



RWTH Aachen, Forschungszentrum Jülich (Germany) 3rd-7th September 2007

CONTENTS

③ OVERVIEW

③ VIRTUALIZATION: CLUSTER OS

③ VIRTUALIZATION: CLUSTER STORAGE

③ FAULT TOLERANCE

③ TOOLS

③ DOES IT WORK?

③ CONCLUSIONS

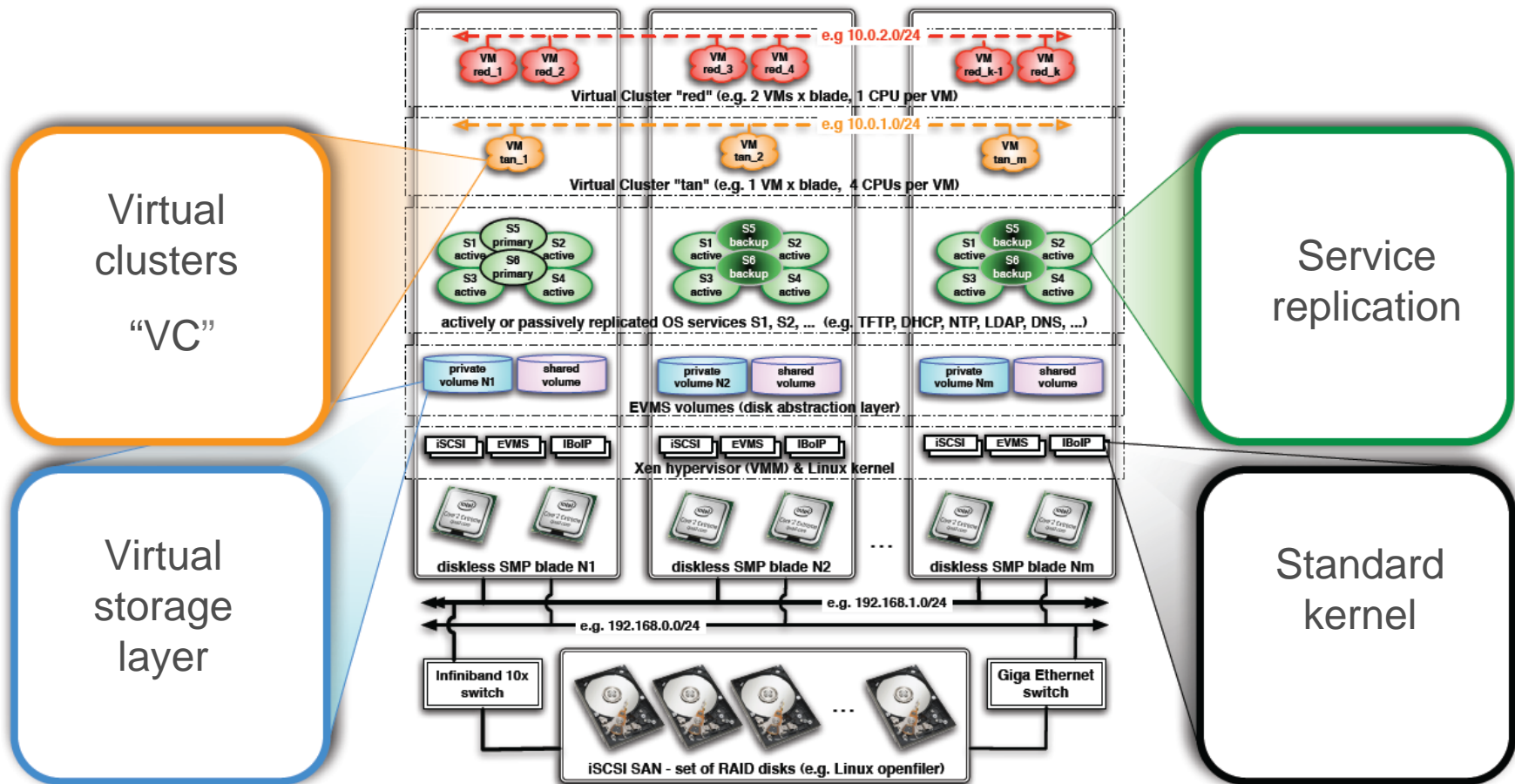


OVERVIEW: COMMON FLAWS

- ③ Cluster configuration is hard to evolve as needs change (reinstallation)
- ③ One cluster → one environment (how to cope with GRIDs, Beowulf, etc?)
- ③ Strong coupling between physical and logical architectures
- ③ Single point of failure (master node)
- ③ Multiple points of failure (local hard disks)
- ③ Resource waste (OS replication for homogeneous nodes)
- ③ Configuration is not reusable (installation procedure must be redone for identical machines)



OVERVIEW: VirtualLinux



... and the tools to manage this



CONTENTS

③ OVERVIEW

③ VIRTUALIZATION: CLUSTER OS

③ VIRTUALIZATION: CLUSTER STORAGE

③ FAULT TOLERANCE

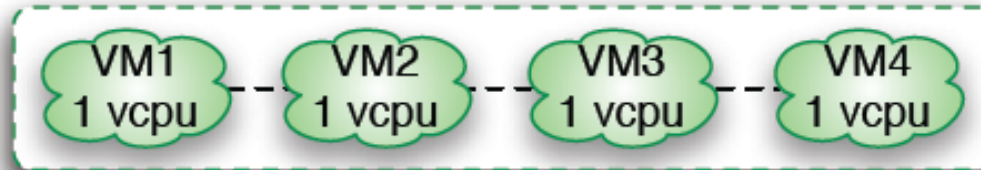
③ TOOLS

③ DOES IT WORK?

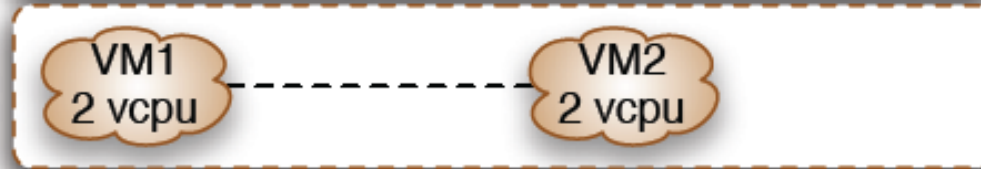
③ CONCLUSIONS



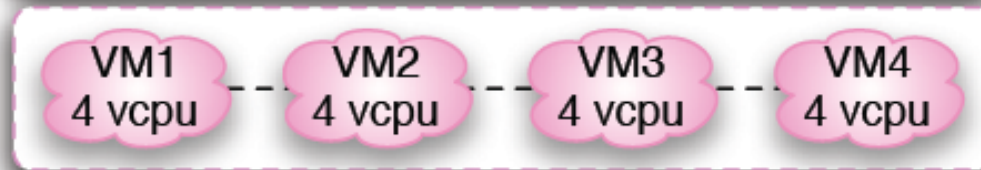
VIRTUALIZATION: CLUSTER OS



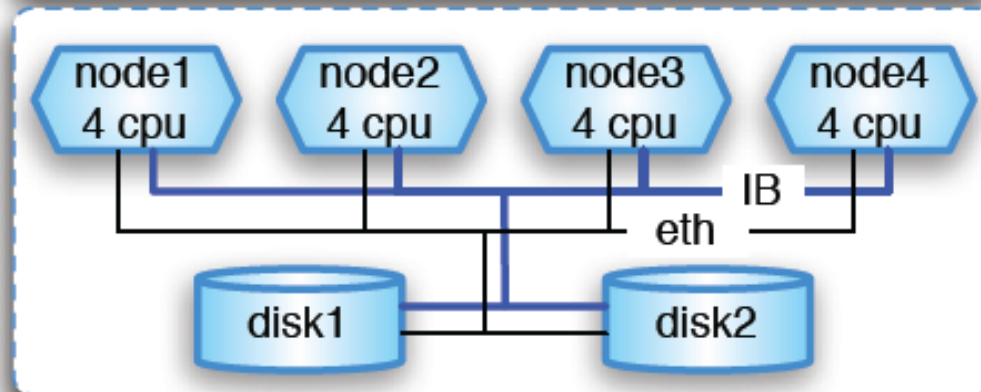
Virtual Cluster "green"
4 VMs x 1 VCPUs
10.0.3.0/24



Virtual Cluster "tan"
2 VMs x 2 VCPUs
10.0.1.0/24



Virtual Cluster "pink"
4VMs x 4VCPUs
10.0.0.0/24



Physical Cluster + external SAN
InfiniBand + Ethernet
4 Nodes x 4 CPUs
Cluster InfiniBand 192.0.0.0/24
Cluster Ethernet 192.0.1.0/24
Internet Gateway 131.1.7.6

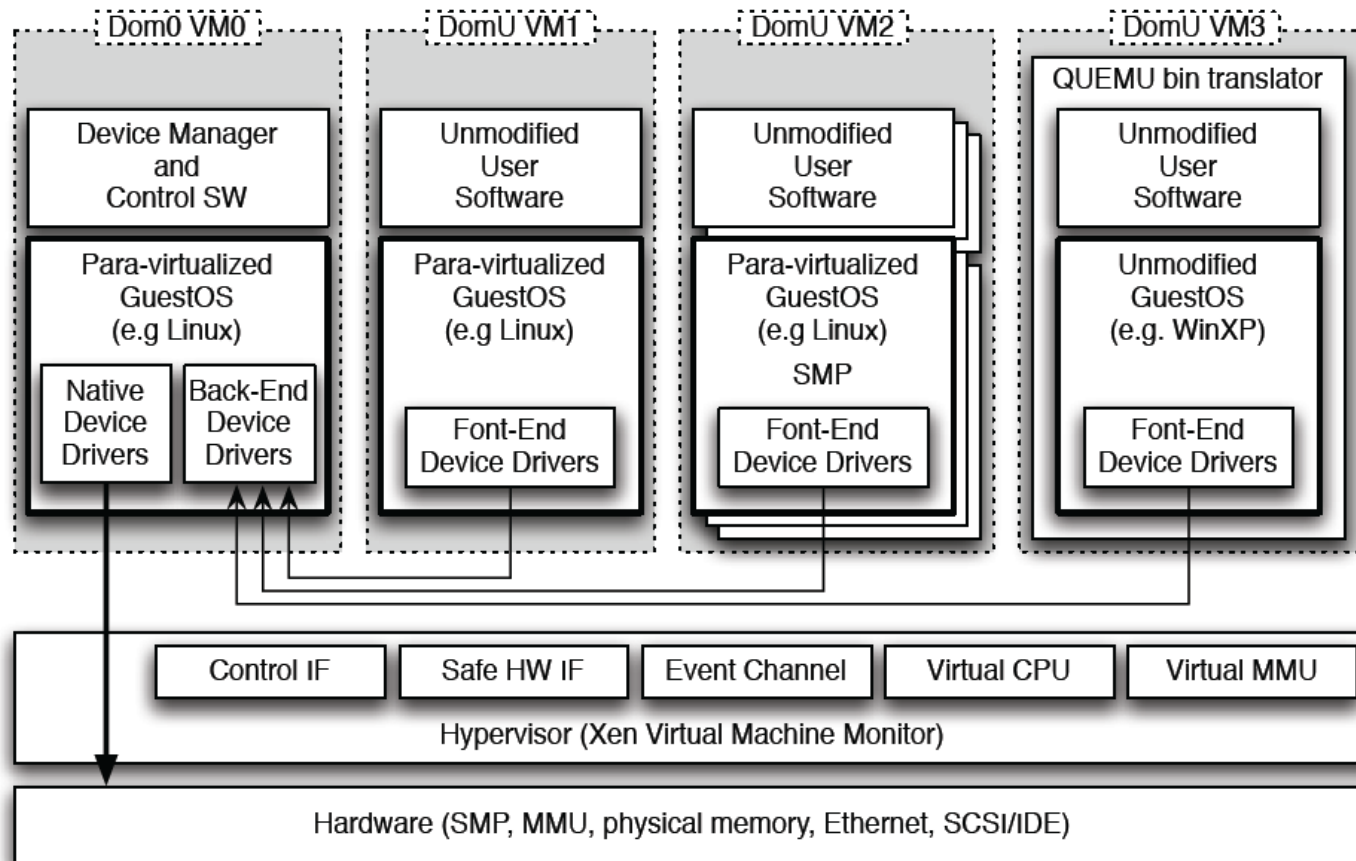


VIRTUALIZATION: CLUSTER OS

- ③ Three layers: physical, host, guest clusters
- ③ The host cluster layer hides and protects hardware details
- ③ Guest clusters are clusters of Virtual Machines (VM) that are totally insulated (Virtual Clusters, VC)
- ③ Key features of VCs:
 - ③ Crashes and instabilities are confined
 - ③ VCs can have different configurations and OS (Beowulf, Grid, RH, Ubuntu, Windows).
 - ③ Guest administrators don't have to be very skilled (limited damage)
 - ③ VC and VMs can be preconfigured and downloaded from repositories
 - ③ VMs within a VC can be moved between physical nodes
 - ③ Possibly, superscalability (I/O bound and CPU bound VCs running on the same physical nodes)



VIRTUALIZATION: CLUSTER OS



A physical node running 3 VMs (2 to 3 VCs nodes. One runs Windows through QEMU). Xen provides the virtualization

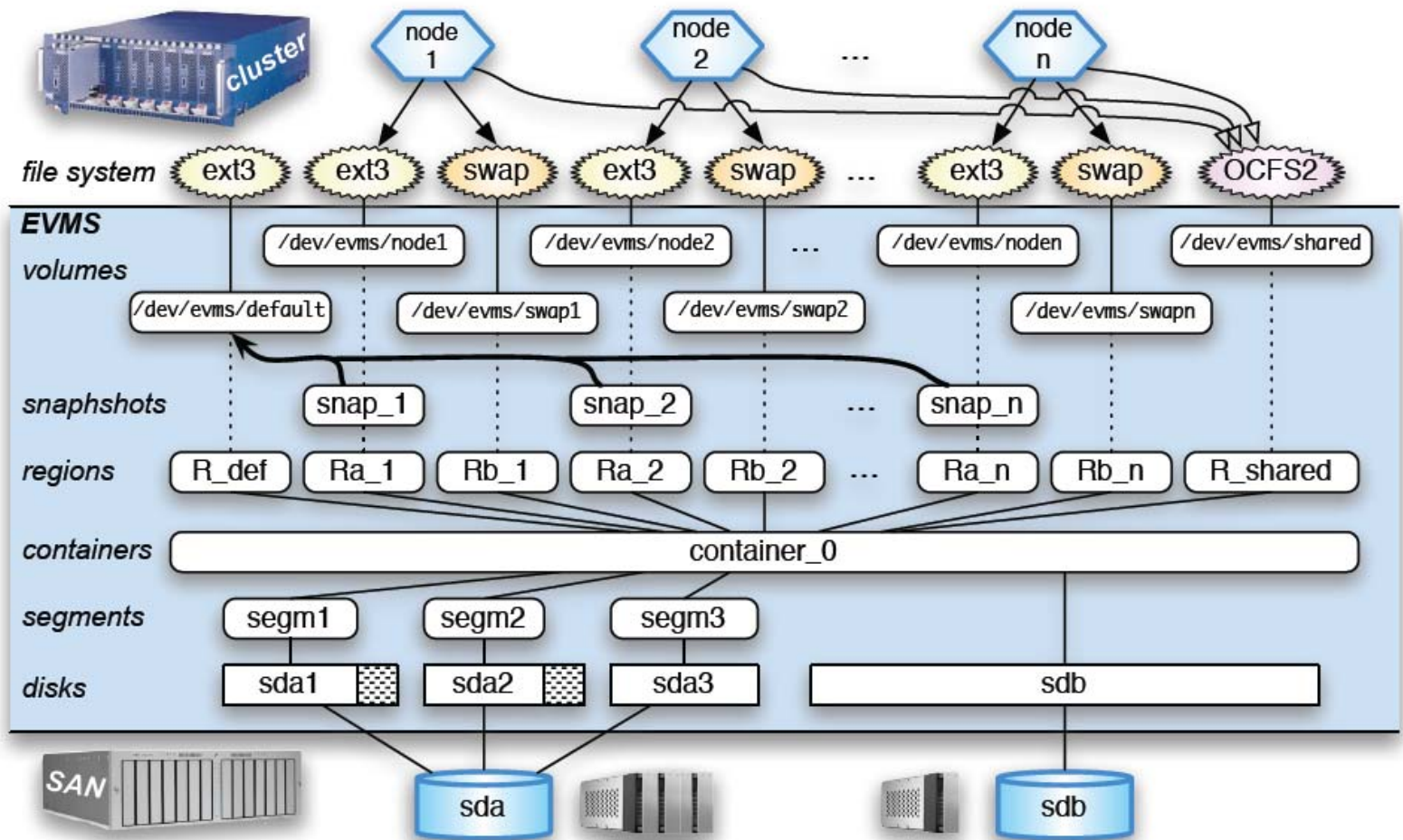


CONTENTS

- ③ OVERVIEW
- ③ VIRTUALIZATION: CLUSTER OS
- ③ VIRTUALIZATION: CLUSTER STORAGE
- ③ FAULT TOLERANCE
- ③ TOOLS
- ③ DOES IT WORK?
- ③ CONCLUSIONS



VIRTUALIZATION: CLUSTER STORAGE



VirtualLinux, EVMS and iSCSI

VIRTUALIZATION: CLUSTER STORAGE

- ③ Storage virtualization in VirtualLinux means EVMS+creative use of snapshots on an iSCSI SAN
- ③ iSCSI provides standard, simple, cheap and cluster wide access to the remote, redundant disks, allowing diskless nodes
- ③ iSCSI can reuse existing infrastructures (GigEth, Infiniband). No additional cost, cabling, adapters or drivers
- ③ EVMS insulates the physical details
- ③ EVMS provides a single, unified system for storage management
- ③ EVMS logical volume naming avoids the use of raw device names (iSCSI devices connected to different nodes would have different names)

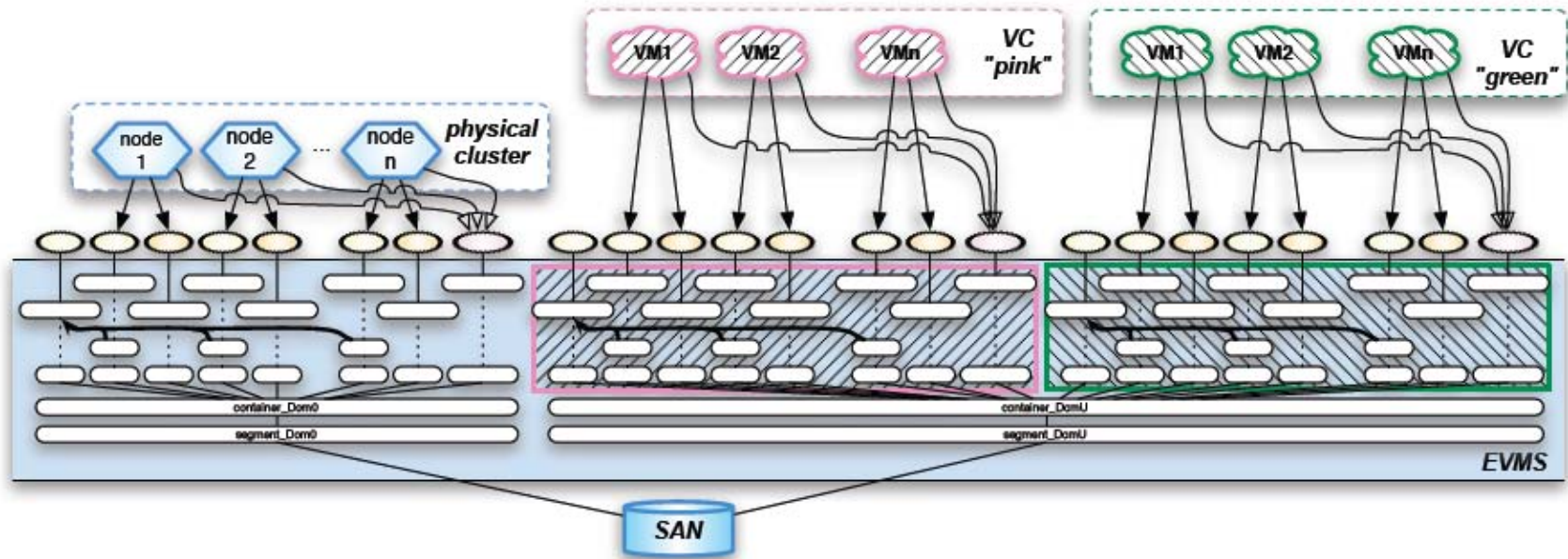


VIRTUALIZATION: CLUSTER STORAGE

- ③ Generally, all VC virtual nodes share the same OS image, with minor differences
- ③ 1000 virtual nodes -> 1000 replicas. Disk and time waste
- ③ The same is true for classic clusters
- ③ Snapshot is traditionally used to backup volumes, **but can be used to create replicas (volumes)**
- ③ Modified snapshot technique:
 - ③ Create the volume and install the OS image
 - ③ Make N snapshots of the original (for N nodes)
 - ③ The snapshots are created using ONLY metadata and pointers to the original copy
 - ③ Creation time of N snapshots almost independent from N (seconds)
 - ③ Modifications are lightweight: a modification to a file just involves the copy and update of the single file
- ③ Memory buffers in kernel limit the N. VirtualLinux extends EVMS semantics (one line of C code) to allow any N



VIRTUALIZATION: CLUSTER STORAGE



As for the physical cluster, each VC node has a VC private virtual shared storage (OCFS2 distributed FS), a private virtual disk and a private virtual swap area. Private virtual disks are functionally identical to traditional local disks. Virtual means that the volume is a snapshot.



CONTENTS

- ③ OVERVIEW
- ③ VIRTUALIZATION: CLUSTER OS
- ③ VIRTUALIZATION: CLUSTER STORAGE
- ③ FAULT TOLERANCE
- ③ TOOLS
- ③ DOES IT WORK?
- ③ CONCLUSIONS



FAULT TOLERANCE

- ③ Traditionally, clusters have a master node. Great risk!
- ③ VirtuaLinux is a masterless architecture
- ③ Active and passive replication (primary-backup) of services
- ③ Service classification:
 - ③ Stateless: Active replication (e.g IB manager). Clients automatically locate the server (broadcasting, etc). Most reactive server replies
 - ③ Stateful: Passive replication (IP Gateway). Only one server exists and is statically known to clients (IP, MAC). Heartbeat with IP takeover between primary and backup servers
 - ③ Node oriented: failure of node specific services is potentially catastrophic for the node only (local SSH, NFS). The rest of the cluster is unaffected. Outside the scope of VirtuaLinux. Scheduler level?
 - ③ Self-healing: fault tolerance is embedded in the service (NTP)



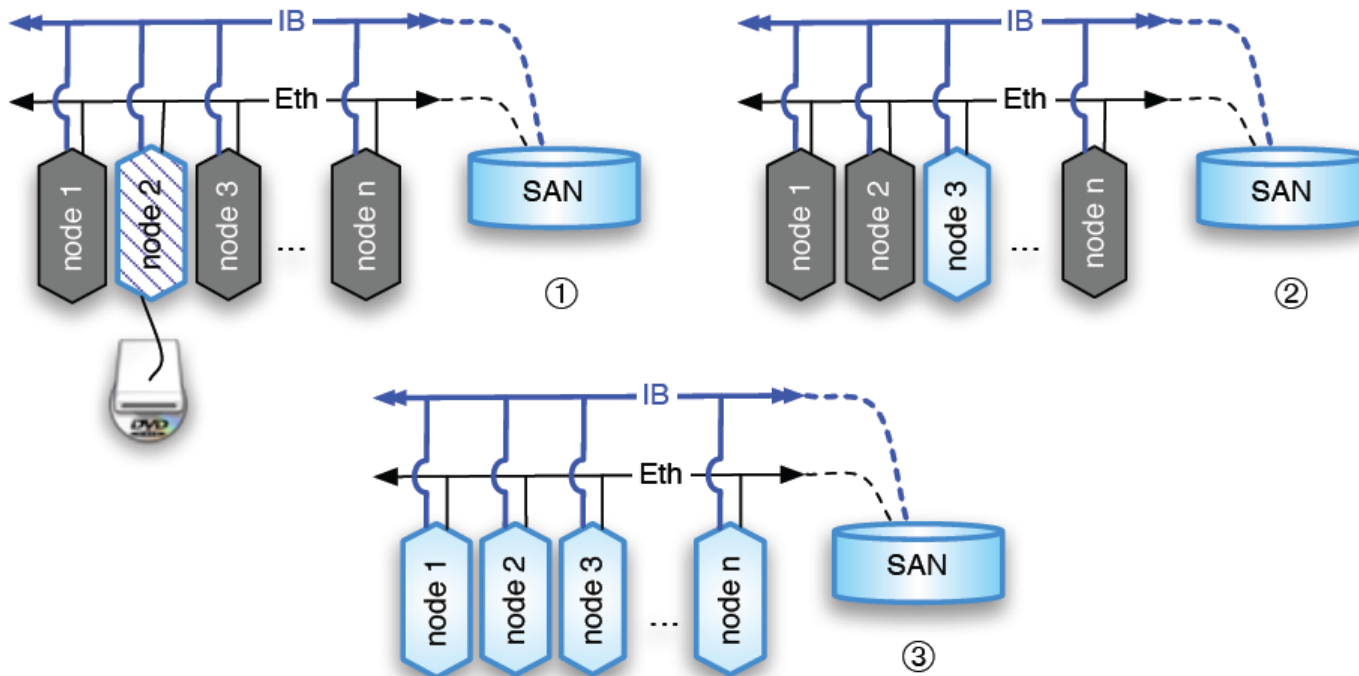
FAULT TOLERANCE

Service	Fault-tolerance	Notes
DHCP	active	Pre-defined map between IP and MAC
TFTP	active	all copies provide the same image
NTP	active	Pre-defined external NTPD fallback via GW
IB manager	active	Stateless service
DNS	active	Cache only DSN
LDAP	service-specific	Service-specific master redundancy
IP GW	passive	Heartbeat on 2 nodes with IP takeover (HA)
Mail	node-oriented	Local node and relays via DNS
SSH/SCP	node-oriented	Pre-defined keys
NFS	node-oriented	Pre-defined configuration
SMB/CIFS	node-oriented	Pre-defined configuration



FAULT TOLERANCE

- ③ How to install a masterless cluster? Chicken and egg problem!
- ③ Solution: Metamaster! A temporary installation node that is then transformed in a standard node
- ③ At the end of the installation, all nodes are identical and provide redundant services. Bye bye master node



CONTENTS

- ③ OVERVIEW
- ③ VIRTUALIZATION: CLUSTER OS
- ③ VIRTUALIZATION: CLUSTER STORAGE
- ③ FAULT TOLERANCE
- ③ TOOLS
- ③ DOES IT WORK?
- ③ CONCLUSIONS



TOOLS

- ③ **VirtualLinux is a set of scripts, independent of the Linux Distribution**
- ③ **VirtualLinux Virtual Cluster Management (VVCM) is a subset of scripts that allow the creation and management of VCs**
- ③ **Main components are:**
 - ③ **Database of the physical and virtual clusters (includes the mapping between physical and virtual clusters, virtual nodes, etc)**
 - ③ **Command line library for the creation, activation and destruction of a VC (VC_Create, VC_Control, VC_Destroy)**
 - ③ **Communication layer (for staging and executing VMs)**
 - ③ **VC start time support for dynamic configuration of network topology and routing policies of the physical nodes**



DOES IT WORK?

Micro-benchmark	Unit	Ub-Dom0	Ub-DomU	CentOS
Simple syscall	usec	0.6305	0.6789	0.0822
Simple open/close	usec	5.0326	4.9424	3.7018
Select on 500 tcp fd's	usec	37.0604	37.0811	75.5373
Signal handler overhead	usec	2.5141	2.6822	1.1841
Protection fault	usec	1.0880	1.2352	0.3145
Pipe latency	usec	20.5622	12.5365	9.5663
Process fork+execve	usec	1211.4000	1092.2000	498.6364
float mul	nsec	1.8400	1.8400	1.8200
float div	nsec	8.0200	8.0300	9.6100
double mul	nsec	1.8400	1.8400	1.8300
double div	nsec	9.8800	9.8800	11.3300
RPC/udp latency localhost	usec	43.5155	29.9752	32.1614
RPC/tcp latency localhost	usec	55.0066	38.7324	40.8672
TCP/IP conn. to localhost	usec	73.7297	57.5417	55.9775
Pipe bandwidth	MB/s	592.3300	1448.7300	956.21

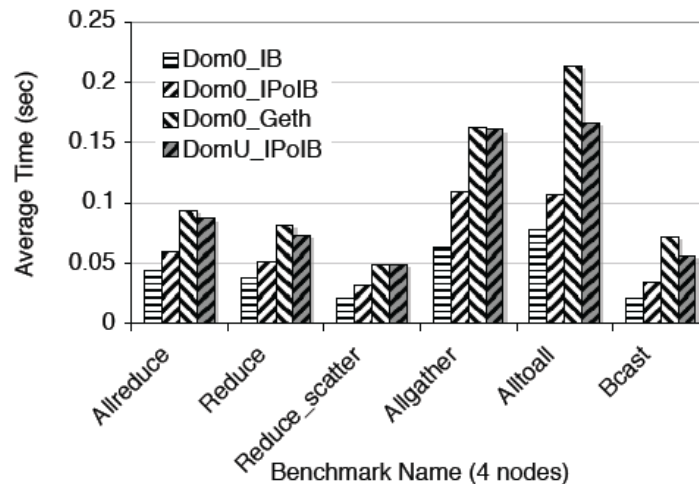
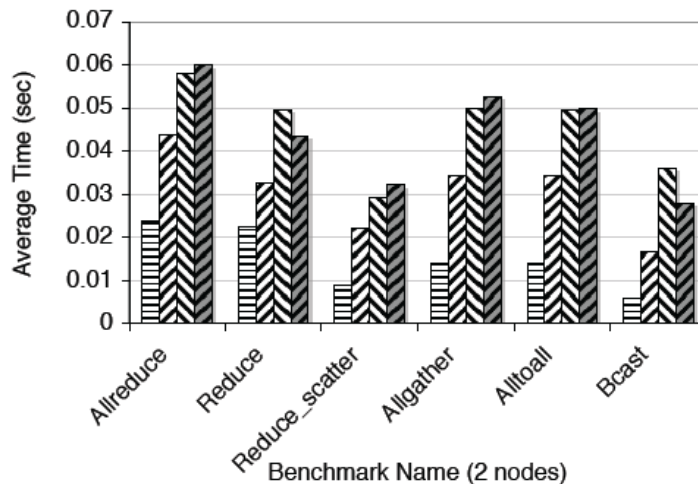
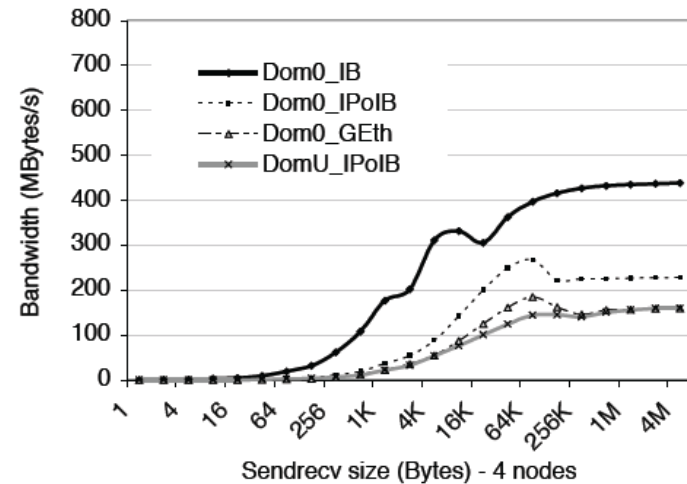
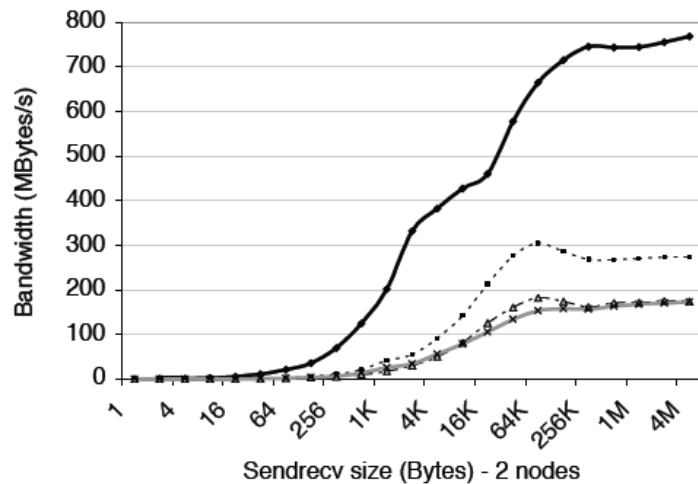
LM microbenchmark:

4 nodes, dual Opteron 280, 8GB ram, GigEth, SDR Infiniband
Environments: Host (Ub-Dom0), Guest (Ub-DomU), not virtualized (CentOS)



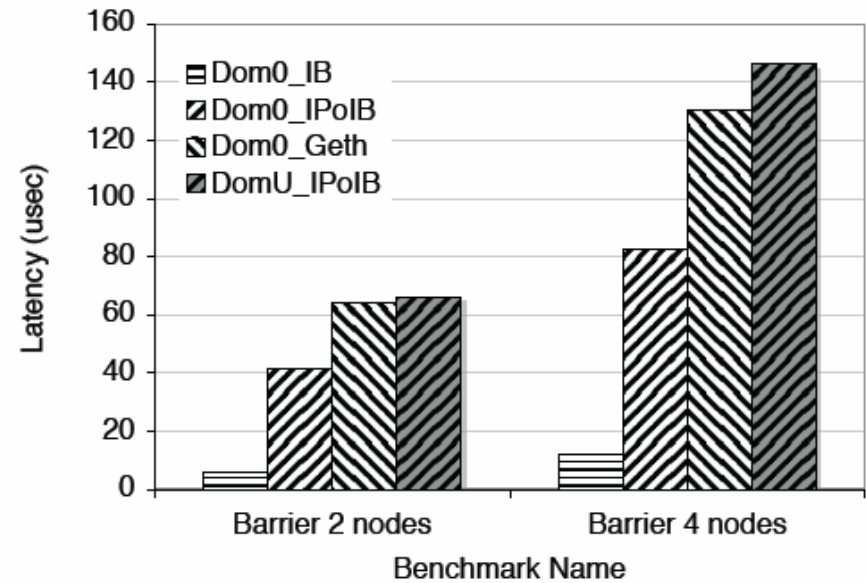
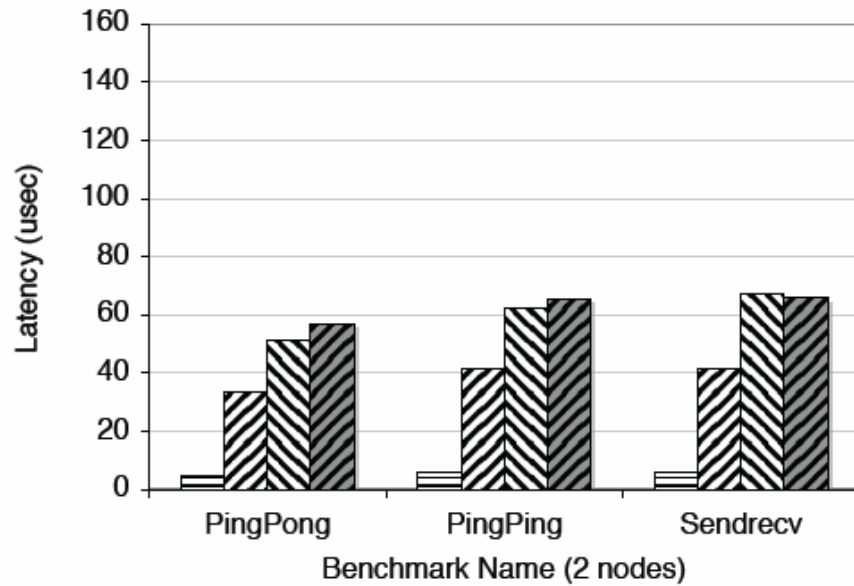
s and system

DOES IT WORK?



Intel MPI Benchmarks (VMAPICH): Bandwidth, collective communications

DOES IT WORK?



Intel MPI Benchmarks (VMAPICH): Latency



DOES IT WORK?

- ③ (Para)Virtualization impact syscalls heavily, but the real cost for computationally intensive applications is not clear
- ③ (Para)Virtualization has no cost for processor instructions (e.g. math)
- ③ (Para)Virtualization significantly penalizes communication bandwidth and latency, but:
 - ③ Current native IB drivers (with user space verbs) cannot be used within VMs
 - ③ TCP is used (IPoIB) within VMs. This is a major source of overhead
 - ③ Most overhead will be eliminated when user-space IB verbs for Xen VM are available (currently under development, see J. Liu, W. Huang, B. Abali and D.K. Panda. High Performance VMM-Bypass I/O in Virtual Machines. USENIX Annual Technical Conference 2006, Boston, MA, May 2006.)



CONTENTS

- ③ OVERVIEW
- ③ VIRTUALIZATION: CLUSTER OS
- ③ VIRTUALIZATION: CLUSTER STORAGE
- ③ FAULT TOLERANCE
- ③ TOOLS
- ③ DOES IT WORK?
- ③ CONCLUSIONS



CONCLUSIONS

- ③ VirtualLinux enables a virtualized, diskless, masterless cluster architecture
- ③ Virtual clusters are insulated from the physical and host cluster, allowing multiple simultaneous VCs on the same hardware
- ③ Storage virtualization coupled with a novel snapshot technique dramatically reduces installation time and permits a centralized management of the VCs
- ③ Storage virtualization permits the backup and restore of entire VCs
- ③ The absence of an individual master and local hard disks replaced by a SAN of redundant arrays avoids single points of failure
- ③ **IB Verbs within VMs are absolutely required for real use of VirtualLinux in a production environment, but they are coming...**



ACKNOWLEDGMENTS

- ④ VirtuaLinux is an ongoing experiment of Eurotech SpA and the HPC laboratory of the Computer Science Department of the University of Pisa
- ④ VirtuaLinux is a open source software under GPL available at <http://virtuallinux.sourceforge.net/>
- ④ VirtuaLinux project has been supported by the initiatives of the LITBIO Consortium, founded within FIRB 2003 grant by MIUR, Italy.
- ④ We are grateful to Peter Kilpatrick for his help in improving the presentation.



THANK YOU FOR YOUR ATTENTION



"On résiste à l'invasion des armées; on ne résiste pas à l'invasion des idées."

Victor Hugo, Paris 1877

