
Architecture Matters: Investigating the Influence of Differential Privacy on Neural Network Design

Felix Morsbach
Institute AIFB
Karlsruhe Institute of Technology
Karlsruhe, Germany
felix.morsbach@kit.edu

Tobias Dehling
Institute AIFB, KASTEL Security Research Labs
Karlsruhe Institute of Technology
Karlsruhe, Germany
dehling@kit.edu

Ali Sunyaev
Institute AIFB, KASTEL Security Research Labs
Karlsruhe Institute of Technology
Karlsruhe, Germany
sunyaev@kit.edu

Abstract

One barrier to more widespread adoption of differentially private neural networks is the entailed accuracy loss. To address this issue, the relationship between neural network architectures and model accuracy under differential privacy constraints needs to be better understood. As a first step, we test whether extant knowledge on architecture design also holds in the differentially private setting. Our findings show that it does not; architectures that perform well without differential privacy, do not necessarily do so with differential privacy. Consequently, extant knowledge on neural network architecture design cannot be seamlessly translated into the differential privacy context. Future research is required to better understand the relationship between neural network architectures and model accuracy to enable better architecture design choices under differential privacy constraints.

1 Introduction

Differential privacy has become the de facto standard for achieving data confidentiality in machine learning settings. The most prominent way to train differentially private neural networks is by clipping gradients in order to limit the impact of each data point and by adding noise to the gradient updates [2]. While gradient clipping is also used in non-private training to avoid overfitting [18], adding noise inherently reduces the accuracy of the machine learning model [10]. This is considered as a trade-off between utility and differential privacy guarantees and is usually measured as accuracy loss, describing the difference between the accuracy with and without differential privacy constraints. Even on comparably simple tasks, such as the CIFAR-10 classification task, the accuracy loss can be significant. This hinders the adoption of differentially private neural networks.

It is well known, that the architecture of a neural network can have a significant influence on the accuracy of the model. For example, expanding the size of a network in depth rather than width is usually said to increase accuracy [7], but also novel architectural features such as residual or dense connections can vastly improve model accuracy [8, 9]. Yet, early works on differentially private neural networks do not seem to account for this in great detail. They either do not include different architectures in their benchmarks [10] or find that the effect is not significant [2]. There is, however, also contrary evidence. For instance, expanding the overall size of the network in a differentially private setting exhibits an inflection point [14]. Increasing the network size beyond this inflection

point reduces model accuracy, which is not the case in the plain setting. Furthermore, the choice of an activation function plays an important role. Bounded activation functions (e. g., tanh) consistently outperform unbounded ones (e. g., ReLU) due to the phenomena of exploding activations during differentially private training [15].

Nevertheless, even with these findings taken into account, there still remains a significant gap between the accuracy of differentially private and non-private neural networks, even on simple tasks. Extant research has already shown that the choice of an appropriate activation function differs between the non-private and differentially private setting and that choosing an activation function appropriate for differentially private neural networks can reduce the incurred accuracy loss [15]. However, there are a multitude of other architectural features, such as the number, ordering, type and configuration of layers, whose impact on model accuracy under differential privacy constraints has not been really understood. Consequently, we investigate whether architectural features other than the activation function of a neural network affect the accuracy loss incurred by differential privacy constraints.

If we can show that more architectural features than the activation function impact model accuracy, this would imply that the accuracy of differentially private models can be improved by carefully tuning a neural network architecture for differential privacy instead of simply copying architectures and hyperparameters that work well without differential privacy constraints [17].

We designed and carried out an experiment to test whether the architecture affects accuracy loss under differential privacy constraints. Our findings show, not only, that the network architecture has an influence on the accuracy loss incurred by differential privacy constraints, but also, that the suitability of architecture choices is sensitive to variations in the targeted level of differential privacy. We elaborate the implications of our findings and argue what future research is necessary to reduce the accuracy loss incurred by differential privacy constraints in order to improve the applicability of privacy-preserving machine learning.

2 Background

Differentially private stochastic gradient descent (DP-SGD) is an algorithm for training neural networks with (ϵ, δ) -differential privacy guarantees [2, 5]. DP-SGD extends classic stochastic gradient descent in two ways. First, gradients are clipped on a per-example basis to a l_2 norm, which is set through a clipping threshold C . Second, random noise is added to the gradient update calibrated via the standard deviation σ , also called noise multiplier. Thus, DP-SGD adds two additional hyperparameters to the training algorithm.

The privacy level (ϵ, δ) of a model trained with DP-SGD is dependent on the batch size, the noise multiplier, the number of epochs, and the number of training examples. Hence, the privacy level for a given set of these hyperparameters can be calculated upfront without the need to actually train the model. Since the privacy level is independent of the architecture of a given model, different architectures can be compared at a fixed privacy level.

3 Experiments

In order to answer our research question, we need to show that the architecture of a neural network does have an effect on the model's accuracy loss under differential privacy constraints. We hypothesize that given a set of neural architectures for a given machine learning task, the architecture which performs best without differential privacy constraints does not necessarily also performs best under differential privacy constraints. Or formulated differently, let there be two neural network architectures A_1 and A_2 and let $U(A)$ denote the accuracy of architecture A without differential privacy and $U_d(A)$ the accuracy with differential privacy. We assume the case exists that $U(A_1) > U(A_2)$ and $U_d(A_1) < U_d(A_2)$.

To test our hypothesis, we conducted an experiment on the standard CIFAR-10 image classification task which contains 60000 color images in 10 classes [11]. We used Tensorflow [1] and the Tensorflow-Privacy extension [3] as machine learning libraries for the implementation. For our experiments we chose 8 convolutional neural network architectures, including prominent examples from the literature such as the LeNet-5 [12], but also architectures from other well-cited works [4, 13, 16] and own creations. The architectures differ mostly in the size and number of convolution,

Table 1: The architectures used for the experiment

#1		#2	
Layer type	Parameters	Layer type	Parameters
Convolution	32 filters of 3x3, ReLU	Convolution	32 filters of 3x3, ReLU
Max-Pooling	2x2	Max-Pooling	2x2
Convolution	64 filters of 3x3, ReLU	Convolution	64 filters of 3x3, ReLU
Fully connected	64 units, ReLU	Max-Pooling	2x2
Softmax	10 units	Convolution	64 filters of 3x3, ReLU
		Fully connected	64 units, ReLU
		Softmax	10 units
#3		#4, [12] but with ReLU	
Layer type	Parameters	Layer type	Parameters
Convolution	32 filters of 3x3, ReLU	Convolution	6 filters of 5x5, ReLU
Max-Pooling	2x2	Avg-Pooling	2x2, stride 2
Convolution	64 filters of 3x3, ReLU	Convolution	16 filters of 5x5, ReLU
Max-Pooling	2x2	Avg-Pooling	2x2, stride 2
Convolution	64 filters of 3x3, ReLU	Fully connected	120 units, ReLU
Fully connected	128 units, ReLU	Fully connected	84 units, ReLU
Softmax	10 units	Softmax	10 units
#5		#6, [4]	
Layer type	Parameters	Layer type	Parameters
Convolution	32 filters of 5x5, ReLU	Convolution	32 filters of 3x3, ReLU
Avg-Pooling	2x2, stride 2	Convolution	32 filters of 3x3, ReLU
Convolution	64 filters of 5x5, ReLU	Max-Pooling	2x2
Avg-Pooling	2x2, stride 2	Dropout	0.25
Fully connected	200 units, ReLU	Convolution	64 filters of 3x3, ReLU
Fully connected	100 units, ReLU	Convolution	64 filters of 3x3, ReLU
Softmax	10 units	Max-Pooling	2x2
		Fully connected	512 units, ReLU
		Dropout	0.5
		Softmax	10 units
#7, [13]		#8, [16]	
Layer type	Parameters	Layer type	Parameters
Convolution	32 filters of 5x5, ReLU	Convolution	32 filters of 5x5, ReLU
Avg-Pooling	2x2, stride 2	Convolution	64 filters of 5x5, ReLU
Convolution	64 filters of 5x5, ReLU	Fully connected	384 units, ReLU
Avg-Pooling	2x2, stride 2	Fully connected	192 units, ReLU
Fully connected	512 units, ReLU	Softmax	10 units
Softmax	10 units		

pooling, and fully-connected layers. We chose convolutional architectures as an initial setting, as their architecture configuration space is more complex and interesting compared to simple feedforward networks with only fully connected layers. See Table 1 for a full definition of the architectures used. We trained each model architecture for 100 epochs with DP-SGD, with a batch size of 250, 5 micro batches, a fixed learning rate of 0.1, a clipping threshold of 1.0 and at two noise multipliers of 0.01 and 0.1.

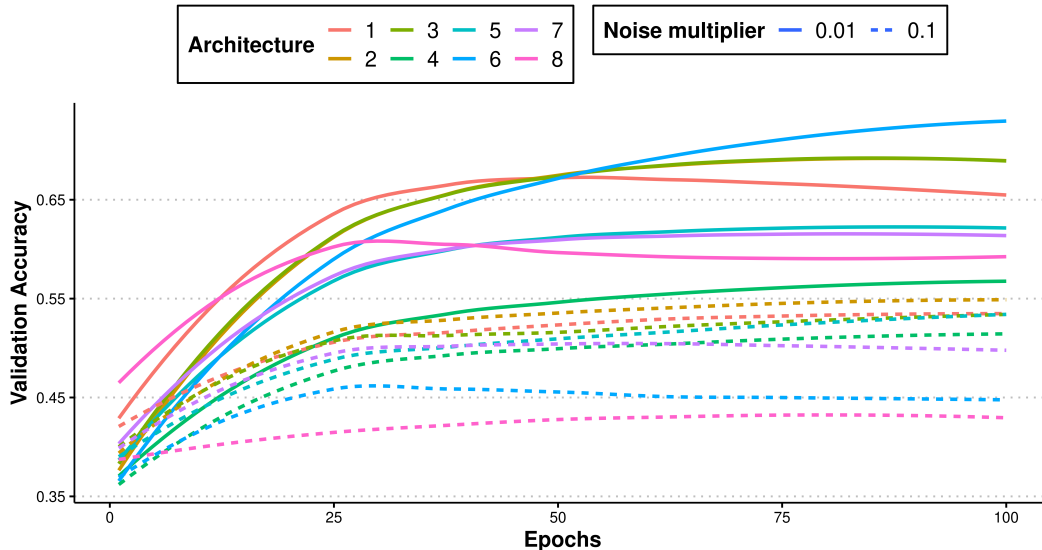


Figure 1: Results of training 8 neural architectures for 100 epochs on the CIFAR-10 image classification task with DP-SGD and a batch size of 250, 5 micro batches, a fixed learning rate of 0.1 and a clipping threshold of 1.0, at two different noise multipliers (0.01 and 0.1). (*Best viewed in color.*)

4 Results

Figure 1 shows the training results of the 8 different architectures on the CIFAR-10 machine learning task, repeated three times in identical settings with two different noise multipliers. We can clearly see that with a low noise multiplier model 6 performs best, but with a high noise multiplier model 6 comes in second to last while model 2 performs best. As the noise multiplier is one of the hyperparameters that determines the privacy level (besides epochs, batch size, and training set size, which are identical for all results), this means that at different privacy levels, different architectures perform best. In turn, this shows that, in a given set of architectures, it is not always the case that a single architecture will perform best for all differential privacy settings. Therefore, it is also not guaranteed that an architecture that performs best without differential privacy necessarily performs best in the differential privacy setting since its performance will be dependent on the targeted privacy level. The raw results and the code used to generate the graph can be found on GitHub¹.

5 Conclusion

In this paper, we investigated whether the architecture of a neural network can influence the accuracy loss incurred through differential privacy constraints. Our findings show that this is the case; neural network architectures that perform well in the non-private setting will not necessarily perform well in the differentially private setting. Furthermore, we found that which architecture performs best also depends on the chosen level of differential privacy.

The implications of our findings are two-fold: First, for research, our findings show that architectures are only comparable at the same privacy level, as the relative ranking of architectures might differ across privacy levels. Second, for practitioners, our findings show that best practices, experiences, or architectures from the non-private setting cannot be easily transferred to the design of differentially private neural networks. Rather, the model architecture has to be designed and tuned for specific differential privacy settings; differential privacy cannot be treated as an afterthought.

Moreover, the design of neural networks is a challenging task that either takes a lot of resources to try many different architectures through neural architecture search [6] or requires a good understanding of neural network design from the modeler. Doing a neural architecture search in the plain setting requires additional computations, but an extensive search will not hurt the final accuracy of the

¹<https://github.com/FMorsbach/ArchitectureMatters>

model. Searching for an architecture or optimizing hyperparameters in the differential privacy setting does, however, consume privacy budget [17]. Therefore, spending more time on architecture search or hyperparameter optimization will decrease the privacy budget available for the actual training, which will decrease the number of epochs available for training; hence, it will probably reduce the accuracy of the final model. As a consequence, neural architecture search is not easily applicable in the differential privacy setting.

Instead, a good understanding of how to design neural network architectures in the differential privacy setting is needed in order to maximize the available privacy budget for the actual training. But as shown by our findings, the experience from the design of neural network architectures without differential privacy constraints can only be transferred to the design of differentially private neural network architectures to a limited degree. Therefore, we need to derive new best practices for the design of architectures tailored specifically for differentially private neural networks, in order to improve the applicability of privacy-preserving machine learning by reducing the accuracy loss incurred by differential privacy constraints.

In our future work, we will derive best practices for the design of neural architectures under differential privacy constraints. We will set up additional experiments across multiple classic machine learning benchmarks. Subsequently, we will analyse the results, derive a set of candidate best practices and test those on different benchmarks. We aim to also incorporate qualitative data from expert interviews to even further advance the best practices for designing neural network architectures that perform well under differential privacy constraints.

Acknowledgments and Disclosure of Funding

This work was supported by funding from the topic Engineering Secure Systems of the Helmholtz Association (HGF) and by KASTEL Security Research Labs.

References

- [1] Martin Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: A system for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, pages 265–283, Savannah, GA, November 2016. USENIX Association.
- [2] Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep Learning with Differential Privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, CCS '16*, pages 308–318, New York, NY, USA, October 2016. Association for Computing Machinery.
- [3] Galen Andrew, Steve Chien, and Nicolas Papernot. TensorFlow Privacy. <https://github.com/tensorflow/privacy>, 2018.
- [4] Pathum Chamikara Mahawaga Arachchige, Peter Bertok, Ibrahim Khalil, Dongxi Liu, Seyit Camtepe, and Mohammed Atiquzzaman. Local differential privacy for deep learning. *IEEE Internet of Things Journal*, 7(7):5827–5842, 2019.
- [5] Cynthia Dwork and Aaron Roth. The Algorithmic Foundations of Differential Privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3–4):211–407, August 2014.
- [6] Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Neural architecture search: A survey. *The Journal of Machine Learning Research*, 20(1):1997–2017, January 2019.
- [7] Ian J. Goodfellow, Yaroslav Bulatov, Julian Ibarz, Sacha Arnoud, and Vinay Shet. Multi-digit Number Recognition from Street View Imagery using Deep Convolutional Neural Networks. In *2nd International Conference on Learning Representations (ICLR)*, 2014.

- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, June 2016.
- [9] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q. Weinberger. Densely Connected Convolutional Networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2261–2269, July 2017.
- [10] Bargav Jayaraman and David Evans. Evaluating Differentially Private Machine Learning in Practice. In *28th USENIX Security Symposium (USENIX Security 19)*, pages 1895–1912, 2019.
- [11] Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. Technical report, University of Toronto, 2009.
- [12] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, November 1998.
- [13] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In Aarti Singh and Jerry Zhu, editors, *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pages 1273–1282, Fort Lauderdale, FL, USA, April 2017. PMLR.
- [14] Nicolas Papernot, Steve Chien, Shuang Song, Abhradeep Thakurta, and Úlfar Erlingsson. Making the Shoe Fit: Architectures, Initializations, and Tuning for Learning with Privacy, September 2019.
- [15] Nicolas Papernot, Abhradeep Thakurta, Shuang Song, Steve Chien, and Úlfar Erlingsson. Tempered Sigmoid Activations for Deep Learning with Differential Privacy. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence*, 2021.
- [16] NhatHai Phan, Xintao Wu, Han Hu, and Dejing Dou. Adaptive laplace mechanism: Differential privacy preservation in deep learning. In *2017 IEEE International Conference on Data Mining (ICDM)*, pages 385–394. IEEE, 2017.
- [17] Koen Lennart van der Veen, Ruben Seggers, Peter Bloem, and Giorgio Patrini. Three Tools for Practical Differential Privacy. *arXiv:1812.02890 [cs, stat]*, December 2018.
- [18] Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why Gradient Clipping Accelerates Training: A Theoretical Justification for Adaptivity. In *Eighth International Conference on Learning Representations*, April 2020.