

Learning with Impartiality to Walk on the Pareto Frontier of Fairness, Privacy, and Utility

MOHAMMAD YAGHINI, PATTY LIU, FRANZISKA BOENISCH, and NICOLAS PAPERNOT, University of Toronto & Vector Institute, Canada

Deploying machine learning (ML) models often requires both fairness and privacy guarantees. Both of these objectives present unique trade-offs with the utility (e.g., accuracy) of the model. However, the mutual interactions between fairness, privacy, and utility are less well-understood. As a result, often only one objective is optimized, while the others are tuned as hyper-parameters. Because they implicitly prioritize certain objectives, such designs bias the model in pernicious, undetectable ways. To address this, we adopt impartiality as a principle: design of ML pipelines should not favor one objective over another. We propose impartially-specified models, which provide us with accurate Pareto frontiers that show the inherent trade-offs between the objectives. Extending two canonical ML frameworks for privacy-preserving learning, we provide two methods (*FairDP-SGD* and *FairPATE*) to train impartially-specified models and recover the Pareto frontier. Through theoretical privacy analysis and a comprehensive empirical study, we provide an answer to the question of where fairness mitigation should be integrated within a privacy-aware ML pipeline.

1 INTRODUCTION

From medical applications [13] to infrastructure planning from census data [7], deploying machine learning (ML) models in critical contexts often requires not only utility (accuracy) guarantees, but also fairness and privacy assurances. Prior attempts to mitigate the tension between fairness, privacy, and utility either attempt to adapt private learning for improved trade-offs with fairness [26, 32, 33], or integrate privacy constraints into bias mitigation methods [14].

We argue that such works obscure the complexity of multi-objective decision making under the veil of “tuning:” they optimize for one metric while specifying arbitrary limits on others. This practice raises fundamental issues. Notably, it puts socially-salient choices at the behest of algorithm designers and engineers. This may result in one objective being relegated to secondary consideration—creating potentially dangerous scenarios, such as introducing additional privacy leakage in the attempt to increase model fairness, or degrading fairness by introducing privacy [9, 25].

Instead of reducing trade-offs between different objectives to a weighted combined metric, we propose algorithms that provide a more informative primitive; such as a trade-off function, or its more general counterpart, a *Pareto frontier*. Richer trade-off representations provide improvements on multiple levels. First, on a technical level, we avoid the pitfall of premature aggregation, and can produce solutions at least as good (if not better) than with a combined weighted metric. Second, we improve the transparency of critical ML models by exposing their inherent trade-offs to the decision makers. At this level, these decision makers, e.g., lawyers, humanities’ experts, elected officials, and the public can make informed decisions regarding the appropriate operating point on the Pareto frontier. Importantly, this does not require them having expertise in optimization or ML. They can use the Pareto frontier over different objectives to choose societally favorable, but still technically feasible, specifications.

Let us illustrate these improvements with Section 1. We showcase a typical Pareto frontier calculated using one of our approaches (Section 4.4) on the UTKFace dataset for a binary prediction task which uses a demographic parity fairness mitigation. Consider two scenarios. In scenario 1, given a differential privacy budget of $\epsilon = 3$, a fairness violation of $\gamma = 0.05$ will allow us to answer 84% of queries posed to the model.

Authors’ address: Mohammad Yaghini, mohammad.yaghini@mail.utoronto.ca; Patty Liu, patty.liu@mail.utoronto.ca; Franziska Boenisch, franziska.boenisch@vectorinstitute.ai; Nicolas Papernot, nicolas.papernot@utoronto.ca, University of Toronto & Vector Institute, Toronto, Ontario, Canada.

Now scenario 2, where we relax our fairness violations constraint to $\gamma = 0.1$. This allows us to answer up to 89% of the queries (+5% improvement), and improve accuracy by 1.5%, while maintaining the same privacy guarantee. For one decision maker, the loss in fairness that results in switching from scenario 1 to 2 may be warranted by the gains in other metrics—while this may be completely unacceptable to another decision maker. The choice between these scenarios is not a technical one but may have important societal consequences. Thus, it should be left to the decision maker and not the algorithm designer.

But how do we obtain such Pareto frontiers? We propose to adopt *Impartiality* as a principle. That is, we design ML pipelines that do not favour one objective (e.g., fairness or privacy) to the detriment of another.

We start our pursuit of impartiality by studying the interactions between our notions of privacy (differential privacy [10]) and fairness (demographic parity [8]). In particular, we characterize the tensions between mitigations developed for these notions. We observe that certain mitigations (for instance, fairness pre-processing) can lead to degraded privacy guarantees. We conclude that ML pipelines naively composed without consideration for these tensions are inherently inefficient. In the spirit of impartial designs, we provide two general fairness mitigation strategies that do not cause additional privacy leakage.

We then turn our attention to more bespoke privacy-preserving ML frameworks. We start with the Differentially Private Stochastic Gradient Descent (DP-SGD) [1] and consider how we can integrate a fairness mitigation within it. Our analysis informs our first algorithm contribution, namely, *FairDP-SGD* which incorporates the aforementioned fairness mitigations. Next, we consider the Private Aggregation of Teacher Ensembles (PATE) [22] which, thanks to its modular design, affords us even more freedom in incorporating a fairness mitigation. A thorough privacy-fairness analysis of our options leads us to our second algorithm contribution, *FairPATE*. We evaluate our models against a suite of non-impartial models on multiple datasets. We find that our impartial designs often produce the most efficient results, and therefore, naturally surface the Pareto frontier—representing the irreconcilable trade-offs between various trustworthiness objectives. In summary, our contributions are as follows:

- (1) We argue for the need for richer representations of the multi-objective ML trustworthiness problem in the form of Pareto frontiers. These impartial representations allow us to instantiate the baseline model specification problem, where a regulatory body uses the Pareto frontier to provide specifications for trustworthy measures. Our discussion focuses on two measures: demographic parity violations γ and a differential privacy bound ϵ .
- (2) To achieve the Pareto frontier, we adopt impartiality between our trustworthiness objectives. As a result, we (a) initiate a study to understand the tensions between demographic parity and differential privacy with a focus on the mutual impact of their mitigations and (b) provide fairness mitigation strategies at zero-cost to privacy. Our analysis includes, to the best of our knowledge, the first privacy analysis of a pre-processed fairness mitigation. We then consider bespoke designs for widely-adopted privacy frameworks. Studying limitations of DP-SGD

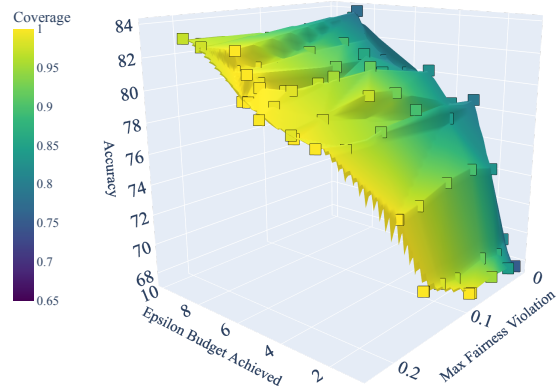


Fig. 1. Percentage of answered queries (coverage) as a function of the maximum privacy budget ϵ and maximum fairness violation γ achieved during training with our FairPATE.

leads us to develop FairDP-SGD. To avoid the limitations of DP-SGD, we consider the PATE framework. After a thorough fairness-privacy analysis we present FairPATE. We provide privacy analysis for both our algorithms.

- (3) We provide interactive¹ Pareto frontiers for various vision tasks, including a medical Chest Xray disease diagnosis dataset (CheXpert) with known fairness issues [3, 24]. Our extensive empirical evaluation compares various in- and pre-processing mitigation techniques. Our results show the inefficiency of non-impartial models: we improve accuracy up to 5% through careful privacy budget consumption model of FairPATE.
- (4) We provide domain-specific best-practices. These not only include fairness and privacy hyper-parameter choices but also, for the first time, the choice of privacy framework and the corresponding point in the ML pipeline at which the mitigation should be inserted (e.g., before or during training).

2 MOTIVATION & BACKGROUND

Acknowledging that machine learning models can pose risks to society—much like pollution-emitting industries pose health risks to the environment—it is a reasonable expectation that a regulatory body should produce technical specifications to curb the societal risks of ML models. This is similar to how the Environment Protection Agency (EPA) in the US, or the European emission standards specify maximum emission ratings.

Problem Setting. Our problem setting is that of *baseline specification*; where a regulatory body needs to decide on the (hyper)parameters for trustworthiness guarantees. Since model builders are always going to optimize for accuracy, the focus of the baseline specification should be on other, societally salient, objectives. In this paper, we focus on fairness and privacy. In particular, we limit our study to the demographic parity notion of fairness guarantees parameterized by $\gamma \in [0, 1]$ which characterizes the maximum tolerable fairness violations. For privacy, we consider differential privacy, parametrized by the privacy budget $\epsilon \in (0, +\infty)$. Ideally, the regulatory body would use a Pareto frontier over the objectives to choose societally favorable, but still technically feasible specifications for their respective parameters.

In the remainder of this section, we provide the necessary background on machine learning, fairness, and privacy. We rigorously define our notions of fairness and privacy in Section 2.1 and Section 2.2, and provide a formal definition of Pareto efficiency in Section 2.3.

ML Background. We assume a classification task where a model $\theta : \mathcal{X} \times \mathcal{Z} \mapsto \mathcal{K}$ maps the features $(\mathbf{x}, z) \in \mathcal{X} \times \mathcal{Z}$ to a label $y \in \mathcal{K}$, where: \mathcal{X} is the domain of non-sensitive attributes, \mathcal{Z} is the domain of the sensitive attribute (as a categorical variable), and \mathcal{K} is the domain of the output label (also categorical). Without loss of generality, we will assume $\mathcal{Z} = [Z]$ (i.e. $\mathcal{Z} = \{1, \dots, Z\}$) and $\mathcal{K} = [K]$.

2.1 Fairness: Demographic Parity

We base our work on the fairness metric of *demographic parity* which requires that ML models produce similar success rates (i.e., rate of predicting a desirable outcome, such as getting a loan) for all sub-populations [8].

We note that in a multi-class setting (i.e., $K > 2$), and even in the binary-class settings where the problem does not admit a reasonable notion of the “desirable outcome”, there can be multiple formulations of the notion of demographic parity (Appendix D). We adopt a natural extension of the well-known binary notion that requires equal rates for any class. Let us first define demographic disparity:

¹Available at https://cleverhans-lab.github.io/impartiality_viz/

The *demographic disparity* $\Gamma(z, k)$ of subgroup z for class k is the difference between the probability of predicting class k for the subgroup z and the probability of the same event for any other subgroup: $\Gamma(z, k) := \mathbb{P}[\hat{Y} = k \mid Z = z] - \mathbb{P}[\hat{Y} = k \mid Z \neq z]$. In practice, we estimate multi-class demographic disparity for class k and subgroup z with: $\hat{\Gamma}(z, k) := \frac{|\{\hat{Y}=k, Z=z\}|}{|\{Z=z\}|} - \frac{|\{\hat{Y}=k, Z \neq z\}|}{|\{Z \neq z\}|}$, where $\hat{Y} = \theta(\mathbf{x}, z)$. We define *demographic parity* when the worst-case demographic disparity between members and non-members for any subgroup, and for any class is bounded:

DEFINITION 1 (γ -DEMPARITY). *For predictions Y with corresponding sensitive attributes Z to satisfy γ -bounded demographic parity (γ -DemParity), it must be that for all z in \mathcal{Z} and for all k in \mathcal{K} , the demographic disparity is at most γ : $\Gamma(z, k) \leq \gamma$.*

2.2 Privacy: Differential Privacy

Differential Privacy (DP) [11] formalizes the intuition that no individual data point should significantly impact the results of an analysis ran on a complete dataset. This allows it to learn properties of the dataset while ensuring individual data points' privacy. More formally, (ϵ, δ) -DP can be expressed as follows:

DEFINITION 2 ((ϵ, δ) -DIFFERENTIAL PRIVACY). *Let $M: \mathcal{D}^* \rightarrow \mathcal{R}$ be a randomized algorithm that satisfies (ϵ, δ) -DP with $\epsilon \in \mathbb{R}_+$ and $\delta \in [0, 1]$ if for all neighboring datasets $D \sim D'$, i.e., datasets that differ in only one data point, and for all possible subsets $R \subseteq \mathcal{R}$ of the result space it must hold that $\mathbb{P}[M(D) \in R] \leq e^\epsilon \cdot \mathbb{P}[M(D') \in R] + \delta$.*

The parameter ϵ bounds the maximal difference between the analysis results on the neighboring datasets while the second parameter δ represents a relaxation of the bound by allowing the results to vary more than the factor e^ϵ . Hence, the total privacy loss is bounded by ϵ with a probability of at least $1 - \delta$ [11]. Note that smaller ϵ correspond to better privacy guarantees for the data.

In ML, there exist two main canonical algorithms to implement DP, first the Private Aggregation of Teacher Ensemble (PATE) [21] an ensemble-based approach for private knowledge transfer, and the Differential Private Stochastic Gradient Descent (DP-SGD) [1]. We present background on both and extend upon them in Section 4.4 and Section 4.1.

PATE. PATE (Figure 3), takes advantage of an unlabeled public data set D_{public} to conserve the privacy of sensitive data D_{private} . Therefore, an ensemble of B teacher models $\{\theta_i\}_{i=1}^B$ is trained using disjoint subsets of D_{private} and their knowledge is transferred to a separate *student* model that can be publicly released. For the knowledge transfer, trained teachers label query data points from D_{public} . The final label of the query is the majority over the vote counts $N(\mathbf{x}) = [n_{i,j}]_{B \times K}$, where K is the number of classes.

PATE estimates the privacy cost of answering queries (i.e. labeling data) through *teachers consensus* with higher consensus revealing less information about individual teachers, and, thereby, consuming less privacy costs. To take advantage of the fact that estimating consensus is less privacy-costly than answering queries, PATE rejects high-cost queries to save on the privacy budget (see Algorithm 1). Both consensus estimation and vote aggregation (answering the query) are noised with $\mathcal{N}(0, \sigma_1^2)$ and $\mathcal{N}(0, \sigma_2^2)$, respectively; where σ_1, σ_2 are tuned for better student accuracy.

Algorithm 1 – Confident-GNMax Aggregator (from [22]) given a query, consensus among teachers is first estimated in a privacy-preserving way to then only reveal confident teacher predictions.

Require: input x , threshold T , noise parameters σ_1 and σ_2

- 1: **if** $\max_j \{ \sum_{i \in [B]} n_{i,j}(x) \} + \mathcal{N}(0, \sigma_1^2) \geq T$
- then**
- 2: **return** $\arg \max_j \{ \sum_{i \in [B]} n_{i,j}(\mathbf{x}) + \mathcal{N}(0, \sigma_2^2) \}$
- 3: **else**
- 4: **return** \perp

DP-SGD. The DP-SGD extends standard stochastic gradient descent (SGD) with two additional steps to implement privacy guarantees. First, the individual data points’ gradients are clipped to a maximum gradient norm bound C . This bounds the gradients’ sensitivity, which ensures that no data points can incur changes to the model above magnitude C . After clipping, Gaussian noise with scale $\mathcal{N}(0, \sigma^2 C^2)$ is added to mini-batches of clipped gradients. The noise distribution has zero mean and standard deviation proportional to a pre-defined noise multiplier σ and the clipping norm C . We detail the DP-SGD algorithm in Algorithm 5 in Appendix C.

To yield tighter privacy bounds, DP-SGD implements a privacy amplification through subsampling [6]: Training data points are sampled into mini-batches with a Poisson sampling per training iteration, in contrast to grouping the entire training data into mini-batches prior to every epoch as done in standard SGD. Hence, the traditional concept of an epoch (as a full training on the entire training data) does not exist in DP-SGD. Instead, each data point is sampled in every iteration according to a given sampling probability. Privacy amplification through subsampling allows to scale down the noise σ by the factor L/N (with L being the expected mini-batch size, N the total number of data points, and $L \ll N$) while still ensuring the same ϵ as with σ [15] which is crucial to the practical performance (privacy-utility trade-offs) of DP-SGD.

2.3 Pareto Efficiency

Let Θ be the set of all feasible baseline ML models with an element $\theta \in \Theta$. A feasible model is one that is achievable through learning (optimization) over a given dataset. Consider I to be the set of measurable trustworthiness objectives such that the *loss value* of objective $i \in I$ is described by $\ell_i(\theta)$. For instance, without loss of generality, $u_{\text{priv}} = \epsilon$ where ϵ is the achieved differential privacy budget from Definition 2, and $l_{\text{fair}} = \hat{\Gamma}(z, k)$ is the demographic parity loss in Section 2.1. We assume lower loss values are desirable.

DEFINITION 3 (PARETO EFFICIENCY). $\theta \in \Theta$ is *Pareto-efficient* if there exists no $\theta' \in \Theta$ such that (a) $\forall i \in I$ we have $\ell_i(\theta') \leq \ell_i(\theta)$, and that (b) for at least one objective $j \in I$ the inequality is strict $\ell_j(\theta') < \ell_j(\theta)$.

If an alternative model θ' exists that satisfies conditions (a) and (b), we say θ' *Pareto dominates* θ . In this case θ is *Pareto-inefficient*, and therefore, is not on *the Pareto frontier*. Intuitively, θ is on the Pareto frontier, if we cannot improve objective i without deteriorating another objective j .

3 RESOLVING TENSIONS BETWEEN FAIRNESS AND PRIVACY

For an impartially-specified model, we would like to consider fairness and privacy on equal footing. However, this is hard to achieve since the types of guarantees and their mitigations are inherently different. Differential privacy is an algorithmic and often data-independent guarantee. Once provided in a pipeline, thanks to its post-processing² property [11], it ensures privacy in the rest of the pipeline. The same cannot be said of fairness mitigations, as most are defined with respect to the model training process (see Figure 2) and are often data-dependent. In this section, we discuss the challenges of implementing fairness into private training through pre-processing. We present a post-processor as an ad-hoc solution concept which can replace pre-processing, or be used alongside it. Since studying in-processing fairness methods requires knowledge of the training procedure, we defer their study until Section 4.

²We note that “post-processing” is an overloaded term. In algorithmic fairness literature, it means that the fairness mitigation happens after a model is trained while in differential privacy, it is a property which states that processing the output of a differentially private mechanism does not incur additional privacy leakage. Unless stated explicitly, we use post-processing in its algorithmic fairness sense.

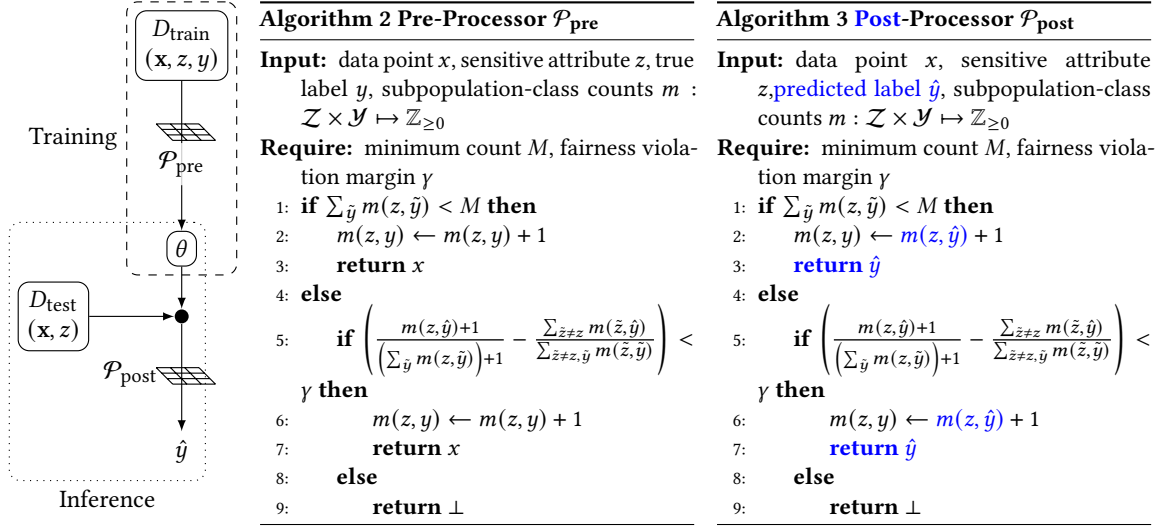


Fig. 2. **Demographic parity pre- and post-processor.** We depict the placement of the fairness mitigation (left). While the pre-processor \mathcal{P}_{pre} (Algorithm 2, middle) operates within training, the post-processor \mathcal{P}_{post} (Algorithm 3, right) is applied at inference time. The subpopulation-class count m refers to the number of data points per-class within each of the subpopulation groups. It is used to empirically estimate the demographic disparity $\hat{\Gamma}(z, k)$, $m : \mathcal{Z} \times \mathcal{K} \mapsto \mathbb{Z}_{\geq 0}$. After a *cold-start phase* (line 1-3) in both algorithms, we start rejecting queries for x if we have too few samples from a given class.

3.1 Shortcomings of Providing Fairness Through Pre-Processing in Private Learning

Fairness Degradation through Private Training. Consider a model θ with no fairness mitigations. Assume we use a demographic parity pre-processor (such as Algorithm 2) to achieve zero demographic parity violations ($\gamma = 0$). If we wish to make the model θ private, we need to introduce randomization through adding a calibrated amount of noise either to the model or its output (see Section 2.2). In either case, noising has the effect of flipping predicted labels, such that the demographic parity violations might increase. Put differently, differential privacy may cause a *label shift* in the predictions, which our earlier fairness mitigation cannot account for.

Privacy Cost of Fairness Pre-Processing. Fairness pre-processing can lead to increased privacy costs during private training. A consequence of differential privacy is the privacy consumption regime [19]: just by observing the data for the purposes of equalizing a fairness measure between subpopulations, we may consume from the privacy budget.³ This budget could otherwise be spent, for instance, on more training passes on data to yield higher accuracy. We formalize this observation in Theorem 3.1 for the case when a universal ordering exists. We defer the proof to Appendix B.

THEOREM 3.1. *Assume the training dataset $D = \{(x, z, y) \mid x \in \mathcal{X}, z \in \mathcal{Z}, y \in \mathcal{Y}\}$ is fed through the demographic parity pre-processor \mathcal{P}_{pre} (as in Algorithm 2 without the minimum subgroup size constraint and applied offline retroactively) following an ordering defined over the input space \mathcal{X} . Let \mathcal{P}_{pre} enforce a maximum violation γ , and $|\mathcal{Z}| = 2$. Suppose now \mathcal{M} is an (ϵ, δ) training mechanism, then $\mathcal{M} \circ \mathcal{P}_{pre}$ is $(K_Y \epsilon, K_Y e^{K_Y \epsilon} \delta)$ -DP where $K_Y = 2 + \left\lceil \frac{2\gamma}{1-\gamma} \right\rceil$.*

³Note that this disadvantage does not hold for fairness post-processing which does not incur additional privacy costs due to the differential privacy post-processing property.

3.2 Improving Fairness-Privacy Interplay Through Additional Pre- and Post-Processing

S1: Fairness Pre-Processor. On its own, Theorem 3.1 is a negative result: at least a subset of pre-processing mitigations will degrade the privacy guarantee. To avoid paying a factor of two (or more) on the privacy analysis, we propose extending the fairness pre-processor to perform an *un-supervised fairness calibration via a public dataset*. By leveraging the information from such a public dataset, one can ensure that demographic parity holds without paying the additional privacy costs. Importantly, since estimating demographic parity does not require a ground-truth label, we can do so in an un-supervised fashion. In this sense, the public dataset serves as an independent calibration set.

S2: Fairness Post-Processor. Our solution for the label shift is to rely on a post-processor. We know, through the differential privacy post-processing property, that this type of mitigation does not cost us any additional privacy budget. More importantly, since the mitigation occurs at the last stage of the ML pipeline (see Figure 2); the fairness guarantee is ensured at test (inference) time in spite of the aforementioned label shift induced by differential privacy.

Our post-processor in Algorithm 3 implements a filtering mechanism that determines which queries to answer given a current estimates of the demographic parity $\Gamma(z, k)$. A query x is rejected if answering it violates the maximum demographic parity violations γ . The algorithm has a cold-start phase where the model releases decisions on all queries until at least M queries for each subgroup z have been answered.

4 IN SEARCH OF IMPARTIAL ALGORITHMS

In designing a trustworthy system, ideally, we want to avoid inefficiencies in every objective. The prior section highlights that certain ways of composing ML pipelines are more efficient than others: while privacy degrades through fairness pre-processing it is not affected by fairness post-processing. This section focuses on in-processing mechanisms, which we left out in the prior section due to their dependence on the training framework. We start our study with DP-SGD—currently the predominant mechanism for privacy-aware machine learning.

4.1 Integrating Fairness into DP-SGD

Having presented the necessary background in DP-SGD in Section 2.2, here, we first propose how one can incorporate a demographic parity fairness in-processor within DP-SGD. Given that DP-SGD builds on optimizing a learning loss in an unconstrained optimization setup, a valid solution consists in incorporating a fairness regularizer [16] within the learning loss. Such a fairness regularizer always needs to estimate a (group) fairness statistic—in our case, the demographic disparity—on the training data. In DP-SGD, this is challenging for two main reasons: 1) Since, in SGD training is conducted in mini-batches, we need to ensure that the mini-batches are representative of the whole dataset and large enough to produce an accurate estimation of the fairness statistic. 2) Using the training data for additional fairness estimation during training consumes additional privacy budget which could otherwise be spent training for more iterations on the data, and thereby, improving model accuracy.

On a first glance, 1) might be easily resolved, for example, by stratified sampling [23] of mini-batches, and using larger mini-batch sizes (or switching to full-batch gradient descent); respectively. However, from a privacy perspective, these changes to the sampling procedure will severely degrade the improved privacy-utility trade-offs from the privacy amplification through subsampling in DP-SGD (see Section 2.2). Instead, we propose addressing 1) through regularization during training and 2) by assessing demographic parity on an unlabeled public dataset, as sketched in **S1** in the prior section. The resulting changes yield our novel FairDP-SGD algorithm.

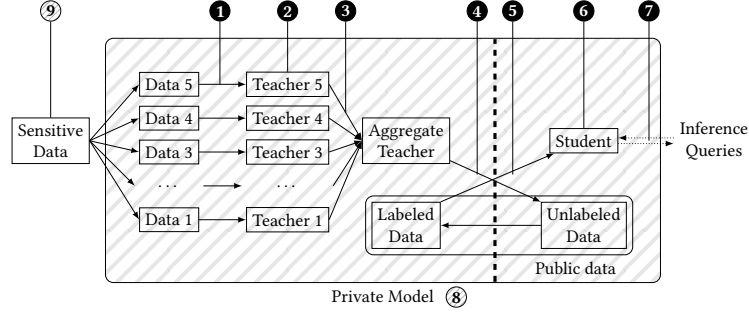


Fig. 3. **Various ways to integrate fairness in PATE:** For teachers: Pre-**1**/In-**2**/Post-Processing **3**. For the student: Pre-**5**/In-**6**/Post-Processing **7**. A fair supervised privacy-preserving algorithm (e.g., FairDP-SGD) replaces the private model (grey stripes) in-processing **8**, while a pre-processor applies to sensitive data directly **9**. Dashed line separates public and private data domains. Our FairPATE’s intervention occurs at **4**.

4.2 Our FairDP-SGD

We introduce FairDP-SGD (Algorithm 6), which extends DP-SGD with a demographic parity fairness regularizer (DPL). In order to avoid paying an extra privacy cost for inferring the fairness measure over the training set, we do so over a public dataset. The resulting demographic parity loss is:

$$\text{DPL}(\theta; X_{\text{public}}) = \max_k \max_z \widehat{\Gamma}(z, k) = \max_k \max_z \left\{ \frac{|\{\hat{Y} = k, Z = z\}|}{|\{Z = z\}|} - \frac{|\{\hat{Y} = k, Z \neq z\}|}{|\{Z \neq z\}|} \right\} \quad (1)$$

where $\hat{Y} = \theta(X_{\text{public}})$ is the prediction of the private model θ on the features X_{public} of the public dataset D_{pub} .

We need to preserve useful gradients when implementing the DPL for it to be effective. Calculating the DPL requires calculating the fairness violation, $\Gamma(z, k)$, which uses the predicted label k . An $\arg \max$ is typically applied to the output to obtain the prediction, but the operation is not differentiable, and hence unsuited to obtain useful gradients. To overcome this limitation, we use a tempered softmax: $\text{softmax}_T(x_i) = \frac{\exp x_i/T}{\sum_i \exp x_i/T}$, where T is the temperature. We set T to a very small value (e.g. 0.01) to make the tempered softmax close to the *argmax* while keeping the overall loss differentiable.

During training, FairDP-SGD adds the DPL to the original loss function and scales it by a regularization factor, λ which balances between both loss functions. Note that in contrast to the fairness violation γ , λ is a *pre-specified* model parameter while γ reports the fairness violation obtained by the final trained model. At inference time, FairDP-SGD also uses the post processor Algorithm 3 to ensure that the results satisfy the specified fairness constraint. The privacy analysis of FairDP-SGD follows that of DP-SGD with Poisson Sampling [1, 35] with no additional privacy costs resulting from fairness assessment on the public data (see **S1** in Section 3.1).

4.3 Integrating Fairness into PATE

Given that FairDP-SGD has only a few degrees of freedom where fairness measures can be implemented, in the next section, we turn our study to the more complex PATE framework. The modular design of PATE (see Figure 3) allows us multiple points of fairness integration which we detail these in the following. These can potentially lead us to achieve stronger improvements to fairness-privacy tradeoff.

Teacher-level ①, ②, and ③. All three designs are non-impartial as they place fairness before privacy mitigation. Since in PATE, privacy is ensured at the level of aggregated teachers and not individual teachers, all the three alternatives can be seen as instances of the fairness pre-processor \mathcal{P}_{pre} in Theorem 3.1 from the final student model’s perspective. As a result, they all suffer from additional privacy leakage; on the level of teacher data ①, model ②, or vote ③.

Student-level ⑤, ⑥, and ⑦. These designs place privacy before fairness, therefore, they are also non-impartial. Thanks to differential privacy post-processing, the privacy budget remains unchanged. However, the drawbacks are in terms of fairness and accuracy. We discussed the former in Section 3.1. Regarding the impact on accuracy, remember that in the query phase of PATE, we incur a much smaller privacy cost for rejecting a query than for answering it (see Section 2.2). Now consider the scenario in ⑤ where a query is labeled but is ultimately rejected due to a fairness violation. In this case, the extra budget incurred for answering the query is wasted. Saving this budget could have allowed us to answer more queries, thus enabling higher student accuracy. Therefore, ⑤ is not Pareto-efficient. We note that our demographic parity post-processor in Algorithm 3 is suitable for ⑦ but, on its own, still inefficient. We demonstrate the inefficiencies of ⑥ and ⑦ empirically in Section 5.1.2.

4.4 Our FairPATE

Having discussed all of the alternatives, it is clear that the only viable for the fairness mitigation option is at ④, namely, at the level of aggregate teacher. This choice reflects the impartiality principle: the fairness mitigation occurs exactly at the point where PATE privacy-preserving mechanisms are implemented—avoiding the aforementioned inefficiencies.

4.4.1 Confident&Fair-GNMax. Our new aggregation mechanism for FairPATE, which we call *Confident&Fair-GNMax* (CF-GNMax) extends PATE’s standard GNMax algorithm (Algorithm 1) with the idea of rejecting queries due to their privacy cost to also rejecting queries due to their disparate impact on fairness. Concretely, CF-GNMax, integrates an additional demographic parity constraint within the aggregator which allows rejecting queries on the basis of fairness, see Algorithm 4. The algorithm checks potential violations of demographic disparity violations and maintains an upper bound γ on them in the course of answering PATE queries (Line 7). The goal is to bound $\Gamma(z, k)$ but which, in practice, needs to be empirically estimated. Concretely, we measure demographic disparity $\widehat{\Gamma}(z, k)$ using the counter $m : \mathcal{Z} \times \mathcal{K} \mapsto \mathbb{Z}_{\geq 0}$ which tracks per-class, per-subgroup decisions.

Care must be taken to produce accurate $\Gamma(z, k)$ estimations: with few samples, $\widehat{\Gamma}(z, k)$ may be a poor estimator of $\Gamma(z, k)$. Therefore, we introduce a *cold-start* stage where there are not yet enough samples to estimate $\widehat{\Gamma}(z)$ accurately. We avoid rejecting queries due to the fairness constraint at this stage. Concretely, we require at least, on average, M samples from the query’s subgroup before we reject a query on the basis of fairness (line 3).

Algorithm 4 – Confident&Fair-GNMax Aggregator

Input: query data point x , sensitive attribute z , predicted class label k , subpopulation subclass counts $m : \mathcal{Z} \times \mathcal{K} \mapsto \mathbb{Z}_{\geq 0}$

Require: minimum count M , threshold T , noise parameters σ_1, σ_2 , fairness violation margin γ

```

1: if  $\max_j \{n_j(x)\} + \mathcal{N}(0, \sigma_1^2) \geq T$  then
2:    $k \leftarrow \arg \max_j \{n_j(x) + \mathcal{N}(0, \sigma_2^2)\}$ 
3:   if  $\sum_{\tilde{k}} m(z, \tilde{k}) < M$  then
4:      $m(z, k) \leftarrow m(z, k) + 1$ 
5:     return  $k$ 
6:   else
7:     if  $\left( \frac{m(z, k) + 1}{(\sum_{\tilde{k}} m(z, \tilde{k})) + 1} - \frac{\sum_{\tilde{z} \neq z} m(\tilde{z}, k)}{\sum_{\tilde{z} \neq z, \tilde{k}} m(\tilde{z}, \tilde{k})} \right) < \gamma$ 
8:        $m(z, k) \leftarrow m(z, k) + 1$ 
9:       return  $k$ 
10:    else
11:      return  $\perp$ 
12:  else
13:    return  $\perp$ 

```

4.4.2 Student-Preprocessor. FairPATE can introduce a label shift in the training data of the student model due to rejecting queries to stay within the fairness constraint. Thus, the student model needs to implement a similar procedure at inference time when answering queries. FairPATE uses a post-processor as the one introduced in Algorithm 3, which mirrors the fairness constraint used to query the teachers while performing inference on the student.

4.4.3 Privacy Analysis. FairPATE’s query phase (CF-GNMAX, Algorithm 4) has two main differences with PATE’s (C-GNMAX, Algorithm 1). First, it consists of two stages: a cold-start stage and a rejection/acceptance stage. Second, in FairPATE, a query may not be answered not only because it is too privacy-costly, but also because answering it would violate the fairness (γ -DemParity) constraint. During the cold-start stage, FairPATE follows the same analysis as PATE’s (Appendix A). When the cold start phase is finished (line 3), we adjust the privacy analysis to account for the fact that in FairPATE, a query q_i can be additionally rejected if answering q_i violates the maximum disparity gap. We can calculate FairPATE’s probability of answering query q_i as:

$$\mathbb{P}[\text{answering } q_i(z, k)] = \begin{cases} 0 & \frac{m(z, k)+1}{(\sum_{\tilde{k}} m(z, \tilde{k})+1)} - \frac{\sum_{\tilde{z} \neq z} m(\tilde{z}, k)}{\sum_{\tilde{z} \neq z, \tilde{k}} m(\tilde{z}, \tilde{k})} > \gamma \\ \tilde{q} & \text{otherwise} \end{cases} \quad (2)$$

where k is the noisy argmax (Line 2), \tilde{q} is calculated using Proposition 1 in Appendix A (as before), and the left side of the condition is simply calculating the new tentative demographic disparity violation $\Gamma(z, k)$ if the query is accepted.

Note that, since the value of the counter $m(z, k)$ is only conditioned on the value of the noisy argmax, by the post-processing property of DP [10], $m(z, k)$ and by extension, Line 7 do not add any additional privacy cost. In other words, rejecting queries on the basis of fairness, does not incur additional privacy cost.

5 EMPIRICAL EVALUATION

We evaluate FairPATE and FairDP-SGD on multiple datasets and derive the Pareto frontiers between fairness, privacy, and accuracy. Based on the Pareto frontiers, we answer the following research questions (**RQs**): **RQ1**: What is the impact of the post-processing strategy (see **S2** in Section 3.2) that filters queries on trade-offs obtained by FairPATE and FairDP-SGD? **RQ2**: Where are fairness mitigations best implemented within the private ML pipeline? **RQ3**: How do FairPATE and FairDP-SGD differ in performance? **RQ4**: Can baseline specification (see Section 2) by a public entity, e.g. a regulatory body, be done without having direct access to a the private data.

Experimental Setup. We evaluate five datasets, namely ColorMNIST [2], CelebA [18], FairFace [17], UTKFace[34], and CheXpert[13]. We refer to Table 2 in Appendix E for details on these datasets. For both FairPATE and FairDP-SGD, we train multiple models to study the trade-offs between model privacy, fairness, and accuracy. We train models with different privacy budget ϵ and fairness specifications. Reminder that in FairPATE, we apply a fairness constraint γ . In FairDP-SGD, we have an unconstrained optimization problem and control the regularization factor λ . In all cases, we make a distinction between specified and achieved privacy budgets and fairness gaps. We report exclusively achieved values of γ and ϵ , as well as model accuracy, and coverage on these models.

5.1 Findings

In this section, we derive Pareto frontiers for FairPATE and FairDP-SGD, as well as several relevant baselines for comparisons. For conciseness, we address our various methods using the visual language adopted in Figure 3. For clarity, we distinguish our main contributions FairPATE (④+⑦) and FairDP-SGD (⑧+⑦).

Setting	Method	ϵ -(\downarrow) Budget	Max (\downarrow) Disparity	Acc. (\uparrow)	Cov. (\uparrow)	ϵ -(\downarrow) Budget	Max (\downarrow) Disparity	Acc. (\uparrow)	Cov. (\uparrow)
		ColorMNIST				UTKFace			
Fair	FairPATE	2.88	0.01	85.6	0.62	8.65	0.01	83.8	0.78
	FairDP-SGD	1.0	0.01	85.4	0.64	8.0	0.01	81.2	0.72
	⑤ + ⑦	2.88	0.01	80.9	0.69	10.0	0.01	82.6	0.82
	⑥ + ⑦	2.88	0.01	84.6	0.63	10.0	0.01	81.4	0.74
Private	FairPATE	1.0	0.10	73.8	1.00	2.0	0.13	74.0	0.98
	FairDP-SGD	1.0	0.10	88.8	1.00	2.0	0.15	75.3	0.99
	⑤ + ⑦	1.0	0.10	73.1	1.00	2.0	0.14	72.3	0.99
	⑥ + ⑦	1.0	0.10	74.2	0.98	2.0	0.15	72.3	0.99
Accurate	FairPATE	2.87	0.10	88.5	0.99	10.0	0.2	82.9	0.97
	FairDP-SGD	2.0	0.10	88.6	0.99	10.0	0.1	81.3	0.96
	⑤ + ⑦	2.88	0.10	88.1	1.0	10.0	0.01	82.6	0.82
	⑥ + ⑦	3.0	0.10	88.5	0.99	10.0	0.15	81.6	0.92

Table 1. **Baseline Comparisons.** ⑤ + ⑦: PATE with pre-processing. ⑥ + ⑦: PATE with in-processing. \downarrow : the lower, the better; \uparrow : the higher, the better. The values are chosen by setting a hard constraint on the variable that we want to optimize for for each objective (**Setting**). Within that constraint, we always optimize for accuracy first and the other variables second. We observe that FairPATE and FairDP-SGD achieve the highest accuracy in most settings. ⑤ + ⑦ always achieves highest coverage when only optimizing for fairness.

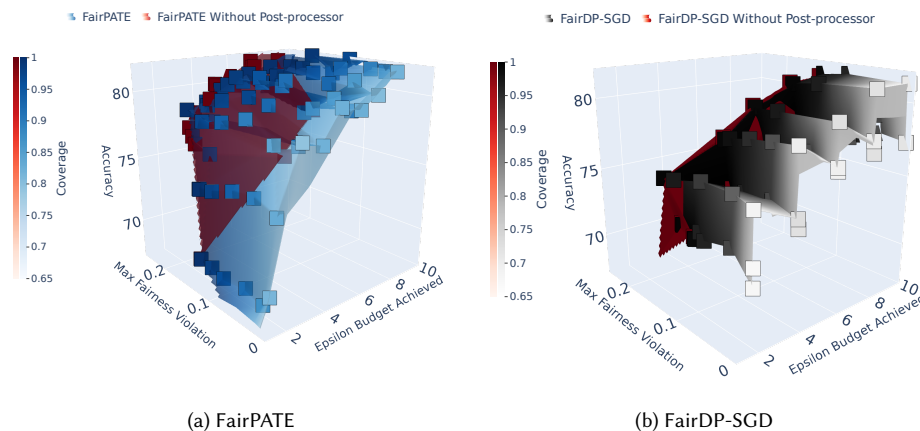


Fig. 4. **FairPATE and FairDP-SGD with vs. without post-processor on UTKFace.** Post-processor helps satisfy small fairness constraints while preserving model accuracy at the cost of answering fewer queries.

5.1.1 RQ1: Post-processing is necessary for ensuring tight fairness gaps. To study the necessity of it, we train FairPATE and FairDP-SGD models without the post-processor ⑦. FairPATE results are shown in Figure 4a. Without the post-processor, results span a smaller range of fairness violations. This is expected as FairPATE introduces a label shift in its training data that should be mirrored in the test data Section 4.4.2, as done by the post-processor. The post-processor, thus, ensures that tighter fairness gaps are feasible. We plot FairDP-SGD results with and without post-processor in Figure 4b. We observe that using demographic parity loss (DPL; see Equation (1)), we can achieve smaller fairness violations but the model utility decreases accordingly. In order to reach very small fairness gaps, we

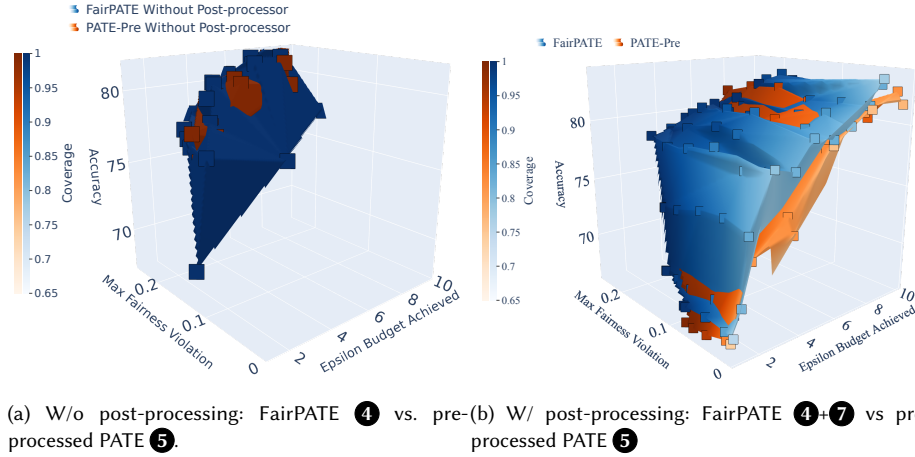


Fig. 5. **FairPATE vs. pre-processed PATE on UTKFace.** FairPATE has higher accuracy than PATE pre-processing in high privacy budget low fairness violation regions, and they have similar accuracy in other regions.

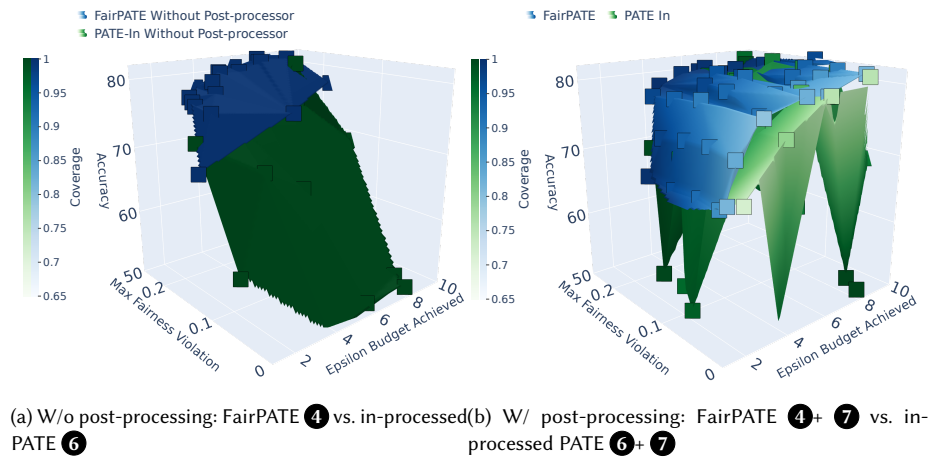


Fig. 6. **FairPATE vs. in-processed PATE on UTKFace.** FairPATE has higher accuracy and coverage than PATE in-processing in low fairness violation regions, and they have similar accuracy in other regions.

lose all utility as the model becomes increasingly inaccurate. The post-processor can preserve utility while satisfying tight fairness constraints at the cost of answering slightly fewer queries.

5.1.2 RQ2: FairPATE Pareto dominates similar designs in most contexts. The closest baseline to FairPATE (④+⑦) is querying with PATE and then applying a pre-processor to the student model (⑤+⑦). The important distinction between the two is that while selecting which queries to answer and which to reject, ④ considers both fairness and privacy at the same time. In contrast, in ⑤, PATE only enforces privacy while the fairness constraint is applied post-hoc. We first compare ④ and ⑤ without the post-processor. The respective student model results are shown in Figure 5a. We observe that the two methods achieve similar performances with higher privacy budget and

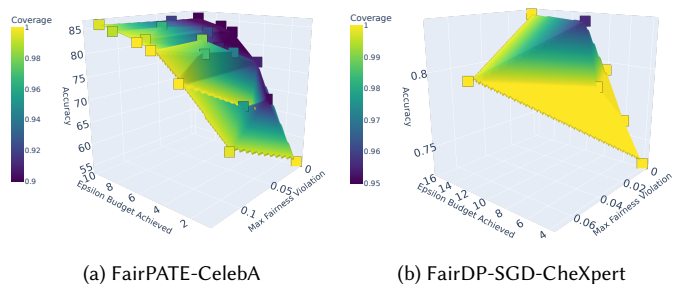


Fig. 8. **FairPATE on CelebA and CheXpert.** The figure plots the model results that are on the Pareto frontier. We observe that in both figures, model accuracy increases with higher privacy budget ϵ , and looser fairness constraints yield higher coverage.

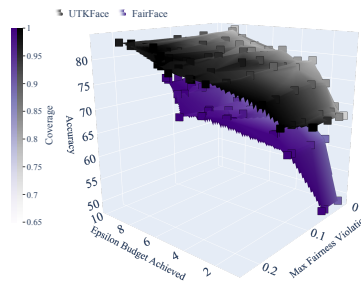


Fig. 9. **Pareto Frontier of FairPATE results on UTKFace and FairFace.** The two surface have very similar shapes despite the differences in accuracy.

higher fairness violation regions. This is because the fairness constraint is too loose to activate FairPATE’s fairness mechanism. In higher privacy budget and low fairness violation regions, FairPATE yields student models of higher accuracy. This is because it is able to answer more queries since it saves privacy costs by rejecting queries due to fairness constraints. Then, we compare the two methods (FairPATE (4) + (7) and (5) + (7)) with the post-processor. Results are shown in Figure 5b. In addition to the previous observations, we also note that the two methods perform similarly in very low privacy budget regions. In these regions, too few queries are answered to leave the cold-start phase and activate the actual fairness-based rejection mechanism.

We also compare FairPATE to (6), where we query with PATE and employ an in-processing method in training the student model. For a controlled experiment, we use the same demographic parity loss as FairDP-SGD for the fairness regularizer, namely DPL Equation (1). The student model results in Figure 6a highlight that models that use in-processing are able to achieve much smaller fairness violations without the post-processor. Yet, utility decreases as fairness violation decreases. With the post-processor (Figure 6b), the two methods perform similarly in larger fairness violation regions, but FairPATE perform better in smaller fairness violation regions with better accuracy and coverage.

5.1.3 RQ3: FairPATE performs better than FairDP-SGD, especially with higher privacy budgets. We compare our two methods, FairPATE and FairDP-SGD in Figure 7. We observe that while they yield similar accuracy in low privacy budget regions, FairPATE provides better accuracy in higher privacy budget regions. Additionally, in low fairness violation regions, FairPATE obtains higher accuracy and higher coverage. This may be an artefact of regularization in FairDP-SGD and could potentially be improved with hyper-tuning. However, this shows that it is easier to derive a Pareto frontier for FairPATE. We attribute this to the fact that performance (accuracy) of PATE (FairPATE) is meaningfully correlated with the number of answered queries (See Appendix F). This allows for more fine-grained control on privacy costs, and thus a smoother Pareto frontier.

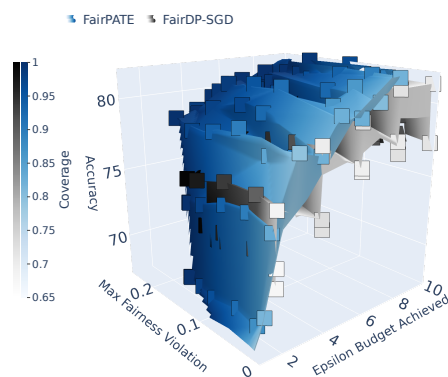


Fig. 7. **FairPATE vs FairDP-SGD on UTKFace.** Our methods yield similar accuracy in low privacy budget or high privacy budget-high fairness violation regions. FairPATE has higher accuracy and higher coverage in high privacy budget-low fairness violation regions.

5.1.4 RQ4: Baseline specification without direct data access

is possible. We set out to explore whether the regulators can still

make good decisions when specifying the baseline model parameters without access to the actual private dataset. Therefore, we compare the Pareto frontier surfaces obtained on different datasets. Figure 8a and Figure 8b plots the Pareto frontier surface from FairPATE on CelebA and CheXpert respectively. Figure 9 plots the Pareto frontier from FairPATE on UTKFace and FairFace. Although the Pareto frontier surfaces show similar trends, the shapes are dataset dependent: different datasets show different trade-offs between the objectives. However, we notice that the Pareto frontier surface shapes on UTKFace and FairFace are very similar. The classification task on both datasets is gender, with race as the sensitive attribute. This shows that a regulator could use the Pareto frontier from a different dataset (which they have access to) to design baseline specifications—as long as the dataset is sufficiently close in terms of classification task.

6 RELATED WORK

Due to the multiplicity of algorithmic fairness notions, as well as privacy; defining a benchmark to study fairness-privacy trade-offs is difficult. In this paper, we focus on discovering the Pareto frontier between demographic parity fairness (a *group fairness* notion [4]) and (central) differential privacy [11]. While these objectives have a significant impact on each other (as we established in Section 3); each has been defined and developed independently of one another.

In contrast, there is a lineage of work that provides new definitions of fairness by characterizing the disparate impact of employing a privacy-aware mechanism [26, 27]. While useful in their own regard, these new definitions do not alleviate the burden of satisfying established notions of fairness, such as demographic parity.

Conceptually, the closest works to our setup are [14, 20] which assume different privacy notions. Both works strive to provide differential privacy (DP) with respect to the sensitive attribute. Jagielski et al. [14] assumes a central notion of DP, while Mozannar et al. [20] assume a *local* DP notion [5]. However, neither of the definitions used provide classical (approximate) differential privacy [11] guarantees with respect to *all features*. Furthermore, algorithms provided in these works, consider linear models and are optimized over tabular data. FairPATE and FairDP-SGD, on the other hand, are scalable deep-learning algorithms. Finally, neither prior work uses their formulations to derive a Pareto frontier.

7 LIMITATIONS & CONCLUSIONS

Trustworthiness in machine learning is inherently a multi-objective endeavour. We acknowledge that as algorithm designers, we are only a part of the decision making process which likely occurs before any human judgement is passed. As such, it is imperative that (i) our design decisions should not limit (human) decision maker choices; and (ii) favour one objective over another. In this paper, we addressed the first challenge by providing a rich trade-off representation between the different objectives (fairness, privacy, and accuracy) in the form of a Pareto frontier. Our answer to the second challenge emerged as a design principle, which we called the impartiality principle. We showed that models that break the impartiality principle are likely not on the Pareto frontier.

Moving forward, the intuition behind our framework is pervasive to different formulations of what it means to be trustworthy. However, our current work assumes demographic parity as the fairness notion. We acknowledge that other group fairness and individual fairness notions are prevalent in the literature. Adopting them would require adaptations of our algorithms. Similarly, we acknowledge the plethora of other privacy notions, including but not limited to, extensions or re-consideration of central (approximate) differential privacy. We leave their study to future work—having shown the value of ML mechanism design that supports obtaining Pareto frontiers.

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] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. 2019. Invariant risk minimization. *arXiv preprint arXiv:1907.02893* (2019).
- [3] Imon Banerjee, Ananth Reddy Bhimireddy, John L. Burns, Leo Anthony Celi, Li-Ching Chen, Ramon Correa, Natalie Dullerud, Marzyeh Ghassemi, Shih-Cheng Huang, Po-Chih Kuo, Matthew P. Lungren, Lyle Palmer, Brandon J. Price, Saptarshi Purkayastha, Ayis Pyrros, Luke Oakden-Rayner, Chima Okechukwu, Laleh Seyyed-Kalantari, Hari Trivedi, Ryan Wang, Zachary Zaiman, Haoran Zhang, and Judy W. Gichoya. [n. d.]. Reading Race: AI Recognises Patient’s Racial Identity In Medical Images. 4, 6 ([n. d.]), e406–e414. [https://doi.org/10.1016/S2589-7500\(22\)00063-2](https://doi.org/10.1016/S2589-7500(22)00063-2) arXiv:2107.10356 [cs, eess]
- [4] Solon Barocas, Moritz Hardt, and Arvind Narayanan. 2019. *Fairness and Machine Learning*. fairmlbook.org. <http://www.fairmlbook.org>.
- [5] Björn Bebensee. [n. d.]. Local Differential Privacy: a tutorial. ([n. d.]). <https://doi.org/10.48550/arXiv.1907.11908>
- [6] Amos Beimel, Shiva Prasad Kasiviswanathan, and Kobbi Nissim. 2010. Bounds on the sample complexity for private learning and private data release. In *Theory of Cryptography Conference*. Springer, 437–454.
- [7] US Census Bureau. 2020. *Formal Privacy Methods for the 2020 Census*. <https://www.census.gov/programs-surveys/decennial-census/decade/2020/planning-management/plan/planning-docs/privacy-methods-2020-census.html> Section: Government.
- [8] Toon Calders and Sicco Verwer. 2010. Three naive Bayes approaches for discrimination-free classification. *Data mining and knowledge discovery* 21, 2 (2010), 277–292.
- [9] Hongyan Chang and Reza Shokri. [n. d.]. On the Privacy Risks of Algorithmic Fairness. *IEEE Computer Society*, 292–303. <https://doi.org/10.1109/EuroSP51992.2021.00028>
- [10] Cynthia Dwork. [n. d.]. Differential privacy. In *Automata, languages and programming* (Berlin, Heidelberg, 2006), Michele Bugliesi, Bart Preneel, Vladimiro Sassone, and Ingo Wegener (Eds.). Springer Berlin Heidelberg, 1–12.
- [11] Cynthia Dwork and Aaron Roth. 2013. The Algorithmic Foundations of Differential Privacy. 9, 3 (2013), 211–407. <https://doi.org/10.1561/04000000042>
- [12] Tom Farrand, Fatemehsadat Mireshghallah, Sahib Singh, and Andrew Trask. 2020. Neither private nor fair: Impact of data imbalance on utility and fairness in differential privacy. In *Proceedings of the 2020 Workshop on Privacy-Preserving Machine Learning in Practice*. 15–19.
- [13] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. 2019. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 590–597.
- [14] Matthew Jagielski, Michael Kearns, Jieming Mao, Alina Oprea, Aaron Roth, Saeed Sharifi-Malvajardi, and Jonathan Ullman. 2019. Differentially private fair learning. In *International Conference on Machine Learning*. PMLR, 3000–3008.
- [15] Peter Kairouz, Brendan McMahan, Shuang Song, Om Thakkar, Abhradeep Thakurta, and Zheng Xu. 2021. Practical and private (deep) learning without sampling or shuffling. In *International Conference on Machine Learning*. PMLR, 5213–5225.
- [16] Toshihiro Kamishima, Shotaro Akaho, and Jun Sakuma. [n. d.]. Fairness-aware Learning through Regularization Approach. In *2011 IEEE 11th International Conference on Data Mining Workshops (2011-12)*. 643–650. <https://doi.org/10.1109/ICDMW.2011.83> ISSN: 2375-9259.
- [17] Kimmo Karkkainen and Jungseock Joo. 2021. FairFace: Face Attribute Dataset for Balanced Race, Gender, and Age for Bias Measurement and Mitigation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 1548–1558.
- [18] Ziwei Liu, Ping Luo, Xiaoang Wang, and Xiaoou Tang. 2015. Deep Learning Face Attributes in the Wild. In *Proceedings of International Conference on Computer Vision (ICCV)*.
- [19] Ilya Mironov. 2017. Renyi Differential Privacy. (2017), 263–275. <https://doi.org/10.1109/CSF.2017.11> arXiv:1702.07476
- [20] Hussein Mozannar, Mesrob Ohannessian, and Nathan Srebro. 2020. Fair learning with private demographic data. In *International Conference on Machine Learning*. PMLR, 7066–7075.
- [21] Nicolas Papernot, Martin Abadi, Ulfar Erlingsson, Ian Goodfellow, and Kunal Talwar. 2016. Semi-supervised knowledge transfer for deep learning from private training data. *arXiv preprint arXiv:1610.05755* (2016).
- [22] Nicolas Papernot, Shuang Song, Ilya Mironov, Ananth Raghunathan, Kunal Talwar, and Úlfar Erlingsson. 2018. Scalable private learning with pate. *arXiv preprint arXiv:1802.08908* (2018).
- [23] Van L. Parsons. 2017. *Stratified Sampling*. John Wiley & Sons, Ltd, 1–11. <https://doi.org/10.1002/9781118445112.stat05999.pub2> arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1002/9781118445112.stat05999.pub2>
- [24] Laleh Seyyed-Kalantari, Guanxiong Liu, Matthew McDermott, Irene Y. Chen, and Marzyeh Ghassemi. [n. d.]. CheXclusion: Fairness gaps in deep chest X-ray classifiers. In *Biocomputing 2021*. WORLD SCIENTIFIC, 232–243. https://doi.org/10.1142/9789811232701_0022
- [25] Vinith M. Suriyakumar, Nicolas Papernot, Anna Goldenberg, and Marzyeh Ghassemi. 2021. Chasing Your Long Tails: Differentially Private Prediction in Health Care Settings. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual Event, Canada) (FAccT ’21)*. Association for Computing Machinery, New York, NY, USA, 723–734. <https://doi.org/10.1145/3442188.3445934>
- [26] Cuong Tran, My H Dinh, Kyle Beiter, and Ferdinando Fioretto. 2021. A Fairness Analysis on Private Aggregation of Teacher Ensembles. *arXiv preprint arXiv:2109.08630* (2021).
- [27] Cuong Tran, My H Dinh, and Ferdinando Fioretto. 2021. Differentially Private Deep Learning under the Fairness Lens. *arXiv preprint arXiv:2106.02674* (2021).

- [28] Cuong Tran, Ferdinando Fioretto, and Pascal Van Hentenryck. 2021. Differentially Private and Fair Deep Learning: A Lagrangian Dual Approach. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- [29] Archit Uniyal, Rakshit Naidu, Sasikanth Kotti, Sahib Singh, Patrik Joslin Kenfack, Fatemehsadat Mireshghallah, and Andrew Trask. 2021. DP-SGD vs PATE: Which Has Less Disparate Impact on Model Accuracy? *arXiv preprint arXiv:2106.12576* (2021).
- [30] Salil Vadhan. 2017. The complexity of differential privacy. *Tutorials on the Foundations of Cryptography: Dedicated to Oded Goldreich* (2017), 347–450.
- [31] Jiaqi Wang, Roei Schuster, Ilya Shumailov, David Lie, and Nicolas Papernot. 2022. In Differential Privacy, There is Truth: On Vote Leakage in Ensemble Private Learning. *arXiv preprint arXiv:2209.10732* (2022).
- [32] Depeng Xu, Wei Du, and Xintao Wu. 2021. Removing disparate impact on model accuracy in differentially private stochastic gradient descent. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 1924–1932.
- [33] Tao Zhang, Tianqing Zhu, Kun Gao, Wanlei Zhou, and S Yu Philip. 2021. Balancing Learning Model Privacy, Fairness, and Accuracy With Early Stopping Criteria. *IEEE Transactions on Neural Networks and Learning Systems* (2021).
- [34] Zhifei Zhang, Yang Song, and Hairong Qi. 2017. Age Progression/Regression by Conditional Adversarial Autoencoder. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.
- [35] Yuqing Zhu and Yu-Xiang Wang. [n. d.]. Poission Subsampled Rényi Differential Privacy. In *Proceedings of the 36th International Conference on Machine Learning* (2019-05-24). PMLR, 7634–7642. <https://proceedings.mlr.press/v97/zhu19c.html> ISSN: 2640-3498.

A STANDARD PATE PRIVACY ANALYSIS

Papernot et al. [22] use Rényi differential privacy (RDP) [19] for accounting of the privacy budget expanded in answering each query. While the true privacy cost for each query is not known, an upperbound is estimated and summed over the course of the query phase. Answering queries stop when a pre-defined budget is exhausted. A student model is then trained on the answered queries.

Theorem A.1 establishes that the upperbound is a function of the probability of *not* answering a query i with the plurality vote i^* . Unsurprisingly, this privacy cost function must tends to zero when the said event is very unlikely (*i.e.*, strong consensus):

THEOREM A.1. *[From [22]] Let \mathcal{M} be a randomized algorithm with (μ_1, ε_1) – RDP and (μ_2, ε_2) – RDP guarantees and suppose that given a dataset D , there exists a likely outcome i^* such that $\Pr[\mathcal{M}(D) \neq i^*] \leq \tilde{q}$. Then the data-dependent Rényi differential privacy for \mathcal{M} of order $\lambda \leq \mu_1, \mu_2$ at D is bounded by a function of $\tilde{q}, \mu_1, \varepsilon_1, \mu_2, \varepsilon_2$, which approaches 0 as $\tilde{q} \rightarrow 0$.*

In practice, Proposition 1 is used to find \tilde{q}_i in Theorem A.1, and μ_1, μ_2 are optimized to achieve the lowest upperbound on the privacy cost of each query for every order λ of RDP.

PROPOSITION 1 (FROM [22]). *For any $i^* \in [m]$, we have $\Pr[\mathcal{M}_\sigma(D) \neq i^*] \leq \frac{1}{2} \sum_{i \neq i^*} \operatorname{erfc}\left(\frac{n_i - n_{i^*}}{2\sigma}\right)$, where erfc is the complementary error function.*

B PRIVACY COST OF PRE-PROCESSING

We provide the proof for Theorem 3.1.

PROOF. We will proceed to show that using a pre-processing that sorts through data following some ordering defined over the whole input space⁴ and, for any given label y , removes the last datapoints (following the ordering) in the majority subclass until it satisfies the γ -constraint will produce datasets at most $2 + K_\gamma = 2 + \left\lceil \frac{2\gamma}{1-\gamma} \right\rceil$ apart. One then applies group privacy to obtain the final claim of the theorem.

Let $D' = D \cup x^*$, and the label of x^* is y^* . We now proceed to analyze how far apart $\mathcal{P}_{\text{pre}}(D)$ and $\mathcal{P}_{\text{pre}}(D')$ can be. First note, they are the same on all labels not y^* , so we need only consider the difference on this label. First, let m be the

⁴An example of such ordering would be to order images based on their pixel values in some specified order of height, width and channel starting by checking the first pixel, then the second pixel, and so on.

size of the minority subclass for label y^* and let $m + c$ be the admissible size of the majority class. That is, we have $\frac{m}{2m+c} - \frac{m+c}{2m+c} < \gamma$. From this we can conclude $c = \lfloor \frac{\gamma}{1-\gamma} 2m \rfloor$. Given this relation between the size of majority class a function of the minority class, we proceed to go through all logical cases to show the maximum difference is as claimed above.

Suppose x^* belongs to the minority subclass for y^* in D . Then we have $m \rightarrow m + 1$ and hence $c \rightarrow \lfloor \frac{\gamma}{1-\gamma} 2(m + 1) \rfloor$. Thus we see $\mathcal{P}_{\text{pre}}(D')$ now admits one more point in the minority class of y^* and at most $1 + \lceil \frac{2\gamma}{1-\gamma} \rceil$ more points to the majority subclass (note we do not replace existing points as we follow the ordering on the input space). Thus the max change between $\mathcal{P}_{\text{pre}}(D)$ and $\mathcal{P}_{\text{pre}}(D')$ is $2 + \lceil \frac{2\gamma}{1-\gamma} \rceil$

Now suppose x^* belongs to majority subclass for y^* in D . In this case we have either x^* appears early enough in the ordering that it now replaces another point in the majority class when applying P , or it is not added. In the former case, this mean we have changed $\mathcal{P}_{\text{pre}}(D)$ by 2: we first removed a point and then added x^* . In the latter case, x^* did not get added into the dataset, more so because of the ordering, $\mathcal{P}_{\text{pre}}(D') = \mathcal{P}_{\text{pre}}(D)$ as the order of points before x^* is still the same. So in this case, once again, the change between $\mathcal{P}_{\text{pre}}(D)$ and $\mathcal{P}_{\text{pre}}(D')$ is less than $2 + \lceil \frac{2\gamma}{1-\gamma} \rceil$.

Thus we have by group privacy (see lemma 2.2 in [30]) that $\mathcal{M} \circ \mathcal{P}_{\text{pre}}$ gives the claimed DP-guarantee, as we set $K_\gamma = 2 + \lceil \frac{2\gamma}{1-\gamma} \rceil$

□

C ADDITIONAL BACKGROUND ON PRIVACY-PRESERVING ML

We include the standard DP-SGD algorithm (Algorithm 5) and FairDP-SGD (Algorithm 6) here for comparison.

Details of the FairDP-SGD algorithm is discussed in Section 4.1.

Algorithm 5 Standard DP-SGD, adapted from [1].

Require: Private training set $D_{\text{prv}} = \{(x_i, y_i) \mid i \in [N_{\text{prv}}]\}$, loss function $\mathcal{L}(\theta, x_i)$, Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

- 1: **Initialize** θ_0 randomly
 - 2: **for** $t \in [T]$ **do**
 - 3: Sample mini-batch L_t with sampling probability L/N ▷ Poisson sampling
 - 4: For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$ ▷ Compute gradient
 - 5: $\tilde{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max\left(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C}\right)$ ▷ Clip gradient
 - 6: $\tilde{\mathbf{g}}_t \leftarrow \frac{1}{|L_t|} (\sum_i \tilde{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$ ▷ Add noise
 - 7: $\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$ ▷ Descent
 - 8: **Output** θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.
-

D FAIRNESS METRICS AND EVALUATIONS

We evaluate and compare different ways to measure the demographic parity gap, $\Gamma(z, k)$. We then select one method to use in our implementations. We explore three different methods that compare the ratio between different sensitive groups to evaluate the chosen fairness metric.

Algorithm 6 FairDP-SGD

Require: Private training set $D_{\text{priv}} = \{(x_i, y_i) \mid i \in [N_{\text{priv}}]\}$, Public calibration set $D_{\text{pub}} = \{(\tilde{x}_i, z_i) \mid i \in [N_{\text{pub}}]\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$, Demographic Parity loss $\text{DPL}(\theta; D_{\text{pub}})$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

- 1: **Initialize** θ_0 randomly
- 2: **for** $t \in [T]$ **do**
- 3: Sample mini-batch L_t with sampling probability L/N ▷ Poisson sampling
- 4: For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \left(\mathcal{L}(\theta_t, x_i) + \lambda \text{DPL}(\theta_t; D_{\text{pub}}) \right)$ ▷ Compute gradient
- 5: $\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max \left(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C} \right)$ ▷ Clip gradient
- 6: $\tilde{\mathbf{g}}_t \leftarrow \frac{1}{|L_t|} \left(\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}) \right)$ ▷ Add noise
- 7: $\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$ ▷ Descent
- 8: **Output** θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

D.1 Demographic Parity Gap Measurements

- (1) Between Groups: This method computes and bounds the maximum difference between two pairs of sensitive groups.

$$\Gamma(z, k) := \max_{\tilde{z}} |\mathbb{P}[\hat{Y} = k | Z = z] - \mathbb{P}[\hat{Y} = k | Z = \tilde{z}]|. \quad (3)$$

- (2) To Overall: This method computes and bounds the difference between each sensitive group and the total of all groups.

$$\Gamma(z, k) := \mathbb{P}[\hat{Y} = k | Z = z] - \mathbb{P}[\hat{Y} = k]. \quad (4)$$

- (3) To Overall Without Double Counting: This method computes and bounds the difference between each sensitive group and the total of all other groups.

$$\Gamma(z, k) := \mathbb{P}[\hat{Y} = k | Z = z] - \mathbb{P}[\hat{Y} = k | Z \neq z]. \quad (5)$$

To compare the three methods, we generate some synthetic data and run queries on them using each method to compare the results.

D.2 Evaluation Results

We first generate synthetic data with two classes and three sensitive groups. The distribution of the generated data is shown below.

Class/Sensitive Group	0	1	2
0	324	420	445
1	287	274	250

D.2.1 By Group. Total number of queries answered = 1661

Class/Sensitive Group	0	1	2
0	315	364	361
1	191	213	217

D.2.2 To Overall. Total number of queries answered = 1832

Class/Sensitive Group	0	1	2
0	324	395	361
1	234	271	247

D.2.3 To Overall Without Double Counting. Total number of queries answered = 1772

Class/Sensitive Group	0	1	2
0	318	371	342
1	218	244	229

D.3 Conclusion

We decide to use the third method, to overall without double counting, as the comparison method. It is a balance between the by group method and the to overall method. We do not want the comparison method to be too strict, because then our algorithm would reject most queries due to fairness. On the other hand, we also do not want it to be too lenient that the fairness constraint is not enforced. One major drawback of the to overall method is that if most of the data is from one sensitive group, then that sensitive group would have too much influence over the overall class label distribution.

E EXPERIMENTAL SETUP

We split each dataset into a training set, an unlabeled set, and a test set. The sizes of these three datasets are determined based on the dataset sizes specified in original PATE [21, 22], and adapted to the difficulties of the prediction tasks. In FairPATE, the training set is further split into equal partitions to train the teacher models. We train as many teachers as possible while still achieving good ensemble accuracy overall. In FairDP-SGD, the whole training set is used to train the private model. The test set is used to evaluate the performance of the final model.

For FairPATE, we first train the teacher ensemble models, then query them with the public dataset, and aggregate their predictions using the FairPATE algorithm. The student model is trained on the public dataset with obtained labels. The model architectures, as well as the parameters used in querying the teacher models are detailed in Table 2 for each dataset, respectively. The model architectures are chosen by referencing what is used in related works for each dataset. FairDP-SGD models are trained with the same model architecture as indicated in the table.

Dataset	Prediction Task	C	Sens. Attr.	SG	Total	U	Model	Number of Teachers	T	σ_1	σ_2
ColorMNIST [2]	Digit	10	Color	2	60 000	1 000	Convolutional Network (Table 3)	200	120	110	20
CelebA [18]	Smiling	2	Gender	2	202 599	9 000	Convolutional Network (Table 4)	150	130	110	10
FairFace [17]	Gender	2	Race	7	97 698	5 000	Pretrained ResNet50	50	30	30	10
UTKFace [34]	Gender	2	Race	5	23 705	1 500	Pretrained ResNet50	100	50	40	15
CheXpert[13]	Disease	2	Race	3	152 847	4 000	Pretrained DenseNet121	50	30	20	10

Table 2. Datasets used for evaluation. Abbreviations: **C**: number of classes in the main task; **SG**: number of sensitive groups; **U**: number of unlabeled samples for the student training. Summary of parameters used in training and querying the teacher models for each dataset. The selection of σ_1 is in accordance with the threshold T . The selection process of σ_2 , is shown in the Appendix E.

We tune the amount of noise injected into the aggregation mechanism of FairPATE by varying the standard deviation of the Gaussian distribution while ensuring the accuracy of the labels produced by the teacher ensemble models to maximize the accuracy of student models.

Layer	Description
Conv2D with ReLU	(3, 20, 5, 1)
Max Pooling	(2, 2)
Conv2D with ReLU	(20, 50, 5, 1)
MaxPool	(2, 2)
Fully Connected 1	(4*4*50, 500)
Fully Connected 2	(500, 10)

Table 3. Convolutional network architecture used in ColorMNIST experiments.

Layer	Description
Conv2D	(3, 64, 3, 1)
Max Pooling	(2, 2)
ReLU	
Conv2D	(64, 128, 3, 1)
Max Pooling	(2, 2)
ReLU	
Conv2D	(128, 256, 3, 1)
Max Pooling	(2, 2)
ReLU	
Conv2D	(256, 512, 3, 1)
Max Pooling	(2, 2)
ReLU	
Fully Connected 1	(14 * 14 * 512, 1024)
Fully Connected 2	(1024, 256)
Fully Connected 2	(256, 2)

Table 4. Convolutional network architecture used in CelebA experiments.

F THE RELATIONSHIP BETWEEN NUMBER OF QUERIES ANSWERED AND STUDENT ACCURACY

We run a set of query experiments to investigate the trade-offs between privacy, fairness, and model utility. We do not train the student model for these querying experiments, but they will be trained later on. Instead, we use the number of queries answered as an estimate of the student model utility since an adequate number of queries needs to be answered to train a student model with good accuracy. In the first set of experiments, we run queries with varying consensus threshold T and fairness violation threshold ρ_{fair} at fixed privacy budget ϵ , and record the number of queries answered. We query the teacher ensemble models with varying privacy budget ϵ and fairness violation threshold ρ_{fair} . For these queries, we measure the maximum fairness violation γ , the achieved ϵ , and the number of queries answered. Using these query results, we also select and plot the points on the Pareto frontier. The results for the UTKFace dataset are shown in Figure 10. The results on the other datasets are found in Figure 12.

Figure 10 (left) plots the trade-offs between the maximum fairness violation γ , the achieved ϵ , and the number of queries answered. As expected, we observe that increasing ϵ allows more queries to be answered. Relaxing ρ_{fair} at fixed ϵ also leads to more queries being answered, although the effect is not as apparent. Additionally, when ϵ is very low, smaller γ is not achievable due to having too few queries answered and the fairness regulation mechanism not being activated as a result.

Figure 10 (right) plots the Pareto frontier of the query results. We plot the privacy constraint, fairness constraint, and the number of queries answered as a 3D plot to better visualize the tension between these different objectives. The figure gives similar insights as the other figure. Another observation is that although smaller γ is achievable when a higher number of queries are answered, at some point the fairness constraint needs to be relaxed in order to answer more queries.

We run an additional set of experiments of querying the teacher ensemble models to investigate the effect of different parameters on the number of queries answered. For these experiments, we run queries with varying consensus threshold T and fairness violation threshold ρ_{fair} at fixed privacy budget ϵ , and record the number of queries answered.

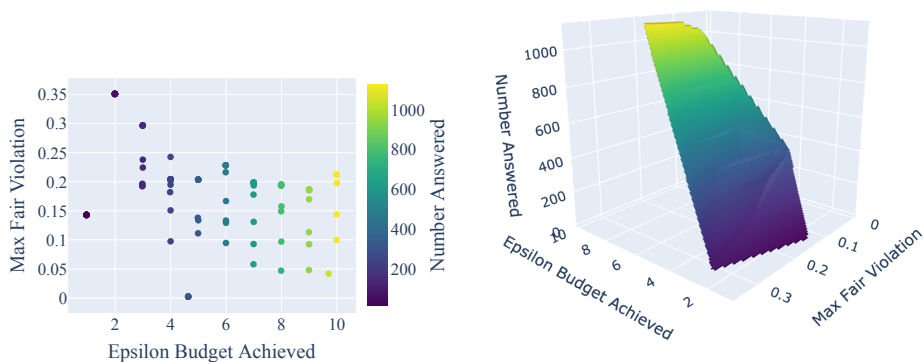


Fig. 10. **Query Experiment Results on UTKFace.** Experimental setup from Table 2. The left figure shows the trade-offs between the maximum ensemble query fairness violation γ_{ens} , the achieved ϵ , and the number of queries answered. The right figure plots the Pareto-frontier. With increasing privacy budget, more queries can be answered. The same holds when loosening the fairness constraint. At small privacy budgets, small fairness constraint might not be achievable.

Appendix F plots the results on UTKFace, and the results on other datasets are in Figure 13. The graph shows the effect of varying the consensus threshold T and fairness violation threshold ρ_{fair} on the number of queries answered. We observe that decreasing T leads to a higher number of queries answered. Similarly, increasing ρ_{fair} to a certain extent also leads to more queries being answered. Once the fairness violation threshold is too large, further relaxing the constraint would not lead to answering more queries, at which point no more queries are rejected due to the fairness constraint. Furthermore, at a fairness constraint of 0, there is a sharp decrease in the number of queries answered. The reason behind this is that if no fairness violation is allowed, no more queries can be answered after the fairness gap reaches 0, as any additional query would break the balance and increase the gap.

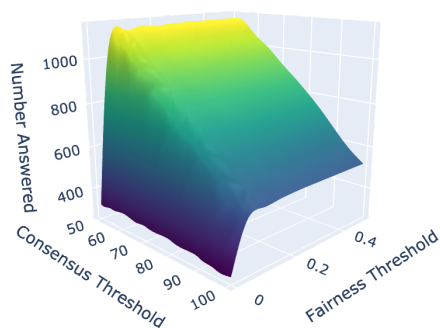


Fig. 11. **Query Experiment Results on UTKFace.** The figure plots the effect of consensus threshold T and fairness threshold $\gamma_{threshold}$ on the number of queries answered. We observe that the number of queries answered increases with smaller T and larger $\gamma_{threshold}$.

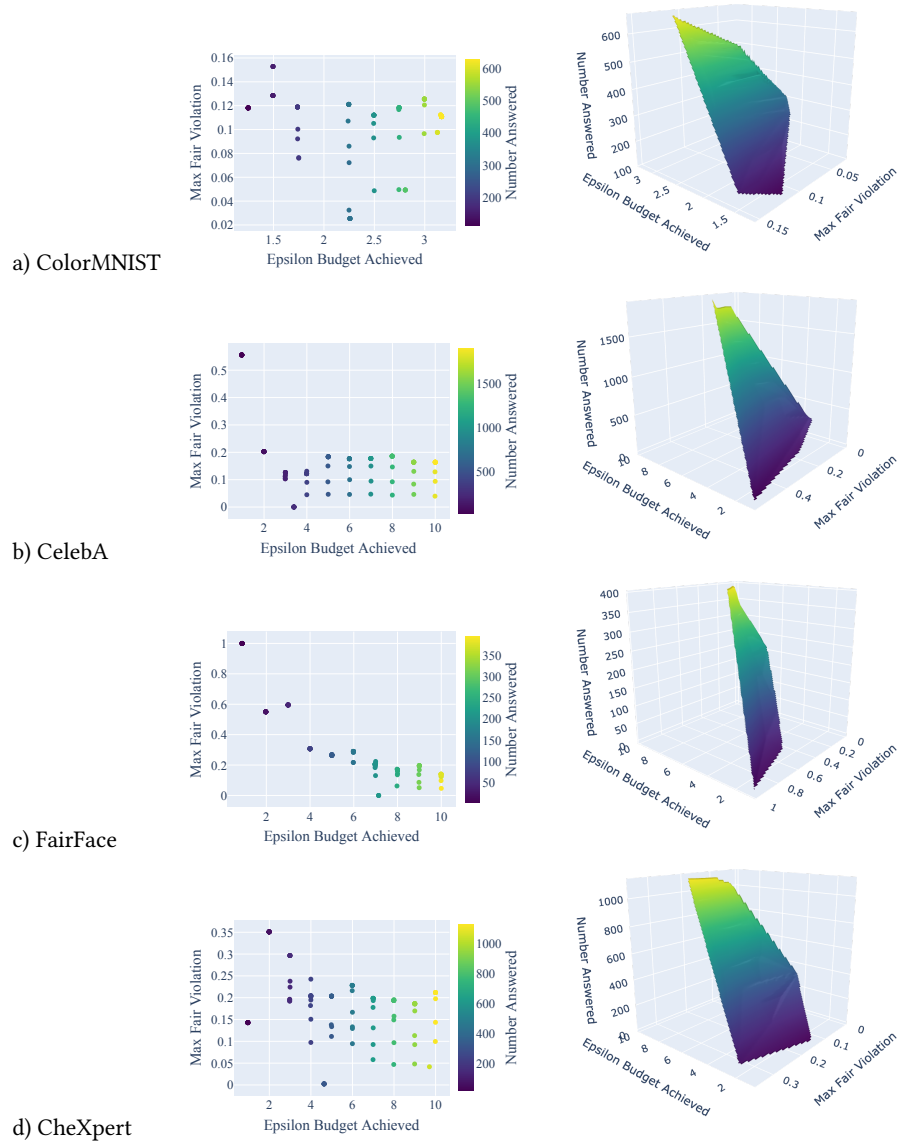


Fig. 12. Query experiments on other datasets. Setup described in Table 2 and discussion is in Appendix F.

G EXTENDED RELATED WORK: INTEGRATING FAIRNESS INTO PRIVATE LEARNING

In the literature, different fairness notions have been implemented within DP-SGD and PATE frameworks.

Fairness and DP-SGD. It has been shown that training with DP-SGD leads to disparate accuracy decrease over different data sub-groups [12, 25]. In particular, model accuracy decreases more for underrepresented data from the tails of the distribution [25]. Farrand et al. [12] presented similar findings and observed that privacy can even have

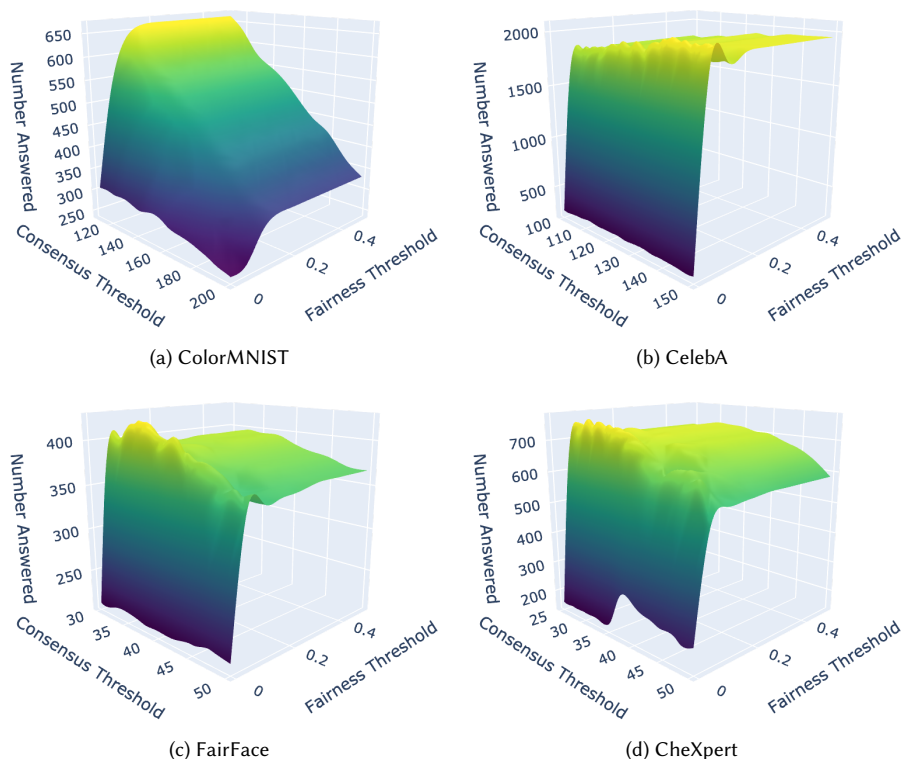


Fig. 13. Query experiments on other datasets. Setup described in Table 2 and discussion is in Appendix F. We found that in order to obtain the best results on student accuracy, some datasets require addition of significant noise σ_1 , which leads to differences in surfaces' shapes.

a negative impact on the model fairness when the training data is only slightly imbalanced. As potential reasons for this, the authors identified the clipping operation in DP-SGD. Since underrepresented data has larger gradients, these gradients are more effected by the clipping operation, and thereby, this data experiences a higher information loss [12]. To limit this effect, Xu et al. [32] proposed adapting the clipping threshold in DP-SGD individually for each sensitive group. They showed how their approach limits the disparate impact of DP-SGD on different groups. However, due to higher information leakage form larger gradients, their method requires larger perturbations. In a similar vein, Zhang et al. [33] propose early stopping to mitigate the negative impact of DP-SGD on model fairness. The authors observe that DP-SGD makes ML model training less stable which they leverage to interrupt training once high-enough fairness is achieved, without a significant loss in accuracy. However, all these methods solely manage fairness as an indirect byproduct of adapting the private training mechanism. Neither of them integrates explicit fairness constraints to yield formal guarantees, such as done in this work.

Tran et al. [28] proposed applying a Lagrangian dual approach for solving the joint optimization of fairness and privacy in ML. Therefore, they rely on a fairness constraint plus adaptive clipping and make the computations of the primal and dual update steps differentially private w.r.t. the considered sensitive attributes. However, their method adds

a significant computational overhead, especially for larger ML models and mini-batch sizes (increase of up to factor 100).

Fairness and PATE. When comparing the fairness impact of DP-SGD and PATE, Uniyal et al. [29] observed that PATE induces lower accuracy parity. The authors reason that this might be because the diversity among the teachers allows to cancel out their individual fairness issues. However, their observations only hold for very small numbers of teachers (10, in contrast to 250 proposed for MNIST in the original PATE paper [21]). This however yields sub-optimal privacy-utility trade-offs since in PATE, stronger privacy guarantees can be obtained when using more teachers which allows for the injection of more noise. In the work closest to ours, Tran et al. [26] study fairness properties of PATE and identified both algorithmic properties of the training (number of teachers, regularizer, privacy noise), and properties of the student data (magnitude of the input norm, and distance to the decision boundary) as factors influencing prediction fairness. To mitigate tensions, they proposed releasing the teacher models' prediction histogram as *soft labels* to train the student model. However, it has been shown that releasing the histograms leaks significant amounts of private information [31], which makes their method leaks privacy above the promised DP guarantees. In contrast, in this work, we integrate fairness in the aggregation process while keeping the teachers' votes private, and, thereby providing the promised privacy guarantees.