

G5208 PCIE5 iDCL 智算服务器

产品可靠性报告

文档版本: V1.0

发布日期: 2025-9-12

文档更新记录

文档版本	发布日期	更新说明
V1.0	2025/9/12	首次发布

文档说明

本文档作为G5208 PCIe5 iDCL测试报告，旨在验证各部件之间可靠性以及产品性能。

硬件环境:

部件大类	品牌型号	FW	数量
处理器	Intel(R) Xeon Platinum 6530	/	2
内存	Samsung M321R8GA0BB0-CQKZJ D5-64G-4800	/	12
硬盘	SK hynix HFS480G3H2X069N -480G-2.5-SATA	410A2Z00	2
NVME硬盘	HUAWEI-ES3500P V6-SSD-3.84T-U.2	1011	1
GPU卡	PALIT GeForce RTX™ 5090 GameRock	98.02.2E.00.AA	8
网卡	Mellanox MCX4121A-ACAT-25G	14.32.1010	1
电源	航嘉-HKS2700D1-2700W	/	4

软件环境:

软件	版本信息
操作系统	Ubuntu 22.04
NVIDIA驱动版本	570.86.16
CUDA版本	V12.8
NCCL版本	2.27
Cutlass版本	3.9
Sysbench	1.0.20
BIOS版本	32.29.03
BMC版本	32.28.02

浸没环境:

Tank参数	详细信息
型号	绿色云图MLD50121YI
可用空间	21U
冷却液温度	35-50°C
冷却水温度	≤35°C
冷却能力	50kw
冷却液型号	LCS200
设备功率	1.5kw
最大运行重量	1000kg
净重	250kg

测试项目清单

测试类别	测试项目	测试目的	测试结果
可靠性测试	兼容性测试	验证G5208 PCIe5 iDLC与主流操作系统适配情况。	PASS
	稳定性测试	测试G5208 PCIe5 iDLC硬件的兼容性，验证出符合新产品出货要求的硬件清单。	PASS
	功耗测试	收集G5208 PCIe5 iDLC整机压测下的满载以及空载功耗情况。	PASS
	温度测试	输出G5208 PCIe5 iDLC满载下各个部件运行温度。	PASS
产品性能测试	GPU性能测试	收集G5208 PCIe5 iDLC的GPU链路速率、带宽，负载运行状态下Ispci检查Inksta，speed速率正常，不存在掉带宽、速率等情况。	PASS
	NCCL带宽测试	通过 GPU all_reduce 通信模型测试平台卡间通讯总线带宽性能；	PASS
	浮点性能测试	通过cutlass工具测试真实场景下GPU浮点运算性能；	PASS
	CPU性能测试	通过sysbench工具测试CPU性能；	PASS
	内存性能测试	通过sysbench工具测试内存带宽性能；	PASS
	硬盘性能测试	通过fio工具测试硬盘读写性能；	PASS

目 录

文档更新记录	i
文档说明	ii
测试项目清单	iii
1.可靠性测试	1
1.1 兼容性测试	1
1.2 稳定性测试	1
1.3 功耗测试	3
1.4 温度测试	4
2.产品性能测试	6
2.1 GPU带宽测试	6
2.2 NCCL带宽测试	8
2.3 浮点性能测试	9
2.4 CPU性能测试	10
2.5 内存性能测试	11
2.6 硬盘性能测试	12

1 可靠性测试

1.1 兼容性测试

测试目的:验证G5208 PCIE5 iDCL与当前主流操作系统适配情况;

测试操作系统:Ubuntu 20.04/Ubuntu22.04/Windows Server 2019;

测试目标:安装主流操作系统, 硬件驱动均可以正常安装且正常运行;

测试结果:操作系统及各硬件的驱动可以正常安装, 驱动来源来自于产品官方网站。

1.2 稳定性测试

测试目的:测试G5208 PCIE5 iDCL硬件的兼容性, 验证出符合出货要求的硬件清单;

测试目标:压测72小时测试途中硬件不出现丢失、死机、非自然重启等不良现象, 在测试环境温度下主要部件温度不超过阈值;

测试工具:

(1)**Stressapptest:**让来自处理器和I/O到内存的数据尽量随机化, 以创造出模拟现实的环境来测试现在的硬件设备是否稳定, 如CPU、内存、硬盘等;

(2)**GPU_BURN:** 一个用于测试图形处理器 (GPU) 性能和稳定性的工具。它利用了现代GPU的计算能力, 通过持续执行繁重的图形运算, 以最大程度地激发GPU的工作负荷。这样做有助于评估GPU的耐久性和稳定性。

测试结果:72小时stressapptest+gpu_burn压力测试，测试过程中无意外造成CPU/GPU等硬件丢失及重启宕机等不良现象。

测试截图:

```
100.0% proc'd: 15662240 (66354 Gflop/s) - 15560116 (66099 Gflop/s) - 15710761 (66710 Gflop/s) - 15708220 (66726 Gflop/s) - 15875079 (67255 Gflop/s) - 15419393 (65484 Gflop/s) - 15497680 (65756 Gflop/s) - 15507723 (65989 Gflop/s) errors: 0 - 0 - 0 - 0 - 0 - 0 - 0 temps: 78 C - 78 C - 76 C - 75 C - 73 C - 79 C
100.0% proc'd: 15662240 (66354 Gflop/s) - 15560116 (66099 Gflop/s) - 15710761 (66710 Gflop/s) - 15708220 (66726 Gflop/s) - 15875079 (67255 Gflop/s) - 15419393 (65484 Gflop/s) - 15497801 (65780 Gflop/s) - 15507723 (65989 Gflop/s) errors: 0 - 0 - 0 - 0 - 0 - 0 - 0 temps: 78 C - 78 C - 76 C - 75 C - 73 C - 79 C
100.0% proc'd: 15662240 (66354 Gflop/s) - 15560116 (66099 Gflop/s) - 15710882 (66709 Gflop/s) - 15708220 (66726 Gflop/s) - 15875079 (67255 Gflop/s) - 15419393 (65484 Gflop/s) - 15497801 (65780 Gflop/s) - 15507723 (65989 Gflop/s) errors: 0 - 0 - 0 - 0 - 0 - 0 - 0 temps: 78 C - 78 C - 76 C - 75 C - 73 C - 79 C
100.0% proc'd: 15662240 (66354 Gflop/s) - 15560116 (66099 Gflop/s) - 15710882 (66709 Gflop/s) - 15708220 (66726 Gflop/s) - 15875079 (67255 Gflop/s) - 15419514 (65490 Gflop/s) - 15497801 (65780 Gflop/s) - 15507723 (65989 Gflop/s) errors: 0 - 0 - 0 - 0 - 0 - 0 - 0 temps: 78 C - 78 C - 76 C - 75 C - 73 C - 79 C
78 C - 75 C
Killing processes with SIGTERM (soft kill)
^reed memory for dev 1
^initted cublas
^reed memory for dev 3
^initted cublas
^reed memory for dev 7
^initted cublas
^reed memory for dev 0
^initted cublas
^reed memory for dev 4
^initted cublas
^reed memory for dev 6
^initted cublas
^reed memory for dev 2
^initted cublas
^reed memory for dev 5
^initted cublas
done

Tested 8 GPUs:
GPU 0: OK
GPU 1: OK
GPU 2: OK
GPU 3: OK
GPU 4: OK
GPU 5: OK
GPU 6: OK
GPU 7: OK

*****Test Summary: Done *****

Log: Seconds remaining: 258930
Log: Seconds remaining: 258920
Log: Seconds remaining: 258910
Log: Seconds remaining: 258900
Log: Seconds remaining: 258890
Log: Seconds remaining: 258880
Log: Seconds remaining: 258870
Log: Seconds remaining: 258860
Log: Seconds remaining: 258850
Log: Seconds remaining: 258840
Log: Seconds remaining: 258830
Log: Seconds remaining: 258820
Log: Seconds remaining: 258810
Log: Seconds remaining: 258800
Log: Seconds remaining: 258790
Log: Seconds remaining: 258780
Log: Seconds remaining: 258770
Log: Seconds remaining: 258760
Log: Seconds remaining: 258750
Log: Seconds remaining: 258740
Log: Seconds remaining: 258730
Log: Seconds remaining: 258720
Log: Seconds remaining: 258710
Log: Seconds remaining: 258700
Log: Seconds remaining: 258690
Log: Seconds remaining: 258680
Log: Seconds remaining: 258670
Log: Seconds remaining: 258660
Log: Seconds remaining: 258650
Log: Seconds remaining: 258640
Log: Seconds remaining: 258630
Log: Seconds remaining: 258620
Log: Seconds remaining: 258610
Log: Seconds remaining: 258600
Log: Pausing worker threads in preparation for power spike (258600 seconds remaining)
Log: Seconds remaining: 258590
Log: Resuming worker threads to cause a power spike (258585 seconds remaining)
```

(1)**DC Recycle**: 进行200次不断电重启，每次重启后获取硬件信息对比，全部正常继续进行reboot操作，对比信息failed则终止DC操作。

测试结果:重启测试200次，测试过程中检测硬件信息均能正常识别。

测试截图:

The 200 loop start:

```
net-check:True
cpu_hw:True
mem_hw:True
disk_hw:True
gpu_hw:True
```

The 200 loop: PASS

1.3 功耗测试

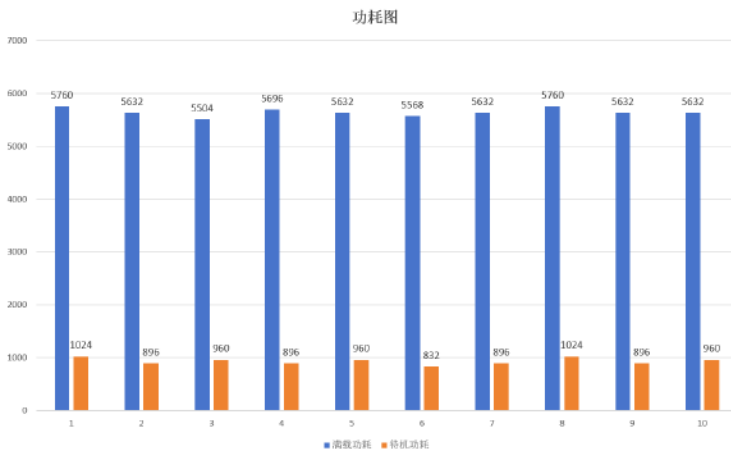
测试目的: 收集G5208 PCIe5 iDCL整机压测下的满载以及空载功耗情况；

测试工具:

(1)**Stressapptest**:让来自处理器和I/O到内存的数据尽量随机化，以创造出模拟现实的环境来测试现在的硬件设备是否稳定，如CPU、内存、硬盘等；

(2)**GPU_BURN**: 一个用于测试图形处理器（GPU）性能和稳定性的工具。它利用了现代GPU的计算能力，通过持续执行繁重的图形运算，以最大程度地激发GPU的工作负荷。这样做有助于评估GPU的耐久性和稳定性。

测试结果:随机抽取10个24小时压测过程中整机满负载&空载下功耗读数测试样例进行统计，结果如下：



1.4 温度测试

测试目的:观察G5208 PCIe5 iDCL满负载压测72小时下GPU温度和状态。

测试工具:

(1)**Stressapptest:** 让来自处理器和I/O到内存的数据尽量随机化，以创造出模拟现实的环境来测试现在的硬件设备是否稳定，如CPU、内存、硬盘等；

(2)**GPU_BURN:** 一个用于测试图形处理器（GPU）性能和稳定性的工具。它利用了现代GPU的计算能力，通过持续执行繁重的图形运算，以最大程度地激发GPU的工作负荷。这样做有助于评估GPU的耐久性和稳定性。

(3)**温度记录仪:** 用于定点记录环境温度。

测试流程: 将整机完全浸没至tank中，尾部为进液处，头部为出液处，将温度记录仪的温感针脚放置进出液处。开机后使用压测程序进行压测，使用监控脚本对GPU的温度信息进行记录。

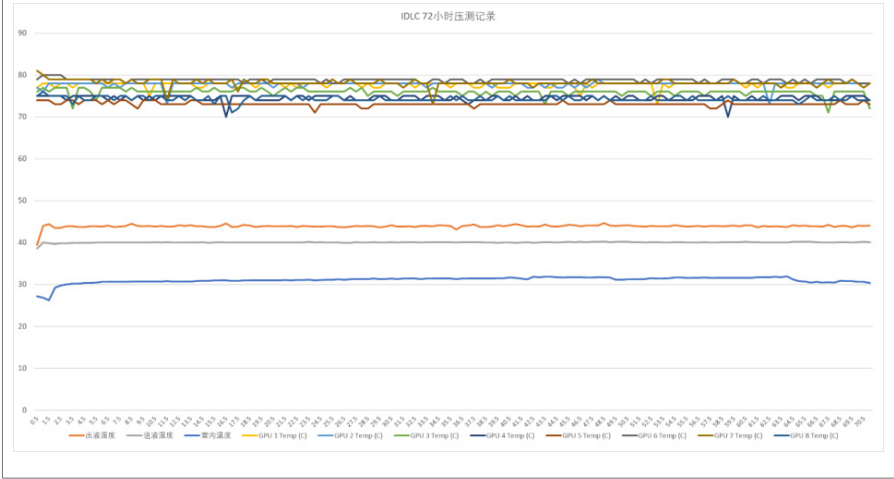
压测结果：

压测时长	72H	测试时间	2025/8/30-2025/9/2
进出液温度	40-42℃/43-44℃	环境温度	26-32℃
测试结果：压测72小时无异常无报错。		测试结论	PASS

测试截图：

```
100.0% proc'd: 15662240 (66354 Gflop/s) - 15560116 (66099 Gflop/s) - 15710761 (66710 Gflop/s) - 15708220 (66726 Gflop/s) - 15875079 (67255 Gflop/s) - 15419393 (65484 Gflop/s) - 15497680 (65756 Gflop/s) - 15507723 (65989 Gflop/s) errors: 0 - 0 - 0 - 0 - 0 - 0 - 0 temps: 78 C - 78 C - 76 C - 75 C - 73 C - 73 C
100.0% proc'd: 15662240 (66354 Gflop/s) - 15560116 (66099 Gflop/s) - 15710761 (66710 Gflop/s) - 15708220 (66726 Gflop/s) - 15875079 (67255 Gflop/s) - 15419393 (65484 Gflop/s) - 15497680 (65756 Gflop/s) - 15507723 (65989 Gflop/s) errors: 0 - 0 - 0 - 0 - 0 - 0 - 0 temps: 78 C - 78 C - 76 C - 75 C - 73 C - 73 C
100.0% proc'd: 15662240 (66354 Gflop/s) - 15560116 (66099 Gflop/s) - 15710882 (66709 Gflop/s) - 15708220 (66726 Gflop/s) - 15875079 (67255 Gflop/s) - 15419393 (65484 Gflop/s) - 15497680 (65756 Gflop/s) - 15507723 (65989 Gflop/s) errors: 0 - 0 - 0 - 0 - 0 - 0 - 0 temps: 78 C - 78 C - 76 C - 75 C - 73 C - 73 C
100.0% proc'd: 15662240 (66354 Gflop/s) - 15560116 (66099 Gflop/s) - 15710882 (66709 Gflop/s) - 15708220 (66726 Gflop/s) - 15875079 (67255 Gflop/s) - 15419514 (65490 Gflop/s) - 15497680 (65756 Gflop/s) - 15507723 (65989 Gflop/s) errors: 0 - 0 - 0 - 0 - 0 - 0 - 0 temps: 78 C - 78 C - 76 C - 75 C - 73 C - 73 C
78 C - 75 C
killing processes with SIGTERM (soft kill)
free memory for dev 1
jinitiated cublas
free memory for dev 3
jinitiated cublas
free memory for dev 7
jinitiated cublas
free memory for dev 0
jinitiated cublas
free memory for dev 4
jinitiated cublas
free memory for dev 6
jinitiated cublas
free memory for dev 2
jinitiated cublas
free memory for dev 5
jinitiated cublas
jone
Tested 8 GPUs:
GPU 0: OK
GPU 1: OK
GPU 2: OK
GPU 3: OK
GPU 4: OK
GPU 5: OK
GPU 6: OK
GPU 7: OK
xxxxxxxxxxxxxxxxxxxxxx
```

进出液温度曲线：



(2)GPU带宽测试:

测试目的:收集G5208 PCIe5 iDCL的GPU P2P和GPU Bandwidth性能, 验证GPU与主机的传输速率是否达到正常标准;

测试工具:CUDA 编译包含的测试样本p2pBandwidthLatencyTest和BandwidthTest 测试程序;

测试结果:系统内存和5090之间的传输通过PCIe5.0进行, 实测可以在主机和GPU之间实现56GB/s传输速率;

p2pBandwidthLatencyTest测试截图:

Unidirectional P2P-Disabled Bandwidth Matrix (GB/s)									
D\0	0	1	2	3	4	5	6	7	
0	1511.17	42.73	42.54	42.55	42.02	42.08	42.07	42.28	
1	42.40	1537.99	42.59	42.59	42.21	42.01	42.08	42.28	
2	42.75	42.42	1556.26	42.57	42.12	42.12	42.42	42.44	
3	42.57	42.02	42.04	1549.09	42.08	42.12	42.11	42.43	
4	42.89	42.14	42.28	42.15	1545.50	42.75	42.97	42.92	
5	41.98	42.11	42.08	42.14	43.12	1539.50	42.79	42.97	
6	42.13	42.10	42.10	42.10	43.11	42.94	1540.97	42.88	
7	42.17	42.08	42.17	42.11	42.84	43.11	42.91	1541.02	
Unidirectional P2P-Enabled Bandwidth (P2P Writes) Matrix (GB/s)									
D\0	0	1	2	3	4	5	6	7	
0	1511.17	43.03	42.87	42.81	42.64	42.38	42.46	42.55	
1	42.80	1539.46	42.90	42.60	42.58	42.35	42.45	42.52	
2	43.03	42.66	1539.41	42.86	42.53	42.51	42.34	42.49	
3	42.91	42.83	42.58	1534.92	42.53	42.52	42.26	42.46	
4	42.28	42.25	42.42	42.33	1537.94	43.32	43.38	43.40	
5	42.15	42.30	42.26	42.30	43.21	1540.97	43.14	43.43	
6	42.47	42.32	42.24	42.25	43.32	43.29	1542.50	43.40	
7	42.47	42.30	42.15	42.27	43.23	43.21	43.24	1540.93	
Bidirectional P2P-Disabled Bandwidth Matrix (GB/s)									
D\0	0	1	2	3	4	5	6	7	
0	1527.30	56.57	56.66	56.66	56.34	56.21	56.31	56.17	
1	56.75	1530.38	56.68	57.02	56.35	56.46	56.17	56.22	
2	56.74	56.61	1537.09	56.66	56.31	56.27	56.23	56.05	
3	56.85	56.63	56.76	1537.09	56.35	56.16	56.02	55.98	
4	55.91	56.74	56.19	56.12	1538.60	56.93	56.93	57.04	
5	56.37	56.59	56.30	56.57	57.76	1541.94	56.76	56.82	
6	56.56	56.59	56.53	56.08	57.00	56.73	1540.12	57.14	
7	56.27	56.23	56.23	56.32	56.94	57.08	56.88	1537.09	
Bidirectional P2P-Enabled Bandwidth Matrix (GB/s)									
D\0	0	1	2	3	4	5	6	7	
0	1522.86	56.81	56.45	56.76	56.11	56.40	56.43	56.55	
1	56.46	1538.60	56.87	56.84	56.15	56.31	56.32	56.22	
2	56.53	56.61	1538.60	56.88	56.12	56.89	56.27	56.42	
3	56.62	56.65	56.41	1537.07	55.87	56.17	56.05	56.82	
4	56.42	56.14	56.12	56.90	1540.12	57.13	56.95	57.20	
5	56.89	56.14	56.59	56.82	56.96	1538.96	57.04	56.96	
6	56.49	56.81	56.27	56.34	56.99	56.92	1537.09	56.72	
7	56.89	56.28	56.39	56.30	57.14	57.13	56.81	1543.10	

P2P-Disabled Latency Matrix (us)									
GPU	0	1	2	3	4	5	6	7	
0	2.10	13.90	14.31	13.86	14.59	14.54	14.53	14.55	
1	13.87	2.11	13.96	14.23	14.61	14.57	14.56	14.50	
2	13.61	13.36	2.09	13.64	14.59	14.54	14.57	14.54	
3	13.53	14.01	14.13	2.13	14.00	14.56	14.60	14.50	
4	14.59	14.26	14.26	14.63	2.08	13.37	13.20	13.26	
5	14.48	14.56	14.56	14.60	13.10	2.09	13.26	13.24	
6	14.57	14.51	14.56	14.30	13.30	13.27	2.06	13.09	
7	14.50	14.36	14.53	14.59	12.94	13.48	13.63	2.08	
P2P-Enabled Latency (P2P Writes) Matrix (us)									
GPU	0	1	2	3	4	5	6	7	
0	2.97	9.38	9.15	9.40	8.39	8.52	8.46	8.52	
1	9.30	2.90	9.26	9.46	8.53	8.52	8.55	8.58	
2	9.17	9.28	2.91	9.53	8.62	8.58	8.60	8.59	
3	9.18	9.29	9.23	2.88	8.61	8.60	8.60	8.64	
4	8.70	8.73	8.61	8.97	2.91	8.27	8.15	8.13	
5	8.61	8.76	8.67	8.92	7.94	2.63	8.18	8.17	
6	8.62	8.82	8.81	9.00	8.12	8.17	2.64	8.27	
7	8.65	8.93	8.85	9.01	8.18	8.24	8.30	2.67	
P2P-Enabled Latency (P2P Reads) Matrix (us)									
GPU	0	1	2	3	4	5	6	7	
0	2.10	14.15	13.91	13.98	14.56	14.46	14.49	14.58	
1	14.32	2.10	13.63	13.62	14.58	14.35	14.51	14.56	
2	14.24	13.45	2.19	13.53	14.50	14.37	14.56	14.54	
3	14.40	13.72	13.99	2.68	14.43	14.30	14.46	14.49	
4	13.62	14.58	14.40	14.57	2.10	13.21	13.08	12.97	
5	14.38	14.48	14.01	14.46	12.95	2.88	13.10	13.39	
6	14.35	13.75	14.54	14.42	12.63	13.11	2.97	12.79	
7	14.63	14.55	14.26	14.62	13.36	13.13	13.18	2.08	
P2P-Enabled Latency (P2P Reads) Matrix (us)									
GPU	0	1	2	3	4	5	6	7	
0	3.80	9.32	9.10	9.31	8.35	8.49	8.43	8.43	
1	9.17	2.85	9.24	9.38	8.49	8.57	8.57	8.58	
2	9.05	9.08	2.89	9.31	8.43	8.50	8.48	8.48	
3	9.27	9.39	9.31	2.97	8.65	8.75	8.74	8.83	
4	8.69	8.72	8.89	8.92	2.60	8.21	8.15	8.11	
5	8.74	8.90	8.85	8.85	8.10	2.66	8.28	8.37	
6	8.65	8.84	8.92	8.98	8.88	8.11	2.60	8.20	
7	8.63	8.83	8.84	8.94	8.98	8.10	8.25	2.67	

Bandwidth:测试截图:

```

root@stone:/home/stone/cuda-samples/Samples/1_Utillities/bandwidthTest# ./bandwidthTest
[CUDA Bandwidth Test] - Startng...
Running on...

Device 0: NVIDIA GeForce RTX 5090
Quick Mode

Host to Device Bandwidth, 1 Device(s)
PINNED Memory Transfers
Transfer Size (Bytes)      Bandwidth(GB/s)
32000000                   56.3

Device to Host Bandwidth, 1 Device(s)
PINNED Memory Transfers
Transfer Size (Bytes)      Bandwidth(GB/s)
32000000                   57.2

Device to Device Bandwidth, 1 Device(s)
PINNED Memory Transfers
Transfer Size (Bytes)      Bandwidth(GB/s)
32000000                   4381.8

Result = PASS

```

2.2 NCCL带宽测试

测试目的:通过 GPU all_reduce 通信模型测试平台卡间通讯总线带宽性能;

测试工具:NCCL

测试结果:设备8卡NCCL总线带宽能达到41GB/s, alltoall总线带宽能达到37GB/s。

测试截图:

all_reduce

```
root@stone:/home/stone/nccl-tests/build# ./all_reduce_perf -b 128M -e 8G -f 2 -g 8
# Collective test starting: all_reduce_perf
# nthread 1 ncpus 8 minBytes 134217728 maxBytes 8589934592 step: 2(factor) warnup iters: 1 iters: 20 agg iters: 1 validation: 1 graph: 0
#
# Using devices
# Rank 0 Group 0 Pid 3526 on stone device 0 [0000:16:00] NVIDIA GeForce RTX 5090
# Rank 1 Group 0 Pid 3526 on stone device 1 [0000:36:00] NVIDIA GeForce RTX 5090
# Rank 2 Group 0 Pid 3526 on stone device 2 [0000:46:00] NVIDIA GeForce RTX 5090
# Rank 3 Group 0 Pid 3526 on stone device 3 [0000:56:00] NVIDIA GeForce RTX 5090
# Rank 4 Group 0 Pid 3526 on stone device 4 [0000:98:00] NVIDIA GeForce RTX 5090
# Rank 5 Group 0 Pid 3526 on stone device 5 [0000:b8:00] NVIDIA GeForce RTX 5090
# Rank 6 Group 0 Pid 3526 on stone device 6 [0000:c8:00] NVIDIA GeForce RTX 5090
# Rank 7 Group 0 Pid 3526 on stone device 7 [0000:d8:00] NVIDIA GeForce RTX 5090
#
# size count type redop root time out-of-place in-place
# (B) (elements) (us) (GB/s) (GB/s) (us) (GB/s) (GB/s)
# 134217728 33554432 float sum -1 5915.0 22.69 39.71 0 5908.9 22.71 39.75 0
# 268435456 67108864 float sum -1 11747 22.85 39.99 0 11732 22.88 40.04 0
# 536870912 134217728 float sum -1 23293 23.05 40.23 0 23254 23.09 40.40 0
# 1073741824 268435456 float sum -1 46212 23.24 40.06 0 46226 23.23 40.05 0
# 2147483648 536870912 float sum -1 91032 23.30 40.92 0 91041 23.30 40.92 0
# 4294967296 1073741824 float sum -1 182619 23.52 41.16 0 182864 23.49 41.10 0
# 8589934592 2147483648 float sum -1 364492 23.57 41.24 0 364576 23.56 41.22 0
# Out of bounds values : 0 OK
# Avg bus bandwidth : 40.5799
# Collective test concluded: all_reduce_perf
```

alltoall

```
root@stone:/home/stone/nccl-tests/build# ./alltoall_perf -b 128M -e 8G -f 2 -g 8
# Collective test starting: alltoall_perf
# nthread 1 ncpus 8 minBytes 134217728 maxBytes 8589934592 step: 2(factor) warnup iters: 1 iters: 20 agg iters: 1 validation: 1 graph: 0
#
# Using devices
# Rank 0 Group 0 Pid 3364 on stone device 0 [0000:16:00] NVIDIA GeForce RTX 5090
# Rank 1 Group 0 Pid 3364 on stone device 1 [0000:36:00] NVIDIA GeForce RTX 5090
# Rank 2 Group 0 Pid 3364 on stone device 2 [0000:46:00] NVIDIA GeForce RTX 5090
# Rank 3 Group 0 Pid 3364 on stone device 3 [0000:56:00] NVIDIA GeForce RTX 5090
# Rank 4 Group 0 Pid 3364 on stone device 4 [0000:98:00] NVIDIA GeForce RTX 5090
# Rank 5 Group 0 Pid 3364 on stone device 5 [0000:b8:00] NVIDIA GeForce RTX 5090
# Rank 6 Group 0 Pid 3364 on stone device 6 [0000:c8:00] NVIDIA GeForce RTX 5090
# Rank 7 Group 0 Pid 3364 on stone device 7 [0000:d8:00] NVIDIA GeForce RTX 5090
#
# size count type redop root time out-of-place in-place
# (B) (elements) (us) (GB/s) (GB/s) (us) (GB/s) (GB/s)
# 134217728 4194304 float none -1 3303.1 40.63 35.55 0 3284.3 40.87 35.76 N/A
# 268435456 8388608 float none -1 6384.7 42.04 36.79 0 6434.1 41.72 36.51 N/A
# 536870912 16777216 float none -1 12591 42.64 37.31 0 12711 42.24 36.96 N/A
# 1073741824 33554432 float none -1 25027 42.90 37.54 0 25299 42.44 37.14 N/A
# 2147483648 67108864 float none -1 40856 43.07 37.60 0 50488 42.60 37.28 N/A
# 4294967296 134217728 float none -1 95552 43.14 37.75 0 100811 42.60 37.28 N/A
# 8589934592 268435456 float none -1 198943 43.18 37.78 0 202060 42.52 37.21 N/A
# Out of bounds values : 0 OK
# Avg bus bandwidth : 37.6382
# Collective test concluded: alltoall_perf
```

2.3 浮点性能测试

测试目的:通过cutlass工具测试真实场景下GPU浮点运算性能

测试工具: cutlass

测试结果: GPU浮点运算测试FP16结果达到409 Tflops

单位:Tflops

GPU	FP16	TF32	FP32	FP64
GPU0	413.51	121.92	76.53	1.95
GPU1	412.23	121.92	76.22	1.94
GPU2	411.15	121.91	76.09	1.94
GPU3	410.93	121.91	76.07	1.94
GPU4	410.46	121.90	76.57	1.94
GPU5	411.46	121.90	76.36	1.94
GPU6	409.60	121.91	76.60	1.94
GPU7	409.25	121.91	76.36	1.94

2.4 CPU性能测试

测试目的:通过sysbench工具测试真实场景下CPU性能

测试工具: sysbench

测试结果: 单核CPU在10秒内完成素数计算任务达到13193次。

测试截图:

```
root@stone:/home/stone/cutlass/build/tools/profiler# sysbench cpu --cpu-max-prime=20000 run
sysbench 1.0.20 (using system LuaJIT 2.1.0-beta3)
```

```
Running the test with following options:
```

```
Number of threads: 1
```

```
Initializing random number generator from current time
```

```
Prime numbers limit: 20000
```

```
Initializing worker threads...
```

```
Threads started!
```

```
CPU speed:
```

```
events per second: 1318.99
```

```
General statistics:
```

```
total time: 10.0006s
```

```
total number of events: 13193
```

```
Latency (ms):
```

```
min: 0.75
```

```
avg: 0.76
```

```
max: 2.15
```

```
95th percentile: 0.78
```

```
sum: 9999.03
```

```
Threads fairness:
```

```
events (avg/stddev): 13193.0000/0.00
```

```
execution time (avg/stddev): 9.9990/0.00
```

2.5 内存性能测试

测试目的:通过sysbench工具测试真实场景下内存性能和延迟

测试工具: sysbench

测试结果: 内存顺序读取模式下的持续带宽可达到30.6GiB/s。

测试截图:

```
root@stone:~# sysbench memory --memory-block-size=1M --memory-total-size=10G --memory-oper=read run
sysbench 1.0.20 (using system LuaJIT 2.1.0-beta3)
```

```
Running the test with following options:
Number of threads: 1
Initializing random number generator from current time
```

```
Running memory speed test with the following options:
block size: 1024KiB
total size: 10240MiB
operation: read
scope: global
```

```
Initializing worker threads...
```

```
Threads started!
```

```
Total operations: 10240 (30637.02 per second)
```

```
10240.00 MiB transferred (30637.02 MiB/sec)
```

```
General statistics:
total time:                0.3328s
total number of events:    10240
```

```
Latency (ms):
min:                    0.03
avg:                    0.03
max:                    0.09
95th percentile:      0.03
sum:                    331.41
```

```
Threads fairness:
events (avg/stddev):    10240.0000/0.00
execution time (avg/stddev): 0.3314/0.00
```


2.6 硬盘性能测试

测试目的:通过fio工具测试真实场景下硬盘顺序读写性能

测试工具: fio

测试结果: 硬盘在顺序读取下的速度可达分别到7437MB/s, 4379 MB/s

测试截图:

```
root@stone:~# fio -filename=/dev/nvme0n1 -direct=1 -iodepth 256 -thread -rw=read -ioengine=libaio -bs=128k -size=200G -numjobs=4 -run
lae=60 -group_reporting -name=test
test: (g=0): rw=read, bs=(R) 128KiB-128KiB, (W) 128KiB-128KiB, (T) 128KiB-128KiB, ioengine=libaio, iodepth=256
...
fio-3.28
Starting 4 threads
Jobs: 4 (f=4): [W(4)][100.0%][r=7128MiB/s][w=57.6k IOPS][eta 00m:00s]
test: (groupid=0, jobs=4): err=0: pid=43291: Fri Sep 5 08:28:49 2025
read: IOPS=56.7k, BW=7093MiB/s (7437MB/s)(4166iB/60010msec)
slat (usec): min=4, max=28972, avg=28.45, stdev=398.53
clat (msec): min=2, max=154, avg=17.99, stdev=4.91
lat (msec): min=2, max=155, avg=18.02, stdev=4.91
clat percentiles (msec):
| 1.00th=[ 8], 5.00th=[ 9], 10.00th=[ 13], 20.00th=[ 16],
| 30.00th=[ 17], 40.00th=[ 18], 50.00th=[ 18], 60.00th=[ 18],
| 70.00th=[ 19], 80.00th=[ 22], 90.00th=[ 25], 95.00th=[ 26],
| 99.00th=[ 30], 99.50th=[ 31], 99.90th=[ 38], 99.95th=[ 42],
| 99.99th=[ 114]
bw ( MiB/s): min=4993, max=8110, per=100.00%, avg=7100.33, stdev=137.68, samples=476
tops
lat (msec): min=39948, max=64880, avg=56802.22, stdev=1101.44, samples=476
lat (msec): 4=6.92%, 10=8.72%, 20=65.93%, 50=25.29%, 100=6.03%
cpu
: usr=5.93%, sys=36.99%, ctx=1181662, majf=1, minf=122809
10 depths : 1=0.1%, 2=0.1%, 4=0.1%, 8=0.1%, 16=9.1%, 32=0.1%, >=64=100.0%
submit : 0=0.0%, 4=100.0%, 8=0.0%, 16=0.0%, 32=0.0%, 64=0.0%, >=64=0.0%
complete : 0=0.0%, 4=100.0%, 8=0.0%, 16=0.0%, 32=0.0%, 64=0.0%, >=64=0.1%
issued rwts: total=3405564,0,0,0 short=0,0,0,0 dropped=0,0,0
latency : target=0, window=0, percentile=100.00%, depth=256

Run status group 0 (all jobs):
READ: bw=7093MiB/s (7437MB/s), 7093MiB/s-7093MiB/s (7437MB/s-7437MB/s), io=4166iB (4460B), run=60018-60018msec

Disk stats (read/write):
nvme0n1: ios=3397396/0, merge=0/0, ticks=48004009/0, in_queue=48004009, util=99.89%
root@stone:~#

root@stone:~# fio -filename=/dev/nvme0n1 -direct=1 -iodepth 256 -thread -rw=write -ioengine=libaio -bs=128k -size=200G -numjobs=4 -run
lae=60 -group_reporting -name=test
test: (g=0): rw=write, bs=(R) 128KiB-128KiB, (W) 128KiB-128KiB, (T) 128KiB-128KiB, ioengine=libaio, iodepth=256
...
fio-3.28
Starting 4 threads
Jobs: 4 (f=4): [W(4)][100.0%][w=4740MiB/s][r=38.0k IOPS][eta 00m:00s]
test: (groupid=0, jobs=4): err=0: pid=43427: Fri Sep 5 08:27:18 2025
write: IOPS=33.4k, BW=4177MiB/s (4379MB/s)(2450iB/60020msec); 0 zone resets
slat (usec): min=5, max=32038, avg=73.51, stdev=663.00
clat (msec): min=5, max=126, avg=30.56, stdev=7.75
lat (msec): min=5, max=126, avg=36.63, stdev=7.80
clat percentiles (msec):
| 1.00th=[ 15], 5.00th=[ 22], 10.00th=[ 25], 20.00th=[ 27],
| 30.00th=[ 29], 40.00th=[ 29], 50.00th=[ 29], 60.00th=[ 30],
| 70.00th=[ 33], 80.00th=[ 34], 90.00th=[ 38], 95.00th=[ 39],
| 99.00th=[ 54], 99.50th=[ 59], 99.90th=[ 89], 99.95th=[ 94],
| 99.99th=[ 111]
bw ( MiB/s): min=2821, max=5806, per=100.00%, avg=4177.57, stdev=177.22, samples=476
tops
lat (msec): min=22574, max=46454, avg=33420.39, stdev=1417.71, samples=476
lat (msec): 10=0.40%, 20=2.21%, 50=91.67%, 100=3.62%, 250=0.02%
cpu
: usr=25.58%, sys=19.86%, ctx=736561, majf=0, minf=163719
10 depths : 1=0.1%, 2=0.1%, 4=0.1%, 8=0.1%, 16=0.1%, 32=0.1%, >=64=100.0%
submit : 0=0.0%, 4=100.0%, 8=0.0%, 16=0.0%, 32=0.0%, 64=0.0%, >=64=0.0%
complete : 0=0.0%, 4=100.0%, 8=0.0%, 16=0.0%, 32=0.0%, 64=0.0%, >=64=0.1%
issued rwts: total=0,2005696,0,0 short=0,0,0,0 dropped=0,0,0
latency : target=0, window=0, percentile=100.00%, depth=256

Run status group 0 (all jobs):
WRITE: bw=4177MiB/s (4379MB/s), 4177MiB/s-4177MiB/s (4379MB/s-4379MB/s), io=2450iB (2630B), run=60028-60028msec

Disk stats (read/write):
nvme0n1: ios=181/2088455, merge=0/0, ticks=3/32702481, in_queue=32702465, util=98.65%
root@stone:~#
```