



英特尔®傲腾™在存储中的应用

姚 舸

南京大学

人工微结构科学与技术协同创新中心

yaoge@nju.edu.cn my@yaoge123.com



e-Science 中心



■ 服务教学科研

- 高性能计算
- 云盘
- 开源镜像
- Git
- 云平台
- 站群
- 信息安全
- 实训平台
- 仪器管理
- 信息发布
- 门禁
- 财务
- 资产



服务



通知公告

· 云盘外链

在外链可以设定精确到分钟的有效期。共享链接用于发布文件，提供下载的。上传链接用于收集文件，提供上传的。可以设定密码、有效期和权限。

🕒 2020-07-04

· 云盘和高性能计算QQ群

高性能计算很早就通过QQ群的形式为大家提供服务和交流的平台，现在云盘也有QQ群啦！云盘QQ群 2343870；高性能计算QQ群 1406661。

🕒 2020-04-13

· 各类服务全球可访问性和延迟监控

中心一直坚持公开运行、数据开放，近年来提供的服务越来越多，现新增了中心各类服务的全球可访问性和延迟信息。

🕒 2020-05-20

· 集群新增AMD CPU和NVIDIA GPU计算资源以及对象存储

新增AMD Rome CPU计算队列7702opa和NVIDIA GPU计算队列722080tiopa/72rtxopa，GPU节点特别适合AI计算。新增支持标准S3协议的对象存储 s3.nju.edu.cn。鼓励用户对比测试，不仅计算费免费，还有奖励

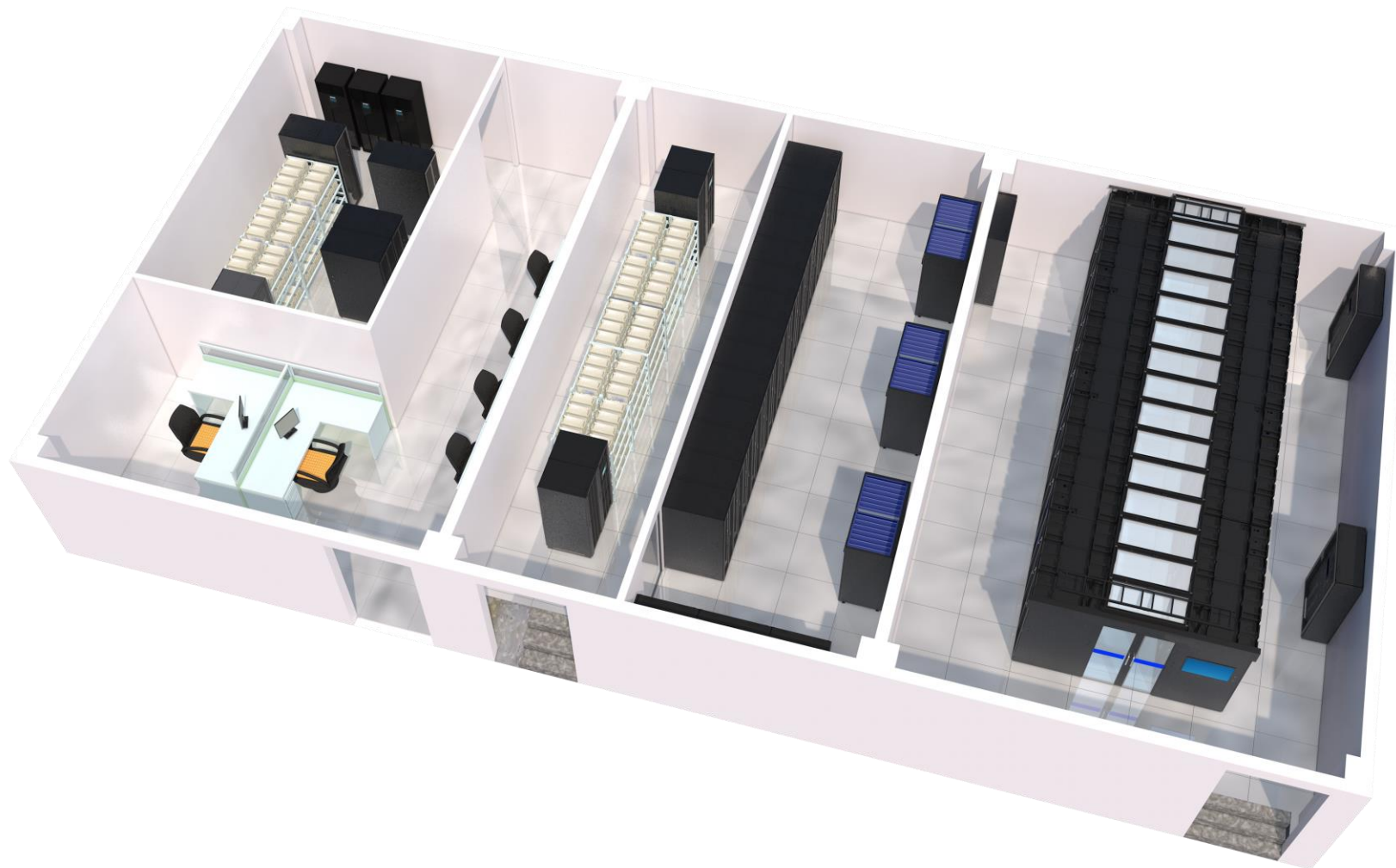
🕒 2020-03-30



基础设施



- 低密度机房
 - 高业务连续性
 - 云平台系统
 - HPC虚拟化
- 高密度机房
 - HPC计算与IO





- 统一机房课题组自行采购管理
 - 采购竞争不充分，硬件与机房环境可能不匹配
 - 学生管理为主，精力不足、更替传承问题
 - 资源使用不均衡，运行效率低



HPC管理



- 统一机房课题组自行采购管理
 - 采购竞争不充分，硬件与机房环境可能不匹配
 - 学生管理为主，精力不足、更替传承问题
 - 资源使用不均衡，运行效率低
- 各个课题组集群融合统一管理
 - 全生命周期管理，协助选型参与购置
 - 自主设计、安装、部署、运维
 - 统一管理、单一集群资源共享
 - 培训讲座，QQ群、微信群、公众号在线服务



品牌多样化



计算

IBM Hewlett Packard Enterprise
Dell intel
Lenovo
SUPERMICR H3C
ASUS 中科曙光 Sugon
inspur
GIGABYTE HUAWEI

存储

IBM Hewlett Packard Enterprise
Dell DDN[®]
AI · BIG DATA · HPC
HUAWEI UNIS WDC
inspur

网络

CISCO Hillstone NETWORKS
H3C BLADE NETWORK TECHNOLOGIES
Dell RUCKUS[™]
an ARRIS company
Voltaire Mellanox TECHNOLOGIES
HUAWEI intel
E Extreme[™]



配置多样化



■ CPU

- Intel Xeon 5150, ~~E5430~~, ~~X5550~~, X7542, ~~E5620~~, E5645, X5650, E5-2643, E5-2660, E5-2692v2, E5-2680v3, E5-2630v4, E5-2640v4, E5-2650v4, E5-2660v4, E5-2680v4, E5-2682v4, E7-4850v4, Gold 6140, Gold 6148, Gold 6248, Platinum 9242, Gold 5220R, Gold 6226R.
- AMD EPYC 7302P, 7702.

■ GPU

- Nvidia Tesla K40, P100, V100(PCIe), V100(SXM2), GeForce RTX 2080 Ti, TITAN RTX.

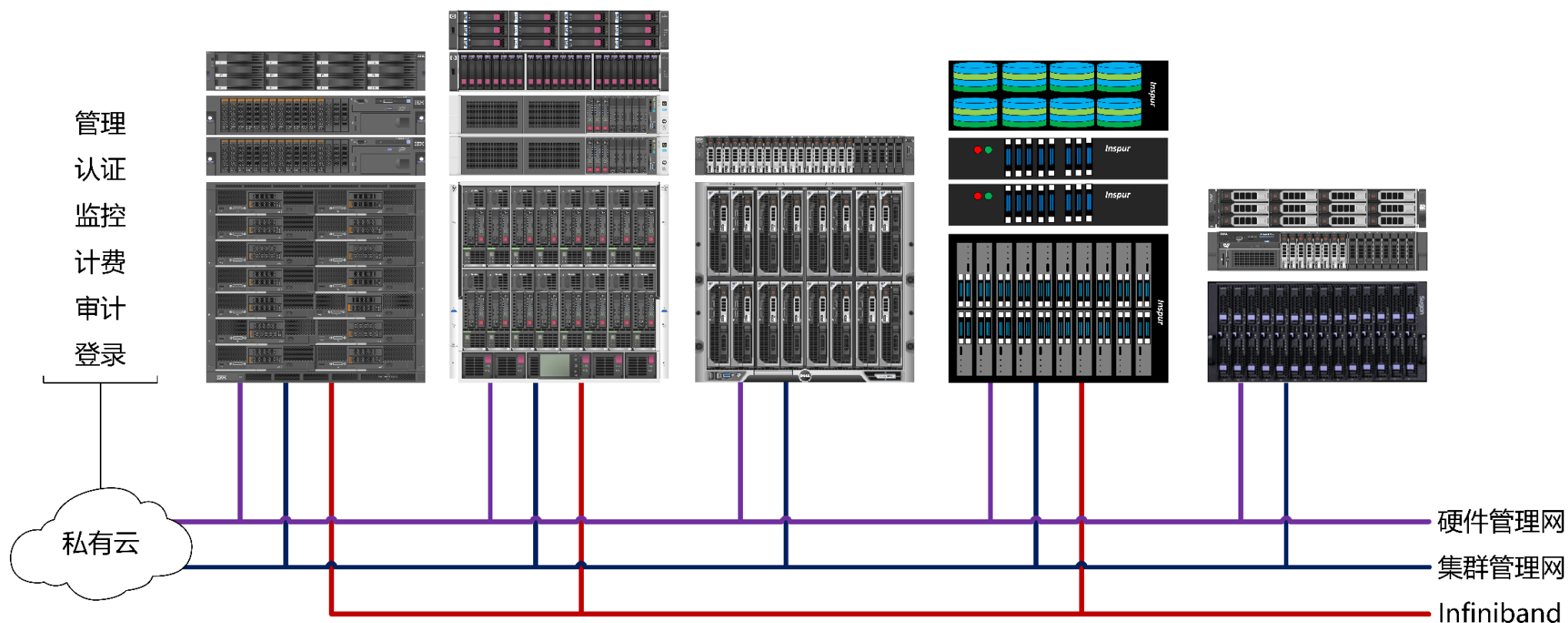
■ NET

- 1/10/25/40/100Gb Ethernet, 10Gb iWARP Ethernet, 20/40/56/100/200Gb InfiniBand, 100Gb Omni-Path Architecture.



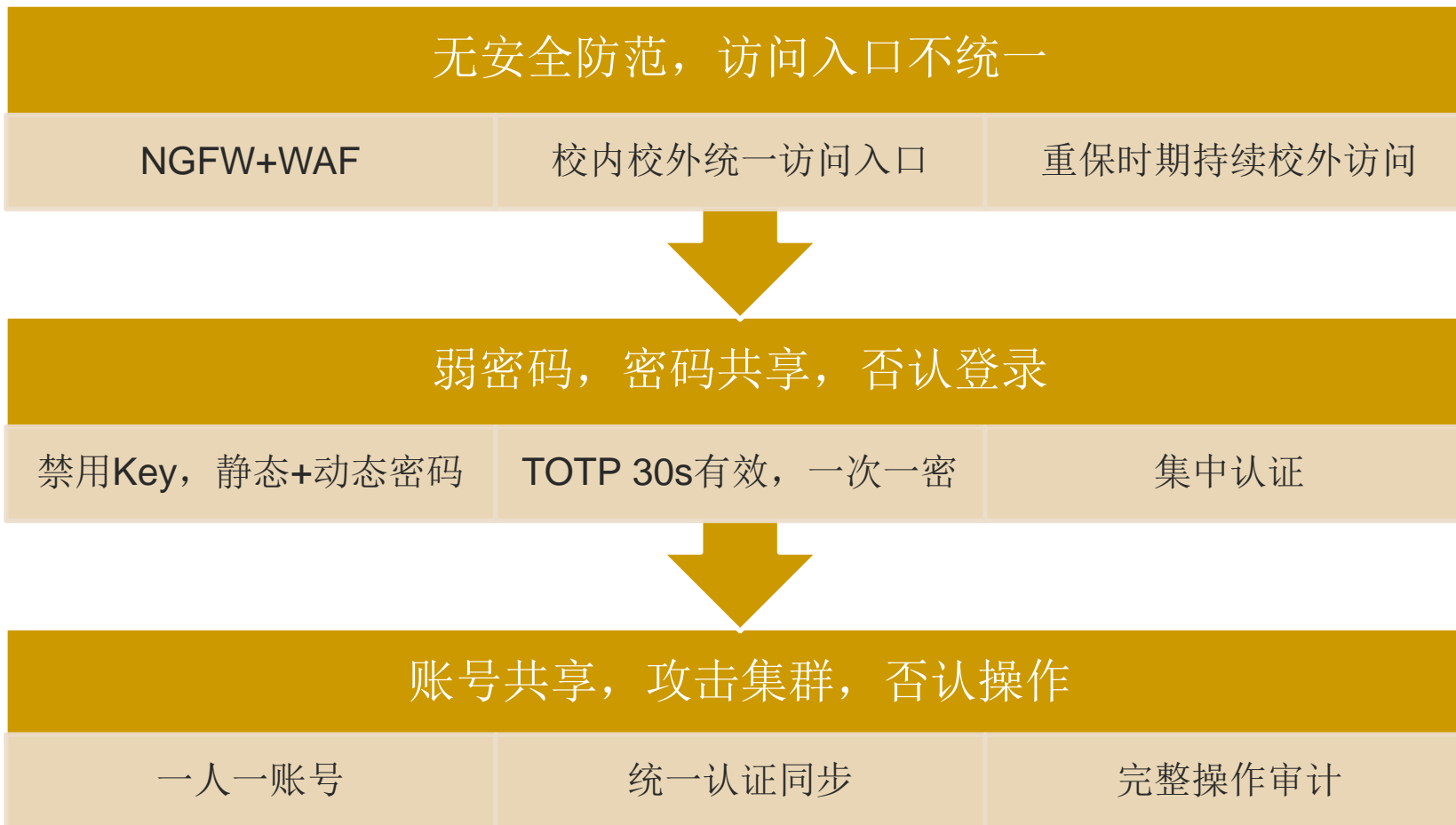
架构多样化

- 计算节点: 619 理论峰值: 661 TFLOPS + 411 TFLOPS
- IO+BB: 9+4 台; 虚机: 12 个; 存储: 4 套; 交换机: 16 ETH、4 IB、2 OPA





安全运营



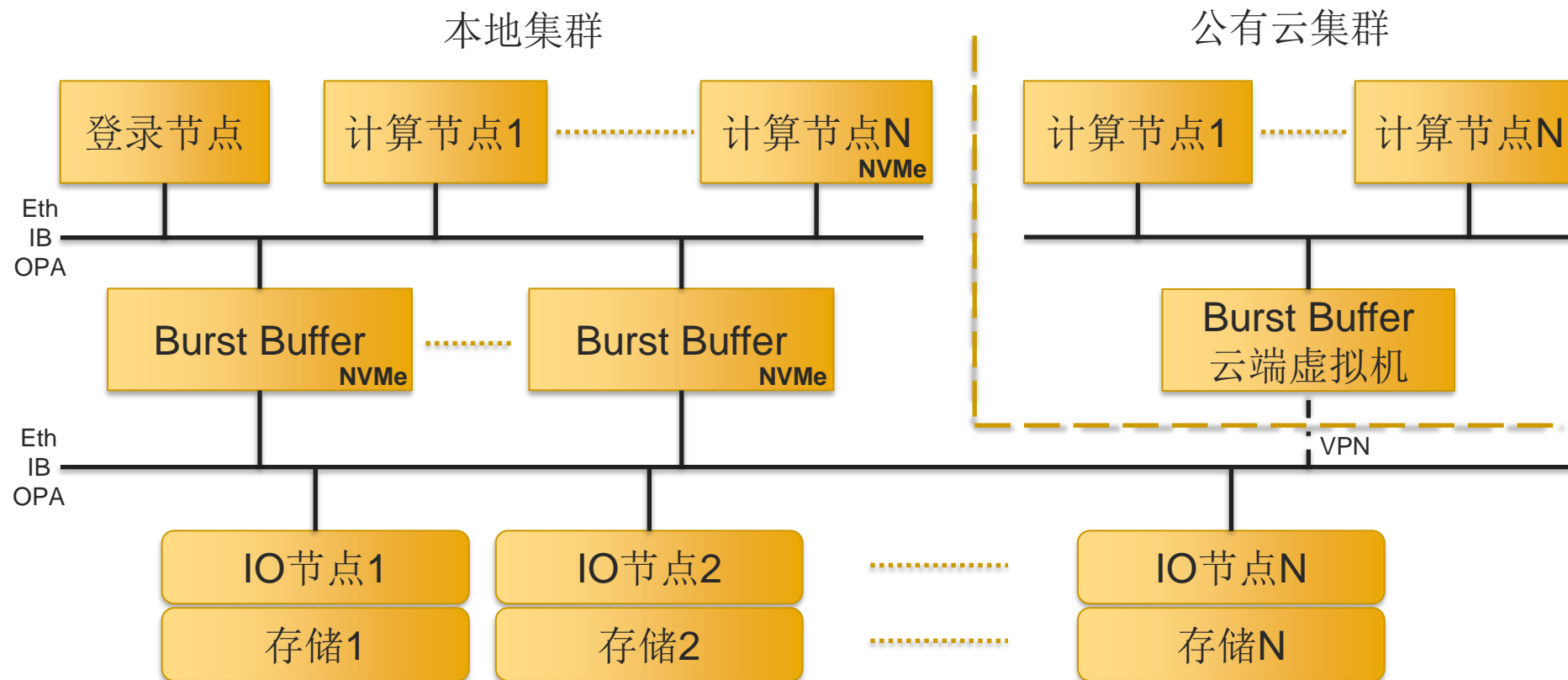


安全运营



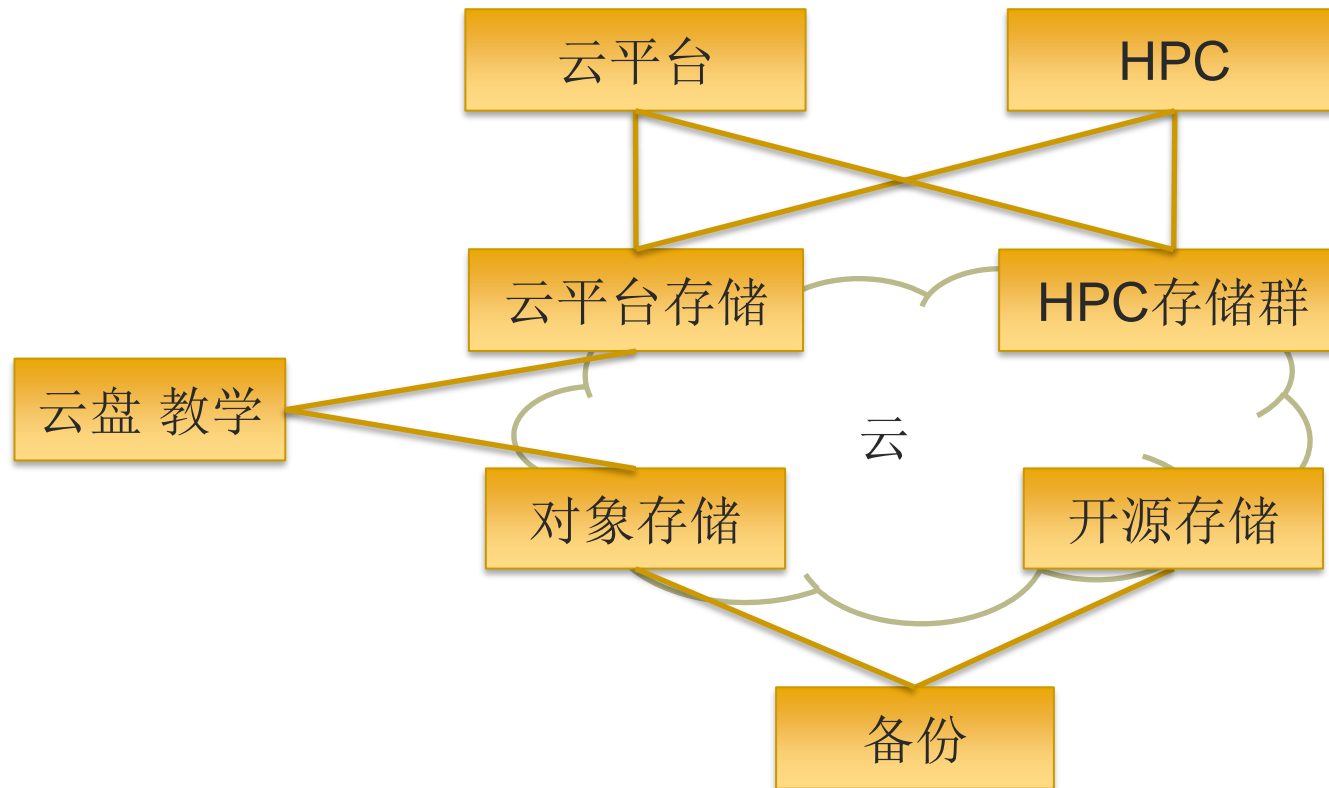


HPC 存储群





存储总览





开源存储 OpenZFS



高性价比，高数据持久性，业务连续性要求不高

- 集成了逻辑卷管理的文件系统
- 代替存储或RAID卡提供数据保护



开源存储 OpenZFS



提升性能

- ARC: 内存中的缓存
- L2ARC: ARC的扩展, 一般采用SSD
- ZIL: 持久性写缓存, 用于意外恢复



开源存储 OpenZFS

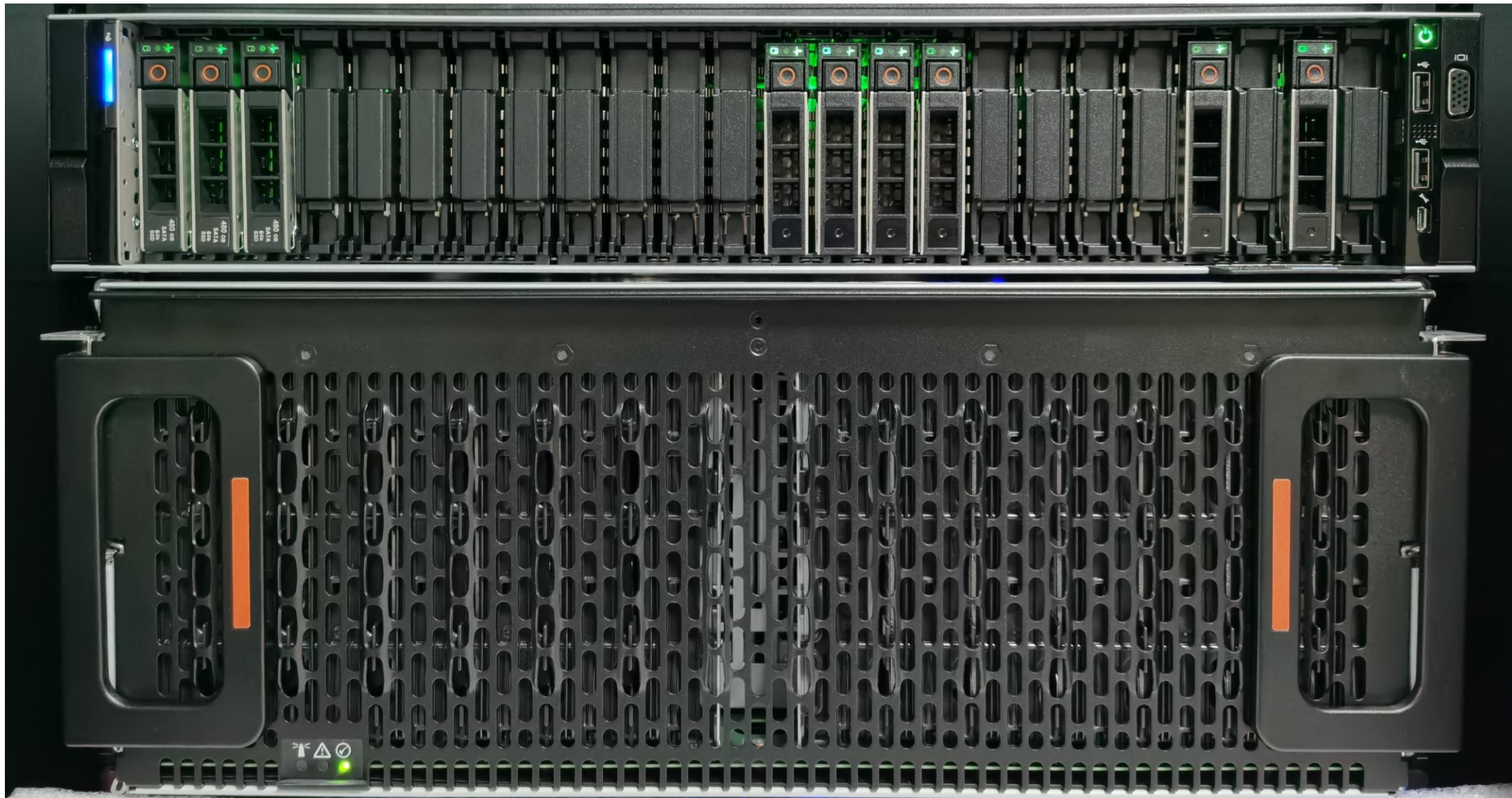


通用硬件

- X86服务器: Dell R740xd, 2*Intel Gold 5218, 12*32GB RAM, Broadcom Dual 25Gb Ethernet, 480GB SSD RAID1, 2*Dell 12Gbps HBA(SAS3508)连接JBOD
- JBOD: HGST Ultrastar Data60 H4060-J, 双控制器, 60*14TB NL-SAS
- 4*Intel SSD P4510 2TB U.2
- 2*Intel Optane SSD P4800X 375GB U.2
- 12*Intel Optane Persistent Memory 128GB



开源存储 OpenZFS



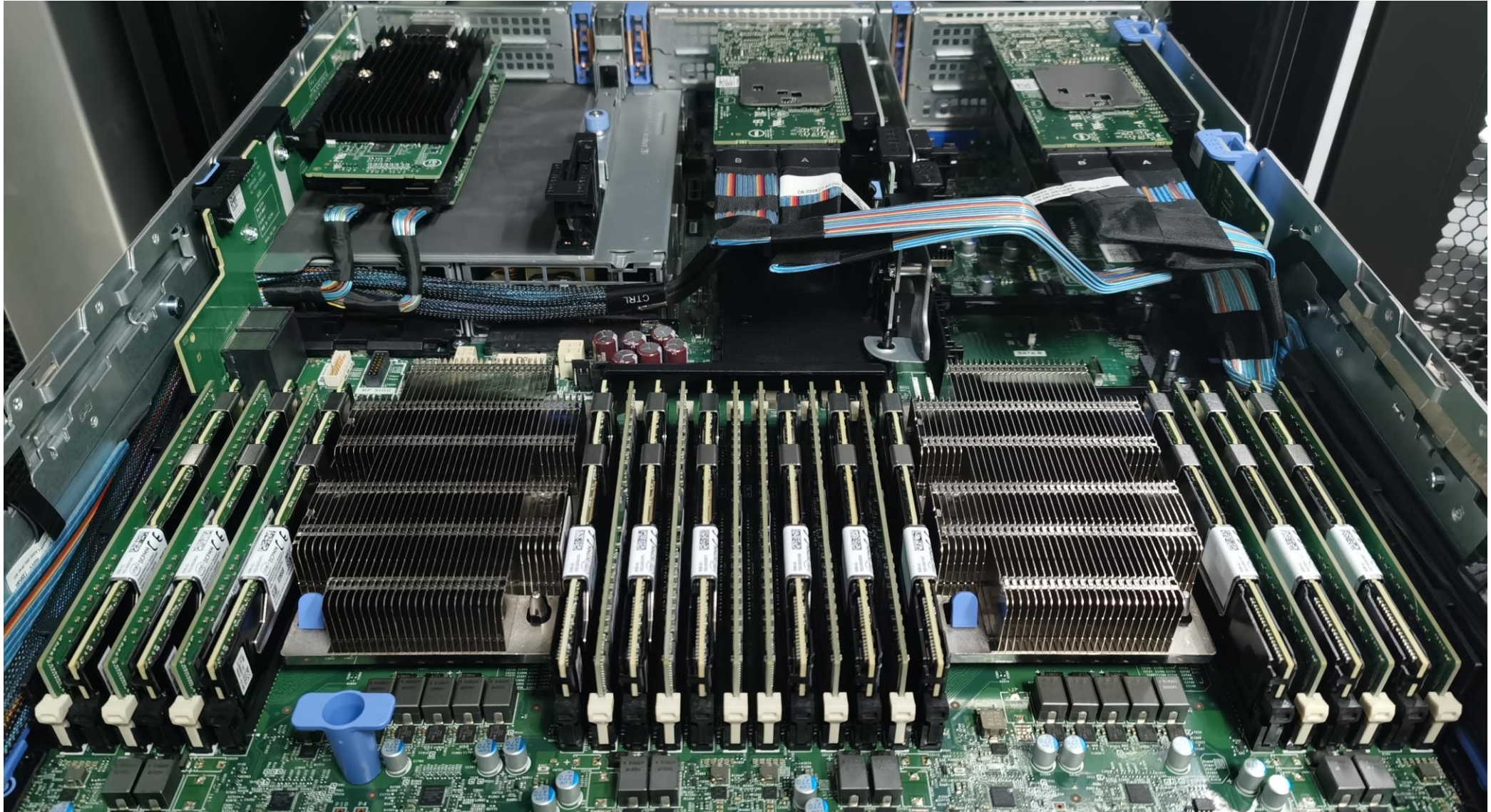


开源存储 OpenZFS





开源存储 OpenZFS





Optane Persistent Memory



MM+AD 混合模式

- 95%设置为内存模式

```
ipmctl create -goal MemoryMode=95
```

- 创建两个内存盘:

```
ndctl create-namespace -m raw
```

```
ndctl create-namespace -m raw
```



开源存储 OpenZFS



- 调整ARC容量

`/sys/module/zfs/parameters/zfs_arc_max`

- 添加L2ARC

`zpool add tank cache nvme0n1 nvme1n1 nvme2n1 nvme3n1`

- 添加ZIL使用镜像保护

`zpool add tank log mirror pmem0 pmem1`

- 异步写也落盘（可选）

`zfs set sync=always tank`



测试环境



■ HDD

- ARC: 247 GB DDR4 RAM
- ZIL: 5* Raidz3 (8+3) 14TB NL-SAS

■ NVMe SSD

- ARC: 247 GB DDR4 RAM
- L2ARC: 2* 2TB Intel SSD P4510
- ZIL: 2* 2TB Intel SSD P4510 Mirror



- NVMe Optane
 - ARC: 1.3 TB Optane MM
 - L2ARC: 4* 2TB Intel SSD P4510
 - ZIL: 2* 375GB Intel Optane SSD P4800X Mirror

- DIMM Optane
 - ARC: 1.3 TB Optane MM
 - L2ARC: 4* 2TB Intel SSD P4510
 - ZIL: 2* 36GB Optane AD Mirror



写入寿命



Product	Capacity	Lifetime Writes
Intel SSD P4510	1 TB	1.92 PBW
Intel SSD P4510	2 TB	2.61 PBW
Intel SSD P4510	4 TB	6.3 PBW
Intel SSD P4510	8 TB	13.88 PBW
Intel SSD P4610	1.6 TB	12.25 PBW
Intel SSD P4610	3.2 TB	21.85 PBW
Intel SSD P4610	6.4 TB	36.54 PBW
Intel Optane SSD P4801X	100 GB	10.9 PBW
Intel Optane SSD P4800X	375 GB	41.0 PBW

写入速度：100MB/s

五年写入：14.69PB



异步顺序写测试



Block Size	Process	HDD	NVMe SSD	NVMe Optane	DIMM Optane	
128k	1	1138 109	1123 111	963 136	1001 124	Bandwidth (MiB/s) Total latency (us)
128k	64	1424 5616	1531 5220	2073 4042	1660 4815	Bandwidth (MiB/s) Total latency (us)
4k	1	149k 6.43	151k 6.36	146k 6.57	148k 6.48	IOPS Total latency (us)
4k	64	347k 183	345k 185	371k 172	375k 170	IOPS Total latency (us)

Ubuntu 20.04.1 LTS, 5.4.0-42-generic, OpenZFS 0.8.3-1ubuntu12.3, fio-3.16
fio ioengine=psync iodepth=1 direct=1 rw=write -sync=0



异步随机写测试



Block Size	Process	HDD	NVMe SSD	NVMe Optane	DIMM Optane	
128k	1	1153 108	1033 120	1234 101	1380 90	Bandwidth (MiB/s) Total latency (us)
128k	64	378 21145	348 22964	345 23159	501 15968	Bandwidth (MiB/s) Total latency (us)
4k	1	7930 125	7203 138	8103 122	10200 97	IOPS Total latency (us)
4k	64	3177 20099	3181 20103	3634 17601	4708 13586	IOPS Total latency (us)

Ubuntu 20.04.1 LTS, 5.4.0-42-generic, OpenZFS 0.8.3-1ubuntu12.3, fio-3.16
fio ioengine=psync iodepth=1 direct=1 rw=randwrite -sync=0



同步顺序写测试



Block Size	Process	HDD	NVMe SSD	NVMe Optane	DIMM Optane	
128k	1	6.72 18590	473 263	407 306	250 499	Bandwidth (MiB/s) Total latency (us)
128k	64	209 38240	1384 5777	1360 5877	974 8213	Bandwidth (MiB/s) Total latency (us)
4k	1	77 12900	10.3k 96	9.66k 103	10.1k 98	IOPS Total latency (us)
4k	64	4089 15640	143k 447	135k 472	90.8k 704	IOPS Total latency (us)

Ubuntu 20.04.1 LTS, 5.4.0-42-generic, OpenZFS 0.8.3-1ubuntu12.3, fio-3.16
fio ioengine=psync iodepth=1 direct=1 rw=write -sync=1



同步随机写测试



Block Size	Process	HDD	NVMe SSD	NVMe Optane	DIMM Optane	
128k	1	5.17 24160	443 281	442 281	253 492	Bandwidth (MiB/s) Total latency (us)
128k	64	71.5 111	436 18351	513 15591	601 13300	Bandwidth (MiB/s) Total latency (us)
4k	1	68 14580	5540 179	5960 166	5636 176	IOPS Total latency (us)
4k	64	971 65820	2980 21467	3605 17741	3800 16833	IOPS Total latency (us)

Ubuntu 20.04.1 LTS, 5.4.0-42-generic, OpenZFS 0.8.3-1ubuntu12.3, fio-3.16
fio ioengine=psync iodepth=1 direct=1 rw=randwrite -sync=1



顺序读测试



Block Size	Process	ARC 2GB HDD	ARC 247GB HDD	DIMM Optane	
128k	1	127 982	2622 47.42	2516 49.41	Bandwidth (MiB/s) Total latency (us)
128k	64	851 9383	3314 2384	3539 2230	Bandwidth (MiB/s) Total latency (us)
4k	1	32.4k 30.34	287k 3.29	288k 3.29	IOPS Total latency (us)
4k	64	218k 29.32	856k 74.08	879k 71.94	IOPS Total latency (us)

Ubuntu 20.04.1 LTS, 5.4.0-42-generic, OpenZFS 0.8.3-1ubuntu12.3, fio-3.16
fio ioengine=psync iodepth=1 direct=1 rw=read



随机读测试



Block Size	Process	ARC 2GB HDD	ARC 247GB HDD	DIMM Optane	
128k	1	11.8 10622	2528 49.04	2345 52.91	Bandwidth (MiB/s) Total latency (us)
128k	64	119 67150	1290 6119	1548 5113	Bandwidth (MiB/s) Total latency (us)
4k	1	97 10252	24.4k 40.62	33.8k 29.26	IOPS Total latency (us)
4k	64	1209 52904	27.9k 2294	35.0k 1777	IOPS Total latency (us)

Ubuntu 20.04.1 LTS, 5.4.0-42-generic, OpenZFS 0.8.3-1ubuntu12.3, fio-3.16
fio ioengine=psync iodepth=1 direct=1 rw=randread



ZFS延迟 (生产环境)



Type	HDD	NVMe Optane	DIMM Optane
Total IO time - Read	27 ms	9 ms	10 ms
Total IO time - Write	16 ms	3 ms	3 ms
Disk IO time - Read	19 ms	8 ms	8 ms
Disk IO time - Write	6 ms	0.77 ms	0.73 ms
Sync queue time - Read	8 ms	0.36 ms	0.77 ms
Sync queue time - Write	1 ms	0.003 ms	0.003 ms
Async queue time - Read	7 ms	1 ms	1 ms
Async queue time - Write	18 ms	3 ms	3 ms



Optane Persistent Memory



优势

- 比内存更低的成本
- 比NAND更长的写入寿命
- MM+AD混合模式下节约用于ZIL的SSD



南京大學