

	Model	DCGAN			WGAN			WGAN-GP		
		5	10	20	5	10	20	5	10	20
% successful	Regular Adam	48.3	68.7	80.0	56.0	84.3	90.3	47.0	64.7	64.7
	Surfing	78.3	98.7	96.3	81.7	97.3	99.3	83.7	95.7	97.3
# iterations	Regular Adam	618	4560	18937	464	1227	3702	463	1915	15445
	Surfing	741	6514	33294	547	1450	4986	564	2394	25991

Table 1: Surfing compared against direct gradient descent over the final trained network. Shown are percentages of “successful” solutions  $\hat{x}_T$  satisfying  $\|\hat{x}_T - x_*\| < 0.01$ , and 75th-percentiles of the total number of gradient descent steps used (across all networks  $G_0, \dots, G_T$  for surfing) until  $\|\hat{x}_T - x_*\| < 0.01$  was reached.

1 We thank the reviewers for carefully reading our paper and providing insightful and constructive comments. We will  
2 respond to each of the concerns that were raised.

3 *Reviewers 1 and 2 both comment on the computational cost of the procedure, compared with running vanilla Adam with*  
4 *multiple random initial points.* We thank the reviewers for raising this important point, which led us to further explore  
5 the computational cost of surfing. In fact, surfing can be performed such that its runtime is close to that of a *single*  
6 initialization of vanilla Adam—the reason is that for the intermediate networks, gradient descent (GD) does not need to  
7 be run until full convergence; the number of GD steps can be quite small and surfing will still succeed.

8 The updated Table 1 illustrates this: Briefly, we re-ran both vanilla Adam and surfing on the DCGAN, WGAN, and  
9 WGAN-GP examples, using the same step size in both methods. We recorded the 75th-percentile of the number of GD  
10 steps  $N$  needed in vanilla Adam to achieve  $\|\hat{x}_T - x_*\| < 0.01$ . We then constrained surfing to use  $N$  total iterations  
11 across networks  $G_0, \dots, G_{99}$ , followed by GD until convergence for the final trained network  $G_{100}$ . The  $N$  steps in  
12 surfing were split across networks  $G_0, \dots, G_{99}$  proportional to a common deterministic schedule, which allotted more  
13 steps to the earlier networks  $G_t$  where the landscape changes more rapidly, and fewer steps to later networks where  
14 this landscape stabilizes. Shown are the success rates and the 75th-percentiles of the total number of GD iterations for  
15 both methods. We see that surfing still has a much higher success rate, at a comparable computational cost to a single  
16 initialization for vanilla Adam. We will update Table 1 of the original manuscript to display this new comparison.

17 *R1: I only have a problem with the way the set  $S(x, \theta, \tau)$  is defined in line 177, since the authors do not require the*  
18 *signs to strictly differ on this set.*  $S(x, \theta, \tau)$  is just the set of neurons that are close to zero before ReLU thresholding.  
19 These are the neurons for which the signs could change after a small change of the network input  $x$ .

20 *R1: Although Algorithm 2 and the empirical algorithm are similar in spirit, lines 1 and 3 in algorithm 2 are crucial for*  
21 *proof of correctness.* Theorem 2 mainly illustrates that the procedure can be formalized, although in its current form the  
22 projected gradient algorithm is not easily implemented.

23 *R1: For the case where  $y = G(z) + \text{noise}$ , where noise has sufficiently low energy, you would expect a local minimum*  
24 *close to  $z$ . Would this not contradict the result of Theorem 3.1?* This case is not covered by Theorem 3.1, because  $y$  is  
25 then correlated with the network parameters. Please see our comment starting on line 157.

26 *R2: I find the paper quite interesting already. To make it even more interesting would involve having a complete*  
27 *theoretical argument establishing the time complexity without the current heuristic.* We agree that a full theoretical  
28 analysis would be preferred. Ultimately we think that something between the simple surfing and projected gradient  
29 surfing methods will be more attractive in both theory and practice.

30 *R3: From my understanding, the first theorem is mainly built on (Hand and Voroninski, 2017), and the second theorem*  
31 *is mainly built on (Bora et al.)* Our analysis builds primarily on Hand and Voroninski. The type of result in Bora et al.  
32 is different, and pertains to properties of near-global minimizers rather than computational procedures for finding them.

33 *R3: For the second theorem, the result implies the deeper the network is, the smaller the delta should be. It would be*  
34 *better to discuss how tight is the analysis, and whether this dependency is necessary in practice.* The dependence of  
35  $\delta$  on network depth comes from upper-bounding the Lipschitz constant of the network  $G(x)$  by  $\prod_{i=1}^d \|W_i\|$ . We do  
36 expect the true Lipschitz constant to increase with network depth in practice. The upper-bound is likely not tight, but it  
37 may be difficult to theoretically improve. The same type of bound was used in Szegedy et al. (2014); Virmaux and  
38 Scaman (2018) which discussed this question in more detail—we will add a discussion of this point to the manuscript.

## 39 References

- 40 Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. (2014). Intriguing  
41 properties of neural networks. In *International Conference on Learning Representations*.
- 42 Virmaux, A. and Scaman, K. (2018). Lipschitz regularity of deep neural networks: Analysis and efficient estimation. In  
43 *Advances in Neural Information Processing Systems*, pages 3835–3844.