### HPC Virtualization: Control Your Stack

Rick Wagner HPC Systems Manager rpwagner@sdsc.edu 4<sup>th</sup> Annual MVAPICH User Group Meeting 2016





### *comet.sdsc.edu* 32 racks of awesomeness









## **Project High Level Goals**

- "... expand the use of high end resources to a much larger and more diverse community
- ... support the entire spectrum of NSF communities
- ... promote a more comprehensive and balanced portfolio
- ... include research communities that are not users of traditional HPC systems."

NSF solicitation 13-528

### The long tail of science needs HPC







# Comet's integrated architecture is a platform for a wide range of computing modalities



## **Virtualization Staffing**

SDSC: Project management, system managements, systems software (Nucleus)



IU: User support, client software (Cloudmesh)







## **Virtual Clusters**

### Goal:

Provide a near bare metal HPC performance and management experience

#### Target Use

Projects that can manage their own cluster, and:

- can't fit our batch environment, and
  - don't want to buy hardware or
  - have bursty or intermittent need





### Single Root I/O Virtualization in HPC

- Problem: Virtualization generally has resulted in significant I/O performance degradation (e.g., excessive DMA interrupts)
- Solution: SR-IOV and Mellanox ConnectX-3 InfiniBand host channel adapters
  - One physical function → multiple virtual functions, each light weight but with its own DMA streams, memory space, interrupts
  - Allows DMA to bypass hypervisor to VMs
- SRIOV enables virtual HPC cluster w/ near-native InfiniBand latency/bandwidth and minimal overhead









## MPI bandwidth slowdown from SR-IOV is at most 1.21 for medium-sized messages & negligible for small & large ones







## MPI latency slowdown from SR-IOV is at most 1.32 for small messages & negligible for large ones



SDSC SAN DIEGO SUPERCOMPUTER CENTER



### WRF Weather Modeling – 2% Overhead with SR-IOV IB

- 96-core (6-node) calculation
- Nearest-neighbor communication
- Scalable algorithms
- SR-IOV incurs modest (15%)
  performance hit
- 2% slower w/ SR-IOV vs native IB!
- Still 20% faster than EC2 Despite 20% slower CPUs





### Quantum ESPRESSO: 28% 8% (!) Overhead







## Selected Technologies: Enabling Major Impact

- KVM—Let's us run virtual machines (all processor features)
- SR-IOV—Makes MPI go fast on VMs
- Rocks—Systems management
- ZFS—Disk image management
- VLANs—Isolate virtual cluster management network
- pkeys—Isolate virtual cluster IB network
- Nucleus—Coordination engine (scheduling, provisioning, status, etc.)
- Client Cloudmesh





## **User-Customized HPC**



SDSC SAN DIEGO SUPERCOMPUTER CENTER





### Accessing Comets Virtual Cluster Capabilities

- REST API
- Command line interface
- Command shell for scripting
- Console Access
- (Portal)

![](_page_14_Picture_6.jpeg)

_	
Usage:	
comet	ll [CLUSTERID] [format=FORMAT]
comet	cluster [CLUSTERID]
	[format=FORMAT]
comet	computeset [COMPUTESETID]
comet	<pre>power on CLUSTERID [count=NUMNODES] [NODESPARAM]</pre>
	[allocation=ALLOCATION]
	[walltime=WALLTIME]
comet	<pre>power (off reboot reset shutdown) CLUSTERID [NODESPARAM]</pre>
comet	console CLUSTERID [COMPUTENODEID]
comet	image list
comet	<pre>image upload [imagename=IMAGENAME] PATHIMAGEFILE</pre>
comet	<pre>image attach IMAGENAME CLUSTERID [COMPUTENODEID]</pre>
comet	<pre>image detach CLUSTERID [COMPUTENODEID]</pre>
comet	node rename CLUSTERID OLDNAME NEWNAME

(ENV)big:client big\$ cm help

Documented commands (type help <command/> ):									
EOF	comet	group	key	open	refresh	server	var		
banner	context	h	launcher	pause	register	shell	verbose		
check	debug	help	limits	portal	reservation	ssh	version		
clear	default	history	list	ру	reset	submit	vm		
cloud	echo	hpc	man	q	rsync	sync			
cluster	exec	image	network	quit	secgroup	timer			
color	flavor	inventory	nova	quota	select	usage			
		•		-					

User does NOT see: Rocks, Slurm, etc.

![](_page_14_Picture_11.jpeg)

![](_page_14_Picture_12.jpeg)

![](_page_15_Picture_0.jpeg)

![](_page_15_Picture_1.jpeg)

![](_page_15_Picture_2.jpeg)

## Cloudmesh

## Hybrid Cloud Management http://cloudmesh.github.io/client/

![](_page_16_Picture_2.jpeg)

laszewski@gmail.com, http://cloudmesh.github.io/client/

![](_page_16_Picture_4.jpeg)

## **Integrated Clouds**

![](_page_17_Figure_1.jpeg)

![](_page_17_Picture_2.jpeg)

laszewski @gmail.com, http://cloudmesh.github.io/client

![](_page_17_Picture_4.jpeg)

## **Cloudmesh Fills the Gap**

- Matches user needs with multiple provider's services.
- Researchers can use one platform to manage their clouds.
- Orchestrates provisioning and allocation of cloud resources.
- Local copy of your cloud data is created, so jobs and VMs are traceable across clouds
- New clouds with similar configurations can be created easily.
- Default attributes allow easy control of cloud artefacts.
- Users can switch easily between clouds.
- Users can switch easily between HPC systems
- cm default cloud=comet
- cm vm boot
- cm default cloud=chameleon
- cm vm boot

![](_page_18_Picture_13.jpeg)

![](_page_18_Picture_15.jpeg)

## Future: Comet Cloudmesh Platform Launchers

• Customizable launchers

COMPLITER CENTER

• Launchers available through commandline or browser Example: Hadoop

\$cm hadoop -n 10 -group=myHadoop

![](_page_19_Figure_4.jpeg)

![](_page_19_Picture_5.jpeg)

## Early Success: Open Science Grid

![](_page_20_Figure_1.jpeg)

SDSC SAN DIEGO SUPERCOMPUTER CENTER

![](_page_20_Picture_3.jpeg)

![](_page_21_Picture_0.jpeg)

## **Control & Responsibility**

- Shared responsibility
  - SDSC: hosting environment
  - Cluster admin & PI: virtual machine stack
- Different from current HPC roles
- Considering an explicit "Acknowledgement of Responsibility" signed by PI & cluster admin
  - Based on SDSC's "Outback Network" agreement
  - Outback is a separate VLAN & IP subnet for user-managed systems

![](_page_22_Picture_8.jpeg)

![](_page_22_Picture_9.jpeg)

## **Innovation in Utilization**

- VMs co-scheduled within the batch system
- Works for automated workflows bringing up virtual compute nodes when needed
- What about responsiveness?
  - Can tune batch policies based on need
  - Reserve some physical nodes for fast launch to some scale
- Reuse existing allocations & accounting process
- Likewise, existing science gateway policies: e.g., community account for cluster

![](_page_23_Picture_8.jpeg)

![](_page_23_Picture_9.jpeg)

## **New Use Case: Training**

- Training
  - Nearly full stack: OS, networking (IP & InfiniBand), applications
  - Any type of cluster: HPC; N-tier web framework; big data
  - Limits: no layer 1 (VLANs, etc.) networking

![](_page_24_Picture_5.jpeg)

![](_page_24_Picture_6.jpeg)

## **New Use Case: Campus Bursting**

- Custom HPC clusters can help campuses extend a familiar environment
- How can a campus get a large compute allocation?
- Need to justify science and compute time
- Single PI proposals from large users?
- Historical campus cluster utilizations
- Don't want to burn out Campus Champion allocations

![](_page_25_Picture_7.jpeg)

![](_page_25_Picture_8.jpeg)

## **And in Other News**

- Check out Singularity: <a href="http://singularity.lbl.gov/">http://singularity.lbl.gov/</a>
- User space containers for HPC
- Deployed on Comet, Gordon, and TSCC (campus cluster)
- Impromptu BoF Wednesday afternoon (location TBD)
- https://github.com/cjprybol/reproducibility-via-singularity

#### Make your research more reproducible with Singularity

Services like GitHub, Bitbucket, and GitLab have democratized access to affordable (or free!) tools that reinforce reproducibility in research, via the code repositories that these services offer. These services make it easier to backup, version control, collaborate on, and distribute code (or any other text-based file). These features make it easier for researchers to write and maintain high-quality code. These features also increase the chance that someone else will review the code for errors or bugs. These reviews can be done by some combination of reading and executing the code. Unfortunately, unless the code is run on the exact same computer, while logged in as the same user, getting the same code to run the same way can be a research project in and of itself.

![](_page_26_Picture_8.jpeg)

![](_page_26_Picture_9.jpeg)

## Singularity

![](_page_27_Figure_1.jpeg)

Linux Kernel

Physical Hardware Layer

### Singularity & Open MPI

![](_page_27_Picture_5.jpeg)

![](_page_27_Picture_6.jpeg)

https://github.com/singularityware

![](_page_27_Picture_8.jpeg)