

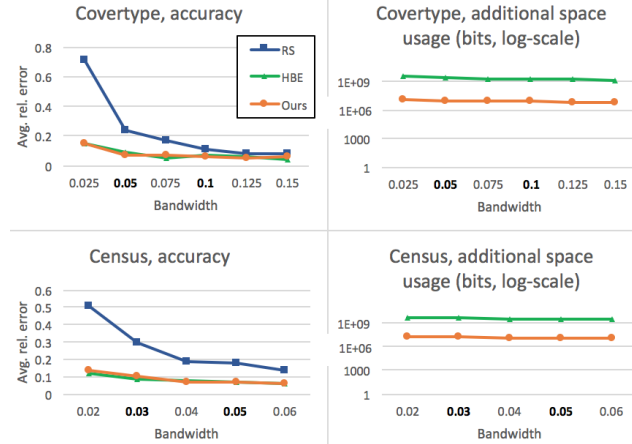
1 We thank the reviewers for their useful feedback and suggestions.

2 **Significance of space improvement:** The need to reduce the space requirements and the preprocessing time of the  
3 [CS17] algorithm was pointed out previously in the literature. For example [SRB+] states that “...straightforward  
4 implementations of the proposed approach require a large amount [...] of preprocessing time and memory to create the  
5 requisite number of hash tables” (Limitations section, pg. 2). This motivated [SRB+] to develop core-set-like methods  
6 for reducing the data size (cf. lines 93–99 in our submission).

7 Regarding theoretical analysis, we emphasize that our result is a direct theoretical improvement to the framework of  
8 [CS17]. While the algorithmic modification is simple (which has the advantage of being practical), the analysis is  
9 rather subtle (cf. Theorem 7). Moreover, it is nearly optimal: As mentioned on page 3, any algorithm must use at least  
10  $N = \Omega(1/(\tau \cdot \epsilon^2))$  data samples to guarantee  $(1 + \epsilon)$ -approximation (see e.g., [PT18]). Thus, the  $1/\sqrt{\tau}$  improvement  
11 that we provide is essentially the best possible, as our storage is proportional to the size  $N$  of the sample data set. We  
12 remark, however, that further improvements could be still possible by storing  $O(N)$  bits instead of memory words.

13 Regarding the quantitative improvement, we note that the lower bound  $\tau$  can be very close to zero in practice, which  
14 makes the  $1/\sqrt{\tau}$  improvement very substantial. Concretely, if in Algorithm 1 we set the budget of kernel computations  
15 per KDE query to  $L$  (recall that analytically  $L \sim 1/\sqrt{\tau}$ ), then our method stores  $L$  hashes per point (in expectation)  
16 while HBE stores  $L^2$ . Since  $L$  is typically of the order of at least dozens or hundreds, the gap in performance can be  
17 large. For example, the plots attached to this file are for  $L = 250$ , resulting in space reduction by a factor of 250.

18 **Additional experiments:** We thank you for the useful  
19 suggestion to study the performance as the bandwidth is  
20 varied. We have run those experiments and representa-  
21 tive results are included here (as space permits). We will  
22 include the full results in the revised version. The plots  
23 shown here are for 250 kernel evaluations per query with  
24 varying bandwidths. The x-axis values in boldface are the  
25 bandwidth values used in our submission. As predicted,  
26 (a) HBE and our method achieve similar error, (b) our  
27 space usage is significantly better (by a factor of 250, as  
28 explained above), (c) as the bandwidth grows and the  
29 KDE values become bounded away from 0 the estimation  
30 task becomes easier, and in particular naïve random sam-  
31 pling improves and converges to the same error as HBE.  
32 We will also include results for other LSHable kernels.



33 **Units of space:** The plots measure space in single bits (however, the labels on the x-axis are accidentally scaled by 100;  
34 we have fixed this). We accept the reviewer’s remark that there are various efficient ways to store each hash, so we redid  
35 the experiments by measuring the *number* of stored hashed points. The resulting plots are essentially scaled versions of  
36 the plots in the paper. In particular, the relative improvement (by a factor of  $L$ ) over HBE is the same. For the sake of  
37 compatibility with our original submission, the plots in this rebuttal still measure space in bits.

38 **Clarifications:** All of the following points raised by the reviewers will be revised or clarified in the final version.  
39 Apologies for any confusion. **R1:** “scheme on page 5” Correct. More accurately for any  $L$ , the expected number of  
40 non-empty bins  $b_j(y)$  is  $\delta L$ , and the estimate in line 152 has expectation  $\tau$ , without having to know it. **R1:** “loss in  
41 variance” The variance affects the error of the estimate, so its impact is indirectly reported on Figs 1–8. We believe the  
42 *relative* error reported on those figures is more natural in our context than the additive squared error (i.e., the empirical  
43 variance), because it directly corresponds to the parameter  $\epsilon$ . **R2:** “stated bound seem to assume  $d = O(1)$ ” Indeed,  
44 the introduction alternates between the number of kernel evaluations and the running time, not always consistently.  
45 Generally, the running time (either of the preprocessing or the kernel evaluation) should have a factor of  $d$  in front of it.  
46 **R2:** “The new bound is only better when  $d < \tau^{1/2}$ ” Indeed, the second column of Table 1 accurately describes the  
47 *preprocessing time* of both algorithms, as opposed to their extra space usage (which can be reduced for both algorithms  
48 by storing pointers not points). Theorem 1 summarizes the bounds correctly, parameterized by the relevant quantities.  
49 **R2:** “additive  $\epsilon\tau$  error guarantee?” This is a subtle point. Our understanding is that without additional assumptions,  
50 such additive error guarantee must require  $1/(\epsilon\tau)^2$  samples by anti-concentration. Whether it can be achieved under  
51 weaker assumptions needs further exploration. **R2:** “strictly better than quadratic?” Correct. **R2:** “query points from  
52 the data set?” Yes. **R3:** “why is the sampling probability chosen as  $\delta = 1/(n\sqrt{\tau})$ ?” Intuitively, so each element appears  
53 in expectation in  $O(1)$  hash tables (cf. lines 71–75). Technically, this equates the variance term of our modification to  
54 that of the original HBE algorithm, ensuring that our variance is at most twice that of HBE (cf. last two lines in proof of  
55 Theorem 7, with  $\beta = 0.5$ ). **R3:** “refer to the proposed one, right?” Correct. **R4:** “what is meant by “typical kernel  
56 value”” Median KDE (which for all the considered settings is within a factor of 2 from the average KDE).