

Disaggregation : future of data centers ?

FRNOG 38 Meeting – Oct. 6th, 2023

Xavier LE VAILLANT 2crsi.com

Disclaimer: The views and opinions expressed in this presentation are those of the speaker and do not necessarily reflect the views or positions of any entities he represent.

First 1U server supporting up to 32 High-End GPUs !





Mona 1.12GG

High-end server based on AMD EPYC[™] 4th Gen CPU with 12x 2.5" SFF Drives

I9-inch - 1U/12x 2.5" ▲MDJ EPYC[™] 4th Gen socket 24x DDR5 @ 4800MHz 2x PCle 5.0 x16 HHHL 1x PCle 5.0 x16 2x OCP 3.0 x16 Air Cooling

Key Features



BREAKING NEWS

> SCAN THE CODE! TO DISCOVER MORE ABOUT THIS PRODUCT



FRNOG 38 Meeting | 2

First 1U server system supporting up to 32 High-End GPUs!



https://gigaio.com/supernode/





55	h user@su	permi	lan1			
U	32					
FPGA	0					
	1					
GPU	0					
	0					
twork	Adapter,	Gen4	(FA400-	4)		
		ROCm	System	Management	Interface	

Xili - NVIDI - Infi

FabreX N

GPU

		===== Co	ncise Inf	0 ====				******		
Temp (DieEdge)	AvgPwr	SCLK	MCLK							
53.0c	42.0W	800Mhz	1600Mhz							010110
57.0c	44.0W	800Mhz	1600Mhz					Re	sNet50 Results	
66.0c	44.0W	800Mhz	1600Mhz						contector neodato	
59.0c	43.0W	800Mhz	1600Mhz							
62.0c	44.0W	800Mhz	1600Mhz		80	ALEE	ACTO	D - 0	00%	
57.0c	42.0W	800Mhz	1600Mhz		- 30	ALCIN	ACTO	/n -	3310	
59.0c	42.0W	800Mhz	1600Mhz							
61.0c	42.0W	800Mhz	1600Mhz							
48.0c	42.0W	800Mhz	1600Mhz							
49.0c	40.0W	800Mhz	1600Mhz	0						
53.0c	40.0W	800Mhz	1600Mhz	26						
51.0c	42.0W	800Mhz	1600Mhz	18						
52.0c	43.0W	800Mhz	1600Mhz	9						
50.0c	40.0W	800Mhz	1600Mhz	E .						
53.0c	42.0W	800Mhz	1600Mhz	-						
53.0c	42.0W	800Mhz	1600Mhz					-		
48.0c	42.0W	800Mhz	1600Mhz				-			
51.0c	42.0W	800Mhz	1600Mhz			-				
59.0c	45.0W	800Mhz	1600Mhz							
55.0c	44.0W	800Mhz	1600Mhz		-					
63.0c	44.0W	800Mhz	1600Mhz							
59.0c	47.0W	800Mhz	1600Mhz		12	3456	789	10 11 12	13 14 15 16 17 18 19 20 21 2	2 23 24 25 26 27 28 29 30 31 32
55.0c	43.0W	800Mhz	1600Mhz							
52.0c	44.0W	800Mhz	1600Mhz						Number of GPUs	
57.0c	44.0W	800Mhz	1600Mhz							
55.0c	42.0W	800Mhz	1600Mhz		-			~ ~		
64.0c	44.0W	800Mhz	1600Mhz	0%	auto	300.0W	0%	0%		
55.0c	42.0W	800Mhz	1600Mhz	0%	auto	300.0W	0%	0%		
58.0c	42.0W	800Mhz	1600Mhz	0%	auto	300.0W	0%	0%		
56.0c	45.0W	800Mhz	1600Mhz	0%	auto	300.0W	0%	0%		
65.0c	45.0W	800Mhz	1600Mhz	0%	auto	300.0W	0%	0%		
60.0c										
	-	Dr. M	oritz Leh	mann	L					
	- 17	@Pro	jectPhys)	(
			Anterson - Sterlandsk							
	0.0	* the uu	alkand I	dot t	o toot	#Fluids	20 0	a tha u	uarldla lardaat #UD	

Over the weekend I got to test #h #GPU server, #GigalOSuperNODE. Here is one of the largest #CFD simulations ever, the Concorde for 1s at 300km/h landing speed. 40 *Billion* cells resolution. 33h runtime on 32 @AMDInstinct MI210, 2TB VRAM.



2Crsi

More computing power... with Accelerators

Systems Using Accelerators on the TOP500





e.g., different research activities



Mechanical Engineering FP32, Vector Engines



NETWORK	lob 1 - Engineering	
		2
		3
	Job 2 – Bioinformatics	5
STODACE		100.
		3
	0.55 (6)	8
ACCELERATORS	0.25 444	8
50 50 50 50 50 50 50 50 50 50		
The second		
	Job 3 – Computer Scie	ence
PERSISTENT MEMORY	Job 3 – Computer Scie	ence
PERSISTENT MEMORY	Job 3 – Computer Scie	ence
PERSISTENT MEMORY	Job 3 – Computer Scie	ence
PERSISTENT MEMORY	Job 3 – Computer Scie	ence

Defines hardware uniquely for each workload

→ The End of Stranded Resources

• The right resources in the right place, on demand

→ Scale-Up and Scale-Out as You Grow

 In-place scaling and selection of servers and accelerators as requirements evolve

⊖ True Heterogeneity

 The right GPUs, servers and accelerators for the job



Composable Disaggregated Infrastructure

Composable Disaggregated Infrastructure (CDI) brings the agility, savings and efficient resource-sharing of the cloud to the management of on-premises equipment.

Using orchestration and high-bandwidth, low-latency fabrics, shared resources can be combined on-demand for shifting workloads. The goal is to get the right ratio for a specific AI training or inference job, change configurations as the workload pipeline changes, and free up expensive GPUs and other accelerators for additional work.

"Software-Defined Hardware"





PCIe-based Fabric



Example of a Rack scale system



PCrsi

PCIe-based Fabric



Example of a Rack scale system



PCIe-based Fabric



Composable Resources





Summary of Data Center Memory Challenges



Decoupling the Memory Controller from the CPU and Providing Options for New Server Architectures to Address Memory Challenges

CXL Forum at FMS2023 - Advancing Data Center Architectures with Memory Tiering, Rambus



Scope of CXL







Source: CXL Forum ISC 23 - The March of Composability – Onward to Memory with CXL, GigalO



Scope of CXL





Memory tiers – Span the Latency Gap



- CXL delivers new expansion options for hot DRAM, with no impact to software applications
- CXL also introduces memory tiering, to the Data Center, much like storage tiering before it
- The industry is now working on software infrastructure to take advantage of these new tiers

Source: CXL Forum at FMS2023 - Advancing Data Center Architectures with Memory Tiering, Rambus



Comprehensive end-to-end CXL solutions



CXL use cases

- Expanders
- Pooling
- Switch
- Accelerators
- Electro-optics
- Re-timers
- Custom Compute
- DPUs / SmartNICs
- SSD Controllers

Source: CXL Forum at FMS 2023 - CXL: Transforming Cloud Infrastructure



CXL vision: optimal resource utilization



Source: CXL Forum at FMS 2023 - CXL: Transforming Cloud Infrastructure



Optical CXL is Required for Scaling



Copper cables struggles to support CXL scaling beyond a few servers



Composability: I/O Fabric (one of) the biggest challenge

Inefficient architecture Low ratio of GPUs to CPUs • 200+ 200 Pool of Servers and Accelerators 150 100 System Efficiency 50 **Typical Limit** Single Server PCle Expansion Disaggregated **Box Scale Rack Scale Row Scale**

Traditional System Model





•

Composability: Benefits

- Improved system utilization by more fully leveraging expensive on-premises assets.
- Flexible hardware profiles create the
 Impossible server
- Pay as you Grow Simplified system expansion and reduced system costs via modular resource-specific nodes
- Reduced Power & Cooling (Sustainability)

Better Managed Life Cycles





Composability: Challenges

• Usage and operational impacts

• Which workloads are most suitable for composability?

• Resource impacts

- Changes (if any) in application code to support composability?
- Will it increase or reduce support requirements?

• Performance impacts

- What about the latency to manage, provision, monitor, and re-claim system resources between jobs?
- Will increased physical distance also add latency?
- Scaling? How far?

• Cost impacts

• Additional network (MPI, I/O, and now PCIe...)?





Thank you !

Q & A