# MOFA: Modular Factorial Design for Hyperparameter Optimization

**Bo Xiong** [1,*], **Yimin Huang**[2], **Hanrong Ye**[2], **Steffen Staab**[1,3], and **Zhenguo Li**[2]

[1]University of Stuttgart, Germany
[2]Huawei Noah's Ark Lab, China
[3]University of Southampton, United Kingdom

## Abstract

This paper presents a novel and lightweight hyperparameter optimization (HPO) method, MOdular FActorial Design (MOFA). MOFA pursues several rounds of HPO, where each round alternates between *exploration* of hyperparameter space by *factorial design* and *exploitation* of evaluation results by *factorial analysis*. Each round first *explores* the configuration space by constructing a low-discrepancy set of hyperparameters that cover this space well while de-correlating hyperparameters, and then *exploits* evaluation results through factorial analysis that determines which hyperparameters should be further explored and which should become fixed in the next round. We prove that the inference of MOFA achieves higher confidence than other sampling schemes. Each individual round is highly parallelizable and hence offers major improvements of efficiency compared to model-based methods. Empirical results show that MOFA achieves better effectiveness and efficiency compared with state-of-the-art methods.

## 1 Introduction

Modern machine learning techniques have achieved promising results in various areas [1, 2, 3, 4]. However, the performance of these models highly depends on the configurations of their hyperparameters [5]. For example, the performance of deep neural networks may fluctuate dramatically under different neural architectures [6, 7]. Different data augmentation policies can lead to different experimental results for an image recognition task [8]. Heuristic tuning by using expert's experiences is a possible solution but only works on simple settings.

To avoid dependence on experts' experiences, automated hyperparameter optimization (HPO) [9, 10] searches for optimal hyperparameters. Currently, there are two lines of works on HPO. (1) *Model-based methods* such as Bayesian Optimization (BO) [11] optimize hyperparameters by learning a surrogate model (e.g. Gaussian process). While there are some more sophisticated approaches such as SMAC [12], BOHB [13] and Reinforcement Learning (RL) [14], all these methods highly depend on the model parametrization and cannot be fully parallelized as their core strategies iteratively pursue improvements. (2) *Model-free methods* (e.g. Random Search) do not depend on any parametric model and can run in fully parallel fashion, as the different hyperparameter configurations are run and evaluated individually [15]. Nevertheless, the sample efficiency of Random Search is not on par with model-based methods to achieve a globally optimal result. Besides, some works studied how to automatically change the search space during optimization [16] and explored identifying hyperparameter importance [17, 18], but these works still rely on the model parametrization.

Factorial designs [19], which deliberate how to adjust parameters and exploit corresponding experimental responses, can be used to improve the sample efficiency while avoiding model parametrization.

---

[*]Work done as an internship at Huawei Noah's Ark Lab

Some HPO methods have begun to make use of factorial designs, such as Latin Hypercubes [20, 21] and Orthogonal Arrays (OAs) [22]. So far, however, factorial designs for HPO have only been applied to explore the hyperparameter space without exploiting possible feedback from the evaluation results returned by running various hyperparameter configurations.

This paper present MOdular FActorial Design (MOFA), a multi-module process that combines factorial design with factorial analysis, which achieves high quality of optimization results while allowing for high efficiency through excellent parallelizability. The main idea of MOFA is to improve *exploration* of hyperparameter sampling with *factorial designs* and improve *exploitation* of evaluation results with *factorial analysis* by identifying the hyperparameter importance and reducing the search space. In each round, MOFA first *explores* the configuration space by constructing a low-discrepancy set of hyperparameters that cover this space well while de-correlating hyperparameters and then *exploits* evaluation results through factorial analysis that determines which hyperparameters should be further explored and which should become fixed in the next round. Specifically, MOFA runs through four modules in each round: Firstly, an Orthogonal Latin Hypercube (OLH)-based sampler ensuring both univariate projection uniformity (low-discrepancy) and orthogonality, which explores the hyperparameter space more efficiently without correlating the hyperparemeters (improving *exploration*). Secondly, a highly parallelized evaluator. Thirdly, a transformer collapsing the OLH performance table into an Orthogonal Array (OA). Finally, factorial analysis narrows down the search space and selects hyperparameters that are most promising for iterative optimization.

To summarize, the main contributions of our paper are:

- We propose MOFA, a novel and lightweight HPO method that shares the advantages of being model-free, parallelizable and sample efficient.
- To our best knowledge, we are the first to exploit factorial design with factorial analysis on the setting of HPO.
- We provide necessary theoretical analysis to show how MOFA controls the worst-case response and its inference reliability.
- Empirical results show that MOFA clearly improves effectiveness and efficiency of HPO compared to state-of-the-art methods.

## 2 Preliminaries

**Latin Hypercubes.** A Latin Hypercube [23] is an $N \times d$ table for an $N$-run experiment with $d$ factors, which is based on the Latin Square in which there is only one point in each row and column of a gridded space. A Latin Hypercube is the generalization of the Latin Square to an arbitrary number of dimensions, whereby each sample is the only one in each axis-aligned hyper-plane containing it. Fig. 1c (left) shows an example of a Latin Hypercube with three factors. A Latin Hypercube has the property of univariate or one-dimensional projection uniformity which means that by projecting an $l$-point design on to any factor we will get $l$ different levels for that factor [24]. This property drastically reduces the number of configurations necessary to achieve a reasonably accurate result.

**Orthogonal Arrays.** An OA is an $N \times d$ table, and each factor has $l$ levels. For any $t$ columns, the different level configurations appear with equal frequency ($t$-dimensional orthogonality). The number $t$ is called the strength. A Latin Hypercube is a special OA of strength one. For example, Fig. 1c (right) is an $OA(9, 3, 3, 2)$, where the number of runs is nine, the number of factors is three, the number of levels is three and the strength is two. In any selected two columns, all possible configurations $(1,1),(1,2),(1,3),(2,1),(2,2),(2,3),(3,1),(3,2),(3,3)$ appear (*complete*) and appear the same number of times (*balanced*). This number is called the *index* ($\lambda$). In OA design, the size of OA should meet the restriction $N = \lambda l^t$, where $l$ should be a prime number. The orthogonality in OA ensures that each factor's main effect can be determined unaffected by interaction with other factors.

**Orthogonal Latin Hypercubes.** A randomly generated Latin Hypercube may be quite structured: the design may not have good univariate projection uniformity or the different factors might be highly correlated [25]. Several criteria such as maximin distance [26] and minimum correlation [27] have been proposed to address these issues. Since OA ensures that the factor's main effect should not be influenced by its interaction effect with others, we use OA to construct Orthogonal Latin Hypercubes (OLH) [24]). Beyond non-orthogonal Latin Hypercubes, OLH ensures both univariate projection uniformity and $t$-dimensional orthogonality, making it more suitable for factorial analysis.
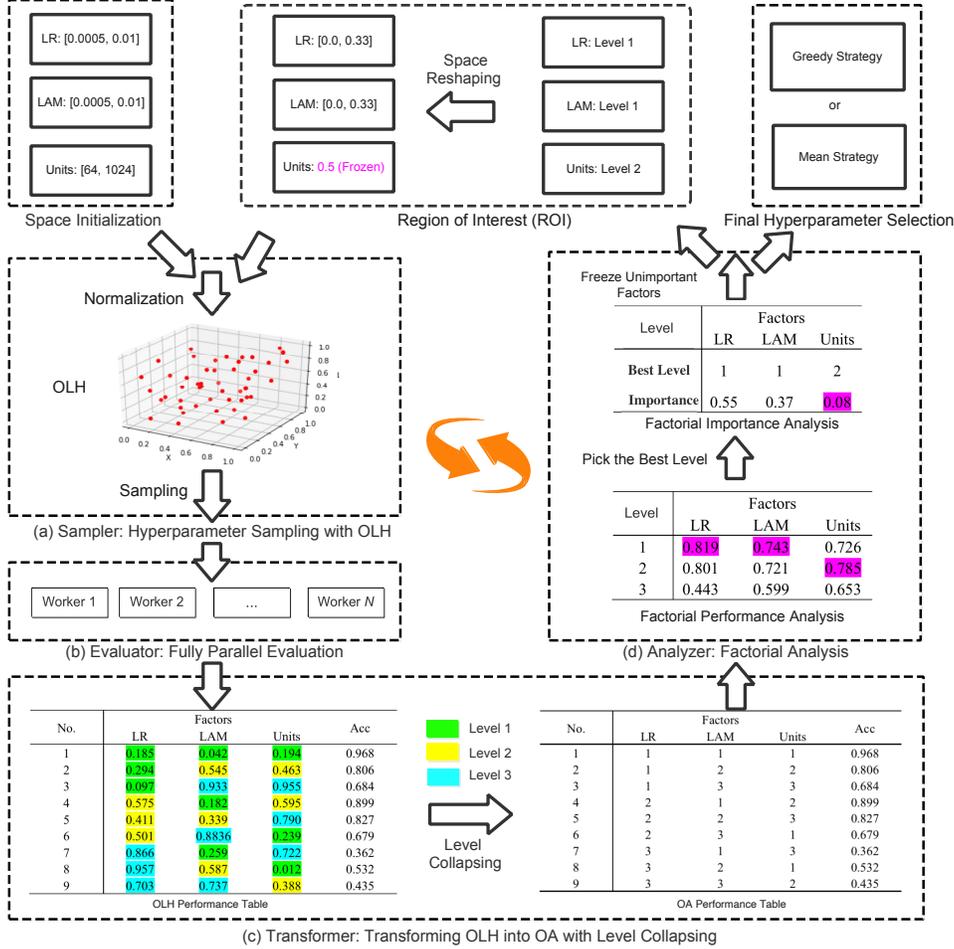
Figure 1: The overview of MOFA. MOFA consists of four modules. (a) an OLH-based sampler for hyperparameter sampling; (b) a paralleled evaluator; (c) a transformer collapsing OLH performance table into an OA. (d) an analyzer narrowing down the search space.

**Space Initialization:** LR: [0.0005, 0.01], LAM: [0.0005, 0.01], Units: [64, 1024]

**Space Reshaping**

**Region of Interest (ROI):** LR: [0.0, 0.33], LAM: [0.0, 0.33], Units: 0.5 (Frozen) — LR: Level 1, LAM: Level 1, Units: Level 2

**Final Hyperparameter Selection:** Greedy Strategy or Mean Strategy

**(a) Sampler: Hyperparameter Sampling with OLH** — Normalization, OLH, Sampling

**(b) Evaluator: Fully Parallel Evaluation** — Worker 1, Worker 2, ..., Worker N

**(d) Analyzer: Factorial Analysis**

Freeze Unimportant Factors

| Level | Factors | | |
|---|---|---|---|
| | LR | LAM | Units |
| Best Level | 1 | 1 | 2 |
| Importance | 0.55 | 0.37 | 0.08 |

Factorial Importance Analysis

Pick the Best Level

| Level | Factors | | |
|---|---|---|---|
| | LR | LAM | Units |
| 1 | 0.819 | 0.743 | 0.726 |
| 2 | 0.801 | 0.721 | 0.785 |
| 3 | 0.443 | 0.599 | 0.653 |

Factorial Performance Analysis

**(c) Transformer: Transforming OLH into OA with Level Collapsing**

| No. | Factors | | | Acc |
|---|---|---|---|---|
| | LR | LAM | Units | |
| 1 | 0.185 | 0.042 | 0.194 | 0.968 |
| 2 | 0.294 | 0.545 | 0.463 | 0.806 |
| 3 | 0.097 | 0.933 | 0.955 | 0.684 |
| 4 | 0.575 | 0.182 | 0.595 | 0.899 |
| 5 | 0.411 | 0.339 | 0.790 | 0.827 |
| 6 | 0.501 | 0.8836 | 0.239 | 0.679 |
| 7 | 0.866 | 0.259 | 0.722 | 0.362 |
| 8 | 0.957 | 0.587 | 0.012 | 0.532 |
| 9 | 0.703 | 0.737 | 0.388 | 0.435 |

OLH Performance Table

Level 1, Level 2, Level 3

Level Collapsing

| No. | Factors | | | Acc |
|---|---|---|---|---|
| | LR | LAM | Units | |
| 1 | 1 | 1 | 1 | 0.968 |
| 2 | 1 | 2 | 2 | 0.806 |
| 3 | 1 | 3 | 3 | 0.684 |
| 4 | 2 | 1 | 2 | 0.899 |
| 5 | 2 | 2 | 3 | 0.827 |
| 6 | 2 | 3 | 1 | 0.679 |
| 7 | 3 | 1 | 3 | 0.362 |
| 8 | 3 | 2 | 1 | 0.532 |
| 9 | 3 | 3 | 2 | 0.435 |

OA Performance Table

## 3    Modular Factorial Design for Hyperparameter Optimization

Fig. 1 shows an overview of MOFA. MOFA starts with a space initialization. Then, MOFA runs through four modules in each iteration. (1) *Sampler* (Fig. 1a): we construct an OLH ensuring both univariate projection uniformity and orthogonality to sample hyperparameter configurations. (2) *Evaluator* (Fig. 1b): the sampled hyperparameter configurations are evaluated in parallel. (3) *Transformer* (Fig. 1c): an OLH performance table is built based on the evaluated results and collapsed into an OA performance table. (4) *Analyzer* (Fig. 1d): factorial performance analysis and factorial importance analysis are conducted to narrow down the search space and select influential hyperparameters.

### 3.1    Sampler: Hyperparameter Sampling with OLH

In the sampler (Fig. 1a), we first normalize the search space for each (discrete or continuous) hyperparameter into $[0, 1]$ so that each hyperparameter is located in the same sampling space. Secondly, we build an OLH that ensures both one-dimensional projection uniformity and $t$-dimensional orthogonality to sample hyperparameter configurations. The one-dimensional projection uniformity makes it explore the search space more efficiently, while orthogonality reduces the interaction between hyperparameters to evaluate their main effects, making the validation results being more suitable for factorial analysis. For a specific HPO task, an OLH with a set of factors and a number of runs must be specified. Each factor corresponds to a hyperparameter. The number of runs should be chosen to meet the restrictions of OA design ($N = \lambda l^t$), where $l$ is the number of level, the index $\lambda$ describes the number of repeats of each level for each factor, $t$ is the strength that describes how many

hyperparameters might simultaneously interact. MOFA offers to balance between a low polynomial and high expressiveness, e.g. a higher strength $t$ provides better decorrelation while the restriction of OA would be stricter, and vice versa. Also, MOFA provides to balance between *exploration* and *exploitation*, e.g. a larger index $\lambda$ makes MOFA *explore* more hyperparameter configurations but consume more budget in a single round, and vice versa. One might also consider using more flexible designs such as near-orthogonal Latin Hypercube [28] to avoid these restrictions, while more flexible design might suffer from some loss of performance. These interesting *trade-off* problems are left for our future research.

## 3.2 Evaluator: Parallel Evaluation

Different from model-based methods that learn a parametric model and update it based on previous experiences, the hyperparameter configurations sampled by OLH can be evaluated in parallel. In practice, we evaluate different hyperparameter configurations asynchronously on multiple workers. We analyze the parallelization in Sec. 5.3.

## 3.3 Transformer: Transforming OLH into OA

After evaluating the sampled hyperparameter configurations, we build an OLH performance table to store all the hyperparameter settings and evaluation results. By default, as depicted in Fig. 1c (left), there are only continuous factor levels. To allow for factorial analysis that only supports discrete levels, we collapse the continuous levels in each column of OLH into $L$ discrete ordered levels (highlighted with colors). Since the OLH meets orthogonality, the collapsed OLH will be an OA of the same size (corresponding to $N$ runs).

## 3.4 Analyzer: Factorial Analysis

**Factorial Performance Analysis:** Based on the OA, we first calculate the *marginal mean* (MM) performance for each factor $F^i$ at level $l$ by $\mathrm{MM}(F^i_l) = \frac{L}{N} \sum_{k \in [1...\frac{N}{L}]} Y(F^i_{lk})$, where $Y(F^i_{lk})$ is the response for $F^i$ at level $l$ on its $k_{th}$ row. Since OA satisfies orthogonality, the MM can be seen as an approximation of the overall performance of each level for each factor, even when some factors are highly correlated (e.g. the *learning rate* and the *learning rate decay*). Then, we pick the level with the best MM performance of each factor for subsequent search with $\mathrm{BL}(F^i) = \mathrm{argmax}_{l \in [1...L]} \mathrm{MM}(F^i_l)$. For example, in Fig. 1d, the best MM performance of factor *LR* is $0.819$, so the corresponding level 1 is selected as the best level of *LR*. Similarly, level 1 and level 2 are selected as the best levels of factor *LAM* and *Units*, respectively.

**Factorial Importance Analysis:** To narrow down the search space, we analyze the importance of each factor. We use the *marginal variance ratio* $\mathrm{MVR}(F^i) = \frac{\mathrm{MV}(F^i)}{\sum_{i=1}^{d} \mathrm{MV}(F^i)}$, to measure the relative importance of each factor, where $\mathrm{MV}(F^i) = \mathbb{E}(\mathrm{MM}(F^i) - \mathbb{E}(\mathrm{MM}(F^i)))^2$ is the *marginal variance* of factor $F^i$, $d$ is the number of factors. Since OA satisfies orthogonality, the $\mathrm{MV}(F^i)$ reflects the stability of results when varying factor $F^i$. Therefore, the MVR reflect its relative stability among all factors, which can be used to measure the relative importance of the factor. Specifically, if results produced by varying a factor are not stable, it indicates that the factor needs to be explored more. Conversely, if a factor is stable enough (the MVR is less than a specified threshold $\beta$), it does not need to be explored. In this case, we directly freeze it to be the current best level (median of the current search space). For example, in Fig. 1d, the importance (MVR) of factor *LR*, *LAM* and *Units* in the current best level are $0.55$, $0.37$ and $0.08$ respectively, and the importance of *Units* is less than the specified threshold $0.1$, so we freeze *Units* to be $0.5$ (the median of search space $[0, 1]$).

## 3.5 Final Hyperparameter Selection

The termination condition can be designed case-by-case, e.g. defined by a a maximum budgets or stop if all the hyperparameters are frozen. Once the iteration ends, two strategies are available to select the final hyperparameter configuration. (1) *Greedy Strategy*: we choose the configuration with the best performance among all of the evaluated experiments as the final hyperparameter setting. (2) *Mean Strategy*: we do a factorial performance analysis for the hyperparameters that have not yet been determined and use the median level of the search space with the largest marginal mean performance

as the final configuration. We combine these two strategies by picking the best hyperparameter configuration among these two strategies.

## 4 Theoretical Analysis

In this section, the theoretical properties of MOFA are analyzed with regard to two aspects. First, we study how MOFA controls the worst-case response. Second, we prove that the inference of MOFA is more reliable. Precisely, if we make a hypothesis test for the inference of MOFA, it will have a higher confidence level.

### 4.1 Worst-case Response

For factorial analysis (Fig. 1d), MOFA aims to compare different levels and choose the best one. The mean statistics $\bar{Y}_i = \sum_{X_j \in \text{level } i} f(X_j)$ is used to estimate the effect of the $i$-th level, $\mathbb{E}_{X \in \text{level } i} f(X)$, where $f(X)$ is the response of the factor vector $X$. Then, according to these estimates, the level with the highest mean statistics is picked as the best level. Although the estimator will approach the effect asymptotically, it can be poor when the dataset $P = (X_1, X_2, \ldots, X_N)$ is small. In this case, one can use the Koksma-Hlawka inequality [29, 30, 31] to control the worst case. The Koksma-Hlawka inequality is stated as Eq.(1).

$$|\mathbb{E}f(X) - \bar{Y}| \leq V(f) \cdot D^*(P) \tag{1}$$

where $V(f)$ is the variation in the sense of Hardy and Krause [32] and $D^*$ is the star discrepancy [33]. The variation $V(f)$ is a constant in our situation. Hence, reducing the star discrepancy of $P$ is a natural way to control the worst-case response.

**Definition 1** *The star discrepancy is defined as follows,*

$$D^*(P) = \sup_{(b_1, b_2 \ldots, b_d) \in [0,1)^d} \left| \frac{A(B; P)}{N} - \lambda_d(B) \right|, \tag{2}$$

*where $\lambda_d$ is the $d$-dimensional Lebesgue measure, $A(B; P)$ is the number of points in $P$ that fall into $B$, and $B = \{(z_1, z_2 \ldots, z_d) \in \mathbb{R}^d \mid 0 \leq z_j < b_j \text{ for } j = 1, 2, \ldots, d\}$.*

Roughly speaking, the discrepancy of $P$ is low if the proportion of points in $P$ falling into an arbitrary set $B$ is close to proportional to the measure of $B$.

As mentioned before, MOFA has univariate projection uniformity. It is exactly the reason why we can control the worst case. Notice that the one-dimensional projections of samples obtained from MOFA are much more evenly distributed than that of random sampling. For intuitive explanation, let $d = 1$, we can see that the star discrepancy of MOFA with size $N$ is at most $1/N$ while the star discrepancy of random samplings of the same size has the order of $1/\sqrt{N}$ [34].

### 4.2 Hypothesis Test

For the MOFA procedure, we compare the mean performance of levels, and choose the level with highest response. To make this inference convincing, it is necessary to show that this level is statistically significantly better than other levels. Thus, a hypothesis test as follows is introduced to analyse this problem,

$$\begin{aligned} H_0 : \ &\mu_1 \geq \mu_2, \ \sigma_1 \text{ and } \sigma_2 \text{ are unknown, but cannot be} \\ &\text{assumed to be equal;} \\ H_a : \ &\mu_1 < \mu_2, \sigma_1 \text{ and } \sigma_2 \text{ do not change,} \end{aligned} \tag{3}$$

where $\mu_i$ and $\sigma_i$ is mean and standard deviation of the effect of the $i$-th level. Let $\bar{y}_i$ and $s_i$ be the estimators of $\mu_i$ and $\sigma_i$, respectively. Without loss of generality, assume $\bar{y}_1 < \bar{y}_2$ and the inference is that the second level is better than the first level. Then, the $p$-value of the hypothesis test (3) is used to measure the confidence of this inference. To show the advantage of MOFA, the following theorem is given.
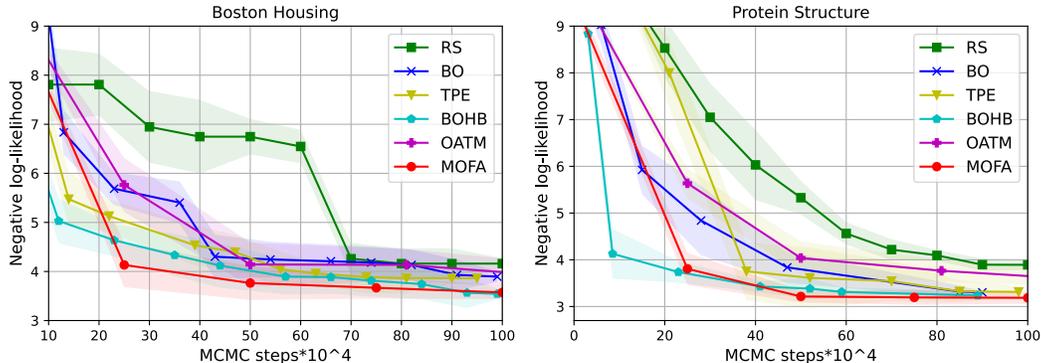
5

Figure 2: The negative log-likelihood of BNN on two different UCI datasets ($\lambda = 1, l = 5, t = 2$).

**Theorem 1** *Let $p_M$ and $p_R$ denote the p-value of the hypothesis test $H_0$ with samples in MOFA and samples in Random Search, respectively. Then, we have $p_M < p_R$.*

From Theorem 1, the confidence level of the inference of MOFA is higher than that of random search.

The Key Point of **Proof:** This comparison of two populations is called the two sample Behrens-Fisher problem. In statistics, Welch's t-test, or unequal variances t-test [35] is a standard solution to this problem. It defines the statistic $t$ by the Eq. (4),

$$ t \quad = \quad (\bar{Y}_1 - \bar{Y}_2) \bigg/ \sqrt{\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2}}. \tag{4}$$

The difficulty of calculating the $p$-value is how to calculate the freedom degree of this $t$ statistic. In MOFA, the factorial importance analysis helps to make this easier, and the sampling strategy helps to improve the confidence level. For details please refer to the supplementary materials.

## 5    Empirical Evaluation

In this section, we comprehensively evaluate MOFA on three different settings. 1) a two-layer Bayesian neural network (BNN); 2) the image classification with ResNet on CIFAR10 and 3) an deep neural networks for EEG-based intention recognition.

### 5.1    Experiment Settings

**Bayesian Neural Networks.** Following [13], we optimize five hyperparameters of a two-layer BNN: the number of units in layer 1 and layer 2, the step length, the length of the burn-in period, and the momentum decay. The initial search intervals for these hyperparameters are set to $[2^4, 2^9]$, $[2^4, 2^9]$, $[10^{-6}, 10^{-1}]$, $[0, 0.8]$ and $[0, 1]$ respectively. To ensure uniformity of hyperparameter search space, we perform log transformation [13] on the first three hyperparameters. We reuse an open-source code[2] to implement the BNN and the baseline methods. Two different UCI datasets, Boston Housing and Protein Structure described in [36] are used for evaluation. The BNN is trained with Markov Chain Monte-Carlo (MCMC) sampling and the number of steps for the MCMC sampling is used as the budgets, we report the negative log-likelihood on the validation data as the final performance.

**ResNet on CIFAR-10.** We evaluate MOFA on ResNet (Wide Residual Networks [37] with Cutout [38] regularization) for image classification on CIFAR10. The hyperparameters optimized include the base learning rate, the interval of the learning rate, the decay rate of the learning rate, the cutout number and the cutout length. The initial search intervals for these hyperparameters are set to $[10^{-10}, 10^{-1}]$ (with log transformation), $[0.01, 1]$, $[0.01, 0.99]$, $[0, 3]$ and $[1, 20]$ respectively. We set the maximum total budget (epochs) to 12000. For BOHB, we set the minimum budget to 5 and maximum budget to 135. 5000 training images are split off as the validation set and the accuracy on the test set is reported for comparison.

---

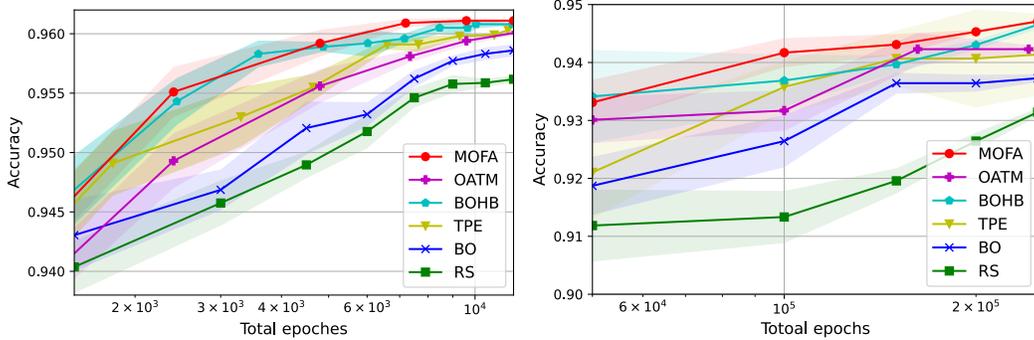[2]https://github.com/automl/HpBandSter/tree/icml_2018

Figure 3: The accuracy of ResNet on CIFAR10 (*left*) and the accuracy of EEG-based intention recognition (*right*). We set the size of OLH to $\lambda = 1, l = 5, t = 2$.

**EEG-based Intention Recognition.** Following [22], we evaluate MOFA on a practical task, EEG-based intention recognition with a deep convolutional neural network, the details of the architecture of the networks is described in [22]. A subset of the dataset (28000 samples) is used in the experiment. The dimension of each sample is 64, which corresponds to 64 channels. The dataset is divided into a training set (80%) and testing set (20%), and we report the models' accuracy on the testing set. We optimize three hyperparameters: the learning rate $lr$, the regularization coefficient $\lambda$, and the number of units in each hidden layer. The initial search intervals for the three hyperparameters is set to $[0.0005, 0.01]$, $[0.0005, 0.01]$ and $[64, 1024]$ with log transformation, respectively.

**Baselines and Implementation Details.** We compare MOFA with Random Search [39], BO [11], TPE [40], BOHB [13] and OATM [22]. For BO and BOHB, the Gaussian Process with Matern kernel function is used as the surrogate model, and expected improvement is used as the aquisition function. As OATM cannot be conducted iteratively, we use OAs with different number of runs for fair comparisons. We set the number of runs of OATM to 25, 50, 81, 121 respectively. For MOFA, we build an OLH with index $\lambda = 1$, level $l = 5$ and strength $t = 2$ for comparison. In this way, the number of runs in each round will be $N = \lambda l^t = 25$. In an ablation study (Section 5.2) we show how $\lambda$, $l$ and $t$ influence the performance by varying the settings of these numbers. All experiments are run 5 times with different random seeds, and the standard deviations are reported.

## 5.2 Hyperparameter Optimization Results

**Overall Results.** Fig. 2 shows the HPO results for BNN within two different UCI datasets. We firstly find that MOFA is significantly better (+12%) than Random Search with any budget. For the Boston Housing dataset, the performance of MOFA is even 40% better than that of Random Search when the budget is limited (<60). Compared with BO and TPE, MOFA improves results by +5% on average, although their performance goes to the same level when increasing budgets. Furthermore, MOFA is even better than BOHB on the Boston Housing dataset and comparable to BOHB on the Protein Structure dataset. One should note that BOHB is still model-based and can not be fully parallelized, while MOFA is a highly parallelizable method. We also find that OATM outperforms Random Search on two datasets and outperforms BO on Boston Housing but it has no advantage over other model-based algorithms like BOHB, which confirms our assumption that inference analysis (*exploitation*) is important to narrow down the search space. Fig. 5.2 shows the performance for HPO on CIFAR10 with ResNet.The results demonstrate that MOFA consistently outperforms baselines, with only BOHB coming close to MOFA. Fig. 5.2 shows the results for HPO on EEG-based intention recognition, showcasing that MOFA outperforms all the baseline solutions including BOHB, improving by +3% and +2% compared to RS and BO, respectively.

**Ablation Study of Varying Index, Level and Strength.** To evaluate the sensitivity of the model's performance to the index, level and strength. We first fix $l = 5$ and $t = 2$ and change the number of $\lambda$ to $1, 2, 3, 4, 5$. Since the number of runs of an OLH is $N = \lambda l^t$, we can get OLHs with runs $25, 50, 75, 100, 125$. For OLH with 25 runs ($\lambda = 1$), we ran it for five rounds, and for OLH with 50 ($\lambda = 2$) runs, we run it for two rounds, while the OLH with 75 ($\lambda = 3$), 100 ($\lambda = 4$), 125 ($\lambda = 5$) runs are only conducted for one round, which means there is no factorial analysis for these three experiments. We further study the the varying number of $l$ and $t$ by fixing $\lambda = 1$. For OLH with $t = 2$, $l$ is set to $3, 5, 7$ respectively, while for OLH with $t = 3$, $l$ is set to $3, 5$. Fig. 5.2 shows the results of MOFA with different numbers of $\lambda$. It clearly shows that OLH ($\lambda = 1$) is better than OLH
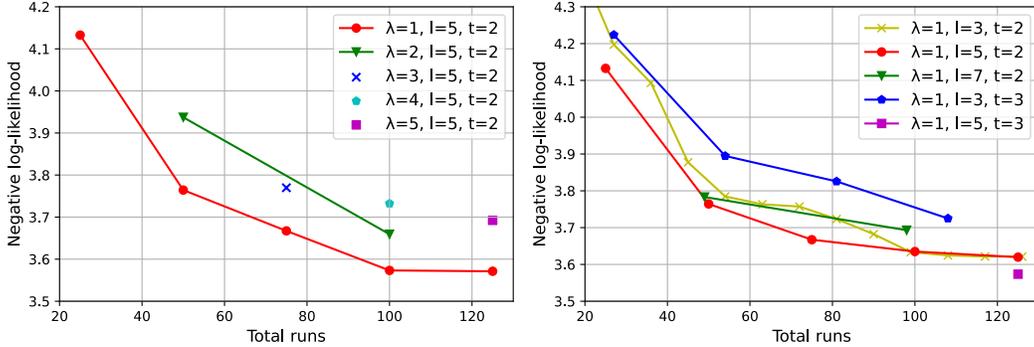
Figure 4: The negative log-likelihood of MOFA with varying $\lambda$ (left) and MOFA with varying $l$ and $t$ (right) on BNN task. The maximum budgets (number of runs) is fixed to $125$.

| Method | Boston | Protein | CIFAR-10 |
|--------|--------|---------|----------|
| RS | 6.756 | 5.234 | 0.938 |
| OA-sampler | 5.133 | 5.125 | 0.944 |
| LH-sampler | 5.213 | 5.084 | 0.943 |
| OLH-sampler | 4.842 | 5.016 | 0.951 |
| MOFA (#1) | 4.232 | 3.824 | 0.955 |
| MOFA (#2) | 3.735 | 3.258 | 0.958 |
| MOFA (#3) | **3.658** | **3.249** | **0.962** |

Table 1: The comparison of MOFA of varying rounds with different sampling schemes.



Figure 5: The GPU resource usage ratio.

($\lambda = 2$). While OLH ($\lambda = 2$) with two rounds outperforms OLH with $\lambda = 4$, though taking the same total runs. In a nutshell, we find that with the same number of runs, increasing the number of iterations boosts the performance, which also confirms the importance of the factorial analysis (reducing the search space). Fig. 5.2 further shows that OLH with $l = 5$ perform better than that with $l = 3$ and $l = 7$, and the increasing strength of OLH does not boost the performance.

**Module Contributions Analysis.** To study the module contributions (e.g. how much gain is contributed by factorial design and how much gain is from factorial analysis), we compare MOFA with different sampling schemes including RS, OA-sampler, LH-sampler and OLH-sampler on three datasets (Boston, Protein and CIFAR-10). Table. 1 clearly shows that OLH-sampler performs better than other sampling schemes, demonstrating the benefits of using OLH (factorial design). While MOFA outperforms OLH-sampler, showcasing the advantages of using factorial analysis.

**Computation Efficiency.** In practice, it is very important to take advantage of parallel resources efficiently, as it can significantly reduce the running time of HPO. We compare the GPU resource usage ratio of MOFA with BO (a fully sequential based method) and BOHB (a semi-parallel HPO method). We ran them on the same GPU server with $8$ GPU cards and report the GPU usage ratio. Fig. 5 shows that MOFA almost makes full use of the GPU resources. While for BO, a large proportion of GPU resources are idle. The resource utilization efficiency of BOHB is higher than BO, but it is not as stable as MOFA.

**MOFA as Initialization.** Another advantage of MOFA is that it can be easily applied to existing HPO methods as a space initialization module. We apply MOFA as an initialization for BO. Specifically, we firstly conduct MOFA for one or two rounds to narrow down the search space and then conduct BO on the new search space. Fig. 6 shows that with initialization of MOFA, BO outperforms the BO that starts from scratch. Also, we find that using MOFA initialization for one round is much better than using it for two rounds. However, though initialized with MOFA, we find that the performance of BO still cannot exceed the original results of MOFA. One might consider to use MOFA as initialization until the search intervals become small enough, and then continue to search better hyperparameter configurations with other HPO methods.

8

## 5.3 Analysis

**Model-Free.** Apparently, MOFA follows the standard pipeline of HPO (exploration & exploitation), but it does not depend on any parametric model. It is lightweight, model-free and easy to implement.

**Parallelization.** Let the time required by MOFA to prepare one iteration of exploration and to exploit the results be denoted by $a$. Assume that the maximal time needed to evaluate a single hyperparameter configuration is $t_e$, where $e$ is the $e_{th}$ iteration, the number of processes working in parallel is $n$, the number of iterations is $E$, and the number of runs in OLH is $N$. Then, the maximum total time spent in parallel mode is $T = \sum_{e=1}^{E} a + (t_e \times N/n)$. Compared to computing efforts of training and validation, the time required by MOFA (or most other HPO methods) for preparation and analysis $a$ is negligible. Thus, we can approximate $T \approx \sum_{e=1}^{E} t_e \times N/n$. Con-



Figure 6: Comparison of BO and BO with MOFA initialization for BNN task.

sidering Figures 2 to 4 and Fig. 5 together, one can recognize that increased parallelism of MOFA directly translates into efficiency gains with a speed-up in the order of $n$.

**Sample Efficiency.** Theoretically, the discrepancy of data points can directly affect efficiency of sampling. MOFA uses an OLH-sampler that ensures better discrepancy than other sampling methods, orthogonality of OLH can further improve sample efficiency by eliminating the influence of correlated hyperparameters. Table 1 shows that the OLH-sampler performs better than other sampling schemes.

**Categorical hyperparameters.** MOFA supports both continuous and discrete hyperparameters. In order to support categorical hyperparameters (e.g. selection of optimizers), one solution is to treat the categorical hyperparameters as one-hot hyperparameters. e.g., we treat the *optimizer* as multiple hyperparameters such as *SGD*, *Adam* and *RMSprop*, while each of them has binary $(0, 1)$ levels.
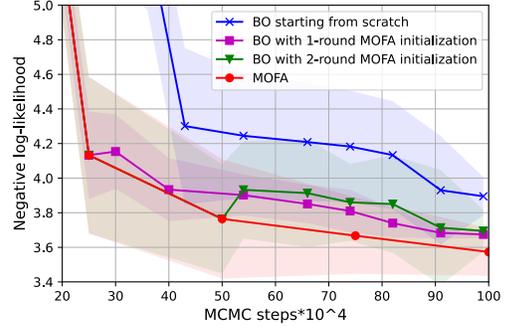
## 6 Related Work

**Model-Based Methods.** Model-based HPO learns a parametric model such as a surrogate model to search optimal hyperparameters guided by previous experiences. For example, BO [11] is a sequential model-based method that balances exploration and exploitation. It is efficient but highly depends on the surrogate model that models the objective function and does not allow for parallelization, since it is a sequential learning-based process. Based on BO, optimization with an early-stopping policy such as Successive Halving (SH) [41], HyperBand [42], BOHB [13] and ASHA [40] were proposed and achieved significant improvements. Different from BO, TPE [40] learns a surrogate model with a graphic structure, which is more suitable in dealing with conditional variables.

**Model-Free Methods.** Model-free HPO tries to tune hyperparameters without any parametric model. Grid Search conducts an exhaustive search on all candidate hyperparameter configurations. It is the most straightforward method but the consumption of computational resources of Grid Search grows exponentially with the increasing number of hyperparameters. Random Search [40] randomly samples a subset of hyperparameter configurations and can assign different budgets to different configurations. However, the randomly selected configurations cannot guarantee to find the global optimum. More recent literature [20, 21] have applied factorial designs into HPO tasks. [43] proposed a fully parallel HPO method based on Latin Hypercube, and provided some theoretical analyses, but they only used it to explore the hyperparameter space without corresponding factorial analyses. For factorial analyses, [22] proposed OATM based on OAs, and conducted orthogonal marginal analysis to select best levels for hyperparameters. However, OATM only supports discrete search space and cannot improve hyperparameter configurations by successive iterations.

## 7 Conclusion and Future Work

This paper presents a novel HPO method, MOFA, which combines the advantages of being model-free, parallelizable and sample efficient. To the best of our knowledge, MOFA is the very first step to

exploit factorial design and analysis on HPO, and more sophisticated techniques in these areas can be explored in the future, e.g., exploring the *trade-off* issues (strength $t$ and index $\lambda$) in OLH designs and exploiting more flexible factorial design such as near-orthogonal Latin Hypercubes [28].

## Broader Impact

Hyperparameter Optimization (HPO) is a pivotal task in automated machine learning. Most of current HPO methods rely heavily on model parametrization and suffer from poor computation efficiency. Rather than learning a parametric sequential-based model, MOFA try to improve the sampling by constructing a low-discrepancy and weak-correlated set of hyperparameters and determining the most promising hyperparameters. MOFA provide a more lightweight solution to HPO, which is model-free and has high efficiency.

One potential issue of MOFA, like many other HPO methods, is that it requires some human knowledge to determine the best size of OLH such like index, level and strength. We advocate peer researchers to look into this to enhance the intelligence of HPO by reducing the human intervention.

## Acknowledgments

## References

[1] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 2818–2826. IEEE Computer Society, 2016.

[2] Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 779–788. IEEE Computer Society, 2016.

[3] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008, 2017.

[4] Bo Xiong, Peng Bao, and Yilin Wu. Learning semantic and relationship joint embedding for author name disambiguation. *Neural Computing and Applications*, pages 1–12, 2020.

[5] Matthias Feurer and Frank Hutter. Hyperparameter optimization. In *Automated Machine Learning*, pages 3–33. Springer, Cham, 2019.

[6] Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: differentiable architecture search. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.

[7] Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 19–34, 2018.

[8] Ekin D. Cubuk, Barret Zoph, Dandelion Mané, Vijay Vasudevan, and Quoc V. Le. Autoaugment: Learning augmentation strategies from data. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 113–123. Computer Vision Foundation / IEEE, 2019.

[9] Quanming Yao, Mengshuo Wang, Yuqiang Chen, Wenyuan Dai, Hu Yi-Qi, Li Yu-Feng, Tu Wei-Wei, Yang Qiang, and Yu Yang. Taking human out of learning applications: A survey on automated machine learning. *arXiv preprint arXiv:1810.13306*, 2018.

[10] Matthias Feurer, Aaron Klein, Katharina Eggensperger, Jost Tobias Springenberg, Manuel Blum, and Frank Hutter. Efficient and robust automated machine learning. In Corinna Cortes, Neil D. Lawrence, Daniel D. Lee, Masashi Sugiyama, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 2962–2970, 2015.

[11] Jonas Mockus, Vytautas Tiesis, and Antanas Zilinskas. The application of bayesian methods for seeking the extremum. *Towards global optimization*, 2(117-129):2, 1978.

[12] Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. Sequential model-based optimization for general algorithm configuration. In *International conference on learning and intelligent optimization*, pages 507–523. Springer, 2011.

[13] Stefan Falkner, Aaron Klein, and Frank Hutter. BOHB: robust and efficient hyperparameter optimization at scale. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 1436–1445. PMLR, 2018.

[14] Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.

[15] Tong Yu and Hong Zhu. Hyper-parameter optimization: A review of algorithms and applications. *arXiv preprint arXiv:2003.05689*, 2020.

[16] Huong Ha, Santu Rana, Sunil Gupta, Thanh Nguyen, Hung Tran-The, and Svetha Venkatesh. Bayesian optimization with unknown search space. *arXiv preprint arXiv:1910.13092*, 2019.

[17] Frank Hutter, Holger Hoos, and Kevin Leyton-Brown. An efficient approach for assessing hyperparameter importance. In *International conference on machine learning*, pages 754–762. PMLR, 2014.

[18] Jan N Van Rijn and Frank Hutter. Hyperparameter importance across datasets. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2367–2376, 2018.

[19] CF Jeff Wu and Michael S Hamada. *Experiments: planning, analysis, and optimization*, volume 552. John Wiley & Sons, 2011.

[20] Dimo Brockhoff, Bernd Bischl, and Tobias Wagner. The impact of initial designs on the performance of matsumoto on the noiseless bbob-2015 testbed: A preliminary study. In *Proceedings of the Companion Publication of the 2015 Annual Conference on Genetic and Evolutionary Computation*, pages 1159–1166, 2015.

[21] Wolfgang Konen, Patrick Koch, Oliver Flasch, Thomas Bartz-Beielstein, Martina Friese, and Boris Naujoks. Tuned data mining: a benchmark study on different tuners. In *Proceedings of the 13th annual conference on Genetic and evolutionary computation*, pages 1995–2002, 2011.

[22] Xiang Zhang, Xiaocong Chen, Lina Yao, Chang Ge, and Manqing Dong. Deep neural network hyperparameter optimization with orthogonal array tuning. In *International Conference on Neural Information Processing*, pages 287–295. Springer, 2019.

[23] MD McKay, RJ Beckman, and WJ Conover. Comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics*, 21(2):239–245, 1979.

[24] Boxin Tang. Orthogonal array-based latin hypercubes. *Journal of the American statistical association*, 88(424):1392–1397, 1993.

[25] V Roshan Joseph and Ying Hung. Orthogonal-maximin latin hypercube designs. *Statistica Sinica*, pages 171–186, 2008.

[26] Jessica Franco. Exploratory designs for computer experiments of complex physical systems simulation. *Theses, Ecole Nationale Supérieure des Mines de Saint-Etienne*, 2008.

[27] Boxin Tang. Selecting latin hypercubes using correlation criteria. *Statistica Sinica*, pages 965–977, 1998.

[28] Thomas M Cioppa and Thomas W Lucas. Efficient nearly orthogonal and space-filling latin hypercubes. *Technometrics*, 49(1):45–55, 2007.

[29] JF Koksma. Een algemeene stelling uit de theorie der gelijkmatige verdeeling modulo 1. *Mathematica B (Zutphen)*, 11(7-11):43, 1942.

[30] Edmund Hlawka. Discrepancy and riemann integration. *Studies in Pure Mathematics*, 3, 1971.

[31] Christoph Aistleitner, Florian Pausinger, Anne Marie Svane, and Robert F Tichy. On functions of bounded variation. *arXiv preprint arXiv:1510.04522*, 2015.

[32] Florian Pausinger and Anne Marie Svane. A koksma–hlawka inequality for general discrepancy systems. *Journal of Complexity*, 31(6):773–797, 2015.

[33] Michael Drmota and Robert F Tichy. *Sequences, discrepancies and applications*. Springer, 2006.

[34] Benjamin Doerr. A lower bound for the discrepancy of a random point set. *Journal of Complexity*, 30(1):16–20, 2014.

[35] Bernard L Welch. The significance of the difference between two means when the population variances are unequal. *Biometrika*, 29(3/4):350–362, 1938.

[36] José Miguel Hernández-Lobato and Ryan P. Adams. Probabilistic backpropagation for scalable learning of bayesian neural networks. In Francis R. Bach and David M. Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 1861–1869. JMLR.org, 2015.

[37] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In Richard C. Wilson, Edwin R. Hancock, and William A. P. Smith, editors, *Proceedings of the British Machine Vision Conference 2016, BMVC 2016, York, UK, September 19-22, 2016*. BMVA Press, 2016.

[38] Terrance DeVries and Graham W Taylor. Improved, regularization of convolutional neural networks with cutout., arxiv. *preprint*, 2017.

[39] James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *The Journal of Machine Learning Research*, 13(1):281–305, 2012.

[40] James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. In John Shawe-Taylor, Richard S. Zemel, Peter L. Bartlett, Fernando C. N. Pereira, and Kilian Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain*, pages 2546–2554, 2011.

[41] Manoj Kumar, George E Dahl, Vijay Vasudevan, and Mohammad Norouzi. Parallel architecture and hyperparameter search via successive halving and classification. *arXiv preprint arXiv:1805.10255*, 2018.

[42] Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband: A novel bandit-based approach to hyperparameter optimization. *The Journal of Machine Learning Research*, 18(1):6765–6816, 2017.

[43] Marie-Liesse Cauwet, Camille Couprie, Julien Dehos, Pauline Luc, Jérémy Rapin, Morgane Rivière, Fabien Teytaud, Olivier Teytaud, and Nicolas Usunier. Fully parallel hyperparameter search: Reshaped space-filling. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 1338–1348. PMLR, 2020.

[44] Art B Owen. Randomly permuted (t, m, s)-nets and (t, s)-sequences. In *Monte Carlo and quasi-Monte Carlo methods in scientific computing*, pages 299–317. Springer, 1995.

## Appendix

## A Proof of Theorem 1

The Welch's t-test is used only when the two population variances are not assumed to be equal and hence must be estimated separately. The t statistic to test whether the population means are different is calculated as:

$$t \quad = \quad (\bar{Y}_1 - \bar{Y}_2) \Big/ \sqrt{\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2}},$$

where $s_i^2$ $(i = 1, 2)$ is the unbiased estimator of the variance of each of the two samples, and is assumed that $s_1 > s_2$ without loss of generality. This statistic can be also used for testing whether one population mean is larger than the other one. We just need to change the two-sided hypothesis test into the one-sided hypothesis test.

In [44], it is proved that

$$\text{Var}(\bar{Y}_M) = \text{Var}(\bar{Y}_R) - c/N + o(N^{-1}),$$

where $c$ is a non-negative constant. Hence, the t statistic $t_M > t_R$ which means $p_M < p_R$ with the same degree of freedom. Note that when the value of t statistic is the same, the p-value $p_M < p_R$ is equivalent to the degree of freedom of MOFA is larger than that of Random sampling. Next, we calculate their d.f.s to check this fact.

The degree of freedom is calculated by the Welch–Satterthwaite equation as follows,

$$\text{d.f.} = \frac{\left(\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2}\right)^2}{\frac{\left(s_1^2/N_1\right)^2}{N_1 - 1} + \frac{\left(s_2^2/N_2\right)^2}{N_2 - 1}}.$$

In our cases, the number of trials in different levels $N_1 = N_2 = N$. Then, we have

$$\text{d.f.} = \frac{(N-1)\left(s_1^2 + s_2^2\right)^2}{\left(s_1^2\right)^2 + \left(s_2^2\right)^2}.$$

We prove the fact that $\text{d.f.}(M) < \text{d.f.}(R)$ through checking the monotonicity of $\left(s_1^2 + s_2^2 - 2c\right)^2 \Big/ \left(\left(s_1^2 - c\right)^2 + \left(s_2^2 - c\right)^2\right)$ w.r.t the constant $c$ denoted by $d(c)$. To check this, we simplify it into

$$d(c) = \frac{\left(s_1^2 - c\right)^2 + \left(s_2^2 - c\right)^2 + 2(s_1^2 - c)(s_2^2 - c)}{\left(s_1^2 - c\right)^2 + \left(s_2^2 - c\right)^2}$$

$$= 1 + 2\left(\frac{s_1^2 - c}{s_2^2 - c} + \frac{s_2^2 - c}{s_1^2 - c}\right)^{-1},$$

which can be seen that it is monotonically decreasing w.r.t $c$.

Next, the p-value is calculated through the cumulative distribution function of $t$ distribution,

$$p = \frac{\int_0^{\frac{d.f.}{t^2 + d.f.}} x^{d.f./2 - 1}(1 - x)^{-1/2} dx}{2 \int_0^1 x^{d.f./2 - 1}(1 - x)^{-1/2} dx}$$

$$\triangleq \frac{\int_0^{\frac{d.f.}{t^2 + d.f.}} g(x) dx}{2 \int_0^1 g(x) dx}.$$

Thus, the p-value can be viewed as the area ratio of $g(x)$. For illustration, let $g(x) = b * x$. Then, no matter how $b$ changes, the triangle area ratio is fixed. Note that in our case, $g(x)$ is a convex function and monotonically decreasing w.r.t the degree of freedom. Since $d.f. > 2$, we know $g(0) = 0$ and $g(1) = \infty$. Hence, the right area of the line $x = \frac{d.f.}{t^2 + d.f.}$ plays a major role in the area ratio. According to the monotonicity, when $c$ comes larger, the line $x = \frac{d.f.}{t^2 + d.f.}$ moves to the left. Consequently, the area ratio becomes smaller. $\qquad \square$