

Hype and trends

Arnaud Legrand, CNRS, University of Grenoble

LIG laboratory, arnaud.legrand@imag.fr

December 9, 2013

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective
There Goes the
Neighborhood

Toward Exascale

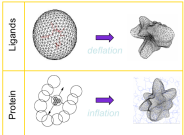
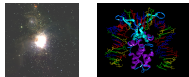
The Mont-Blanc
Project
The Deep
Project
Programming
and Application
Challenges
There Goes the
Neighborhood
Neighborhood
on Large
Systems

- 1 **Virtualization**
 - How Virtualization Changed the Grid Perspective
 - There Goes the Neighborhood
- 2 **Toward Exascale**
 - The Mont-Blanc Project
 - The Deep Project
 - Programming and Application Challenges
 - There Goes the Neighborhood
 - Neighborhood on Large Systems

Some Particularly Challenging Computations

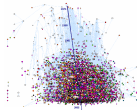
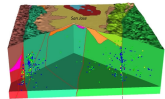
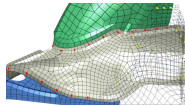
Science

- ▶ Global climate modeling
- ▶ Astrophysical modeling
- ▶ Biology (genomics; protein folding; drug design)
- ▶ Computational Chemistry
- ▶ Computational Material Sciences and Nanosciences



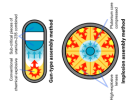
Engineering

- ▶ Crash simulation
- ▶ Semiconductor design
- ▶ Earthquake and structural modeling
- ▶ Computation fluid dynamics (airplane design)
- ▶ Combustion (engine design)



Business and Humanities

- ▶ Financial and Economic modeling
- ▶ Transaction processing, web services and search engines
- ▶ Social Networking



Defense

- ▶ Nuclear weapons – tested by simulations
- ▶ Cryptography

Performance in Scientific Computations

Scientific Problems are Large

- ▶ The finer the Mesh, the better the Prediction: need more points for quality
Forecast prediction: hundreds of km: one day ahead; 1 week ahead: kilometers
- ▶ Some intrinsically large problems (cosmology, atom studies, etc)

We want the result quickly

- ▶ Need to run numerous experiments to find the one invalidating the theory

↪ Computer systems devoted to science: the biggest existing ones

- ▶ Large amount of interconnected processing units
- ▶ High bandwidth, low latency Networks (never rely on the Internet!)

Why would Business need Computers

Initially, no need for performance

- ▶ Business computations seldom extend beyond ordinary rational arithmetic (unless when science is involved in business)
- ▶ Many desktop usage \leadsto the business uses computers without relying on them
- ▶ Computer systems distributed iff the company is: interconnect business units

And then came the Internet

- ▶ Some company relying on the Internet emerged (eBay, amazon, google)
- ▶ Computers naturally play a central role in their business plan
- ▶ Cannot afford to loose clients \leadsto **High Availability Computing**

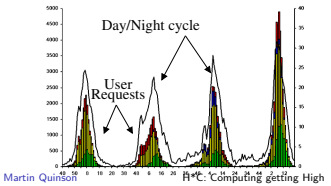
Why would Business need Computers

Initially, no need for performance

- ▶ Business computations seldom extend beyond ordinary rational arithmetic (unless when science is involved in business)
- ▶ Many desktop usage ~ the business uses computers without relying on them
- ▶ Computer systems distributed iff the company is: interconnect business units

And then came the Internet

- ▶ Some company relying on the Internet emerged (eBay, amazon, google)
- ▶ Computers naturally play a central role in their business plan
- ▶ Cannot afford to loose clients ~ **High Availability Computing**
- ▶ But load is very changing



Why would Business need Computers

Initially, no need for performance

- ▶ Business computations seldom extend beyond ordinary rational arithmetic (unless when science is involved in business)
- ▶ Many desktop usage ~ the business uses computers without relying on them
- ▶ Computer systems distributed iff the company is: interconnect business units

And then came the Internet

- ▶ Some company relying on the Internet emerged (eBay, amazon, google)
- ▶ Computers naturally play a central role in their business plan
- ▶ Cannot afford to loose clients ~ **High Availability Computing**
- ▶ But load is very changing ~ Servers dimensioned for flash crowds



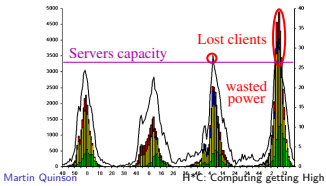
Why would Business need Computers

Initially, no need for performance

- ▶ Business computations seldom extend beyond ordinary rational arithmetic (unless when science is involved in business)
- ▶ Many desktop usage ~ the business uses computers without relying on them
- ▶ Computer systems distributed iff the company is: interconnect business units

And then came the Internet

- ▶ Some company relying on the Internet emerged (eBay, amazon, google)
- ▶ Computers naturally play a central role in their business plan
- ▶ Cannot afford to loose clients ~ **High Availability Computing**
- ▶ But load is very changing ~ Servers dimensioned for flash crowds



Amazon idea

- ▶ Rent unused power to others!
- ▶ Computers better amortized
Buy bigger ones, loose no client
- ▶ Infrastructure as a Service (IaaS)
- ▶ **Highly Cost-Efficient Computing**

Here Come the Clouds

Client Incentives

- ▶ IT maintenance burden assumed by external specialists
- ▶ **Pay only used power:** rent a server 1h, send computations *in the cloud*, enjoy
This is called **Elastic Computing**
- ▶ The created need revealed very profound: everyone wants it now
- ▶ Clients even want to rent OS+apps (PaaS) or software (SaaS)

Virtualization

- ▶ **Installing an OS:** \approx one hour. Not flexible enough.
- ▶ **Rent virtual machines instead:** overprovisioning and other optimizations

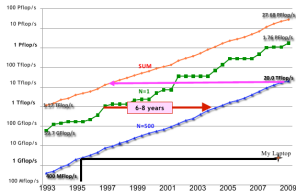
The Data Centers Growth

- ▶ Scale allows Cost Cuttings, as always. Motivation for big DC already existed
- ▶ Clouds removes the wastes due to over-dimensioning
- ⇒ **Corporate Data Centers become as big as Scientific Supercomputers!**
- ▶ ... and share the same difficulties. The twins are technically reconciled 😊

How big are these machines?

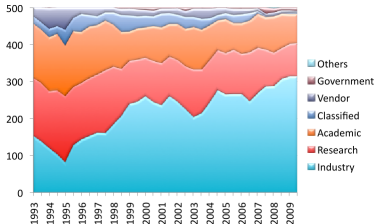
There is an International Ranking

- ▶ TOP500: updated twice a year since 1993
- ▶ Computational power growth: Exponential
- ▶ My laptop is a 10 years old supercomputer! (and my phone is a 10 years old desktop)



Machine usage

- ▶ 60% used by the industry
- ▶ The industry does science for sure
- ▶ But the increase is now due to clouds
- ▶ Some of this machines are classified HPC and Cloud don't need to argue: The big players are intelligences :)



42: The Answer to the Ultimate Question of Life, the Universe, and Everything

On Monday November 14th 2011, the Top 500 Supercomputer list was updated

Rank	Site	Computer/Year	Vendor	Cores	Rmax	Rpeak	Power
1	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 VI-Ifx 2.0GHz, Tofu interconnect / 2011	Fujitsu	705024	10510.00	11280.38	12659.9
2	National Supercomputing Center in Tianjin China	NUDT YH MPP, Xeon X5670 6C 2.93 GHz, NVIDIA 2050 / 2010	NUDT	186368	2566.00	4701.00	4040.0
3	DOE/SC/Oak Ridge National Laboratory United States	Cray XT5-HE Opteron 6-core 2.6 GHz / 2009	Cray Inc.	224162	1759.00	2331.00	6950.0
...							
42	Amazon Web Services United States	Amazon EC2 Cluster, Xeon 8C 2.60GHz, 10G Ethernet / 2011	Self-made	17024	240.09	354.10	??

From <http://perspectives.mvdirona.com>

42: The Answer to the Ultimate Question of Life, the Universe, and Everything

On Monday November 14th 2011, the Top 500 Supercomputer list was updated

Rank	Site	Computer/Year	Vendor	Cores	Rmax	Rpeak	Power
1	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 VI-Ifx 2.0GHz, Tofu interconnect / 2011	Fujitsu	705024	10510.00	11280.38	12659.9
2	National Supercomputing Center in Tianjin China	NUDT YH MPP, Xeon X5670 6C 2.93 GHz, NVIDIA 2050 / 2010	NUDT	186368	2566.00	4701.00	4040.0
3	DOE/SC/Oak Ridge National Laboratory United States	Cray XT5-HE Opteron 6-core 2.6 GHz / 2009	Cray Inc.	224162	1759.00	2331.00	6950.0
...							
42	Amazon Web Services United States	Amazon EC2 Cluster, Xeon 8C 2.60GHz, 10G Ethernet / 2011	Self-made	17024	240.09	354.10	??

► Virtualization Tax is Now Affordable

*When Cray 1 supercomputer was announced in 1976, it didn't even use virtual memory. It was believed at the time that only real-mode memory access could deliver the performance needed. Now **virtual memory** in a **guest operating system** running under a **hypervisor**.*

From <http://perspectives.mvdirona.com>

42: The Answer to the Ultimate Question of Life, the Universe, and Everything

On Monday November 14th 2011, the Top 500 Supercomputer list was updated

Rank	Site	Computer/Year	Vendor	Cores	Rmax	Rpeak	Power
1	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 VI-Ifx 2.0GHz, Tofu interconnect / 2011	Fujitsu	705024	10510.00	11280.38	12659.9
2	National Supercomputing Center in Tianjin China	NUDT YH MPP, Xeon X5670 6C 2.93 GHz, NVIDIA 2050 / 2010	NUDT	186368	2566.00	4701.00	4040.0
3	DOE/SC/Oak Ridge National Laboratory United States	Cray XT5-HE Opteron 6-core 2.6 GHz / 2009	Cray Inc.	224162	1759.00	2331.00	6950.0
...							
42	Amazon Web Services United States	Amazon EC2 Cluster, Xeon 8C 2.60GHz, 10G Ethernet / 2011	Self-made	17024	240.09	354.10	??

- ▶ Virtualization Tax is Now Affordable
- ▶ Commodity Networks can Compete with IB, Myrinet, etc.

*This is the only Top500 entrant below number 128 on the list that is not running either Infiniband or a proprietary, purpose-built network. This result at #42 is an **all Ethernet network** showing that a commodity network, if done right, can produce industry leading performance numbers.*

*What's the secret? **10Gbps** directly the host is the first part. The second is full **non-blocking networking fabric** (clos network) where all systems can communicate at full line rate at the same time.*

From <http://perspectives.mvdirona.com>

42: The Answer to the Ultimate Question of Life, the Universe, and Everything

On Monday November 14th 2011, the Top 500 Supercomputer list was updated

Rank	Site	Computer/Year	Vendor	Cores	Rmax	Rpeak	Power
1	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 VI-Ilfx 2.0GHz, Tofu interconnect / 2011	Fujitsu	705024	10510.00	11280.38	12659.9
2	National Supercomputing Center in Tianjin China	NUDT YH MPP, Xeon X5670 6C 2.93 GHz, NVIDIA 2050 / 2010	NUDT	186368	2566.00	4701.00	4040.0
3	DOE/SC/Oak Ridge National Laboratory United States	Cray XT5-HE Opteron 6-core 2.6 GHz / 2009	Cray Inc.	224162	1759.00	2331.00	6950.0
...							
42	Amazon Web Services United States	Amazon EC2 Cluster, Xeon 8C 2.60GHz, 10G Ethernet / 2011	Self-made	17024	240.09	354.10	??

- ▶ Virtualization Tax is Now Affordable
- ▶ Commodity Networks can Compete with IB, Myrinet, etc.
- ▶ Anyone can own a Supercomputer for an hour
You can have a top50 supercomputer for under \$2,600/hour

From <http://perspectives.mvdirona.com>

42: The Answer to the Ultimate Question of Life, the Universe, and Everything

On Monday November 14th 2011, the Top 500 Supercomputer list was updated

Rank	Site	Computer/Year	Vendor	Cores	Rmax	Rpeak	Power
1	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 VI-Ifx 2.0GHz, Tofu interconnect / 2011	Fujitsu	705024	10510.00	11280.38	12659.9
2	National Supercomputing Center in Tianjin China	NUDT YH MPP, Xeon X5670 6C 2.93 GHz, NVIDIA 2050 / 2010	NUDT	186368	2566.00	4701.00	4040.0
3	DOE/SC/Oak Ridge National Laboratory United States	Cray XT5-HE Opteron 6-core 2.6 GHz / 2009	Cray Inc.	224162	1759.00	2331.00	6950.0
...							
42	Amazon Web Services United States	Amazon EC2 Cluster, Xeon 8C 2.60GHz, 10G Ethernet / 2011	Self-made	17024	240.09	354.10	??

- ▶ Virtualization Tax is Now Affordable
- ▶ Commodity Networks can Compete with IB, Myrinet, etc.
- ▶ Anyone can own a Supercomputer for an hour
- ▶ This configuration has not been re-evaluated since and is thus now ranked 165

From <http://perspectives.mvdirona.com>

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective
There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project
The Deep
Project
Programming
and Application
Challenges
There Goes the
Neighborhood
Neighborhood
on Large
Systems

- 1 **Virtualization**
 - How Virtualization Changed the Grid Perspective
 - There Goes the Neighborhood
- 2 **Toward Exascale**
 - The Mont-Blanc Project
 - The Deep Project
 - Programming and Application Challenges
 - There Goes the Neighborhood
 - Neighborhood on Large Systems

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems

Dynamic scheduling of virtual machines, scalability and fault tolerance are still the issues!

Adrien Lèbre, Flavien Quesnel
ASCOLA Research Group
Ecole des Mines de Nantes

Hype and trends

A. Legrand

Virtualization

**How
Virtualization
Changed the
Grid Perspective**

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

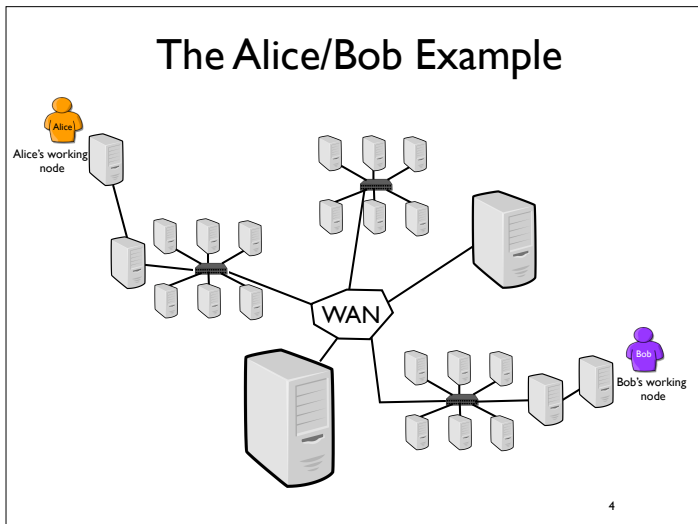
There Goes the
Neighborhood

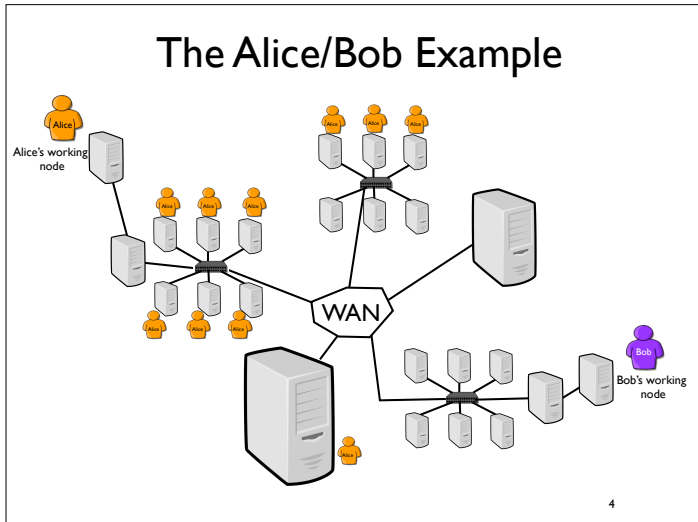
Neighborhood
on Large
Systems

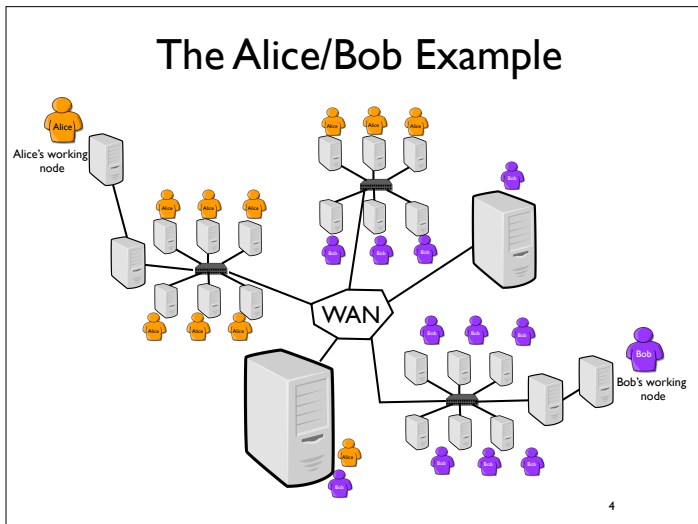
How Virtualization Changed The Grid Perspective

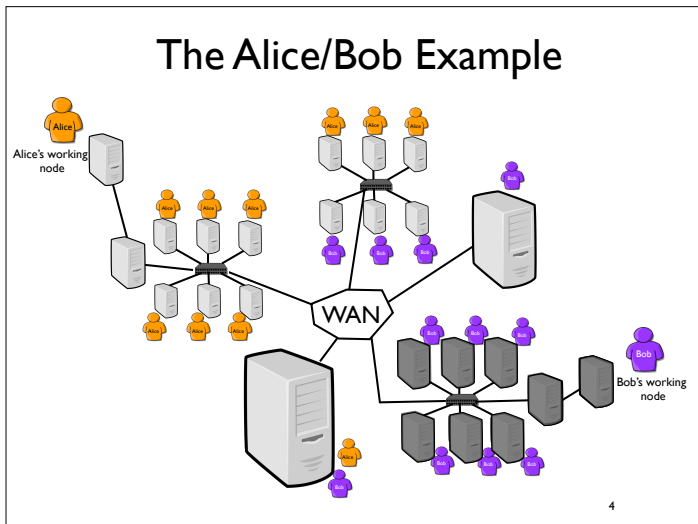
xxx Computing

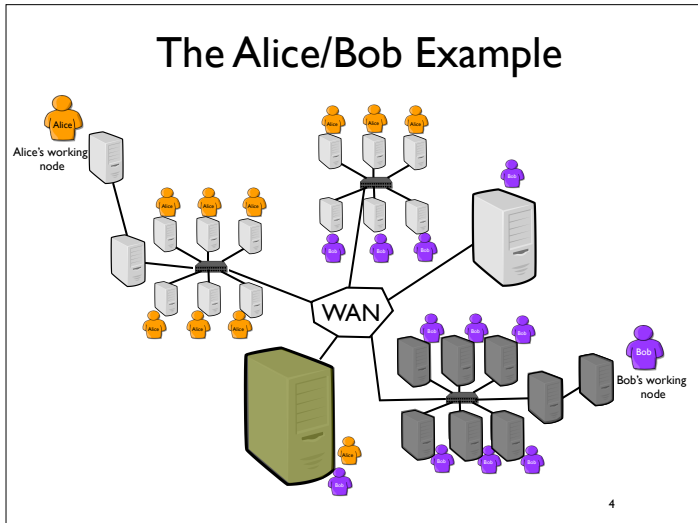
- xxx as Distributed
(Cluster / Grid / Desktop / “Hive” / Cloud / Sky / ...)
- A common objective
provide computing resources (both hardware and software)
in a flexible, transparent, secure, ... way



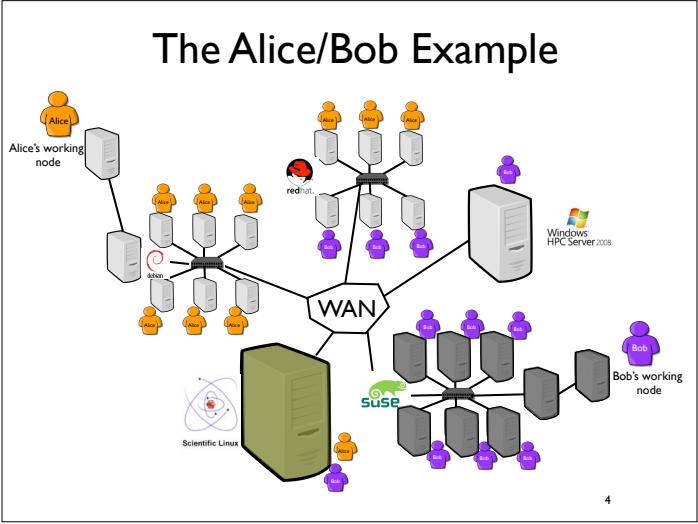








The Alice/Bob Example



Hype and trends

A. Legrand

Virtualization

How Virtualization Changed the Grid Perspective

There Goes the Neighborhood

Toward Exascale

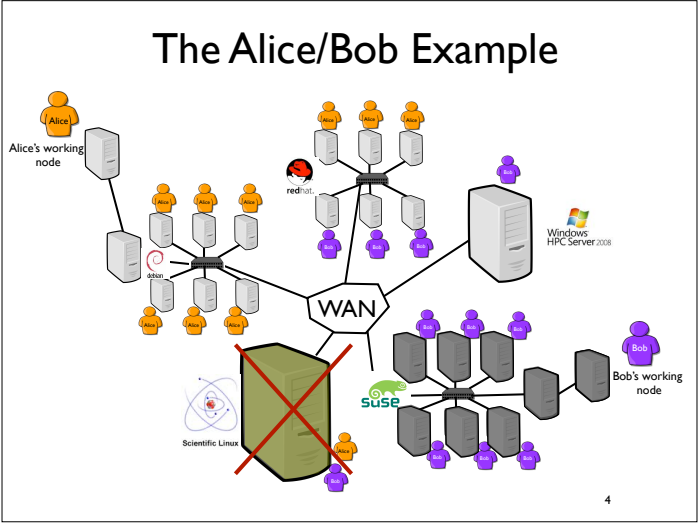
The Mont-Blanc Project

The Deep Project

Programming and Application Challenges

There Goes the Neighborhood

Neighborhood on Large Systems



Hype and trends

A. Legrand

Virtualization

How Virtualization Changed the Grid Perspective

There Goes the Neighborhood

Toward Exascale

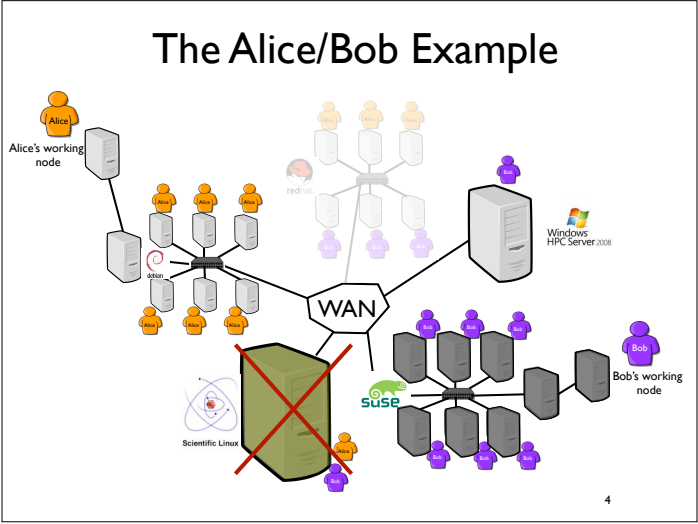
The Mont-Blanc Project

The Deep Project

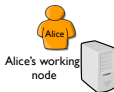
Programming and Application Challenges

There Goes the Neighborhood

Neighborhood on Large Systems



What a Grid!?!

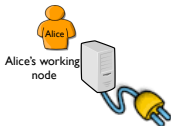


Resource booking (based on user's estimates)
Security concerns (job isolation)
Heterogeneity concerns (hardware and software)
Scheduling limitations (a job cannot be easily relocated)
Fault tolerance issues

...

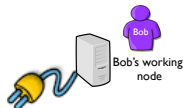


What a Grid!?!

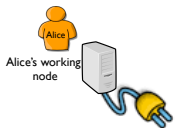


- Resource booking (based on user's estimates)
- Security concerns (job isolation)
- Heterogeneity concerns (hardware and software)
- Scheduling limitations (a job cannot be easily relocated)
- Fault tolerance issues

...



What a Grid!?!

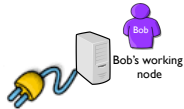


Resource

A lot of progress has been done since the 90's and several proposals partially addressed these concerns.

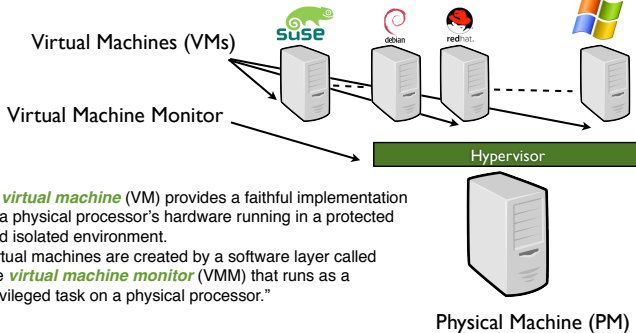
However none of them is mature enough and Strong limitations still persist !

...



Here Comes System Virtualization

- One to multiple OSES on a physical node thanks to a hypervisor (an operating system of OSES)

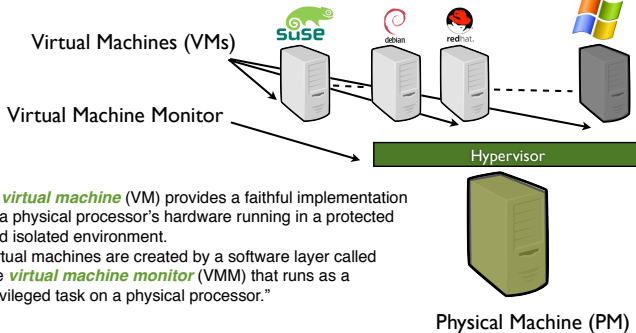


“A *virtual machine* (VM) provides a faithful implementation of a physical processor’s hardware running in a protected and isolated environment.

Virtual machines are created by a software layer called the *virtual machine monitor* (VMM) that runs as a privileged task on a physical processor.”

Here Comes System Virtualization

- One to multiple OSES on a physical node thanks to a hypervisor (an operating system of OSES)



“A *virtual machine* (VM) provides a faithful implementation of a physical processor’s hardware running in a protected and isolated environment.

Virtual machines are created by a software layer called the *virtual machine monitor* (VMM) that runs as a privileged task on a physical processor.”

Virtualization History

- Proposed in the 60's by IBM

More than 70 publications between 66 and 73

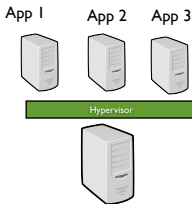
*“Virtual Machines have finally arrived. Dismissed for a number of years as merely academic curiosities, **they are now seen as cost-effective techniques for organizing computer systems resources to provide extraordinary system flexibility and support for certain unique applications**”.*

Goldberg, Survey of Virtual Machine Research, 1974

Virtualization History

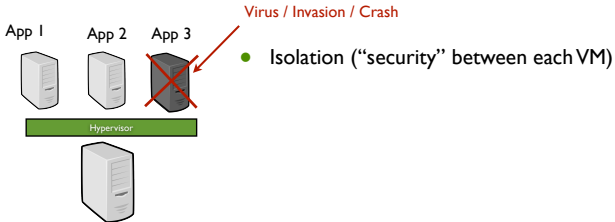
- The 80's
 - No real improvements
 - Virtualization seems given up
- End of the 90's:
 - HLL-VM : High-Level Language VM
 - Java and its famous JVM!
 - Virtual Server: Exploit for Web hosting
(Linux `chroot` / containers)
 - Revival of System Virtualization approach (VmWare/Xen)
 - Hard or soft partitioning of SMP/Numa Server

VM Capabilities

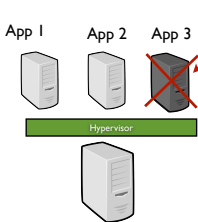


- Isolation (“security” between each VM)

VM Capabilities



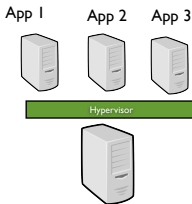
VM Capabilities



Virus / Invasion / Crash

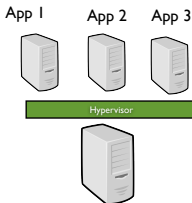
- Isolation (“security” between each VM)
- Snapshotting (a VM can be easily resumed from its latest consistent state)

VM Capabilities



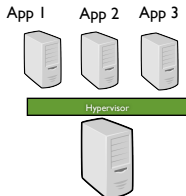
- Isolation (“security” between each VM)
- Snapshotting (a VM can be easily resumed from its latest consistent state)

VM Capabilities

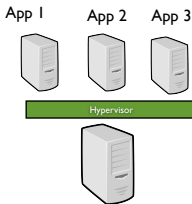


- Suspend/Resume

- Isolation (“security” between each VM)
- Snapshotting (a VM can be easily resumed from its latest consistent state)

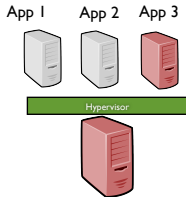


VM Capabilities

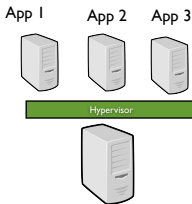


- Suspend/Resume

- Isolation (“security” between each VM)
- Snapshotting (a VM can be easily resumed from its latest consistent state)

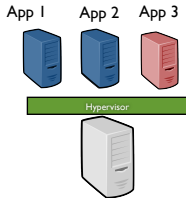


VM Capabilities

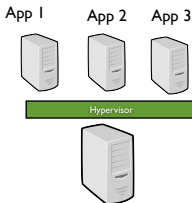


- Suspend/Resume

- Isolation (“security” between each VM)
- Snapshotting (a VM can be easily resumed from its latest consistent state)

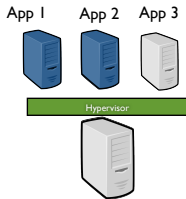


VM Capabilities

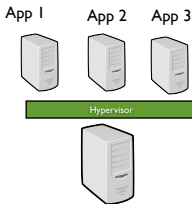


- Suspend/Resume

- Isolation (“security” between each VM)
- Snapshotting (a VM can be easily resumed from its latest consistent state)

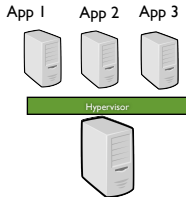


VM Capabilities

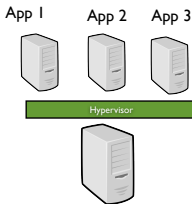


- Suspend/Resume

- Isolation (“security” between each VM)
- Snapshotting (a VM can be easily resumed from its latest consistent state)

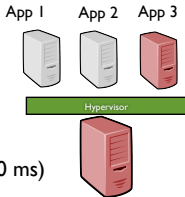


VM Capabilities



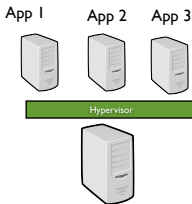
- Isolation (“security” between each VM)
- Snapshotting (a VM can be easily resumed from its latest consistent state)

- Suspend/Resume



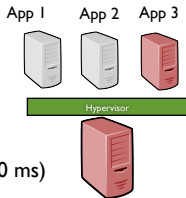
- Live migration (negligible downtime ~ 60 ms)

VM Capabilities



- Isolation (“security” between each VM)
- Snapshotting (a VM can be easily resumed from its latest consistent state)

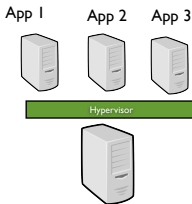
- Suspend/Resume



- Live migration (negligible downtime ~ 60 ms)

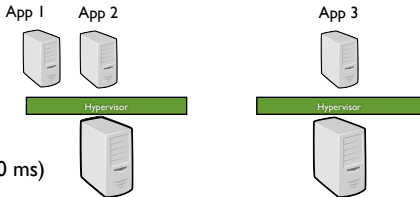


VM Capabilities



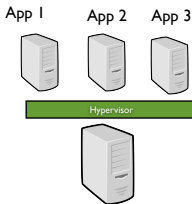
- Isolation (“security” between each VM)
- Snapshotting (a VM can be easily resumed from its latest consistent state)

- Suspend/Resume



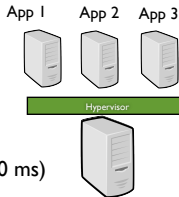
- Live migration (negligible downtime ~ 60 ms)

VM Capabilities



- Isolation (“security” between each VM)
- Snapshotting (a VM can be easily resumed from its latest consistent state)

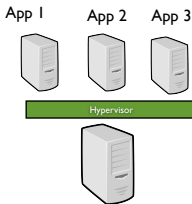
- Suspend/Resume



- Live migration (negligible downtime ~ 60 ms)

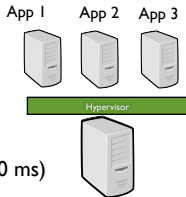


VM Capabilities

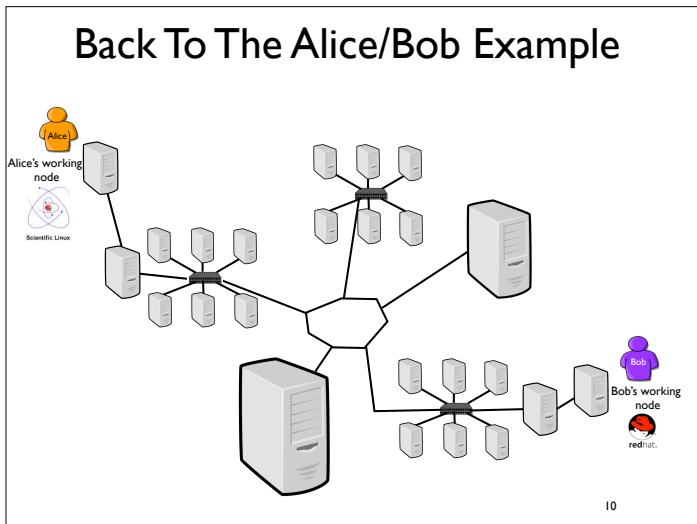


- Isolation (“security” between each VM)
- Snapshotting (a VM can be easily resumed from its latest consistent state)

- Suspend/Resume



- Live migration (negligible downtime ~ 60 ms)



Hype and trends

A. Legrand

Virtualization

How Virtualization Changed the Grid Perspective

There Goes the Neighborhood

Toward Exascale

The Mont-Blanc Project

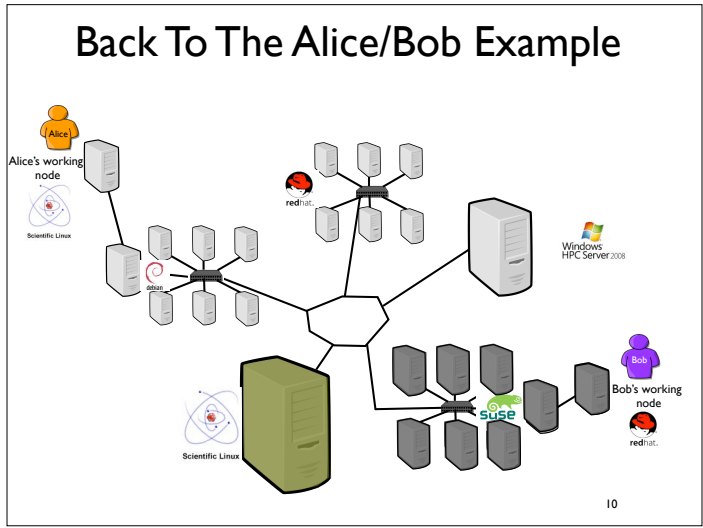
The Deep Project

Programming and Application Challenges

There Goes the Neighborhood

Neighborhood on Large Systems

Back To The Alice/Bob Example



Courtesy of Adrien Lèbre (2010)

Hype and trends

A. Legrand

Virtualization

How Virtualization Changed the Grid Perspective

There Goes the Neighborhood

Toward Exascale

The Mont-Blanc Project

The Deep Project

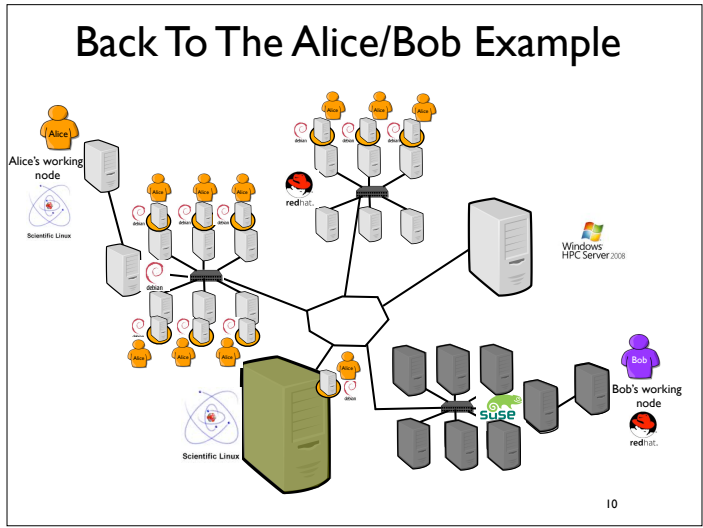
Programming and Application Challenges

There Goes the Neighborhood

Neighborhood on Large Systems

Neighborhood on Large Systems

Back To The Alice/Bob Example



Courtesy of Adrien Lèbre (2010)

Hype and trends

A. Legrand

Virtualization

How Virtualization Changed the Grid Perspective

There Goes the Neighborhood

Toward Exascale

The Mont-Blanc Project

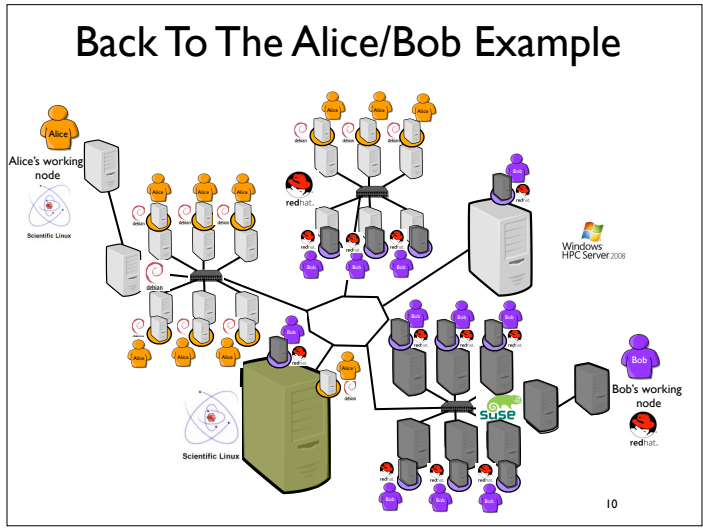
The Deep Project

Programming and Application Challenges

There Goes the Neighborhood

Neighborhood on Large Systems

Back To The Alice/Bob Example



Hype and trends

A. Legrand

Virtualization

How Virtualization Changed the Grid Perspective

There Goes the Neighborhood

Toward Exascale

The Mont-Blanc Project

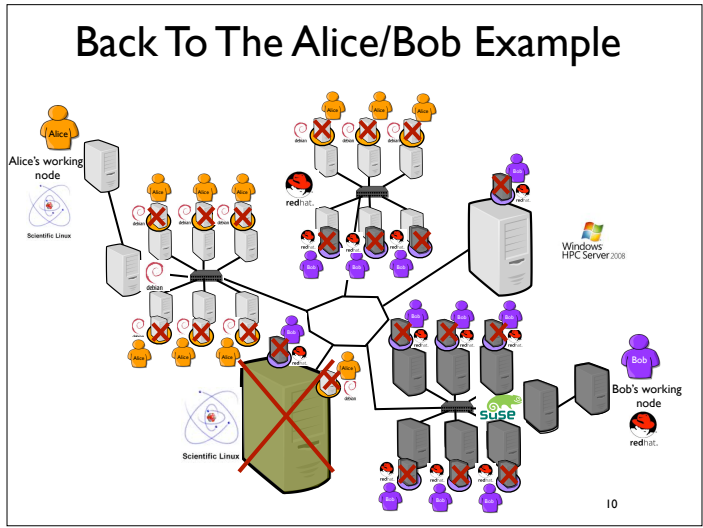
The Deep Project

Programming and Application Challenges

There Goes the Neighborhood

Neighborhood on Large Systems

Back To The Alice/Bob Example



Hype and trends

A. Legrand

Virtualization

How Virtualization Changed the Grid Perspective

There Goes the Neighborhood

Toward Exascale

The Mont-Blanc Project

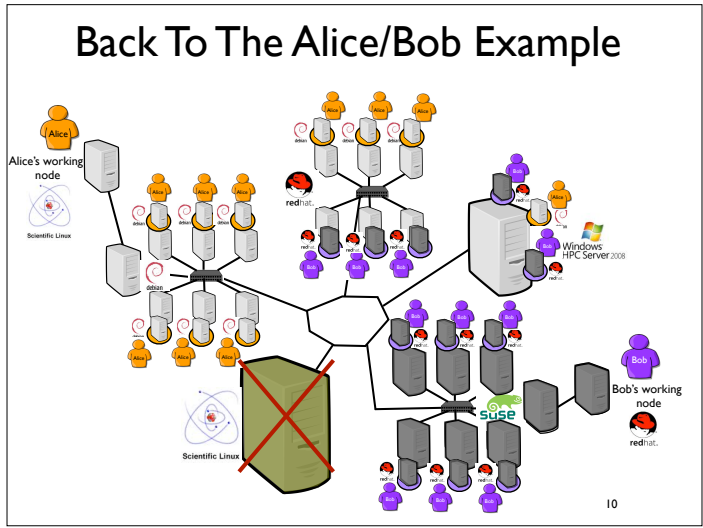
The Deep Project

Programming and Application Challenges

There Goes the Neighborhood

Neighborhood on Large Systems

Back To The Alice/Bob Example



Hype and trends

A. Legrand

Virtualization

How Virtualization Changed the Grid Perspective

There Goes the Neighborhood

Toward Exascale

The Mont-Blanc Project

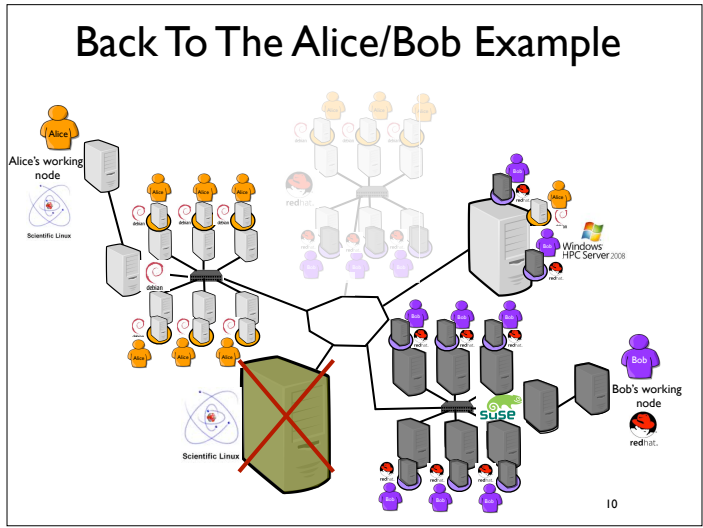
The Deep Project

Programming and Application Challenges

There Goes the Neighborhood

Neighborhood on Large Systems

Back To The Alice/Bob Example



Hype and trends

A. Legrand

Virtualization

How Virtualization Changed the Grid Perspective

There Goes the Neighborhood

Toward Exascale

The Mont-Blanc Project

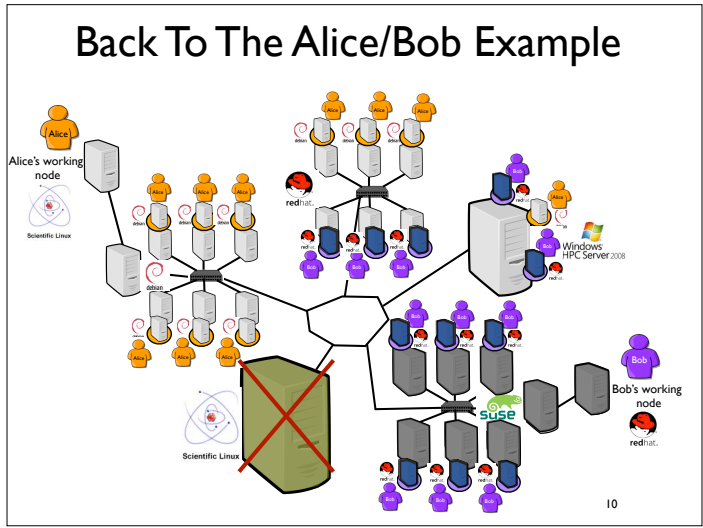
The Deep Project

Programming and Application Challenges

There Goes the Neighborhood

Neighborhood on Large Systems

Back To The Alice/Bob Example



Hype and trends

A. Legrand

Virtualization

How Virtualization Changed the Grid Perspective

There Goes the Neighborhood

Toward Exascale

The Mont-Blanc Project

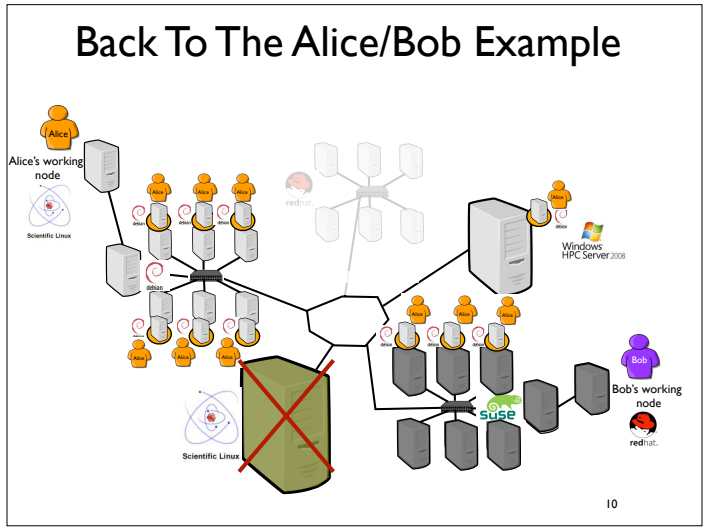
The Deep Project

Programming and Application Challenges

There Goes the Neighborhood

Neighborhood on Large Systems

Back To The Alice/Bob Example



xxxx Computing

- xxxx as Utility

“We will probably see the spread of *computer utilities*, which, like present electric and telephone utilities, will service individual homes and offices across the country”

xxxx Computing

- xxxx as Utility

“We will probably see the spread of *computer utilities*, which, like present electric and telephone utilities, will service individual homes and offices across the country”

Len Kleinrock, 1960

credits: I. Foster

1961, Prof. John McCarthy
... nree Point Checklist

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems

Focus on dynamical scheduling concerns

What can be done thanks to VM capabilities

Context

Job scheduling strategies for clusters/grids:
static allocation of resources / "user-intrusive"

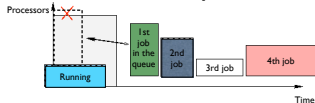
Based on user estimates (time/resources)
For a bounded amount of time
(e.g. 4 nodes for 2 hours)

Resources are reassigned at the end
of the slot without considering real
needs of applications
*(in the worst case, running applications can
be simply withdrawn from resources, i.e. G5K
best effort mode)*

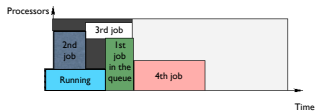
⇒ Coarse-grain exploitation
of the architecture

Context

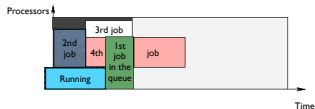
- Batch scheduler policies: closed to FCFS



Jobs arrive in the queue
and have to be scheduled.



FCFS + Easy backfilling
Jobs 2 and 3 have been backfilled.
Some resources are unused (dark areas)



Easy backfilling with preemption
The 4th job can be started without
impacting the first one.
A small piece of resources is still unused.

⇒ consolidation and preemption to finely exploit
distributed resources

Consolidation and Preemption

- Few schedulers include preemption mechanisms based on checkpointing solutions:
 - 🤨 Strongly middleware/OS dependent
 - 🤨 Still not consider application resource changes
- SSI approaches include both consolidation and preemption of processes:
 - 🤨 Strongly middleware/OS dependent
 - 🤨 SSI developments are tedious (most of them have been given up)
- Exploit all VM capabilities (start/stop - suspend/resume - migrate)

Consolidation and Preemption

- The Entropy proposal

F. Hermenier, Ph.D. in CS (University of Nantes / 2009)
Use of Live migration capability to finely exploit cluster
resources [Hermenier et al. 09]

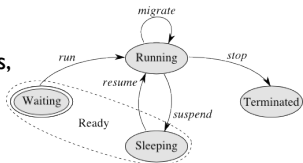
Generalization: the Cluster-Wide Context Switch concept
[Hermenier et al. 10]

- Use case - energy concerns in Datacenters

Cluster-Wide Context Switch

- General idea: manipulate **vjobs** instead of jobs (by encapsulating each submitted job in one or several VMs)

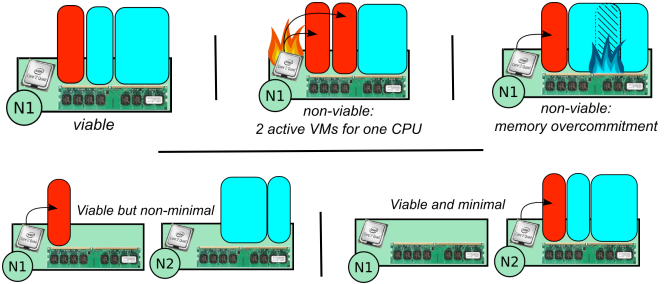
- In a similar way of usual processes, each vjob is in a particular state:



- A cluster-wide context switch (a set of VM context switches) enables to efficiently rebalance the cluster according to the: scheduler objectives / available resources / waiting vjobs queue

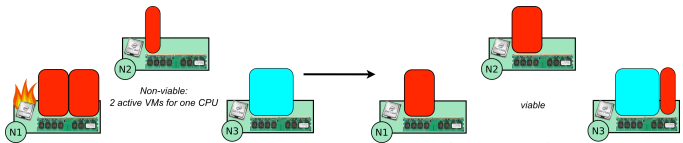
The Entropy Proposal

- To finely exploit resources (efficiency and energy constraints)
- Find the “right” mapping between VM needs and resources provided by PM



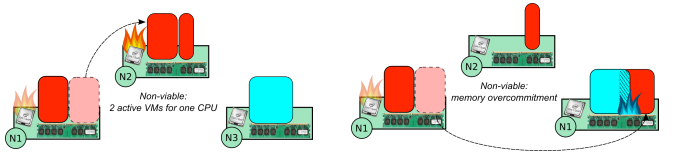
credits: F. Hermenier, Mines Nantes

The Entropy Proposal



Current Status

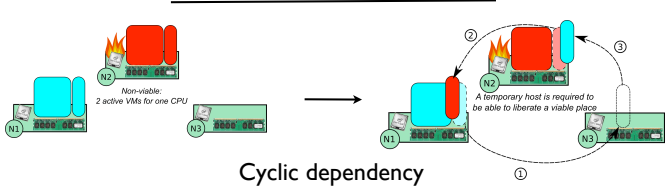
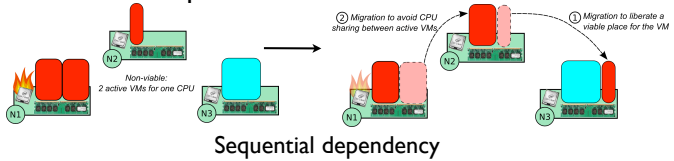
Correct Status



Non-viable manipulations

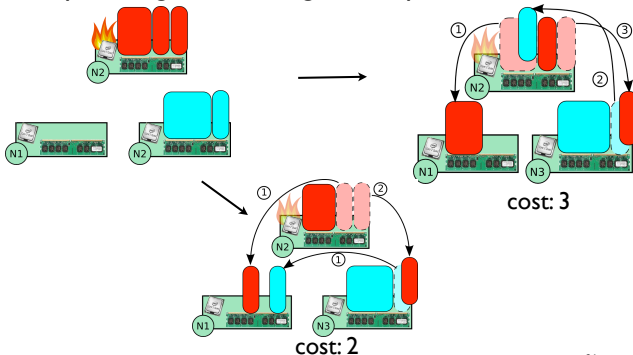
The Entropy Proposal

- Order VM Operations



The Entropy Proposal

- Optimizing the reconfiguration process

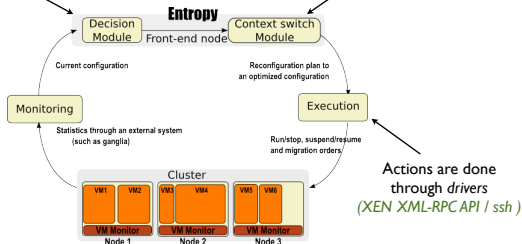


The Entropy Proposal

- The big picture: an autonomic model

Scheduling algorithm: select the jobs to run
(*objectives/strategies defined by administrators*)

Compute an efficient reconfiguration plan to reach the expected configuration
(*through the Choco constraint solver*)



- <http://entropy.gforge.inria.fr>, irc.freenode.net #entropy,

The Entropy Proposal

- To sum up



An autonomic framework to make the implementation of vjobs scheduling policies easier

Strength: composition of constraints
Developed since 2006 (ANR SelfXL / MyCloud, ANR Emergence, 10 persons)



“Prix de la croissance verte numérique” in 2009



Scalability of both computation and execution of the reconfiguration plan



Work in progress

Performance/scalability/...

Is consolidation really painless?

(12 January 2010) From

http://alan.blog-city.com/has_amazon_ec2_become_over_subscribed.htm

- ▶ Amazon in the early days was fantastic.

*Instances started up within a **couple of minutes**, they rarely had any problems and even their **SMALL INSTANCE** was **strong enough** to power even the moderately used MySQL database. For a good 20 months, all was well in the Amazon world, with really no need for concern or complaint.*

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems

Is consolidation really painless?

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems

(12 January 2010) From

http://alan.blog-city.com/has_amazon_ec2_become_over_subscribed.htm

- ▶ Amazon in the early days was fantastic.
- ▶ Neighborhood isn't what it use to be

Noisy Neighbors: *A **quick termination** and a **new spin** up would usually, through the laws of randomness, have us in a quiet neighborhood where we could do what we needed.*

*As time went on, and our load increased, the real **usefulness** of the **SMALL** instances, soon **disappeared** with us pretty much writing off any real production use of them. This is a shame, as many of our web servers are not CPU intensive, just I/O instensive.*

***Moving up to the "High-CPU Medium Instance"** as our base image has given us some of that early-pioneer feeling that we are indeed getting the intended throughput that we expect from an instance.*

Is consolidation really painless?

(12 January 2010) From

http://alan.blog-city.com/has_amazon_ec2_become_over_subscribed.htm

- ▶ Amazon in the early days was fantastic.
- ▶ Neighborhood isn't what it use to be
- ▶ The commute is such a drag

However, in the last month of two, we've even noticed that these "High-CPU Medium Instance" have been suffering a similar fate of the Small instances.

In normal circumstances, a ping between two internal nodes within Amazon is around the 0.3ms level, with the odd ping reporting a whopping 7ms ever 30 or so packets.

When our instances appear to be dying or at least shaky, then this network latency jumps up to a whopping 7241ms.

**Under extreme load, the virtual operating system
is not able to process the network queue.**

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems

Is consolidation really painless?

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems

(12 January 2010) From

http://alan.blog-city.com/has_amazon_ec2_become_over_subscribed.htm

- ▶ Amazon in the early days was fantastic.
- ▶ Neighborhood isn't what it use to be
- ▶ The commute is such a drag
- ▶ Different road surfaces

In one particular "fire fighting mode", we spent an hour literally spinning up new instances and terminating them until we found ourselves on a node that actually responded to our network traffic.

Not all the Amazon instances are equal in terms of the underlying hardware, and depending on which processor you get allocated can make a huge difference to the performance of your running instance.

So not only should we check for the CPU we are running on, we now must also take note of the network performance before we can safely push an instance into production.

This is not what cloud computing is all about.

Cloud versus Cloud



Too complex: do I need to become a sys admin?

What is the best programming model, what are the tools I need to make effective use of them?

It costs too much! And what if Amazon raises prices?

Performance: especially I/O, especially Big Data!

Custom user environments!
On-demand access!
Elastic computing!
Isolation!
Capital expense -> operational expense!



Cloud Storage Basics

- **Ephemeral/Transient Storage**
 - Local virtual disk attached to an instance
 - Persists only for the lifetime of an instance
 - Included in the cost of an instance
 - Varying capacity, e.g., 160 GB-48 TB on AWS
- **Persistent attached storage**
 - Block storage volumes that can be attached to an instance
 - Lifetime independent of a particular instance, can be mounted by many
 - Price based on space and time used
 - E.g., AWS Elastic Block Storage (EBS), Azure drives
- **Storage Clouds**
 - Data storage as binary objects (BLOBs)
 - Price differs based levels of service, e.g., access time or reliability, space used and time
 - E.g., AWS Simple Storage Service (S3), AWS Glacier, Azure BLOBs, Google Cloud Storage

Streaming Applications



12/1/13

 www.nimbusproject.org

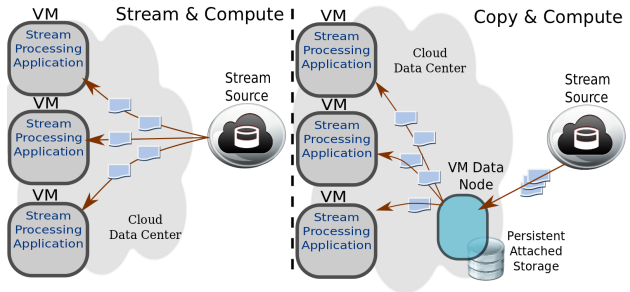
- Repeatedly apply an operation to a stream of data (time events)
- Examples:
 - Virtual Observatories: OOI, Forest project at ANL, IFC
 - Experiment processing: STAR, APS
- Requirements:
 - An “always-on” service
 - Real-time event-based data stream processing capabilities
 - Highly volatile need for data distribution and processing

4

ATLAS Data Analysis



Streaming Scenarios



Streaming Scenarios (2)

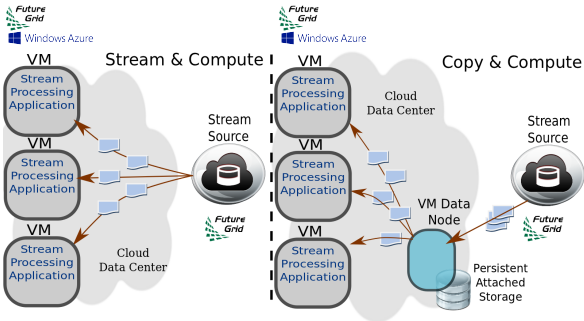
Stream&Compute (SC)

- Simpler model with fewer moving parts
- Potentially better response time
- Overlap computation and communication (potentially faster)
- Uses ephemeral storage (potentially cheaper)

Copy&Compute (CC)

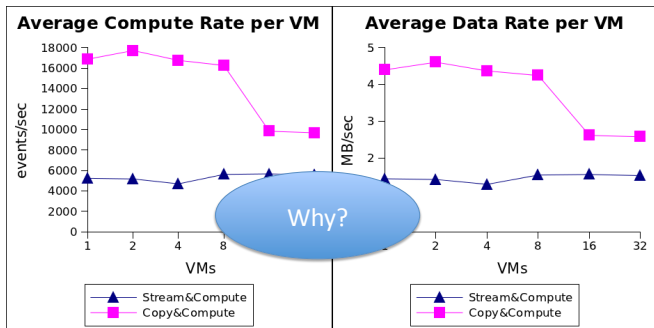
- Independent of network saturation
- Persistent storage: less liable to data loss

Experimental Configuration

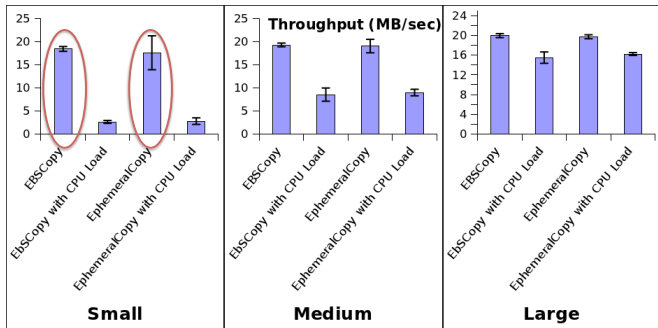


- *Compute rate: events processed per time unit*
- *Data rate: amount of data acquired per time unit*

SC versus CC (Azure)



Data Throughput vs CPU Load



Cost

$$TotalCost = \frac{TotalEvents}{CompRate_{Total}} * (N_{VMsData} + N_{VMsComp}) * VM_{Cost} + Storage_{Cost}$$

- Cost of instance: ~\$0.1 per hour
- Cost of storage: ~\$0.1 per 1GB month
- In our case (320M events & 5 GB attached storage)
 - Stream&Compute: \$1.33
 - Copy&Compute: \$0.48
 - Overall: SC is 2.77 times more expensive

Conclusions

- To stream or not to stream?
 - Not to stream!
 - Difference of ~4x in performance and ~3x in cost
- Amplification of virtualization performance trade-offs in the presence of remote traffic
- Hypervisor design
 - Need for controlled allocation of CPU to I/O processing
- *Paper: Tudoran et al., "Evaluating Streaming Strategies for Event Processing across Infrastructure Clouds", submitted to CCGrid*

Recap

Types and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems

- ▶ Virtualization changed the grid perspective because it solved many of heterogeneity, isolation and fault tolerance issues.
- ▶ Virtualization changed classical batch scheduling issue because preemption helps.
- ▶ Now, focus on consolidation and energy minimization.

Recap

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems

- ▶ Virtualization changed the grid perspective because it solved many of heterogeneity, isolation and fault tolerance issues.
- ▶ Virtualization changed classical batch scheduling issue because preemption helps.
- ▶ Now, focus on consolidation and energy minimization.
- ▶ Remember EC2 has been #42 on Top500:
 - ▶ Commodity Networks can Compete with IB, Myrinet, etc.
 - ▶ No power consumption was reported...
 - ▶ The worldwide demand for data center power in 2005 was equivalent to the output of about 17 1,000-megawatt power plants (1% of world electricity consumption in 2005).
 - ▶ Google continuously uses enough electricity to power 200,000 homes.
The average energy consumption on the level of a typical user, is about 180 watt-hours a month (a 60-watt light bulb for three hours).

<http://www.nytimes.com/2011/09/09/technology/google-details-and-defends-its-use-of-electricity.html>

- ▶ Partly because of the 2008 recession, power consumption by data centers hasn't grown at expected rates.

http://www.nytimes.com/2011/08/01/technology/data-centers-using-less-power-than-forecast-report-says.html?_r=1

html?_r=1

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective
There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project
The Deep
Project
Programming
and Application
Challenges
There Goes the
Neighborhood
Neighborhood
on Large
Systems

- 1 Virtualization
 - How Virtualization Changed the Grid Perspective
 - There Goes the Neighborhood
- 2 Toward Exascale
 - The Mont-Blanc Project
 - The Deep Project
 - Programming and Application Challenges
 - There Goes the Neighborhood
 - Neighborhood on Large Systems

Tianhe-2 (MilkyWay-2)

Hype and trends

A. Legrand

Virtualization

How

Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

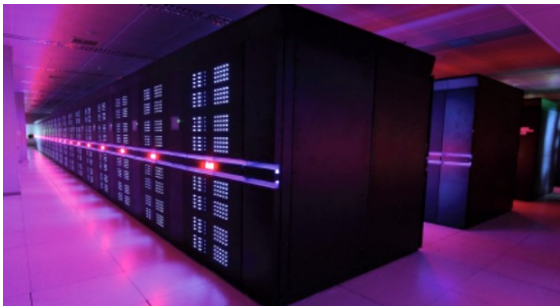
The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems



- ▶ Simulation, analysis, and government security applications
- ▶ 16,000 computer nodes, each comprising two Intel Ivy Bridge Xeon processors and three Xeon Phi chips for a total of 3,120,000 cores
- ▶ 33.8 PFlops (Peak=54.9 PFlops)
- ▶ **17.8 GW!!!**

Tianhe-2 (MilkyWay-2)

Hype and trends

A. Legrand

Virtualization

How

Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems



- ▶ Simulation, analysis, and government security applications
- ▶ 16,000 computer nodes, each comprising two Intel Ivy Bridge Xeon processors and three Xeon Phi chips for a total of 3,120,000 cores
- ▶ 33.8 PFlops (Peak=54.9 PFlops)
- ▶ **17.8 GW!!!**

Toward Exascale

Hype and trends

A. Legrand

Virtualization

How Virtualization Changed the Grid Perspective

There Goes the Neighborhood

Toward Exascale

The Mont-Blanc Project

The Deep Project

Programming and Application Challenges

There Goes the Neighborhood

Neighborhood on Large Systems

Green500 Rank	MFLOPS/W	Site*	Computer*	Total Power (kW)
1	2026.48	IBM - Rochester	BlueGene/Q, Power BQC 16C 1.60 GHz, Custom	85.12
2	2026.48	IBM Thomas J. Watson Research Center	BlueGene/Q, Power BQC 16C 1.60 GHz, Custom	85.12
3	1996.09	IBM - Rochester	BlueGene/Q, Power BQC 16C 1.60 GHz, Custom	170.25
4	1988.56	DOE/NNSA/LLNL	BlueGene/Q, Power BQC 16C 1.60 GHz, Custom	340.50
5	1689.86	IBM Thomas J. Watson Research Center	NNSA/SC Blue Gene/Q Prototype 1	38.67
6	1378.32	Nagasaki University	DEGIMA Cluster, Intel i5, ATI Radeon GPU, Infiniband QDR	47.05
7	1266.26	Barcelona Supercomputing Center	Bullx B505, Xeon E5649 6C 2.53GHz, Infiniband QDR, NVIDIA 2090	81.50
8	1010.11	TGCC / GENCI	Curie Hybrid Nodes - Bullx B505, Nvidia M2090, Xeon E5640 2.67 GHz, Infiniband QDR	108.80
9	963.70	Institute of Process Engineering, Chinese Academy of Sciences	Mole-8.5 Cluster, Xeon X5520 4C 2.27 GHz, Infiniband QDR, NVIDIA 2050	515.20
10	958.35	GSIC Center, Tokyo Institute of Technology	HP ProLiant SL390s G7 Xeon 6C X5670, Nvidia GPU, Linux/Windows	1243.80

- ▶ Exponential improvements at the rate of one order of magnitude every 3 years: One petaflops was achieved in 2008, one exaflops is expected in 2020.
- ▶ Based on a 20 MW power budget, which is already very high, this requires an efficiency of 50 GFLOPS/Watt.
- ▶ However, the current leader in energy efficiency (IBM BlueGene/Q) achieves only 2.0 GFLOPS / Watt. Thus, a 25× improvement is required.

Where does the power go?

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective
There Goes the
Neighborhood

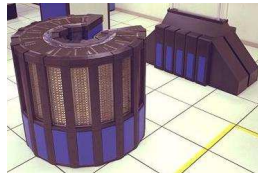
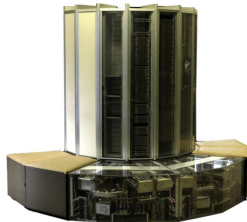
Toward Exascale

The Mont-Blanc
Project
The Deep
Project
Programming
and Application
Challenges
There Goes the
Neighborhood
Neighborhood
on Large
Systems

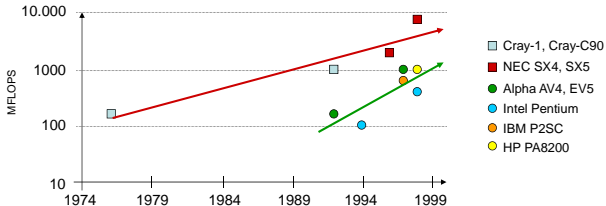
- ▶ In current systems the processors consume a lion's share of the energy approximately 43% or more.
- ▶ The remaining energy is used to power up the memories, the interconnection network, and the storage system.
- ▶ Furthermore, a significant fraction is wasted in power supply overheads, and in thermal dissipation (cooling), which do not contribute to performance at all.

In the beginning ... there were only supercomputers

- Built to order
 - Very few of them
- Special purpose hardware
 - Very expensive
- Control Data, Convex, ...
- Cray-1
 - 1975, 160 MFLOPS
 - 80 units, 5-8 M\$
- Cray X-MP
 - 1982, 800 MFLOPS
- Cray-2
 - 1985, 1.9 GFLOPS
- Cray Y-MP
 - 1988, 2.6 GFLOPS
- Fortran+vectorizing compilers



The Killer Microprocessors



- Microprocessors killed the Vector supercomputers
 - They were not faster ...
 - ... but they were significantly cheaper and greener
- Need 10 microprocessors to achieve the performance of 1 Vector CPU
 - SIMD vs. MIMD programming paradigms

Then, commodity took over special purpose



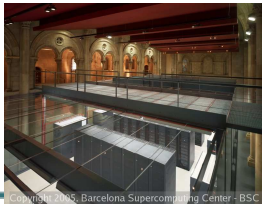
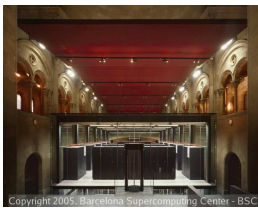
- ASCI Red, Sandia
 - 1997, 1 Tflops (Linpack),
 - 9298 cores @ 200 Mhz
 - 1.2 Tbytes
 - Intel Pentium Pro
 - Upgraded to Pentium II Xeon, 1999, 3.1 Tflops

- ASCI White, LLNL
 - 2001, 7.3 TFLOPS
 - 8192 proc. @ 375 Mhz,
 - 6 Tbytes
 - (3+3) Mwats
 - IBM Power 3

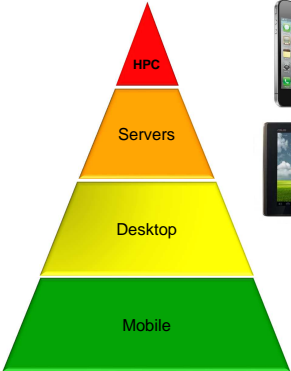
Message-Passing Programming Models

Finally, commodity hardware + commodity software

- MareNostrum
 - Nov 2004, #4 Top500
 - 20 Tflops, Linpack
 - IBM PowerPC 970 FX
 - Blade enclosure
 - Myrinet + 1 GbE network
 - SuSe Linux

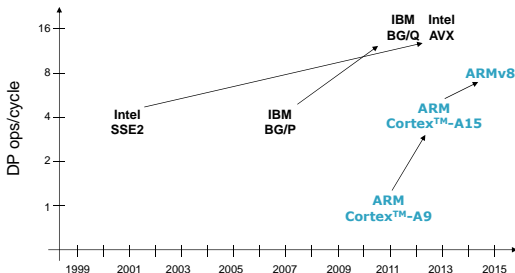


The next step in the commodity chain



- Total cores in Jun'12 Top500
 - 13.5 Mcores
- Tablets sold in Q4 2011
 - 27 Mtablets
- Smartphones sold Q4 2011
 - > 100 Mphones

ARM Processor improvements in DP FLOPS



- IBM BG/Q and Intel AVX implement DP in 256-bit SIMD
 - 8 DP ops / cycle
- ARM quickly moved from optional floating-point to state-of-the-art
 - ARMv8 ISA introduces DP in the NEON instruction set (128-bit SIMD)

Hype and trends

A. Legrand

Virtualization

How Virtualization Changed the Grid Perspective
There Goes the Neighborhood

Toward Exascale

The Mont-Blanc Project

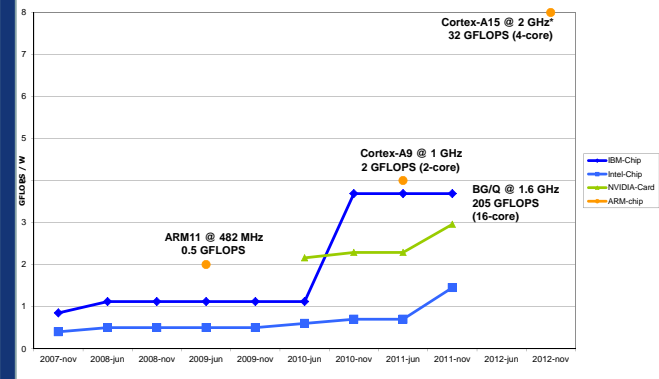
The Deep Project

Programming and Application Challenges

There Goes the Neighborhood

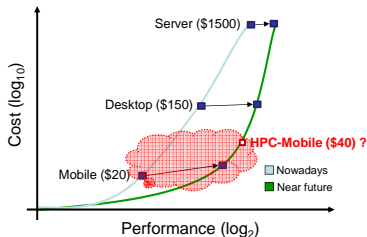
Neighborhood on Large Systems

ARM processor efficiency vs. IBM / Intel / Nvidia



* Based on ARM Cortex-A9 @ 2GHz power consumption on 45nm, not an ARM commitment

Are the “Killer Mobiles™” coming?



- Where is the sweet spot? Maybe in the low-end ...
 - Today ~ 1:8 ratio in performance, 1:100 ratio in cost
 - Tomorrow ~ 1:2 ratio in performance, still 1:100 in cost ?
- The same reason why microprocessors killed supercomputers
 - Not so much performance ... but much lower cost, and power

Hype and trends

A. Legrand

Virtualization

How Virtualization Changed the Grid Perspective

There Goes the Neighborhood

Toward Exascale

The Mont-Blanc Project

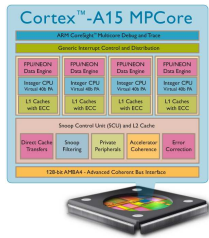
The Deep Project

Programming and Application Challenges

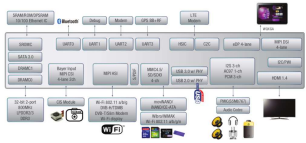
There Goes the Neighborhood

Neighborhood on Large Systems

Killer mobile™ example: Samsung Exynos 5450 *



- 4-core ARM Cortex-A15 @ 2 GHz
 - 32 GFLOPS
- 8-core ARM Mali T685
 - 168 GFLOPS*
- Dual channel DDR3 memory controller
- All in a low-power mobile socket



* Data from web sources, not an ARM or Samsung commitment

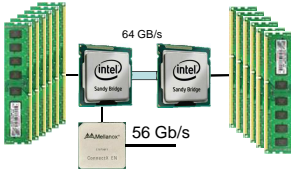


Are we building BlueGene again?

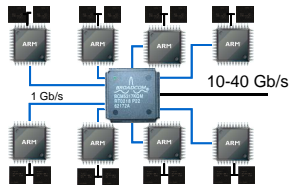
- Yes ...
 - Exploit Pollack's Rule in presence of abundant parallelism
 - Many small cores vs. Single fast core
- ... and No
 - Heterogeneous computing
 - On-chip GPU
 - Commodity vs. Special purpose
 - Higher volume
 - Many vendors
 - Lower cost
 - Lots of room for improvement
 - No SIMD / vectors yet ...
 - Build on Europe's embedded strengths



Can we achieve competitive performance?

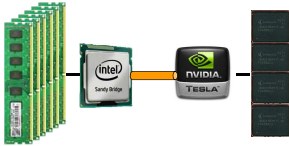
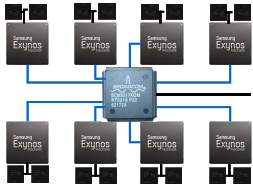


- 2-socket Intel Sandy Bridge
 - 370 GFLOPS
 - 1 address space
 - 44 MB on-chip memory
 - 136 GB/s
 - 64 GB/s intra-node (2 x QPI)




- 8-socket ARM Cortex A-15
 - 256 GFLOPS
 - 8 address spaces
 - 16 MB on-chip memory
 - 102 GB/s
 - 1 Gb/s intra-node (1 GbE)

Can we achieve competitive performance?

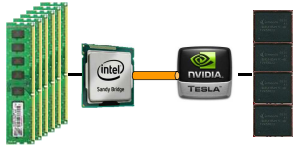



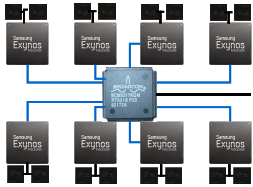
- Sandy Bridge + Nvidia K20
 - 1685 GFLOPS
 - 2 address spaces
 - 32 GB/s between CPU-GPU
 - 16x PCIe 3.0
 - 68 + 192 GB/s

- 8-socket Exynos 5450
 - 1600 GFLOPS
 - 16 address spaces
 - 12.8 GB/s between CPU-GPU
 - Shared memory
 - 102 GB/s

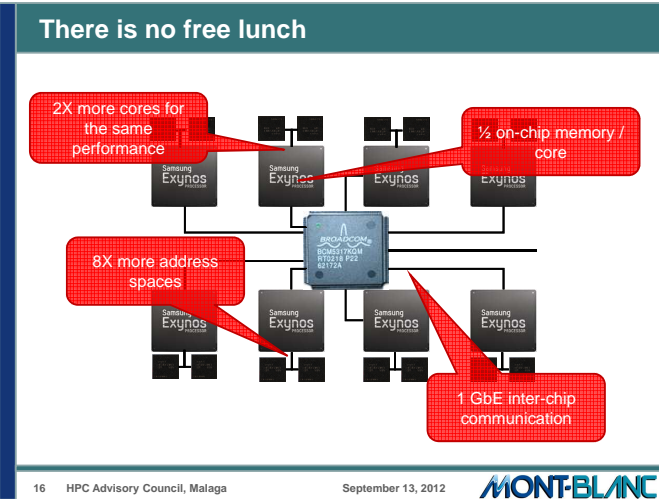
14 HPC Advisory Council, Malaga
September 13, 2012


Then, what is so good about it?

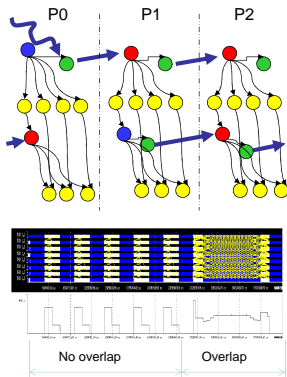




- Sandy Bridge + Nvidia K20
 - > \$3000
 - > 400 Watt
- 8-socket Exynos 5450
 - < \$200
 - < 100 Watt

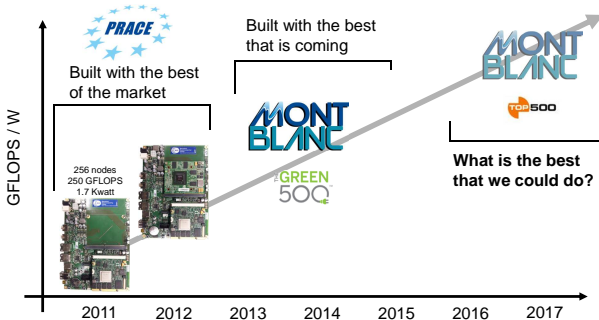


OmpSs runtime layer manages architecture complexity



- Programmer exposed a simple architecture
- Task graph provides lookahead
 - Exploit knowledge about the future
- Automatically handle all of the architecture challenges
 - Strong scalability
 - Multiple address spaces
 - Low cache size
 - Low interconnect bandwidth
- Enjoy the positive aspects
 - Energy efficiency
 - Low cost

A big challenge, and a huge opportunity for Europe



- Prototypes are critical to accelerate software development
 - System software stack + applications

Possible Architecture: 200 PFlops with 10MWatt

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

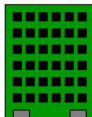
There Goes the
Neighborhood

Neighborhood
on Large
Systems

- ▶ On a power envelope of 10 Watts, this implies that each multi-core chip must achieve 600 GFLOPS of peak performance.
- ▶ If we assume 8 GFLOPS processors (2 GHz, 4 operations per cycle), this requires 75 cores per chip, consuming 0.15 Watts / core.
- ▶ As a reference, the current dual-core ARM Cortex A9 consumes 1.9 Watts at 2 GHz and uses 6.7 mm². The 800 Mhz version consumes only 0.5 Watts and 4.6 mm². That is, 0.25 Watts per processor, quite close to the target 0.15 Watts required.
- ▶ We are much closer to the target in this direction, than using today's high-end processors.



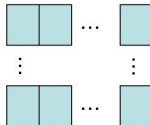
Multi-core chip:
60 GFLOPS /W
10 Watts
600 GFLOPS
8 GFLOPS / core
75 cores / chip
0.15 Watts / core



Compute node:
36 chips
2,700 cores
22 TFLOPS
1.000 Watts / node



Rack:
42 compute nodes
1.512 chips
86.400 cores
0.9 PFLOPS
50 Kwatts / rack



Exascale system:
225 racks
16.800 nodes
604.800 chips
4.5 M cores
200 PFLOPS
10 MWatts

Projection for Exascale

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

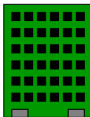
There Goes the
Neighborhood

Neighborhood
on Large
Systems

- ▶ $1000 \text{ Pflops} / 20\text{MWatt} = 10\text{GFlops} / \text{Watt} \rightsquigarrow 200 \text{ 8 Gflops core} / \text{chips and } 0.05\text{Watt} / \text{core!!!}$



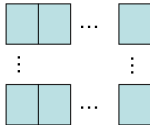
Multi-core chip:
150 GFLOPS / W
10 Watts
1.5 TFLOPS
8 GFLOPS / core
200 cores / chip
0.05 Watts / core



Compute node:
36 chips
7,200 cores
58 TFLOPS
1,000 Watts / node



Rack:
42 compute nodes
1.512 chips
302,400 cores
2.5 PFLOPS
50 Kwatts / rack



Exaflop system:
400 racks
16,800 nodes
604,800 chips
12 M cores
1,000 PFLOPS
20 MWatts

- ▶ Require new memory architecture, network, ...

Thermal dissipation

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective
There Goes the
Neighborhood

Toward Exascale

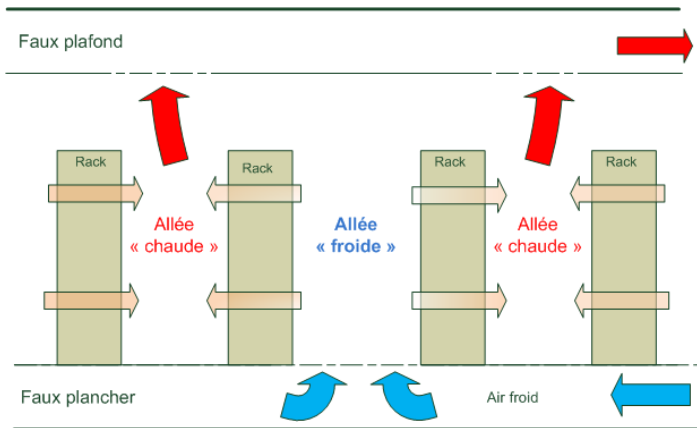
The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems



Thermal dissipation

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

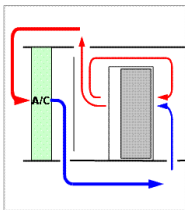
Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems

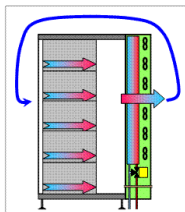
Air-cooled (a)

PUE = 1.9



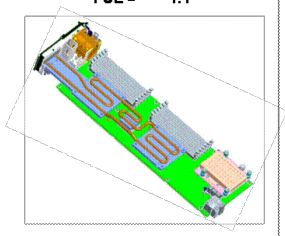
Water-cooled doors (b)

PUE = 1.5



Direct-Liquid-cooling (c)

PUE = 1.1



Thermal dissipation

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

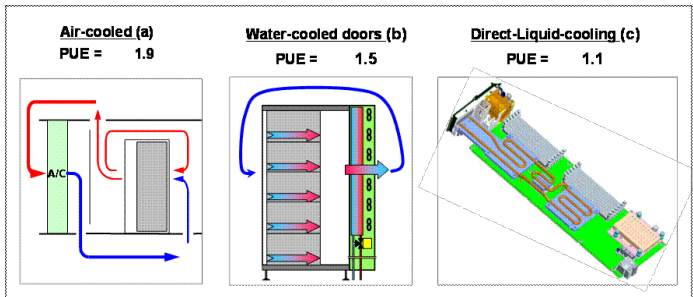
The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems



The use of low-power embedded technologies will have significant implications on the thermal characteristics of the system, which will require re-evaluating these cooling methods, and maybe proposing new ones.

Interconnect

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective
There Goes the
Neighborhood

Toward Exascale

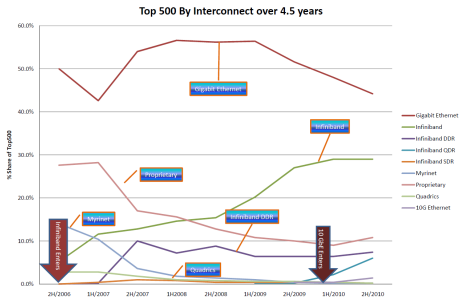
The Mont-Blanc
Project
The Deep
Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems

- ▶ Current HPC systems are characterized by either the large scale integration of low-power embedded devices, or clusters of commodity x86 servers (with increasing use of GPU acceleration).
- ▶ The interconnect for such systems are either based on proprietary technology or on widely available switch network technology such as Infiniband or Ethernet.



- ▶ For the majority of HPC cluster systems in the Top100, the network of choice is Infiniband primarily due to the performance and price
- ▶ For more than a third of systems within the Top500 the dominant interconnect is Ethernet
- ▶ Ethernet is a standard low-power interface that can be used as the method of interconnection
- ▶ Storage and Network Convergence

Toward Exascale

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective
There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems

Green500 Rank	MFLOPS/W	Site*	Computer*	Total Power (kW)
1	4,503.17	GSIC Center, Tokyo Institute of Technology	TSUBAME-KFC - LX 1U-4GPU/104Re-1G Cluster, Intel Xeon E5-2620v2 6C 2.100GHz, Infiniband FDR, NVIDIA K20x	27.78
2	3,631.96	Cambridge University	Wilkes - Dell T620 Cluster, Intel Xeon E5-2630v2 6C 2.600GHz, Infiniband FDR, NVIDIA K20	52.62
3	3,517.84	Center for Computational Sciences, University of Tsukuba	HA-PACS TCA - Cray 3623G4-SM Cluster, Intel Xeon E5-2680v2 10C 2.800GHz, Infiniband QDR, NVIDIA K20x	78.77
4	3,185.91	Swiss National Supercomputing Centre (CSCS)	Piz Daint - Cray XC30, Xeon E5-2670 8C 2.600GHz, Aries interconnect, NVIDIA K20x Level 3 measurement data available	1,753.66
5	3,130.95	ROMEO HPC Center - Champagne-Ardenne	romeo - Bull R421-E3 Cluster, Intel Xeon E5-2650v2 8C 2.600GHz, Infiniband FDR, NVIDIA K20x	81.41
6	3,068.71	GSIC Center, Tokyo Institute of Technology	TSUBAME 2.5 - Cluster Platform SL390s G7, Xeon X5670 6C 2.930GHz, Infiniband QDR, NVIDIA K20x	922.54
7	2,702.16	University of Arizona	iDataPlex DX360M4, Intel Xeon E5-2650v2 8C 2.600GHz, Infiniband FDR14, NVIDIA K20x	53.62
8	2,629.10	Max-Planck-Gesellschaft MPI/IPP	iDataPlex DX360M4, Intel Xeon E5-2680v2 10C 2.800GHz, Infiniband, NVIDIA K20x	269.94
9	2,629.10	Financial Institution	iDataPlex DX360M4, Intel Xeon E5-2680v2 10C 2.800GHz, Infiniband, NVIDIA K20x	55.62
10	2,358.69	CSIRO	CSIRO GPU Cluster - Nitro G16 3GPU, Xeon E5-2650 8C 2.000GHz, Infiniband FDR, Nvidia K20m	71.01

- ▶ Greenest supercomputer = 4.5GFlops/W
- ▶ Fastest supercomputer = 1.8GFlops/W
- ▶ Accelerators (GPU, Xeon Phi) are becoming more and more common.

The DEEP Exascale project

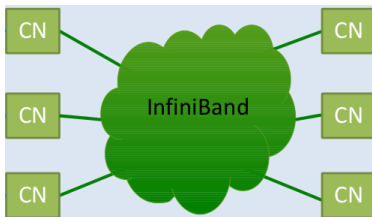
- DEEP: Dynamical Exascale Entry Platform
- one of the three **Exascale** projects funded by the EU: **DEEP, CRESTA and MONTBLANC**.
- It involves 16 partners from 8 different countries and is coordinated by the **Jülich Supercomputing Centre**.
- The project is a two-fold approach to the exascale challenge:

Hardware a novel supercomputing architecture: instead of adding accelerator cards to cluster nodes, an accelerator cluster, called Booster, will complement a conventional HPC system and increase its performance.

Software a matching software stack that includes programming models, libraries and performance tools adapted to its architecture.

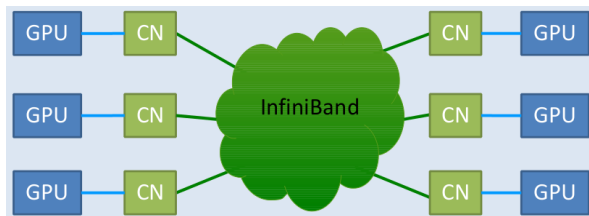
- It will enable unprecedented scalability.

DEEP's architecture



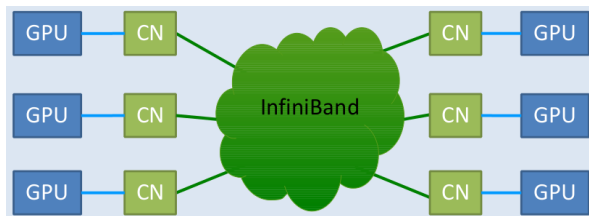
Today's Clusters: CPU nodes connected by Infiniband

DEEP's architecture



Accelerators **statically** attached to CPU

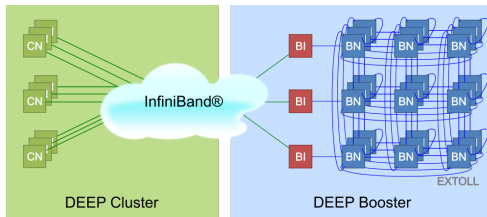
DEEP's architecture



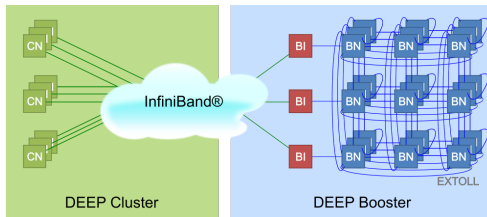
Accelerators **statically** attached to CPU

- Ideally: accelerator cluster + CPU cluster
- Problem: accelerators cannot run autonomously

DEEP's architecture



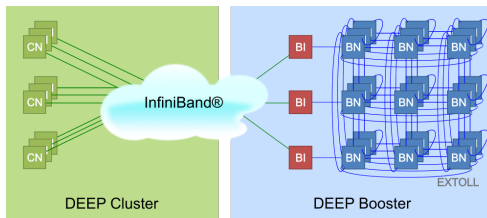
DEEP's architecture



Booster: accelerator cluster

- Intel MIC processors
- EXTOLL network – developed at University of Heidelberg

DEEP's architecture



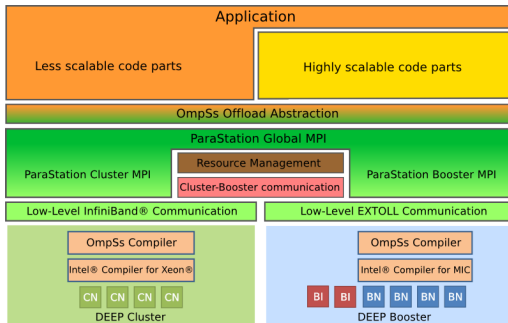
Booster: accelerator cluster

- Intel MIC processors
- EXTOLL network – developed at University of Heidelberg

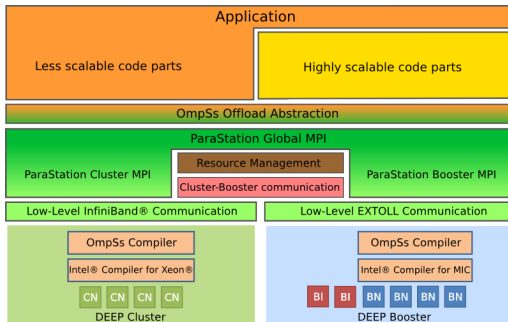
Cluster

- Intel Xeon processors
- Infiniband connect by Mellanox

DEEP's software

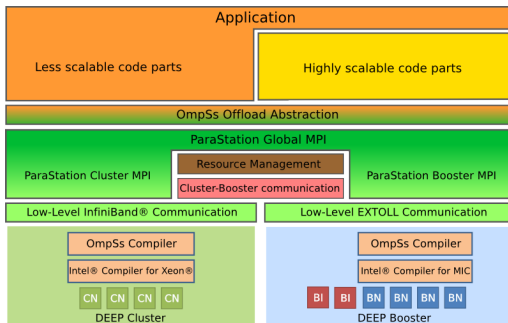


DEEP's software



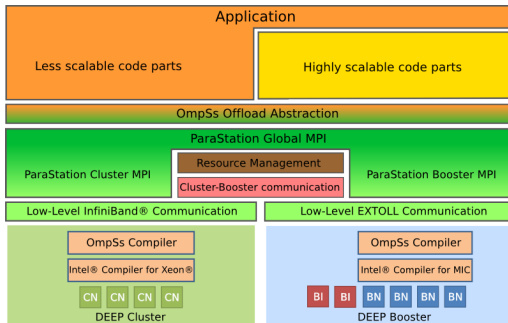
OmpSs developed by the Barcelona Supercomputing Center, allows to decompose applications into tasks sent to the Cluster or the Booster efficiently.

DEEP's software



Parastation MPI will allow to run traditional applications both on the Cluster and the Booster, for this Parastation MPI will be extended to support the Booster and its 3D torus topology.

DEEP's software



All software will be completely reworked in order to be optimized for DEEP.

System software

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

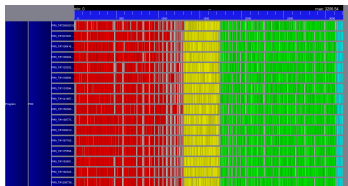
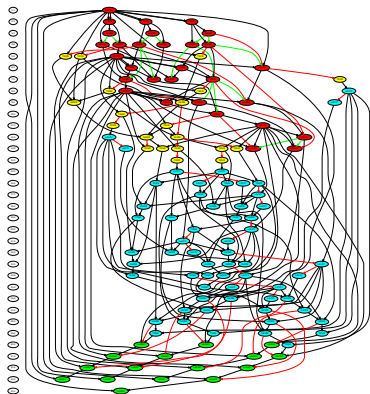
The Deep
Project

Programming
and Application
Challenges

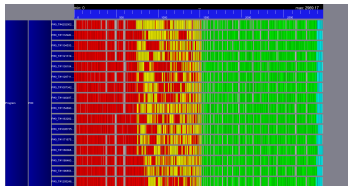
There Goes the
Neighborhood

Neighborhood
on Large
Systems

Classical MPI programs have static load balancing and synchronizations. Hence, they exhibit load imbalance when scale increases.



(a) With synchronization between each stage.



(b) With interleaved stages.

Figure 5: Execution traces of the `DGEMV` routine with a 5000-by-5000 matrix and $NB = 250$ on a 16-cores architecture.

“High Performance Matrix Inversion Based on LU Factorization for Multicore Architectures”.

Dongarra, Faverge, Ltaief, Luszczek. 4th Workshop on Many-Task Computing on Grids and

Failure management.

Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems

As the size of new supercomputers scales to tens of thousands of sockets, the mean time between failures (MTBF) is decreasing to just several hours and long executions need some kind of fault tolerance method to survive failures → a lot of attention on failure management.

Hype and trends

A. Legrand

- Virtualization
 - How Virtualization Changed the Grid Perspective
 - There Goes the Neighborhood
- Toward Exascale
 - The Mont-Blanc Project
 - The Deep Project
 - Programming and Application Challenges
 - There Goes the Neighborhood
 - Neighborhood on Large Systems

Performance tip!



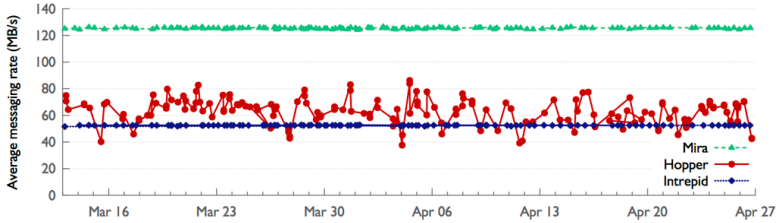
Cray machines



IBM machines

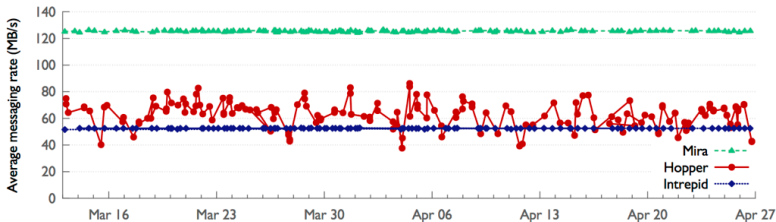
Performance variability

Average messaging rates for batch jobs running a laser-plasma interaction code



Performance variability

Average messaging rates for batch jobs running a laser-plasma interaction code



Total number of bytes sent on the network
Time spent sending the messages



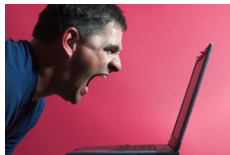
Leads to several problems ...

- Individual jobs run slower:
 - More time to complete science simulations
 - Increased wait time in job queues
 - Inefficient use of machine time allocation/core-hours
- Overall lower throughput
- Increase energy usage/costs



Also affects software development

- Debugging performance issues
- Quantifying the effect of various software changes on performance
 - code changes
 - compiler/software stack changes
- Requesting time for a batch job
- Writing allocation proposals



Setup: Machines

- Hopper: a Cray XE6 at LBNL
 - 2.1 GHz Optrons, 1.28 Petaflop/s
 - 3D Torus, 4 X, Z and 2 Y links, 9.4 GB/s
- Intrepid: an IBM Blue Gene/P at ANL
 - 0.85 GHz PowerPC 450, 0.56 Petaflop/s
 - 3D Torus, 0.425 GB/s
- Mira: an IBM Blue Gene/Q at ANL
 - 1.6 GHz PowerPC A2, 10 Petaflop/s
 - 5D Torus, 2 GB/s



Hype and trends

A. Legrand

Virtualization

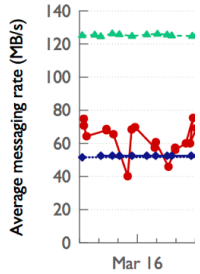
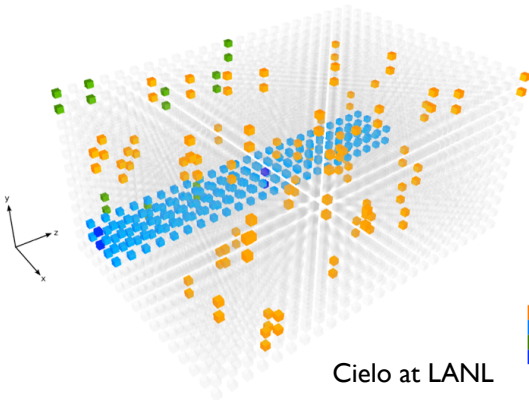
How Virtualization Changed the Grid Perspective
There Goes the Neighborhood

Toward Exascale

The Mont-Blanc Project
The Deep Project
Programming and Application Challenges

There Goes the Neighborhood
Neighborhood on Large Systems

Focus on Cray XE

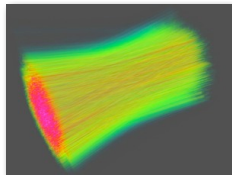


Cielo at LANL

- DVS
- vis
- login/NFS
- vis login

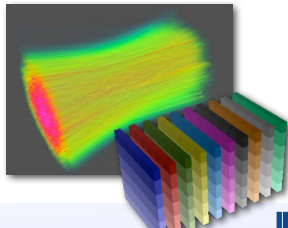
Setup: Application

- pF3D: a highly scalable communication-heavy code
 - used to study laser beam and plasma interactions
- Balanced computation and communication across MPI processes
- 3D virtual process grid
 - 1D FFTs in X and Y direction
 - Send-receives in Z direction



Setup: Application

- pF3D: a highly scalable communication-heavy code
 - used to study laser beam and plasma interactions
- **Balanced computation and communication across MPI processes**
- **3D virtual process grid**
 - ID FFTs in X and Y direction
 - Send-receives in Z direction



Data collection for the paper

- One or more runs on each machine every day: 512 nodes
- Information collected:
 - pF3D stats: messaging rate, time spent in different phases
 - queue status: running jobs and their placement
 - mpiP profiles: time spent in MPI operations



Data collection for the paper

- One or more runs on each machine every day: 512 nodes
- Information collected:
 - pF3D stats: messaging rate, time spent in different phases
 - queue status: running jobs and their placement
 - mpiP profiles: time spent in MPI operations

Machine	Run No.	No. of nodes	No. of cores	No. of jobs	Period		Process Topology	Domain $n_x \times n_y \times n_z$	(x,y) FFT msg. (kB)	Adv. msg. (kB)
					From	To				
Hopper		512	8,192	153	Mar, 2013	Apr, 2013	$32 \times 16 \times 16$	$128 \times 128 \times 8$	4, 8	384
Intrepid		512	2,048	102	Mar, 2013	Apr, 2013	$32 \times 16 \times 4$	$128 \times 128 \times 8$	4, 8	384
Mira		512	8,192	116	Mar, 2013	Apr, 2013	$32 \times 16 \times 16$	$128 \times 128 \times 8$	4, 8	384



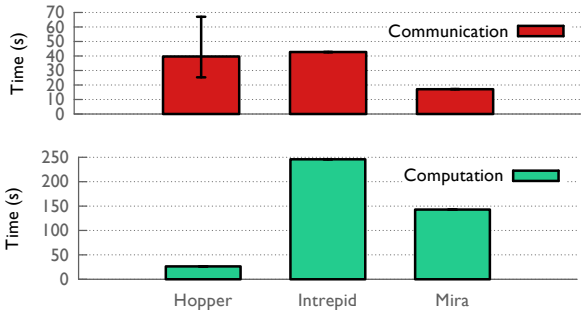
Hype and trends

A. Legrand

- Virtualization
- How Virtualization Changed the Grid Perspective
- There Goes the Neighborhood
- Toward Exascale
- The Mont-Blanc Project
- The Deep Project
- Programming and Application Challenges
- There Goes the Neighborhood
- Neighborhood on Large Systems

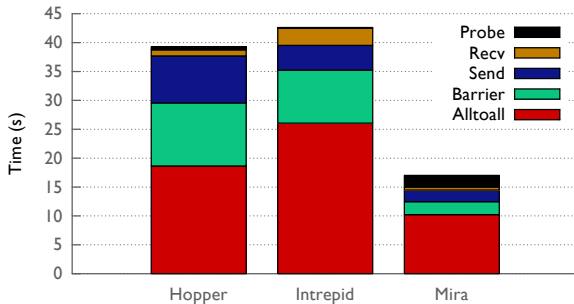
pF3D characterization

Time spent in communication and computation in pF3D



pF3D characterization

Time spent in MPI calls on 512 nodes



Hype and trends

A. Legrand

Virtualization

How
Virtualization
Changed the
Grid Perspective

There Goes the
Neighborhood

Toward Exascale

The Mont-Blanc
Project

The Deep
Project

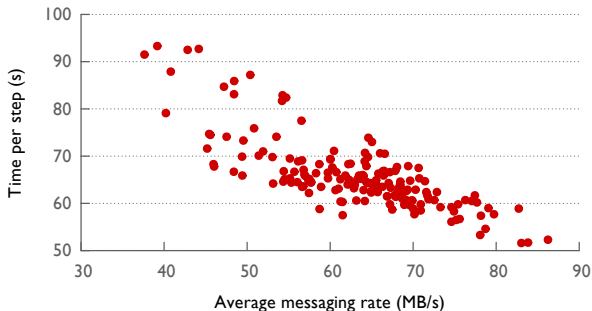
Programming
and Application
Challenges

There Goes the
Neighborhood

Neighborhood
on Large
Systems

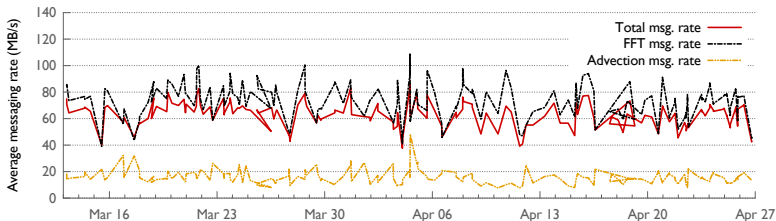
Communication in pF3D

Time versus messaging rate on Hopper



Communication in pF3D

Overall, FFT and advection messaging rates for pF3D on Hopper



Sources of variability

- Operating system noise (OS jitter)
 - OS daemons running on some cores of each node
- Placement/location of the allocated nodes for the job (Allocation shape)
- Contention for shared resources (Inter-job contention)
 - Sharing network links with other jobs



Hype and trends

A. Legrand

Virtualization

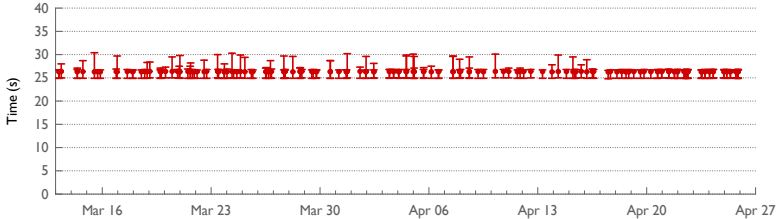
How Virtualization Changed the Grid Perspective
There Goes the Neighborhood

Toward Exascale

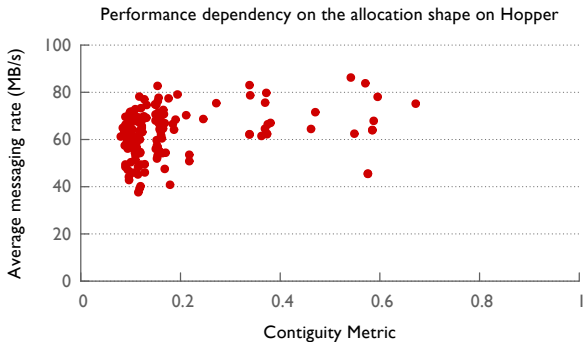
The Mont-Blanc Project
The Deep Project
Programming and Application Challenges
There Goes the Neighborhood
Neighborhood on Large Systems

OS jitter

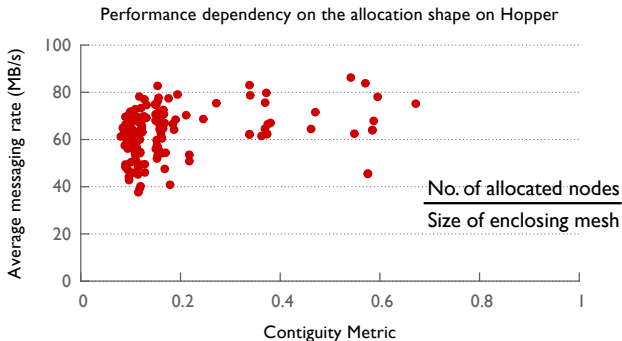
Variation in computation time within a job on Hopper



Degree of fragmentation



Degree of fragmentation



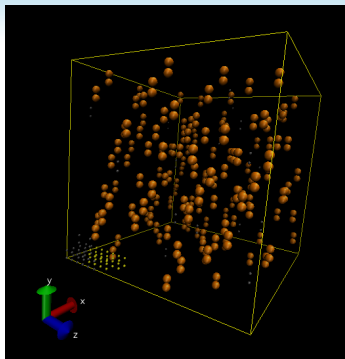


Overview

- What is running on Blue Waters?
- What are the issues and what to do about them?
 - Scalability
 - Runtime consistency
 - Other job interference
 - IO
 - Congestion Protection
 - Interrupts

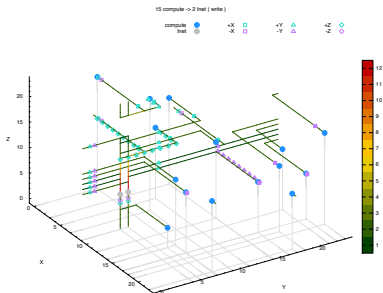
IO

- LNETs scattered across the torus (orange colored geminis).
- Specific OSTs served by specific LNETs (not a full fat tree for the IB between OSTs and LNETs).
- IO is “topology sensitive”.



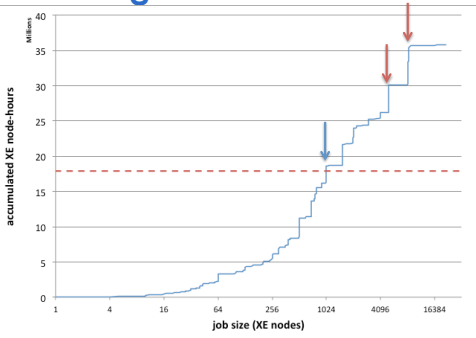


Routing of IO write



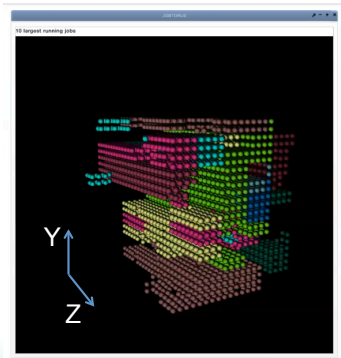
- 15 compute geminis (●) (30 nodes) writing to files served by a LNET pair (●).
- Color scale is the number of convergent routes on the link.

XE Usage in the last 3 months



- 50% of usage is 1,024 nodes or larger.
- Two teams using 5,000 and 8,192 nodes.
- During Friendly User period, several teams sustained runs at full system.
- Nothing prevents users from submitting very large jobs and priority goes to larger jobs.
- Average expansion factor for large jobs < 10.

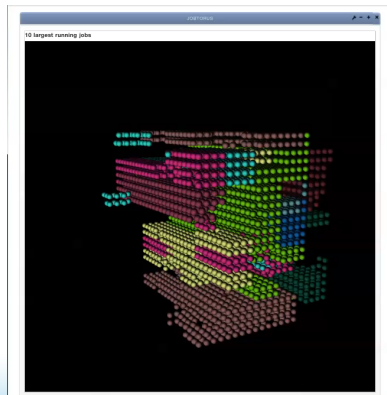
TorusView of 10 largest running jobs



- Relatively compact allocations.
- Some scattered clustering.
- Lots of concave shapes.
- Not showing all the small jobs filling in the rest of the torus.

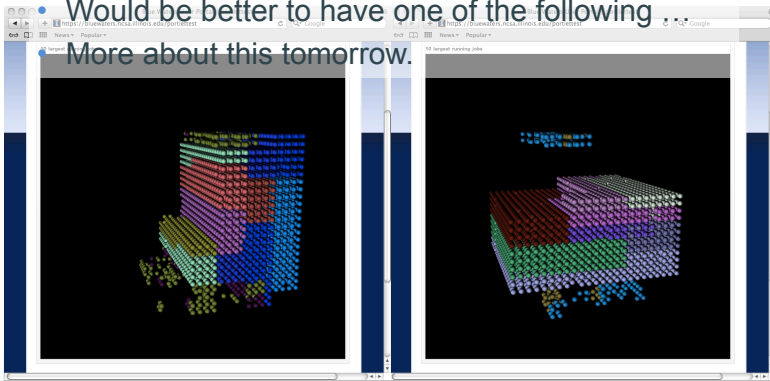
TorusView of 10 largest running jobs

- Allocations shift planes as the end of the Z direction is hit.
- Voids where larger job allocations wrap around smaller ones.



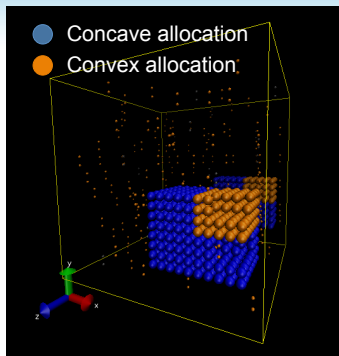
Better nid allocation

- Would be better to have one of the following ...
- More about this tomorrow.



Impact of nid allocation

- Job – Job interaction
 - Analysis of key application communication intensity and sensitivity
 - 20% slowdown typical, 100% or more possible.

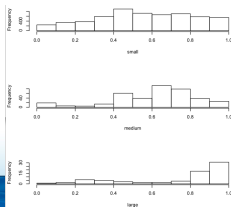
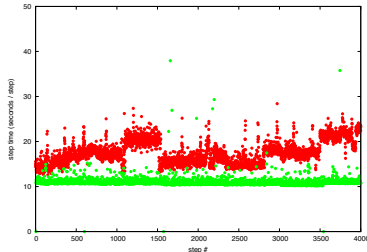


Communication	MILC	NAMD	NWCHEM	PSDNS	WRF
Intensive	2	2	3	2	1
Sensitive	2	3	1	2	1

1 – low 3 – high
as viewed by convex app.

Impact of poor nid allocation - Consistency

- Two jobs (8,192 nodes) with nearly same nid allocation (s10_8972n). Red job affected by other workload communicating through the region.



- Run time variation - poor wallclock accuracy (padding wallclock).

Congestion Protection

- To avoid data loss, traffic injection is throttled for a period of time, when reaching a point where forward progress is stalling. Throttling is applied and removed until congestion is cleared.
- System monitors percentage of time that traffic trying to enter the network from the nodes and percentage of time network tiles are stalled.
- Fortunately not a common occurrence. It does happen, typically in bursts.
- Can happen with node-node (MPI, PGAS) or node-LNET (IO) traffic.
- Many-to-one and long-path patterns.
- Libraries and user can control node injection as a precaution.
- In CP reports, flit rates represent data arriving at the node from the interconnection network.

```

Node
----
2220440 CentOS6.Linux.      2048    21496    46466    16100:45    19:41:40
2220442 CentOS6.Linux.      2048    8115     46466    16100:45    18:37:02
2219384 nsm2                2000    --       43448    01:58:31    18:02:09
2220893 psolve              2000    45732    4752    17:12:34    17:36:30
2219759 wal_tmd_biq_g       1536    --       12648    07:19:16
2219818 nuchan              1000    --       3745    13:08:50    18:02:07
2220448 nuchan              1000    4328749    32745    17:00:22    18:15:32
2219439 wa_aperture_nia     768    --       12648    13:20:04
2219512 nsm2                700    --       43964    18:35:55
--
-----
Top bandwidth Applications
-----
#
#  apid 2219386 usercid 43448 numtile 2000 agname          nsm2 flits/sec Total 3075
#  apid 2219439 usercid 32745 numtile 1000 agname          nuchan flits/sec Total 2742
#  apid 2220442 usercid 46466 numtile 2048 agname          CentOS6.Linux. flits/sec Total 3715
#  apid 2220440 usercid 46466 numtile 2048 agname          CentOS6.Linux. flits/sec Total 2491
#  apid 2219117 usercid 43864 numtile 700 agname          nsm2 flits/sec Total 2371
#  apid 2219119 usercid 42864 numtile 700 agname          nsm2 flits/sec Total 2075
#  apid 2219759 usercid 12960 numtile 1536 agname          wal_tmd_biq_g flits/sec Total 2071
#  apid 2219514 usercid 43864 numtile 700 agname          nsm2 flits/sec Total 1762
#  apid 2220444 usercid 12960 numtile 512 agname          wa_aperture_nia flits/sec Total 1266
#  apid 2217219 usercid 47396 numtile 500 agname          python flits/sec Total 1269
--
-----
Congestion Candidate COMPUTE nodes
-----
#
#  c0b-10c1a2d1 404851 flits/sec (pid 18481); apid 2220477 usercid 14394 numtile 32 agname nsm_script.sh
#  c0b-10c1a2d0 401956 flits/sec (pid 23034); apid 2219894 usercid 14394 numtile 32 agname nsm_script.sh
#  c0b-10c1a2d3 28418 flits/sec (pid 5798); apid 2219756 usercid 14394 numtile 32 agname nsm_script.sh
#  c0b-10c1a2d1 28218 flits/sec (pid 25847); apid 2219672 usercid 15077 numtile 44 agname nsm_saw
#  c0b-10c1a2d1 22254 flits/sec (pid 8024); apid 2219756 usercid 14394 numtile 32 agname nsm_script.sh
#  c0b-10c1a2d0 20193 flits/sec (pid 24813); apid 2219672 usercid 15077 numtile 44 agname nsm_saw
#  c0b-10c1a2d0 19141 flits/sec (pid 8094); apid 2219756 usercid 14394 numtile 32 agname nsm_script.sh
#  c0b-10c1a2d0 18794 flits/sec (pid 8129); apid 2219756 usercid 14394 numtile 32 agname nsm_script.sh
#  c0b-10c1a2d1 18273 flits/sec (pid 5818); apid 2219756 usercid 14394 numtile 32 agname nsm_script.sh
#  c0b-10c1a2d0 17453 flits/sec (pid 5814); apid 2219756 usercid 14394 numtile 32 agname nsm_script.sh
-----
Top 100 Congestion Candidate Nodes (414 compute nodes: 134930785 flits/s, 590 service nodes: 1257373796 flits/s)
-----
#
#  c0b-10c1a2d0 4128749 flits/sec nid 12810; apid 2220440 usercid 32745 numtile 1000 agname nuchan
#  c0b-10c1a2d3 3386708 flits/sec nid 12810; apid 2220440 usercid 32745 numtile 1000 agname nuchan
#  c0b-10c1a2d1 3051520 flits/sec nid 15484; apid 2220440 usercid 32745 numtile 1000 agname nuchan
#  c0b-10c1a2d0 2633807 flits/sec nid 17984; apid 2220440 usercid 32745 numtile 1000 agname nuchan
#  c0b-10c1a2d1 2612123 flits/sec nid 15484; apid 2220440 usercid 32745 numtile 1000 agname nuchan
#  c0b-10c1a2d3 2779003 flits/sec nid 12847; apid 2220440 usercid 32745 numtile 1000 agname nuchan
#  c0b-10c1a2d0 2737704 flits/sec nid 13044; apid 2220440 usercid 32745 numtile 1000 agname nuchan
#  c0b-10c1a2d0 2623978 flits/sec nid 15336; apid 2220440 usercid 32745 numtile 1000 agname nuchan
#  c0b-10c1a2d0 2619900 flits/sec nid 19030; apid 2220440 usercid 32745 numtile 1000 agname nuchan
#  c0b-10c1a2d3 2648270 flits/sec nid 15489; apid 2220440 usercid 32745 numtile 1000 agname nuchan

```

Hype and trends

A. Legrand

Virtualization

How Virtualization Changed the Grid Perspective

There Goes the Neighborhood

Toward Exascale

The Mont-Blanc Project

The Deep Project

Programming and Application Challenges

There Goes the Neighborhood

Neighborhood on Large Systems

Congestion Protection Analysis

- Look at application to node relation.
- wrf listed as top application and the top 10 nodes are wrf nodes.
- nwchem running at same time (listed #4).
- The OVIS state of the network data should help here.

