

1 We thank all the reviewers for their insightful and encouraging feedback.

2 **Initial points of acquisition function optimization** Following the implementation details of Spearmint<sup>1</sup>, COMBO  
3 adopts the idea of *spray points*<sup>2</sup> to promote exploitation, which is the same method as in the reference pointed by **R2**.  
4 Due to the discrete nature of COMBO’s search space, the implementation detail is slightly different. In Spearmint,  
5 spray points are points around the best-evaluation point, after perturbing by a zero-mean Gaussian with user-specified  
6 variance (e.g. 0.001<sup>2</sup>). In contrast, in COMBO’s combinatorial graphs, we have spray vertices. Spray vertices are  
7 randomly chosen vertices neighboring the best evaluated vertex, such that the shortest path distance to the best vertex is  
8 less than or equal to a user-specified distance (e.g. 2). As **R2** suggested, this heuristic promotes exploitation.

9 Using random vertices for exploration is similar to Spearmint. In Spearmint we use Sobol sequences, whereas in  
10 COMBO random vertices are chosen uniformly. For initial points for the acquisition function optimization, we probe  
11 acquisition values on 20,020 points consisted of 20 spray vertices and 20,000 random vertices, similar to Spearmint.  
12 The best 20 from 20,020 points are used as initial points for further optimization. In the camera-ready version we will  
13 give the full array of details regarding the exploitation-exploration trade-off, including the references pointed by **R2**.

14 **Computational complexity** In addition to the surrogate model fitting and acquisition function optimization in  
15 ordinary Bayesian optimization methods, COMBO has a one-off pre-processing step of eigendecomposition to compute  
16 Fourier basis. The complexity of the one-off eigendecomposition is *linear* with respect to the number of variables.

17 *Surrogate model fitting.* 1-dimensional slice sampling is used with the Gibbs sampler. The entire sampling procedure has  
18  $O(k^2 d^2) + O(kd^3)$  complexity. Specifically, the 1-dimensional slice sampling requires (a)  $O(kd^2)$  multiplications for  
19  $k$  variables and  $d$  evaluations and (b) solving a linear system of cubic complexity  $O(d^3)$ . The  $O(kd^2)$  multiplications,  
20 specifically, are due to the Kronecker product kernel that the diffusion kernel yields on a Graph Cartesian product.

21 *Acquisition function optimization.* We observe that the heuristic of local optimization with initial points selected from  
22 the pool of spray vertices and random vertices (similar to the description given by **R2**) works better than more complex  
23 solutions with simulated annealing as stated in the paper. We leave the trade-off between optimization and efficiency of  
24 other methods, such as, evolutionary search and genetic algorithm as future work.

25 **Error measure** The error measure in the experiment is the standard error. The evaluations in all experiments except  
26 NAS are noise-free. Thus, small yet still different from 0 standard errors indicate different optima.

27 **Future works and limitations.** COMBO currently handles continuous variables only if they are discretized. Mixing  
28 combinatorial and continuous variables is challenging for Bayesian Optimization, not only for designing an appropriate  
29 surrogate model but also for the optimization of the acquisition function.

30 Besides addressing the aforementioned limitations, bridging the well-established graph theory with combinatorial  
31 optimization opens up lots of interesting research directions (**R1**). Examples are irregular structures not amenable to the  
32 graph Cartesian product or even learning the graph structure. Further, exploring problems, like joint neural architecture  
33 and hyperparameter search, traditionally dominated by low sample efficiency methods (e.g., evolutionary search) is  
34 of great interest: the superior uncertainty quantification of the Bayesian Optimization surrogate model, like Gaussian  
35 Processes, would bring great advantages.

36 As pointed out by reviewers, we hope the graph theory framework on combinatorial optimization introduced by  
37 COMBO will attract attention to the under-explored problem of combinatorial Bayesian optimization. We appreciate the  
38 meticulous comments from all reviewers on the structure, presentation and typos. We will include all the clarifications  
39 in the camera-ready version and release all code, experiments and results upon acceptance.

---

<sup>1</sup>J. Snoek, H. Larochelle, R.P. Adams. Practical Bayesian Optimization of Machine Learning Algorithms. NeurIPS 2012.  
(<https://github.com/JasperSnoek/spearmint>)

<sup>2</sup><https://github.com/JasperSnoek/spearmint/blob/b37a541be1ea035f82c7c82bbd93f5b4320e7d91/spearmint/spearmint/chooser/GPEIOptChooser.py#L235>