

Experience With Composable GPUs

Paul Walker

Composability at Durham - Background

- Some context
- Pilot / test systems from 2 manufacturers
- Interfaces for admin / users
- Experience

Cosma

- Cosma 5 – legacy system due for upgrade
- Cosma 7 – 449 Skylake nodes
- Cosma8 – 528 Epyc nodes
- Hardware lab



Background

Key points in evolution towards distributed computing resources:

- Storage – 1980s: NAS, 1990s: SAN, FibreChannel, Lustre
- Processing – 1990s, 2000s: Message Passing Interface etc.
- GPU – 2000s to present: SLI, NVLink, composability
- Memory – Now(ish): Compute Express Link (CXL)

Move from tightly coupled, monolithic systems toward increasingly flexible, composable architectures.

Each resource type can be optimized and allocated independently.

Liquid Composable GPU

- 3 nodes – formerly 4
- Dell R6525 / R7525, AMD Epyc, 1TB / 4TB RAM
- 3 x Nvidia A100, 40GB
- NVME
- PCIe switch fabric



172.16.188.1 / Connected

Node ID:

9



System



EXPAND ALL



COLLAPSE ALL



Fabric 72 Fabric ID : 72

+ CREATE GROUP



0



0



0



0



0



0



cosma8 Group ID : 1



1



0



8



0



0



0



mad04 Machine ID : 2



RESTART



ON



CPU

1



GPU

1



SSD

8



NIC

0



FPGA

0



MEM

0

pcpu0

gpu2 (A100 PCIe 40G...

scm1 (NVM Expre...

scm5 (NVM Expre...

scm9 (NVM Expre...

scm13 (NVM Expr...

scm17 (NVM Expr...

scm21 (NVM Expr...

scm25 (NVM Expr...



mad05 Machine ID : 3



1



1



8



0



0



0



mad06 Machine ID : 4



1



1



8



0



0



0

 RE-CENTER



← BACK

cosma8 ▼mad06 ▼

Machine Edit: mad06

✕ CANCEL↺ REPROGRAM & REBOOT✓ REPROGRAM

Group Free Pool



CPUs

1

pcpu2

Gen3

x0

ON



SSDs

8

scm0

Gen3

x4

intel



scm4

Gen3

x4

intel



scm8

Gen3

x4

intel



scm12

Gen3

x4

intel



scm16

Gen3

x4

intel



scm20

Gen3

x4

intel



scm24

Gen3

x4

intel



scm28

Gen3

x4

intel



Assigned



CPUs

1 / 1

pcpu3

Gen3

x0

ON



GPUs

1 / 10

gpu0

Gen3

x16

Tesla_A100



SSDs

8 / 10

scm3

Gen3

x4

intel



scm7

Gen3

x4

intel



scm11

Gen3

x4

intel



scm15

Gen3

x4

intel



scm19

Gen3

x4

intel



scm23

Gen3

x4

intel



scm27

Gen3

x4

intel



scm31

Gen3

x4

intel



Liquid - Experience

- Good points
- Not so good points

```

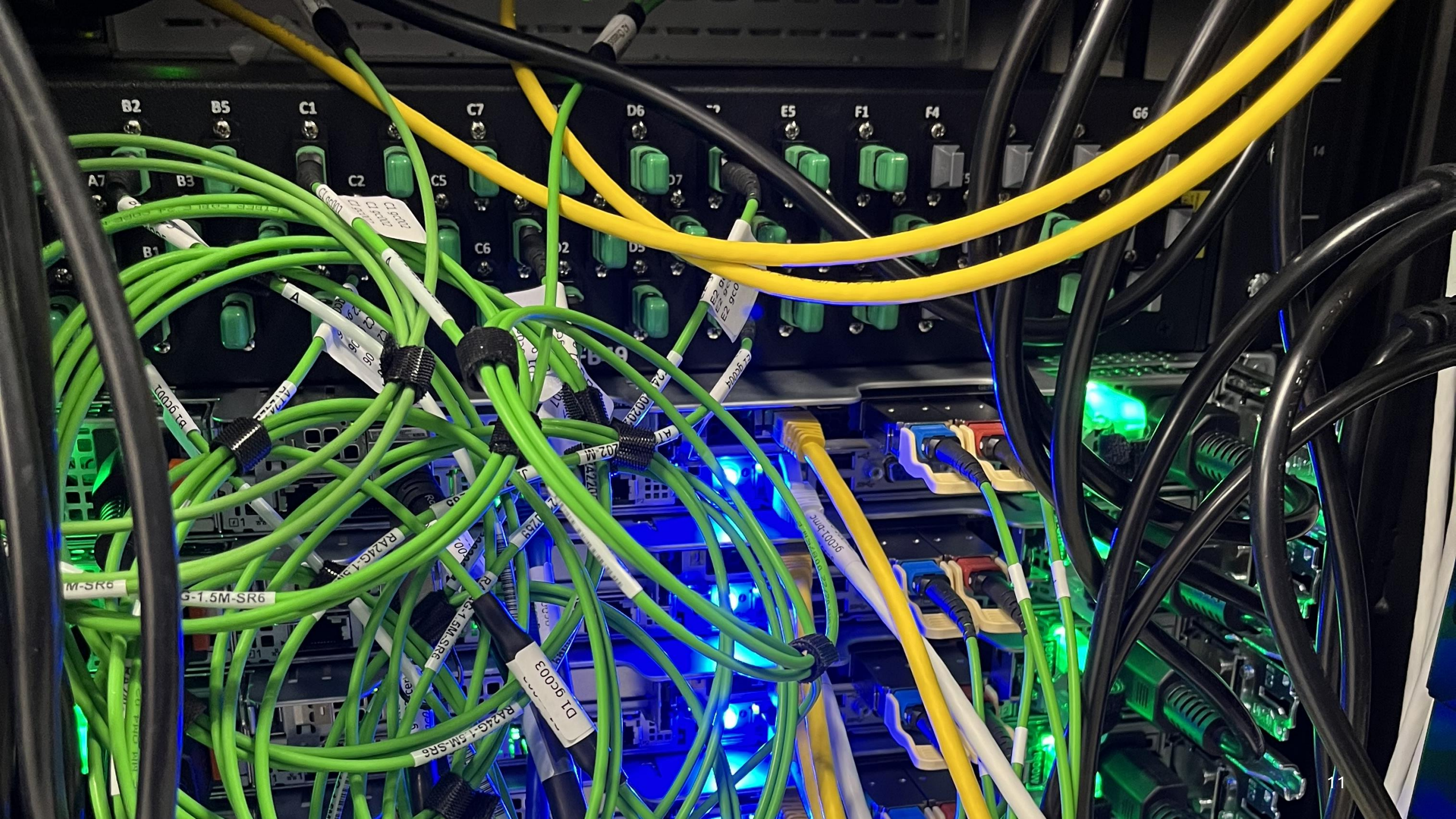
+-----+-----+-----+-----+-----+-----+
| NVIDIA-SMI 570.86.10              Driver Version: 570.86.10      CUDA Version: 12.8     |
+-----+-----+-----+-----+-----+-----+
| GPU  Name                Persistence-M | Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp   Perf          Pwr:Usage/Cap |      Memory-Usage | GPU-Util  Compute M. |
|                                           | MIG M.         |
+-----+-----+-----+-----+-----+-----+
|   0   NVIDIA A100-PCIE-40GB         Off | 00000000:F8:00.0 Off |   0          Default |
| N/A   30C    P0              33W / 250W |  4MiB / 40960MiB |      0%      Disabled |
+-----+-----+-----+-----+-----+-----+

+-----+-----+-----+-----+-----+-----+
| Processes:                         |
| GPU   GI    CI          PID    Type   Process name                      GPU Memory |
|      ID     ID                                   |             Usage   |
+-----+-----+-----+-----+-----+-----+
| No running processes found         |
+-----+-----+-----+-----+-----+-----+

```

CerIO Composable GPU

- 8 nodes
- Dell R660, Intel Sapphire Rapids, 2TB RAM
- 8 x Nvidia A30, 24GB
- Novel switchless fibre fabric







CerIO Composable GPU

```
[cerio@login0a cli]$ ./cerio device list
```

NodeName	Device	Slot	PCIeDevice	PCIeVendor	PCIeClass	State
gc-chas-l	gpu1	105	GA100GL [A30 PCIe]	NVIDIA Corporation	3D controller	Inserted
gc-chas-l	gpu2	104	GA100GL [A30 PCIe]	NVIDIA Corporation	3D controller	Inserted
gc-chas-l	gpu3	101	GA100GL [A30 PCIe]	NVIDIA Corporation	3D controller	Inserted
gc-chas-l	gpu4	102	GA100GL [A30 PCIe]	NVIDIA Corporation	3D controller	Inserted
gc-chas-r	gpu1	105	GA100GL [A30 PCIe]	NVIDIA Corporation	3D controller	Inserted
gc-chas-r	gpu2	104	GA100GL [A30 PCIe]	NVIDIA Corporation	3D controller	Inserted
gc-chas-r	gpu3	101	GA100GL [A30 PCIe]	NVIDIA Corporation	3D controller	Inserted
gc-chas-r	gpu4	102	GA100GL [A30 PCIe]	NVIDIA Corporation	3D controller	Inserted

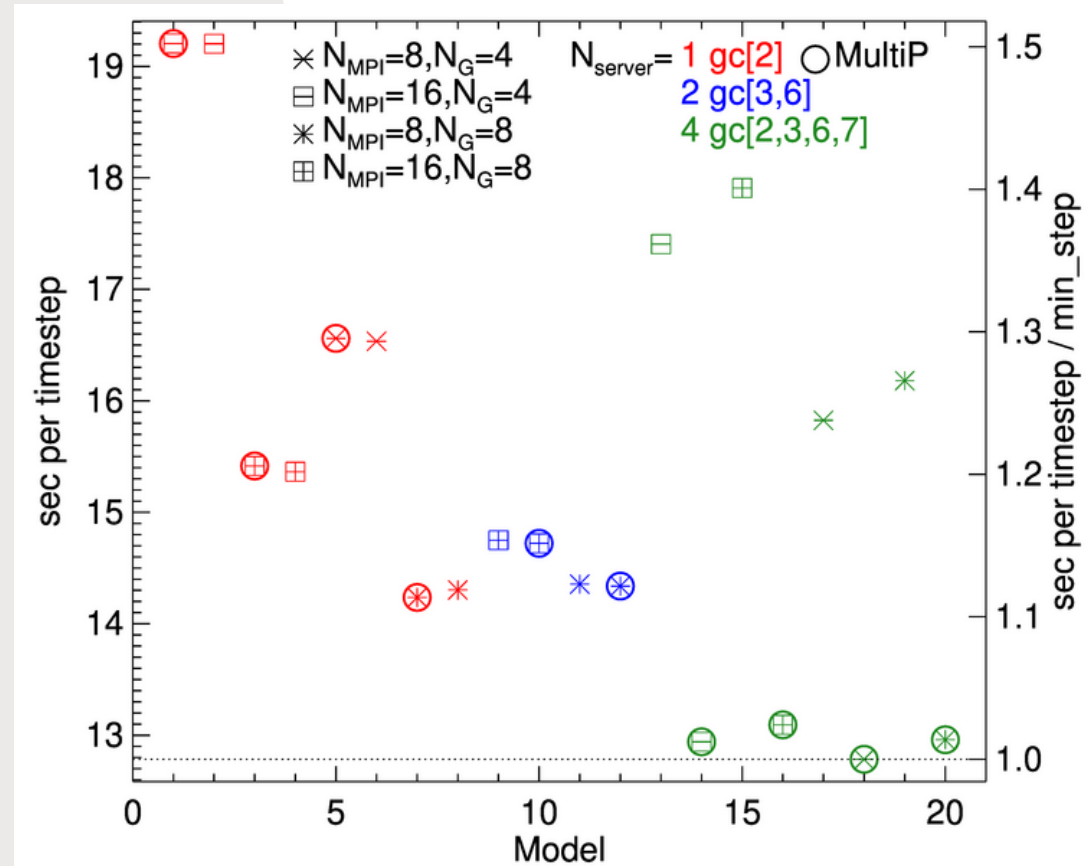
```
[cerio@login0a cli]$ ./cerio node list
```

Name	NodeID	SerialNumber	IpAddress	Version	CommState	UUID
gc-chas-l	70	S24030022	172.17.186.111	ENG-DASH32-C1.2.0-223	true	
gc-chas-r	71	S24030023	172.17.186.112	ENG-DASH32-C1.2.0-223	true	
gc001-cerio	69	S24030021	172.17.186.121	ENG-DASH31-C1.2.0-223	true	
gc002-cerio	68	S24030020	172.17.186.122	ENG-DASH31-C1.2.0-223	true	
gc003-cerio	67	S24030019	172.17.186.123	ENG-DASH31-C1.2.0-223	true	
gc004-cerio	65	S24030017	172.17.186.124	ENG-DASH31-C1.2.0-223	true	
gc005-cerio	64	S24030016	172.17.186.125	ENG-DASH31-C1.2.0-223	true	
gc006-cerio	63	S24030015	172.17.186.126	ENG-DASH31-C1.2.0-223	true	
gc007-cerio	62	S24030014	172.17.186.127	ENG-DASH31-C1.2.0-223	true	
gc008-cerio	61	S24030013	172.17.186.128	ENG-DASH31-C1.2.0-223	true	

```
[cerio@login0a cli]$ ./cerio comp list
```

CompUID	Host	HostComp	DSP	Bus(Pri/Sec/Sub)	Chassis	ChassisComp	ChassisMode	Device	DevState	Persist	Type
GC003	gc003-cerio	Activated	0	0xae/0xb0/0xb0	gc-chas-r	Activated	device	gpu1	Inserted	disabled	static
GC003	gc003-cerio	Activated	1	0xae/0xb1/0xb1	gc-chas-r	Activated	device	gpu2	Inserted	disabled	static
GC006	gc006-cerio	Activated	0	0xae/0xb0/0xb0	gc-chas-l	Activated	device	gpu3	Inserted	disabled	static
GC006	gc006-cerio	Activated	1	0xae/0xb1/0xb1	gc-chas-l	Activated	device	gpu4	Inserted	disabled	static
GC007	gc007-cerio	Activated	0	0xae/0xb0/0xb0	gc-chas-r	Activated	device	gpu3	Inserted	disabled	static
GC007	gc007-cerio	Activated	1	0xae/0xb1/0xb1	gc-chas-r	Activated	device	gpu4	Inserted	disabled	static
GC002	gc002-cerio	Activated	0	0xae/0xb0/0xb0	gc-chas-l	Activated	device	gpu1	Inserted	disabled	static
GC002	gc002-cerio	Activated	1	0xae/0xb1/0xb1	gc-chas-l	Activated	device	gpu2	Inserted	disabled	static

CerIO Benchmarking



CerIO - Experience

- Good points
- Not so good points

Conclusion - Wishlist

- Early days for technology but shows promise
- Current systems have issues, as to be expected
- Dynamic composition
- Integration with SLURM
- More profiling / benchmarking