



Master Informatics Eng.

2020/21

A.J.Proença

The move from multicore to manycore *(online)*
(some slides are borrowed)

From multicore to manycore: key issues



Lessons from Intel KNL to fit many cores into a single chip

- lower the compute capability of each core, but not too much
 - P54C (KNC) and Atom Silvermont (KNL) are too slow
- use a scalable interconnection network fabric on-chip (NoC)
 - to minimize shared-cache/memory access latencies
 - to provide enough data communication bandwidth
 - to minimize traffic bottlenecks
- group cores in clusters to improve the quality of the NoC
- reduced cache size/levels with strong impact on performance
- smaller fabrication processes (KNL: 14nm; Icelake: 10nm; Apple M1: 5nm)
- mix general-purpose PUs with application-oriented modules:
GPUs for vector computing, NNP for tensor computing, ...
- move to MCM/chiplets: simpler chips have better wafer production yield

Interconnect Fundamentals



- **Networks-on-chip:** an adapted follow-up of interconnection systems to link servers in supercomputers
- Key parameters that define a **NoC**:
 - **topology:** defines how the nodes and links are connected, namely all possible paths a message can take through the network
 - **routing algorithm:** selects the specific path a message will take from source to destination
 - **flow control protocol:** determines how a message actually traverses the assigned route
 - **router micro architecture:** implements the routing and flow control protocols and critically shapes its circuits

Manycore chips/packages: an overview



Key server chips/packages that addresses those issues:

- Intel: the Xeon Processor Scalable family
- AMD: the Epyc Zen family
- Sunway: the SX260x0 family
- ARM: the ARMv8 server-level competitors
 - Marvell ThunderX family
 - Fujitsu A64FX Arm chip
 - Neoverse N1 hyperscale reference design
 - Ampere Altra Arm Processor
 - Amazon Graviton
 - Huawei HiSilicon Kunpeng 920
- Cerebras: a Wafer Scale Engine
- Apple (*not server...*): the SoC approach (*no chipllets!*)

Manycore chips/packages: an overview



Key server chips/packages that addresses those issues:

– **Intel: the Xeon Processor Scalable family**

– AMD: the Epyc Zen family

– Sunway: the SX260x0 family

– ARM: the ARMv8 server-level competitors

- Marvell ThunderX family

- Fujitsu A64FX Arm chip

- Neoverse N1 hyperscale reference design

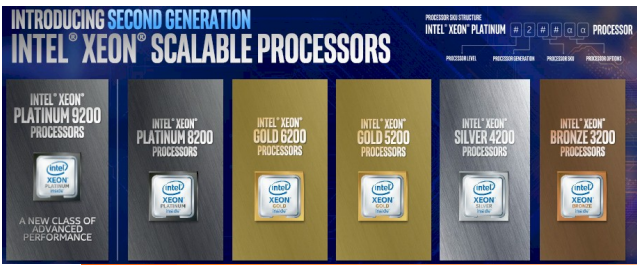
- Ampere Altra Arm Processor

- Amazon Graviton

- Huawei HiSilicon Kunpeng 920

– Cerebras: a Wafer Scale Engine

– Apple (*not server*): the SoC approach (*no chiplets!*)



Intel Xeon Scalable Processor

(formerly code-named Skylake-SP)

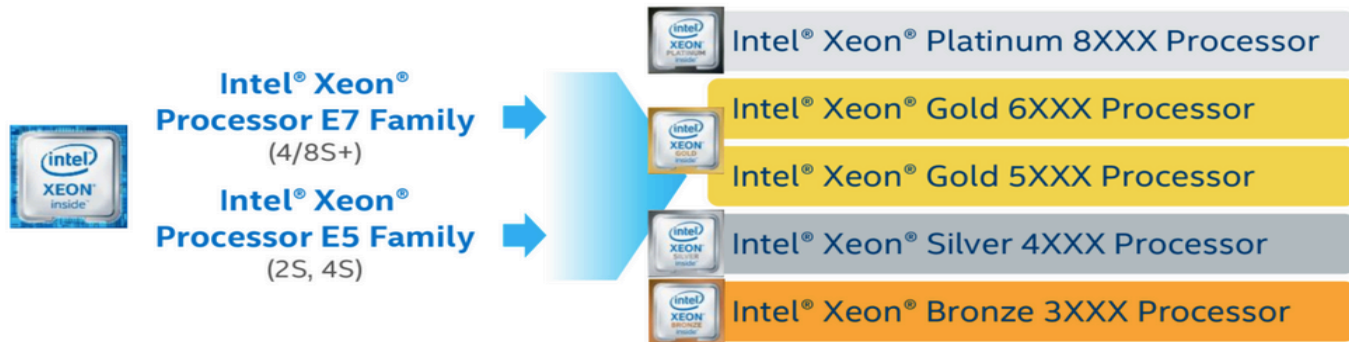
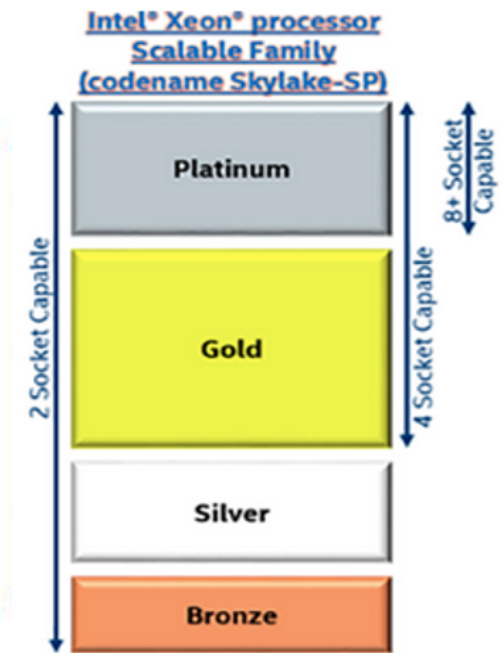


First Generation Intel® Xeon® Scalable Processor

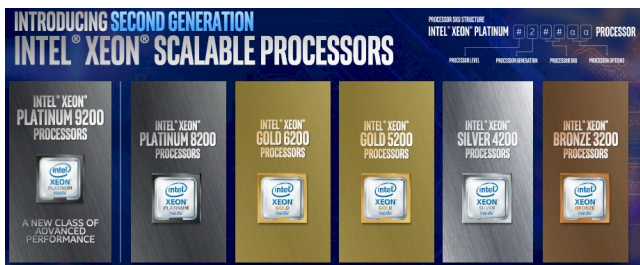
Introduced in July 2017

- Skylake-SP core microarchitecture with data center specific enhancements
- Intel® AVX-512 with 32 DP flops per cycle per core
- Data center optimized cache hierarchy – 1MB L2 per core, non-inclusive L3
- New Intel® Mesh architecture
- Enhanced 6 channel memory subsystem
- 48 lanes of PCIe Gen3 with integrated DMA, NTB, and VMD devices
- New Intel® Ultra Path Interconnect (Intel® UPI)

Features	Intel® Xeon® Scalable Processor
Cores and Threads Per CPU	Up to 28 cores and 56 threads
Last-level Cache (LLC)	Up to 38.5 MB (non-inclusive)
QPI/UPI Speed (GT/s)	Up to 3x UPI @ 10.4 GT/s
PCIe® Lanes/ Controllers	Up to 48 / 12 / PCIe 3.0 (2.5, 5, 8 GT/s)
Memory Population	Up to 6 channels of up to 2 RDIMMs, LRDIMMs, or 3DS LRDIMMs
Max Memory Speed	Up to 2666 MHz



AJP



PU generations for HPC: Intel Xeon Processor Scalable Family

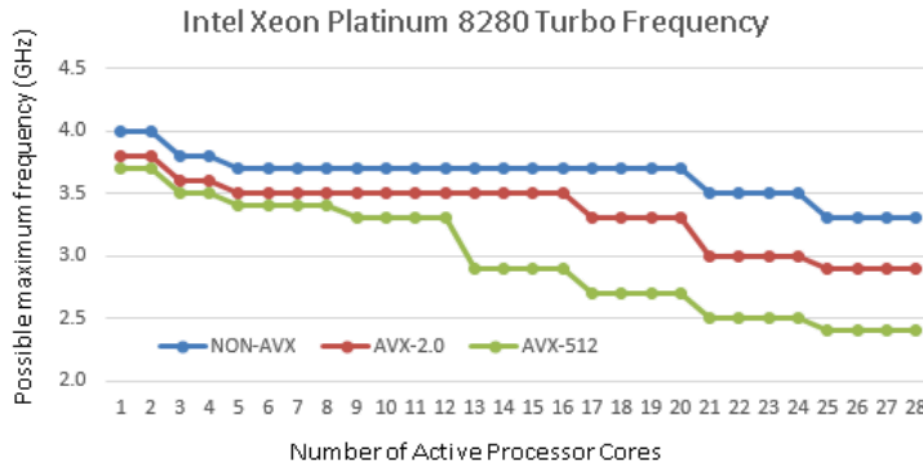
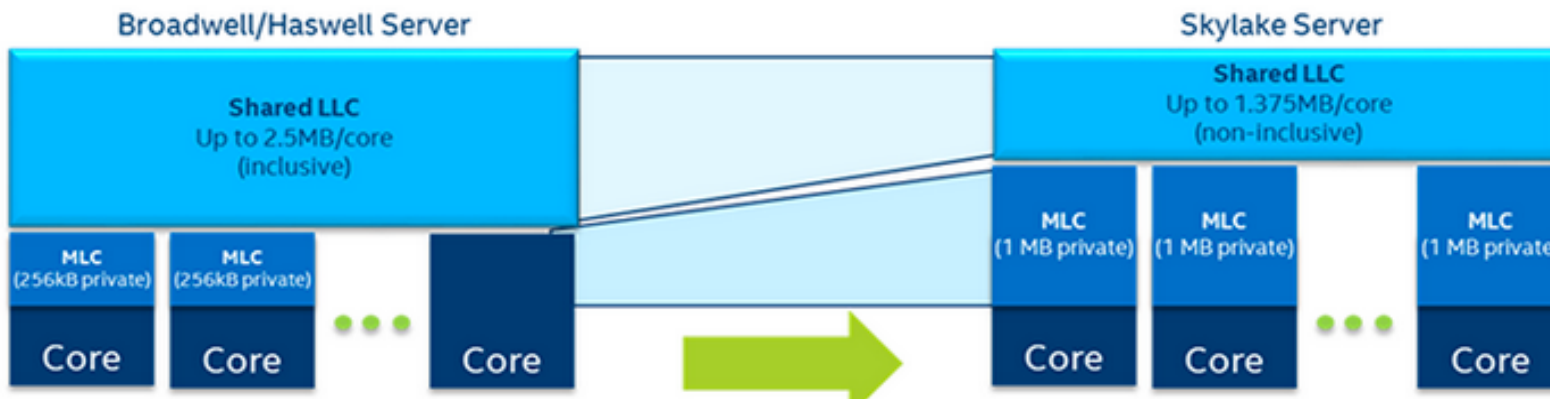


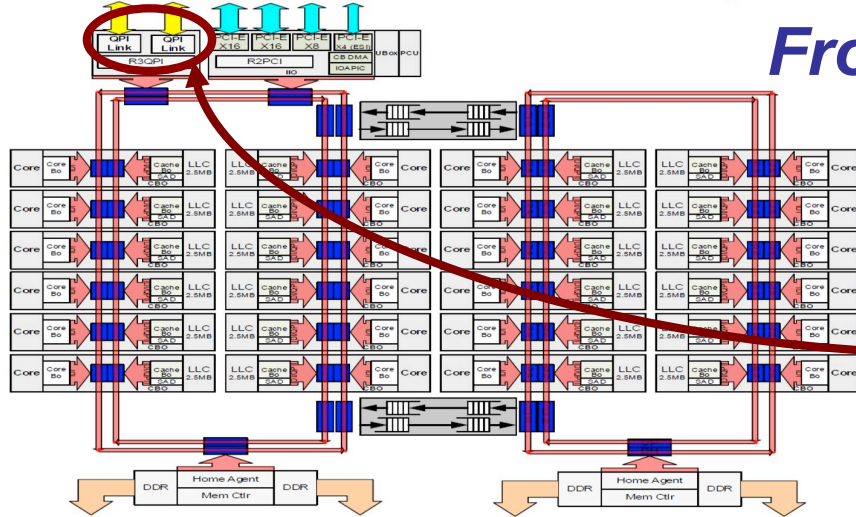
Figure 1 Intel Xeon Platinum 8280 Turbo Frequency

Grantley Platform		Purley Platform	
Intel® Microarchitecture Codenamed Haswell		Intel® Microarchitecture Codenamed Skylake	
Haswell	Broadwell	Skylake-SP	Cascade Lake-SP
22nm	14nm	14nm	14nm
New Micro-architecture		New Micro-architecture	
Features		Cascade Lake CPU	
Cores and Threads		Up to 28 Cores and 56 Threads	
Last-level Cache		Up to 38.5 MB (non-inclusive)	
UPI Speed (GT/s)		Up to 3x UPI @ 10.4 GT/s	
PCIe* 3.0 Lanes		Up to 48 lanes with 12 controllers	
Memory Speed		Up to 6 channels @ up to 2666 MHz	

Cache Hierarchy Changes



From Broadwell to Skylake (server): the move from ring to mesh



UPI required for dual-socket
(Ultra Path Interconnect)

Broadwell

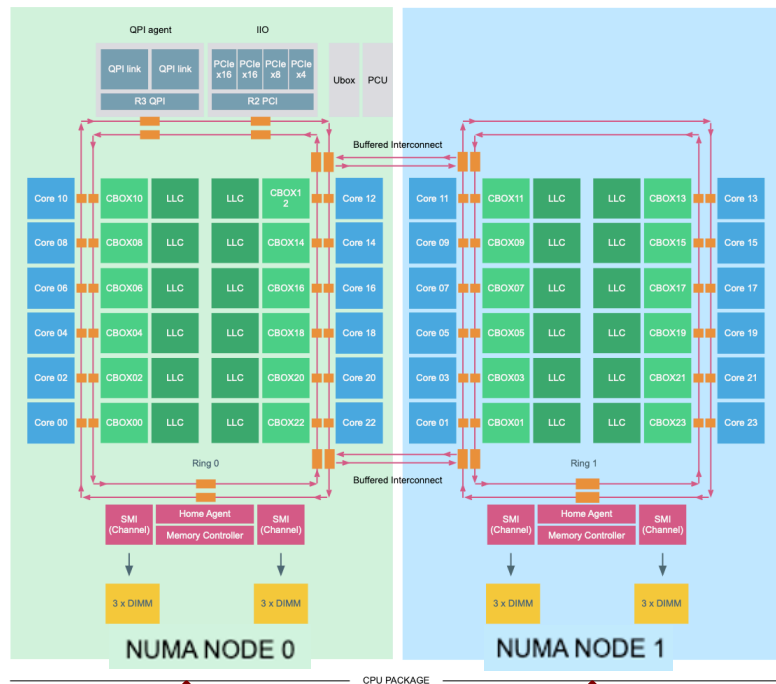
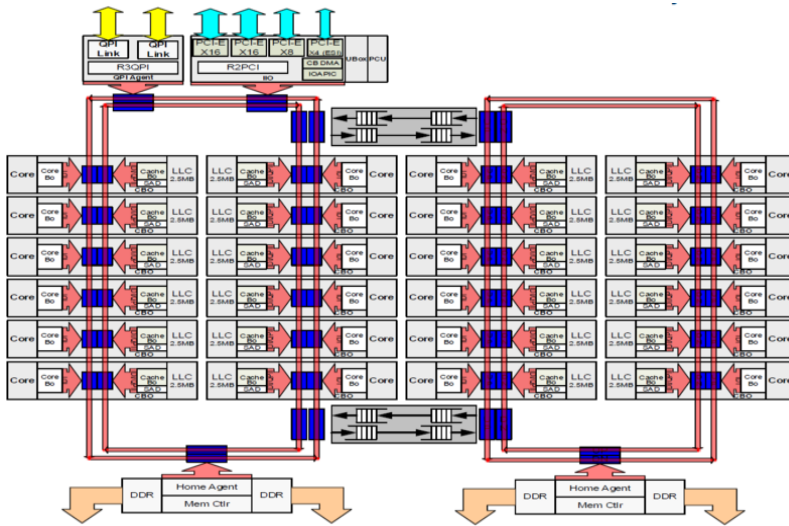
ring interconnection does not scale for large #cores

Skylake (server)

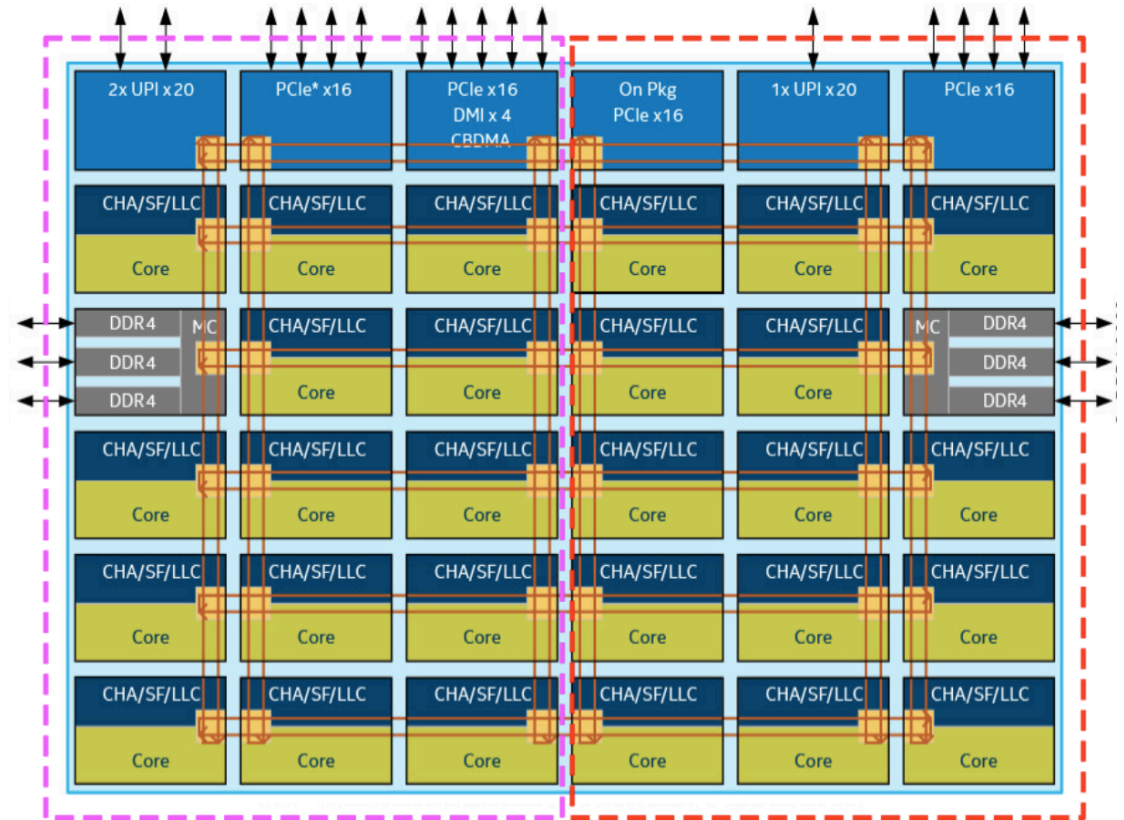
(mesh follows KNL)



From Broadwell to Skylake (server): Sub-NUMA Clusters



↪ Broadwell ↩

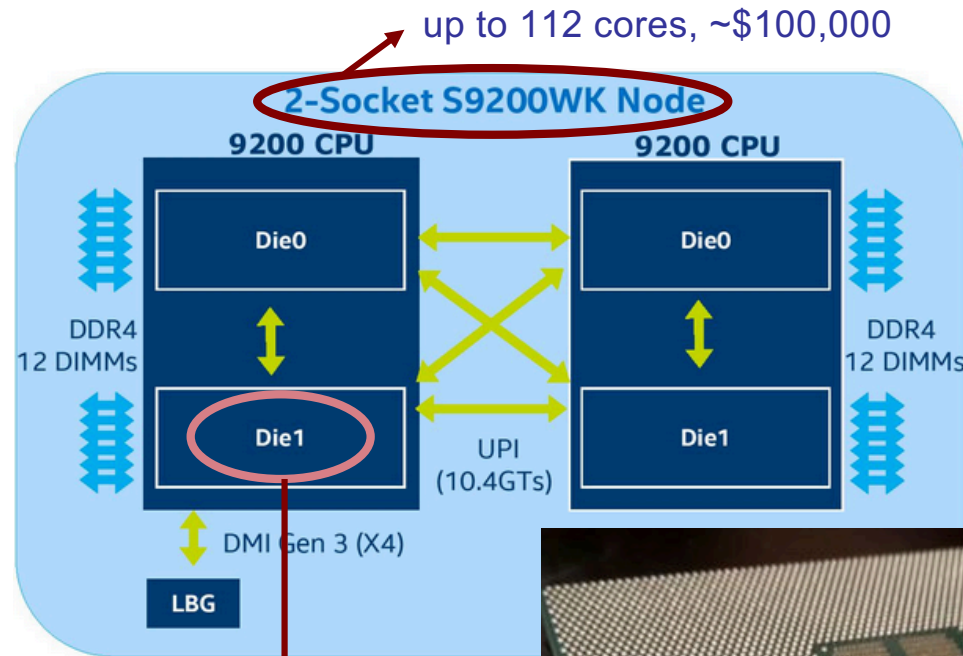
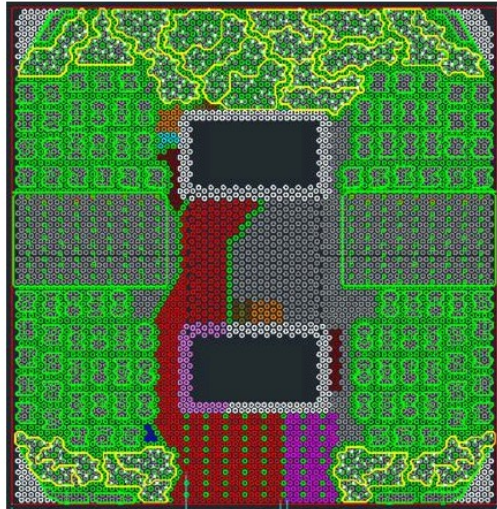


SNC Domain 0 **SNC Domain 1**

↪ Skylake ↩

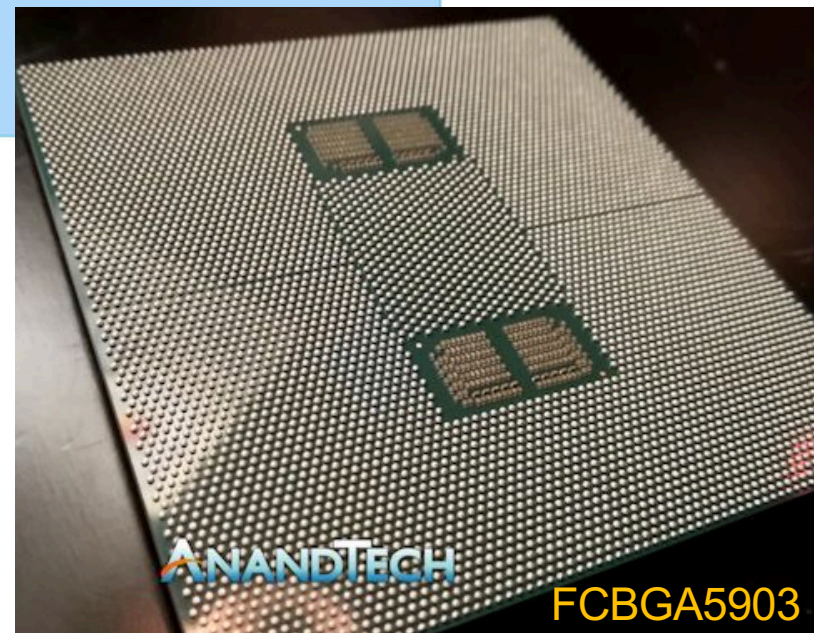
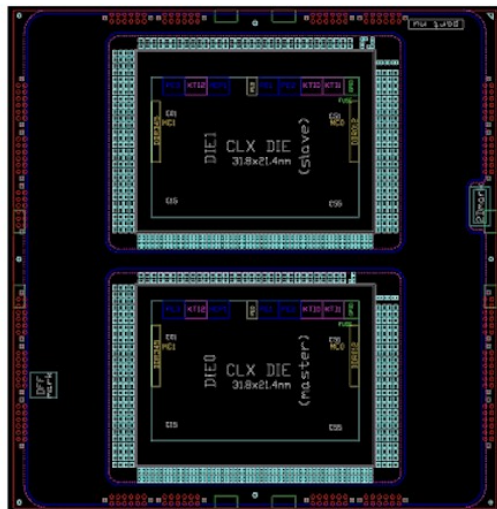


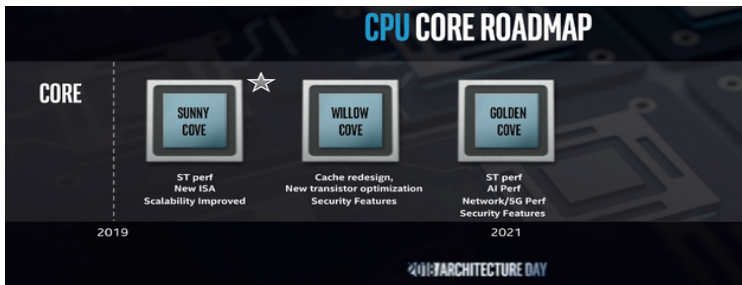
Intel Xeon Scalable Processor: the Advanced Performance (AP) device



8200 die, up to 28 cores

Two 8200 dies on the same package





Sunny Cove Microarchitecture

(larger L2 & L3 in servers)

SUNNYCOVE MICROARCHITECTURE

	HASWELL	SKY LAKE	ICE LAKE
L1 Data Cache	32KB	32KB	48KB
L2 Cache	256KB	256KB	512KB
L2 TLB	1024	1536 16 (1G)	2048 (4k) Shared 1024 for 2M/4M 1024 for 1G
µop Cache	1.5K µops	1.5K µops	2.25K µops
OoO Window	182	224	352
In-Flight Loads	72	72	128
In-Flight Stores	42	56	72

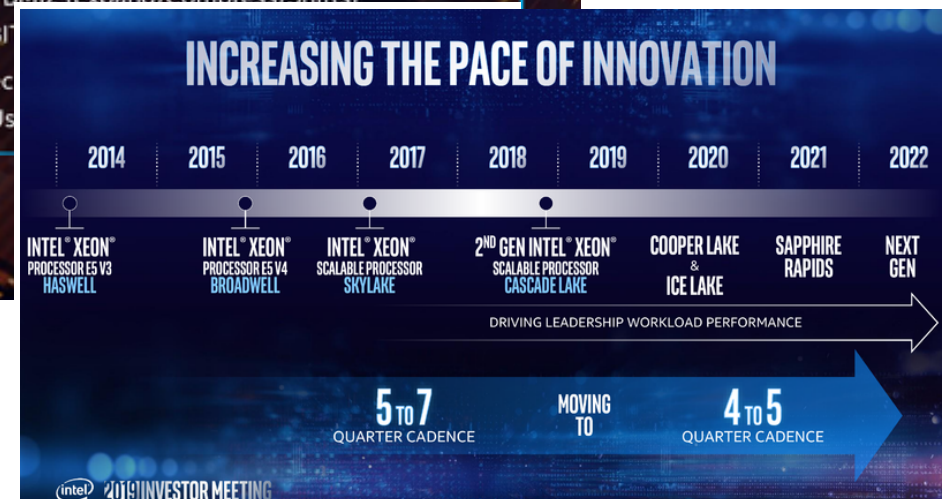
NEW CAPABILITIES

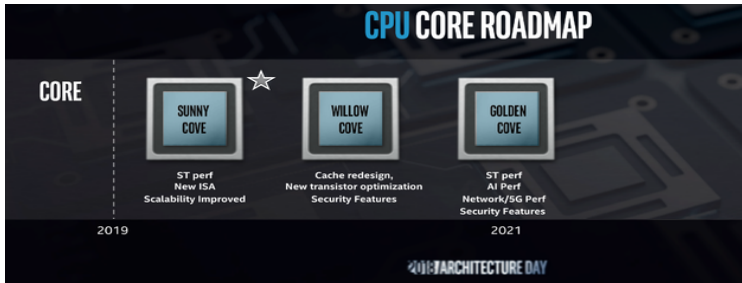
New Instructions for Crypto Performance

- Big Number Arithmetic (IFMA)
- Vector AES
- Vector Carryless Multiply
- Galois Field
- SHA

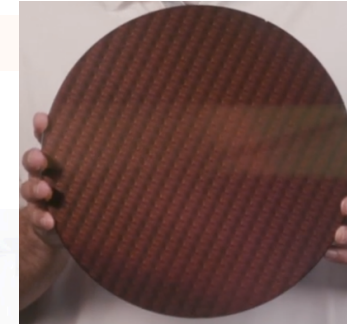
Additional Vector Capabilities

- DLBoost – Inference Acceleration
- VBMI (Permute/Shifts)
- VBMI2 (Expand/Compress/Shifts)





*Next generation:
Willow Cove & Tiger Lake*



Tiger Lake Architecture

**WILLOW COVE
SCALAR
ARCHITECTURE**



**TIGER LAKE
SOC**

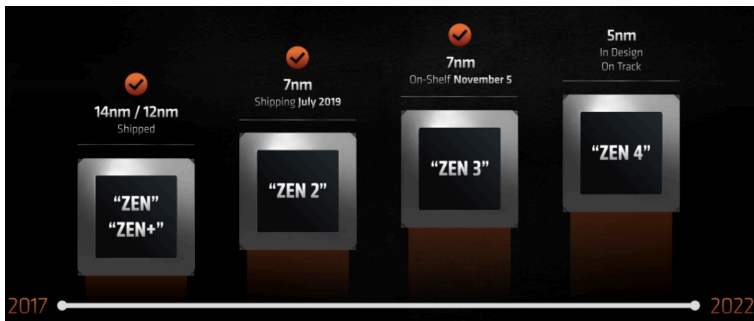
No server version yet...

Manycore chips/packages: an overview

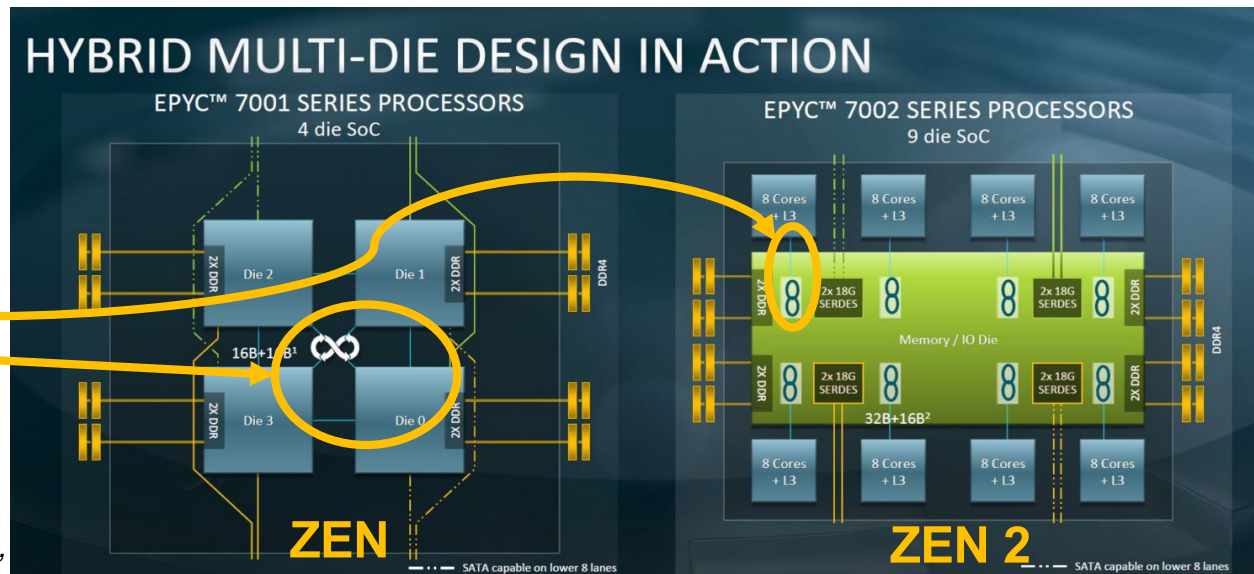
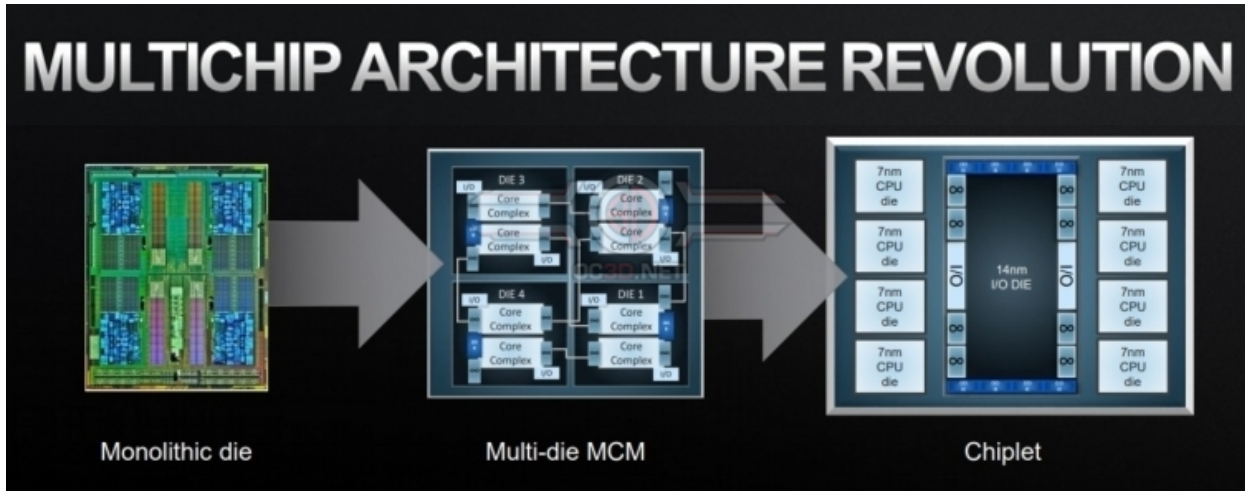


Key server chips/packages that addresses those issues:

- Intel: the Xeon Processor Scalable family
- **AMD: the Epyc Zen family**
- Sunway: the SX260x0 family
- ARM: the ARMv8 server-level competitors
 - Marvell ThunderX family
 - Fujitsu A64FX Arm chip
 - Neoverse N1 hyperscale reference design
 - Ampere Altra Arm Processor
 - Amazon Graviton
 - Huawei HiSilicon Kunpeng 920
- Cerebras: a Wafer Scale Engine
- Apple (*not server*): the SoC approach (*no chiplets!*)



Key Intel Xeon competitor: AMD Epyc (Zen, Zen 2, 3, 4)

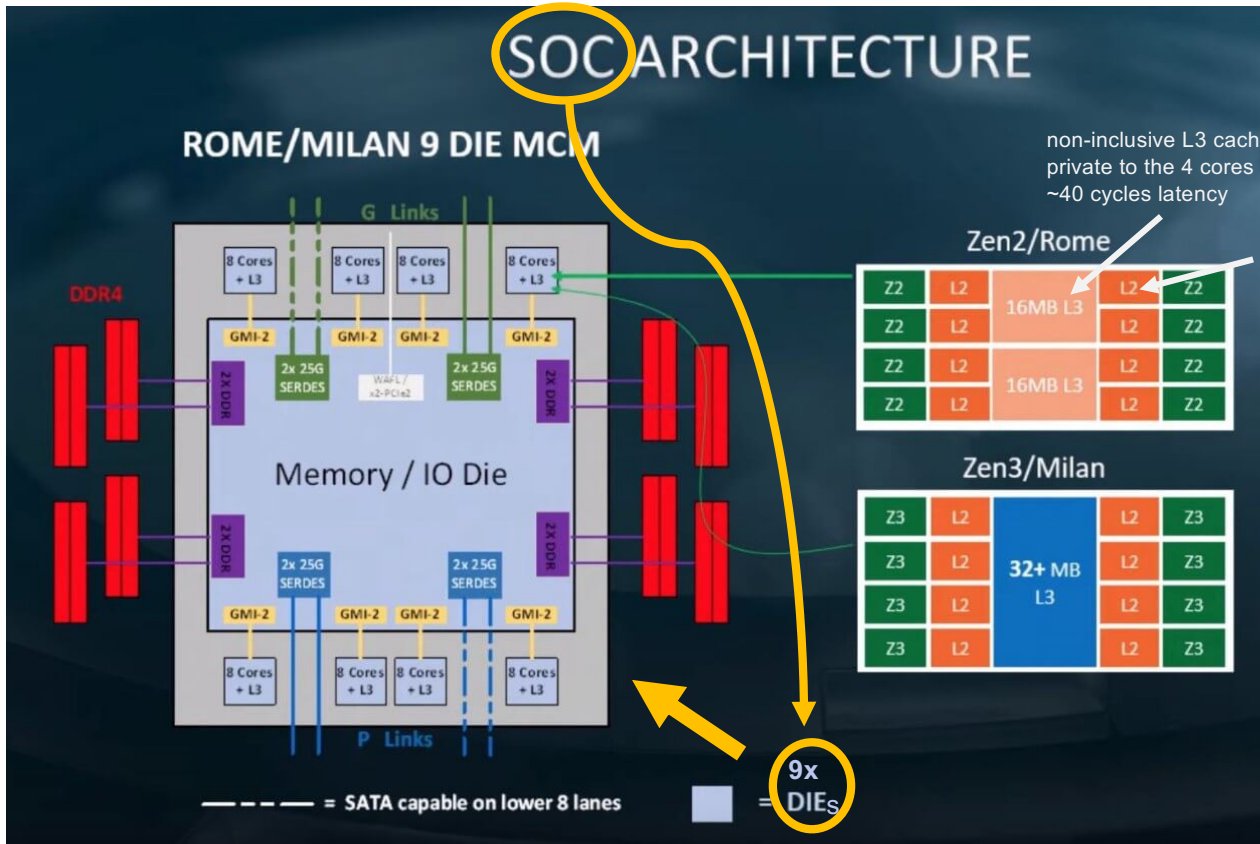


Infinity:
AMD interconnection fabric,
a superset of HyperTransport

AJProença, Advanced Architectures,



AMD Epyc: from Zen 2 (Rome) to Zen 3 (Milan)



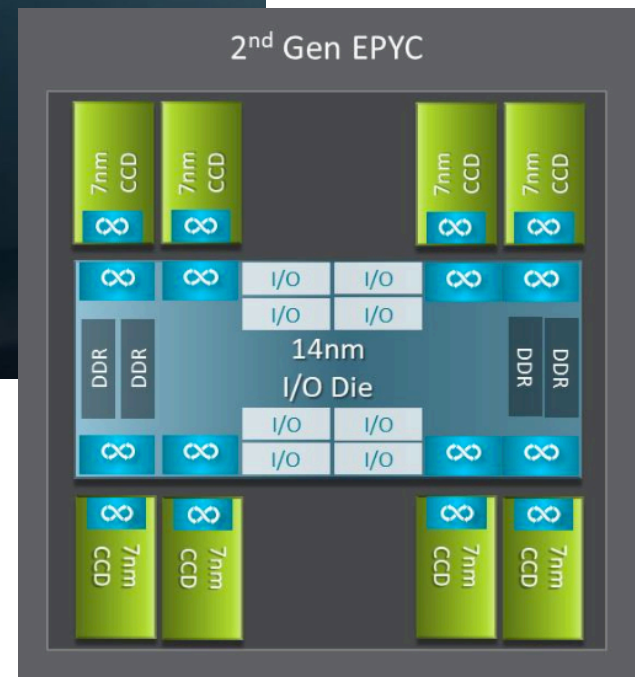
non-inclusive L3 cache private to the 4 cores ~40 cycles latency

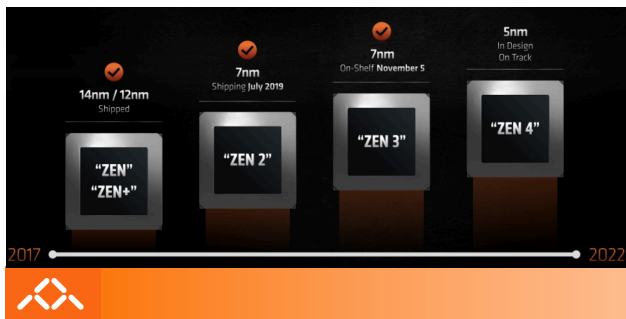


512 KiB L2 cache/core 12-cycles latency



8 CCD dies at 7 nm
1 I/O die at 14 nm

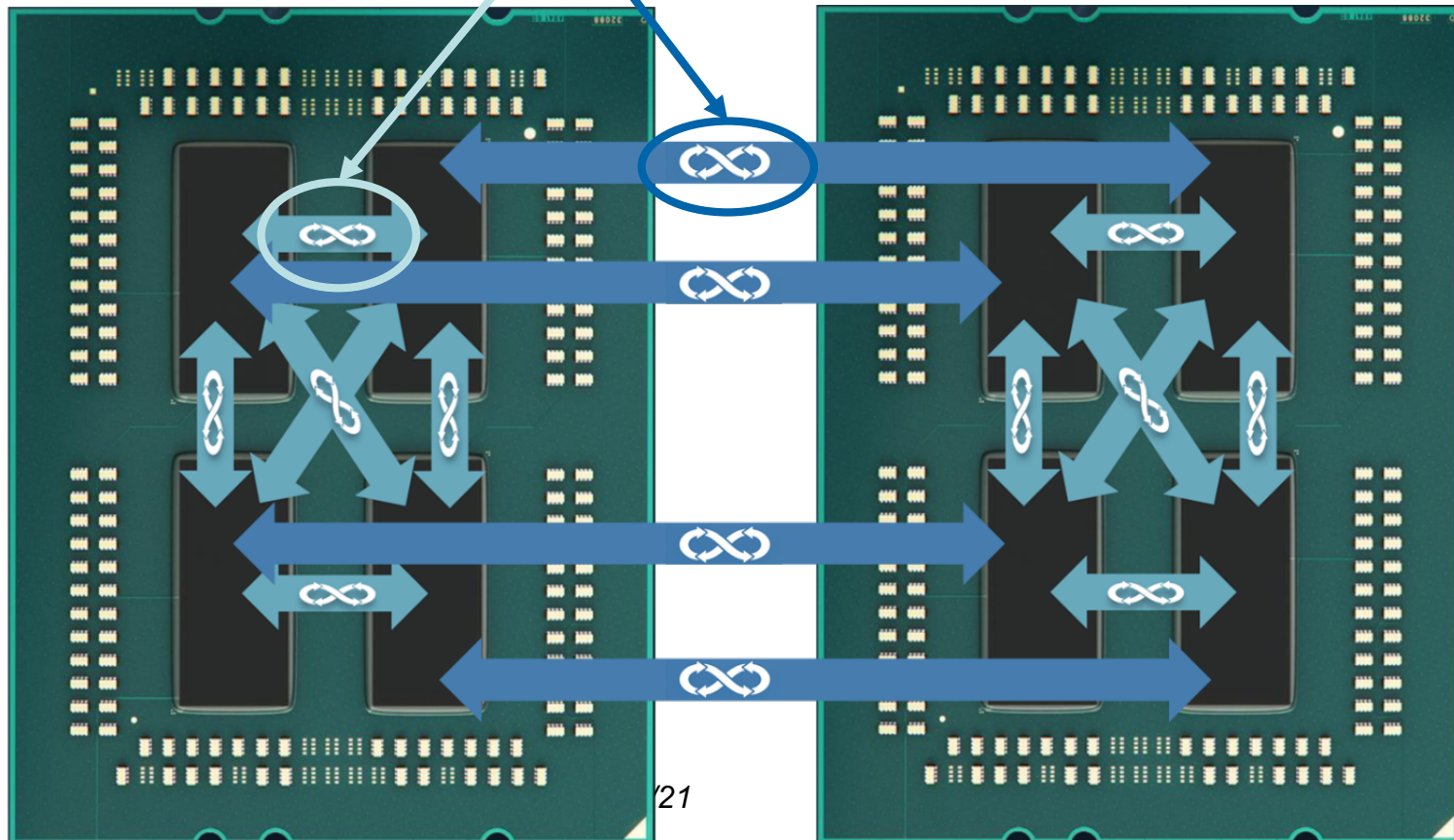


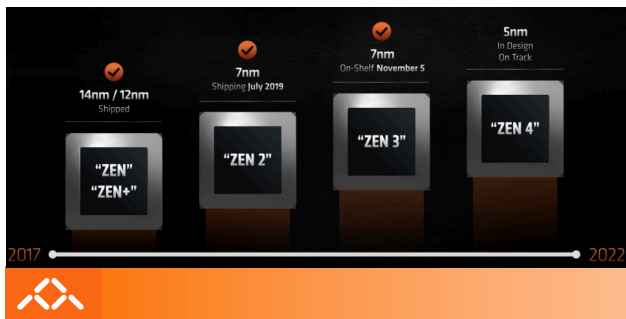


The AMD Infinity Fabric

Infinity Fabric (IF)

- AMD system interconnect architecture, a 256-wide bi-directional crossbar:
 - **Infinity Fabric On-Package (IFOP)**: die-to-die communication in same package
 - **Infinity Fabric InterSocket (IFIS)**: for package-to-package communications

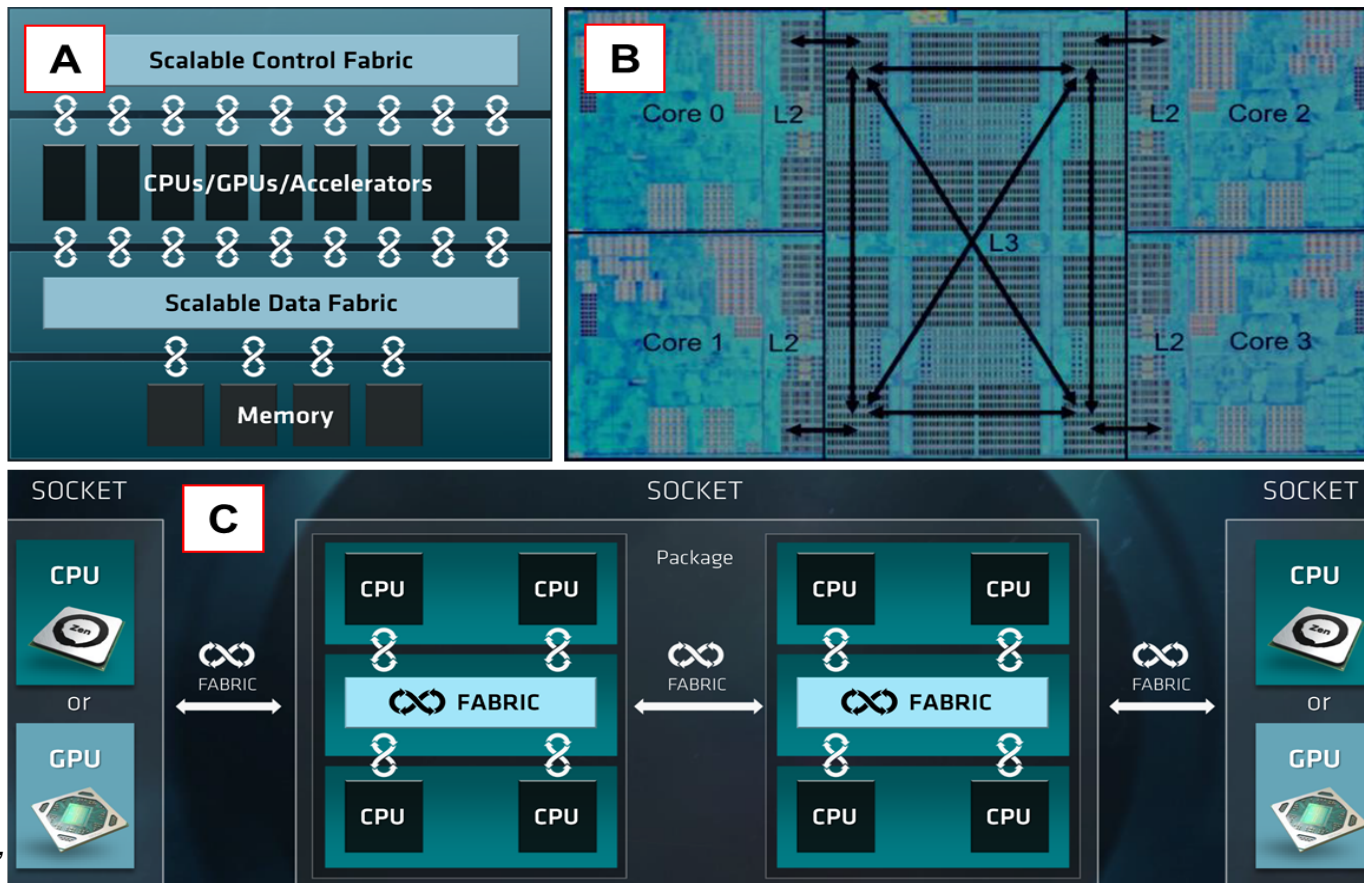




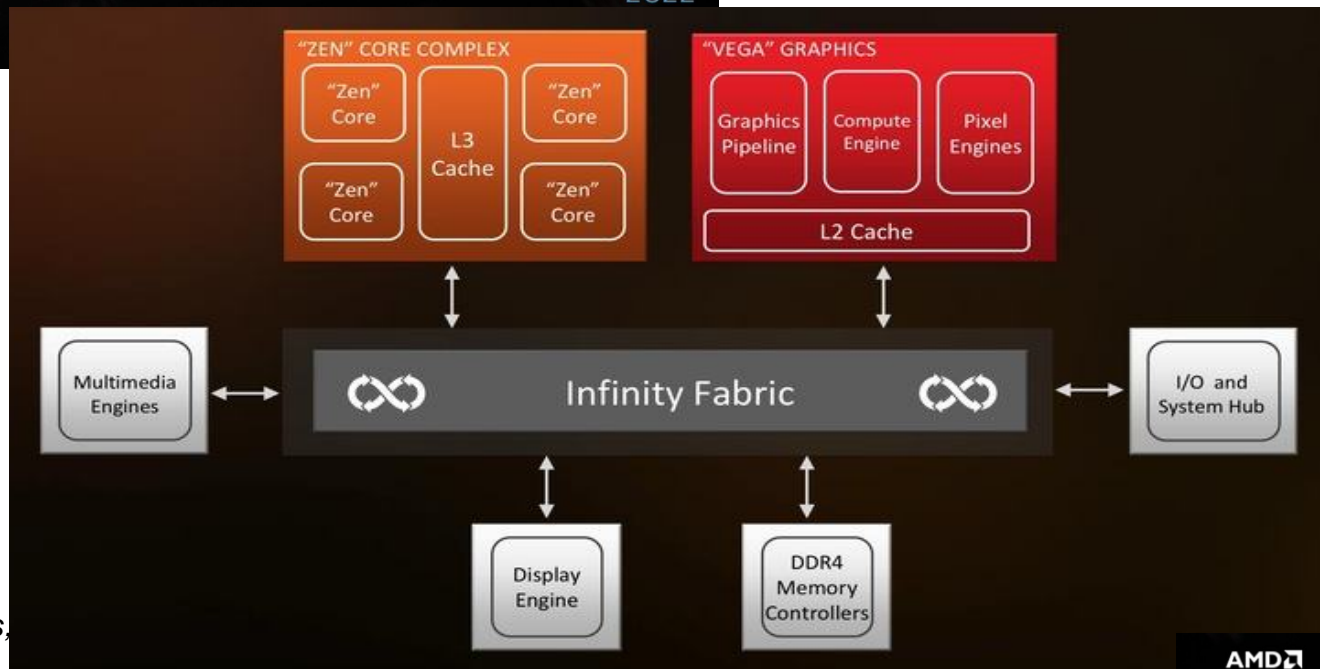
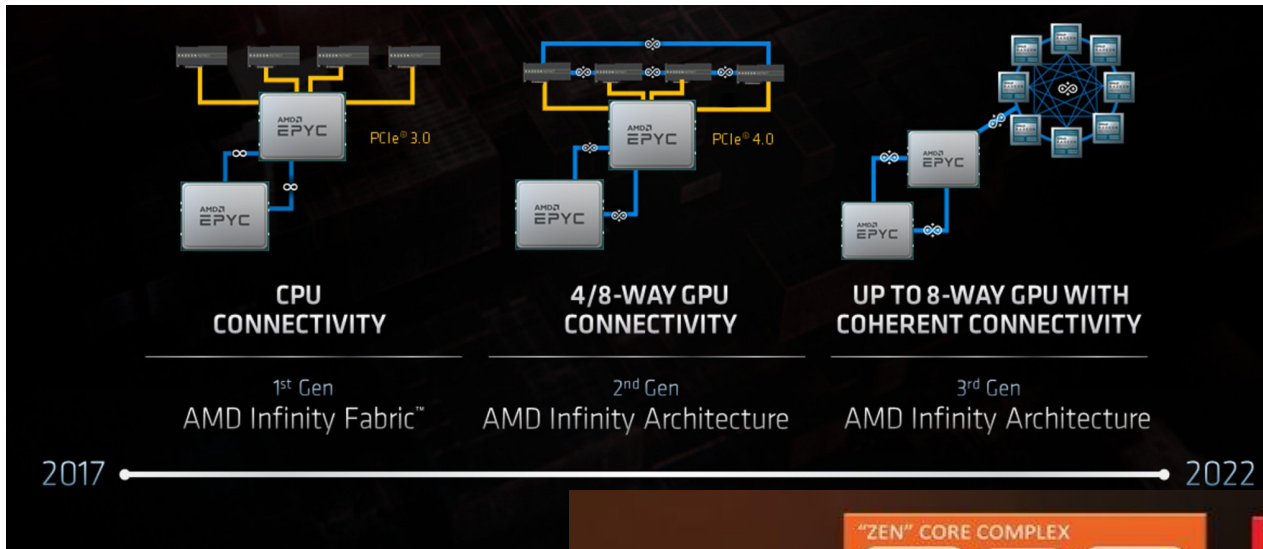
The Infinity Fabric scaling

The **Infinity Fabric** scaling:

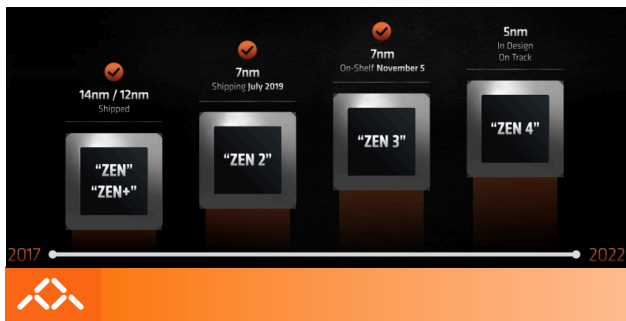
- A. Scalable Control Fabric (SCF)** connects compute elements in CCX (**C**ore **C**omple**X**)
- B. System Data Fabric (SDF)** coherent communications among caches and memory
- C. SCF and SDF** scale between dies on an MCM and between sockets



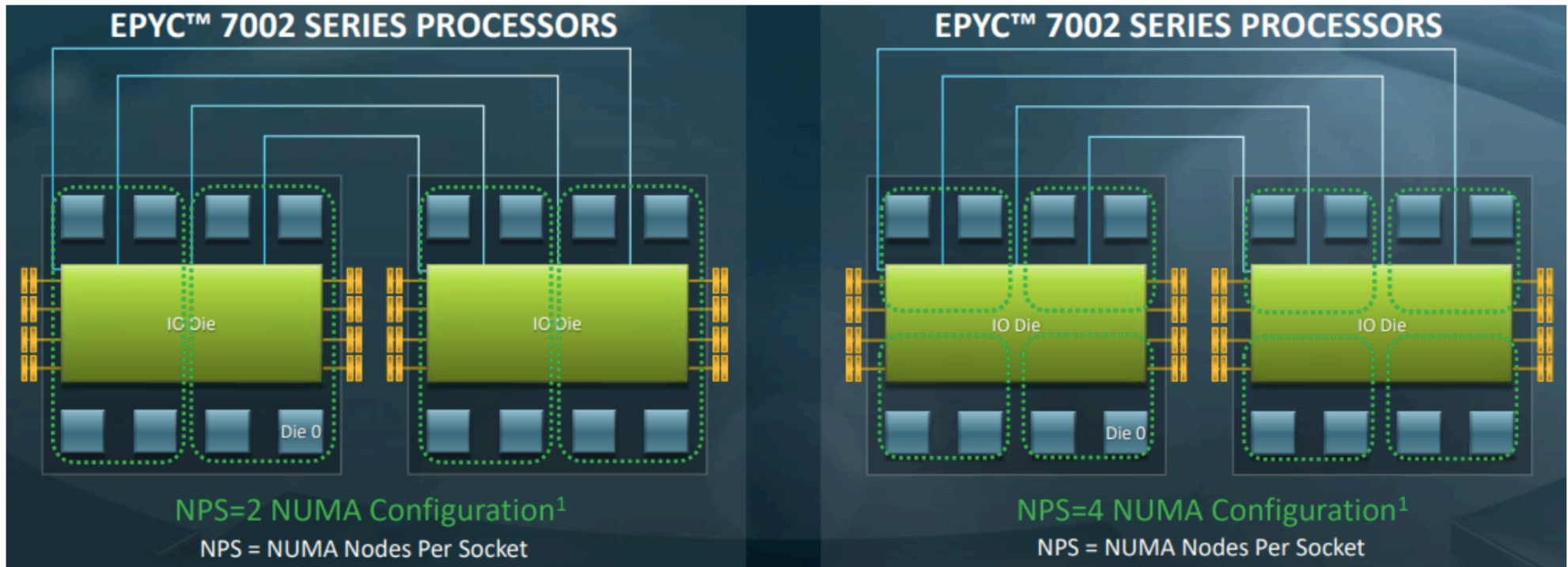
AMD Infinity Architecture Roadmap

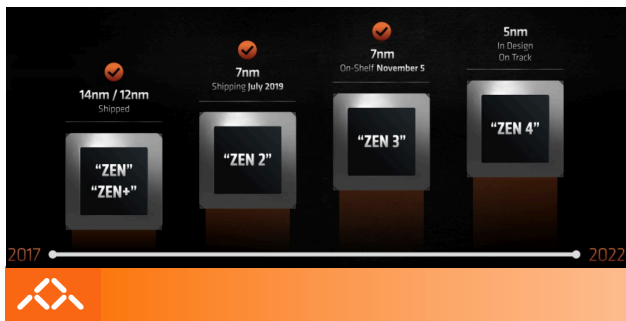


AJProença, Advanced Architectures,



NUMA domains in Epyc MCM





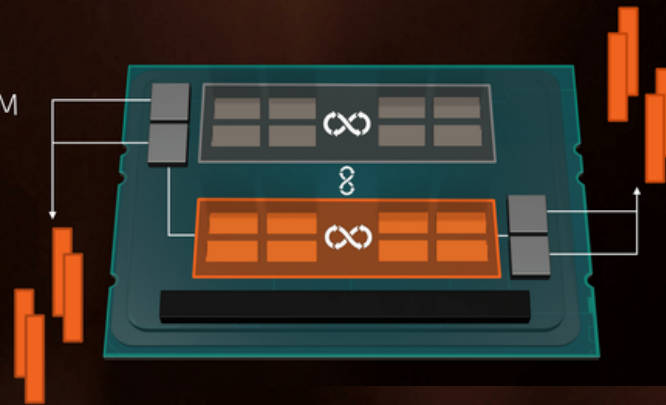
Memory modes in Epyc

CONTROLLING THE MEMORY

Distributed Mode (UMA)

Transactions spread evenly across DRAM

For apps that prefer **WIDE** DRAM access



17 AMD SIGGRAPH17 Tech Day | Confidential - Under Embargo Until 8/10, 9:00am EDT

CONTROLLING THE MEMORY

Local Mode (NUMA)

Transactions in die-local memory



For apps that prefer **FAST** DRAM access

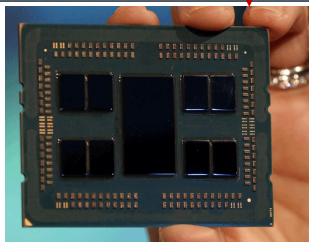


18 AMD SIGGRAPH17 Tech Day | Confidential - Under Embargo Until 8/10, 9:00am EDT

Intel Xeon vs. AMD Epyc

(both launched 1st half 2019)

	AMD EPYC 7742 (Rome)		Intel Xeon Platinum 8280 (Skylake-SP)
	<p>7 nm, I/O 14 nm 9-die package</p> <p>2.25 GHz – 3.4 GHz</p> 		<p>14 nm single-die, 2 SNC</p> <p>2.7 GHz – 4.0 GHz AVX512: 1.8 – 3.7 GHz</p> 
Cores / Threads	64c / 128t	2-way SMT	28c / 56t
L2/L3 Cache	512 KiB/core / 256 MB		1 MiB/core / 38.5 MB
Max Memory / Bandwidth	4 TB / 190.7 GiB/s		1 TB / 131.13 GiB/s
Memory Channels	8		6
PCI-E Lanes	128x PCI-E 4.0		48x PCI-E 3.0



Skylake-AP
 14 nm
 2-die 82xx, 56 cores
 12 memory channels

Intel Xeon vs. AMD Epyc

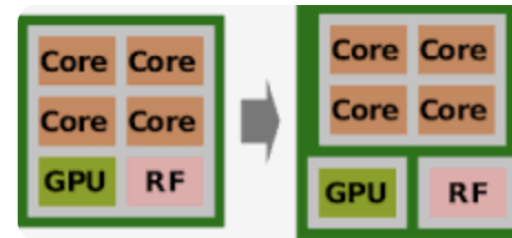
		Microarchitecture Comparison					
Mobile Architecture		Skylake	Cannon Lake	Sunny Cove*	Zen	Zen 2	
	L1-D Cache	32 KiB/core 8-way	32 KiB/core 8-way	48 KiB/core 12-way	32 KiB/core 8-way	32 KiB/core 8-way	
	L1-I Cache	32 KiB/core 8-way	32 KiB/core 8-way	32 KiB/core 8-way	64 KiB/core 4-way	32 KiB/core 8-way	
Skylake-(server):	L2 Cache	256 KiB/core 4-way	256 KiB/core 4-way	512 KiB/core 8-way	512 KiB/core 8-way	512 KiB/core 8-way	
1 MiB/core 16-way	L3 Cache	2 MiB/core 16-way	2 MiB/core 16-way	2 MiB/core 16-way	2 MiB/core	4 MiB/core	
1.375 MiB /core 11-way	L3 Cache Type	Inclusive Non-Inclusive	Inclusive	Inclusive	Non-Inclusive	Non-Inclusive	
	Decode	4 + 1	4 + 1	4 + 1	4	4	
	uOP Cache	1.5k	1.5k	2.25k	2k	4k	
	Reorder Buffer	224	224	352	192	224	
	Execution Ports	8	8	10	10	11	
	AGUs	2 + 1	2 + 1	2 + 2	1 + 1	2 + 1	
	AVX-512	-	1 x FMA	1 x FMA	-	-	

* Sunny Cove numbers for Client. Server will have different L2/L3 cache and FMA, like Skylake

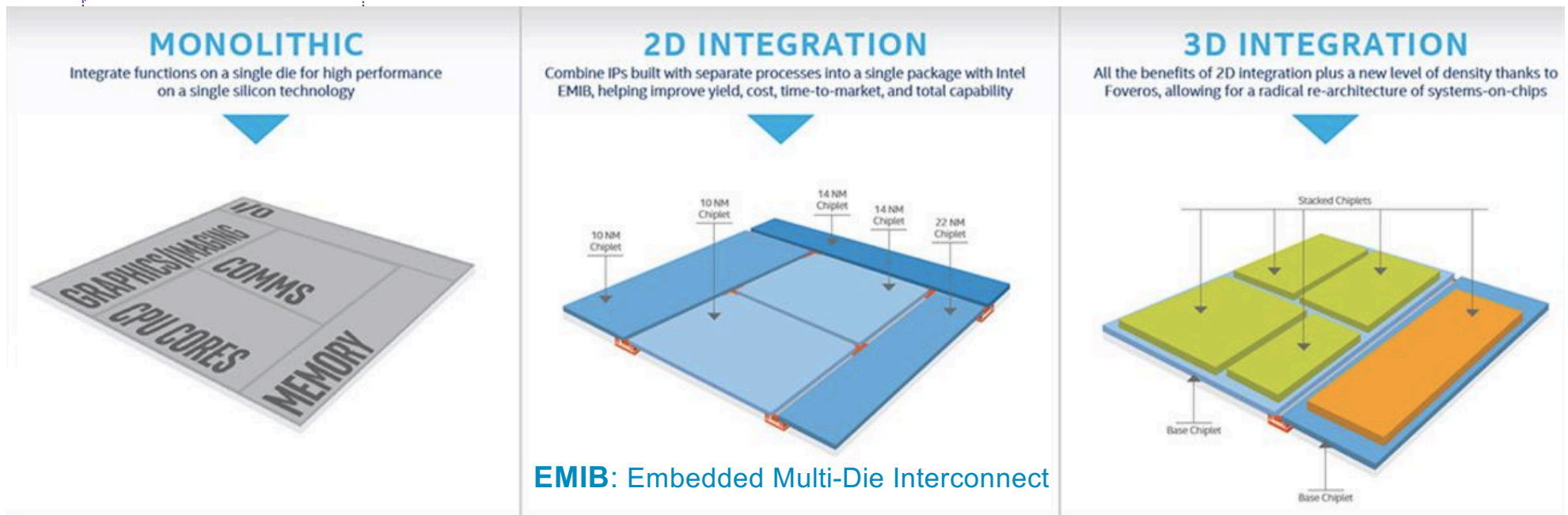
The move at Intel from single-die SoC to stacked chipllets



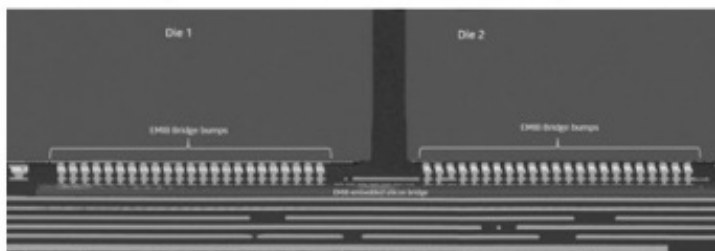
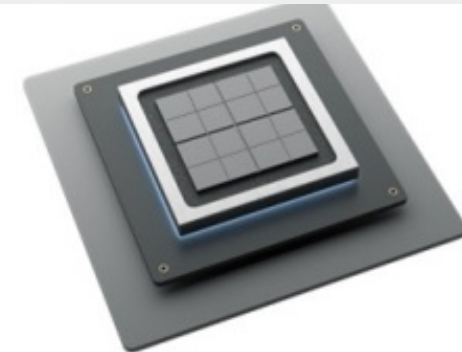
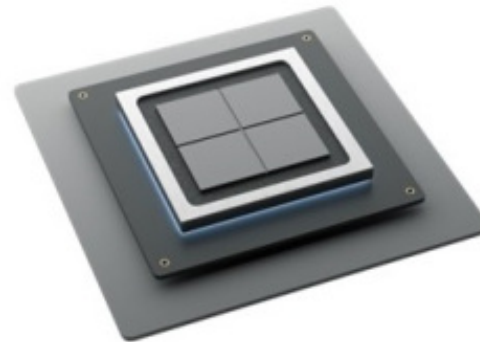
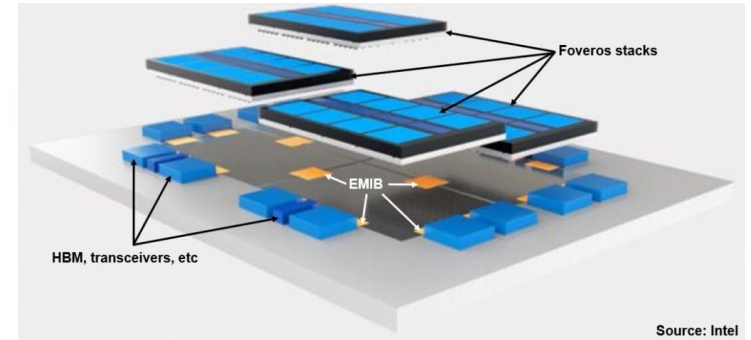
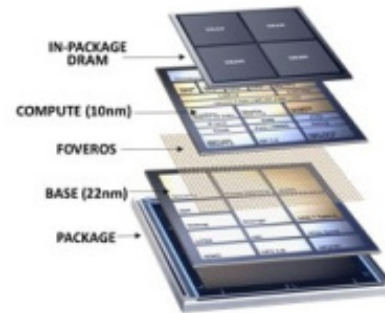
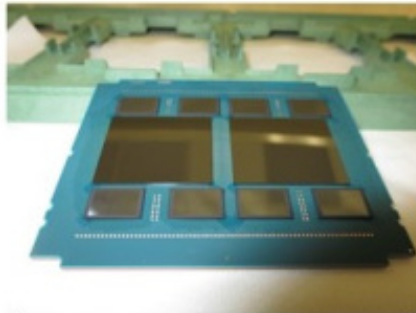
A **chipllet** is an integrated circuit block that has been specifically designed to work with other similar **chipllets** to form larger more complex chips. In such chips, a system is subdivided into functional circuit blocks, called "**chipllets**", that are often made of reusable IP blocks. Mar 27, 2020



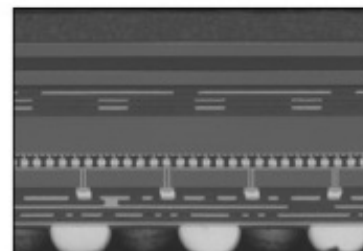
[en.wikichip.org](https://en.wikichip.org/wiki/chipllet) > wiki > chipllet



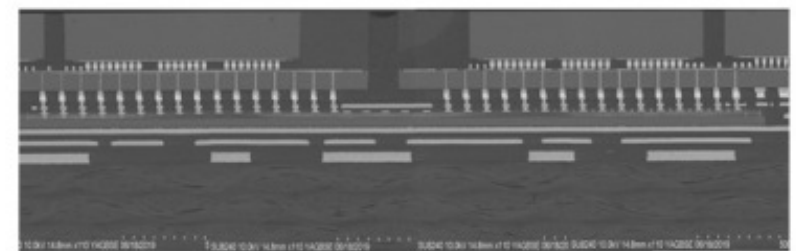
Advanced package architectures at Intel



EMIB

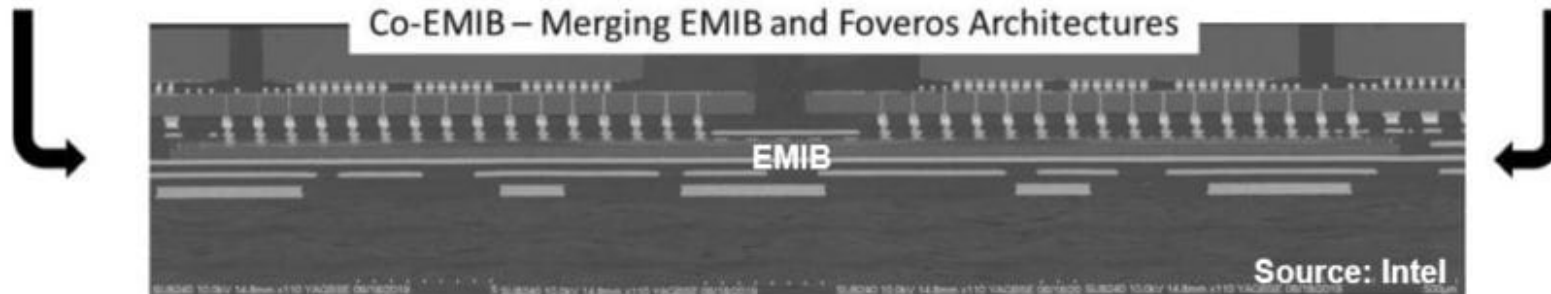
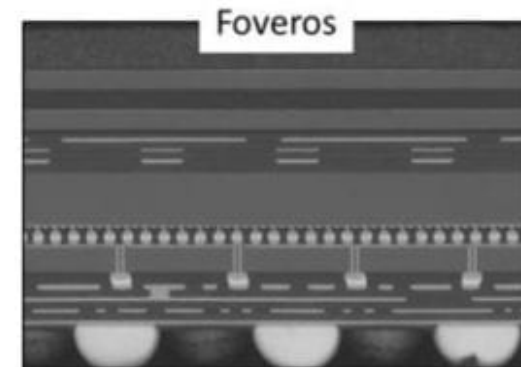
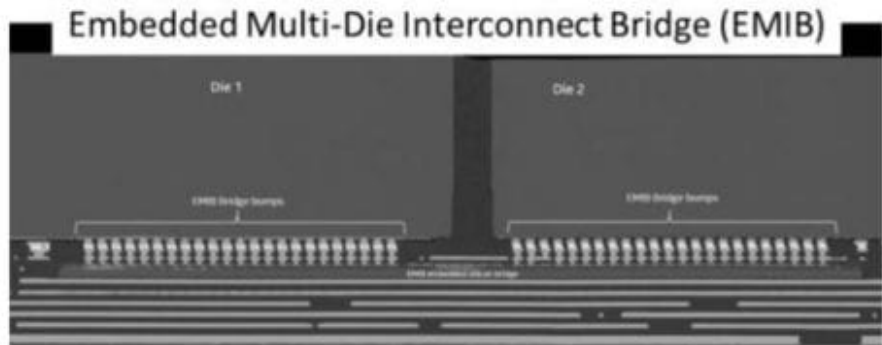
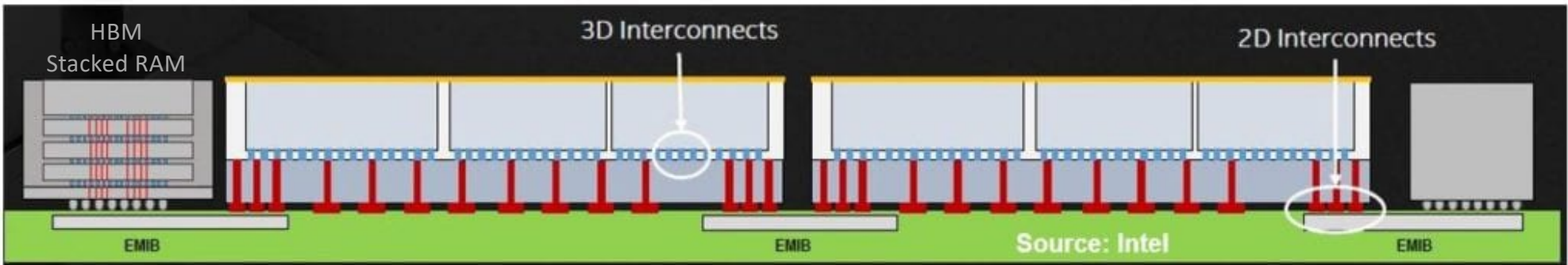
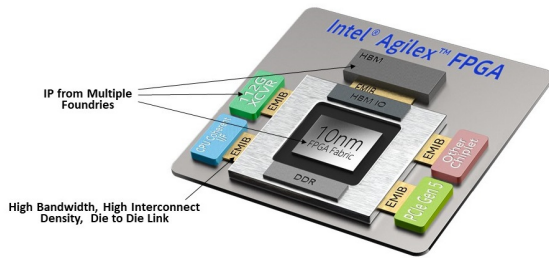


Foveros



Co-EMIB

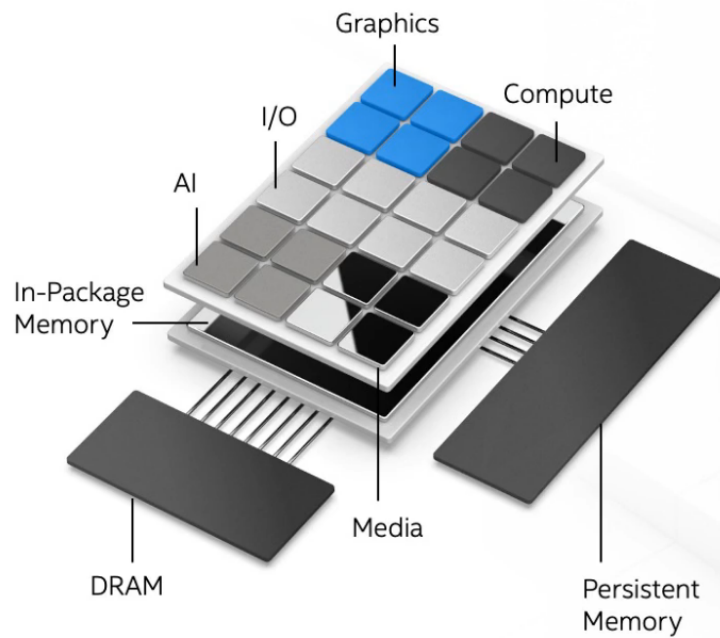
Advanced package architectures at Intel



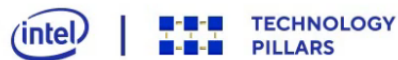
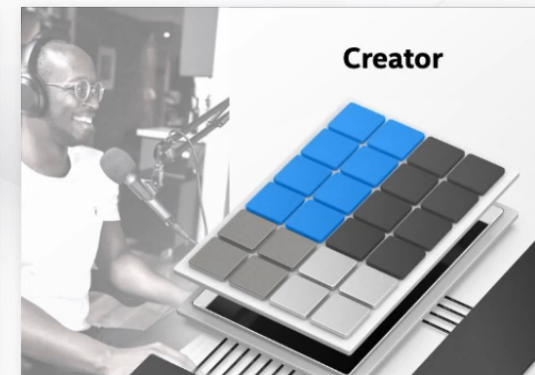
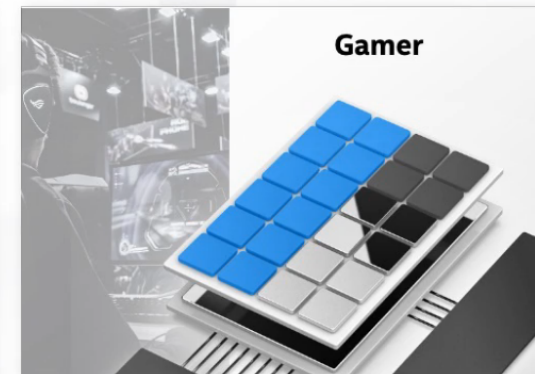
A long term vision at Intel



Purpose Built Client



Long Term Vision



Under embargo until August 13th, 2020 at 6:00 a.m. Pacific Time.

Architecture Day **2020**

Manycore chips/packages: an overview



Key server chips/packages that addresses those issues:

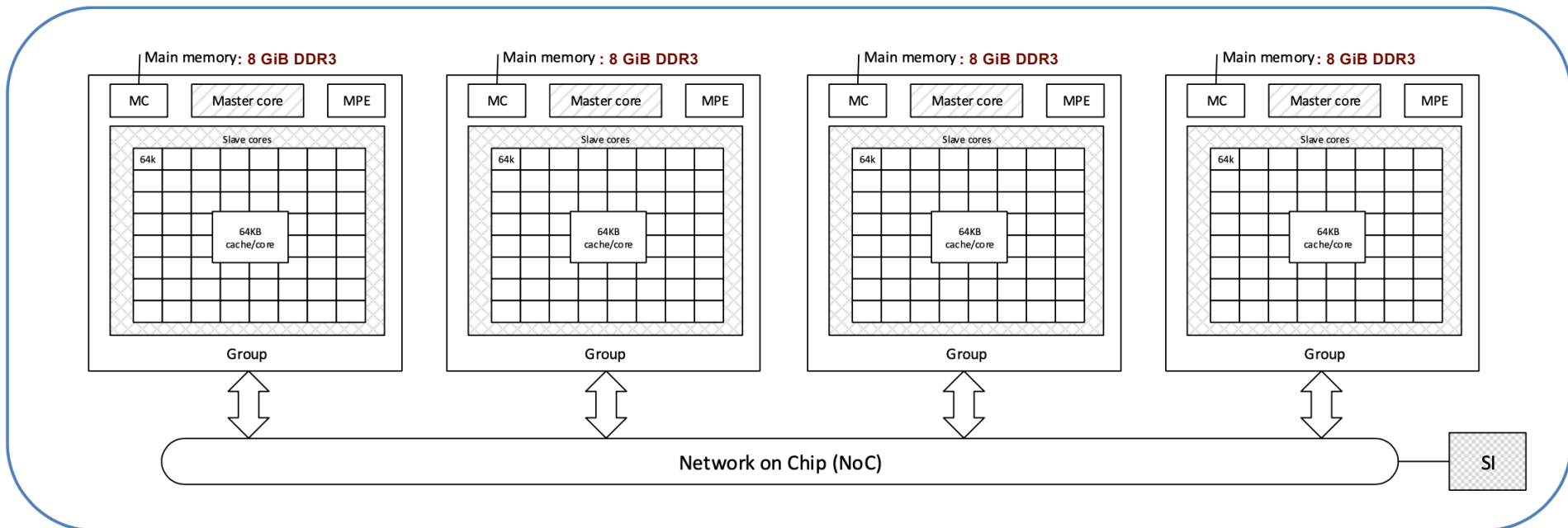
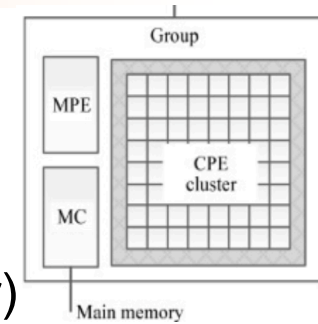
- Intel: the Xeon Processor Scalable family
- AMD: the Epyc Zen family
- **Sunway: the SX260x0 family**
- ARM: the ARMv8 server-level competitors
 - Marvell ThunderX family
 - Fujitsu A64FX Arm chip
 - Neoverse N1 hyperscale reference design
 - Ampere Altra Arm Processor
 - Amazon Graviton
 - Huawei HiSilicon Kunpeng 920
- Cerebras: a Wafer Scale Engine
- Apple (*not server*): the SoC approach (*no chiplets!*)



The Shen Wei SW 26010 in SunWay TaihuLight (#1 in June'16 TOP500)

Shen Wei SW 26010 (260 cores):

- 4 core groups (CG, as a SNC), connected via a NoC
- each CG has a Management Processing Element (MPE) and a Memory Controller (MC) and a **8x8 mesh** of 64 Computing Processing Elements (CPE Cluster)
- a CPE is a 64-bit RISC OoO (*out-of-order*) core w/ a **256-bit vector unit**, no SMT, 16 KiB L1 instruction cache, and 64 KiB Scratch Pad Memory (**not L2 cache**)



Manycore chips/packages: an overview



Key server chips/packages that addresses those issues:

- Intel: the Xeon Processor Scalable family
- AMD: the Epyc Zen family
- Sunway: the SX260x0 family
- **ARM: the ARMv8 server-level competitors**
 - Marvell ThunderX family
 - Fujitsu A64FX Arm chip
 - Neoverse N1 hyperscale reference design
 - Ampere Altra Arm Processor
 - Amazon Graviton
 - Huawei HiSilicon Kunpeng 920
- Cerebras: a Wafer Scale Engine
- Apple (*not server*): the SoC approach (*no chiplets!*)

Support for
at least dual-socket



WIKIPEDIA
The Free Encyclopedia



ARM brand: a bit of history...

ARM architecture

Current owner of Arm Holdings: NVidia

From Wikipedia, the free encyclopedia

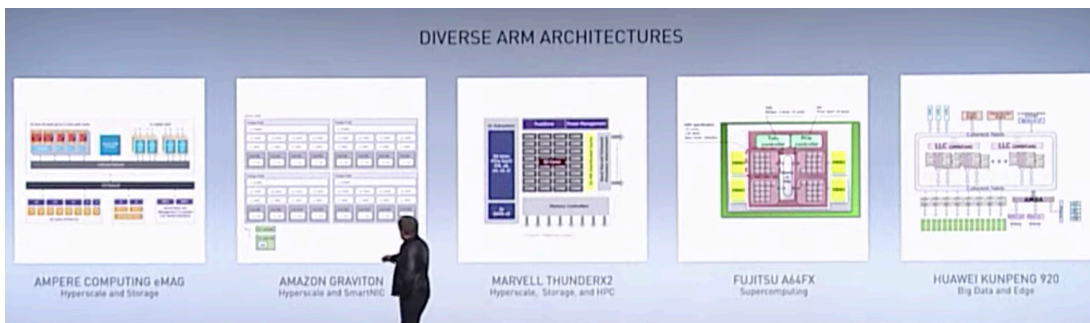
ARM, previously **Advanced RISC Machine**, originally **Acorn RISC Machine**, is a family of [reduced instruction set computing](#) (RISC) architectures for computer processors, configured for various environments. [Arm Holdings](#) develops the architecture and licenses it to other companies, who design their own products that implement one of those architectures—including [systems-on-chips](#) (SoC) and [systems-on-modules](#) (SoM) that incorporate memory, interfaces, radios, etc. It also designs [cores](#) that implement this [instruction set](#) and licenses these designs to a number of companies that incorporate those core designs into their own products.

Processors that have a RISC architecture typically require fewer [transistors](#) than those with a [complex instruction set computing](#) (CISC) architecture (such as the [x86](#) processors found in most [personal computers](#)), which [improve cost, power consumption, and heat dissipation](#). These

ARM architectures

The ARM logo

Designer	Arm Holdings
Bits	32-bit, 64-bit
Introduced	1985; 34 years ago
Design	RISC
Type	Register-Register
Branching	Condition code, compare and branch
Open	Proprietary



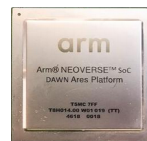
HPCs with ARMv8: server-level competitors



1. Marvell ThunderX product family



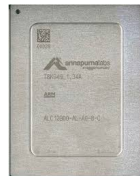
2. Fujitsu A64FX Arm chip



3. Neoverse N1 hyperscale reference design



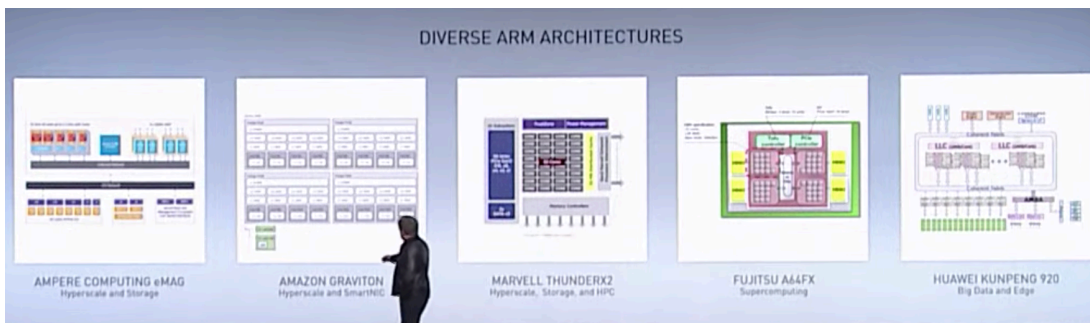
4. Ampere Altra Arm Processor



5. Amazon Graviton



6. Huawei HiSilicon Kunpeng 920



HPCs with ARMv8: server-level competitors



1. Marvell ThunderX product family



2. Fujitsu A64FX Arm chip



3. Neoverse N1 hyperscale reference design



4. Ampere Altra Arm Processor



5. Amazon Graviton



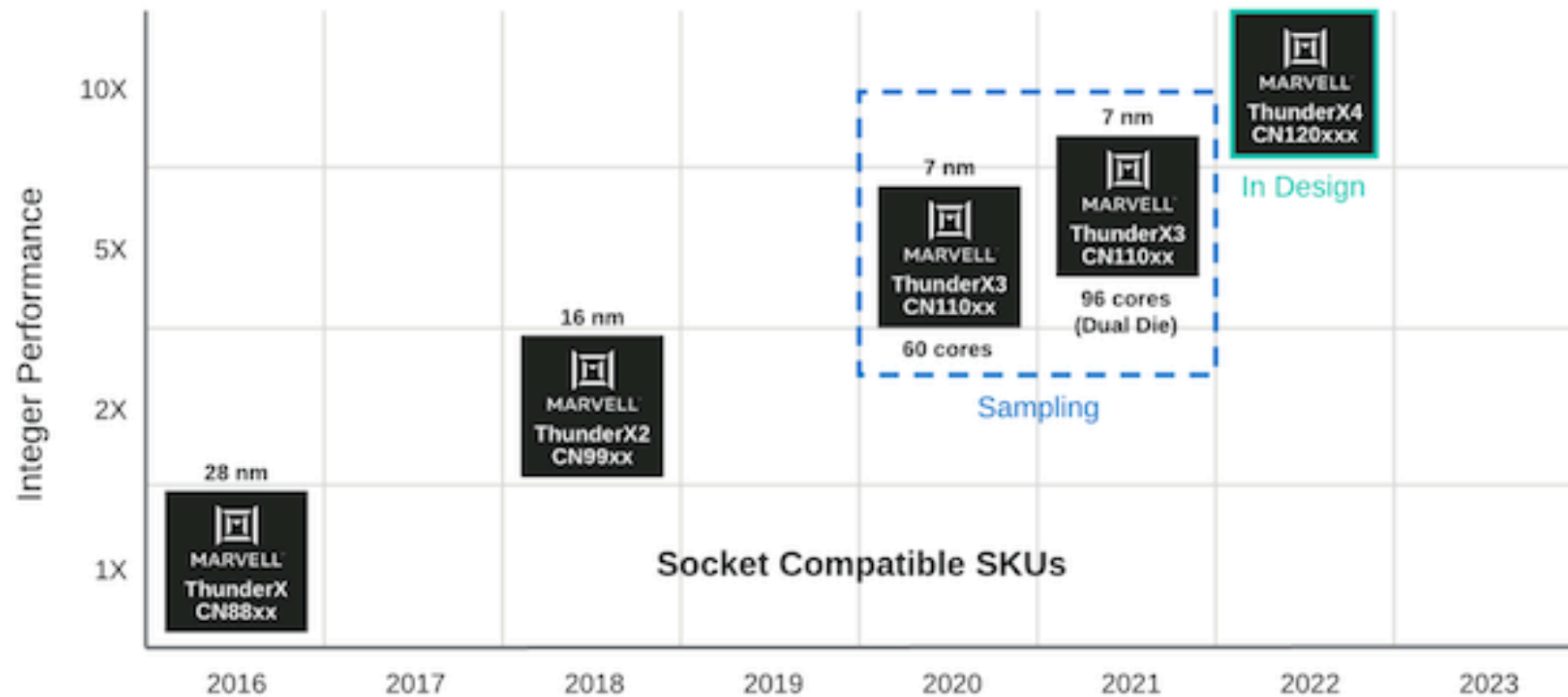
6. Huawei HiSilicon Kunpeng 920



Marvell Server Processor Roadmap



Marvell server processor roadmap



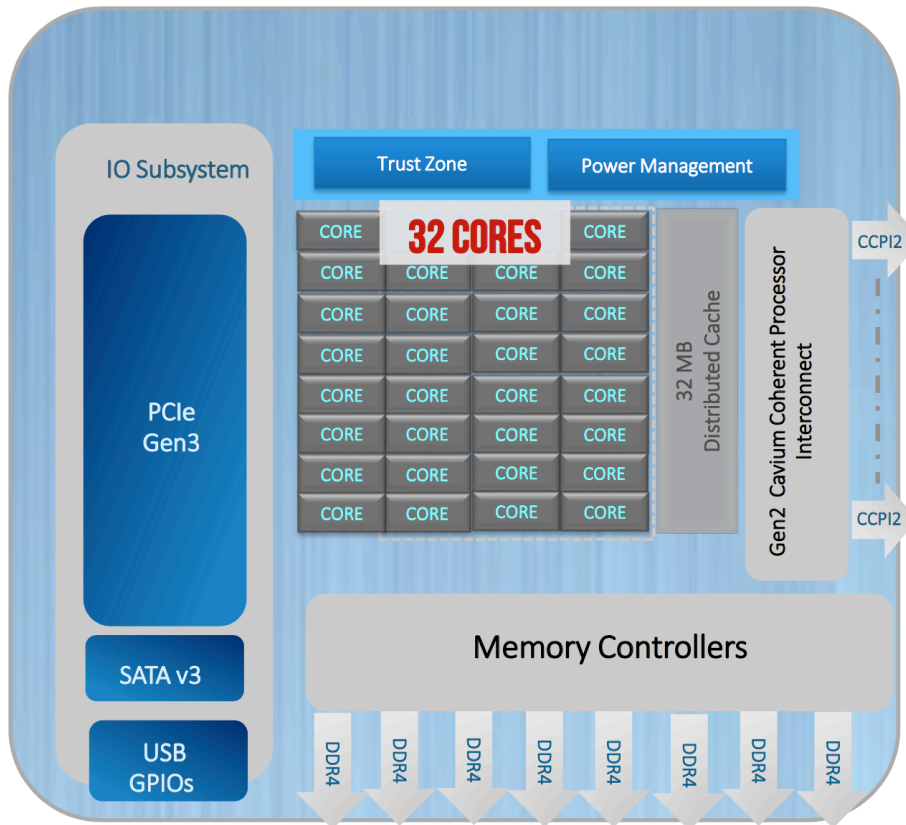
© 2020 Marvell. All rights reserved.

3



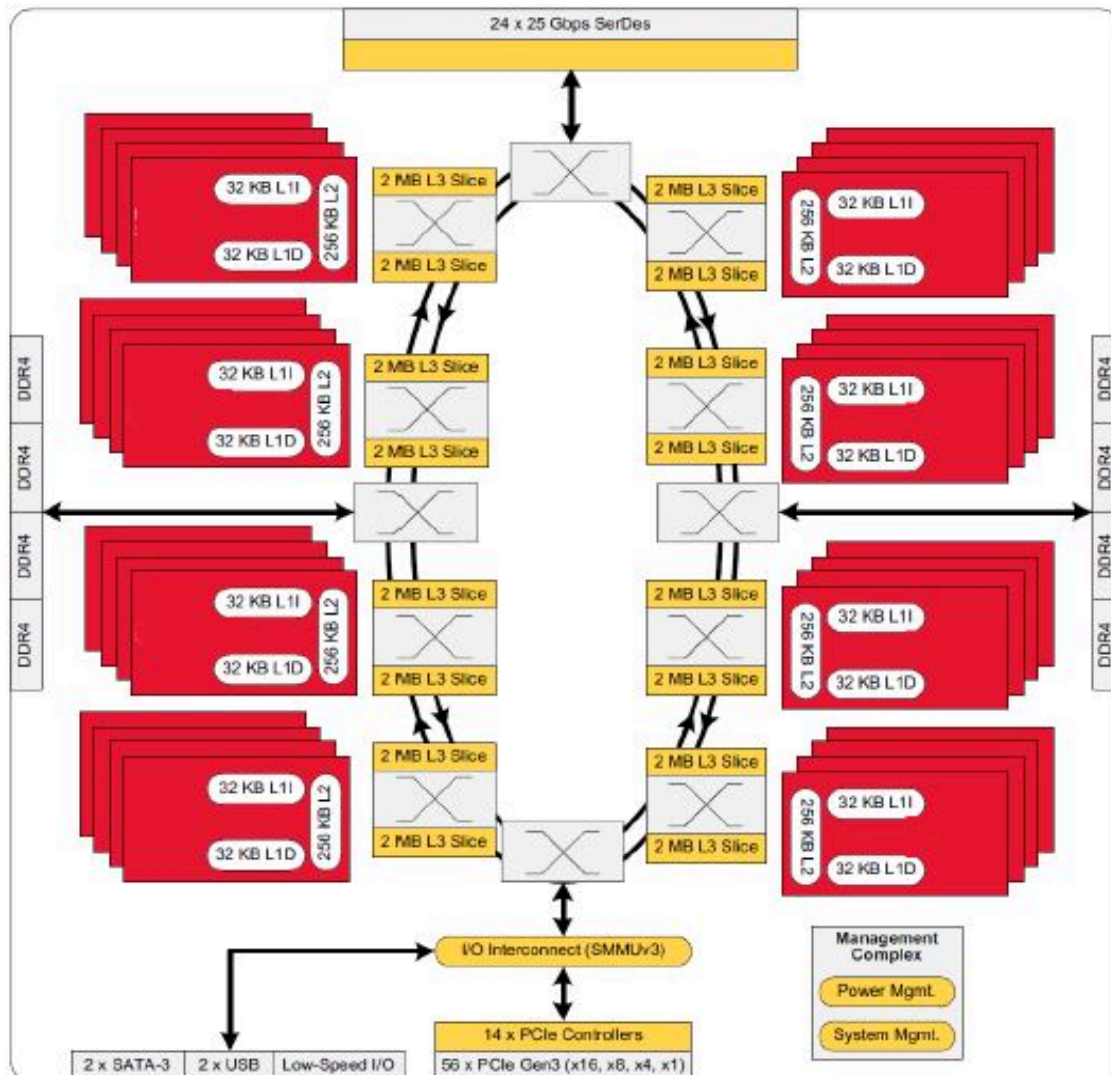
The Marvell/Cavium ThunderX2

THUNDERX2[®] Family Key Features

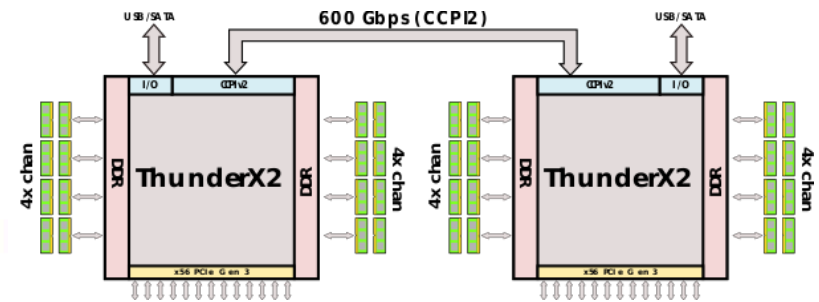


- Up to 32 custom Armv8.1 cores, up to 2.5GHz
- Full OoO, 1, 2, 4 threads per core
- 1S and 2S Configuration
- Up to 8 DDR4-2667 Memory Controllers, 1 & 2 DPC
- Up to 56 lanes of PCIe, 14 PCIe controllers
- Full SoC: Integrated SATAv3 USB3 and GPIOs
- Server class RAS & Virtualization
- Extensive Power Management
- LGA and BGA for most flexibility
- 40+ SKUs
- Volume SKU List Price: \$1795 (180W) to \$800 (75W)

The Marvell/Cavium ThunderX2 architecture block diagram

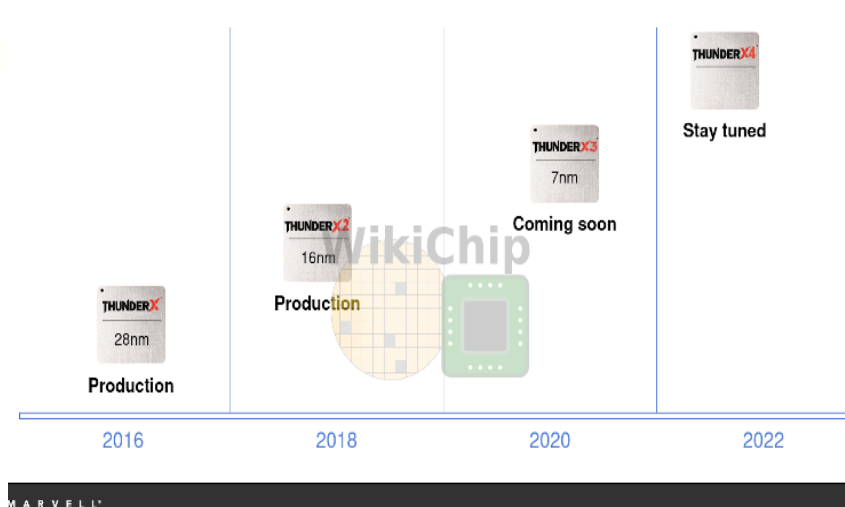


Scalability



ThunderX[®] roadmap

Driving >2X generational performance improvement

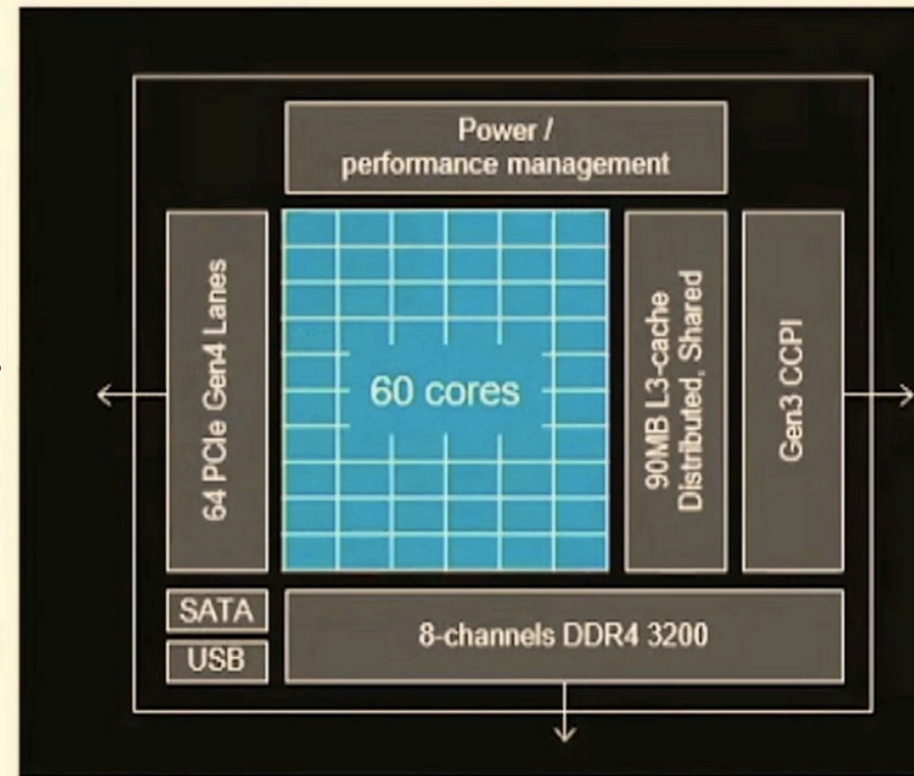




Next generation: ThunderX3

ThunderX3™ overview

- Single die: Up to 60 cores
- Dual die: Up to 96 cores
- Arm v8.3 with select v8.4/v8.5 features
- 30% single thread gain at equal frequency over ThunderX2
- Up to four threads per core
- High bandwidth switched ring interconnect
- Up to 8 DDR4-3200 channels
- Single die: 2X-3X perf over ThunderX2 at equal power
 - Further gains from dual die
- Up to 64 PCIe Gen4, 16 PCIe controllers
- Fine grain power monitoring/management
- TSMC 7nm



© 2020 Marvell. All rights reserved.

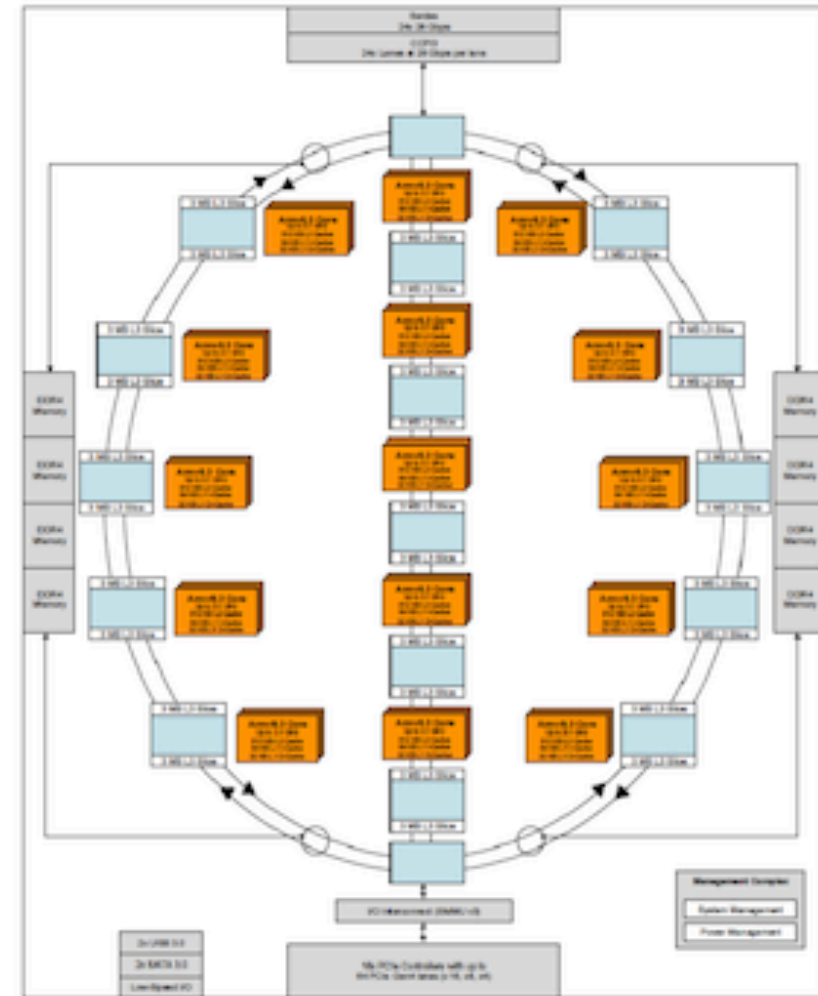
4



ThunderX3: L3 Cache and Interconnect

L3-Cache and interconnect

- Cores / L3-caches organized as switched rings
- DDR channels, I/O tap into rings
- L3-cache organized as tiles that are cache line striped
 - 1 ½ MB per core
 - No notion of L3 cache affinity to cores
 - Good for shared text and shared data
- Exclusive L3-cache – filled on evict from L2-cache
- Snoop based coherence with snoop filters
 - Single socket and two socket



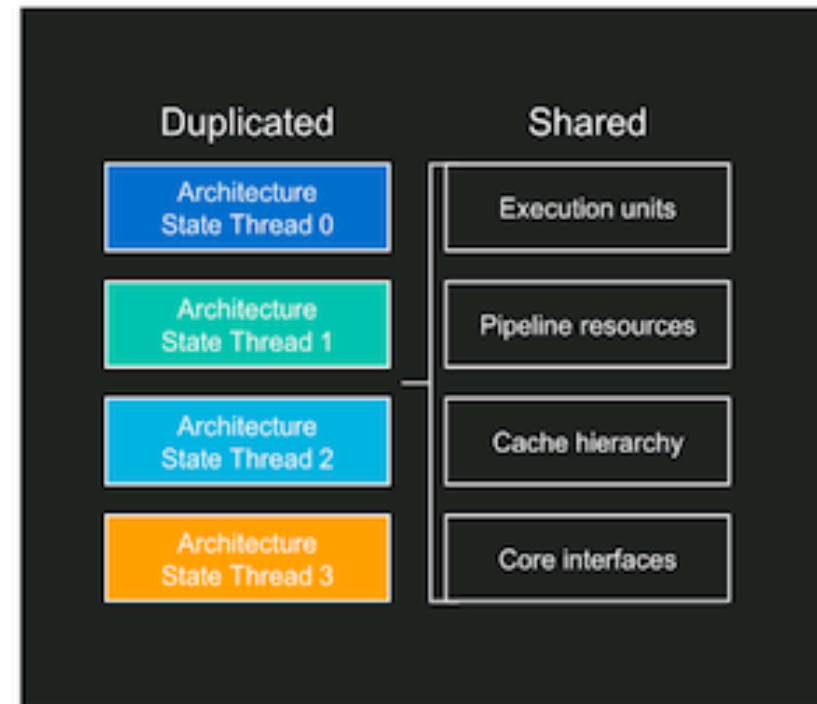
© 2020 Marvell. All rights reserved.

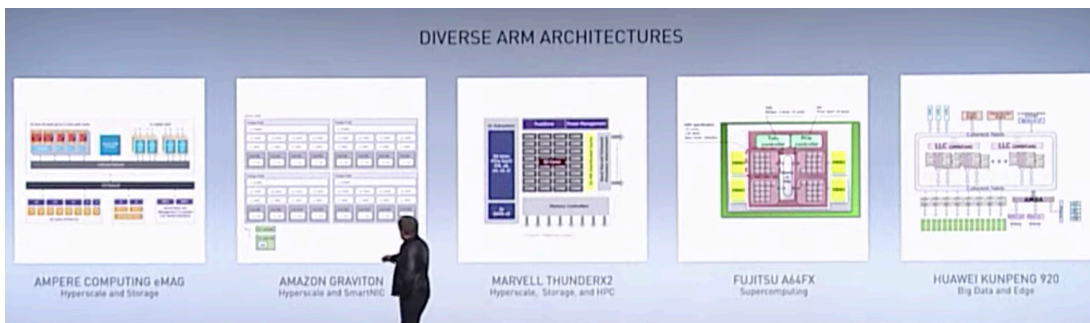


4-way SMT in ThunderX3

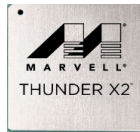
Multithread execution

- Four hardware threads per core
- Each thread includes full copy of Arm architecture state
- Threads share core pipeline resources
- To OS each thread appears as a regular Arm CPU
 - So four CPUs per core
- Area impact of 4-way SMT relative to no SMT: ~5%
- ThunderX3 has 60 cores / 240 threads per die





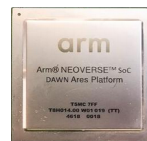
HPCs with ARMv8: server-level competitors



1. Marvell ThunderX product family



2. Fujitsu A64FX Arm chip



3. Neoverse N1 hyperscale reference design



4. Ampere Altra Arm Processor



5. Amazon Graviton

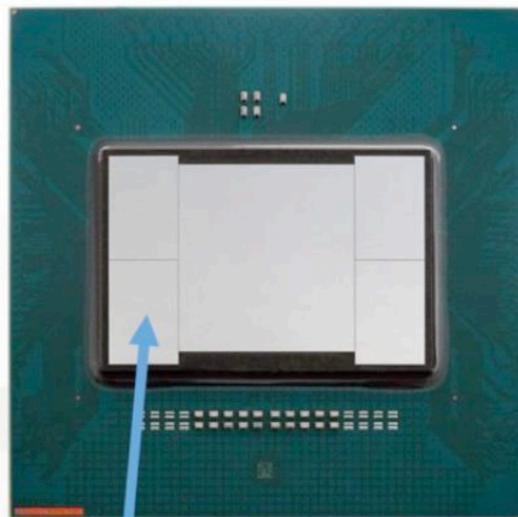


6. Huawei HiSilicon Kunpeng 920

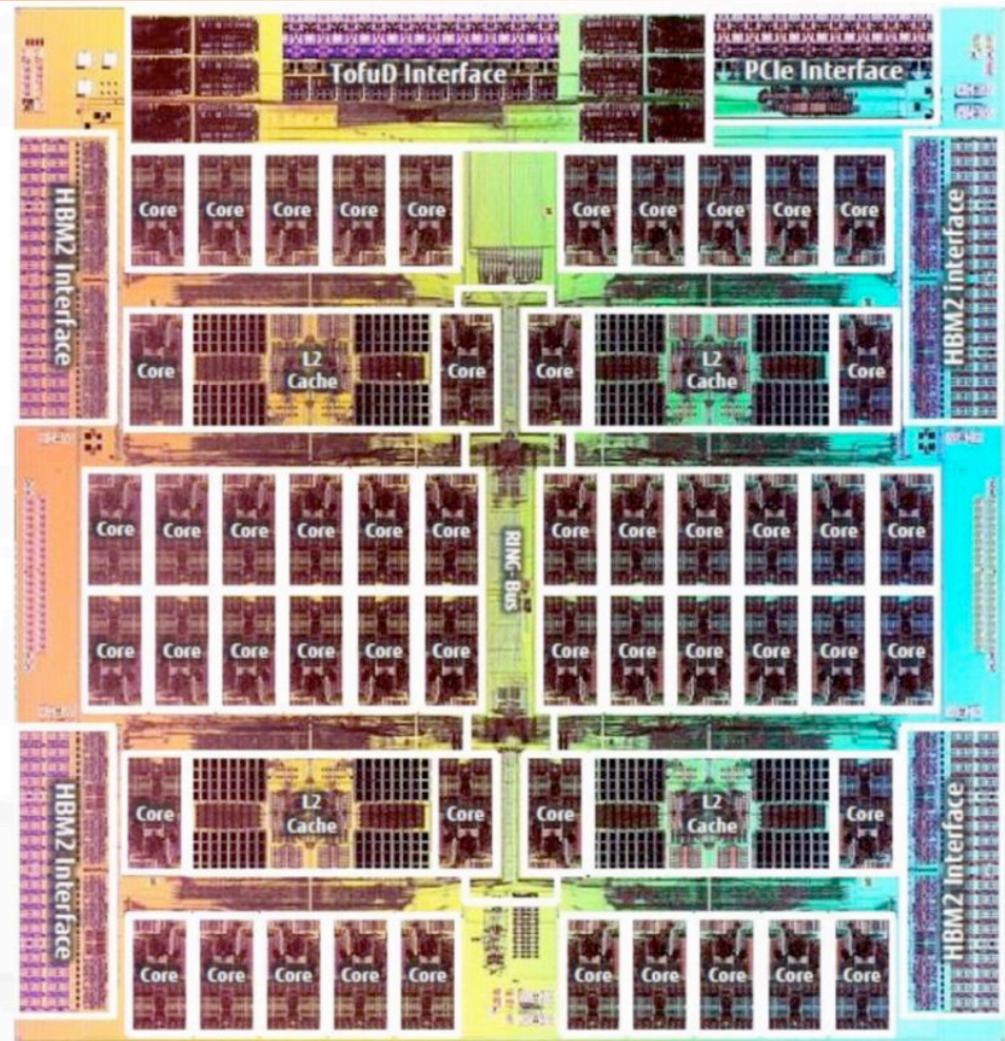


Fujitsu's A64FX ARM Chip

- TSMC 7nm FinFET
- CoWoS technologies for HBM2



HBM2



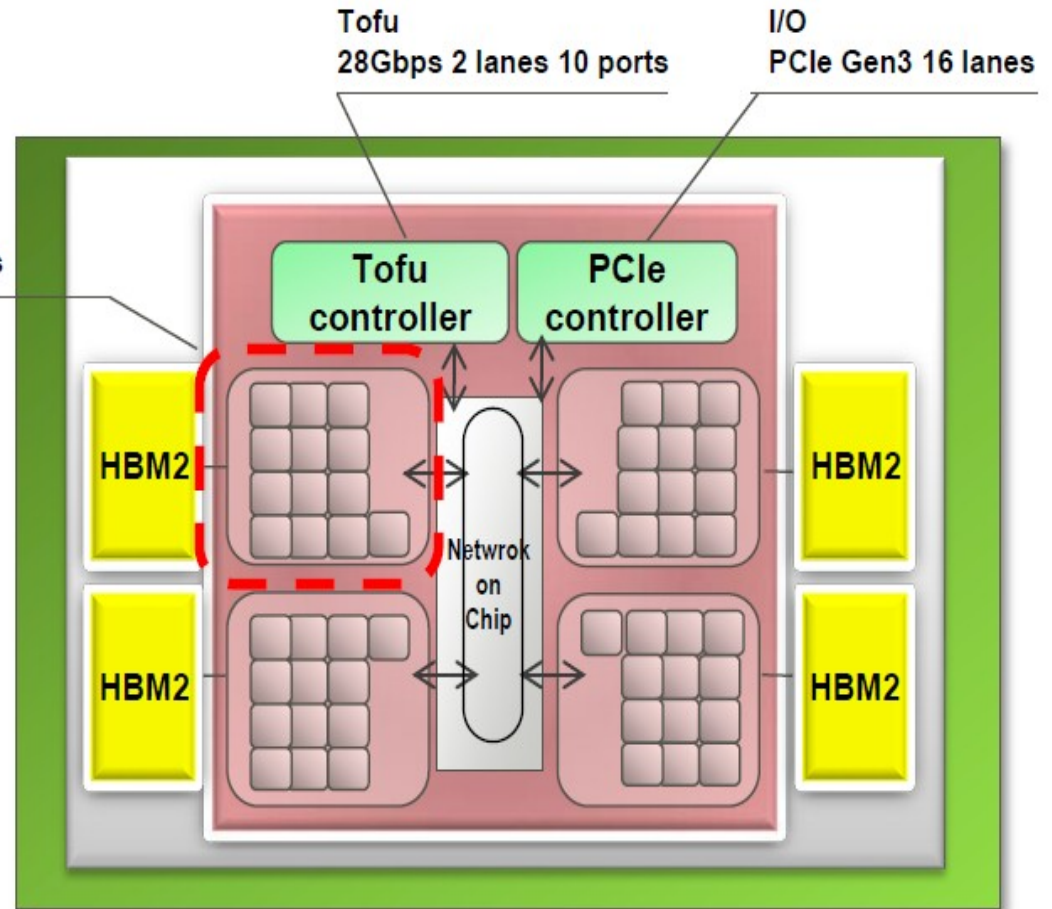


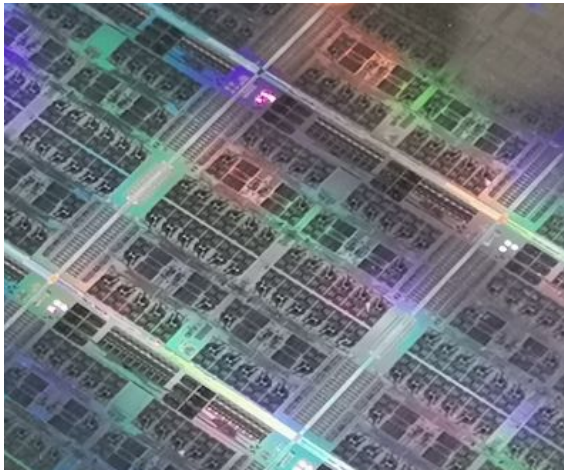
Fujitsu's A64FX Arm Chip: 48+4 cores

A64FX Arm

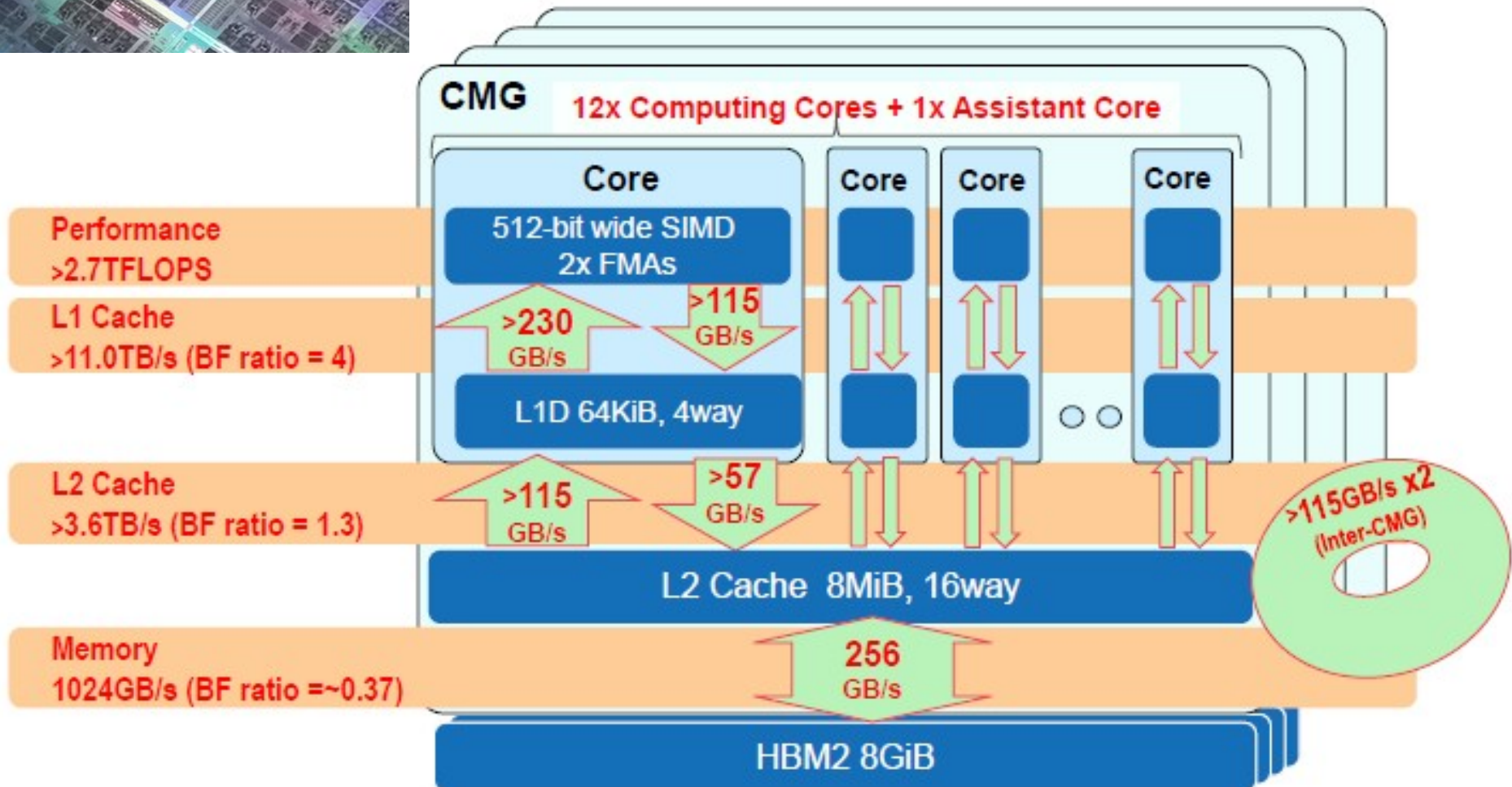
- Feature size: 7 nm
- Armv8.2-A spec with 512-bit SVE extensions
- HP math and a dot-product engine
- 4 core-memory groups
- NoC: a double ring bus
- cores in CMG linked by a crossbar to L2 cache & to HBM2 mem controller
- 8 MiB L2 cache; no L3 cache
- a Tofu-D controller on the die
- #1 TOP500 since Jun'20 uses A64FX package

CMG specification
13 cores
L2\$ 8MiB
Mem 8GiB, 256GB/s





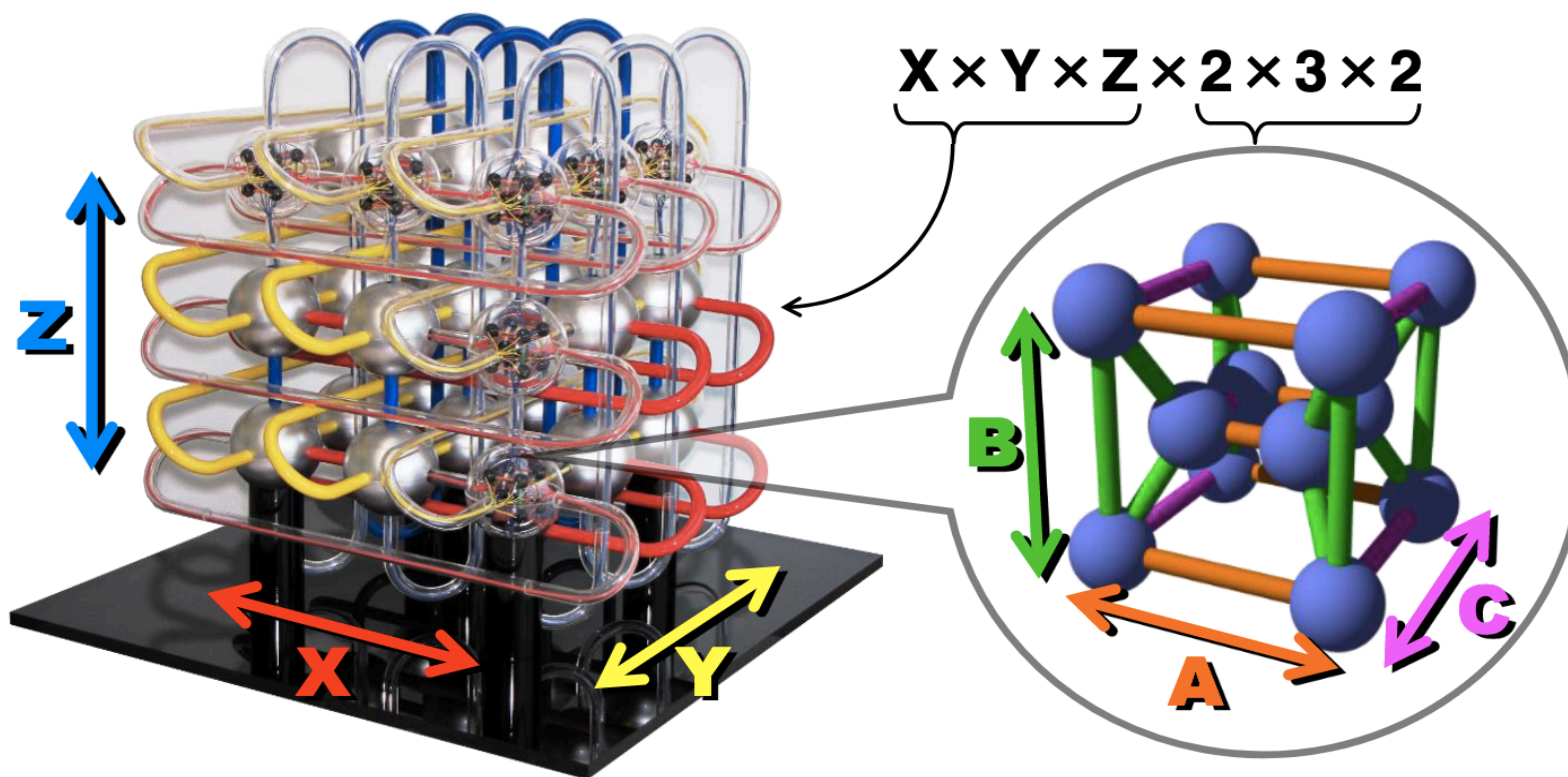
Block diagram of the A64FX chip



6D Mesh/Torus Network

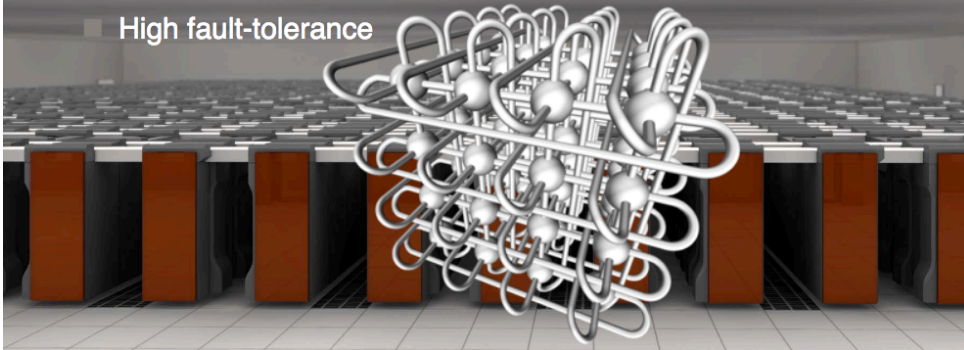


- Six coordinate axes: X, Y, Z, A, B, C
 - X, Y, Z: the size varies according to the system configuration
 - A, B, C: the size is fixed to $2 \times 3 \times 2$
- Tofu stands for “torus fusion”: $(X, Y, Z) \times (A, B, C)$



■ Tofu: Fujitsu's original 6D mesh/torus interconnect

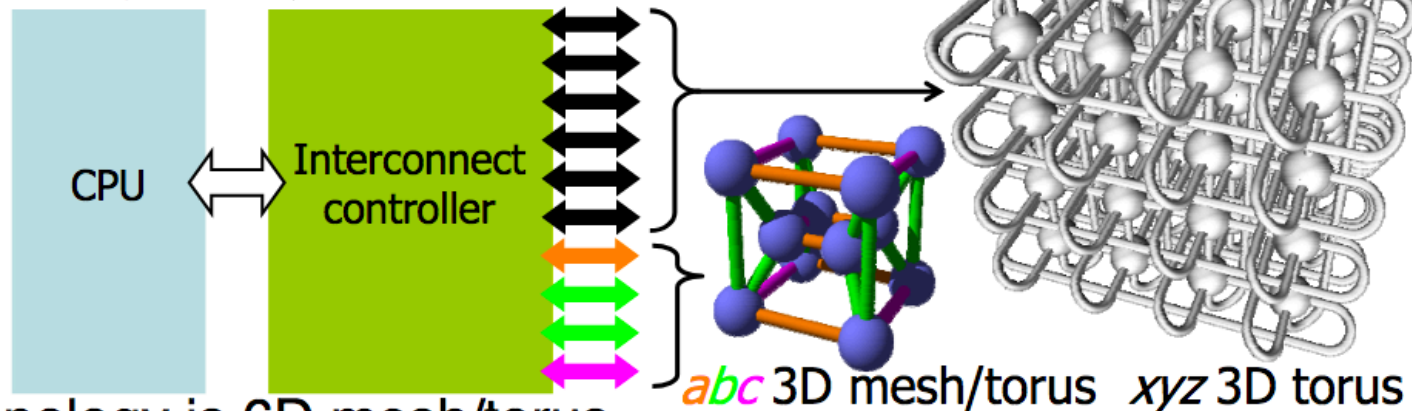
- High communication performance
- High system scalability
- High fault-tolerance



Tofu-D: 6D mesh/torus interconnect



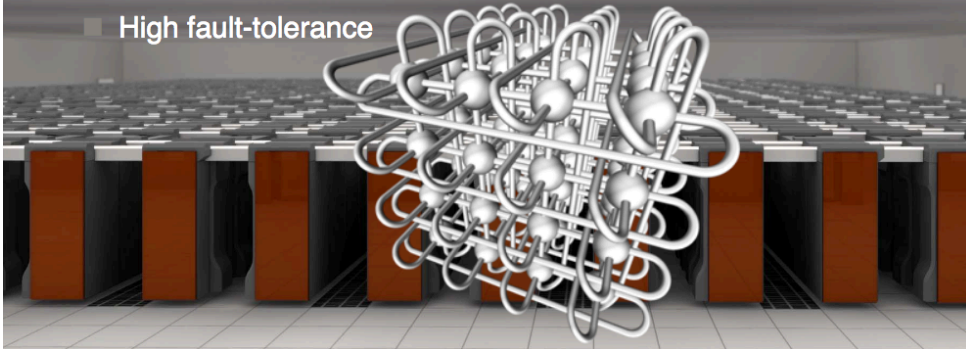
- 6 links \Rightarrow Scalable xyz 3D torus
- 4 links \Rightarrow Fixed size abc 3D mesh/torus
 - $|a|=2, |b|=3, |c|=2 \Rightarrow 12$ nodes



- Total topology is 6D mesh/torus
 - Cartesian product of xyz and abc mesh/torus

■ Tofu: Fujitsu's original 6D mesh/torus interconnect

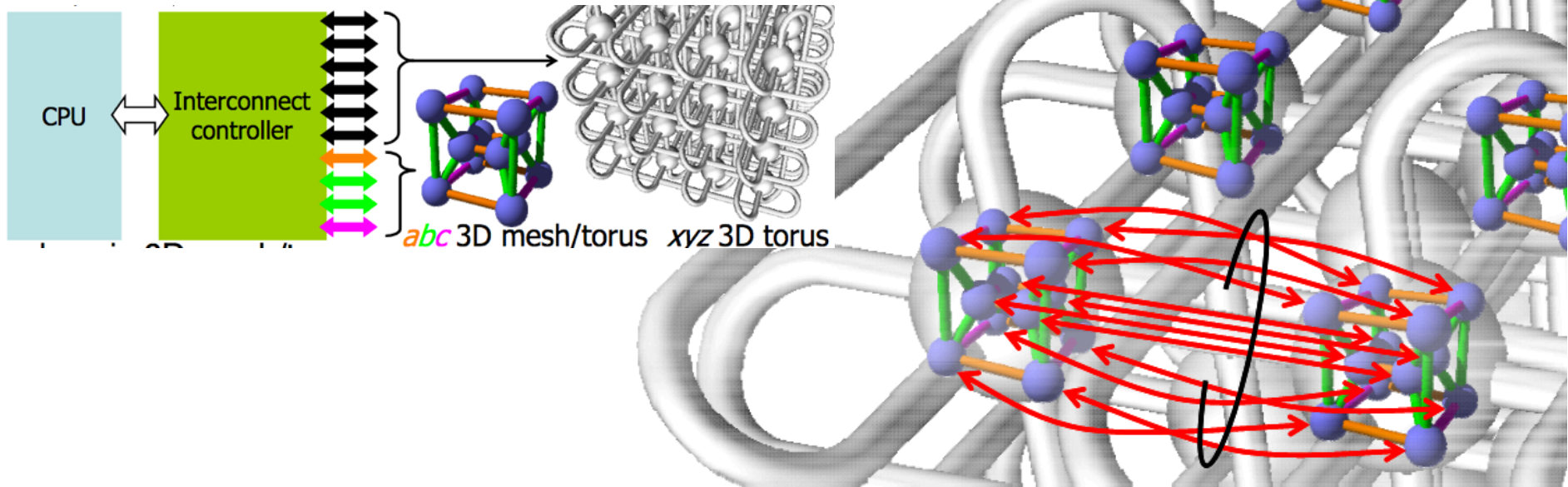
- High communication performance
- High system scalability
- High fault-tolerance

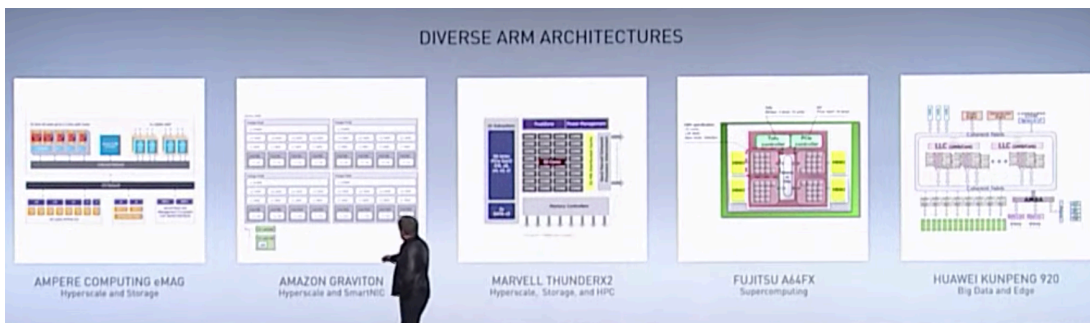


Tofu-D:
6D mesh/torus interconnect



■ Each pair of adjacent *abc* mesh/torus is interconnected with twelve links





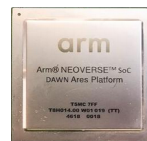
HPCs with ARMv8: server-level competitors



1. Marvell ThunderX product family



2. Fujitsu A64FX Arm chip



3. Neoverse N1 hyperscale



4. Ampere Altra Arm Processor



5. Amazon Graviton



6. Huawei HiSilicon Kunpeng 920

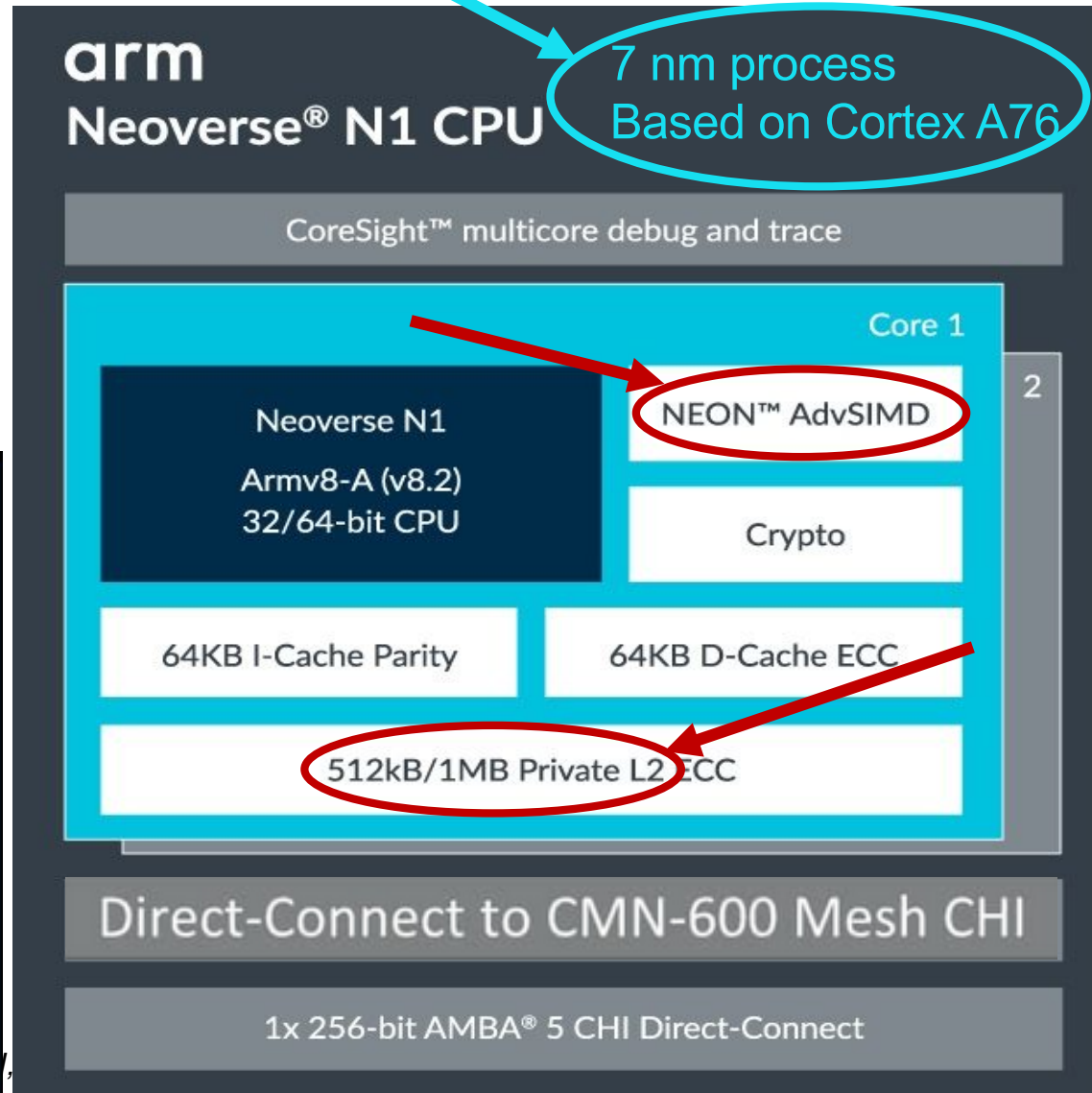
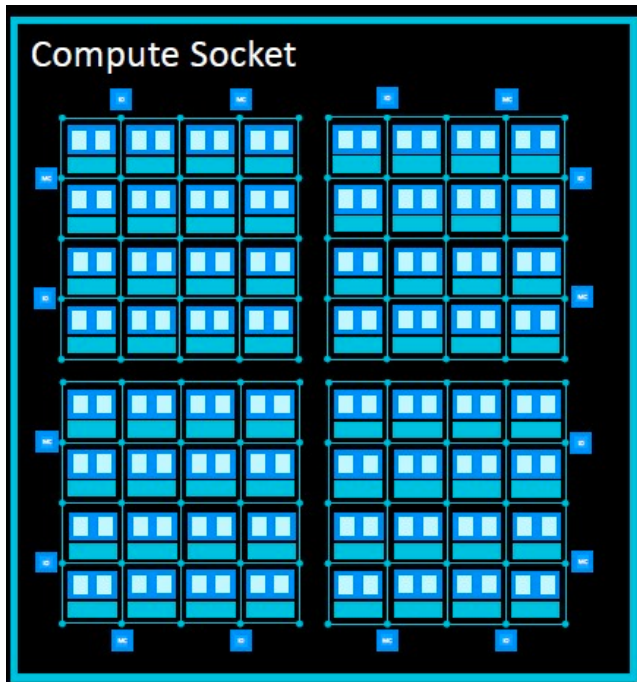


Arm Neoverse N1

(announced Feb'19)

Hyperscale
Datacenter

150W and beyond
64-128 core



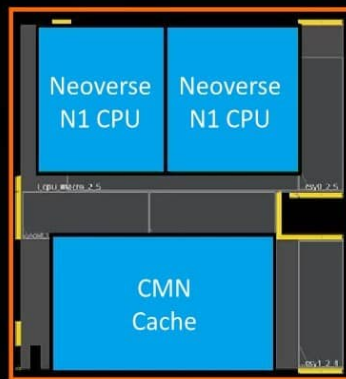


Arm Neoverse N1

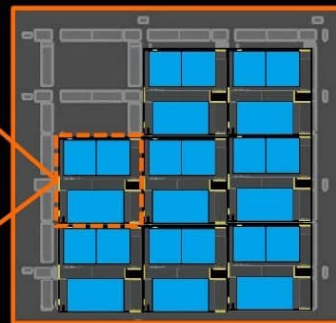


Building hyperscale compute in 7nm

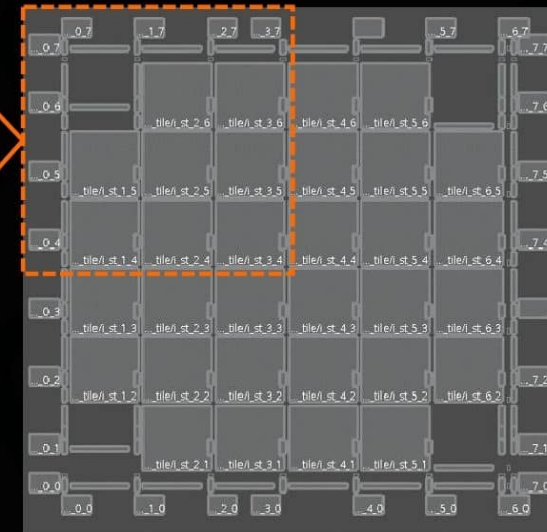
CPU Tile



Super Tile



Top-level Mesh



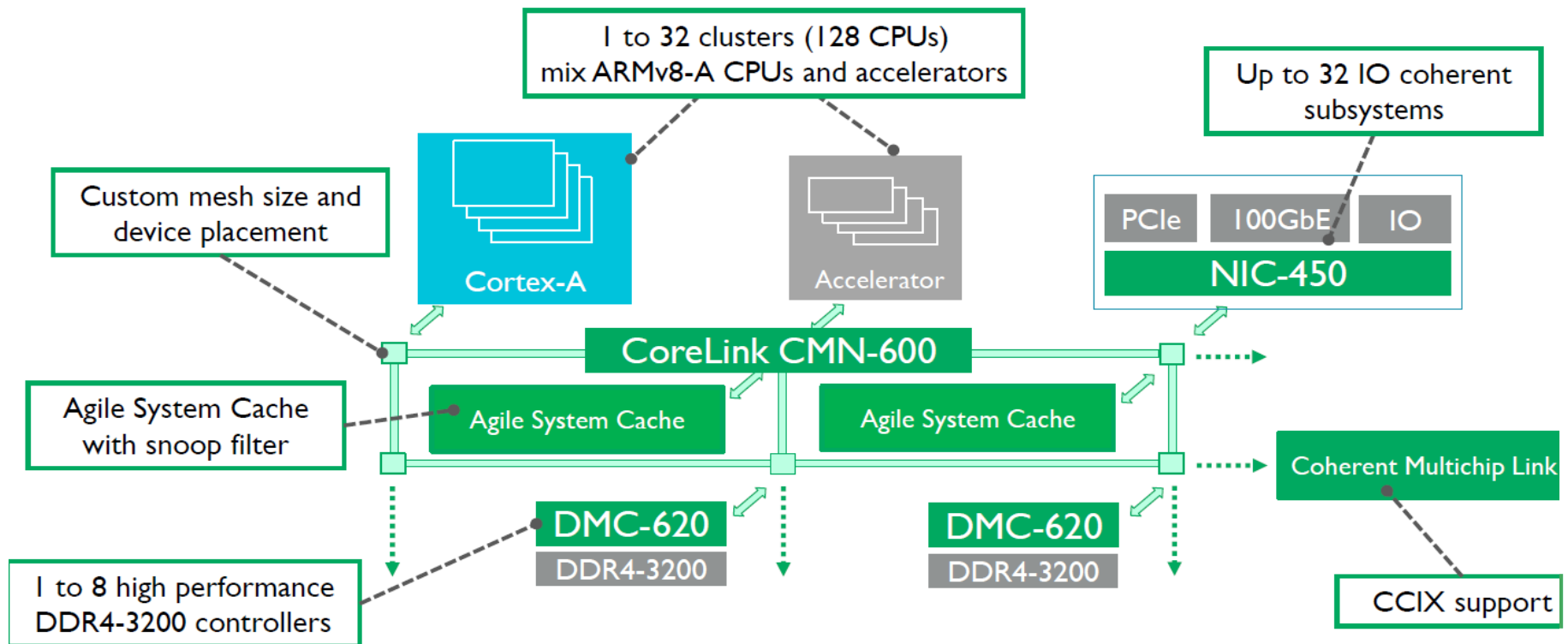
1.8-2.2GHz+
Coherent Mesh Network
CMN

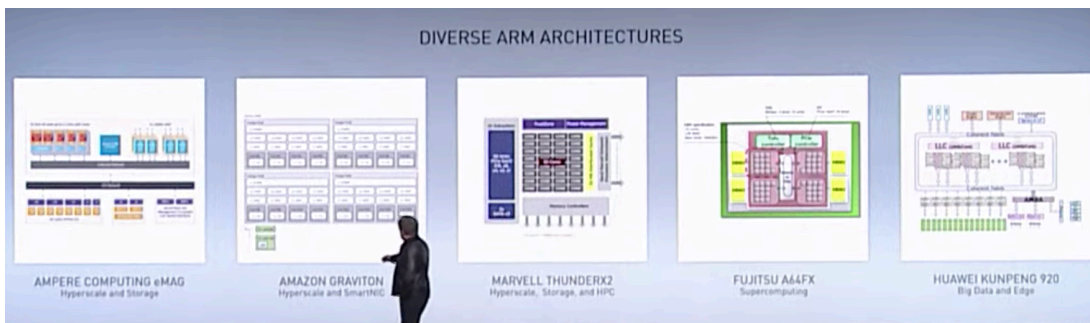


CMN-600 in Arm Neoverse (a SoC interconnect IP)

CMN: Coherent Mesh Network

New scalable coherent mesh architecture





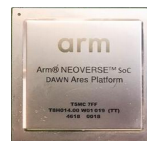
HPCs with ARMv8: server-level competitors



1. Marvell ThunderX product family



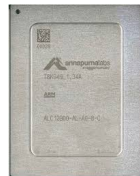
2. Fujitsu A64FX Arm chip



3. Neoverse N1 hyperscale reference design



4. Ampere Altra Arm Processor



5. Amazon Graviton



6. Huawei HiSilicon Kunpeng 920



Ampere Altra family: a SoC based on Neoverse N1

(announced Mar'20)

Ampere™ Altra™ processor complex

80 64-bit Arm CPU cores @ 3.0 GHz Turbo

- 4-Wide superscalar aggressive out-of-order execution
- Single threaded cores for performance and security isolation

Arm v8.2+ features

Large Cache, all with ECC Protection

- 64 KB L1 I/D-cache per core
- 1 MB L2 cache per core
- 32 MB system level cache

2x 128 SIMD Units

int8 and fp16 for ML Inference performance

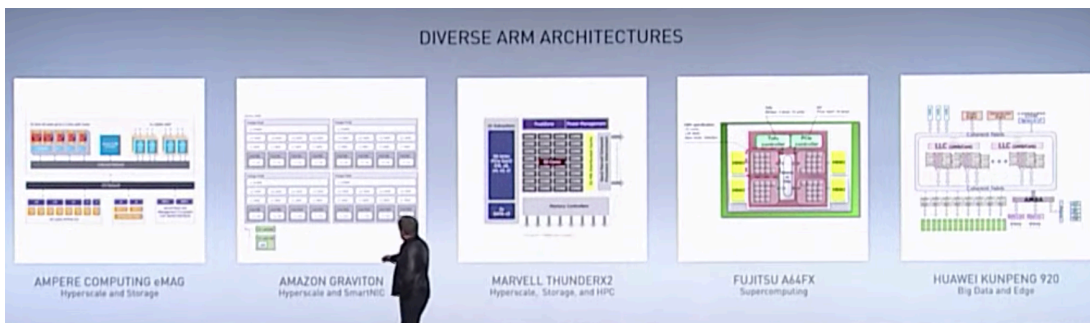
8 72-bit DDR4-3200 channels exceeding 200 GB/s per socket

Embargo: March 3, 2020 (6:00 AM Pacific time)



Ampere® Altra™ Max
7nm

- Up to 128 Cores



HPCs with ARMv8: server-level competitors



1. Marvell ThunderX product family



2. Fujitsu A64FX Arm chip



3. Neoverse N1 hyperscale reference design



4. Ampere Altra Arm Processor



5. Amazon Graviton



6. Huawei HiSilicon Kunpeng 920



Amazon Web Services (AWS): Graviton2 with Arm Neoverse N1 cores



AWS Graviton2

- 7 nm process
- up to 64 Arm Neoverse N1 core
- 1 MiB L2 per core
- shared 32 MiB L3
- without SMT
- single-socket

AWS Designing a 32-Core Arm Neoverse N1 CPU for Cloud Servers

by [Anton Shilov](#) on December 2, 2019 1:00 PM EST

Posted in [Servers](#) [CPUs](#) [Arm](#) [Amazon](#) [AWS](#) [Neoverse N1](#)



Graviton1 Processor



First Arm-based processor in major cloud



Built on 64-bit ARM Neoverse cores with AWS-designed 16 nm silicon



Up to 16 vCPUs, 10 Gbps enhanced networking, 3.5 Gbps EBS bandwidth

Graviton2 Processor



Built with 64-bit Arm Neoverse cores with AWS-designed 7 nm silicon process



Up to 64 vCPUs, 25 Gbps enhanced networking, 18 Gbps EBS bandwidth



7x performance,
4x compute cores,
5x faster memory



Graviton2 vs. AMD & Intel

	m6g	m5a	m5n
CPU Platform	Graviton2	EPYC 7571	Xeon Platinum 8259CL
vCPUs		64	
Cores Per Socket	64	32	24 (16 instantiated)
SMT	-	2-way	2-way
CPU Sockets	1	1	2
Frequencies	2.5GHz	2.5-2.9GHz	2.9-3.2GHz
Architecture	Arm v8.2	x86-64 + AVX2	x86-64 + AVX512
µarchitecture	Neoverse N1	Zen 1 / Naples	Cascade Lake
L1I Cache	64KB	64KB	32KB
L1D Cache	64KB	32KB	32KB
L2 Cache	1MB	512KB	1MB
L3 Cache	32MB shared	8MB shared per 4-core CCX	35.75MB shared per socket
Memory Channels	8x DDR4-3200	8x DDR-2666 (2x per NUMA-node)	6x DDR4-2933 per socket
NUMA Nodes	1	4	2
DRAM		256GB	

SNC

AJProença,

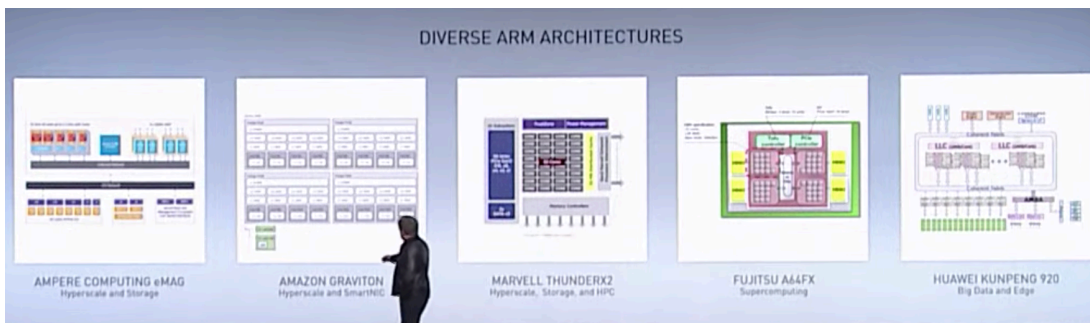
Comparison of major ARM servers



	Marvell ThunderX3 110xx	Cavium ThunderX2 9980-2200	Ampere Altra Q80-33	Amazon Graviton2
Process Technology	TSMC 7nm	TSMC 16 nm	TSMC 7 nm	TSMC 7nm
Die Type	Monolithic <i>or</i> Dual-Die MCM	Monolithic	Monolithic	Monolithic
Micro-architecture	Triton	Vulcan	Neoverse N1 (Ares)	
Cores	60 (1 Die) Switched 3x Ring 96 (2 Die)	32 Ring bus	80 Mesh	64 Mesh
Threads	240 (1 Die) 384 (2 Die)	128	80	64
Max. number of sockets	2	2	2	1
Base Frequency	?	2.2 GHz	-	-
Turbo Frequency	3.1 GHz	2.5 GHz	3.3 GHz	2.5 GHz
L3 Cache	90MB	32 MB	32 MB	32 MB
DRAM	8-Channel DDR4-3200	8-Channel DDR4-2667	8-Channel DDR4-3200	8-Channel DDR4-3200

SMT →





HPCs with ARMv8: server-level competitors



1. Marvell ThunderX product family



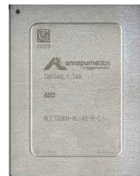
2. Fujitsu A64FX Arm chip



3. Neoverse N1 hyperscale reference design



4. Ampere Altra Arm Processor



5. Amazon Graviton



6. Huawei HiSilicon Kunpeng 920

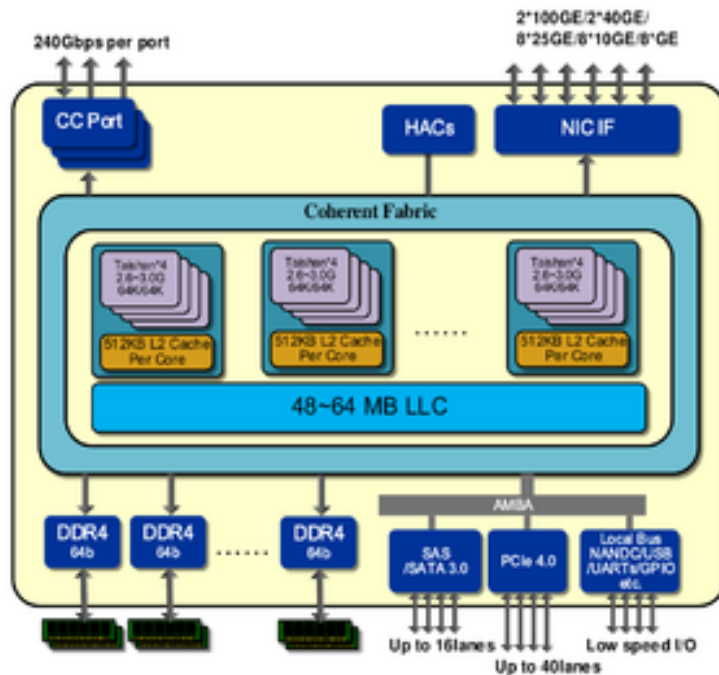


The Huawei Kunpeng 920 (previously known as HiSilicon Hi1620)

(launched in 2019)

Huawei Kunpeng 920 is based on TaiShan V110 core, a semi-custom ARM Cortex-A72

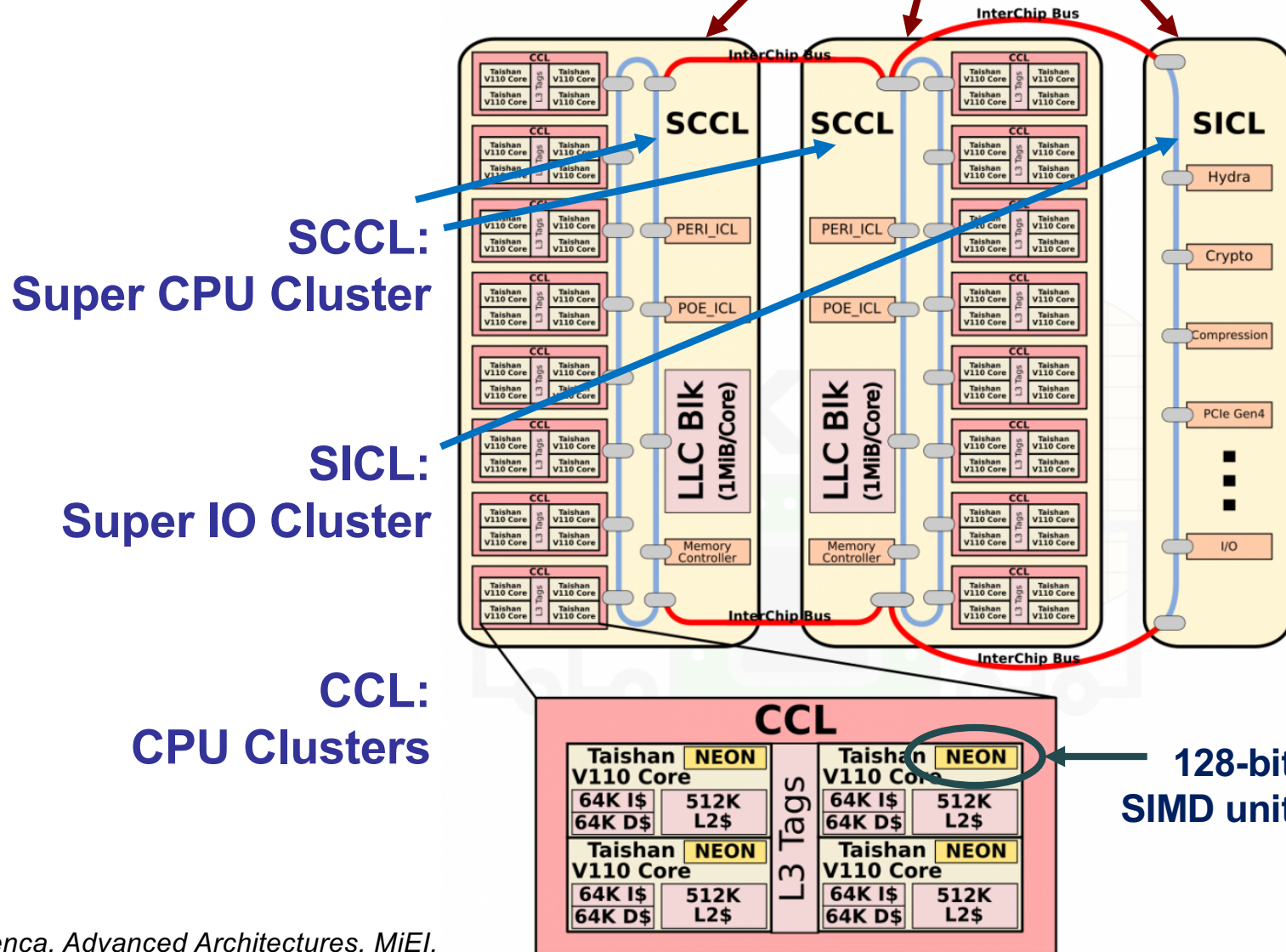
Hi1620 Specifications Overview



CPU core	Up to 64 ARMv 8.2 cores, 3.0 GHz, 48-bit physical address 4 issue OoO superscalar design 64 KB L1 I Cache and 64 KB L1 D cache 128-bit SIMD unit
L2 cache	512 KB private per core, 24 MB total
L3 cache	48 MB shared for all (1 MB/core), Partitioned
Memory	8-channel DDR4-2400/2666/2933/3200 16 ranks/channel, 1DPC and 2DPC configurations x4/x8 support ECC, SDDC, DDDC
PCIe	40 lanes of PCIe Gen4.0 16x
Integrated I/O	8 lanes of ETH, Combo MACs, supporting 2 x 100GE, 2 x 40GE, 8 x 25GE/10GE, 10 x GE, supporting SR-IOV RoCEv2/RoCEv1 x4 USB 3.0 x8 SAS 3.0 x2 SATA 3.0
Crypto engine	AES, DES/3DES, MD5, SHA1, SHA2, HMAC, CMAC Up to 100 Gbit/s
Compression	GZIP, LZS, LZ4 Up to 40 Gbit/s (compress)/100 Gbit/s (decompression)
RAID	RAID5/6, DIF, XOR, PQ acceleration
CCIX	Cache coherency interface for accelerator, like Xilinx FPGA World's 1st CCIX solution
Scale-up	Coherent SMP interface for 2P/4P 3*240Gbps bandwidth
Power	TDP ~150 W (48C 2.6 GHz)

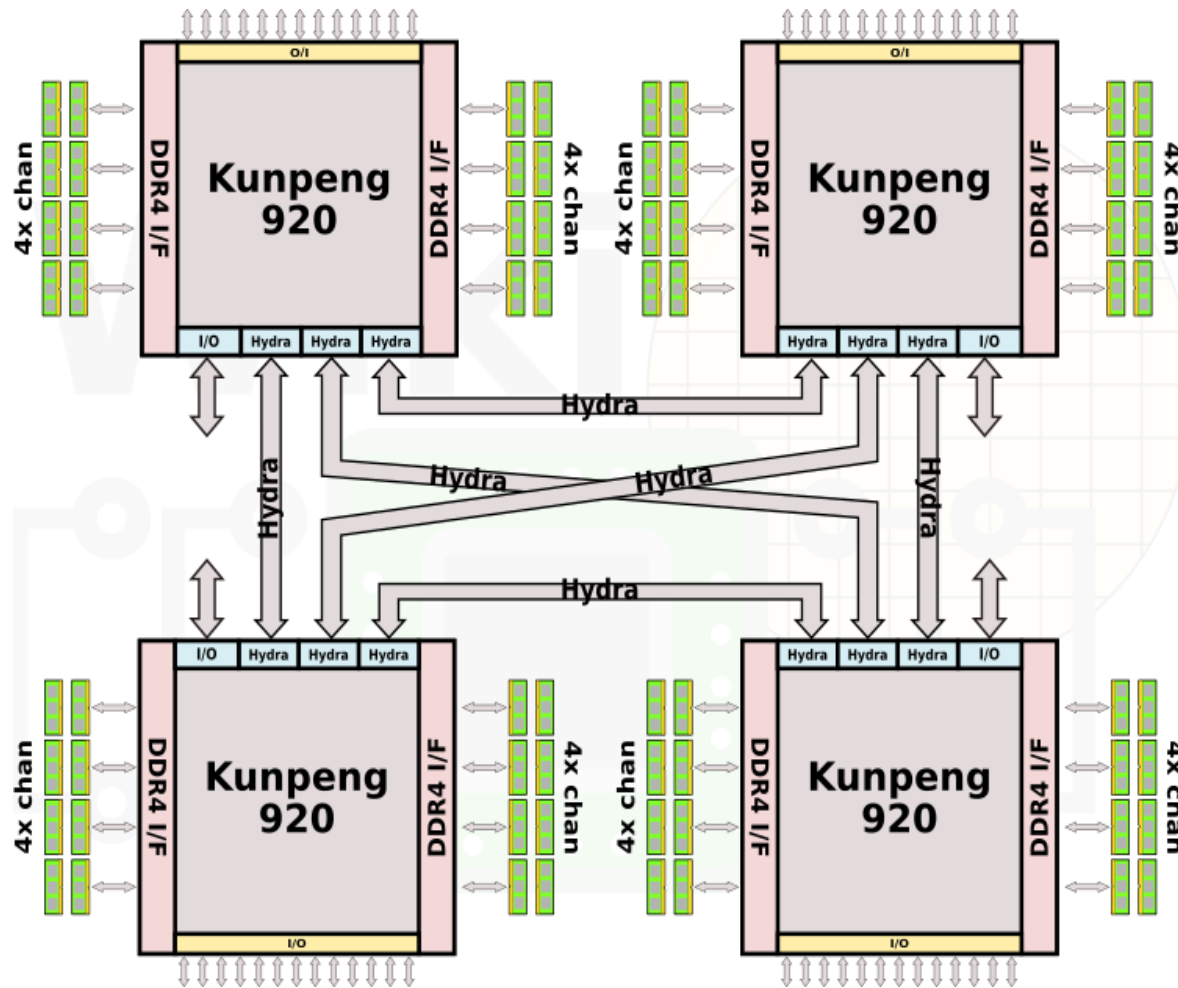


The Huawei Kunpeng 920: a multi-chip 48-64 cores





The Huawei Kunpeng 920: multi-socket support



Next-Gen Kunpeng 930:

- higher-performance core
- SMT support
- Arm SVE (vector comp)
- expected in 2021



Manycore chips/packages: an overview



Key server chips/packages that addresses those issues:

- Intel: the Xeon Processor Scalable family
- AMD: the Epyc Zen family
- Sunway: the SX260x0 family
- ARM: the ARMv8 server-level competitors
 - Marvell ThunderX family
 - Fujitsu A64FX Arm chip
 - Neoverse N1 hyperscale reference design
 - Ampere Altra Arm Processor
 - Amazon Graviton
 - Huawei HiSilicon Kunpeng 920
- **Cerebras: a Wafer Scale Engine**
- **Apple** *(no chipleths!)*: **the SoC approach** *(no chipleths!)*



Cerebras: a Wafer Scale Engine (WSE) (Aug'19)



Cerebras Wafer Scale Engine (WSE):
the largest chip ever built

46,225 mm² chip

56x larger than the biggest GPU ever made

400,000 core

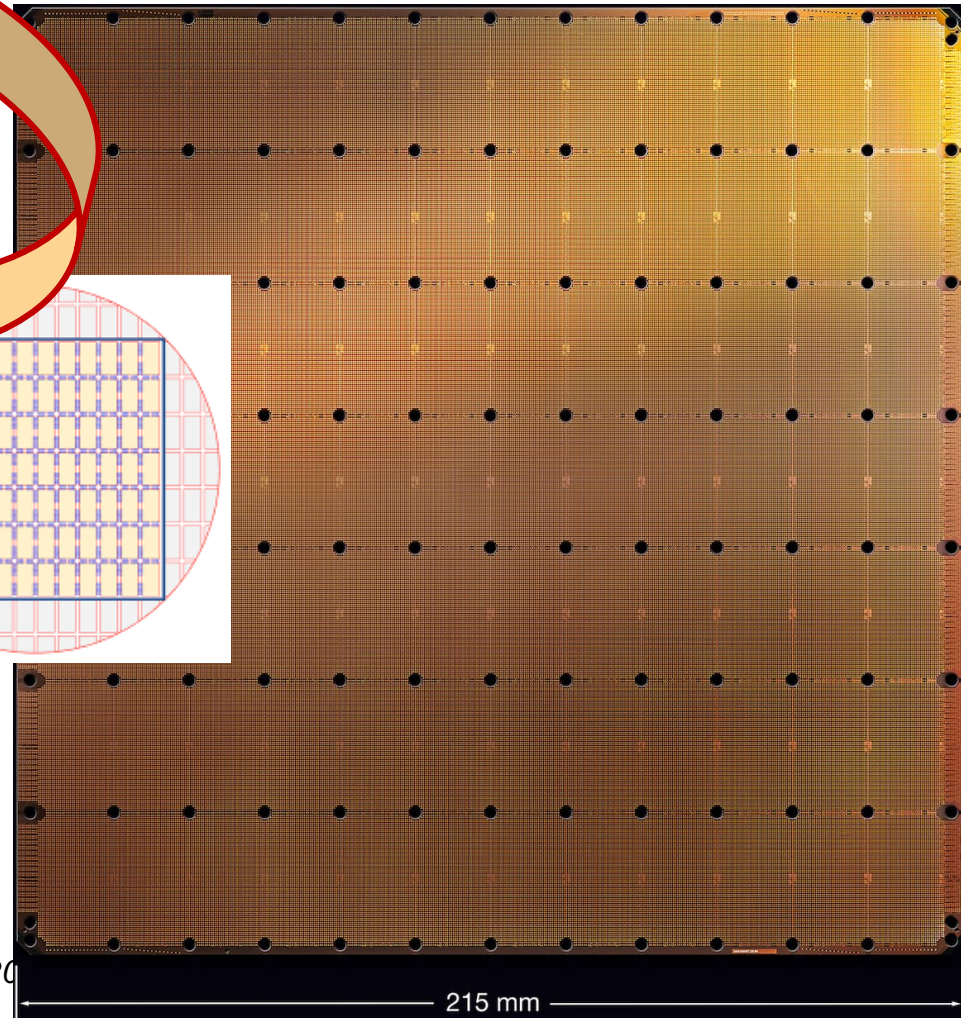
78x more cores

18 GB on-chip SRAM

3000x more on-chip memory

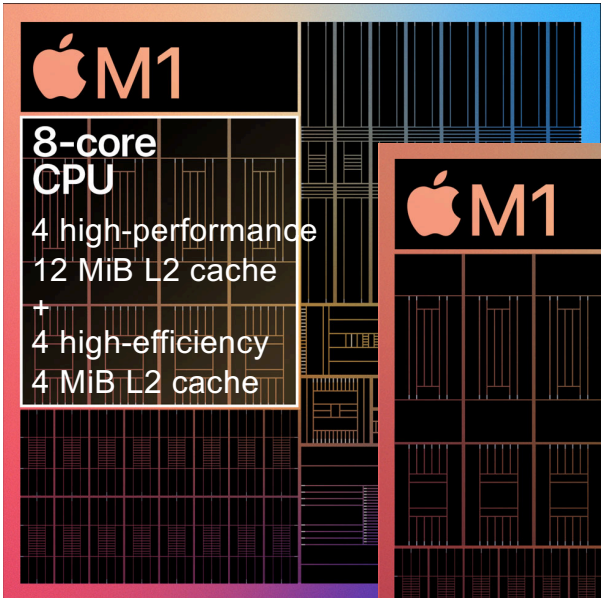
100 Pb/s interconnect

33,000x more bandwidth

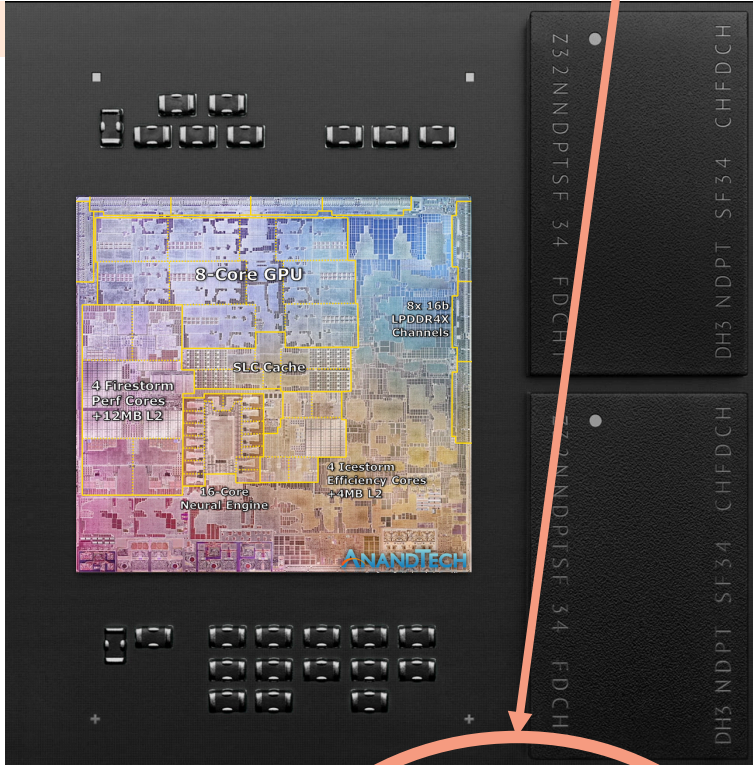
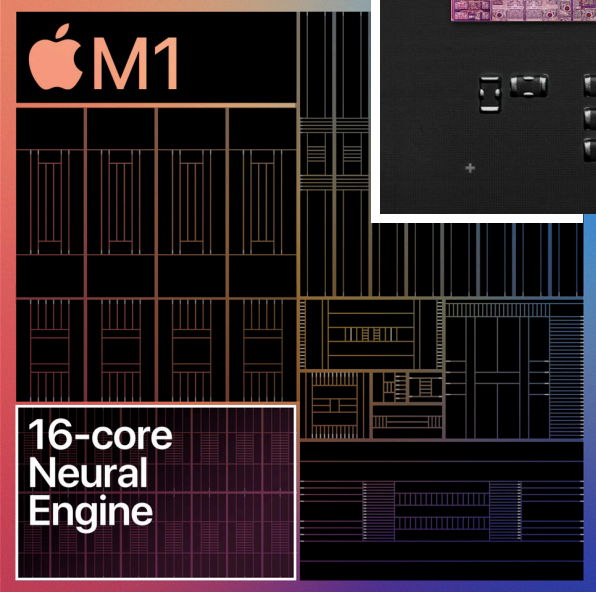
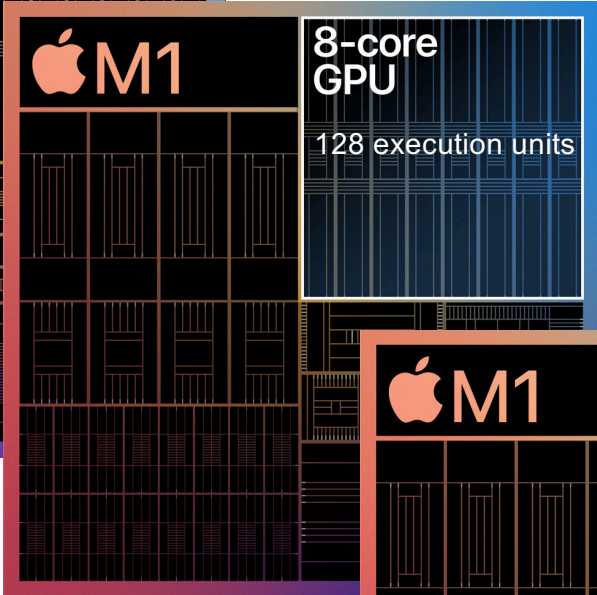




Apple M1 SoC

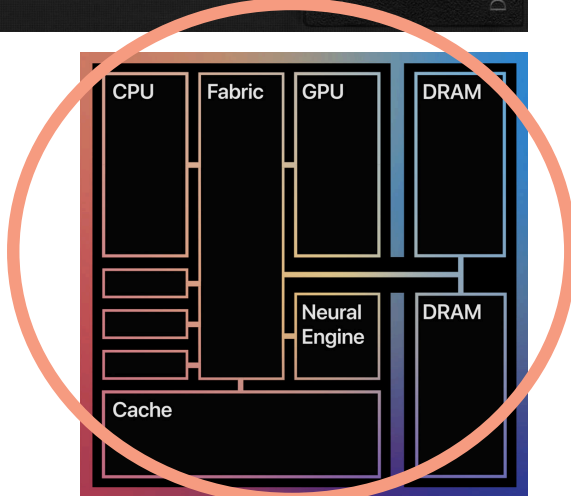


5-nanometer process



Apple claims SoC is better:

- shorter latencies
- better silicon process (5 nm)
- better wafer fabrication (increased yield)



turing lecture

Nobel equivalent in Computer Science

DOI:10.1145/3282307

Innovations like domain-specific hardware, enhanced security, open instruction sets, and agile chip development will lead the way.

BY JOHN L. HENNESSY AND DAVID A. PATTERSON

A New Golden Age for Computer Architecture

WE BEGAN OUR Turing Lecture June 4, 2018¹¹ with a review of computer architecture since the 1960s. In addition to that review, here, we highlight current challenges and identify future opportunities, projecting another golden age for the field of computer architecture in the next decade, much like the 1980s when we did the research that led to our award, delivering gains in cost, energy, and security, as well as performance.

“Those who cannot remember the past are condemned to repeat it.”
—George Santayana, 1905



engineers, including ACM A.M. Turing Award laureate Fred Brooks, Jr., thought they could create a single ISA that would efficiently unify all four of these ISA bases.

They needed a technical solution

And now?