# Resource Management in Legion[*]

## Steve J. Chapin, Dimitrios Katramatos, John Karpovich, and Andrew Grimshaw

*Department of Computer Science, School of Engineering & Applied Science,*
*University of Virginia, Charlottesville, VA 22903–2442,*
*{chapin,dk3x,karp,grimshaw}@virginia.edu*

**Abstract**

The recent development of gigabit networking technology, combined with the proliferation of low-cost, high-performance microprocessors, has given rise to metacomputing environments. These environments can combine many thousands of hosts, from hundreds of administrative domains, connected by transnational and worldwide networks. Managing the resources in such a system is a complex task, but is necessary to efficiently and economically execute user programs.

In this paper, we describe the resource management portions of the Legion metacomputing system, including the basic model and its implementation. These mechanisms are flexible both in their support for system-level resource management but also in their adaptability for user-level scheduling policies. We show this by implementing a simple scheduling policy and demonstrating how it can be adapted to more complex algorithms.

*Keywords:* parallel and distributed systems, task scheduling, resource management, autonomy

## 1 Introduction

The recent development of gigabit networking technology, combined with the proliferation of low-cost, high-performance microprocessors, has given rise to metacomputing environments. These environments can combine many thousands of hosts, from hundreds of administrative domains, connected by local,

transnational, and world-wide networks. Managing the resources in such a system is a complex task, but is necessary to efficiently and economically execute user programs. The Legion project is developing metacomputing software, and in this paper, we will describe the resource management subsystem of Legion. In particular, we will describe the Legion scheduling model, our implementation of the model, and the use of these mechanisms to support user-level scheduling.

Legion [6] is an object-oriented metacomputing environment, intended to connect many thousands, perhaps millions, of hosts ranging from PCs to massively parallel supercomputers. Such a system will manage millions to billions of objects. To be successful, Legion will require much more than simply ganging computers together via gigabit channels—a sound software infrastructure must allow users to write and run applications in an easy-to-use, transparent fashion. Furthermore, the software must unite machines from thousands of administrative domains into a single coherent system. This requires extensive support for autonomy, so that we can assure administrators that they retain control over their local resources.

In a sense, then, we have two goals which can often be at odds: users want to ensure that their programs receive the best treatment, while administrators want to ensure that their systems are safe, secure, and available for their priority users. Legion provides a methodology allowing each group to express their desires, and the system acts as a mediator to find a resource allocation that is acceptable to both parties. With such a system in place, users may neither know nor care whether their jobs are running across the hall or across the country. Administrators can offer excess cycles to the Legion system, or even set up workstation farms selling cycles as a commodity, secure in the knowledge that their local access and use policies will be respected.

Legion achieves this vision through a flexible, modular approach to scheduling support. Throughout the paper, we will refer to the current implementation of Legion, or the default behavior. This is because Legion is fundamentally a set of interface definitions for an object system, and our prototype is only one implementation that manifests those interfaces. We fully expect others to reimplement or augment portions of the system, reflecting their needs for specific functionality. For scheduling, as in other cases, we provide reasonable default policies and allow users and system administrators to customize behavior to meet their needs and desires. Our mechanisms have cost that scales with capability—the effort required to implement a simple policy is low, and rises slowly, scaling commensurately with the complexity of the policy being implemented. This continuum is provided through a substrate rich in functionality that simplifies the implementation of scheduling algorithms.

Section 2 describes the Legion metacomputing system, and section 3 outlines

the resource management subsystem. We develop a Scheduler using Legion resource management in section 4, and describe other resource management systems for metacomputing in section 5. Finally, we give concluding remarks in section 6.

## 2  Legion

The Legion design encompasses ten basic objectives: site autonomy, support for heterogeneity, extensibility, ease-of-use, parallel processing to achieve performance, fault tolerance, scalability, security, multi-language support, and global naming. These objectives are described in greater depth in Grimshaw et al. [6]. Resource Management is concerned primarily with autonomy and heterogeneity, although other issues certainly play a role.

Supporting heterogeneity requires Legion to accommodate vastly differing computing capabilities among constituent machines, including differences in architectures, operating systems, and installed software. Such support is important to run complex distributed computations, such as a weather forecasting and visualization program—portions of the computation may be best suited for vector supercomputers, message-passing architectures, or graphics workstations. Autonomy means that each site has the freedom to have heterogeneous resources, define local policies, and refuse to run jobs from remote sites. Users have the freedom to choose where they would like their jobs to run, and to decline an unsatisfactory choice made by the system.

The resulting Legion design contains a set of *core* objects, without which the system cannot function, a subset of which are shown in figure 1. These objects are critical to resource management in that they provide the basic resources to be managed, and the infrastructure to support management. Between core objects and user objects lie *service* objects—objects which improve system performance, but are not truly essential to system operation. Examples of service objects include caches for object implementations, file objects, and the resource management infrastructure.

In the remainder of this section, we will examine the core objects and their role in resource management. For a complete discussion of the Legion Core Objects, see [10]. We will defer discussion of the service objects until section 3.
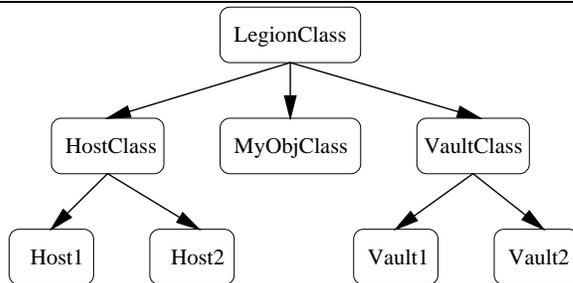
3

Fig. 1. The Legion Core Object Hierarchy

*2.1 Legion Core Objects*

Class objects (e.g. HostClass, LegionClass) in Legion serve two functions. As in other object-oriented systems, Classes define the types of their instances. In Legion, Classes are also active entities, and act as managers for their instances. Thus, a Class is the final authority in matters pertaining to its instances, including object placement. The Class defines the create_instance() method, which is responsible for placing an instance on a viable host. create_instance takes an optional argument suggesting a placement, which is necessary to implement external Schedulers. In the absence of this argument, the Class makes a quick (and almost certainly non-optimal) placement decision.[1]

The two remaining core objects represent the basic resource types in Legion: Hosts and Vaults.[2] Each has a corresponding guardian object class. Host Objects encapsulate machine capabilities (e.g., a processor and its associated memory) and are responsible for instantiating objects on the processor. In this way, the host acts as an arbiter for the machine's capabilities. Our current Host Objects represent single-host systems (both uniprocessor and multiprocessor shared memory machines), although this is not a requirement of the model. We are working with the Globus project and the NSF PACI centers to implement generic functionality that will allow Host Objects to interact with queue management systems such as LoadLeveler and Condor.

To support scheduling, Hosts grant reservations for future service. The exact form of the reservation depends upon the Host Object implementation, but they must be non-forgeable tokens; the Host Object must recognize these tokens when they are passed in with service requests. It is not necessary for any

---

[1] The current default is to place the object "here," i.e. using the class's Host and Vault, if possible.

[2] We are developing Network Objects to encapsulate host connectivity and interconnection.

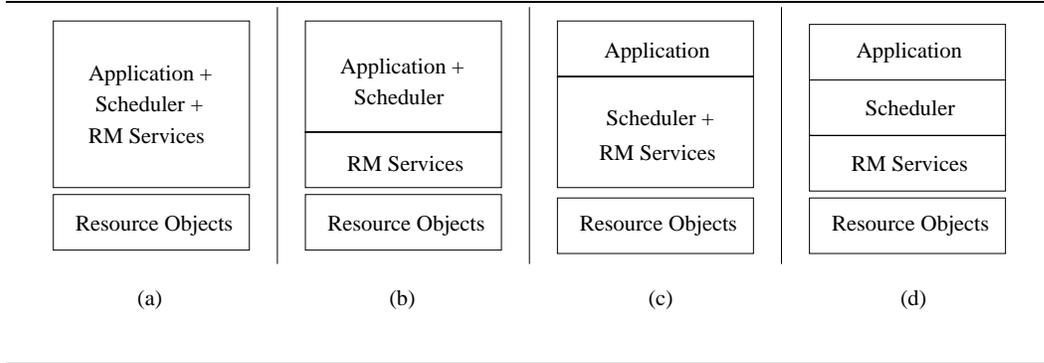| Application +<br>Scheduler +<br>RM Services | | Application +<br>Scheduler | | Application | | Application |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | | RM Services | | Scheduler +<br>RM Services | | Scheduler |
| | | | | | | RM Services |
| Resource Objects | | Resource Objects | | Resource Objects | | Resource Objects |
| (a) | | (b) | | (c) | | (d) |

Fig. 2. Choices in Resource Management Layering

other object in the system to be able to decode the reservation token. Our current implementation of reservations encodes both the Host and the Vault which will be used for execution of the object. Vaults are the generic storage abstraction in Legion. To be executed, a Legion object must have a vault to hold its Object Persistent Representation (OPR). The OPR holds the persistent state of the object, and is used for migration and for shutdown/restart purposes.

Hosts also contain a mechanism for defining event triggers—this allows a host to, e.g., initiate object migration if its load rises above a threshold. Conceptually, triggers are guarded statements which raise events if the guard evaluates to a boolean true. These events are handled by the Reflective Graph and Event (RGE) mechanisms in all Legion objects. RGE is described in detail in [14,15]; for our purposes, it is sufficient to note that this capability exists.

## 3 Resource Management Infrastructure (RMI)

Our philosophy of scheduling is that it is a negotiation of service between autonomous agents, one acting on the part of the application (consumer) and one on behalf of the resource or system (provider). This approach has been validated by both our own past history [4,8] and the more recent work of groups such as the AppLeS project at UCSD [1]. These negotiating agents can either be the principals themselves (objects or programs), or Schedulers and intermediaries acting on their behalfs. Scheduling in Legion is never of a dictatorial nature; requests are made of resource guardians, who have final authority over what requests are honored.

Figure 2 shows several different layering schemes that can naturally arise in metasystems. In part (a), the application does it all, negotiating directly with resources and making placement decisions. In part (b), the application still

5

makes its own placement decision, but uses the provided Resource Management services to negotiate with system resources. Part (c) shows an application taking advantage of a combined placement and negotiation module, such as was provided in MESSIAHS [4]. The most flexible layering scheme, shown in part (d), performs each of these functions in a separate module. Without loss of generality, we will write in terms of the third layering scheme, with the understanding that the Scheduler may be combined with other layers, thus producing one of the simpler layering schemes. Any of these layerings is possible in Legion; the choice of which to use is up to the individual application writer.

Legion provides simple, generic default Schedulers that offer the classic "90%" solution—they do an adequate job, but can easily be outperformed by Schedulers with special knowledge of the application. Application writers can take advantage of the resource management infrastructure, described below, to write per-application or application-type-specific user-level Schedulers. We are working with Weissman's group at UTSA [16] to develop Schedulers for broad classes of applications with similar structures (e.g. 5-point stencils).

Our resource management model, shown in figure 3, supports our scheduling philosophy by allowing user-defined Schedulers to interact with the infrastructure. The components of the model are the basic resources (hosts and vaults), the information database (the Collection), the schedule implementor (the Enactor), and an execution Monitor. Before we examine each component in detail, we will examine their interactions at a higher level. Note that figure 3 and the following discussion are intended to detail the logical components and steps involved in the scheduling process. Again, this description conforms to our implementation of the interfaces; others are free to substitute their own modules—for example, several components may be combined (e.g. the Scheduler or Enactor and the Monitor) for efficiency. The steps in object placement are as follows:

  (i)  The Collection is populated with information describing the resources.
 (ii)  The Scheduler queries the Collection, and
(iii)  based on the result and knowledge of the application, computes a mapping of objects to resources. This application-specific knowledge can either be implicit (in the case of an application-specific Scheduler), or can be acquired from the application's classes.
 (iv)  This mapping is passed to the Enactor, which
  (v)  invokes methods on hosts and vaults to
 (vi)  obtain reservations from the resources named in the mapping.
(vii)  After obtaining reservations, the Enactor consults with the Scheduler to confirm the schedule, and
(viii) after receiving approval from the Scheduler,
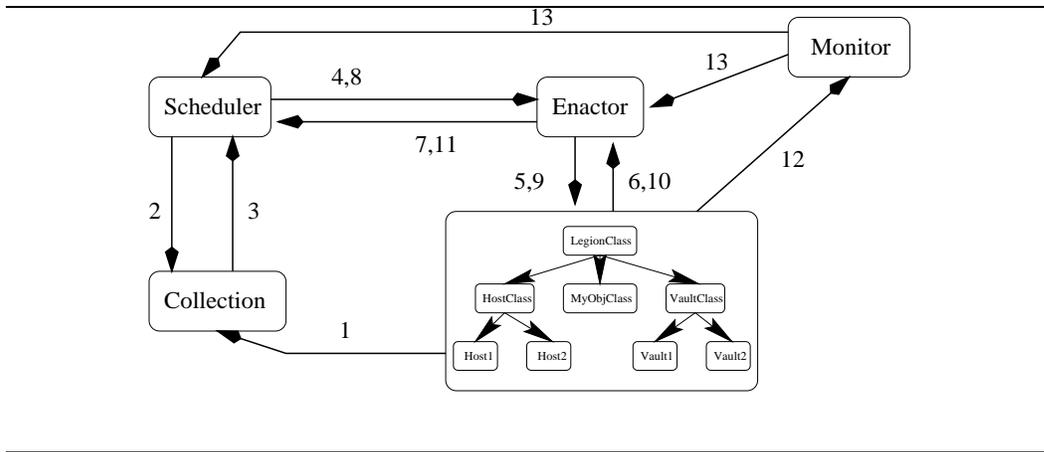 (ix)  attempts to instantiate the objects through member function calls on the

Fig. 3. Use of the Resource Management Infrastructure

     appropriate class objects.
  (x)  The class objects report success/failure codes, and
 (xi)  the Enactor returns the result to the Scheduler.
(xii)  If, during execution, a resource decides that the object needs to be migrated, it performs an outcall to a Monitor,
(xiii)  Which notifies the Scheduler and Enactor that rescheduling should be performed.

The remainder of this section examines each of the components in greater detail.

## 3.1   Host and Vault Objects

The resource management interface for the Host object appears in table 1. There are three broad groups of functions: reservation management, object management, and information reporting.

| Reservation Management | Process Management | Info. Reporting |
|---|---|---|
| make_reservation() | startObject() | get_compatible_vaults() |
| check_reservation() | killObject() | vault_OK() |
| cancel_reservation() | deactivateObject() | |

Table 1
Host Object Resource Management Interface

The reservation functions are used by the Enactor to obtain a reservation token for each subpart of a schedule. When asked for a reservation, the Host is responsible for ensuring that the vault is reachable, that sufficient resources

7

are available, and that its local placement policy permits instantiating the object.

In addition to the information reporting methods listed above, the Host also supports the attribute database included in all Legion objects. These information reporting methods for Host Objects allow us to build Collections using a pull model—the Collection can query the host to determine its current state. All Legion objects include an extensible attribute database, the contents of which are determined by the type of the object. Host objects populate their attributes with information describing their current state, including architecture, operating system, load, available memory, etc. Future versions of host objects will export scheduling policy information so that user-level Schedulers can better determine whether particular hosts are good candidates for object placement.

The Host Object reassesses its local state periodically, and repopulates its attributes. If a push model[3] is being used, it will then deposit information into its known Collection(s). The flexibility of Legion object attribute databases allows the Host Object to export a rich set of information, well beyond the minimal "architecture, OS, and load average" information used by most current scheduling algorithms. For example, the Host could export information such as the amount charged per CPU cycle consumed, domains from which it refuses to accept object instantiation requests, or a description of its willingness to accept extra jobs based on the time of day. This kind of information can help Schedulers to make better choices at the outset, thus avoiding the computation of subtly nonfeasible schedules.

The current implementation of Vault Objects does not contain dynamic state to the degree that Host Objects do. Vaults, therefore, only participate in the scheduling process at the start, when they verify that they are compatible with a host. They may, in the future, be differentiated by the amount of storage available, cost per byte, security policy, etc.

*3.2   The Collection*

The Collection acts as a repository for information describing the state of the resources comprising the system. Each record is stored as a set of Legion object attributes. As seen in figure 4, Collections provide methods to join (with an optional installment of initial descriptive information) and update records, thus facilitating a push model for data. The security facilities of Legion authenticate the caller to be sure that it is allowed to update the data in the

---

[3] Our current default is a push model, although we are implementing intermediate agents while will pull data from hosts and push it into collections.

```
int JoinCollection(LOID joiner);
int JoinCollection(LOID joiner, LinkedList <Uval_ObjAttribute>);
int LeaveCollection(LegionLOID leaver);
int QueryCollection(String Query, &CollectionData result);
int UpdateCollectionEntry(LOID constituent, LinkedList <Uval_ObjAttribute>);
```

Fig. 4. Collection Interface

| | | |
|---|---|---|
| *int-binop* | $\Rightarrow$ | + \| $-$ \| / \| * \| mod \| & \| \| \| |
| | | max \| min |
| *int-expr* | $\Rightarrow$ | *int-expr int-binop int-expr* \| |
| | | (*int-expr*) \| *integer* \| |
| | | int(*float-expr*) \| *id* |
| *string-expr* | $\Rightarrow$ | *string-expr* + *string-expr* \| |
| | | (*string-expr*) \| *string* \| *id* |
| *float-binop* | $\Rightarrow$ | + \| $-$ \| / \| * \| max \| min |
| *float-expr* | $\Rightarrow$ | *float-expr float-binop float-expr* \| |
| | | (*float-expr*) \| *float* \| |
| | | float(*int-expr*) \| *id* |
| *comp* | $\Rightarrow$ | < \| > \| = \| >= \| <= \| <> |
| *bool-binop* | $\Rightarrow$ | and \| or \| xor |
| *bool-expr* | $\Rightarrow$ | *bool-expr bool-binop bool-expr* \| |
| | | not *bool-expr* \| |
| | | *int-expr comp int-expr* \| |
| | | *float-expr comp float-expr* \| |
| | | *string-expr comp string-expr* \| |
| | | match(*string-expr, string-expr*) \| |
| | | (*bool-expr*) \| true \| false \| *id* |

Fig. 5. Grammar for Collection Query Language

Collection. As noted earlier, Collections may also pull data from resources. Users, or their agents, obtain information about resources by issuing queries to a Collection. A Collection query is string conforming to the grammar in figure 5, which is largely the same as that used in our earlier work [3]. This grammar allows typical operations (field matching, semantic comparisons, and boolean combinations of terms). Identifiers refer to attribute names within a particular record, and are of the form $AttributeName.

For example, to find all hosts that run the IRIX operating system version 5.x, one could use the regular expression matching feature for strings and query

as follows:

match($host_os_name, "IRIX") and match($host_os_name, "5\..*")

In its current implementation, the Collection is a passive database of static information, queried by Schedulers. We plan to extend Collections to support function injection—the ability for users to install code to dynamically compute new description information and integrate it with the already existing description information for a resource. This capability is especially important to users of the Network Weather Service [17], which predicts future resource availability based on statistical analysis of past behavior.

An important use of Collections is to structure resources within the Legion system. Having a few, global, Collections will prohibit the scalability we wish to achieve. Therefore, Collections may receive data from, and send data to, other Collections. This allows us to have a Collection for each administrative domain, and to combine Collections in other Collections. This is analogous to the hierarchical structuring of scheduling modules in [4], and we expect to see the same scalability benefits realized there.

*3.3  The Scheduler and Schedules*

The Scheduler computes the mapping of objects to resources. At a minimum, the Scheduler knows how many instances of each class must be started. Application-specific Schedulers may implicitly have more extensive knowledge about the resource requirements of the individual objects, and any Scheduler may query the object classes to determine such information (e.g., the available implementations, or memory or communication requirements). The Scheduler obtains resource description information by querying the Collection, and then computes a mapping of object instances to resources. This mapping is passed on to the Enactor for implementation. It is not our intent to directly develop more than a few widely-applicable Schedulers; we leave that task to experts in the field of designing scheduling algorithms. Our job is to build mechanisms that assist them in their task.

Schedules must be passed between Schedulers and Enactors. A graphical representation for a Schedule appears in figure 6. Each Schedule has at least one Master Schedule, and each Master Schedule may have a list of Variant Schedules associated with it. Both master and variant schedules contain a list of mappings, with each mapping having the type (Class LOID $\rightarrow$ (Host LOID x vault LOID)). Each mapping indicates that an instance of the class should be started on the indicated (host, vault) pair. In the future, this mapping process may also select from among the available implementations of an object as well.
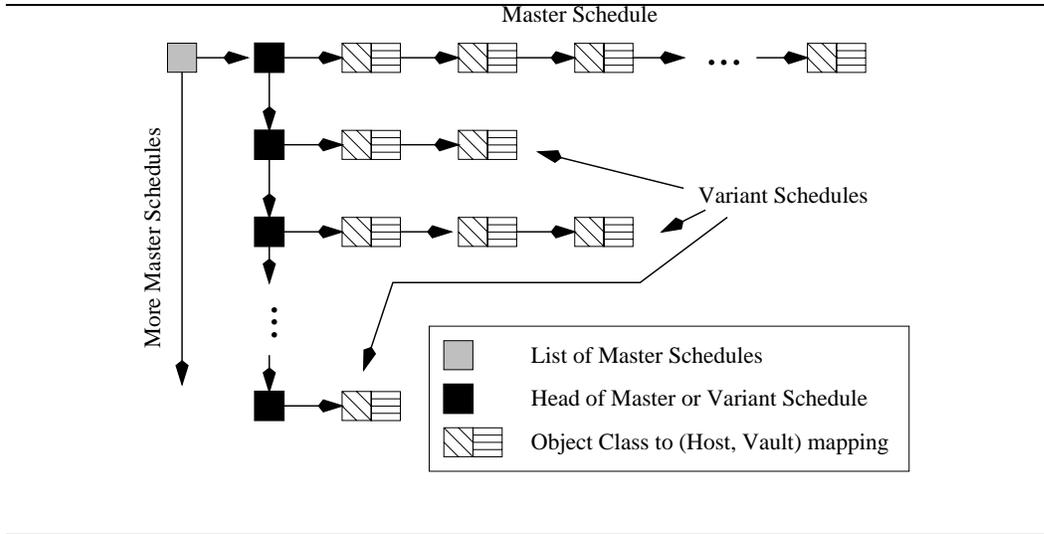
Fig. 6. The Schedule data structure

There are three important data types for interacting with the Enactor: the LegionScheduleFeedback, LegionScheduleList, and LegionScheduleRequestList. A LegionScheduleList is simply a single schedule (e.g. a Master or Variant schedule). A LegionScheduleRequestList is the entire data structure shown in figure 6. LegionScheduleFeedback is returned by the Enactor, and contains the original LegionScheduleRequestList and feedback information indicating whether the reservations were successfully made, and if so, which schedule succeeded.

## 3.4 The Enactor

The pertinent portion of the Enactor interface appears in figure 7. A Scheduler first passes in the entire set of schedules to the make_reservations() call, and waits for feedback. If all schedules failed, the Enactor may (but is not required to) report whether the failure was due to an inability to obtain resources, a malformed schedule, or other failure. If any schedule succeeded, the Scheduler can then use the enact_schedule() call to request that the Enactor instantiate objects on the reserved resources, or the cancel_reservations() method to release the resources.

We have mentioned master and variant schedules, but have not explained how they are used by the Enactor. Each entry in the variant schedule is a single-object mapping, and replaces one entry in the master schedule. If all mappings in the master schedule succeed, then scheduling is complete. If not, then a Variant schedule is selected that contains a new entry for the failed mapping. This Variant may also have different mappings for other instances,

11

```
&LegionScheduleFeedback make_reservations(&LegionScheduleList);
int cancel_reservations(&LegionScheduleRequestList);
&LegionScheduleRequestList enact_schedule(&LegionScheduleRequestList);
```

Fig. 7. Enactor Interface

which may have succeeded in the Master schedule. Implementing the Variant schedule entails making new reservations for items in the Variant schedule and canceling any corresponding reservations from the Master schedule. Our default Schedulers and Enactor work together to structure the Variant schedules so as to avoid reservation thrashing (the canceling and subsequent remaking of the same reservation).

As mentioned earlier, Class objects implement a create_instance() method. This method has an optional argument containing an LOID and a reservation token. Use of the optional argument allows directed placement of objects, which is necessary to implement externally computed schedules. The Class object is still responsible for checking the placement for validity and conformance to local policy, but the Class does not have to go through the standard placement steps.

### 3.5 Application Monitoring

As noted earlier, Legion provides an event-based notification mechanism via its RGE model [14]. Using this mechanism, the Enactor can register an outcall to the host objects; this outcall will be performed when a trigger's guard evaluates to true. There is no explicitly-defined interface for this functionality, as it is implicit in the use of RGE facilities. If desired, the Enactor or Scheduler can perform the monitoring, with the outcall registered appropriately.

## 4  Examples of Use

We now give an example of a Scheduler that uses our resource management infrastructure. While it does not take advantage of any application-specific knowledge, it does serve to demonstrate some of the flexibility of the mechanisms. We start with a simple random policy, and demonstrate how to build a "smarter" Scheduler based on the simple random policy. This improved Scheduler provides a template for building Schedulers with more complex placement algorithms. We then discuss our plans for building more sophisticated Schedulers with application and domain-specific knowledge.

```
Generate_Random_Placement(ObjectClass list) {
       for each ObjectClass 𝒪 in the list, do {
              query the class for available implementations
              query Collection for hosts matching available implementations
              k = the number of instances of this class desired
              for i := 1 to k, do {
                     pick a host ℋ at random
                     extract list of compatible vaults from ℋ
                     randomly pick a compatible vault 𝒱
                     append the target (ℋ, 𝒱) to the master schedule
              }      }
       return the master schedule
}
```

Fig. 8. Pseudocode for random placement

For the sake of brevity and presentation, we have omitted the full source code in favor of pseudocode. The source code is contained in release 1.4 of the Legion system, first made available in September 1998. The current release of the Legion software is available from [9].

*4.1  Random Scheduling*

The Random Scheduling Policy, as the name implies, randomly selects from the available resources that appear to be able to run the task. There is no consideration of load, speed, memory contention, communication patterns, or other factors that might affect the completion time of the task. The goal here is simplicity, not performance.

Pseudocode for our random schedule generator in figure 8. The Generate_Random_Placement() function is called with a list of classes for which instantiation is desired. The Scheduler iterates over this list, and executes the following steps for each item. First, the Scheduler extracts the list of available implementations from the attribute list of the class of the object it is to instantiate. The Scheduler then queries the Collection for matching hosts, and picks a matching host at random. After extracting that host's list of compatible vaults from the description returned by the Collection, the Scheduler randomly selects a vault. This (host, vault) pair is added to the master schedule. This pair selection is done once for each instance desired for this class.

Note that this algorithm only builds one master schedule, and does not take

advantage of the variant schedule feature, nor does it calculate multiple schedules. The Scheduler could call this function multiple times to generate additional master schedules. This is not efficient, nor will it necessarily generate a near-optimal schedule, but it is simple and easy. This is, in fact, the equivalent of the default schedule generator for Legion Classes in releases prior to 1.4.

After generating the mapping, the Scheduler must interact with the Enactor to determine if the placement was successful. Although not shown in figure 8, the simple implementation passes a single master schedule to the Enactor via the make_reservations() and enact_schedule() methods, and reports the success or failure of that call back to the object that invoked the Scheduler. No attempt is currently made to generate other placements, although a more sophisticated Scheduler would certainly do so.

*4.2  Improved Random Scheduling (IRS)*

There are many possible improvements on our random placement algorithm, both for efficiency of calculation and for efficacy of the generated schedule. The improvement we focus on is not in the basic algorithm; the IRS still selects a random host and vault pair. Rather, we will compute multiple schedules and accommodate negative feedback from the Enactor. The pseudocode for IRS is in figures 9 and 10.

The improved version generates $n$ random mappings for each object class, and then constructs $n$ schedules out of them. The Scheduler could just as easily build $n$ schedules through calls to the original generator function, but IRS does fewer lookups in the Collection. Note also that, because this is random placement, we do not consider dependencies between objects in the placement. A more sophisticated Scheduler would take this into account either when generating the individual instance mappings or when combining instance mappings into a schedule.

The Wrapper function has three global variables that limit the number of times it will try to generate schedules, the number of times it will attempt to enact each schedule, and the number of variant schedules generated per call to the generation function. [4] Again, this is a simple-minded approach to solving the problem, but serves to demonstrate how one could construct a richer Scheduler.

------

[4] We realize that the value returned from the generator and passed to the Enactor should be a list of master schedules; we take liberty with the types in the pseudocode for the sake of brevity.

```
IRS_Generate_Placement(ObjectClass list, int n) {
        for each ObjectClass O in the list, do {
                query the class for available implementations
                query Collection for hosts matching available implementations
                k = the number of instances of this object desired
                for l := 1 to n, do {
                        for i := 1 to k, do {
                                pick a host H at random
                                extract list of compatible vaults from H
                                randomly pick a compatible vault V
                                append the target (H, V) to the list for this instance
        }       }       }
        master schedule = first item from each object instance list
        for l := 2 to n, do {
                select the l^{th} component of the list for each object instance
                construct a list of all that do not appear in the master list
                append to list of variant schedules
        }
        return the master schedule
}
```

Fig. 9. Pseudocode for the IRS Placement Generator

```
IRS_Wrapper(ObjectClass list) {
        for i in 1 to SchedTryLimit, do {
                sched = IRS_Generate_Placement(ObjectClass List, NSched);
                for j in 1 to EnactTryLimit, do {
                        if (make_reservations(sched) succeeded) {
                                if (enact_placement(sched) succeeded) {
                                        return success;
        }       }       }       }
        return failure;
}
```

Fig. 10. Pseudocode for the IRS Wrapper

*4.3  Specialized Policies*

We are in the process of defining and implementing specialized placement
policies for structured multi-object applications. Examples of these applica-

tions include MPI-based or PVM-based simulations, parameter space studies, and other modeling applications. Applications in these domains quite often exhibit predictable communication patterns, both in terms of the compute/communication cycle and in the source and destination of the communication. For example, we are working with the DoD MSRC in Stennis, Mississippi to develop a Scheduler for an MPI-based ocean simulation which uses nearest-neighbor communication within a 2-D grid.

## 5  Related Work

The Globus project [5] is also building metacomputing infrastructure. At a high level, their scheduling model closely resembles that of Legion, as we first presented it at the 1997 Legion Winter Workshop [2]. There is a rough correspondence between Globus Resource Brokers and Legion Schedulers; Globus Information Services and Legion Collections; Globus Co-allocators and Legion Enactors; and Globus GRAMs and Legion Host Objects. However, there are substantial differences in realization of the model, due primarily to two features of Legion not found in Globus: the object-oriented programming model and strong support for local autonomy among member sites. Legion achieves its goals with a "whole-cloth" design, while Globus presents a "sum-of-services" architecture layered over pre-existing components. Globus has the advantage of a faster path to maturity, while Legion encompasses functionality not present in Globus.

There are many software systems for managing a locally-distributed multi-computer, including Condor [11] and LoadLeveler [13]. These systems are typically Queue Management Systems intended for use with homogeneous resource pools. While extremely well-suited to what they do, they do not map well onto wide-area environments, where heterogeneity, multiple administrative domains, and communications irregularities dramatically complicate the job of resource management. Indeed, these types of systems are complementary to a metasystem, and we will incorporate them into Legion by developing specialized Host Objects to act as mediators between the queuing systems and Legion at large.

SmartNet [7] provides scheduling frameworks for heterogeneous resources. It is intended for use in dedicated environments, such as the suite of resources available at a supercomputer center. Unlike Legion, SmartNet is not intended for large-scale systems spanning administrative domains. Thus, SmartNet could be used within a Legion system by developing a specialized Host Object, similar to the Condor and LoadLeveler Host Objects mentioned earlier. IBM's DRMS [12] also provides scheduling frameworks, in this case targeted towards reconfigurable applications. The DRMS components serve functions similar

to those of the Legion RMI, but like SmartNet, DRMS is not designed for wide-area metacomputing systems.

## 6    Conclusions and Future Work

This paper has described the resource management facilities in the Legion metacomputing environment. We have examined the components of the RM subsystem, presented their functionality, and described the interfaces of each component. Using these interfaces, we have implemented sample Schedulers, including a simple random Scheduler and a more sophisticated, but still random, Scheduler. These sample Schedulers point the way to building more complex and sophisticated Schedulers for real-world applications.

We are in the process of benchmarking the current system so that we can measure the improvement in performance as we develop more intelligent Schedulers. We expect to incorporate Network Objects as a core Legion resource in late 1998 or early 1999. The object interfaces will evolve in response to need—as we work with our research partners who are developing scheduling algorithms, we will enrich both the content and capability of the Resource Management Infrastructure and the Legion core objects.

## References

[1] F. Berman and R. Wolski. Scheduling from the perspective of the application. In *Proceedings of the 5th International Symposium on High-Performance Distributed Computing (HPDC-5)*. IEEE, August 1996.

[2] S. Chapin and J. Karpovich. Resource Management in Legion. Legion Winter Workshop. http://www.cs.virginia.edu/˜legion/WinterWorkshop/slides/ Resource_Management/, January, 1997.

[3] S. Chapin and E. Spafford. Support for Implementing Scheduling Algorithms Using MESSIAHS. *Scientific Programming*, 3:325–340, 1994. special issue on Operating System Support for Massively Parallel Computer Architectures.

[4] S. J. Chapin. Distributed Scheduling Support in the Presence of Autonomy. In *Proceedings of the 4th Heterogeneous Computing Workshop, IPPS*, pages 22–29, April 1995. Santa Barbara, CA.

[5] I. Foster and C. Kesselman. Globus: A metacomputing infrastructure toolkit. *International Journal of Supercomputer Applications*, to appear.

[6] A. S. Grimshaw, Wm. A. Wulf, and the Legion Team. The legion vision of a worldwide virtual computer. *Communications of the ACM*, 40(1), January 1997.

[7] D. Hensgen, L. Moore, T. Kidd, R. Freund, E. Keith, M. Kussow, J. Lima, and M. Campbell. Adding rescheduling to and integrating condor with smartnet. In *Proceedings of the 4th Heterogeneous Computing Workshop*, pages 4–11. IEEE, 1995.

[8] J. Karpovich. Support for object placement in wide area heterogeneous distributed systems. Technical Report CS-96-03, Dept. of Computer Science, University of Virginia, January 1996.

[9] Legion main web page. http://legion.virginia.edu.

[10] M. J. Lewis and A. S. Grimshaw. The core legion object model. In *Proceedings of the 5th International Symposium on High-Performance Distributed Computing (HPDC-5)*. IEEE, August 1996.

[11] M. Litzkow, M. Livny, and M. W. Mutka. Condor—A Hunter of Idle Workstations. In *Proceedings of the International Conference on Distributed Computing Systems*, pages 104–111, June 1988.

[12] J. E. Moreira and V. K. Naik. Dynamic resource management on distributed systems using reconfigurable applications. *IBM Journal of Research & Development*, 41(3), 1997.

[13] A. Prenneis, Jr. Loadleveler: Workload management for parallel and distributed computing environments. In *Proceedings of Supercomputing Europe (SUPEUR)*, October 1996.

[14] A. Nguyen-Tuong, S. J. Chapin and A. S. Grimshaw. Designing Generic and Reusable ORB Extensions for a Wide-Area Distributed System. The Eighth IEEE International Symposium on High-Performance Distributed Computing (HPDC), poster session, July, 1998.

[15] C. L. Viles, M. J. Lewis, A. J. Ferrari, A. Nguyen-Tuong, and A. S. Grimshaw. Enabling flexiblity in the legion run-time library. In *Proceedings of the International Conference on Parallel and Distributed Processing Techniques and Applications (PDPTA'97)*, pages 265–274, June 1997.

[16] J. Weissman and X. Zhao. Scheduling parallel applications in distributed networks. *Journal of Cluster Computing*, to appear.

[17] R. Wolski. Dynamically forecasting network performance to support dynamic scheduling using the network weather service. In *Proceedings of the 6th International Symposium on High-Performance Distributed Computing (HPDC-6)*, August 1997.