



ISC

High Performance
The HPC Event.

HPC IODC Workshop

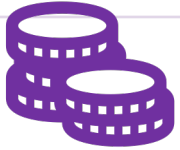
HPC on Public Clouds Opportunities, Challenges & Learnings

Vinay Gaonkar, Co-founder, Kmesh, Inc
vinay@kmesh.io



Kmesh.io

Public Clouds Provide Attractive Infrastructure for HPC



Lower Cost to Innovate:

Inexpensive to try new things and innovate



Flexibility:

Newer GPU, amounts of shared memory, and processor



Choice:

Multiple node types as needed



On-Demand:

Spin up only as needed



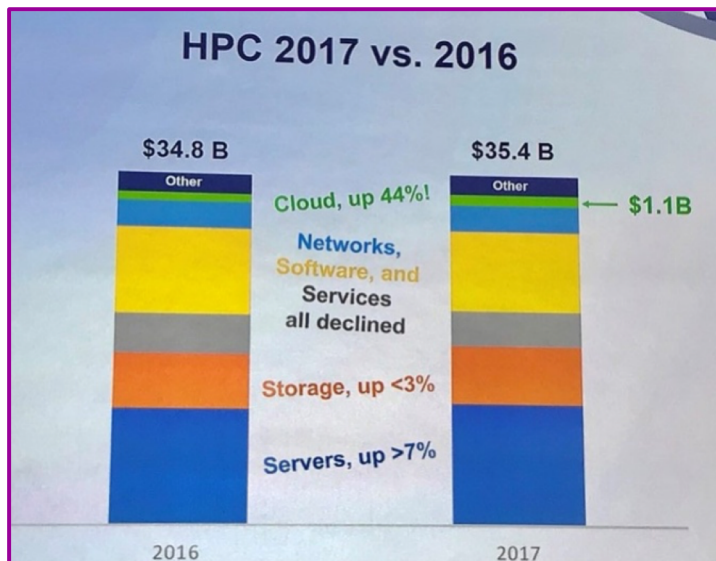
Extendable:

Bursting as needed
POCs



Top-Performing Technologies:

Latest Intel CPUs or NVIDIA GPUs



64 percent of HPC sites now run at least some of their work in public cloud.

– Hyperion Research
(<https://www.top500.org/news/cloud-computing-in-hpc-surges/>)

Public Cloud Opportunity : Attractive Cost Model

On-Prem vs Cloud Comparison(2015)¹

Location	Cost (Per Core Hour)
On Prem (<i>P/C 100% utilized and servers 100% utilized</i>)	\$0.014
On Prem (<i>P/C 50% utilized and servers 100% utilized</i>)	\$0.021
On Prem (<i>P/C 50% utilized and servers 85% utilized</i>)	\$0.023
AWS On Demand	\$0.0419
AWS Reserved 3 yr	\$0.0237
AWS Spot Instance	\$0.0108
GCP On Demand	\$0.03545
GCP with Sustained Use Discount	\$0.024815
GCP Preemptable Instance	\$0.0075

¹ <https://www.internet2.edu/blogs/detail/14114>

Public Cloud Opportunity : Attractive Cost Model

On-Prem (/core/hour)

ReScale ²	Northwestern University ³	University of Maine ⁴
\$0.12 \$100,000 /month for 1600 cores	\$0.04 (Generic HPC) \$0.16 (4xP100 Tesla GPU)	\$0.42 (HPC node)

Public Clouds

Cloud	Instance/VM	Processor	CPU	Memory (GB)	NW performance (Gbps)	Storage (GB)	Total Storage Capacity (TB)	Cost/Hour	Cost/CPU core/hour	Monthly Cost (1600 cores)
AWS	i3.metal	Intel Xeon Broadwell	72	512	25	8x1900 NVMe (included)	15.2	\$4.99	\$0.07	\$159,744.00
AWS	p3dn.24xlarge	Intel® Xeon® Skylake 8xNVIDIA Tesla V100 GPUs	96	768	100	1800 (Included)	1.8	\$31.22	\$0.33	\$374,616.00
Azure	Standard_HC44rs	Intel Xeon Platinum 8168	44	352		Networked SSDs		\$4.35	\$0.10	\$113,883.05
Azure	NC24 v3	Intel Xeon Broadwell 4xNVIDIA Tesla V100 GPUs	24	448		Networked SSDs		\$16.81	\$0.70	\$806,668.80

¹ <https://resources.rescale.com/the-real-cost-of-high-performance-computing/>

² <https://www.it.northwestern.edu/bin/docs/research/quest-full-access.pdf>

³ <https://acg.umaine.edu/pricing/fee-structure/>

Public Cloud Challenge: 'Infinite' Cloud Resources

Multiple Clouds

Which cloud?



Plethora of Instances

Which instance(s)?



Sizes for Linux virtual machines in Azure

11/13/2018 • 2 minutes to read • Contributors all

This article describes the available sizes and options for the Azure virtual machines you can use to run your Linux apps and workloads. It also provides deployment considerations to be aware of when you're planning to use these resources. This article is also available for [Windows virtual machines](#).

Type	Sizes	Description
General purpose	B, Dsv3, Dv3, DSv2, Dv2, Av2, DC	Balanced CPU-to-memory ratio. Ideal for testing and development, small to medium databases, and low to medium traffic web servers.
Compute optimized	Fsv2, Fs, F	High CPU-to-memory ratio. Good for medium traffic web servers, network appliances, batch processes, and application servers.
Memory optimized	Esv3, Ev3, M, GS, G, DSv2, Dv2	High memory-to-CPU ratio. Great for relational database servers, medium to large caches, and in-memory analytics.
Storage optimized	Lsv2, Ls	High disk throughput and IO ideal for Big Data, SQL, NoSQL databases, data warehousing and large transactional databases.
GPU	NV, NVv2, NC, NCv2, NCv3, ND, NDv2 (Preview)	Specialized virtual machines targeted for heavy graphic rendering and video editing, as well as model training and inferencing (ND) with deep learning. Available with single or multiple GPUs.
High performance compute	H	Our fastest and most powerful CPU virtual machines with optional high-throughput network interfaces (RDMA).

Many Storage Options

Which storage(s)?



Block Blobs Scalable object storage for documents, videos, pictures, and unstructured text or binary data. Choose from Hot, Cool, or Archive tiers. Prices for locally redundant storage (LRS) Archive Block Blob start from: \$0.002 /GB per month See Pricing >	Azure Data Lake Storage Combines the power of a Hadoop compatible file system with integrated hierarchical namespace with the massive scale and economy of Azure Blob Storage to help speed your transition from proof of concept to production. Prices for LRS start from: \$0.001 /GB per month See Pricing >
Managed Disks Persistent, secured disks that support simple and scalable virtual machine deployment. Designed for 99.999% availability. Choose Premium (SSD) Disks for low latency and high throughput. Prices for Standard Managed Disks start from: \$1.54 per month See Pricing >	Files Fully managed file shares in the cloud, accessible via standard Server Message Block (SMB) protocol. Enables sharing files between applications using Windows APIs or REST API. Prices for LRS File storage start from: \$0.060 /GB per month See Pricing >

Instances Comparison : IOR Tests

Specs are only guidance. Real performance matter!

← Specification → Real Performance →

Model	vCPU	Mem (GiB)	Networking Performance (Gbps)	Ephemeral	Ephemeral Storage	Dedicated EBS Bandwidth (Mbps)	\$ Per Hour	\$/GiB (Monthly)	RecordSize	Read MB/s	Write MB/s
i3.xlarge	4	30.5	Up to 10	1 x 1.9 NVMe SSD	1946	850	\$0.31	0.12	Lustre=128K ZFS=128K	601	440
i3.2xlarge	8	61	Up to 10	1 x 1.9 NVMe SSD	1946	1,700	\$0.62	0.24		621	620
c5d.4xlarge	16	32	Up to 10	1 x 400 NVMe SSD	400	3,500	\$0.77	1.44		837	387
m5d.4xlarge	16	64	Up to 10	2 x 300 NVMe SSD	600	3,500	\$0.90	1.13		621	302
i3.4xlarge	16	122	Up to 10	1 x 1.9 NVMe SSD	3891	850	\$1.25	0.24		621	622
i3.xlarge	4	30.5	Up to 10	1 x 1.9 NVMe SSD	1946	850	\$0.31	0.12	Lustre=1MBK ZFS=128K	597	438
i3.2xlarge	8	61	Up to 10	1 x 1.9 NVMe SSD	1946	1,700	\$0.62	0.24		621	620
c5d.4xlarge	16	32	Up to 10	1 x 400 NVMe SSD	400	3,500	\$0.77	1.44		840	387
m5d.4xlarge	16	64	Up to 10	2 x 300 NVMe SSD	600	3,500	\$0.90	1.13		621	303
i3.4xlarge	16	122	Up to 10	1 x 1.9 NVMe SSD	3891	850	\$1.25	0.24		621	622

of slots = 8

- Instance cost != better IOR performance
- Write performance depends on the storage size
- Larger ZFS transfer size is better for IOR BW. Optimal ZFS transfer = 128k

Instances Comparison : MDTEST

Model	vCPU	Memory (GiB)	\$/Hour	Directory creation	Directory stat	Directory removal	File creation	File stat	File read	File removal	Tree creation	Tree removal	CPU Utilization
i3.xlarge	4	30.5	\$ 0.31	4,827.42	17,709.11	6,363.59	6,660.80	9,780.99	9,267.85	9,737.73	2,276.53	2,601.46	77.10%
i3.2xlarge	8	61	\$ 0.62	5,876.86	19,147.13	8,928.47	8,168.89	10,642.64	10,460.31	11,516.76	2,412.13	2,570.56	67.70%
c5d.4xlarge	16	32	\$ 0.77	4,034.69	30,153.35	14,768.59	13,628.94	16,969.22	17,973.16	18,410.46	4,152.59	4,449.19	19.60%
m5d.4xlarge	16	64	\$ 0.90	9,861.89	25,202.48	12,248.51	11,300.76	13,988.44	13,963.48	15,866.89	3,431.61	3,821.88	32.80%
i3.4xlarge	16	122	\$ 1.25	6,990.27	19,614.65	8,810.93	8,400.30	10,692.49	10,716.49	11,489.88	2,530.25	2,685.22	20.40%

Command Used: `mpiexec <host> /usr/bin/mdtest -l 10 -i 2 -z5 -b 5 -i 3 -d /mnt/kmesh/mdtest5/`

Cross-cloud Comparison : Instances with NVMe

VDbench Comparison

		i3.2xlarge (AWS)			L16s_v2 (Azure)			%Diff		
		8 vCPU	61GiB	1x1900NVMe SSD	16vCPU	128Gib	2x1900GiB NVMe SSD			
		\$0.624/hr			\$1.248/hr					
	Transfer Size (Bytes)	Total IOPS	Total MiB/s	Latency (MS)	Total IOPS	Total MiB/s	Latency (MS)	IOPS	MiB/s	Latency
100% Read	4K	49,024.80	191.50	1.32	36,650.50	238.30	2.59	-25%	24%	196%
0% Read	4K	11,129.80	43.47	5.79	8,483.60	54.04	10.20	-24%	24%	176%
100% Read	32K	34,292.60	1,071.64	1.89	18,147.50	958.76	5.68	-47%	-11%	300%
0% Read	32K	19,004.20	593.89	3.34	12,367.00	646.21	7.87	-35%	9%	235%
100% Read	128K	16,148.70	2,018.58	4.11	9,242.80	1,950.44	11.09	-43%	-3%	270%
0% Read	128K	5,862.70	732.84	10.77	3,745.60	770.21	24.02	-36%	5%	223%
100% Read	256K	9,263.70	2,315.92	6.94	5,607.70	2,333.38	16.98	-39%	1%	245%
0% Read	256K	4,383.10	1,095.79	14.30	2,152.30	874.68	39.62	-51%	-20%	277%

Cross-cloud Comparison : Instances with Persistent Drives

VDBench comparison

			EBS:c5.4xlarge (AWS) 16 vCPU 32 GiB Memory \$0.680/hour			PD:Standard F16s(Azure) 16 vCPU 32 GiB Memory \$0.796/hour			Diff %		
R/W Ratio	Seek %	Transfer Size (Bytes)	Total IOPS	Total MiB/s	Latency (MS)	Total IOPS	Total MiB/s	Latency (MS)	Total IOPS	Total MiB/s	Latency (MS)
100% Read	Random	4K	7385.4	28.85	4.327	3035.5	11.86	10.521	-59%	-59%	243%
100% Read	50	4K	5298.3	20.7	6.034	2158.1	8.43	14.812	-59%	-59%	245%
100% Read	Sequential	4K	83160.3	324.84	0.384	24244.1	94.7	1.323	-71%	-71%	345%
75% Read	Random	4K	6303.8	24.62	5.07	2728.1	10.66	11.709	-57%	-57%	231%
75% Read	50	4K	5666.8	22.14	5.64	2404.6	9.39	13.288	-58%	-58%	236%
75% Read	Sequential	4K	55336.6	216.16	0.577	36259	141.64	0.88	-34%	-34%	153%
50% Read	Random	4K	5642.1	22.04	5.664	2595.7	10.14	12.308	-54%	-54%	217%
50% Read	50	4K	6312.8	24.66	5.062	2735.6	10.69	11.678	-57%	-57%	231%
50% Read	Sequential	4K	51685.5	201.9	0.617	19642.7	76.73	1.622	-62%	-62%	263%
0% Read	Random	4K	3717.5	14.52	8.599	2162.4	8.45	14.776	-42%	-42%	172%
0% Read	50	4K	6757.7	26.4	4.727	3232.7	12.63	9.876	-52%	-52%	209%
0% Read	Sequential	4K	23989.9	93.71	1.329	9932.1	38.8	3.213	-59%	-59%	242%
100% Read	Random	32K	7186.6	224.58	4.447	2146.5	67.08	14.89	-70%	-70%	335%
100% Read	50	32K	3748.7	117.15	8.531	1374.6	42.96	23.262	-63%	-63%	273%
100% Read	Sequential	32K	12354.7	386.08	2.589	2900.2	90.63	11.018	-77%	-77%	426%
75% Read	Random	32K	7787.8	243.37	4.1	2805.5	87.67	11.385	-64%	-64%	278%
75% Read	50	32K	3602.3	112.57	8.873	1470	45.94	21.737	-59%	-59%	245%
75% Read	Sequential	32K	10752.1	336	2.972	6563.9	205.12	4.869	-39%	-39%	164%
50% Read	Random	32K	7074.3	221.07	4.512	3615.4	112.98	8.826	-49%	-49%	196%
50% Read	50	32K	3991.5	124.73	8.005	1677.3	52.41	19.052	-58%	-58%	238%
50% Read	Sequential	32K	10928.2	341.51	2.922	5197.5	162.42	6.149	-52%	-52%	210%
0% Read	Random	32K	5505.1	172.03	5.796	3138.7	98.09	10.164	-43%	-43%	175%
0% Read	50	32K	7130.7	222.84	4.471	3852.4	120.39	8.276	-46%	-46%	185%
0% Read	Sequential	32K	13303.6	415.74	2.393	7543.9	235.75	4.227	-43%	-43%	177%
100% Read	Random	128K	2570.1	321.26	12.445	1297	162.12	24.654	-50%	-50%	198%
100% Read	50	128K	1793.4	224.18	17.837	949.3	118.66	33.692	-47%	-47%	189%
100% Read	Sequential	128K	3300.5	412.56	9.693	1744.8	218.1	18.33	-47%	-47%	189%
75% Read	Random	128K	2801.3	350.16	11.405	1593.4	199.18	20.053	-43%	-43%	176%
75% Read	50	128K	2132.3	266.54	14.989	1124.9	140.62	28.418	-47%	-47%	190%
75% Read	Sequential	128K	3282.9	410.36	9.733	1915.4	239.42	16.687	-42%	-42%	171%
50% Read	Random	128K	2536.5	317.06	12.589	1484.9	185.62	21.512	-41%	-41%	171%
50% Read	50	128K	2187	273.38	14.603	1241.7	155.21	25.735	-43%	-43%	176%
50% Read	Sequential	128K	3287.8	410.98	9.711	1961.1	245.13	16.289	-40%	-40%	168%
0% Read	Random	128K	2251.4	281.42	14.173	1100.3	137.54	29.032	-51%	-51%	205%
0% Read	50	128K	2588.8	323.6	12.322	1224.9	153.11	26.074	-53%	-53%	212%
0% Read	Sequential	128K	3336.8	417.1	9.553	1591.5	198.94	20.061	-52%	-52%	210%

CPU utilization matters

% CPU Idle c5.4xlarge (AWS)	% CPU Idle F16s (Azure)	CPU Idle Relative Diff
79.73	89.08	9.35
73.29	80.72	7.43
72.69	80.50	7.81
76.25	80.81	4.56
79.77	88.31	8.54
73.13	81.26	8.13
72.62	79.29	6.67
82.54	83.14	0.60
79.33	86.89	7.56
72.73	79.59	6.86
71.68	78.13	6.45
83.85	83.49	-0.36

Cloud infrastructure is different.

1

CPU Cycles : Generally under-utilized
Networking & storage : Over-subscribed

Inter-Instance traffic and “networked” storage performance depends heavily on networking

2

Performance specs assume extreme ideal conditions

Noisy neighbor is a real problem on the cloud
Higher price does not translate to better performance

3

Bare-metal != “physical systems”

Networking - shared and over-subscribed
Networked storage - bottleneck

Public Cloud Challenge: No Comprehensive Benchmark

IO-500

This is an intermediate list based on corrected calculations¹⁾ for the IO-500 ranked list²⁾ from November 2018 (from SC 2018).

Please see also the 10 node challenge ranked list.

The list shows the best result for a given combination of system/institution/filesystem.

#	information							io500		
	institution	system	storage vendor	filesystem type	client nodes	client total procs	data	score	bw	md
									GiB/s	kIOP/s
1	Oak Ridge National Laboratory	Summit	IBM	Spectrum Scale	504	1008	zip	330.56	88.20	1238.93
2	University of Cambridge	Data Accelerator	Dell EMC	Lustre	528	4224	zip	158.71	71.40	352.75
3	Korea Institute of Science and Technology Information (KISTI)	NURION	DDN	IME	2048	4096	zip	156.91	554.23	44.43
4	JCAHPC	Oakforest-PACS	DDN	IME	2048	16384	zip	137.78	560.10	33.89
5	WekaIO	WekaIO	WekaIO		17	935	zip	78.37	37.39	164.26
6	KAUST	ShaheenII	Cray	DataWarp	1024	8192	zip	77.37	496.81	12.05
7	University of Cambridge	Data Accelerator	Dell EMC	BeeGFS	184	5888	zip	74.58	58.81	94.57
8	Google	Exascaler on GCP	Google	Lustre	120	960	zip	47.23	23.06	96.74
9	JCAHPC	Oakforest-PACS	DDN	Lustre	256	8192	zip	42.18	20.04	88.78
10	KAUST	ShaheenII	Cray	Lustre	1000	16000		41.00*	54.17	31.03*
11	JSC	JURON	ThinkparQ	BeeGFS	8	64		35.77*	14.24	89.81*
12	DKRZ	Mistral	Seagate	Lustre	100	1000		32.15	22.77	45.39
13	DDN	Bancholab	DDN	Lustre	10	240	zip	31.50	6.33	156.69
14	IBM	Sonasad	IBM	Spectrum Scale	10	10	zip	24.24	4.57	128.61
15	Clemson University	ofsdev	Dell	BeeGFS	32	32	zip	22.02	8.55	56.68
16	Fraunhofer	Seislab	ThinkparQ	BeeGFS	24	24		16.96	5.13	56.14
17	Joint Institute for Nuclear Research	Govorun	RSC	Lustre	24	192	zip	12.08	3.34	43.65
18	PNNL	EMSL		Lustre	126	252		11.12	4.88	25.33

*** Google & DDN #8 ***

```
io500_info_data_storage_type=NVMe
io500_info_filesystem=Lustre
io500_info_filesystem_vendor=DDN
io500_info_filesystem_version='Exascaler 5.0 PRE - Lustre 2.11-56 (Master)'
io500_info_institute_name=Google
io500_info_metadata_storage_type=NVMe
io500_info_num_client_nodes=120
io500_info_num_data_server_nodes=120
io500_info_num_data_storage_devices=480
io500_info_num_metadata_server_nodes=4
io500_info_num_metadata_storage_devices=16
io500_info_procs_per_node=8
io500_info_storage_age_in_months=0
io500_info_storage_install_date=11/18
io500_info_storage_interface=NVMe
io500_info_storage_network='16Gbps Ethernet'
io500_info_system_name='Lustre on Google Cloud Platform'
io500_info_whatever='Medium Size configuration part of development between Google and DDN'
io500_ior_cmd=/scratch/io-500-dev-master/bin/ior
io500_ior_easy_params='-t 2048k -b 60000m -F'
io500_ior_easy_size=60000
io500_ior_hard_other_options=
io500_ior_hard_writes_per_proc=42000
io500_mdreal_cmd=/scratch/io-500-dev-master/bin/md-real-io
io500_mdreal_params='-P=5000 -I=1000'
io500_mdtest_cmd=/scratch/io-500-dev-master/bin/mdtest
```

<https://www.vi4io.org/io500/start>

IO-500

IO⁵⁰⁰

This is an intermediate list based on corrected calculations¹⁾ for the IO-500 ranked list²⁾ from November 2018 (from SC 2018).

Please see also the 10 node challenge ranked list.

The list shows the best result for a given combination of system/institution/filesystem.

#	information							io500		
	institution	system	storage vendor	filesystem type	client nodes	client total procs	data	score	bw	md
									GiB/s	kIOP/s
1	Oak Ridge National Laboratory	Summit	IBM	Spectrum Scale	504	1008	zip	330.56	88.20	1238.93
2	University of Cambridge	Data Accelerator	Dell EMC	Lustre	528	4224	zip	158.71	71.40	352.75
3	Korea Institute of Science and Technology Information (KISTI)	NURION	DDN	IME	2048	4096	zip	156.91	554.23	44.43
4	JCAHPC	Oakforest-PACS	DDN	IME	2048	16384	zip	137.78	560.10	33.89
5	WekaIO	WekaIO	WekaIO		17	935	zip	78.37	37.39	164.26
6	KAUST	ShaheenII	Cray	DataWarp	1024	8192	zip	77.37	496.81	12.05
7	University of Cambridge	Data Accelerator	Dell EMC	BeeGFS	184	5888	zip	74.58	58.81	94.57
8	Kmesh	Kmesh on AWS	Kmesh	Lustre	84	224	zip	56.77	16.61	17.31
9	JCAHPC	Oakforest-PACS	DDN	Lustre	256	8192	zip	42.18	20.04	88.78
10	KAUST	ShaheenII	Cray	Lustre	1000	16000		41.00*	54.17	31.03*
11	JSC	JURON	ThinkparQ	BeeGFS	8	64		35.77*	14.24	89.81*
12	DKRZ	Mistral	Seagate	Lustre	100	1000		32.15	22.77	45.39
13	DDN	Bancholab	DDN	Lustre	10	240	zip	31.50	6.33	156.69
14	IBM	Sonasad	IBM	Spectrum Scale	10	10	zip	24.24	4.57	128.61
15	Clemson University	ofsdev	Dell	BeeGFS	32	32	zip	22.02	8.55	56.68
16	Fraunhofer	Seislab	ThinkparQ	BeeGFS	24	24		16.96	5.13	56.14
17	Joint Institute for Nuclear Research	Govorun	RSC	Lustre	24	192	zip	12.08	3.34	43.65
18	PNNL	EMSL		Lustre	126	252		11.12	4.88	25.33

* Kmesh's #8+ AWS Run *

```
io500_info_data_storage_type=NVMe
io500_info_filesystem=Lustre
io500_info_filesystem_vendor=Kmesh
io500_info_filesystem_version='Kmesh 1.0 - Lustre 2.11'
io500_info_institute_name=Kmesh
io500_info_metadata_storage_type=NVMe
io500_info_num_client_nodes=25
io500_info_num_data_server_nodes=25
io500_info_num_data_storage_devices=159
io500_info_num_metadata_server_nodes=25
io500_info_num_metadata_storage_devices=25
io500_info_procs_per_node=8
io500_info_storage_age_in_months=0
io500_info_storage_install_date=04/18
io500_info_storage_interface=NVMe
io500_info_storage_network='25Gbps Ethernet'
io500_info_system_name='Lustre on Amazon Web Services(AWS)'
io500_info_whatever='Extrapolation of a small config of Kmesh on AWS'
```

+ Not submitted yet

Are these different?

Kmesh #8 Run

Instance/Region	Number of MDS	Number of OSS	Number of Storage Devices (MDS)	Number of Storage Devices (OSS)	Client Number of Processors	Score BW	Score IOPS	Total Score	Cost/hour	Total Cost
i3.metal MDT0/OOS0..MDT25/OSS25 AWS us-east-1 region	25	25	25	149	224	19	17	56.77	\$175.23	\$246.13

\$250

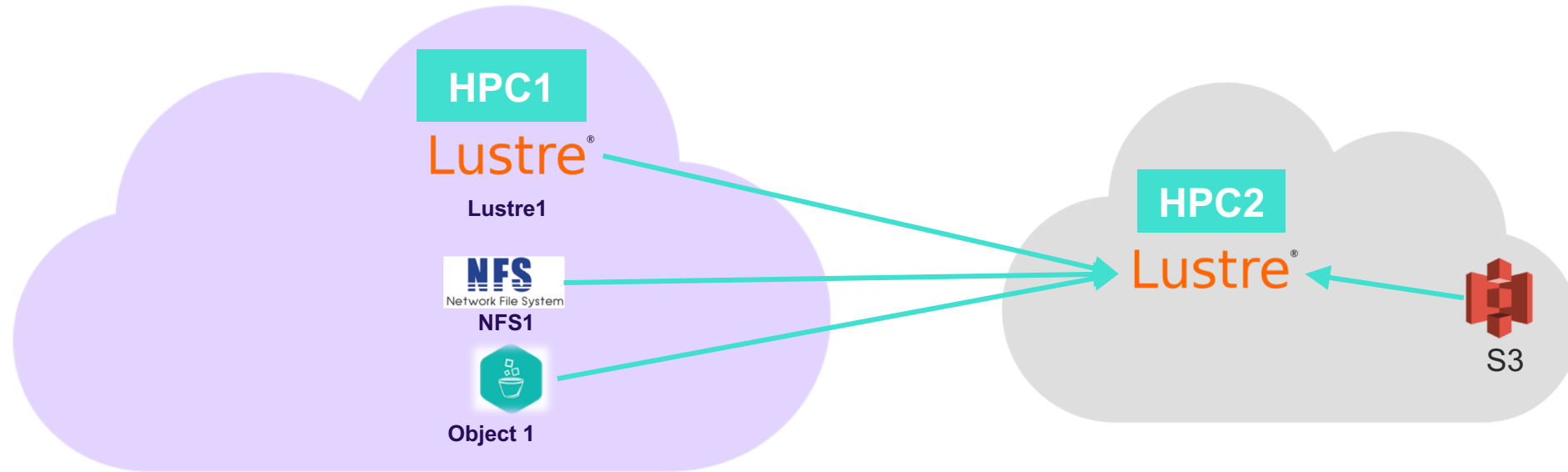
Google & DDN #8

- ✓ No cost information
- ✓ No details on VM/Instance information
- ✓ No GCP Region information

\$?

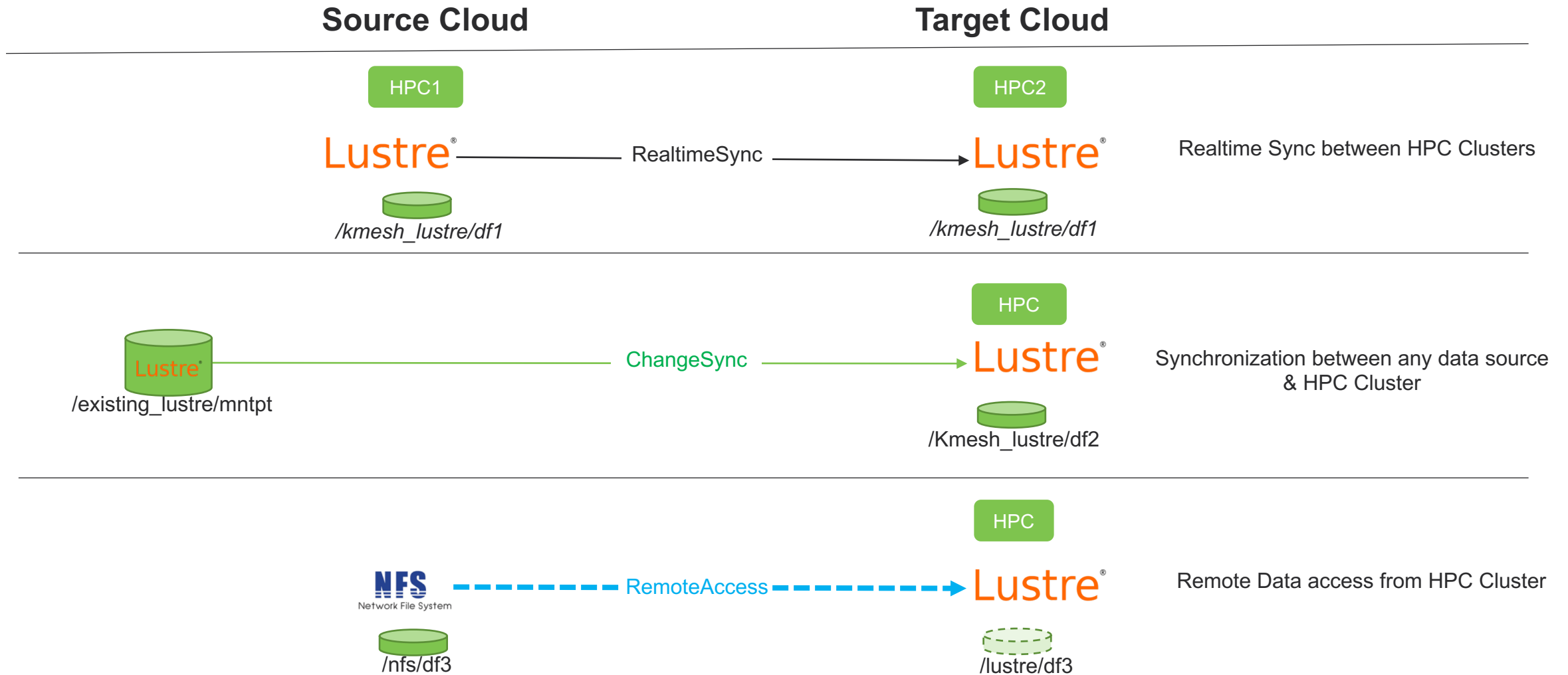
❖ IO 500 does not contain any of this information

Public Cloud Challenge: Data Synchronization



- One-time copy
- Continuous synchronization
- Real-Time synchronization
- Variety of sources

Potential HPC Cloud Bursting Models



Cloud Do's and Don'ts

- 👍 Do size HPC system primarily based on *end-to-end* network bandwidth
- 👍 Do consider that cloud infrastructure can go down
- 👍 Do extensive benchmark for your application
(IO500 need to incorporate 'cloud' elements 🤔)
- 👎 Don't expect same 1-to-1 replacement for on-premise system
- 👎 Don't buy more CPU core than you required
- 👎 Don't ignore data synchronization complexities