

ARGONNE NATIONAL LABORATORY  
9700 South Cass Avenue  
Argonne, Illinois 60439

**Machine-Learning-Based Load Balancing for Community  
Ice Code Component in CESM**

**Prasanna Balaprakash<sup>1,2</sup>, Yuri Alexeev<sup>2</sup>, Sheri Mickelson<sup>1</sup>, Sven  
Leyffer<sup>1</sup>, Robert Jacob<sup>1</sup>, and Anthony Craig<sup>3</sup>**

Mathematics and Computer Science Division

Preprint ANL/MCS-P4070-0413

January 2015

---

<sup>2</sup> Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, IL 60439

<sup>3</sup> Leadership Computing Facility, Argonne National Laboratory, Argonne, IL 60439

<sup>4</sup> UCAR, Seattle, WA

# Machine-Learning-Based Load Balancing for Community Ice Code Component in CESM

Prasanna Balaprakash<sup>1,2</sup>, Yuri Alexeev<sup>2</sup>, Sheri A. Mickelson<sup>1</sup>, Sven Leyffer<sup>1</sup>,  
Robert Jacob<sup>1</sup>, and Anthony Craig<sup>3</sup>

<sup>1</sup> Mathematics and Computer Science Division, Argonne National Laboratory,  
Argonne, IL

<sup>2</sup> Leadership Computing Facility, Argonne National Laboratory, Argonne, IL

<sup>3</sup> UCAR, Seattle, WA

**Abstract.** Load balancing scientific codes on massively parallel architectures is becoming an increasingly challenging task. In this paper, we focus on the Community Earth System Model, a widely used climate modeling code. It comprises six components each of which exhibits different scalability patterns. Previously, an analytical performance model has been used to find optimal load-balancing parameter configurations for each component. Nevertheless, for the Community Ice Code component, the analytical performance model is too restrictive to capture its scalability patterns. We therefore developed machine-learning-based load-balancing algorithm. It involves fitting a surrogate model to a small number of load-balancing configurations and their corresponding runtimes. This model is then used to find high-quality parameter configurations. Compared with the current practice of expert-knowledge-based enumeration over feasible configurations, the machine-learning-based load-balancing algorithm requires six times fewer evaluations to find the optimal configuration.

## 1 Introduction

The Community Earth System Model (CESM) is one of the most widely used climate models in the world. Results from this model are a major part of the Intergovernmental Panel on Climate Change assessment reports [1]. CESM1.1.1 consists of six model components—atmosphere, ocean, sea-ice (CICE), land, river, and land-ice models—that communicate through a coupler. Each of the CESM model components has different scalability patterns and performance characteristics. In this paper, we focus on static load-balancing of computation, which is usually simple to implement with negligible overhead, making it suitable for “fine-grained” parallelism consisting of many small tasks. Previously, the load-balancing problem has been formulated as a mixed-integer nonlinear optimization problem and solved by using the optimization solver MINOTAUR [2]. This is a heuristic method that consists of gathering benchmarking data, calibrating a performance model using the data, and making decisions about optimal allocation by using the model. The performance model predicts the execution time of the program running in parallel as a function of problem size and

the number of processors employed. Nonetheless, several challenges in intramodel load balancing for the CICE computations occur only where sea ice is located and the sun is shining. This restriction presents a load-balance problem because processors are allocated across the entire Earth grid and several locations on the grid that do not have any sea ice [3]. The poor fit of the CICE results in inefficient processor allocations to all components—incorrect allocation of the CICE affects all other allocations because the total number of processors available to components is a fixed number. This is the primary motivation for us to develop sophisticated approaches for load balancing the CICE component of the CESM.

Recently, machine-learning methods [4] have received considerable attention for tuning performance of large scientific codes and kernels on high-performance computing systems. In particular, supervised machine-learning tries to learn the relationship between the input and the output of an unknown response function by fitting a model from few representative examples. When the model is accurate enough, it can predict the output at new unseen inputs, which provides numerous benefits, in particular when the evaluation becomes expensive.

In this paper, we present a machine-learning-based approach for static load-balancing problems, and we apply it to find high-quality parameter configurations for load balancing the CICE component of the CESM on IBM Blue Gene/P (BG/P). The novelty of the proposed algorithm consists of iteratively using the model to choose configurations with shorter predicted runtime for evaluation on the target architecture. We emphasize, however, that the algorithm is general and not specific to the CESM and/or BG/P. The contributions of the paper are as follows: 1) a machine-learning-based algorithm for static load-balancing problem, 2) deployment of a machine-learning method as a diagnostic tool for analyzing the sensitivity of the load-balancing parameters on the execution time, 3) empirical analysis of several state-of-the-art machine-learning methods for modeling the relationship between the load-balancing parameters and their corresponding execution time, and 4) 6x savings in core-hour usage for load balancing the CICE component of the CESM on BG/P.

## 2 The CICE Component on BG/P

For the CICE component, we need to find the optimal load-balancing parameter configuration  $x^*$  with the shortest the runtime ( $f^*$ ) for task counts  $\in \{80, 128, 160, 256, 320, 376, 512, 640, 800, 1024\}$ . The task count corresponds to number of MPI tasks; the number of OpenMP threads per MPI task is set to four because of memory restrictions on BG/P. The CICE component comprises six parameters. Three integer parameters, namely, maximum number of CICE blocks, `max.block`; the size of a CICE block in the first and second horizontal dimensions `block.x` and `block.y` respectively. Two categorical parameters that determine the decomposition strategy, `decomp.typ`  $\in \{\text{blkrobin, roundrobin, spacecurve, blkcart, cartesian}\}$  and `decomp.set`  $\in \{\text{null, slenderX1, slenderX1}\}$ . A binary parameter `mask.h`  $\in \{0,1\}$  that specifies to run the code with or without halo.

**Table 1.** Decomposition strategies and their corresponding `block.x` and `block.y` sizes

| <code>decomp.set</code> | <code>decomp.typ</code> | <code>block.x</code> | <code>block.y</code> |
|-------------------------|-------------------------|----------------------|----------------------|
| null                    | blkrobin, blkcart       | 1, 2, 4, 8           | 24, 48, 96,          |
|                         | roundrobin, spacecurve  |                      | 192, 3840            |
| slenderX1               | cartesian               | 4, 5, 8, 10          | 4 6 8 12             |
| slenderX2               |                         |                      |                      |

The constraints that define a feasible set  $\mathcal{D}$  of configurations are as follows. The parameter `max.block`  $\in \{1, 2, 3, 4, 5, 6, 7, 8, 10, 11, 12, 13, 14, 15, 16, 20, 24, 26, 30, 32, 40, 48, 64\}$  is determined by computing  $(\text{CICE\_X\_Grid\_Size} \times \text{CICE\_Y\_Grid\_Size}) / (\text{block.x} \times \text{block.y} \times \text{task count})$ . The feasible values for `decomp.set`, `decomp.typ`, `block.x`, and `block.y` are constrained as shown in Table 1. The decomposition strategies have different rules, and not all combinations of block sizes are possible. The `blkcart` method must have a multiple of four blocks per compute core. The `spacecurve` method must have 2, 3, and 5 only in `max.block`. The `slenderX1` method requires that the `block.x` multiplied by the task count divide evenly into the CICE X grid size. The value of `block.y` must also be divisible by the CICE Y grid size. The `slenderX2` method requires that the `block.x` multiplied by the task count be divisible by the CICE X grid size multiplied by 2. The decomposition also requires that the `block.y` multiplied by 2 divide evenly into the CICE Y grid size.

### 3 Machine-Learning Based Load-Balancing Algorithm

Given a set of training data  $\{(x_1, y_1), \dots, (x_l, y_l)\}$ , where  $x_i \in \mathcal{D}$  and  $y_i = f(x_i) \in \mathbb{R}$  are the load-balancing parameter configuration and its corresponding runtime, respectively, the supervised machine-learning approach includes finding a surrogate function  $h$  for the expensive  $f$  such that the difference between  $f(x_i)$  and  $h(x_i)$  is minimal for  $\forall i \in \{1, \dots, l\}$ . The function  $h$ , which is an empirical performance model, can be used to predict the runtimes of unevaluated  $x' \in \mathcal{D}$ . The key idea behind the machine-learning-based load-balancing algorithm is iteratively using the model to choose configurations with shorter predicted runtime for evaluation and retrain the model with the evaluated configurations.

The pseudo-code is shown in Algorithm 1. The symbols  $\cup$  and  $-$  denote set union and difference operators, respectively. Given a task count  $c$ , a pool  $\mathcal{X}_p$  of unevaluated configurations of task count  $c$ , the maximum number  $n_{\max}$  of allowed evaluations, and initial sample size  $n_s$ , the algorithm proceeds in two phases: parallel initialization phase and sequential iterative phase. In the initialization phase, the algorithm first samples  $n_s$  configurations at random and evaluates them in parallel to obtain their corresponding runtimes. A supervised learning method uses these points as a training set to build a predictive model. The sequential iterative phase consists of predicting the runtimes of all remaining unevaluated configurations using the model, evaluating the configuration with shortest predicted runtime, and retraining the model with the evaluation

---

**Algorithm 1** Pseudo-code for the machine-learning-based load-balancing algorithm

---

**Input:** task count  $c$ , configuration pool  $\mathcal{X}_p$  of task count  $c$ , max evaluations  $n_{\max}$ , initial sample size  $n_s$

```
1  $\mathcal{X}_{\text{out}} \leftarrow$  sample  $\min\{n_s, n_{\max}\}$  distinct configurations from  $\mathcal{X}_p$ 
2  $\mathcal{Y}_{\text{out}} \leftarrow$  Evaluate.Parallel( $c, \mathcal{X}_{\text{out}}$ )
3  $\mathcal{M} \leftarrow$  fit( $\mathcal{X}_{\text{out}}, \mathcal{Y}_{\text{out}}$ )
4  $\mathcal{X}_p \leftarrow \mathcal{X}_p - \mathcal{X}_{\text{out}}$ 
5 for  $i \leftarrow n_s + 1$  to  $n_{\max}$  do
6    $\mathcal{Y}_p \leftarrow$  predict( $\mathcal{M}, \mathcal{X}_p$ )
7    $x_i \leftarrow x \in \mathcal{X}_p$  with the shortest runtime in  $\mathcal{Y}_p$ 
8    $y_i \leftarrow$  Evaluate( $c, x_i$ )
9   retrain  $\mathcal{M}$  with  $(x_i, y_i)$ 
10   $\mathcal{X}_{\text{out}} \leftarrow \mathcal{X}_{\text{out}} \cup x_i$ ;  $\mathcal{Y}_{\text{out}} \leftarrow \mathcal{Y}_{\text{out}} \cup y_i$ 
11   $\mathcal{X}_p \leftarrow \mathcal{X}_p - x_i$ 
12 end for
```

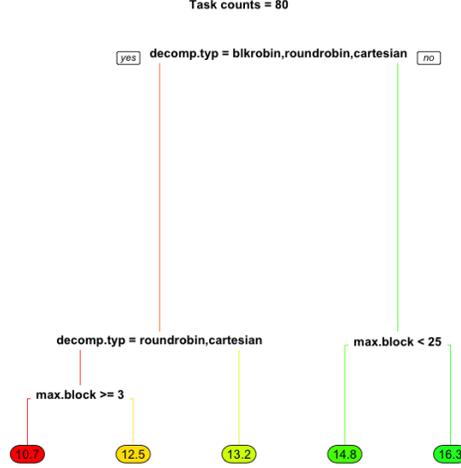
**Output:**  $x \in \mathcal{X}_{\text{out}}$  with the shortest runtime in  $\mathcal{Y}_{\text{out}}$

---

results. Without loss of generality, Algorithm 1 can be run in parallel for each task count  $c \in C$ . Because the best supervised learning algorithm depends on the relationship between the input and output, we test four state-of-the-art machine-learning algorithms as candidates for Algorithm 1: random forest (RF) [5], support vector machines (SVM) [6], Gaussian process regression (GP) [7], and neural networks (NN) [8].

RF belongs to the class of recursive partitioning methods [9]. They are widely used tools for predictive modeling in many scientific fields. These methods recursively partition the multi-dimensional input space  $\mathcal{D}'$  of training points into a number of hyper rectangles. The partition is done in such a way that input configurations with similar outputs fall within the same rectangle. The partition gives rise to a set of if-else rules that can be represented as a decision tree. For each hyper rectangle, a constant value is assigned—typically this is an average over the output values that fall within the given hyper rectangle. An example tree which is obtained on the CICE component data is shown in Figure 1. Given an unseen input  $x^* \in \mathcal{D}^* \subset \mathcal{D}$ , the algorithm uses the if-else rule to find the leaf and returns the corresponding constant value as the predicted value. RF uses a collection of regression trees, where each tree is obtained by the principle of recursive partitioning. For each tree generation, the algorithm takes a subsample of random points from the given training set. The subsample is either a bootstrap sample of the same size drawn with replacement or a subset of smaller size, drawn without replacement. Due to the randomness in the sampling, each subsample differs from each other. Given that each individual tree is built on the subsample, it can differ significantly from other trees. For a given  $x^*$ , each tree can make a prediction with respect to its own subsample. The power of RF comes from the aggregation of predicted output values from different trees and

the natural way of handling the categorical parameters. Consequently, it can deal with large dimensional inputs even in the presence of complex interactions and non-linearity.



**Fig. 1.** Illustration of a decision tree obtained via recursive partitioning on CICE component data for the task count 80.

SVM for nonlinear regression consists of mapping the given  $\mathcal{D}'$  of the training points into a high dimensional feature space and performing linear regression in the feature space:

$$g(\mathcal{D}') = \langle w \cdot \psi(\mathcal{D}') \rangle + b, \quad (1)$$

where  $\psi : \mathbb{R}^n \rightarrow \mathcal{F}$  being the nonlinear transformation,  $b$  being the bias term, and  $w \in \mathcal{F}$ . Finding  $g(\mathcal{D}')$  consists in specifying a loss function that need to be optimized and a kernel function  $k(\cdot)$  for nonlinearity transformation  $\psi$ . For the former, we use  $\epsilon$  intensive-loss function in which zero penalty is added to the loss function when predicted value of a training point is within  $\epsilon$  from its observed value. For the latter, we use the widely used Gaussian radial basis function kernel. Now, Equation 1 can be written as follows:

$$g(\mathcal{D}') = \sum_{i=1}^l \alpha_i \times [k(x_i, x_1), \dots, k(x_i, x_l)] + b, \quad (2)$$

$$k(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|_2^2}{2\sigma^2}\right), \quad (3)$$

where coefficients  $\alpha_i$  can be found by solving  $\epsilon$  intensive-loss function,  $\|\mathbf{x} - \mathbf{x}'\|_2^2$  is squared Euclidean distance that decreases with an increase in dissimilarity between  $x_i$  and  $x_j$ , and  $\sigma$  is a parameter of the kernel.

GP follows a probabilistic approach for regression. Given a training data of  $l$  points, GP assumes that  $\mathbf{Y} = [y_1, \dots, y_l]$  as a sample from a  $l$ -variate Gaussian distribution. For an unseen input  $x^*$ , the probability  $p(y^*|\mathbf{Y})$  follows the Gaussian distribution  $\mathcal{N}$  with a user defined kernel function  $k(\cdot)$ :

$$y^*|\mathbf{Y} \sim \mathcal{N}(\mathbf{K}_* \mathbf{K}^{-1} \mathbf{Y}, \mathbf{K}_{**} - \mathbf{K}_* \mathbf{K}^{-1} \mathbf{K}_*^{\mathbf{T}}), \quad (4)$$

where

$$\mathbf{K} = \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{bmatrix} \quad (5)$$

$$\begin{aligned} \mathbf{K}_* &= [k(x_*, x_1), \dots, k(x_*, x_n)], \\ \mathbf{K}_{**} &= k(x_*, x_*). \end{aligned}$$

Note that  $\mathbf{T}$  represents matrix transpose operation. For  $k(\cdot)$ , we use the Gaussian radial basis function as in Equation 3. The predicted value  $\hat{y}_*$  and variance  $var(y^*)$  of  $y^*$  are given by the parameters of  $\mathcal{N}$ :

$$\begin{aligned} \hat{y}^* &= \mathbf{K}_* \mathbf{K}^{-1} \mathbf{Y}, \\ var(y^*) &= \mathbf{K}_{**} - \mathbf{K}_* \mathbf{K}^{-1} \mathbf{K}_*^{\mathbf{T}}. \end{aligned} \quad (6)$$

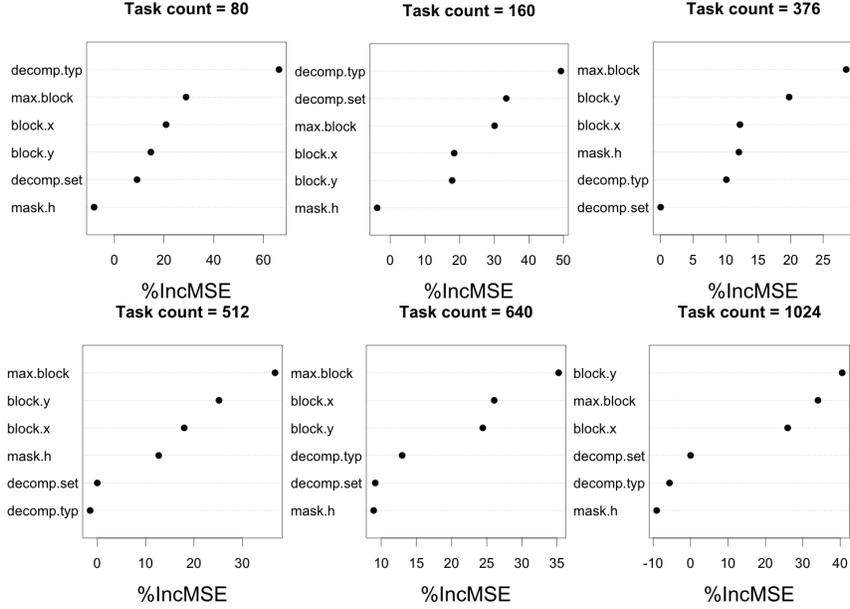
NN is a classical and one of most widely used supervised learning approaches. We focus on a single-hidden-layer neural network, an effective variant that comprises one input layer, one hidden layers, and one output layer. The nonlinear regression performed by NN can be written as follows:

$$\mathbf{Y} = h(\mathcal{D}') = \mathbf{B}\varphi(\mathbf{A}\mathcal{D}' + \mathbf{a}) + \mathbf{b},$$

where  $\mathbf{A}$  and is the matrix of weights and bias vector for the first layer (between input and hidden layer) and  $\mathbf{B}$  and  $\mathbf{b}$  are the weight matrix and the bias vector of the second layer (between hidden and output layer). The function  $\varphi$  denotes an element wise nonlinearity. The training in neural network consists in adapting all the weights and biases  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{a}$ , and  $\mathbf{b}$  to their optimal values for the training set  $\{(x_1, y_1), \dots, (x_l, y_l)\}$ . The optimization problem consists in minimizing the squared reconstruction error  $\sum_{i=1}^l \|h(x_i) - y_i\|^2$  and it can be solved effectively with back-propagation algorithm.

## 4 Experimental Results

We evaluated the effectiveness of the proposed load-balancing algorithm with the four machine-learning methods. In addition, we include two approaches in the comparison: Expert-knowledge-based enumeration (EE) and random search

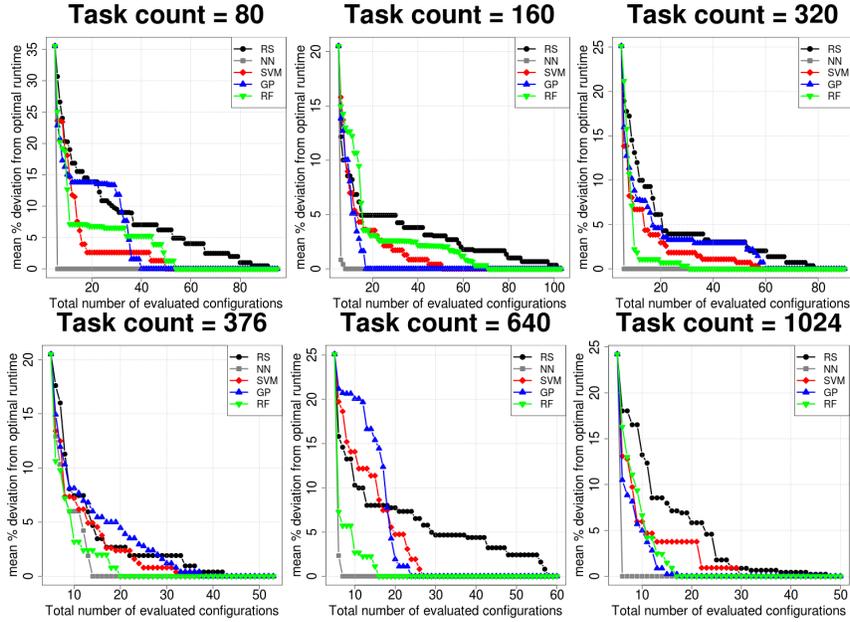


**Fig. 2.** Sensitivity analysis of the load-balancing parameters on the runtime of the CICE component for different task counts. For each parameter, the plot shows the percentage increase in mean squared error (%IncMSE) when the values of the corresponding parameter gets imputed.

(RS). EE is the current practice for finding the optimal load-balancing configuration for the CICE component of the CESM. In addition to the application-specific constraints, expert knowledge of the code and the architecture were used to prune the feasible set of configurations  $\mathcal{D}$  for the CICE component. As a result, for each task count  $c$ , there are 50 to 60 ( $|\mathcal{D}_c|$ ) feasible configurations; in total, for all the 10 task counts, there are  $|\mathcal{D}| = 653$  parameter configurations. This method evaluates all 653 parameter configurations. Moreover, we followed the current practice for defining the runtime  $f(x)$  for  $x$ : the code was run twice with the same  $x$  and the shortest runtime was taken as  $f(x)$ . In RS, for each task count  $c$ , parameter configurations were sampled at random without replacement from  $\mathcal{D}_c$  and were evaluated. To minimize the impact of randomness involved in the initialization procedure of Algorithm 1 and in the five approaches, we repeated all of them 10 times, each with a different random seed. Moreover, we stored the runtime of each configuration from EE in a lookup table and reused the results for running all other algorithms. For Algorithm 1, for each task count  $c$ ,  $\mathcal{D}_c$  obtained in the EE approach was given as the configuration pool  $\mathcal{X}_p$ , and the initial sample size  $n_s$  was set to 5. The approaches were implemented and run in the **R** programming language and environment [10] version 2.15.2 using the **nnet** (NN), **kernlab** (SVM, GP), and **randomForest** (RF) packages. The default parameter values were used for each method. Experiments were carried out on Intrepid, a BG/P supercomputer at Argonne.

**Sensitivity analysis:** First, we present an empirical analysis to explain why the previously proposed analytical performance model fails to predict the runtime of the CICE component and why distinct models may be constructed for each task count. For this purpose, we used the RF method to analyze the impact of each load-balancing parameter on the resulting runtimes. For the training data, we randomly sampled 50% of the data (parameter configuration and runtimes) obtained with EE approach. An RF model was fitted on this training set. The mean squared error (MSE) on the original training set is given by  $\frac{\sum_{i=1}^l (f(x_i) - \hat{f}(x_i))^2}{l}$ , where  $l$  is the number of training points, and  $f(x_i)$  and  $\hat{f}(x_i)$  are the original and predicted runtime value of parameter configuration  $x_i$ , respectively. In order to assess the impact of a parameter  $m$ , the values of  $m$  in the training set were randomly permuted. Again, an RF model was fitted on this imputed training set, and the mean squared error was computed. If a parameter  $m$  is important, then permuting the values of  $m$  should affect the prediction accuracy significantly and eventually increase the mean squared error. The results are shown in Figure 2. We observe that the trend in the parameter importance is not the same over all the task counts. For task counts up to 320, `decomp.set` and/or `decomp.type` have a strong impact on the runtimes; for large task counts, they become relatively less important—`max.block`, `block.x`, and `block.y` have a strong impact on the runtime. For 1024, only `max.block`, `block.x`, and `block.y` have an impact on the runtime; the other three parameters have negative %IncMSE, suggesting that they do not affect the runtime. In summary, the impact of parameter values on the runtimes and the type of nonlinear interactions between them change with an increase in the task counts. The previously developed analytical model does not take this effect into account for the CICE component, and consequently it falls short in runtime prediction. Moreover, if we build a single model for all task counts with task count being an input to the model, we might lose these task-count-specific interactions, thus affecting the runtime quality of the obtained configurations.

**Comparison between variants:** With EE as a baseline, we next examined the effectiveness of the five approaches in finding the optimal load-balancing configuration for the CICE component. As a measure of the effectiveness of each variant, we use the percentage deviation from the optimal runtime (*%dev*). Given a variant  $v$  and task count  $c$ , this is given by  $\frac{f_v^c - f_{opt}^c}{f_{opt}^c} \times 100$ , where  $f_v^c$  is the shortest runtime obtained by variant  $v$  and  $f_{opt}^c$  is the optimal runtime obtained from EE. Because we repeated each method 10 times to reduce the impact of randomness, we consider the mean percentage deviation from the optimal runtime of a variant as *%dev* averaged over 10 repetitions. We also used a statistical t-test to check whether the observed differences in the *%dev* of the variants are significant. Figure 3 shows the comparison between the approaches. The results show that RS requires almost the same number of evaluations as does EE for all task counts. These results indicate that the problem of finding high-quality configurations is not an easy task; clearly, we need more sophisticated approaches to find high-quality configurations within fewer evaluations. The variants of Algorithm 1 obtain optimal configurations with fewer evaluations, and they outperform RS.



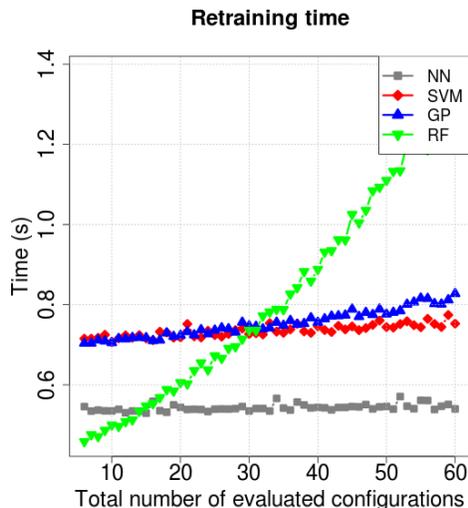
**Fig. 3.** Comparison between approaches for different task counts of the CICE component. The lines represent the mean percentage deviation from the optimal runtime as a function of the number of evaluated configurations.

NN completely dominates all other variants and RS. The key advantage of NN comes from its requiring less than 10 evaluations to find the optimal parameter configuration on 9 out of 10 task counts—only on  $c = 376$ , does it require 15 evaluations.

In Table 2, we analyze  $\%dev$  of each variant, when it is allowed to perform only 10 evaluations (for machine-learning variants this corresponds to five evaluations after the initialization). The results show that mean  $\%dev$  of NN is zero and it lower than all other variants. For all but one task counts, the observed differences are significant in statistical sense. NN fails to find optimal runtime for  $c = 376$ , where it is 6% away from the optimal runtime and it is comparable to other approaches.

As soon as a new evaluation becomes available, each machine-learning variants is retrained on all the available input-output pairs. This is the most computationally expensive part in the iterative phase of Algorithm 1. In Figure 4, we analyze the retraining time required by the machine-learning variants after each evaluation. The reported time is an average time over all repetitions and task count. The results show that NN outperforms all other variants in retraining time. The time remains fairly constant throughout with an average of 0.5 seconds. This can be attributed the effective back propagation algorithm adopted in the underlying optimization routine. For GP and SVM, there is a slight increase in retraining time. Nonetheless, the retraining time of RF increases linearly with an increase in the number of training points suggesting that it might not be suit-

able for sequential learning with large number of points. Note that there exists some advanced algorithm-specific techniques to avoid retraining from scratch, however, none of the machine-learning methods adopts such technique in our study. Furthermore, in all these algorithms, the time to predict an unseen input  $x^*$  is negligible (in the order of milli to micro seconds) because they belong to a class of eager learning algorithms as opposed to lazy learning algorithms where a model is built only when  $x^*$  needs to be predicted.



**Fig. 4.** Time taken by various machine-learning methods in Algorithm 1 for retraining after each evaluation.

## 5 Related work

Compared with dynamic strategies [11,12,13,14,15,16], static load-balancing approaches have received relatively less attention in the literature. The problem of static load-balancing can be formulated as a graph-partitioning problem that belongs to a class of  $\mathcal{NP}$ -hard problem for which finding optimal solution is computationally hard. Many efficient algorithms are developed to tackle this problem in operations research community and are used for static load-balancing. These algorithms can be grouped into geometry-based algorithms, graph-based algorithms, and partitioning algorithms [17]. In [18], the authors carried out an experimental comparison of eleven static load-balancing algorithms for heterogeneous distributed computing systems. They showed the relatively simple Min-Min heuristic performs well in comparison to the other techniques such as simulated annealing and genetic algorithms, and tabu search. However, the state-of-the-art high-performing algorithms comprises hybrid algorithms, multilevel

**Table 2.** Mean percentage deviation from the optimal runtime averaged over 10 replications with the maximum budget of 10 evaluations

| Task count | NN    | RF            | GP            | SVM           | RS            |
|------------|-------|---------------|---------------|---------------|---------------|
| 80         | 0.000 | <i>12.668</i> | <i>15.032</i> | <i>18.089</i> | <i>20.246</i> |
| 128        | 0.000 | 3.269         | <i>7.620</i>  | <i>5.177</i>  | <i>12.846</i> |
| 160        | 0.000 | <i>12.649</i> | <i>8.050</i>  | <i>6.989</i>  | <i>8.563</i>  |
| 256        | 0.000 | 4.575         | <i>8.468</i>  | <i>7.340</i>  | <i>10.024</i> |
| 320        | 0.000 | <i>2.208</i>  | <i>8.818</i>  | <i>6.709</i>  | <i>13.105</i> |
| 376        | 6.005 | <b>3.186</b>  | 8.132         | 7.206         | 7.456         |
| 512        | 0.000 | <i>10.269</i> | <i>11.794</i> | <i>6.472</i>  | <i>9.090</i>  |
| 640        | 0.000 | 2.674         | <i>20.058</i> | <i>14.072</i> | <i>10.292</i> |
| 800        | 0.000 | <i>7.435</i>  | <i>5.996</i>  | <i>6.182</i>  | <i>8.770</i>  |
| 1024       | 0.000 | <i>6.645</i>  | <i>4.985</i>  | 5.966         | <i>13.241</i> |

Note: The value is typeset in *italics* (bold) when a variant is significantly worse (better) than NN according to a *t*-test with significance (alpha) level 0.05.

approaches, and parallel implementations of the above algorithms [17]. We refer the reader to [17,19] for a survey for static load balancing approaches. Recently, in [20], a genetic algorithm was adopted for tasks scheduling and load balancing in heterogeneous parallel multiprocessor system. Nonetheless, the domain-specific constraints of the CICE component make the search problem hard and prevents the straightforward adoption of heuristic search algorithms [21]. In order to handle these constraints effectively, the search algorithms need a sophisticated constraint-handling mechanism; consequently they lose generality and become problem-specific.

The idea of using machine learning in load-balancing has received considerable attention for dynamic strategies. Examples include neural network [22], decision tree [23], and reinforcement learning approaches [24]. However, to the best of our knowledge, the adoption of machine-learning approaches for application and architecture specific static load-balancing has not been investigated before. Finally, this is the first work on the use of machine learning approaches for analyzing the sensitivity of the load-balancing parameters.

## 6 Summary and Outlook

We developed a machine-learning-based approach for static load-balancing problem and applied it for load balancing the CICE component of the CESM running on BG/P. We deployed a machine-learning method as a diagnostic tool for analyzing the sensitivity of the load-balancing parameters on the runtime and provided an explanation for inadequacy of the analytical performance model. The main contribution of the paper is the development and empirical analysis of the machine-learning-based algorithm that allowed us to load balance the CICE component of the CESM on BG/P with significant savings in core-hour usage. Compared to the current practice of expert-knowledge-based enumeration

over feasible parameter configurations, we showed that the proposed algorithm requires 6x fewer evaluations to find the optimal load-balancing configurations.

A inherent limitation of our algorithm consists in the sequential evaluation of parameter configurations that will affect the wall clock time. To address this issue, we will develop unsupervised learning methods to partition the feasible set into a number of similar groups and learning those regions in parallel. To that end, we will investigate parallel machine-learning algorithms. Since the inefficient processor allocations of CICE component can affect overall scaling of the CESM, we will use the proposed approach and assess the overall performance of the CESM. Furthermore, two projects, Climate-Science Computational End Station Development and Attributing Changes in the Risk of Extreme Weather and Climate, granted computational time on ALCF's BG/P and Q supercomputers under the DOE INCITE program will directly benefit from this work. We are planning to investigate the effectiveness of the proposed algorithm for load-balancing various climate simulations in these projects.

## Acknowledgments

This work was supported by the U.S. Department of Energy, Office of Science, Advanced Scientific Computing Research, under Contract DE-AC02-06CH11357. An award of computer time was provided by the Innovative and Novel Computational Impact on Theory and Experiment (INCITE) program. This research used resources of the Argonne Leadership Computing Facility at Argonne National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under contract DE-AC02-06CH11357.

## References

1. B. Metz, O. Davidson, P. Bosch, R. Dave, and L. Meyer, "Contribution of working group III to the fourth assessment report of the Intergovernmental Panel on Climate Change," 2007.
2. "MINOTAUR: A toolkit for MINLP," [http://wiki.mcs.anl.gov/minotaur/index.php/Main\\_Page](http://wiki.mcs.anl.gov/minotaur/index.php/Main_Page).
3. 2013, <http://www.cesm.ucar.edu/events/ws.2012/Presentations/SEWG2/craig.pdf>.
4. C. M. Bishop *et al.*, *Pattern Recognition and Machine Learning*. Springer, New York, 2006, vol. 1.
5. L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
6. M. A. Hearst, S. Dumais, E. Osman, J. Platt, and B. Scholkopf, "Support vector machines," *Intelligent Systems and Their Applications, IEEE*, vol. 13, no. 4, pp. 18–28, 1998.
7. C. E. Rasmussen and C. K. Williams, "Gaussian processes for machine learning (adaptive computation and machine learning)," 2005.
8. S. Haykin, *Neural Networks: A Comprehensive Foundation*, 1st ed. Upper Saddle River, NJ, USA: Prentice Hall PTR, 1994.
9. E. J. Atkinson and T. M. Therneau, "An introduction to recursive partitioning using the rpart routines," *Rochester: Mayo Foundation*, 2000.

10. R Core Team, *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria, 2013. [Online]. Available: <http://www.r-project.org>
11. L. V. Kale and S. Krishnan, *CHARM++: a portable concurrent object oriented system based on C++*. ACM, 1993, vol. 28, no. 10.
12. K. Barker, A. Chernikov, N. Chrisochoides, and K. Pingali, "A load balancing framework for adaptive and asynchronous applications," *Parallel and Distributed Systems, IEEE Transactions on*, vol. 15, no. 2, pp. 183–192, 2004.
13. K. J. Barker and N. P. Chrisochoides, "An evaluation of a framework for the dynamic load balancing of highly adaptive and irregular parallel applications," in *Proceedings of the 2003 ACM/IEEE conference on Supercomputing*. ACM, 2003, p. 45.
14. C. Huang, G. Zheng, L. Kalé, and S. Kumar, "Performance evaluation of adaptive mpi," in *Proceedings of the eleventh ACM SIGPLAN symposium on Principles and practice of parallel programming*. ACM, 2006, pp. 12–21.
15. C. Boneti, R. Gioiosa, F. J. Cazorla, and M. Valero, "A dynamic scheduler for balancing hpc applications," in *Proceedings of the 2008 ACM/IEEE conference on Supercomputing*. IEEE Press, 2008, p. 41.
16. R. Sharma and P. Kanungo, "Dynamic load balancing algorithm for heterogeneous multi-core processors cluster," in *Communication Systems and Network Technologies (CSNT), 2014 Fourth International Conference on*. IEEE, 2014, pp. 288–292.
17. Y. Hu and R. Blake, "Load balancing for unstructured mesh applications," *Parallel and Distributed Computing Practices*, vol. 2, no. 3, pp. 117–148, 1999.
18. T. D. Braun *et al.*, "A comparison of eleven static heuristics for mapping a class of independent tasks onto heterogeneous distributed computing systems," *Journal of Parallel and Distributed Computing*, vol. 61, no. 6, pp. 810–837, 2001.
19. S. Ichikawa and S. Yamashita, "Static load balancing of parallel PDE solver for distributed computing environment," in *Proc. 13th Int'l Conf. Parallel and Distributed Computing Systems*, 2000, pp. 399–405.
20. M. Effatparvar and M. Garshasbi, "A genetic algorithm for static load balancing in parallel heterogeneous systems," *Procedia-Social and Behavioral Sciences*, vol. 129, pp. 358–364, 2014.
21. P. Balaprakash, S. M. Wild, and P. D. Hovland, "Can search algorithms save large-scale automatic performance tuning?" in *Int. Conf. on Computational Science*, 2011.
22. Y. Jia and J.-Z. Sun, "A load balance service based on probabilistic neural network," in *International Conference on Machine Learning and Cybernetics*, vol. 3. IEEE, 2003, pp. 1333–1336.
23. M. A. Dantas and A. R. Pinto, "A load balancing approach based on a genetic machine learning algorithm," in *19th International Symposium on High Performance Computing Systems and Applications (HPCS 2005)*. IEEE, 2005, pp. 124–130.
24. T. Helmy and S. Shahab, "Machine learning-based adaptive load balancing framework for distributed object computing," in *Advances in Grid and Pervasive Computing*. Springer, 2006, pp. 488–497.

The submitted manuscript has been created by the UChicago Argonne, LLC, Operator of Argonne National Laboratory (Argonne) under Contracts No. DE-AC02-06CH11357 and DE-FG02-05ER25694 with the U.S. Department of Energy. The U.S. Government retains for itself, and others acting on its behalf, a paid-up, nonexclusive, irrevocable worldwide license in said article to reproduce, prepare derivative works, distribute copies to the public, and perform publicly and display publicly, by or on behalf of the Government. The NCAR is sponsored by the National Science Foundation.