The Virtual Institute for I/O and the IO-500

Julian Kunkel¹, Jay Lofstead², John Bent³

¹ Deutsches Klimarechenzentrum (DKRZ) ² Sandia National Laboratory

Contact: kunkel@dkrz.de

³ Seagate Government Solutions



INTRODUCTION

The research community in high-performance computing is organized loosely. There are many distinct resources such as homepages of research groups and benchmarks. The Virtual Institute for I/O aims to provide a hub for the community and particularly newcomers to find relevant information in many directions. Additionally, we host the **high-performance storage list**. Similarly to the top500, it contains information about supercomputers and their storage systems. Additionally, in the community, we are working on standardizing an I/O benchmark.

This poster introduces the Virtual Institute for I/O, the high-performance storage list and the effort for the IO-500 which are unfunded com-

COMMUNITY CONTENT OF THE VI4IO WIKI

Worldwide research groups that address high-performance I/O including:

- A taglist for available knowledge
- Research products such as file systems
- Ongoing research projects

	Virtuar	Institute				Search Recent Changes
ou are here: Virtu	al Institute for I/O » Gr	oups » Research » [DE/UHAM			
DE/FZJ	DE/HLRS	DE/JGU	DE/LRZ	DE/SCC	DE/TUD	DE/UHAM
ES/BSC	ES/HPC4EAS	ES/UPM	FR/CEA	FR/INRIA	GR/FORTH	US/ANL
US/SNL	US/UCSC					
ins grou hoi	name: Scientific Comp ionality: DE tittution: Universität Han up head: Prof. Dr. Thoma mepage: https://wr.inf cations: https://wr.inf cations: \$3,5668558.	iburg is Ludwig ormatik.uni-hamburg.de/ ormatik.uni-hamburg.de/ 9.9742932		e, PVFS2		 DE/UHAM Products Projects Past projects Similar groups

V4 Virtual Institute for I/O ou are here: Virtual Institute for I/O » Tools » Benchmarks » IOI **IOR** IOR Vsage
 Example Output layers: POSIX, MPI-IO, HDF5, NCMP webpage: 🕥 https://github.com/LLNL/i You can set the API to be used using the -a flag. A multitude of other commandline options is do Here are a few important one a set the api to one of: POSIX, MPIIO, HDFS5 or NCMP -N number of tasks _i number of repetitions of the whole test Edit Command line used: ./ior -a POSI lachine: Linux hostname Test 0 started: Mon Mar 21 16:19:50 2010

Relevant I/O related tools and benchmarks.

HPSL 2017

The current list contains 33 sites:

#	Site		Supercomputer			Storage	
	Name	nationality	Name	compute_peak	memory_capacity	Name	capacity ↑
				in PFLOPs	in TiB		in PiE
1	LANL	US	Trinity	11.08	1,919.03	Lustre	72.83
2	DKRZ	DE	Mistral	3.12	204.00	Lustre02 Lustre01 HPSS	52.00
3	LLNL	US	Sequoia	20.10	1,364.24	Grove	48.85
4	RIKEN	JP	K Computer	10.62	1,136.87	Lustre FEFS	39.77
5	NCAR	USA	Cheyenne	5.33	184.40	HPSS GPFS	37.00
6	NERSC	US	Cori Phase I	4.90	204.00	Lustre	30.00
7	ORNL	US	Titan	27.10	645.74	Spider 2	28.00
8	NCSA	US	Blue Waters	13.40	1,500.00	HPSS Lustre	26.40
9	JCAHPC	JP	Oakforest-PACS	24.91	836.09	Lustre Burst Buffer	24.10
10	CINECA	IT	Marconi A2 Fermi	12.93	413.97	GPFS GPFS	23.71
11	ANL	US	Mira	10.00	698.49	GPFS	21.32
12	JSC	DE	Juqueen	5.90	407.45	HPSS JUST	20.30
13	JAMSTEC	JP	Earth Simulator	1.31	291.04	Home Data Work Archive	19.62
14	KMA	KR	Miri	2.90	0.00	Lustre	19.2
15	NSCC	CN	TaihuLight	125.00	1,191.44	Sunway	17.70
16	AFRL	US	Thunder	5.61	406.54	Lustre	15.54
17	KAUST	SAU	Shaheen II	7.20	718.50	Lustre HPSS	15.28
18	LRZ	DE	SuperMUC Phase 2	3.58	176.44	GPFS	15.00
19	NASA	US	Pleiades	4.97	603.90	Lustre	14.2
20	NSCG	CN	Tianhe-2 Tianhe-1A	59.60	1,169.61	Tianhe-2 H2FS Tianhe-2 Lustre Lustre	14.18
21	TACC	US	Stampede	9.60	245.56	Lustre	12.43
22	ERDC DSRC	US	Topaz	4.57	401.63	Lustre	10.66
23	HLRS	DE	Hazel Hen	7.40	876.75	HPSS Lustre	8.8
24	TEP	FR	Pangea	6.71	49.11	Lustre	8.17
25	GSIC	JP	Tsubame 2.5	5.76	67.67	Lustre	6.93
26	ENI	IT	HPC2	4.60	0.00	GPFS	6.6
27	PGS	US	Abel	5.37	531.14	Lustre	5.3
28	Nagoya University	JP	PRIMEHPC	3.20	83.67	Lustre	5.3
29	ECMWF	UK	Cray XC40	4.25	0.00	HPSS Lustre	5.3
30	ARL	US	Excalibur	3.70	385.63	Lustre	4.09
31	EPCC	UK	Archer	2.55	0.00	Lustre	3.9
32	PNL	US	Cascade	3.40	167.35	Lustre	2.4
33	CSCS	CHE	Piz Daint	7.79	153.70	Lustre	2.2

munity projects.

THE VIRTUAL INSTITUTE FOR I/O

Goals of the Virtual Institute for I/O (VI4IO) are

- Provide a platform for I/O researchers and enthusiasts for exchanging information
- Foster training and international collaboration in the field of high-performance I/O
- Track/encourage the deployment of large storage systems by hosting information about high-performance storage systems

The philosophical cornerstones of VI4IO are:

- Treat contributors/participants equally
- Allow free participation without any fee inclusive to all
- Independent of vendors/research facilities

OPEN ORGANIZATION

The organization uses a wiki as central hub

- Registered users can edit the content
- Mayor changes should be discussed on the contribute mailing list
- Tag clouds link between similar entities
- Supported by mailing lists, e.g.:
 - Call-for-papers

HIGH-PERFORMANCE STORAGE LIST

The High-Performance storage list contains the characteristics of site, supercomputer and connected storage (see the box about the system model). The list shown in the box on the right is sortable on the metric of choice. It allows to add/remove metrics (see the list next). Graphs are created based on selectable grouping.

Metrics: Most metrics can be determined without measurement and describe hardware and software characteristics that should be well known to the site and vendor. A few metrics cover actually observed metadata and I/O performance, in this case the measurement procedure must be clear. The list stores data entered in the wiki into a database and converts data to a base unit.

The following list of supported metrics includes a description:

Institution

- institution: The abbreviation of the institution. Note that systems are linked together based on year and institution.
- year: The year for which the data is valid.
- nationality: The international abbreviation for the nationality of the institution.
- web page: The web page of the institution.
- energy consumption: The overall energy consumption of the datacenter.
- power usage effectiveness: The PUE of the datacenter.

Supercomputer

- institution: see above
- year: see above
- vendor: The vendor of the supercomputer.
- software: A list of keywords with relevant software components, e.g., which file system, parallelization software.
- installation: This is the date when the supercomputer has been installed. Multi-phase installations should appear with their last upgrade date.
- compute peak: The theoretical peak performance in FLOPs. • node count: The number of nodes. • total cores: The total number of available cores • memory capacity: The available memory capacity in Bytes. • memory bandwidth: The sum of the theoretical memory bandwidth available in B/s. • memory per node: The memory capacity per node. • application domain: A list of the main (scientific) domains that use this supercomputer. • applications: A list of the main applications (if known). • energy consumption: The energy consumption of the supercomputer (without storage) in Watts – this does not take the PUE into account. • interconnect: A list of keywords about the interconnect. • processor: A list of keywords specifying the processor. • graph500: The achieved performance according to the graph 500 list. This is not the position in the list, as this may change over time. • graph500 problem scale: according to the graph 500 list. • top500: The achieved performance according to the top 500. • green500: The achieved efficiency according to the green 500. • architecture: A list of keywords covering the system architecture, e.g., i386 64, GPGPU

storage type: all, local storage, shared storage, tape archive, MAIL aggregation: no. sum, avg. max, min – ! removed non-reducable column + energy consumption, power usage effectiveness, initial facility costs, annual staff cos nationality supercomput + node count, total cores, energy consumption, graph500, graph500 problem scale, top500, green500, memory bandwidth, life time annual procurement costs, annual facility update costs, annual too compute peak, memory capacity + energy consumption, drives, cache size, slots, peak, metadata rate, sustained write, sustained read, servers, hdds, ssds, life time annual procurement costs, annual facility update costs, annual too supercomputer storage storage supercompute node count capacity sum sum PiB AFRL, ANL, ARL, ERDC hunder, Mira, Excalibur, Topaz ustre, GPFS, Lustre, Lustre DSRC, LANL, LLNL, NASA rinity, Sequoia, Pleiades, Blue Lustre, Grove, Lustre, HPS NCSA, NERSC, ORNL, PGS Vaters, Cori Phase I, Titan, Abel ustre, Lustre, Spider 2, Lustre PNL. TACC Cascade, Stamped Lustre, Lustre 42338 DKRZ, HLRS, JSC, LR Mistral, Hazel Hen, Juquee tre02, Lustre01, HPSS, HPS Lustre, HPSS, JUST, GPFS SuperMUC Phase 2 96272 95.75 GSIC, JAMSTEC, JCAHF Tsubame 2.5, Earth Simulato Lustre, Home, Data, Work Dakforest-PACS, PRIMEHPC, M Vagova University, RIKE Archive, Lustre, Burst Buffe Lustre, Lustre, FEFS Computer HPSS, GPFS USA 37.00 NCAR Chevenne TaihuLight, Tianhe-2, Tianhe-1A vay, Tianhe-2 H2FS, Tianhe-3 NSCC, NSCG 56960 31.94 Lustre, Lustre 1500 CINECA, ENI GPFS, GPFS, GPFS 30.38 Marconi A2, Fermi, HPC 19.27 KMA Lustre SAU 6174 15.28 KAUST Shaheen II Lustre, HPSS ECMWF, EPCC Cray XC40, Arche HPSS, Lustre, Lustre 9.24 TEP 4608 8.17 Pangea Lustre 5772 CSCS Piz Daint CHE 2.22 Lustre 534808 637.49 adarea metrics + energy_consumption, power_usage_effectiveness, initial_facility_costs, annual_staff_costs

Storage

- institution: see left
- year: see left
- type: The type of the storage, i.e., tape archive/shared storage
- installation: see left
- energy consumption: The energy consumption of the storage part in Watts this does not take the PUE into account.
- capacity: The effective capacity that is available to users. It includes overhead of erasure (RAID) coding and potential hot/cold spares. This value can be easily derived from the number of available storage devices that support the listed file system.
- interconnect: A list of keywords about the interconnect.
- drives: The total number of tape drives for a nearline tape/-MAID archive.
- cache size: The amount of storage cache in a nearline HSM.
- slots: The number of slots in a nearline tape/MAID archive to hold media.
- vendor: The vendor of the storage hardware. • software: A list of keywords specifying the software further. • hardware: A list of keywords specifying the hardware further. • peak: The theoretical peak performance of the storage system. The value is the performance that could theoretically be achieved when transferring data between clients and storage. It is limited by 1) the aggregated network throughput between client and servers, 2) the aggregated (RAID) controller throughput, 3) the network topology. • metadata rate: Metadata throughput. The value can be determined using any I/O benchmark of choice that ensures that client-side and server-side caches are overwhelmed. • sustained write: Best I/O throughput ever measured when accessing files. The read and write values can be determined using any I/O benchmark of choice that ensures that clientside and server-side caches are overwhelmed. • *sustained read*: (see the description for write) • servers: The number of storage servers of the storage system. • hdds: The number of HDDs that belong to the storage system. • ssds: The number of SSDs that belong to the storage system.

- Announcements
- Contributions / suggestions

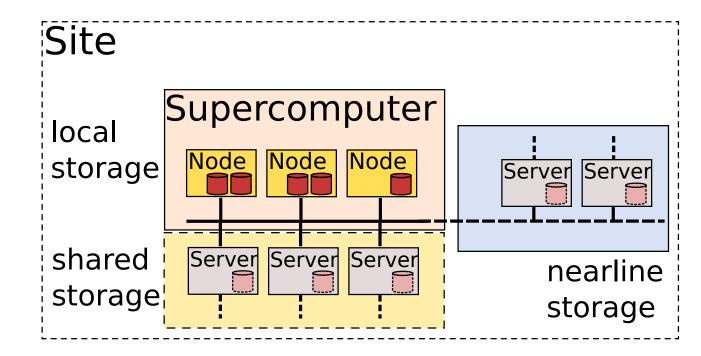
HPSL System Model

The system model describes how characteristics are assigned to components. Storage is difficulty to assign to a single component as it is often shared across supercomputers, therefore, a component based model is used.

Supported components:

- Site: Describes the facility
- Supercomputer: A system
- Storage (shared, local or nearline storage)

Conceptual example:



The web page allows the creation of a topology for the facility to indicate the relation between the components. An example for DKRZ: **Measurement procedure for sustained performance:** Compared to other lists (TOP500, Green500) that have a clear measurement process, the rules for determining sustained performance for the HPSL are relaxed due to the complexity of I/O benchmarks, However it must be clarified how the measurement has been conducted.

IO-500 EFFORT

We are discussing the creation of a benchmark to compare facilities and storage systems. This challenge is explored on our task page: https://www.vi4io.org/std/io500 and mailing list.

Goals for the benchmark:

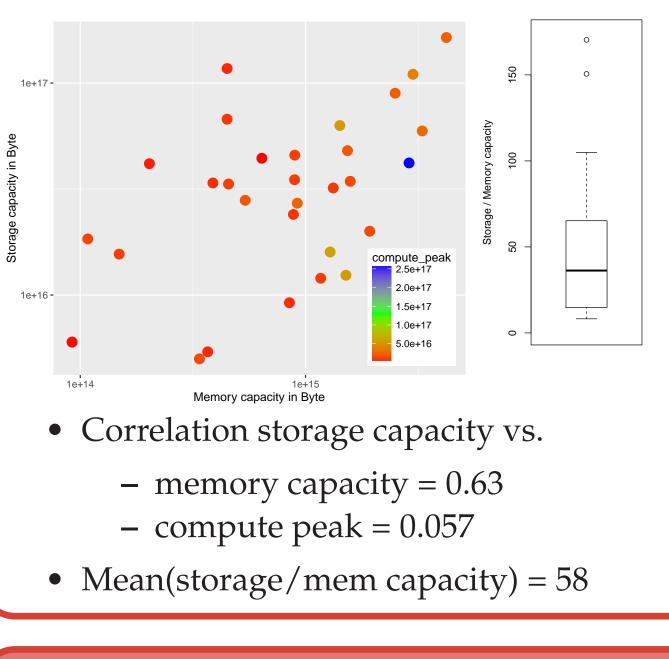
• Capture user-experienced performance



- ality, web_page
- supercomputer: vendor, software, installation, application_domain, applications, interconnect, processor, architecture, memory_per_node
 storage: type, installation, interconnect, vendor, software, hardware
 filter by: US, DE, JP, USA, CN, IT, KR, SAU, UK, FR, CHE, none

DERIVED ANALYSIS

With the collected data many in-depth analysis becomes possible, for example, the relationship between storage and memory capacity:



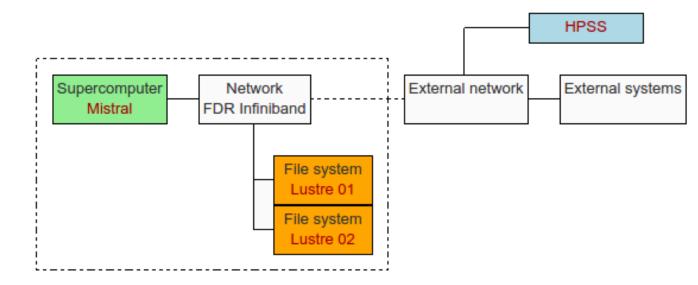
ONGOING WORK

Support standardization efforts
 – IO-500 benchmark

Site characteristics

Кеу	Value	Rf				
Institution	DKRZ					
Nationality	DE					
Web page	https://www.dkrz.de/					
Energy consumption	2.00 MW					
Power usage effectiveness	1.04					

System architecture



DKRZ hosts the Mistral supercomputer which is tightly coupled with two Lustre file systems. We have some small compute and supporting infrastructure that may access the storage of Mistral and a large HPSS system.

Supercomputers

Mistral

Storage systems

- HPSS
- Lustre02
- Lustre01

- Reported performance is representative for:
 - IOEasy: Applications with well optimized I/O patterns
 - IOHard: Applications that require a random workload
 - IOMD: Usage that depends on metadata/small objects

Challenges:

- Representative: for optimized, naive I/O heavy workloads; and small objects
- Inclusive: cover various storage technology and non-POSIX APIs
- Trustworthy: representative results and prevent cheating
- Cheap: easy to run and short benchmarking time (in the order of minutes)

Strategy:

- Build on existing benchmarks
- Plugin systems should allow for alternative storage technology
- Start by reporting one metric per benchmark, decide later about a single number



- Lossy compression interfaces
- IO-500 agenda:
 - June'17, proposal for benchmark
 - Benchmark runs on Top-500 sites
 - Nov'17, SC presentation of results
- Extending benchmarks, HPSL sitesSupport training and teaching for storage

VI4IO AND YOU

Content is under open licenses. You are welcome to join the mailing lists or participate!



https://vi4io.org