1824

# **Brokering in UNICORE**

### John Brooke and Donal Fellows University of Manchester

john.brooke@manchester.ac.uk

donal.k.fellows@manchester.ac.uk

Unicore Summit 11-12 October, 2005 Sophia Antipolis, France.





### Grids as Virtual Organizations

- •Used in paper Anatomy of the Grid (Foster, Kesselman, Tuecke)
- •" ... Grid concept is coordinated resource sharing in dynamic, multi-institutional virtual organizations ..."
- •The link to the Power Grid concept is that the Power Grid is resource sharing by producers to provide a unified service to consumers.

•A large unresolved question is how do Virtual Organisations federate across security boundaries (e.g. firewalls) and organisational boundaries (resource allocation).



The analogy with a power grid

•The power grid delivers electrical power in the form of a wave (A/C wave)

•The form of the wave can change over the Grid but there is a universal (scalar) measure of power, Power = voltage x current.

•This universal measure facilitates the underlying economy of the power grid. Since it is indifferent to the way the power is produced (gas, coal, hydro etc...) different production centres can all switch into the same Grid.

•To define the abstractions necessary for a Computational Grid we MUST understand what we mean by computational resource.

MANCHESTER 1824



The University of Manchester

### An abstract space for jobcosting

- Define a job as a vector of computational resources
  - (r1,r2,...,rn)
- A Grid resource advertises a cost function for each resource
  - (cl,c2,...,cn)
- Cost function takes vector argument to produce job cost

• (r1\*c1 + r2\*c2 + ... + rn\*cn)





# A Dual Job-Space

MANCHESTER 1824

- •Thus we have a space of "requests" defined as a vector space of the computational needs of users over a Grid. For many jobs most of the entries in the vector will be null.
- •We have another space of "services" who can produce "cost vectors" for costing for the user jobs (providing they can accommodate them).
- •This is an example of a dual vector space.
- •A strictly defined dual space is probably too rigid but can provide a basis for simulations.
- •The abstract job requirements will need to be agreed. It may be a task for a broker to translate a job specification to a "user job" for a given Grid node.





Combining the strengths of UMIST and The Victoria University of Manchester

**UNIC** 



- Computational jobs ask questions about the internal structure of the provider of computational power in a manner that an electrically powered device does not.
- For example, do we require specific compilers, libraries, disk resource, visualization servers?
- What if it goes wrong, do we get support? If we transfer data and methods of analysis over the Internet is it secure?
- A resource broker for high performance computation is a different order of complexity to a broker for an electricity supplier.

# Emergent behaviour

MANCHESTER 1824

- Given this complexity, self-sustaining global Grids are likely to emerge rather than be planned.
- Planned Grids can be important for specific tasks, the LCG is an example. They are not required to be self-sustaining and questions of accounting and resource transfer are not of central interest.
- Self-sustaining (in a financial sense) Grids must have ways of accounting for and costing resource.
- This allows for trading mechanisms between suppliers of computational resource analogous to trading in energy and power Grids



The University of Manchester

MANCHESTER 1824

### **Resource Requestor and Provider Spaces**

- Resource requestor space (RR), in terms of what the user wants: e.g. Relocatable Weather Model, 10<sup>6</sup> points, 24 hours, full topography.
- Resource Provider space (RP), 128 processors, Origin 3000 architecture, 40 Gigabytes Memory, 1000 Gigabytes disk space, 100 Mb/s connection.
- We may even forward on requests from one resource provider to another, recasting of O3000 job in terms of IA64 cluster, gives different resource set.
- Linkage and staging of different stages of workflow require environmental support, a hosting environment.



# What is Computational Brokering?

• Many Definitions!

MANCHESTER 1824

The University of Manchester

- Selection of where to run a job
- Selection of what jobs to run on a resource
- Selection of when to run a job
- Discovery of how to run a job
- Generalized user representative agent

• We Focus on "Where to Run a Job"



### Why Broker over a Grid?

- Discover Resources
  - Users should not know all places to run jobs
- Discover Costs

MANCHESTER

- A typical pricing model is complex
- Sites may wish to conceal information from competitors
- Discover Queuing and Running Times
  - Some systems are very fast, but use batch queues
  - Some systems are slower due to competition, but offer immediate execution
- Improve Resource Utilization
- Account for Partial Resource Availability



The University of Manchester

### Grid Brokering





# Why Distributed Brokering?

- Sites Know the State of their Resources Best
- Sites Can Conceal their Resource Configuration
- Different VOs Need Different Selection Algorithms
  - Preferred site sets will vary
  - Different applications have different performance characteristics
- Divide-and-Conquer

MANCHESTER



The University of Manchester

### **Distributed Brokering**



The Victoria University of Manchester

#### **UNIC**<br/> **RE**

# The UNICORE Broker

• Flexible

MANCHESTER

- Can issue multiple alternative offers per resource
- Can issue offers for multiple resources
- Distributed
  - Delegate brokering requests to other brokers
- Powerful
  - Workflow brokering
  - Supports application-specific brokering
- Interoperable
  - Supports transparent brokering across both UNICORE and Globus grids



The University of Manchester

### **UNICORE** Architecture







The University of Manchester

### **UNICORE Broker Internal Architecture**





The University of Manchester

#### **Unicore/GS** Architecture





The University of Manchester

#### **UniGrids Broker Service Architecture**







### **UniGrids Grid Brokering Architecture**



20



- Open Grid Software Architecture
  - Open standard under development
  - Both business and academia
- Fundamental Architecture of Future Grids
  - Based on web-services
    - Service-oriented
  - Covers many areas of basic grid systems
- Major Aspect is Execution Management

The University of Manchester

MANCHESTER 1824

### **OGSA Execution Management**



NIC



# OGSA Resource Selection

- Part of OGSA Execution Management
- Focuses on Choosing Which Service to Use
  - Brokering

MANCHESTER

- Super-scheduling
- Partitioned Selection System
  - CSG suggests where to place atomic tasks
  - EPS creates overall plans for workflows
    - Based on cost, time to execute, etc.



The University of Manchester

MANCHESTER 1824

### **OGSA Resource Selection Architecture**

- JM asks EPS for ways to run a workflow
- EPS asks CSG for ways to instantiate atomic tasks
- CSG returns set of candidate executions
  - Where, how, etc.
- EPS picks good combinations for overall workflow
  - Task combination is non-trivial
  - Uses info-services to discover additional costs of combination
- JM uses the resulting plan(s)
- Very similar to UNICORE Broker architecture





- UNICORE Broker is Flexible Agent Supporting Rich Distributed Grid Brokering
  - Workflows
  - Application-specific brokering
  - Works with multiple low-level grids
  - Has been deployed in grid across multiple organizations
- UniGrids Broker is Next Evolutionary Phase
  - Driving OGSA standards
  - Leveraging low-level grid information systems



### Grid Resource Description Problem

- Two Independent Grid Systems
  - Unicore (http://www.unicore.org/)
  - Globus (http://www.globus.org/)
- Both Need to Describe Systems that run Compute Jobs
- Very Different Description Languages
  - Unicore's Resource model, part of the AJO Framework
  - Globus's GLUE Schema (DataTAG, iVGDL) for GT2 and GT3
- For interoperability, we want to take a Unicore job and run it on Globus resources
- Therefore, we need to translate the Job's Resource Requirements between the two Systems



### Two Types of Integration

- Right here, Right now:
  - Integrate existing features such as the way Unicore and Globus currently describe hardware resources
  - Best done by evolution, preserving much of the character of the legacy system components
- The future:

MANCHESTER 1824

- Integrating future features such as the way Unicore and Globus will describe Software Resources
- Best done by revolution, introducing a new system, reached by consensus between the two teams of architects

#### Methodology for translation



MANCHESTER 1824 UNIC

#### Methodology for translation



- Develop an ontology for the Unicore resource terminology
- Develop an ontology for the Globus resource terminology
- Map concepts in the Unicore ontology to concepts in the Globus ontology
- We assume a consensus between the concepts in Unicore and GLUE

MANCHESTER 1824



#### Methodology for translation service

- Address Data Transformation Issues for Translating Attributes
- Find a technology that has these characteristics:
  - can model the two ontologies
  - has support for linking abstract concepts to code fragments
  - easily allows someone to update mappings
  - is appropriate for a video conferencing setting
  - writes modelling information to a file format that can be used by other applications
- Use the data files created by the application to run the translator service.

MANCHESTER 1824



#### An Ontology Building Life-cycle









# LUE: Modelling resources

MANCHESTER 1824



Combining the strengths of UMIST and The Victoria University of Manchester

#### **UNIC**®RE

MANCHESTER 1824

### **GLUE: Marking up transcripts**

PC Pack 4 vB2.9h - Protocol Tool on me108 [GLUE v2] - IDI XI Knowledgebase Edit Protocol Markers view Tool Window Hep-😅 😂 🖻 🧮 🗐 🗐 🎒 Erose 📑 Concept 📩 Tosk 📑 Atrioutes and Volues Computine Element. A computing element represents an entry point into a queuing system. There is one computing element per queue. Queuing systems with multiple queues are represented by creating one computing element per queue. The information associated with a computing element is limited only to information relevant to the quene. All information about the physical resources access by a queue are represented by the Cluster • information element. Claster. A cluster is a container that groups together subclusters, or nodes. Subcluster elements represent "homogeneous" collections of computational nodes, while a nodes represent unique nodes, such as head nodes, or individual computing nodes. A cluster may be referenced by more then one computing element.  $\frac{SubClaster}{}$ . A subcluster represents a "homogeneous" collection of nodes, where the homogeneity is defined . by a collection whose required node attributes all have the same value. For example, a subcluster represents a set of nodes with the same CPU, memory, OS, network interfaces, etc. Strictly speaking, subclusters are not necessary, but they provide a convenient way of representing useful collections of nodes. A subcluster captures a node count and the set of attributes for which homogeneous values are being asserted. 235. Represents a physical computing element. This element characterizes the physical configuration of a computing node, including processors, software, storage elements, etc. Concept — Tesk — A li buliesen d Value Clue Component: marked up os (components) -GueComputing Element of arked up as 'Computing Element Give Queuring System, marked up as 'queuring system' physical resolutions access (O, O) Markup GUEV2



The University of Manchester

### **GLUE:** Provenance Information

RC Pack Annotation Tool	
File Help	0.000000
日本  本市市区  今日  京都  東田公園  キャギロ	
$f \propto x' + I = x \approx  x  \approx  x  =  x $	
Glue Load 1min	
Native Element: Host.SMPLoad.Load1min Definition: 1-minute average processor availability for an SMP node (multi CPU), wh between the available CPUs and the average runable task count during th Source: Glue Computing Element Schema version 1.1 FINAL; revised 12 March 20	ich is the difference at time X 100 003



- GLUE has container classes that include "Computing Element", "Cluster", "Subcluster" and "Host". From the heading "Representing Information", the GLUE document indicates:
  - "...hosts are composed into sub-clusters, sub-clusters are grouped into clusters, and then computing elements refer to one or more clusters."



• These container objects may hold any number optional auxiliary classes that actually describe the GRID features.

MANCHESTER 1824



# **GLUE:** Auxiliary Classes

 The documentation provides few details about the nature of a Host other than that it is a "physical computing element". Much of the meaning for Host has to be derived from what it might contain. Consider the following two valid definitions:



• A Host is a physical computing element characterized by Main Memory, a Benchmark, a Network Adapter and an Operating System

• A Host is a physical computing element characterized by an Architecture, a Processor and an Operating System.

MANCHESTER 1824



- Unicore and GLUE have different philosophies for describing resources :-(
- In Unicore, the resources are described in terms of resource requests
- In GLUE, resources are described in terms of the availability of resources.

### **Compatible Concepts**

Unicore Ontology	GLUE Ontology
Network Performance	Glue SI00 Benchmark
Network Performance	Glue SF00 Benchmark
Floating Point Performance	Glue SI00 Benchmark
Floating Point Performance	Glue SF00 Benchmark
Data Processing Performance	Glue S100 Benchmark
Data Processing Performance	Glue SF00 Benchmark
Maximum Memory Capacity Request	Host Virtual Main Memory Available
Maximum Memory Capacity Request	Subcluster Virtual Main Memory Available
Minimum Memory Capacity Request	Host RAM Main Memory Available
Minimum Memory Capacity Request	Subcluster RAM Main Memory Available
Priority Value	Priority



The University of Manchester

#### **Translation Service Prototype**







### Questions?

donal.k.fellows@manchester.ac.uk

Combining the strengths of UMIST and The Victoria University of Manchester