

1 We thank the reviewers for their careful reading and constructive comments. We feel that the reviews are largely positive.  
2 In the remainder, we want to address some of the issues raised, and we will address them in detail in the revision. We  
3 will also release the code for the numerical experiments.

4 "Why would one care about the variable sketches, rather than the fixed sketch", "It's fine if this paper's method and  
5 analysis is not the best for chronological and historical reasons, however please state of the art clear and provide a  
6 reference." *This reference appeared after our work. It is based on the IHS with a fixed SRHT embedding, and their  
7 convergence rate appears currently the best known (in the asymptotic sense) for the SRHT. Differently, we emphasize  
8 that our analysis provides an **exact and closed-form formula of the convergence rate**. From a practical standpoint,  
9 their optimal fixed sketch algorithm critically relies on a momentum term, and this can be sensitive to noise and  
10 rounding. Without momentum, their fixed sketch algorithm has worse convergence rate than ours.*

11 "Did you experiment applying the truncated Walsh-Hadamard transform when using SRHT ?", "No numerical com-  
12 parisons against state-of-the-art randomized methods are given.", "Could you [...] empirically investigate how [many  
13 refreshing steps] could be skipped?" *We'll include these additional experiments in the final version, and comparisons to  
14 state-of-the-art methods. In particular, our algorithm is faster than the pre-conditioned conjugate gradient method [24].  
15 Further, it's more robust to noise in the gradients compared to the aforementioned fixed sketch algorithm. Although we  
16 do not have theoretical guarantees for skipping refreshing the embedding, we observe in practice that refreshing at  
17 each iteration can be omitted at the cost of convergence rate. Detailed numerical results will be included.*

18 "Practical performance improvement by using orthogonal transforms is slight", "The benefit of orthogonal occurs as  
19  $\xi \rightarrow 1$  [but we] are interested with  $\xi$  as small as possible, close to  $\gamma$ .", "[The algorithm] might be slower [when]  $\varepsilon$  cannot  
20 really be treated as a constant." *We emphasize one of our key findings, that is, the SRHT has the remarkable benefit of a  
21 fast projection method compared to Gaussian embeddings, along with always improving the convergence rate. For  
22  $\gamma = 0.1$ ,  $\xi = 0.5$ , the limiting convergence ratio between  $\rho_h$  and  $\rho_g$  is about 60%. So orthogonal still has benefit. Even  
23 with  $\xi$  being close to  $\gamma$ , the ratio is still strictly less than one. Indeed, we improve theoretical time complexity when  $\varepsilon$  is  
24 treated as a constant, which is a reasonable setting, and we will discuss more general cases in the final version.*

25 "Could you please precisely give the references claiming that  $m \approx d \log(d)$  is prescribed for state-of-the-art algorithms  
26 and for which algorithm?" *Up to constant factors, the authors of [24] originally prescribed  $m \geq d^2$  (see Lemma 1).  
27 Improved concentration bounds on the SRHT [27] can be used to improve this lower bound to  $d \log d$ . See also [6],  
28 Thm 3.1, where the bound  $m \geq C \log d [\sqrt{d} + \sqrt{\log n}]^2$  is stated.*

29 "'For SRHT, we use the optimal step sizes'. I thought it was not proven to be optimal for SRHT. Did I miss something?"  
30 *We will make clearer that this step size is optimal conditional on  $\beta = 0$ .*

31 "A comment on how to deal in the case where  $n$  is not a power of 2." *One can use padding with zeros, which increases  
32 the value of  $n$ . This slightly increases the convergence rate. Or one can take a random subset of coordinates of SRHT,  
33 which empirically does not increase convergence rate, but is somewhat slower to compute.*

34 "How was equation (2) with inverse of sketched Hessian solved?" *In practice, the fastest method is to solve approximately  
35 the linear system with an iterative solver such as conjugate gradient.*

36 "Can this analysis be also extended to accommodate count-sketch or any other sparse sketching matrices?" *We rely on  
37 recent results in random matrix theory (RMT) which, to our knowledge, have not been derived for sparse embeddings.  
38 But there is recent work that analyzes both SRHT/Haar and some sparse embeddings in the same asymptotic framework,  
39 for PCA: [arxiv.org/abs/2005.00511](https://arxiv.org/abs/2005.00511). It may work here but possibly with stronger assumptions on the data.*

40 "The SRHT in this paper incorporates an additional permutation.", "The optimality of IHS is only proved for Haar  
41 transform." *We do consider this additional permutation as we leverage recent results from the RMT. We cannot think of  
42 drawbacks of the extra permutation (computational or otherwise). For proving optimality for the SRHT, we would need  
43 results currently unknown in the RMT. Please see Remark A.1 for more details.*

44 "'SRHT [...] contains less randomness, but is more structured and faster to generate' than Haar matrix." *SRHT relies on  
45 a random permutation and  $n$  sign-flips, while constructing Haar matrices usually needs order  $n^2$  random Gaussian  
46 variables. Thanks to the matrix decomposition of the SRHT, multiplications are faster to perform via FFT.*

47 "Analysis is asymptotic and holds in the infinite limit only.", "... not very standard for sketching results", "Paper assumes  
48 that the ratio  $d/n = \gamma$  is fixed.", "The language of the paper oversells things [...] without qualifying that everything is  
49 asymptotic." *The asymptotic framework is a good fit as one will only use sketching when the dataset is large. It also  
50 leads to clean theoretical results. Finite-sample results sometimes hide large constants. Our results are down to the  
51 constant and thus can be used easily by practitioners. The asymptotic result is also powerful enough to illustrate that  
52 Hadamard projection is superior to Gaussian projection. Moreover, the asymptotic results agree with simulations well  
53 with a few thousands samples. We will make our claims more precise and use asymptotically optimal wherever needed.*