

# Predictive Performance Modeling for Distributed Computing using Black-Box Monitoring and Machine Learning

Carl Witt<sup>1,\*</sup>, Marc Bux, Wladislaw Gusew, Ulf Leser

*Humboldt-Universität zu Berlin, 6, Unter den Linden, Berlin, GER.*

---

## Abstract

In many domains, the previous decade was characterized by increasing data volumes and growing complexity of computational workloads, creating new demands for highly data-parallel computing in distributed systems. Effective operation of these systems is challenging when facing uncertainties about the performance of jobs and tasks under varying resource configurations, e. g., for scheduling and resource allocation. We survey predictive performance modeling (PPM) approaches to estimate performance metrics such as execution duration, required memory or wait times of future jobs and tasks based on past performance observations. We focus on non-intrusive methods, i. e., methods that can be applied to any workload without modification, since the workload is usually a black-box from the perspective of the systems managing the computational infrastructure. We classify and compare sources of performance variation, predicted performance metrics, required training data, use cases, and the underlying prediction techniques. We conclude by identifying several open problems and pressing research needs in the field.

*Keywords:* Distributed Computing, Resource Management, Machine Learning, Black-Box Monitoring

---

## 1. Introduction

Growing computational demands in science and industry are strong drivers of innovations in distributed computing, initiating, for instance, the development of grid infrastructures or cloud computing [1, 2]. One of the fundamental problems occurring when performing complex computations in a distributed environment is scheduling: Given a set of resources (compute nodes, network, storage) and a workload, a scheduler has to decide when and which resources (nodes) are to be assigned to which units of work (jobs or tasks). Scheduling decisions are typically guided by a goal, such as minimization of runtime or of resource usage, and often have to meet additional constraints, such as resource limits of the individual nodes or data dependencies between tasks [3]. To achieve its given goals as good as possible under the specified constraints, a scheduler requires precise estimates about the expected performance of a given job or task for a given input on given resources. Only with such knowledge it is possible to take decisions which consider their future implications, which is a pre-requisite for achieving near-optimal schedules [4].

Unfortunately, such estimates are very difficult to obtain in practice [5, 6, 7, 8]. In most systems, the scheduler has no knowledge about the particular computation performed by a job or task, i. e., it treats everything as black-box. Accordingly, approaches involving analysis of source codes [9, 10] or

hand-crafted analytical models for specific operators or workloads [11, 12] are not applicable. Instead, the only available type of information are usually a few observations of performance for varying problem sizes and resource configurations extracted from log files of past executions of the same or a similar workload. The challenge is to derive precise predictions from these sparse samples.

Due to these difficulties, many practical schedulers completely disregard such estimations [13], creating schedules that potentially are far from optimal [14, 13, 15]. Other systems delegate the task of estimating resource consumption to the user, demanding a deep understanding of the workload, the input data, and the infrastructure at hand, often leading to inaccurate estimates [16, 17].

In this paper, we survey predictive performance modeling (PPM) approaches for execution durations, resource consumption, and wait times for black-box jobs and tasks based on historical data. Such approaches extrapolate historical observations to future situations without requiring detailed knowledge about the internals of the workload. Their basic idea is to model relationships between workload and execution environment based on historical performance observations, using, for instance, regression [18] or time series analysis [19].

Research on PPM dates back at least to the late eighties, motivated by load balancing in a distributed system [8]. Later, cluster computing constituted the focus of the research, followed by a series of relevant works in the context of grid computing [1]. The development to this point was characterized by an increasing consideration of system heterogeneity, moving from single-core systems over clusters of homogeneous machines to data centers encompassing a battery of different node

---

\*Corresponding author

*Email addresses:* wittcarl@informatik.hu-berlin.de (Carl Witt),  
buxmarcn@informatik.hu-berlin.de (Marc Bux),  
gusewly@informatik.hu-berlin.de (Wladislaw Gusew),  
leser@informatik.hu-berlin.de (Ulf Leser)

<sup>1</sup>ORCID: 0000-0002-7918-605X

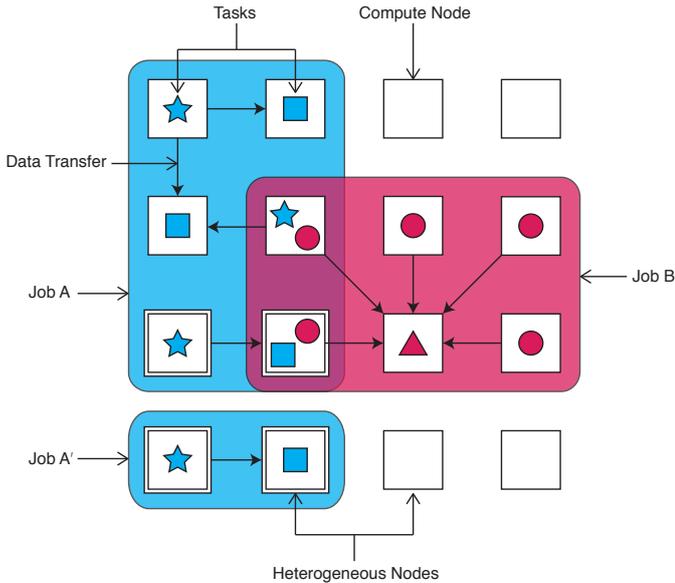


Figure 1: Overview of the distributed resource and workload model. Jobs have different resource usage patterns (computation and communication), occur at different scales (Job A and Job A'), utilize compute nodes with heterogeneous hardware specifications, and experience resource contention on a shared node (e. g., tasks  $\star$  and  $\bullet$ ).

types. System heterogeneity was further increased through the recent introduction of cloud computing, where resources are typically virtualized and available in various configurations. On the workload side, big data workloads [20, 21, 22] and complex analysis pipelines [23, 24] added to the diversity of the processing tasks, increasing the demand for other resources than CPU, like disk I/O, and the complexity of the resource usage patterns.

We are aware of only a few previous surveys in related fields. [25, 26] reviewed performance modeling, but focused on analytical models, breaking the black-box assumption. [27, 28] survey the general field of distributed resource management but without special attention to PPM. [29] surveyed machine learning methods for predicting the performance of algorithms but only considered isolated programs on a single machine, whereas we target methods for complex workloads on a distributed environment. [30] reviews techniques to model the performance of relational queries using queuing theory, but focuses on relational, transactional workloads.

This survey is structured as follows. Section 2 provides background on predictive resource management in a distributed system. Section 3 describes our comparison scheme for PPM approaches. Section 4 presents approaches to forecast the behavior of a single task on a single node. In Section 5, we review performance models for complex workloads (jobs) on a set of nodes. In Section 3.2 we review techniques for black-box performance monitoring. We conclude in Section 7 with open questions and possible research directions.

## 2. Background and Scope

In this work, we survey methods for predictive performance modeling (PPM) to improve resource efficiency and user sat-

isfaction in distributed systems. The predictions are provided either to a system component such as the scheduler or a user facing a similar decisions, e. g., the amount of resources to request for a job. In both cases, decision making is driven by certain constraints and goals, such as fairness, deadlines, or costs. To make valid decisions under a predefined optimization goal, PPM can be used to compare alternative resource allocations, job placement, starting time, resource types, etc.

In the following subsections, we first give an abstract model for the most important resources within a distributed execution environment and for a workload. We then outline exemplary scheduling use cases and their information needs.

### 2.1. Abstract View on Resources and Workloads

We model a distributed computing system as a set of *compute nodes* that are connected via a network. Each node encompasses various devices, such as CPUs, memory, disks, and network controllers. The nodes that comprise a system can vary in their available resources, featuring, for instance, different numbers of cores or different sizes of main memory. The network may also connect nodes at varying speeds, e. g., across racks in a cluster or across sites in a grid. Note that many modern system architectures virtualize resources to improve maintainability and resource usage [31].

Users submit *jobs* to the system for execution. A job works on some input data and may have parameters, which may affect the amount of work to be done and its resource usage patterns. Some jobs can be broken down to *tasks* which represent the smallest unit of work for the system. Each task is executed on a single node, but nevertheless may employ parallelism by utilizing multiple threads.

Distributed systems, and thus their resources, are usually shared by different users. Different policies for sharing resources exist. In a system that is *space-shared*, a set of resources is assigned exclusively to one job, and this assignment remains until the job finishes. In a space-shared system, jobs are typically submitted to one or more queues, implementing system properties like fairness or job priorities [32]. In a *time-shared* system, multiple jobs can run on overlapping resources at the same time. This implies that they potentially compete for system resources, giving rise to resource contention [33]. However, from a scheduling point of view, these policies are not fundamentally different, as in both cases the amount of resources and the placement need to be decided such that some optimization goal is achieved.

Our model of distributed computing covers various scenarios, including cluster, grid, and cloud computing, high performance computing and big data workloads. A very generic model is sufficient because of the versatility of the black-box approach. Although these paradigms differ strongly in their organizational aspects such as resource ownership, pricing, access, etc. black-box modeling techniques are applicable to all of them because they rely on patterns in the observed performance data rather than models of the internals of a system or workload.

## 2.2. Scheduling Use Cases

In this section, we briefly review some of the use cases that arise when using PPM to support a scheduler in a distributed system, which is a frequent motivation for research in this field. A scheduler is the component of a distributed system which assigns resources to jobs and tasks in a way that aims at optimizing some criterion, such as throughput, time-to-finish, energy consumption, or even monetary cost, e. g., when renting resources in a cloud.

Scheduling heuristics exist for various scenarios, such as running independent batch jobs by various users on a cluster [34, 13], or for executing workflows of interdependent tasks [35, 36]. Except for extremely constrained scenarios, scheduling problems are difficult to solve and usually NP-hard [37]. Accordingly, there exists a rich literature on heuristics and strategies for scheduling [38, 39, 40, 41]. The focus of our survey are methods to obtain performance predictions that support informed decision making. Such performance estimates are fundamental for many approaches to scheduling; in the following, we describe, for illustration, three popular techniques all relying on accurate cost estimates.

Backfilling is a scheduling technique to improve the utilization and response time of a system [42]. In a space-shared compute cluster, jobs are often scheduled according to a first-come-first-served policy. A backfilling policy allows short jobs to take advantage of resources left idle by larger jobs. To make sure that a job really fits a gap between reservations, predictions of job execution duration are required (e. g. [13]).

Task graph scheduling is the problem of executing computational tasks in compliance with their control and data dependencies [43]. The problem again is NP-hard, which has led to the development of numerous scheduling heuristics, like the popular Heterogeneous Earliest Finish Time heuristic [36]. These heuristics typically rely on estimated execution and file transfer times [44].

Site selection denotes the problem of choosing, for a given job, from various eligible sites the one that best fits the scheduling goal, depending on properties like required data transfer times, expected execution times, or the current load on the sites [45]. Site selection was intensively studied in grid computing, motivated by the availability of various clusters for a given job [46]. Another application area are systems that chose the best from a set of available clouds [47]. Clearly, the expected runtime of a job on a site (or cloud) is an important parameter in any such decision method.

## 3. Comparing Black-Box Performance Prediction Methods

In this section, we describe the four distinction criteria along which we compare methods for PPM: (1) Methods can be distinguished by the unit of prediction: Some predict the performance of single tasks, while others predict performance of entire jobs. (2) Another distinction lies in the metrics which are predicted, such as wall-clock time or peak memory usage. (3) Different methods take different properties of the system and

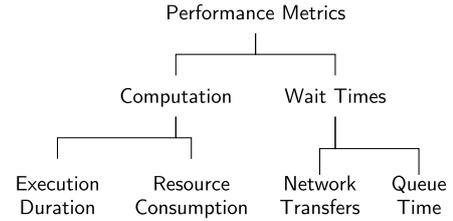


Figure 2: The performance metrics capturing the different notions of performance in a distributed system.

tasks into consideration, i.e., they differ in the principal performance factors which their prediction model covers. (4) Approaches can be classified according to the prediction method they employ, such as classification or regression. In the following, we explain each of the criteria in more detail.

### 3.1. Workload Granularity

We divide approaches in two categories depending on the target of analysis. The first category comprises task models (Section 4), i. e., methods that predict the performance of a single node executing a single task. Approaches in this category may rely on various performance information sources, such as hardware performance counters, the operating system, and binary instrumentation. In the context of distributed computing, the relevance of task performance prediction arises in various scenarios, e. g., as input to scheduling heuristics for several tasks with dependencies [36, 48].

The second category of approaches comprises job models (Section 5), which assess the collective performance of a set of nodes working together to solve a workload of tasks. A standard use case for job-level performance prediction is batch scheduling in a high-performance computing system. Such systems are typically based on per-job information from previous execution, like the total amount of the requested resources or the overall execution duration.

### 3.2. Performance Metrics

Methods for PPM differ in the metric of performance they try to assess, i.e., the concrete aspect of performance predicted. We categorize these into metrics related to computation and metrics related to wait times (see Figure 2).

The most commonly predicted metrics are related to computation, a term which itself subsumes different aspects. The most prominent one is *execution duration* (wall-clock time), i. e., the time that elapses between the beginning and the completion of a task or a job. It is typically correlated to other metrics, such as input size [7], program arguments [49], or hardware characteristics [50]. Estimated wall-clock times constitute important information for schedulers, for instance to estimate at what point-in-time a given node becomes free again, or to estimate the critical path through the tasks of a job [43]. A second perspective on computation is the *amount of resources* required to complete a task or job. For instance, predicting the number of CPU instructions can be used as an intermediate metric to estimate wall-clock time [51]. Other important aspects

Table 1: Principal Performance Factors

Abbrv.	Performance Factor
W	Workload: resource usage patterns
H	Heterogeneity: type of resources
S	Scale: amount of resources and problem size
C	Contention: performance degradation caused by resource sharing

of computation are the peak amount of memory used and the amount of disk I/O or network traffic caused. Both are critical measures for a scheduler to prevent spilling from main memory to disk and contention of processes competing for constrained data exchange channels [5, 52].

A second important category of performance metrics covers *wait times*, resulting from activities supplementary to the actual computing. Especially the *duration of file transfers* during a task or job execution has been subject of various works. Transfer times depend mainly on the amount of data to transfer, the bandwidth of the communication channel, the current load on this channel, and the possibility to use multiple channels in parallel [53, 54, 55]. Thus, transfer times largely depend on properties of a system which are not under the exclusive control of the scheduler, making analytical models difficult to employ; instead, predictions are usually based on real-time network probing [56, 57, 55, 53]. Transfer time predictions are important to estimate when all necessary data for a scheduled task will be available, or to find the fastest route to data in case multiple replica exist, or to choose a node for executing a task which has the fastest access to the data this task requires. Specific research has been devoted to the communication costs of programs written in certain paradigms, e. g., Message Passing Interface [58] or MapReduce [20]. We do not cover approaches specific to these paradigms, as the focus of this survey are black-box approaches. Finally, considerable wait times often occur when a scheduled job has to wait for resources to become free. This effect is usually captured by measuring the *time spent in a queue* waiting for resources [59]. A scheduler can utilize predicted queue times to select among multiple available computation sites.

### 3.3. Principal Performance Factors

Distributed systems are complex and exhibit multiple effects which may cause variations in the performance metrics of jobs or tasks. Methods for predicting performance metrics can be classified by the concrete set of properties of such systems which they take into account, either explicitly or implicitly. The four most influential performance factors are resource usage patterns, resource heterogeneity, scale, and resource contention.

The first and most fundamental factor is the *resource usage pattern* of the job or task that is to be scheduled. Different jobs or tasks might stress different subsystems of the system: A typical distinction is between compute-intensive jobs, demanding mostly CPU cycles, and data-intensive jobs, requiring mostly

main memory and fast data transfer [56, 59, 60, 34, 61]. In the black-box scenario we consider in this survey, one has limited insight into what a program does, but one may observe its resource usage over time and its interaction with the system. The challenge is that the behavior of a program, and thus its resource usage patterns, can depend in complex ways on the input data and the program arguments. For example, the execution duration of a scientific simulation might be dominated by the number of scenarios considered, which could be specified as application parameters.

A second important performance factor is *resource heterogeneity* in the target system, denoting varying amounts or kinds of resources across compute nodes. For instance, running a given task on a CPU with higher clock rate is expected to be faster than on a node with a slower CPUs; however, the effect depends much on the resource usage pattern of the task, i.e., whether it is compute-intensive or not. Resource heterogeneity is present in many real-world systems, often induced by extensions of systems over time where each extension utilizes the best machines affordable at the time of the expansion. Resource heterogeneity can also be present by design to allow systems to offer different levels of service [62]. Schedulers critically depend on information about resource heterogeneity to prevent load imbalance [63, 64]. Performance prediction in heterogeneous environments ideally provides a means to normalize observations with respect to hardware, such that predictions for other kinds of hardware can be extrapolated from prior performance observations.

A third performance factor is the *scale of the problem*. This includes both the amount of resources assigned to an application, e. g., the number of nodes, and the size of the problem, e. g., the amount of input data. The scaling behavior of an application is affected by various factors, such as communication patterns and contention or degree of parallelism. Since the execution times of jobs and tasks typically do not scale linearly with node (or core) count, predicting scale-out behavior is important to sensibly trade off resource demands and application performance [65]. Requesting as few resources as possible is both beneficial for optimizing resource usage, since it avoids idle resources, and for the application itself, since in space-sharing systems queue times can be expected to increase quickly with the amount of requested resources [65, 59].

A fourth performance factor is *resource contention* arising from simultaneous resource usage of different jobs or tasks. Resources in a distributed environment are often shared and the load imposed by other applications can affect the performance of a task. For instance, given the performance of a task on an idle node, the performance on a loaded node can be predicted using time series forecasting methods [66]. Another approach to predict performance degradation due to resource sharing is based on the current load on shared CPU resources, as revealed by hardware performance counters [67, 68].

### 3.4. Prediction Method

PPM methods predict a certain performance metric for a job or a task based on information about certain performance factors, typically extracted from log files of previous runs. In this

Table 2: Prediction Methods

Abbrv.	Prediction Method
CL	Classification
LL	Local Learning
R	Regression
TS	Time Series

section, we summarize the four most common techniques for actually computing predictions, which are classification, local learning, regression, and time series analysis. Note that a detailed description of the technical basis of each of these methods are beyond the scope of our work; the interested readers might refer to [69] for an introduction in classification, regression and local learning, and to [19] for time series analysis.

In the following, we focus on application possibilities of these methods to PPM and their specific advantages. Regarding the disadvantages, three apparent issues in the context of PPM are (1) since observations correspond to task or job executions, training data is potentially very expensive to collect and may only be sparsely available. (2) System workload is difficult to model and its characteristics are likely to change over time [70], such that a trained model might become obsolete. (3) It can be harder to explain insufficient prediction accuracy and find ways to improve it than for analytical or simulation models.

*Classification* refers to techniques which assign one of a (small) set of classes to a problem instance at hand. Applied to performance prediction, these classes essentially are discretized ranges of the performance metric addressed, and the problem instance is the triple of the distributed system at hand, the job or task to be executed, and the input data. Thus, approaches based on classification are fundamentally based on the assumption that tasks or jobs can be partitioned into categories with similar behavior regarding the performance metric. Predictions can be made based on historical runs, textual job descriptions, or properties of the user or organization submitting a job [34, 71, 6, 13, 60]. Classes may be derived data-driven by clustering historical runs [8] or may be specified manually based on expert knowledge. One advantage of the classification approach is that it is relatively simple and has been shown to work reasonably well using only easily collectable job and task metadata [60, 6]. A trained classification model can also provide insights into the properties of a system’s workload.

*Local learning*<sup>2</sup> approaches predict the performance of a job or task to be the same as that of the most similar job or task they have seen in the past. In a more general setting, they also might interpolate between the most similar jobs or tasks to predict the performance metric of interest. Local learning requires a function to compute the similarity of jobs and is based on the hypothesis that jobs similar under this function exhibit similar performance metrics. In contrast to classification approaches, they do not require the definition of distinct classes of

jobs. Local approaches are more flexible than classification because they do not require performance metrics to be consistent within classes, and the model can be adapted without retraining by adding new instances or discarding older ones [72].

*Regression* denotes a broad class of methods which map observations into a numerical space and then fit a function over the observations which approximates a given target variable. In our setting, the observations are performance factors of the job or task to execute, and the target variable is the performance metric to predict. Regression methods usually treat performance factors as numerical variables and consider only a limited set of functions; the by-far most popular methods express a target variable as linear combination of the independent variables. Thus, a trained regression model can also provide explicit insights about what performance factors affect a metric of interest most severely.

All previous methods are based on the assumption that the observations they are based on (for classifying jobs, for performing regression, for comparing jobs) are independent from each other. In contrast, *time series analysis* focusses on the temporal development of observations, i.e., they generate predictions based on performance trends and patterns over time. Thus, the temporal order of observations is the primary source of information, assuming that the past development of performance factors dominates future developments. These models lend naturally to performance prediction in systems whose resource usage change over time, which for the other methods requires additional steps, such as retraining the model. Another advantage is that performance variation from any source can be captured, as long as the effect is exposed in a predictable pattern over time.

#### 4. Task Performance Models

In this section, we survey concrete PPM methods at the task level, that is, performance prediction for single programs running on a single node. Such methods are important for schedulers as they allow taking into account the performance variation resulting from placing a task at different compute nodes, or the performance degradation resulting from sharing resources with another task. The presentation of the approaches is structured according the principal performance factors, as summarized in Table 3. An approach may occur in multiple of the following subsections, where we consider a paper’s contribution relevant to the overall topic of modeling that specific performance factor.

##### 4.1. Resource Usage Pattern Modeling at the Task Level

Task resource usage patterns can be modeled in different ways. One way is to characterize them implicitly by observing similarities in behavior and performance of programs. This leads to classification or local learning methods based on reusing performance observations from similar tasks. Another approach is to model resource usage across repeated executions of the same task. This leads to time series methods that predict the performance of a task based on the last executions of a task.

<sup>2</sup>also referred to as instance-based learning, non-parametric regression, and nearest neighbor regression

Table 3: Task Performance Models. The principal performance factors (W,H,S,C) and methods (Mthds.) are abbreviated according to Tables 1 and 2.

References	Context	W H S C	Input	Output	Mthds.
Marin et al. [73]	Explore interactions between an application and a platform	• • • •	CPU speed and cache sizes, small input instances, application binaries	Execution duration, L1, L2 and TLB cache miss counts	R
Hoste et al. [74]	Select platform that yields best performance	• • • •	Binary instrumentation results, benchmark scores	Execution duration (relative)	LL
Yadwadkar et al. [75]	Cloud instance selection	• • • •	Benchmark results of VM instances, representative task test run	Execution duration	R
Ferreira da Silva et al. [52]	Predict resource consumption of workflow tasks		Size of the input files	Execution duration, peak memory usage, disk usage	R
Chatzopoulos et al. [76]	Assess scalability, detect bottlenecks	• • • •	Sampled durations and stalled cycles for different core counts	Execution duration	R
Govindan et al. [77]	Virtual machine consolidation	• • • •	Memory access pattern, degradation table	Execution duration relative to execution on an idle node	LL
Gao et al. [78]	Effective scheduling in heterogeneous grids	• • • •	Execution duration, predicted execution duration of other running tasks	Execution duration	TS
Devarakonda and Iyer [79]	Load balancing in distributed systems	• • • •	Execution duration, file I/O, peak memory	Same as input	TS
Dinda[66, 80]	Deadline-constrained scheduling in shared distributed systems	• • • •	CPU load, nominal execution duration	Execution duration	TS
Zhang et al.[81]	Scheduling on distributed infrastructures	• • • •	CPU load, nominal execution duration	Execution duration	TS
Wolski et al.[56, 82, 83]	Dynamic scheduling in shared distributed systems	• • • •	CPU load, network throughput	Same as input	TS
Yang et al.[84, 85]	Adaptive scheduling on multi-user grids	• • • •	CPU load	Same as input	TS
Matsunaga and Fortes [5]	Supporting job schedulers in a heterogeneous cloud or grid environment	• • • •	Hardware specifications and benchmark results, input size	Execution duration, resident set size, size of output produced	R
Iverson et al. [7]	DAG scheduling in heterogeneous cluster	• • • •	Numerical description of the input, benchmark results	Execution duration	LL
Zhao et al. [68]	Improve resource utilization while “providing QoS guarantees”	• • • •	Benchmark results	Execution duration (relative)	R

Resource usage patterns can also be captured in a fine-grained manner by instrumenting a program or by monitoring hardware performance counters of a CPU. Such methods are also used to predict the performance degradation of tasks due to resource contention, which is covered separately in Section 4.4.

One of the earliest time series approaches to PPM was proposed by Devarakonda and Iyer [8]. In their approach, the performance observations are represented as points in a three-dimensional space comprising the dimensions execution duration, peak memory usage, and file I/O. Using the  $k$ -means algorithm, clusters of observations are determined. Subsequently, a Markov model is assembled where each state corresponds to one cluster and transition probabilities are derived from the order of observations. Based on the previous performance of a task, the performance of the next invocation of a program is predicted, by averaging over cluster centers according to transition probability.

Marin and Mellor-Crummey [73] instrument the binary of a task to observe its memory access patterns. These patterns are quantified using memory reuse distance histograms, which are based on “the number of unique memory locations accessed between a pair of accesses to a particular data item”. The program to predict is executed on small input instances of different sizes, and the bins of the resulting histograms are subjected to regression to describe how the memory access patterns change with input size.

Hoste et al. [74] use a local learning approach to predict execution durations. They introduce the concept of benchmark space to describe application behavior. The coordinates of an application are defined by 47 characteristics, such as instruction mix and branch predictability, which are independent of the CPU’s microarchitecture (e. g., its cache size and branch predictor size). These characteristics are obtained for a variety of representative tasks by means of binary instrumentation. A new task is profiled in the same way as the representative tasks before its execution to determine its coordinates in benchmark space and thus to quantify its behavioral similarity to the representative tasks. Performance measurements of benchmarks that are close to the task in benchmark space are used to predict the performance of the task on different nodes. Like Marin et al. [73], the authors focus on predicting the best suited node rather than absolute performance values. Thus, they evaluate their method by comparing predicted and actual task performance rankings of nodes rather than the differences between predicted and actual execution duration.

Zhang et al. [81] developed a method to provide a continuous prediction of CPU load, in contrast to common time-series-based approaches which are limited to a one- or a multiple-step-ahead forecast. They propose to fit a polynomial curve of second or third order on the CPU load time series using least squares regression. While this method is able to capture long-term developments in CPU load, it is slow to react to sudden turning points. To address this issue, the authors propose a means of predicting turning points by comparing the latest CPU load tendencies (i. e., differences between current subsequent measurements) against patterns that occurred in the past and that were associated with a turning point. If such a pattern can

be found, an imminent turning point is deemed likely and the forecast is derived from the historical pattern as opposed to the polynomial curve.

A seminal work for PPM is the Network Weather Service (NWS) [56] developed by Wolski et al. It is one of the pioneering works in distributed resource usage monitoring and prediction, which lays the foundation for PPM at higher levels, e. g., to guide scheduling decisions. The NWS is a distributed performance forecasting framework that continuously measures and forecasts CPU and network performance. It uses various time series analysis methods [82], which are applied simultaneously. The method which has exhibited the smallest prediction error so far is used for the next forecast. To keep the intrusiveness of the system low, the frequency at which new forecasts are generated is adjusted automatically, based on the accuracy of earlier forecasts. The provided bandwidth measurements and forecasts have been extensively used in other performance prediction approaches. Forecasts of CPU load capture the aggregated short-term resource usage patterns of all tasks running on a node, which can be used to estimate the execution time of an additional task, taking into account the background load. Refer to Section 4.4 for details on approaches that focus on the resource contention performance factor. To minimize prediction errors when forecasting grid performance, Wu et al. suggested to extend the autoregressive estimator by applying Kalman and Savitzky-Golay filters to measurement data to reduce noise [86]. Furthermore, they propose to adaptively adjust the amount of considered measurements to minimize the prediction error, similar to the adaptive sliding window estimator in the NWS.

#### 4.2. Heterogeneity Modeling at the Task Level

Performance variation from hardware heterogeneity arises in various scenarios. In a compute cluster, older hardware might be operated side-by-side with newer hardware. In a cloud computing scenario, customers choose from various instance types and in a virtualized setup, users create virtual machines with custom resource limits. Predicting the performance of a task on different node types is useful for a scheduler to avoid load imbalance or to prioritize tasks. In contrast to Section 4.1, this section focuses on modeling the effects of heterogeneity in the compute infrastructure, i. e., the compute nodes, rather than the workload.

Iverson et al. [7] propose a local learning approach to predict the execution duration of a task on a given node, based on the input size and the hardware characteristics, expressed as benchmark scores. To be able to reuse observations made on one node for similar nodes, the concept of machine space is introduced. A node’s benchmark scores define its coordinates in machine space and performance observations from nodes that are close in machine space are considered reusable across these nodes. To avoid high dimensionality in the machine space, a distance-preserving dimensionality reduction is applied. To decrease the influence of extreme data points, a constant fraction of the largest and smallest observations is discarded. The work generally assumes that resources are space-shared, i. e., there is no background load on the nodes.

Yadwadkar et al. [75] predict the execution duration of a task as a function of the chosen instance types in a cloud computing scenario. Their method is based on an offline profiling phase in which, for each instance type, the performance of a set of benchmarks is measured. In addition to the measured mean and tail performance, the task’s load on the CPU, memory, network, disk, and operating system is measured. The combination of a task’s performance and resource usage is termed its “fingerprint”. The collected profiling data is used to train a random forest that learns a mapping from a benchmark’s fingerprint on a pair of instance types to a fingerprint on another instance type. In the online phase, users submit a “representative” task which is executed and profiled on a pair of instance types to then extrapolate performance on all instance types and deliver recommendations on performance-cost tradeoffs. Although the method works well given a representative task, it does not account for varying input sizes or scaling behavior (problem size) or resource contention.

Marin and Mellor-Crummey [73] use resource usage patterns, specifically memory access patterns, to predict the relative performance of a task on different node types. Given the cache size of a target architecture and the input size, the cache miss counts are predicted. Similar techniques are used to predict the cost of the computation, which can be combined with the memory hierarchy latency to predict execution duration. Since exact prediction of execution duration is difficult, the authors focus on ranking the compute nodes according to their suitability for a given task.

Matsunaga and Fortes [5] build a regression tree for each program based on input size and hardware characteristics, e. g., CPU speed, cache size, amount of memory, and disk benchmark scores. Furthermore, the location of the input data (e. g., on the node or in a network file system) is considered. Based on these data, they predict execution duration, output file size, and resident set size of the task. They report that in a few situations where resource performance is non-linear, e. g., a network file system under load, local learning or support vector machines outperform the regression tree, which overall still performs best. Note that output file sizes and memory consumption are seldomly considered in the literature, Devarakonda and Iyer’s early work [8] being one of the few examples. Ferreira da Silva et al. [52] take a similar approach, but use a density-based clustering to partition the observations. For each cluster of observations, the pearson correlation coefficient is computed to test on a linear relationship. If there is one, execution duration, peak memory usage, and disk usage are predicted by regressing on the input file size, otherwise the mean observed resource usage is returned.

Venkataraman et al. [23] consider heterogeneity across instance types in a cloud scenario. They compare the execution duration of several tasks across instance types and use the insights to learn an task-specific regression model that relates execution duration to the input size and the number of nodes. A separate model has to be learned for each instance type. Before investing the cost of training a model for each instance type, they can optimize the costs for a given deadline or budget by choosing the appropriate instance type.

#### 4.3. Scale Modeling at the Task Level

As described in Section 3.3, we consider as the scale factor both to the number of CPU cores assigned to a task and the amount of input data. The former is important in the context of malleable tasks, i. e., in cases where more CPU cores can be assigned to reduce execution duration. Considering input sizes is important for resource consumption prediction, such as execution duration and peak memory usage.

Zhao et al. [68] predict the execution duration of a task in relationship to the number of threads by applying a thread resource contention model (Section 4.4). Each thread is treated like an independent single-threaded instance of the task, contending for CPU resources. The model takes into account additional CPU cycles incurred by resource contention and thread synchronization, as measured via profiling. The predicted total work is divided by the number of threads to obtain an execution duration estimate. They observe that inter-thread contention for resources dominates the performance variation, whereas thread creation costs and data sharing benefits are negligible. However, contention effects were only noticeable on a CPU with a relatively small cache, whereas on a CPU with eight cores and 24 MB shared cache, the PARSEC benchmarks [87] exhibited only a few percent of performance degradation due to inter-thread contention.

Chatzopoulos et al. [76] predict task execution duration as a function of the number of cores used. They use hardware performance counters to measure the number of wasted (“stalled”) CPU cycles due to various reasons, such as waiting for the cache to fetch a line. The key insights are that stalled cycles correlate highly with execution duration, but exhibit scaling trends even before they affect measured runtime. This allows better extrapolation compared to extrapolating execution durations. To fit a model, runtimes and stalled cycles are measured on different numbers of cores and then extrapolated by selecting the best fit from selected rational, cubic, and polynomial functions. When fitting the function, more weight is given to the measurements at higher numbers of cores. The method works well for scale-out predictions on similar platforms but has its limitations when it comes to different architectures, e. g., due to other performance counters or non-uniform memory access.

Matsunaga and Fortes [5] study both the effects of varying input sizes and numbers of threads for two bioinformatics tasks. They first collect training data by running the tasks on various input sizes and thread counts and then compare several machine learning algorithms in terms of prediction accuracy. They find that the PQR2 regression tree yields the lowest average percentage error and is also able to capture non-linear effects of execution duration with regard to input size and the number of CPU cores.

#### 4.4. Contention Modeling at the Task Level

Due to the trend of increasing numbers of cores per CPU, contention between tasks for shared resources has received much attention. Early approaches are based on operating-system-level measurement of resource load and task resource demand.

Binary instrumentation can be used to measure task memory access patterns, which can be used to predict performance degradation due to contention for CPU caches. Another approach is to profile the tasks' sensitivity for contention to predict the aggregate memory pressure and the according performance degradation for co-running tasks.

Dinda [66] uses autoregressive time series models to predict background load, i. e., CPU utilization by other processes, on a node. Load predictions are combined with a task's nominal execution duration, measured on an idle machine, to predict the increase in execution duration on the loaded node. The predictions have been used to determine suitable hosts for executing CPU-bound tasks in a distributed system of homogeneous nodes [80]. A solution for predicting nominal execution durations, however, is not proposed.

Yang et al. [84] assume the same notion of background load like Dinda [66]. They also take into account the variation of the background load. To assign a task, the scheduler favors nodes that have both low background load and expose little variation in background load. Both quantities can be observed and fed into standard time series predictors, but the authors observe that standard prediction scheme do not sufficiently emphasize the latest few measurements. Hence, they propose several simple strategies based on the current tendency, i. e., the difference between the latest and second-to-latest measurements. The authors found that their strategies outperform not only the other proposed strategies, but also the predictions of the Network Weather Service.

Gao et al. [78] estimate execution duration based on previous invocations and on the amount of tasks that are already running on a node. Starting from a scenario where only one program is executed, *node curves* are derived that relate the average execution duration of a task to the number of tasks already running on the node. This model is then extended to an arbitrary number of task types. The determined execution duration estimates are employed by a sampling-based scheduler, which leaves room for (re-)exploring node performance, while also exploiting favorable task-machine assignments.

Govindan et al. [77] predict the increase in execution duration due to CPU cache contention. A program that occupies cache space and bandwidth, the synthetic cache loader, is used to induce synthetic pressure on the CPU's memory subsystem. The synthetic cache loader can be configured to occupy a given number of sets and ways in a set-associative cache. It is assumed that a limited number (e. g., 256) of synthetic cache loader configurations approximately covers all possible memory access patterns. In a classification-like manner, each task is mapped to the most similar cache loader configuration, for which degradation behavior is known. Once for every node type, performance degradation of all pairwise cache loader colocations is measured. To extend to more than two co-runners, a reduction scheme is proposed that maps the pressure of two co-runners to a single, more aggressive configuration of the cache loader. This reduction is applied recursively to predict performance degradation with up to four co-runners in the experiments.

Zhao et al. [68] also follow a memory-centric approach us-

ing the same basic aggregate pressure approach as in [88, 77]. They distinguish between CPU cache space consumption and cache bandwidth consumption. This is motivated by the observation that performance degradation can be modeled as a piecewise linear function of the aggregate pressure on both cache space and cache bandwidth. To measure the cache space and bandwidth consumption, hardware performance counters are used. To find the piecewise linear function that best models the execution duration of a task, all models in a set of user-defined models are evaluated. To reduce the costs of trying all possible model forms, a task-independent model is first created that optimizes the functional form, whereas task-specific parameters for that model are fitted in a second step.

## 5. Job Performance Models

In this section, we survey PPM approaches at the job level, that is, prediction methods for complex workloads utilizing several nodes connected through a network. Like for tasks, we structure the review of performance models along the principal performance factors from Section 3.3.

### 5.1. Resource Usage Pattern Modeling at the Job Level

Compared to the task level, resource usage patterns at the job level include communication between nodes. In addition, a job resource usage pattern depends on the resource usage patterns of its tasks, leading to potentially much more complex patterns. Job resource usage patterns can sometimes be inferred from job metadata, such as the name of the application or the submitting user. The assumption is that jobs with similar metadata will have similar resource usage patterns and thus similar performance metrics. Historically, analyses of supercomputer workloads showed that categorizing jobs according to their metadata reduces the variance among the observed runtimes [94]. Not very surprisingly, invocations of the same job, submitted by the same user on the same number of requested processors have been observed to expose much less runtime variation than the set of all jobs observed in a cluster. Gibbons [60] verified this observation and used it to predict job execution duration by simply averaging runtimes of prior jobs with the same metadata. Downey [59] makes a similar observation in a San Diego Supercomputer Center workload, where job runtimes approximately followed a log-uniform distribution, with different parameters for the queues for short, medium, and long jobs.

Smith et al. [97] build on the observations by Gibbons [60] to predict job execution durations from scheduler logs. They coined the term *template* for a set of metadata attributes used to partition the observations into groups. The execution duration for a new job is predicted by looking up prior observations in the according category and aggregating them either by using the arithmetic mean or by regression over the number of assigned processors. A genetic algorithm is used to determine the metadata features to group by and the aggregation method. Instead of searching for a single best template, they also propose to search for "template sets" that consist of one to ten templates,

Table 4: Job Performance Models. The principal performance factors (W,H,S,C) and methods (Mthds.) are abbreviated according to Tables 1 and 2.

References	Context	W H S C	Input	Output	Mthds.
Delimitrou et al. [34]	Data center scheduling	• • • •	Short test runs of the application	Performance metrics under different resource configurations	CL
Smith [89]	Site selection in a compute grid	• • • •	Job metadata, requested processors, maximum run time	Execution duration, queue time, transfer duration	LL
Arslan et al. [54]	Optimizing end-to-end data transfer parameters for a set of files	• • • •	Historical data transfer times, parameters, and background traffic probe data	Transfer duration	R, LL
Venkataraman et al. [23]	Cloud resource configuration	• • • •	Input size, number of machines	Execution duration	R
Pumma et al. [90]	Scientific computing	• • • •	Input size, hardware usage during probe execution	Execution duration	CR
Alipourfard et al. [91]	Cloud resource configuration	• • • •	Job and initial cloud configuration	Optimized Configuration	R
Gaussier et al. [15]	Provide runtime estimates for a backfilling scheduler	• • • •	User’s job history, job submission time	Execution duration	R
Lee et al. [49]	Tune application parameters for a platform	• • • •	Application parameters, sampled execution runs	Execution duration	R
Dobber et al. [92]	Increase robustness of applications in the grid	• • • •	Previous job durations	Execution duration	TS
Tsafir et al. [13]	Backfilling in a HPC context	• • • •	History of user estimated and observed runtimes	Execution duration	TS
Sanjay and Vadhiyar [93]	Site selection in a grid environment	• • • •	CPU and network benchmark results	Execution duration	R
Cunha et al. [47]	Site selection (cloud or private cluster)	• • • •	Metadata, requested Processors, queue state, submission time	Execution duration, queue time	LL
Pfeiffer and Wright [50]	“high-performance (HPC) system acquisitions”	• • • •	CPU and network characteristics, HPCC benchmark results	Execution duration	R
Li et al. [46]	Site selection (grid)	• • • •	Metadata, executable name, requested processors, maximum run time, submission time, queue state	Execution duration, queue time	LL

Table 4: (continued) Job Performance Models.

References	Context	W H S C	Input	Output	Mthds.
Downey [59]	Site selection	. . . ●	Number of running jobs, their age, and occupied number of processors, requested number of processors	Queue time	R
Barnes et al. [65]	Improve cluster efficiency	. . ● .	Application parameters, number of processors	Execution duration	R
Gibbons [60]	Provide estimates to scheduler	. . ● .	Executable name, user name, assigned number of processors, current job age	Execution duration	CL, R
Wu et al. [86]	Scheduling and load balancing on the grid	. . . ●	CPU load history	CPU load	TS
Brevik and Nurmi [95, 96]	Estimating queueing delays for batch jobs	. . . ●	Queue time history	Queue time	TS

each giving one prediction with a confidence interval for a new job. Li et al. [45] found that averaging all templates' predictions can improve prediction accuracy. The template approach later has been extended by considering more job attributes, e. g., the topology of a workflow-structured application [98], predicting other resource types, e. g., memory usage [99], and by using other prediction methods within the categories, e. g., time series predictors [45, 100].

Tsafrir et al. [13] predict the runtime of a job as the average of the runtimes of the last two submitted jobs of the same user. This is a time series approach that aims more at the resource usage pattern of the user rather than the job. The authors compare several variations of the above scheme and conclude that no method works best for all systems and workloads. Other time series-based approaches include Sonmez et al. [100], who show that the accuracy of time series predictors for job runtime and queue time predictions can be boosted when building separate models for each user, computing site, or both. Gaussier et al. [15] propose a set of features similar to the above time series features for a linear regression model. They include a variety of information about the user's job history, e. g., the execution durations of the last three jobs and the time since the last job of that user was completed. Their cost function penalizes the underestimation of execution duration stronger than overprediction because backfilling schedulers need to take corrective actions when jobs exceed their reservation, in the worst case even killing the job.

Dobber et al. [92] compare several time series predictors for execution duration prediction. The proposed adaptive exponential smoothing technique improves upon the basic exponential smoothing predictor implemented in the Network Weather Service [56] (see Section 5.4) by continuously adjusting the smoothing parameter based on the observed forecast error. This way, the weight of the latest measurement is increased if, for instance, the last forecast was found to be inaccurate. They also

propose a dynamic exponential smoothing method, which combines a sliding window approach with the adaptive exponential smoothing predictor. Results were found to be comparable or even preferable to the Network Weather Service, due to an increased robustness to peaks and changes in performance.

Recently, more sophisticated classification methods have been applied to job performance prediction. The performance predictions of Matsunaga et al. [5] and Dwyer et al. [101] are based on regression trees, a technique that combines classification and regression [102]. Both predict different metrics, but identify regression trees as most accurate among a wide range of machine learning techniques. This contradicts the observations by Smith et al. [97] and Tsafrir [13], who favor simple predictors like the mean over a regression-based aggregation of observations.

Pumma et al. [90] use an approach that completely ignores metadata. Instead, each job is run for a short while and profiled using platform independent low level features derived from hardware performance counters. These include for instance instruction mix, branch predictability, and address distances between successive cache accesses. Using a pre-built classification scheme based on seven major types ("Dwarfs") of scientific applications [103], the job is assigned to a job class that implies its resource usage pattern. The C4.5 decision tree is trained using manually labelled benchmarks from three benchmark suites [104, 105, 106]. For each job class, linear regression is used to model runtime as a function of input size and observed resource usage patterns (the measured low level characteristics).

Delimitrou et al. [34] use singular value decomposition to classify an incoming job based on a short test run of the application. The idea is to extrapolate the performance information for untested configurations from other applications with similar resource usage patterns, as revealed by the test runs. The assumption here is again that the workload can be partitioned into jobs that expose similar behavior under the principal per-

formance factors.

Lee et al. [49, 107] analyze the impact of application parameters on the execution duration of two high performance scientific applications. These applications have large parameter spaces comprising both parameters that define the workload, e. g., working set size per processor, and those that define how to execute the workload, e. g., which broadcast method to use. Such parameters clearly affect the resource usage pattern, and through the latter parameters, their work gains a performance tuning aspect. Moreover, a statistical analysis is incorporated, which indicates non-linear relationships between parameters and execution duration. The authors thus choose neural networks and piecewise polynomial regression for modeling the relationship.

### 5.2. Heterogeneity Modeling at the Job Level

Heterogeneity at the job level is more complex than at the task level, because it may affect different parts of a job. At the task level, hardware heterogeneity, i. e., hardware specifications affect the entire task. At the job level, parts of the job may run on different nodes, and these parts depend on the scheduler, its awareness of hardware differences and its ability to incorporate them in its decisions. A typical problem in grid computing is to choose a site to run a job on, under the premise that resource characteristics of the sites are heterogeneous. A variant of the site selection problem is computing in a hybrid cloud, in which, for instance, some of the compute resources are owned by the organization and additional resources can be acquired via cloud computing on demand. The estimation of job finish times depends on predictions of queue time, execution duration, and file transfer duration.<sup>3</sup>

Sanjay and Vadhiyar [93] model execution duration of a parallel application as the sum of computation duration and non-overlapped communication duration. Both durations are modeled by regression, using a mixture of 77 polynomial and logarithmic functions. They address cross-platform performance prediction by starting from a regression model learned on a reference platform, conducting dedicated small-scale experiments on the reference and target platform and using the observed performance ratios to scale the coefficients of the reference model.

Cunha et al. [47] consider performance differences between a local cluster and cloud resources. They assume that executing in the cloud slows down execution by a constant, application-dependent factor, for which they propose a linear model and an empirical model, based on relative performance of eight applications on three clusters and on three cloud platforms. Cunha et al. acknowledge that it is difficult to predict relative performance because of various factors, such as the network performance. To alleviate this issue, they use a “cutoff function” that favors executing on the local cluster in case of high uncertainty associated with the predictions.

Pfeiffer and Wright [50] assume a homogeneous cluster, but build models that regress on the performance characteristics of

a node type and thus have the possibility to predict performance for different (but homogeneous) clusters. They also express job execution duration as the sum of computation time and communication time that is not overlapped with computation. The former is modeled as a linear combination of CPU and memory benchmark results, whereas the latter is expressed as a linear combination of interconnect bandwidth and latency benchmark results. Both fits are performed using non-negative least squares regression, on a subset of the benchmarks automatically selected by backward elimination [108].

Alipourfard et al. [91] propose a search method to select cloud configurations that optimize the performance of a job. A configuration comprises the number of virtual machine instances, number of cores, CPU speed, RAM per core, disk count and speed, and network capacity. A gaussian process model is used to estimate the uncertainty associated with a configuration and to select the one that has the largest potential for increasing performance. The configuration evaluated and the model is updated until expected improvements fall below a threshold. A limitation of the method is its reliance on representative workloads. Since the search procedure is too slow to be applied to a job at hand, it is applied offline to representative workloads instead. This is particularly problematic with respect to changes in input size, which the method does not model explicitly. This makes the method a perfect fit recurring jobs, which often operate on data sets of similar size.

### 5.3. Scale Modeling at the Job Level

The scale of a job refers both to the problem size, e. g., amount of input data, and the amount of resources assigned to the job. Compared to tasks, the effects of scale can be much more pronounced and thus require careful decision making. For instance, choosing the right amount of resources is important because both undersizing and oversizing can have adverse effects. Allocating too many resources might result in an inefficient utilization, e. g., resources left idle although reserved, or cause additional parallelization overheads that diminish or even outweigh performance gains. Also, wait time in a cluster resource queue usually increases with the amount of requested resources. Too few resources on the other hand may result in a missed deadline or the failure of the job. Since scale-out is often modeled as a linear or sub-linear function, regression methods are a natural fit to this problem. Another option is to include the amount of resources, e. g., the number of requested cores into the distance metric of a local learning approach [89, 46].

Barnes et al. [65] regress execution duration both on application parameters and processor count. Here, the application parameters describe the problem size, e. g., the resolution of a fluid dynamics simulation. They try both a linear and quadratic form for the processor count and pick the better fit. All variables are transformed by taking their logarithm before regression to avoid the large execution durations dominating the short ones in term of influencing the model coefficients. They observe that regressing separately on the time spent for computation and communication works better, as these two quantities tend to scale differently.

<sup>3</sup>We address queue time estimation in Section 5.4, as queue times are a result of jobs contending for resources.

Venkataraman et al. [23] model application execution duration as a function of input size and the number of nodes. They argue that execution duration can be modeled as the sum of constant costs, the costs of parallelizable work, a logarithmic communication cost for aggregating data in a tree pattern, and linear per-machine overheads. They use a non-negative least squares regression to find the coefficients of these cost factors given some runs of an application at different scales.

Arslan et al. [54] predict the duration of large file transfers on the network, focusing on three parameters: amount of data, concurrency, and parallelism, where concurrency denotes the transfer of several files in parallel and parallelism refers to transferring the blocks of a single file with several TCP streams in parallel. The authors apply local learning to retrieve the most similar observations to the current network load situation as found by historical probe data. Then, the relationship between the transfer parameters concurrency, and parallelism is found using polynomial regression. These models are finally used to predict optimal settings for a transfer of a set of files. A similar approach is proposed by Liu et al. [109], who also consider the two file transfer parameters concurrency and parallelism. Based on large Globus GridFTP logs, they engineer additional features, such as the load on the two endpoints. Finally, they build a regression model that relates the transfer configuration to its duration.

#### 5.4. Contention Modeling at the Job Level

Job resource contention commonly arises in two ways: queue times for compute jobs and background load on the network for file transfer jobs. We cover these scenarios in the following subsections.

##### 5.4.1. Queue Time Prediction

On space-shared clusters, jobs queue up for resources and the scheduler decides which job to run next. For a metascheduler, i. e., a scheduler that submits a job to one of several sites, predictions of queue time are important, as they add up to the job execution duration. Queue time can be predicted based on the queued jobs, their predicted execution durations, predicted new job arrivals, and the scheduling policy applied to the queue.

Downey [59] considers the case of a simple first-come-first-served scheduler to predict wait time for the job at the head of the queue based on the number of processors it has requested. He first fits a log-uniform distribution to the execution durations of the jobs in the queues for short, medium, and long jobs. Given the current age of a job, the probability of running for another  $t$  minutes can be estimated. Using this, the probability of a number of processors becoming free during the next  $t$  minutes can be estimated. Smith et al. [97] applied the template method (see Section 5.1) they developed for runtime predictions to queue time prediction and report improvements over Downey’s method. They also consider more schedulers, including a backfilling scheduler. Li et al. [46] combine the metadata-based classification of Smith and an idea similar to Downey’s approach in a local learning approach. The state of the resource is characterized by attributes like the number of

currently running jobs and their estimated runtimes. A genetic algorithm is used to optimize various hyperparameters.

Brevik and Nurmi et al. [95, 96] propose a time series-based approach for queue time prediction. The first component is a conservative quantile estimator for distributions, which is applied to queue waiting times, disregarding the temporal order of observations. The second component is a change point detector. The assumption is that queue times are usually consistent for an extended period of time before a change causes entrance into another consistent phase. A change point is assumed after observing a fixed number of queue times beyond the .95 quantile of the current distribution of queue times. Similar to Smith’s template approach, queue time observations are grouped according to their job metadata.

Sonmez et al. [100] conducted a comparative evaluation of estimators for job execution duration and for queue time on grid job logs. This includes various time series methods and Nurmi’s method [96]. They confirm that partitioning observations, e. g., by site, user or both, improves predictions, although they do not observe strong improvements in the scheduling quality in a trace-based simulation. This contradicts other studies [13, 4], emphasizing the difficulty of comparing different systems under various workloads.

##### 5.4.2. File Transfer Duration Prediction

Resource contention also plays an important role for file transfers. The replica selection problem arises when multiple copies (replicas) of a file reside in the distributed environment. A scheduler may want to select the one with shortest transfer duration to a specific endpoint. Here, resource contention has to be considered in the form of background load on the network and the endpoints.

Faerman et al. [55] address the problem of predicting file transfer durations based on network bandwidth probes. The problem is that, to keep probing overhead low, the probes are usually small in comparison to actual files. Instead of just dividing file size by estimated bandwidth, the authors regress transfer duration on bandwidth probe measurements (a single coefficient regression model), which presumably captures the overhead of larger file transfers. Swamy and Wolski [110] also map probed bandwidth values to file transfer durations. Instead of a regression model, they map the cumulative distribution function of the bandwidth probe values  $CDF_B$  to the cumulative distribution of file transfer durations  $CDF_D$ . The forecast is obtained by looking up  $CDF_D^{-1}$  for  $CDF_B(B)$ . The bandwidth value  $B$  is suggested to be obtained from a Network Weather Service forecast. Vazhkudai and Schopf [111] show that also considering disk load values improved predictions. A comprehensive study considering disk load and several additional factors was recently conducted by Liu et al. [109].

Qiao et al. [112] consider the impact of the frequency of bandwidth probes. They claim that for small files, a higher-resolution signal is needed, whereas for large files, a coarse-grained observation history will suffice. They evaluate numerous network packet traces with bandwidth utilization data captured on WANs and LANs in combination with smoothing techniques to adapt the time scale of the bandwidth utilization time

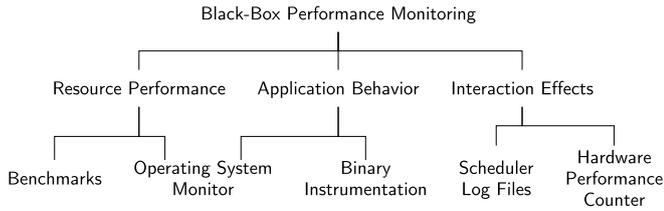


Figure 3: Methods for performance monitoring that do not require modification of the workload or execution system.

series. Subsequently, they apply various prediction methods to this data, including the mean and last value alongside autoregression approaches and nonlinear schemes. They find that often there is a sweet spot in the resolution of the bandwidth probe signal which improves predictability, but the optimal resolution depends on unidentified criteria. Notably, simple models are shown to be competitive with more complex prediction schemes, a result that has already been observed by Vazhkudai and Schopf [111].

## 6. Obtaining Performance Indicators

All methods presented in this survey rely on performance measurements of jobs or tasks that do not require modifications to the workload. We thus give a brief overview of common tools and methods for black-box performance monitoring. To observe performance of jobs and tasks in a non-intrusive way, various tools and methods are available. These produce data at different levels of granularity and at different levels of overhead. As shown in Figure 3, one can distinguish approaches that aim at measuring the behavior of an application, the performance of the resources, and the interaction between both.<sup>4</sup> In the following, we briefly explain each approach, give usage examples, cite exemplary tools and discuss their overhead.

*Benchmarks* can be used to measure the performance of a node with respect to some representative task. Measuring such task-specific performance is important because resource consumption may depend on complex interactions, e. g., between CPU, memory and mass storage. To compare the performance of two nodes with respect to some task, their benchmark scores are likely more useful than their hardware specifications, given that the task to schedule is similar to the benchmark task. Popular benchmark suites are SPEC [114] and PARSEC [87]. As benchmarks are usually only executed once, their overhead is negligible. *Load probes* are a related technique, which aim at the current, transient performance of a subsystem, usually the network. A small amount of data is sent in order to measure the current available bandwidth and latency. This technique was used for instance in the Network Weather Service [56] and many others [115, 55]. With respect to overhead, probe sizes closer to the actual transfer size are more representative, but also incur higher overhead [55].

<sup>4</sup>Another common source of information are application logs, but these require information on how to parse, interpret, and use the relevant information [113], which is not available in a black-box scenario.

*Operating systems* provide means to inspect the current load on a node’s devices and all running processes. In a virtualized environment, the hypervisor plays a similar role as the operating system and can be used to collect performance measurements [67]. In a Linux environment, the virtual `proc` file system exposes performance information, e. g., about the CPU and memory usage of individual processes. Another option is to monitor the task’s interaction with the operating system by intercepting system calls. This can be used, e. g., to measure the amount of data a task reads and writes, which is an aspect of application behavior, rather than resource performance. Tools like `uptime`, `vmstat`, and `iostat` are useful to determine the current load of CPU or mass storage devices [56, 111]. As an example of a higher-level tool, Kickstart [116] can be used to wrap the execution of any task and collect comprehensive performance information. The overhead of polling operating system performance metrics is low, while intercepting system calls for tasks can be expensive, depending on how much information is collected.

*Binary instrumentation* is a technique that modifies the executable of a program to collect information about its behavior while it executes. In general, it is challenging to obtain hardware-independent information about a black-box task, because its observed performance is always a product of the task behavior and the resources it interacts with. Using binary instrumentation, memory access patterns can be characterized in a microarchitecture independent manner [117], giving some low-level insights into program behavior. For instance, Marin et al. [73] use binary instrumentation to measure memory access patterns to select a CPU with appropriate cache size. Typical tools for binary instrumentation are PEBIL [118] and Dyninst [119]. This technique has relatively high overhead and the execution duration of the binary can increase considerably.

*Scheduler log files* are a widely used source for training prediction models. They usually comprise metadata and performance data, such as the submitting user, the application name and the observed resource consumption of a job. Since high-level performance metrics like the job execution duration are a joint product of resource usage pattern and hardware characteristics, we classify log files as a monitoring technique for interaction effects. Comprehensive logs from large-scale systems in science and industry are occasionally published, e. g., the parallel workloads archive [120] or Google cluster traces [121]. Each application can also log its own performance data, and a tailored analysis can yield valuable information [113], but since it cannot be exploited in an application-independent way, such approaches are out of the scope for this work. Since the scheduler usually collects accounting information anyways during operation, log data comes virtually for free.

*Hardware performance counters* provide low level insights into task behavior. Performance counters are dedicated registers in a CPU that can be used for measuring events, such as cache misses, branch mispredictions, or the total number of CPU cycles consumed by a process. They are often used to reveal the performance degradation due to sharing resources with another task, e. g., by measuring an increase in cache misses or cycles per instruction compared to executing with exclusive access to

Table 5: Open topics in PPM research. Benefits are abbreviated as prediction accuracy (A), model generalizability (G), and overhead reduction (O).

Domain	Opportunities	Benefits
Prediction	Online Learning	A,G,O
	Active Learning	A,O
	Multi-resource models	A
	Temporal resource usage	A
	Multi-level parallelism	A
	Continuous prediction	A
Data collection	Resource-agnostic workload patterns	A,G,O
	Monitoring data integration	A,G
Decision making	Holistic approaches	A,G,O
	Sufficient accuracy models	O

resources. For instance, Zhao et al. [68] and Dwyer et al. [101] predict the slowdown experienced by two applications contending for shared CPU resources based on performance counter observations. As the configuration of these hardware performance counters is microarchitecture-specific, typically a high-level interface like PAPI [122] or LIKWID [123] is used for interacting with performance counters. Because measurement takes place in hardware, hardware performance counters incur only negligible overhead.

## 7. Open Issues and Directions for Research

The reviewed literature demonstrates the long-standing and evolving research in predictive performance modeling. In this section, we identify future trends and research opportunities in the field. We tag research opportunities according to three classes of improvements. Most research aims at improving the prediction accuracy (A) delivered by a predictor in a certain scenario. A second, often opposing goal is model generalizability (G), i. e., the ability of a predictor to perform well in a large range of scenarios. Thirdly, collecting training data for predictors can incur significant costs. Therefore, reducing the overhead (O) associated to integrating PPM into a system is another research goal. An overview of the research opportunities discussed in the following is given in Table 5.

### 7.1. Prediction

Often, PPM methods involve an offline phase in which performance measurements are taken, models are trained, and hyperparameters are chosen. On the upside, this approach allows for costly collection of high quality training data (since it is usually performed only once) and keeps the complexity of the production system low. On the downside, this approach cannot directly benefit from data collected during system operation. Performance may degrade over time, as the system or its workload changes. Thus, we see a major research opportunity in approaching training and prediction in an *online* fashion, i. e., collecting training data and updating prediction models during system operation.

A major challenge in PPM using machine learning is that training data sets can be expensive to obtain, as each observation might come at significant costs. In addition to using data that is generated during system operation, *active learning* [124] approaches and sampling strategies [125, 126] might reduce data collection overhead and increase accuracy by achieving better coverage of the feature space.

In the literature, models for predicting execution duration dominate. Recently however, other resource types such as memory [127, 128] shift more into focus. A research opportunity would be to investigate and predict the *interdependent usage of resources* rather than predicting the usage of each resource type separately [5, 8]. Predictions regarding disk and memory usage are essential to more efficiently provision resources, e. g., in a cloud scenario [129]. Predictions for resources other than the CPU gain importance as I/O- and memory-intensive workloads arise [130] in distributed computing.

A related research direction is the exploration of *temporal variations*. Resource usage of tasks and jobs usually varies across time and predicting the usage of multiple resources as a (multi-dimensional) time series rather than as their aggregate metrics could lead to more sophisticated schemes to improve resource utilization and reduce contention between jobs and tasks. Phase analysis [131, 132] is a related research field that explores the resource usage of programs on a cycle-accurate level, but we see a research gap in exploring the use of several resources in a black-box fashion and across various time scales.

Large-scale computing involves *multiple levels of parallelism*, ranging from multi-core architectures over clusters to grid and federated cloud [133] scenarios. This flexibility offers new trade-offs with severe performance implications, for instance when having to choose between fewer machines with more cores and more machines with fewer cores. Resources like memory, network, and disks per compute node contribute to complex performance behavior [91]. Predicting performance for various resource configurations and granularities of parallelization is especially pressing in cloud scenarios, where users have to choose from a large range of compute node types [23].

Another challenge is handling errors from PPM. While all methods strive for high average accuracy, maximum errors can still be high. To mitigate those effects and to avoid catastrophic decisions, schemes to correct decisions [134, 135] and update predictions [34] are emerging. Another possible research direction is to extend single predictions before making a decision with *continuous predictions* that account for newly arriving monitoring information and provide a decision maker with updated, more accurate predictions.

### 7.2. Data Collection

A largely unexplored aspect is the collection of *resource-agnostic* workload patterns. Observed performance is usually a product of a task’s or job’s resource usage patterns and a specific resource configuration. For example, the number of cache misses is a product of the application’s memory access pattern and the cache size. To predict performance across various resource configurations, it would be helpful to quantify application behavior in a resource-independent way. Approaches that

collect such information in a black-box fashion exist [117, 73] but incur high runtime overheads and cover only low-level performance metrics. Detecting and accounting for more general interaction effects between workload and resources would support PPM by reducing the amount of data that needs to be collected and would help to generalize better to other resource configurations.

Various technologies have transformed the landscape of distributed computing during the last years. For instance, the cgroups kernel feature in Linux operating systems now allows for better monitoring and control of the resource usage of processes or groups of processes. On the networking side, software defined networking has introduced new opportunities to monitor and thus predict resource usage. Hardware performance counters allow for low-level characterization of application performance. Aggregating this low-level information to more descriptive metrics and integrating the data of various sources to characterize an application’s resource usage is both an engineering and a research challenge.

### 7.3. Decision Making

There is a dichotomy in the literature between scheduling and prediction: scheduling research often takes the availability of performance estimates for granted and research on PPM seldomly explores its impact on scheduling. We see a large research potential in the exploration of *holistic approaches* to prediction, data collection, and decision making. For instance, improving prediction quality could be added the objectives of a scheduler. By integrating scheduling and prediction, systems that continuously improve accuracy, adopt to changing workloads, and make optimal use of performance observations during system operation become feasible.

In this context, the question of *sufficient prediction accuracy* arises, i. e., how scheduling quality is affected by prediction accuracy. This relationship has been empirically evaluated before [60, 4, 13], but the results are ambiguous and still lack a theoretical foundation. Further research in this area might enable models of sufficient accuracy, that allow for better trade-offs between data collection efforts, the resulting prediction accuracy, and the returns of improved decision making.

## 8. Conclusion

Predictive performance modeling methods are important to support distributed computing at growing scales and levels of automation, under increasing hardware complexity and workload diversity. We surveyed approaches that provide performance predictions to improve the utilization, throughput, and other efficiency metrics of distributed computing systems. The approaches were compared in various dimensions, most importantly the principal performance factors they account for, i. e., resource usage patterns, resource heterogeneity, problem scale, and resource contention. From this perspective, we reviewed and classified the literature according to prediction model, input and output data, and use cases.

In conclusion, machine learning methods are an approach to the problem of efficiently managing and planning resources in

distributed systems under various factors of uncertainty. They offer greater versatility and reduced development costs compared to simulation and analytical performance models that are tailor-made for specific applications and infrastructures. The key idea is to put a stronger emphasis on observed performance factors as opposed to the performance factors anticipated by an analytical model or simulation, which offers greater flexibility at the price of reduced prediction accuracy. However, we believe that machine learning approaches are necessary to make PPM available in practice for a large range of applications and computing platforms. Further developments in PPM and their integration into scheduling algorithms promise improved resource utilization, reduced execution times and costs, and more satisfied users.

## Acknowledgements

Carl Witt, Marc Bux, and Wladislaw Gusew have received funding by the Deutsche Forschungsgemeinschaft (DFG) through the SOAMED graduate school (GRK 1651).

## References

### References

- [1] I. T. Foster, Y. Zhao, I. Raicu, S. Lu, Cloud Computing and Grid Computing 360-degree compared, in: Grid Computing Environments Workshop, IEEE, 2008.
- [2] H. Hussain, S. U. R. Malik, A. Hameed, S. U. Khan, G. Bickler, N. Min-Allah, M. B. Qureshi, L. Zhang, W. Yongji, N. Ghani, J. Kolodziej, A. Y. Zomaya, C.-Z. Xu, P. Balaji, A. Vishnu, F. Pintel, J. E. Pecero, D. Kli-azovich, P. Bouvry, H. Li, L. Wang, D. Chen, A. Rayes, A survey on resource allocation in high performance distributed computing systems, *Parallel Computing* 39 (11) (2013) 709–736.
- [3] T. L. Casavant, J. G. Kuhl, A taxonomy of scheduling in general-purpose distributed computing systems, *IEEE Transactions on Software Engineering* 14 (2) (1988) 141–154.
- [4] H. Casanova, A. Legrand, D. Zagorodnov, F. Berman, Heuristics for scheduling parameter sweep applications in grid environments, 9th Heterogeneous Computing Workshop (HCW 2000) (2000) 349–363.
- [5] A. Matsunaga, J. A. B. Fortes, On the Use of Machine Learning to Predict the Time and Resources Consumed by Applications, in: *IEEE/ACM International Conference on Cluster, Cloud and Grid Computing*, IEEE, 2010, pp. 495–504.
- [6] W. Smith, I. T. Foster, V. E. Taylor, Predicting application run times using historical information, *Job Scheduling Strategies for Parallel Processing*.
- [7] M. A. Iverson, F. Özgüner, L. Potter, Statistical Prediction of Task Execution Times through Analytic Benchmarking for Scheduling in a Heterogeneous Environment, *IEEE Transactions on Computers* 48 (12) (1999) 1374–1379.
- [8] M. V. Devarakonda, R. K. Iyer, Predictability of Process Resource Usage: A Measurement-Based Study on UNIX, *IEEE Trans. Software Eng.* 15 (12) (1989) 1579–1586.
- [9] G. R. Nudd, D. J. Kerbyson, E. Papaefstathiou, S. C. Perry, J. S. Harper, D. V. Wilcox, Pace—A Toolset for the Performance Prediction of Parallel and Distributed Systems, *International Journal of High Performance Computing Applications* 14 (3) (2000) 228–251.
- [10] V. E. Taylor, X. Wu, R. Stevens, Prophecy: an infrastructure for performance analysis and modeling of parallel and grid applications, *ACM SIGMETRICS Performance Evaluation Review* 30 (4) (2003) 13–18.
- [11] M. Khan, Y. Jin, M. Li, Y. Xiang, C. Jiang, Hadoop Performance Modeling for Job Estimation and Resource Provisioning, *IEEE Transactions on Parallel and Distributed Systems* 27 (2) (2016) 441–454.

- [12] A. Ganapathi, H. A. Kuno, U. Dayal, J. L. Wiener, A. Fox, M. I. Jordan, D. A. Patterson, Predicting Multiple Metrics for Queries, ICDE.
- [13] D. Tsafirir, Y. Etsion, D. G. Feitelson, Backfilling Using System-Generated Predictions Rather than User Runtime Estimates., IEEE Transactions on Parallel and Distributed Systems 18 (6) (2007) 789–803.
- [14] E. Agullo, O. Beaumont, L. Eyraud-Dubois, S. Kumar, Are Static Schedules so Bad? A Case Study on Cholesky Factorization, IPDPS (2016) 1021–1030.
- [15] É. Gaussier, D. Glesser, V. Reis, D. Trystram, Improving backfilling by using machine learning to predict running times., SC (2015) 64–10.
- [16] W. Tang, N. Desai, D. Buettner, Z. Lan, Analyzing and adjusting user runtime estimates to improve job scheduling on the Blue Gene/P, in: Proceedings of the 2010 IEEE International Symposium on Parallel and Distributed Processing, IPDPS 2010, Illinois Institute of Technology, Chicago, United States, IEEE, 2010, pp. 1–11.
- [17] J. M. Ramírez-Alcaraz, A. Tchernykh, R. Yahyapour, U. Schwiegelshohn, A. Quezada-Pina, J. L. González-García, A. Hiraes-Carbajal, Job Allocation Strategies with User Run Time Estimates for Online Scheduling in Hierarchical Grids, Journal of Grid Computing 9 (1) (2011) 95–116.
- [18] F. E. Harrell, Regression Modeling Strategies, With Applications to Linear Models, Logistic and Ordinal Regression, and Survival Analysis, Springer, 2015.
- [19] W. W. S. Wei, Time Series Analysis, Univariate and Multivariate Methods, Pearson Addison Wesley, 2006.
- [20] J. Dean, S. Ghemawat, MapReduce: simplified data processing on large clusters, Communications of the ACM 51 (1) (2008) 107–113.
- [21] A. Gandomi, M. Haider, Beyond the hype: Big data concepts, methods, and analytics, International Journal of Information Management 35 (2) (2015) 137–144.
- [22] I. A. T. Hashem, I. Yaqoob, N. B. Anuar, S. Mokhtar, A. Gani, S. U. Khan, The rise of "big data" on cloud computing - Review and open research issues., Inf. Syst. 47 (2015) 98–115.
- [23] S. Venkataraman, Z. Yang, M. J. Franklin, B. Recht, I. Stoica, Ernest - Efficient Performance Prediction for Large-Scale Advanced Analytics, in: 13th USENIX Symposium on Networked Systems Design and Implementation, 2016.
- [24] C. S. Liew, M. P. Atkinson, M. Galea, T. F. Ang, P. Martin, J. I. V. Hemert, Scientific Workflows: Moving Across Paradigms, ACM Computing Surveys (CSUR) 49 (4) (2017) 66–39.
- [25] S. Pillana, I. Brandic, S. Benkner, A survey of the state of the art in performance modeling and prediction of parallel and distributed computing systems, International Journal of Computational Intelligence Research 4 (1).
- [26] S. Seneviratne, S. Witharana, A survey on methodologies for runtime prediction on grid environments, in: International Conference on Information and Automation for Sustainability, The University of Sydney, Sydney, Australia, 2014.
- [27] B. Jennings, R. Stadler, Resource Management in Clouds - Survey and Research Challenges., J. Network Syst. Manage. 23 (3) (2014) 567–619.
- [28] M. B. Qureshi, M. M. Dehnavi, N. Min-Allah, M. S. Qureshi, H. Husain, I. Rentifis, N. Tziritas, T. Loukopoulos, S. U. Khan, C.-Z. Xu, A. Y. Zomaya, Survey on Grid Resource Allocation Mechanisms, Journal of Grid Computing 12 (2) (2014) 399–441.
- [29] F. Hutter, L. Xu, H. H. Hoos, K. Leyton-Brown, Algorithm runtime prediction: Methods & evaluation, Artificial Intelligence 206 (2014) 79–111.
- [30] R. Osman, W. J. Knottenbelt, Database system performance evaluation models: A survey, Performance Evaluation 69 (10) (2012) 471–493.
- [31] P. Barham, B. Dragovic, K. Fraser, S. Hand, T. Harris, A. Ho, R. Neugebauer, I. Pratt, A. Warfield, P. Barham, B. Dragovic, K. Fraser, S. Hand, A. Ho, R. Neugebauer, I. Pratt, A. Warfield, Xen and the art of virtualization, Vol. 37, ACM, 2003.
- [32] M. Hovestadt, O. Kao, A. Keller, A. Streit, Scheduling in HPC Resource Management Systems: Queuing vs. Planning, in: Computer Performance Engineering, Springer Berlin Heidelberg, Berlin, Heidelberg, 2003, pp. 1–20.
- [33] Y. Zhu, L. M. Ni, A survey on grid scheduling systems, Tech. rep., Shanghai Jiao Tong University (2003).
- [34] C. Delimitrou, C. Kozyrakis, Quasar: resource-efficient and QoS-aware cluster management., ASPLOS (2014) 127–144.
- [35] R. Ferreira da Silva, R. Filgueira, I. Pietri, M. Jiang, R. Sakellariou, E. Deelman, A characterization of workflow management systems for extreme-scale applications, Future Generation Computing Systems.
- [36] H. Topcuoglu, S. Hariri, M.-Y. Wu, Performance-effective and low-complexity task scheduling for heterogeneous computing, IEEE Transactions on Parallel and Distributed Systems 13 (3) (2002) 260–274.
- [37] J. K. Lenstra, A. H. G. Rinnooy Kan, P. Brucker, Complexity of Machine Scheduling Problems, in: Studies in Integer Programming, Elsevier, 1977, pp. 343–362.
- [38] F. Dong, S. G. Akl, Scheduling algorithms for grid computing: State of the art and open problems, Tech. Rep. 2006-504, Queen's University (2006).
- [39] S. Smanchat, K. Viriyapant, Taxonomies of workflow scheduling problem and techniques in the cloud, Future Generation Computing Systems 52 (2015) 1–12.
- [40] Z.-H. Zhan, X.-F. Liu, Y.-J. Gong, J. Zhang, H. S.-H. Chung, Y. Li, Cloud Computing Resource Scheduling and a Survey of Its Evolutionary Approaches, ACM Computing Surveys (CSUR) 47 (4) (2015) 1–33.
- [41] R. V. Lopes, D. Menascé, A Taxonomy of Job Scheduling on Distributed Computing Systems, IEEE Transactions on Parallel and Distributed Systems 27 (12) (2016) 3412–3428.
- [42] A. W. Mu'alem, D. G. Feitelson, Utilization, Predictability, Workloads, and User Runtime Estimates in Scheduling the IBM SP2 with Backfilling, IEEE Transactions on Parallel and Distributed Systems 12 (6) (2001) 529–543.
- [43] Y.-K. Kwok, I. Ahmad, Static scheduling algorithms for allocating directed task graphs to multiprocessors, ACM Computing Surveys 31 (4) (1999) 406–471.
- [44] J. Blythe, S. Jain, E. Deelman, Y. Gil, K. Vahi, A. Mandal, K. Kennedy, Task scheduling strategies for workflow-based applications in grids., CCGRID (2005) 759–767.
- [45] H. Li, D. Groep, J. Templon, L. Wolters, Predicting job start times on clusters, in: International Symposium on Cluster Computing and the Grid, Leiden University, Leiden, Netherlands, IEEE, 2004, pp. 301–308.
- [46] H. Li, D. Groep, L. Wolters, Mining performance data for metascheduling decision support in the Grid, Future Generation Computer Systems 23 (1) (2007) 92–99.
- [47] R. L. F. Cunha, E. R. Rodrigues, L. P. Tizzei, M. A. S. Netto, Job placement advisor based on turnaround predictions for HPC hybrid clouds, Future Generation Computer Systems 67 (2017) 35–46.
- [48] J. Liu, E. Pacitti, P. Valduriez, M. Mattoso, A Survey of Data-Intensive Scientific Workflow Management, Journal of Grid Computing.
- [49] B. C. Lee, D. M. Brooks, B. R. de Supinski, M. Schulz, K. Singh, S. A. McKee, Methods of Inference and Learning for Performance Modeling of Parallel Applications, PPOPP.
- [50] W. Pfeiffer, N. J. Wright, Modeling and predicting application performance on parallel computers using HPC challenge benchmarks, IEEE International Symposium on Parallel and Distributed Processing (2008) 1–12.
- [51] M. Kuperberg, K. Krogmann, R. Reussner, Performance Prediction for Black-Box Components Using Reengineered Parametric Behaviour Models, in: Component-Based Software Engineering, Springer Berlin Heidelberg, Berlin, Heidelberg, 2008, pp. 48–63.
- [52] R. Ferreira da Silva, G. Juve, M. Rynge, E. Deelman, M. Livny, Online Task Resource Consumption Prediction for Scientific Workflows, Parallel Processing Letters 25 (03) (2015) 1541003.
- [53] S. Vazhkudai, J. M. Schopf, Using Regression Techniques to Predict Large Data Transfers, International Journal of High Performance Computing Applications 17 (3) (2003) 249–268.
- [54] E. Arslan, K. Guner, T. Kosar, HARP - predictive transfer optimization based on historical analysis and real-time probing., in: the International Conference for High Performance Computing, Networking, Storage and Analysis, 2016.
- [55] M. Faerman, A. Su, R. Wolski, F. Berman, Adaptive performance prediction for distributed data-intensive applications, in: ACM International Conference on Supercomputing, ACM Press, New York, New York, USA, 1999.
- [56] R. Wolski, N. T. Spring, J. Hayes, The network weather service: a distributed resource performance forecasting service for metacomputing., Future Generation Computer Systems (1999) 757–768.

- [57] D. M. Swamy, R. Wolski, Multivariate resource performance forecasting in the network weather service., SC.
- [58] W. Gropp, E. Lusk, A. Skjellum, Using MPI: portable parallel programming with the message-passing interface, MIT Press, 1999.
- [59] A. B. Downey, Predicting Queue Times on Space-Sharing Parallel Computers., IPPS (1997) 209–218.
- [60] R. Gibbons, A Historical Application Profiler for Use by Parallel Schedulers., Job Scheduling Strategies for Parallel Processing 1291 (Chapter 3) (1997) 58–77.
- [61] T. Miu, P. Missier, Predicting the Execution Time of Workflow Activities Based on Their Input Features, in: High Performance Computing, Networking, Storage and Analysis, IEEE, 2012, pp. 64–72.
- [62] F. Ahmad, S. T. Chakradhar, A. Raghunathan, T. N. Vijaykumar, Tarazu - Optimizing MapReduce on Heterogeneous Clusters, *AcM Sigplan Notices* 47 (4) (2012) 61–74.
- [63] D. Cheng, J. Rao, Y. Guo, X. Zhou, Improving MapReduce performance in heterogeneous environments with adaptive task tuning, in: the 15th International Middleware Conference, ACM Press, New York, New York, USA, 2014, pp. 97–108.
- [64] M. Bux, U. Leser, DynamicCloudSim: Simulating heterogeneity in computational clouds, *Future Generation Computer Systems*.
- [65] B. J. Barnes, B. Rountree, D. K. Lowenthal, J. Reeves, B. de Supinski, M. Schulz, A regression-based approach to scalability prediction, in: International Conference on Supercomputing, ACM, New York, New York, USA, 2008, pp. 368–377.
- [66] P. A. Dinda, Online prediction of the running time of tasks, 10th IEEE International Symposium on High Performance Distributed Computing (2002) 383–394.
- [67] Y. Koh, R. C. Knauerhase, P. Brett, M. Bowman, Z. Wen, C. Pu, An Analysis of Performance Interference Effects in Virtual Environments., *ISPASS* (2007) 200–209.
- [68] J. Zhao, H. Cui, J. Xue, X. Feng, Predicting Cross-Core Performance Interference on Multicore Processors with Regression Analysis, *IEEE Transactions on Parallel and Distributed Systems* 27 (5) (2016) 1443–1456.
- [69] T. Hastie, R. Tibshirani, J. Friedman, The Elements of Statistical Learning, Data Mining, Inference, and Prediction, Second Edition, Springer Science & Business Media, 2009.
- [70] A. B. Downey, D. G. Feitelson, The elusive goal of workload characterization., *SIGMETRICS Performance Evaluation Review* () 26 (4) (1999) 14–29.
- [71] L. T. Yang, X. Ma, F. Mueller, Cross-Platform Performance Prediction of Parallel Applications Using Partial Execution., SC.
- [72] N. H. Kapadia, J. A. B. Fortes, C. E. Brodley, Predictive Application-Performance Modeling in a Computational Grid Environment., in: International Symposium on High-Performance Distributed Computing, 1999.
- [73] G. Marin, J. Mellor-Crummey, Cross-Architecture Performance Predictions for Scientific Applications Using Parameterized Models, *ACM SIGMETRICS Performance Evaluation Review* 32 (1) (2004) 2.
- [74] K. Hoste, A. Phansalkar, L. Eeckhout, A. Georges, L. K. John, K. De Bosschere, Performance Prediction based on Inherent Program Similarity, in: International Conference on Extending Database Technology, ACM Press, New York, New York, USA, 2006, p. 114.
- [75] N. J. Yadwadkar, B. Hariharan, J. E. Gonzalez, B. Smith, R. H. Katz, Selecting the Best VM Across Multiple Public Clouds: A Data-Driven Performance Modeling Approach, ACM, New York, New York, USA, 2017.
- [76] G. Chatzopoulos, A. Dragojevic, R. Guerraoui, ESTIMA - extrapolating scalability of in-memory applications., *PPOPP 12-16-March-2016* (8) (2017) 1–11.
- [77] S. Govindan, J. Liu, A. Kansal, A. Sivasubramaniam, Cuanta - quantifying effects of shared on-chip resource interference for consolidated virtual machines., *Symposium on Cloud Computing* (2011) 22–14.
- [78] Y. Gao, H. Rong, J. Z. Huang, Adaptive grid job scheduling with genetic algorithms, *Future Generation Computer Systems* 21 (1) (2005) 151–161. doi:10.1016/j.future.2004.09.033.
- [79] M. V. Devarakonda, R. K. Iyer, Predictability of Process Resource Usage: A Measurement-Based Study on UNIX, *IEEE Transactions on Software Engineering* 15 (12) (1989) 1579–1586. doi:10.1109/32.58769.
- [80] P. A. Dinda, A prediction-based real-time scheduling advisor, in: Proceedings - International Parallel and Distributed Processing Symposium, IPDPS 2002, Northwestern University, Evanston, United States, IEEE, 2002, pp. 88–95.
- [81] Y. Zhang, W. Sun, Y. Inoguchi, Predict task running time in grid environments based on CPU load predictions, *Future Generation Computer Systems* 24 (6) (2008) 489–497.
- [82] R. Wolski, Dynamically forecasting network performance using the Network Weather Service, *Cluster Computing*, IEEE International Conference on 1 (1) (1998) 119–132.
- [83] R. Wolski, N. Spring, J. Hayes, Predicting the CPU availability of time-shared Unix systems on the computational grid, *Cluster Computing*, IEEE International Conference on 3 (4) (2000) 293–301.
- [84] L. Yang, I. T. Foster, J. M. Schopf, Homeostatic and Tendency-Based CPU Load Predictions, IPDPS.
- [85] L. Yang, J. M. Schopf, I. T. Foster, Conservative Scheduling, in: the 2003 ACM/IEEE conference, ACM Press, New York, New York, USA, 2003, p. 31.
- [86] Y. Wu, K. H. 0001, Y. Yuan, W. Zheng, Adaptive Workload Prediction of Grid Performance in Confidence Windows., *IEEE Trans. Parallel Distrib. Syst.*
- [87] C. Bienia, Benchmarking Modern Multiprocessors, Ph.D. thesis, Princeton University, Princeton University (2011).
- [88] J. Mars, L. Tang, R. Hundt, K. Skadron, M. Lou Soffa, Bubble-Up - increasing utilization in modern warehouse scale computers via sensible co-locations., MICRO.
- [89] W. Smith, Prediction Services for Distributed Computing, IPDPS (2007) 1–10.
- [90] S. Puma, W.-c. Feng, P. Phunchongharn, S. Chapeland, T. Achalakul, A runtime estimation framework for ALICE, *Future Generation Computer Systems* 72 (2017) 65–77.
- [91] O. Alipourfard, H. H. Liu, J. Chen, S. Venkataraman, M. Yu, M. Zhang, CherryPick - Adaptively Unearthing the Best Cloud Configurations for Big Data Analytics., NSDI.
- [92] M. Dobber, R. van der Mei, G. Koole, A prediction method for job runtimes on shared processors: Survey, statistical analysis and new avenues, *Performance Evaluation* 64 (2007) 755–781.
- [93] H. A. Sanjay, S. Vadhiyar, Performance modeling of parallel applications for grid scheduling, *Journal of Parallel and Distributed Computing* 68 (8) (2008) 1135–1145.
- [94] D. G. Feitelson, B. Nitzberg, Job characteristics of a production parallel scientific workload on the NASA Ames iPSC/860, *Economics of Grids* 949 (Chapter 19) (1995) 337–360.
- [95] J. Brevik, D. Nurmi, R. Wolski, Predicting bounds on queuing delay for batch-scheduled parallel machines., *PPOPP* (2006) 110–118.
- [96] D. Nurmi, J. Brevik, R. Wolski, QBETS: queue bounds estimation from time series, in: International Conference on Job scheduling strategies for parallel processing, California State University Long Beach, Springer-Verlag, 2007.
- [97] W. Smith, V. E. Taylor, I. T. Foster, Using Run-Time Predictions to Estimate Queue Wait Times and Improve Scheduler Performance., *Job Scheduling Strategies for Parallel Processing* (1999) 202–219.
- [98] F. Nadeem, T. Fahringer, Optimizing execution time predictions of scientific workflow applications in the Grid through evolutionary programming, *Future Generation Computer Systems* 29 (4) (2013) 926–935.
- [99] R. M. Piro, A. Guarise, G. Patania, A. Werbrouck, Using historical accounting information to predict the resource usage of grid jobs, *Future Generation Computer Systems* 25 (5) (2009) 499–510.
- [100] O. Sonmez, N. Yigitbasi, A. Iosup, D. Epema, Trace-based evaluation of job runtime and queue wait time predictions in grids, in: ACM International Symposium on High Performance Distributed Computing, ACM Press, New York, New York, USA, 2009, pp. 111–120.
- [101] T. Dwyer, A. Fedorova, S. Blagodurov, M. Roth, F. Gaud, J. Pei, A practical method for estimating performance degradation on multicore processors, and its application to HPC workloads, in: High Performance Computing, Networking, Storage and Analysis, International Conference on, Simon Fraser University, IEEE Computer Society Press, 2012, pp. 1–11.
- [102] L. Rokach, O. Maimon, Data Mining with Decision Trees - Theory and Applications. 2nd Edition, Series in Machine Perception and Artificial Intelligence 81.

- [103] K. A. Asanovic, R. A. Bodik, B. C. A. Catanzaro, J. J. A. Gebis, P. A. Husbands, K. A. Keutzer, D. A. A. Patterson, W. L. A. Plishker, J. A. Shalf, S. W. A. Williams, K. A. A. Yelick, The Landscape of Parallel Computing Research: A View from Berkeley, Tech. Rep. UCB/EECS-2006-183 (2006).
- [104] S. Che, M. Boyer, J. Meng, D. Tarjan, J. W. Sheaffer, S.-H. Lee, K. Skadron, Rodinia: A benchmark suite for heterogeneous computing, in: 2009 IEEE International Symposium on Workload Characterization (IISWC), IEEE, 2009, pp. 44–54.
- [105] D. H. Bailey, E. Barszcz, J. T. Barton, D. S. Browning, R. L. Carter, L. Dagum, R. A. Fatoohi, P. O. Frederickson, T. A. Lasinski, R. S. Schreiber, H. D. Simon, V. Venkatakrisnan, S. K. Weeratunga, The Nas Parallel Benchmarks, *International Journal of Supercomputer Applications and High Performance Computing* 5 (3) (1991) 63–73.
- [106] A. Kaiser, TORCH Computational Reference Kernels - A Testbed for Computer Science Research, Lawrence Berkeley National Laboratory.
- [107] E. Ipek, B. R. de Supinski, M. Schulz, S. A. McKee, An Approach to Performance Prediction for Parallel Applications, in: *Computer Performance Engineering*, Springer Berlin Heidelberg, Berlin, Heidelberg, 2005, pp. 196–205.
- [108] N. R. Draper, H. Smith, *Applied regression analysis, Wiley series in probability and mathematical statistics*, Wiley, 1998.
- [109] Z. Liu, P. Balaprakash, R. Kettimuthu, I. T. Foster, Explaining Wide Area Data Transfer Performance, *HPDC* (2017) 167–178.
- [110] M. Swamy, R. Wolski, Multivariate resource performance forecasting in the network weather service, in: *SC '02: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, IEEE, 2002.
- [111] S. Vazhkudai, J. Schopf, Using regression techniques to predict large data transfers, *The International Journal of High Performance Computing Applications* 17 (3) (2003) 249–268.
- [112] Y. Qiao, J. A. Skicewicz, P. A. Dinda, An Empirical Study of the Multi-scale Predictability of Network Traffic, in: *International Symposium on High-Performance Distributed Computing*, 2004.
- [113] J. Tan, X. Pan, S. Kavulya, R. Gandhi, P. Narasimhan, SALSA - Analyzing Logs as StAte Machines, in: *USENIX Workshop on the Analysis of System Logs*, 2008.
- [114] J. L. Henning, SPEC CPU2006 benchmark descriptions, *ACM SIGARCH Computer Architecture News* 34 (4) (2006) 1–17.
- [115] Y. Yuan, Y. Wu, G. Yang, W. Zheng, Adaptive hybrid model for long term load prediction in computational grid, in: *IEEE International Symposium on Cluster Computing and the Grid*, Tsinghua University, Beijing, China, 2008, pp. 340–347.
- [116] G. Juve, B. Tovar, R. Ferreira da Silva, D. Król, D. Thain, E. Deelman, W. E. Allcock, M. Livny, Practical Resource Monitoring for Robust High Throughput Computing., in: *IEEE International Conference on Cluster Computing*, 2015, pp. 650–657.
- [117] K. Hoste, L. Eeckhout, Microarchitecture-Independent Workload Characterization., *IEEE micro* 27 (3) (2007) 63–72.
- [118] M. Laurenzano, M. M. Tikir, L. Carrington, A. Snaveley, PEBIL - Efficient static binary instrumentation for Linux., *ISPASS*.
- [119] A. R. Bernat, B. P. Miller, Anywhere, any-time binary instrumentation., in: *SIGPLAN-SIGSOFT workshop on Program analysis for software tools*, ACM Press, New York, New York, USA, 2011, p. 9.
- [120] D. G. Feitelson, D. Tsafir, D. Krakov, Experience with using the Parallel Workloads Archive, *Journal of Parallel and Distributed Computing* 74 (10) (2014) 2967–2982.
- [121] C. Reiss, A. Tumanov, G. R. Ganger, R. H. Katz, M. A. Kozuch, Heterogeneity and dynamicity of clouds at scale, in: *ACM Symposium on Cloud Computing*, ACM Press, New York, New York, USA, 2012, pp. 1–13.
- [122] J. Dongarra, K. London, S. Moore, P. Mucci, D. Terpstra, Using papi for hardware performance monitoring on linux systems, in: *LCI '01: Proceedings of the 2nd Conference on Linux Clusters: The HPC Revolution*, Vol. 5, Linux Clusters Institute, 2001.
- [123] J. Treibig, G. Hager, G. Wellein, LIKWID: A Lightweight Performance-Oriented Tool Suite for x86 Multicore Environments, in: *International Conference on Parallel Processing*, IEEE, 2010, pp. 207–216.
- [124] F. Olsson, A literature survey of active machine learning in the context of natural language processing, Tech. rep., Swedish Institute of Computer Science (2009).
- [125] A. Sarkar, J. Guo, N. Siegmund, S. Apel, K. Czarnecki, Cost-Efficient Sampling for Performance Prediction of Configurable Systems, 2015 30th IEEE/ACM International Conference on Automated Software Engineering (ASE) (2015) 342–352.
- [126] A. Singh, A. Rao, S. Purawat, I. Altintas, A machine learning approach for modular workflow performance prediction, in: *the 12th Workshop, ACM Press, New York, New York, USA, 2017*, pp. 1–11.
- [127] E. Rodrigues, R. L. F. Cunha, M. A. S. Netto, M. Spriggs, Helping HPC Users Specify Job Memory Requirements via Machine Learning, in: *Proceedings of HUST 2016: 3rd International Workshop on HPC User Support Tools - Held in conjunction with SC 2016: The International Conference for High Performance Computing, Networking, Storage and Analysis*, IBM Research, Yorktown Heights, United States, IEEE, 2017, pp. 6–13.
- [128] C. Reiss, Understanding Memory Configurations for In-Memory Analytics, Ph.D. thesis (2016).
- [129] G. Juve, Resource management for scientific workflows, Ph.D. thesis, University of Southern California, University of Southern California (Jan. 2012).
- [130] E. Pennisi, Will Computers Crash Genomics?, *science* 331 (6018) (2011) 666–668.
- [131] T. Sherwood, S. Sair, B. Calder, Phase Tracking and Prediction., *ISCA* (2003) 336–347.
- [132] W. Zhang, J. Li, Y. Li, H. Chen, Multilevel Phase Analysis, *ACM Transactions on Embedded Computing Systems (TECS)* 14 (2) (2015) 31–29.
- [133] R. C. Coutinho, L. Drummond, Y. Frota, Optimizing virtual machine allocation for parallel scientific workflows in federated clouds, *Future Generation Computer Systems*.
- [134] N. Mishra, J. D. Lafferty, H. Hoffmann, ESP: A Machine Learning Approach to Predicting Application Interference, *ICAC* (2017) 125–134.
- [135] M. Jeon, Y. He, H. Kim, S. Elnikety, S. Rixner, A. L. Cox, TPC: Target-driven parallelism combining prediction and correction to reduce tail latency in interactive services, in: *International Conference on Architectural Support for Programming Languages and Operating Systems - ASPLOS*, Microsoft Research, Redmond, United States, ACM Press, New York, New York, USA, 2016, pp. 129–141.