

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. It hosts the **comprehensive data center list (CDCL)**. Similarly to the top500, it contains information about supercomputers and their storage systems.

I/O benchmarking, particularly, the intercomparison of measured performance between sites is tricky as there are more hardware components involved and configurations to take into account. Therefore, together with the community, we standardized an HPC I/O benchmark, the **IO-500** benchmark, for which the first list had been released during supercomputing in Nov. 2017. Such a benchmark is also useful to assess the impact of system issues like the **Melt-down** and **Spectre* bugs**.

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

IO-500 EFFORT

Together with the community, we created the IO-500 benchmark to compare storage systems.

Goals for the benchmark:

- Capture user-experienced performance
- Reported performance is representative for:
 - IOEasy: Applications with well optimized I/O patterns
 - IOHard: Applications that require a random workload
 - MDEasy: Metadata/small objects
 - MDHard: Small files (3901 bytes) in a shared directory
 - Find: Finding relevant objects based on patterns

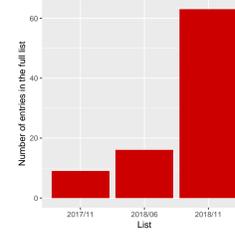


Fig. 1: Growth of the list

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)

Benefit for the community beyond the IO-500:

- Support the development of benchmarks that are used (IO-500 builds on standard benchmarks)
- Feed back best practices of tool usage (e.g., find) and benchmarks
- Aid detailed comparison of individual system characteristics while having a **ranked list**
- Share best-practices to obtain good performance

IO-500 LIST NOV 2018

There are several lists available, the **full list** contains all the results submitted for comparison while entries can enter a **ranked list** upon user choice and only one solution per system.

There are several ranked lists with awards that stimulate aspects of I/O system development:

- The **IO-500 award** for the fastest system in the ranked IO-500 list
- The **10 node challenge** fosters performance for small-scale runs

Further awards may follow.

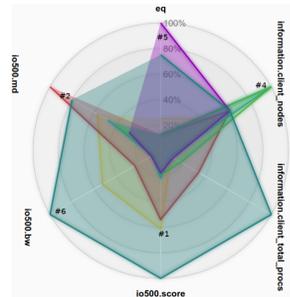
The ranked IO-500 list:

#	Information				IO500	
	institution	system	storage vendor	filesystem type	client nodes	total score
1	Oak Ridge National Laboratory	Summit	IBM Spectrum Scale	LuFS	1008	330.56
2	University of Cambridge	Data Accelerator	Dell EMC	LuFS	528	42.24
3	Korea Institute of Science and Technology Information (KISTI)	NURION	DDN	IME	2048	4096
4	JCAHPC	Oakforest-PACS	DDN	IME	2048	16384
5	WekaIO	WekaIO	WekaIO		17	935
6	KAUST	Shahereh	Cray	DataWarp	1024	8192
7	University of Cambridge	Data Accelerator	Dell EMC	BeefFS	184	5888
8	Google	Exascale on GCP	Google	LuFS	120	960
9	JCAHPC	Oakforest-PACS	DDN	IME	256	8192
10	KAUST	Shahereh	Cray	LuFS	1000	16000
11	JSC	JURON	Thinkstor	BeefFS	8	64
12	DKRZ	Melita	Seagate	LuFS	100	1000

Flexible equations It supports equations to compute derived metrics, here `mdtest.easy_write / client_nodes` for the full list:

#	evaluation	institution	system	storage vendor	filesystem type	client nodes	client total score	data score	io500
1	24.17	JSC	JURON	Thinkstor	BeefFS	8	64	35.77	14.24
2	19.97	DDN	BeefFS	DDN	LuFS	10	240	31.50	6.20
3	11.34	Clemson University	obdev	Dell	BeefFS	16	84	12.51	2.10
4	10.36	Clemson University	obdev	Dell	BeefFS	16	128	11.01	1.90
5	8.48	Clemson University	obdev	Dell	BeefFS	10	80	10.17	2.32
6	6.87	WekaIO	WekaIO	WekaIO		10	700	58.25	27.05
7	6.09	Oak Ridge National Laboratory	Summit	IBM	Spectrum Scale	504	1008	330.56	1238.93
8	5.72	IBM	Spectrum Scale	IBM	Spectrum Scale	10	10	24.24	4.57
9	5.27	WekaIO	WekaIO	WekaIO		17	935	78.37	37.39
10	5.19	Clemson University	obdev	Dell	BeefFS	10	30	13.12	3.17

The system displays a net graph for further distinguishing the best systems:



As we can see, this can be used to create arbitrary new rankings and investigate the data.

All results are available The individual submission scripts and results for the benchmarks are preserved and can be accessed. The data is also available as CSV file for offline analysis.

DATA CENTER LIST

The comprehensive data center list with its system model describes how characteristics are assigned to components. Storage is difficult to assign to a single component as it is often shared across supercomputers, therefore, a flexible component based model is used.

Supported components:

- Site: Describes the facility
- Supercomputer: A system
- Storage system
- Nodes
- Network
- Building

The schema is under active development – we aim to describe data center characteristics. The web page allows the creation of a topology for the facility to indicate the relation between the components – ultimately multiple views will be created to show, e.g.:

- Logical network connectivity
- Physical layout in racks
- Building organization

Metrics: Most metrics can be determined without measurement and describe hardware and software characteristics that should be 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 defined. The list stores data entered in the wiki into a database and converts data to a base unit.

The following is an example of the schema for the DKRZ system:

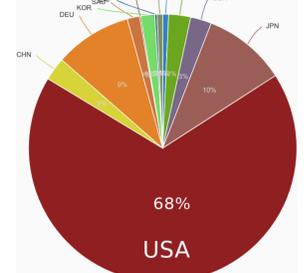
The rules for determining performance are relaxed due to the complexity of I/O measurements, but this is augmented by the IO-500.

HPSL 2019

The current list contains 43 sites:

#	site.institution	site.storage.system.net.capacity.in.PiB	site.supercomputer.compute.peak.in.PFLOPS	site.supercomputer.memory.capacity.in.TB
1	Oak Ridge National Laboratory	250.04	220.64	3511.66
2	National Energy Research Scientific Computing Center	197.65	37.71	887.53
3	Lawrence Livermore National Laboratory	72.83	11.08	2193.93
4	German Climate Computing Center	52.00	3.89	683.80
5	Lawrence Livermore National Laboratory	48.85	20.10	1800.00
6	RKIN Advanced Institute for Computational Science	30.77	16.62	1250.00
7	National Center for Atmospheric Research	37.00	5.33	202.75
8	National Center for Supercomputing Applications	27.80	13.40	1649.27
9	Global Scientific Information and Computing Center	25.84	17.89	275.98
10	Aust Center for Advanced HPC	24.10	24.91	919.29

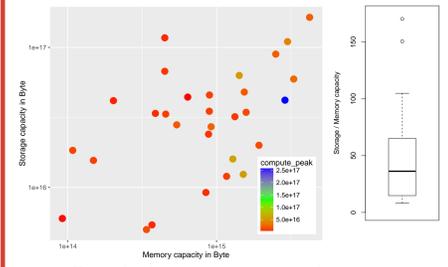
Various views are possible, moreover, it supports flexible data aggregation, e.g., capacity by country:



site.nationality	site.storage.system.net.capacity.in.PiB	site.supercomputer.compute.peak.in.PFLOPS	site.supercomputer.memory.capacity.in.TB	site.supercomputer.nodes
USA	776.95	358.90	13728.07	344706
JPN	114.66	57.93	2897.27	98254
DEU	189.50	23.05	2595.32	48232
CHN	32.16	184.80	2596.00	56990
ITA	30.38	17.53	455.17	1500
GBR	29.02	13.05	371.01	485
KOR	19.27	2.90	0.00	0
SAU	16.96	7.20	790.00	6174
FRN	6.17	6.71	64.00	4008
CHE	7.73	28.32	631.00	6751
ZAF	2.98	0.00	8.86	100
AUT	1.81	0.68	42.18	0

DERIVED ANALYSIS

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



- Correlation storage capacity vs.
 - memory capacity = 0.63
 - compute peak = 0.057
- Mean(storage/mem capacity) = 58

ONGOING WORK

- IO-500:
 - Clarified execution rules
 - Procedures to adapt IO-500
 - Integration of optional benchmarks
 - Continuous integration deployment including performance regression
 - Finalize vendor engagement program
- VI4IO standardization efforts
 - Data center representation
 - Next-generation interfaces (NGI)
- CDCL list
 - Extended schema and alternative views
 - More CDCL sites
 - Better link between IO-500 and CDCL
- Support training and teaching for storage

VI4IO, IO-500, AND YOU

You are welcome to join the mailing lists or our slack channel and participate!



Join us on Slack:

The content is under open licenses.

More details on:

- <https://vi4io.org>
- <http://io-500.org>

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
 - Announcements
 - Contributions / suggestions

COMMUNITY CONTENT

The wiki covers A) worldwide research groups that address high-performance I/O including:

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

Everyone is welcome to add (own) group(s)!

B) Relevant I/O related tools and benchmarks

C) Comprehensive Data Center List (see the other boxes)