

# SSSD23 multi tenant file system consolidation

![](_page_0_Picture_2.jpeg)

jochen.zeller@sva.de

# ABOUT US

![](_page_1_Figure_1.jpeg)

A SHORT COMPANY INTRODUCTION  $\frac{1}{21.03.2023 / 2}$ 

# WHAT MAKES US UNIQUE?

![](_page_2_Figure_1.jpeg)

![](_page_2_Picture_2.jpeg)

# / SO MANY CLUSTER AND FILE SYSTEMS

![](_page_3_Figure_1.jpeg)

- 2 x AIX VM as NSD server per cluster
- SVC as backend storage, SVC vdisk mirror
- Many file systems per cluster
- Not every client mounts all file systems
- Scale clients run applications on RHEL

| cluster | file systems | capacity (sum) |
|---------|--------------|----------------|
| 8       | 116          | 78TB           |

The task:

- Replace IBM POWER by x86
- Replace IBM AIX by RHEL
- Replace SVC
- The separation of data must remain

![](_page_3_Picture_13.jpeg)

# / GENERAL CONDITIONS

- Pure Storage, FlashArray, is the customers successor for SVC
- Mirroring, high availability between two data centers, is necessary
- Performance requirements: 2GB/sec overall is fine
- Are virtual NSD servers an option, like the AIX VMs were before?
  - 16 NSD server VMs are cheaper than 16 rack servers
- Storage compression would be nice, current SVC compression is 1:2
- Spectrum Scale clients access the file system by Ethernet / tcpip
- NFS shares would be great

21.03.2023

![](_page_4_Picture_9.jpeg)

# / PLANNING WORKSHOP

- Let's start with a vision: We will have only one file system
- With a single file system, we just need two or four NSD servers
  - File system consolidation with separation means a remote cluster setup
- What about the storage?
  - Alternatives to Pure Storage? IBM?
  - ESS is to expensive for this small configuration with low performance requirements
  - IBM FS5200 could be a great (and cheap?) solution
    - Mirroring with HyperSwap (no additional license cost)
      - Reduces the Spectrum Scale licenses by half compared to Spectrum Scale mirroring
    - FCMs with compression

![](_page_5_Picture_11.jpeg)

![](_page_5_Picture_12.jpeg)

![](_page_5_Picture_13.jpeg)

# / THE HARDWARE SOLUTION

- 4 NSD server + 1 quorum VM
  - 2 x 32Gbit FC
  - 2 x 25Gbit Ethernet
- 2 x IBM FS5200 in Hyper Swap cluster
  - 70TB @ 12 x 9,6TB FMC
  - 140TB with 1:2 compression, 120TB assigned
- SAN for storage access and hyperswap (already present)

![](_page_6_Figure_8.jpeg)

![](_page_6_Picture_9.jpeg)

### / HYPERSWAP AND MULTIPATHING

- Both FS5200 build a cluster and both provide the same LUNs
- The paths of a LUN are on one site "active" and on the other site "enabled"
- Host multipathing uses ALUA <u>A</u>symmetric <u>Logical Unit A</u>ccess
- Hyperswap can change the active site for a LUN automatically

![](_page_7_Picture_5.jpeg)

![](_page_7_Figure_6.jpeg)

# / SPECTRUM SCALE AND HYPERSWAP – DATA ACCESS PATH

![](_page_8_Figure_1.jpeg)

# / PERFORMANCE

- FS5200 HyperSwap has a reduced write performance compared to a single FS5200
  - Host writes to one FS5200, this must copy the data to the second FS5200. When the second FS5200 reports "sync", the IO is reported by the first FS5200 as sync to the host.

![](_page_9_Figure_3.jpeg)

![](_page_9_Figure_4.jpeg)

# / TO THE ACTUAL TOPIC: MULTI TENANCY / REMOTE MOUNT

![](_page_10_Figure_1.jpeg)

## / REMOTE FILESET LOOK AND FEEL

- Mount one (or more) remote file system
- Share access of one or more filesets
- Content of root fileset is always visible / accessible
- In our case, a "1s" to the file system root /gpfs shows all directories c101,c102,c103, ... c199

**# ls** 

- cl01 cl02 cl03 cl04 cl05
- Because all subdirectories in the cIXY directories are filesets (fs01, fs02, ...), they are only visible if they are shared by the source cluster.

# cd cl02 # ls

fs01 fs02 fs09

![](_page_11_Picture_11.jpeg)

![](_page_11_Picture_12.jpeg)

# / HOW TO SETUP THE REMOTE CLUSTER

- **THE** question: many multi node clusters or even more single node clusters?
- The problem with multi node cluster:
  - You have to take care of the cluster quorum
  - This means that at least three nodes are not equal to the others
- Decision: single node cluster

![](_page_12_Picture_6.jpeg)

• Why: The client cluster setup and remote relationship will be automated with Ansible anyway. Whether execute 8 times or more than 100 times does not make much difference. By the way:

#### Q5.6:

What is the limit of remote clusters that can join a local cluster? A5.6:

There is not really a limit. The smallest cluster possible is a single node cluster, which means that 16,383 clusters can join a local cluster (16384 - 1).

https://www.ibm.com/docs/en/STXKQY/gpfsclustersfaq.html#inclusters

### / MULTI TENANT FILE SYSTEM CONSOLIDATION. DONE.

![](_page_13_Figure_1.jpeg)

![](_page_13_Picture_2.jpeg)

# / A TECHNICAL AWESOME PROJECT

MCOT – multiple connections over TCPIP Cluster Export Services of AFM for data migration of Hor Berline o FCM hardware compression

data migrat

![](_page_14_Picture_2.jpeg)

#### / SUMMARY

- IBM FS5200 he has very good price-performance ratio (NVMe, FCM) for small configurations
- HyperSwap is a very good and cost neutral solution for two data center mirroring
- Remote fileset access is a very good feature for consolidation and multi tenancy, and certainly also for other use cases
- Thanks to Ansible automation

The good thing about Spectrum Scale is that you can design and adjust a lot. The bad thing about Spectrum Scale is that you can design and adjust a lot.

![](_page_15_Picture_6.jpeg)

![](_page_15_Picture_7.jpeg)

## / THE END!

![](_page_16_Figure_1.jpeg)

![](_page_16_Picture_2.jpeg)

# / CONTACT

![](_page_17_Picture_1.jpeg)

#### JOCHEN ZELLER

IT Architect

Technical Leader IBM Spectrum Scale

| Phone.: | +49 151 180 256 77   |
|---------|----------------------|
| Mail:   | jochen.zeller@sva.de |

![](_page_17_Picture_6.jpeg)