

Using Run-Time Predictions to Estimate Queue Wait Times and Improve Scheduler Performance

Warren Smith^{1,2}, Valerie Taylor², and Ian Foster¹

¹ {wsmith, foster}@mcs.anl.gov

Mathematics and Computer Science Division

Argonne National Laboratory

Argonne, IL 60439

<http://www.mcs.anl.gov>

² taylor@ece.nwu.edu

Electrical and Computer Engineering Department

Northwestern University

Evanston, IL 60208

<http://www.ece.nwu.edu>

Abstract. On many computers, a request to run a job is not serviced immediately but instead is placed in a queue and serviced only when resources are released by preceding jobs. In this paper, we build on run-time prediction techniques that we developed in previous research to explore two problems. The first problem is to predict how long applications will wait in a queue until they receive resources. We develop run-time estimates that result in more accurate wait-time predictions than other run-time prediction techniques. The second problem we investigate is improving scheduling performance. We use run-time predictions to improve the performance of the least-work-first and backfill scheduling algorithms. We find that using our run-time predictor results in lower mean wait times for the workloads with higher offered loads and for the backfill scheduling algorithm.

1 Introduction

On many high-performance computers, a request to execute an application is not serviced immediately but instead is placed in a queue and serviced only when resources are released by running applications. We examine two separate problems in this environment. First, we predict how long applications will wait until they execute. These estimates of queue wait times are useful to guide resource selection when several systems are available [7], to co-allocate resources from multiple systems [2], to schedule other activities, and so forth. Our technique for predicting queue wait times is to use predictions of application execution times along with the scheduling algorithms to simulate the actions made by a scheduler and determine when applications will begin to execute.

We performed queue wait-time prediction and scheduling experiments using four workloads and three scheduling algorithms. The workloads were recorded

from an IBM SP at Argonne National Laboratory, an IBM SP at the Cornell Theory Center, and an Intel Paragon at the San Diego Supercomputing Center. The scheduling algorithms are first-come first-served (FCFS), least-work-first (LWF), and backfill. We find that there is a built-in error of 34 to 43 percent when predicting wait times of the LWF algorithm and a smaller built-in error of 3 to 4 percent for the backfill algorithm. We also find that more accurate run-time predictions result in more accurate wait-time predictions. Specifically, using our run-time prediction technique instead of maximum run times or the run-time prediction techniques of Gibbons [8] or Downey [3] improves run-time prediction error by 39 to 92 percent and this improves wait-time prediction performance by 16 to 89 percent.

Second, we improve the performance of the LWF and backfill scheduling algorithms by using our run-time predictions. These algorithms use run-time predictions when making scheduling decisions, and we therefore expect that more accurate run-time predictions will improve scheduling performance. Using our run-time predictions using the same four workloads describe above, we find that the accuracy of the run-time predictions has a minimal effect on the utilization of the systems we are simulating. We also find that using our run-time predictors result in mean wait times that are within 22 percent of the mean wait times that are obtained if the scheduler exactly knows the run times of all of the applications. When comparing the different predictors, our run-time predictor results in 2 to 67 percent smaller mean wait times for the workload with the highest offered load. No prediction technique clearly outperforms the other techniques when the offered load is low.

The next section summarizes our approach to predicting application run times and the approaches of other researchers. Section 3 describes our queue wait-time prediction technique and presents performance data. Section 4 presents the results of using our run-time predictions in the LWF and backfill scheduling algorithms. Section 5 presents our conclusions.

2 Predicting Application Run Times

This section briefly describes our prediction technique, discusses the scheduling algorithms used with this work, and reviews the run-time prediction techniques of two other researchers, Gibbons [8] and Downey [3]. For further details, see [12].

2.1 Our Run-Time Prediction Technique

Our general approach to predicting application run times is to derive run-time predictions from historical information of previous “similar” runs. This approach is based on the observation [12,3,5,8] that similar applications are more likely to have similar run times than applications that have nothing in common. We address the issues of how to define similar and how to generate predictions from similar past applications.

Defining Similarity A difficulty in developing prediction techniques based on similarity is that two jobs can be compared in many ways. For example, we can compare the application name, submitting user name, executable arguments, submission time, and number of nodes requested. In this work, we are restricted to those values recorded in workload traces obtained from various supercomputer centers. However, because the techniques that we propose are based on the automatic discovery of efficient similarity criteria, we believe that they will apply even if quite different information is available.

Table 1. Characteristics of the trace data used in our studies.

Workload Name	System	Number of Nodes	Location	Number of Requests	Mean Run Time (minutes)
ANL ¹	IBM SP2	120	ANL	7994	97.75
CTC	IBM SP2	512	CTC	13217	171.14
SDSC95	Intel Paragon	400	SDSC	22885	108.21
SDSC96	Intel Paragon	400	SDSC	22337	166.98

The workload traces that we consider are described in Table 1; they originate from Argonne National Laboratory (ANL), the Cornell Theory Center (CTC), and the San Diego Supercomputer Center (SDSC). Table 2 summarizes the information provided in these traces. Text in a field indicates that a particular trace contains the information in question; in the case of “Type,” “Queue,” or “Class” the text specifies the categories in question.

The general approach to defining similarity that we as well as Downey and Gibbons take is to use characteristics such as those presented in Table 2 to define *templates* that identify a set of *categories* to which jobs can be assigned. For example, the template (q, u) specifies that jobs are to be partitioned by *queue* and *user*; on the SDSC Paragon, this template generates categories such as $(q16m, wsmith)$, $(q64l, wsmith)$, and $(q16m, foster)$.

We find that using characteristics 1–6 of Table 2 in the manner just described works reasonably well. On the other hand, the number of nodes is an essentially continuous parameter, and so we prefer to introduce an additional parameter into our templates, namely, a “node range size” that defines what ranges of requested number of nodes are used to decide whether applications are similar. For example, the template $(u, n=4)$ specifies a node range size of 4 and generates categories $(wsmith, 1-4 \text{ nodes})$ and $(wsmith, 5-8 \text{ nodes})$.

Once a set of templates has been defined (using a search process described later), we can categorize a set of applications (e.g., the workloads of Table 1)

¹ Because of an error when the trace was recorded, the ANL trace does not include one-third of the requests actually made to the system. To compensate, we reduced the number of nodes on the machine from 120 to 80 when performing simulations.

Table 2. Characteristics recorded in workloads. The column “Abbr” indicates abbreviations used in subsequent discussion.

	Abbr	Characteristic	Argonne	Cornell	SDSC
1	t	Type	batch, interactive	serial, parallel, pvm3	
2	q	Queue			29 to 35 queues
3	c	Class		DSI/PIOFS	
4	u	User	Y	Y	Y
5	s	Loadleveler script		Y	
6	e	Executable	Y		
7	a	Arguments	Y		
8	na	Network adaptor		Y	
9	n	Number of nodes	Y	Y	Y
10		Maximum run time	Y	Y	
11		Submission time	Y	Y	Y
12		Start time	Y	Y	Y
13		Run time	Y	Y	Y

by assigning each application to those categories that match its characteristics. Categories need not be disjoint, and hence the same job can occur in several categories. If two jobs fall into the same category, they are judged similar; those that do not coincide in any category are judged dissimilar.

Generating Predictions We now consider the question of how we generate run-time predictions. The input to this process is a set of templates and a workload for which run-time predictions are required. In addition to the characteristics described in the preceding section, running time, maximum history, type of data to store, and prediction type are also defined for each template. The running time is how long an application has been running when a prediction is made; we use this characteristic by forming a prediction from a category using only the data points that have an execution time less than this running time. The maximum history indicates the maximum number of data points to store in each category generated from a template. The type of data is either an actual run time or a relative run time. A relative run time incorporates information about user-supplied run time estimates by storing the ratio of the actual run time to the user-supplied estimate (the maximum run times provided by the ANL and CTC workloads). The prediction type determines how a run-time prediction is made from the data in each category generated from a template. In our previous work, we considered four prediction types: a mean, a linear regression, an inverse regression, and a logarithmic regression [13,4]. We found that the mean is the single best predictor, so this study uses only the mean to form predictions.

The output from this process is a set of run-time predictions and associated confidence intervals. (A confidence interval is an interval centered on the run-time prediction within which the actual run time is expected to appear some specified percentage of the time.) The basic algorithm is described below and comprises three steps: initialization, prediction, and incorporation of historical information.

1. Define T , the set of templates to be used, and initialize C , the (initially empty) set of categories.
2. At the time each application a begins to execute:
 - (a) Apply the templates in T to the characteristics of a to identify the categories C_a into which the application may fall.
 - (b) Eliminate from C_a all categories that are not in C or that cannot provide a valid prediction (i.e., do not have enough data points).
 - (c) For each category remaining in C_a , compute a run-time estimate and a confidence interval for the estimate.
 - (d) If C_a is not empty, select the estimate with the smallest confidence interval as the run-time prediction for the application.
3. At the time each application a completes execution:
 - (a) Identify the set C_a of categories into which the application falls. These categories may or may not exist in C .
 - (b) For each category $c_i \in C_a$
 - i. If $c_i \notin C$, create c_i in C .
 - ii. If $|c_i| = \text{maximum history}(c_i)$, remove the oldest point in c_i .
 - iii. Insert a into c_i .

Note that steps 2 and 3 operate asynchronously, since historical information for a job cannot be incorporated until the job finishes. Hence, our algorithm suffers from an initial ramp-up during which there is insufficient information in C to make predictions. This deficiency could be corrected by using a training set to initialize C .

The use of maximum histories in step 3(b) of our algorithm allows us to control the amount of historical information used when making predictions and the amount of storage space needed to store historical information. A small maximum history means that less historical information is stored, and hence only more recent events are used to make predictions.

Template Definition and Search We use search techniques to identify good templates for a particular workload; this approach is in contrast to the strategies of Gibbons and Downey, who use a fixed set of templates. While the number of application characteristics included in our traces is relatively small, the fact that effective template sets may contain many templates means that an exhaustive search is impractical. Our previous work compared greedy and genetic algorithm searches and found that genetic algorithm searches outperform greedy searches. Therefore, we use only genetic algorithm searches in this work.

Genetic algorithms are a probabilistic technique for exploring large search spaces, in which the concept of cross-over from biology is used to improve efficiency relative to purely random search [10]. A genetic algorithm evolves individuals over a series of generations. The process for each generation consists of evaluating the fitness of each individual in the population, selecting which individuals will be mated to produce the next generation, mating the individuals, and mutating the resulting individuals to produce the next generation. The process then repeats until a stopping condition is met. The stopping condition we use is that a fixed number of generations have been processed. There are many different variations to this process, and we will next describe the variations we used.

Our individuals represent template sets. Each template set consists of between 1 and 10 templates, and we encode the following information in binary form for each template:

1. Whether a mean or one of the three regressions is used to produce a prediction.
2. Whether absolute or relative run times are used.
3. Whether each of the binary characteristics associated with the workload in question is enabled.
4. Whether node information should be used and, if so, the range size from 1 to 512 in powers of 2.
5. Whether the amount of history stored in each category should be limited and, if so, the limit between 2 and 65536 in powers of 2.

A fitness function is used to compute the fitness of each individual and therefore its chance to reproduce. The fitness function should be selected so that the most desirable individuals have higher fitness and produce more offspring, but the diversity of the population must be maintained by not giving the best individuals overwhelming representation in succeeding generations. In our genetic algorithm, we wish to minimize the prediction error and maintain a range of individual fitnesses regardless of whether the range in errors is large or small. The fitness function we use to accomplish this goal is

$$F_{min} + \frac{E_{max} - E}{E_{max} - E_{min}} \times (F_{max} - F_{min}),$$

where E is the error of the individual (template set), E_{min} and E_{max} are the minimum and maximum errors of individuals in the generation, and F_{min} and F_{max} are the desired minimum and maximum fitnesses desired. We choose $F_{max} = 4F_{min}$.

We use a common technique called stochastic sampling with replacement to select which individuals will mate to produce the next generation. In this technique, each parent is selected from the individuals by selecting individual i with probability $\frac{F_i}{\sum F}$.

The mating or crossover process is accomplished by randomly selecting pairs of individuals to mate and replacing each pair by their children in the new population. The crossover of two individuals proceeds in a slightly nonstandard way

because our chromosomes are not fixed length but a multiple of the number of bits used to represent each template. Two children are produced from each crossover by randomly selecting a template i and a position p in the template from the first individual $T_1 = t_{1,1}, \dots, t_{1,n}$ and randomly selecting a template j in the second individual $T_2 = t_{2,1}, \dots, t_{2,m}$ so that the resulting individuals will not have more than 10 templates. The new individuals are then $\tilde{T}_1 = t_{1,1}, \dots, t_{1,i-1}, n_1, t_{2,j+1}, \dots, t_{2,m}$ and $\tilde{T}_2 = t_{2,1} \dots t_{2,j-1}, n_2, t_{1,i+1}, \dots, t_{1,n}$. If there are b bits used to represent each template, n_1 is the first p bits of $t_{1,i}$ concatenated with the last $b - p$ bits of $t_{2,j}$, and n_2 is the first p bits of $t_{2,j}$ concatenated with the last $b - p$ bits of $t_{1,i}$.

In addition to using crossover to produce the individuals of the next generation, we also use a process called elitism whereby the best individuals in each generation survive unmutated to the next generation. We use crossover to produce all but two individuals for each new generation and use elitism to select the last two individuals for each new generation. The individuals resulting from the crossover process are mutated to help maintain a diversity in the population. Each bit representing the individuals is flipped with a probability of .01.

Run-Time Prediction Experiments We use run time predictions to predict queue wait times and improve the performance of scheduling algorithms. Therefore, we need to determine what workloads to search over to find the best template sets to use. We have already described four sets of trace data that were recorded from supercomputers. Next, we describe the three scheduling algorithms we consider.

We use the first-come first-served (FCFS), least-work-first (LWF), and backfill scheduling algorithms in this work. In the FCFS algorithm, applications are given resources in the order in which they arrive. The application at the head of the queue runs whenever enough nodes become free. The LWF algorithm also tries to execute applications in order, but the applications are ordered in increasing order using estimates of the amount of work (number of nodes multiplied by estimated wallclock execution time) the application will perform.

The backfill algorithm is a variant of the FCFS algorithm. The difference is that the backfill algorithm allows an application to run before it would in FCFS order if it will not delay the execution of applications ahead of it in the queue (those that arrived before it). When the backfill algorithm tries to schedule applications, it examines every application in the queue, in order of arrival time. If an application can run (there are enough free nodes and running the application will not delay the starting times of applications ahead of it in the queue), it is started. If an application cannot run, nodes are “reserved” for it at the earliest possible time. This reservation is only to make sure that applications behind it in the queue do not delay it; the application may actually start before the reservation time.

Each scheduling algorithm predicts application run times at different times when predicting queue wait times for the jobs in each trace. When predicting queue wait times, we predict the wait time of an application when it is sub-

mitted. A wait-time prediction in this case requires run-time predictions of all applications in the system so the run-time prediction workload contains predictions for all running and queued jobs every time an application is submitted. We insert data points for an application into our historical database as soon as each application completes. To try to find the optimal template set to use to predict execution times, we use a workload for each algorithm/trace pair and search over each of these 12 workloads separately.

When using run-time predictions while scheduling, run-time predictions are also made at different times for each algorithm/trace pair, and we attempt to find the optimal template sets to use for each pair. The FCFS algorithm does not use run-time predictions when scheduling, so we only consider the LWF and backfill algorithms here. For the LWF algorithm, all waiting applications are predicted whenever the scheduling algorithm attempts to start an application (when any application is enqueued or finishes). This occurs because the LWF algorithm needs to find the waiting application that will use the least work. For the backfill algorithm, all running and waiting applications are predicted whenever the scheduling algorithm attempts to start an application (when any application is enqueued or finishes).

We generate our run-time prediction workloads for scheduling using maximum run times as run-time predictions. We note that using maximum run times will produce predictions and insertions slightly different from those produced when the LWF and backfill algorithms use other run-time predictions. Nevertheless, we believe that these run-time prediction workloads are representative of the predictions and insertions that will be made when scheduling using other run-time predictors.

2.2 Related Work

Gibbons [8,9] also uses historical information to predict the run times of parallel applications. His technique differs from ours principally in that he uses a fixed set of templates and different characteristics to define templates. He uses the six templates/predictor combinations listed in Table 3. The running time (`runtime`) characteristic indicates how long an application has been executing when a prediction is made for the application. Gibbons produces predictions by examining categories derived from the templates listed in Table 3, in the order listed, until a category that can provide a valid prediction is found. This prediction is then used as the run-time prediction.

The set of templates listed in Table 3 results because Gibbons uses templates of (`u,e`), (`e`), and () with subtemplates in each template. The subtemplates add the characteristics `n` and `runtime`. Gibbons also uses the requested number of nodes slightly differently from the way we do: rather than having equal-sized ranges specified by a parameter, as we do, he defines the fixed set of exponential ranges 1, 2-3, 4-7, 8-15, and so on.

Another difference between Gibbons's technique and ours is how he performs a linear regression on the data in the categories (`u,e`), (`e`), and (). These categories are used only if one of their subcategories cannot provide a valid

Table 3. Templates used by Gibbons for run-time prediction.

Number	Template	Predictor
1	(u, e, n, rtime)	mean
2	(u, e)	linear regression
3	(e, n, rtime)	mean
4	(e)	linear regression
5	(n, rtime)	mean
6	()	linear regression

prediction. A weighted linear regression is performed on the mean number of nodes and the mean run time of each subcategory that contains data, with each pair weighted by the inverse of the variance of the run times in their subcategory.

Downey [3] uses a different technique to predict the execution time of parallel applications. His procedure is to categorize all applications in the workload, then model the cumulative distribution functions of the run times in each category, and finally use these functions to predict application run times. Downey categorizes applications using the queues that applications are submitted to, although he does state that other characteristics can be used in this categorization. In fact, Downey’s prediction technique within a category can be used with our technique for finding the best characteristics to use to categorize applications.

Downey observed that the cumulative distributions of the execution times of the jobs in the workloads he examined can be modeled relatively accurately by using a logarithmic function: $\beta_0 + \beta_1 \ln t$. Once the distribution functions are calculated, he uses two different techniques to produce a run-time prediction. The first technique uses the median lifetime given that an application has executed for a time units. If one assumes the logarithmic model for the cumulative distribution, this equation is

$$\sqrt{ae^{\frac{1.0-\beta_0}{\beta_1}}}.$$

The second technique uses the conditional average lifetime

$$\frac{t_{max} - a}{\log t_{max} - \log a}$$

with $t_{max} = e^{(1.0-\beta_0)/\beta_1}$.

3 Predicting Queue Wait Times

We use the run-time predictions described in the preceding section to predict queue wait times. Our technique is to perform a scheduling simulation using the predicted run times as the run times of the applications. This will then provide predictions of when applications will start to execute. Specifically, we simulate

the FCFS, LWF, and backfill scheduling algorithms and predict the wait time for each application as it is submitted to the scheduler. The accuracy of using various run-time predictors is shown in Table 4 through Table 9.

Table 4 shows the wait-time prediction performance when actual run times are used during prediction. No data is shown for the FCFS algorithm because there is no error when computing wait-time predictors in this case: later-arriving jobs do not affect the start times of the jobs that are currently in the queue. For the LWF and backfill scheduling algorithms, wait-time prediction error does occur because later arriving jobs can affect when the jobs currently in the queue can run. As one can see in the table, the wait-time prediction error for the LWF algorithm is between 34 and 43 percent. For the backfill scheduling algorithm, there is a smaller error of 3 to 4 percent. This error is higher for the LWF algorithm because later arriving jobs that wish to perform smaller amounts of work move to the head of the queue. Any error for the backfill algorithm seems unexpected at first, but errors in wait-time prediction can occur because scheduling is performed using maximum run times. For example, a job J_2 arriving in the queue can start ahead of an already queued job J_1 because the scheduler does not believe job J_1 can use those nodes; a running job finishes early, and job J_1 would start except that the job J_2 is using nodes that it needs. This example results in a wait-time prediction error for job J_1 before job J_2 arrives in the queue.

Table 4. Wait-time prediction performance using actual run times.

Workload	Scheduling Algorithm	Mean Error (minutes)	Percentage of Mean Wait Time
ANL	LWF	37.14	43
ANL	Backfill	5.84	3
CTC	LWF	4.05	39
CTC	Backfill	2.62	10
SDSC95	LWF	5.83	39
SDSC95	Backfill	1.12	4
SDSC96	LWF	3.32	42
SDSC96	Backfill	0.30	3

Table 5 shows the wait-time prediction errors while using maximum run times as run-time predictions. Maximum run times are used to predict run times in scheduling systems such as EASY [11]. These predictions are provided in the ANL and CTC workload and are implied in the SDSC workloads because each of the queues in the two SDSC workload has maximum limits on resource usage. To derive maximum run times for the SDSC workloads, we determine the longest running job in each queue and use that as the maximum run time for all jobs in that queue. The wait-time prediction error when using actual run times as run-time predictors is 59 to 99 percent better than the wait-time prediction error of

the LWF and backfill algorithms when using maximum run times as the run-time predictor.

Table 5. Wait-time prediction performance using maximum run times.

Workload	Scheduling Algorithm	Mean Error (minutes)	Percentage of Mean Wait Time
ANL	FCFS	996.67	186
ANL	LWF	97.12	112
ANL	Backfill	429.05	242
CTC	FCFS	125.36	128
CTC	LWF	9.86	94
CTC	Backfill	51.16	190
SDSC95	FCFS	162.72	295
SDSC95	LWF	28.56	191
SDSC95	Backfill	93.81	333
SDSC96	FCFS	47.83	288
SDSC96	LWF	14.19	180
SDSC96	Backfill	39.66	350

Table 6 shows that our run-time prediction technique results in wait-time prediction errors that are from 34 to 77 percent of mean wait times. Our data also shows that run-time prediction errors that are from 33 to 73 percent of mean application run times. The best wait-time prediction performance occurs for the ANL and CTC workloads and the worst for the SDSC96 workload. This is the opposite of what we expect from the run-time prediction errors. The most accurate run-time predictions are for the SDSC96 workload and the least accurate are for the CTC workload. The results imply that accurate run-time predictions are not the only factor that determines the accuracy of wait-time predictions.

The results when using our run-time predictor also show that the mean wait time prediction error is 20 percent better to 62 percent worse than when predicting wait times for the LWF algorithm using actual run times. Finally, using our run-time predictor results in 42 to 88 percent better wait time predictions than when using maximum run times as the run-time predictors.

Table 7 shows the wait-time prediction errors when using Gibbons's run time predictor. Our run-time prediction errors are between 39 and 68 percent better than Gibbons's and our wait-time prediction errors are between 13 and 83 percent better.

Tables 8 and 9 show the wait-time prediction error when using Downey's conditional average and conditional median predictors. The wait-time prediction errors we achieve when using our run-time predictor are between 19 and 87 percent better than these errors and our run-time prediction errors are between 42 and 92 percent better.

Table 6. Wait-time prediction performance using our run-time predictor.

Workload	Scheduling Algorithm	Mean Error (minutes)	Percentage of Mean Wait Time
ANL	FCFS	161.49	30
ANL	LWF	44.75	51
ANL	Backfill	75.55	43
CTC	FCFS	30.84	31
CTC	LWF	5.74	55
CTC	Backfill	11.37	42
SDSC95	FCFS	20.34	37
SDSC95	LWF	8.72	58
SDSC95	Backfill	12.49	44
SDSC96	FCFS	9.74	59
SDSC96	LWF	4.66	59
SDSC96	Backfill	5.03	44

In summary, the results show that our run-time predictor is more accurate than maximum run times, Gibbons’s predictor, or Downey’s predictors. The results also show that the wait-time prediction errors are smaller when our run-time predictor is used. Clearly, there is a correlation between wait-time prediction error and run-time prediction error.

Table 7. Wait-time prediction performance using Gibbons’s run-time predictor.

Workload	Scheduling Algorithm	Mean Error (minutes)	Percentage of Mean Wait Time
ANL	FCFS	350.86	66
ANL	LWF	76.23	91
ANL	Backfill	94.01	53
CTC	FCFS	81.45	83
CTC	LWF	32.34	309
CTC	Backfill	13.57	50
SDSC95	FCFS	54.37	99
SDSC95	LWF	11.60	78
SDSC95	Backfill	20.27	72
SDSC96	FCFS	22.36	135
SDSC96	LWF	6.88	87
SDSC96	Backfill	17.31	153

Table 8. Wait-time prediction performance using Downey’s conditional average run-time predictor.

Workload	Scheduling Algorithm	Mean Error (minutes)	Percentage of Mean Wait Time
ANL	FCFS	443.45	83
ANL	LWF	232.24	277
ANL	Backfill	339.10	191
CTC	FCFS	65.22	66
CTC	LWF	14.78	141
CTC	Backfill	17.22	64
SDSC95	FCFS	187.73	340
SDSC95	LWF	35.84	240
SDSC95	Backfill	62.96	223
SDSC96	FCFS	83.62	503
SDSC96	LWF	28.42	361
SDSC96	Backfill	47.11	415

Table 9. Wait-time prediction performance using Downey’s conditional median run-time predictor.

Workload	Scheduling Algorithm	Mean Error (minutes)	Percentage of Mean Wait Time
ANL	FCFS	534.71	100
ANL	LWF	254.91	304
ANL	Backfill	410.57	232
CTC	FCFS	83.33	85
CTC	LWF	15.47	148
CTC	Backfill	19.35	72
SDSC95	FCFS	62.67	114
SDSC95	LWF	18.28	122
SDSC95	Backfill	27.52	98
SDSC96	FCFS	34.23	206
SDSC96	LWF	12.65	161
SDSC96	Backfill	20.70	183

4 Improving Scheduler Performance

Our second application of run-time predictions is to improve the performance of the LWF and backfill scheduling algorithms. Table 10 shows the performance of the scheduling algorithms when the actual run times are used as run-time predictors. This is the best performance we can expect in each case and serves as an upper bound on scheduling performance.

Table 10. Scheduling performance using actual run times.

Workload	Scheduling Algorithm	Utilization (percent)	Mean Wait Time (minutes)
ANL	LWF	70.34	61.20
ANL	Backfill	71.04	142.45
CTC	LWF	51.28	11.15
CTC	Backfill	51.28	23.75
SDSC95	LWF	41.14	14.48
SDSC95	Backfill	41.14	21.98
SDSC96	LWF	46.79	6.80
SDSC96	Backfill	46.79	10.42

Table 11 shows the performance of using maximum run times as run time predictions in terms of average utilization and mean wait time. The scheduling performance when using the maximum run times can once again be considered an upper bound for comparison. When comparing this data to the data in Table 10, one can see that the maximum run times are an inaccurate predictor but this fact does not affect the utilization of the simulated parallel computers. Predicting run times with actual run-times when scheduling results in 3 to 27 percent lower mean wait times, except in one case where using maximum run times results in 6 percent lower mean wait times. The effect of accurate run-time predictions is highest for the ANL workload which has the largest offered load.

Table 12 shows the performance of using our run-time prediction technique when scheduling. The run-time prediction error in this case is between 23 and 93 percent of mean run times, slightly worse than the results when predicting run-times for wait-time prediction. This worse performance is due to more predictions being performed. First, more predictions are made of applications before they begin executing; and these predictions do not have information about how long an application has executed. Second, more predictions are made of long-running applications, the applications that contribute the largest errors to the mean errors.

Our run-time prediction technique results in mean wait times that are 5 percent better to 4 percent worse than when using actual run times as predictions for the least-work-first algorithm. For the backfill algorithm, mean wait times when using our run-time predictor are 11 to 22 percent worse. These results can

Table 11. Scheduling performance using maximum run times.

Workload	Scheduling Algorithm	Utilization (percent)	Mean Wait Time (minutes)
ANL	LWF	70.70	83.81
ANL	Backfill	71.04	177.14
CTC	LWF	51.28	10.48
CTC	Backfill	51.28	26.86
SDSC95	LWF	41.14	14.95
SDSC95	Backfill	41.14	28.20
SDSC96	LWF	46.79	7.88
SDSC96	Backfill	46.79	11.34

be understood by noticing that the backfill algorithm requires more accurate run-time predictions than LWF. LWF just needs to know if applications are “big” or “small,” and small errors do not greatly affect performance. The performance of the backfill algorithm depends on accurate run-time predictions because it tries to fit applications into time/space slots.

Table 12. Scheduling performance using our run-time prediction technique.

Workload	Scheduling Algorithm	Utilization (percent)	Mean Wait Time (minutes)
ANL	LWF	70.28	78.22
ANL	Backfill	71.04	148.77
CTC	LWF	51.28	13.40
CTC	Backfill	51.28	22.54
SDSC95	LWF	41.14	16.19
SDSC95	Backfill	41.14	22.17
SDSC96	LWF	46.79	7.79
SDSC96	Backfill	46.79	10.10

When comparing our run-time prediction technique to using maximum run times, our technique has a minimal effect on the utilization of the systems, but it does decrease the mean wait time in six of the eight experiments. Table 13 through Table 15 show the performance of the scheduling algorithms when using Gibbons’s and Downey’s run-time predictors. The results indicate that once again, using our run-time predictor does not produce greater utilization. The results also show that our run-time predictor results in 13 to 50 percent lower mean wait times for the ANL workload, but there is no clearly better run-time predictor for the other three workloads. The ANL workload has much larger mean wait times and higher utilization (greater offered load) than the other workloads (particularly the SDSC workloads). This may indicate that greater

prediction accuracy of our technique when scheduling becomes “hard”. To test this hypothesis, we compressed the interarrival time of applications by a factor of two for both SDSC workloads and then simulated these two new workloads. We found that our run-time predictor results in mean wait times that are 8 percent better on average, but are 43 percent lower to 31 percent higher than mean wait times than obtained when using Gibbons’s or Downey’s techniques.

The results also show that of Gibbons’s and Downey’s run-time predictors, Downey’s conditional average is the worst predictor and Gibbons’s predictor is the most accurate. The data shows that our run-time predictor is between 2 and 86 percent better than the other predictors, except for the CTC workload. For this workload, our predictor is the worst. This may be explained by the limited template searches we performed for that workload (because of time constraints). The accuracy of the run-time predictions for the CTC workload carries over to the mean wait times of the scheduling algorithms when using the various run-time predictors: our mean wait times are the worst.

Table 13. Scheduling performance using Gibbons’ run-time prediction technique.

Workload	Scheduling Algorithm	Utilization (percent)	Mean Wait Time (minutes)
ANL	LWF	70.72	90.36
ANL	Backfill	71.04	181.38
CTC	LWF	51.28	11.04
CTC	Backfill	51.28	27.31
SDSC95	LWF	41.14	15.99
SDSC95	Backfill	41.14	24.83
SDSC96	LWF	46.79	7.51
SDSC96	Backfill	46.79	10.82

5 Conclusions

In this work, we apply predictions of application run times to two separate scheduling problems. The problems are predicting how long applications will wait in queues before executing and improving the performance of scheduling algorithms. Our technique for predicting application run times is to derive a prediction for an application from the run times of previous applications judged similar by a template of key job characteristics. The novelty of our approach lies in the use of search techniques to find the best templates. For the workloads considered in this work, our searches found templates that result in run-time prediction errors that are significantly better than those of other researchers or using user-supplied maximum run times.

We predict queue wait times by using run-time predictions and the algorithms used by schedulers. These two factors are used to simulate scheduling

Table 14. Scheduling performance using Downey’s conditional average run-time predictor.

Workload	Scheduling Algorithm	Utilization (percent)	Mean Wait Time (minutes)
ANL	LWF	71.04	154.76
ANL	Backfill	70.88	246.40
CTC	LWF	51.28	9.87
CTC	Backfill	51.28	14.45
SDSC95	LWF	41.14	16.22
SDSC95	Backfill	41.14	20.37
SDSC96	LWF	46.79	7.88
SDSC96	Backfill	46.79	8.25

Table 15. Scheduling performance using Downey’s conditional median run-time predictor.

Workload	Scheduling Algorithm	Utilization (percent)	Mean Wait Time (minutes)
ANL	LWF	71.04	154.76
ANL	Backfill	71.04	207.17
CTC	LWF	51.28	11.54
CTC	Backfill	51.28	16.72
SDSC95	LWF	41.14	16.36
SDSC95	Backfill	41.14	19.56
SDSC96	LWF	46.79	7.80
SDSC96	Backfill	46.79	8.02

algorithms and decide when applications will execute. Estimates of queue wait times are useful to guide resource selection when several systems are available, to co-allocate resources from multiple systems, to schedule other activities, and so forth. This technique results in a wait-time prediction error of between 33 and 73 percent of mean wait times when using our run-time predictors. This error is significantly better than when using the run-time predictors of Gibbons, Downey, or user-supplied maximum run times. We also find that even if we predict application run times with no error, the wait-time prediction error for the least-work-first algorithm is significant (34 to 43 percent of mean wait times).

We improve the performance of the least-work-first and backfill scheduling algorithms by using our run-time predictions when scheduling. We find that the utilization of the parallel computers we simulate does not vary greatly when using different run-time predictors. We also find that using our run-time predictions does improve the mean wait times in general. In particular, our more accurate run-time predictors have the largest effect on mean wait time for the ANL workload, which has the highest utilization. In this workload, the mean wait times are between 7 and 67 percent lower when using our run-time predictions than when using other run-time predictions. We also find that on average, the mean wait time when using our predictor is within 8 percent of the mean wait time that would occur if the scheduler knows the exact run times of the applications. The mean wait time when using our technique ranges from 5 percent better to 22 percent worse than when scheduling with the actual run times.

In future work, we will investigate an alternative method for predicting queue wait times. This method will use the current state of the scheduling system (number of applications in each queue, time of day, etc.) and historical information on queue wait times during similar past states to predict queue wait times. We hope this technique will improve wait-time prediction error, particularly for the LWF algorithm, which has a large built-in error using the technique presented here. Further, we will expand our work in using run-time prediction techniques for scheduling to the problem of combining queue-based scheduling and reservations. Reservations are one way to co-allocate resources in metacomputing systems [1,6,2,7]. Support for resource co-allocation is crucial to large-scale applications that require resources from more than one parallel computer.

Acknowledgments

We thank the Mathematics and Computer Science Division of Argonne National Laboratory, the Cornell Theory Center, and the San Diego Supercomputer Center for providing us with the trace data used in this work. We also thank Richard Gibbons for providing us with his run-time prediction code and Allen Downey for assisting in the comparison with his work.

This work was supported by the Mathematical, Information, and Computational Sciences Division subprogram of the Office of Computational and Technology Research, U.S. Department of Energy, under Contract W-31-109-Eng-38, and by a NSF Young Investigator award under Grant CCR-9215482.

References

1. C. Catlett and L. Smarr. Metacomputing. *Communications of the ACM*, 35(6):44–52, 1992.
2. K. Czajkowski, I. Foster, N. Karonis, C. Kesselman, S. Martin, W. Smith, and S. Tuecke. A Resource Management Architecture for Metasystems. *Lecture Notes on Computer Science*, 1998.
3. Allen Downey. Predicting Queue Times on Space-Sharing Parallel Computers. In *International Parallel Processing Symposium*, 1997.
4. N. R. Draper and H. Smith. *Applied Regression Analysis, 2nd Edition*. John Wiley and Sons, 1981.
5. Dror Feitelson and Bill Nitzberg. Job Characteristics of a Production Parallel Scientific Workload on the NASA Ames iPSC/860. *Lecture Notes on Computer Science*, 949:337–360, 1995.
6. Ian Foster and Carl Kesselman. Globus: A Metacomputing Infrastructure Toolkit. *International Journal of Supercomputing Applications*, 11(2):115–128, 1997.
7. Ian Foster and Carl Kesselman, editors. *The Grid: Blueprint for a New Computing Infrastructure*. Morgan Kaufmann, 1999.
8. Richard Gibbons. A Historical Application Profiler for Use by Parallel Schedulers. *Lecture Notes on Computer Science*, 1297:58–75, 1997.
9. Richard Gibbons. A Historical Profiler for Use by Parallel Schedulers. Master's thesis, University of Toronto, 1997.
10. David E. Goldberg. *Genetic Algorithms in Search, Optimization, and Machine Learning*. Addison-Wesley, 1989.
11. David A. Lifka. The ANL/IBM SP Scheduling System. *Lecture Notes on Computer Science*, 949:295–303, 1995.
12. Warren Smith, Ian Foster, and Valerie Taylor. Predicting Application Run Times Using Historical Information. *Lecture Notes on Computer Science*, 1459:122–142, 1998.
13. Neil Weiss and Matthew Hassett. *Introductory Statistics*. Addison-Wesley, 1982.