

# **On-demand File Systems**

Production use, use cases and lessons learned

Mehmet Soysal Steinbuch Centre for Computing (SCC)



#### www.kit.edu



#### **Overview**

- Motivation / Background
- Integration into HPC batch system
- Generic Benchmarks
- Use Cases
  - OpenFOAM, NAStJA, FLAPS
- Summary and Conclusion

### **Motivation / Background**



- ADA-FS project within SPPEXA
  - ZDV@Mainz, ZIH@Dresden, SCC@KIT
  - Tackling the I/O bottleneck on HPC Systems
  - Use of on-demand file systems (ODFS)
- At SCC we have:
  - traditionally node-local storage (HDD, SSD, NVMe ...)
  - heavy I/O use cases
- We needed a solution to handle heavy I/O until use cases are optimized

# Motivation / Background (1)



#### ForHLR II

- Compute nodes with 400GB SATA-SSD
- ~550/450 MB/s R/W per SSD

#### Horeka

- Compute nodes with 1TB NVMe-SSD
- ~3000/1000 MB/s R/W per NVMe



#### Motivation / Background

#### Integration into HPC batch system

Generic Benchmarks

Use Cases

OpenFOAM, NAStJA, ML/DL (Gromacs)

Summary and Conclusion

#### **Using Global Parallel File System**





### **Using On-demand File System**





#### **Integration into Batch System**

![](_page_7_Picture_1.jpeg)

BeeGFS / BeeOND is used as ODFS for production use

Other FS planned, e.g., GekkoFS<sup>1</sup>...

- Only on-demand minimal changes to system usage
  - Request on job submission "--constraint BEEOND"
- Different implementations
  - MOAB + (SLURM) + data staging
  - SLURM + Burst Buffer
  - SLURM + Pro-/Epilog
- Loopback device as storage

1) Vef et al. "GekkoFS - A Temporary Distributed File System for HPC Applications"

# Integration using Job Pro-/Epilog

![](_page_8_Picture_1.jpeg)

![](_page_8_Figure_2.jpeg)

Simple integration

- ODFS start/stop in pro/-epilogue
- No data staging (User must move data within job)
- Very robust

# **Extending Slurm Burst Buffer**

![](_page_9_Picture_1.jpeg)

![](_page_9_Figure_2.jpeg)

- Extending the Burst Buffer plugin
  - Only proof of concept
- Creation of a reservation for a burst buffer

1) Bachelorthesis : Valentin Voigt "Entwicklung eines on-demand Burst-Buffer-Plugins für HPC-Batch-Systeme"

### **Moab and Data Staging**

![](_page_10_Picture_1.jpeg)

![](_page_10_Figure_2.jpeg)

MOAB's data staging feature extended

Job is split in three sub-jobs

MOAB offers special feature to inherit a resource allocation

#### 12/35 November 3, 2021 Mehmet Soysal – On-demand File Systems: lessons learned

#### Loopback device

![](_page_11_Figure_3.jpeg)

Yamamoto et al., Analysis and elimination of client evictions on a large scale lustre based file system. Presentation at the Lustre User Group 15 (LUG'15), 2015.

![](_page_11_Picture_5.jpeg)

#### **Summary – Integration**

![](_page_12_Picture_1.jpeg)

- The simplest solution is usually the best (Pro-/Epilogue)
- One storage target per node
- Using loopback device as storage target
- Pre-create directories with different stripe settings

![](_page_13_Picture_0.jpeg)

- Motivation / Background
- Integration into HPC batch system

#### Generic Benchmarks

- Use Cases
  - OpenFOAM, NAStJA, ML/DL (Gromacs)
- Summary and Conclusion

# Karlsruhe Institute of Technology

# **Generic Benchmarks**

- Start/Stop time
- Read/Write Throughput
- Metadata performance
  - Multiple Metadata Server (MDS)
- Topology awareness

# Startup and Shutdown of BeeGFS

![](_page_15_Picture_1.jpeg)

| Nodes        | 8    | 16   | 32   | 64   | 128   | 256   |
|--------------|------|------|------|------|-------|-------|
| Startup (s)  | 10.2 | 16,7 | 29,3 | 56,5 | 152,1 | 222,4 |
| Shutdown (s) | 11,9 | 12,1 | 9,4  | 15,9 | 36,1  | 81,1  |

Results above with beeond script using pdsh (clean start/stop)
Probably to long for very short jobs

- With optimization possible to start on 256 nodes within a minute<sup>1</sup>
- Shutdown can be done within seconds ("stoplocal")

1) Mehmet Soysal "Speeding up beeond startup" https://groups.google.com/g/fhgfs-user/c/g8ysFS35Ucs/m/E-RtxKyiCAAJ

### **R/W Throughput – ForHLR & Horeka**

![](_page_16_Picture_1.jpeg)

![](_page_16_Figure_2.jpeg)

17/35 2. November 2021 Mehmet Soysal – On-demand File Systems: lessons learned

Steinbuch Centre for Computing (SCC) Karlsruhe Institute of Technology (KIT)

#### **Metadata Performance (Horeka)**

![](_page_17_Picture_1.jpeg)

![](_page_17_Figure_2.jpeg)

# **Topology awareness**

![](_page_18_Picture_1.jpeg)

What if we do not have a fully-non blocking interconnect?

- Large Systems have a "weaker" interconnect (Dragonfly, Tofu, Hypercube ...)
- Cost of interconnect
- BeeGFS offers storage pools

#### Storage pools for topology awareness

# Small Island ForHLR II

![](_page_19_Picture_1.jpeg)

![](_page_19_Figure_2.jpeg)

![](_page_20_Picture_0.jpeg)

#### **Storage Pool per Leaf Switch**

![](_page_20_Figure_2.jpeg)

# **Comparison w/o Storage Pools**

![](_page_21_Picture_1.jpeg)

![](_page_21_Figure_2.jpeg)

### Summary – Benchmarks

![](_page_22_Picture_1.jpeg)

- BeeGFS scales well
- Bandwidth scales almost linear
- MDS performance scales with multiple MDS
- BeeGFS has the same "weakness" like Lustre
  - MDS bottleneck
  - Working in same directory => slow performance
  - Need to choose right strip size/count
  - Care about file size / stripcount / OST capacity

![](_page_23_Picture_0.jpeg)

- Motivation / Background
- Integration into HPC batch system
- Generic Benchmarks
- Use Cases
  - OpenFOAM, NAStJA, FLAPS
- Summary and Conclusion

![](_page_24_Picture_0.jpeg)

#### **Use Cases**

#### OpenFoam<sup>1</sup>

- 240 nodes (20 Cores)
- ~450k files / 120 GB per snapshot
- NAStJA<sup>2</sup> Stencil Code solver
  - 240 nodes (20 Cores)
  - 4800 files / 4800MB per snapshot

#### FLAPS<sup>3</sup> - dynamic particle swarm optimization

- ML/DL + Gromacs
- 50 nodes
- Heavy metadata

- 1) Zirwes et al. "Automated Code Generation for Maximizing Performance of Detailed Chemistry Calcu- lations in OpenFOAM". In: High Performance Computing in Science and Engineering '17.
- 2) Berghoff et al. "Massively Parallel Stencil Code Solver with Autonomous Adaptive Block Distribution".In: IEEE Transactions on Parallel and Distributed Systems 29.10 (2018)
- 3) Weiel et al. "Dynamic particle swarm optimization of biomolecular simulation parameters with flexible objective functions." Nat Mach Intell 3, https://doi.org/10.1038/s42256-021-00366-3

# **OpenFOAM – Lustre Server Unix Load**

- 5 short runs
- Average unix load of MDS/OSS
- Whole system is getting "laggy"

![](_page_25_Figure_4.jpeg)

![](_page_25_Picture_5.jpeg)

# **OpenFOAM – Runtime per Timestep**

![](_page_26_Picture_1.jpeg)

- Blue line average of 5 runs
- 100 timesteps ~ 30 min
- Snapshot every 5 timesteps
- Less variation when using ODFS (black bars)
- Simulation takes slightly more time when using ODFS

![](_page_26_Figure_7.jpeg)

# **NAStJA – Runtime per Timesteps**

![](_page_27_Picture_1.jpeg)

- Snapshot every timestep
- Almost no variation when using ODFS
- Very high interference on ODFS

![](_page_27_Figure_5.jpeg)

![](_page_28_Picture_0.jpeg)

### **FLAPS – One Generation**

![](_page_28_Figure_2.jpeg)

![](_page_28_Figure_3.jpeg)

![](_page_29_Picture_0.jpeg)

#### **FLAPS – Two Generations**

![](_page_29_Figure_2.jpeg)

### Summary – Use Cases

![](_page_30_Picture_1.jpeg)

Effect on application hard to predict

- Some running slower and some faster<sup>1</sup>
- Metadata performance more important
- Easy to use for users no code changes needed
  - Data staging is handled by users
- Some use cases are enforced to use ODFS

1) Versick et al. "Performance gains in an ESM using parallel ad-hoc file systems" EGU 2020

![](_page_31_Picture_0.jpeg)

- Motivation / Background
- Integration into HPC batch system
- Generic Benchmarks
- Use Cases
  - OpenFOAM, NAStJA, FLAPS
- Summary and Conclusion

# **Summary & Conclusion (1)**

![](_page_32_Picture_1.jpeg)

- ODFS (BeeGFS) is easy to integrate
- Startup time is acceptable
- SSD wear leveling should be observed
  - We did not notice any problems
- Loopback device improves performance
  - Especially with small files
- BeeGFS scales well (Metadata performance ok)
  - Multi MDS
  - Topology awareness with storage pools possible

# Summary & Conclusion (2)

![](_page_33_Picture_1.jpeg)

- Impact on application hard to predict
- Observation: Add one additional node to job
  - Use this additional node for MDS
  - mpirun –nolocal
- Observation: Metadata performance more important than bandwidth
  - Added options for multi metadata servers
- User handles staging staging
  - Concurrent data staging possible with less interference<sup>1</sup>

#### Advantages of ODFS outweigh any disadvantages

Most important: Load reduction of the global parallel file system

1) Soysal et al. "Using On-Demand File Systems in HPC Environments" HPCS 2019 DOI: 10.1109/HPCS48598.2019.9188216

![](_page_34_Picture_0.jpeg)

# Thanks for your attention

Questions ?

Scripts available: https://github.com/mehsoy/ODFS-tools

35/35 November 3, 2021 Mehmet Soysal – On-demand File Systems: lessons learned

Steinbuch Centre for Computing (SCC) Karlsruhe Institute of Technology (KIT)

#### **Backup Slides**

![](_page_35_Picture_1.jpeg)

![](_page_35_Figure_2.jpeg)