

---

## Supplementary material for: BO-HB: Robust and Efficient Hyperparameter Optimization at Scale

---

Stefan Falkner<sup>1</sup> Aaron Klein<sup>1</sup> Frank Hutter<sup>1</sup>

### A. Available Software

To promote reproducible science and enable other researchers to use our method, we provide an open-source implementation of BOHB and Hyperband. It is available under <https://github.com/automl/HpBandSter>. The benchmarks and our scripts used to produce the data shown in the paper can be found in the *icml\_2018* branch.

### B. Comparison to other Combinations of Bayesian optimization and Hyperband

Here we discuss the differences between our method and the related approaches of Bertrand et al. (2017) and Wang et al. (2018) in more detail. We note that these works are independent and concurrent; our work extends our preliminary short workshop papers at NIPS 2017 (Falkner et al., 2017) and ICLR 2018 (Falkner et al., 2018).

While the general idea of combining Hyperband and Bayesian optimization by Bertrand et al. (2017) is the same as in our work, they use a Gaussian process for modeling the performance. The budget is modeled like any other dimension of the search space, without any special treatment. Based on our experience with Fabolas (Klein et al., 2017), we expect that the squared exponential kernel might not extrapolate well, which would hinder performance. Also, the small evaluation provided by Bertrand et al. (2017) does not allow strong conclusions about the performance of their method.

Wang et al. (2018) also independently combined TPE and Hyperband, but in a slightly different way than we did. In their method, TPE is used as a subroutine in every iteration of Hyperband. In particular, a new model is built from scratch at the beginning of every SuccessiveHalving run (Algorithm 3, line 8 in Wang et al. (2018)). This means that in later iterations of the algorithm, the model does not

---

<sup>1</sup>Department of Computer Science, University of Freiburg, Freiburg, Germany. Correspondence to: Stefan Falkner <sfalkner@informatik.uni-freiburg.de>.

benefit from any of the evaluations in previous iterations. In contrast, BOHB collects all evaluations on all budgets and uses the largest budget with enough evaluations (admittedly a heuristic, but we would argue a reasonable one) as a base for future evaluations. This way, BOHB aggregates more knowledge into its models for the different budgets as the optimization progresses. We believe this to be a crucial part of the strong performance of our method. Empirically, Wang et al. (2018) did not achieve the consistent and large speedups across a wide range of applications BOHB achieved in our experiments.

### C. Successive Halving

SuccessiveHalving is a simple heuristic to allocate more resources to promising candidates. For completeness, we provide pseudo code for it in Algorithm 1. It is initialized with a set of configurations, a minimum and maximum budget, and a scaling parameter  $\eta$ . In the first stage all configurations are evaluated on the smallest budget (line 3). The losses are then sorted and only the best  $1/\eta$  configurations are kept in the set  $C$  (line 4). For the following stage, the budget is increased by a factor of  $\eta$  (line 5). This is repeated until the maximum budget for a single configuration is reached (line 2). Within Hyperband, the budgets are chosen such that all SuccessiveHalving executions require a similar total budget.

---

**Algorithm 1:** Pseudocode for SuccessiveHalving used by Hyperband as a subroutine.

---

```
input : initial budget  $b_0$ , maximum budget  $b_{max}$ , set of  $n$  configurations  $C = \{c_1, c_2, \dots, c_n\}$ 
1  $b = b_0$ 
2 while  $b \leq b_{max}$  do
3    $L = \{\tilde{f}(c, b) : c \in C\}$ 
4    $C = \text{top}_k(C, L, \lfloor |C|/\eta \rfloor)$ 
5    $b = \eta \cdot b$ 
```

---

### D. Details on the Kernel Density Estimator

We used the MultivariateKDE from the statsmodels package (Seabold & Perktold, 2010), which constructs a factorized

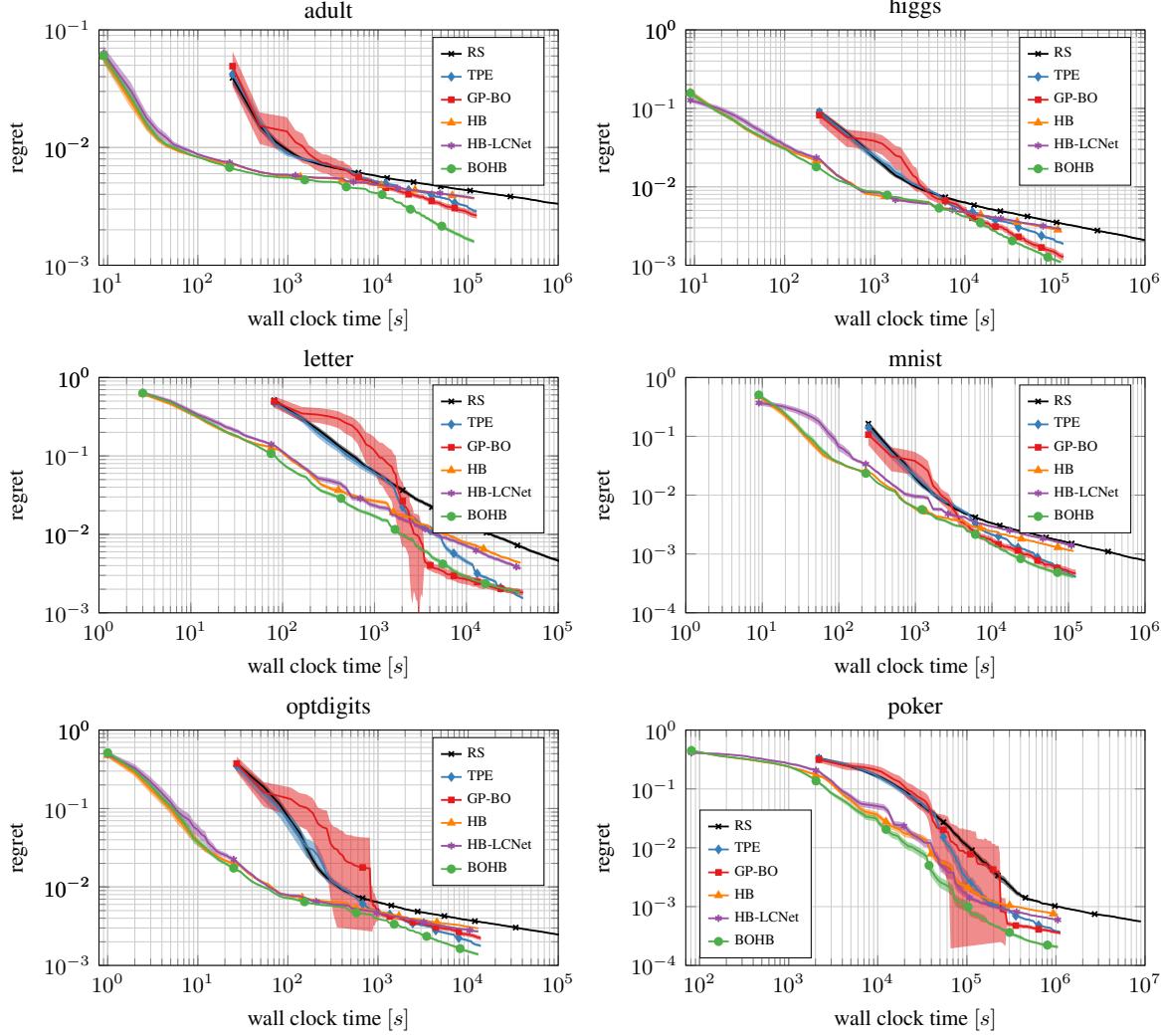


Figure 1. Mean performance on the surrogates for all six datasets. As uncertainties, we show the standard error of the mean based on 512 runs (except for GP-BO, which has only 50 runs).

kernel, with a one-dimensional kernel for each dimension. Note that using this product of 1-d kernels differs from the original TPE, which uses a pdf that is the product of 1-d pdfs. For the continuous parameters a Gaussian kernel is used, whereas the Aitchison-Aitken kernel is the default for categorical parameters. We used Scott's rule for efficient bandwidth estimation, as preliminary experiments with maximum-likelihood based bandwidth selection did not yield better performance but caused a significant overhead.

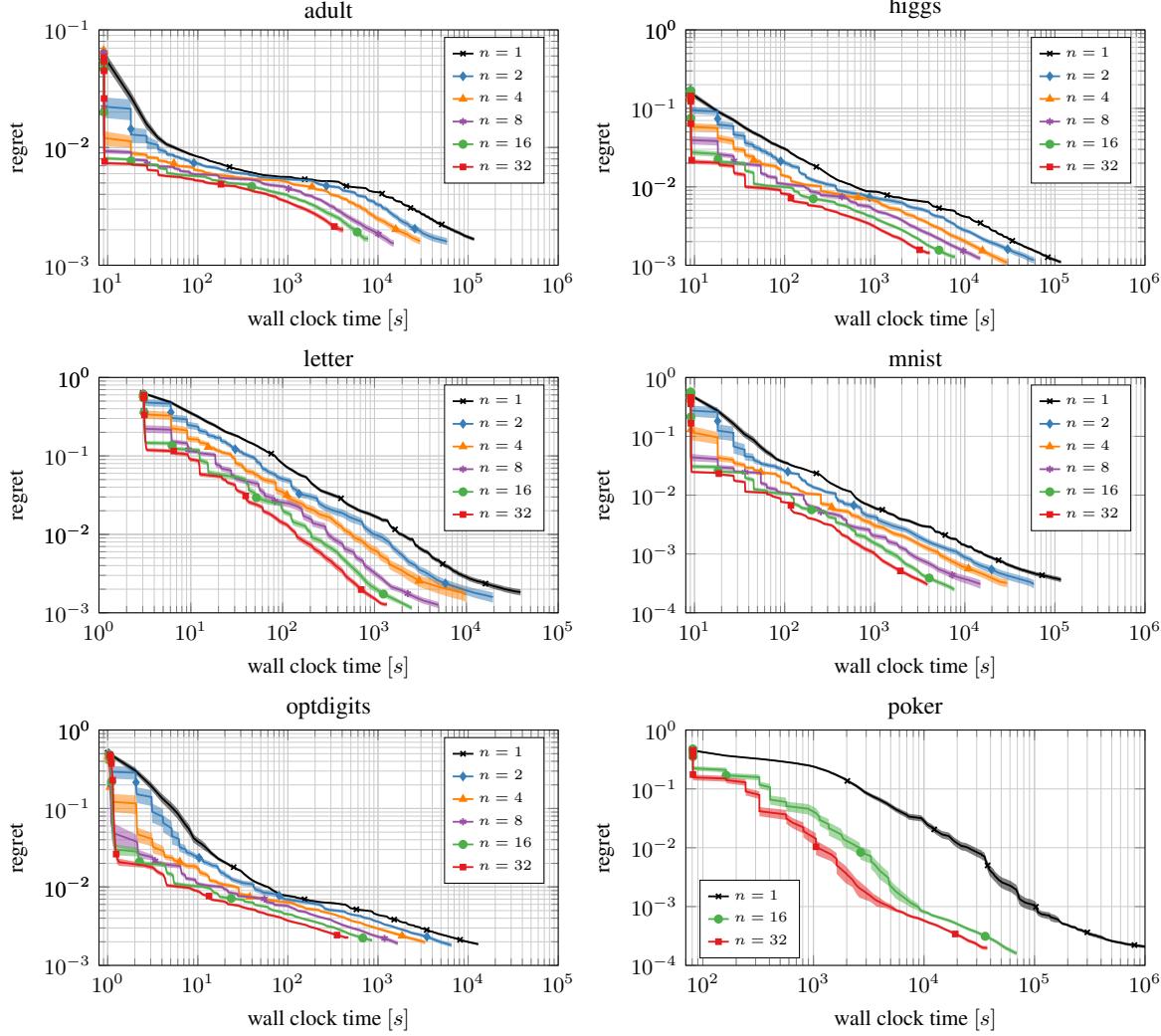
## E. Performance of all methods on all surrogates

Figure 1 shows the performance of all methods we evaluated on all our surrogate benchmarks. Random search is clearly

the worst optimizer across all datasets when the budget is large enough for GP-BO and TPE to leverage their model. Hyperband and the two methods based on it (HB-LCNet) and BOHB improve much more quickly due to the smaller budgets used. On all surrogate benchmarks, BOHB starts to outperform HB after the first couple of iterations (sometimes even earlier, e.g., on dataset letter). The same dataset also shows that traditional BO methods can still have an advantage for very large budgets, since in these late stages of the optimization process the low fidelity evaluations of BOHB can cause a constant overhead without any gain.

## F. Performance of parallel runs

Figure 2 shows the performance of BOHB when run in parallel on all our surrogate benchmarks. The speed-ups



*Figure 2.* Mean performance on the surrogates for all six datasets with different numbers of workers  $n$ . As uncertainties, we show the standard error of the mean based on 128 runs. Because we simulated them in real time to capture the true behavior, poker is too expensive to evaluate with less than 16 workers within a day.

are quite consistent, and almost linear for a small number of workers (2-8). For more workers, more random configurations are evaluated in parallel before the first model is built, which degrades performance. But even for 32 workers, linear speedups are possible (see, e.g., dataset letter, for reaching a regret of  $2 \times 10^{-3}$ ).

We note that in order to carry out this evaluation of parallel performance, we actually simulated the parallel optimization by making each worker wait for the given budget before returning the corresponding performance value of our surrogate benchmark. (The case of one worker is an exception, where we can simply reconstruct the trajectory because all configurations are evaluated serially.) By using this approach in connection with threads, each evaluation of a parallel algorithm still only used 1 CPU, but the run

actually ran in real time. For this reason, we decided to not evaluate all possible numbers of workers for dataset poker, for which each run with less than 16 workers would have taken more than a day, and we do not expect any different behavior compared to the other datasets.

## G. Evaluating the hyperparameters of BOHB

In this section, we evaluate the importance of the individual hyperparameters of BOHB, namely the number of samples used to optimize the acquisition function (Figure 3), the fraction of purely random configuration  $\rho$  (Figure 4), the scaling parameter  $\eta$  (Figure 5), and the bandwidth factor used to encourage exploration (Figure 6).

Additionally, we want to discuss the importance of  $\eta$ ,  $b_{min}$

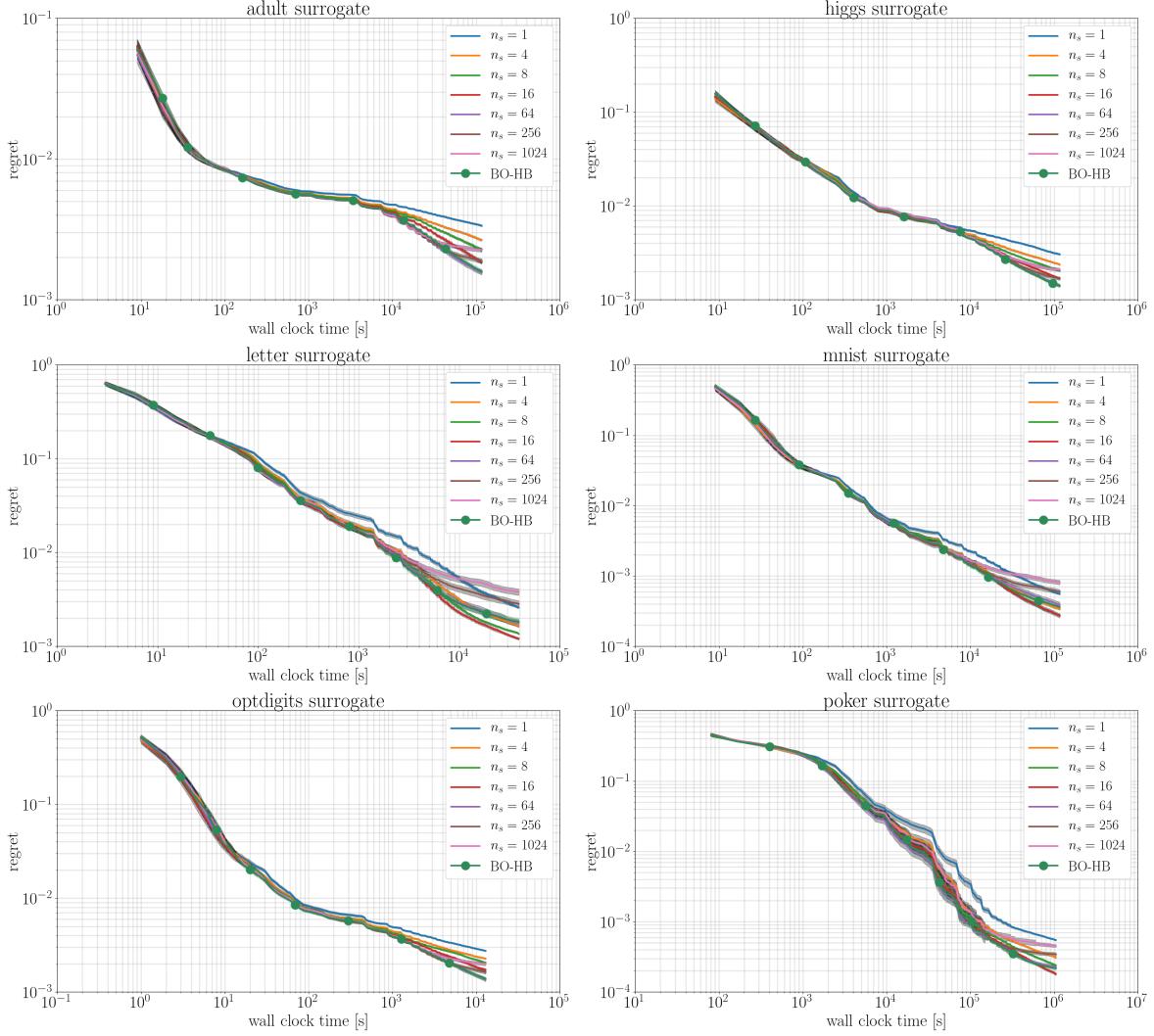


Figure 3. Performance on the surrogates for all six datasets for different number of samples

and  $b_{max}$  already present in HB. The parameter  $\eta$  controls how aggressively SH cuts down the budget and the number of configurations evaluated. Like HB (Li et al., 2017), BOHB is also quite insensitive to this choice in a reasonable range. For our experiments, we use the same default value ( $\eta = 3$ ) for HB and BOHB.

More important for the optimization are  $b_{min}$  and  $b_{max}$ , which are problem specific and inputs to both HB and BOHB. While the maximum budget is often naturally defined, or is constrained by compute resources, the situation for the minimum budget is often different. To get substantial speedups, an evaluation with a budget of  $b_{min}$  should contain some information about the quality of a configuration with larger budgets; for example, when subsampling the data, the smallest subset should not be one datum, but rather enough points to fit a meaningful model. This requires knowledge about the benchmark and the algorithm

being optimized.

## H. Counting Ones

We now show results for the counting ones function for different dimensions. Figure 7 shows the mean performance of all applicable methods in  $d = 8, 16, 32$  and  $64$  dimensions for a budget of 8192 full function evaluations.

We draw the following conclusions from the results:

1. Despite its simple definition, this problem is quite challenging for the methods we applied to it. RS and HB both suffer from the fact that drawing configurations at random performs quite poorly in this space. The model-based approaches SMAC and TPE performed substantially better, especially with large budgets. They required a larger number of samples before converging

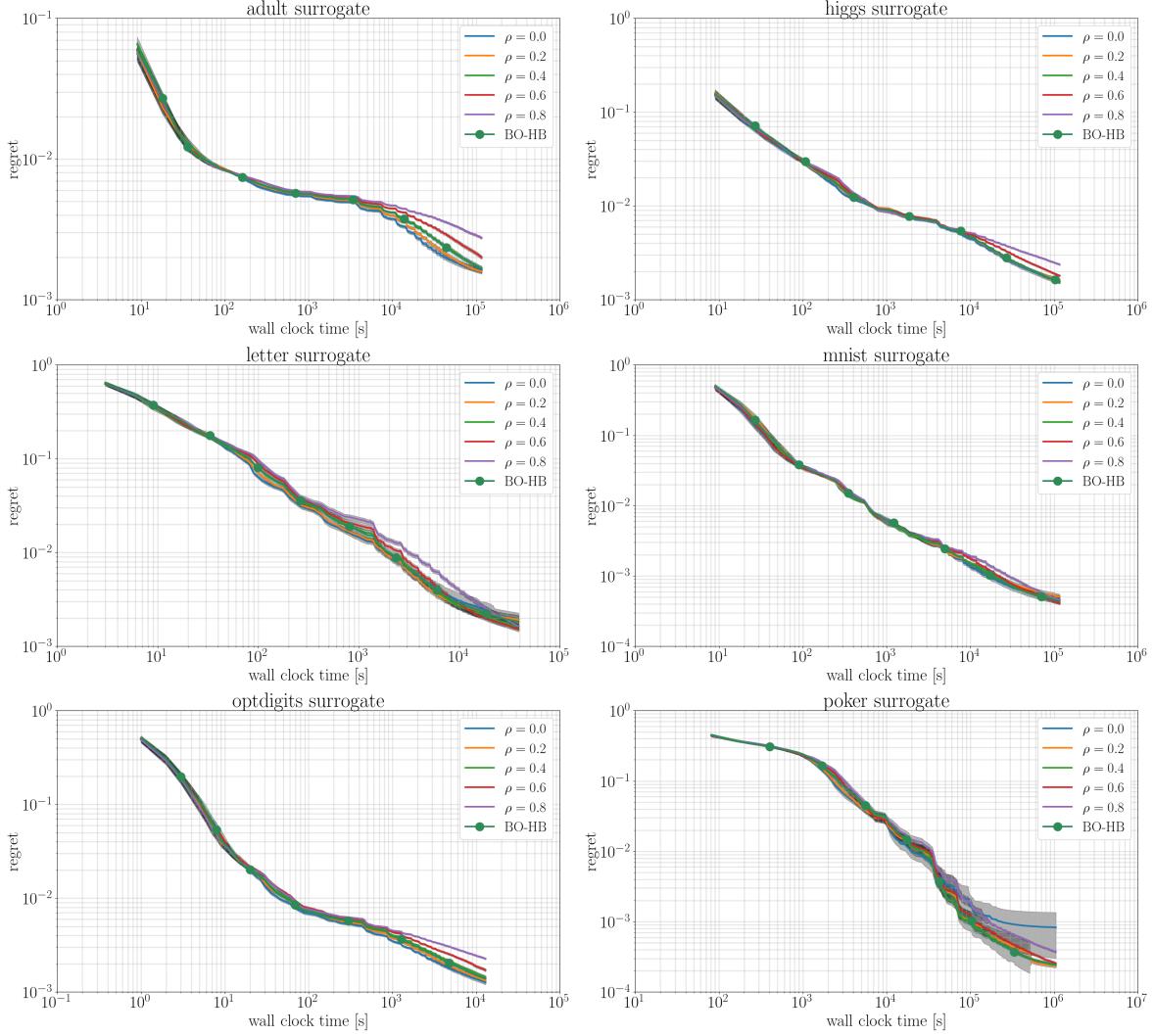


Figure 4. Performance on the surrogates for all six datasets for different random fractions

to the true optimum than BOHB. However, we would like to mention that SMAC and TPE treated the problem as a blackbox optimization problem; the results for SMAC could likely be improved by treating individual samples as “instances” and using SMAC’s intensification mechanism to reject poor configurations based on few samples and evaluate promising configurations with more samples.

- BOHB struggles in the very high dimensional case. We attribute this to the fact that the noise is substantially higher in this case, such that larger budgets are required to build a good model. Therefore, given a large enough budget, BOHB’s evaluations on small budgets lead to a constant overhead over only using the more reliable evaluations on larger budgets. Since the optimization problem is perfectly separable (there are no interaction effects between any dimensions), we also expect

TPE’s univariate KDE to perform better than BOHB’s multivariate one.

## I. Surrogates

### I.1. Constructing the Surrogates

To build a surrogate, we sampled 10 000 random configurations for each dataset, trained them for 50 epochs, and recorded their classification error after each epoch, along with their total training time. We fitted two independent random forests that predict these two quantities as a function of the hyperparameter configuration used. This enabled us to predict the classification error as a function of time with sufficient accuracy. As almost all networks converged within the 50 epochs, we extend the curves by the last obtained value if the budget would allow for more epochs.

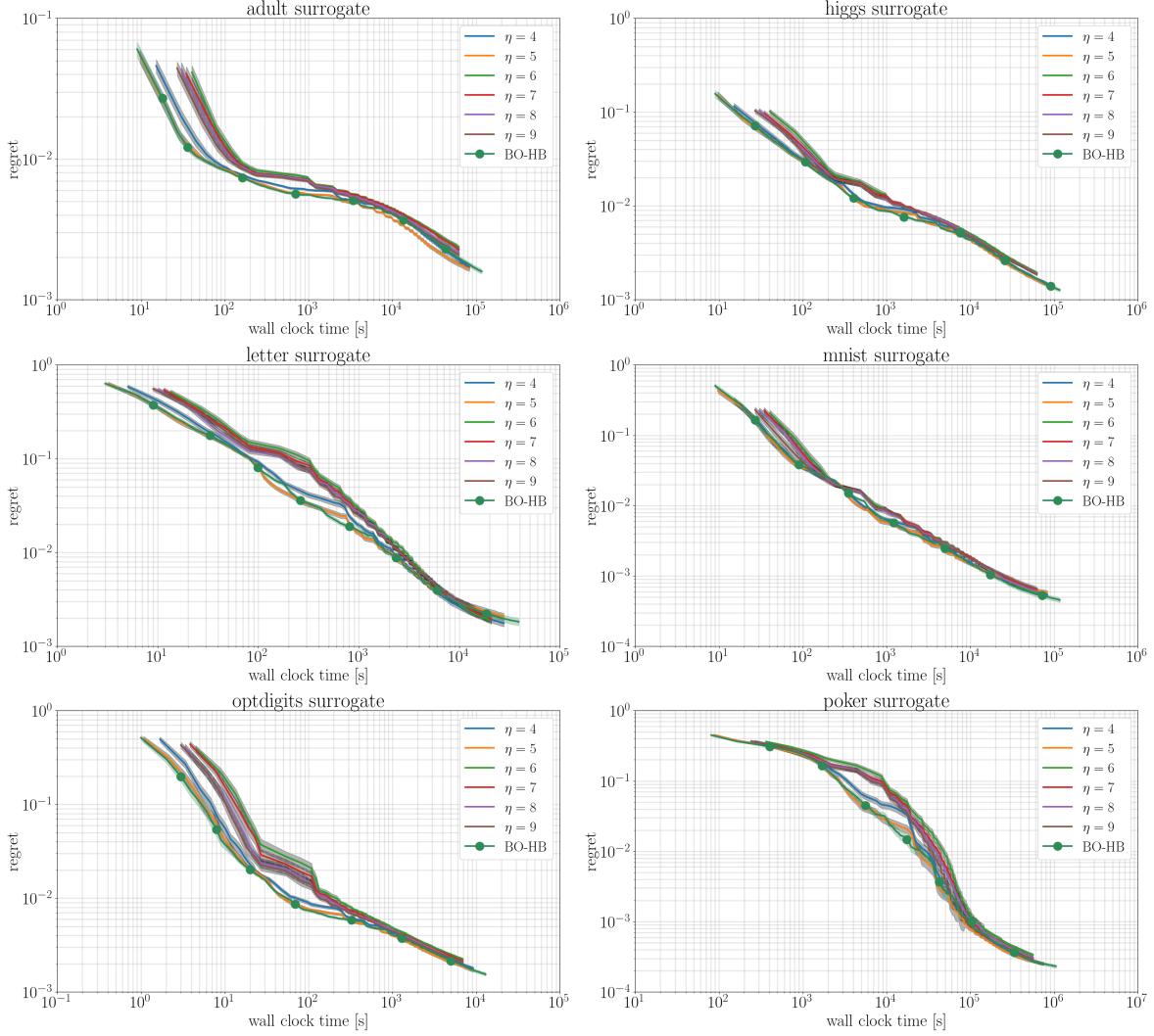


Figure 5. Performance on the surrogates for all six datasets for different values of  $\eta$ .

The surrogates enable cheap benchmarking, allowing us to run each algorithm 256 times. Since evaluating a configuration with the random forest is inexpensive, we used a global optimizer (differential evolution) to find the true optimum. We allowed the optimizer 10 000 iterations which should be sufficient to find the true optimum.

Besides these positive aspects of benchmarking with surrogates, there are also some drawbacks that we want to mention explicitly:

- (a) There is no guarantee that the surrogate actually reflects the important properties of the true benchmark.
- (b) The presented results show the optimized classification error on the validation set used during training. There is no test performance that could indicate overfitting.
- (c) Training with stochastic gradient descent is an inher-

ently noisy process, i.e. two evaluations of the same configuration can result in different performances. This is not at all reflected by our surrogates, making them a potentially easier to optimize than the true benchmark they are based on.

- (d) By fixing the budgets (see below) and having deterministic surrogates, the global minima might be the result of some small fluctuations in the classification error in the surrogates' training data. That means that the surrogate's minimizer might not be the true minimizer of the real benchmark.

None of these downsides necessarily have substantial implications for comparing different optimizers; they simply show that the surrogate benchmarks are not perfect models for the real benchmark they mimic. Nevertheless, we believe that, especially for development of novel algorithms,

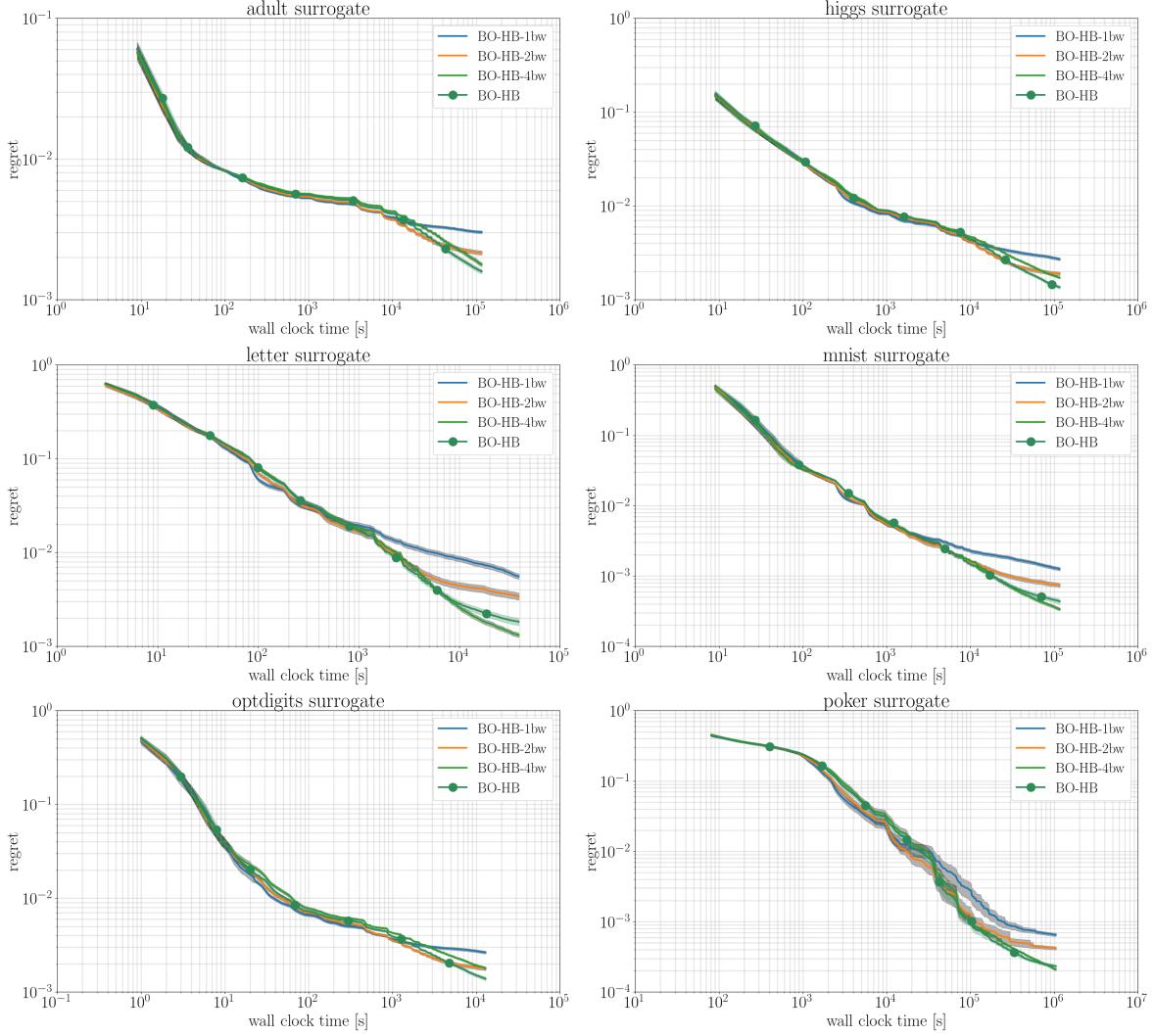


Figure 6. Performance on the surrogates for all six datasets for different bandwidth factors.

Table 1. The hyperparameters and architecture choices for the fully connected networks.

Hyperparameter	Range	Log-transform
batch size	$[2^3, 2^8]$	yes
dropout rate	$[0, 0.5]$	no
initial learning rate	$[10^{-6}, 10^{-2}]$	yes
exponential decay factor	$[-0.185, 0]$	no
# hidden layers	$\{1, 2, 3, 4, 5\}$	no
# units per layer	$[2^4, 2^8]$	yes

the positive aspects outweigh the negative ones.

## I.2. Determining the budgets

To choose the largest budget for training, we looked at the best configuration as predicted by the surrogate and its

Table 2. The budgets used by HB and BOHB; random search and TPE only used the last budget

Dataset	Budgets in seconds for HB and BOHB
Adult	9, 27, 81, 243
Higgs	9, 27, 81, 243
Letter	3, 9, 27, 81
Poker	81, 243, 729, 2187

training time. We chose the closest power of 3 (because we also used  $\eta = 3$  for HB and BOHB) to achieve that performance. We chose the smallest budget for HB such that most configurations had finished at least one epoch. Table 2 lists the budgets used for all datasets.

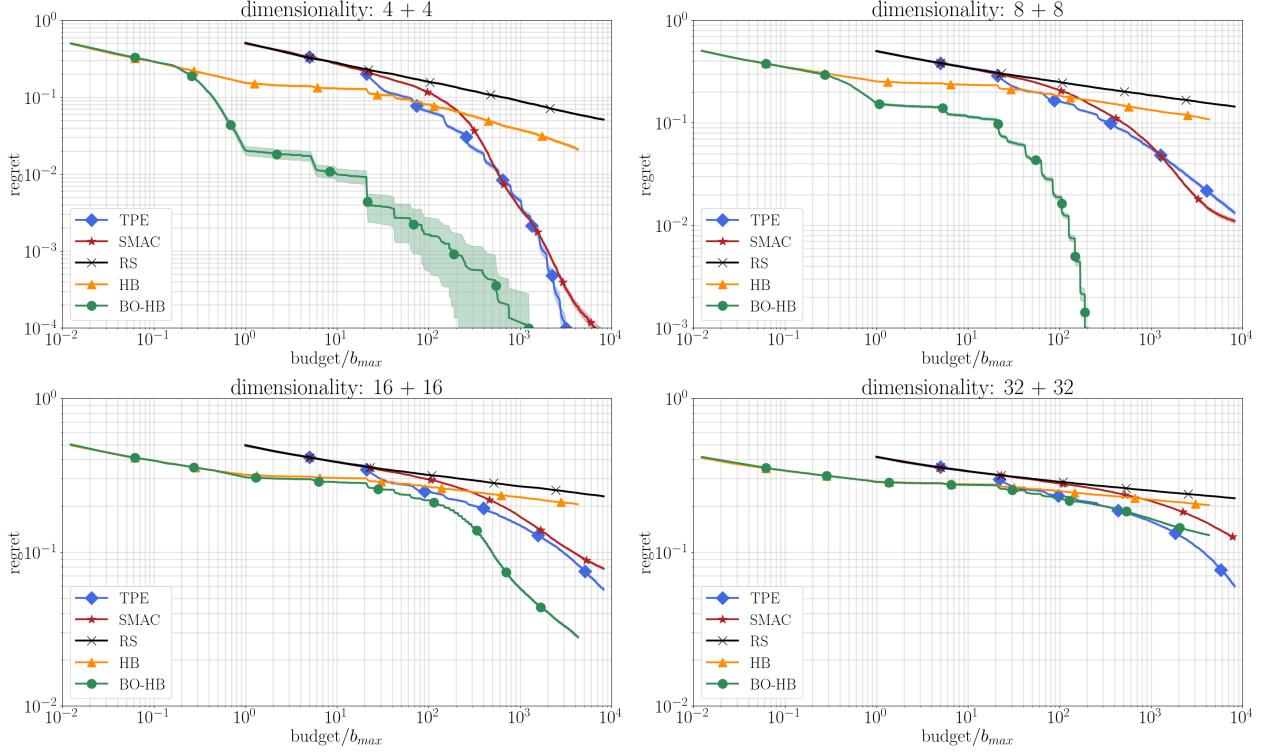


Figure 7. Mean performance of BOHB, HB, TPE, SMAC and RS on the mixed domain counting ones function with different dimensions. As uncertainties, we show the standard error of the mean based on 512 runs.

Table 3. The hyperparameters for the Bayesian neural network task.

Hyperparameter	Range	Log-transform
# units layer 1	[ $2^4, 2^9$ ]	yes
# units layer 2	[ $2^4, 2^9$ ]	yes
step length	[ $10^{-6}, 10^{-1}$ ]	yes
burn in	[0, .8]	no
momentum decay	[0, 1]	no

Table 4. The hyperparameters for the PPO Cartpole task.

Hyperparameter	Range	Log-transform
# units layer 1	[ $2^3, 2^7$ ]	yes
# units layer 2	[ $2^3, 2^7$ ]	yes
batch size	[ $2^3, 2^8$ ]	yes
learning rate	[ $10^{-7}, 10^{-1}$ ]	yes
discount	[0, 1]	no
likelihood ratio clipping	[0, 1]	no
entropy regularization	[0, 1]	no

## J. Bayesian Neural Networks

We optimized the hyperparameters described in Table 3 for a Bayesian neural network trained with SGHMC on two UCI regression datasets: Boston Housing and Protein Structure. The budget for this benchmark was the number of steps for the MCMC sampler. We set the minimum budget to 500 steps and the maximum budget to 10000 steps. After sampling 100 parameter vectors, we computed the log-likelihood on the validation dataset by averaging the predictive mean and variances of the individual models. The performance of all methods for both datasets is shown in Figure 8.

## K. Reinforcement Learning

Table 4 shows the hyperparameters we optimized for the PPO Cartpole task.

## References

Bertrand, H., Ardon, R., Perrot, M., and Bloch, I. Hyperparameter optimization of deep neural networks: Combining hyperband with Bayesian model selection. *Proceedings of Conférence sur l'Apprentissage Automatique (CAP 2017)*, 2017.

Falkner, S., Klein, A., and Hutter, F. Combining hyperband and Bayesian optimization. In *NIPS 2017 Bayesian*

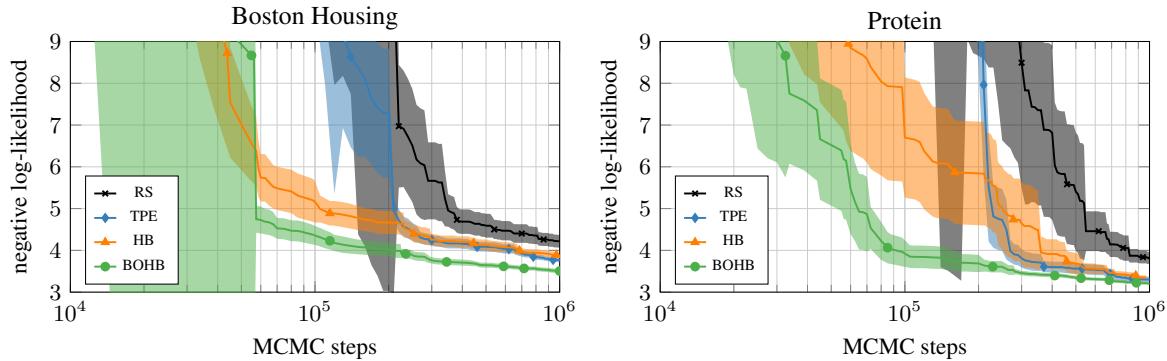


Figure 8. Mean performance of TPE, RS, HB and BOHB for optimizing the 5 hyperparameters of a Bayesian neural network on two different UCI datasets. As uncertainties, we show the standard error of the mean based on 50 runs.

*Optimization Workshop, December 2017.*

Falkner, S., Klein, A., and Hutter, F. Practical hyperparameter optimization for deep learning. In *ICLR 2018 Workshop Track*, 2018.

Klein, A., Falkner, S., Bartels, S., Hennig, P., and Hutter, F. Fast Bayesian optimization of machine learning hyperparameters on large datasets. In *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2017.

Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A., and Talwalkar, A. Hyperband: Bandit-based configuration evaluation for hyperparameter optimization. In *Proceedings of the International Conference on Learning Representations (ICLR'17)*, 2017. Published online: [iclr.cc](http://iclr.cc).

Seabold, S. and Perktold, J. Statsmodels: Econometric and statistical modeling with python. In *9th Python in Science Conference*, 2010.

Wang, J., Xu, J., and Wang, X. Combination of hyperband and bayesian optimization for hyperparameter optimization in deep learning. *arXiv preprint arxiv:1801.01596*, 01 2018.