

1 **R1, 7/10** *Although the maximum speedups over [18] in matrix inversion is quite sizable (27x in Fig 1), the gain over*  
2 *the much simpler, standard baseline seems moderate (2.7x to 4.3x). This potentially weakens the significance of the*  
3 *contributions, unless the proposed approach has advantages in other dimensions compared to the baseline.*

4 Indeed, the approach has advantages in other dimensions. FastH attains 5x to 27x speed-up for orthogonal matrices as  
5 used by [4,6,7,8,10,14,15,16,18], i.e. we can expect 5x to 27x if we use FastH in the referenced approaches (Figure 3).  
6 In particular, we found FastH to be so much faster than previous methods, that it can speed up even matrix inversion,  
7 something that was not possible with previous methods. We updated the introduction to reflect this distinction.

8 **R4, 5/10** *The cost of searching proper  $k$  is  $O(d^3)$  which is 'negligible' which makes me confused about which is not*  
9 *negligible in designed algorithm.*

10 Thanks for raising this concern. The default  $k = \text{batch\_size}$  used in our experiments works well and one does not need  
11 to search for  $k$ . However, one might be able to improve slightly by trying different  $k$ . This only needs to be done **once**  
12 for a given hardware setup. This differs to our algorithm, FastH, which is used **every step during training**, i.e., FastH  
13 is used  $10^5$  times while the search is done 1 time. We rephrased a few sentences in section 3.3 to clarify this confusion.

14 *The result in section 4.2, the sequential method spend too much time on inner products isn't mentioned in section 3.*

15 We opted to use Section 3 to introduce our FastH algorithm, while clarifying the mentioned issue in the introduction  
16 (L22-28) and the background section (L66-70). We believe the current organization allows for a sharper separation  
17 between previous work and our proposed algorithm FastH.

18 *There is a gap between theoretical time-complexity and empirical time-complexity which makes the analysis of time-*  
19 *complexity in section 3 can't support the effectiveness of the designed algorithm.*

20 We believe there is a misconception here. “FastH retains the **same desirable time complexity** as the sequential  
21 algorithm from [18] while reducing the number of sequential operations” (introduction L31-33). In other words, FastH  
22 is  $27x$  faster than [18] due to less sequential work, **not** due to a difference in time complexity.

23 **R2, 6/10** *Limited applicability: seems the technique applies only to layers whose #input neurons = #output neurons.*

24 Thanks for raising this concern. Our technique does apply when the number of input neurons  $n$  is different to number of  
25 output neurons  $m$ , that is, for a regular linear layer with weight matrix  $W \in \mathbb{R}^{n \times m}$ . The weight matrix has a singular  
26 value decomposition  $W = U\Sigma V^T$  for orthogonal  $U, V$  where  $U \in \mathbb{R}^{n \times n}$ ,  $V \in \mathbb{R}^{m \times m}$  and diagonal  $\Sigma \in \mathbb{R}^{n \times m}$ .  
27 FastH works for both  $U$  and  $V$ . This can furthermore be extended to attain semi-orthogonal  $W$ .<sup>1</sup> We added a paragraph  
28 in the subsection concerning extensions with the hopes that it clarifies this confusion.

29 *The empirical methodology is a bit problematic since it does not explore deep learning at all, but rather on the time to*  
30 *do a single step and the cost of various matrix operations.*

31 Thanks for raising this concern. The use of Householder matrices in deep learning has received much attention in  
32 previous work, e.g., [6,10,14,16,18]. FastH computes **exactly the same** as the algorithm used by [6,10,14,16,18];  
33 repeating their experiments with FastH would thus attain the same results, albeit faster. Since we believe [6,10,14,16,18]  
34 adequately demonstrate the usefulness of Householder matrices for deep learning, we found the additional value of  
35 more such experiments were not that high. Furthermore, there are additional benefits to studying the performance of  
36 single operators as opposed to end-to-end deep learning experiments. Time complexity is a more transparent measure  
37 to investigate than the validation loss of deep learning models. Such measure shows the benefits of our approach  
38 **irrespective** of the architecture, optimizer, loss function and the many hyperparameters of complex networks.

39 *Novelty: There is very little in the way of a fundamentally new idea. ... . The authors simply adjust the tool to the job.*

40 As evidence against the claimed lack of novelty, we present three articles that would benefit from a  $O(dm^2)$  parallel  
41 algorithm but did not “simply adjust the tool to the job.” [18] realized the issue with sequential computation and  
42 suggested a parallel  $O(d^3)$  algorithm. [10] uses the related CWY decomposition for gradient computations, without  
43 realizing WY can increase GPU utilization. They instead speed up by using fewer Householder matrices. An  
44 article contemporaneous<sup>2</sup> to our submission (<https://arxiv.org/abs/2004.08675>) addresses parallel Householder  
45 products with the related CWY decomposition, but attains a  $O(d^3)$  bound (see their table 2, serial complexity is  $L^3$ ).

46 **R3, 7/10.** *I would like to ask the authors to provide more details regarding the experimental setup to help reproducibility.*

47 We will soon open-source “neuralsvd.py” from the supplementary material, which we updated to run our main experi-  
48 ment and draw Figure 3. We also updated “README.txt” to contain more details regarding the experimental setup.

<sup>1</sup>Semi-orthogonal means  $W^T W = I \neq W W^T$ . This is true if  $n > m$  and  $\Sigma_{ii} = 1$ .

<sup>2</sup>Contemporaneous, as per NeurIPS guidelines, refers to a work published less than two months before the submission deadline of which the authors are retroactively not aware.