

Efficient Hyperparameter Optimization for Differentially Private Deep Learning

Aman Priyanshu

Manipal Institute of Technology
aman.priyanshu@learner.manipal.edu

Fatemehsadat Miresghallah
University of California, San Diego
fmireshg@eng.ucsd.edu

Rakshit Naidu

Carnegie Mellon University
rnamakal@andrew.cmu.edu

Mohammad Malekzadeh
Imperial College London
m.malekzadeh@imperial.ac.uk

ABSTRACT

Tuning the *hyperparameters* in the differentially private stochastic gradient descent (DPSGD) is a fundamental challenge. Unlike the typical SGD, private datasets cannot be used many times for hyperparameter search in DPSGD; *e.g.*, via a grid search. Therefore, there is an essential need for algorithms that, within a given search space, can find near-optimal hyperparameters for the best achievable privacy-utility tradeoffs efficiently. We formulate this problem into a general optimization framework for establishing a desirable privacy-utility tradeoff, and systematically study three cost-effective algorithms for being used in the proposed framework: *evolutionary*, *Bayesian*, and *reinforcement learning*. Our experiments, for hyperparameter tuning in DPSGD conducted on MNIST and CIFAR-10 datasets, show that these three algorithms significantly outperform the widely used *grid search* baseline. As this paper offers a first-of-a-kind framework for hyperparameter tuning in DPSGD, we discuss existing challenges and open directions for future studies. As we believe our work has implications to be utilized in the pipeline of private deep learning, we open-source our code at <https://github.com/AmanPriyanshu/DP-HyperparamTuning>.

1 INTRODUCTION

Deep neural networks (DNNs) [9] can learn very useful patterns from large multi-dimensional datasets, enabling motivational applications; *e.g.*, in health [12, 21]. However, large amounts of training data are required for not only learning the near-optimal DNN parameters for the underlying task, but also for finding the right set of hyperparameters that enable appropriate learning. For a task defined on public datasets, the same data can be reused as many times as we wish. But, as every reuse of the available data comes at a price of some privacy loss, *hyperparameter tuning* has been a fundamental challenge for tasks defined on private datasets.

Differential Privacy (DP) [8] provides strong guarantees for the individuals participating in private datasets. DP restricts the maximum contribution of each sample on the result of a computation on the private dataset. Differentially-private stochastic gradient descent (DPSGD) [1] is a widely accepted algorithm for training DNNs on private datasets, where zero-mean Gaussian noise, with a predefined variance, is added to the clipped gradients computed for each sample in the training dataset at each iteration. Noisy gradients often result in a degraded accuracy for the trained DNN.

Previous works look at two variants: (1) optimizing privacy parameters of a private model for achieving comparable performance

to a non-private model and (2) providing privacy guarantees to reach moderate performance [17]. However, in practice, both hyperparameters and privacy parameters need to be optimized within the user-specified privacy budget. Thus, in this paper, we propose a systematic study for learning hyperparameters faster (constrained by a privacy budget) and with less privacy cost through four different optimization algorithms.

Although there is a wide range of hyperparameters that one can choose from in DPSGD (*e.g.*, noise multiplier, clipping factor, batch size, learning rate, etc.), in this paper, we specifically focus on two important hyperparameters: *noise multiplier* σ (*i.e.*, the standard deviation of the Gaussian noise) and *learning rate* η . We optimize for these two parameters specifically as the epsilon (ϵ *i.e.* privacy leakage) and validation loss (minimizing the validation loss) are highly dependent on the chosen values for σ and η , respectively. To this end, we study three cost-effective algorithms: *evolutionary*, *Bayesian*, and *reinforcement learning* and compare the results with the *grid search* base-line. We also display consistent results across these techniques and provide insights on which algorithm could pave the path towards better hyperparameter tuning in DPSGD. We also open-source our code to enable future practitioners' optimize for specific hyperparameters, according to their requirements.

2 RELATED WORK

The most widely used methods for hyperparameter tuning in deep learning are manual search, random search, and grid search [16]. Manual search refers to the manual tuning of hyperparameters by individuals experimenting on a deep neural network. This method is frequently used due to prior experiences and intuition. On the other hand, random search provides a path towards hyperparameter space exploration. However, it is non-exhaustive and may not be able to discover high-performing hyperparameters. Thus, grid search is utilized to provide a sufficient exploration within a restricted search space. Although, due to its non-adaptive nature (*i.e.*, the hyper-parameter sets selected to be evaluated are not selected using already available results), it utilizes abundant resources and requires significant computational time.

Not all common practices in training deep models are always directly applicable when we apply DPSGD. For instance, [13] shows that while *ReLU* is the most common activation function in conventional deep learning, for DPSGD a bounded activation function, such as *tanh* or tempered *sigmoid*, is more efficient. Also, [17] argues that while SGD hyperparameters and DP parameters are considered

independent in the original DPSGD paper [1], this is not a reasonable assumption. For example, by choosing a smaller batch size, we can achieve better privacy, but to keep accuracy well, we thus need to also reduce the learning rate accordingly. However, smaller learning rates usually slow down the convergence [4], thus we need more epochs and hence more privacy loss in DPSGD. Therefore, [17] suggests using a public dataset first, to find an appropriate DNN architecture with an optimized set of hyperparameters, and then train the model on the private dataset. However, [17] did not propose any method for searching over DPSGD hyperparameters.

Model selection in multivariate linear regression under the constraint of differential privacy is studied in [11]; based on penalized least squares and likelihood. Especially, [11] reports that under differential privacy, the procedure of model selection becomes more sensitive to the tuning parameters. Moreover, the appropriate choice of tuning parameters requires some additional information in the data, and it is mentioned as a future topic in [11] to develop differentially private methods to estimate these hyperparameters. [14] introduces AdaClip, which achieves the same privacy guarantee with much less added noise by using coordinate-wise adaptive clipping of the gradient. As the convergence of DPSGD depends on the variance of the gradient, AdaClip also improves the accuracy of the trained model. While such an adaptive clipping provides better tradeoffs, it needs to estimate the variance and thus introduces four new hyperparameters by itself, making our problem more complicated. Similarly, [2] introduces a method for adaptively tuning the clipping threshold to track a given “quantile” of the update norm distribution during training. Again, this method also needs to tune a new hyperparameter in the range of $[0, 1]$.

The DPareto algorithm is proposed in [3], where Bayesian optimization is used for hyperparameter tuning. The paper describes the Pareto Front and empirically validates its application using different neural network architectures across two datasets. Their study uses the multi-objective Bayesian optimizer to find the best hyperparameters, utilizing hypervolume to find the relative merit of different objectives. On the other hand, we use a single-objective Bayesian optimizer for our study, which reduces computational costs and aims to optimize the reward function defined by us.

[19] uses reinforcement learning to efficiently tune hyperparameters needed for quantization of deep neural networks and find the bit-widths for weights of each layer that would provide optimal computation-accuracy trade-off.

3 METHODOLOGIES

3.1 Problem Formulation

We consider the problem of training a DNN with a fixed architecture (i.e., the number, type, and size of each layer) using DPSGD. Let D_{train} denotes the training set and D_{valid} denote the validation set. Let $\mathcal{H} = \{h_1, \dots, h_N\}$ denotes the set of N hyperparameters that are used during training on D_{train} and have impact on both validation loss (val_loss) and privacy loss in DP (ϵ), on D_{valid} . To provide a general but customizable framework, we define *reward* as a weighted linear combination of val_loss and ϵ :

$$reward = \alpha^U \cdot (e^{-val_loss}) + \alpha^P \cdot (e^{-\epsilon}). \quad (1)$$

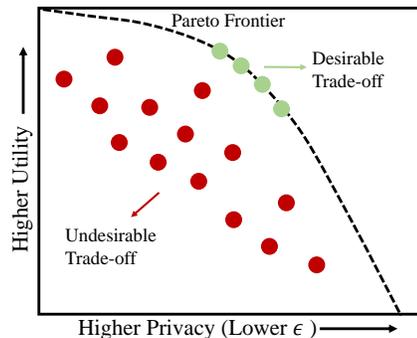


Figure 1: The Pareto frontier of potential privacy-utility trade-offs. Each point (green/red) shows the trade-off by choosing a specific set of values for the hyperparameters.

We use regularizers α^U and $\alpha^P \in [0, 1]$ to control the importance of utility and privacy, respectively (to control the privacy-utility trade-off). In our proposed framework, we first set these α regularizers, and then start searching for the optimal hyperparameters in \mathcal{H} using the algorithms explained in the following section. In this paper, we consider $\mathcal{H} = \{\sigma, \eta\}$, where σ denotes the noise multiplier and η denotes the learning rate; in DPSGD. Our aim for the following experiments remains to optimize the reward given by Equation (1). Notice that, in practice, the value of α^U and α^P depends on the requirements of the underlying task.

3.2 Evolutionary optimization

Evolutionary optimization algorithms provide an opportunity to explore as well as exploit the hyperparameter search space [6, 20]. Its random initialization and mutational attributes allow it to take advantage of the random search optimization algorithm. In contrast, its adaptive nature enables it to exploit critical values, which give better results. In our implementation, we encoded each hyperparameter as a gene, a set of hyperparameters made up the genome of an individual, i.e., the experiment. The range and precision for each hyperparameter are predetermined, allowing the optimization algorithm to search within a limited search space. The initial population is determined by random sampling of hyperparameters from this search space. Once aggregated, the population is trained, and a fitness score or reward is measured using 1. Subsequently, each generation is formed using selection, cross-over, and mutation based on the individuals with the highest fitness from the previous generation. The methodology can be optimized for high exploitation, thereby reducing resource wastage [20].

3.3 Bayesian Optimization

Bayesian optimization treats neural network training and performance as a black-box function. It combines prior experience with the black-box neural network with sample information to approximate the function distribution using the Bayesian formula [18]. Based on this estimation of the function distribution, optimal values can be extrapolated. The estimate distribution is effectively a probabilistic model for the function, which exploits this model to decide where to explore the function next while integrating out uncertainty [15]. This methodology allows one to find the minima

Table 1: Comparison of different methods based on the best-achieved reward and the average time required to attain this reward for the search space on (A) CIFAR-10 and (B) MNIST.

Method	Time (in hours)	Best Reward (in %)	Accuracy (in %)	Epsilon (ϵ)
(A) CIFAR-10				
Grid Search	150.020	51.406	44.936	0.600
Evolutionary	11.064	52.044	37.999	0.599
Bayesian	49.636	51.846	43.864	0.581
Reinforcement	52.971	52.398	44.884	0.590
(B) MNIST				
Grid Search	43.712	72.260	89.133	0.683
Evolutionary	5.250	72.615	73.745	0.175
Bayesian	2.853	73.385	81.562	0.349
Reinforcement	31.165	74.906	75.022	0.240

of complex non-convex functions with relatively few evaluations. However, this is only due to our assumption that the function is drawn from a Gaussian process prior. Our experiments utilized Hyperopt, a Sequential Model-Based Optimization (SMBO) that provides high performance at a low computational budget [5].

3.4 Reinforcement Learning

Evolutionary optimization and Bayesian optimization provide both adequate methodologies for search space exploration and exploitation. However, this classical problem can also be dealt with by reinforcement learning. In our application of this method, we begin by initializing a regression network capable of estimating the reward output of training on a particular set of hyperparameters. We start by sampling a random collection of hyperparameters used to train the DPSGD model and obtain the reward to fit the regression network. We then proceed to extract the estimated reward of the entire search space for our hyperparameter tuning. The best-performing hyperparameters are obtained from this estimation. These hyperparameters are mutated to hyperparameters in near proximity to them for the next episode, thereby allowing the model to exploit values that may give high performance. Subsequently, in the following episodes, we select a certain percent of experiments based on the reward estimate of the regression network; in contrast, the others continue to be randomly sampled. This is determined based on the epsilon-decreasing strategy, where the value of exploration-exploitation-epsilon decreases as the experiment progresses. This methodology allows us to estimate the hyperparameter-reward function and verify the proximal search space of high-performing hyperparameters, giving us generalized results.

4 EVALUATION

For the experiments in this section, we used a Tesla P100 16GB as GPU, with 13GB RAM Intel Xeon as CPU for our experiments. Note that the random seed is fixed across all experiments for uniformity and reproducibility purposes. In the rest of this section, we will discuss the benchmarks used and the results of each experiment.

4.1 Benchmarks

To assess and analyze the effectiveness of optimization algorithms across both CIFAR-10 and MNIST datasets, we use Grid-Search on a similar search complexity as the other methods. We display the computational time taken, best reward achieved, and its respective accuracy and epsilon value in Table 1. Grid search displays a poor understanding of the epsilon-accuracy as it is not adaptive in nature. It achieves a reward of 72.2% and 51.4% on the MNIST and CIFAR-10 datasets, respectively.

4.2 Optimization Algorithms

As described in earlier sections, we ran our experiments over three distinct optimization algorithms. Here, we observe that although Reinforcement Learning provides the highest performance, it comes at the expense of computational time. On the other hand, Evolutionary Algorithms and Bayesian Optimization provide consistent results with respect to computational time and performance.

Additionally, we can see from Table 1 that although Grid Search returns a highly accurate model, it is compensated by the high privacy leakage that occurs due to it. Contrary, adaptive optimization algorithms can leverage previous samples to search for a better privacy-utility tradeoff, allowing them to achieve high rewards. We see that Evolutionary Optimization achieves the lowest epsilon value for the MNIST dataset. In contrast, Bayesian Optimization achieves the same for the CIFAR-10 dataset.

4.3 Evaluating Computation Time for Satisfactory Reward

Although finding the best reward is our goal, we also evaluate the computational time required by each algorithm to achieve the maximum reward attained by Grid Search. The time consumed is calculated based on the time taken for an optimization algorithm to achieve a reward equal to or greater than the baseline reward. Here, baseline reward refers to the highest reward achieved by the Grid Search algorithm. We display these results in Figure 2. It can clearly be seen that grid search is much more time-consuming in nature, compared to other optimization algorithms.

From Figure 2, we can observe that Bayesian and evolutionary optimization algorithms display consistency and uniformity across datasets. Whereas, reinforcement-learning-based optimization fails to generalize the time taken to achieve a given reward. This is due to the exploration-exploitation trade-off mechanism applied by the aforementioned method. An appropriate exploration-exploitation-epsilon must be selected if the user intends to have a highly exploitative model for optimization.

4.4 Evaluation Sample-Specific Privacy

In this subsection, we look at sample-specific privacy. We count the total number of times that each sample has revisited for each algorithm during optimization. As training on any sample at a given time can lead to sensitive information leakage, we consider this an essential feature for privacy evaluation. In Figure 3, we display a bar graph representing the number of times any sample in the training set is visited before every model achieves its respective highest reward. These evaluations continue to display the impractical nature of Grid Search and validate its high privacy leakage. On

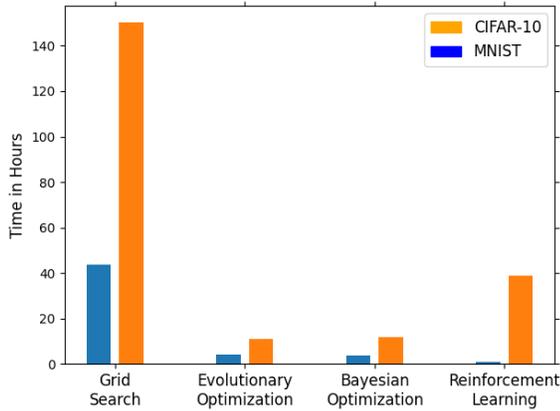


Figure 2: The time taken by different methodologies to achieve reward greater than or equal to baseline rewards.

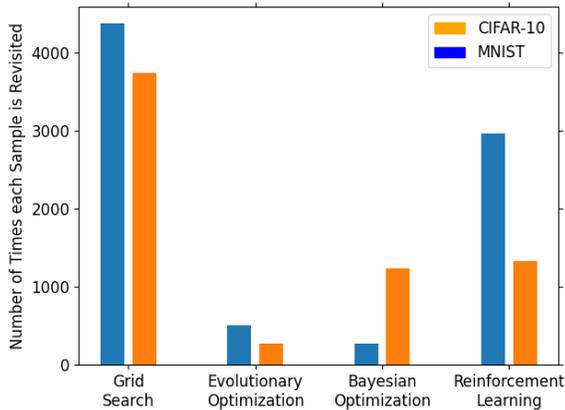


Figure 3: The number of times each sample is re-visited before the best reward is achieved.

the other hand, Bayesian and Evolutionary Optimization continues to give expected results compared to the prior. The Reinforcement Learning approach again displays high variance across datasets, giving it a weaker generalization capacity.

4.5 Learning and Convergence Analysis of Reinforcement Learning Approach

We further study the behavior of the reinforcement learning approach in our analysis through Figure 4. As the model learns over newer data every epoch, we can observe that we continue to adapt the expected reward accordingly. The model makes more resolute changes to the reward estimate nearest to the previous global maxima. This allows the model to exploit better performing hyperparameter values, allowing it to restrict search within a high-performing area. However, this also supplements that as the number of epochs increases, estimates which remain unexplored do not change in value. Therefore, an experimentally robust exploration-exploitation-epsilon must be selected for generalized results.

On comparing the two datasets, we can see the more complex nature of the search space for the MNIST dataset. This can be attributed to the simple complexity of the dataset, which allows learning over different hyperparameters. However, the CIFAR-10 dataset is much more complex, leading to only the most optimal hyperparameters being highlighted.

5 CONCLUSION AND FUTURE WORK

In this paper, we discussed different methodologies for hyperparameter tuning for the private training of deep neural networks using DPSGD algorithm. We proposed a novel, customizable reward function that allows users to define a single objective function for establishing their desired privacy-utility tradeoff. We quantified, compared, and analyzed the methods of grid search (as the baseline), Bayesian optimization, evolutionary optimization, and reinforcement learning, across two datasets, CIFAR-10, and MNIST. We observed that Bayesian and evolutionary optimization behave similarly in terms of the privacy-utility trade-off point they provide, and how efficiently they find it. Reinforcement learning, however, provides a more desirable trade-off but with varying efficiencies across datasets. All three methods perform much better than the baseline grid search algorithm. We believe that our work serves as a valuable resource for privacy-preserving ML practitioners, developers, and researchers for hyperparameter tuning.

For future work, one can use our proposed method alongside that of [7, 10], where a portion of the privacy budget is allocated to finding the appropriate learning rate on the private dataset. Another direction is to extend our proposed method to tune other hyperparameters in DPSGD, and even the network architecture and non-linear activation functions that are used.

Acknowledgement

Mohammad Malekzadeh was partially supported by the UK EPSRC (grant no. EP/T023600/1) within the CHIST-ERA program.

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] Galen Andrew, Om Thakkar, H Brendan McMahan, and Swaroop Ramaswamy. 2019. Differentially private learning with adaptive clipping. *arXiv preprint arXiv:1905.03871* (2019).
- [3] Brendan Avent, Javier González, Tom Diethe, Andrei Paleyes, and Borja Balle. 2020. Automatic Discovery of Privacy–Utility Pareto Fronts. *Proceedings on Privacy Enhancing Technologies* 2020, 4 (2020), 5–23.
- [4] Yoshua Bengio. 2012. Practical recommendations for gradient-based training of deep architectures. In *Neural networks: Tricks of the trade*. Springer, 437–478.
- [5] James Bergstra, Dan Yamins, and David D. Cox. Hyperopt: A Python Library for Optimizing the Hyperparameters of Machine Learning Algorithms. (????).
- [6] Erik Bochinski, Tobias Senst, and Thomas Sikora. 2017. Hyper-parameter optimization for convolutional neural network committees based on evolutionary algorithms. In *2017 IEEE International Conference on Image Processing (ICIP)*. 3924–3928. <https://doi.org/10.1109/ICIP.2017.8297018>
- [7] Chen Chen and Jaewoo Lee. 2020. Stochastic adaptive line search for differentially private optimization. In *2020 IEEE International Conference on Big Data (Big Data)*. IEEE, 1011–1020.
- [8] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. 2006. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*. Springer, 265–284.
- [9] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep learning*. MIT press.
- [10] Jaewoo Lee and Daniel Kifer. 2018. Concentrated differentially private gradient descent with adaptive per-iteration privacy budget. In *Proceedings of the 24th*

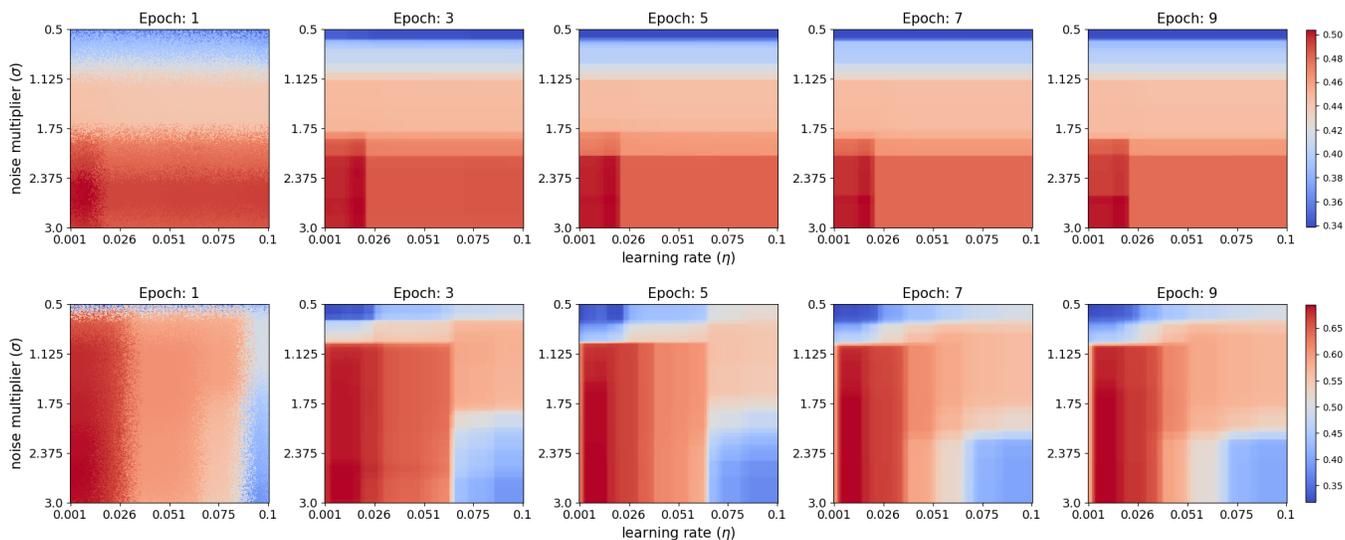


Figure 4: Convergence of best reward estimate of the ANN model, based on input *noise_multiplier* (σ , Y-Axis) and *learning_rate* (η , X-Axis). We display model convergence for the CIFAR-10 dataset on top and the MNIST dataset at the bottom.

- ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 1656–1665.
- [11] Jing Lei, Anne-Sophie Charest, Aleksandra Slavkovic, Adam Smith, and Stephen Fienberg. 2018. Differentially private model selection with penalized and constrained likelihood. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 181, 3 (2018), 609–633.
- [12] Mohammad Malekzadeh, Burak Hasircioglu, Nitish Mital, Kunal Katarya, Mehmet Emre Ozfatura, and Deniz Gündüz. 2021. Dopamine: Differentially Private Federated Learning on Medical Data. *The Second AAAI Workshop on Privacy-Preserving Artificial Intelligence (PPAI-21)* (2021).
- [13] Nicolas Papernot, Abhradeep Thakurta, Shuang Song, Steve Chien, and Ulfar Erlingsson. 2021. Tempered Sigmoid Activations for Deep Learning with Differential Privacy. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- [14] Venkatesh Pichapati, Ananda Theertha Suresh, Felix X Yu, Sashank J Reddi, and Sanjiv Kumar. 2019. AdaClip: Adaptive clipping for private SGD. *arXiv preprint arXiv:1908.07643* (2019).
- [15] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. 2012. Practical Bayesian Optimization of Machine Learning Algorithms. In *Advances in Neural Information Processing Systems*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger (Eds.), Vol. 25. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2012/file/05311655a15b75fab86956663e1819cd-Paper.pdf>
- [16] Florian Tramer and Dan Boneh. 2021. Differentially Private Learning Needs Better Features (or Much More Data). In *International Conference on Learning Representations (ICLR)*.
- [17] Koen Lennart van der Veen, Ruben Seggers, Peter Bloem, and Giorgio Patrini. 2018. Three tools for practical differential privacy. *arXiv preprint arXiv:1812.02890* (2018).
- [18] Jia Wu, Xiu-Yun Chen, Hao Zhang, Li-Dong Xiong, Hang Lei, and Si-Hao Deng. 2019. Hyperparameter Optimization for Machine Learning Models Based on Bayesian Optimization. *Journal of Electronic Science and Technology* 17, 1 (2019), 26–40. <https://doi.org/10.11989/JEST.1674-862X.80904120>
- [19] Amir Yazdanbakhsh, Ahmed T Elthakeb, Prannoy Pilligundla, F Mireshghallah, and Hadi Esmailzadeh. 2018. Releq: An automatic reinforcement learning approach for deep quantization of neural networks. *arXiv preprint arXiv:1811.01704* 1, 2 (2018).
- [20] Steven R. Young, Derek C. Rose, Thomas P. Karnowski, Seung-Hwan Lim, and Robert M. Patton. 2015. Optimizing Deep Learning Hyper-Parameters through an Evolutionary Algorithm. In *Proceedings of the Workshop on Machine Learning in High-Performance Computing Environments (MLHPC '15)*. Association for Computing Machinery, New York, NY, USA, Article 4, 5 pages. <https://doi.org/10.1145/2834892.2834896>
- [21] Alexander Ziller, Dmitrii Usynin, Nicolas Remerscheid, Moritz Knolle, Marcus Makowski, Rickmer Braren, Daniel Rueckert, and Georgios Kaissis. 2021. Differentially private federated deep learning for multi-site medical image segmentation. *arXiv preprint arXiv:2107.02586* (2021).