

Computer Architecture

Hot Chips

Intel & AMD 64/Zen 3

© 2022 IEEE International Solid-State Circuits Conference

Dr Jeff Drobman

website



drjeffsoftware.com/classroom.html

email



jeffrey.drobman@csun.edu



Index

- ❖ Intel Roadmap [Hot Chips](#) → slide 3
- ❖ AMD Roadmap [ISSCC](#) → slide 13
- ❖ AMD RADN3 GPU → slide 33



Section



Intel Roadmap

Hot Chips 34

Executive Summary

- Intel has a rich history of foundational process innovations in pursuit of Moore's Law.
- Advanced packaging gives architects and designers new tools in their pursuit of Moore's Law.
- Intel has a full pipeline of research that gives us the confidence of maintaining Moore's Law.
- All considered, numerous options are available to designers and architects in their continued mission to deliver Moore's Law



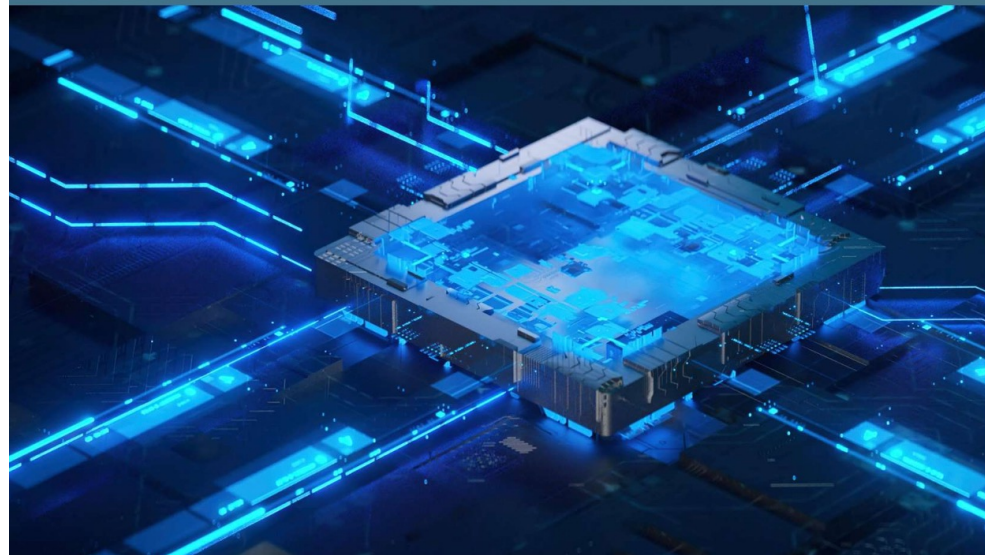
By Dr. Ann Kelleher

*Executive Vice President and General Manager
of Technology Development*

Hot Chips 34

intel newsroom

[Newsroom Home](#) | [Newsroom Search](#) | [International Sites](#) | [Email Opt-in](#)



A New Era of Chipmaking to Meet the World's Demand for Compute

Intel CEO Pat Gelsinger details how advanced compute and packaging are needed to meet the world's insatiable demand for compute and implement fully immersive digital experiences at Hot Chips 34.



Hot Chips 34



- **Meteor Lake, Arrow Lake and Lunar Lake processors** will transform personal computers with tile-based chip designs that create efficiencies in manufacturing, power and performance. This is done through discrete CPU, GPU, SoC and I/O tiles stacked in 3D configurations using Intel's Foveros interconnect technology. This platform transformation is reinforced by industry support for the open Universal Chiplet Interconnect Express (UCIe™) specification enabling chiplets designed and manufactured on different process technologies by different vendors to work together when integrated with advanced packaging technologies.
- **Intel Data Center GPU, code-named Ponte Vecchio**, was built to address the compute density across high performance computing (HPC) and AI supercomputing workloads. It also takes full advantage of Intel's open software model, using OneAPI to simplify API abstractions and cross-architecture programming. Ponte Vecchio is composed of several complex designs that manifest in tiles, connected using a combination of embedded multi-die interconnect bridge (EMIB) and Foveros advanced packaging technologies. The high-speed MDFI interconnect allows the package to scale up to two stacks, allowing a single package to contain more than 100 billion transistors.



Hot Chips 34

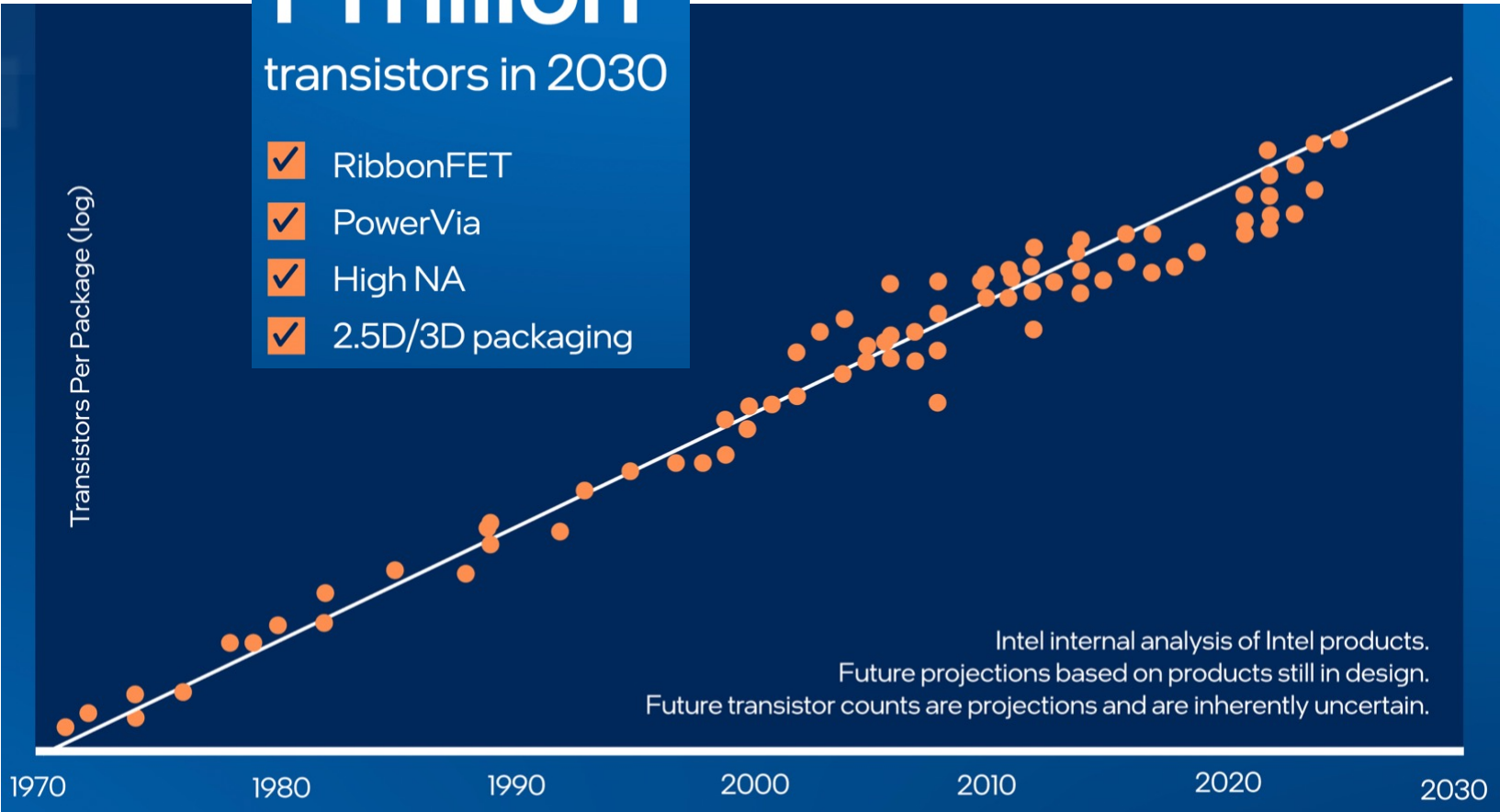
- **Xeon D-2700 and 1700 series** are designed to address edge use cases for 5G, IoT, enterprise and cloud applications, with special consideration to the power and space constraints that are common in many real-world implementations. These chips are also examples of tile-based design, including state-of-the-art compute cores, 100G Ethernet with flexible packet processor, inline crypto acceleration, time coordinated computing (TCC), time-sensitive networking (TSN) and built-in optimization for AI processes.
- **FPGA technology** continues to be a powerful and flexible tool for hardware acceleration, with particular promise for radio frequency (RF) applications. Intel has identified new efficiencies by integrating digital and analog chiplets, as well as chiplets from different process nodes and foundries, cutting development time and maximizing flexibility for developers. Intel will share the results of its chiplet-based approach in the near future.

Hot Chips 34

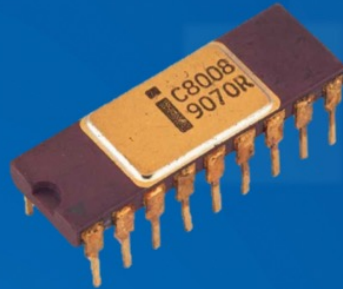
Aspiring to
1 Trillion
transistors in 2030

- ✓ RibbonFET
- ✓ PowerVia
- ✓ High NA
- ✓ 2.5D/3D packaging

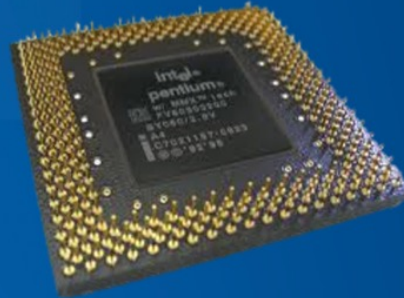
Transistors Per Package (log)



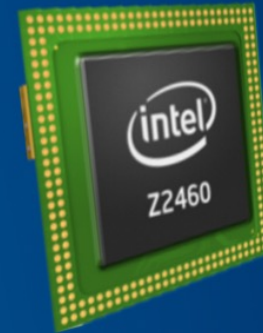
Hot Chips 34



Leadframe /
Wirebond



Flip Chip
Ceramic

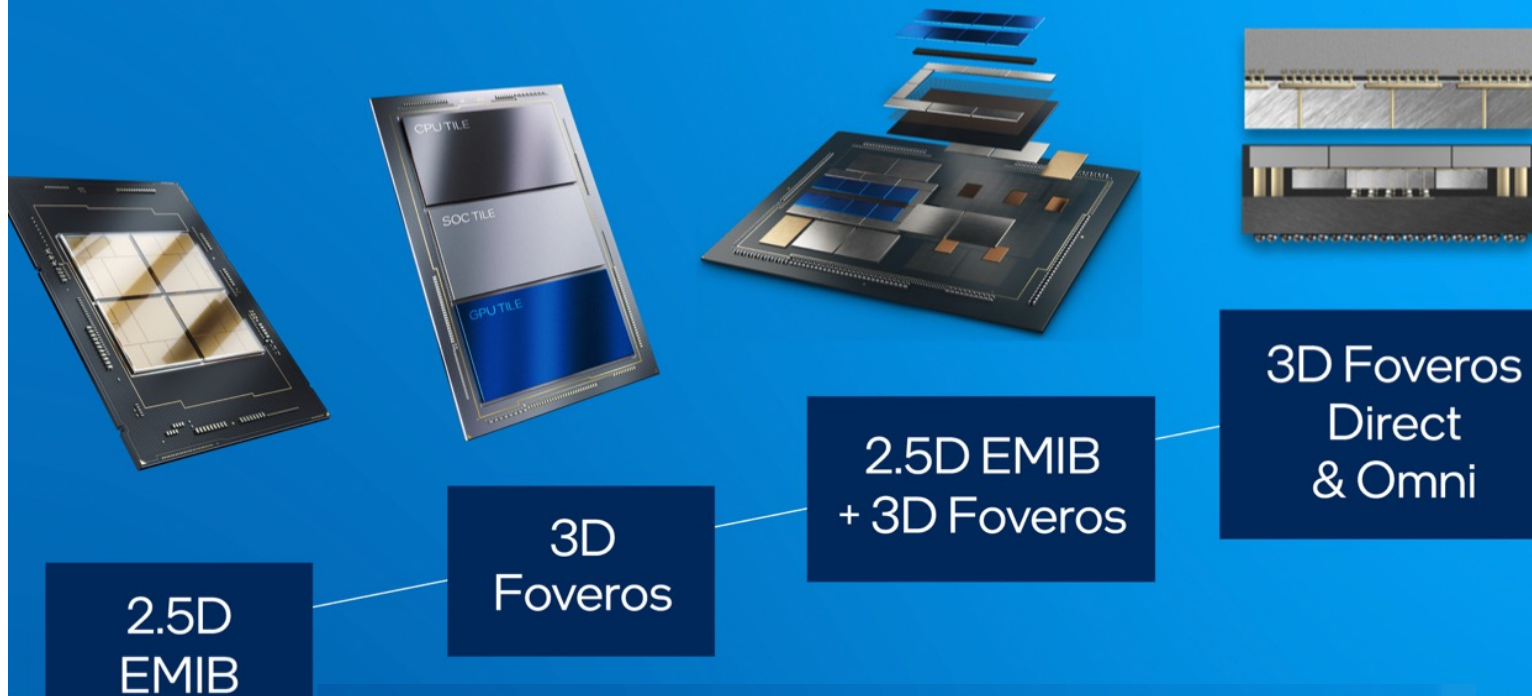


Flip Chip
Organic &
Multi Chip Pkg

Package main function:
provide power and signaling
from motherboard to die

Hot Chips 34

Advanced packaging era

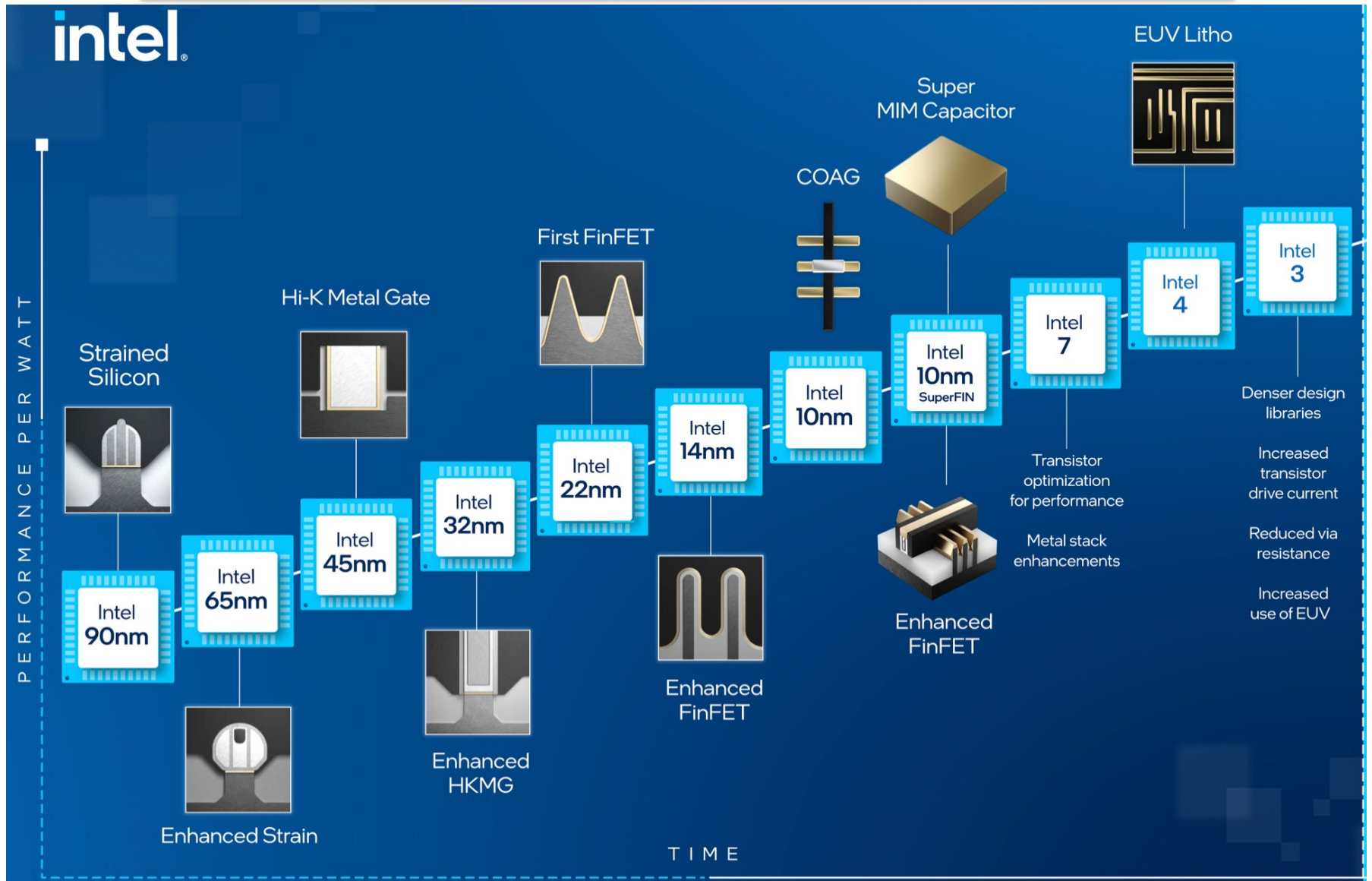


Added Package value:

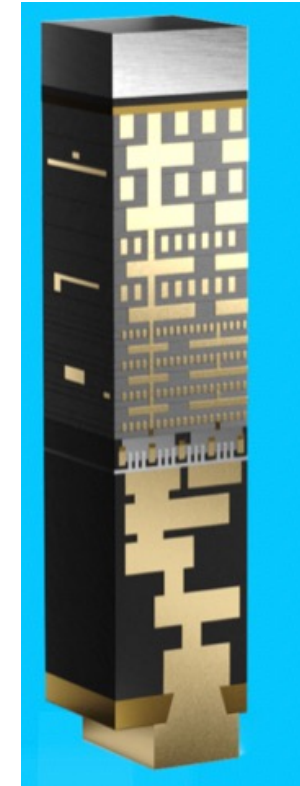
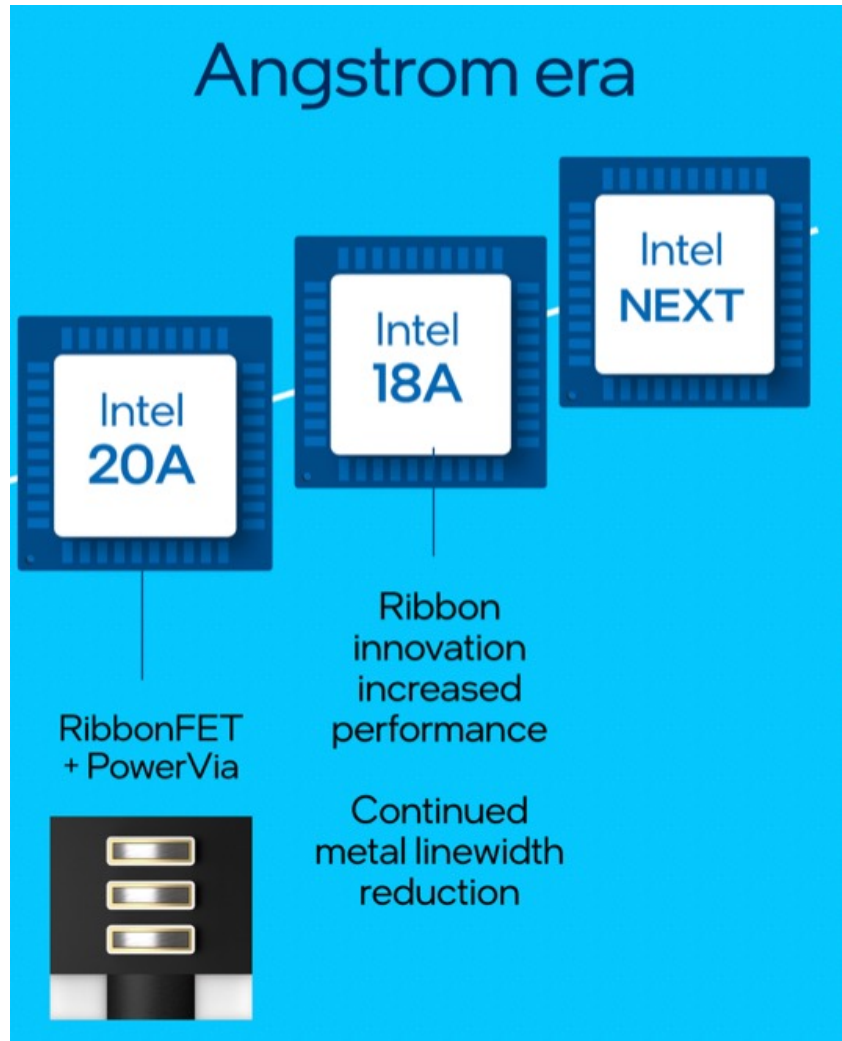
high density interconnects that enable larger die complexes from multiple process nodes



Hot Chips 34



Hot Chips 34



Section

AMD Roadmap

ISSCC

AMD

View Online



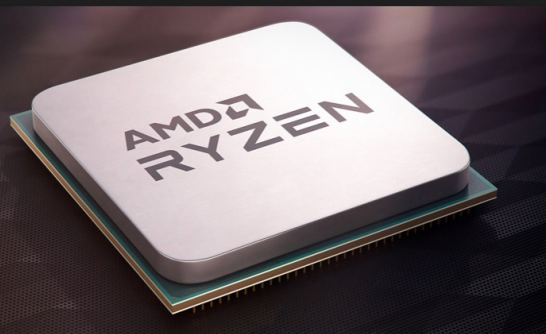
FROM A MODEST BEGINNING TO LEADING THE WORLD, OUR LEGACY CONTINUES

May marks a new phase here at our HQ in Austin, Texas!

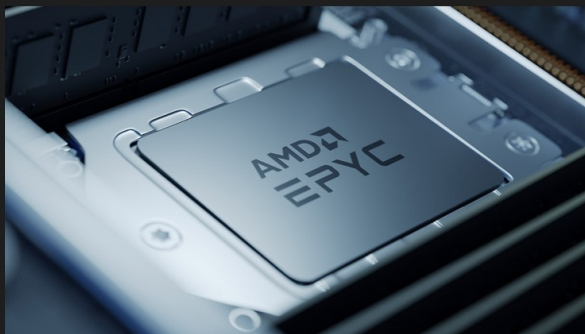
We're happy to see the "[Memory Lane](#)" discussions becoming more popular. Filled with vintage assets (i.e., ads and images) from the good old days the "[Name the AMDer!](#)" gained 136 views and 7 replies in April! Thank you [jeanpierrevelly](#), [DrJeffD](#), [john_springer](#), [donmc71](#), [Rick_Marz](#), and [petegasperini](#) for your answers!

AMD CPU/GPU

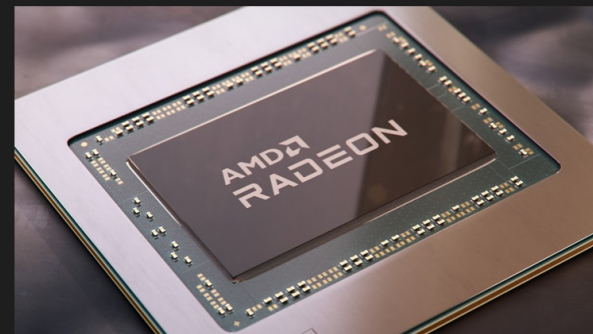
Ultimate Performance Technology



AMD Ryzen™ Processors



AMD EPYC™ Processors



AMD Radeon™ Graphics

Discrete Radeon GPU

(not so *discreet*)

“Zen 3” Market Segments



**AMD RYZEN™ 5000
SERIES MOBILE
PROCESSORS**



**AMD RYZEN™ 5000
SERIES DESKTOP
PROCESSORS**



**3RD GEN AMD EPYC™
SERVER PROCESSORS**

- Single CPU core across laptop, desktop, and server
- AMD 2nd Generation TSMC 7nm FinFET CPU
- Need to balance performance and power efficiency

AMD Process Roadmap

HIGH PERFORMANCE MOMENTUM



2017

2022



Zen uArch

THE EVOLUTION OF "ZEN"

RESOURCE	"ZEN"	"ZEN 2"	"ZEN 3"
Issue width	10	11	16
INT reg	168	180	192
INT sched	84	92	96
FP reg	160	160	160
ROB	192	224	256
FADD, FMUL, FMA	3/4/5	3/3/5	3/3/4
FP width	128	256	256
L1 BTB	256	512	1024
L2 BTB	4k	7k	6.5k

CORE COMPARISON



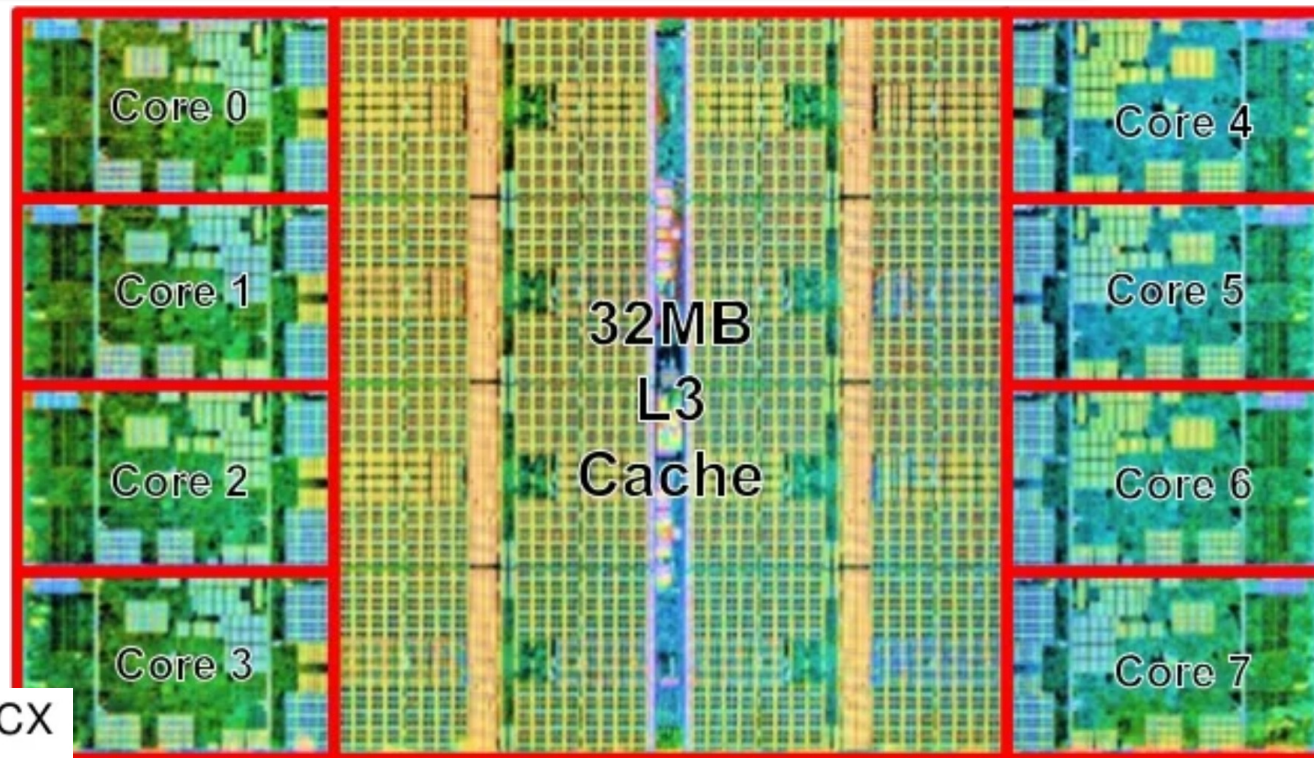
Zen uArch

THE EVOLUTION OF "ZEN"

RESOURCE	"ZEN"	"ZEN 2"	"ZEN 3"
LDQ	72	72	72
STQ	44	48	64
Micro-Op-cache	2k	4k	4k
L1 Icache	64k	32k	32k
L1 Dcache	32k	32k	32k
L2 cache	512k	512k	512k
L3 cache/core	2M	4M	4M
L2 TLB size	1.5k	2k	2k
L2 TLB latency	8	6	6
L2 latency	12	12	12
L3 latency	35	39	46

CACHE COMPARISON

AMD Zen 3 CCX



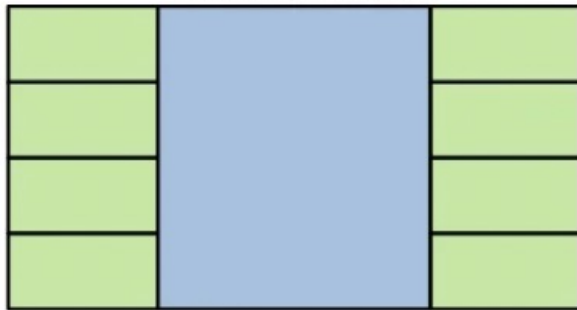
- Increased 4→8 cores per CCX
- Same max. 4MB of L3 per core
- Single CCX per chiplet
- 2x L3 cache directly accessible per core

**“Zen 3”: AMD 2nd Generation 7nm x86-64
Microprocessor Core**

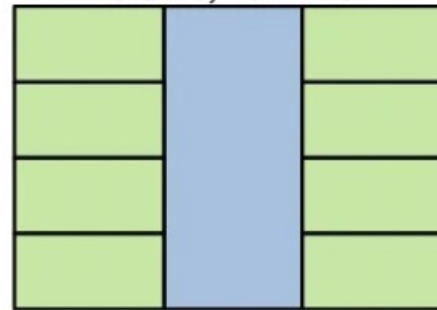
AMD Zen 3 CCX

“Zen 3” CCX Configs

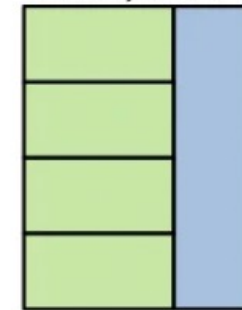
Performance-Desktop/Workstation/Server
8 Cores, 32MB L3



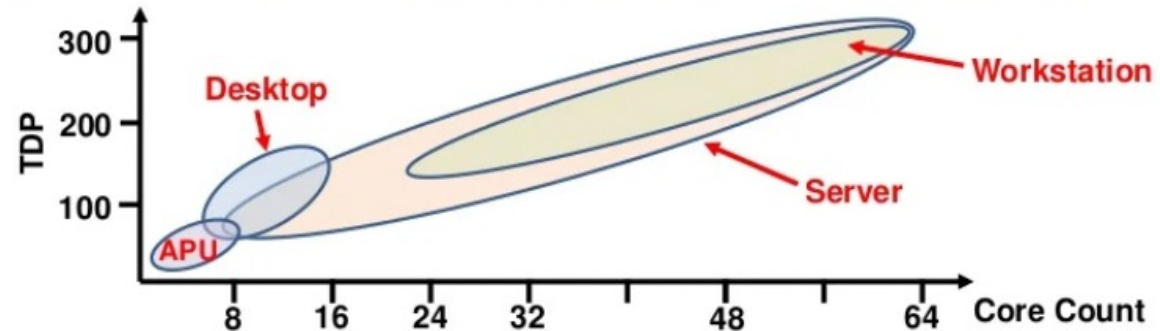
APU
8 Core, 16MB L3



Value
4 Core, 8MB L3



- Multiple CCXs can be placed to create up to 64-core products
- “Zen 3” CPU core can span from approx. 15-320W TDP



© 2022 IEEE International Solid-State Circuits Conference

2.7: “Zen 3”: AMD 2nd Generation 7nm x86-64 Microprocessor Core

9 of 32

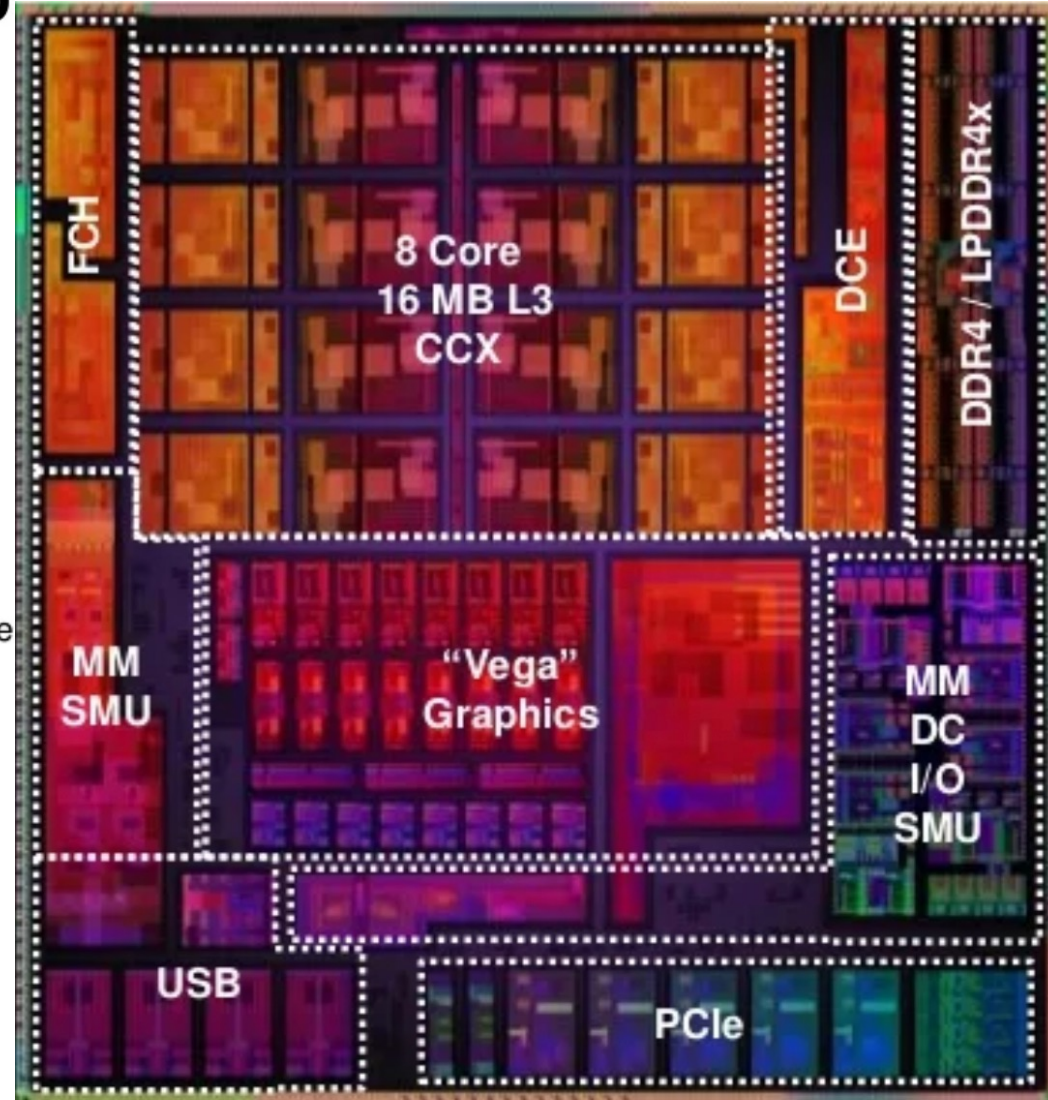
**“Zen 3”: AMD 2nd Generation 7nm x86-64
Microprocessor Core**

AMD APU

APU Monolithic Chip

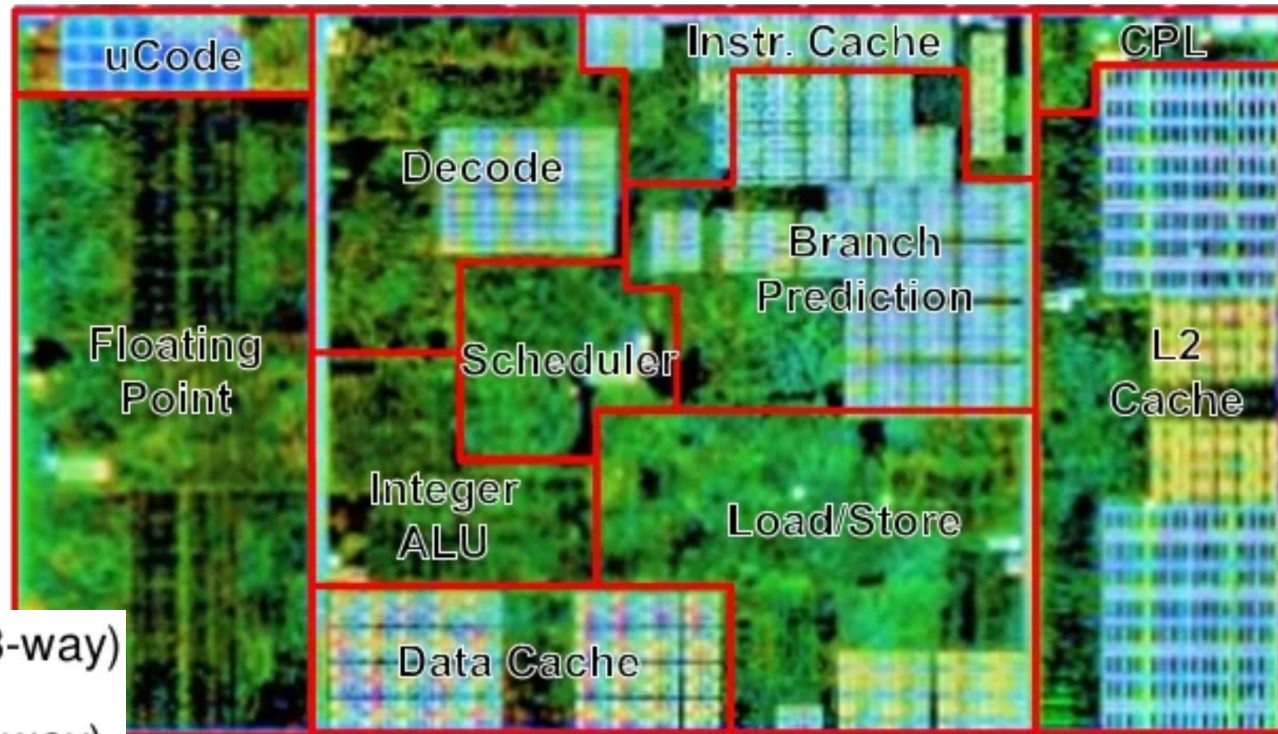
Key chip features

- 8 core, 16MB L3 CCX
- 8 compute-unit (“Vega”) graphics
- 2 memory controllers
 - DDR4 up to 3200 MT/s
 - LPDDR4x up to 4266 MT/s
- Multimedia (MM) engines
 - 2nd Gen Video Codec¹, 3rd Gen Audio ACP
- 2nd Gen display controller (DC)
- I/O controllers
 - PCIe® Gen4, USB-C, USB-3.1, USB-2.0, NVMe
- System management unit (SMU)
- Fusion controller hub (FCH)



180 mm² 10.7B Transistors (7nm)

Core Functional Units



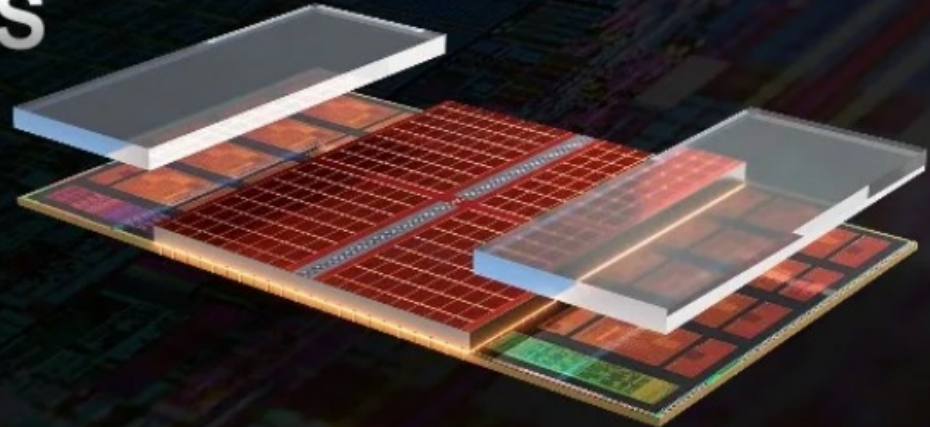
- 512kB private L2 cache (8-way)
- 32kB instruction cache (8-way)
- 32kB data cache (8-way)
- Comprised of ~30 sub-blocks
- Chip pervasive logic (CPL)
 - Clock/test unit



AMD 3D V-Cache



HETEROGENEOUS INTEGRATION WITH 3D PACKAGING

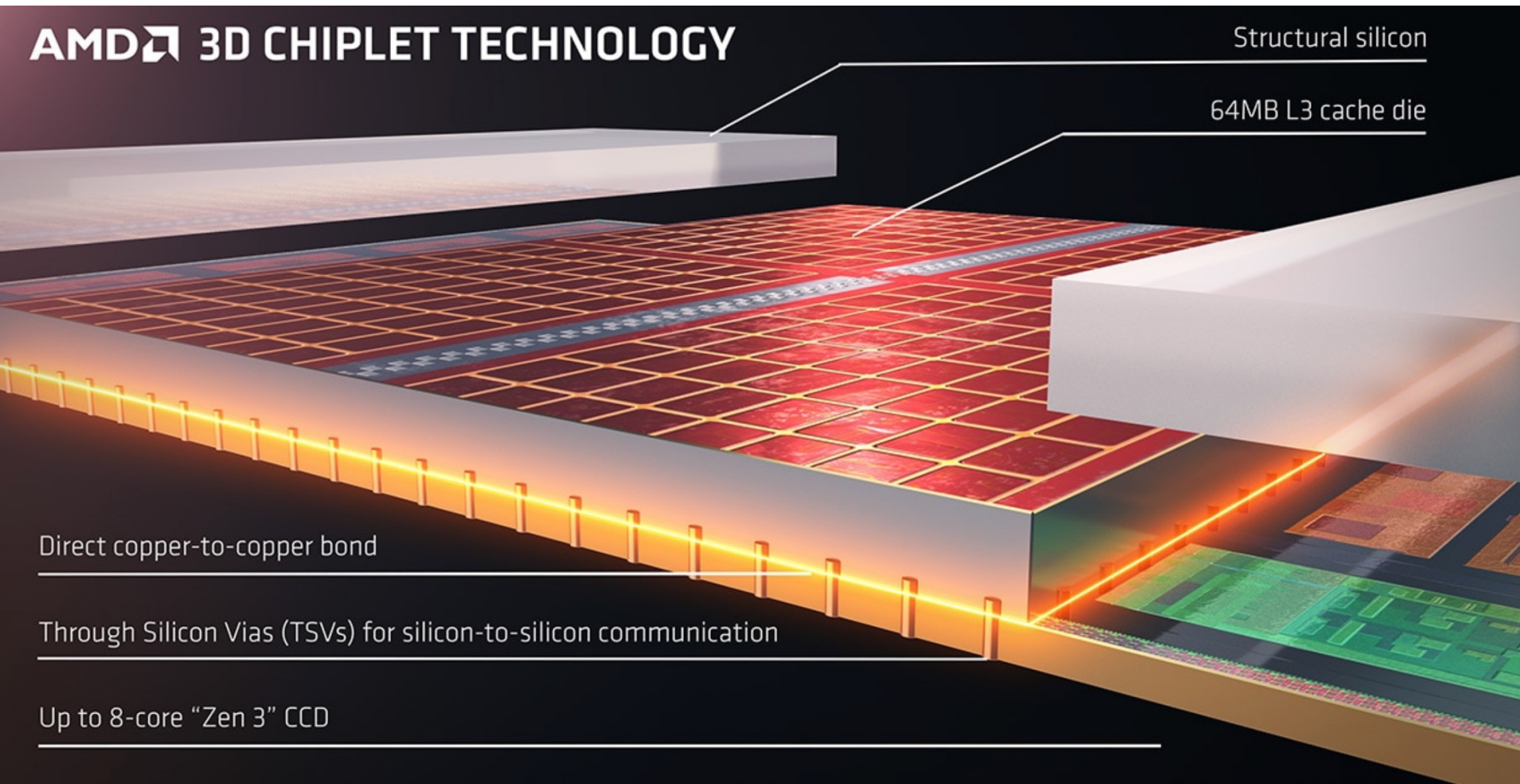


Heterogeneous Integration with 3D Packaging

Explore an overview of AMD 3D V-Cache™ technology and the performance uplift this groundbreaking tech provides.

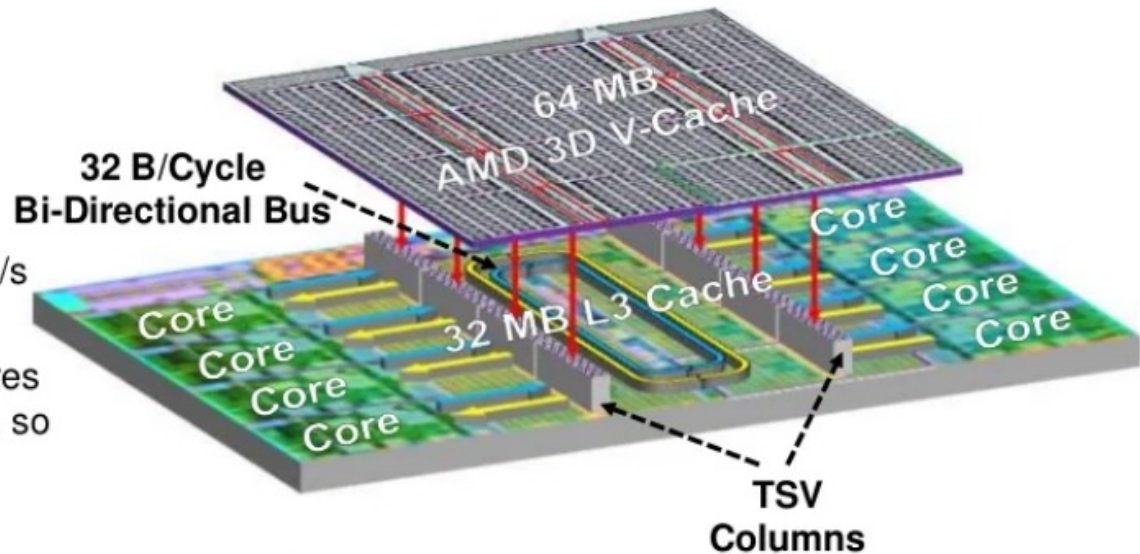
AMD 3D V-Cache

AMD 3D CHIPLET TECHNOLOGY



AMD 3D V-Cache Ready

- Two columns of TSVs on left/right side of the L3 cache
- AMD 3D V-Cache extends L3 Cache capacity by 64MB (3x)
- Total inter-die bandwidth: >2 TB/s
- All control and routing to the cores is implemented on the base die, so AMD 3D V-Cache can be completely focused on density



**“Zen 3”: AMD 2nd Generation 7nm x86-64
Microprocessor Core**

¹AMD, Fort Collins, CO, ²AMD, Santa Clara, CA, ³AMD, Austin, TX



AMD Zen 3 uArch

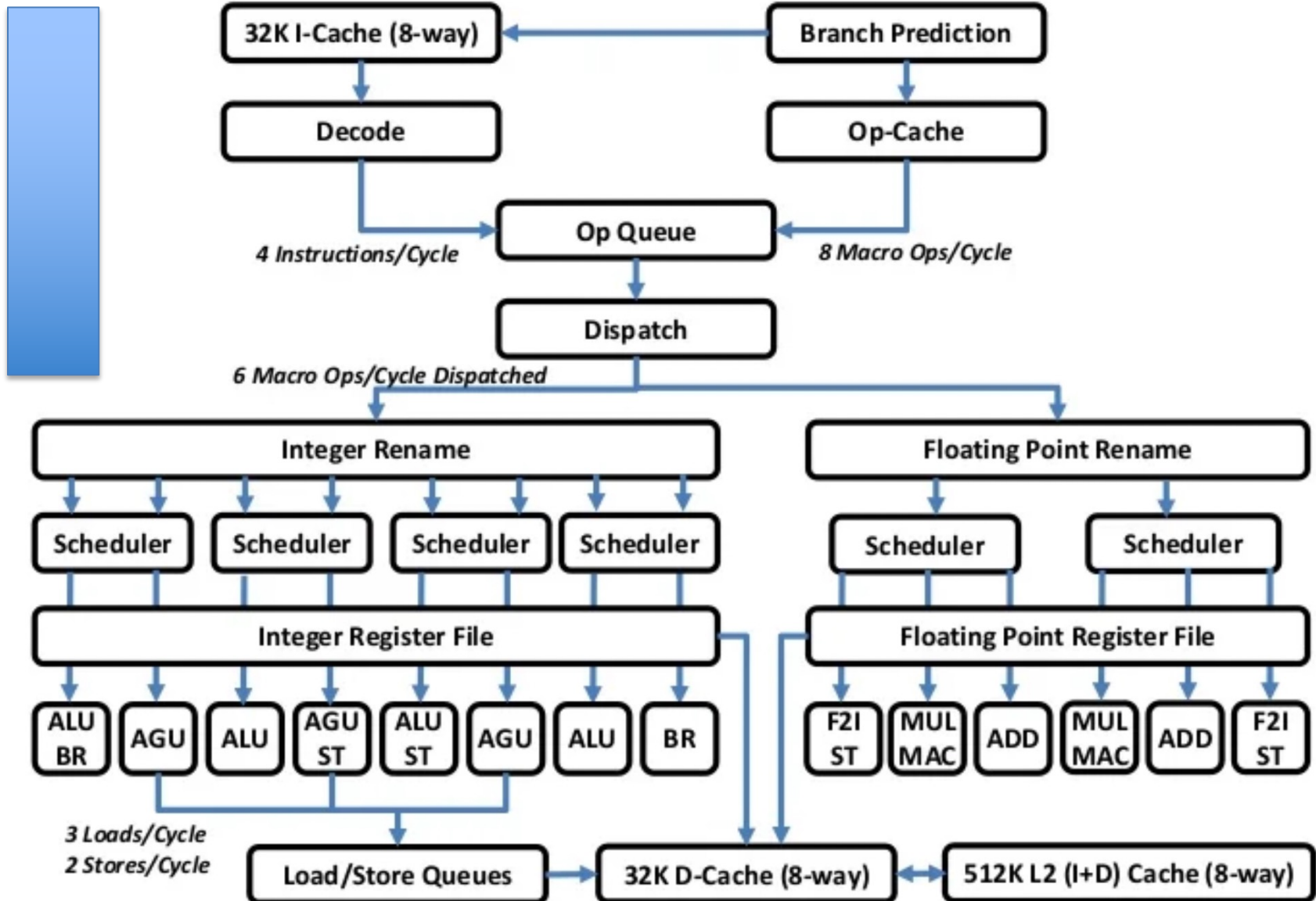
Major changes from “Zen 2”

- L1 BTB: 512→1024 entries
- Improved branch pred. bandwidth
- Int issue width: 7→10
- Reorder buffer: 224→256 entries
- FP issue width: 4→6
- FMAC latency: 5→4 cycles
- LD/ST bandwidth: 2/1→3/2
- TLB table walkers: 4→6





AMD Zen 3 uArch





AMD Benchmarks

© 2022 IEEE International Solid-State Circuits Conference

Server Performance

Benchmark	Config	"Zen 2"	"Zen 3"	Uplift
SPECint@2017	64 Cores [7763 vs. 7H12] ¹	717	854	+19%
	32 Cores [75F3 vs. 7532] ²	444	596	+34%
SPECfp@2017	64 Cores [7763 vs. 7H12] ³	543	651	+20%
	32 Cores [75F3 vs. 7532] ²	434	546	+26%
SPECjbb@2017	64 Cores [7763 vs. 7H12] ⁴	249k	314k	+26%





AMD Benchmarks

Client Performance

Single-Thread

Benchmark	Segment	"Zen 2"	"Zen 3"	Uplift
Cinebench R20	Desktop [5950X vs. 3900XT] ^{1,5}	546 (4.7 GHz)	640 (4.9 GHz)	+17%
	Mobile [5800U vs. 5600U] ²	474 (4.4Ghz)	551 (4.6 GHz)	+16%

Multi-Thread

Benchmark	Config	"Zen 2"	"Zen 3"	Uplift
Cinebench R20	8 Cores [5800U vs. 4800U] ²	3218	3655	+14%
PCMark 10	8 Cores [5800U vs. 4800U] ³	5081	6074	+20%
PCMark Apps	8 Cores [5800U vs. 4800U] ³	8663	10663	+23%

Major Gaming Uplifts with "Zen 3": +26% on average

Benchmark Game	Config	Uplift
CS:GO™ (DirectX® 9)	12 Cores	+46%
PUBG™ (DirectX® 11)		+33%
DOTA™ (Vulkan®)	[3900XT vs. 5900X] ⁴	+24%
F1™ 2019 (DirectX® 12)		+24%
Battlefield™ V (DirectX® 12)		+5%

Benchmark Game	Config	Uplift
League of Legends™ (DirectX® 11)	12 Cores	+50%
Shadow of the Tomb Raider™ (DirectX® 12)		+28%
Far Cry™ New Dawn (DirectX® 11)	[3900XT vs. 5900X] ⁴	+22%
Ashes of the Singularity™ (Vulkan®)		+19%
Total War™ : Three Kingdoms (DirectX® 11)		+6%



Tools & SDKs

AMD Zen Software Studio

AMD Optimizing C/C++ and Fortran Compilers (“AOCC”)

— The AOCC compiler system is a high performance, production software generation tool optimized for AMD processors based on the AMD “Zen” core architecture.

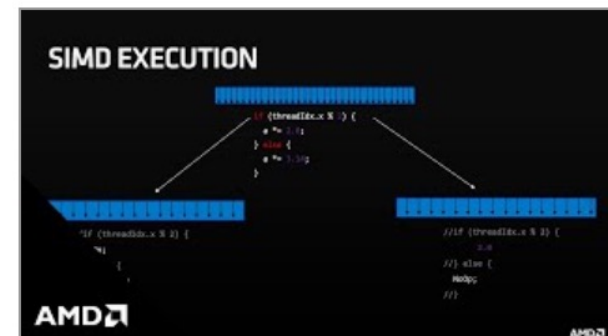
AMD μ Prof — AMD μ Prof is a suite of powerful tools that help developers optimize software for performance and power, optimized for AMD processors based on the AMD “Zen” core architecture.

AMD Optimizing CPU Libraries (“AOCL”) — AOCL is a set of numerical libraries optimized for AMD processors based on the AMD “Zen” core architecture.

Other SDKs and Tools

Tools for DMTF DASH — DASH (Desktop and mobile Architecture for System Hardware) is a client management standard released by the DMTF (Distributed Management Task Force). DASH is a web services based standard for secure out-of-band and remote management of desktops and mobile systems. Client systems that support out-of-band management help IT administrators perform tasks independent of the power state of the machine or the state of the operating system.

AMD Ryzen™ Master Monitoring SDK — The AMD Ryzen™ Master Monitoring SDK is a public distribution



Porting CUDA to HIP

In the final video of the series, presenter Nicholas Malaya...

GPU Programming Software

In this video, presenter Damon McDougall summarizes the various Compilers,...

GPU Programming Concepts (Part 3)

In this video, presenter Noel Chalmers concludes the discussion on...



Section

AMD RADN3 GPU



AMD RDNA3

12-4-22

AMD INSTINCT™ MI200 SERIES



AMD INSTINCT™
MI200 OAM
MI250, MI250X



AMD INSTINCT™
MI210 PCIe®
COMING SOON



AMD RDNA3

12-4-22

AMD
EPYC | INSTINCT



AMD RDNA3

12-4-22

AVAILABLE
DECEMBER 13, 2022



AMD Radeon™ RX 7900 XT

RDNA™ 3 Compute Units	84
GDDR6 Memory	20 GB
Memory Bus	320 bit
Shader Flops	52 TFLOPS
Display Port	2.1

\$899 SEP

AMD Radeon™ RX 7900 XTX

RDNA™ 3 Compute Units	96
GDDR6 Memory	24 GB
Memory Bus	384 bit
Shader Flops	61 TFLOPS
Display Port	2.1

\$999 SEP

Nvidia RTX 4080

SM Count	76
GDDR6 Memory	16 GB
Memory Bus	256 bit
Shader Flops	49 TFLOPS
Display Port	1.4

\$1,199 SEP

AMD RDNA3

12-4-22

LEADERSHIP DESIGN

EXQUISITE DESIGN THROUGH RELENTLESS ATTENTION TO DETAIL

NEW GPU thermal interface material for high performance and maximum reliability

Die-cast aluminum backplate that improves PCB rigidity

Ultra-soft MOSFET thermal interface material

NEW Dispensed GDDR6 thermal interface material

Premium die-cast aluminum shroud

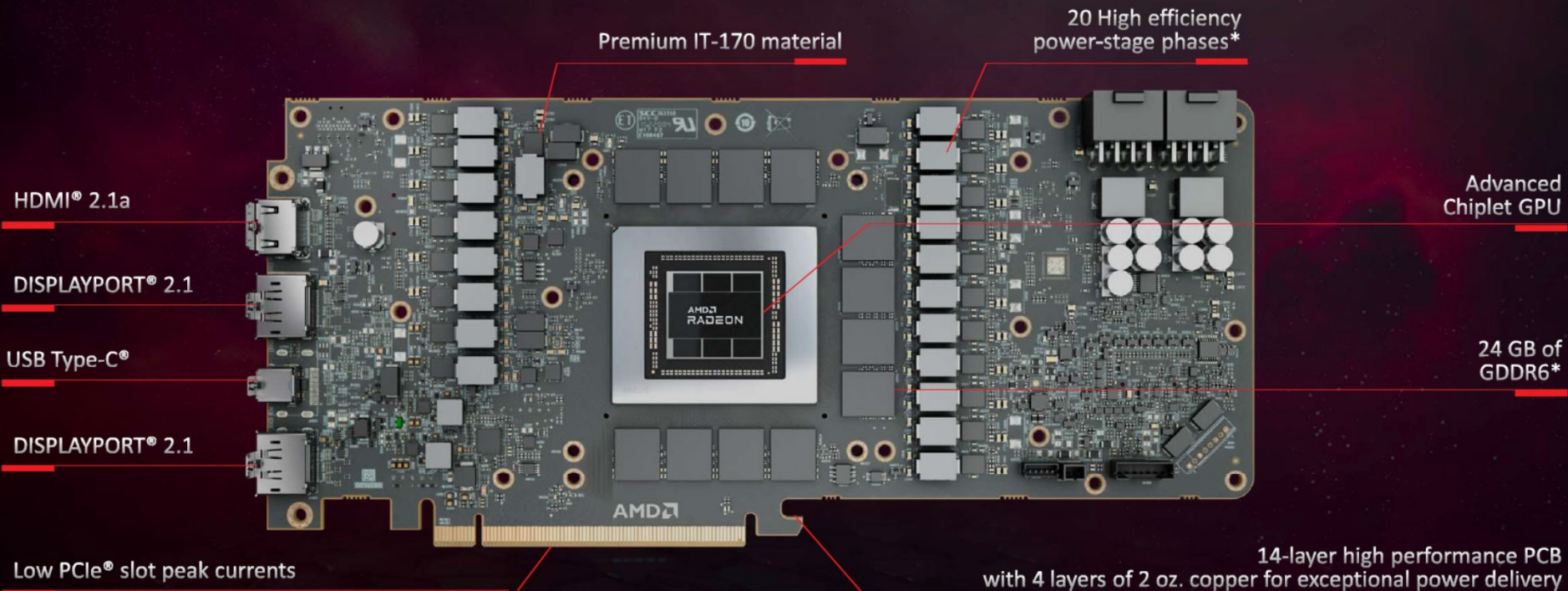
10% Larger vapor chamber optimized for maximum performance*

AMD RDNA3

12-4-22

INTELLIGENT ENGINEERING

A REFINED DESIGN THAT DELIVERS INCREDIBLE PERFORMANCE & EFFICIENCY



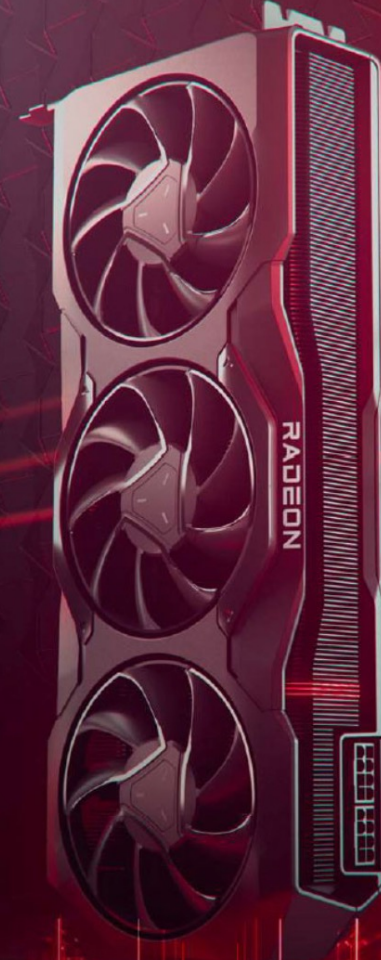
AMD RDNA3

12-4-22

BREAKTHROUGH PERFORMANCE

RADEON™ RX 7900 SERIES

	AMD RADEON™ RX 7900 XT	AMD RADEON™ RX 7900 XTX	NVIDIA RTX 4080
AMD RDNA™ 3 COMPUTE UNITS	84	96	-
STREAM PROCESSORS	5376	6144	-
GAME CLOCK	2.0 GHZ	2.3 GHZ	-
BOOST CLOCK (UP TO)	2.4 GHZ	2.5 GHZ	-
2 ND GENERATION INFINITY CACHE™	80 MB	96 MB	-
GDDR6 MEMORY	20 GB	24 GB	16 GB
MEMORY BUS WIDTH	320 BIT	384 BIT	256 BIT
TOTAL BOARD POWER	300 W	355 W	320 W

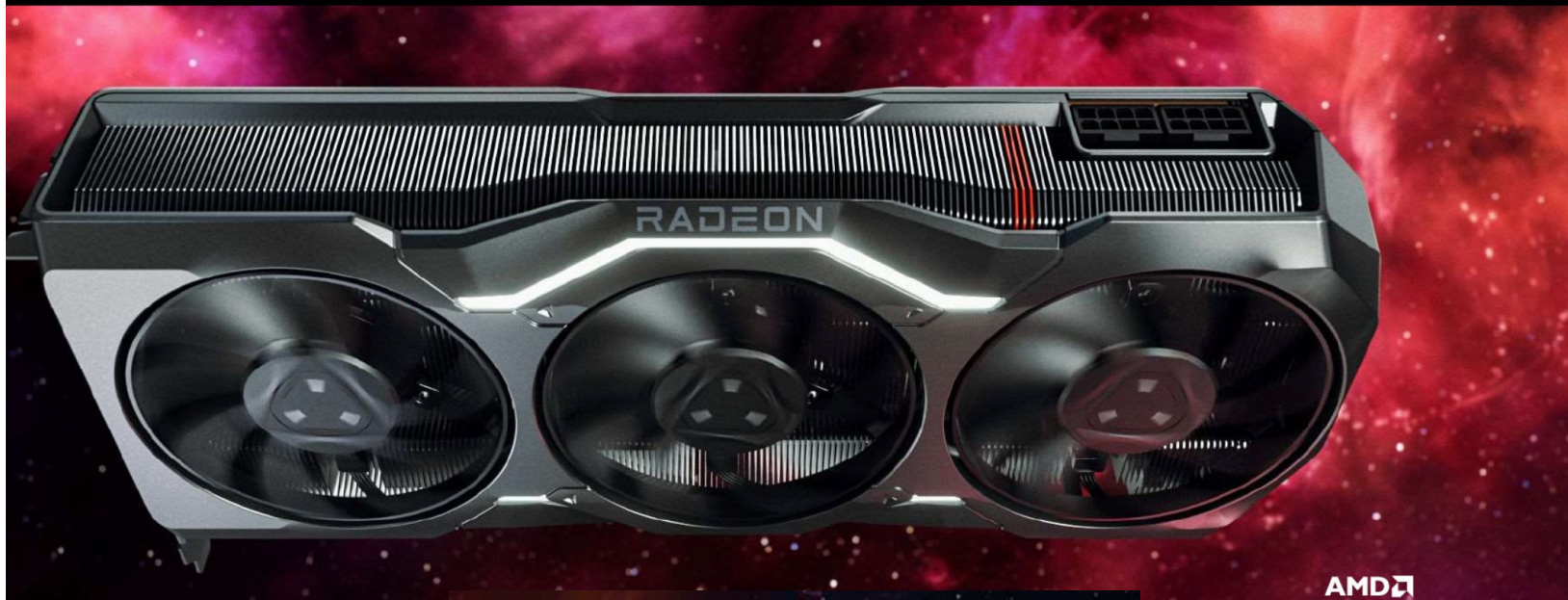


AMD RDNA3

Embargoed Until Nov. 14, 2022, at 9am ET

12-4-22

INTRODUCING **AMD RADEON™ RX 7900 SERIES** THE MOST ADVANCED GRAPHICS FOR GAMERS & CREATORS



RX 7900 Series Deep Dive

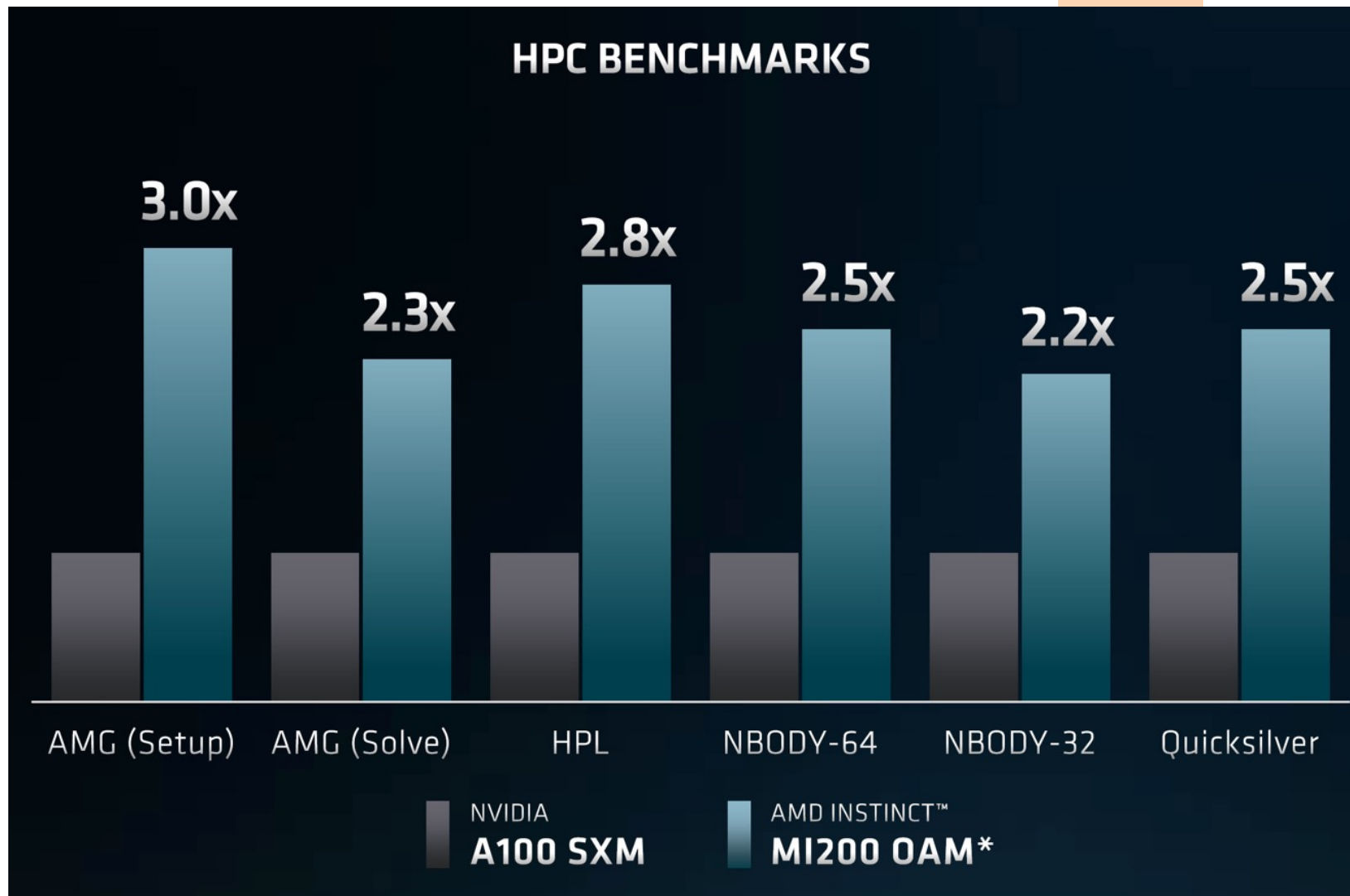
Scott Olschewsky

Director of Product Management



AMD vs Nvidia A100

12-4-22



SHATTERING PERFORMANCE BARRIERS IN HPC & AI

PEAK PERFORMANCE	A100	MI200*	INSTINCT™ ADVANTAGE
FP64 VECTOR	9.7 TF	47.9 TF	4.9X
FP32 VECTOR	19.5 TF	47.9 TF	2.5X
FP64 MATRIX	19.5 TF	95.7 TF	4.9X
FP32 MATRIX	N/A	95.7 TF	N/A
FP16, BF16 MATRIX	312 TF	383 TF	1.2X
MEMORY SIZE	80 GB	128 GB	1.6X
MEMORY BANDWIDTH	2.0 TB/s	3.2 TB/s	1.6X



AMD and Nvidia GPU Specifications

Graphics Card	RX 7900 XTX	RX 7900 XT	RX 6950 XT	RTX 4090
Architecture	Navi 31	Navi 31	Navi 21	Ada Lovelace
Process Technology	TSMC N5 + N6	TSMC N5 + N6	TSMC N7	TSMC N5
Transistors (Billion)	58 (45.7 + 6x 2.05)	56 (45.7 + 5x 2.05)	26.8	76
Die size (mm ²)	300 + 222	300 + 185	519	608
CUs / SMs	96	84	80	128
SPs / Cores (Shaders)	6144 (12288)	5376 (10752)	5120	16384
Tensor / Matrix Cores	?	?	?	512



AMD RDNA3

12-4-22

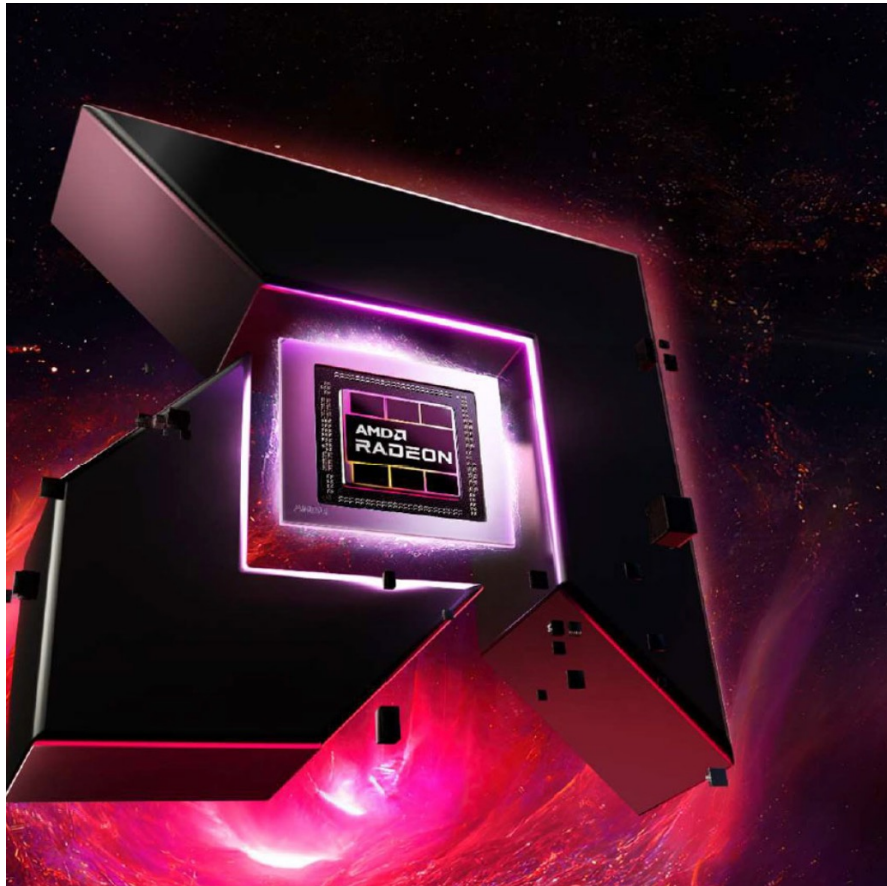
AMD and Nvidia GPU Specifications

Graphics Card	RX 7900 XTX	RX 7900 XT	RX 6950 XT
Ray Tracing "Cores"	96	84	80
Boost Clock (MHz)	2500	2400	2310
VRAM Speed (Gbps)	20	20	18
VRAM (GB)	24	20	16
VRAM Bus Width	384	320	256
L2 / Infinity Cache	96	80	128
ROPs	192	192	128



AMD RDNA3

12-4-22



AMD RDNA 3

UP TO
54%
MORE PERFORMANCE PER WATT

165%
MORE TRANSISTORS PER MM²

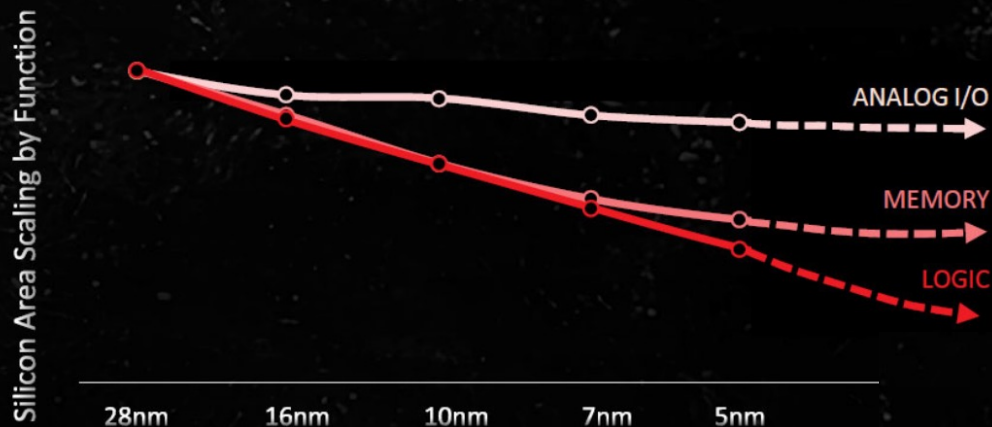
5.3 TB/S
WORLD'S FASTEST
INTERCONNECT

AMD RDNA3

CHIPLET TECHNOLOGY OUR MOTIVATION

12-4-22

Limited and Divergent Scale Factors



Density gains diminishing

Increasing Costs^[1]

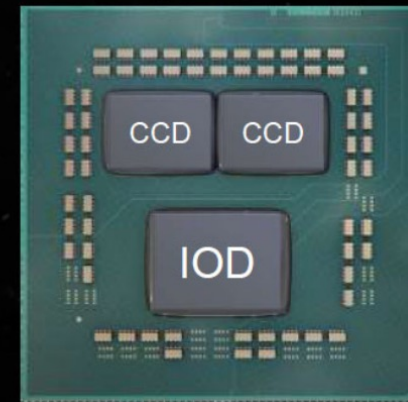
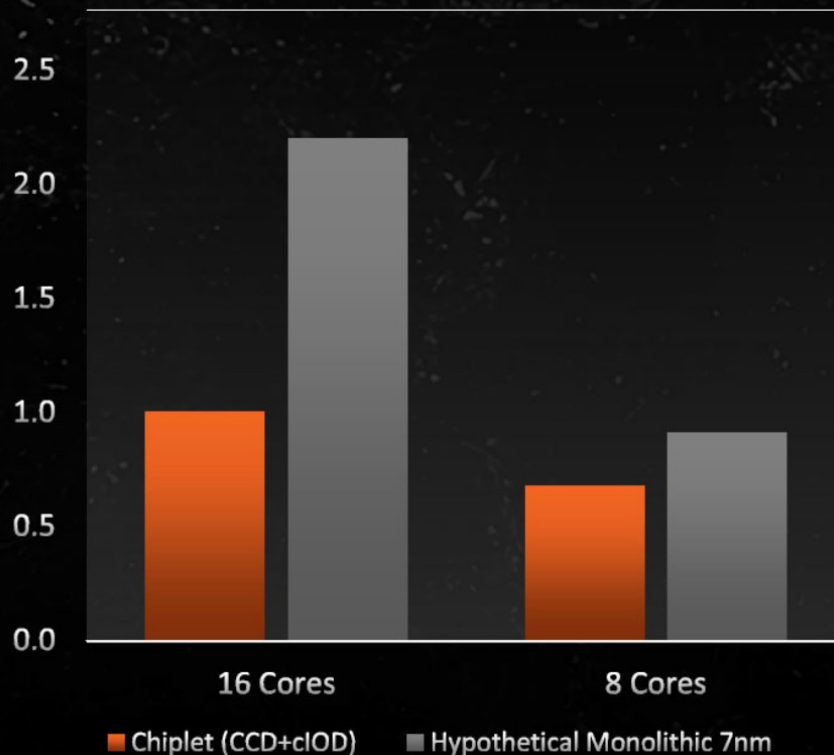


Costs Increasing

12-4-22

CHIPLET TECHNOLOGY BENEFIT FOR AMD RYZEN™

Normalized Die Cost [1]



Great cost benefit vs. monolithic
Linear cost with core count



AMD RDNA3

CHIPLET TECHNOLOGY

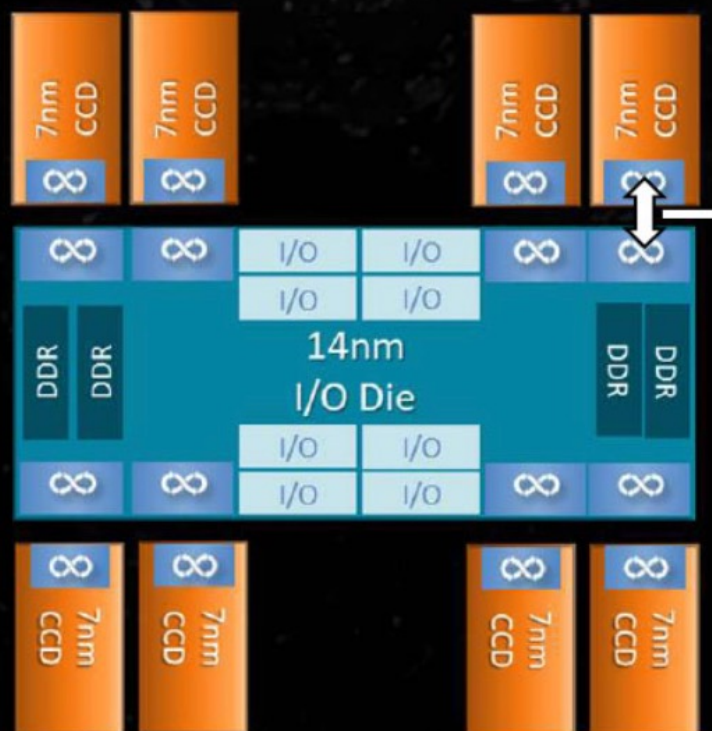
CAN IT WORK FOR GRAPHICS?

12-4-22

Traditional Monolithic



EPYC CPU Server



100's of signals

AMD RDNA3

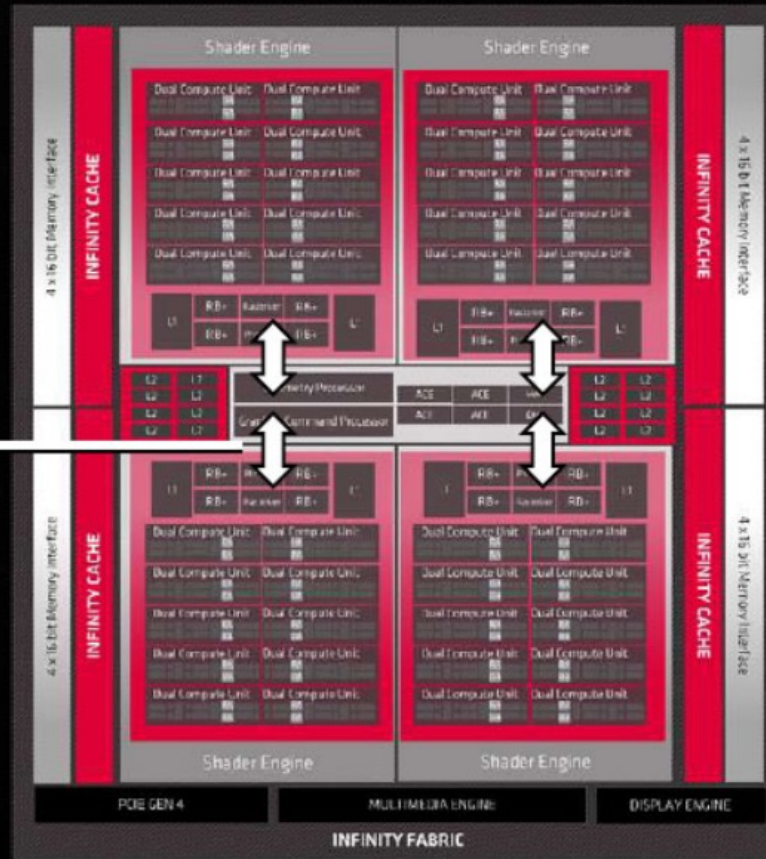
CHIPLET TECHNOLOGY

CAN IT WORK FOR GRAPHICS?

12-4-22

"Navi21" GPU

10's of
1000's of
signals

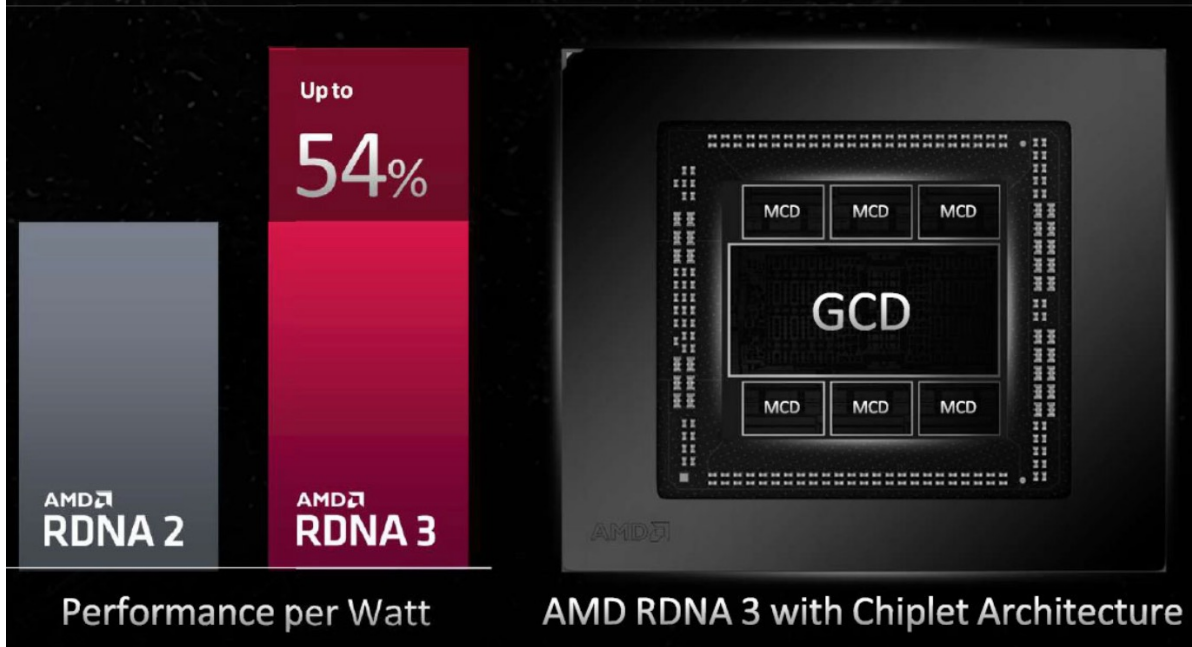


- GPU shader engines require massive amounts of connectivity compared to CPUs

AMD RDNA3

12-4-22

CHIPLET TECHNOLOGY SUMMARY



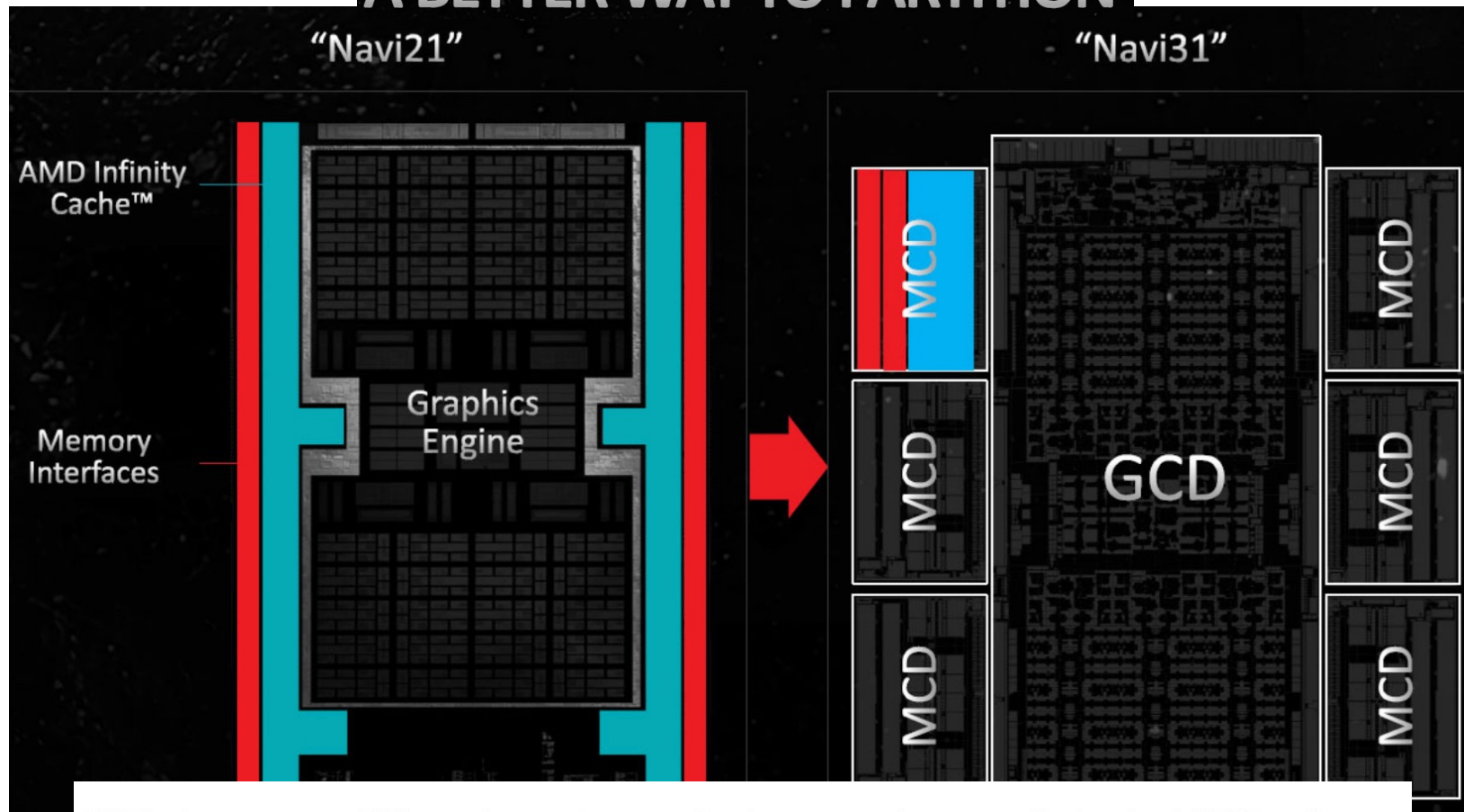
- Chiplet architecture with advanced packaging is the future,
- AMD has leveraged our leadership chiplet expertise to deliver the first chiplet-based gaming GPU
- Massive 5.3TB/s bandwidth with innovative Infinity Links on High Performance Fanout
- Negligible overheads in latency and power enable leadership performance/Watt

AMD RDNA3

CHIPLET TECHNOLOGY

A BETTER WAY TO PARTITION

12-4-22



GPUs have very different requirements. Large caches can help, but GPUs also really like having gobs of memory bandwidth to feed all the GPU cores. For example, even the beastly EPYC 9654 with a 12-channel DDR5 configuration 'only' delivers up to 460.8 GB/s of bandwidth. The fastest graphics cards like the RTX 4090 can easily double that.

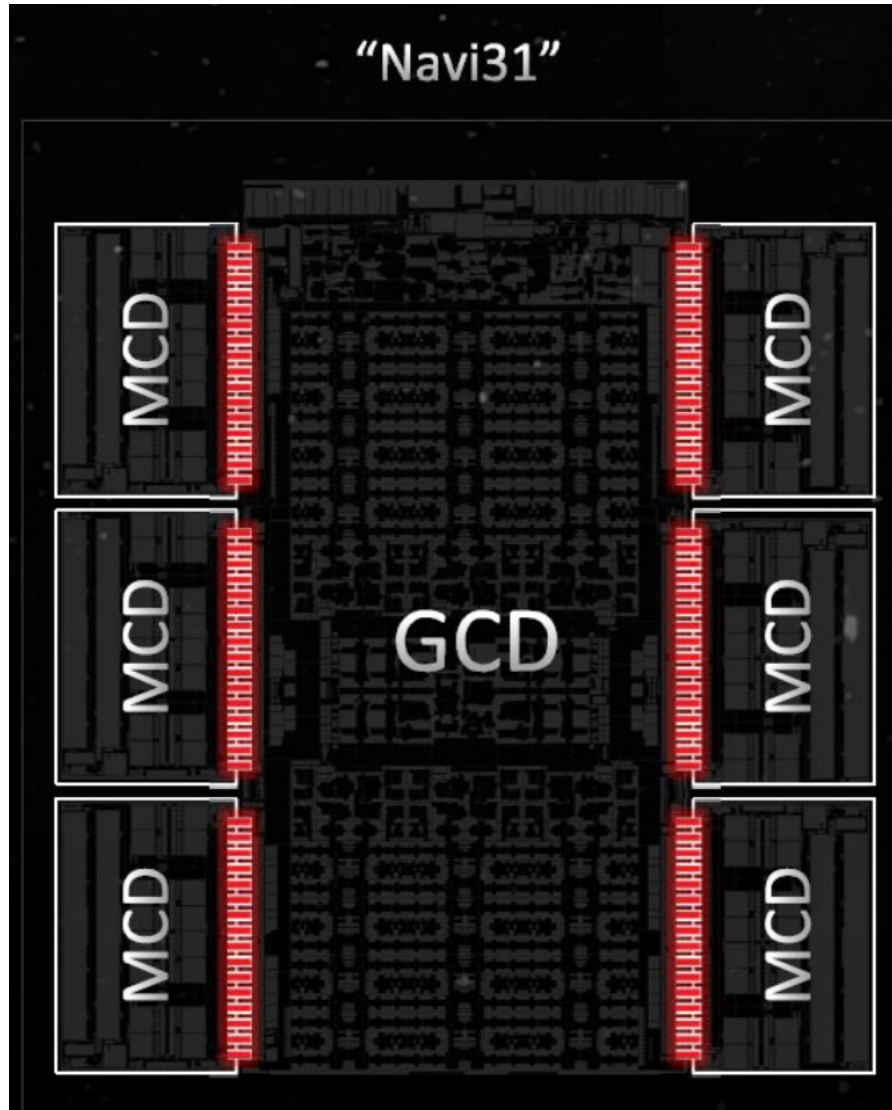
AMD RDNA3

12-4-22

CHIPLET TECHNOLOGY HOW TO CONNECT THE CHIPLETS?

Graphics
Compute
Die

Memory
Cache
Die





AMD RDNA 3 GPU Architecture Deep Dive: The Ryzen Moment for GPUs

By [Jarred Walton](#) published 20 days ago

12-4-22

Swimming with the next generation GPUs



Jarred Walton

Jarred Walton is a senior editor at Tom's Hardware focusing on everything GPU. He has been working as a tech journalist since 2004, writing for AnandTech, Maximum PC, and PC Gamer. From the first S3 Virge '3D decelerators' to today's GPUs, Jarred keeps up with all the latest graphics trends and is the one to ask about game performance.

One surprising piece of information is that the entire MI200 chip contains up to 58 billion transistors. That's certainly a lot, but Nvidia's A100 by comparison has 54.2 billion transistors in a single GPU core. Unless we've missed something, that means the total size of the MI200 chips are roughly the same size as Nvidia's A100, except they potentially pack a lot more compute performance into that area.

12-4-22



Jarred Walton

Jarred Walton
GPU. He has
AnandTech, I
decelerators'
trends and is

Fundamentally, MI200 appears to use an updated and enhanced version of the GPU that powered the MI100 — AMD calls the architecture CDNA2, similar to the RDNA2 vs. RDNA shift on the consumer side. MI100 had 120 CDNA CUs (compute units) and 7680 streaming processors. MI100 used TSMC's N7 fabrication node, and also supported up to 32GB of HBM2 memory clocked at 1.2 Gbps. MI200 takes the ball and runs with it, boosting all of the key performance metrics.

(Image credit: AMD)

The first major change relative to the MI100 comes in the use of a multi-die package. This is basically taking the same chiplet approach that AMD used in its recent Zen 2 and Zen 3 CPUs and applying that to GPUs, though with some enhancements. The two CDNA dies (that's "Compute DNA", as opposed to the graphics-focused RDNA used in consumer GPUs) are linked together via an Infinity Fabric, with 25 Gbps links providing up to 100 GBps of bi-directional bandwidth between the GPUs. There are eight available links in the MI200 OAM (OCP Accelerator Module, where OCP is "Open Compute Platform") configuration, yielding 800 GBps of bandwidth between the two chiplets.



Jarred Walton

Jarred Wal
GPU. He h
AnandTech
decelerati
trends and

Navi 31 consists of two core pieces, the Graphics Compute Die (GCD) and the Memory Cache Dies (MCDs). There are similarities to what AMD has done with its Zen 2/3/4 CPUs, but everything has been adapted to fit the needs of the graphics world.

For Zen 2 and later CPUs, AMD uses an Input/Output Die (IOD) that connects to system memory and provides all of the necessary functionality for things like the PCIe Express interface, USB ports, and more recently (Zen 4) graphics and video functionality. The IOD then connects to one or more Core Compute Dies (CCDs) — alternatively "Core Complex Dies," depending on the day of the week) via AMD's Infinity Fabric, and the CCDs contain the CPU cores, cache, and other elements.

A key point in the design is that typical general computing algorithms — the stuff that runs on the CPU cores — will mostly fit within the various L1/L2/L3 caches. Modern CPUs up through Zen 4 only have two 64-bit memory channels for system RAM (though [EPYC Genoa server processors](#) can have up to twelve DDR5 channels).



Jarred Walton

Jarred Wal
 GPU. He h
 AnandTech
 deceleratc
 trends and

Putting that into numbers, MI100 was the first GPU to provide over 10 TFLOPS of FP64 vector compute. With its higher clocks, dual-GPUs, and doubled FP64 rates, the MI200 has a peak FP64 vector rate of 47.9 TFLOPS — AMD was quick to point out that this represents a 4.9X increase over the Nvidia A100 FP64 vector rates.

MI200 also adds FP64 matrix support, with a peak rate that's double the vector unit rate: 95.7 TFLOPS. Again, by way of comparison, the Nvidia A100 FP64 vector performance is 19.5 TFLOPS. That's on paper, of course, so we need to see how that translates into the real world. AMD claims performance is around three times as fast as the A100 in several workloads, though it's difficult to say if that will be the case across all workloads.

On the FP16 side of things, the performance isn't quite as high. Nvidia's A100 has 312 TFLOPS of FP16/BF16 compute, compared to 383 TFLOPS for the MI200, but Nvidia also has sparsity. Basically, sparsity allows the GPU to skip some operations, specifically multiplication by zero (which, so my math teacher taught me, is always zero). Sparsity can potentially double the compute performance of the A100, so there should be some use cases where Nvidia maintains the lead.

AMD RDNA3

12-4-22



Jarred Walton

Jarred Walton is a senior GPU. He has been working on graphics for over 10 years.

Graphics
Compute
Die

The GCD houses all the Compute Units (CUs) along with other core functionality like video codec hardware, display interfaces, and the PCIe connection. The Navi 31 GCD has up to 96 CUs, which is where the typical graphics processing occurs. But it also has an Infinity Fabric along the top and bottom edges (linked via some sort of bus to the rest of the chip) that then connects to the MCDs.

Memory
Cache
Die

The MCDs, as the name implies (Memory Cache Dies) primarily contain the large L3 cache blocks (Infinity Cache), plus the physical GDDR6 memory interface. They also need to contain Infinity Fabric links to connect to the GCD, which you can see in the die shot along the center facing edge of the MCDs.

GCD will use TSMC's N5 node, and will pack 45.7 billion transistors into a 300mm² die. The MCDs meanwhile are built on TSMC's N6 node, each packing 2.05 billion transistors on a chip that's only 37mm² in size. Cache and external interfaces are some of the elements of modern processors that scale the worst, and we can see that overall the GCD averages 152.3 million transistors per mm², while the MCDs only average 55.4 million transistors per mm².

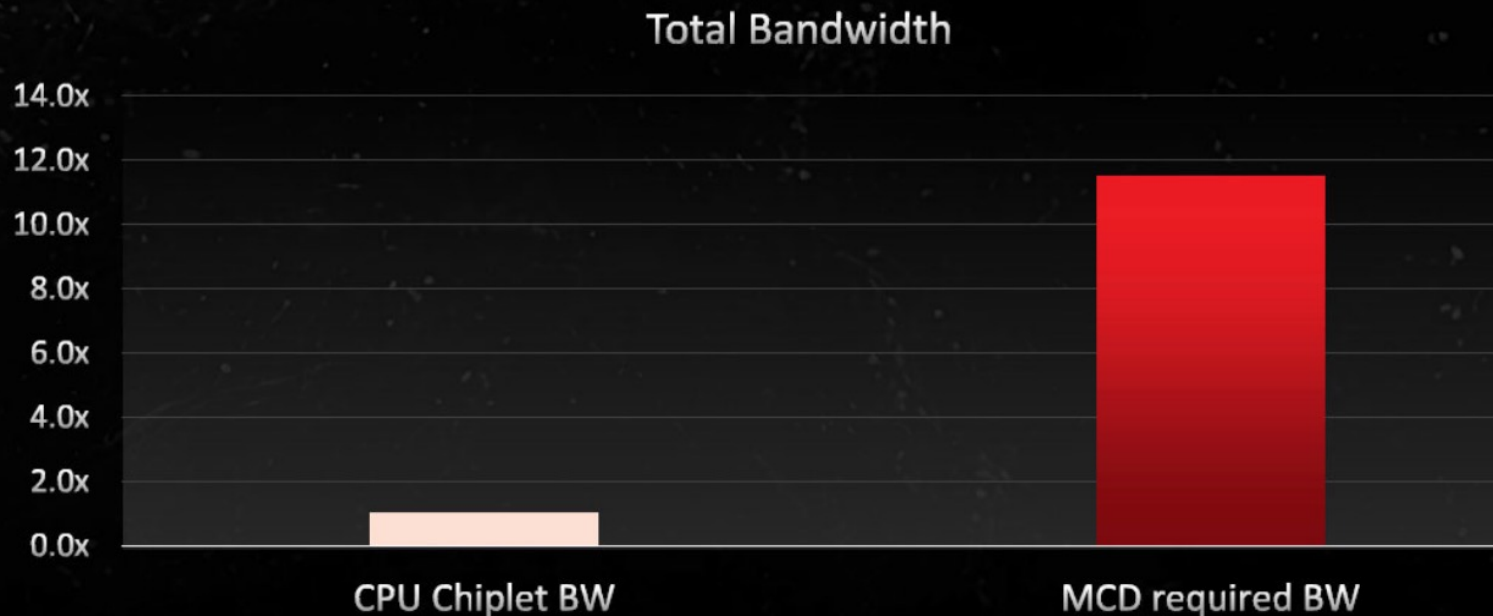


AMD RDNA3

12-4-22

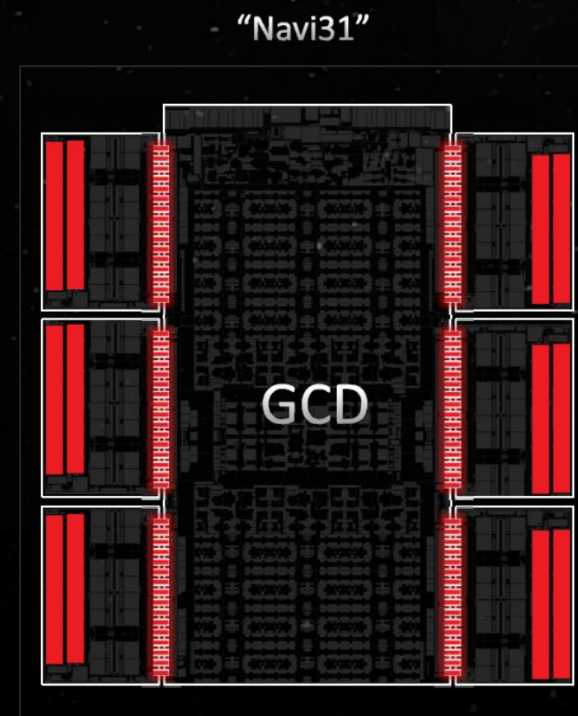
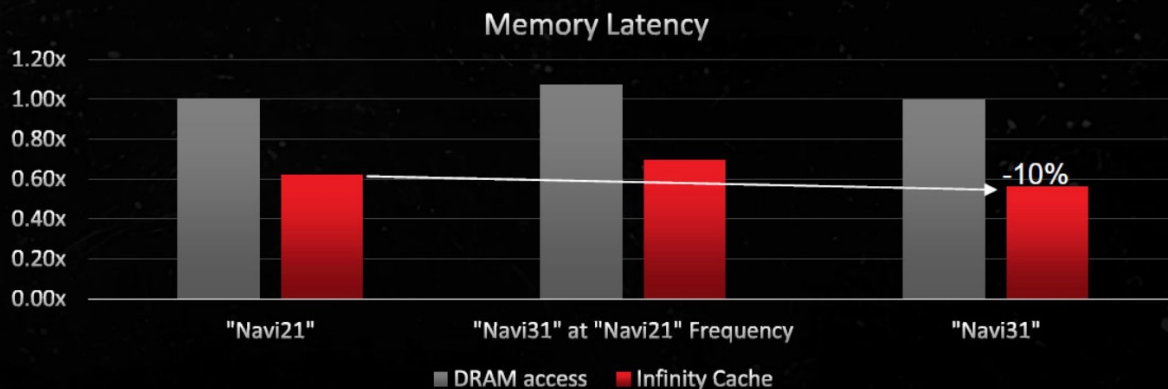
CHIPLET TECHNOLOGY HOW TO CONNECT THE CHIPLETS?

- GCD-MCD partitioning is great, but the bandwidth requirements are still extremely high
- Over 10X what a CPU CCD requires in EPYC
- Breakthrough Advanced packaging and a new interface is required:
 - **High Performance Fanout and Infinity links**



CHIPLET TECHNOLOGY INFINITY FANOUT LINKS MEMORY LATENCY

- The Infinity Link chiplet interfaces costs a modest amount of latency vs. on-die
- We eliminate this latency with higher clock rates
 - Base Infinity Fabric clock by +43%
 - Gfx game clock +18%
- The common case of Infinity cache hit is ~10% lower latency on "Navi31"

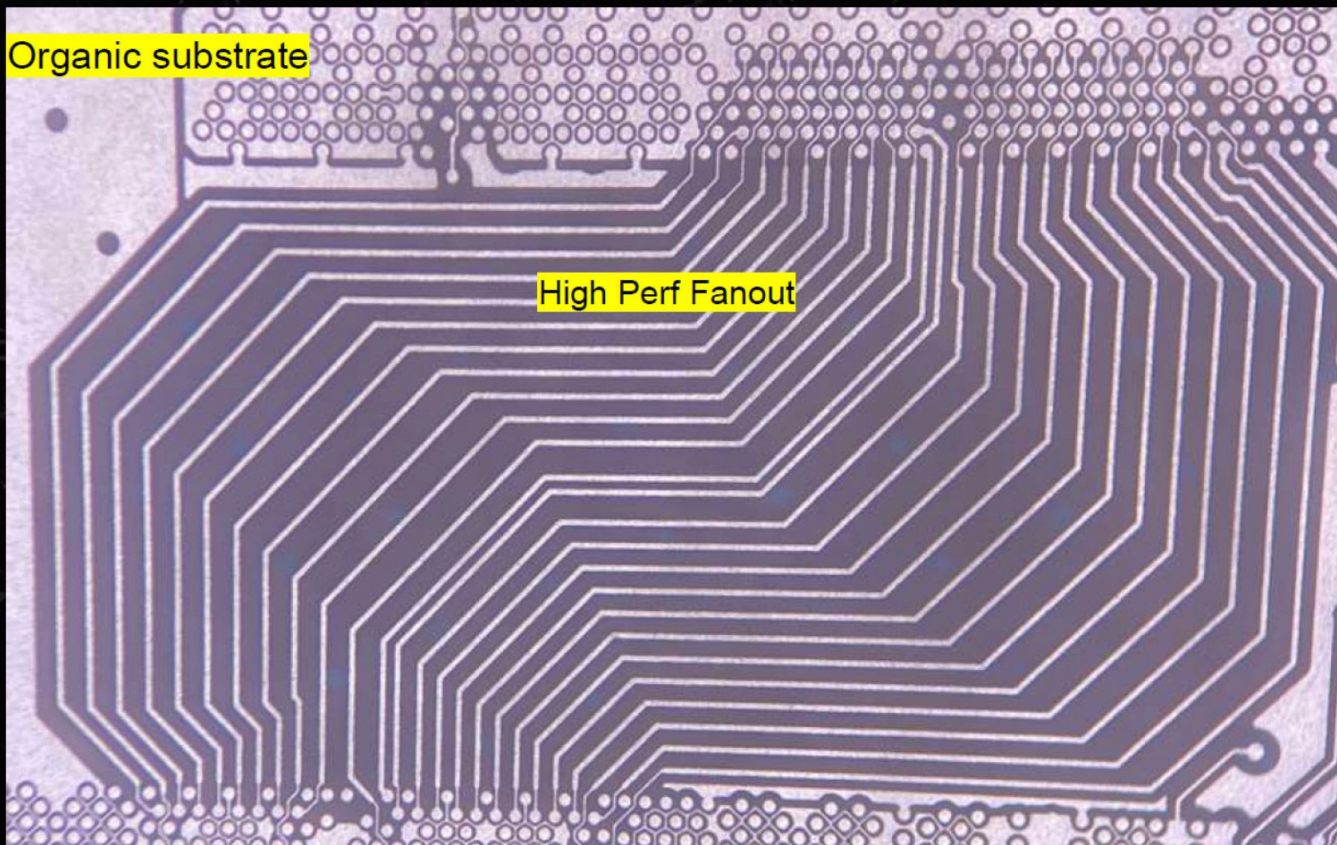


12-4-22

CHIPLET TECHNOLOGY

HIGH PERFORMANCE FANOUT

INTERCONNECT



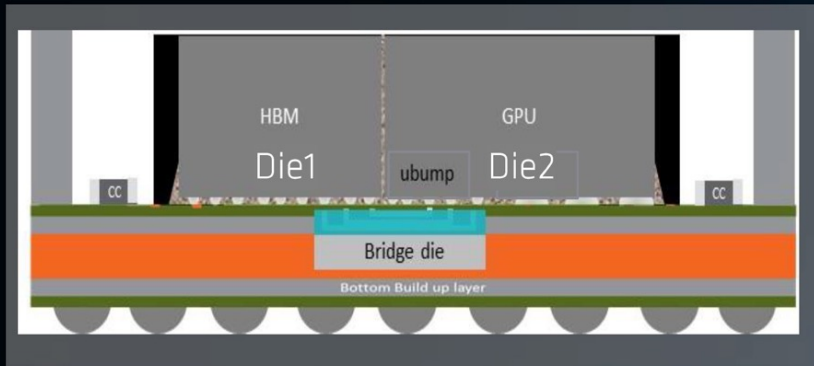
25 wires on organic substrate
compared to 50 wires on High
Performance Fanout

Images approximately to scale

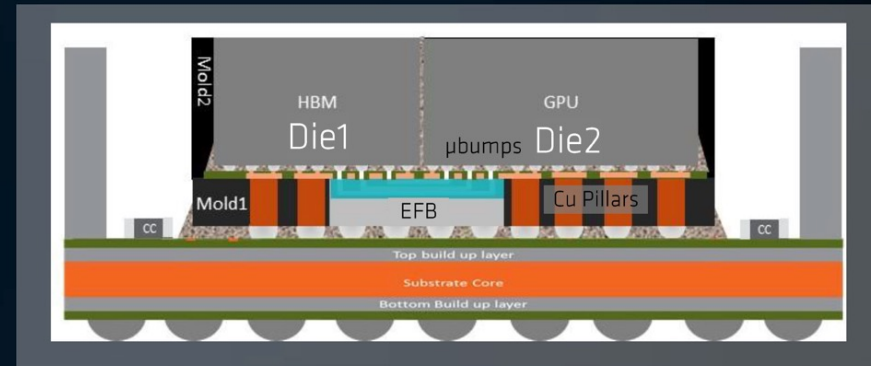
12-4-22

2.5D “BRIDGE” ARCHITECTURE LANDSCAPE

Substrate Embedded 2.5D



Elevated Fanout Bridge 2.5D





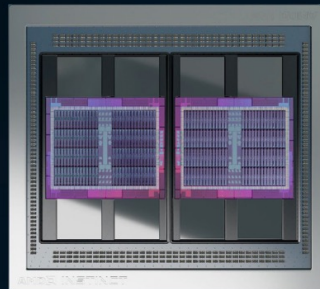
AMD RDNA3

12-4-22

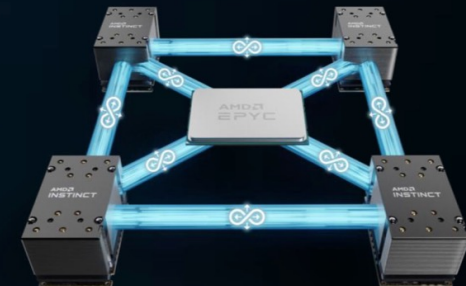
AMD INSTINCT™ MI200 SERIES

AMD
CDNA 2

WORKLOAD-OPTIMIZED
COMPUTE ARCHITECTURE



FIRST MULTI-DIE GPU



3RD GEN AMD INFINITY
ARCHITECTURE

AMD RDNA3

12-4-22



AMD INSTINCT™ MI200 SERIES WORLD'S MOST ADVANCED DATA CENTER ACCELERATOR

UP
TO **58B**

Transistors in 6nm

UP
TO **220**

Compute Units

UP
TO **880**

2nd Gen Matrix Cores

UP
TO **128**

GB HBM2E @ 3.2 TB/s

AMD RDNA3

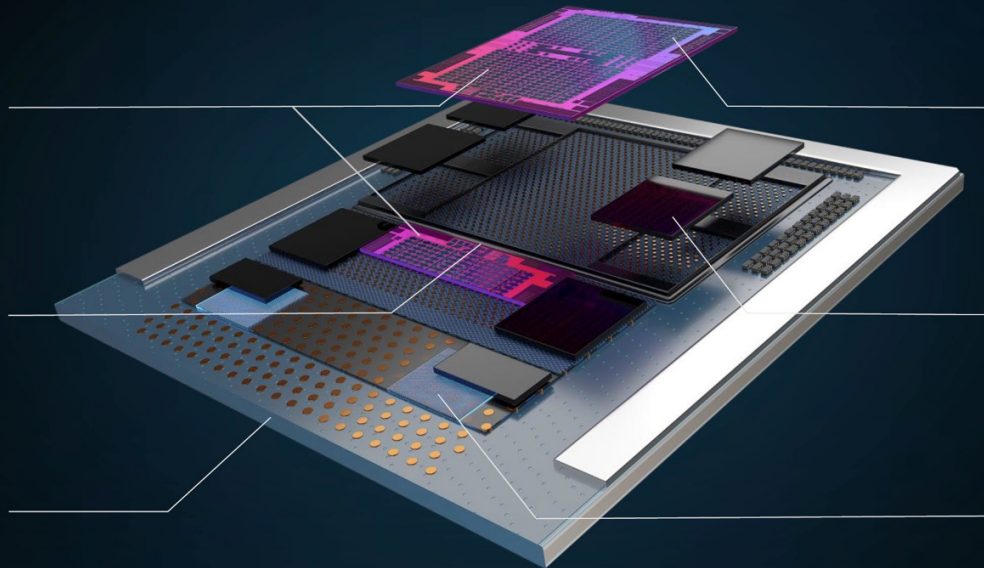
12-4-22

AMD INSTINCT™ MI200 SERIES KEY INNOVATIONS

TWO
AMD CDNA™2 DIES

ULTRA HIGH BANDWIDTH
DIE INTERCONNECT

COHERENT CPU-TO-GPU
INTERCONNECT



2ND GEN MATRIX
CORES FOR HPC & AI

EIGHT STACKS
OF HBM2E

2.5D ELEVATED
FANOUT BRIDGE (EFB)

AMD INSTINCT™ MI200 OAM SERIES



AMD RDNA3

12-4-22

AMD's High Performance Fanout Interconnect

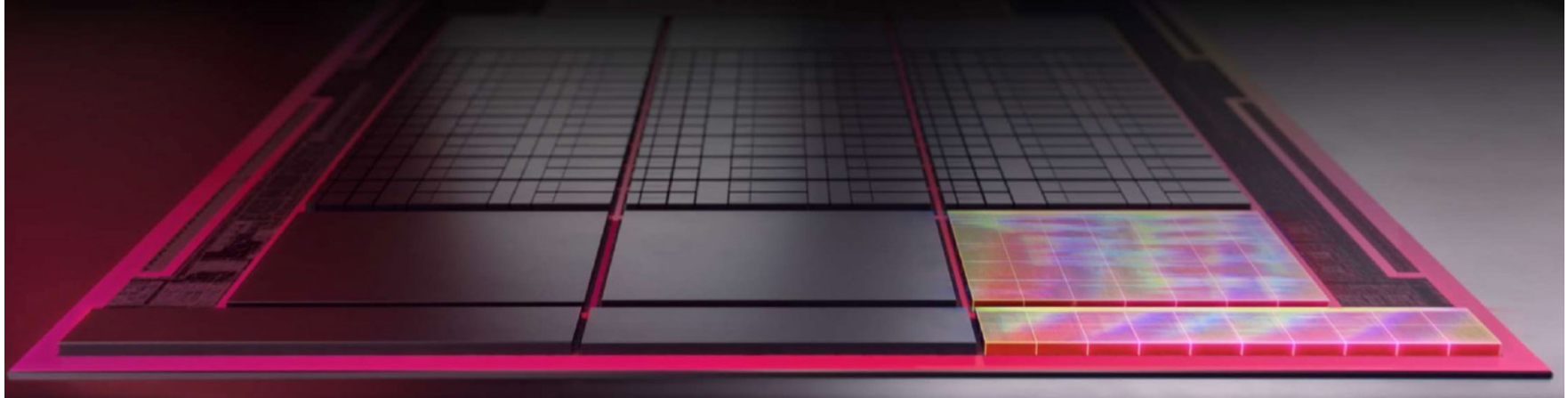
Interconnect	Picojoules per Bit (pJ/b)	
On-die	0.1	
Foveros	0.2	TSMC
EMIB	0.3	Intel
UCIe	0.25-0.5	
Infinity Fabric (Navi 31)	0.4	AMD
TSMC CoWoS	0.56	TSMC
Bunch of Wires (BoW)	0.5-0.7	
Infinity Fabric (Zen 4)	???	AMD
NVLink-C2C	1.3	Nvidia
Infinity Fabric (Zen 3)	1.5 (?)	AMD



AMD RDNA3

12-4-22

INTRODUCING **AMD RADIANCE DISPLAY™ ENGINE** MAXIMUM FIDELITY WITH AMD RADEON™ RX 7900 SERIES GRAPHICS



FULL REC2020
COVERAGE

12-BIT HDR
COLOR

UP TO
68 BILLION COLORS

AMD RDNA3

12-4-22

HIGH PERFORMANCE RENDERING

A SMARTER APPROACH TO REALTIME LIGHTING AND SHADOWS

96FPS

FSR ON

ASSASSIN'S
CREED
VALHALLA

AMD
RDNA 3

AMD
Software
Adrenalin Edition

AMD
FidelityFX



AMD RDNA3

12-4-22

TRUE 8K GAMING

NEXT-GENERATION EXPERIENCES ONLY AMD RADEON™ CAN DELIVER

8K Max Settings, FSR

RADEON™
RX 7900 XTX

73 fps

Uncharted: Legacy of Thieves Collection
(Ultra + FSR)

96 fps

Assassin's Creed Valhalla
(Very High + FSR)

149 fps

Death Stranding
(Very High + FSR + Ultrawide)

190 fps

Call of Duty: Modern Warfare 2
(Extreme + FSR + Ultrawide)

12-4-22

THREE WAYS TO OVERCLOCK FOR BEGINNERS TO ENTHUSIASTS

PERFORMANCE TUNING PRESETS

One click to turn on
Rage Mode

PRE-SET TGP AND FAN

PERFORMANCE TUNING AUTOMATIC

Varying frequency
for simple overclocking

PRE-SET TGP AND FAN

PERFORMANCE TUNING MANUAL

Fully manual overclocking
for customization enthusiasts

FREEDOM TO CUSTOMIZE ALL VARIABLES

AMD
RADEON
Boost

ON 80FPS
OFF 60FPS

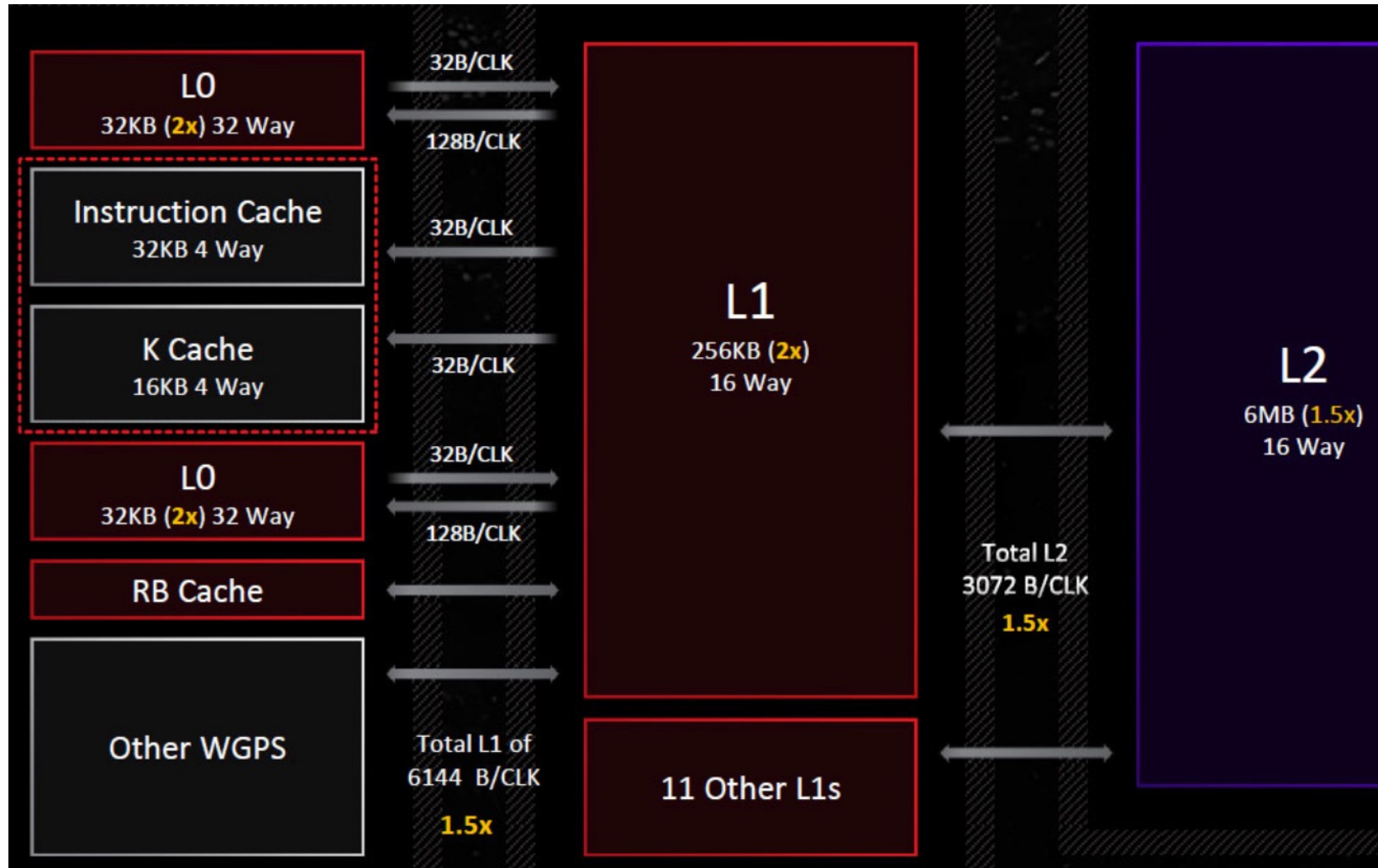
BORDERLANDS 3

AMD
RADEON
Anti-Lag

AMD RDNA3

OPTIMIZED AND BALANCED CACHE SYSTEM

12-4-22

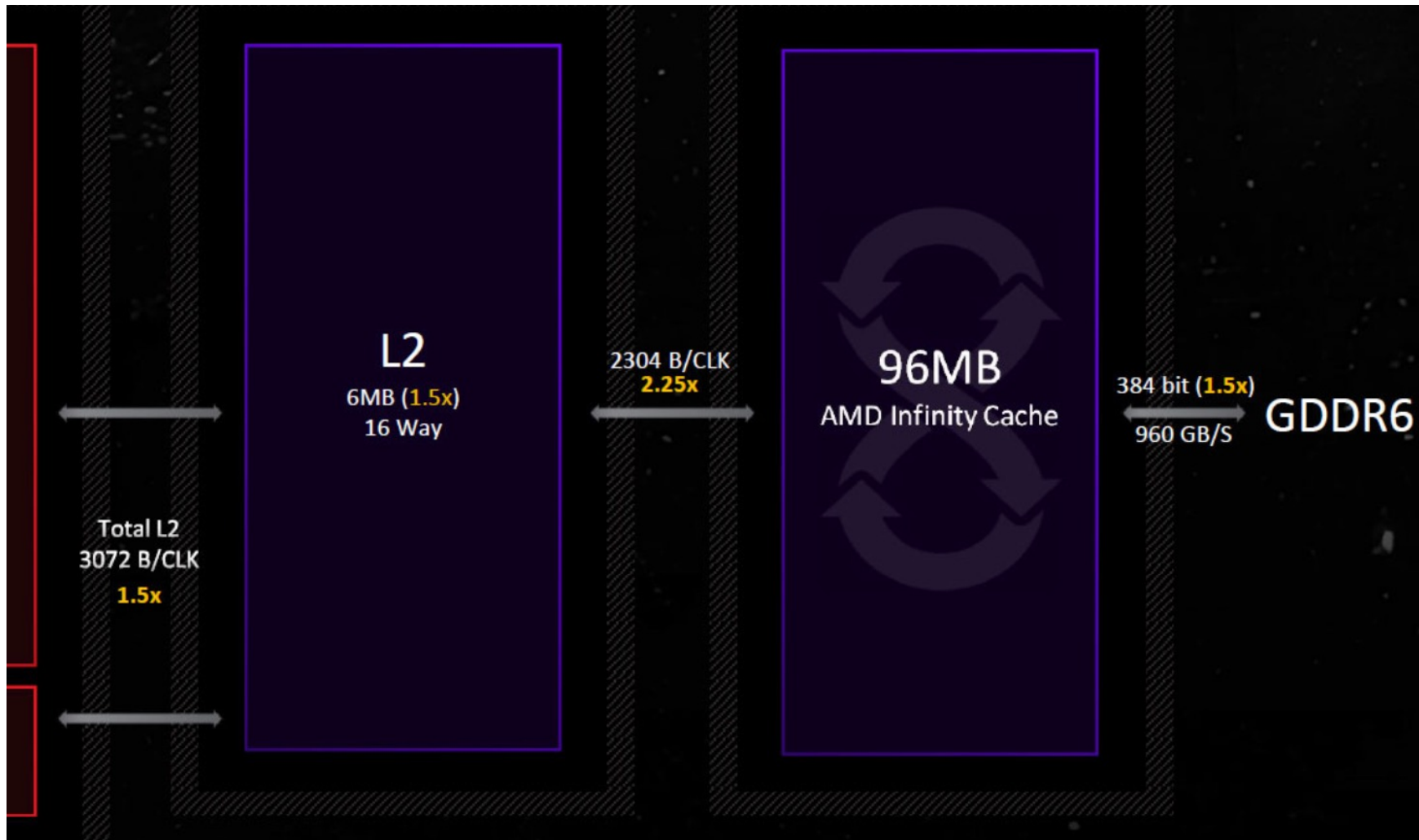




AMD RDNA3

OPTIMIZED AND BALANCED CACHE SYSTEM

12-4-22



AMD RDNA3

12-4-22

BREAKTHROUGH PERFORMANCE

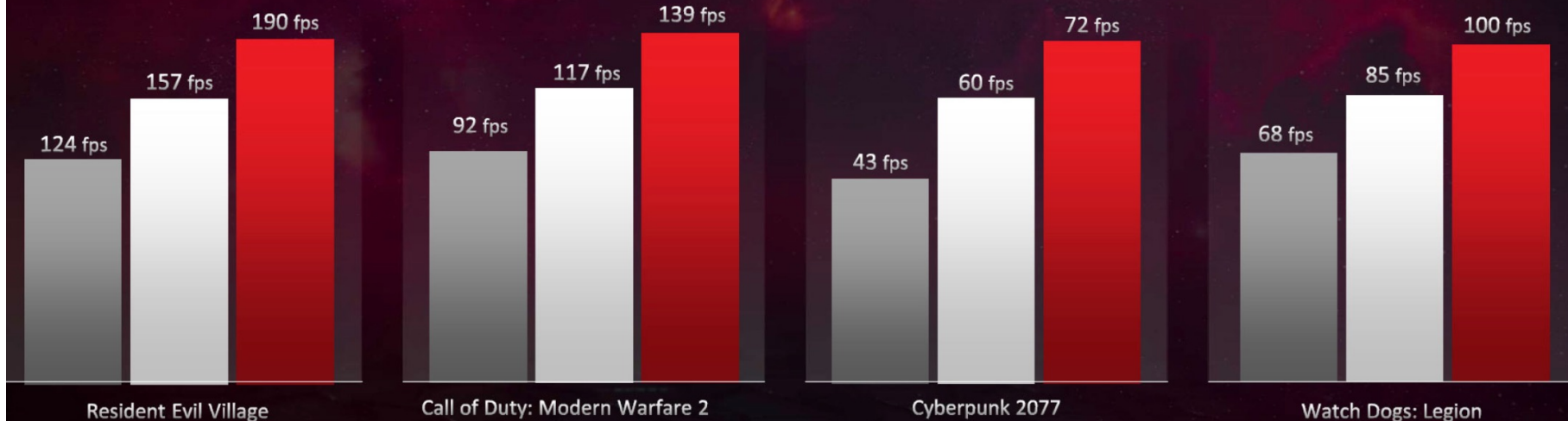
UP TO 67% MORE PERFORMANCE VS RX 6950 XT IN RASTERIZATION

4K Max Settings (Up To)

RADEON™
RX 6950 XT

RADEON™
RX 7900 XT

RADEON™
RX 7900 XTX





AMD RDNA3

12-4-22

BREAKTHROUGH PERFORMANCE

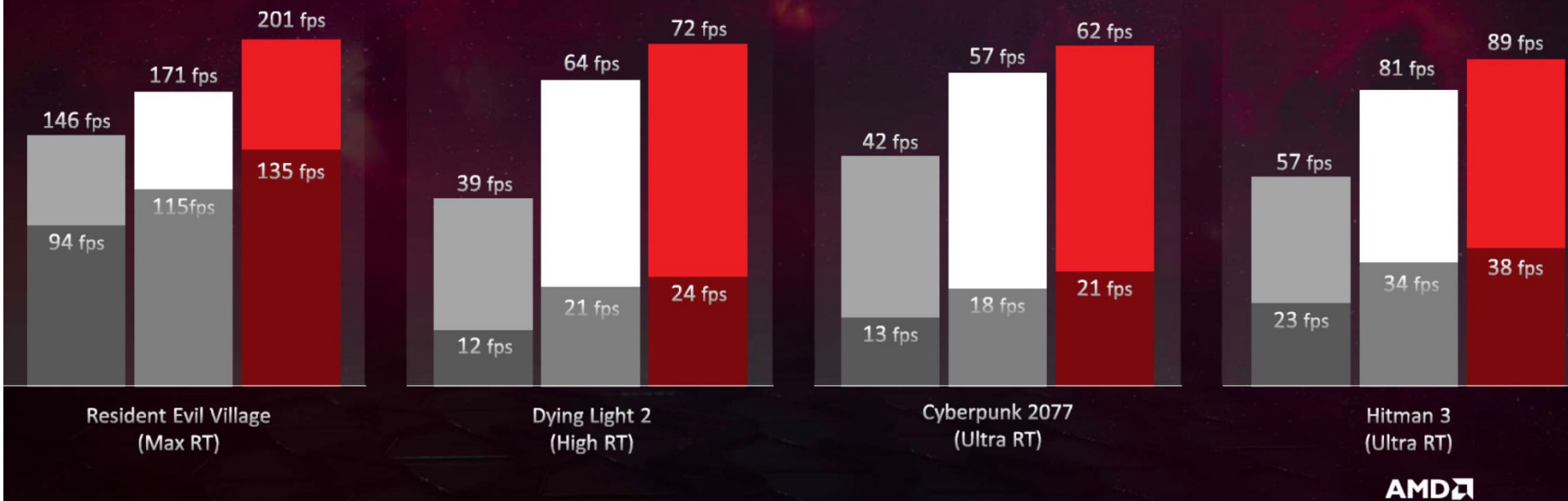
UP TO 82% MORE PERFORMANCE VS RX 6950 XT IN RAYTRACING

4K Max Settings (Up To), Raytracing, FSR On/Off

RADEON™
RX 6950 XT
FSR ON / OFF

RADEON™
RX 7900 XT
FSR ON / OFF

RADEON™
RX 7900 XTX
FSR ON / OFF

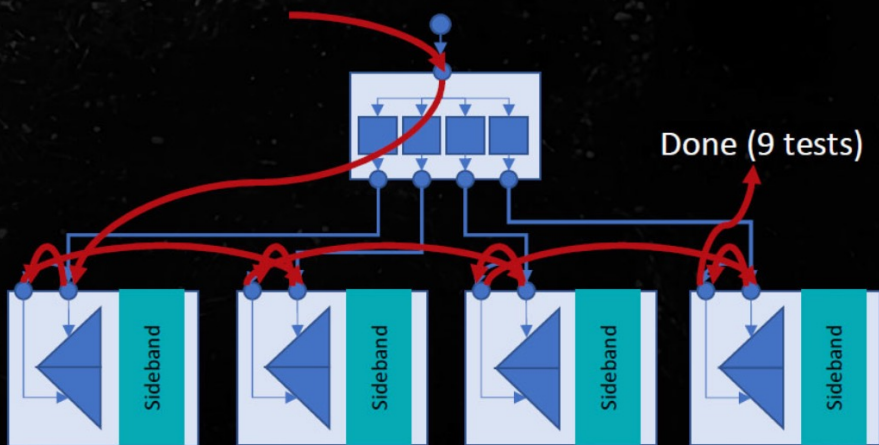


AMD RDNA™ 3 2ND GENERATION RAY TRACING

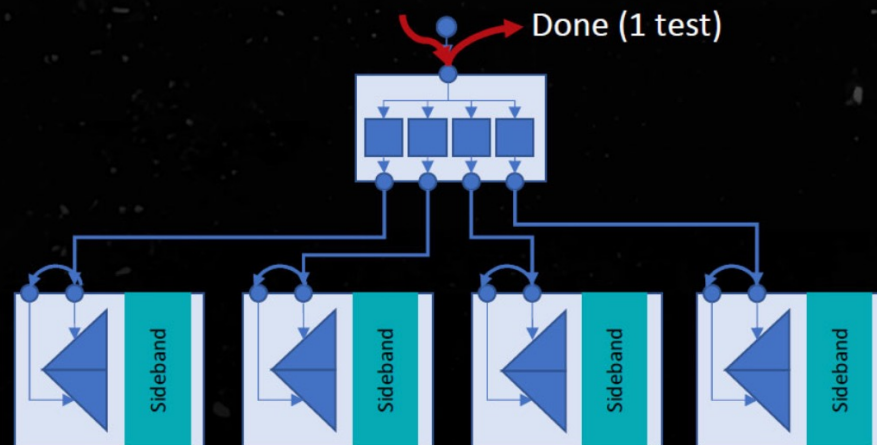
Hardware managed DXR Ray Flags

- Add Geometry Flags in BVH Nodes and Instance Ray Flags to the Node Pointer
- Hardware support for DXR Ray Flags reduces the required instruction count by ~15 per traversal loop iteration
- Use of DXR Ray Flags can improve performance due to triangle/subtree culling, and triangle/procedural geometry skipping by reducing the number of traversal iterations required

RDNA2 Late Culling

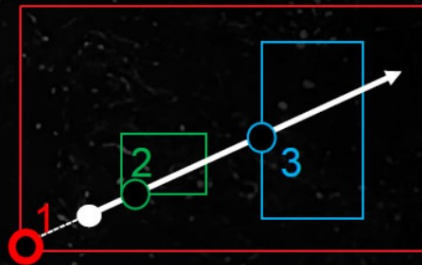


RDNA3 Early Subtree Culling

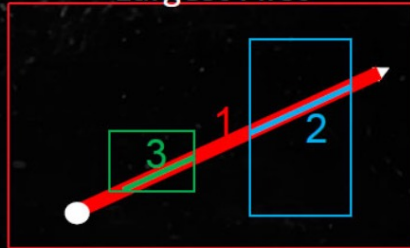


AMD RDNA™ 3 2ND GENERATION RAY TRACING

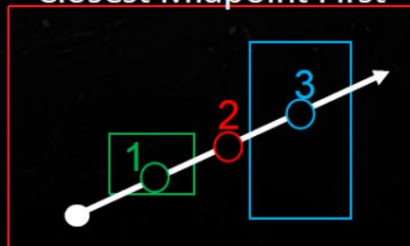
Closest First



Largest First



Closest Midpoint First



Extracting efficiency from each Ray

- Closest First

- Nodes are intersected in order of the closest intersection
- Good generic sorting heuristic

New hardware for specialized box sorting modes to improve performance by reducing traversal iterations for different ray types

- Largest First

- Nodes that have a larger overlap with the ray are intersected first
- Optimized for shadow rays / Terminate on First Hit

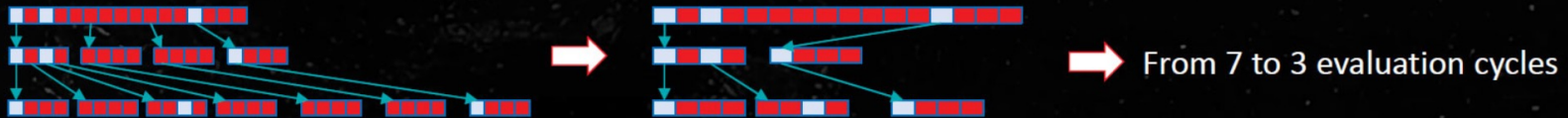
- Closest Midpoint

- Nodes are ordered from the closest midpoint of the intersection interval to the furthest
- Optimized for reflection rays / Closest Hit

AMD RDNA™ 3 2ND GENERATION RAY TRACING

RT Ray Material Shading creates many divergent memory requests patterns

- A New two step scheduling algorithm to discard empty ray quads and optimizes the cycles per ray



Hardware Stack Management Optimizations

- New `ds_bvh_stack_rtn` reduces required VALU & LDS instructions by ~50 per traversal iteration
- Reduced vector memory bus activity by a factor of 4 for stack updates per traversal iteration

Radeon™ RX 7900 series is equipped with 1.5x VGPR

- Enables up to 1.5x Rays in flight to hide latency and extract heavier ALU and intersection logic utilization
- RT traversal and shading performance with recursive DXR pipelines and large number of material shader usage

Additional uplift with larger caches and improved hit rates for complex scenes with secondary ray traversal and shading

Achieving up to a 1.8x performance uplift from config, frequency and features on heavy RT workloads



AMD RDNA3

12-4-22

CP, GEOMETRY AND PIXEL PIPE ADVANCES

Multi-Draw-Indirect Accelerator (MDIA)

- Up to 2.3x performance improvement of MultiDrawIndirect and MultiDrawIndexIndirect Execution
- Reduce CPU API and Driver overhead by accelerating gathering and parsing of Multi Draw command data

Native hardware support for 12 Primitive/Ck through culling – 50% Increase

- Up to 24 Vertices/Ck to support hardware-based culling for most meshes at peak rates
- Native 2x fixed function primitive culling hardware to remove SW based culling overhead

Configuration support for 50% more rasterized performance/clock

- Up to 6 Peak Primitives and 192 Peak Pixels of Rasterization/ Clock **+50%**

Random Order Opaque exports

- Replace large reorder buffer with scoreboard to enable non-overlapping or opaque Pixel Shader result posting OOO
- Enables divergent opaque Pixel Waves that finish early to post exports and free resources
- Improved shader resource utilization with active work

Pixel Wait Sync – Finer grain dependency management

- Producer: Previous Pixel Shader or fixed function writes → Consumer: Subsequent Pixel Shader or fixed Function dependent reads
- New sync mechanism stalls consumer at latest 3D pipe stage until producer are done writing
- Enables consumer geometry stages to execute ahead of producer completion to remove delays and start up latencies

12-4-22

RDNA™3: PREMIUM ADVANCED GRAPHICS INDUSTRY DEFINING CHIPLET ARCHITECTURE

Advanced **Chiplet** Design

- Disruptive Architecture vs Monolithic
- 5nm high performance Graphics Die
- 6nm Memory Cache Dies (MCD)
- Advanced technology – 1st in gaming

Architected to exceed 3Ghz – Industry 1st

- 61 TFLOPs Boost FP32 – ~2.7x increase
- ~1.54x Perf/Watt for a 3rd Generation

New **ALUs** Instructions and Throughput

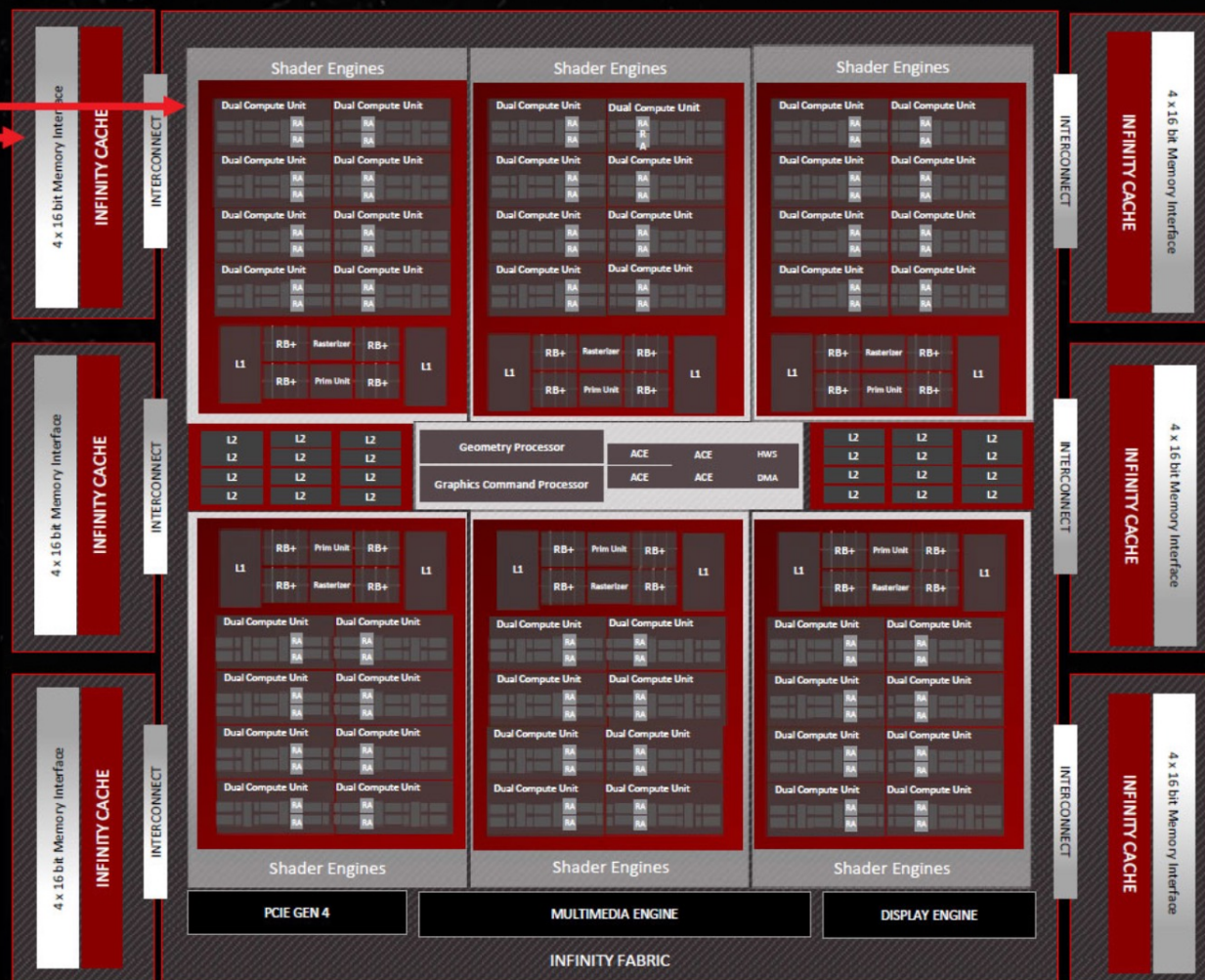
- Up to 2x ALU rates plus BF16 support
- New Instructions for effective utilization

Optimized & Balanced **Cache System**

- 96 MB 2nd Gen Infinity Cache
- 6MB L2 Cache – 50% Increase
- 3MB L1 Cache – 300% Increase
- 3MB L0 Cache - 240% Increase

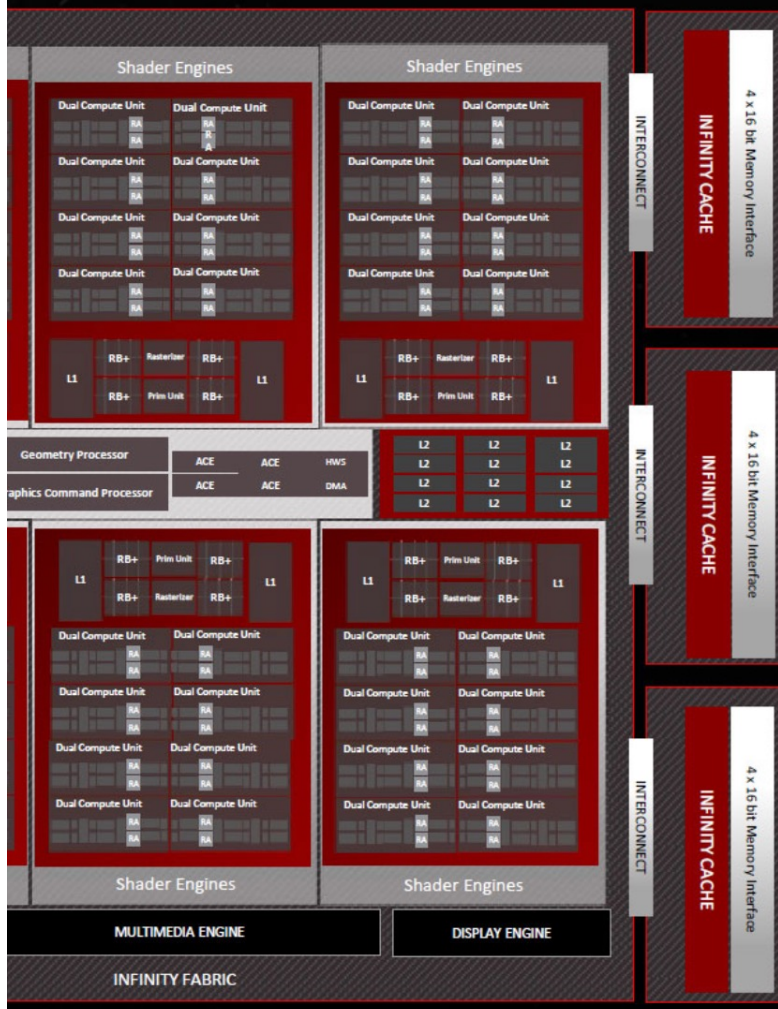
2nd Generation **Ray Tracing**

- RT Features for Performance & Efficiency
- Larger Caches for Complex RT Workloads
- Up to **1.8x** RT performance @2.5GHz



12-4-22

AMD ADVANCED GRAPHICS RDNA3 CHIPLET ARCHITECTURE



Flexible CP & Geometry Pipe

- Multi Draw Indirect Accelerator (MDIA)
- 12 Primitive/Clock – 50% Increase
- 2x Hardware Prim/Vert Cull Rates

Advances in Pixel Pipe

- 6 Prims Rasterized/Clock – 50% Increase
- 192 Pixels/Clock – 50% Increase
- Random Order Opaque exports
- Pixel Wait Sync

High Speed GDDR6 Memory

- Up to 384b @ 20 Gbps – 960 GB/s
- Up to 24 GB of GDDR6

AMD Radiance Display™ Engine

- DisplayPort™ 2.1 & HDMI 2.1a
- 12 bit/channel for up to 68 billion colors

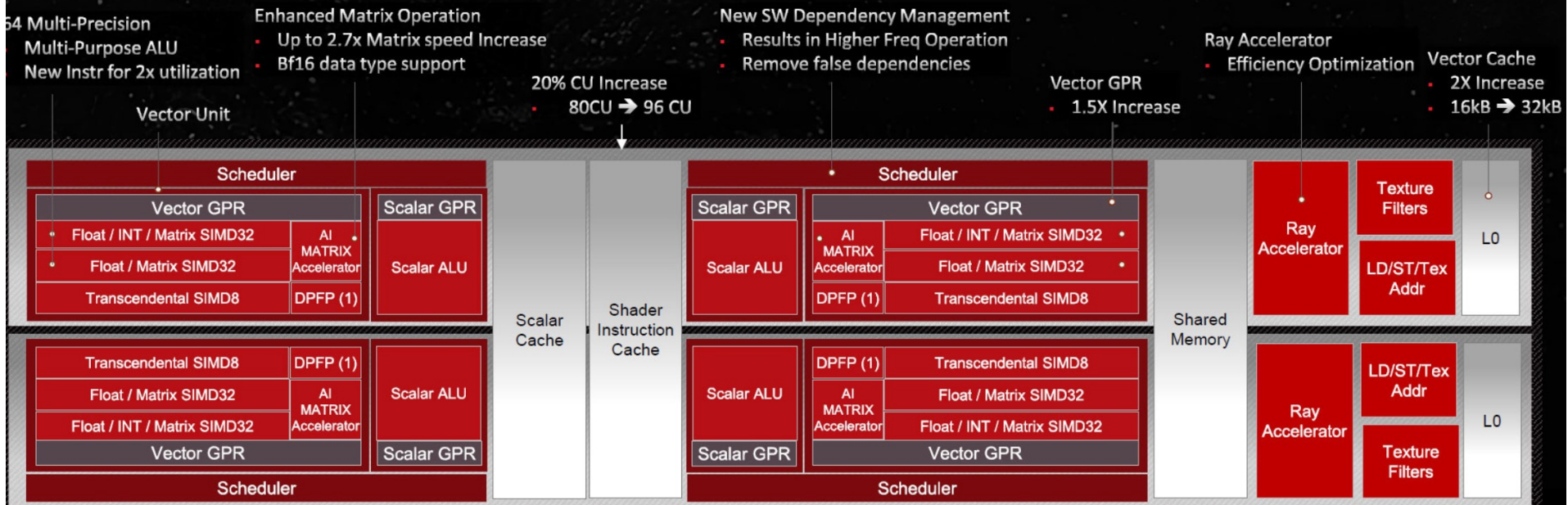
New Dual Media Engine

- Simultaneous Encode/Decode (AVC/HEVC)
- 8k60 AV1 Encode/Decode
- AI Enhanced Video Decode

Leading Features

- Full DirectX12 Ultimate
- Fidelity FX Super Resolution
- AMD Advantage Smart Technologies

THE ENHANCED COMPUTE UNIT PAIR

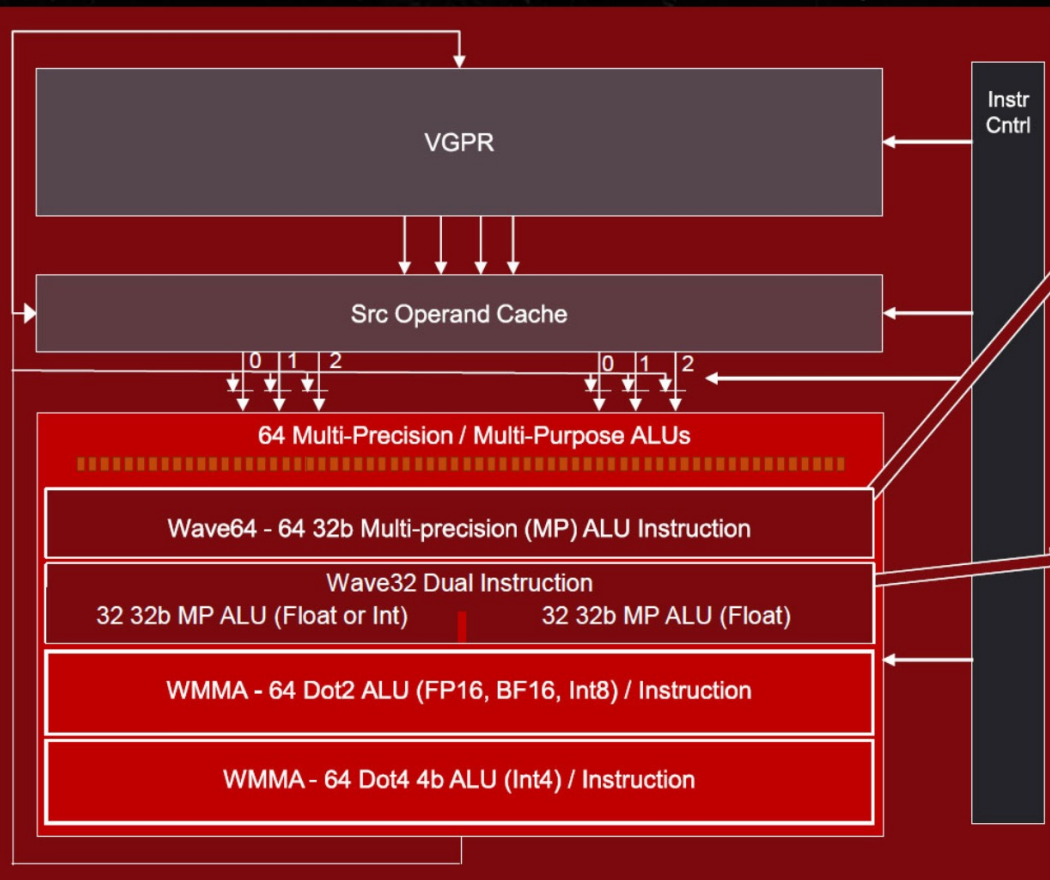


Enhanced CU delivers approximately 17.4% architectural improvement clock for clock

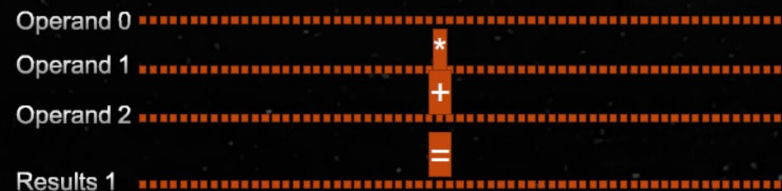
AMD RDNA3

12-4-22

VECTOR UNIT AS 1 SIMD64 OR 2 SIMD32



One Clock Wave 64 FMA Instructions

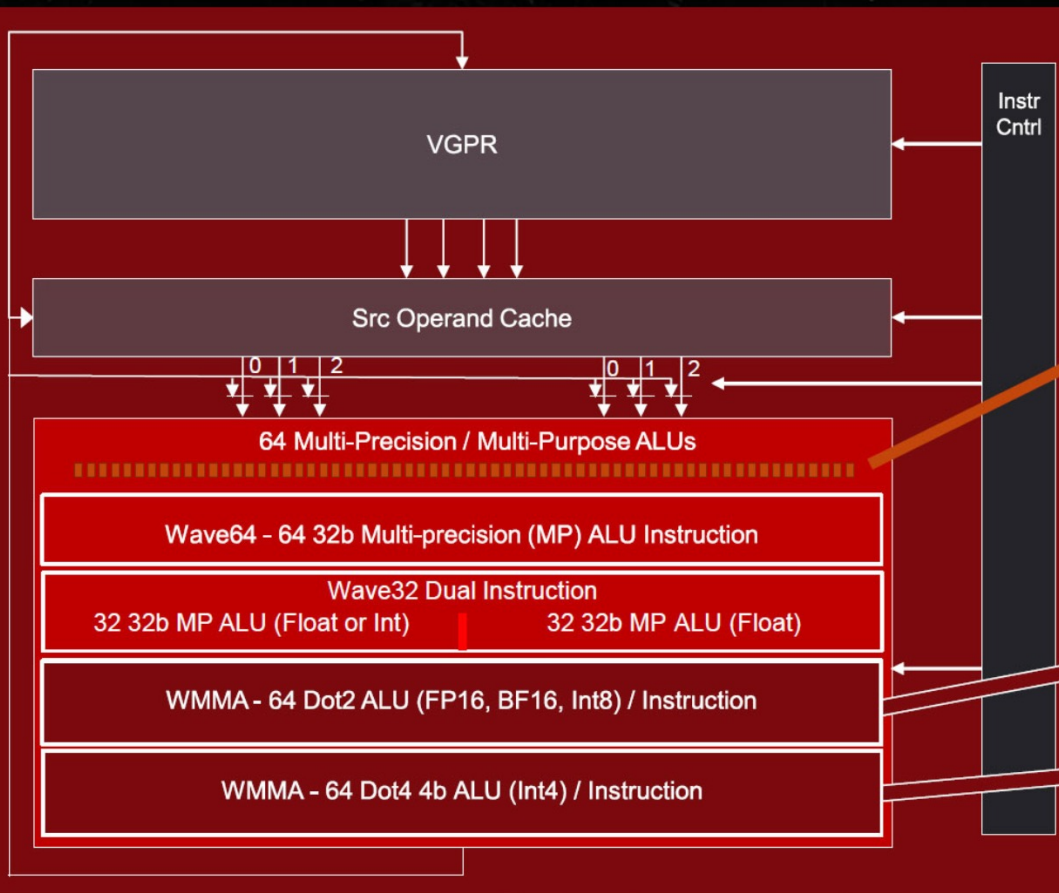


Dual issue Wave32 one clock Instructions



Instructions have same 1 clock issue and dependency rules of RDNA Architectures

VECTOR UNIT AS MATRIX ACCELERATOR



Wave Matrix Multiply Accumulate



WMMA - Wave Matrix Multiply Accumulate

64 Vector 32b float ALU employed

64 Dot2 Ops/cycle - FP16, BF16, & Int8

$$C_{00} += A_{00} \times B_{00} + A_{01} \times B_{10}$$

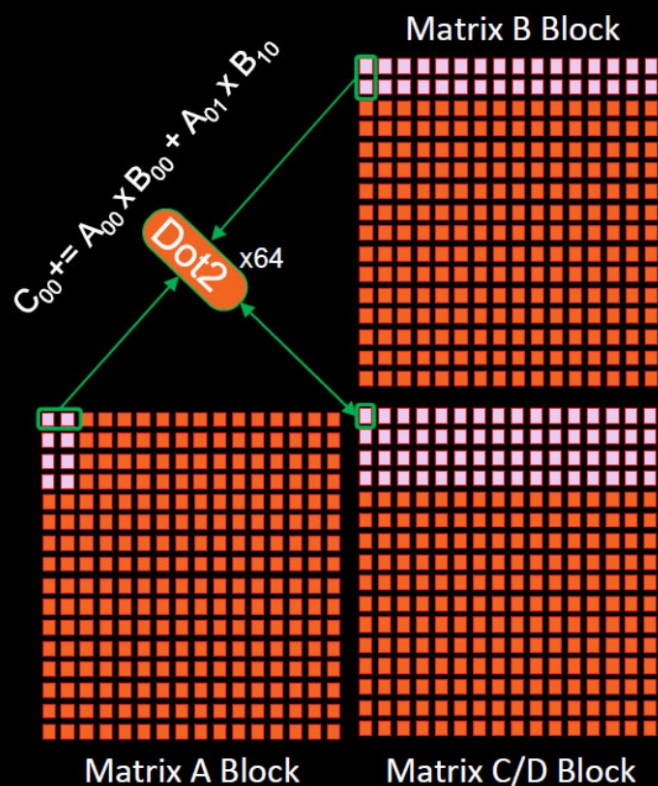
64 Dot4 Ops/cycle - Int4

$$C_{00} += A_{00} \times B_{00} + A_{01} \times B_{10} + A_{02} \times B_{20} + A_{03} \times B_{30}$$

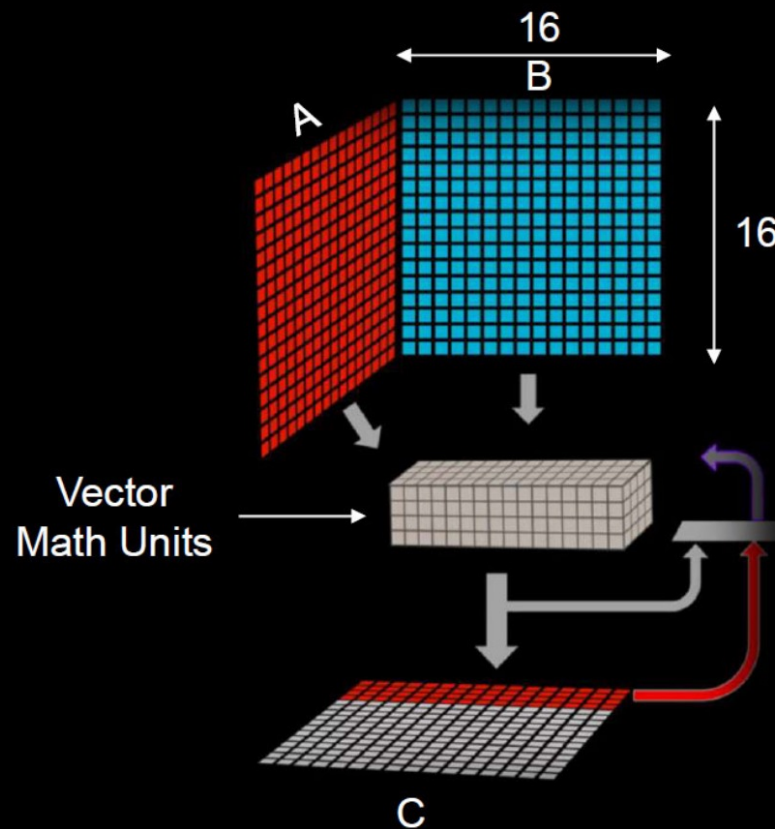
12-4-22

ONE VECTOR UNIT

WMMA MATRIX MULTIPLICATION



64 Dot2 per cycle (256 Ops/cycle)
32 cycle: 2048 Dot2 Operations





GPU Software



ROCm™ Open Software Platform for GPU Compute

DRIVING MAINSTREAM ADOPTION & ECOSYSTEM ENABLEMENT

Applications	HPC Apps		ML Frameworks	
Cluster Deployment	Singularity	SLURM	Docker	Kubernetes
Tools	Debugger	Profiler, Tracer	System Valid.	System Mgmt.
Portability Frameworks	Kokkos	RAJA	GridTools	ONNX
Math Libraries	RNG, FFT	Sparse	BLAS, Eigen	MIOpen
Scale-Out Comm. Libraries	OpenMPI	UCX	MPICH	RCCL
Programming Models	OpenMP	HIP	OpenCL™	Python
Processors	CPU + GPU			

2021: ROCm 4.0
Production-Ready HPC & ML Stack



AMD
ROCm 5.0

EXPANDING SUPPORT & ACCESS

OPTIMIZING PERFORMANCE

ENABLING DEVELOPER SUCCESS