensions of Tracking Health and Taking Care. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW2 (2020), 1–24.
- [8] Arpan Bhattacharjee, Shahriar Badsha, and Shamik Sengupta. 2021. Personalized privacy preservation for smart grid. In *2021 IEEE International Smart Cities Conference (ISC2)*. IEEE, 1–7.
- [9] Jason Brownlee. 2019. How to Develop a CNN for MNIST Handwritten Digit Classification. URL <https://machinelearningmastery.com/...>
- [10] Christopher A Choquette-Choo, Natalie Dullerud, Adam Dziedzic, Yunxiang Zhang, Somesh Jha, Nicolas Papernot, and Xiao Wang. 2021. CaPC Learning: Confidential and Private Collaborative Learning. *arXiv preprint arXiv:2102.05188* (2021).
- [11] Fatemeh Deldar and Mahdi Abadi. 2019. PDP-SAG: Personalized privacy protection in moving objects databases by combining differential privacy and sensitive attribute generalization. *Ieee Access* 7 (2019), 85887–85902.
- [12] Cynthia Dwork. 2006. Differential privacy. In *International Colloquium on Automata, Languages, and Programming*. Springer, 1–12.
- [13] Cynthia Dwork. 2008. Differential privacy: A survey of results. In *International conference on theory and applications of models of computation*. Springer, 1–19.
- [14] Cynthia Dwork, Aaron Roth, et al. 2014. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science* 9, 3-4 (2014), 211–407.
- [15] Hamid Ebadi, David Sands, and Gerardo Schneider. 2015. Differential privacy: Now it’s getting personal. *Acm Sigplan Notices* 50, 1 (2015), 69–81.
- [16] Vitaly Feldman and Tijana Zmcic. 2021. Individual privacy accounting via a renyi filter. In *Thirty-Fifth Conference on Neural Information Processing Systems*.
- [17] Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. 2015. Model inversion attacks that exploit confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*. 1322–1333.
- [18] Karan Ganju, Qi Wang, Wei Yang, Carl A Gunter, and Nikita Borisov. 2018. Property inference attacks on fully connected neural networks using permutation invariant representations. In *Proceedings of the 2018 ACM SIGSAC conference on computer and communications security*. 619–633.
- [19] Chris Jay Hoofnagle and Jennifer M Urban. 2014. Alan Westin’s privacy homo economicus. *Wake Forest L. Rev.* 49 (2014), 261.
- [20] Carlos Jensen, Colin Potts, and Christian Jensen. 2005. Privacy practices of Internet users: Self-reports versus observed behavior. *International Journal of Human-Computer Studies* 63, 1-2 (2005), 203–227.
- [21] James Jordon, Jinsung Yoon, and Mihaela van der Schaar. 2019. Differentially Private Bagging: Improved utility and cheaper privacy than subsample-and-aggregate. In *Advances in Neural Information Processing Systems*. pp. 4323–4332.
- [22] Zach Jorgensen, Ting Yu, and Graham Cormode. 2015. Conservative or liberal? Personalized differential privacy. In *2015 IEEE 31st international conference on data engineering*. IEEE, 1023–1034.
- [23] Ronny Kohavi and Barry Becker. 1996. Adult data set. *UCI machine learning repository* 5 (1996), 2093.
- [24] Jan Kolter and Günther Pernul. 2009. Generating user-understandable privacy preferences. In *2009 International Conference on Availability, Reliability and Security*. IEEE, 299–306.
- [25] Yann LeCun and Corinna Cortes. 2010. MNIST handwritten digit database. <http://yann.lecun.com/exdb/mnist/>. <http://yann.lecun.com/exdb/mnist/>
- [26] Haoran Li, Li Xiong, Zhanglong Ji, and Xiaoqian Jiang. 2017. Partitioning-based mechanisms under personalized differential privacy. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 615–627.
- [27] Ilya Mironov. 2017. Rényi differential privacy. In *2017 IEEE 30th Computer Security Foundations Symposium (CSF)*. IEEE, 263–275.
- [28] Takeru Miyato, Shin-Ichi Maeda, Masanori Koyama, and Shin Ishii. 2019. Virtual Adversarial Training: A Regularization Method for Supervised and Semi-Supervised Learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41, 8 (2019), 1979–1993. <https://doi.org/10.1109/TPAMI.2018.2858821>
- [29] Milad Nasr, Shuang Songi, Abhradeep Thakurta, Nicolas Papernot, and Nicholas Carlin. 2021. Adversary instantiation: Lower bounds for differentially private machine learning. In *2021 IEEE Symposium on Security and Privacy (SP)*. IEEE, 866–882.
- [30] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. 2011. Reading digits in natural images with unsupervised feature learning. (2011).
- [31] Ben Niu, Yahong Chen, Boyang Wang, Jin Cao, and Fenghua Li. 2020. Utility-aware Exponential Mechanism for Personalized Differential Privacy. In *2020 IEEE Wireless Communications and Networking Conference (WCNC)*. IEEE, 1–6.
- [32] Ben Niu, Yahong Chen, Boyang Wang, Zhibo Wang, Fenghua Li, and Jin Cao. 2021. AdaPDP: Adaptive personalized differential privacy. In *IEEE INFOCOM 2021-IEEE Conference on Computer Communications*. IEEE, 1–10.
- [33] Nicolas Papernot, Martin Abadi, Úlfar Erlingsson, Ian Goodfellow, and Kunal Talwar. 2017. Semi-supervised Knowledge Transfer for Deep Learning from Private Training Data. In *Proceedings of the International Conference on Learning Representations*. <https://arxiv.org/abs/1610.05755>
- [34] Nicolas Papernot, Shuang Song, Ilya Mironov, Ananth Raghunathan, Kunal Talwar, and Úlfar Erlingsson. 2018. Scalable Private Learning with PATE. In *International Conference on Learning Representations (ICLR)*. <https://arxiv.org/abs/1802.08908>
- [35] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2017. Membership inference attacks against machine learning models. In *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE, 3–18.
- [36] Shuang Song, Kamalika Chaudhuri, and Anand D. Sarwate. 2013. Stochastic gradient descent with differentially private updates. In *2013 IEEE Global Conference on Signal and Information Processing*. 245–248. <https://doi.org/10.1109/GlobalSIP.2013.6736861>
- [37] Peter Sörries, Claudia Müller-Birn, Katrin Glinka, Franziska Boenisch, Marian Margraf, Sabine Sayegh-Jodehl, and Matthias Rose. 2021. Privacy Needs Reflection: Conceptual Design Rationales for Privacy-Preserving Explanation User Interfaces. *Mensch und Computer 2021-Workshopband* (2021).
- [38] Humphrey Taylor. 2003. Most people are “privacy pragmatists” who, while concerned about privacy, will sometimes trade it off for other benefits. *The Harris Poll* 17, 19 (2003), 44.
- [39] Maximilian Teltzrow and Alfred Kobsa. 2004. Impacts of user privacy preferences on personalized systems. In *Designing personalized user experiences in eCommerce*. Springer, 315–332.
- [40] Florian Tramèr and Dan Boneh. 2020. Differentially private learning needs better features (or much more data). *arXiv preprint arXiv:2011.11660* (2020).
- [41] Guido Van Rossum and Fred L. Drake. 2009. *Python 3 Reference Manual*. CreateSpace, Scotts Valley, CA.
- [42] Sandra Wachter, Brent Mittelstadt, and Chris Russell. 2017. Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harv. J. L. & Tech.* 31 (2017), 841.
- [43] Yu-Xiang Wang. 2019. Per-instance Differential Privacy. *Journal of Privacy and Confidentiality* 9, 1 (2019).
- [44] Zhibo Wang, Jiahui Hu, Ruizhao Lv, Jian Wei, Qian Wang, Dejun Yang, and Hairong Qi. 2018. Personalized privacy-preserving task allocation for mobile crowdsensing. *IEEE Transactions on Mobile Computing* 18, 6 (2018), 1330–1341.
- [45] Aiping Xiong, Tianhao Wang, Ninghui Li, and Somesh Jha. 2020. Towards Effective Differential Privacy Communication for Users’ Data Sharing Decision and Comprehension. In *2020 IEEE Symposium on Security and Privacy (SP)*. IEEE, 392–410.
- [46] Da Yu, Gautam Kamath, Janardhan Kulkarni, Tie-Yan Liu, Jian Yin, and Huishuai Zhang. 2022. Individual Privacy Accounting for Differentially Private Stochastic Gradient Descent. *arXiv preprint arXiv:2206.02617* (2022).
- [47] Xin-Yuan Zhang, Liu-Sheng Huang, Shao-Wei Wang, Zhen-Yu Zhu, and Hong-Li Xu. 2016. Personalized Differential Privacy Preserving Data Aggregation for Smart Homes. In *3rd International Conference on Wireless Communication and Sensor Networks (WCSN 2016)*. Atlantis Press, 203–209.
- [48] Yuqing Zhu, Xiang Yu, Manmohan Chandraker, and Yu-Xiang Wang. 2020. Private-kNN: Practical Differential Privacy for Computer Vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

A ADDITIONAL BACKGROUND ON RDP AND PATE

This section supplements Section 2 by providing formalizations of Rényi divergence, RDP composition, and the tight bound.

Rényi Differential Privacy

Definition 9 (cf. [27], Def. 3). Let P and Q be two probability distributions over \mathcal{D} . Rényi divergence of order $\alpha \in \mathbb{R}_+ \setminus \{1\}$ for P and Q can be defined as:

$$\mathbb{D}_\alpha [P \parallel Q] := \frac{1}{\alpha - 1} \cdot \log \mathbb{E}_{x \sim Q} \left[\left(\frac{P(x)}{Q(x)} \right)^\alpha \right], \quad (13)$$

where $x \sim Q$ expresses that samples $x \in \mathcal{D}$ follow the probability distribution Q .

Composition under RDP can be expressed as:

Lemma 3 (cf. [27], Prop. 1). Let $\mathcal{R}_1, \mathcal{R}_2$ be arbitrary result spaces. Let further $M_1: \mathcal{D}^* \rightarrow \mathcal{R}_1, M_2: \mathcal{D}^* \rightarrow \mathcal{R}_2$ be mechanisms that satisfy (α, ϵ_1) and (α, ϵ_2) -RDP, respectively. Then, the composition $M_3(D) \mapsto (M_1(D), M_2(D))$ satisfies $(\alpha, \epsilon_1 + \epsilon_2)$ -RDP.

Note that Lemma 3 also holds for adaptive sequential composition as shown in [27].

RDP guarantees can be transformed into DP guarantees as follows:

Lemma 4 (cf. [27], Prop. 3). Let $M: \mathcal{D}^* \rightarrow \mathcal{R}$ be an (α, ϵ) -RDP mechanism. Then, M also satisfies (ϵ', δ) -DP with

$$\epsilon' = \epsilon + \frac{\ln 1/\delta}{\alpha - 1} \quad (14)$$

for all $\delta \in (0, 1]$.

Lemma 4 is proved in [27].

Tight Bound Privacy Analysis of PATE

The tight bound for PATE can be defined as follows.

Lemma 5 (cf. [34], Thm. 6). Let M simultaneously satisfy (α_1, ϵ_1) -RDP and (α_2, ϵ_2) -RDP. Both RDP bounds can be computed by applying the loose bound for two different alpha values. Suppose that $1 \geq q \geq \mathbb{P}[M(D) \neq j^*]$ holds for a likely teacher voting j^* . Additionally suppose that $\alpha \leq \alpha_1$ and $q \leq \exp((\alpha_2 - 1) \cdot \epsilon_2) / \left(\frac{\alpha_1}{\alpha_1 - 1} \cdot \frac{\alpha_2}{\alpha_2 - 1} \right)^{\alpha_2}$. Then, M satisfies (α, ϵ) -RDP for any neighboring dataset D' of D with

$$\epsilon = \frac{1}{\alpha - 1} \cdot \log((1 - q) \cdot A + q \cdot B), \quad (15)$$

where A and B are defined as follows:

$$A := \left(\frac{1 - q}{1 - (q \cdot e^{\epsilon_2})^{\frac{\alpha_2 - 1}{\alpha_2}}} \right)^{\alpha - 1}, \quad (16)$$

$$B := \left(\frac{e^{\epsilon_1}}{q^{\frac{1}{\alpha_1 - 1}}} \right)^{\alpha - 1}. \quad (17)$$

This holds since according to [34], Prop. 7, for a GNMax aggregator M with parameter σ and for any class $j^* \in \mathcal{Y}$ the following

statement applies:

$$\mathbb{P}[M(D) \neq j^*] \leq \frac{1}{2} \sum_{j \neq j^*} \operatorname{erfc} \left(\frac{n_{j^*} - n_j}{2\sigma} \right) \quad (18)$$

where $\operatorname{erfc}(\cdot)$ denotes the complementary error function defined by:

$$\operatorname{erfc}(a) := \frac{2}{\sqrt{\pi}} \int_a^\infty e^{-t^2} dt. \quad (19)$$

See [34] for the proofs.

B THE VANISHING-MECHANISM

Our *vanishing* mechanism keeps the independent partitioning of the original PATE approach and implements individualized privacy by having teachers participate in more or fewer votings according to their training data points' privacy budget. Therefore, in vanishing, data points with the same privacy budget have to be allocated to the same teacher (s). We call data points with the same privacy budget a *privacy group* g_j . Teachers trained on privacy groups with higher privacy requirements (lower budgets) contribute to fewer votings, whereas teachers in lower-requirement groups contribute to more votings. We implement vanishing by randomly sampling teachers for participating in given voting according to their data points' privacy requirements. To be able to apply the same magnitude of privacy noise for each voting, we make sure that the number of teachers sampled per voting stays constant, see Algorithm 3. The vanishing mechanism is also visualized in Figure 4c.

Algorithm 3: Select teacher models for voting in the **vanishing** method.

Input: Privacy budget ϵ_j for each privacy group g_j ,
 $j = 1, \dots, G$, each teacher model t_i .

Result: Participation s_i for each teacher t_i .

```

1 for Each teacher  $t_i$  do
2    $s_i \leftarrow 0$ ; /*Initialize participation*/
3  $\epsilon_{max} \leftarrow \max_{j=1}^G \epsilon_j$ ;
4 for Each privacy group  $g_j$  do
5    $S \leftarrow$  randomly select  $\frac{\epsilon_j}{\epsilon_{max}}$  teachers from group  $g_j$ ;
6   for Each teacher  $t_i$  in  $S$  do
7      $s_i \leftarrow 1$ ; /*Update participation*/
8   end
9 end

```

We call the resulting aggregation method *vanishing GNMax* (vGNMax). Its vote count mechanism can be defined as follows:

Definition 10 (Vanishing Vote Count). Let $t_i: \mathcal{X} \rightarrow \mathcal{Y}$ be the i -th out of $k \in \mathbb{N}$ teachers. Let further $N \in \mathbb{N}$ be the number of sensitive data points and $m_i \in \{0, 1\}^N$ a mapping that describes which points are learned by t_i . Moreover, let $s_i \in \{0, 1\}$ be the current participation of t_i . The vanishing vote count $\hat{n}: \mathcal{Y} \times \mathcal{X} \rightarrow \mathbb{N}$ of any class $j \in \mathcal{Y}$ for any unlabeled public data point $x \in \mathcal{X}$ is defined as

$$\hat{n}_j(x) := \sum_{i=1}^k s_i \cdot \mathbb{1}(t_i(x) = j). \quad (20)$$

For example, if there are 2 groups of teachers where the higher privacy budget is twice the lower privacy budget, then the teachers with the higher privacy budget always vote while the teachers from the lower privacy budget participate only in half of the votings, and the number of teachers for given voting is 3/4 of the total number of teachers.

Privacy Analysis for Vanishing

Proposition 4 (Vanishing Sensitivity). *Let $d^{(i)} \in \mathcal{D}$ be a sensitive data point learned by teacher $t_i \in \{t_1, \dots, t_k\}$. Let $s := (s_1, \dots, s_k) \in \{0, 1\}^k$ be the selection of teachers that participate in the current voting. Then, the individual sensitivity of the vanishing vote count, regarding $d^{(i)}$, is:*

$$\Delta_{\text{vanishing}, d}^{(i)} = s_i. \quad (21)$$

PROOF. In vanishing PATE, every data point only influences the vote of one teacher, which, in the worst-case results in two vote counts being changed. However, in contrast to non-individualized PATE, privacy is only spent if the teacher corresponding to $d^{(i)}$ participates in the current voting. \square

Note that the vanishing mechanism could potentially benefit from the privacy amplification by subsampling [4]. However, how to combine the data-dependent RDP, as used in PATE, with the subsampling mechanism remains an open problem [48] which is outside of the scope of this work.

The vanishing approach does not change PATE hyperparameters (σ , σ_T , and T), in contrast to the upsampling method. The sensitive data has to be grouped budget-wise before being provided to the teachers. The votes are scaled so that the total sum of the votes is equal to the total number of teachers.

C IMPLEMENTATION DETAILS OF THE INDIVIDUALIZED VARIANTS OF PATE

This section contains details on the implementation of the individualized GNMax variants which are left out in Sections 4 and 5 for the sake of brevity.

Hyperparameter Search

The goal of the practical implementation of our variants of PATE is to ensure that their parameters align. Thus, the optimization of PATE hyperparameters, *i.e.* number of teachers k , noise standard deviation for consensus σ_T , threshold for consensus T , and the noise standard deviation for label creation σ , have to be adjusted so that the different variants of PATE are comparable.

We show how the parameters used in variants of individualized PATE: numbers of duplications for upsampling, participation frequencies for vanishing, and teacher weights for weighting, influence the individual loose bound through individual sensitivities. For example, the teachers’ weights for the weighting scheme translate directly to the teacher’s sensitivities, which are set for the privacy analysis. The same holds for the duplication factor in upsampling and the participation frequencies for vanishing.

The privacy costs depend not only on the loose bound but also on the data-dependent tight bound and on the currently optimal RDP order(s). Therefore, it does not suffice to set the individual

parameters or sensitivities proportional to the individual budgets. To find adequate parameters, we conduct experiments to analyze the relation between individual sensitivities and resulting privacy costs in Figure 3. The figure reports results for sensitivities, which are set according to the duplication, participation, and weighting factors of PATE. The goal is to relate individual sensitivities by adjusting parameters so that all privacy budgets exhaust approximately at the same time. To enable comparisons among the different variants, we describe the parameters by corresponding individual sensitivities. We randomly divide the sensitive data into two equally sized groups, one with higher and one with lower individual sensitivity. The parameters are adjusted so that a ratio of c to 1 is achieved for individual sensitivities with each $c \in \{2, \dots, 9\}$. We conduct this experiment on the MNIST dataset and train ten different ensembles. Each ensemble is then used for five different voting processes. We perform 4,000 votings in each process to compare the different cost growths over time.

Details of Upsampling

Upsampling PATE extends the training data for teachers by duplicates. Figure 3 shows that the ratio of privacy costs approaches the ratio of their corresponding budgets after some votings. Therefore, numbers of duplicates should align to the relation of privacy budgets. This can be achieved by initializing the numbers of duplicates as the different budgets and then scaling them up equally until each of them reaches an integer with some desired precision. Note that a higher precision might lead to very high numbers of duplicates if not all budgets are multiples of each other as in our experiments.

Details of Vanishing

Vanishing PATE differentiates privacy on teacher-level by avoiding participation in some votings. Therefore, sensitive data has to be grouped budget-wise and then be given to teachers *s.t.* all data points in a teacher have (almost) the same privacy budget. Afterward, the participation frequencies have to be set according to the lowest budget of each teacher. Figure 3 suggests that the frequencies corresponding to two different budgets should have a relation that is at least quadratic to the relation of their corresponding budgets. In our experiments we used a relation that equals the relation of budgets to the power of four since the costs of data with different budgets are closer after a few votings before they approach constant ratios as for upsampling and weighting. So we initialized frequencies to the corresponding budgets, squared them, and finally divided them by the highest frequency so that frequencies were probabilities and the highest one was 100%.

To be comparable to the other individualized variants, the voting accuracy of vanishing PATE should be retained by decreasing the noise intensity according to the smaller number of voting teachers. A weaker noise entails higher privacy costs for participating teachers’ data. Experiments showed that the privacy costs of all data is lower if the number of participating is stable over the whole voting process. Therefore, our implementation maintains a stable number of participating teachers by selecting random alternations that are changed periodically. More precisely, randomly selected sets of teachers participate periodically in votings where the period aligns to their frequency and equal periods are shifted to achieve a

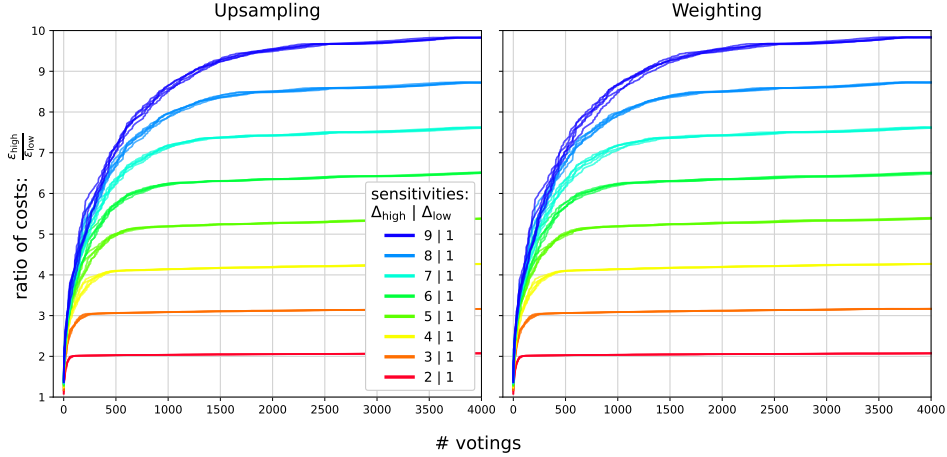


Figure 3: Tuning parameters for individualized variants of PATE. Privacy cost relation between two equal-sized groups of sensitive data (high and low sensitivity) shown over 4,000 votings on the MNIST dataset for both individualized Confidential-GNMax variants. All costs are given in (ϵ, δ) -DP for $\delta = 10^{-5}$. Each of ten teacher ensembles is used to vote five times for all labels in the public dataset that is shuffled differently for each combination. After some votings, the ratio of different costs almost remains constant at the ratio of corresponding sensitivities. Note that the sensitivities shown are only proportions of sensitivities. This means that for upsampling and weighting, the *plain* sensitivities (that correspond to the number of duplications and the teachers’ weights) are scaled while for vanishing, *average* sensitivities (*i.e.* participation frequencies), are scaled and depicted. We use the figure to find what sensitivity values should be set for our two balanced (with the same number of data points) privacy groups. For a given ratio of ϵ -s and the desired number of votings (generated labels), we find the optimal ratios of sensitivities that directly correspond to the hyperparameters of our variants of PATE. For example, if the ratio of the privacy budgets $\frac{\epsilon_{high}}{\epsilon_{low}}$ is 3, then the corresponding ratio of sensitivities for the upsampling method should be 3 (the points from the lower privacy group are not duplicated - sampled once, while the points from the higher privacy group should be duplicated twice - with the total number of 3 points).

stable number of participating teachers per voting. After some votings, new sets of teachers with identical frequencies are randomly sampled to reduce the risk of biases that could be introduced into labels by cliques of teachers with similar knowledge.

Details of Weighting

Following Figure 3, weights can be set to their corresponding budgets divided by the average budget. Thus, all hyperparameters can remain unchanged while their optimization regarding the accuracy of teachers and voting still holds.

Setting Parameters for Individualized Variants

We show how to set the parameters of the individualized variants of PATE so that different privacy budgets are exhausted at approximately the same time. Figure 3 visualizes the relation between the parameters of our individualized Confidential-GNMax variants and the resulting individualized privacy costs according to tight bounds over time. We observe that uGNMax and wGNMax behave very similarly and their cost ratios stay almost constant after a few votings. Contrary, vGNMax needs more votings to lower the gain of its cost ratio. For uGNMax and wGNMax, the cost ratio according to the tight bound seems to be approximately equal to the ratio of sensitivities, whereas the cost ratio approaches the square root of the ratio of sensitivities for vGNMax. Therefore, in our experiments,

we adjust the individualization parameters (duplications, participation frequencies, and weights) so that the resulting sensitivities relate to the actual budgets.

E.g. let $\epsilon_1, \epsilon_2 \in \mathbb{R}_+$ be two DP budgets with $\epsilon_2 = c \cdot \epsilon_1$ for any $c > 1$. Then, for the uGNMax, the duplications u_2 of points having the higher budget ϵ_2 are set to $u_2 := c \cdot u_1$ where u_1 is the number of duplications for points having the lower budget. For vGNMax, the participation frequency s_1 of teachers trained on points having the lower budget ϵ_1 is set to $s_1 := 1/c^2$ while the frequency s_2 of teachers trained on points having the higher budget is always one². Finally, for the wGNMax, the weights of teachers are set to the corresponding budgets and then normalized so that the sum of all weights equals the number of teachers. Thus, $w_1 := \epsilon_1/\bar{w}$ and $w_2 := \epsilon_2/\bar{w}$ where \bar{w} is the average weight of all teachers.

D VISUALIZATION OF METHODS

We visualize the methods in Figure 4.

E ADDITIONAL RESULTS

Results of experiments on the Adult income dataset as well as more comprehensive results of MNIST experiments are presented on the following pages.

²We set $s_1 := 1/c^4$ for vGNMax in our experiments so that different privacy budgets exhausted approximately at the same time.

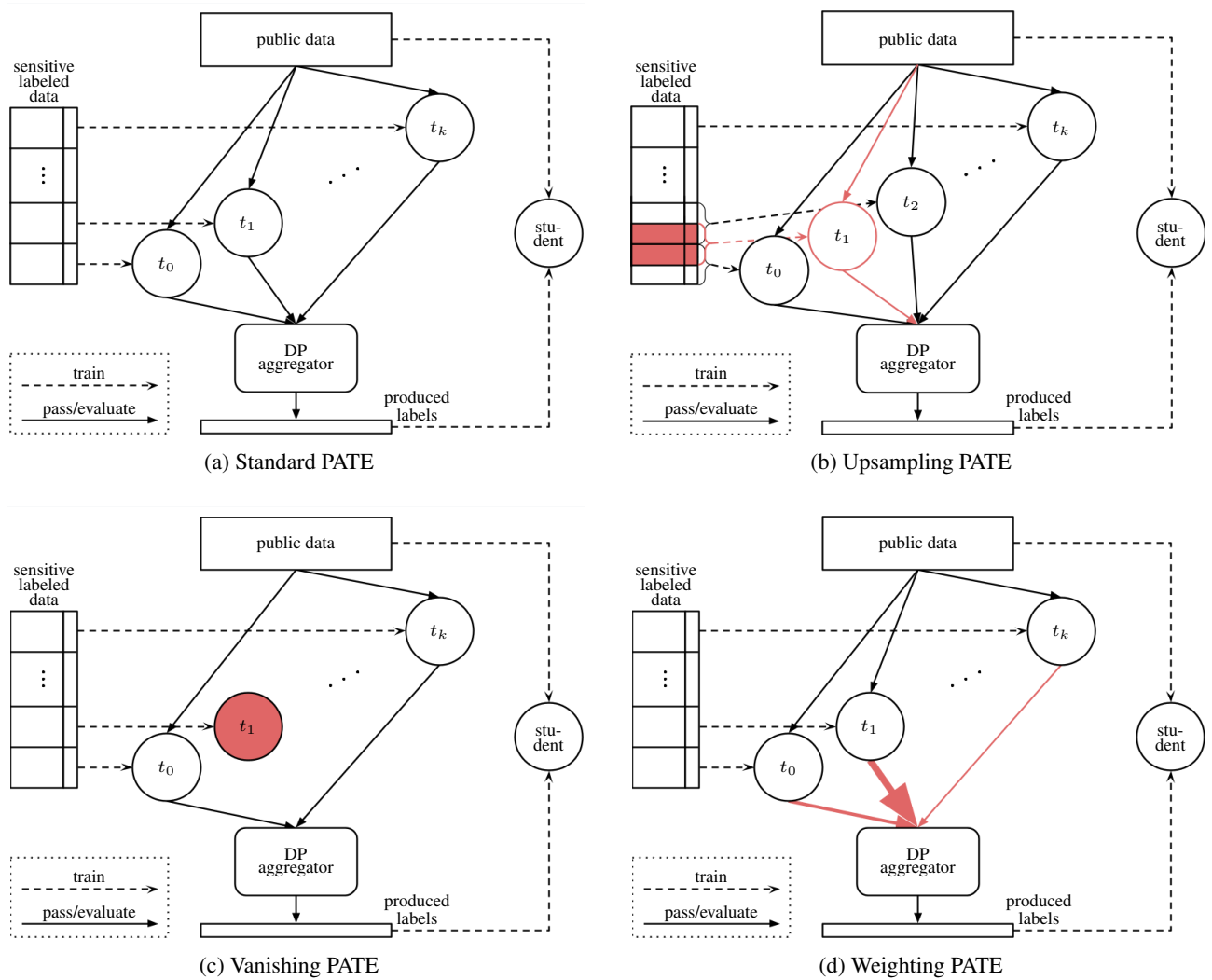


Figure 4: PATE variants. Teachers t_0, t_1, \dots, t_k are trained on partitions of sensitive labeled data. Afterward, public unlabeled data is given to the teacher ensemble whose votes are aggregated *s.t.* labels are produced. Finally, the student is trained on the public data with produced labels. The individualized variants modify this procedure to individually adjust the influence of sensitive data on produced labels. Sensitive data may be used to train multiple teachers (b). Instead, teachers may avoid participating in some votings (c), or be weighted differently (d).

higher budget in ϵ	25% ratio		50% ratio		75% ratio	
	U	W	U	W	U	W
log 4	140	139	202	203	272	273
log 8	198	198	346	349	541	543
log 16	264	259	530	530	868	872
baseline	88					

Table 7: Number of labels generated per individualization (Adult). Results computed over five voting processes for different budget distributions for Upsampling and Weighting. Non-individualized experiments with the lower group’s privacy budget log 2 serve as baselines.

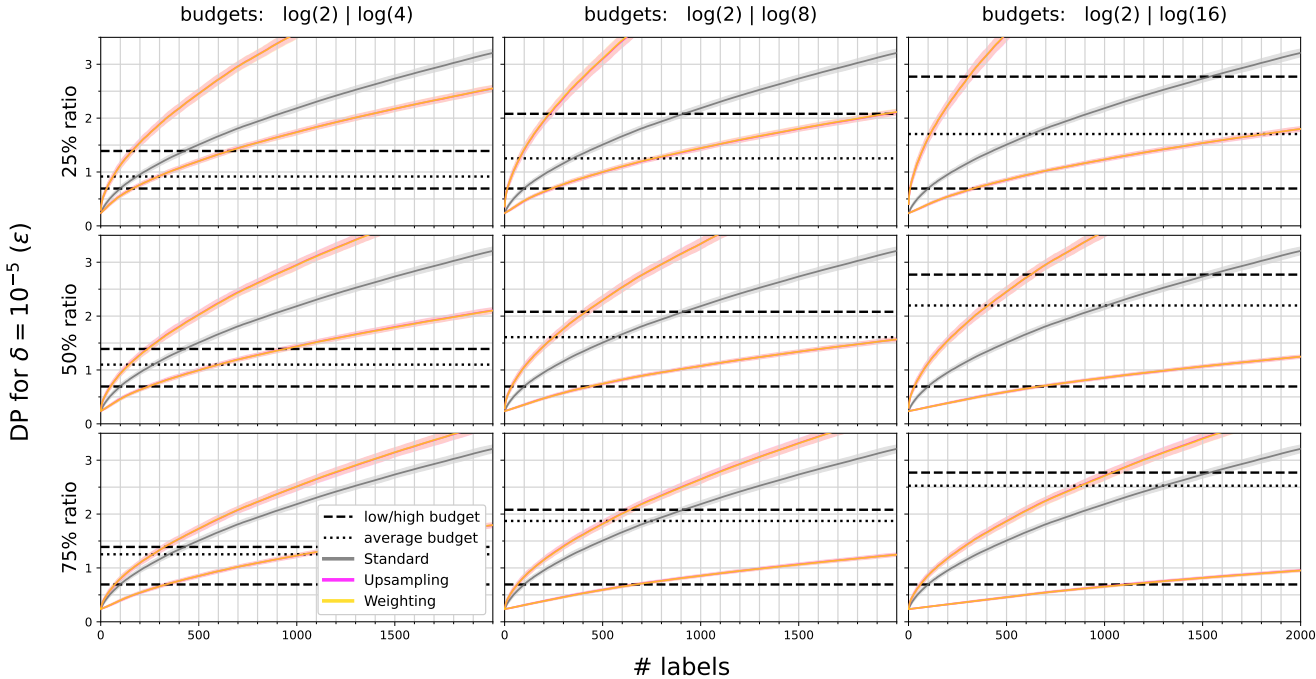


Figure 5: Privacy cost history (MNIST). Costs for generating the first 2,000 labels. Results are averaged over five voting processes by ten teacher ensembles, each for different budget distributions and GNMax variants upsampling and weighting. Privacy costs and budgets are given in (ϵ, δ) -DP for $\delta = 10^{-5}$. Ratios indicate the proportion of data with the higher budget. Costs are listed per group of data points sharing the same budget.

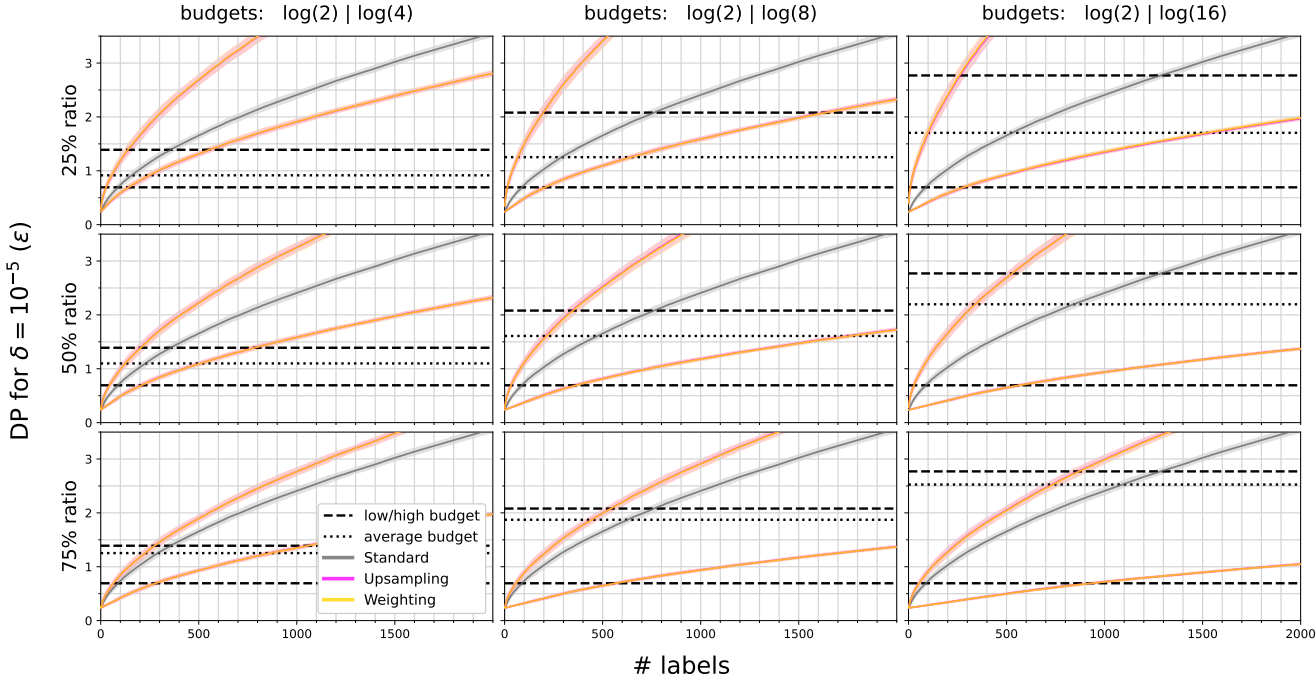


Figure 6: Privacy cost history (Adult). Costs for generating the first 2,000 labels. Results are averaged over five voting processes by ten teacher ensembles, each for different budget distributions and GNMax variants upsampling and weighting. Privacy costs and budgets are given in (ϵ, δ) -DP for $\delta = 10^{-5}$. Ratios indicate proportion of data with the data have the higher budget. Costs are given per group of points that share the same budget.

higher budget in ϵ	25% ratio		50% ratio		75% ratio	
	U	W	U	W	U	W
log 4	92.32	93.26	92.52	93.08	94.38	93.70
log 8	93.12	88.94	94.48	94.68	96.20	95.74
log 16	93.96	90.24	96.38	96.32	96.90	96.60
<i>baseline</i>	88.70					

Table 8: Student accuracy per individualization (MNIST). Results for Upsampling and Weighting based on the generated labels (see Table 5). Non-individualized experiments with the lower group’s privacy budget log 2 serve as a baseline.

higher budget in ϵ	25% ratio		50% ratio		75% ratio	
	U	W	U	W	U	W
log 4	81.02	80.87	81.76	81.76	82.16	82.26
log 8	81.79	81.67	82.52	82.60	82.87	82.89
log 16	82.30	82.25	82.82	82.84	83.07	83.04
<i>baseline</i>	79.85					

Table 9: Student accuracy per individualization (Adult). Results depict the average accuracies, computed over five voting processes for different budget distributions for Upsampling and Weighting. Non-individualized experiments with the lower group’s privacy budget log 2 serve as a baseline.

higher budget in ϵ	25% ratio		50% ratio		75% ratio		100% ratio	
	low	high	low	high	low	high	low	high
log 4	90.08	56.05	87.75	61.70	85.58	66.24	83.62	69.68
log 8	87.74	61.73	83.57	69.84	79.99	75.46	76.86	79.44
log 16	85.58	66.12	79.93	75.48	75.48	81.01	71.57	84.79
<i>baseline</i>	(92.55, 48.82)							

(a) Teacher Accuracy.

higher budget in ϵ	25% ratio		50% ratio		75% ratio		100% ratio	
	low	high	low	high	low	high	low	high
log 4	95.21	55.99	93.02	64.11	91.24	68.93	88.59	74.11
log 8	93.19	64.11	88.71	73.91	84.42	81.14	81.41	85.23
log 16	91.02	68.55	84.47	81.36	79.99	86.49	75.32	90.22
<i>baseline</i>	(97.13, 46.29)							

(b) Voting Accuracy.

Table 10: Per-group teacher and voting accuracy (in %) for unbalanced individualization (Adult). Results reported for low-income (majority) and high-income (underrepresented) class as an average over 50 different students trained through five voting processes by ten teacher ensembles using upsampling. The ratios specify what proportion of the underrepresented class was assigned the higher privacy budget. The remaining data obtained a privacy budget of log 2. Non-individualized experiments with all data points having a privacy budget of log 2 serve as the baseline. As the proportion of data points from the underrepresented class that obtain a higher privacy budget (or their respective budget) increases, we observe an increase in accuracy on this class. At the same time, the accuracy on the majority class decreases.

higher budget in ϵ	# produced labels			
	25% ratio	50% ratio	75% ratio	100% ratio
log 4	90	95	101	109
log 8	93	108	132	162
log 16	96	129	172	225
<i>baseline</i>	88			

Table 11: Labels generated per unbalanced individualization (Adult). Results depict the average over five voting processes by ten teacher ensembles using upsampling. Ratios indicate the proportion of the underrepresented class with the indicated higher budgets. The remaining data receives a privacy budget of log 2. Non-individualized experiments with a uniform privacy budget of log 2 serve as a baseline.