

Predicting Cross-Architecture Performance of Parallel Programs

Daniel Nichols[†], Alexander Movsesyan[†], Jae-Seung Yeom^{*}, Abhik Sarkar^{*}, Daniel Milroy^{*},
Tapasya Patki^{*}, Abhinav Bhatele[†]

[†]*Department of Computer Science, University of Maryland*

^{*}*Center for Applied Scientific Computing, Lawrence Livermore National Laboratory*

E-mail: [†]dnicho@umd.edu, [†]bhatele@cs.umd.edu

Abstract—A variety of hardware architectures, both CPUs and GPUs, are used today to build supercomputers and parallel clusters. Often times, users can choose which hardware platform they want to run on. Modern scientific workflows have multiple computational tasks, and each task may be better suited for a different architecture in terms of performance. Deciding where to run an application or workflow task is not straightforward because of the complexity of applications, and hardware architectures, which makes performance predictions challenging. Hence, modeling the performance of scientific applications across a variety of architectures is important for achieving the best performance. In this paper, we present a machine learning based methodology to model the relative performance of applications across multiple architectures using hardware performance counters. Our machine learning model can predict the relative performance of an application with a mean absolute error of 0.11, and can be used effectively to make performance-aware and multi-architecture scheduling decisions, reducing makespan by up to 20%.

Index Terms—performance modeling, architectures, machine learning, multi-cluster scheduling

I. MOTIVATION

An increasing number of scientific workloads are being expressed as workflows with sets of computational tasks and dependencies between them [1], [2]. These workflows typically involve ensembles of tasks (jobs) in a pipeline that run different codes such as simulations, uncertainty quantification analysis, and machine learning training. As applications become more portable due to the emergence of portable programming models [3], package managers [4], and containerization techniques, different tasks or jobs might be better suited for different hardware architectures. Given these portable workflows and the increasingly heterogeneous set of computing resources available to end users today, it is important to develop capabilities to efficiently place these tasks on the most efficient resources available.

Different tasks or applications in a workflow can be assigned to different architectures if users have access to a variety of compute nodes via a multi-resource job scheduler, which is becoming increasingly common, both in data centers and HPC facilities. As a result, the demand for such multi-resource schedulers [5] is emerging. In an ideal setting, scheduler can automatically decide the most suitable architecture for different jobs in terms of performance. This can remove the

user from the decision making process and let a system scheduler decide what hardware to run an application on. However, in practice, this requires being able to predict the performance of incoming jobs across diverse architectures. This is a complex problem that would involve developing models for understanding the performance of scientific applications across diverse architectures.

Cross-architecture performance modeling is a challenging problem because application execution times are dependent on several factors with non-trivial relationships to performance. The performance depends on how well the application’s behavior aligns with the properties of the hardware it is running on. These hardware properties, such as peak flop/s, memory bandwidth, and cache sizes are easy to obtain, however, the behavior of the application is non-trivial to model. Application performance can depend on a number of characteristics such as arithmetic intensity, memory loads/stores, branching behavior, I/O, and many more. Characterizing these and using them to model performance on a diverse set of architectures is challenging due to the number of contributing factors and complexity of the relationship.

In this paper, we propose a solution to the cross-architecture performance modeling task by training a machine learning model to predict the relative performance of an application across a set of architectures given performance counters of the application from one architecture. In order to accomplish this, we collect a data set of application runs from four different HPC systems with different architectures and measure a hand selected set of performance counters. These counters, along with the recorded execution times, are used to train a regression model to predict relative performance vectors. Additionally, we demonstrate the generalizability of our model by evaluating it on a set of applications it has not seen before. Finally, we demonstrate the makespan improvement from using this model in a multi-resource scheduling simulation.

In this paper we make the following contributions:

- The Multi-Platform HPC (MP-HPC) dataset of hardware performance counters for a wide variety of scientific applications recorded on four different HPC systems.
- A regression model that can predict the relative performance of an application across multiple systems with a mean absolute error of 0.11.

- A qualitative comparison of the importance of different counters in cross-architecture performance modeling.
- A demonstration of the potential makespan improvement if our model is used to assist scheduling decisions in a multi-resource scheduler.

II. BACKGROUND

In this section, we provide background on performance profiling, relative performance vectors, and regression modeling.

A. Performance Profiling

When studying performance related aspects of an application, performance profiling tools are often used to collect data. These tools record profiling metrics such as wall time during an application execution, and often attribute those metric values to different regions of the application’s code.

A popular performance profiling tool for parallel programs is HPCToolkit [6]. It is a sampling-based tool that can collect numerous operating system and hardware counters and attribute them to nodes on a calling context tree. It can record counters such as cache misses, floating point operations, branch instructions, etc. While many tools can record this type of data, HPCToolkit has been demonstrated to be more accurate than the others with relatively low overheads [7].

Typically, this analysis is done through HPCToolkit’s graphical interface, hpcviewer, making studying trends in large numbers of profiles very difficult [8]. The Python library, Hatchet [9], solves this problem by providing a programmatic interface to the profiles produced by HPCToolkit and other popular profilers. Additionally, it provides extensive functionality for calling context tree pruning and analysis through pandas Dataframe operations.

B. Regression Modeling

Traditional machine learning (ML) is often tasked as learning to predict some output given an input to the machine learning model. When the output is discrete it is called *classification*. On the other hand, when the output is continuous it is called *regression*. The latter of these training objectives is used in the modeling in this paper.

When training a regression model, it is necessary to have a dataset, \mathcal{D} , of existing data where outputs are known. The amount of data, $|\mathcal{D}|$, needed is dependent on the model, features, and problem complexity, however, typical regression tasks can require thousands to tens of thousands of training samples. This data, along with its corresponding outputs, is used to optimize the model’s predictions respective to some learning objective. Common learning objectives in regression are to minimize *mean absolute error*, *mean squared error*, coefficient of determination (R^2), etc. on a testing data set. The testing data set is separate from the data that the model was trained on, so that reported values do not include overfitting, i.e. the model does not memorize the data set. Often times the function being minimized is additively combined with a regularization term to reduce model complexity.

$$\theta^* = \min_{\theta} \mathcal{L}(x; \theta) + \Omega(\theta)$$

Here $\mathcal{L}(x; \theta)$ is the loss function at data sample $x \in \mathcal{D}$ parameterized by θ . It is intended to model predictive capacity of the model such as with mean absolute error. The second term, $\Omega(\theta)$, is the regularization term, which models the complexity of the machine learning model. Penalizing model complexity helps prevent overfitting of the data set.

III. RELATED WORK

Below, we present related work in the areas of machine learning based performance modeling and cross-architecture performance modeling. We further discuss how our work differs from and builds on top of existing work.

A. Performance Modeling using Machine Learning

Performance modeling is a well studied research area with lots of literature surrounding analytical and statistical models. Recently, with the increase in machine learning innovations, there has been a large focus on the latter. Machine learning can help model complex relationships between applications and their final performance. It has been used to model job runtimes [10]–[13], variability [14]–[16], power consumption [17], and many other things [18]–[20].

These models are often used to study and understand complex relationships between applications and their performance. Malakar et al. [18] compare the capability of various different machine learning methods on modeling performance. Furthermore, Zhou et al. [21] demonstrate how to extrapolate models from small scale runs to larger scale runs. Many works also use these models in downstream tasks to improve performance. In [22] standard machine learning techniques such as k -Nearest-Neighbors and XBoost are used to model MPI collective performance and inform auto-tuning decisions. This fits into the broader study of using machine learning models to more efficiently explore the combinatorial search space in auto-tuning [12], [23], [24]. Similar to this paper, there are other works that use machine learning models to make informed scheduling decisions on HPC systems such as to reduce variability [15] or avoid IO bottlenecks [25].

B. Cross-Architecture Performance Modeling

Two works by Ardalani et al. [26], [27] focus specifically on cross-architecture performance modeling. More specifically, they consider predicting performance across architectures. Like our work these use expert derived counters to model the computational behaviour of an application. However, they only focus on mapping sequential C code performance to GPU performance. They do not look at a multiple architectures or a wide variety of applications and they only consider single functions rather than entire application binaries. Another similar paper by Yang et al. [28] introduces a model for predicting performance of parallel applications between two architectures using cumulative averages and a filter model. This work differs from ours in that it requires running the application on both architectures to make predictions and it only considers CPU architectures. These lines of work, [26], [27] and [28], do not explore any potential uses of their models

in applications such as multi-resource scheduling. Some works have used heuristics and machine learning models to do resource placements for tasks in workflows [29], [30]. These have used test runs and search space pruning to find optimal resource sets. However, none of them use cross-architecture predictive models to inform their resource selection.

IV. OVERVIEW OF OUR METHODOLOGY

We first provide an overview of our methodology to predict the relative performance of an application across a set of architectures given performance counters of the application from one architecture. This includes two things – the data collection phase and the model training phase (Figure 1). In the first phase, we collect performance profiles for a variety of applications running on N different HPC systems with different architectures and record a hand-selected set of performance counters. These counters, along with the recorded execution times, are used to train a regression model to predict relative performance in the second phase.

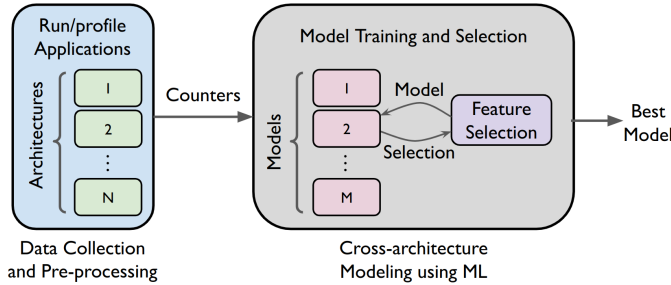


Fig. 1. Overview of data collection and machine learning pipeline. Applications are profiled on several architectures and performance counters are collected for training the model. Model and feature selection are done iteratively until the best set is selected.

Since our goal is to predict performance on other architectures relative to a baseline on one architecture, we introduce the term *Relative Performance Vector* (RPV) that encodes the relative performance of an application across several architectures. To define RPV, let us consider a set of applications A , corresponding input problems I_A , and systems S . For a particular application and input problem pair $(a, i) \in A \times I_A$ executed on N systems in S we can define the *Relative Performance Vector* as $\text{rpv} : (A, I_A) \times S \mapsto \mathbb{R}^N$ such that $\text{rpv}(a, i, s)$ is the vector of the performance of (a, i) across all platforms relative to that on system s . Here we assume that (a, i) can run on all the systems in S . For example, consider running an application-input pair (TestApp, “-s 5”) on systems X , Y , and Z . If the application runs in ten minutes on system X , eight minutes on system Y , and 21 minutes on system Z , then the performance vector relative to X would be:

$$\text{rpv}(\text{TestApp}, \text{“-s 5”}, X) = \begin{bmatrix} 1.0 \\ 0.8 \\ 2.1 \end{bmatrix}$$

Application being run Time on X relative to X Time on Y relative to X
System run on Time on Z relative to X
Input Arguments

We also define $\text{rpv}(\cdot, \cdot, \min)$ and $\text{rpv}(\cdot, \cdot, \max)$ as the performance vectors relative to the systems where lowest and highest performance is obtained, respectively. The RPV provides a concise, mathematical representation for relative performance across systems that can be used in our further downstream modeling tasks.

In order to model the mapping $\text{rpv} : (A, I_A) \times S \mapsto \mathbb{R}^N$, we need a large number of input and output data to train on. This requires a large number of samples in the $(A, I_A) \times S$ space. To collect these, we profile a variety of applications at several of their inputs on several architectures. These runs provide hardware counters that may provide insight into an application’s behavior for many application, input, and architecture tuples.

We use the counters collected during profiling to form the MP-HPC dataset and, in turn, use this dataset for the machine learning (ML) component (second phase). The ML component uses the profiled counters from a particular architecture to predict the relative performance vector across a set of systems. We try different ML models and feature sets to identify the best performing model. This model is exported and used in downstream relative performance prediction tasks such as cross-architecture scheduling.

V. DATA COLLECTION AND PRE-PROCESSING

In this section, we provide details of how we generated the MP-HPC dataset used for our modeling problem. We describe the process of running and profiling the applications, and collecting the performance metrics.

A. Scientific Applications

In order to model the relative performance of applications run on an HPC machine, we need to collect performance data from applications that are typically run on these machines. We accomplish this by running a set of applications, benchmarks, and proxy applications from the ECP Proxy Applications Suite [31] and E4S Test Suite [32]. These are chosen because they are designed to be representative of actual workloads on HPC systems, but are simpler to build and run than full scientific applications.

Table II lists the applications used in our data set. There are twenty applications in total, and eleven of them have GPU support. The GPU support comes from a variety of libraries such as OpenMP, Kokkos [33], RAJA [3], and native CUDA or HIP. Each application is paired with different input configurations when run, in order to test different problems and problem sizes. We build and install all of the applications with their default build settings provided in their respective Spack [4] packages.

B. Architecture Descriptions

We run each application-input pair on four different machines with different architectures. These are listed in Table I. There are two Intel Xeon based, CPU-only machines and two GPU-based machines. The first GPU machine uses IBM

TABLE I

OVERVIEW OF THE FOUR ARCHITECTURES WE COLLECT PERFORMANCE DATA ON. THERE ARE TWO CPU ONLY SYSTEMS AND TWO CPU+GPU SYSTEMS. THE CPUS SPAN THREE VENDORS: INTEL, IBM, AND AMD, WHILE THE GPUS ORIGINATE FROM TWO VENDORS: NVIDIA AND AMD.

System	CPU Type	CPU cores/node	CPU Clock Rate (GHz)	GPU Type	GPUs/node
<i>Quartz</i>	Intel Xeon E5-2695 v4	36	2.1	—	—
<i>Ruby</i>	Intel Xeon CLX-8276	56	2.2	—	—
<i>Lassen</i>	IBM Power9	44	3.5	NVIDIA V100	4
<i>Corona</i>	AMD Rome	48	2.8	AMD MI50	8

TABLE II

THE APPLICATIONS USED IN OUR STUDY ALONG WITH A BRIEF DESCRIPTION OF WHAT EACH APPLICATION DOES AND WHETHER IT SUPPORTS RUNNING ON A GPU.

Application	Description	GPU
AMG	Algebraic multigrid solver	✓
CANDLE	Deep learning models for cancer studies	✓
CoMD	Molecular dynamics and materials science algorithms	
CosmoFlow	3D convolutional neural network for astrological studies	✓
CRADL	Multiphysics and ALE hydrodynamics	✓
Ember	Communication patterns	
ExaMiniMD	Molecular dynamics simulations	✓
Laghos	FEM for compressible gas dynamics	✓
miniFE	Unstructured implicit FEM codes	✓
miniGAN	Generative Adversarial Neural Network training	✓
miniQMC	Real space quantum Monte Carlo algorithms	✓
miniTri	Triangle based data analytics algorithms	
miniVite	Graph community detection	
DeepCam	Climate segmentation benchmark	✓
Nekbone	High-order, incompressible Navier-Stokes solver	
PICSARlite	Particle-in-Cell simulation	
SW4lite	Seismic wave simulation	✓
SWFFT	Distributed-memory parallel 3D FFT	
Thornado-mini	Radiative transfer solver in multi-group, two-moment estimations	
XSBBench	Monte Carlo neutronics simulations	

Power9 CPUs and NVIDIA V100 GPUs, while the second uses AMD Rome CPUs and AMD MI50 GPUs.

On each of these systems the applications are run in three configurations – on one core, on one node using all the cores, and on two nodes. The one-core runs use one GPU if applicable. MPI is used for the one and two node runs to make use of all the cores and GPUs on the node. Some applications only support run configurations with square or power of two MPI processes and are, thus, run on the nearest number of ranks possible to one or two nodes. If an application does not support running on a GPU, we run it on the CPU only and use comparable CPU counters. If an application does support running on a GPU, then only GPU counters are collected. During these runs, HPCToolkit [6] (with CUPTI [34] on

NVIDIA GPUs or rocProfiler [35] on AMD) is used to record the application counters, and after the application run is complete, Hatchet [9] is used to parse these counters from the HPCToolkit output. For multi-process and multi-GPU runs, we record the mean value of the counters across all processes. The final results from all runs are then collected into a Pandas dataframe for use in the later tasks.

C. Details of Recorded Hardware Counters

To understand the varied computational characteristics of different applications in Table II, we record several hardware counters during the application runs. Table III lists the counters recorded on each architecture in our data set. Counter names are not consistent across different architectures and they may also represent slightly different data. However, we have tried to identify similar counters that model the same underlying performance characteristics that affect final performance. Most of these counters fit into one of three categories: control flow, data intensity, or I/O. These categories capture the main performance characteristics of applications across different architectures. Broadly speaking, applications with more complex control flow will fair better on CPUs, which are geared towards latency. On the other hand, applications with more data intensity generally benefit on throughput-geared GPUs.

D. Preparing the Final Dataset

Using the counters listed on the right of Table III we compute a set of derived values as the final features in the data set. These features are detailed on the left of Table III. The instruction related counters branch, store, load, single FP, double FP, and integer arithmetic are all computed to be ratios of the total number of instructions (note that the feature *arithmetic intensity* refers to the ratio of arithmetic instructions, not the conventional flop-to-bandwidth ratio). This normalizes the values across runs, which may have drastically different numbers of total instructions. The remaining eight features are normalized by subtracting that feature’s mean to center its values and dividing them by its standard deviation. We additionally include whether the run was from a GPU or not, how many nodes, and how many cores the run used. The architecture feature is a one-hot-encoded vector encoding what architecture the counters were collected on. In the context of this paper, that is four separate features that are used to denote whether the run is from Quartz, Ruby, Corona, or Lassen.

The final MP-HPC dataset has 21 columns and 11,312 rows. Each row represents a run of an application-input pair for

TABLE III
LIST OF FINAL FEATURES IN THE COLLECTED DATA SET AND THE SOURCE COUNTERS/VALUES THEY ARE DERIVED FROM. WE COMBINE DERIVED VALUES FROM THE RECORDED COUNTERS AND META-DATA ABOUT THE RUN CONFIGURATION.

Feature	Description	Source Counters & Values			
		Quartz	Ruby	Lassen	Corona
Branch Intensity	Ratio of branch instructions to total instructions	PAPI_BR_INS	PAPI_BR_INS	cf_executed	–
Store Intensity	Ratio of store instructions to total instructions	PAPI_SR_INS	PAPI_SR_INS	inst_executed_local_stores, inst_executed_global_stores	LDSInsts, GDSInsts
Load Intensity	Ratio of load instructions to total instructions	PAPI_LD_INS	PAPI_LD_INS	inst_executed_local_loads, inst_executed_global_loads	LDSInsts, GDSInsts
Single FP Intensity	Ratio of single precision FP instructions to total instructions	PAPI_SP_OPS	PAPI_SP_OPS	flop_count_dp	VALUInsts, SALUInsts
Double FP Intensity	Ratio of double precision FP instructions to total instructions	PAPI_DP_OPS	PAPI_DP_OPS	flop_count_sp	VALUInsts, SALUInsts
Arithmetic Intensity	Ratio of integer arithmetic instructions to total instructions	bdw_ep::ARITH	clx::ARITH	inst_integer	–
L1 Load Misses	L1 cache load misses	PAPI_L1_LDM	PAPI_L1_LDM	local_load_requests, local_hit_rate	–
L1 Store Misses	L1 cache store misses	PAPI_L1_STM	PAPI_L1_STM	local_store_requests, local_hit_rate	–
L2 Load Misses	L2 cache load misses	PAPI_L2_LDM	PAPI_L2_LDM	gld_efficiency	TCC_MISS_sum, TCC_EA_RDREQ
L2 Store Misses	L2 cache store misses	PAPI_L2_STM	PAPI_L2_STM	gst_efficiency	TCC_MISS_sum, TCC_EA_WRREQ
IO Bytes Written	Bytes written to IO	IO	IO	IO	IO
IO Bytes Read	Bytes read from IO	IO	IO	IO	IO
Extended Page Table	Extended page table size	EPT	EPT	EPT	EPT
Memory Stalls	Memory stalls	PAPI_MEM_SCY	PAPI_MEM_SCY	GINST:STL_ANY	MemUnitStalled
Nodes	Number of nodes	Run Configuration	Run Configuration	Run Configuration	Run Configuration
Cores	Number of cores	Run Configuration	Run Configuration	Run Configuration	Run Configuration
Uses GPU	1 if counters from GPU; 0 otherwise	0	0	1 if app uses GPU	1 if app uses GPU
Architecture	one-hot-encoded vector for what architecture these counters were recorded on	(1 0 0 0)	(0 1 0 0)	(0 0 1 0)	(0 0 0 1)

a specific number of MPI processes on a single architecture. The columns are derived from the counters collected during the run and meta-data about the run (see Table III).

VI. MACHINE LEARNING BASED MODELING

Next, we present our methodology for training the machine learning models, evaluating their performance, and identifying the best models and features.

A. Training the XGBoost Regression Model

Now that we have a data set of counters from applications and the corresponding relative performance vectors across a set of architectures, we want to use machine learning to predict the relative performance vectors given counters from one of the architectures. In order to learn how to predict relative performance vectors, we use the XGBoost (eXtreme Gradient Boosting) regression model [36]. This model is an ensemble of decision trees that are additively combined to make final predictions. If $\hat{y}_i \in \mathbb{R}$ is the predicted regression value of the model, then it can be computed as,

$$\hat{y}_i = \sum_{k=1}^K f_k(x_i) \quad \text{where } f_k \in \mathcal{F}$$

number of trees $\rightarrow K$ regression tree k space of regression trees $\rightarrow \mathcal{F}$

predicted value $\rightarrow \hat{y}_i$

As described in Section II-B, we can add a regularized objective function to avoid over-fitting:

$$\mathcal{L}(\hat{y}_i) = \sum_{i=1} \underbrace{l(\hat{y}_i, y_i)}_{\text{training loss}} + \sum_k \underbrace{\Omega(f_k)}_{\text{convex loss function}} \quad (1)$$

predicted value $\rightarrow \hat{y}_i$ complexity of tree f_k $\rightarrow \Omega(f_k)$

Since this is parameterized by functions ($f_k \in \mathcal{F}$), it cannot be optimized using typical optimization methods. Thus, *gradient tree boosting* greedily adds in the best functions throughout training iterations by selecting the f_t that minimizes Equation 1 the most. These f_t can be additively combined into a new loss function as,

$$\mathcal{L}^{(t)} = \sum_{i=1} l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t)$$

training iteration $\rightarrow t$ tree that minimizes $\mathcal{L}^{(t-1)}$ $\rightarrow f_t$

This can be optimized using second-order approximations and standard convex minimization methods. XGBoost implements this gradient tree boosting method alongside a number of state-of-the-art techniques for tree splitting and pruning. Additionally, it provides efficient implementations that can scale to large numbers of data samples and run on GPUs.

It is a state-of-the-art machine learning algorithm for learning on tabular data.

In order to train an XGBoost regressor, we use its publicly available Python library at version 1.7.1 [36]. We train the model on a CPU on the Ruby system. Training the XGBoost model takes on the order of tens of seconds on average. The model takes the features from Section V-D and predicts the relative performance across all four architectures as a vector. Mean absolute error (MAE) is used as the minimization objective during training. During this training, 10% of the data is set aside as a testing data set, while the other 90% is shown to the model as a training data set. While training on the training data set, the data is further split into five folds as part of k-fold cross-validation. The model is trained on four out of the five folds at a time, while the other is used as validation. This is done for all five combinations and the average MAE is reported.

We additionally train several other common machine learning regressors to compare the quality of the XGBoost model to other state-of-the-art methods. For this, we include linear regression and decision forests. These are implemented from the scikit-learn Python library [37]. As with XGBoost these are trained with a 90-10 train-test split and 5-fold cross-validation. We also test against *mean* prediction as a baseline for the ML models. This regressor guesses the mean RPV in the training set for all samples in the test set.

B. Model and Feature Selection

To select the best model and feature set, we first train all the models on all the features. After training we select the best set of features using those reported by XGBoost and the decision forest, since these models expose feature importances. These features are then used to re-train all the models again.

In order to measure the feature importances of the trained model, we use XGBoost to easily recover importance values. XGBoost, in its Python framework, computes feature importances during training and exposes them in its model interface. It calculates them based on the average gain across all decision splits in the trees. During training a tree will add splits on a feature to improve its predictive performance. The improvement in performance from this split is called the gain. When there are multiple regression targets the gain is averaged over each output.

For any given feature if we average the gain from all the splits on that particular feature in a tree, then we can compute the importance of the feature for that tree. This includes all the splits in XGBoost's sets of trees. Finally, we can compute this value for all of the features in the data set to retrieve a feature importance vector.

With this method of calculating feature importances we can expose the relative contribution of each feature to the model's predictions. A higher importance indicates that that feature contributes more to the models performance than other lower scored features. Decision tree feature importances can also be calculated based on the frequency and coverage of splits for a feature, however, these can be biased towards features with

a large number of unique values and numeric features. Both of these are present in our data set, so we elect to use the average gain.

Since the data set has a relatively small number of features, the feature selection will likely have negligible impact on model training time. However, discovering the most impactful features gives insight into what is most necessary in predicting cross-architecture performance. Additionally, it allows us to collect less features in future implementations of this methodology. This is a considerable optimization as data collection is the most time and resource intensive portion of our machine learning pipeline.

C. Evaluation Metrics

We evaluate the model's performance using two different metrics: *Mean Absolute Error* (MAE) and *Same Order Score* (SOS). The MAE encodes the average magnitude of error in the relative performance predictions. This measure provides a value that is easy to reason about regarding predictive performance. An MAE of 0.1 means that the model predicts the relative performance of applications within ± 0.1 on average across each vector.

$$\text{MAE} = \frac{1}{|\mathcal{D}_{\text{rpv}}|} \sum_{i=1}^{|\mathcal{D}_{\text{rpv}}|} \|\text{rpv}_i - \widehat{\text{rpv}}_i\|_1$$

data set of relative performance vectors (points to \mathcal{D}_{rpv})
predicted rpv for run i (points to $\widehat{\text{rpv}}_i$)
rpv for run i (points to rpv_i)
number of architectures (points to d)

The SOS metric denotes the fraction of samples where the model predicts the relative performance vector in the correct order. We define two vectors \mathbf{a} and \mathbf{b} as being in the same order if the i -th elements \mathbf{a}_i and \mathbf{b}_i are both the n -th largest in their respective vector, for all i . The SOS is then defined as the fraction of predicted relative performance vectors that are in the same order as their respective true relative performance vector. This is shown in the equation below with the indicator function being used to count the relative performance vectors where the ordering is preserved.

$$\text{SOS} = \frac{1}{|\mathcal{D}_{\text{rpv}}|} \sum_{i=1}^{|\mathcal{D}_{\text{rpv}}|} \mathbb{1}\{\text{ranks}(\text{rpv}_i) = \text{ranks}(\widehat{\text{rpv}}_i)\}$$

ordering of relative performance vectors (points to $\text{ranks}(\text{rpv}_i)$ and $\text{ranks}(\widehat{\text{rpv}}_i)$)

This metric shows how well the model understands the ordering of performance on different architectures, but ignores the magnitude of its predictions. Thus, the SOS combined with MAE gives reasonable insight into how well the model is predicting relative performance vectors. Both of these metrics are computed over the testing set for data samples that the model has not seen before.

VII. SCHEDULING EXPERIMENT

Once a model is trained to predict relative performance vectors, it can be used to make informed cross-architecture scheduling decisions. We test this capability in our trained model by simulating a multi-resource scheduling environment.

We create a workload of 50,000 jobs randomly sampled from our existing data set with replacement. These are scheduled using the First-Come-First-Serve with EASY backfilling scheduling algorithm (FCFS+EASY) [38] presented in Algorithm 1. This algorithm uses the `Machine` function to assign a job to one of these machines: Quartz, Ruby, Lassen, or Corona. If the machine cannot satisfy the resource requirement of a job (the number of nodes it needs), then the job is reserved at the earliest possible time or backfilled, otherwise, it is run immediately. The function `Start(j, m)` represents running job j on machine m . We use the observed run times on each machine from the data set to determine how long the job would run for simulation purposes.

Algorithm 1 Multi-resource Scheduling Algorithm using FCFS+EASY. This standard algorithm queues jobs using policy \mathcal{R}_1 (FCFS in our case) and the EASY backfilling algorithm parameterized by the policy \mathcal{R}_2 (FCFS in our case). The function `Machine` is used to assign jobs to machines. The symbol \setminus represents the set minus operation.

Input: $Q \leftarrow$ queue of jobs
 $\mathcal{R}_1 \leftarrow$ Queue ordering policy
 $\mathcal{R}_2 \leftarrow$ Backfill ordering policy
 $M \leftarrow$ Set of machines used for multi-resource scheduling
`Machine(j, i, M)` \leftarrow Function that assigns jobs to machines

```

1:  $i \leftarrow 0$ 
2: sort  $Q$  according to  $\mathcal{R}_1$ 
3: for job  $j \in Q$  do
4:   if  $j$  can start now then
5:     pop  $j$  from  $Q$ 
6:     Start( $j, \text{Machine}(j, i, M)$ )
7:      $i \leftarrow i + 1$ 
8:   else
9:     Reserve  $j$  at earliest possible time
10:     $L \leftarrow Q \setminus \{j\}$ 
11:    sort  $L$  according to  $\mathcal{R}_2$ 
12:    for job  $j' \in L$  do
13:      if  $j'$  can start now without delaying  $j$  then
14:        pop  $j'$  from  $L$  and  $Q$ 
15:        Start( $j', \text{Machine}(j', i, M)$ )
16:         $i \leftarrow i + 1$ 

```

We run this scheduling simulation with four different implementations of the `Machine` function that represent different assignment strategies: Round-Robin, Random, User+RR, and Model-based. These different strategies expose the common interface for scheduling, `Machine(j, i, M)`, where j is the job to schedule, i is the index of j in the queue, and M is the set of machines considered for multi-resource scheduling. Depending on the algorithm, some of these arguments are not used. The Round-Robin placement places jobs on machines in a round-robin fashion rotating between machines for each consecutive job. The Random placement uniformly selects a random machine among the four to run on. The

User+RR placement mimics typical user behavior by running GPU-enabled applications on GPU systems and CPU-only applications on CPU-only systems. A round-robin scheme is used to decide which GPU system to use for GPU-enabled applications, and likewise for CPU-only applications.

Finally, Algorithm 2 shows the Model-based assignment strategy, which uses an ML-based model to pick the fastest machine for each job and run it there. If the machine cannot satisfy the resource requirement of the job, then the strategy picks the next fastest and so on. We implement this scheduling simulation in Python using our data set to get run time information for jobs. The nodes available on each machine reflect the number available on the actual machines. This is not meant to substitute rigorous scheduling simulation studies but only to demonstrate how such ML models can be used in a production multi-resource job scheduler.

Algorithm 2 Model-based strategy to decide which machine to assign a job to based on the predicted relative performance.

Input: $j \leftarrow$ Job to schedule
 $i \leftarrow$ Index of j in queue
 $M \leftarrow$ Set of machines used for multi-resource scheduling

```

1: function MachineModel-based( $j, i, M$ )
2:    $rpv \leftarrow \text{Model}(j)$ 
3:    $m \leftarrow \text{argmax}_{s \in M} rpv$ 
4:   if all  $s \in M$  are full then
5:     return  $m$ 
6:   else
7:      $M' \leftarrow M$ 
8:     while  $m$  is full do
9:        $M' \leftarrow M' \setminus \{m\}$ 
10:       $m \leftarrow \text{argmax}_{s \in M'} rpv$ 
11:   return  $m$ 

```

A. Evaluation Metrics

When evaluating the efficiency of our scheduling algorithm we are concerned with performance from the perspective of individual jobs as well as the scheduler as a whole. Users will hope to see a faster turnaround time from job submission to completion for their jobs, while system administrators may look at the job throughput of a given scheduler to measure its performance. To quantify both of these metrics, we use *average bounded slowdown* and *makespan*.

The average bounded slowdown represents the average slowdown of a set of jobs with a fixed bound to prevent overpenalizing extremely short jobs. Slowdown is the ratio of the sum of execution time and wait time to the execution time (not including the wait time). This provides a per-job evaluation metric to see how much each job is affected by the scheduling algorithm. The bounded slowdown can be calculated as shown below:

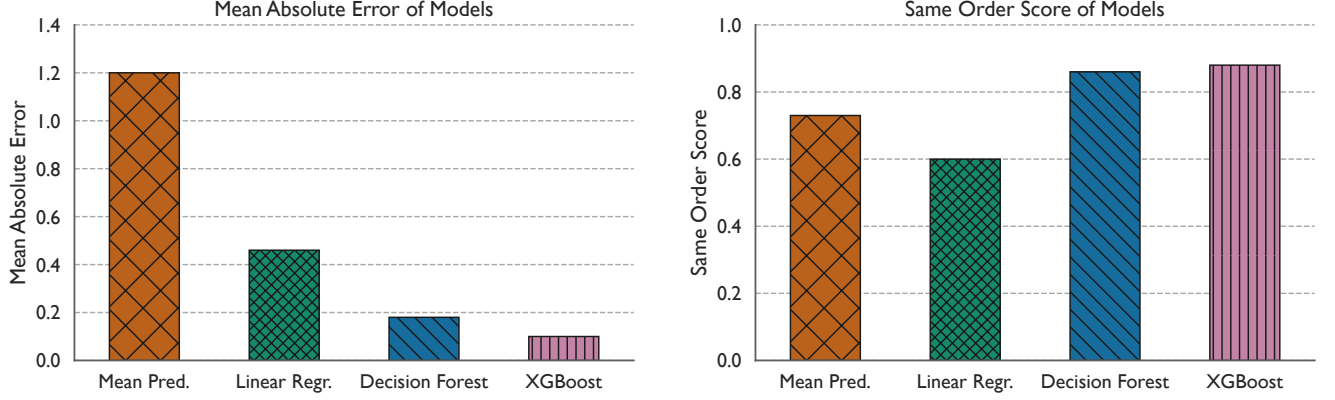


Fig. 2. The MAE (left) and SOS (right) of each machine learning model over the testing data set after training (Lower MAE is better, and higher SOS is better.) XGBoost outperforms the other models with an MAE of 0.11.

$$\text{BoundedSlowdown}(j) = \max \left(\frac{\text{wait time of } j + \text{run time of } j}{\max(p_j, \tau)}, 1 \right)$$

batch job → j (blue box)
wait time of j → w_j (green box)
run time of j → p_j (red box)
small time interval to prevent overpenalizing short jobs → τ (purple box)

We use $\tau = 10$ in our evaluation. Using the equation above, we can compute the average bounded slowdown over a set of jobs J as,

$$\overline{\text{BoundedSlowdown}}(J) = \frac{1}{|J|} \sum_{j \in J} \text{BoundedSlowdown}(j)$$

set of jobs → J (purple box)

Using the set of jobs J , we can also define the makespan as the time from the first job submission to the time when the last job finishes. This measures the amount of time it takes for a scheduler to complete executing all the jobs in a workload, and is commonly used to compare different machine assignment algorithms over fixed workloads.

$$\text{makespan}(J) = \left(\max_{j \in J} (w_j + r_j + p_j) \right) - \left(\min_{j \in J} r_j \right)$$

wait time of j → w_j (green box)
start time of j → r_j (orange box)
duration of j → p_j (red box)
set of jobs → J (purple box)

Both of these metrics are computed for each machine assignment algorithm across our workload. We compare them between all the machine assignment algorithms to observe the benefit from cross-architecture performance modeling. For each metric, a lower value indicates better performance.

VIII. RESULTS

In this section, we present the results from our training of regression models and scheduling experiments.

A. Evaluation of ML Models

Figure 2 shows the mean absolute error and same-order-score of each model on the testing data set. We observe that XGBoost performs the best for both of these metrics.

The XGBoost model scores a MAE of 0.11 (see Figure 2, left). This signifies that the model can use counter values for a job recorded on one architecture and predict its relative performance on other architectures within 0.11 on average. This is an 81.6% improvement over guessing the mean relative performance vector from the data. From this we can infer that the model is not simply guessing according to the distribution of the runtime data, but is rather correlating counter data with its performance prediction.

The linear and decision forest models perform better than guessing the mean, but obtain worse MAE values than XGBoost. Decision forest scores the closest to XGBoost likely since they are both ensembles of decision trees. However, XGBoost implements boosting alongside a number of other pruning techniques that strengthen its prediction.

We observe similar performance from XGBoost on the SOS metric where it is the best model (see Figure 2, right). It is able to predict the relative performance vector in the correct architecture order in 86% of samples in the testing set. This means that XGBoost is able to predict the fastest and slowest architectures for a particular application and input in a large number of scenarios, which is a valuable result to a user who is likely trying to avoid the slowest architecture and run on the fastest. Additionally, if the system with the fastest architecture is busy, then the user can select the next fastest and so on. As with the MAE metric, the decision forest has similar, but worse performance than XGBoost. Unlike with MAE, the linear model performs the worst on the SOS metric. This suggests that using the average machine order is a better predictor than the linear model.

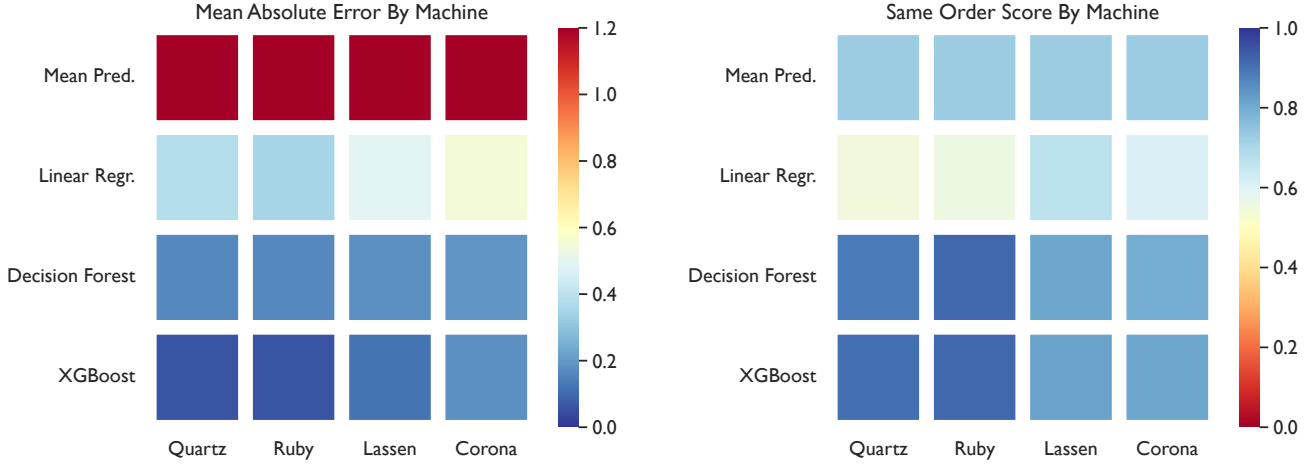


Fig. 3. The MAE (left) and SOS (right) of each model when predicting using counters data from one particular machine. For instance, the bottom right of the left visualization represents the MAE when predicting relative performance vectors with XGBoost and counters data from Corona. Lower MAE, or more blue, is better, and higher SOS, or more blue, is better.

B. Ablation Study

Here we study the effects on modeling performance when removing certain features and/or data from the training set. Figure 3 further details how well the models perform when the counters for only one architecture are used for modeling performance. In both visualizations in Figure 3, the “mean” prediction row is constant, since the mean relative performance vector is independent of the input features. Figure 3 (left) shows the MAE scores for each ML model. Each cell in the heatmap shows the MAE for one ML model and using the counters from one specific architecture. For instance, the bottom right of the left visualization represents the MAE when predicting relative performance vectors with XGBoost and counters data from Corona. We observe the same trends as Figure 2 where XGBoost has the best MAE. However, we notice that counters data from Ruby lead to a lower MAE and, thus, better predicted relative performance vectors. In fact, using counters from the two CPU systems, Ruby and Quartz, generally leads to better MAE. This same trend continues for the SOS metric in Figure 3 (right).

The fact that counters recorded on CPU machines lead to better predictions on average is an important observation for using this model in practice. CPU machines are generally less expensive and more readily available. Users can run their code on them and get predictions from the model for less available or more expensive resources, such as GPUs. Additionally, users can obtain an estimate of the speedup from running on a given architecture without actually having access to or being capable of running on that architecture. For instance, if a particular application does not support AMD GPUs a user could estimate the performance increase/decrease if they were to implement AMD GPU support.

We hypothesize that the CPU performance metrics give better predictions due to the maturity of CPU performance counters and the profiling tools used to record them. CPU

performance counters have been used extensively and the difficulties in recording them accurately have been well studied. On the other hand, GPU profiling, particularly for AMD, is a relatively new feature in HPCToolkit and the counters may not be as reliable as those recorded on a CPU.

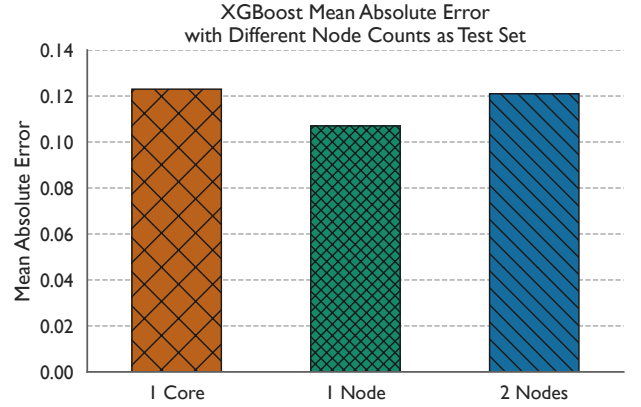


Fig. 4. The MAE of XGBoost when data for a specific resource count is removed from the training set and used for evaluation. The model performs best at predicting 1 node performance when trained on 1 core and 2 nodes data. Note that all scores are relatively low.

Figure 4 shows the performance of XGBoost when trained on data from two of the three resources amounts (1 core, 1 node, and 2 nodes), and evaluated on the third. We observe that predicting the one node relative performance vectors gives the best MAE. It is unclear whether this is because modeling the one node performance is easier or that the one core and two node data is more representative. Regardless, all three node counts score very close to 0.11 MAE, which is still a significant result.

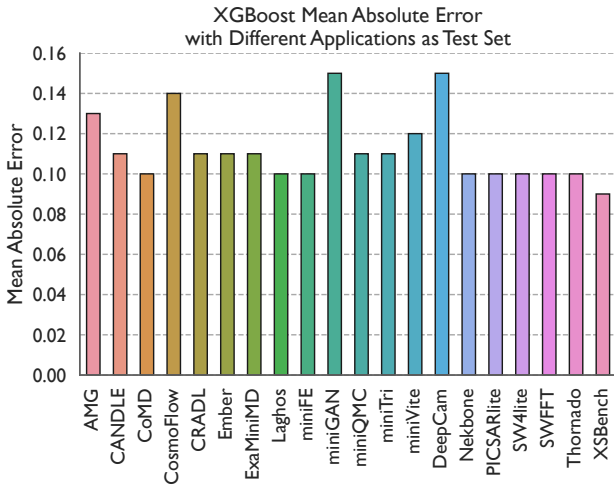


Fig. 5. The MAE of XGBoost when a specific application is removed from the training set and used for evaluation. Results are generally strong (low MAE values) across all applications.

Additionally, we can study the performance of XGBoost when trained on all but one application and evaluated on the removed application in Figure 5. Again, we see that the model performs well across all applications. However, it does notably perform worse for the ML and Python-based applications. This is possibly due to more noise and/or complicated software stacks involved in running each of these applications. These applications also tend to depend on more libraries and have more dependencies than the other applications.

C. Feature Importances

Figure 6 shows the feature importances for the XGBoost model. The most important feature is the ratio of the number of branch instructions to that of total instructions. This feature captures the control flow complexity of a program as those with more branch instructions have a more complex control flow. Since programs with more control flow generally perform worse on GPUs, the model likely uses this feature to make CPU-GPU predictions.

Next we see that the ratio of integer and single precision FP arithmetic to the total number of instructions are the next most important features in prediction. These provide insight into the data throughput of the model. In this case, applications with higher data intensity are more likely to perform better on the GPU as they are designed for high throughput data-parallel computation. These two features combined with the branching intensity make sense as the three most important features as they help the model predict relative performance between CPUs and GPUs, which is where we see the largest performance differences in the data.

The next three most important features are Ruby, Lassen, and Uses GPU, which detail where the counters were collected. This is necessary for the model to predict the relative performance vector and is likely why these are the next three most important features. We also observe that the L2 store

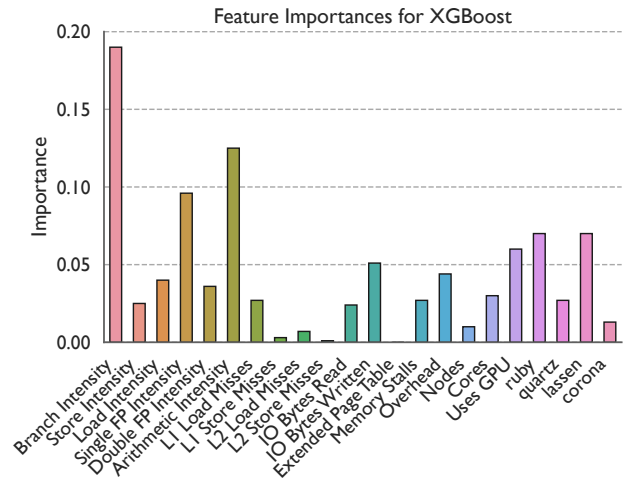


Fig. 6. Importances of each feature in the XGBoost model. A higher feature importance value suggests that it is more influential in the decision making of the model. The branch instructions intensity is the most important feature followed by the integer and single-precision floating point arithmetic intensity.

misses and extended page table features are not used in the prediction, so we can remove these during feature selection.

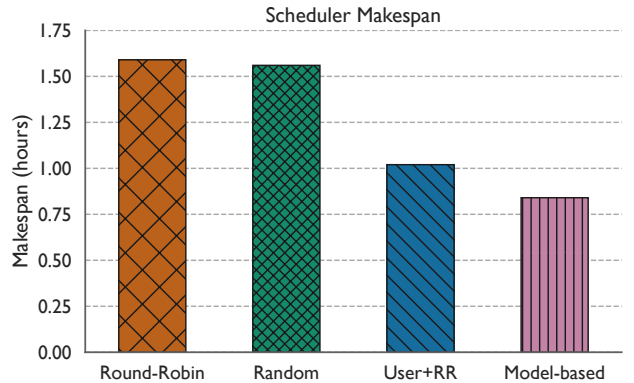


Fig. 7. The makespan of each machine selection algorithm in the scheduling simulation. Lower is better.

D. Evaluation of Scheduling Simulations

Figures 7 and 8 show the results from the scheduling simulation. Figure 7 lists the makespan for the scheduler with each machine assignment algorithm. The Model-based machine assignment method gives the lowest makespan at 0.87 hours meaning it is able to finish the job workload in a shorter amount of time than the others. Placing jobs on the most efficient resource helps improve the makespan by allowing jobs to finish sooner. The next best method is the User+RR placement algorithm. This method represents how users submit jobs to the scheduler with only the limited knowledge of the performance of their applications across machines. This is

followed by the Round-Robin and Random placement methods that perform the worst.

Figure 8 shows the average bounded-slowdown for each machine placement method. The slowdown measures the ratio of wait time and run time to just run time. As with makespan, the Model-based assignment performs the best compared to the other algorithms.

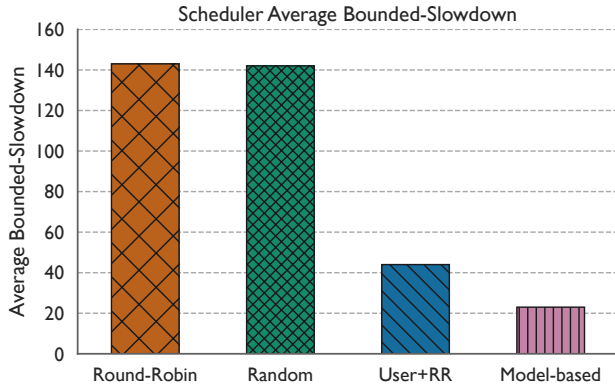


Fig. 8. The average bounded-slowdown of each machine selection algorithm in the scheduling simulation. Lower is better.

IX. CONCLUSION

The convergence of traditional HPC and new simulation, analysis, and data-science approaches provides unprecedented opportunities for scientific discovery, but also creates workflows that are more complex than ever before. These workflows often combine many applications with vastly different performance requirements that are best handled by certain types of computing hardware. Meanwhile, HPC centers and cloud platforms offer various types of computing resources to satisfy diverse needs. In this work, we study one of the many capabilities workflow users need to effectively utilize such resources: cross-architecture performance modeling. We collect the MP-HPC dataset of hardware counters across several different architectures for numerous scientific applications. We create expert derived features from these counters and train a machine learning model to predict relative performance vectors across a set of architectures with a MAE of 0.11. We further showcase how this can be used to efficiently schedule jobs across a heterogeneous set of resources.

ACKNOWLEDGMENT

This material is based upon work supported in part by the National Science Foundation under Grant No. 2047120. This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory (LLNL) under Contract DE-AC52-07NA27344 (LLNL-CONF-855652). This work was supported in part by LLNL LDRD projects 23-ERD-045 and 24-SI-005.

REFERENCES

- [1] H. I. e. a. Ingólfsson, “Machine learning-driven multiscale modeling reveals lipid-dependent dynamics of ras signaling proteins,” in *Proceedings of the National Academy of Sciences of the United States of America*, vol. 119,1, 2022.
- [2] D. H. Ahn, X. Zhang, J. Mast, S. Herbein, F. Di Natale, D. Kirshner, S. A. Jacobs, I. Karlin, D. J. Milroy, B. De Supinski, B. Van Essen, J. Allen, and F. C. Lightstone, “Scalable composition and analysis techniques for massive scientific workflows,” in *2022 IEEE 18th International Conference on e-Science (e-Science)*, 2022, pp. 32–43.
- [3] R. D. Hornung and J. A. Keasler, “The RAJA Portability Layer: Overview and Status,” Lawrence Livermore National Laboratory, Tech. Rep. LLNL-TR-661403, Sep. 2014.
- [4] T. Gamblin, M. LeGendre, M. R. Collette, G. L. Lee, A. Moody, B. R. de Supinski, and S. Futral, “The spack package manager: bringing order to hpc software chaos,” in *SC15: International Conference for High-Performance Computing, Networking, Storage and Analysis*. Los Alamitos, CA, USA: IEEE Computer Society, nov 2015. [Online]. Available: <https://doi.ieeecomputersociety.org/10.1145/2807591.2807623>
- [5] D. H. Ahn, N. Bass, A. Chu, J. Garlick, M. Grondona, S. Herbein, J. Koning, T. Patki, T. R. W. Scogland, B. Springmeyer, and M. Taufer, “Flux: Overcoming scheduling challenges for exascale workflows,” in *2018 IEEE/ACM Workflows in Support of Large-Scale Science (WORKS)*, 2018, pp. 10–19.
- [6] L. Adhianto, S. Banerjee, M. Fagan, M. Krentel, G. Marin, J. Mellor-Crummey, and N. R. Tallent, “Hpc toolkit: Tools for performance analysis of optimized parallel programs,” *Concurrency and Computation: Practice and Experience*, vol. 22, no. 6, pp. 685–701, 2010.
- [7] O. Cankur and A. Bhatele, “Comparative evaluation of call graph generation by profiling tools,” in *High Performance Computing*, A.-L. Varbanescu, A. Bhatele, P. Luszczek, and B. Marc, Eds. Cham: Springer International Publishing, 2022, pp. 213–232.
- [8] A. Bergel, A. Bhatele, D. Boehme, P. Gralka, K. Griffin, M.-A. Hermanns, D. Okanovic, O. Pearce, and T. Vierjahn, “Visual analytics challenges in analyzing calling context trees,” in *Programming and Performance Visualization Tools*, ser. Lecture Notes in Computer Science, vol. 11027, Apr. 2019. [Online]. Available: https://link.springer.com/chapter/10.1007/978-3-030-17872-7_14
- [9] A. Bhatele, S. Brink, and T. Gamblin, “Hatchet: Pruning the overgrowth in parallel profiles,” in *Proceedings of the ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis*, ser. SC ’19, Nov. 2019, ILNL-CONF-772402. [Online]. Available: <http://doi.acm.org/10.1145/3295500.3356219>
- [10] L. Zhou, X. Zhang, W. Yang, Y. Han, F. Wang, Y. Wu, and J. Yu, “Prep: Predicting job runtime with job running path on supercomputers,” in *Proceedings of the 50th International Conference on Parallel Processing*, ser. ICPP ’21. New York, NY, USA: Association for Computing Machinery, 2021. [Online]. Available: <https://doi.org/10.1145/3472456.3473521>
- [11] M. R. Wyatt, S. Herbein, T. Gamblin, A. Moody, D. H. Ahn, and M. Taufer, “Prionn: Predicting runtime and io using neural networks,” in *Proceedings of the 47th International Conference on Parallel Processing*, ser. ICPP ’18. New York, NY, USA: Association for Computing Machinery, 2018. [Online]. Available: <https://doi.org/10.1145/3225058.3225091>
- [12] H. Menon, A. Bhatele, and T. Gamblin, “Auto-tuning parameter choices using bayesian optimization,” in *Proceedings of the IEEE International Parallel & Distributed Processing Symposium*, ser. IPDPS ’20. IEEE Computer Society, May 2020.
- [13] A. Marathe, R. Anirudh, N. Jain, A. Bhatele, J. Thiagarajan, B. Kailkhura, J.-S. Yeom, B. Rountree, and T. Gamblin, “Performance modeling under resource constraints using deep transfer learning,” in *Proceedings of the ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis*, ser. SC ’17. IEEE Computer Society, Nov. 2017. [Online]. Available: <http://doi.acm.org/10.1145/3126908.3126969>
- [14] A. Bhatele, J. J. Thiagarajan, T. Groves, R. Anirudh, S. A. Smith, B. Cook, and D. K. Lowenthal, “The case of performance variability on dragonfly-based systems,” in *Proceedings of the IEEE International Parallel & Distributed Processing Symposium*, ser. IPDPS ’20. IEEE Computer Society, May 2020.

- [15] D. Nichols, A. Marathe, K. Shoga, T. Gamblin, and A. Bhatele, "Resource utilization aware job scheduling to mitigate performance variability," in *Proceedings of the IEEE International Parallel & Distributed Processing Symposium*, ser. IPDPS '22. IEEE Computer Society, May 2022.
- [16] A. Bhatele, A. R. Titus, J. J. Thiagarajan, N. Jain, T. Gamblin, P.-T. Bremer, M. Schulz, and L. V. Kale, "Identifying the culprits behind network congestion," in *Proceedings of the IEEE International Parallel & Distributed Processing Symposium*, ser. IPDPS '15. IEEE Computer Society, May 2015. [Online]. Available: <http://doi.ieeecomputersociety.org/10.1109/IPDPS.2015.92>
- [17] A. Borghesi, A. Bartolini, M. Lombardi, M. Milano, and L. Benini, "Predictive modeling for job power consumption in hpc systems," in *High Performance Computing*, J. M. Kunkel, P. Balaji, and J. Dongarra, Eds. Cham: Springer International Publishing, 2016, pp. 181–199.
- [18] P. Malakar, P. Balaprakash, V. Vishwanath, V. Morozov, and K. Kumar, "Benchmarking machine learning methods for performance modeling of scientific applications," in *2018 IEEE/ACM Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems (PMBS)*, 2018, pp. 33–44.
- [19] J. J. Thiagarajan, N. Jain, R. Anirudh, A. Giménez, R. Sridhar, A. Marathe, T. Wang, M. Emani, A. Bhatele, and T. Gamblin, "Bootstrapping parameter space exploration for fast tuning," in *Proceedings of the International Conference on Supercomputing*, ser. ICS '18, Jun. 2018. [Online]. Available: <http://doi.acm.org/10.1145/3205289.3205321>
- [20] J. J. Thiagarajan, R. Anirudh, B. Kailkhura, N. Jain, T. Islam, A. Bhatele, J.-S. Yeom, and T. Gamblin, "PADDLE: Performance analysis using a data-driven learning environment," in *Proceedings of the IEEE International Parallel & Distributed Processing Symposium*, ser. IPDPS '18. IEEE Computer Society, May 2018. [Online]. Available: <http://doi.ieeecomputersociety.org/10.1109/IPDPS.2018.00088>
- [21] W. Zhou, J. Zhang, J. Sun, and G. Sun, "Using small-scale history data to predict large-scale performance of hpc application," in *2020 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, 2020, pp. 787–795.
- [22] S. Hunold, A. Bhatele, G. Bosilca, and P. Knees, "Predicting mpi collective communication performance using machine learning," in *2020 IEEE International Conference on Cluster Computing (CLUSTER)*, 2020, pp. 259–269.
- [23] P. Balaprakash, J. Dongarra, T. Gamblin, M. Hall, J. K. Hollingsworth, B. Norris, and R. Vuduc, "Autotuning in high-performance computing applications," *Proceedings of the IEEE*, vol. 106, no. 11, pp. 2068–2083, 2018.
- [24] Y. Cho, J. W. Demmel, J. King, X. S. Li, Y. Liu, and H. Luo, "Harnessing the crowd for autotuning high-performance computing applications," in *2023 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2023, pp. 635–645.
- [25] M. R. Wyatt, S. Herbein, K. Shoga, T. Gamblin, and M. Taufer, "Canario: Sounding the alarm on io-related performance degradation," in *2020 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2020, pp. 73–83.
- [26] N. Ardalani, U. Thakker, A. Albarghouthi, and K. Sankaralingam, "A static analysis-based cross-architecture performance prediction using machine learning," 2019.
- [27] N. Ardalani, C. Lestourgeon, K. Sankaralingam, and X. Zhu, "Cross-architecture performance prediction (xapp) using cpu code to predict gpu performance," in *Proceedings of the 48th International Symposium on Microarchitecture*, ser. MICRO-48. New York, NY, USA: Association for Computing Machinery, 2015, p. 725–737. [Online]. Available: <https://doi.org/10.1145/2830772.2830780>
- [28] L. Yang, X. Ma, and F. Mueller, "Cross-platform performance prediction of parallel applications using partial execution," in *SC '05: Proceedings of the 2005 ACM/IEEE Conference on Supercomputing*, 2005, pp. 40–40.
- [29] L. L. Nesi, L. M. Schnorr, and A. Legrand, "Multi-phase task-based hpc applications: Quickly learning how to run fast," in *2022 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2022, pp. 357–367.
- [30] B. Tovar, B. Lyons, K. Mohrman, B. Sly-Delgado, K. Lannon, and D. Thain, "Dynamic task shaping for high throughput data analysis applications in high energy physics," *IPDPS International Parallel and Distributed Processing Symposium*. [Online]. Available: <https://par.nsf.gov/biblio/10356916>
- [31] "Ecp proxy applications," <https://proxyapps.exascaleproject.org/>, accessed: 2023-09-30.
- [32] "The extreme-scale scientific software stack," <https://e4s-project.github.io/index.html>, accessed: 2023-09-30.
- [33] C. R. Trott, D. Lebrun-Grandié, D. Arndt, J. Ciesko, V. Dang, N. Ellingwood, R. Gayatri, E. Harvey, D. S. Hollman, D. Ibanez, N. Liber, J. Madson, J. Miles, D. Poliakoff, A. Powell, S. Rajamanickam, M. Simberg, D. Sunderland, B. Turcksin, and J. Wilke, "Kokkos 3: Programming model extensions for the exascale era," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 4, pp. 805–817, 2022.
- [34] "Cupti," accessed: 2023-09-30. [Online]. Available: <https://docs.nvidia.com/cuda/cupti/index.html>
- [35] "rocpfprofiler," accessed: 2023-09-30. [Online]. Available: <https://rocm.docs.amd.com/projects/rocpfprofiler/en/latest/rocpfprof.html>
- [36] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD '16. New York, NY, USA: Association for Computing Machinery, 2016, p. 785–794. [Online]. Available: <https://doi.org/10.1145/2939672.2939785>
- [37] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [38] J. Lelong, V. Reis, and D. Trystram, "Tuning easy-backfilling queues," in *Job Scheduling Strategies for Parallel Processing*, D. Klusáček, W. Cirne, and N. Desai, Eds. Cham: Springer International Publishing, 2018, pp. 43–61.