

Peter Münch

T/L Test and Development – Elastic Storage for Linux on System z

IBM Research & Development Germany



# Elastic Storage

## for Linux on IBM System z





# Session objectives

- This presentation introduces the Elastic Storage, based on General Parallel File System technology that will be available for Linux on IBM System z. Understand the concepts of Elastic Storage and which functions will be available for Linux on System z. Learn how you can integrate and benefit from the Elastic Storage in a Linux on System z environment. Finally, get your first impression in a live demo of Elastic Storage.



# Trademarks

**The following are trademarks of the International Business Machines Corporation in the United States and/or other countries.**

AIX*	FlashSystem	Storwize*	Tivoli*
DB2*	IBM*	System p*	WebSphere*
DS8000*	IBM (logo)*	System x*	XIV*
ECKD	MQSeries*	System z*	z/VM*

\* Registered trademarks of IBM Corporation

**The following are trademarks or registered trademarks of other companies.**

Adobe, the Adobe logo, PostScript, and the PostScript logo are either registered trademarks or trademarks of Adobe Systems Incorporated in the United States, and/or other countries. Cell Broadband Engine is a trademark of Sony Computer Entertainment, Inc. in the United States, other countries, or both and is used under license therefrom.

Intel, Intel logo, Intel Inside, Intel Inside logo, Intel Centrino, Intel Centrino logo, Celeron, Intel Xeon, Intel SpeedStep, Itanium, and Pentium are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

IT Infrastructure Library is a registered trademark of the Central Computer and Telecommunications Agency which is now part of the Office of Government Commerce.

ITIL is a registered trademark, and a registered community trademark of the Office of Government Commerce, and is registered in the U.S. Patent and Trademark Office.

Java and all Java based trademarks and logos are trademarks or registered trademarks of Oracle and/or its affiliates.

Linear Tape-Open, LTO, the LTO Logo, Ultrium, and the Ultrium logo are trademarks of HP, IBM Corp. and Quantum in the U.S. and

Linux is a registered trademark of Linus Torvalds in the United States, other countries, or both.

Microsoft, Windows, Windows NT, and the Windows logo are trademarks of Microsoft Corporation in the United States, other countries, or both.

OpenStack is a trademark of OpenStack LLC. The OpenStack trademark policy is available on the [OpenStack website](#).

TEALEAF is a registered trademark of Tealeaf, an IBM Company.

Windows Server and the Windows logo are trademarks of the Microsoft group of countries.

Worklight is a trademark or registered trademark of Worklight, an IBM Company.

UNIX is a registered trademark of The Open Group in the United States and other countries.

\* Other product and service names might be trademarks of IBM or other companies.

## Notes:

Performance is in Internal Throughput Rate (ITR) ratio based on measurements and projections using standard IBM benchmarks in a controlled environment. The actual throughput that any user will experience will vary depending upon considerations such as the amount of multiprogramming in the user's job stream, the I/O configuration, the storage configuration, and the workload processed. Therefore, no assurance can be given that an individual user will achieve throughput improvements equivalent to the performance ratios stated here.

IBM hardware products are manufactured from new parts, or new and serviceable used parts. Regardless, our warranty terms apply.

All customer examples cited or described in this presentation are presented as illustrations of the manner in which some customers have used IBM products and the results they may have achieved. Actual environmental costs and performance characteristics will vary depending on individual customer configurations and conditions.

This publication was produced in the United States. IBM may not offer the products, services or features discussed in this document in other countries, and the information may be subject to change without notice. Consult your local IBM business contact for information on the product or services available in your area.

All statements regarding IBM's future direction and intent are subject to change or withdrawal without notice, and represent goals and objectives only.

Information about non-IBM products is obtained from the manufacturers of those products or their published announcements. IBM has not tested those products and cannot confirm the performance, compatibility, or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

Prices subject to change without notice. Contact your IBM representative or Business Partner for the most current pricing in your geography.

This information provides only general descriptions of the types and portions of workloads that are eligible for execution on Specialty Engines (e.g, zIIPs, ZAAPs, and IFLs) ("SEs"). IBM authorizes customers to use IBM SE only to execute the processing of Eligible Workloads of specific Programs expressly authorized by IBM as specified in the "Authorized Use Table for IBM Machines" provided at [www.ibm.com/systems/support/machine\\_warranties/machine\\_code/aut.html](http://www.ibm.com/systems/support/machine_warranties/machine_code/aut.html) ("AUT"). No other workload processing is authorized for execution on an SE. IBM offers SE at a lower price than General Processors/Central Processors because customers are authorized to use SEs only to process certain types and/or amounts of workloads as specified by IBM in the AUT.

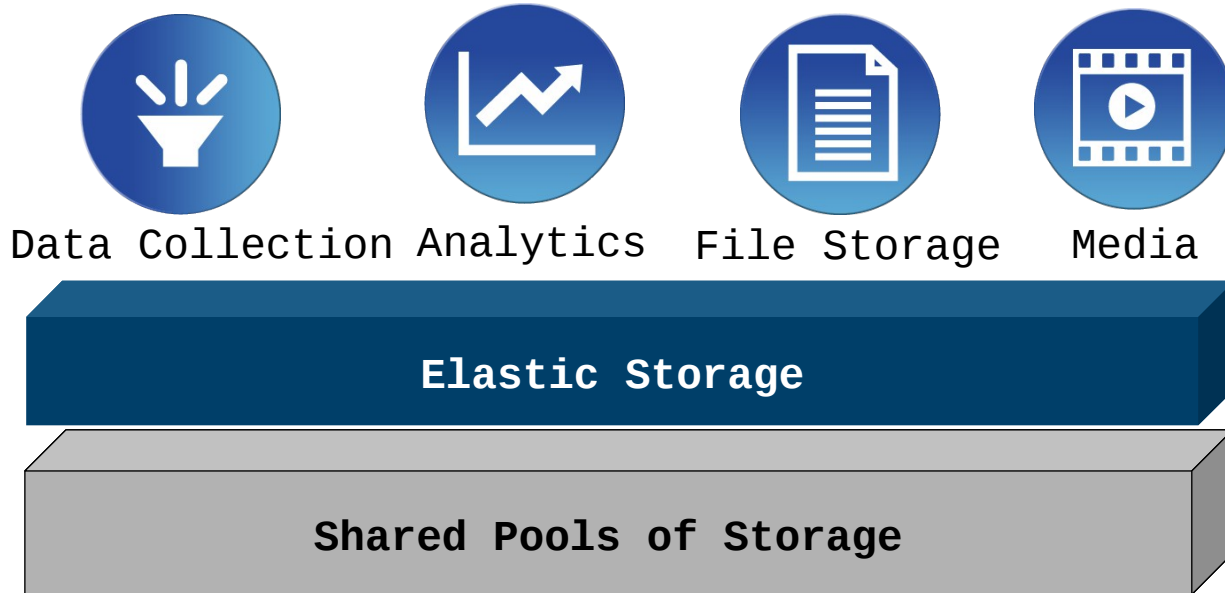


# Agenda

- Elastic Storage - General overview
- Elastic Storage for Linux on System z
  - Overview Version 1
  - Usage scenarios
    - WebSphere AppServer
    - WebSphere MQ
  - Outlook
- Quick Install Guide
- Demo

# Elastic Storage

*Provides fast data access and simple, cost effective data management*



- Streamline Data access
- Centralize Storage Management
- Improve Data Availability



# Clustered and Distributed File Systems

- Clustered file systems
  - File system shared by being simultaneously mounted on multiple servers accessing the same storage
  - Examples: Elastic Storage, Oracle Cluster File System (OCFS2), Global File System (GFS2)
- Distributed file systems
  - File system is accessed through a network protocol and do not share block level access to the same storage
  - Examples: NFS, OpenAFS, CIFS

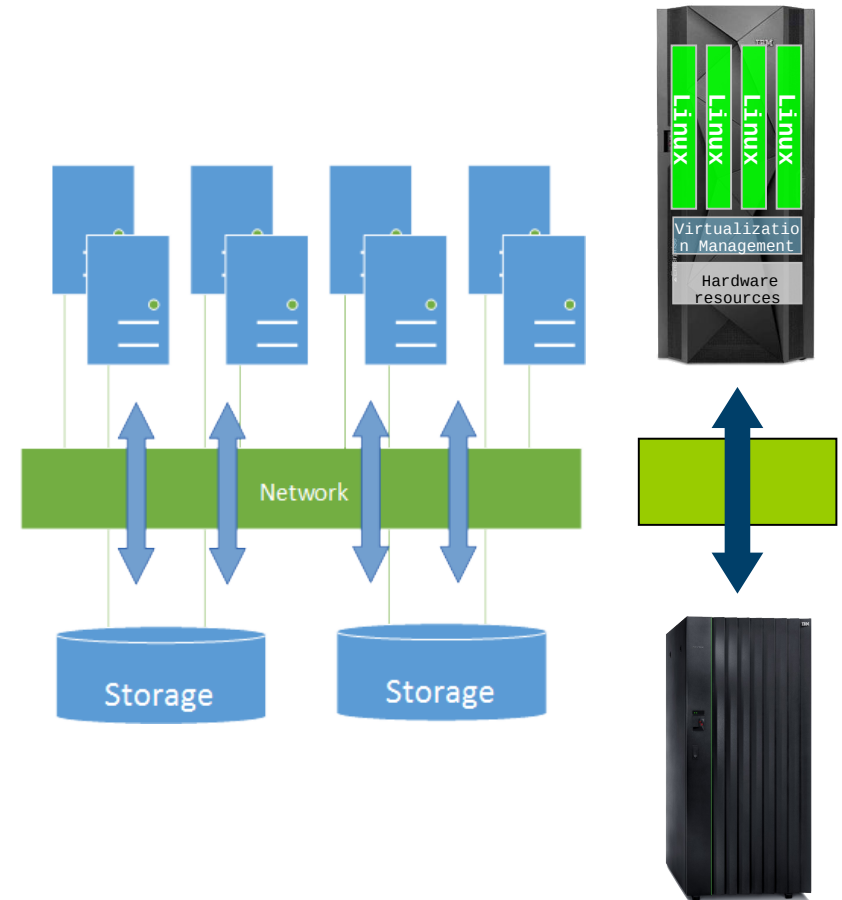
Available for Linux for System z:

- SUSE Linux Enterprise Server
  - Oracle Cluster File system (OCFS2)
- Red Hat Enterprise Linux
  - GFS2 (via Sine Nomine Associates)



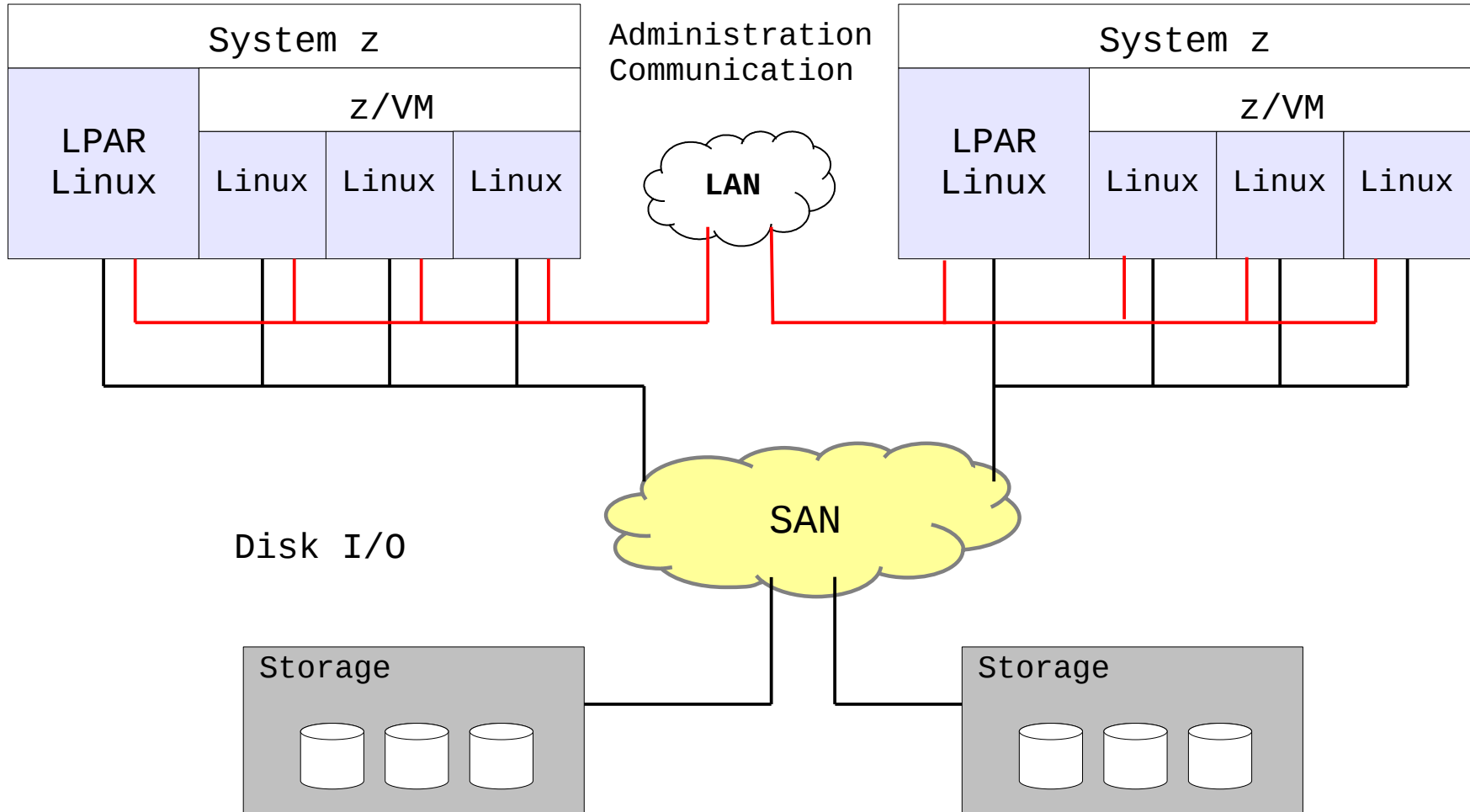
# What is Elastic Storage?

- IBM's shared disk, parallel cluster file system
- **Cluster:** 1 to 16,384\* nodes, fast reliable communication, common admin domain
- **Shared disk:** all data and metadata on storage devices accessible from any node through block I/O interface ("disk": any kind of block storage device)
- **Parallel:** data and metadata flow from all of the nodes to all of the disks in parallel.



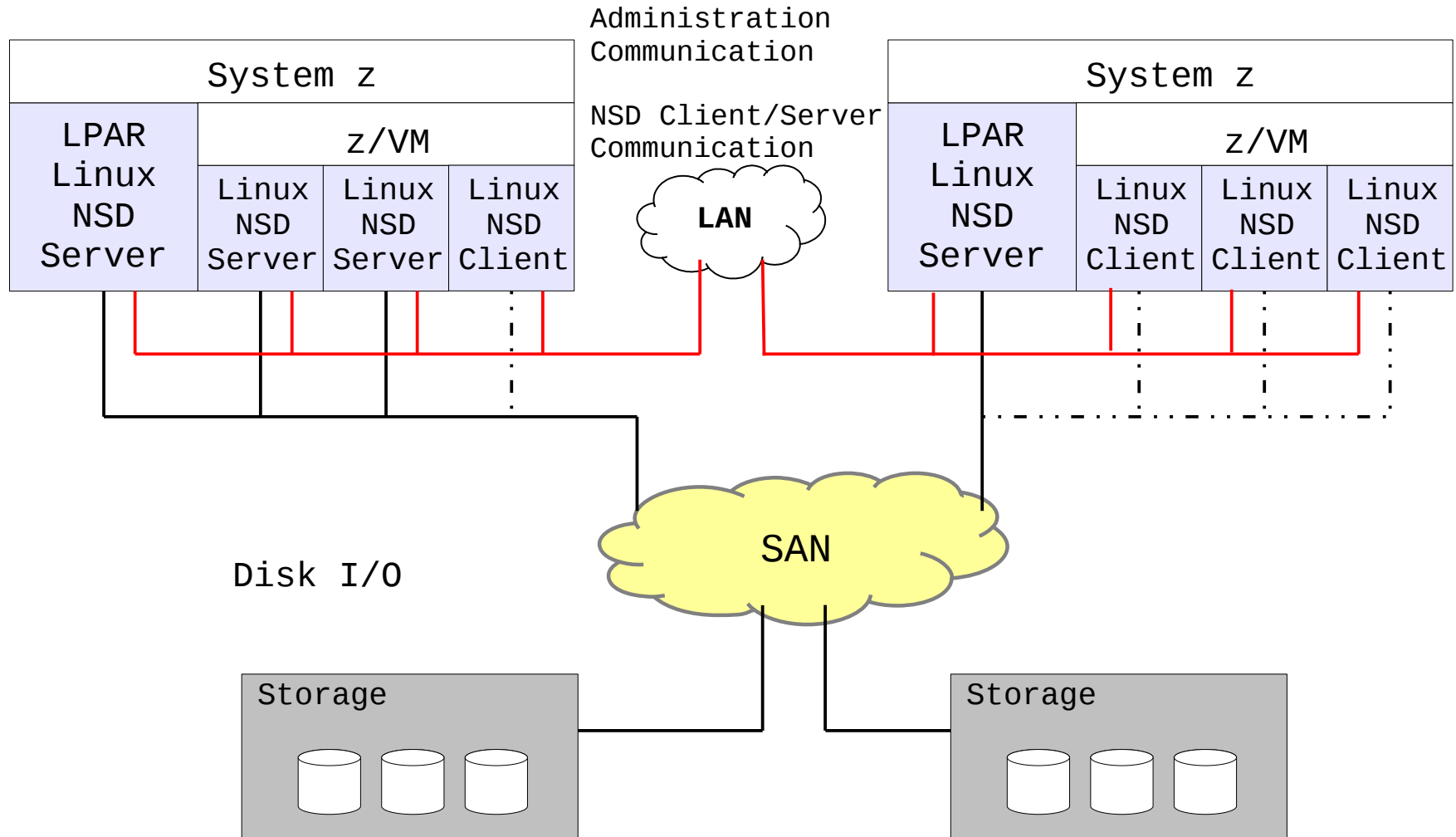
\*largest cluster in production as of August 2014  
Is LRZ SuperMUC 9400 Nodes of x86\_64

# Shared Disk (SAN) Model





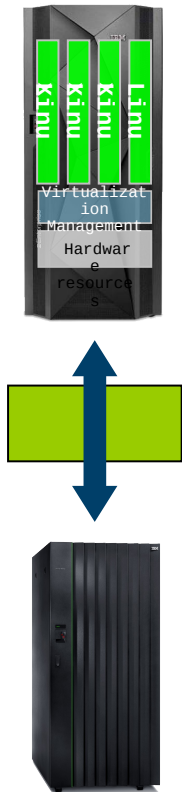
# Network Shared Disk (NSD) Model





# Elastic Storage Features & Applications

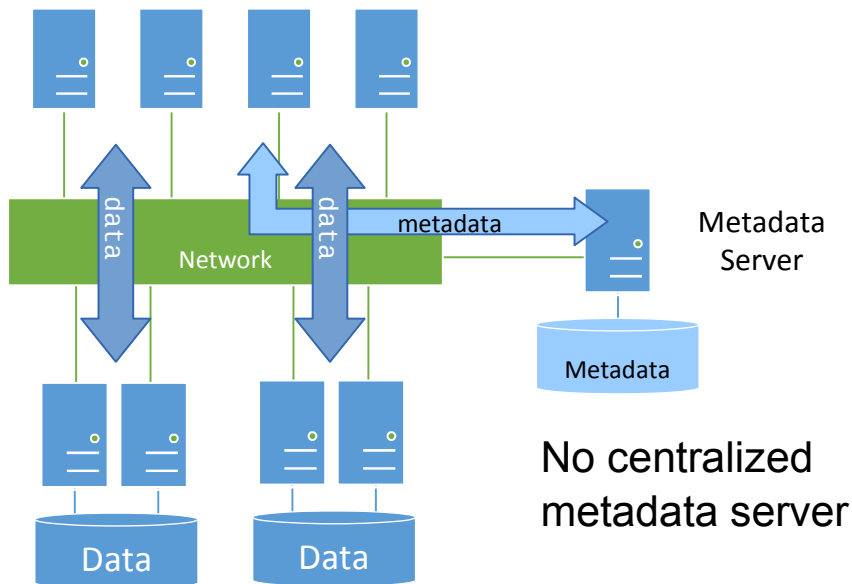
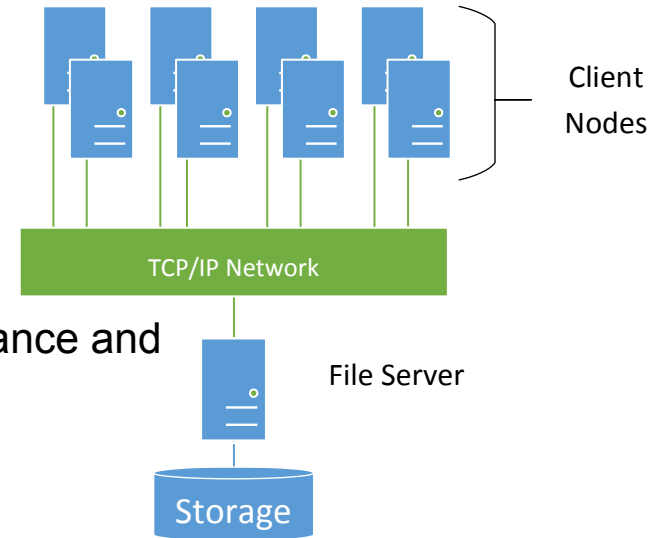
- Standard file system interface with POSIX semantics
  - Metadata on shared storage
  - Distributed locking for read/write semantics
- Highly scalable
  - High capacity (up to  $2^{99}$  bytes file system size, up to  $2^{63}$  files per file system)
  - High throughput (TB/s)
  - Wide striping
  - Large block size (up to 16MB)
  - Multiple nodes write in parallel
- Advanced data management
  - ILM (storage pools), Snapshots
  - Backup HSM (DMAPI)
  - Remote replication, WAN caching
- High availability
  - Fault tolerance (node, disk failures)
  - On-line system management (add/remove nodes, disks, ...)



# What Elastic Storage is NOT

Not a client-server file system  
like NFS, CIFS or AFS

No single-server performance and  
bottleneck scaling limits



No centralized  
metadata server



# Elastic Storage for Linux on System z

Based on GPFS Express Edition 4.1



## Elastic Storage for Linux on System z – Version 4.1

- Linux instances in LPAR mode or on z/VM, on the same or different CECs
  - Elastic Storage has no dependency on a specific version of z/VM
- Up to 32 cluster nodes with same or mixed Linux distributions/releases
- Heterogeneous clusters with client nodes without local storage access running on AIX, Linux on Power and Linux on x86
- Support for ECKD-based and FCP-based storage
- Support for IBM System Storage DS8000 Series, IBM Storwize V7000 Disk Systems, IBM XIV Storage Systems and IBM FlashSystem Systems, SVC
- Supported workloads are IBM WebSphere Application Server, IBM WebSphere MQ or similar workloads

*The Express Edition does not include features, therefore IBM is planning to offer enhanced functionality in future versions of Elastic Storage for Linux on System z.*

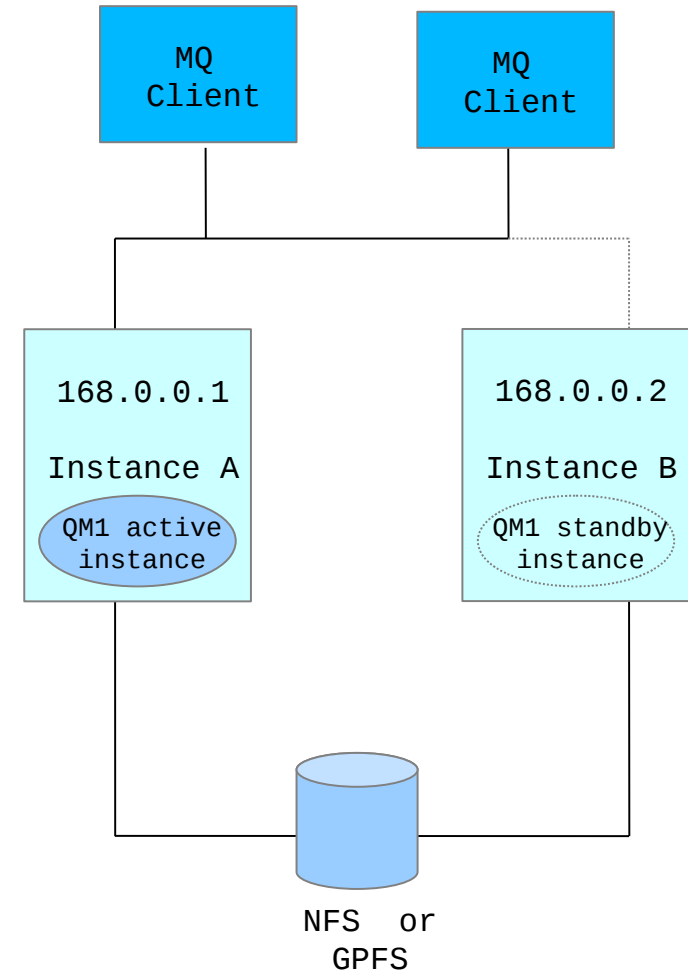


## Elastic Storage for Linux on System z – Version 1 (cont'd)

- Minimum supported Linux distributions:
  - SUSE Linux Enterprise Server (SLES) 11 SP3 + Maintweb-Update
  - Red Hat Enterprise Linux (RHEL) 6.5 + Errata Update
  - Red Hat Enterprise Linux (RHEL) 7.0
- While Elastic Storage V1 for Linux on System z does not support all functionality available for other platforms, this gap will be closed with the next updates.
- Elastic Storage for Linux on System z is part of the mainstream development, all future enhancements of Elastic Storage will become available for Linux on System z.

# Use Case for WebSphere MQ Multi-Instance Queue Manager (MIQM)

- High availability configuration of WebSphere MQ with two instances of the queue manager running on different servers, and either instance can be active.
  - A shared file system is required on networked storage, such as a NFS, or a cluster file system such as Elastic Storage

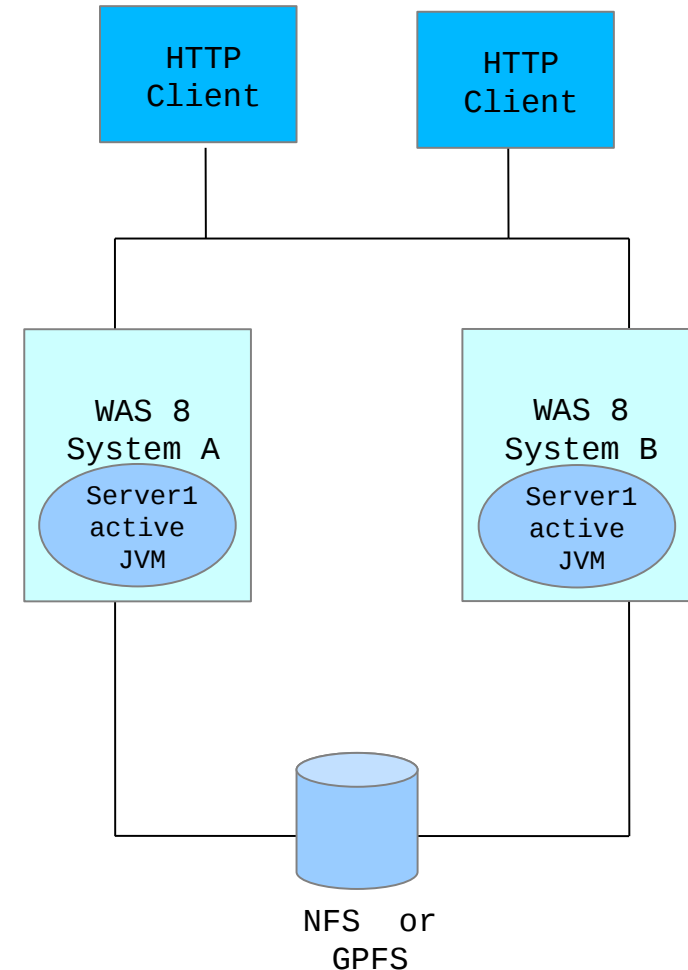




# Use Case for WebSphere AppServer HA Cluster

## HA Cluster

- High availability configuration of WebSphere AppServer with two instances of the application running on different servers, and both instances are active.
  - A shared file system is required for transaction logs on networked storage, such as a NFS, or a cluster file system such as Elastic Storage







# Outlook

- Multi-Cluster support
- Stretch-Cluster support (20, 40, 100, 200km for active/active DR configurations)
- Active File Management (AFM) / Information Lifecycle Management (ILM)
- AFM for active/backup configurations for clients not basing on hardware-based cross-site data replication (HA and DR)
- Tivoli Storage Manager (both backup and Hierarchical Storage Management (HSM))
- Support for heterogeneous clusters (Linux on System x,p,z)
- Encryption
- Support for other storage servers



# Quick Install Guide

## Elastic Storage for Linux on System z

Based on GPFS Express Edition 4.1



## Prerequisites Linux Distribution and Storage Hardware

- Supported Linux Distribution

Distribution	Minimum level	Kernel
SLES 11	SUSE Linux Enterprise Server 11 SP3 + Maintweb Update or later maintenance update or Service Pack	3.0.101-0.15-default
RHEL 6	Red Hat Enterprise Linux 6.5 + Errata Update RHSA-2014-0328 or later miner update	2.6.32-431.11.2.el6
RHEL 7		3.10.0-123.el7

- Supported Storage System
  - DS8000, XIV, V7000 and FlashSystem
- Elastic Storage has no dependency on a specific version of z/VM



# Software Prerequisites

- Additional Kernel Parameter

- set the following kernel parameters in */etc/zipl.conf* when booting the kernel
- `vmalloc = 4096G`
- `user_mode = home`

```
# cat /etc/zipl.conf  
Parameters = "... vmalloc=4096G user_mode=home ..."
```

- Passwordless communication between nodes of GPFS cluster
- Cluster system time coordination via NTP or equivalent
- Required kernel development packages to be installed on at least one system to build the kernel modules



# Exchange ssh keys between all GPFS nodes

- Passwordless access between all GPFS nodes is a prerequisite
- Exchange ssh key from one node to all other nodes
  - Create ssh-keys at node1:

```
# cd .ssh  
# ./ssh-keygen #hit return by all questions
```

- Copy ssh keys to authorized\_keys at node1:

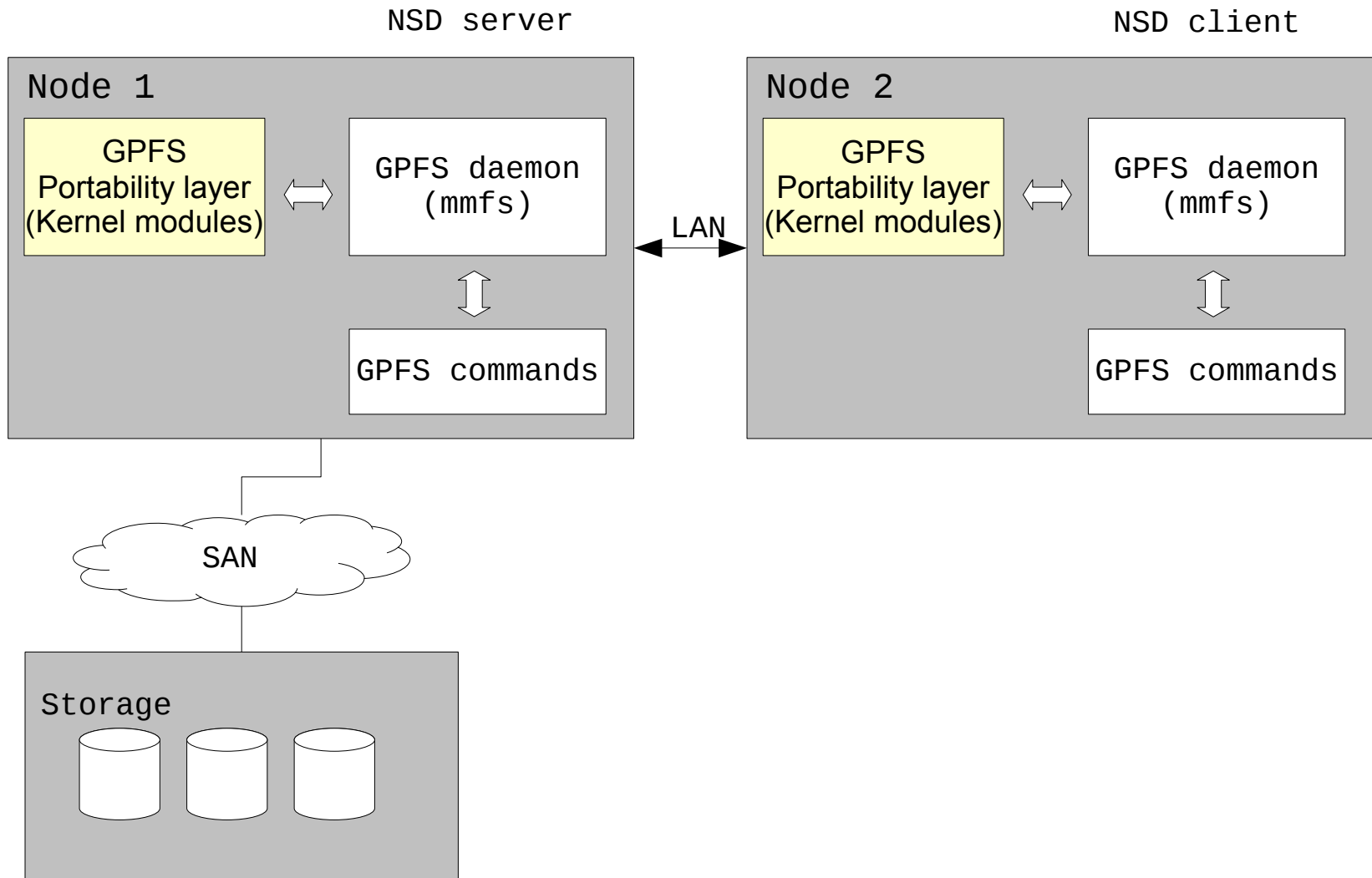
```
# cat id_rsa.pub >> authorized_keys  
# ssh localhost  
# ssh node1  
# ssh node1.domain.com
```

Copy id\_rsa.pub to other nodes

```
# ssh-copy-id -i /root/.ssh/id_rsa.pub root@node2
```

- Do ssh connects from each node to each other node and localhost (with and without the domain name)

# Overview





# Install GPFS product

- Install GPFS product RPM packages on all nodes of the cluster
  - Packages name: gpfs.\*.rpm
- GPFS product files can be found after installation at
  - /usr/lpp/mmfs
- Build the GPFS kernel modules (portability layer) e.g. development system

```
# cd /usr/lpp/mmfs/src/  
# make Autoconfig  
# make World  
# make InstallImages
```

- Build an rpm (make rpm) and install this rpm on all related nodes
- Reboot all nodes



# Plan for GPFS Cluster

- Create a NodeFile to define the role of the nodes (FS Manager): e.g. nodes.file

```
node1:quorum-manager:
node2:quorum-manager:
node3:quorum:
node4::
```

- Create a stanza file to define Network Shared Disks (NSD) to be used by GPFS file systems : e.g. nsd.file

```
%nsd: device=/dev/mpatha
      nsd=NSD_1
      servers=node1,node2
      usage=dataAndMetadata
%nsd: device=/dev/mpathb
      nsd=NSD_2
      servers=node1
      usage=dataAndMetadata
%nsd: device=/dev/mpathc
      nsd=NSD_3
      servers=node1
      usage=dataAndMetadata
```





# Quick Install Guide

- Create a GPFS cluster

– -A options: Start GPFS daemons automatically when nodes come up

```
node1# mmcrcluster -N nodes.file -C cluster1 -r /usr/bin/ssh  
-R /usr/bin/scp -A
```

- Change the type of GPFS license associated with the nodes

```
node1# mmchlicense server --accept -N node1,node2,node3  
node1# mmchlicense client --accept -N node4
```

- Start the GPFS cluster on all nodes

```
node1# mmstartup -a
```



## Quick Install Guide (cont'd)

- Get information about the previously activated GPFS cluster

```
node1# mmfsccluster
GPFS cluster information
=====
GPFS cluster name:          cluster1
GPFS cluster id:           18000255686092070264
GPFS UID domain:          cluster1.domain.com
Remote shell command:      /usr/bin/ssh
Remote file copy command:  /usr/bin/scp
Repository type:           CCR

GPFS cluster configuration servers:
-----
Primary server:    node1.domain.com (not in use)
Secondary server:  (none)

Node  Daemon node name      IP address      Admin node name  Designation
-----
1     node1.domain.com         10.20.80.86    node1.domain.com quorum-manager
2     node2.domain.com         10.20.80.87    node1.domain.com quorum-manager
3     node3.domain.com         10.20.80.88    node1.domain.com quorum
4     node4.domain.com         10.20.80.89    node1.domain.com
```

## Quick Install Guide (cont'd)

- Get information about the status of the GPFS cluster

```
node1# mmgetstate -a
Node number  Node name      GPFS state
-----
          1      node1      active
          2      node2      active
          3      node3      active
          4      node4      active
```

- Create Network Shared Disks used by GPFS

```
node1# mmcrnsd -F nsd.file
```

- Create an GPFS file system

— -A option: File system will be mounted when GPFS daemon starts

```
node1# mmcrfs esfs1 -F nsd.file -T /elastic_storage -A yes
node1# mmcrfs esfs2 "NSD_4;NSD_5" -T /elastic_storage2 -A yes
```



# Quick Install Guide

- Retrieve information about the Network Shared Disks

```
node1# mmlsnsd
```

File system	Disk name	NSD servers
-----	-----	-----
esfs1	NSD_1	node1.domain.com,node2.domain.com
esfs1	NSD_2	node1.domain.com
esfs1	NSD_3	node1.domain.com

- Mount all GPFS file systems on all nodes in the cluster

```
node1# mmmount all -a
```



# Manage GPFS Cluster: useful commands

- Manage Elastic Storage Cluster / Node
  - mmcrcluster, mmchcluster, mmlscluster
  - mmstartup, mmshutdown
  - mmchlicense
  - mmaddnode, mmchnode, mmdelnode, mmlsnode
- Manage Network Shared Disks (NSD)
  - mmcrnsd, mmchnsd, mmdelnsd, mmlsnsd
- Manage Elastic Storage Filesystem
  - mmcrfs, mmchfs, mmdelfs, mmlsfs
  - mmcrsnapshot, mmdelsnapshot, mmlssnapshot
  - mmadddisk, mmchdisk, mmdeldisk, mmlsdisk



# Resources

- **ibm.com:**

[ibm.com/systems/platformcomputing/products/gpfs/](http://ibm.com/systems/platformcomputing/products/gpfs/)

- **Public Wiki:**

[ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/General Parallel File System \(GPFS\)](http://ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/General_Parallel_File_System_(GPFS))

- **IBM Knowledge Center:**

[ibm.com/support/knowledgecenter/SSFKCN/gpfs\\_welcome.html?lang=en](http://ibm.com/support/knowledgecenter/SSFKCN/gpfs_welcome.html?lang=en)

- **Data sheet: IBM General Parallel File System (GPFS) Version 4.1**

[ibm.com/common/ssi/cgi-bin/ssialias?subtype=SP&infotype=PM&appname=STGE\\_DC\\_ZQ\\_USEN&htmlfid=DCD12374USEN&attachment=DCD12374USEN.PDF](http://ibm.com/common/ssi/cgi-bin/ssialias?subtype=SP&infotype=PM&appname=STGE_DC_ZQ_USEN&htmlfid=DCD12374USEN&attachment=DCD12374USEN.PDF)

# Questions?



***Peter Münch***

*T/L Test and  
Development*

*Elastic Storage for  
Linux on System z*

*Hechtsheimer Strasse 2  
55131 Mainz, Germany*

*Phone +49 (0)6131-84-3462  
muenchp@de.ibm.com*