# Grid Computing – Technology, Trends & Attributes

Mitesh Agarwal IT Architect – HPC Solutions Sun Microsystems mitesh.agarwal@sun.com

http://sun.com/grid





# Agenda

- What is a Grid?
- How Does it Work?
- Underpinning Technologies
- Types of Grids
- Trends
- Attributes
- Discussions, Q&A





# What Is Grid Computing?

- A hardware and software infrastructure that connects distributed computers, storage devices, databases and software applications through a network, and is managed by distributed resource management software
- A way of using many resources to perform many kinds of tasks, accessible from many places by many people



 A universal computing infrastructure that builds on the power of the Net and enables more efficient computation, collaboration, and communication



### Grid Computing Explained

#### Problem-solving through resource pooling in virtual systems:

Virtualization of	Resources into a dynamic, single compute resource from federated assets			
Transparent scalability of	CPU cycles, storage			
Access that is	Dependable, consistent, pervasive, inexpens			



#### Why use Grid Computing?

- Reduce Costs
- Increase Productivity
- Increase Quality
- Shorten Time to Market
- Transparent Growth
- Do things you couldn't do before





#### The Strategic Need for Compute Power Across Diverse Industries

Industry	Computationally Intensive activities		
Life Sciences	Genetic sequencing, database queries		
Electronic Design	Simulations, verifications, regression testing		
Financial Services	Risk and portfolio analysis, simulations		
Automotive Manufacturing	Crash testing simulations, stress testing, aerodynamics modeling		
Scientific Research	Large computational problems, collaboration		
Oil and Gas			
Exploration	Simulations, seismic analysis, visualization		
Digital Content			
Creation	Frame rendering		



#### Levels of Grid Computing



#### Department Grids

- Single user community
- Single organization

#### **Enterprise Grids**

- Multiple user communities
- Single organization

#### Global Grids

- Multiple user communities
- Multiple organizations

#### Academia & Research Business



#### Grid Enabling Technologies Across all levels of deployment





### Grid Engine: Sun's Key Grid Technology

- Software to manage jobs & resources on distributed systems
- Well-established software
  - Initially developed in 1994, acquired by Sun in 2000
  - Core technology powers over 8000 Grids worldwide
- N1 Grid Engine 6 :
  - Enterprise-scale components on top of open-source core technology
  - Fully supported by Sun on: Solaris, Linux, Mac OS X, HP-UX, AIX, Irix (Windows: end of 2004)



# Examples of Grids

Cycle-stealing: use resources when unoccupied

- "desktop grid"
- datacenter sharing

Compute Farms: racks of dedicated nodes

- High-throughput
- "Beowulf"

HPC Clusters: clumps of large SMP systemsMassively-parallel jobs, eg, weather, astro



#### Or any combination of these



![](_page_10_Picture_0.jpeg)

### N1 Grid Engine Overview

![](_page_10_Picture_2.jpeg)

#### **Resource Management**

Selection of Jobs Job-based policies for resource ROI, SLA, QoS, etc User-based policies for sharing, ranking, by departments, projects, user groups, etc

Selection of Resources
System characteristics: CPU, memory, OS, patches, etc.
Status of systems: avail. mem, load, free disk space, etc.
Status of other resources: licenses, shared storage, other software, etc.

![](_page_11_Picture_0.jpeg)

### N1 Grid Engine Overview

![](_page_11_Picture_2.jpeg)

#### **Resource Control**

Control of jobs Customizable action methods for Suspend, Resume, Kill, Migrate, Restart Manual or automated via triggers

#### *Control of resources* Regulate load on systems based upon resource value thresholds Control access to systems via permissions, time/date, jobtype

![](_page_12_Picture_0.jpeg)

### N1 Grid Engine Overview

![](_page_12_Picture_2.jpeg)

#### **Resource** Accounting

Accounting of jobs Current resource consumption always monitored Total detailed consumption recorded at end of job Includes record of user, department, project, etc,

Accounting of resources Utilization of all resources recorded (eg, CPU, memory, license, etc) Usage accounting kept for users, departments, projects, etc.

![](_page_13_Picture_0.jpeg)

![](_page_13_Figure_1.jpeg)

![](_page_14_Picture_0.jpeg)

#### Grid Engine Architecture

![](_page_14_Figure_2.jpeg)

![](_page_15_Picture_0.jpeg)

#### **Batch Job Execution**

#### job script

- any script, eg, #!/bin/csh #!/bin/perl
- can do pre-work, eg, preprocessing
- may stage in/out binary/datasets, if desired
- invokes application with any options

applications run without modification: 1) use job script, or 2) submit binary directly

% qsub -1 memfree=1GB, arch=linux, jobtype=analysis myjob.sh arg1 arg2

Master

host

sent to exec host at execution

Spooled onto

master host

at submission

Exec host

![](_page_16_Picture_0.jpeg)

### Interactive Applications

![](_page_16_Figure_2.jpeg)

![](_page_17_Picture_0.jpeg)

#### Resources

#### THE HEART OF GRID ENGINE MANAGEMENT

Per Host, eg,

• load\_avg

• mem\_free

• OS/patch-level

Global, eg,

- floating licenses
- shared storage

Built-in and custom resources

- Static resources: *strings, numbers, boolean*
- Countable resources: *eg, licenses, MB of memory/disk*
- Measured resources: value provided through Load Sensor

Resources used for

- Job resource request: *job A needs 1 license and 1GB*
- Action thresholds: *suspend jobs if load\_avg* > 1.5
- Resource preference: *send jobs to hosts with least load; out of those, choose hosts with most free memory*

![](_page_18_Picture_0.jpeg)

### Parallel and Checkpointing Environments

Environment: a **set of hosts and associated procedures** that are used to support **parallel or checkpointing applications** 

applications must inherently support parallel/checkpointing execution

![](_page_18_Picture_4.jpeg)

![](_page_19_Picture_0.jpeg)

![](_page_19_Figure_1.jpeg)

![](_page_20_Picture_0.jpeg)

#### N1GE 6 Scheduler Next-generation Functionality

• Sophisticated Policy system

![](_page_20_Picture_3.jpeg)

- Align resource allocation with business priorities
- Ability to implement SLA, QoS guarantees
- Ensure greatest ROI for your assets
- Look-ahead planning
  - Resource Reservation / Preemption / Backfilling
    - Can reserve any resource, eg memory, CPU, license, disk
  - Ensure that important jobs get the resources they need, while maximizing overall utilization

![](_page_21_Picture_0.jpeg)

#### 1) Ticket-based Policies

Resource allocation aligned to business priorities

- policy basis includes: department, project, user groups, cumulative utilization, category priority
- powerful, flexible, tunable, easy to configure

![](_page_21_Figure_5.jpeg)

![](_page_22_Picture_0.jpeg)

#### Ticket policy: Share Tree

![](_page_22_Figure_2.jpeg)

![](_page_23_Picture_0.jpeg)

# Ticket Policy: Strict Order

Objective: strict determination of job dispatch order

- jobs grouped according to type or grouping
- jobs from certain groups go ahead of others, regardless of submit order or current running jobs

Job Submit order 1) Project A: job1 2) Project C: job1 3) Project B: job1 4) Project C: job2 5) Project A: job2 6) Project C: job3 7) Project B: job2 8) Project A: job3 9) Project C: job4 10) Project C: job5 11) Project B: job3

![](_page_23_Figure_6.jpeg)

![](_page_23_Picture_7.jpeg)

- Job Dispatch order 1) Project A: job1
- 2) Project A: job2
- 3) Project A: job2
- 4) Project B: job1
- 5) Project B: job2
- 6) Project B: job3
- 7) Project C: job1
- 8) Project C: job2
- 9) Project C: job3
- 10) Project C: job4
- 11) Project C: job5

![](_page_24_Picture_0.jpeg)

### Ticket Policy: Preferential Order

Objective: preferred hierarchy of groups of jobs

- jobs grouped according to projects, users, departments
- groups are given certain target percentage of resources
- scheduler tries to implement targets globally, eg, Project B jobs might go before Project A jobs if Project B is using less than target share at that time

![](_page_24_Figure_6.jpeg)

![](_page_25_Picture_0.jpeg)

# Unifying Ticket Policies

Share Tree policy average entitlements over time

Number of tickets affects priority of job

Sum of

all ticket

contributions

Functional policy *fixed entitlements* 

Override policy fine-tuning, temporary changes

![](_page_26_Picture_0.jpeg)

# 2) Urgency Policies

- Resource-based urgency: gives priority to
  - expensive asset (SW license, \$\$\$ server)
  - demanding job (multi-CPU, lots of memory)
- Deadline-based urgency
  - Enables SLA-type policies
- Wait-time urgency
  - Enables quality-of-service guarantees

![](_page_27_Picture_0.jpeg)

# 3) POSIX priority Sub-policy

- Priority value assigned to jobs with a simple number between -1023 and 1024
- Used to implement site-specific custom priorities (eg, external co-scheduler)
- Also enables admin override of any other policies

![](_page_28_Picture_0.jpeg)

# **Combining Policies**

- Final priority of jobs calculated based on unifying the three sub-policies
- Each policy normalized to 0.0 < N < 1.0 before combining using weight factors
- Allows hierarchy of policies

$$\mathbf{prio} = \mathbf{W}_{urg} \times \mathbf{N}_{urg} + \mathbf{W}_{tix} \times \mathbf{N}_{tix} + \mathbf{W}_{psx} \times \mathbf{N}_{psx}$$

$$N_{urg}$$
 = normalized Urgency  
 $N_{tix}$  = normalized Tickets  
 $N_{psx}$  = normalized Posix  
W = weighting factors

![](_page_29_Picture_0.jpeg)

![](_page_29_Figure_1.jpeg)

Jobs, with resource requirements

- number indicates priority
- length represents duration of

va au iiva mant

![](_page_30_Picture_0.jpeg)

#### Grid Resources Representation

![](_page_30_Figure_2.jpeg)

![](_page_31_Picture_0.jpeg)

![](_page_31_Figure_1.jpeg)

![](_page_32_Picture_0.jpeg)

### Scheduling with Resource Reservation

![](_page_32_Figure_2.jpeg)

![](_page_33_Picture_0.jpeg)

### Resources Reservation with Backfilling

![](_page_33_Figure_2.jpeg)

![](_page_34_Picture_0.jpeg)

#### N1GE 6 Reporting Analysis / Monitoring / Accounting

- Module for doing analysis, monitoring, accounting reports, etc.
  - Fine-grained resource recording

![](_page_34_Figure_4.jpeg)

- Stored in RDBMS in well-defined schema
- provides built-in capability for analysis, reporting, chargeback, etc
- Web-based console tool provided for generating reports, queries, etc.
- Standard SQL access for 3<sup>rd</sup> party tools

![](_page_35_Picture_0.jpeg)

#### Example Analysis DATA IMPORTED INTO STAROFFICE

![](_page_35_Figure_2.jpeg)

![](_page_36_Picture_0.jpeg)

110011101

1010110010

# N1GE 6 Enhanced Performance

- Berkeley DB spooling
- Multi-threaded Master Daemon
- New communication system
- Scalability goals:
  - 10,000 unique hosts
  - 500,000 unique jobs

![](_page_37_Picture_0.jpeg)

N1GE 6 Other Features Standards Compliance

• DRMAA 1.0

![](_page_37_Picture_3.jpeg)

- API to submit, control, monitor jobs
- Audience: developers who want to "grid-enable" their applications, without customizing to particular DRM/Grid system
- C-binding, Java-binding (future: Perl)
- XML-based monitoring
  - Status command 'qstat' outputs full status info in well-defined XML DTD

![](_page_38_Picture_0.jpeg)

# Grid Engine

#### Useful Links

- http://www.sun.com/grid
  - Success stories
  - Web-based training and testing
  - Grid Certification course information
  - Partner information
- http://www.sun.com/gridware
  - Sun Grid Engine products
- http://gridengine.sunsource.net
  - Open Source Project
  - Mailing Lists for users, developers (get answers from SGE experts)
  - FAQs, HOWTOs, documentation for users and developers
  - Courtesy Binaries for all major Unix platforms

![](_page_38_Picture_15.jpeg)

![](_page_39_Picture_0.jpeg)

#### The Grid Architecture Dilemma: Scale Vertically or Scale Horizontally?

#### Scale Vertically:

- Parallel applications: OpenMP
- Large Shared Memory
- Top Performance
- Higher acquisition cost
- Lower development and management complexity & cost

![](_page_39_Picture_8.jpeg)

#### <u>The Deciding</u> <u>Factor</u>

What do the workloads require?

#### Scale Horizontally:

- Serial and parallel applications: MPI
- Throughput
- Lower acquisition cost
- Higher development and management complexity & cost

![](_page_40_Picture_0.jpeg)

#### **Processor Choice**

SI	MP	x86
Large Shared Memory	X	
32 bit x86 code		X
64 bit x86 code		X
Mixed x86 code		X
Single Threaded		X
Multi threaded	X	

![](_page_41_Picture_0.jpeg)

# Node Selection

- Understand your workload characteristics
  - Vertical or horizontal scaled
  - Interconnect requirements
  - Processor requirements
    - 32 bit v 64 bit
    - SPARC or x86 code
    - Single threaded, multi-threaded
  - Application requirements/restrictions

![](_page_42_Picture_0.jpeg)

#### **Qmon: N1 Grid Engine's GUI**

![](_page_42_Figure_2.jpeg)

![](_page_43_Picture_0.jpeg)

# **Grid Engine Portal**

- Web interface to Sun Grid Engine & SGE-EE
- Highly secure internet access to Grid applications
- Submit and Monitor SGE jobs (no Unix skills required)
- Securely upload input files to the Portal Server
- Securely download output files to local workstation
- View X-Windows based applications using VNC
- Register new applications to the Grid in minutes
- Dynamic Project Directory Creation
- Accounting capabilities
- Runs on Sun ONE Portal Server

![](_page_44_Picture_0.jpeg)

#### Grid Engine Portal (Web Interface)

🛱 Sun ONE Portal Server 6.0 - Mozilla				
, Eile Edit View Go Bookmarks Tools Window Help		—		
Content Layout				
lob List				
<u>Blast Demo</u>				
	Submit new job			
Project List		<u>_?@x</u>		
<u>Rasmol Demo</u>	edit	delete		
<u>X-Display ACT Demo</u>	<u>edit</u>	<u>delete</u>		
<u>Bio Java Demo</u>	<u>edit</u>	<u>delete</u>		
<u>TCoffee Demo</u>	<u>edit</u>	<u>delete</u>		
<u>ClustalW Demo</u>	<u>edit</u>	<u>delete</u>		
<u>FastDNAml Demo</u>	<u>edit</u>	<u>delete</u>		
<u>NAMD PSFGen Demo</u>	<u>edit</u>	<u>delete</u>		
<u>PAML Demo</u>	<u>edit</u>	delete		
<u>Wise2 Demo</u>	edit	<u>delete</u>		
<u>PHYLIP Demo</u>	edit	delete		
<u>Fasta Application</u>	<u>edit</u>	<u>delete</u>		
LOOPP Application	<u>edit</u>	<u>delete</u>		
DOWSER Application	<u>edit</u>	<u>delete</u>		
<u>Hmmer Convert</u>	<u>edit</u>	<u>delete</u>		
<u>Hmmer Build</u>	<u>edit</u>	<u>delete</u>		
<u>Hmmer Calibrate Index Emit</u>	<u>edit</u>	<u>delete</u>		
<u>Hmmer Search Pfam</u>	edit	<u>delete</u>		
<u>Hmmer Align</u>	edit	<u>delete</u>		
<u>Hmmer Fetch</u>	<u>edit</u>	delete		
<u>Blast Demo</u>	<u>edit</u>	<u>delete</u>		
<u>ReadSeq</u>	<u>edit</u>	<u>delete</u>		
	Cre	ate new project		

![](_page_45_Picture_0.jpeg)

# What Grid is Not

It's *not* futuristic

Grid technology is:

- Here now
- Real

Sun grid solutions are:

- Based on solid technology
- Ready to be delivered today!

![](_page_46_Picture_0.jpeg)

### What Grid is Not

It's *not* new technology

- The evolution of grid has been ongoing for many years
- Sun has been an active participant in the growth and development of grid technology
- Sun has been assisting customers deploy grid technology for several years

![](_page_47_Picture_0.jpeg)

### Sun's Philosophy "All Grid, All the Time"

Data Grids

Compute Grids

Graphics Grids

![](_page_47_Picture_5.jpeg)

![](_page_48_Picture_0.jpeg)

### Complete Grid Stack

Grid Management	Applications Sun N1 Grid Engine				theme		SSIONAL
Cluster Management	Sun Control Station				Ψαμαζο Α	VIAILAS	l, Prote
Operating System	<b>s red</b> hat.	Suse	SOLARIS"	SOLARIS"			chitectura
Processor	XEON		MD Contraction		ولر		pport, Ar
Interconnect	Gigabit	Myrinet	Infini	SunFire Link band			KS, Su ervices

![](_page_49_Picture_0.jpeg)

# Compute Grid: Systems to Scale Horizontally

Low Cost Building blocks for Compute Grid Services

![](_page_49_Picture_3.jpeg)

- V20z / V40z: 32bit x86, Solaris, Linux
- Compute Grid Rack (V20z)
- V210 / V240: 64bit SPARC, Solaris
- Small 4 8 way SMP Servers
  - V480/490
  - V880/890

![](_page_49_Picture_10.jpeg)

![](_page_49_Picture_11.jpeg)

Sun Proprietary/Confidential: Internal Use Only

![](_page_49_Picture_13.jpeg)

# Horizontal Compute Grid Rack

![](_page_50_Figure_1.jpeg)

#### Compute Grid: Systems to Scale Vertical Survives Low Complexity Building blocks for Compute Grid Services

![](_page_51_Picture_1.jpeg)

- Large SMP Servers
  - SF V1280/E2900
  - 4800/4900
  - SF6800/E6900
  - SF12K/E20K
    - SF15K/E25K

Sun Proprietary/Confidential: Internal Use Only

![](_page_51_Picture_9.jpeg)

![](_page_52_Picture_0.jpeg)

### Data Grid: HPC SAN

![](_page_52_Figure_2.jpeg)

![](_page_53_Picture_0.jpeg)

### Graphics Grid:

Access for More Users to Visualization Services at Required Visual Quality and Performance Levels

![](_page_53_Figure_3.jpeg)

![](_page_54_Picture_0.jpeg)

### The Sun Ranch

- Sun microprocessor design
- Large scale EDA computing using Cluster Grid approach
- Best design practices for highly productive computing

9-yrs compute time/day! 98% CPU usage 24/7/365! 10,000+ UltraSPARC CPUs in 3 grids 100% SPARC,Solaris,Storedge 30+ HA-clusters 520TB A5x00 & T3 DiskArrays Gigabit networking 300 SunRays 250 EDA CAD Tools Tape Robot Backups 15Ksqft space at 200Watts/sqft

![](_page_54_Picture_6.jpeg)

![](_page_55_Picture_0.jpeg)

#### Sun Differentiation: Innovation Through IP Ownership

![](_page_55_Figure_2.jpeg)

![](_page_56_Picture_0.jpeg)

Mitesh Agarwal IT Architect Sun Microsystems mitesh.agarwal@sun.com

http://sun.com/grid

![](_page_56_Picture_3.jpeg)