

Support for Implementing Scheduling Algorithms Using MESSIAHS*

Steve J. Chapin
Dept. of Math. & Computer Science
Kent State University
Kent, OH 44242-0001
sjc@cs.kent.edu

Eugene H. Spafford
Dept. of Computer Sciences
Purdue University
West Lafayette, IN 47907-1398
spaf@cs.purdue.edu

Abstract

The MESSIAHS project is investigating mechanisms that support task placement in heterogeneous, distributed, autonomous systems. MESSIAHS provides a substrate on which scheduling algorithms can be implemented. These mechanisms were designed to support diverse task placement and load balancing algorithms.

As part of this work, we have constructed an interface layer to the underlying mechanisms. This includes the MESSIAHS Interface Language (MIL) and a library of function calls for constructing distributed schedulers. This paper gives an overview of MESSIAHS, describes two sample interface layers in detail, and gives example implementations of well-known algorithms from the literature built using these layers.

1 Introduction

Recent initiatives in high-speed, heterogeneous computing have spurred renewed interest in large-scale distributed systems, and the desire for better utilization of existing resources has contributed to this movement. A typical departmental computing environment already has a substantial investment in computing equipment, including dozens or hundreds of workstations. Studies have shown that the utilization of this equipment can be as low as 30% of capacity [1, 2].

A solution to this problem is to conglomerate the separate processors into a distributed system, and to recursively join the distributed systems into larger systems to further expand the computational power of the whole. Large-scale distributed systems can have a combined computing power outperforming that of supercomputers [3].

A central element of effective utilization of such systems is task scheduling. Task scheduling has two components: macro-scheduling (also defined as global scheduling [4]) and task

*This work was sponsored in part by NASA GSRP grant number NGT 50919.

assignment [5]) and micro-scheduling (or local scheduling [4]). Macro-scheduling chooses where to run a process, while micro-scheduling selects which eligible process to execute next on a particular processor. All further uses of the term *scheduling* in this paper refer to macro-scheduling.

The MESSIAHS¹ system [6, 7, 8, 9] provides a set of mechanisms that facilitate scheduling in distributed, heterogeneous, autonomous systems. For our purposes, distributed, or loosely-coupled, systems communicate via message passing rather than a shared memory bus. Heterogeneous systems may have different instruction set architectures, data formats, and attached devices. All policy decisions in autonomous systems are made locally. Our vision of distributed systems includes all three attributes, connecting machines of different architectures with individual administrative authorities via a communications network. Section 2 gives precise definitions of autonomy and heterogeneity.

It is vital that a system for distributed computation support autonomy because of the prevailing decentralization of computing resources. There is usually no longer a single, authoritative controlling entity for the computers in a large organization. A scientist may control a few of his own machines, and his department may have administrative control over several such sets of machines. That department may be part of a regional site, which is, in turn, part of a national organization. No single entity, from the scientist up to the national organization, has complete control over all the computers it may wish to use. An example of such usage is when two research organizations pool their resources to solve a common problem.

Heterogeneity is important because it yields the most cost-effective and efficient method for performing some computations. A large computation might have certain pieces best suited for execution on a supercomputer, while other parts might run best on a hypercube or a graphics workstation. If the distributed system is restricted to only using one architecture, the computation will suffer needless delay. In other cases, tasks such as text processing or high-level language interpretation may be independent of any single architecture.

Our subsequent uses of the terms *distributed system* or *system* refer to a distributed, autonomous, heterogeneous system, and *node* refers to an individual machine within an autonomous system. Our definition of system includes a single machine, as well as two homogeneous workstations communicating via a local-area network. This definition also encompasses systems as complex as thousands of machines, including personal computers, workstations, file and computation servers, and supercomputers, spread among several remote sites and connected by a wide-area network.

Within this distributed system, each individual system has its own policy for deciding when to accept or remove tasks. The local administrator defines this policy, which is implemented over the MESSIAHS mechanisms via an interface layer. The interface layer provides a virtual machine interface; the underlying mechanism can be presented to algorithm writers in various ways. The language described here provides an interface that is easy to use, yet powerful enough to implement a wide variety of scheduling algorithms. Primitive operations are supplied to access system and task state information, manipulate tasks, and control the behavior of the local system.

This approach is distinct from that taken in distributed programming systems such as

¹Mechanisms Effecting Scheduling Support In Autonomous, Heterogeneous Systems

PVM [10] in which the program distribution is visible, and even forced upon, the programmer. The MESSIAHS approach more closely reflects that taken in Condor [11], which schedules processes invisibly for the programmer. Program distribution is under the control of the autonomous system, and therefore the administrators, rather than the programmer.

Portions of this paper discuss implementation issues of the MESSIAHS prototype. The prototype was written in C for SunOS 4.1, and runs on both Sun 3 and Sun SPARC architectures.

2 The MESSIAHS Architecture

MESSIAHS supports task placement in distributed systems with hierarchical structure based on administrative domains, modeled by directed acyclic graphs. Multiple subordinate systems can be combined into an encapsulating system, yielding the hierarchical structure. The nodes of the graph represent the autonomous systems, and edges indicate encapsulation. The graph is directed downward; the edges are directed from encapsulating nodes, or parents, to subordinate nodes, or children. Children of the same parent are siblings². The *neighbors* of a system are its children, parents, and siblings.

Figure 1 shows an example distributed system based on the Purdue University Computer Sciences department. In the example, the Computer Sciences department contains two administrative domains, Cypress and General. Cypress in turn encapsulates the research machines belonging to the Cypress project, and General contains the general purpose servers for the department. Bredbeddle and Percival are children of Cypress, and therefore are siblings.

Each component system runs a *scheduling module* that implements the local scheduling policy and manages administrative aspects of the system. These modules exchange data sets describing the status of the systems. On demand, the modules also process scheduling requests, which contain a description of the task for which the sender requests service.

Each module only exchanges status information with modules running on its neighbors. Because of the hierarchical structure of the system, some nodes might be invisible to other nodes. In the example system from figure 1, Arthur receives information updates only from Nyneve and General, and sees no information that can be directly related to Percival or Bredbeddle. The capabilities of Bredbeddle and Percival are subsumed and combined within General's state advertisement.

Individual systems enjoy *execution autonomy*, *communication autonomy*, *design autonomy*, and *administrative autonomy* as defined in [9, 13, 14]. Execution autonomy means that each system decides whether it will honor a request to execute a task; each system also has the right to revoke a task that it had previously accepted. Communication autonomy means that each system decides the content and frequency of state advertisements, and what other messages it sends. A system is not required to advertise all its capabilities, nor is it required to respond to messages from other systems. Design autonomy gives the architects of a system freedom to design and construct it without regard to existing systems, yielding heterogeneous systems.

²These correspond to the formal definitions of father, son, and brother found in [12].

Administrative autonomy means that each system can have its own usage policies and behavioral characteristics, independent of any others. In particular, a local system can run in a manner counterproductive to a global optimum. In the usual case, scheduling modules will cooperate, but administrators must be free to set their local policies or they will not participate in the distributed system. Both [11] and [1] note that users are willing to execute remote jobs on their workstations if the scheduling policy places higher priority on local jobs.

Figure 2 displays the structure of a MESSIAHS scheduling module. The machine-dependent layer handles raw data exchange between scheduling modules, collects the local state information, and interacts with the task manipulation mechanisms specific to the local operating system. The abstract data and task management layer provides an abstract interface for the machine-dependent operations to the data reporting layer. The shaded layer, data reporting and task manipulation, is the focus of this paper. This layer presents the user with the interface to the MESSIAHS mechanisms. The administrator supplies the topmost layer, which embodies the scheduling policy for the system.

MESSIAHS does not determine policy; the three layers provide mechanisms to implement scheduling policies. The interface layer is the administrator's vehicle for expressing and implementing the local policy through the MESSIAHS mechanisms. Sections 4 and 5 describe two interface layers, but we next examine the two lower levels upon which the interface layer is built. This will provide a frame of reference for discussion of the interface layer.

2.1 The Machine-Dependent Layer

The machine-dependent layer provides the interface in table 1 to the management layer of the module. The prototype does not implement those functions marked with a †. As noted earlier, discussion of implementation details pertains to the prototype running on SunOS 4.1.

The functions divide into three main groups: data collection, message passing, and task management. The data collection routines gather information that forms the system description for the local host. The message-passing routines implement abstract message exchange between modules. The task management routines provide access to the underlying operating system process manipulation primitives.

The data collection operations are implemented using the `kvm_open()`, `kvm_read()`, `kvm_nlist()`, and `kvm_close()` routines that access kernel state in SunOS 4.1. The `collect_process_data()` function collects information on the number of processes in the ready queue, and the percentage of processor utilization. `collect_memory_data()` determines how much of the physical memory is in use. `collect_disk_data()` finds the amount of public free space on a system, typically in the `/tmp` directory on SunOS. `collect_network_data()` determines the average round-trip time between a host and its neighboring systems within the graph.

An alternative data collection implementation could use the `rstat()` call, which uses the Remote Procedure Call (RPC) mechanisms of SunOS to query a daemon that monitors the kernel state. However, the `rstatd` daemon does not provide information on physical memory statistics or communication time estimates, which are required to implement the mechanisms. Use of `rstat()` and `rstatd` also involves communication and context-switching overhead.

The message passing routines use the SunOS socket abstraction for communication and

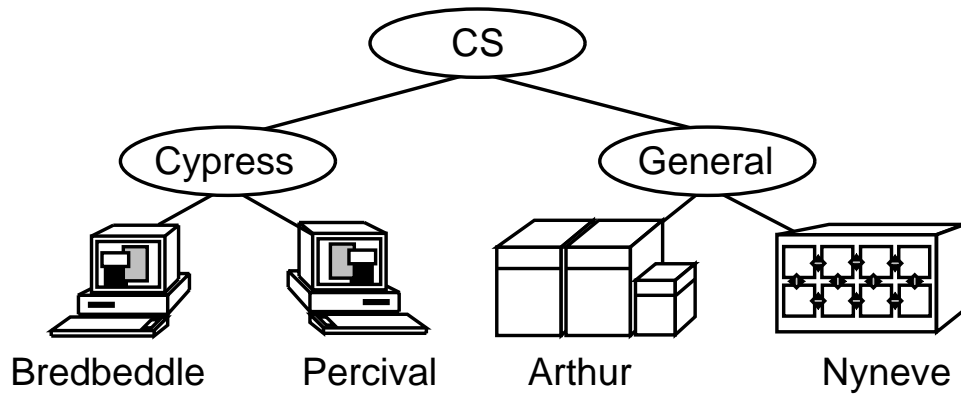


Figure 1: a sample distributed system

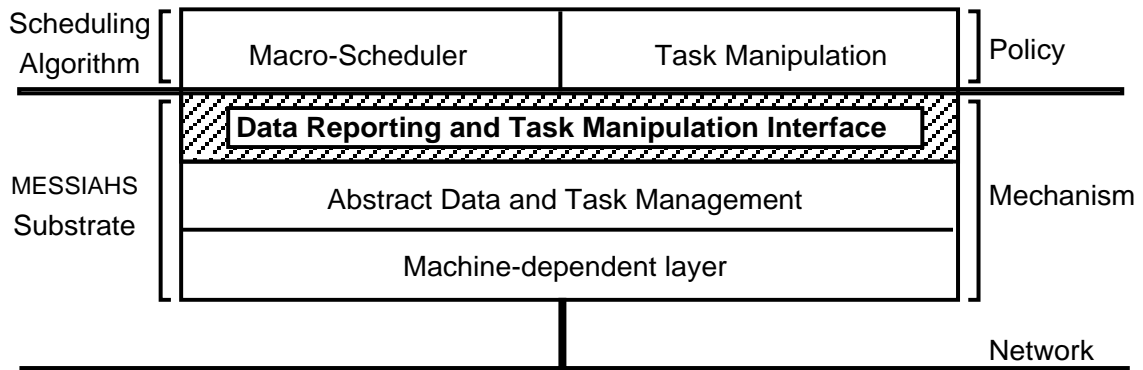


Figure 2: structure of a MESSIAHS scheduling module

Table 1: Functions in the machine-dependent layer

Purpose	Function Name	Description
data collection	collect_process_data	collect data regarding the number of processes and load statistics
	collect_memory_data	collect data on available and total memory
	collect_disk_data	collect data on available temporary disk space
	collect_network_data	collect data on inter-module communication time
message passing	get_message	receive a message from the network
	send_message	send a message over the network
task management	suspend_task	pause a running task
	resume_task	continue executing a suspend task
	kill_task	halt execution of a task and remove it from the system
	checkpoint_task [†]	save the state of a task
	migrate_task [†]	checkpoints a task and moves it to a target host
	revert_task [†]	returns a task to its originating system

the User Datagram Protocol (UDP) to exchange information between modules. UDP was chosen because it provides an unreliable datagram protocol, which is the minimum level of service required for the update and control channels. The message passing routines encode the data using the XDR standard for external data representation.

The task manipulation primitives use the SunOS `kill()` system call, which sends a software interrupt, called a *signal*, to a process. The signals used are `SIGSTOP`, which pauses a process, `SIGCONT`, which resumes a paused process, and `SIGKILL`, which terminates a process. The task migration primitive is not implemented in the prototype, but is a stub procedure for later completion.

2.2 Abstract Data and Communication Management

The middle layer in figure 2 comprises the abstract data and task manipulation functions. These functions use the basic mechanism provided by the machine-dependent layer to construct higher-level semantic operations. For example, the `send_sr()` routine, which sends a schedule request to a neighbor, is implemented using the `send_message()` function. Table 2 lists the abstract data and task management functions.

The message-passing functions construct a message from the pertinent data and use the `send_message()` function to communicate with a neighboring module. There is one `send` routine for each message type.

MESSIAHS maintains two hash tables containing description vectors: one table containing description vectors of foreign tasks executing on the local host and another table for description vectors of neighboring systems. The hash tables use double hashing as described in Knuth [15, pp. 521–526] for efficiency. The `sys_lookup()` and `task_lookup()` routines search the tables for a particular task or system. The `sys_first()`, `sys_next()`, `task_first()`, and `task_next()` routines iterate over the tables, returning successive description vectors with each call.

The event manipulation routines provide access to the internal event queues used by the module. The `register_event()` function inserts a timed event into the timeout queue, and the `enqueue()` and `dequeue()` routines allow direct manipulation of the queues. The `set timeout` routines enqueue timeout events of particular types, and the `set period` functions set the timeout periods for the various timers in MESSIAHS. If a timeout period is set to 0, the associated timer is disabled. Input timeouts occur when a neighbor has not sent a status message to the local host within the timeout period. Output timeouts indicate that the local host should advertise its state to its neighbors. Recalculation timeouts cause the local host to recompute its update vectors. When a revocation timeout occurs, the host checks its state to see if tasks should be revoked.

3 Support for Scheduling Policies

As seen in figure 2, the scheduling policy is implemented over the interface layer. Through the interface layer, MESSIAHS either directly provides or supports five mechanisms that can be used to construct scheduling policies. These five mechanisms are system description, decision filters, task revocation, data combination and condensation, and node configuration and behavior customization.

Table 2: Functions in the abstract data and communication layer

Purpose	Function Name	Description
data exchange	send_sr	send a schedule request message
	send_sa	send a schedule accept message
	send_sd	send a schedule deny message
	send_trq	send a task request message
	send_ta	send a task accept message
	send_td	send a task deny message
	send_trv	send a task revoke message
	send_ssq	send a system status query
	send_ssv	send a system status vector
	send_tsq	send a task status query
	send_tsv	send a task status vector
	send_jr	send a join request
	send_jd	send a join deny
	description	sys_lookup
vector access	sys_first	return the first neighbor from the system hash table
	sys_next	return the next neighbor from the system hash table
	task_lookup	find the TDV for a task in the task hash table
	task_first	return the first task from the task hash table
events	task_next	return the next task from the task hash table
	register_event	insert an event into the timeout event queue
	enqueue_event	enqueue an event
	dequeue_event	dequeue an event
	new_queue	allocate an event queue
	qempty	check if a queue is empty
	set_input_timeout	enqueue an input timeout
	set_output_timeout	enqueue an output timeout
	set_recalc_timeout	enqueue a recalculation timeout
	set_revoke_timeout	enqueue a revocation timeout
	set_oto_period	set the output timeout period
set_ito_period	set the input timeout period	
set_rcto_period	set the recalculation timeout period	
set_rvto_period	set the revocation timeout period	

3.1 Intrinsic Mechanisms

MESSIAHS uses a mechanism called *description vectors* to characterize available resources and requests for resources. A system description vector, or SDV, lists the capabilities of an autonomous system and comprises the state advertised between systems. A task description vector, or TDV, describes the resources required by a computational job. Description vectors contain a fixed portion that is optimized for task placement support, and an extensible portion that administrators can use to implement new scheduling policies or to extend the basic descriptions of requirements or abilities.

To determine the basis for the fixed portion of the description vector, we reviewed 18 algorithms from the existing scheduling literature [2–5, 7, 13, 17–20, 22–24, 26–28, 30, 31]. Table 3 depicts the resulting data set. We found that only two characteristics—processor speed and inter-processor communication time estimates—were used by more than four algorithms. Therefore, we included processor speed estimates in the description vector and provide a mechanism to determine inter-system communication time. We also augmented SDVs with other data items that we expect will be useful to writers of future scheduling algorithms. This data supports the common case, as represented by the surveyed algorithms, while the extension mechanism allows the inclusion of special-purpose data.

The `address` and `module` fields uniquely identify a scheduling module: the address specifies a machine, and the module indicates which module on that machine. MESSIAHS allows multiple modules to run on a single machine (see [7]). The `nsys` field indicates how many systems the vector represents; just as a distributed system encapsulates multiple subordinate systems, the description vector for a system contains information describing its component systems. The `n tasks`, `n active tasks`, and `n suspended tasks` list the number of total tasks, running tasks, and suspended tasks for the system. The `willingness` gives the rough probability that the system will accept a new task, and `loadave` estimates the computational load on the entire system. The `Procclass` field is an array of records describing statistical measures of the processor utilization, processor speed, free memory, and disk space.

Execution autonomy mandates the ability to remove a task from execution on the local system. Aborting a running task fulfills the autonomy requirements, but does not support load-balancing algorithms based on process migration. Therefore, MESSIAHS includes mechanisms to kill, checkpoint, suspend, resume, and migrate jobs.

In support of administrative and communication autonomy, tunable parameters affect the general behavior of the node. These parameters are independent of any single scheduling policy, and effect all polices running on the node. These four parameters are listed in table 4.

The `recalc_timeout` field and `revocation_timeout` fields determine how often prescribed events occur. The `SPECint92` and `SPECfp92` are measures of processor speed using the SPEC benchmark suite [30]. The SPEC benchmark suite consists of applications-oriented programs, specifically selected to represent real-world workloads.

The machine architecture type (e.g. SPARC or VAX) does not appear as a universal parameter because many jobs are architecture independent. For example, text formatting requests require the presence of a particular text processing package, but do not depend on the underlying architecture.

Table 3: fixed portion of a system description vector

field name	purpose
address	address of the system
module	id of module on this system
nsys	number of systems described by the vector
ntasks	total number of tasks currently accepted by the system
nactivetasks	number of active tasks running on the system
nsuspendedtasks	number of inactive tasks waiting on the system
willingness	desire of the system to take on new tasks
loadave	an estimate of the load average for the entire system
Procclass	information on the different classes of processors in the system

Table 4: general state parameters

parameter	purpose
recalc_timeout	suggested period, in seconds, between recalculations of the local system description
revocation_timeout	period, in seconds, between checks for possible task revocation
SPECint92	the integer performance rating of the node, per the SPEC integer benchmark.
SPECfp92	the floating point performance rating of the node, per the SPEC floating point benchmark.

3.2 Supported Mechanisms

MESSIAHS supports the use of *filters* to implement scheduling policies. Decision filters take two description vectors as input, and return an integer value denoting how well the two vectors match according to the local policy. Larger values indicate closer matches. Scheduling modules employ filters to determine where to attempt scheduling a task (including on the local node), and what tasks are eligible for migration or revocation.

MESSIAHS allows multiple scheduling policies to operate within the system simultaneously, and a single node can support two or more scheduling policies. For example, batch queues for text processing, remote compilation, and remote program execution could all coexist within the same distributed system, each with its own individual scheduling policy. The administrator for each node could determine whether that node would participate as a server for any or all of the services.

Communication autonomy requires that the local policy control the flow of information out of a system. This mandates a mechanism to combine and compact the data set, and to allow the advertisement of restricted sets of information. In addition, data condensation is essential to avoid arbitrary limits on scaling the mechanisms. If systems concatenated all the data describing subordinate systems, the resources required to transmit and process a description vector would soon outstrip the capabilities of many networks and processors.

Unfortunately, some information loss is unavoidable if data compression takes place. Recall that in our example system, Arthur has no first-hand information about Bredbeddle or Percival. Therefore, Arthur might misdirect scheduling requests to General, based on the union of Percival's and Bredbeddle's abilities. For example, if Percival had 100 megabytes of free disk space and 4 megabytes of memory, while Bredbeddle had 10 megabytes of disk space and 32 megabytes of memory, the scheduling module on Arthur might mistakenly think that resources were available to execute a task requiring 16 megabytes of memory and 50 megabytes of disk space. These misdirected requests cause a small efficiency loss, but no tasks will be misscheduled as a result.

4 The Language

The shaded interface layer shown in figure 2 provides scheduling algorithms with access to lower-level mechanisms. We have chosen to provide two interface layers: a simple programming language, similar to that used in Univers [31], and a library of high-level language functions. This section describes the MESSIAHS Interface Language (MIL), and the next section describes the library of function calls.

MIL contains direct support for dynamic scheduling algorithms, without precluding support for static algorithms. Static algorithms consider only the system topography, not the state, when calculating the mapping. Dynamic algorithms take the current system state as input, therefore the resultant mapping depends on the state (see [4]). Figure 3 depicts the structure of an MIL program. The grammars for deriving the various rules, along with explanations of their semantics, appear in the rest of this section.

```
begin state
  <node state rules>
end
begin combining
  <data combination rules>
end
begin schedfilter
  <sched request filter rules>
end
begin taskfilter
  <task request filter rules>
end
begin revokefilter
  <revocation filter rules>
end
begin revokerules
  <revocation rules>
end
```

Figure 3: MIL specification template

4.1 Expressions and Types

MIL defines four basic types for data values: integers (INT), booleans (BOOL), floats (FLOAT), and strings (STRING). Integers can be written in decimal or in hexadecimal. Booleans have either the value `true` or `false`. Floats are two decimal digit sequences separated by a decimal point, e.g. 123.45. Strings are a sequence of characters delimited by quotation marks (").

Identifiers are a dollar sign followed by either a single word, or two words separated by a period. The latter case specifies fields within description vectors. The legal vectors are the received task description (`task`), the description of a task already executing on the system (`loctask`), the system description of a neighboring system (`sys`), the description of the local node (`me`), and the description being constructed by data combination (`out`). `loctask` is used for the task request filter and the revocation filter. `sys` is used for the data combination rules and the schedule request filter. `out` is used only for the data combination rules, and `me` can appear in any of the combination rules, filtering, or task revocation sections.

The following grammar defines the expression types used by the language. This grammar only derives expressions of the base types; in particular, there is no access to the `Procclass` field of the SDV with MIL.

$$\begin{aligned}
 \textit{int-binop} &\rightarrow + \mid - \mid / \mid * \mid \textit{mod} \mid \& \mid \mathbf{!} \mid \\
 &\quad \textit{max} \mid \textit{min} \\
 \textit{int-expr} &\rightarrow \textit{int-expr} \textit{int-binop} \textit{int-expr} \mid \\
 &\quad (\textit{int-expr}) \mid \textit{integer} \mid \\
 &\quad \textit{int}(\textit{float-expr}) \mid \textit{id} \\
 \\
 \textit{string-expr} &\rightarrow \textit{string-expr} + \textit{string-expr} \mid \\
 &\quad (\textit{string-expr}) \mid \textit{string} \mid \textit{id} \\
 \\
 \textit{float-binop} &\rightarrow + \mid - \mid / \mid * \mid \textit{max} \mid \textit{min} \\
 \textit{float-expr} &\rightarrow \textit{float-expr} \textit{float-binop} \textit{float-expr} \mid \\
 &\quad (\textit{float-expr}) \mid \textit{float} \mid \\
 &\quad \textit{float}(\textit{int-expr}) \mid \textit{id} \\
 \\
 \textit{comp} &\rightarrow < \mid > \mid = \mid >= \mid <= \mid <> \\
 \textit{bool-binop} &\rightarrow \textit{and} \mid \textit{or} \mid \textit{xor} \\
 \textit{bool-expr} &\rightarrow \textit{bool-expr} \textit{bool-binop} \textit{bool-expr} \mid \\
 &\quad \textit{not} \textit{bool-expr} \mid \\
 &\quad \textit{int-expr} \textit{comp} \textit{int-expr} \mid \\
 &\quad \textit{float-expr} \textit{comp} \textit{float-expr} \mid \\
 &\quad \textit{string-expr} \textit{comp} \textit{string-expr} \mid \\
 &\quad \textit{match}(\textit{string-expr}, \textit{string-expr}) \mid \\
 &\quad (\textit{bool-expr}) \mid \textit{true} \mid \textit{false} \mid \textit{id}
 \end{aligned}$$

4.2 Access to Intrinsic Mechanisms

MIL includes five task manipulation primitives: `kill`, `suspend`, `wake`, `migrate`, and `revert`. Other operations, such as process checkpointing, are available in the lower-level mechanisms, but are not explicitly included in the language. `kill` aborts a task, discards any interim results, and frees system resources used by the task. `suspend` temporarily blocks a running task. `wake` resumes a suspended task. `migrate` checkpoints a task and attempts to schedule it on neighboring systems. `revert` checkpoints the task and returns it to the originating system for rescheduling. Task revocation rules take the following form, using a boolean guard to determine when to take an action.

$$\begin{array}{ll}
 \textit{task-action} & \rightarrow \textit{kill} \mid \\
 & \textit{suspend} \mid \\
 & \textit{wake} \mid \\
 & \textit{migrate} \mid \\
 & \textit{revert} \\
 \textit{revocation-rule} & \rightarrow \textit{bool-expr} : \textit{task-action} ;
 \end{array}$$

The node state section is a list of types, identifiers, and constant values. Node state declarations are parameters that affect system state. Unlike the extension variables, they do not directly appear in the system description vector. The four node state parameters are `specint92`, `specfp92`, `recalc_timeout`, and `revocation_timeout`. The `specint92` and `specfp92` parameters list the speed of the host in terms of the SPEC benchmarks [30]. The `recalc_timeout` and `revocation_timeout` parameters determine the timeout periods for the associated events.

4.3 Filters and Data Combination

In MIL, a filter is a series of guarded statements, similar to combining rules. In place of an *action*, filters define integer expressions,

$$\textit{filter-stmt} \rightarrow \textit{bool-expr} : \textit{int-expr} ;$$

A return value of 0 indicates that there is no match. A negative value indicates an error, and a positive value measures the affinity of the two vectors. As noted earlier, higher values indicate a better match. If the guard expression uses an undefined variable, the guard evaluates to `false`. If the integer expression references an undefined variable, the filter returns `-1`, indicating an error. With appropriate extension variables and guards, a single scheduling module can serve multiple scheduling policies as stated in section 3.2.

MIL provides a mechanism to combine description vectors. To support communication autonomy, this mechanism allows the administrator to write rules specifying operations to coalesce the data.

```

int-action      →  discard | set int-expr
float-action   →  discard | set float-expr
bool-action    →  discard | set bool-expr
string-action →  discard | set string-expr

combining-rule →  int id bool-expr:
                    int-action ; |
                    float id bool-expr:
                    float-action ; |
                    string id bool-expr:
                    string-action ; |
                    bool id bool-expr:
                    bool-action ;

```

The boolean expression acts as a guard, and the action is performed for a particular (*type*, *identifier*) pair if the value of the guard is **true**. Administrators may supply multiple rules for the same pair. If multiple rules exist, the module evaluates them in the order written, performing the action corresponding to the first guard that evaluates to **true**.

If no matching rule is found for a pair, the identifier is discarded. Explicit discarding of data items, via the **discard** action, fulfills the constraint of communication autonomy. The **set value** action assigns *value* to the current pair in the outgoing description vector. An error in evaluating a guard automatically evaluates to **false**. If the evaluation of an action expression causes a run-time error, e.g. a division by 0, the action converts to **discard**.

4.4 Specification Evaluation

The extension and node state rules are interpreted when the specification is first loaded. The data combination rules are applied when a recalculation timeout occurs. When a revocation timeout occurs, the module passes once through the list of revocation rules, repeatedly evaluating each one until its guards return **false**. If the guard evaluates to **true**, the revocation filter is applied to the appropriate list of tasks to provide a target for the revocation action. If no task matches, the module moves on to the next rule in the list.

When a scheduling request arrives, the module iterates over the list of available systems, evaluating the request filter rules in-order until a guard that evaluates to **true** is found, or the rules are exhausted. If no matching rule is found, 0 is returned. If a rule is found, its value is returned as the suitability ranking for that system. The module follows a similar procedure for task requests, iterating over the set of available tasks.

4.5 A Small Example

Figure 4 shows a simple MIL specification for a SPARC IPC participating in a distributed L^AT_EX text-processing system. Line 1 in the node state section sets the period for SDV recalculation at 60 seconds. Every minute, each participating system will compute its SDV and forward updates to its neighbors.

The SDV extension variable `hasLaTeX` is true if the system has \LaTeX available and wishes to act as a formatting server. Clients requesting \LaTeX processing set the `needsLaTeX` variable to `true` in their task description vector. The combining rule in line 2 sets the outgoing `hasLaTeX` variable if any of the incoming description vectors have it set, and the rule on line 3 sets the `hasLaTeX` variable for the local hosts. Hosts providing the \LaTeX service would use line 3; hosts not providing the service would use line 2 to propagate advertisements by other hosts.

The scheduling filter rule in line 4 compares the available system vectors to the incoming task vector, accepts servers with load averages of less than five, and ranks the systems based on their load average. The guard would fail for a neighbor that had not set the `hasLaTeX` variable, and return `false`.

5 A Library of Function Calls

This section describes a library of function calls, called a *scheduling toolkit* that provides access to the underlying mechanism. The toolkit provides access to the functions in the low and middle layers as well as the functions listed in table 5.

The `send_Uvec()`, `send_Dvec()`, and `send_Svec()` functions send update vectors to a system's parents, children, and siblings, respectively.

As shown in figure 5, statistics vectors (`statvec`) are components of the `procclass` structure, which are used to condense the advertised state information for a virtual system. Processors are grouped into *process classes* on a logarithmic scale, based on their computation speed. The `statvec` fields represent multiple processors using statistical descriptions of their capabilities. Processor speed was chosen as the grouping factor because research of the existing scheduling algorithms indicates that processor speed is the primary consideration for task placement (see chapter 2 of [7]). The SPEC ratings were chosen as the default speed rating because they are the most widely available benchmark for both integer and floating point performance. Other measures of speed can be included through the extension mechanism.

The `merge_statvec()` function merges two statistics vectors, and `merge_procclass()` merges two processor classes into one. The `merge_SDV()` function provides a default mechanism for merging two system description vectors into one. The functions in figure 5 are used to implement MIL, described in the previous section.

The programmer uses the toolkit to write a set of event handlers. These handlers comprise the scheduling policy. MESSIAHS predefines the set of handlers listed in table 6, which may be overloaded by the administrator to create a new policy.

6 Example Algorithms Using MIL

In addition to the simple \LaTeX batch processing system described earlier, we present two applications built using MIL. The first demonstrates the task revocation facility as used by a general-purpose distributed batch system. The second implements a load-balancing algorithm.

```

begin state
1.   int    $recalc_period    60;
end
begin combining
2.   bool   $out.hasLaTeX
      $sys.hasLaTeX:  set true;
3.   bool   $out.hasLaTeX
      $sys.address == $me.address:
          set true;
end
begin schedfilter
4.   $task.needsLaTeX and $sys.hasLaTeX
      and int($sys.loadave) < 5 :
          6 - int($sys.loadave);
end

```

Figure 4: a simple MIL specification

Table 5: Functions in the MESSIAHS toolkit

Purpose	Function Name	Description
data exchange	send_Uvec	send the U update vector to a parent
	send_Dvec	send the D update vector to a child
	send_Svec	send the U update vector to a sibling
description vector manipulation	merge_SDV	merges two SDVs into one
	merge_statvec	merge two statistics vectors into one
miscellaneous	merge_procclass	merge two procclass sets into one
	mk_sid_sb	return a printable form of the system identification number
	Log	produce output in the error log
	pLog	produce output in the error log, including operating-system specific error messages

```

struct statvec {
    float min, max, mean, stddev, total;
};

typedef struct statvec Statvec;

struct procclass {
    bit32    nsys;          /* # of machines in this class */
    Statvec  qlen;         /* run queue statistics */
    Statvec  busy;         /* load on cpu (percentage) */
    Statvec  physmem;      /* total physical memory */
    Statvec  freemem;      /* available memory */
    Statvec  specint92;    /* ratings for SPECint 92 */
    Statvec  specfp92;     /* ratings for SPECfp 92 */
    Statvec  freedisk;     /* public disk space (/tmp) stats */
};

typedef struct procclass Procclass;

#define SDV_NPROCCLASS 7
#define SDV_MAXUSERDEF 2048 /* multiple of 2 for cksum */

struct SDV {
    SysId    sid;          /* Autonomous System ID */
    bit32    nsys;         /* number of total systems */
    bit32    ntasks;       /* number of total tasks */
    bit32    nactivetasks; /* number of active tasks */
    bit32    nsuspendedtasks; /* number of suspended tasks */
    float    willingness;  /* probability of taking on
                           /* a new task
    float    global_load;  /* global load average
    Procclass procs[SDV_NPROCCLASS]; /* information on the
                                   /* different classes of procs
                                   /* in the autonomous system
    bit32    userdeflen;    /* length of user-defined data
    bit8     userdef[SDV_MAXUSERDEF]; /* user defined data
};

typedef struct SDV Sdv;

```

Figure 5: MESSIAHS data structures

6.1 Distributed Batch

The MITRE distributed batch [1], Condor [11], and Remote Unix [2] systems support general-purpose distributed processing for machines running the UNIX operating system. Figure 6 lists a short specification file for a SPARC IPC participating in a distributed batching system. The state rules (lines 1–4) give the speed ratings for an IPC and the recalculation and revocation timeout periods.

The combining rules in lines 5 and 6 ensure that the processor type variable, `proctype`, contains the string `:SPARC` and that the operating system variable `OSname` contains the string `:SunOS4.1`. Lines 7 and 8 propagate incoming processor and operating system names.

The example schedule request filter (lines 9 and 10) computes a rating function in the range [0, 200] for the local system, and [0, 400] for remote systems. The scheduling request rules ensure that the processor type and operating system match, and assign a priority to a match based on the system load average. Because there is no provision for requesting tasks from a busy system, the section for task request rules is empty.

Hosts participating in the batch system preserve autonomy by varying the parameters of the schedule request filter. For example, tasks submitted by a local user can be given higher priority by basing the rating function on the source address of the task.

The task revocation rules (lines 12 and 13) determine, based on the computational load on the node, whether active tasks should be suspended, or whether suspended tasks should be returned to execution. The `true` guard in the revocation filter rule (line 10) matches any available task, and the value portion of the rule assigns an equal priority to all tasks under consideration.

6.2 Load Balancing

Several researchers have investigated load balancing and sharing policies for distributed systems, such as those described in [32], [33], and [34].

The *greedy load-sharing algorithm* [32], makes decisions based on a local optimum. When a user submits a task for execution, the receiving system attempts to place the task with a less busy neighbor, according to a weighting function. If no suitable neighbor is found, the task is accepted for local execution.

The suggested weighting function to determine if a task should be placed remotely is $f(n) = n \text{ div } 3$, where n is the number of tasks currently executing on the local system. The algorithm searches for neighbors whose advertised load is less than or equal to one-third the local load. Because the greedy algorithm depends on local state, it is dynamic.

The policy specification in figure 7 implements a variant of the greedy algorithm. The original algorithm used a limited probing strategy to collect the set of candidates for task reception. The version in figure 7 sets the recalculation and retransmission periods low (line 1), and depends on the SDV dissemination mechanism to determine the candidate systems.

The combination rules (lines 2 and 3) set the `$minload` field to be the minimum of the load advertised by neighbors and the local load. The filter assigns a low priority to local execution (line 4), and rates the neighboring systems on a scale of two through 100 (line 5). Any eligible neighbor takes precedence over local execution, but if the resultant candidate

```

    begin state
1.      float $SPECint92      13.8;
2.      float $SPECfp92      11.1;
3.      int   $recalc_period  30;
4.      int   $revocation_period 30;
    end
    begin combining
5.      string $out.proctype not match($out.proctype, "SPARC"):
        set $out.proctype + ":SPARC";
6.      string $out.OSname   not match($out.OSname, "SunOS4.1"):
        set $out.OSname + ":SunOS4.1";
7.      string $out.proctype
        not match($out.proctype, $sys.proctype):
        set $out.proctype + $sys.proctype;
8.      string $out.OSname   not match($out.OSname, $sys.OSname):
        set $out.OSname + $sys.OSname;
    end
    begin schedfilter
9.      $sys.address == $me.address and
        match($sys.proctype, $task.proctype) and
        match($sys.OSname, $task.OSname):
        max(200 - (100 * int($sys.loadave)), 0);
10.     match($sys.proctype, $task.proctype) and
        match($sys.OSname, $task.OSname):
        max(400 - (100 * int($sys.loadave)), 0);
    end
    begin revokefilter
11.     true: 1;
    end
    begin revokerules
12.     $me.loadave > 2.0 and $me.nactivetasks > 2: suspend;
13.     $me.loadave < 1.0 and $me.nuspendedtasks > 0: wake;
    end

```

Figure 6: MIL remote execution specification

Table 6: Predefined event handlers in MESSIAHS

Function Name	Corresponding Event
handle_msg_sr	sched request message
handle_msg_sa	sched accept message
handle_msg_sd	sched deny message
handle_msg_trq	task request message
handle_msg_ta	task accept message
handle_msg_td	task deny message
handle_msg_trv	task revoke message
handle_msg_ssq	system status query message
handle_msg_ssv	system status vector message
handle_msg_tsq	task status query message
handle_msg_tsv	task status vector message
handle_msg_jr	join request message
handle_msg_jd	join deny message
handle_input_timeout	input timeout
handle_output_timeout	output timeout
handle_recalc_timeout	recalculation timeout
handle_revoke_timeout	revocation timeout

```

begin state
1.   int $recalc_period 5;
end
begin combining
2.   int $out.minload ($sys.address == $me.address):
      set min($out.minload, $me.ntasks);

3.   int $out.minload true:
      set min($out.minload, $sys.minload);
end
begin schedfilter
4.   $sys.address == $me.address: 1;

5.   $sys.minload <= ($me.ntasks / 3):
      max(100 - $sys.minload, 2);
end

```

Figure 7: MIL specification for greedy load sharing

set is empty, the local system executes the task.

The greedy algorithm has no provision for task revocation; any tasks accepted run to completion. Thus, systems using the depicted specification yield some execution autonomy in the spirit of cooperation.

7 Example Algorithms Using the Toolkit

As stated earlier, MESSIAHS contains a set of event handlers which may be overloaded by the administrator to create a new policy. For example, the MESSIAHS prototype includes a default handler for schedule request message events. The administrator customizes the scheduling policy by writing a filter routine.

This section presents the implementation of three scheduling algorithms using the toolkit. Figure 8 lists the code for Arrival Balanced Scheduling [17], figure 9 lists the code for the greedy algorithm, and figure 10 lists the code for the BOS algorithm [16].

The Greedy algorithm was described in section 6. Arrival Balanced Scheduling assigns a task to the processor that will complete it first, as estimated by the scheduling host. The estimated runtime of the task, the current load on the target host, and the speed of the target host are used to estimate the finishing time of the task. The BOS algorithm employs a simple centralized scheme using a heuristic approach to schedule a task force on a set of homogeneous processors. The algorithm generates an initial mapping, then uses a bounded probabilistic approach to move towards the optimal solution.

The implementations of three algorithms demonstrate that the underlying mechanisms are easy to use and are flexible enough to support a wide variety of algorithms. The longest of the three algorithms, BOS, represents less than one-half of one percent of the code for the scheduling support module. Writing a new algorithm involves editing a code skeleton and inserting the algorithm code in a C `switch` statement. This process takes only a few minutes for a programmer familiar with the MESSIAHS code. In contrast, writing a scheduler from scratch, including data collection, data communication, and task management would take man-months of effort.

This ratio of schedule code size to support code size is consistent with that seen in other distributed scheduling support systems, such as Condor. However, MESSIAHS has ease-of-use advantages because of its separation of mechanism and policy, and because of its support for customizable scheduling policies.

Performance measurements were taken for each of the three algorithms, based on simulated tasks [7, chapter 6]. These results indicate, but do not prove, that the overhead incurred by use of the prototype is minor, typically less than 10% for dynamic algorithms and less than 40% for static algorithms. The 40% slowdown for a static algorithm may be acceptable in some environments, because the MESSIAHS version of the algorithm works in an environment the original static algorithm could not.

In addition, it appears that the MESSIAHS mechanisms perform better as the ratio of inter-task delay to update frequency increases. This increased ratio means that update information travels farther within the distributed system between task arrivals, and thus the scheduling modules are working with more up-to-date information.

```
nt = 0;
for (i = 0; i < SDV_NPROCCLASS; i++) {
    pp = &(psdv->procs[i]);
    if (pp->nsys > 0) {
        float ps, la, d;

        la = pp->qlen.min;
        ps = pp->specint92.mean;
        d = (la + 1) * ptdv->runtime * ptdv->specint92;
        value = (int) (1000 * ps / d);
        if (value > nt) {
            nt = value;
        }
    }
}
return nt;
```

Figure 8: Toolkit implementation of the ABS algorithm

```
if (sidmatch(psdv->sid, pmysdv->sid)) {
    return 1;
} else if (psdv->global_load <= (pmysdv->ntasks / 3)) {
    gl = (int) psdv->global_load;
    nt = (int) psdv->ntasks;
    value = ((100 - gl) * 1000) + (999 - nt);
    return(value);
} else {
    return 0;
}
```

Figure 9: Toolkit implementation of the greedy algorithm

```
i = 100000;

if (is_sibling(pste)) {
    return (0);
}

/* add in 'self' and 'target' */
if (pste == pmyste) {
    i -= ((pmysdv->ntasks + 1) * (pmysdv->ntasks + 1));
} else {
    i -= (pmysdv->ntasks * pmysdv->ntasks);
    gl = (pste->sdv.ntasks + 1) * (pste->sdv.ntasks + 1);
    nt = MAX(pste->sdv.nsys, 1);
    i -= (gl / nt);
}

for (pshte = sys_first(); pshte != (Shte *) NULL;
     pshte = sys_next(pshte)) {

    if (pshte->entry != pste) {
        nt = MAX(pshte->entry->sdv.nsys, 1);
        gl = pshte->entry->sdv.ntasks *
            pshte->entry->sdv.ntasks;
        i -= (gl / nt);
    }
}

return i;
```

Figure 10: Toolkit implementation of the BOS algorithm

8 Concluding Remarks

The mechanisms provided by the MESSIAHS system, MIL and the scheduling toolkit support global task scheduling and load sharing in scalable distributed systems. These mechanisms also protect the autonomy of the individual systems, while uniting heterogeneous machines into a coherent distributed system.

The language presented here is simple and expressive. It addresses two neglected areas of distributed scheduling, heterogeneity and autonomy. MIL supports a broad range of existing scheduling algorithms, while enabling rapid development, prototyping, and analysis of new policies.

Because of its simplicity, MIL is somewhat limited. It cannot store history and has no control flow or looping constructs. Because of this, scheduling algorithms that accept multiple tasks and a set of system descriptions as input cannot be expressed precisely using this language. MIL also assumes that neighbors can be trusted to tell the truth in their SDV advertisements, and depends on a model of timely information exchange.

A more complex approach that addresses these limitations, implemented as a set of library calls for high-level languages, is the scheduling toolkit described in section 5. The toolkit is a more complex interface to the underlying mechanisms than MIL is, but is also more expressive and efficient than MIL. Algorithms developed using MIL can be implemented and refined using the toolkit. Preliminary performance results obtained from the toolkit demonstrated that overhead of less than 10% is achievable for dynamic scheduling algorithms.

The prototype continues to evolve. The existing task environment is incompletely defined; in particular, the performance results were obtained using simulated tasks. The primary focus of current research on the prototype is to add support for task migration and execution, while still preserving as much autonomy as possible.

In summary, MESSIAHS embodies mechanisms supporting task placement in distributed, heterogeneous, autonomous systems. This support includes extensible mechanisms for implementing the local scheduling policy. This paper briefly described the MESSIAHS scheduling support mechanisms, defined a simple language and a library of function calls for constructing schedulers, and gave sample implementations of representative scheduling policies using these tools.

References

- [1] C. A. Gantz, R. D. Silverman, and S. J. Stuart. A Distributed Batching System for Parallel Processing. *Software-Practice and Experience*, 19, 1989.
- [2] M. J. Litzkow. Remote UNIX: Turning Idle Workstations Into Cycle Servers. In *USENIX Summer Conference*, pages 381–384, 2560 Ninth Street, Suite 215, Berkeley, CA 94710, 1987. USENIX Association.
- [3] A. H. Karp, K. Miura, and H. Simon. 1992 Gordon Bell Prize Winners. *IEEE Computer*, 26(1):77–82, January 1993.
- [4] T. L. Casavant and J. G. Kuhl. A Taxonomy of Scheduling in General-Purpose Distributed Computing Systems. *IEEE Transactions on Software Engineering*, 14(2):141–154, February 1988.

- [5] V. M. Lo. Heuristic Algorithms for Task Assignment in Distributed Systems. *IEEE Transactions on Computers*, 37(11):1384–1397, November 1988.
- [6] S. J. Chapin and E. H. Spafford. Constructing Distributed Schedulers with the MESSIAHS Interface Language. In *27th Hawaii International Conference on Systems Sciences*, volume 2, pages 425–434, Maui, Hawaii, January 1994.
- [7] S. J. Chapin. Scheduling Support Mechanisms for Autonomous, Heterogeneous, Distributed Systems. Ph.D. Dissertation, Purdue University, 1993.
- [8] S. J. Chapin and E. H. Spafford. Scheduling Support for an Internetwork of Heterogeneous, Autonomous Processors. Technical Report TR-92-006, Department of Computer Sciences, Purdue University, West Lafayette, IN, January 1992.
- [9] S. J. Chapin and E. H. Spafford. An Overview of the MESSIAHS Distributed Scheduling Support System. Technical Report TR-93-011 (supercedes TR-93-004), Department of Computer Sciences, Purdue University, West Lafayette, IN, January 1993.
- [10] A. Geist, A. Beguelin, J. Dongarra, W. Jiang, R. Manchek, and V. Sunderam. *PVM 3.0 User's Guide and Reference Manual*. Oak Ridge National Laboratory, Oak Ridge, TN 37831, February 1993.
- [11] A. Bricker, M. Litzkow, and M. Livny. Condor Technical Summary. Technical Report 1069, Department of Computer Science, University of Wisconsin-Madison, January 1992.
- [12] A. V. Aho, J. E. Hopcroft, and J. D. Ullman. *The Design and Analysis of Computer Algorithms*. Addison-Wesley Publishing Company, Reading, MA, 1974. ISBN 0-201-00029-6.
- [13] F. Eliassen and J. Veijalainen. Language Support for Multidatabase Transactions in a Cooperative, Autonomous Environment. In *TENCON '87*, pages 277–281, Seoul, 1987. IEEE Regional Conference.
- [14] W. Du, A. K. Elmagarmid, Y. Leu, and S. D. Ostermann. Effects of Local Autonomy on Global Concurrency Control in Heterogeneous Distributed Database Systems. In *Second International Conference on Data and Knowledge Systems for Manufacturing and Engineering*, pages 113–120. IEEE, 1989.
- [15] D. E. Knuth. *The Art of Computer Programming, Volume III: Searching and Sorting*. Addison-Wesley, 1973. ISBN 0-201-03803-X.
- [16] G. J. Bergmann and J. M. Jagadeesh. An MIMD Parallel Processing Programming System with Automatic Resource Allocation. In *Proceedings of the ISMM International Workshop on Parallel Computing*, pages 301–304, Trani, Italy, September 10–13 1991.
- [17] B. A. Blake. Assignment of Independent Tasks to Minimize Completion Time. *Software-Practice and Experience*, 22(9):723–734, September 1992.
- [18] F. Bonomi. On Job Assignment for a Parallel System of Processor Sharing Queues. *IEEE Transactions on Computers*, 39(7):858–869, July 1990.
- [19] F. Bonomi and A. Kumar. Adaptive Optimal Load Balancing in a Nonhomogeneous Multi-server System with a Central Job Scheduler. *IEEE Transactions on Computers*, 39(10):1232–1250, October 1990.

- [20] A. Drexl. Job-Processor-Scheduling für heterogene Computernetzwerke (Job-Processor Scheduling for Heterogeneous Computer Networks). *Wirtschaftsinformatik*, 31(4):345–351, August 1990.
- [21] A. Ghafoor and I. Ahmad. An Efficient Model of Dynamic Task Scheduling for Distributed Systems. In *Computer Software and Applications Conference*, pages 442–447. IEEE, 1990.
- [22] D. Hochbaum and D. Shmoys. A Polynomial Approximation Scheme for Scheduling on Uniform Processors: Using the Dual Approximation Approach. *SIAM Journal of Computing*, 17(3):539–551, June 1988.
- [23] C. C. Hsu, S. D. Wang, and T. S. Kuo. Minimization of Task Turnaround Time for Distributed Systems. In *Proceedings of the 13th Annual International Computer Software and Applications Conference*, 1989.
- [24] C. C. Price and M. A. Salama. Scheduling of Precedence-Constrained Tasks on Multiprocessors. *Computer Journal*, 33(3):219–229, June 1990.
- [25] S. Ramakrishnan, I. H. Cho, and L. Dunning. A Close Look at Task Assignment in Distributed Systems. In *INFOCOM '91*, pages 806–812, Miami, FL, April 1991. IEEE.
- [26] V. Sarkar and J. Hennessy. Partitioning Parallel Programs for Macro-Dataflow. In *ACM Conference on Lisp and Functional Programming*, pages 202–211, August 1986.
- [27] Vivek Sarkar and John Hennessy. Compile-time Partitioning and Scheduling of Parallel Programs. *SIGPLAN Notices*, 21(7):17–26, July 1986.
- [28] H. S. Stone. Multiprocessor Scheduling with the Aid of Network Flow Algorithms. *IEEE Transactions on Software Engineering*, SE-3(1):85–93, January 1977.
- [29] C. M. Wang and S. D. Wang. Structured Partitioning of Concurrent Programs for Execution on Multiprocessors. *Parallel Computing*, 16:41–57, 1990.
- [30] Standard Performance Evaluation Corporation. *The SPEC Newsletter*, published quarterly.
- [31] M. Bowman, L. L. Peterson, and A. Yeatts. Univers: An Attribute-Based Name Server. *Software-Practice and Experience*, 20(4):403–424, April 1990.
- [32] S. Chowdhury. The Greedy Load Sharing Algorithm. *Journal of Parallel and Distributed Computing*, 9:93–99, 1990.
- [33] D. L. Eager, E. D. Lazowska, and J. Zahorjan. A Comparison of Receiver-Initiated and Sender-Initiated Dynamic Load Sharing. Technical Report 85-04-01, University of Washington, Department of Computer Science, April 1985.
- [34] M. F. Pucci. Design Considerations for Process Migration and Automatic Load Balancing. Technical report, Bell Communications Research, 1988.