
Contents

1 Metacomputing: Harnessing Informal Supercomputers	1
<i>Mark Baker and Geoffrey Fox</i>	
1.1 General Introduction	1
1.1.1 Why do we need Metacomputing ?	2
1.1.2 What is a Metacomputer?	3
1.1.3 The Parts of a Metacomputer	3
1.2 The Evolution of Metacomputing	4
1.2.1 Introduction	4
1.2.2 Some Early Examples	5
1.3 Metacomputer Design Objectives and Issues	8
1.3.1 General Principles	9
1.3.2 Underlying Hardware and Software Infrastructure	9
1.3.3 Middleware - The Metacomputing Environment	10
1.4 Metacomputing Projects	13
1.4.1 Introduction	13
1.4.2 Globus	13
1.4.3 Legion	17
1.4.4 WebFlow	22
1.5 Emerging Metacomputing Environments	28
1.5.1 Introduction	28
1.5.2 Metacomputing Environments	28
1.5.3 Metacomputing Interfaces	34
1.5.4 Summary	35
1.6 Summary and Conclusions	35
1.6.1 Introduction	35
1.6.2 Summary of the Reviewed Metacomputing Environments	35
	i

1.6.3	Some Observations	36
1.6.4	Metacomputing Trends	37
1.6.5	The Impact of Metacomputing	37
1.7	Bibliography	38

Metacomputing: Harnessing Informal Supercomputers

MARK BAKER[†] AND GEOFFREY FOX[‡]

[†]Division of Computer Science
University of Portsmouth
Southsea, Hants, PO4 8JF

[‡]NPAC at Syracuse University
Syracuse, NY, 13244

Email {*Mark.Baker@port.ac.uk*, *gef@npac.syr.edu*}

1.1 General Introduction

The term metacomputing, when first encountered, seems a rather strange and typically 'geek' word. Its origin is believed to have been the CASA project, one of several US Gigabit testbeds around in 1989. Larry Smarr, the NCSA Director, is generally accredited with popularising the term thereafter.

A search through an ordinary dictionary to try and decypher the term would, at the time of writing, be fruitless. So, what does metacomputing mean? There seems to be many, sometimes conflicting, interpretations. Perhaps one should refer back to the Greek word 'meta' first to help understand the full word. Among the many meanings of 'meta', one will often find references to 'sharing' and 'action in common'. These words are the key to understanding the concept of metacomputing. Using these terms one can interpret metacomputing and understand it to be computers sharing and acting together to solve some common problem.

At this point, to reduce potential confusion, it is worth distinguishing between a parallel computer and a metacomputer. The key difference is the behaviour of individual computational nodes. A metacomputer is a dynamic environment that has some informal pool of nodes that can join or leave the environment whenever they desire. The nodes can be viewed as independent machines. So, in a slightly

confusing sense, a parallel computer, such as an IBM SP2, can be viewed as a 'metacomputer in a box'. Whereas, an SMP parallel computer, such as Tera MTA or Sun Enterprise 10000, cannot. The difference is that individual computational nodes in an SMP are not independent.

More recently Catlett and Smarr [1] have related the term metacomputing to "the use of powerful computing resources transparently available to the user via a networked environment". Their view is that a metacomputer is a networked virtual supercomputer. To an extent our usage of the term metacomputing still holds true to this definition apart from the need to explicitly refer to 'powerful computing resources'. Today's typical desktop computing resources can be viewed as powerful resources of yesterday.

The steps necessary to realise a metacomputer include:

- The integration of individual software and hardware resources into combined networked resource
- The implementation of middleware to provide a transparent view of the resources available
- The development and optimisation of distributed applications to take advantage of the resources.

1.1.1 Why do we need Metacomputing ?

The short answer to this question is that our computational needs are infinite, whereas our financial resources are finite. As we are all well aware, the sophisticated applications that we run on our desktops today seem to require more and more computing resources with every new revision. This trend is likely to continue as users and developers demand, for example, additional functionality or more realistic simulations. The net result is that from desktops to parallel supercomputers, users will always want more and more powerful computers. This is where metacomputing comes into its own. Why not try and utilise the potentially hundreds of thousands of computers that are inter-connected in some unified way? The realisation of this concept is the essence of metacomputing. It should be mentioned here that seamless access to remote resources is an uncontroversial topic. Whereas linking remote resources together to say execute a parallel application is more contentious as there are overheads involved. There are some situations where there is a strong case, for example, when linked resources must be geographically distributed such as in filtering and visualisation of scientific data. In general, it is fair to say that metacomputing comes with an efficiency penalty and it is usually better to run separate jobs on components of the metacomputer. However, linking of self contained applications is of growing importance in science (multi-disciplinary applications) and industry (linking different components of an organisation together).

It is not too futuristic to envisage that at some stage in the not so distant future individuals, be they engineers, scientists, students, health-care workers or business

persons, will be able to access the computing resources that they need to run their particular application with the same ease that we today switch on a light or turn on a kitchen appliance. Their application may be the simulation of the fluid flow around the after end of a ship, image processing of NMR scans, a Monte-Carlo financial simulation for a stockbroker, or a final year project for a student dissertation. It should be noted that these metacomputing resources can be accessed via a remote laptop or desktop [2].

1.1.2 What is a Metacomputer?

The simplest analogy to help describe a metacomputer is the electricity grid. When you turned on the power to your computer or switched on your television, you probably did not think about the original source of the electricity to drive these appliances. It was not necessary for you to select a generator with adequate capacity; or consider the gauge of wire used to connect the outlet or whether the power lines are underground or on pylons. Basically you were using a national power grid sophisticated enough to route the electrons across hundreds of miles, and easy enough for a child to use. In the same manner a metacomputer is a similarly easy-to-use assembly of networked computers that can work together to tackle a single task or set of problems.

It is not surprising, therefore, that the terms 'The Grid' and 'Computational Grids' are being used to describe a universal source of computing power [3]. A Grid can be viewed as the means to provide pervasive access to advanced computational resources, databases, sensors and people. It is believed that it will allow a new class of applications to emerge and will have a major impact on our approach to computing in the 21st century. For our purposes, Computational Grids are equivalent to metacomputing environments.

Metacomputing encompasses the following broad categories:

- Seamless access to high-performance resources.
- 'Parameter' studies (embarrassingly parallel application, see FAFNER, Section 1.2.2).
- The linkage of scientific instruments, analysis system, archival storage, visualisation (so called 4-way metacomputing, see I-WAY, Section 1.2.2).
- The general complex linkage of N distributed components.

1.1.3 The Parts of a Metacomputer

A metacomputer is a virtual computer architecture. Its constituent components are individually not important. The key concept is how these components work together as a unified resource. On an abstract level the metacomputer consists of the following components:

- *Processors and Memory*

The most obvious component of any computer system is the microprocessor that provides its computational power. The metacomputer will consist of an array of processors. Associated with each processor will be some dynamic memory.

- *Networks and Communications Software*

The physical connections between computers turn then from a collection of individual machines into an inter-connected network. The link between machines could, for example be via: modems, ISDN, standard Ethernet, FDDI, ATM or a myriad of other networking technologies. Networks with high bandwidth and low latency that provide rapid and reliable connections between the machines are the most favoured. To actually communicate over these physical connections it is necessary to have some communications software running. This software bridges all of the gaps, between different computers, between computers and people, even between different people.

- *Virtual Environment*

Given that we have an inter-connected, communicating, network of computers, processors with memory - there needs to be something like an operating system that can be used to configure, manage and maintain the metacomputing environment. This virtual environment needs to span the extent of the metacomputer and make it usable by both administrators and individual users. Such an environment will enable machines and/or instruments that may be located in the same building, or separated by thousands of miles, to appear as one system.

- *Remote Data Access and Retrieval*

In a metacomputing environment there is the potential that multiple supercomputers performing at Gflop/s are interacting with each other across national or international networks streaming GBytes of data in and out of secondary storage. This is a major challenge for any metacomputing environment. The challenge will become ever greater as new data-intensive applications are designed and deployed.

1.2 The Evolution of Metacomputing

1.2.1 Introduction

In this section we describe two early metacomputing projects that were in the vanguard of this type of technology. The projects differ in many ways, but both

had to overcome a number of similar hurdles, including communications, resource management and the manipulation of remote data, to be able to work efficiently and effectively. The two projects also attempted to provide metacomputing resources from opposite end of the computing spectrum. Whereas FAFNER [4] was capable of running on any workstation with more than 4 MBytes of memory, I-WAY [5] on the other hand was a means of unifying the resources of large supercomputing centres.

1.2.2 Some Early Examples

FAFNER

Public key cryptographic systems use two keys: a public and private key. A user must keep their private key a secret, but the public key is publicly known. Public and private keys are mathematically related, so that a message encrypted with a recipient's public key, can only be decrypted by their private key. The RSA algorithm [6] is an example of a public key algorithm. It is named after its developers Rivest, Shamir, and Adleman, who invented the algorithm at MIT in 1978.

The RSA keys are generated mathematically in part by combining prime numbers. The security of RSA is based on the premise that it is very difficult to factor extremely large numbers, in particular those with hundreds of digits. RSA keys use either 154 or 512-digit keys. The usage of this type of cryptographic technology has led to integer factorisation becoming an active research area. To keep abreast of the state-of-the-art in factoring, RSA Data Security Inc. initiated the RSA Factoring Challenge in March 1991. The Factoring Challenge provides a test-bed for factoring implementations and provides one of the largest collections of factoring results from many different experts worldwide.

Factoring is computationally very expensive. For this reason parallel factoring algorithms have been developed so that factoring can be distributed over a network of computational resources. The algorithms used are trivially parallel and require no communications after the initial set up. With this set up, it is possible that many contributors can provide a small part of a larger factoring effort. Early efforts relied on email to distribute and receive factoring code and information. More recently, in 1995, a consortium led by Bellcore Labs., Syracuse University and Co-Operating Systems started a project of factoring via the Web, known as FAFNER.

FAFNER was set up to factor RSA130 using a new numerical technique called the Number Field Sieve (NFS) factoring method using computational Web servers. The consortium produced a Web interface to NFS. A contributor then uses a Web-form to invoke server-side CGI scripts written in Perl. Contributors could, from one set of Web pages, access a wide range of support services for the sieving step of the factorisation: NFS software distribution, project documentation, anonymous user registration, dissemination of sieving tasks, collection of relations, relation archival services, and real-time sieving status reports. The CGI scripts produced supported cluster management, directing individual sieving workstations through appropriate day/night sleep cycles to minimise the impact on their owners. Contributors down-

loaded and built a sieving software daemon. This then became their Web client that used HTTP protocol to GET values from and POST the resulting relations back to a CGI script on the Web server.

Three factors combined to make this approach succeed.

1. The NFS implementation allowed even single workstations with 4 Mbytes to perform useful work using small bounds and a small sieve.
2. FAFNER supported anonymous registration - users could contribute their hardware resources to the sieving effort without revealing their identity to anyone other than the local server administrator.
3. A consortium of sites was recruited to run the CGI script package locally, forming a hierarchical network of RSA130 Web servers which reduced the potential administration bottleneck and allowed sieving to proceed around the clock with minimal human intervention.

The FAFNER project won an award in TeraFlop challenge at SC95 in San Diego. It paved the way for a wave of Web-based metacomputing projects, some of these are described in Sections 1.4 and 1.5.

I-WAY

The Information Wide Area Year (I-WAY) was an experimental high-performance network linking many high-performance computers and advanced visualisation environments. The I-WAY project was conceived in early 1995 with the idea not to build a network but to integrate existing high-bandwidth networks with telephone systems. The virtual environments, datasets, and computers used resided at seventeen different US sites and were connected by ten networks of varying bandwidths and protocols, using different routing and switching technologies.

The network was based on ATM technology, which at the time was an emerging standard. This network provided the wide-area backbone for various experimental networking activities at SC95, supporting both TCP/IP over ATM and direct ATM-oriented protocols.

To help standardise the I-WAY software interface and management, the key sites installed point-of-presence (I-POP) computers to serve as their gateways to the I-WAY. The I-POP machines were UNIX workstations configured uniformly and possessed a standard software environment called I-Soft. I-Soft helped overcome issues such as heterogeneity, scalability, performance, and security. The I-POP machines were the gateways into each site participating in the I-WAY project.

The I-POP machine provided uniform authentication, resource reservation, process creation, and communication functions across I-WAY resources. Each I-POP machine was accessible via the Internet and operated within its site's firewall. It also had an ATM interface that allowed monitoring and potential management of the site's ATM switch.

For the purpose of managing its resources efficiently and effectively, the I-WAY project developed a resource scheduler known as the Computational Resource Broker (CRB). The CRB basically consisted of user-to-CRB and CRB-to-local-scheduler protocols. The actual CRB implementation was structured in terms of a single central scheduler and multiple local scheduler daemons - one per I-POP machine. The central scheduler maintained queues of jobs and tables representing the state of local machines, allocating jobs to machine and maintaining state information on the AFS file system.

Security was a major feature of the I-WAY project. An emphasis was made on providing a uniform authentication environment. Authentication to I-POPs was handled by using a telnet client modified to use Kerberos authentication and encryption. In addition, the CRB acted as an authentication proxy, performing subsequent authentication to I-WAY resources on a user's behalf.

I-WAY used AFS to provide a shared file repository for software and scheduler information. An AFS cell was set up and made accessible from only I-POPs. To move data between machines where AFS was unavailable, a version of remote copy (`ircp`) was adapted for I-WAY.

To support user-level tools, a low-level communications library, Nexus, was adapted to execute in the I-WAY environment. Nexus supported automatic configuration mechanisms that enabled it to choose the appropriate configuration depending on the technology being used. For example, communications via TCP/IP or AAL5 when using the Internet or ATM. The MPICH and CAVEcomm libraries were also extended to use Nexus.

The I-WAY project was application driven and defined five types of applications:

- Supercomputer - Supercomputing
- Remote Resource - Virtual Reality
- Virtual Reality - Virtual Reality
- Multi-Supercomputer - multi-Virtual Reality
- Video, Web, GII-Windows

The I-WAY project was successfully demonstrated at SC'95 in San Diego. The I-POP machine was shown to simplify the configuration, usage and management of this type of wide-area computational test-bed. I-Soft was a success in terms that most applications ran, most of the time. More importantly, the experiences and software developed as part of the I-WAY project have been fed into the Globus project described in Section 1.4.2.

A Summary of Early Experiences

The projects described in this section both attempted to produce metacomputing environments by integrating hardware from opposite ends of the computing spectrum. FAFNER was a ubiquitous system that would work on any platform where a

Web server could be run. Typically its clients were at the low-end of the computing performance spectrum. Whereas I-WAY, unified the resources at supercomputing sites. The two projects also differed in the types of applications that could utilise their environments. FAFNER was tailored to a particular factoring application that was in itself trivially parallel and was not dependent on a fast interconnect. I-WAY, on the other hand, was designed to cope with a range of diverse high-performance applications that typically needed a fast interconnect. Both projects, in their way, lacked scalability. For example, FAFNER was dependent on quite a lot of human intervention to distribute and collect sieving results and I-WAY was limited by the design of components that made up I-POP and I-Soft.

FAFNER lacked a number of features that would now be considered obvious. For example, every client had to compile, link and run a FAFNER daemon in order to contribute to the factoring exercise. Today, one would probably download an already set up and configured Java applet. FAFNER was really a means of task-farming a large number of fine-grain computations. Individual computational tasks were unable to communicate with one-and-other, or with their parent Web-server. Again, today perhaps, using technology such as Java RMI, tasks would register themselves, ask for work, co-ordinate their computation, deliver results and so on with even less human intervention or interaction.

Likewise, with I-WAY a number of features would today seem inappropriate. The installation of an I-POP platform made it easier to set up I-WAY services in a uniform manner, but meant that each site needed to be specially set up to participate in I-WAY. In addition, the I-POP platform created one, of many, single-points-of-failure in the design of the I-WAY. Even though this was not reported to be a problem, the failure of an I-POP would mean that a site would drop out of the I-WAY environment. Today many of the services provided by the I-POP and I-Soft would be available on all the participating machines at a particular site.

Regardless of the aforementioned features of both FAFNER and I-WAY, both projects were highly successful. Each project was in the vanguard of metacomputing and has helped pave the way for many of the succeeding projects. In particular, FAFNER was the forerunner of projects such as WebFlow (described in Section 1.4.4) and the I-WAY software, I-Soft, was very influential on the approach used to design the components employed in the Globus Metacomputing Toolkit (described in Section 1.4.2).

1.3 Metacomputer Design Objectives and Issues

In this section we lay out and discuss the basic criteria required by all wide area distributed environments or a metacomputer. In the first part of this section we outline the underlying hardware and software technologies potentially being used. We then move on to discuss the necessary attributes of the middleware that creates the virtual environment we call a metacomputer.

1.3.1 General Principles

In attempting to facilitate the collaboration of multiple organisation running diverse autonomous heterogeneous resources a number of basic principles should be followed so that the metacomputing environment:

- Does not interfere with the existing site administration or autonomy.
- Does not compromise existing security of users or remote sites.
- Does not need to replace existing operating systems, network protocols or services.
- Allows remote sites to join or leave the environment whenever they choose
- Does not mandate the programming paradigms, languages, tools or libraries that a user wants.
- Provides a reliable and fault tolerance infrastructure with no single point of failure.
- Provides support for heterogeneous components.
- Uses standards, existing technologies, and is able to interact with legacy applications,
- Provides appropriate synchronisation and component program linkage.

1.3.2 Underlying Hardware and Software Infrastructure

As one would expect, a metacomputing environment must be able to operate on top of the whole spectrum of current and emerging hardware and software technologies. An obvious analogy is the Web. Users of the Web do not care if the server they are accessing is on a UNIX or NT platform. They are probably unaware that they are using HTTP on top of TCP/IP, and they certainly do not want to know that they are accessing a database supported by a parallel computer, such as an IBM SP2, or an SMP, such as the SGI Origin 2000. From the client browser's point-of-view, they 'just' want their requests to Web services handled quickly and efficiently. In the same way a user of a metacomputer does not want to be fussed with details of its underlying hardware and software infrastructure. A user is really only interested in submitting their application to the appropriate resources and getting correct results back in a timely fashion.

An ideal metacomputing environment will therefore provide access to the available resources in a seamless manner such that physical discontinuities such as differences between platforms, network protocols, administrative boundaries become completely transparent. In essence the metacomputing middleware turns a radically heterogeneous environment into a virtual homogeneous one.

1.3.3 Middleware - The Metacomputing Environment

In this section we outline and describe the idealised design features that are required by a metacomputing system to provide users with a seamless computing environment.

Administrative Hierarchy

An Administrative Hierarchy is the way that each metacomputing environment divides itself up to cope with a potentially global extent. For example, DCE uses cells and DNS has a hierarchical namespace. The reasons why this category is important stems from the administrative need to provide resources on autonomous systems on a global basis. The Administrative Hierarchy determines how administrative information flows through the metacomputer. For example, how does the resource manager find its resources ? Does it interrogate one global database of resources or a hierarchy of servers, or perhaps servers configured in some peer-related manner?

Communication Services

The communication needs of applications using a metacomputing environment are diverse, they can range from reliable point-to-point to unreliable multicast communications. The communications infrastructure needs to support protocols that are used for bulk-data transport, streaming data, group communications and those used by distributed objects.

This communication services provide the basic mechanisms needed by the metacomputing environment to transport administrative and user data. The network services used also provide the metacomputer with important Quality of Service parameters, such as latency, bandwidth, reliability, fault-tolerance and jitter control. Typically the network services will be built from a relatively low-level communication API that can be used to support a wide range of high-level communication libraries and protocols. These mechanisms provide the means to implement a wide range of communications methodologies, including RPC, DSM, stream-based and multicast.

Directory/Registration Services

A metacomputer is a dynamic environment where the location and type of services available are constantly changing. A major goal is to make all resources accessible to any process in the system, without regard to the relative location of the resource user. It is necessary to provide mechanisms to enable a rich environment in which information about metacomputing is reliably and easily obtained by those services requesting the information. The Registration and Directory Services components provide the mechanisms for registering and obtaining information about the metacomputer structure, resources, services and status.

Processes, Threads and Concurrency Control

The term process originates in the literature on the design of operating systems and is generally considered as a unit of resource allocation both for CPU and memory. The advent of shared memory multi-processors brought about the provision of *Light Weight Processes* or *Threads*. The name thread comes from the expression 'thread of control'. Modern OSs, like NT, permit an OS process to have multiple threads of control. With regards to metacomputers, this category is related to the granularity of control provided by the environment to its applications. Of particular interest is the methodology used to share data and maintain its consistency when multiple processes or threads have concurrent access to it.

Time and Clocks

Time is an important concept in all systems. Firstly, time is an entity that we wish to measure accurately, as it may be a record of when a particular transaction occurred. Or, if two or more computer clocks are synchronised it can be used to measure the interval when two or more events occurred. Secondly, algorithms have been developed that depend on clock synchronisation. These algorithms may, for example be used for maintaining the consistency of distributed data or as part of the Kerberos authentication protocol.

Naming Services

In any distributed system, names are used to refer to a wide variety of resources such as computers, services or data objects. The Naming Service provides a uniform name space across the complete metacomputing environment. Typical naming services are provided by the international X.500 naming scheme or DNS, the Internet's scheme.

Distributed Filesystems and Caching

Distributed applications, more often than not, require access to files distributed among many servers. A distributed filesystem is therefore a key component in a distributed system. From an applications point of view it is important that a distributed file system can provide a uniform global namespace, support a range of file I/O protocols, require little or no program modification and provide means that enable performance optimisations to be implemented, such as the usage of caches.

Security and Authorisation

Any distributed system involves all four aspects of security: *confidentiality* - prevents disclosure of data; *integrity* - prevents tampering with data, *authentication* - verifying identity and *accountability* - knowing whom to blame. Security within a metacomputing environment is a complex issue requiring diverse resources autonomously administered to interact in a manner that does not impact on the usability of the resources or introduce security holes in individual systems or the

environments as a whole. A security infrastructure is key to the success or failure of a metacomputing environment.

System Status and Fault Tolerance

There is a very high likelihood that some component in a metacomputing environment will fail. To provide a reliable and robust environment it is important that a means of monitoring resources and applications is provided. For example, if a particular platform goes out-of-service, it is important that no further jobs are scheduled on it until it becomes in-service again. In addition, jobs that were running on the system when it crashed should be re-run when it is available again or re-scheduled onto an alternative system. To accomplish this task, tools that monitor resources and application need to be deployed. So, when a platform is unavailable, information is passed to the directory services, or perhaps, when a job crashes, some part of the system reschedules that job to run again.

Resource Management and Scheduling

The management of processor time, memory, network, storage and other components in a distributed system is clearly very important. The overall aim is to efficiently and effectively schedule the applications that need to utilise the available resources in the metacomputing environment. From a user's point of view, resource management and scheduling should be almost transparent, their interaction with it being confined to a manipulating mechanism for submitting their application. It is important in a metacomputing environment that a resource management and scheduling service can interact with those that may be installed locally. For example, it may be necessary to operate in conjunction with LSF, Codine, or Condor at different remote sites.

Programming Tools and Paradigms

Ideally, every user will want to use a diverse range of programming paradigms and tools with which to develop, debug, test, profile, run and monitor their distributed application. A metacomputing environment should include interfaces, APIs and conversion tools so as to provide a rich development environment. Common scientific languages such as C, C++ and Fortran should be available, as should message passing interfaces like MPI and PVM. A range of programming paradigms should be supported, such as message passing and distributed shared memory. In addition a suite of numerical and other commonly used libraries should be available.

User and Administrative GUI

The interfaces to the services and resources available should be intuitive and easy to use. In addition they should work on range of different platforms and operating systems.

Availability

Earlier in this section we mentioned the need to provide middleware that provided heterogeneous support. In particular we are concerned about issues, such as does a particular resource management system work on a particular operating system, or perhaps will the communication services run on top of particular network architecture such as Novell or SNA. The issues that relate to this category are those that relate to the portability of the software services provided by the metacomputing environment. The metacomputing software should be either be easily 'ported' on to a range of commonly used platforms or should use technologies that enable it to be platform neutral, in a similar manner to Java Byte-code.

1.4 Metacomputing Projects

1.4.1 Introduction

In this section we map the techniques and technologies that three representative current metacomputing environments use with the aid of the design objectives and issues laid out in the previous section. The main purpose of this template is to help the reader review the methodologies used by each project. The three projects reviewed in this section are: Globus from Argonne National Laboratory, Argonne National Laboratory Legion from the University of Virginia University of Virginia and WebFlow from Syracuse University Syracuse University. The reasons that these three particular project to be chosen were:

- **Globus** - provides a toolkit based on a set of existing components with which to build a metacomputing environment.
- **Legion** - provides a high-level unified object model out of new and existing components to build a metasystem.
- **WebFlow** - provides a Web-based metacomputing environment.

1.4.2 Globus

Introduction

Globus [8] [9], provides a software infrastructure that enables applications to handle distributed, heterogeneous computing resources as a single virtual machine. The Globus project is a US multi-institutional research effort that seeks to enable the construction of computational grids. A computational grid, in this context, is a hardware and software infrastructure that provides dependable, consistent, and pervasive access to high-end computational capabilities, despite the geographical distribution of both resources and users. A central element of the Globus system is the Globus Metacomputing Toolkit (GMT), which defines the basic services and capabilities required to construct a computational grid. The toolkit consists of a

set of components that implement basic services, such as security, resource location, resource management and communications.

It is necessary for computational grids to support a wide variety of applications and programming paradigms. Consequently, rather than providing a uniform programming model, such as the object-oriented model, the GMT provides a bag of services from which developers of specific tools or applications can use to meet their own particular needs. This methodology is only possible when the services are distinct and have well-defined interfaces (API), that can be incorporated into applications or tools in an incremental fashion.

Globus is constructed as a layered architecture in which high-level global services are built upon essential low-level core local services. The Globus toolkit is modular, and an application can exploit Globus features such as resource management or information infrastructure without using the Globus communication libraries.

The GMT currently consists of the following:

- Resource allocation and process management (GRAM).
- Unicast and multicast communications services (Nexus).
- Authentication and related security services (GSI).
- Distributed access to structure and state information (MDS).
- Monitoring of health and status of system components (HBM).
- Remote access to data via sequential and parallel interfaces (GASS).
- Construction, caching and location of executables (GEM)

Administrative Hierarchy

Globus has no obvious administrative hierarchy. Every Globus-enabled resource is a peer of every other enabled resource.

Communication Services

Communication services within Globus are provided by Nexus, a communication library that is designed specifically to operate in grid environment. Nexus is distinguished by its support for multi-method communication, providing an application a single API to a wide range of communication protocols and characteristics. Nexus defines a relatively low-level communication API that can be used to support a wide range of high-level communication libraries and languages. Nexus communication services are used extensively in other parts of the Globus toolkit.

Directory/Registration Services

The Globus Metacomputing Directory Service (MDS) provides information about the status of Globus system components. MDS is part of the information infrastructure of the GMT and is capable of storing static and dynamic information about

the status of a metacomputing environment. MDS uses a Lightweight Directory Access Protocol [10] (LDAP) server that can store metacomputing-specific objects. LDAP is a streamlined version of the X.500 directory service. The MDS houses information pertaining to the potential computing resources, their specifications, and their current availability.

Processes, Threads and Concurrency Control

Globus works at the process level. The Nexus API can be used to construct communication primitives between threads. There is no concurrency control in Globus.

Time and Clocks

Globus does not mandate the usage of a particular time service and relies on those already used at each site.

Naming Services

Globus makes extensive usage of LDAP as well as DNS and X.500.

Distributed Filesystems and Caching

The Globus system currently provides three interfaces for remote access of user data:

- Global Access to Secondary Storage (GASS) provides basic access to remote files. Operations supported include remote read, remote write and append.
- Remote I/O - The RIO library implements a distributed implementation of the MPI-IO, parallel I/O API.
- Globus Executable Management (GEM) enables loading and executing a remote file through GRAM using GASS caching calls.

The Remote I/O for Metasystems (RIO) library provides basic mechanisms for tools and applications that require high-performance access to data located in remote, potentially parallel file systems. RIO implements the Abstract I/O (ADIO) device interface specification, which defines basic I/O functionality that can be used to implement a variety of higher-level I/O libraries. ROMIO has adopted the parallel I/O interface defined by the MPI forum in MPI-IO and hence allows any program already using MPI-IO to work without unchanged in a wide-area environment. The RIO library has been developed as part of the GMT, although it can also be used independently. ROMIO can be used with Nexus communications, GSI security and MDS to provide configuration information.

The GMT data movement and access service, GASS, defines a global name space via URLs, allows access to remote files via standard I/O interfaces and provides specialised support for data-movement in a wide-area environment. GASS addresses

bandwidth management issues associated with repeated access to remote files by providing a file cache: where a 'local' copy of remote file can be stored. Files are moved in to and out of the cache when a file is opened or closed by an application. GASS uses a simple locking protocol for local concurrency control, but does not implement a wide-area cache coherency mechanism.

Security and Authorisation

Globus uses an authentication system known as the Generic Security Service API (GSI) using an implementation of the Secure Sockets Layer. This system uses the RSA encryption algorithm and the associated public and private keys. The GSI authentication relies on an X509 certificate, provided by the user in their directory, that identifies them to the system. This certificate includes information about the duration of the permissions, the RSA public key, and the signature of the Certificate Authority (CA). With the certificate is the user's private key. The certificates can be created only by the CA, who reviews the X509 certificate request submitted by the user, and accepts or denies it according to an established policy.

System status and Fault Tolerance

Globus provides a range of basic services designed to enable the construction of application specific fault recovery mechanisms. In Globus it is currently assumed detection of a fault is a necessary prerequisite to fault recovery or fault tolerance. The main fault detection service in Globus is the Heartbeat Monitor (GHM) that enables a process to be monitored and periodic heartbeats to be sent to one or more monitors. The Nexus communication library also provides support for fault detection.

Resource Management and Scheduling

The Globus Resource Allocation Manager (GRAM) is the lowest level of Globus architecture. GRAM allows jobs to run remotely and provides an API for submitting, monitoring, and terminating jobs. GRAM provides the local component for resource management and is responsible for the set of resources operating under the same site-specific allocation policy. Such a policy will often be implemented by a local resource management package, such as LSF, Codine, or Condor.

GRAM is responsible for:

- Parsing and processing the Resource Specification Language (RSL) specifications that outline job requests. The request specifies resource selection, job process creation, and job control. This is accomplished by either denying the request or creating one or more processes (jobs) to satisfy the request. The RSL is a structured language that can be used to define resource requirements and parameters by a user.
- Enabling remote monitoring and managing of jobs already created.

- Updating MDS with information regarding the availability of the resources it manages.

Programming Tools and Paradigms

Globus currently supports: MPI, Java, Compositional C++, Simple RPC and Perl. There are on-going efforts to add a Sockets API, an IDL, Legion and Netsolve.

User and Administrative GUI

Globus makes extensive usage of the Web and command line interfaces for administration. for example, LDAP can be browsed via the Web. There are also a growing number of Java components that can be used with Globus.

Availability

Globus is available on most versions of UNIX and is currently being developed for NT.

1.4.3 Legion

Introduction

Legion [11] [12], is an object-based metasystem developed at the University of Virginia. Legion provides the software infrastructure so that a system of heterogeneous, geographically distributed, high-performance machines can interact seamlessly. Legion attempts to provide a user, at their workstations, with a single, coherent, virtual machine. The Legion system is organised by classes and metaclasses (classes of classes).

In Legion:

- *Everything is an object* - Objects represent all hardware and software components. Each object is an active process that responds to method invocations from other objects within the system. Legion defines an API for object interaction, but not the programming language or communication protocol.
- *Classes manage their instances* - Every Legion object is defined and managed by its own active class object. Class objects are given system-level capabilities; they can create new instances, schedule them for execution, activate or deactivate an object as well as provide state information to client objects.
- *Users can define their own classes* - As in other object-oriented systems users can override or redefine the functionality of a class. This feature allows functionality to be added or removed to meet a user's needs.
- *Core objects* - Legion defines the API to a set of core objects that support the basic services needed by the metasystem.

Legions has the following set of core object types:

- *Classes and Metaclasses* - Classes can be considered managers and policy makers. Metaclasses are classes of classes.
- *Host objects* - Host objects are abstractions of processing resources, they may represent a single processor or multiple hosts and processors.
- *Vault objects* - Vault objects represents persistent storage, but only for the purpose of maintaining the state of Object Persistent Representation (OPR).
- *Implementation Objects and Caches* - Implementation objects hide the storage details of object implementations and can be thought of as equivalent to executable files in UNIX. Implementation cache objects provide objects with a cache of frequently used data.
- *Binding Agents* - A Binding agent maps object IDs to physical address. Binding agents can cache bindings and organise themselves in hierarchies and software combining trees.
- *Context objects and Context spaces* - Context objects map context names to Legion object IDs, allowing users to name objects with arbitrary-length string names. Context spaces consist of directed graphs of context objects that name and organise information.

A Legion object is an instance of its class. Objects are independent, active and capable of communicating with each other via unordered non-blocking calls. Like other object-oriented systems, the set of methods of an object describes its interface. The Legion interfaces are described in an Interface Definition Language (IDL) .

A Legion object can be in one of two different states, active or inert. An active object runs as a process that is ready to accept function invocations. An inert object is represented by an Object Persistent Representation (OPR) . An OPR is an image of the object which resides on some stable storage, this is analogous to a process that has been swapped-out to disk. In a similar manner an OPR contains state information that enables the object to be reactivated.

Legion implements a three-tiered naming system.

1. Users refer to objects using human-readable strings, called context names.
2. Context objects map context names to LOIDs (Legion object identifiers), which are location-independent identifiers that include an RSA public key.
3. A LOID is mapped to an LOA (Legion object address) for communication. A LOA is a physical address (or set of addresses in the case of a replicated object) that contains sufficient information to allow other objects to communicate with the object (e.g., an IP address and port number pair).

Administrative Hierarchy

Legion has no obvious administrative hierarchy. Objects distributed about the Legion environment are peers to one another.

Communication Services

Legion uses standard TCP/IP to support communications between objects. Every Legion object is linked with a UNIX sockets-based delivery layer, called the Modular Message Passing System (MMPS) .

Directory/Registration Services

A Binding agent in Legion maps LOIDs to LOAs. A LOID/LOA pair is called a binding. Binding agents can cache bindings and organise themselves in hierarchies and software combining trees.

Processes, Threads and Concurrency Control

Currently Legion has one process per active object and objects communicate via MMPS. There is no concurrency control included in Legion.

Time and Clocks

Legion does not mandate the usage of a particular time service and relies on those already used at each site.

Naming Services

Legion Context objects map context names to LOIDs, allowing users to name objects with arbitrary-length string names. A LOID is mapped to an LOA for communication purposes. A LOA consists of an IP address and port number. It is assumed that Legion uses DNS to translate names to IP addresses.

The Context Manager is a Java GUI that can be used to manage context space. Context space is organised into a series of sub-contexts (also called contexts) and each context contains context names of various Legion objects. In the Context Manager all context-related objects such as contexts, file objects, and objects are represented by icons that can be manipulated. Basic context manager commands are Move, Alias, Get Interface, Get Attributes, Destroy, Activate and Deactivate.

Distributed Filesystems and Caching

Legion provides a virtual file system that spans all the machines in a Legion system. I/O support is provided via a set of library functions with UNIX-like file and stream operations to read, write and seek. These functions provide location independent, secure, access to context space and to 'files' in the system. Different users can also use the virtual file system to collaborate, sharing data files and even accessing the same running computations.

Legion has a special core object called a vault object. This represents persistent storage, but only for the purpose of maintaining the state of OPRs. The vault object may manage a portion of a UNIX filesystem, or a set of databases.

Security and Authorisation

Legion does not require any special privileges from the host systems that run it. The Legion security model is oriented towards protecting objects and object communication. Objects are accessed and manipulated via method calls; an object's rights are centred in its capabilities to make those calls. The user determines the security policy for an object by defining the object's rights and the method calls they allow. Once this is done, Legion provides the basic mechanism for enforcing that policy.

Every object in Legion supports a special member function called `MayI`. An object with no security will have a null `MayI`. All method invocations to an object must first pass through `MayI` before the target member function is invoked. If the caller has the appropriate rights for the target method, `MayI` allows that method invocation to proceed.

To make rights available to a potential caller, the owner of an object gives it a certificate listing the rights granted. When the caller invokes a method on the object, it presents the appropriate certificate to `MayI`, which then checks the scope and authenticity of the certificate. Alternatively, the owner of an object can permanently assign a set of rights to a particular caller or group. `MayI` is responsible for confirming the identity of a caller and its membership of an authorised group, followed by comparing the rights authorised with the rights required for the method call.

To provide secure communication, every Legion object has a public key pair; the public key is part of the object's name. Objects can use the public key of a target object to encrypt their communications to it. Likewise, an object's private key can be used to sign messages. This ensures authentication and integrity. This integration of public keys into object names eliminates the need for a certification authority. If an intruder tries to tamper with the public key of a known object, the intruder will create a new and unknown name.

System Status and Fault Tolerance

Legion does not mandate any fault-tolerance policies, applications are responsible for selecting the level they need. Fault tolerance will be built into generic base classes and applications will be able to invoke methods that provide the functionality that they require. Legion will support object reflection, replication and check pointing for the purposes of fault tolerance.

Resource Management and Scheduling

Host objects represent processors, and more than one may run on each computing resource. Host objects create and manage processes for active Legion objects. Classes invoke the member functions on host objects in order to activate instances on the computing resources that the hosts represent. Legion provides resource owners with the ability to initiate, manage and control and kill their resources.

The Legion-scheduling module consists of three components:

- A resource state information database (*Collection*)
 - A module which maps request to resources (*Scheduler*)
 - An agent responsible for implementing the schedules (*Enactor*).
1. The *Collection* interacts with resource objects to collect state information describing the system.
 2. The *Scheduler* queries the *Collection* to determine a set of available resources that match the Scheduler's requirements.
 3. After computing a schedule, or set of desired schedules, the *Scheduler* passes a list of schedules to the *Enactor* for implementation.
 4. The *Enactor* then makes reservations with the individual resources and reports the results to the Scheduler.
 5. Upon approval by the *Scheduler*, the *Enactor* place objects on the hosts, and monitors their status.

Host objects can be adapted to different environments to suit user needs. For example, a host object may provide an interface to the underlying resource management system, such as LSF, Codine, or Condor.

Programming Tools and Paradigms

Legion supports MPL (Mentat Programming Language) and BFS (Basic Fortran Support). MPL is a parallel C++ language. Legion is written in MPL. BFS is a set of pseudo-comments for Fortran and a pre-processor that gives the Fortran programmer access to Legion objects.

Object Wrapping is used in Legion for encapsulating existing legacy codes into objects. It is possible to encapsulate a PVM, HPF, or shared memory threaded application in a Legion object. Legion also provides a complete emulation of both PVM and MPI with user libraries for C, C++ and Fortran. Legion also supports Java.

User and Administrative GUI

Legion has a command-line and graphical user interface. The Legion GUI, known as the Context Manager, is a Java application that runs context-related commands. The Context Manager uses icons to represent different parts of context space (file objects, sub-contexts, etc.) and runs most context-related commands. The Context Manager can be run from the command-line of any platform compatible with the Java Development Kit (JDK) 1.1.3. In addition there is a Windows 95 client application, called the Legion Server, that allows users to run the Context Manager from Windows 95.

Availability

Legion is available on: x86/Alpha (Linux), Solaris (SPARC), AIX (RS/6000), IRIX (SGI), DEC UNIX (Alpha) and Cray T90.

1.4.4 WebFlow

Introduction

WebFlow [13] [14], is a computational extension of the Web model that can act as a framework for the wide-area distributed computing and metacomputing. The main goal of the WebFlow design was to build a seamless framework for publishing and reusing computational modules on the Web so that end-users, via a Web browsers, can engage in composing distributed applications using WebFlow modules as visual components and editors as visual authoring tools. Webflow has a three-tier Java-based architecture that can be considered a visual dataflow system. The front-end uses applets for authoring, visualisation and control of the environment. WebFlow uses servlet-based middleware layer to manage and interact with backend modules such as legacy codes for databases or high performance simulations.

Webflow is analogous to the Web. Web pages can be compared to WebFlow modules and hyperlinks that connect Web pages to inter-modular dataflow channels. WebFlow content developers build and publish modules by attaching them to Web servers. Application integrators use visual tools to link outputs of the source modules with inputs of the destination modules, thereby forming distributed computational graphs (or compute-webs) and publishing them as composite WebFlow modules. A user activates these compute-webs by clicking suitable hyperlinks, or customise the computation either in terms of available parameters or by employing some high-level commodity tools for visual graph authoring.

The high performance backend tier is implemented using Globus toolkit:

- The Metacomputing Directory Services (MDS) to used map and identify resources,
- The Globus Resource Allocation Manager (GRAM) is used to allocate resources.
- The Global Access to Secondary Storage (GASS) is used for a high performance data transfer.

WebFlow can be regarded as a high level, visual user interface and job broker for Globus.

With WebFlow new applications can be composed dynamically from reusable components just by clicking on visual module icons, dragging them into the active WebFlow editor area, and linking by drawing the required connection lines. The modules are executed using Globus components combined with the pervasive commodity services where native high performance versions are not available.

The prototype WebFlow system is based on a mesh of Java enhanced Web Servers (Apache), running servlets that manage and co-ordinate distributed computation. This management infrastructure is implemented by three servlets: *Session Manager*, *Module Manager*, and *Connection Manager*. These servlets use URL addresses and can offer dynamic information about their services and current state. Each management servlet can communicate with others via sockets. The servlets are persistent and application independent.

Future implementations of WebFlow will use emerging standards for distributed objects, and take advantage of commercial technologies, such as the CORBA as the base distributed object model.

Administrative Hierarchy

WebFlow has no obvious administrative hierarchy. A WebFlow node is a Web server with unique URL address, and it is a peer to other nodes.

Communication Services

WebFlow communication services are built on multiple protocols. Applet-Web Server communication uses HTTP, Server-to-Server communication is currently implemented using TCP/IP, soon to be replaced by IIOP. The module developer chooses communications between a backend module and its front-end control panel (Java applet): typically it is either TCP/IP or IIOP. The modules exchange data (serialised Java objects) via input and output ports. Originally, the port-to-port connection was implemented using TCP/IP. This model is being changed now. The module is a Java Bean, and it interacts with other modules via Java events over IIOP. The data flow paradigm with port-to-port communication is the default that enables users to visually compose an application from independent modules. However, the user is not restricted to this model. A module can be a high performance application to be run on a multiprocessor machine with intra-module communications using any communication service (for example MPI) available on the target system. Also, the modules can interact with each other via remote methods invocation (Java events over IIOP).

Also, WebFlow supports multiple protocols for file transfer, ranging from HTTP to FTP to IIOP to Globus GASS. The user chooses the protocol to be used depending on the file, performance and security requirements.

Directory/Registration Services

WebFlow does not define its own directory services. The usage of the CORBA naming services and interface repository is planned. It should be noted that WebFlow typically is used in conjunction with Globus and will its directory services (MDS).

Processes, Threads and Concurrency Control

Each module runs as separate Java threads, and all modules are running concurrently. It is the user's responsibility to synchronise modules (for example, the user may want the module to block on receiving the input data).

Time and Clocks

WebFlow does not mandate the usage of a particular time service and relies on those already used at each site.

Naming Services

Currently, no specialised naming service other than DNS is used. CORBA services are planned.

Distributed Filesystems and Caching

WebFlow does not offer 'a native' distributed file system or support caching. This is left to the user. As a part of WebFlow distribution there is a file browser module, that allows the user to browse and select files accessible by the host Web server. The selected files can then be sent to a desired destination using HTTP, IIOP, FTP or GASS. For example, a WebFlow module that serves as the Globus GRAM proxy takes the name of the input file and URL of the GRAM contact as input. This information is sufficient to stage the input file on the target machine and retrieve the output file using GASS over FTP.

Security and Authorisation

WebFlow requires two security levels: secure Web transactions between client and the middle tier, and secure access to the backend resources. Secure Web transactions in WebFlow are based on TLS 1.0 and modelled after the AKENTI system, and secure access to the backend resources is delegated to the backend service providers, such as Globus. The secure access to resources directly controlled by WebFlow has not yet been addressed.

System Status and Fault Tolerance

The original implementation of WebFlow did not address these issues. The new WebFlow middle-tier will use CORBA mechanisms to provide fault tolerance, including a heartbeat monitor. In the backend, WebFlow relies on services provided by the backend service provider.

Resource Management and Scheduling

WebFlow delegates the resource management and scheduling to the metacomputing toolkit (Globus) and/or a local resource management package such as PBS or CONDOR.

Design Objective	Globus	Legion	Webflow
Admin. Hierarchy	Peer	Peer	Peer
Comms Service	Nexus - Low-Level	MMPS - Sockets-based	Hierarical - Sockets+MPI
Dir/Reg Services	MDS - LDAP	Via Binding agent	MDS - LDAP
Processes	Process-based	Object/process-based	Process-based
Clock	Not specified	Not specified	Not specified
Naming Services	LDAP + DNS/X.500	Context Manger + DNS	LDAP + DNS
Filesystems & caching	GASS + ROMIO	Custom Legion filesystem	GASS
Security	GSI (RSA + X.509 certs)	Object-based with RSA	SSL
Fault Tolerance	Heart-beat monitor	Not available yet	None
Resource Management	GRAM + RSL + Local	Host object + Local	GRAM-based
Prog. Paradigms	Many and varied	MPL, BFS + wrappers	MPI
User Interfaces	GUI + command-line	GUI + command-line	Applet-based GUI
Availability	Most UNIX	Most UNIX	Most UNIX and NT

Table 1.1. Metacomputing Functionality Matrix

Programming Tools and Paradigms

WebFlow modules are Java objects. Object wrapping is used in WebFlow for encapsulating existing codes into objects. WebFlow test applications include modules with encapsulated Fortran, Fortran with MPI, HPF, C with MPI, Pascal as well as Java.

User and Administrative GUI

WebFlow offers a visual-authoring tool, implemented as a Java applet that allows the user to compose a (meta-) application from pre-existing modules. In addition, a developer can use a simple API to build a custom graphical user interface.

Availability

Since WebFlow is implemented in Java, it runs on all platforms that support JVM. So far it has been tested on Solaris, IRIX, and Windows NT.

Summary and Conclusions

Introduction

In this section we have attempted to lay out the functionality and features of three representative metacomputer architectures with the design criteria we outlined in Section 1.3. This task in itself has been rather difficult as it has been necessary to map the developer's terminology for components within their environments to those used more commonly in distributed computing. In the final part of this section we summarise the functionality of each environment and conclude by making some observations about the approaches each environment uses.

Functionality Matrix

In the functionality matrix, shown in Table 1, we outline the components within each metacomputing environment that deals with our design criteria.

From table 1 it can be seen that:

- *Administrative Hierarchy* - All three environments use a peer-based administrative hierarchy, which makes the services they provide globally scalable and reduces potential administrative bottlenecks and single-points of failure.
- *Communications Service* - Globus uses Nexus to provide its underlying communications services, whereas Legion and WebFlow use sockets-based protocol.
- *Directory/Registration Services* - Both Globus and Webflow use the commodity LDAP service, whereas Legion uses a custom binding agent.
- *Processes* - all three environments are process based - but each has the ability to encompass threads and enable consistency control.
- *Clock* - The three environments do not require special timing services.
- *Naming services* - Globus and WebFlow use LDAP in conjunction with DNS, where as Legion uses a custom context manager in conjunction with DNS.
- *Filesystems and caching* - Globus makes extensive use of the remote access tool GASS and the parallel I/O interface ROMIO. Legion has a custom global filesystem and the ability to interface with other I/O system. WebFlow utilises GASS to provide file system services.
- *Security* - All three environments use RSA in some form. Globus uses GSI to provide its security services. Legion uses an Object based system where every object has a security method `MayI`. WebFlow uses SSL for security purposes.
- *Fault Tolerance* - Of the three environments, only Globus provide tool - the heartbeat monitor.
- *Resource Management* - Globus implements an extensive resource management and scheduling system, GRAM. Legion has the concept of host object for local resource management. Both Globus and Legion have interfaces to other resource management system, such as Codine and LSF. WebFlow utilises the services of GRAM.
- *Programming Paradigms* - All three environments provide a raft of tools and utilities to support various programming paradigms.
- *User Interfaces* - Globus and Legion provide both command-line and GUI interfaces. WebFlow uses just a GUI.
- *Availability* - All three environments are available on most UNIX platforms.

Some Observations

Globus is constructed as a layered architecture in which high-level global services are built upon essential low-level core local services. The Globus toolkit is modular. This means that an application can exploit an array of features without needing to implement all of them. Globus can be viewed as a metacomputing framework based on a set of APIs to the underlying services. Even though Globus provides the services needed to build a metacomputer, the Globus framework allows alternative local services to be used if desired. For example, the GRAM API allows alternative resource management systems to be utilised, such as Condor or NQE.

Abstracting the services into a set of standard APIs has a number of advantages. These include:

- The underlying services can be changed without affecting applications that use them,
- This type of layered approach simplifies the design of a rather complicated system,
- It encourages developers of tools and services, they need to support only one API making their development and testing cycle shorter and cheaper.

Globus provides application developers with a pragmatic means of implementing a range of services to provide a wide-area application execution environment.

Legion takes a very different approach to provide a metacomputing environment, it encapsulates all its components as objects. The methodology used has all the normal advantages of an object-oriented approach, such as, data abstraction, encapsulation, inheritance and polymorphism.

It can be argued that many aspects of this object-oriented approach potentially makes it ideal for designing and implementing a complex environment such as a metacomputer. For example, Legion's security mechanism, where each object uses RSA keys and a `MayI` method, seems straightforward and more natural than security mechanisms used in many other environments. In addition, the set of methods associated with each object naturally becomes its external interface and hence its API.

Using an object-oriented methodology in Legion does not come without a raft of problems. It is not obvious how best to encapsulate non object-oriented programming paradigms, such as message passing or distributed shared memory. In addition, the majority of real-world computing services have procedural interfaces and it is necessary to produce object-oriented wrappers to interface these services to Legion. For example, the APIs to DNS or resource management systems such as Condor or Codine are procedural.

WebFlow takes a different approach to both Globus and Legion. It is implemented in a hybrid manner using a three-tier architecture that encompasses both the Web and third party backend services. This approach has a number of advantages, including the ability to 'plug-in' a diverse set of backend services. For

example, currently many of these services are supplied by the Globus Metacomputing Toolkit, but they could be replaced with components from CORBA or Legion. WebFlow also has the advantage that it is more portable and can be installed anywhere a Web server supporting servlets is capable of running.

1.5 Emerging Metacomputing Environments

1.5.1 Introduction

There are a large number and diverse range of emerging distributed systems currently being developed. These systems range from metacomputing frameworks to application testbeds, and from collaborative environments to batch submission mechanisms. In this section we briefly describe and reference a few of the better known systems. The aim is to bring the reader's attention not only some of the large number of diverse projects that exist, but also to detail the different approaches used to solve the inherent problems encountered.

1.5.2 Metacomputing Environments

Arcade (ICASE, NASA Langley Research Center)

Arcade [15] [16] is a Java-based framework that uses the Web to provide support for a team to collaboratively design, execute and monitor multi-disciplinary applications on a distributed heterogeneous system. This framework is suitable for applications such as the multi-disciplinary design optimisation of an aircraft. Arcade applications can consist of multiple heterogeneous modules executing on independent distributed resources while interacting with each other to solve the overall design problem.

BAYANIHAN (MIT Laboratory for Computer Science)

BAYANIHAN [17] [18] is a software framework that uses Java and HORB [19] - a distributed object library. The framework allows users to co-operate in solving computational problems by using their Web browser to volunteer their computers' processing power. HORB uses something akin to RMI to pass data between objects. BAYANIHAN also provides a PVM interface for parallel computing.

Bond (Computer Sciences Department, Purdue University)

Bond [20] [21] is a metacomponent architecture for network computing on a grid of autonomous nodes. The Bond environment is composed of metaobjects, which consist of information and network object pairs. Reflection mechanisms allow objects to discover each other's properties and core agents and services provide scheduling, user-level resource management, security, fault-tolerance and other functions.

COVISE (Computing Center, University of Stuttgart)

COVISE [22] [23] (Collaborative Visualization and Simulation Environment) is an extendable distributed software environment used to integrate simulations, post-processing and visualisation functionalities in a seamless manner. An application in COVISE is divided into a number of processing steps, which are represented by modules that can be implemented as separate processes and arbitrarily spread across different heterogeneous platforms.

Charlotte (Department of Computer Science, New York University)

Charlotte [24] [25] is a distributed shared memory environment for parallel programming over the Web. Charlotte is built on top of Java and does not rely on any native code. The runtime system of Charlotte uses eager scheduling and two-phase idempotent execution strategy to provide system load balancing and fault tolerance. In a Charlotte program, the data is logically partitioned into private and shared segments, the former data is local and private and the later is shared and distributed.

Distributed Batch Controller (Computer Science, University of Wisconsin)

DBC [26] [27] is a system that processes data using widely distributed computational resources controlled by the Condor Resource Management (CRM) system. CRM scavenges idle CPU cycles from UNIX workstations.

DISCworld (Department of Computer Science, University of Adelaide)

DISCWorld [28] [29] forms a serverless middleware layer that provides a software infrastructure to support legacy applications by encapsulating them as explicitly named services with well-defined interfaces and actions. Key technologies for making this architecture possible are Java and the use of Remote Method Invocation (RMI) mechanisms.

DOCT (San Diego Supercomputer Center)

The Distributed Object Computation Testbed [30] (DOCT) is a broad research and development effort to unify access methodologies across multiple sites. The DOCT environment federates heterogeneous resources (data, high-performance computers, and networks) administered by different authorities to achieve common objectives, such as in support of a government agency mission. DOCT uses an object-oriented approach, emphasising retention and tracking of all data sets in the environment to manage complex documents comprised of text, images, and multimedia files.

EROPPA (Genias Benelux, The Netherlands)

The objective of EROPPA [31] (Environment for Remote Operations on Post Production Application) was to design, implement and test methodologies for remote

and/or distributed access to 3D graphics applications that run on high performance facilities. The project was driven by real market requirements and the needs of companies in the post production area.

HARNESS (Oak Ridge National Laboratory)

The Heterogeneous Adaptable Reconfigurable Networked SystemS [32] [33] (HARNESS) is an experimental metacomputing framework built around the services of a customisable and reconfigurable Distributed Virtual Machine (DVM). HARNESS defines a flexible kernel and views a DVM as a set of components connected by the use of a shared registry, which can be implemented in a distributed fault tolerant manner. Any particular kernel component thus derives its identity from this distributed registry. The flexibility of service components comes from the way the kernel supplies DVM services by allowing components, which implement services, to be created and installed dynamically. HARNESS uses the micro-kernel approach, where services are added as needed to provide the functionality that the users require.

Hector (Mississippi State University's NSF Engineering Research Center)

Hector [34] [35] provides networks of workstations with support for a widely accepted parallel programming environment (MPI), task migration, resource allocation, fault tolerance.

IceT (Emory University Dept. of Mathematics and Computer Science)

IceT [36] [37] is a framework for collaborative and high performance distributed computing built on top of the Java programming substrate. IceT allows dynamic merging and splitting of virtual environments, multi-user awareness, portability of processes and data across virtual machines and a provisional framework for multi-user programs. This environment is similar to PVM, but using the JVM. IceT supports a message-passing interface much like PVM.

JavaNow (Fermi National Accelerator Lab. and Argonne National Lab.)

The Java Network of Workstations [38] framework (JavaNOW) provides a mechanism where large, complex computational tasks can be broken down and executed in a parallel fashion across a network. The concept of JavaNOW is similar to PVM, however, the JavaNOW interface is based on logical distributed associative shared memory instead of message passing. The interface for JavaNOW is similar to the shared memory model of Linda. It allows users to store complete objects in shared memory.

JavaParty (University of Karlsruhe)

JavaParty [39] [40] is a system for distributed parallel programming in heterogeneous workstation clusters. JavaParty extends Java by adding the concept of remote

objects. The JavaParty consists of a pre-processor and a runtime-system to support the system.

JAVELIN (Dep. of Computer Science University of California, Santa Barbara)

The central idea in the Javelin [41] [42] architecture is a computing broker, which collects and monitors resources, and services requests for these resources. A user interacts with the broker in two ways, by registering or requesting resources. The broker matches clients with host machines for running computations. A key design feature of this architecture is that the host software runs completely at the user-level, and only requires that registered machines have a Web browser, which can run untrusted code, such as Java applets.

Jini (Sun Microsystems)

Jini [43] is an attempt by Sun to create a new distributed computing architecture. Jini is an object-oriented framework that embodies a model of how devices and software inter-operates as well as how distributed systems function. The infrastructure consists of two main components: 'Discovery and Join' and 'Lookup'. Discovery and Join is a mechanism whereby a device or application identifies itself to the network. The mechanism is two-phase. First, the entity broadcasts a discovery package, this contains sufficient information to enable the network to start a dialog with the entity that has just joined. Second, once acknowledged, the entity can now join by sending a message containing details about its own characteristics. Lookup is a component that stores information about Jini registered devices and applications. Clients using Jini use Lookup to find the services that they wish to access.

The distributed programming model used by Jini promotes three technologies: Leasing, Distributed Transactions and Distributed Events. Leasing is where an object negotiates the usage of a service for a period of time. Communication within Jini is based on Remote Method Invocation (RMI). The Jini distributed programming model is based on the JavaSpaces model, which is in itself based on Linda.

KnittingFactory (Department of Computer Science, New York University)

KnittingFactory [44] is an infrastructure that supports building Web-based parallel and collaborative applications. KnittingFactory is based around communicating Java applets that interact by using a special Directory service and Class server. KnittingFactory extends Charlotte to provide its programming and collaborative environment.

Metacomputer OnLine (Paderborn Center for Parallel Computing)

The Metacomputer Online [45] [46] (MOL) project aims to integrate existing software modules in an open extensible environment. The MOL architecture consists of three general modules, these support programming environments (PVM, MPI, and PARIX), resource management and access systems (Codine, NQS, PBS and CCS),

and supporting tools (GUIs, plus tools such as the task migrator MARS, WARP performance predictor and WAMM).

NEOS (Argonne National Lab.)

The Network-Enabled Optimisation System [47] [48] (NEOS) is an Internet-based service environment for the solution of optimisation problems. The main components of NEOS are the NEOS Guide and Server. A user of NEOS only needs to describe the optimisation problem; all additional information required by the optimisation solver, is determined automatically.

Netsolve (University of Tennessee at Knoxville)

NetSolve [49] [50] is a client/server application designed to solve computational science problems in a distributed environment. The Netsolve system is based around loosely coupled distributed systems, connected via a LAN or WAN. Netsolve clients can be written in C and Fortran, use Matlab or the Web to interact with the server. A Netsolve server can use any scientific package to provide its computational software. Communications within Netsolve is via Sockets. Good performance is ensured by a load-balancing policy that enables NetSolve to use the computational resources available as efficiently as possible. NetSolve offers the ability to search for computational resources on a network, choose the best one available, solve a problem (with retry for fault-tolerance), and return the answer to the user.

The Nile Project (University of Texas)

The Nile Project [51] is developing a distributed computing solution for the CLEO High Energy Physics experiment, based at Cornell's electron storage ring accelerator facility. The goal is to provide a self-managing, fault-tolerant, heterogeneous system with access to a distributed database in excess of 100 TBytes. These resources are spread across the United States and Canada at 24 collaborating institutions. The Nile metacomputing software is based on CORBA.

Ninf (Japanese Consortium)

The Network Infrastructure [52] [53] for global computing (Ninf) is a client-server-based system that allows access to multiple remote compute and database servers. Ninf clients can semi-transparently access remote computational resources from languages such as C and Fortran. A programmer is able to build a global computing application by using the Ninf remote libraries as its components, without being aware of the complexities of the underlying system they are programming.

Ninja (Dept. of Computer Science, Berkeley)

The Ninja [54] project aims to develop a software infrastructure to support the next generation of Internet-based applications. Ninja has the concept of a service, an Internet-accessible application that is scalable, fault-tolerant and highly available.

Ninja will enable the development of a suite of interoperable and immediately accessible Internet-based services across a spectrum of devices ranging from PCs and workstations to Cellphones and Personal Digital Assistants.

NWS (UCSD)

The Network Weather Service [55] [56] (NWS) is a system that takes periodic climatic measurements from distributed networked resources, and uses numerical models to dynamically generate forecasts of future meteorological conditions.

PARDIS (Department of Computer Science, Indiana University)

PARDIS [57] [58] is based on CORBA in that it allows the programmer to construct meta-applications without concern for component location, heterogeneity of resources, or data translation and marshalling in communication between them. PARDIS supports SPMD objects representing data-parallel computations. These objects are implemented as a collaboration of computing threads capable of directly interacting with the PARDIS ORB - the entity responsible for brokering requests between clients and servers. The PARDIS ORB interacts with parallel applications through a run-time system interface implemented by the underlying the application software package. PARDIS can be used to interface directly parallel packages, based on different run-time system approaches, such as the POOMA library and High Performance C++.

PROMENVIR (Parallel Application Centre, Southampton)

PROMENVIR [59] (PRObabilistic MEchanical desigN enVIRonment) is a meta-computing tool for performing stochastic analysis of generic physical systems. The tool provides a user with the framework for running a stochastic Monte Carlo (MC) simulation, using a preferred deterministic solver (e.g. NASTRAN), then provides statistical analysis tools to analyse the results. A fundamental component of the package is the Advanced Parallel Scheduler (APS), that acts as a meta-application manager and orchestrates resource usage by an application. The APS creates daemons on remote hosts, which are responsible for submitting jobs, copying files and communicating load information back to the master workstation. It is capable of initiating UNIX processes directly, but can also be used to control and submit meta-applications via conventional load-sharing software, such as LSF, Codine, or Condor..

SNIFE (University of Tennessee at Knoxville)

SNIFE [60] [61] (Scalable Networked Information Processing Environment) is a metacomputing system that aims to provide a reliable, secure, fault tolerant environment for distributed applications and data stores across the Internet. The system combines global naming and replication of both processing and data to support large-scale information processing applications. SNIFE contains seven major com-

ponents: Metadata servers, file servers, per-host SNIPE daemons, client libraries, resource managers, 'playgrounds', and consoles.

Symera (NCSA)

NCSA Symera [62] [63] (Symbiotic Extensible Resource Architecture) is a distributed-object system based on Microsoft's Distributed Component Object Model (DCOM). Symera is designed to support both sequential and parallel applications. The Symera management system consists of an NT service that is installed on each platform in a cluster. The management system hosts objects that allocate resources, schedules jobs, implements fault tolerance as well as object activation and migration. Symera is written in C++ that conforms to the Win32 standard.

WAMM (CNUCE, Italy)

WAMM [64] [65] (Wide Area Metacomputer Manager) is a graphical tool, built on top of PVM. It provides users with a GUI to assist in tasks such as: host add, check, removal, process management, compilation on remote hosts, remote commands. All functions are accessible via menus and buttons with a geographical view of the system. The hosts are grouped following a tree structure. The root node, corresponding to a WAN, can contain MAN and LAN. The system allows the remote execution of UNIX commands as well as application compilation. WAMM has been set up and tested between a number of Italian research labs.

1.5.3 Metacomputing Interfaces

UNICORE (German Consortium and ECWMF)

The Uniform Interface to Computing Resources [66] (UNICORE) project aims to deliver software that allows users to submit jobs to remote high performance computing resources without having to learn details of the target operating system, data storage conventions and techniques, or administrative policies and procedures at the target site. The user interface is based on Java and Web browser technology. A Network Job Supervisor (NJS) at each UNICORE site interprets an Abstract Job Object, which is generated by the user interface, manages submitted jobs and any associated job data. NJS can inter-operate with resource management systems, such as, Cray NQE, IBM Load Leveler and Codine.

Websubmit (NIST)

WebSubmit [67] is a Web-based framework that aims to provide seamless access to applications on a collection of heterogeneous computing resources. WebSubmit has a strong emphasis on security and uses the Secure Sockets Layer protocol for user authentication. A user, when validated by a WebSubmit authority, is given access to a group of application modules. Each application module is presented as an HTML form; this form is filled out and submitted to the server, which then processes the request and executes the desired tasks on the specified remote system.

1.5.4 Summary

The projects described in this section are a cross-section of those currently undertaken. It is interesting to note that all are using Java and the Web as the communications infrastructure. It is also evident that Java has revolutionised the shape and characteristic of the software environments for heterogeneous distributed systems. It seems that the developers of distributed systems no longer have to focus on aspects such as portability and heterogeneity, by using Java they seem able to concentrate on designing and implementing functional distributed environments. It is not clear, among the raft of projects listed in this section, which environments will succeed. However, each project, in its own way, is contributing to our knowledge of how to design, build and implement efficient and effective distributed virtual environments.

1.6 Summary and Conclusions

1.6.1 Introduction

In this chapter we have attempted to describe and discuss many aspects of metacomputers. We started off by discussing why there is a need for such environments. We then moved on to describe two early metacomputing projects. Here we also outlined some of the benefits and experiences learned. Having set the scene, we then laid out a design template to map out the critical services that a metacomputing environment needs to encompass. Then, using this template, we mapped the services of three differing environments onto it. This mapping made comparing and contrasting the services that each metacomputing environment provided clearer to understand. Having described three fairly mature environments, we then briefly described some thirty-odd emerging distributed environments and tools. Finally, here, we summarise what we have discovered whilst researching this chapter and conclude by making a few predictions about metacomputing environments of the future.

1.6.2 Summary of the Reviewed Metacomputing Environments

Globus is constructed as a layered architecture in which high-level global services are built upon essential low-level core local services. The Globus toolkit is modular and as such an application can exploit an array of features without needing to implement all of them. Globus can be viewed as a metacomputing framework based on a set of APIs to the underlying services. Globus provides application developers with a pragmatic means of implementing a range of services to provide a wide-area application execution environment.

Legion takes a very different approach to provide a metacomputing environment, it encapsulates all its components as objects. The methodology used has all the normal advantages of an object-oriented approach, such as, data abstraction, encapsulation, inheritance and polymorphism. It can be argued that many aspects

of this object-oriented approach potentially makes it ideal for designing and implementing a complex environment such as a metacomputer. However, using an object-oriented methodology does not come without a raft of problems, many of these are tied-up with the need for Legion to interact with legacy applications and services. In addition, as Legion is written in MPL, it is necessary to 'port' MPL onto each platform before Legion can be installed.

WebFlow takes a different approach to both Globus and Legion. It is implemented in a hybrid manner using a three-tier architecture that encompasses both the Web and third party backend services. This approach has a number of advantages, including the ability to 'plug-in' a diverse set of backend services. For example, currently many of these services are supplied by the Globus toolkit, but they could be replaced with components from CORBA or Legion. WebFlow also has the advantage that it more portable and can be installed anywhere a Web server supporting servlets is capable of running.

So, in summary, we believe that all three environments have their merits. Fundamentally, the Globus Metacomputing Toolkit is currently the most comprehensive attempt to provide a metacomputing environment. The Globus team have taken a very pragmatic approach to providing the services that are needed in a metacomputer. The design methodology they have used - abstracting the services of some underlying entity into a well thought out API - will give the project longevity, as the entities that provide the service can be updated without changing the fundamental service API. In addition Globus uses existing standard commodity software components to provide many of its services, for example LDAP, X.509 and RSA. This has a number of beneficial implications, including code reuse and avoiding the necessity to create all the services from scratch.

Alternatively, Legion is a very ambitious and impressive project. We believe the object-oriented approach they have taken has a lot of merit. A fundamental flaw with Legion currently is the reliance on MPL. If Legion were written in Java, which no doubt the University of Virginia is seriously contemplating, then we would have much more faith in Legion's longevity. Also, perhaps, we would question the use of this system as opposed to one based on the well-known standard CORBA.

WebFlow is still basically an experimental prototype system that is being used to explore a range of new and emerging technologies. It has much merit, particularly in its comprehensive GUI frontend and its ability to utilise standard backend components designed by other organisations.

1.6.3 Some Observations

The Java programming language successfully addresses several key issues that plague the development of distributed environments, such as heterogeneity and security. It also removes the need to install programs remotely, the minimum execution environment is a Java-enabled Web browser. Java has become a prime candidate for building distributed environments.

In a metacomputing environment it is not possible to mandate the types of

services or programming paradigms that particular users or organisations must use. A metacomputer needs to provide extensible interfaces to any service desired.

Providing adequate security in a metacomputer is a complex issue. A careful balance needs to be maintained between the usability of an environment and security mechanisms utilised. The security methods must not inhibit the usage of an environment but it must ensure that the resources are secure from malicious intruders.

1.6.4 Metacomputing Trends

It is very difficult to predict the future. In a field such as computing, the technological advances are moving very fast. Windows of opportunity for ideas and products seem to open and close in the seeming 'blink of the eye'. However, some trends are evident.

Java, its related technologies and growing repository of tools and utilities, is having a huge impact on the growth and development of metacomputing environments. From a relatively slow start, the development of metacomputers is accelerating fast with the advent of these new and emerging technologies. It is very hard to ignore the presence of the sleepy giant CORBA in the background. We believe that frameworks incorporating CORBA services will be very influential on the design of metacomputing environments in the future.

Whatever technology or computing paradigm becomes influential or most popular, it can be guaranteed that at some stage in the future its star will wane. Historically, in the computing field, this fact can be repeatedly observed. The lesson from this observation must therefore be drawn that, in the long term, backing only one technology can be an expensive mistake. The framework that provides a metacomputing environment must be adaptable, malleable and extensible. As technology and fashions change it is crucial that a metacomputing environment evolves with them.

1.6.5 The Impact of Metacomputing

Metacomputing is not only a computing paradigm for just providing computational resources for supercomputing-sized parallel applications. It is an infrastructure that can bond and unify globally remote and diverse resources ranging from meteorological sensors to data-vaults, and from parallel supercomputers to personal digital organisers. As such, it will provide pervasive services to all users that need them.

Larry Smarr in *The GRID: Blueprint for a New Computing Infrastructure* [3] observes that metacomputing has serious social consequences and is going to have as revolutionary an effect as railroads did in the American mid-West in the early 19th century. But that instead of a 30 - 40 year lead-time to see its effects, its impact is going to be much faster. He concludes, that the effects of computational grids are going to change the world so quickly that mankind will struggle to react and change in the face of the challenges and issues they present.

So, at some stage in the future, our computing-needs will be satisfied in the

same pervasive and ubiquitous manner that we use the electricity power-grid. The analogies with the generation and delivery of electricity are hard to ignore and the implications are enormous.

Acknowledgements

The authors wish to thank Ian Foster (ANL) and Tom Haupt (Syracuse) for information and useful suggestions about their projects. We would also like to thank Wolfgang Gentzsch (Genias) for early access to a FGCS Special Issue on *Metacomputing* [68]. The authors would like to thank Kate Dingley, Tony Kalus, John Rosbottom and Rose Rayner for proof reading the copy of this chapter.

1.7 Bibliography

- [1] C. Catlett and L. Smarr, *Metacomputing*, Communications of the ACM, 35(6):44-52, 1992.
- [2] Desktop Access to Remote Resources - <http://www-fp.mcs.anl.gov/~gregor/datorr/>
- [3] *The GRID: Blueprint for a New Computing Infrastructure*, Edited by I. Foster and C. Kesselman, Morgan Kaufmann Publishers, Inc, 1998, ISBN 1-55860-475-8
- [4] FAFNER - <http://www.npac.syr.edu/factoring.html>
- [5] I-WAY - <http://146.137.96.14/>
- [6] RSA - <http://www.rsa.com/>
- [7] I. Foster, J. Geisler, W. Nickless, W. Smith, S. Tuecke, *Software Infrastructure for the I-WAY Metacomputing Experiment*, Concurrency: Practice and Experience, 10(7):567-581, 1998.
- [8] Globus - <http://www.globus.org/>
- [9] I. Foster and C. Kesselman, *The Globus Project: A Status Report*. Proc. IPPS/SPDP '98 Heterogeneous Computing Workshop, pg. 4-18, 1998.
- [10] W. Yeong, T. Howes and S. Kille, *Lightweight Directory Access Protocol*, RFC 1777, 28/03/95, Draft Standard.
- [11] Legion - <http://legion.virginia.edu/>
- [12] A. Grimshaw, W. Wulf, et al. *The Legion Vision of a Worldwide Virtual Computer*, Communications of the ACM January 1997, Vol. 40, No 1
- [13] WebFlow - <http://osprey7.npac.syr.edu:1998/iwt98/products/webflow/>

- [14] T. Haupt, E. Akarsu, G. Fox and W. Furmanski, *Web based metacomputing*, Special Issue on Metacomputing to appear "Future Generation Computer Systems", North Holland, Due for publication in early 1999.
- [15] ARCADE - <http://www.cs.odu.edu/~ppvm/>
- [16] K. Maly, P. Vangala, and M. Zubair, *JAVADC: A Web-Java Based Environment to Run and Monitor Parallel Distributed Applications*, Technical Report, Old Dominion University, 1997.
- [17] BAYANIHAN - <http://www.cag.lcs.mit.edu/bayanihan/>
- [18] L. Sarmenta, S. Hirano, and S. Ward, *Towards Bayanihan: Building an Extensible Framework for Volunteer Computing Using Java*, ACM 1998 Workshop on Java for High-Performance Network Computing. (submitted)
- [19] S. Hirano, *HORB: Extended execution of Java Programs*, Proceedings of the 1st International Conference on World-Wide Computing and its Applications (WWCA97), March 1997, <http://ring.etl.go.jp/openlab/horb/>
- [20] BOND - <http://bond.cs.purdue.edu/>
- [21] L. Bölön, *Bond Objects – A White Paper*, Department of Computer Sciences, Purdue University CSD-TR No. 1998-002.
- [22] COVISE - <http://www-ks.rus.uni-stuttgart.de/PROJ/suntrec/covise.html>
- [23] A. Wierse, U. Lang, R. Rühle, *A System Architecture for Data-oriented Visualization* Database Workshop, Visualization '93, San Jose to appear in "Proceedings of the IEEE Workshop on Database for Data Visualization", Lecture Notes in Computer Science, Volume 871, Lee, Grinstein (Eds.), Springer Verlag.
- [24] Charlotte - <http://www.cs.nyu.edu/milan/publications/pdcs96/>
- [25] A. Baratloo, M. Karaul, Z. Kedem, P. Wyckoff, *Charlotte: Metacomputing on the Web*, Proc. of the 9th International Conference on Parallel and Distributed Computing Systems, September 1996.
- [26] CONDOR - <http://www.cs.wisc.edu/condor/>
- [27] M. Livny and R. Raman, *High-Throughput Resource Management, The GRID: Blueprint for a New Computing Infrastructure*, Editors, I. Foster and C. Kesselman, Morgan Kaufmann, 1998, ISBN 1-55860-475-8
- [28] DISCWorld - <http://www.dhpc.adelaide.edu.au/>

- [29] K. Hawick and F. Vaughan, *Discworld - distributed information systems cloud of high performance computing resources* - design discussion document, Computer Science Department, University of Adelaide, DHPC Working Note, March 1997.
- [30] DOCT - <http://www.sdsc.edu/DOCT/>
- [31] EROPPA - <http://www.telecom.ntua.gr/eroppa/>
- [32] HARNESS - <http://www.epm.ornl.gov/harness>
- [33] G. Geist, *Harness: The Next Generation Beyond PVM*, proceedings of the 5th European PVM/MPI Users' Group Meeting, Lecture Notes in Computing Science, Vol 1479, pp. 74-82, ISBN 3-540-65041-5
- [34] Hector - <http://www.erc.msstate.edu/>
- [35] S. Russ, et. al *The Hector Parallel Run-Time Environment*, IEEE Transactions on Parallel and Distributed Systems Vol. 9(11), Nov. 1998, pp. 1104-1112
- [36] IceT - <http://www.mathcs.emory.edu/~gray/>
- [37] P. Gray and V. Sunderam, *IceT: Distributed Computing and Java*, Concurrency, Practice and Experience, ed. Geoffrey C. Fox, Vol. 9 (11), pp 1161-1168, Nov. 1997.
- [38] JavaNOW - <http://www.jhpc.org/javanow.htm>
- [39] JavaParty - <http://www.ipd.ira.uka.de/JavaParty/>
- [40] M. Philippsen and M. Zenger, *JavaParty - Transparent Remote Objects in Java*, Concurrency: Practice and Experience, Vol. 9 (11), pp. 1225-1242, November 1997.
- [41] Javelin - <http://www.cs.ucsb.edu/research/javelin/>
- [42] B. Christiansen, P. Cappello, M. Ionescu, M. Neary, K. Schauser, and D. Wu, *Javelin: Internet-Based Parallel Computing Using Java*, ACM Workshop on Java for Science and Engineering Computation, 1997.
- [43] JINI - <http://www.java.sun.com/products/jini/>
- [44] KnittingFactory - <http://www.cs.nyu.edu/>
- [45] Metacompter OnLine - <http://www.uni-paderborn.de/pc2/projects/mol/>
- [46] F. Ramme, T. Römke, J. Simon, *The MOL Project: An Open Extensible Metacomputer* Proc. Heterogenous Computing Workshop HCW 97, IEEE Computer Society Press, 1997, pp. 17-31.

- [47] NEOS - <http://www.mcs.anl.gov/home/otc/>
- [48] J. Czyzyk, M. Mesnier and J. More *The Network-Enabled Optimization System (NEOS) Server*, OTC Technical Report, February, 1997.
- [49] NetSolve - <http://www.cs.utk.edu/~casanova/NetSolve/>
- [50] H. Casanova and J. Dongarra, *Netsolve: A Network Server for Solving Computational Science Problems*, Technical Report CS-95-313, University of Tennessee, November 1995.
- [51] The Nile Project <http://www.nile.utexas.edu/>
- [52] Ninf - <http://ninf.etl.go.jp/>
- [53] M. Sato, H. Nakada, S. Sekiguchi, S. Matsuoka, U. Nagashima, and H. Takagi, *Ninf: A Network based Information Library for a Global World-Wide Computing Infrastructure*, Lecture Notes in Computer Science, High-Performance Computing and Networking, pages pp. 491-502, 1997.
- [54] Ninja - <http://ninja.cs.berkeley.edu>
- [55] NWS - <http://nws.npaci.edu/>
- [56] R. Wolski, N. Spring, and J. Hayes, *The Network Weather Service: A Distributed Resource Performance Forecasting Service for Metacomputing*, UCSD Technical Report Number TR-CS98-599, September, 1998.
- [57] PARDIS - <http://www.cs.indiana.edu/hyplan/kksiazek/papers.html>
- [58] K. Keahey and D. Gannon, *PARDIS: A Parallel Approach to CORBA*, Proceedings of the 6th IEEE International Symposium on High Performance Distributed Computation, August 1997.
- [59] N. Floros, K. Meacham, J. Papay and M. Surridge, *Predictive parallel scheduling for meta-applications*, Special Issue on *Metacomputing* in Future Generation Computer Systems, North Holland, Due for publication in early 1999.
- [60] SNIPE - <http://www.netlib.org/SNIPE/>
- [61] G. Fagg, K. Moore, J. Dongarra and A. Geist, *Scalable Networked Information Processing Environment (SNIPE)*, Proceeding of SuperComputing 97, San Jose, CA., November 1997.
- [62] NCSA Symera - <http://symera.ncsa.uiuc.edu/>
- [63] P. Flanigan and J. Karim, *NCSA Symera: Distributed parallel-processing using DCOM*, Dr. Dobb's Journal November 1998.

- [64] WAMM - <http://miles.cnuce.cnr.it/pp/wamm/>
- [65] R. Baraglia, G. Faieta, M. Formica, D. Laforenza. *WAMM: A Visual Interface for Managing Metacomputers*, EuroPVM'95. Ecole Normale, Supérieure de Lyon, Lyon, France, September 14-15, 1995, pp. 137-142.
- [66] UNICORE - <http://www.fz-juelich.de/unicore/>
- [67] WebSubmit - <http://www.boulder.nist.gov/websubmit/>
- [68] *Metacomputing*, Editor Wolfgang Gentzsch, Future Generation Computer Systems, Due for publication in early 1999.