

Aum HPC Processor

Development under
National
Supercomputing Mission



Sanjay Wandhekar
Senior Director, HoD HPC
Technologies Group, C-DAC
sanjayw@cdac.in

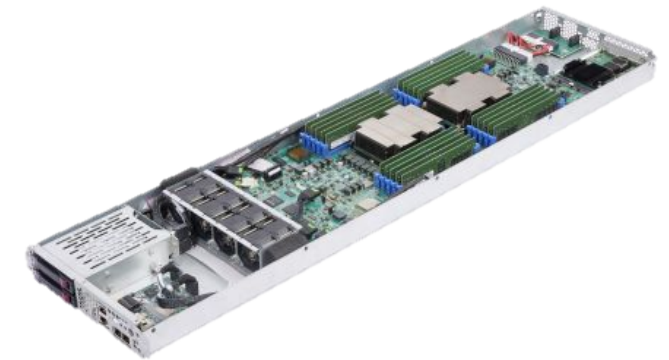
सुस्वागतम्
நல்வரவு
வழிவந்த
సుసాగతం
සුච්චාගතම
সুসাগতম
ಸುಸಾಗತ
സുസാഗതം
सुसगतम
सुआगतम
सुस्वागतम्
خوش آمدید



National Supercomputing Mission (NSM)

Brief about National Supercomputing Mission

- **Supercomputing infrastructure** in the country
- **Indigenous Supercomputing ecosystem** in phased manner: From “Assembly” to “Manufacturing” to “Design and Manufacturing” of Supercomputers
 - Servers
 - HPC network
 - Software stack
 - HPC Processor
 - Liquid cooling technologies
- **Supercomputing Applications** of National interest
- **Human resources** for applications development and HPC maintenance





Motivation for India's own HPC Processor - AUM

- Processor architecture suitable for both HPC & General Purpose Computing- Extracting maximum application level performance
- Energy efficiency: Arm Architecture
- Capability building with bargaining power
- Immunity from possible export restrictions to India in Future
- Technological sovereignty: Designed and Engineered in India
- Security (Back doors etc.) : Highest priority for strategic sectors



HPC Processor development program

- Develop a competitive HPC Processor for HPC, AI and server market
- Develop a complete ecosystem leveraging open source components
 - Open source software ecosystem
 - Reference Boards
 - Reference server designs
- Build a Pilot HPC system with > 1 PF compute power
- Be ready with Exascale system design and subsystems based on AUM Processor
- Industry Collaboration – Design of SoC, Server designs, Deploy and market solutions based on AUM Processor
- Targeted for both HPC and Cloud market
- Planned to be available in 2024

Some of the Architectural decisions

- Best Efficiency
 - Memory Bandwidth
 - Easy to optimize (Vector Size)
 - Superior Application level program / Watt
- Better I/O for Data Access
 - HBM and DDR
 - Many PCIe5 Lanes
 - CXL for Coherent accelerators
- No Competition with specialized devices like GPUs – Keep Provision of GPUs for specialized applications
- Security Features Provision

- Superior Application level program / Watt -> Increase Memory sub-system performance
- Need Much better Bytes/Flop performance – Target > 0.5 Byte/Flop



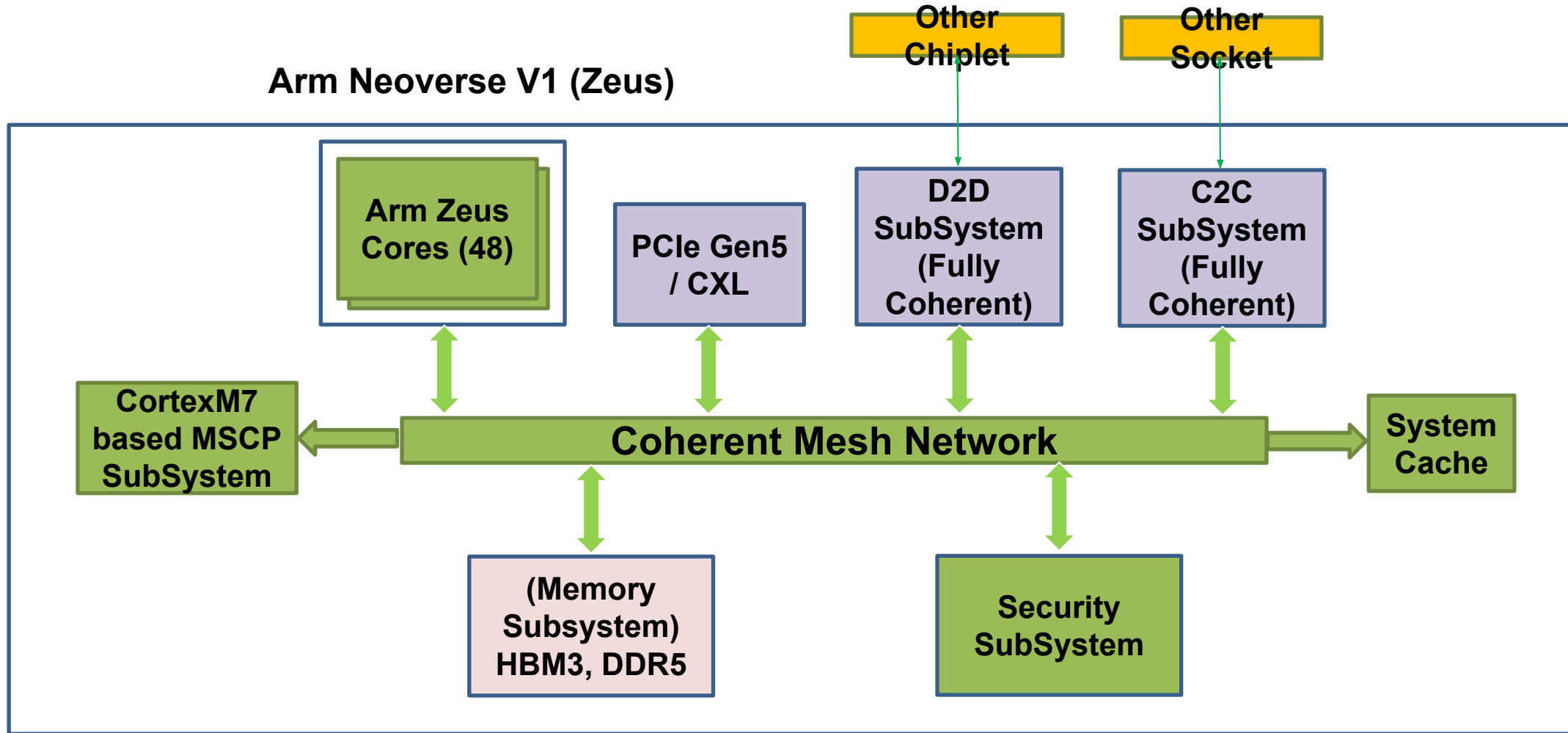
High Performance Conjugate Gradients (HPCG) Benchmark

Rank	Site	Computer	HPL Rmax (Pflop/s)	Bytes/Flop	HPCG (Pflop/s)	Fraction of Peak
1	RIKEN Center for Computational Science Japan	Supercomputer Fugaku — A64FX 48C 2.2GHz, Tofu D	442.01	0.3	16	3.00%
2	DOE/SC/ORNL USA	Summit — IBM POWER9 22C 3.07GHz, Dual-rail Mellanox EDR Infiniband, NVIDIA Volta GV100	148.6	<0.2	2.926	1.50%
3	DOE/SC/LBNL/NERSC United States	Perlmutter — AMD EPYC 7763 64C 2.45GHz, Slingshot-10, NVIDIA A100 SXM4 40 GB	64.59	<0.2	1.905	2.10%
4	DOE/NNSA/LLNL USA	Sierra — IBM POWER9 22C 3.1GHz, Dual-rail Mellanox EDR Infiniband, NVIDIA Volta GV100	94.64	<0.2	1.796	1.40%

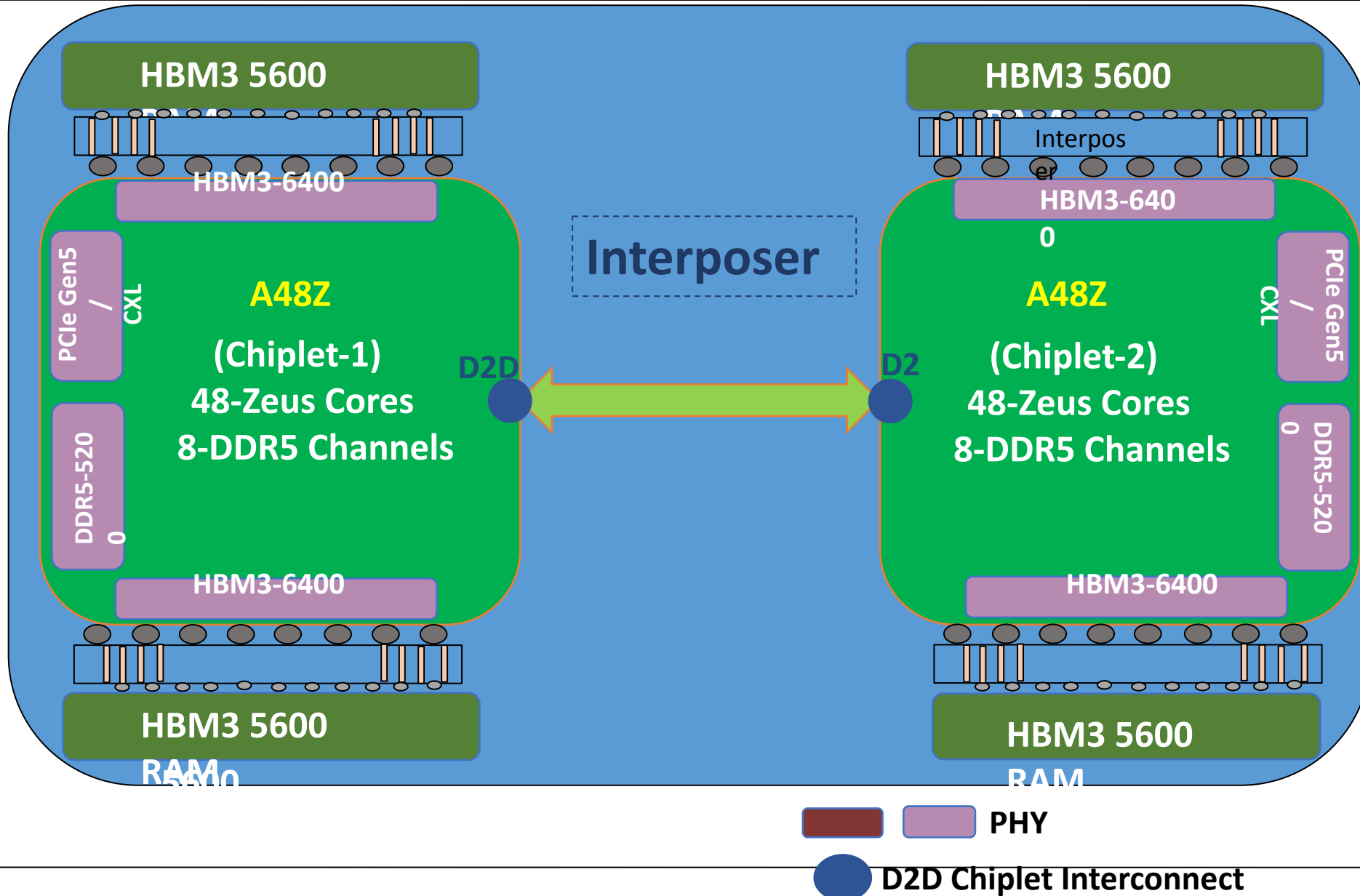
K-Computer: Bytes/Flop = 0.5
HPCG – 5.2% of Peak

Superior Application level program
Better Bytes/Flop i.e. higher Memory B/w

C-DAC HPC SoC (A48Z) Block Diagram (48-Cores)



C-DAC AUM ॐ Microprocessor - 96 Cores





AUM - HPC Processor Development



- 96 core HPC Processor
 - ARM 8.4 architecture
 - 96MB L2 cache, 96MB System cache
 - 8 channel 5200 Mhz DDR5 memory
 - **64 GB HBM3** 5600Mhz memory
 - **PCIe5 64/128 Lanes** – CXL support for coherent Accelerators/ NIC
 - SMP support up to 2 sockets
 - Security Features - Secure boot and Crypto support
 - 5nm Technology Node, Chiplet based architecture, 2-Chiplets, 96-Cores and up to 96-GB HBM3 memory in a socket
- Dual socket Server design with up to 4 Industry standard GPU accelerators – Both HPC and AI applications (CPU Only node ~ 10 TF/Node)
- Indigenous Software eco-system for Aum Processor leveraging open source eco-system



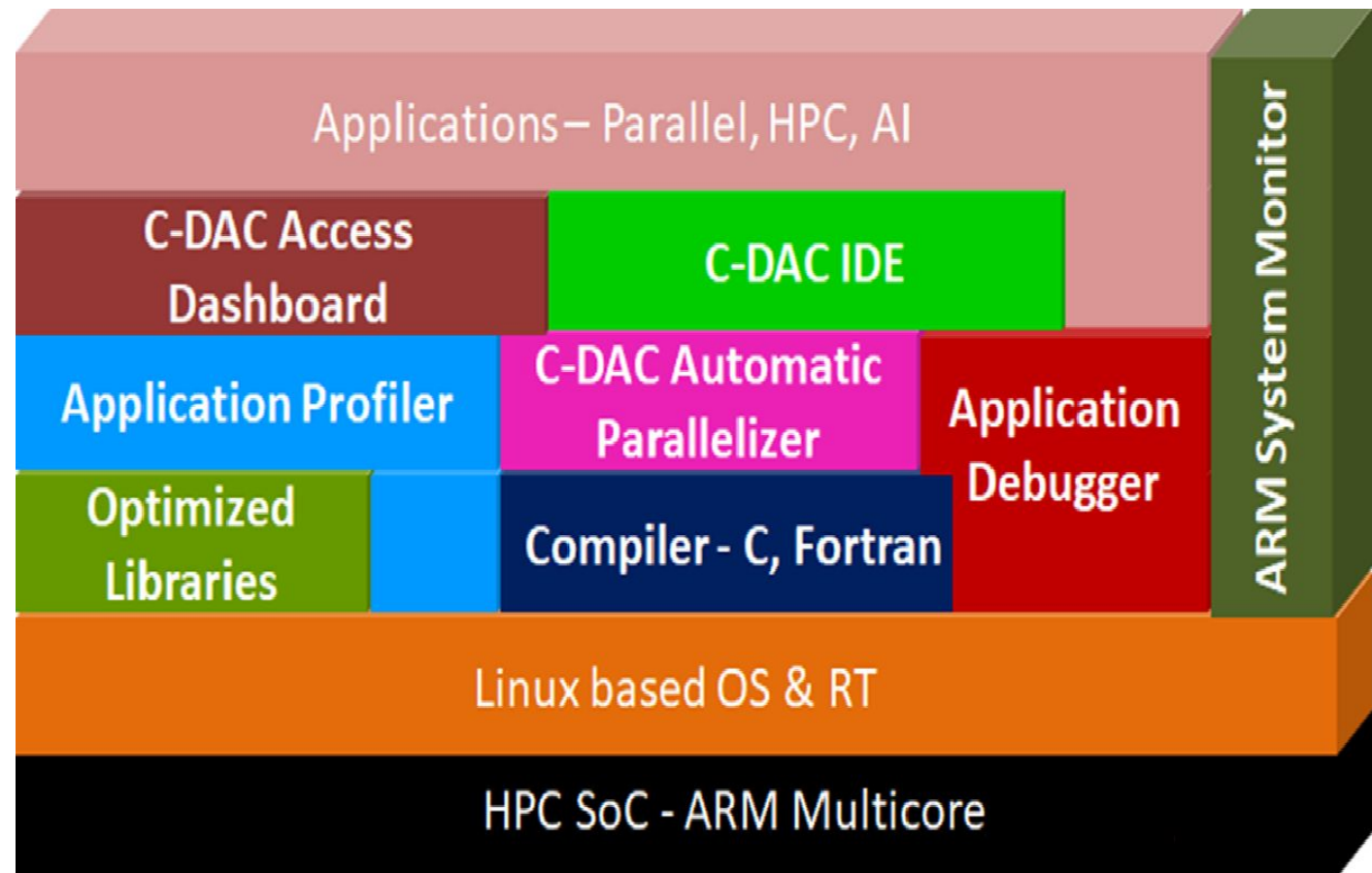
Specification Comparison



	Fujitsu A64FX	C-DAC AUM HPC Processor
Fabrication Technology	7nm FF TSMC	5nm FF
Core Configuration	(48+4)-Cores, 2.2 GHz (typical)	96-Cores, 3.0 GHz (typical) 3.5+ GHz (turbo)
DDR Configuration	No DDR	16-Channels (32 bit) DDR5-5200 BW = 332.8 GB/s
HBM	32-GB HBM2 (4-Controllers) BW = 1 TB/s	64-GB HBM3 (4-Controllers) BW = 2.87 TB/s
PCIe	16 PCIe Gen3 Lanes	64 PCIe Gen5 Lanes
Power	Not Known	300 W (TDP)
Performance (DP)	2.7 TFLOPS per socket	4.6+ TFLOPS per socket
Bytes/FLOPS	0.38	0.7

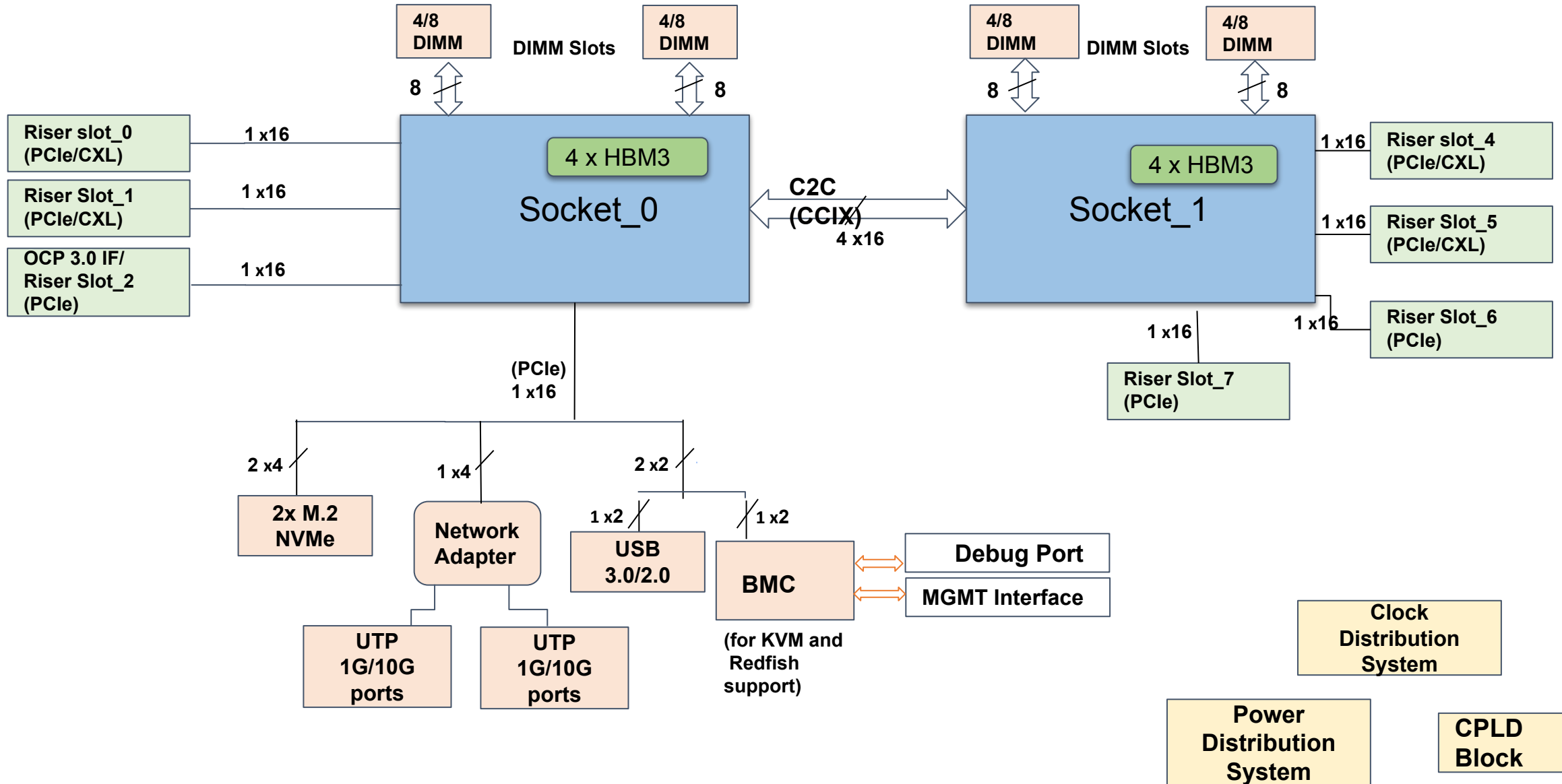
- **System Software, Dev Tools and Utilities**

- HPC Compiler (C and Fortran, multicore + accelerator, HPC & AI applications)
- IDE for HPC & AI applications on ARM system supporting multiple parallel paradigms
- Automatic Parallelizer generate parallel code for multicore / accelerators
- Application Debugger & Profiler
- Optimized Math-AI Libraries
- ARM System Monitor and Utilities
- Parallel Runtime System
- Secure Access Interface

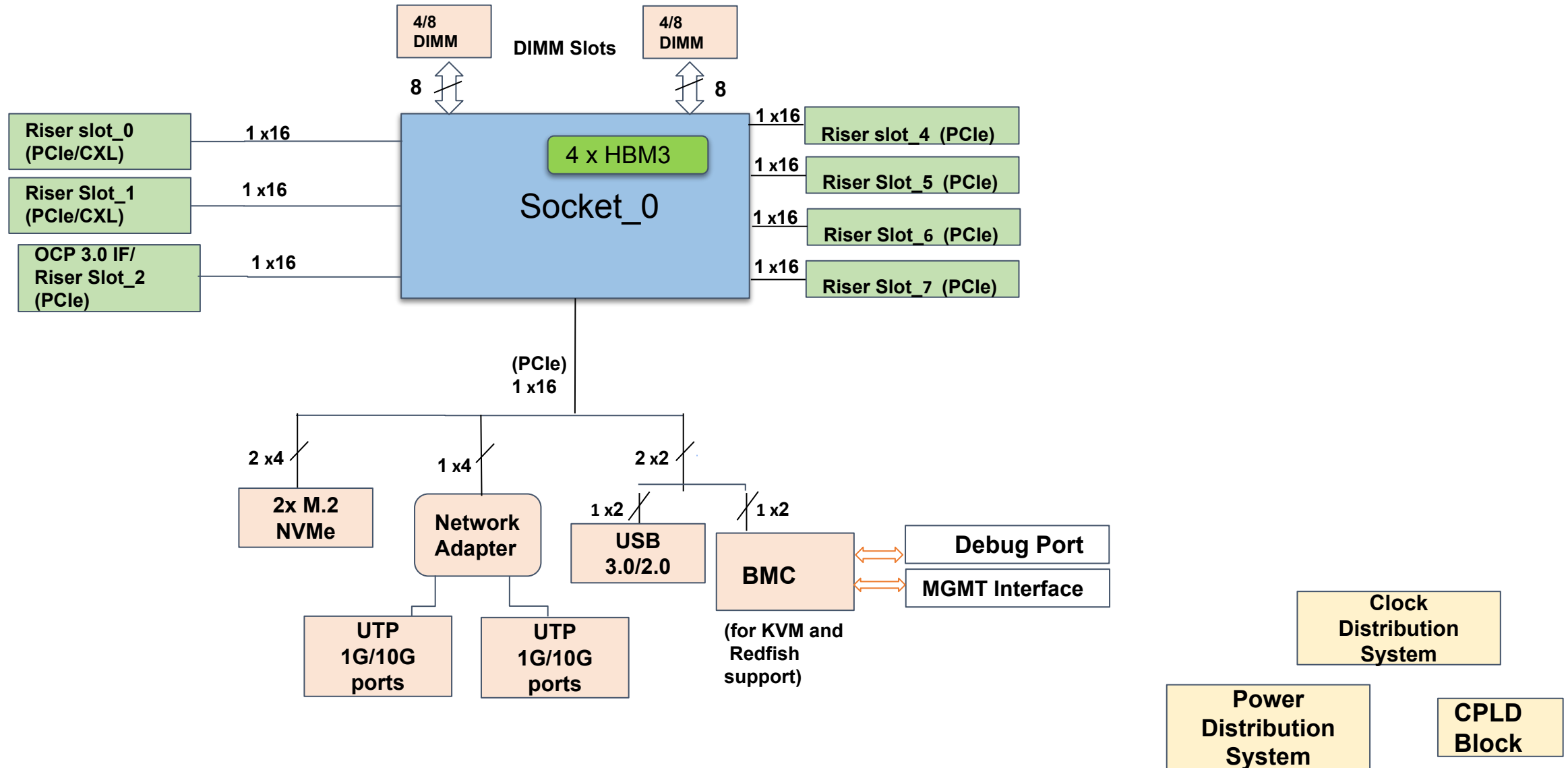


- **HPC System Middleware**

Dual Socket Compute Node: ANANTA



Single Socket Compute Node: ANANTA





Summary Aum Processor

- Competitive HPC Processor for HPC, AI and server market
- Address Strategic requirements
- Complete ecosystem
 - Open source software ecosystem
 - Reference Boards
 - Reference server designs – Derivatives as per market requirements
 - Industry partners OEMs/ODMs/Solution providers
- Towards Indigenous Exascale system including Processors
- Targeted market HPC/AI, Cloud, storage, edge computing
- Planned to be available in 2024



Thank You



Market Comparison



	Ampere Altra	SiPearl Rhea	C-DAC ॐ
Cores	80 Arm Neoverse N1 (Ares) Cores	72 Arm Neoverse V1 (Zeus) Cores	96 Arm Neoverse V1 (Zeus) Cores
Cache	L1: 64 KB L1 I / 64 KB L1 D Cache per core L2: 1 MB L2 cache per core System Cache: 32 MB	System Cache: 128 MB	L1: 64KB I-Cache / 64KB D-Cache per Core L2: 1MB Unified Cache per Core System Cache: 96MB, Snoop Filter: 192MB
Frequency	3.0 (base), 3.3 GHz (turbo)	2.5 GHz (base), 3.0 GHz (turbo)	3.0 Ghz
HBM	No HBM	96GB of HBM2E	96GB of HBM3
DDR	up to 4 TB per socket.	4 Channels DDR5	16 Channles DDR5
PCI	128 PCIe4 Lanes	104 Lanes PCIe5: - Upto 64-lanes for coherent connectivity - Remaining lanes as PCIe5	128 Lanes PCIe5: - Upto 64-lanes for coherent connectivity - Remaining lanes as PCIe5/CXL
TDP	250 W	320W	280 - 320 W
Node	TSMC 5nm	TSMC 6nm	TSMC 5nm
Package		2.5 D	2.5D
Release Year	2020	2022 /23	2023 / 24