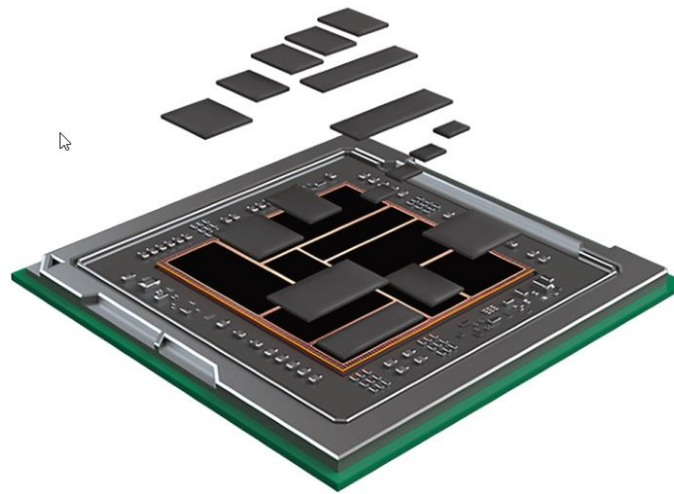


Beyond Moore's Law: Leveraging Advanced Packaging Technology

Cătălin Ciobanu, Transilvania University of Braşov (UnitBV)



Transilvania
University
of Braşov



Moore's Law



Gordon Moore (1929 – 2023)

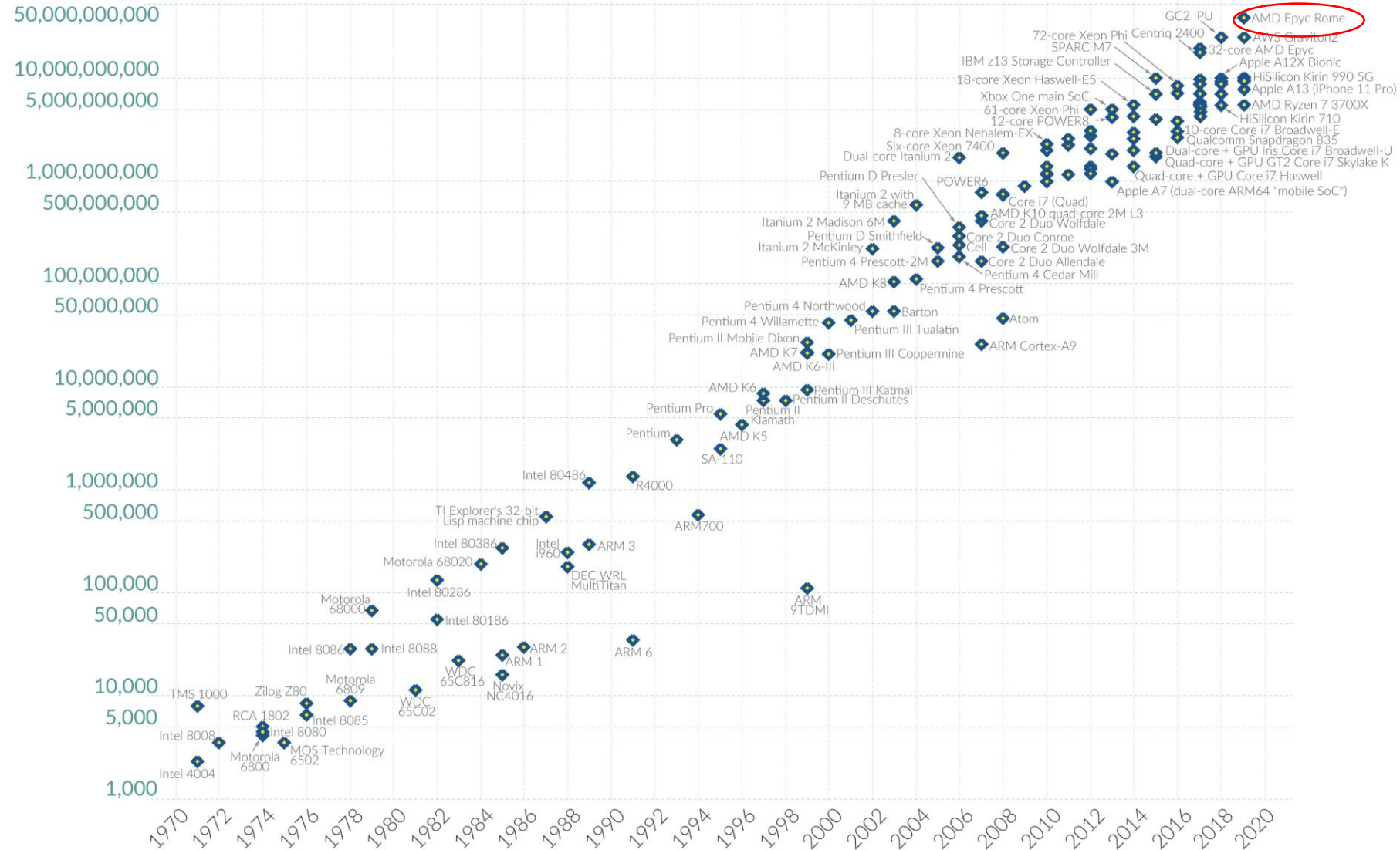
- Prof. Carver Mead popularized the phrase „Moore's Law”
- „Moore's law is the observation that the number of transistors in an integrated circuit (IC) doubles about every two year”
- Several factors contributing to this exponential behavior (1975 IEEE International Electron Devices Meeting)
 - Metal-oxide semiconductor (MOS) technology
 - **The exponential rate of increase in die sizes** coupled with a decrease in defective densities
 - **Finer minimum dimension**



Moore's Law: The number of transistors on microchips doubles every two years

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers.

Transistor count



Data source: Wikipedia (wikipedia.org/wiki/Transistor_count) OurWorldinData.org – Research and data to make progress against the world's largest problems. Licensed under CC-BY by the authors Hannah Ritchie and Max Roser. Sursa: https://en.wikipedia.org/wiki/Moore%27s_law

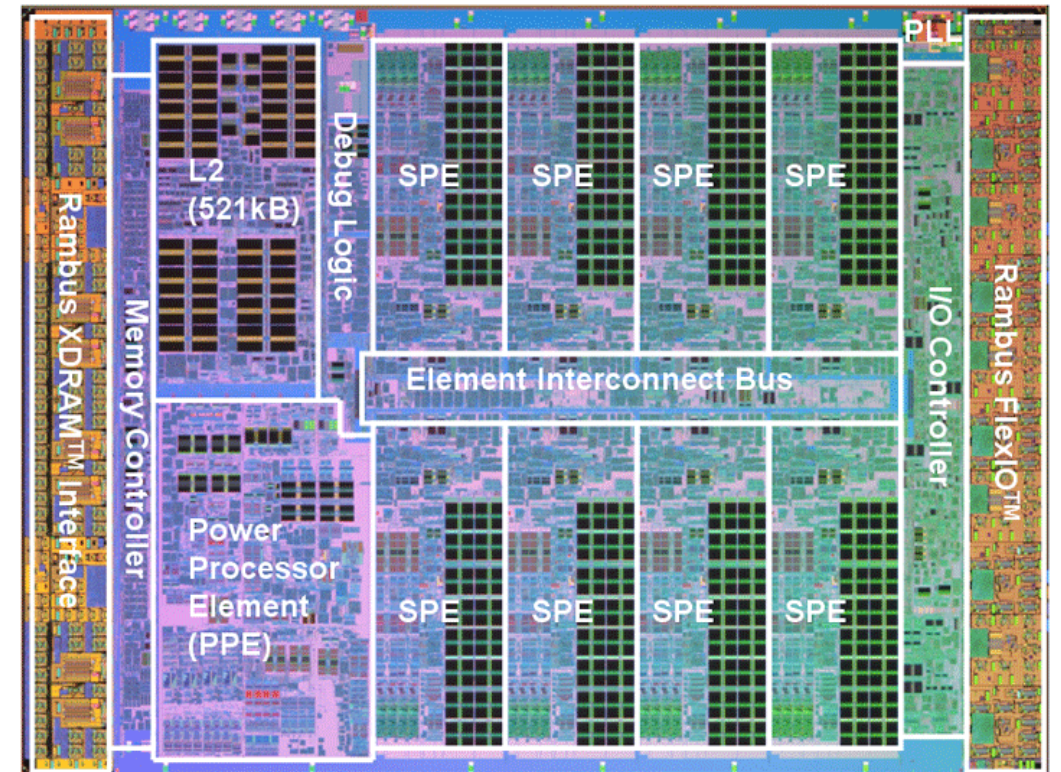
Beyond Moore's Law

- ❑ New **challenges** in semiconductor industry – multiple **constraints in advanced chip design**
 - ❑ **Power** limitations
 - ❑ **Thermal** limitations
 - ❑ **Limited** Instructions Per Clock (**IPC**) gains in new generations of processors
- ❑ Industry **trends**
 - ❑ **Multicore** designs
 - ❑ Wide use of **accelerators**
- ❑ Wafer costs for each new technology node are rising
 - ❑ **Monolithic** designs are **prohibitively expensive**
 - ❑ **Lower yields** inherent to large chips
- ❑ Reticle limit for future High-NA (Numerical Aperture) Extreme Ultraviolet (EUV) Lithography will be halved
 - ❑ **856mm²** -> **429mm²** due to the use of an amorphous lens array¹

1: <https://en.wikichip.org/wiki/mask>

CELL Broadband Engine (2006) - Heterogeneous Multicore

- ❑ Power Processor Element (PPE) for control tasks
- ❑ Synergistic Processor Elements (SPE) for data-intensive processing



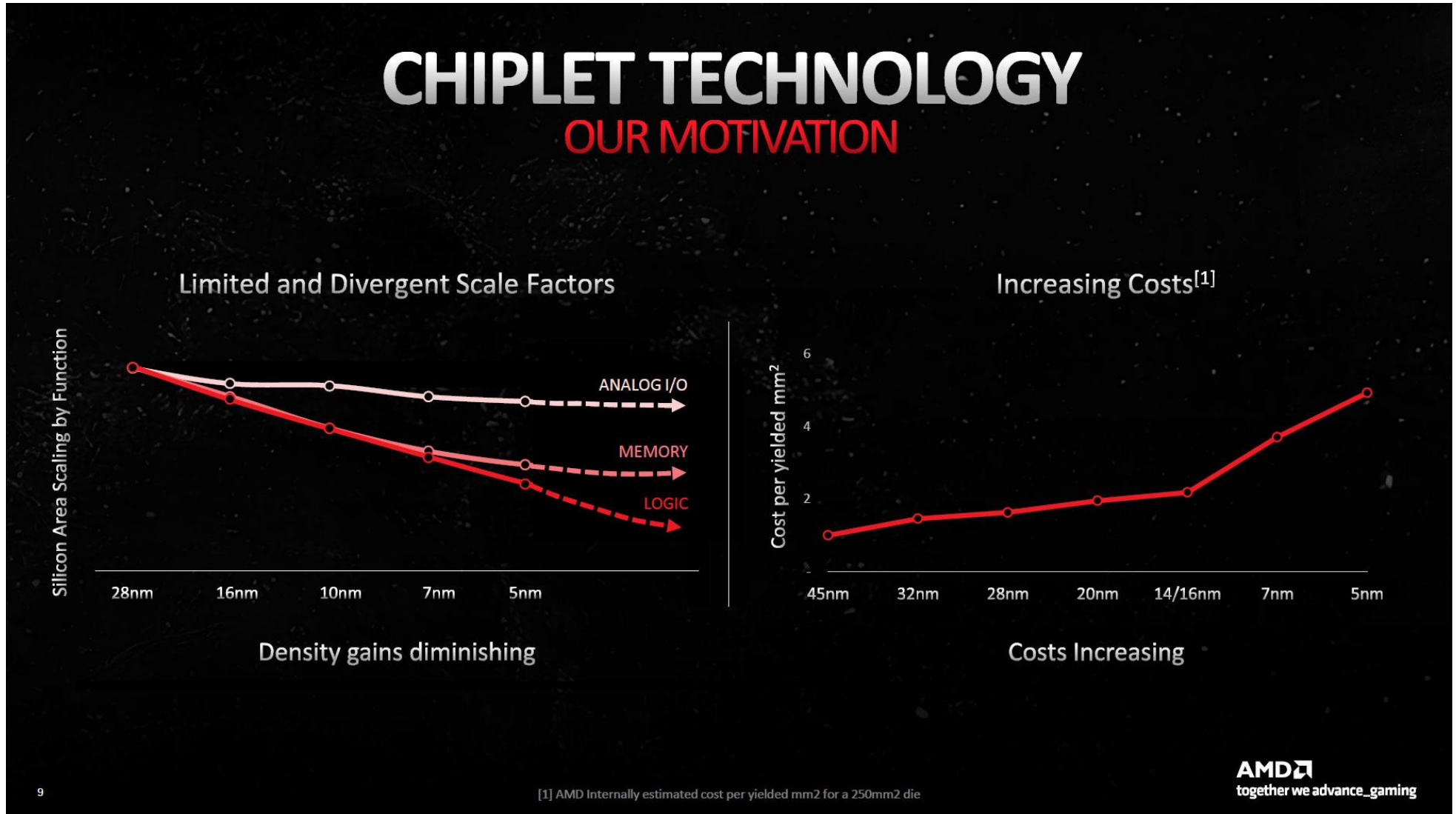
Solution – Advanced packaging - Chiplet technology

- ❑ “Chiplet technology is a microelectronics design and manufacturing approach where multiple smaller dies or chiplets are combined into a single package, with each of chiplets performing a specific function.”¹
- ❑ **Multiple chiplets** – connected to form a SoC-like solution
 - ❑ **Disaggregate** large designs
 - ❑ **Avoid** the silicon **reticle size limitation**
- ❑ **Multiple advantages** when using chiplets
 - ❑ **Flexibility**
 - ❑ **Scalability**
 - ❑ **Cost savings** – off-the-shelf chiplets
 - ❑ **Mixing functions from different process nodes**



Chipelets Motivation

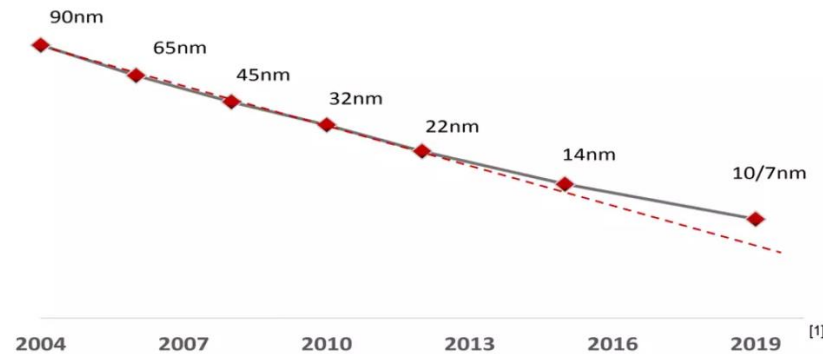
Chiplets Motivation



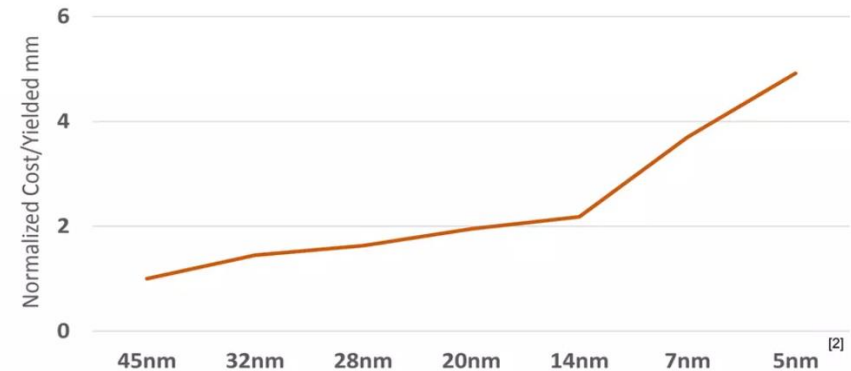
Moore's Law Keeps Slowing

Barriers to Large Cache Size

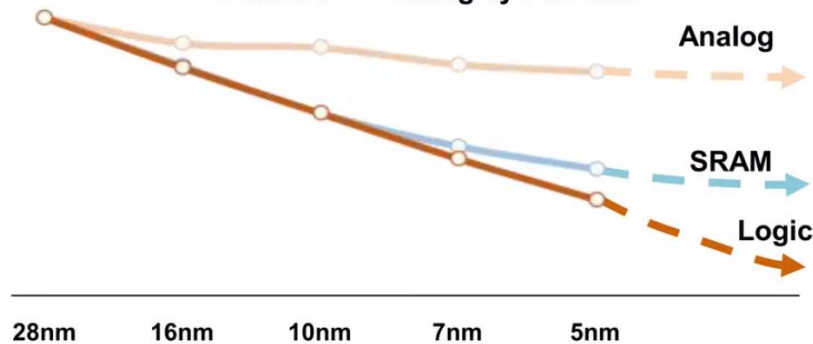
MOORE'S LAW KEEPS SLOWING



WHILE COSTS CONTINUE TO INCREASE



Silicon Area Scaling by Function



■ **SRAM scaling & cost = barriers to large on-die caches**

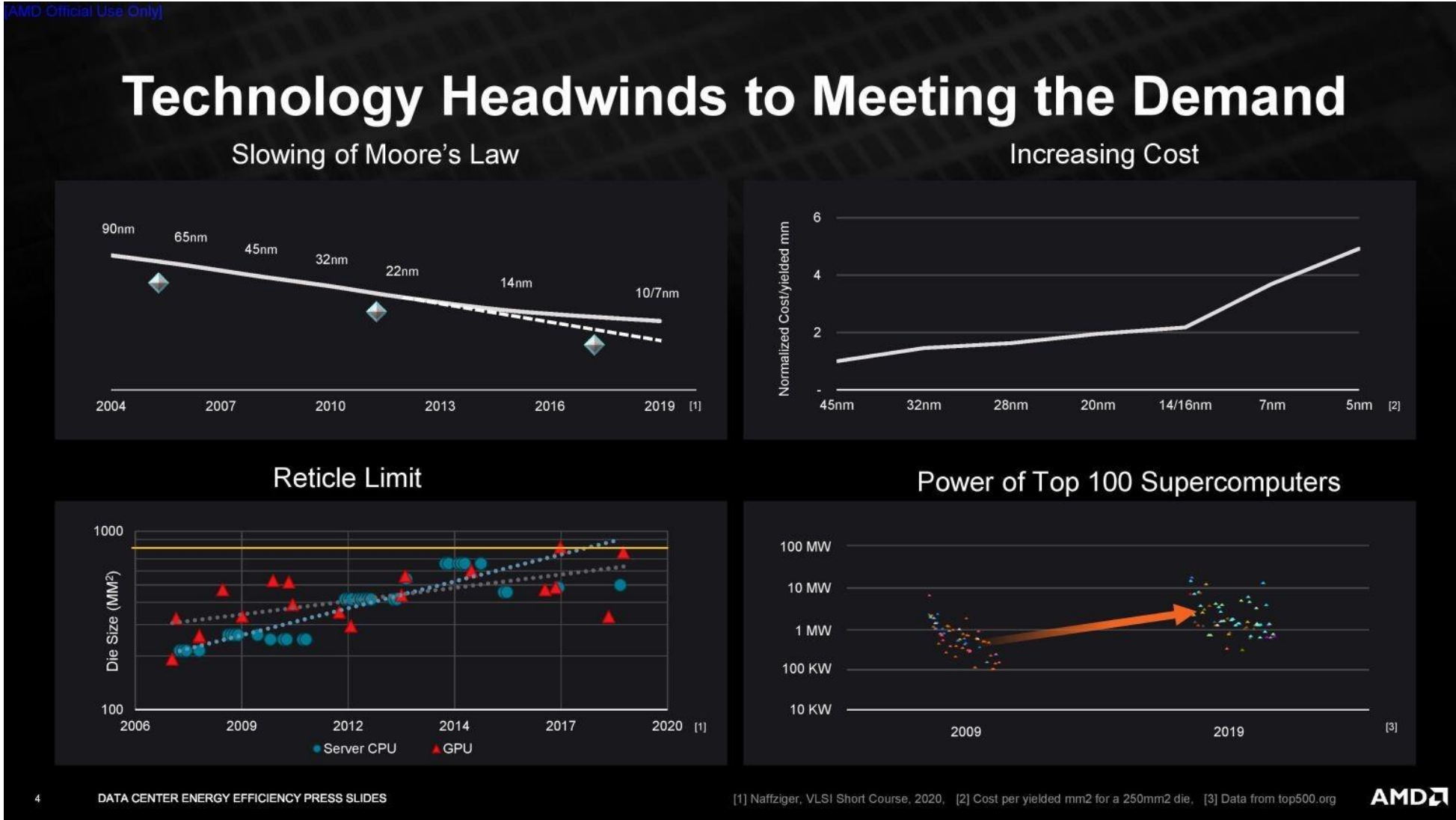
- Cost
- Product flexibility
- Latency

} Chipllets

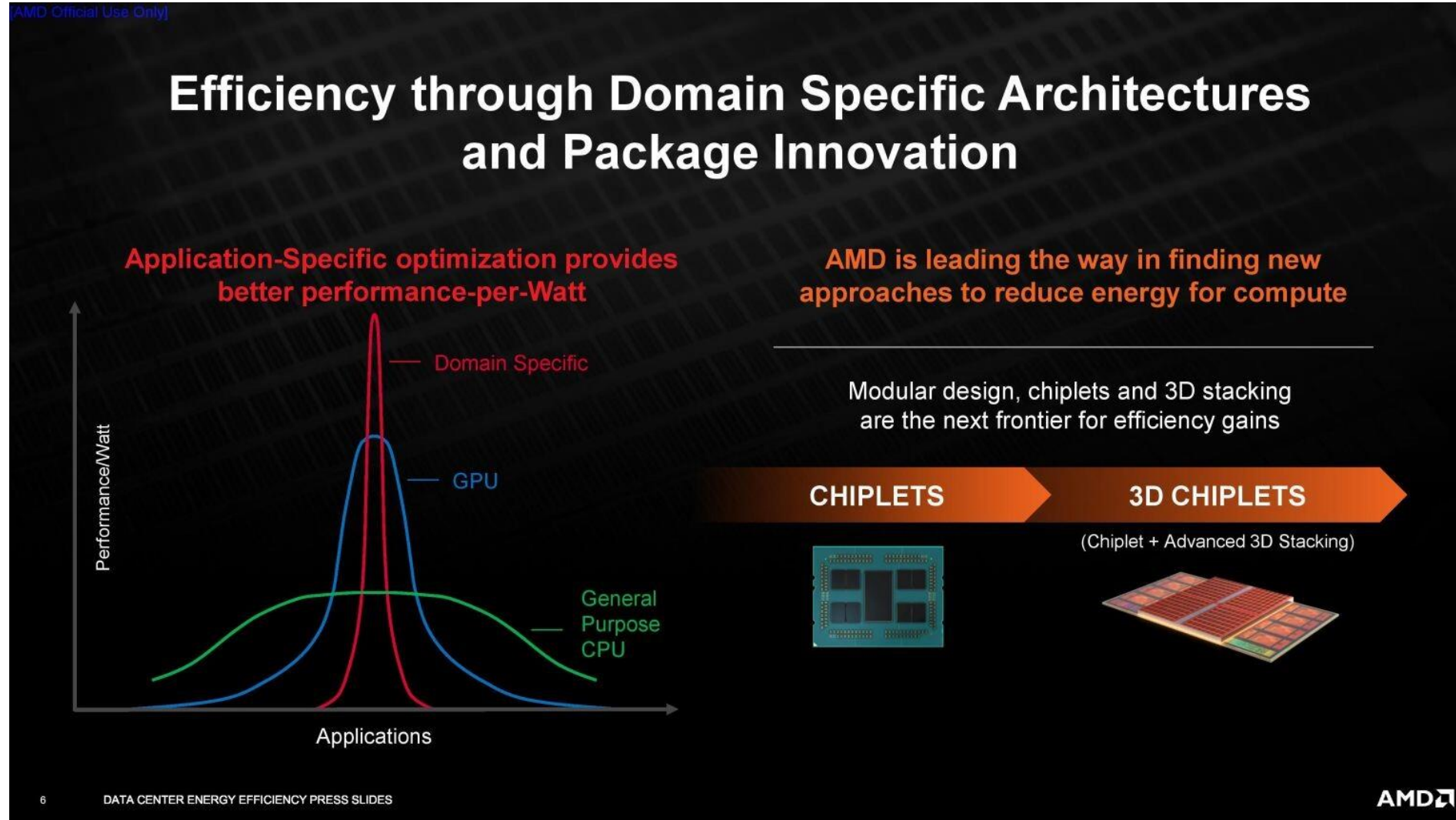
© 2022 IEEE International Solid-State Circuits Conference

[1] Naffziger, VLSI Short Course, 2020, [2] Cost per yielded mm² for a 250mm² die
26.4: 3D V-Cache: The Implementation of a Hybrid-Bonded 64MB Stacked Cache for a 7nm x86-64 CPU

Slowing of Moore's Law + Increasing Cost



Application-Specific optimization, Modular design: chiplets + 3D Stacking

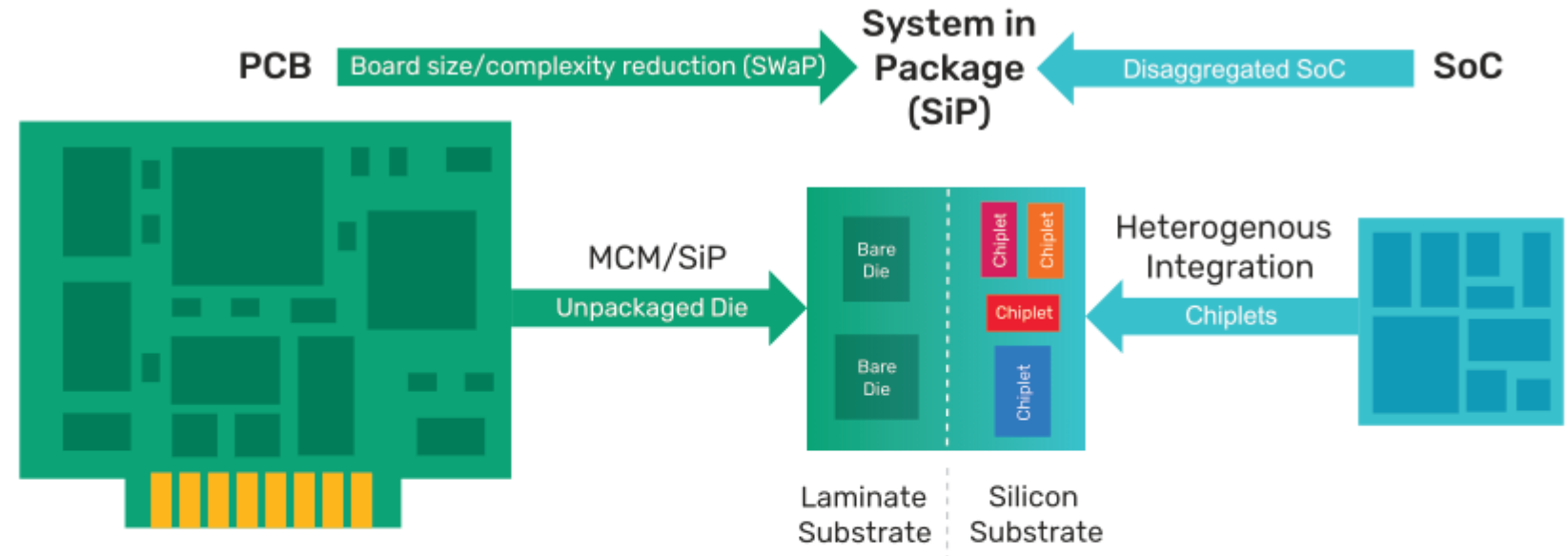




Chiplet Technology

Chiplet Technology and Heterogeneous Integration

- ❑ **Transition** from layout of laminate substrates to silicon substrate
- ❑ **Electrical and thermal analysis challenges**
- ❑ Multi-chip Modules (**MCM**) date back to the 1960s
- ❑ System in Package (**SiP**) began to replace the term MCM in the late 1990s
- ❑ **2.5D IC packaging** - Silicon substrates – high density, using through-silicon vias (TSVs)



PCB to MCM/SiP Benefits

- ▶ Smaller footprint
- ▶ PCB simplification
- ▶ Higher bandwidth
- ▶ Lower power

SoC to HI Benefits

- ▶ Reduced NRE costs
- ▶ Shorter time to market
- ▶ Larger than reticle size designs
- ▶ More flexible IP use-model

Figure 1: Heterogeneous integration

Chiplet Technology Challenges

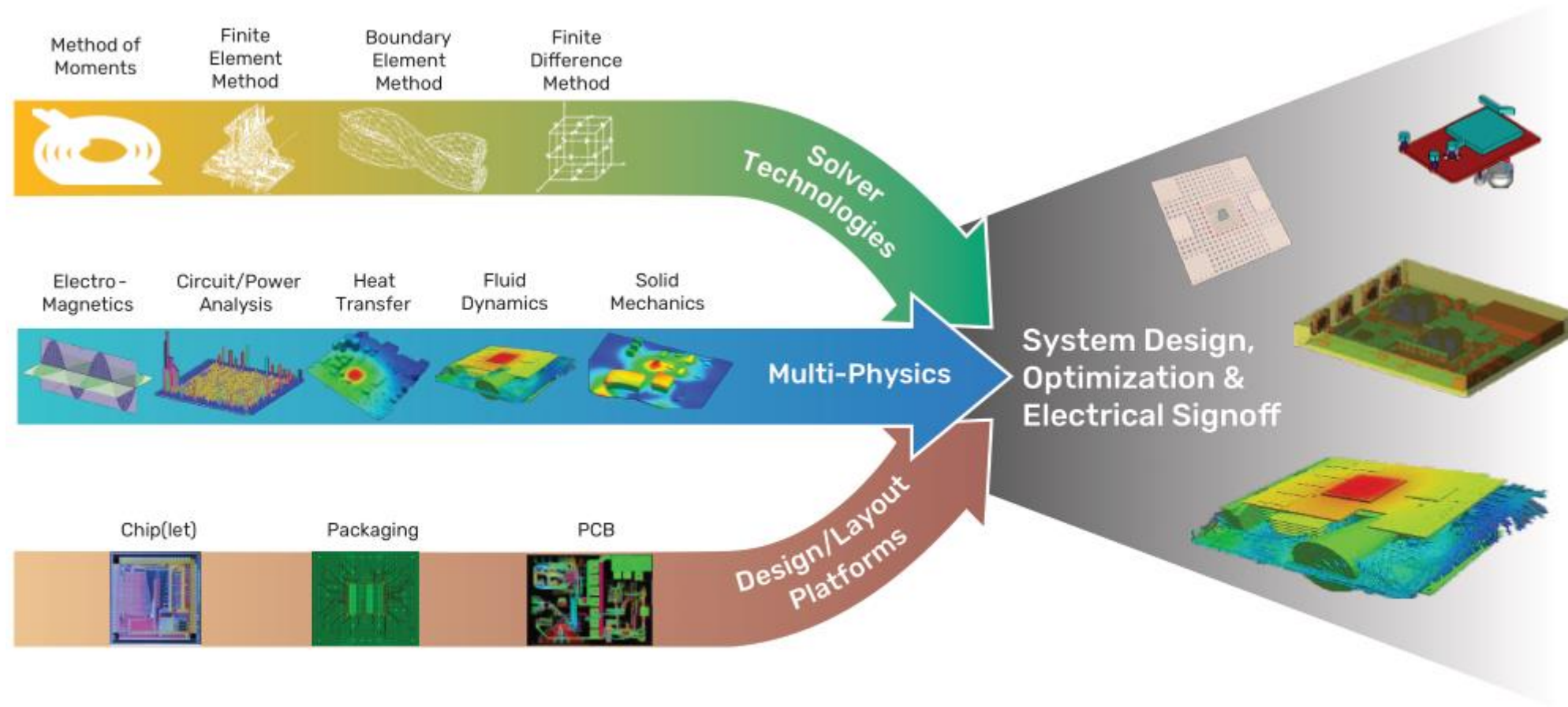
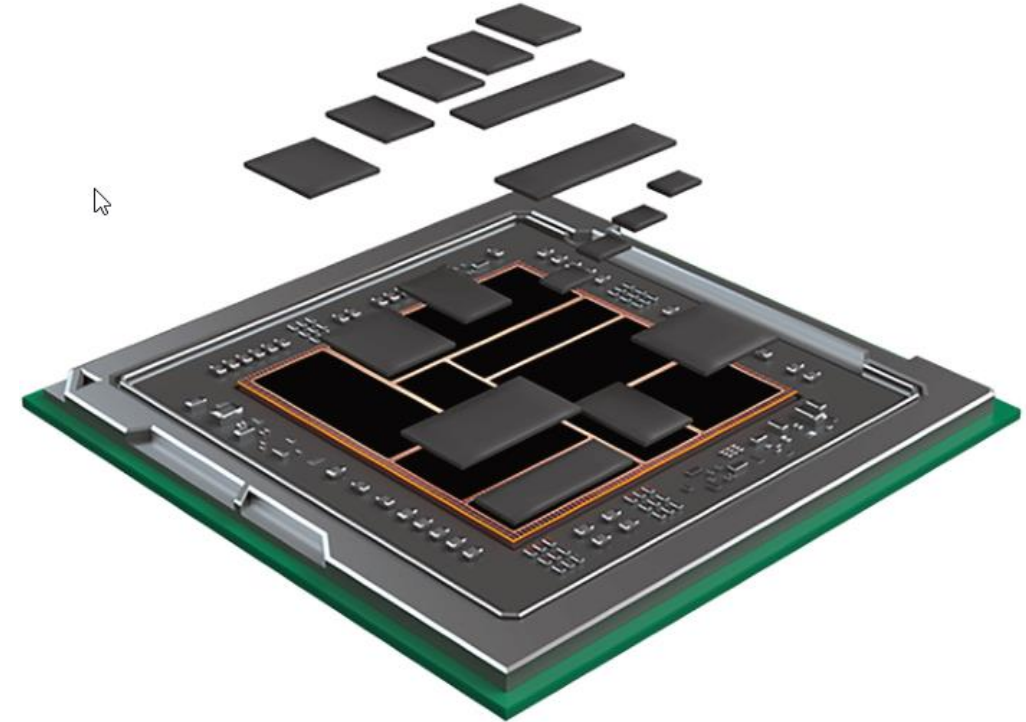


Figure 2: System-level electrothermal analysis

Chiplets Heterogeneous Integration (HI)

- Similar to designing a small PCB
- Each chiplet built with a common/known communication interface
 - PCIe
 - HBM
 - AIB
- Lower development cost – modular integration
- Lower manufacturing costs – purchasing known-good die (KGD)
- Cost advantage – volume manufacturing when reusing the same chiplets in many designs
- Many vendors exploring this space
 - Intel CO-EMIB – EMIB + Foveros in the same package
 - Intel Omni-Directional Interconnect (ODI) – horizontal communication (similar to EMIB) or vertically using TSVs (similar to Foveros)
 - TSMC's Chip-on-Wafer-on-Substrate (CoWoS)



Advanced Packaging Technology Evolution

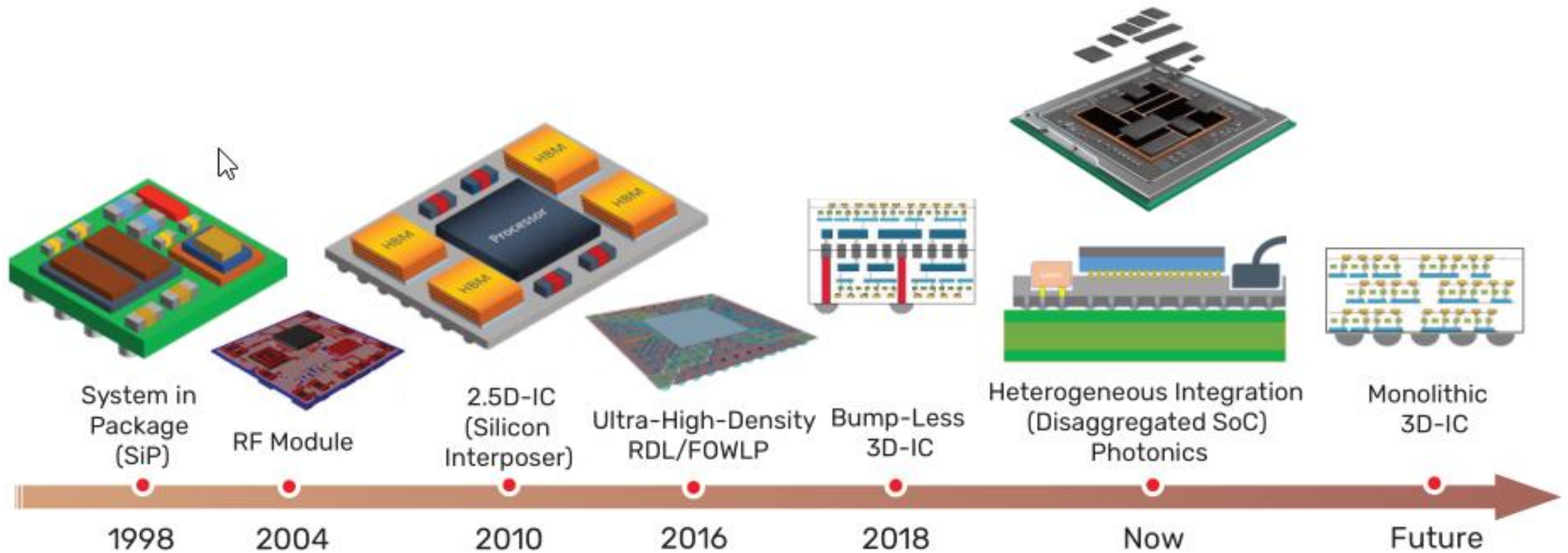






Figure 4: Evolution of advanced multi-chip(let) packaging technologies

Industry Shift Towards Multi-Die SoCs

- ❑ Shift fueled by several converging trends
 - ❑ Some SoCs – too big for manufacturability
 - ❑ Some SoCs require different process nodes for optimal cost/perf
 - ❑ Desire for enhanced product scalability and composability
- ❑ Lack of design ecosystem – pause/postpone multi-die projects
- ❑ Early adopters developed proprietary die-to-die interfaces
 - ❑ Limits the ability to assemble dies from different vendors
- ❑ Solution: standardized die-to-die interconnects
 - ❑ **Optical Interface Forum (OIF)** – The XSR and USR physical layer specifications optimized for die-to-die connectivity
 - ❑ **Chips Alliance** – The **AIB** specification which was originally introduced by Intel
 - ❑ **Open Compute Platform (OCP)** – The OpenHBI and Bunch-of-Wires (BOW) specifications optimized for different use cases
 - ❑ **Unified Chiplet Interconnect Express (UCIe)** – A comprehensive die-to-die interconnect specification covering multiple use cases and a complete protocol stack

Die-to-die interconnect standards

Alliance					
Standard	XSR	BOW	OHBI	AIB	UCle
Data Rate	112G / 224G	8G / 16G	8G / 16G	6G	16G / 32G
Protocol	Not Defined	Not Defined	Not Defined	Not Defined	Streaming, PCIe, CXL
Package Types	2D	2D, 2.5D	2D, 2.5D	Bridge	2D, 2.5D, Bridge
Target Applications	Optical Networking (CPO/NPO)	Cost sensitive aggregation cases	High density scale and split cases for data center	Mil-aero ecosystem	Scale & Split w/ streaming Aggregation w/ PCIe/CXL

Source: <https://www.synopsys.com/designware-ip/technical-bulletin/ucle-multi-die-socs.html>

Synopsys Multi-Die SoCs Gaining Strength with Introduction of UCle

UCle: Companies Join Forces – standardized die-to-die interconnect

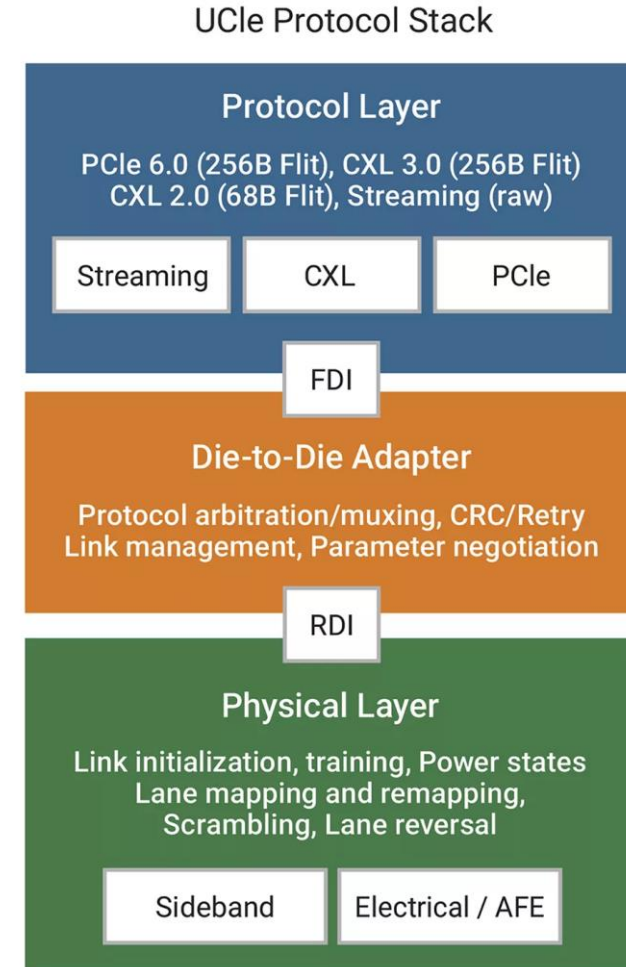


Source: <https://www.synopsys.com/designware-ip/technical-bulletin/ucie-multi-die-socs.html>

Synopsys Multi-Die SoCs Gaining Strength with Introduction of UCle

UCIe: a Complete Stack for Die-to-die interconnect

- ❑ Supports data rates from 8Gbps/pin 16Gbps/pin
 - ❑ Expected to support 32GBps/pin
- ❑ Supports all types of package technology
 - ❑ UCIe for advanced packages (silicon interposer, silicon bridge or RDL fanout)
 - ❑ UCIe for standard packages (organic substrate or laminate)
 - ❑ Both options share the same architecture and protocols
- ❑ Very competitive performance advantages to multi-die SoC designers
 - ❑ High energy efficiency (pJ/b)
 - ❑ High edge usage efficiency (Tbps/mm)
 - ❑ Low latency (ns)
- ❑ UCIe – compelling roadmap
 - ❑ Higher data rates
 - ❑ New protocols
 - ❑ 3d packaging
 - ❑ Security
 - ❑ Testability

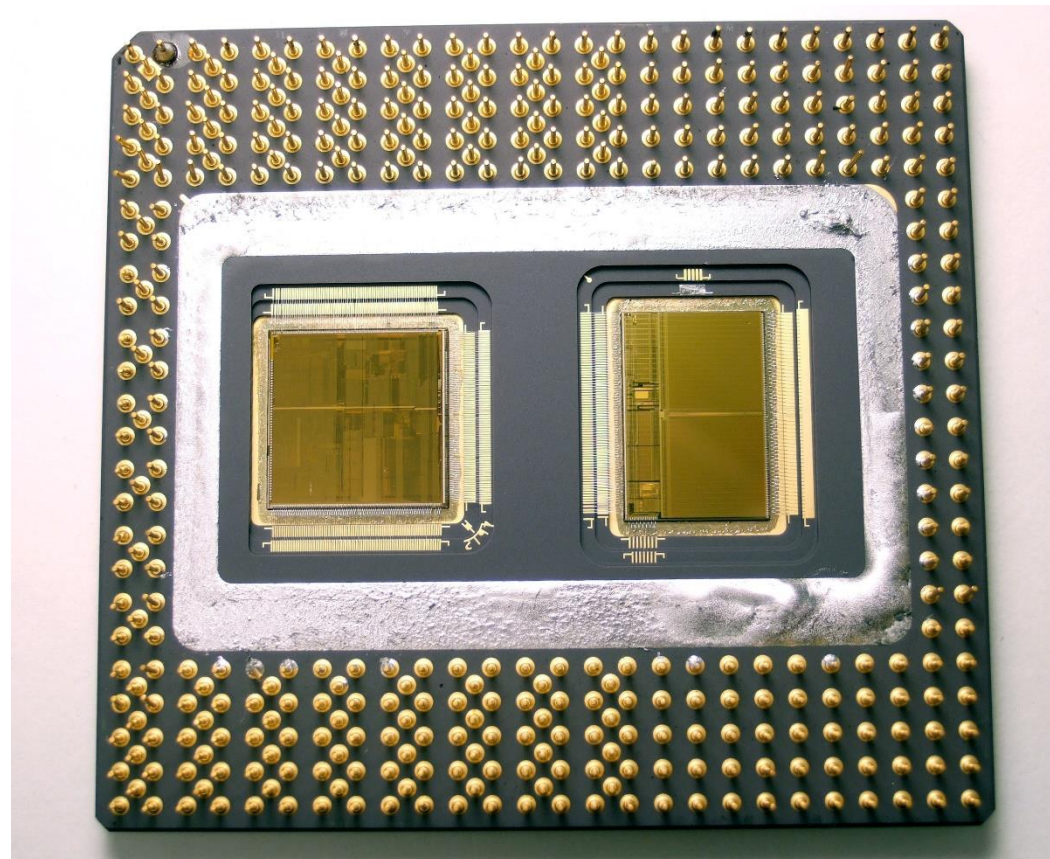


Brief History of Multi Chip Modules

Intel Pentium Pro Multi-Chip Module (MCM) (1995)

Pentium Pro is packaged in Multi-Chip Module (MCM) – “on package cache”

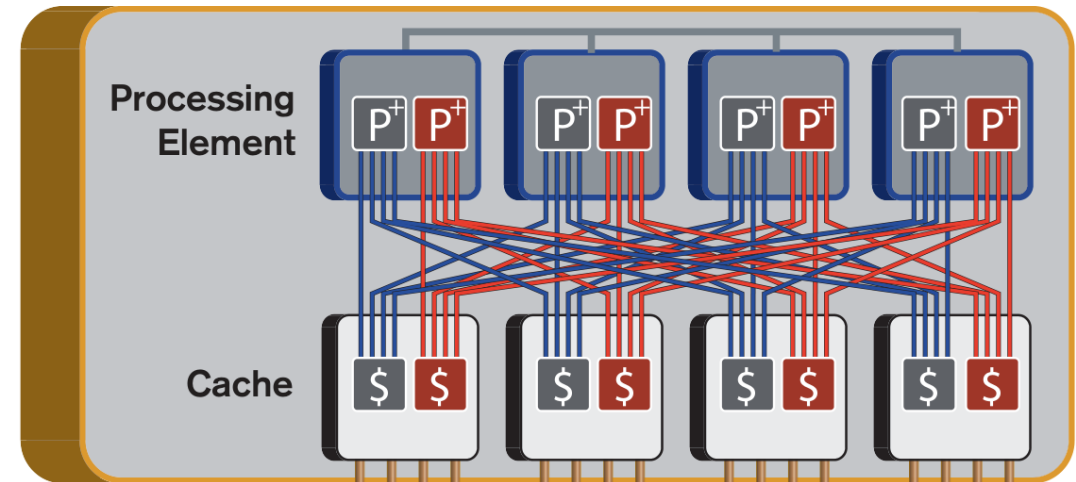
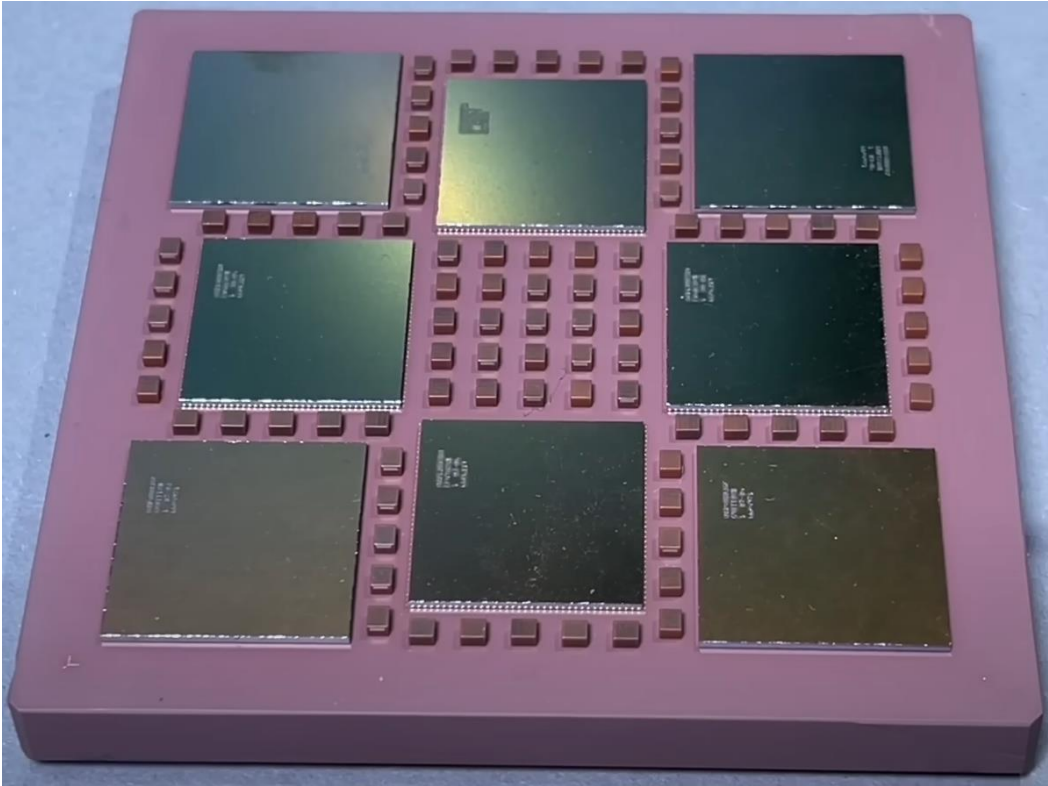
- ❑ Separate L2 die in the package, same speed as CPU core
- ❑ The dies are connected to the package using conventional wire bonding
- ❑ Up to two 512KB cache dies
- ❑ 0.35 μm to 0.50 μm



Cray X1 (2003)

8-chip Multi-Chip Module (MCM)

- ❑ 4 processor chips
- ❑ 4 custom streaming cache chips



Source: CPU Galaxy - CRAY X1 Multi Chip Module - Monster CPU with 8 Cores – Teardown, Feb 18, 2023 <https://www.youtube.com/watch?v=QdteJTicqDE>
https://www.craysupercomputers.com/downloads/CrayX1/CrayX1E_Datasheet.pdf

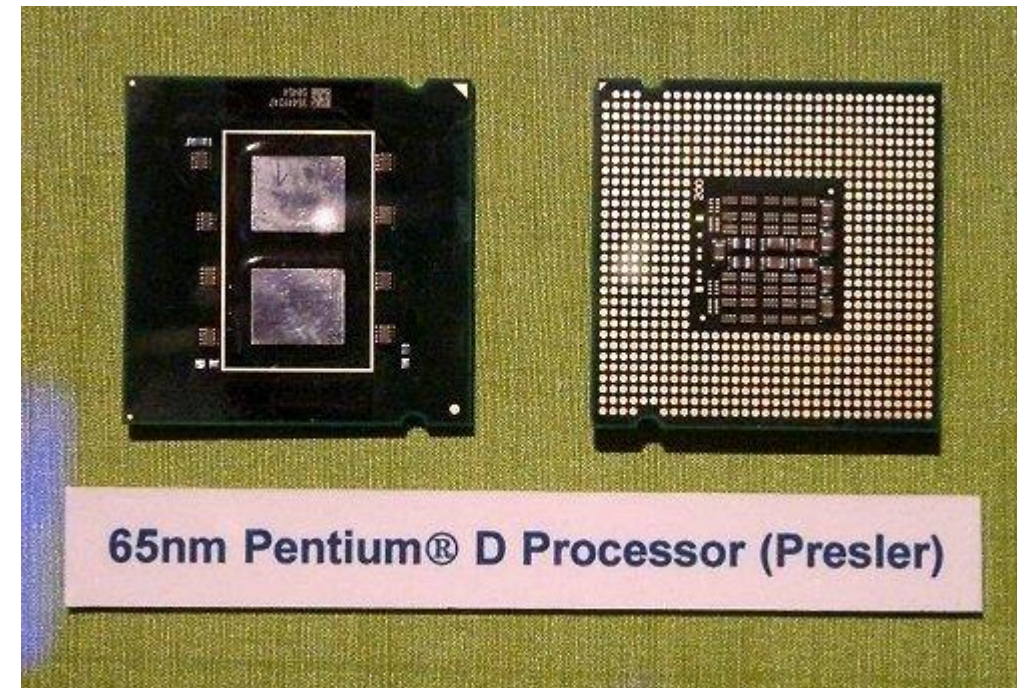
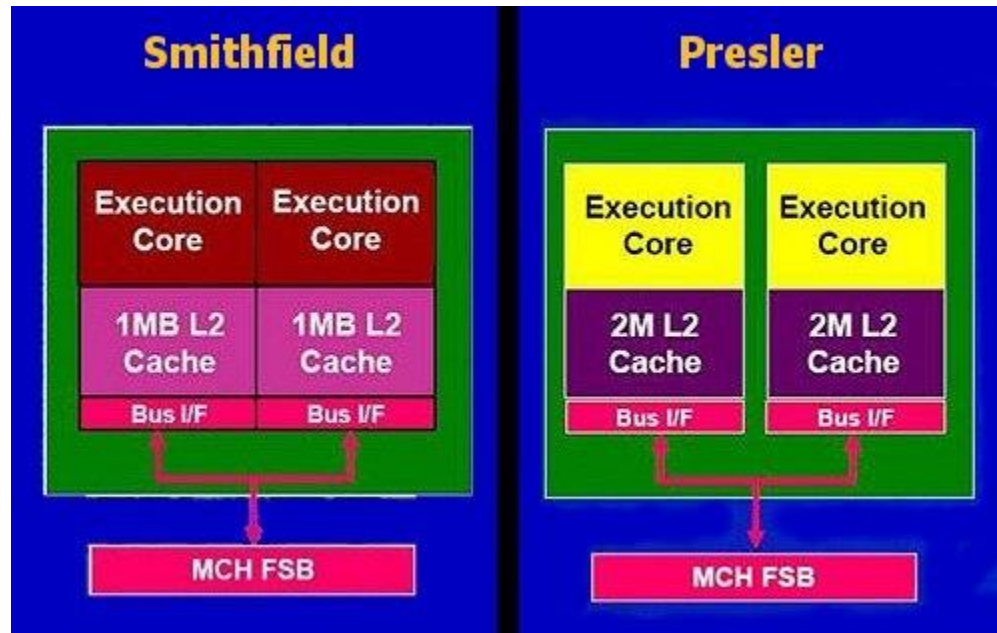
Intel Pentium D Presler (2006)

Smithfield Pentium D 800 series had a single die with two cores

- ❑ Two single-core dies placed on the same substrate
- ❑ Reduced cache available to each core
- ❑ 90nm, 206 mm²

Presler Pentium D 900 series used a Multi-Chip Module (MCM)

- ❑ Two single-core dies placed on the same substrate
- ❑ Each die could be sold as a Pentium 4 CPU
- ❑ Reduced cost due to higher yields
- ❑ 65nm, 162 mm² for both cores



AMD Zen Embraces Chiplets, 3D Vcache

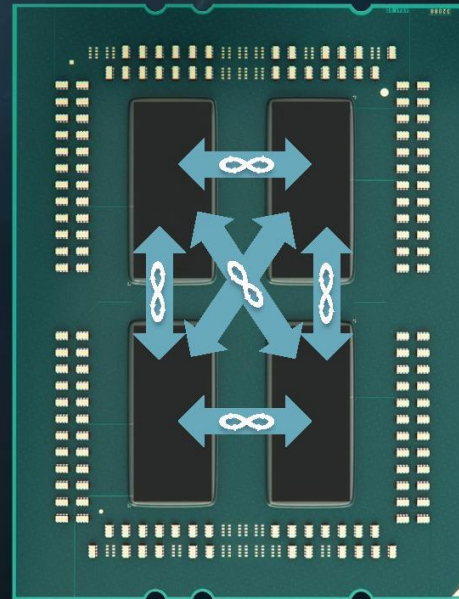
AMD EPYC Naples (2017)

- ❑ Four 14nm compute dies, each with 8 cores
- ❑ Each die connected directly to the others via Infinity Fabric

INFINITY FABRIC: DIE-TO-DIE INTERCONNECT

- Fully connected Coherent Infinity Fabric within socket
- Optimized low power, low latency MCM links between die
- 42GB/sec bi-dir BW per link, ~2pJ/bit TDP
- Single-ended, low power zero transmission
- Infinity Control Fabric connected between dies

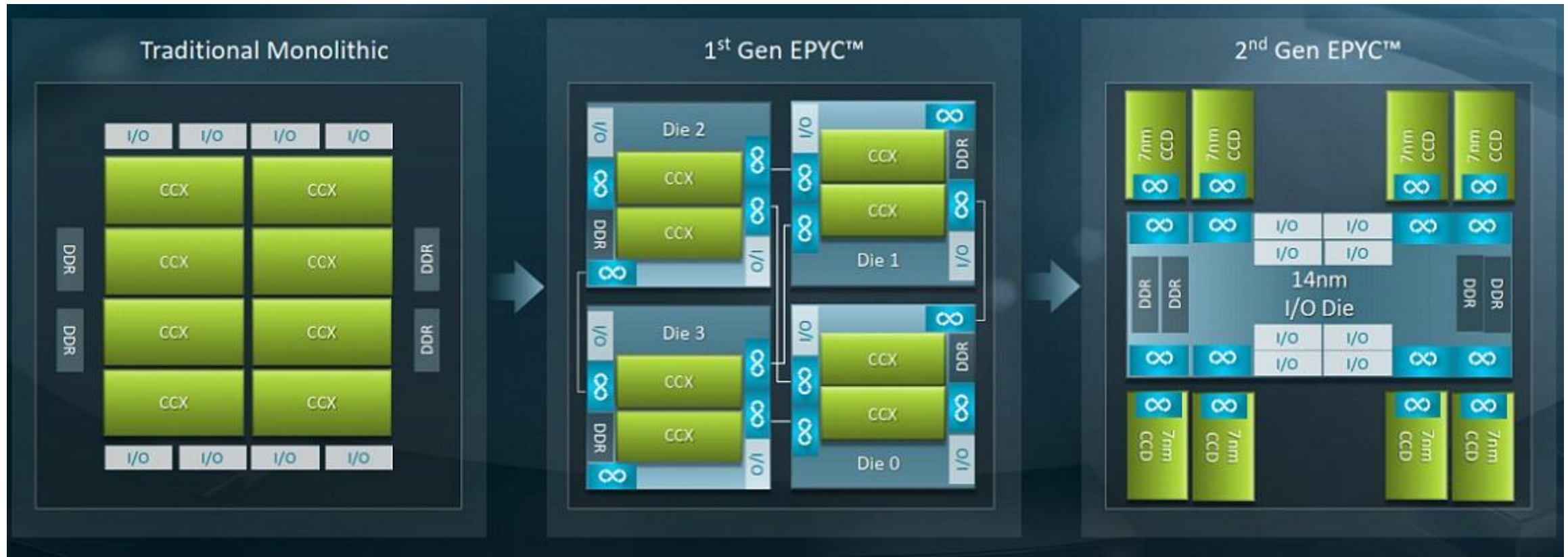
Purpose-built MCM links optimized for power, bandwidth, and latency



AMD 2nd Generation EPYC Rome (2019)

2nd gen. EPYC with Zen2 architecture

- ❑ 9 7nm dies
- ❑ 8 7nm Complex Core Die (CCD) chiplets, 8 cores each
- ❑ A 14nm IO die
- ❑ Connected via second-gen infinity fabric



AMD Zen2 Rome Chiplets

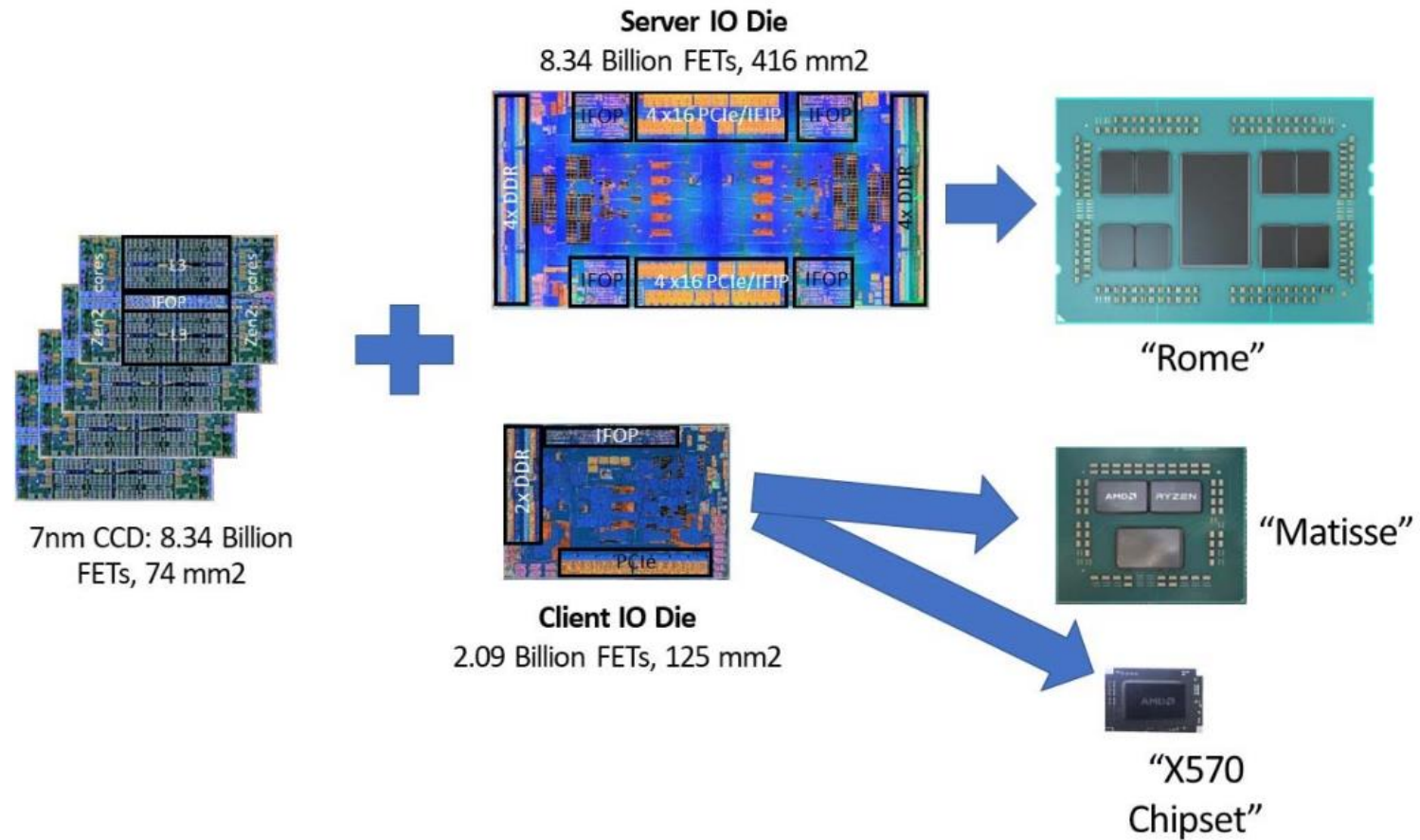


Figure 2.2.1: Three heterogeneous technology chiplets leveraged to many products and markets.

AMD Zen2 Cost-performance scalability

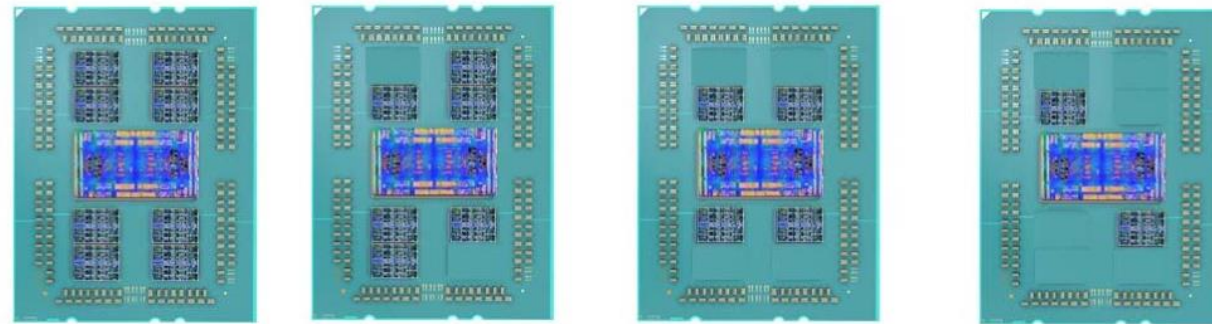
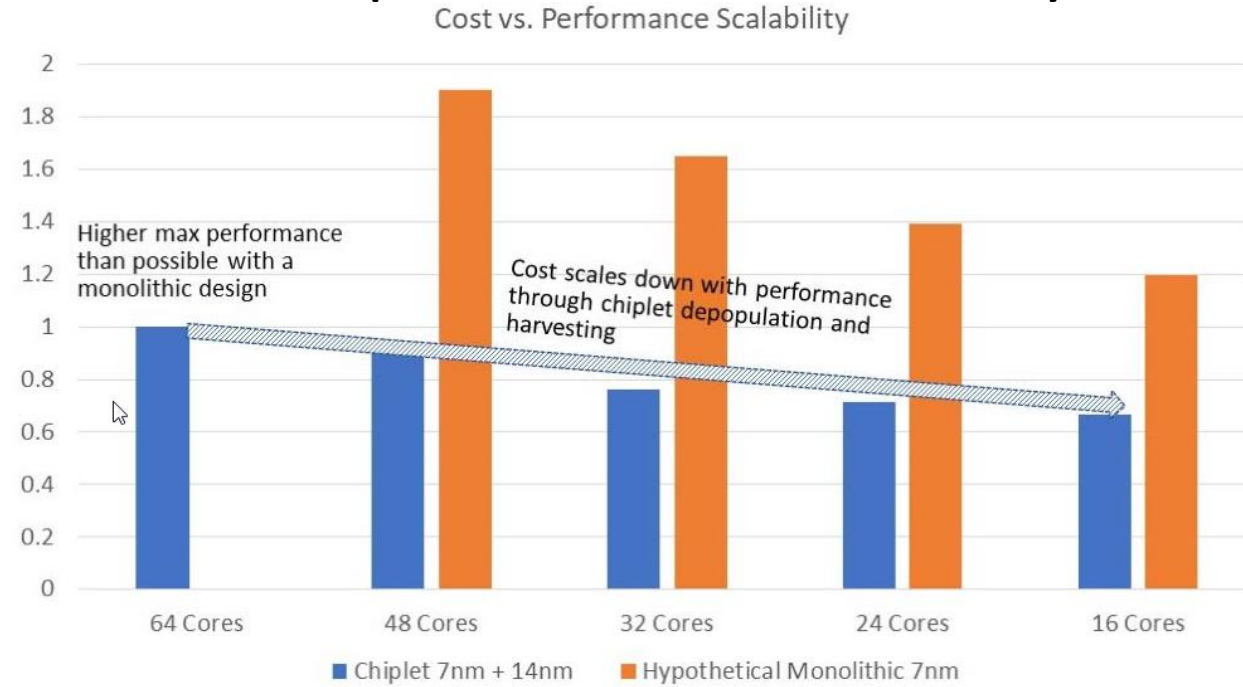


Figure 2.2.2: Cost-performance scalability with chiplet design.

AMD Zen2 Infinity On-Package (IFOP) SerDes Architecture

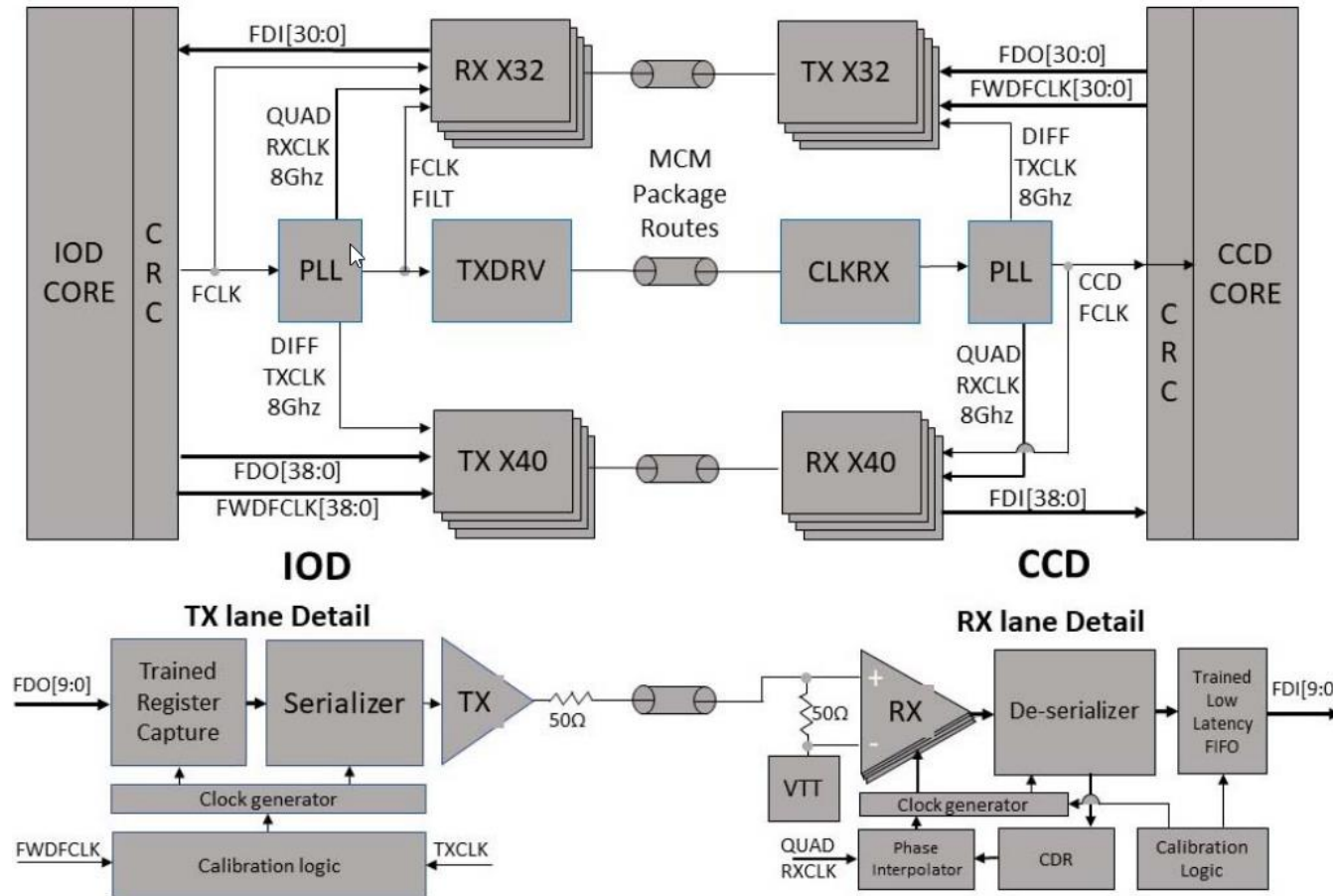


Figure 2.2.3: Infinity Fabric On-Package (IFOP) SerDes architecture.

AMD Zen2 Infinity On-Package (IFOP) SerDes Architecture

“Rome”

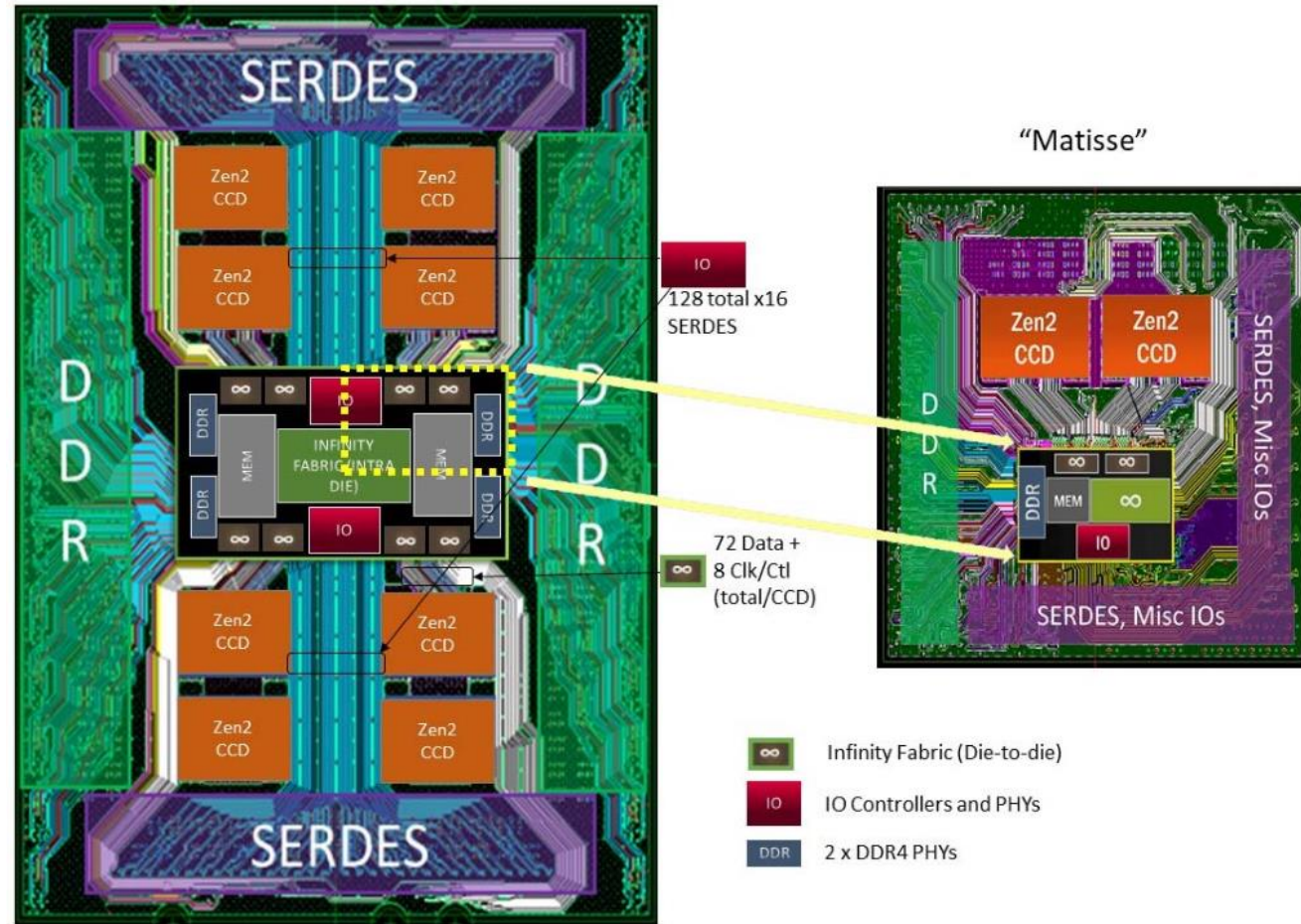


Figure 2.2.4: 'Rome' and 'Matisse' package design and IOD leverage.

AMD Naples vs. Rome

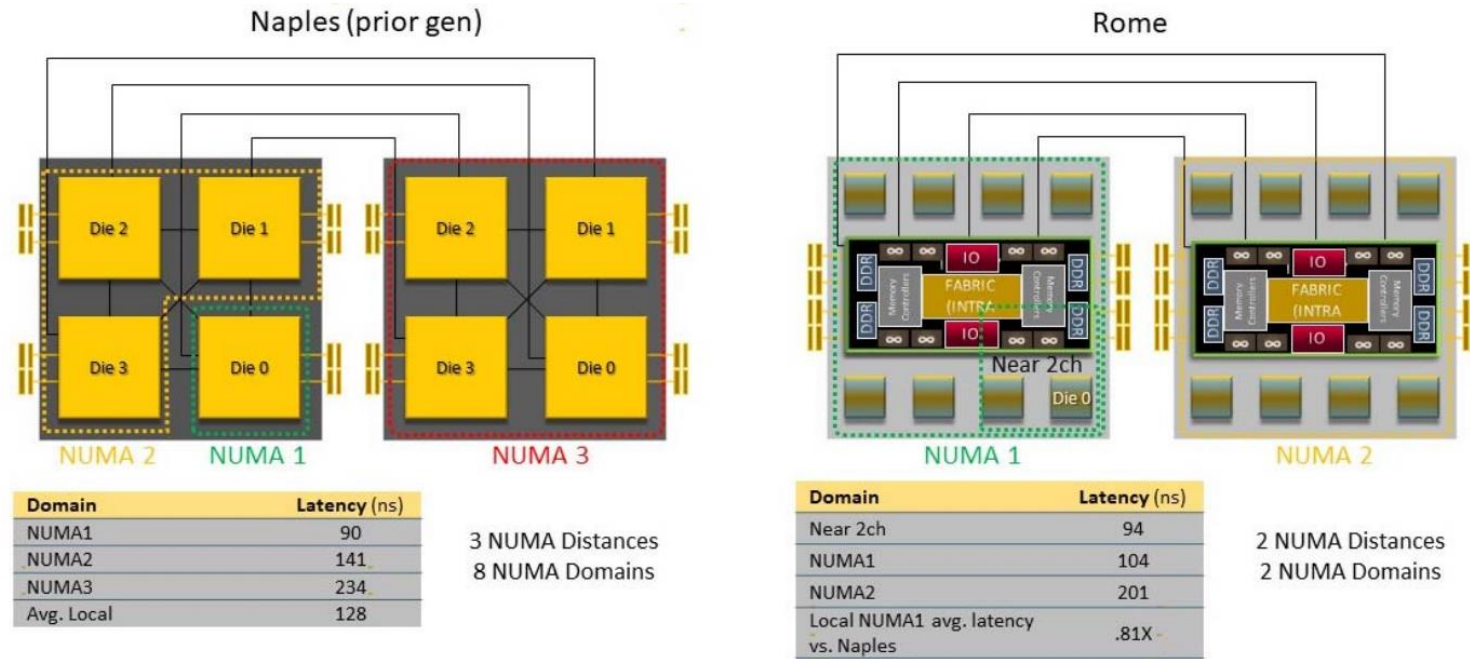


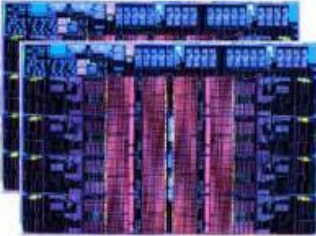

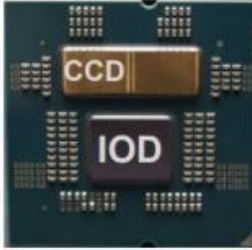
Figure 2.2.6: 'Rome' central IOD reduces the number of NUMA domains and distances for much improved memory latency attributes relative to it's predecessor.

AMD Zen4 line-up

Product Configurations

Desktop

1 or 2 "Zen 4" CCDs + 6nm IO Die = Ryzen™ 7000 Series

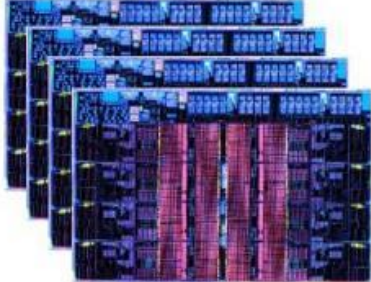
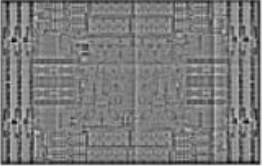

12.4 mm x 9.5 mm
~3.37 Billion Transistors

6-16 Cores
Boost up to 5.7 GHz

Up to 16.5 B transistors per socket

Server

2-12 "Zen 4" CCDs + 6nm IO Die = EPYC™ 9004 Series

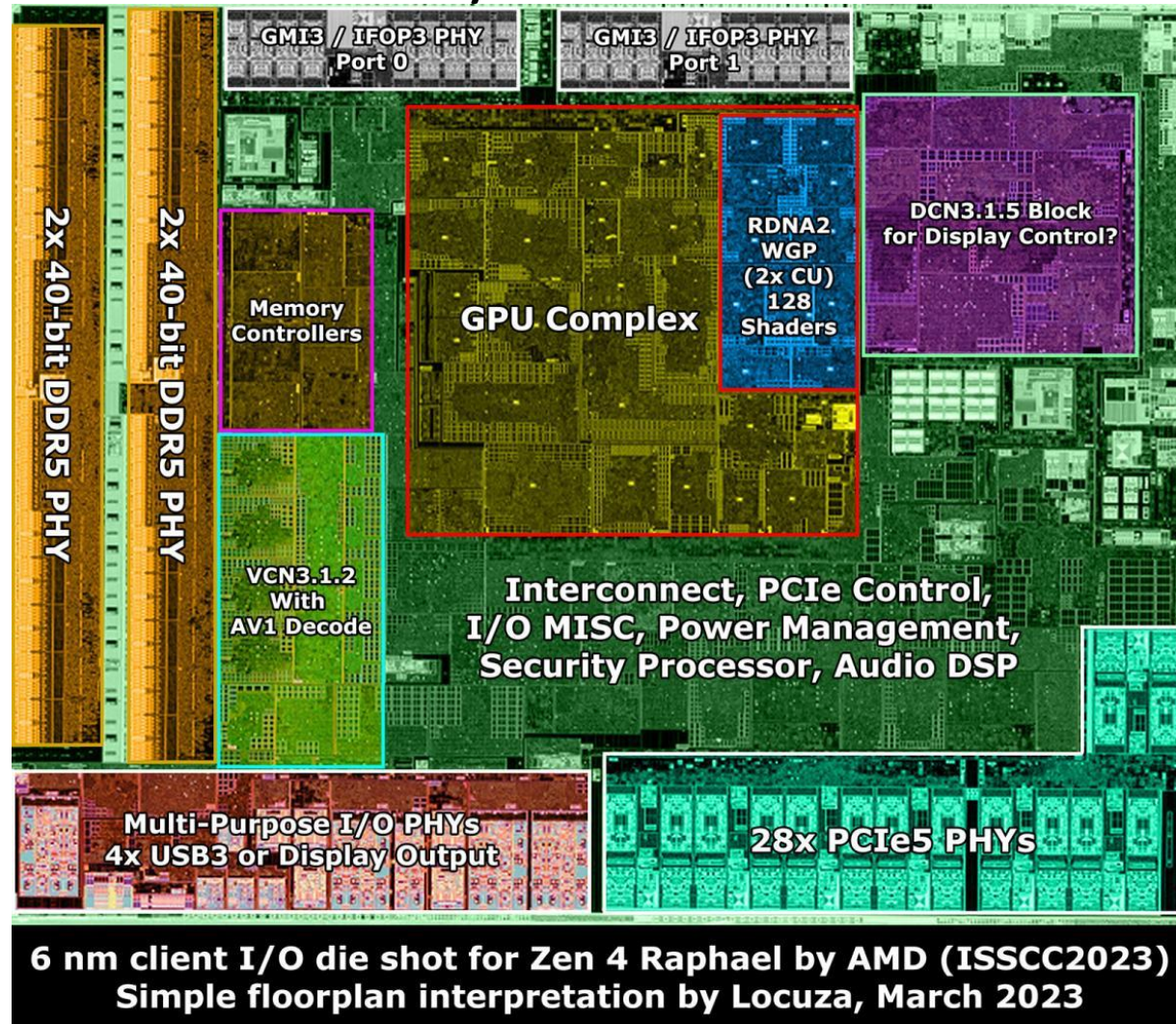




24.8 mm x 15.6 mm
~11 Billion Transistors

Up to 90.2B transistors per socket

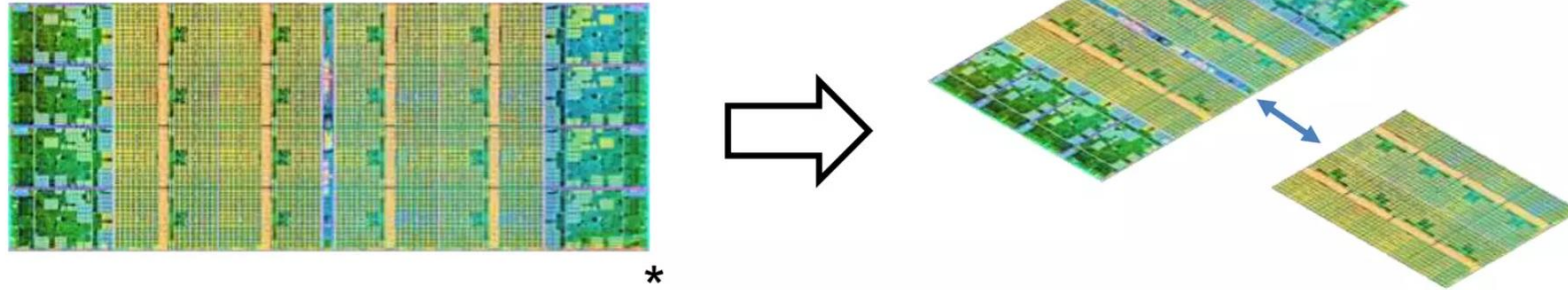
Product Positioning	Core Count	Default TDP	cTDP
Core Performance	16-48	320-360 W	400 W
Core Density	64-96	280-360 W	400 W
Total Cost of Ownership	16-32	200-280 W	400 W

AMD Ryzen Zen4 IO die



2.5D chiplets, 3D chiplets

More than Moore

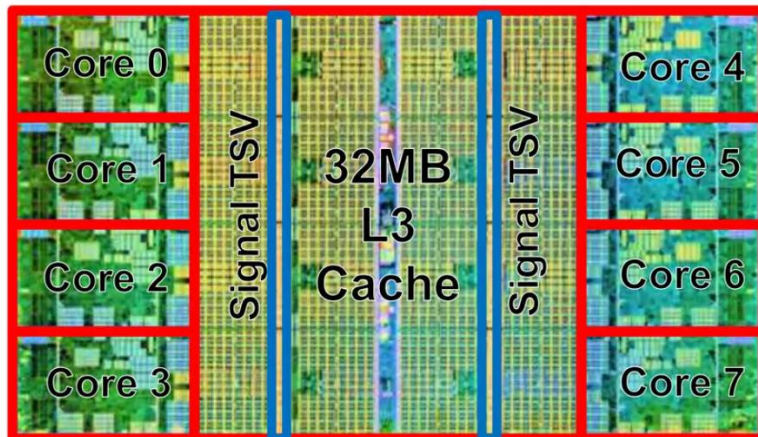
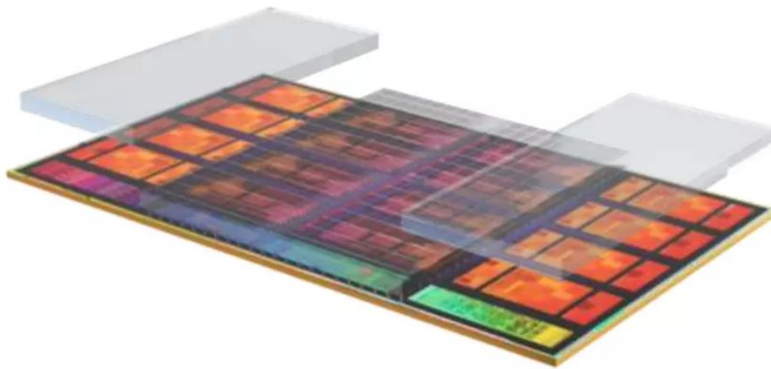


- **2.5D chiplets can provide product flexibility and reduce cost**
- **However, 3D can be even better!**
 - Improves effective memory latency
 - Reduces long datapath and I/O's dynamic powers
 - Fits more transistors within a given package cavity size

*Hypothetical processor with large cache

AMD Ryzen Zen3 3D V-cache

AMD 3D V-Cache™ Components: CCD

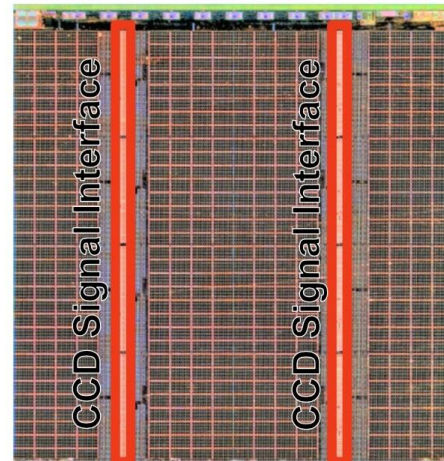
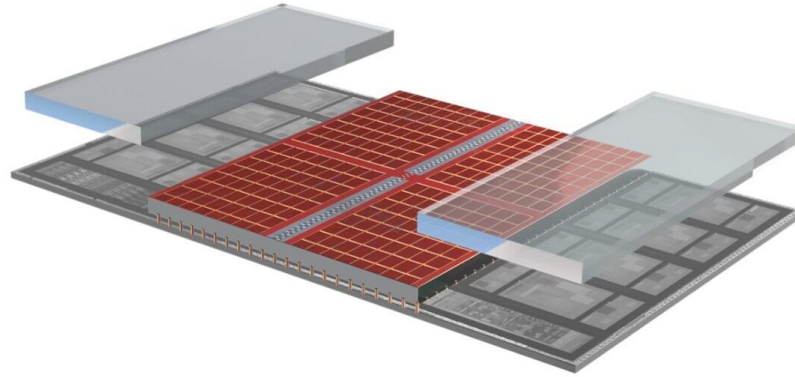


- “Zen 3” x86-64 CPU Core Complex Die (CCD)
- TSMC 7nm technology
- 8 cores per Core Complex (CCX)
- 32MB shared L3 Cache
- +19%¹ IPC (Ave) vs. “Zen 2”
- 81mm²
- AMD 3D V-Cache™ support integrated from Day 1

¹SEE ENDNOTES: R5K-003

AMD Ryzen Zen3 3D V-cache

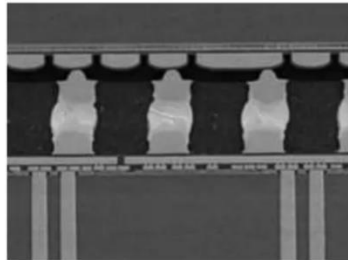
AMD 3D V-Cache™ Components: L3D



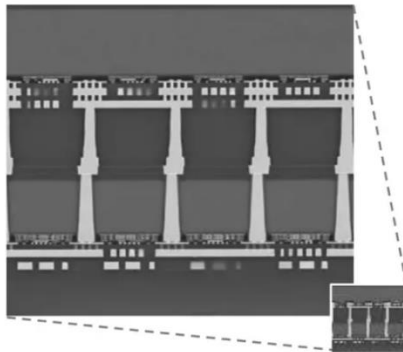
- **AMD 3D V-Cache™ extended L3 Die (L3D)**
- **TSMC 7nm FinFET Technology**
- **13 layers Cu + 1 layer Al metal stack**
- **64MB L3 Cache Extension**
- **41mm²**

AMD 3D V-cache Interface

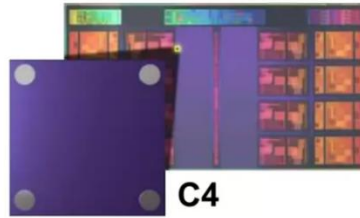
Micro Bump vs. Hybrid Bond



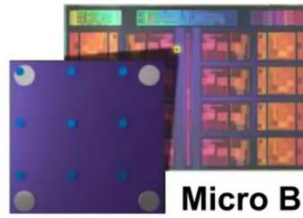
Micro Bump 3D



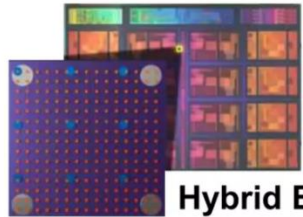
Hybrid Bond 3D



C4



Micro Bump 3D



Hybrid Bond 3D ^[1]

■ Compared to Micro Bump 3D solutions, Hybrid Bond offers

- >15x interconnect density
- >3x interconnect energy efficiency
- Superior thermal conductance

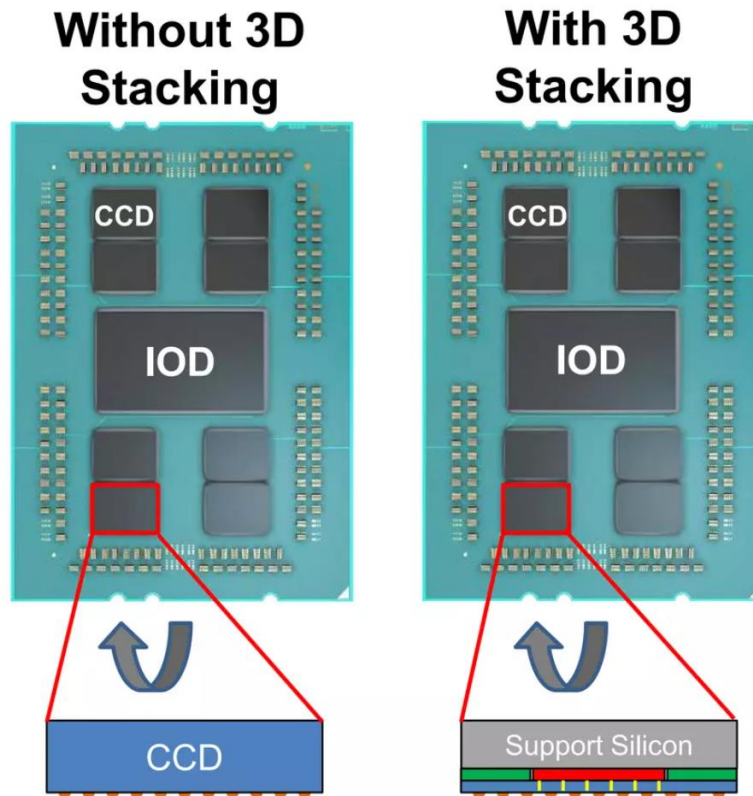
SEE ENDNOTES: EPYC-027

C4 and Micro Bump 3D illustrations are hypothetical

[1] Swaminathan, Hot Chips Tutorial, 2021

AMD 3D V-cache Interface

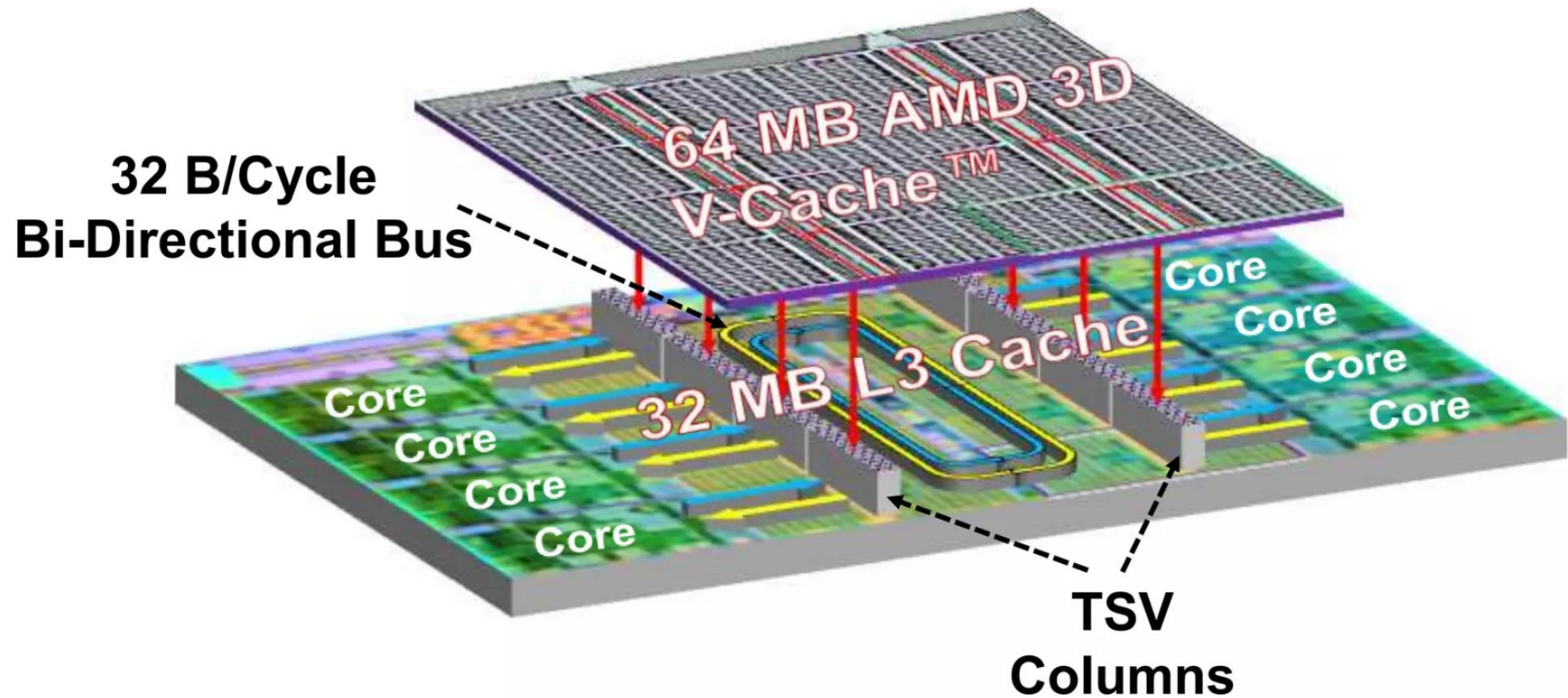
Server Configurations



- **AMD 3rd Gen EPYC™ Server CPU**
 - Up to 8 "Zen 3" CCDs
 - 1 I/O Die (IOD)
- **AMD 3rd Gen EPYC™ Server CPU with AMD 3D V-Cache™**
 - Up to 8 thinned CCDs + L3Ds
 - Support silicon added to match 2D CCD Z-height
- **Both designs compatible with the same package**

AMD 3D V-cache Interface

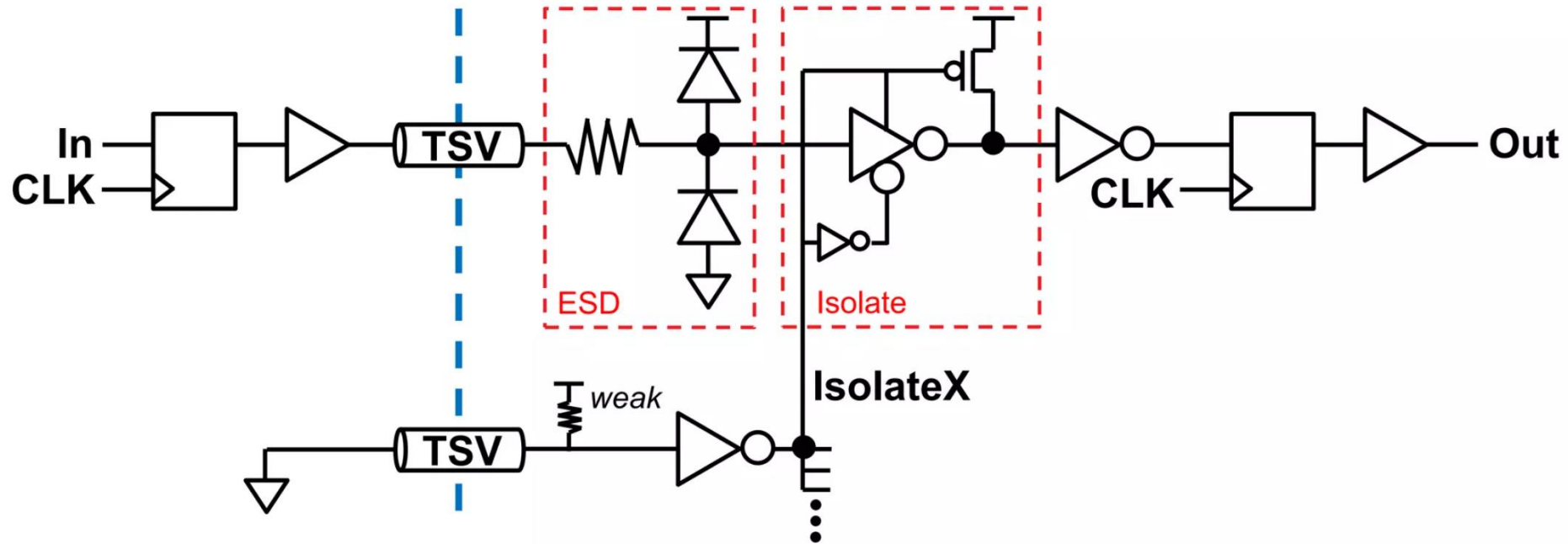
Cache Interface Illustration



[5]

AMD Hybrid-Bonded 64MB Stacked Cache

3D Signal Interface



- **Simple digital signal interface between dies**
 - Enabled by HB technology's low parasitics

AMD 64MB Stacked Cache Performance

Server Performance



Intel Tiles with EMIB and Foveros

Intel Lakefield (2020): Foveros Die Interface (FDI) die stacking

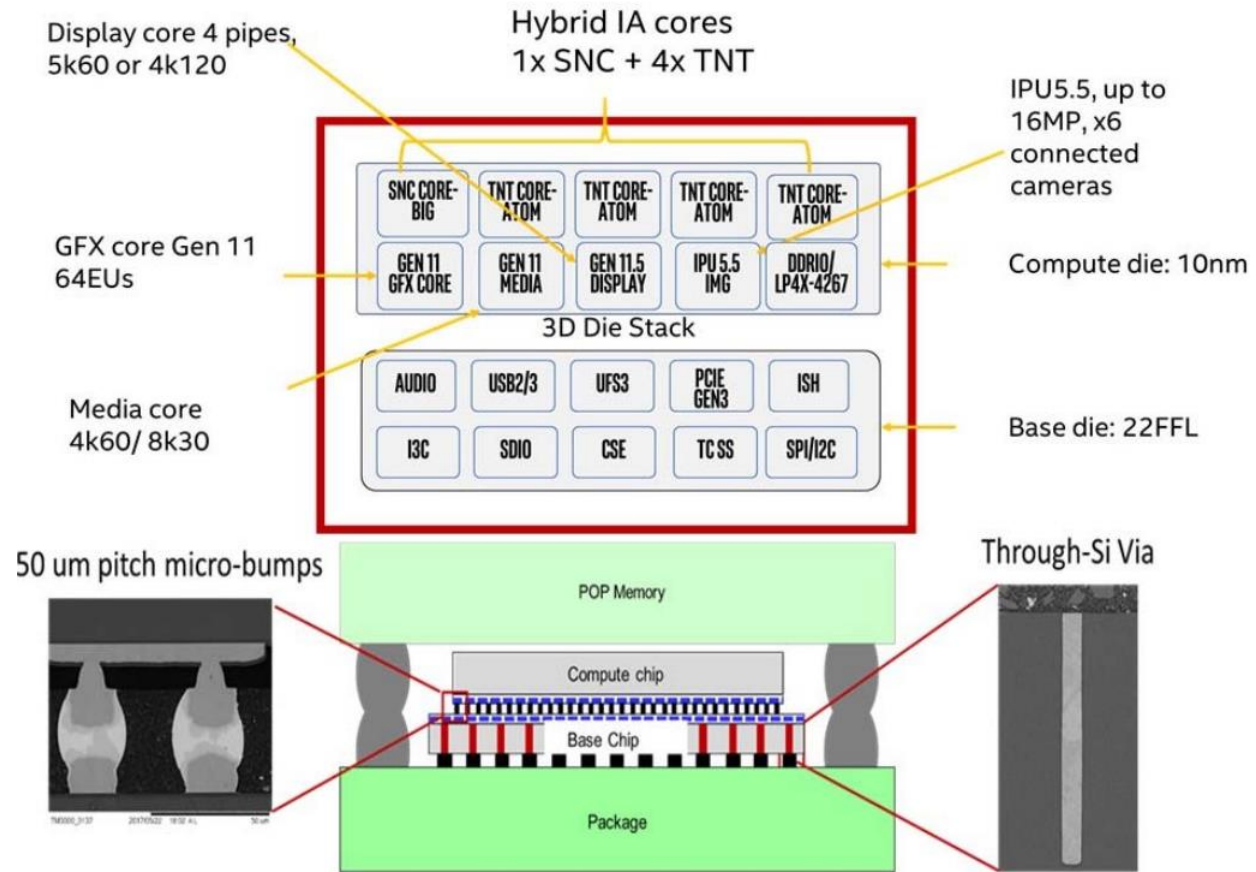


Figure 8.1.1: System partitioning with 3D stacking across compute and base die.

Intel Lakefield (2020): Foveros Die Interface (FDI) die stacking

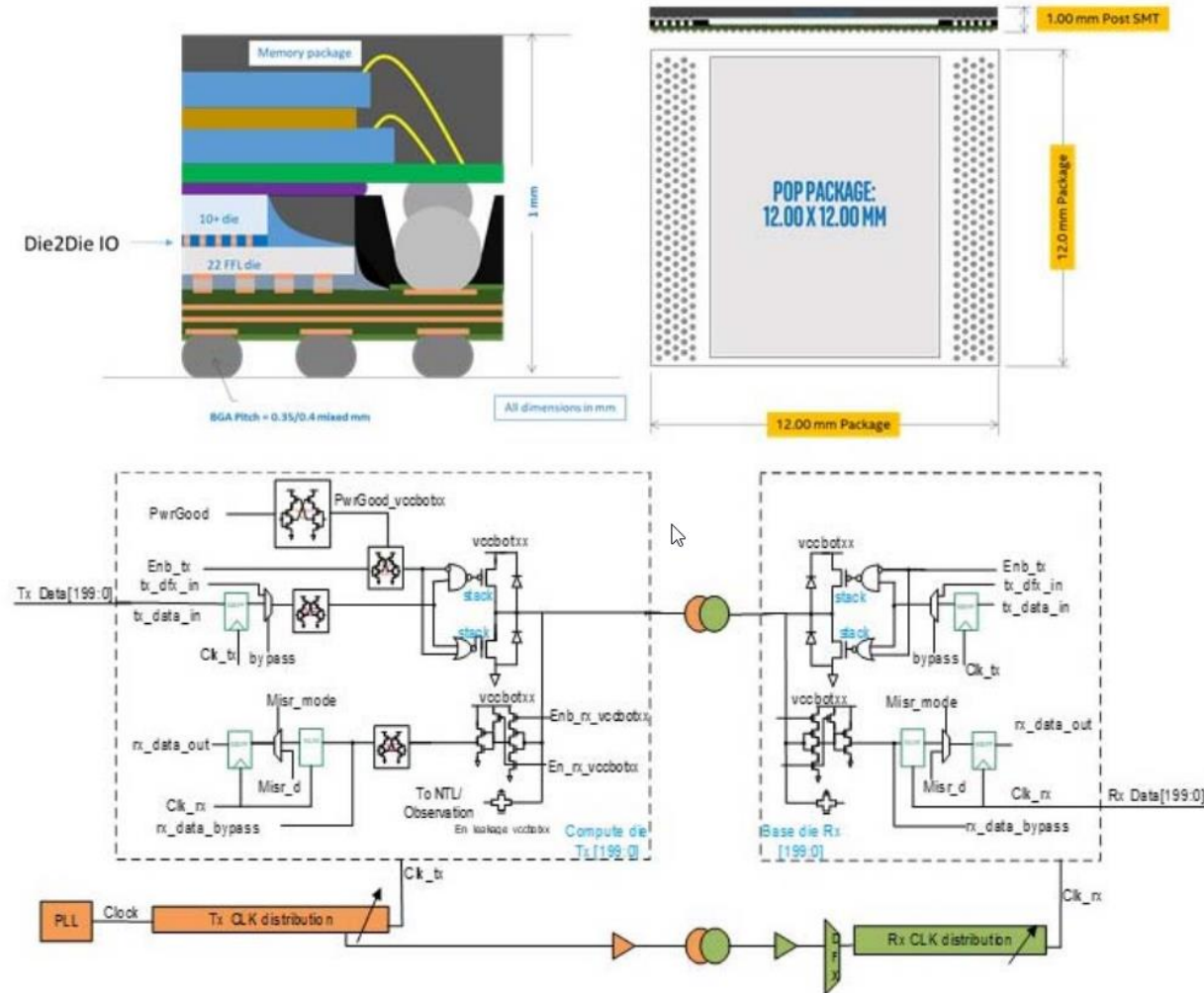
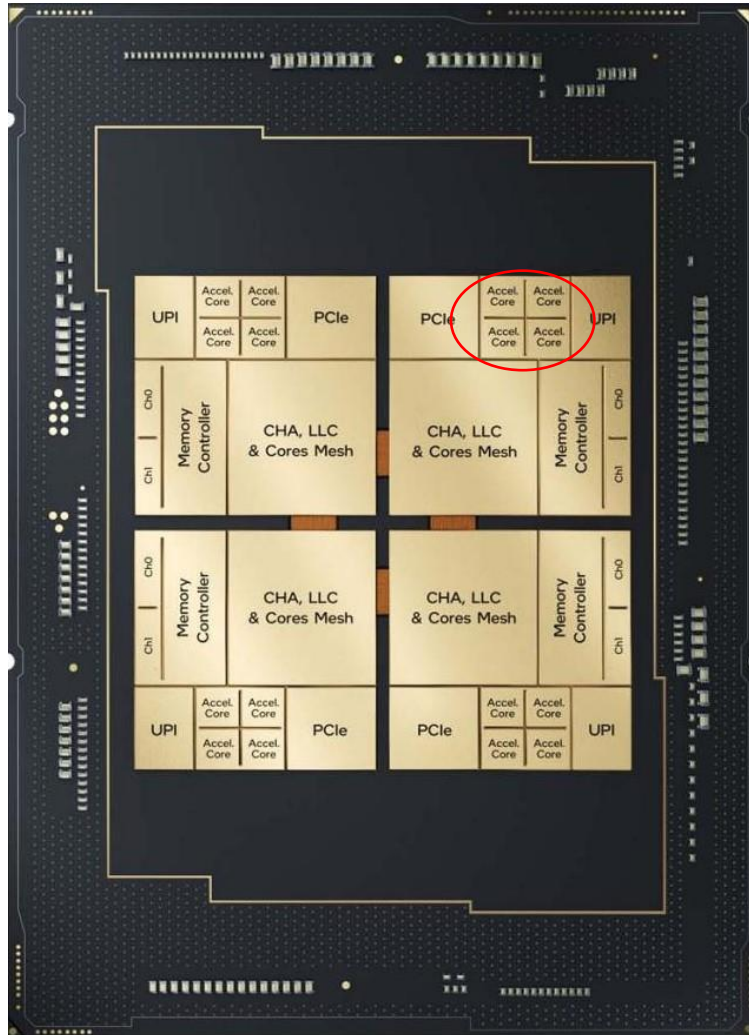


Figure 8.1.3: Simplified Die2Die IO circuits.

Intel Sapphire Rapids – Multi-Tile Design



Intel Sapphire Rapids (2023)



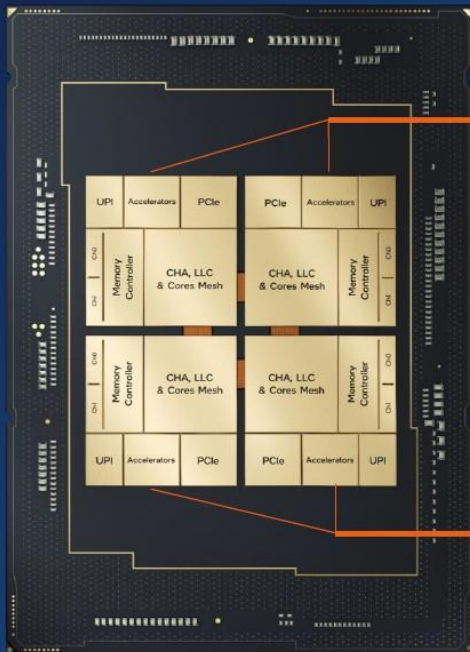
Sapphire Rapids Accelerators

- Intel QuickAssist Technology (QAT)
 - Faster compression and encryption
- Intel Dynamic Load Balancer (DLB)
 - Load balancing
 - Queue management
 - Packer prioritization
- In-Memory Analytics Accelerator (IAA)
 - In-memory databases
 - Big data analytics
- Data Streaming Accelerator (DSA)
 - High performance data copy and transformation

Intel Sapphire Rapids (2023)

Enabling the Intel® Accelerator Engines

Tools for developers to take advantage and deploy today



IAA, DSA,
QAT, DLB

IAA, DSA,
QAT, DLB



Intel® Advanced Matrix Extensions (Intel® AMX)

- TensorFlow
- PyTorch
- ONNX Runtime
- OpenVINO
- oneDNN (Intel oneAPI)



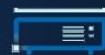
Intel® Advanced Vector Extensions (Intel® AVX) for vRAN

- FlexRAN
- Data Plane dev Kit (DPDK)*



Intel® In-memory Analytics Accelerator (Intel® IAA)

- Intel Query Processing Library



Intel® Data Streaming Accelerator (Intel® DSA)

- Storage Perf Dev Kit (SPDK)*
- Data Plane Dev Kit (DPDK)*



Intel® QuickAssist Technology (Intel® QAT)

- QATzip* (Intel lib)
- OpenSSL**
- Boring SSL



Intel® Dynamic Load Balancer (Intel® DLB)

- VPP IPsec
- Data Plane Dev Kit (DPDK)*



Intel Sapphire Rapids (2023)

A Higher Performance Server Architecture

Benefits of Intel® Accelerator Engines

Intel® Advanced
Matrix Extensions
(Intel® AMX)

Up to

8.6x

higher speech recognition
inference performance
with built-in AMX BFI6 vs.
FP32

Intel® Dynamic
Load Balancer
(Intel® DLB)

Up to

96%

lower latency
at the same throughput for
Istio-Envoy Ingress with Intel®
DLB vs. software for Istio
Ingress gateway

Intel® Data
Streaming
Accelerator
(Intel® DSA)

Up to

1.7x

higher IOPs for SPDK-
NVMe with built-in Intel®
DSA vs. ISA-L software

Intel® In-Memory
Analytics
Accelerator
(Intel® IAA)

Up to

2.1x

higher RocksDB
performance with Intel® IAA
vs Ztsd software

Intel® QuickAssist
Technology
(Intel® QAT)

Up to

84%

fewer cores to achieve
same connections/s on
NGINX with built-in QAT
vs. out-of-box software

Accelerators Enable Step Function Performance Beyond Base Architecture

Intel Sapphire Rapids – EMIB

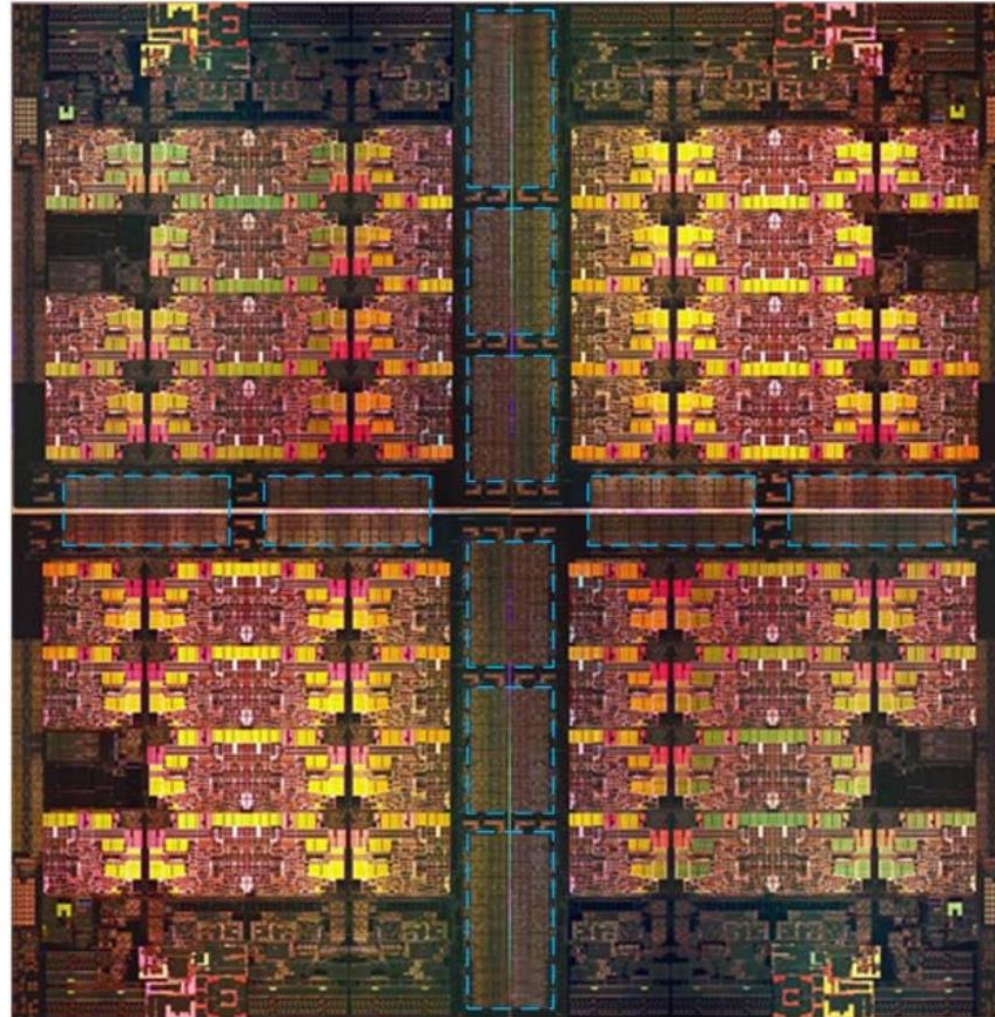


Figure 2.2.7: Die photo of left and right die arranged in 2x2 quasi-monolithic configuration. EMIB placement highlighted in blue.

Intel Sapphire Rapids Multi-Die Fabric (MDF) IO

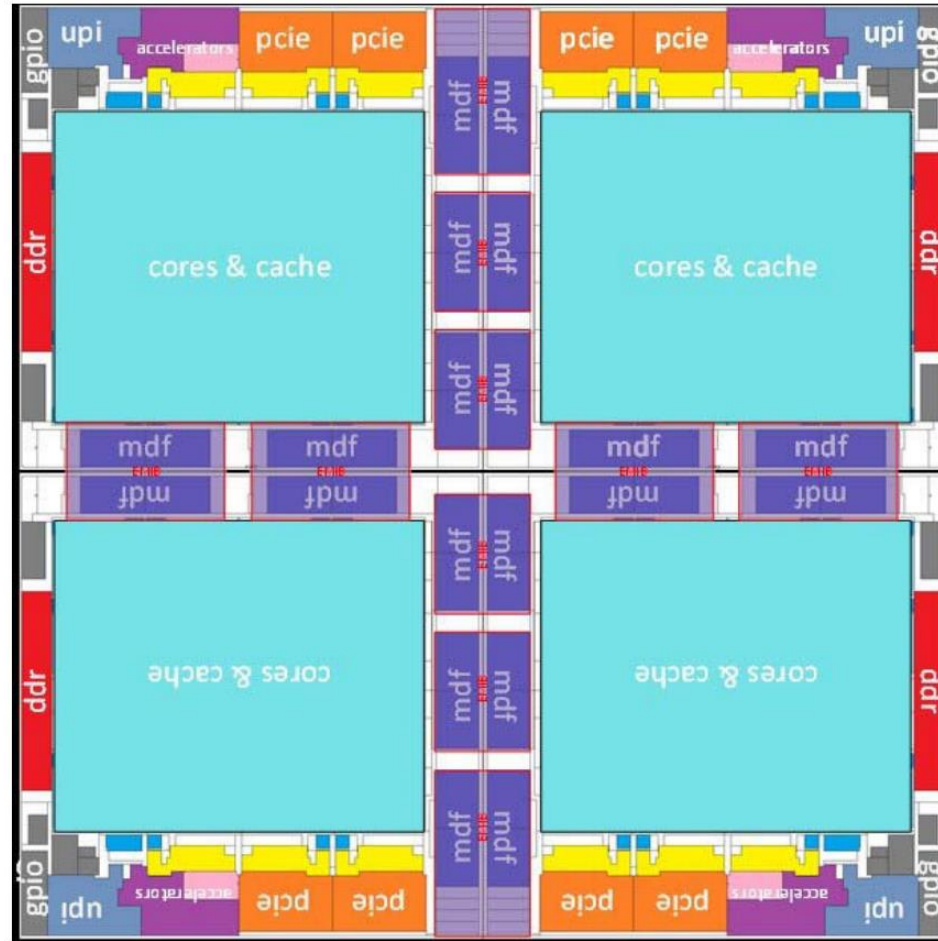


Figure 2.2.1: Die floorplan in 2x2 quasi-monolithic configuration. EMIB highlighted at die-to-die interfaces.

Intel Sapphire Rapids Cross-die timing model

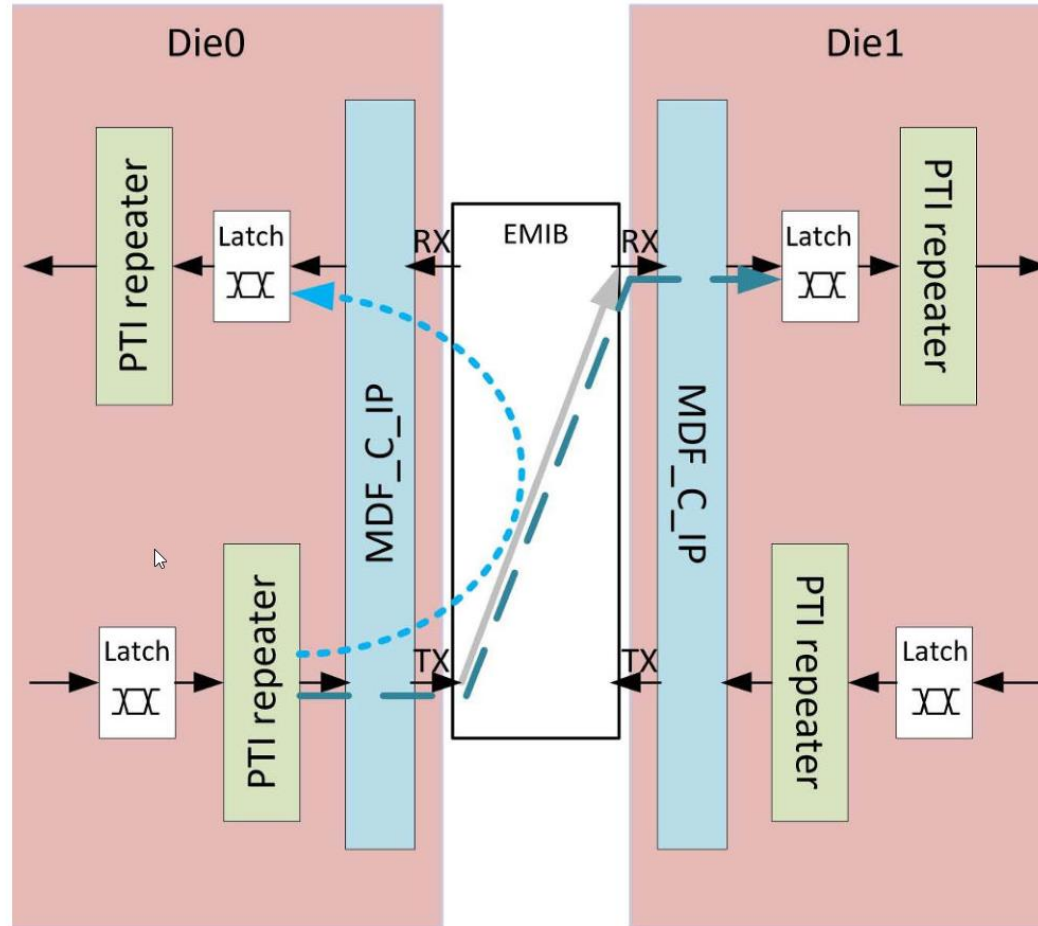


Figure 2.2.2: A cross-die timing model, or a single-die loopback model.

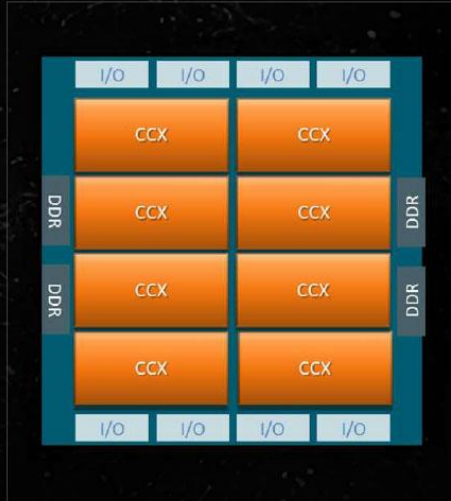


Chipelets-augmented GPUs

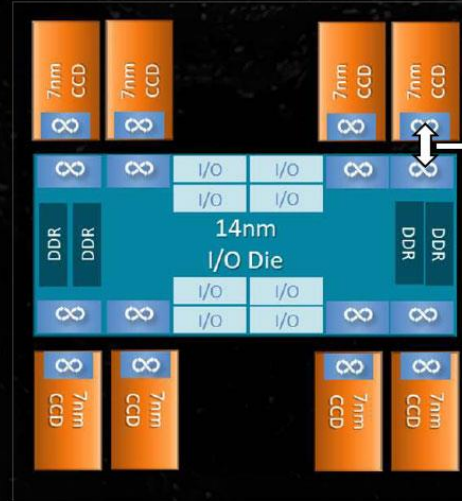
AMD RDNA3 – Chiplets for GPUs

CHIPLET TECHNOLOGY CAN IT WORK FOR GRAPHICS?

Traditional Monolithic

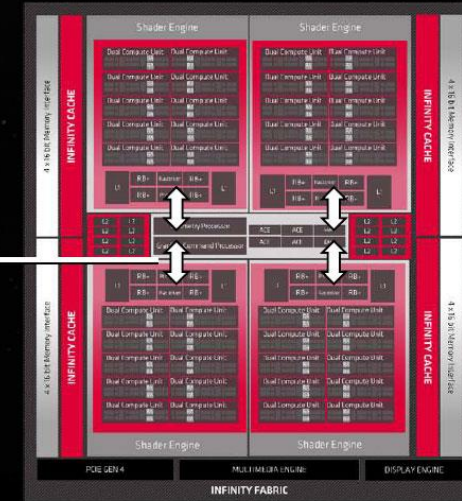


EPYC CPU Server



100's of signals

"Navi21" GPU



10's of 1000's of signals

- Chiplets enabled use of advanced nodes where they benefit CPU performance but mature nodes for IO and interfaces
- High speed organic package links meet CPU Bandwidth requirements

- GPU shader engines require massive amounts of connectivity compared to CPUs
- A different approach is required

AMD RDNA3 – Graphics Compute Die and Memory Cache Die

- Each MCD has 16MB of cache



The diagram shows a square chiplet die layout. In the center is a large square labeled 'GCD' (Graphics Compute Die) with a 5nm process technology. Surrounding the GCD are six smaller squares labeled 'MCD' (Memory Cache Die) with a 6nm process technology, arranged in two rows of three. The entire die is surrounded by a grid of pins. The AMD logo is visible in the bottom left corner of the die image.

Chiplets

The right process technology for the right job

300 mm ²	6x37 mm ²
5nm GCD	6nm MCD

AMD RDNA3

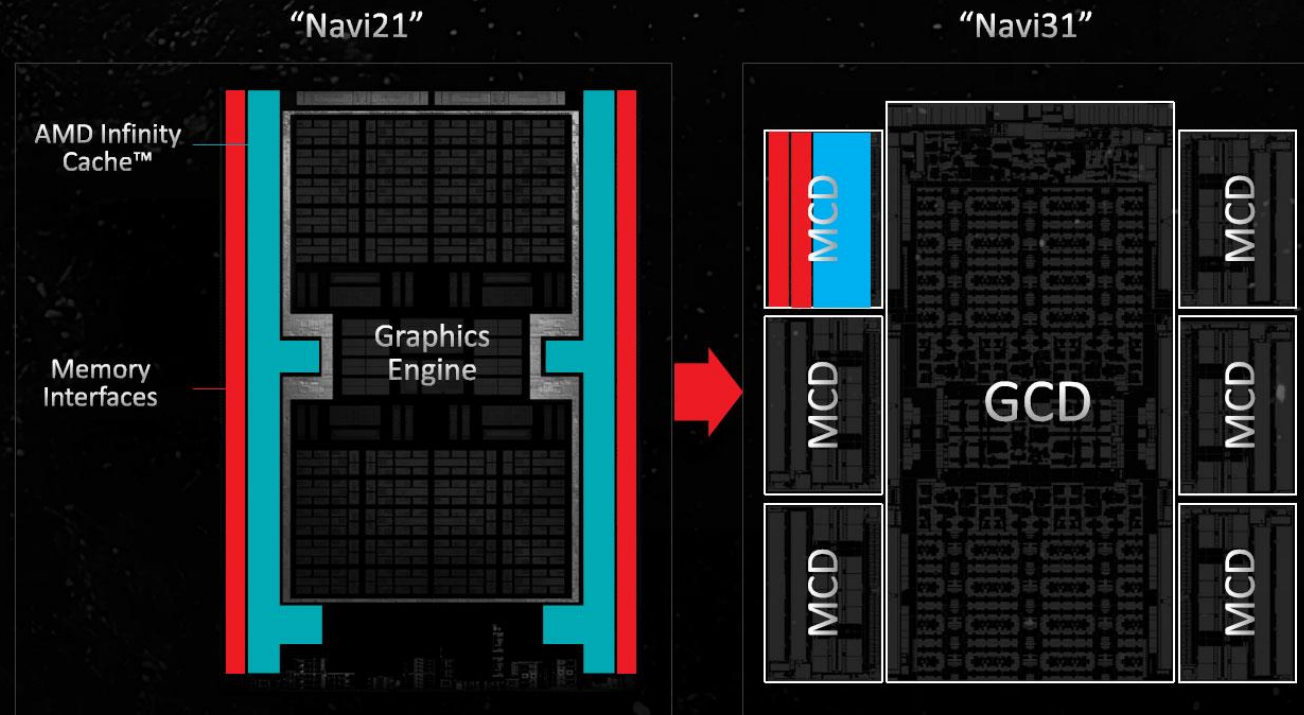
CHIPLET TECHNOLOGY A BETTER WAY TO PARTITION

The graphics engine is what benefits from advanced N5 technology

- AMD Infinity Cache™ critical to performance but barely shrinks into N5
- GDDR6 interfaces are also large and won't shrink at all

Split those poorly scaling components off as a chiplet and shrink the GFx core into N5

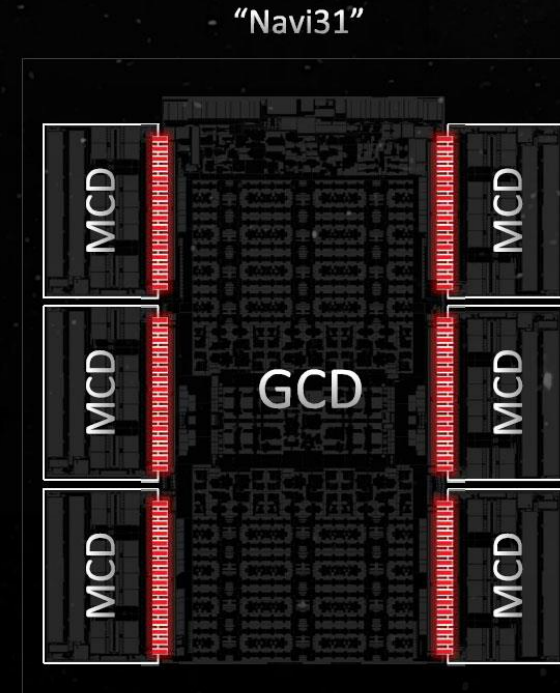
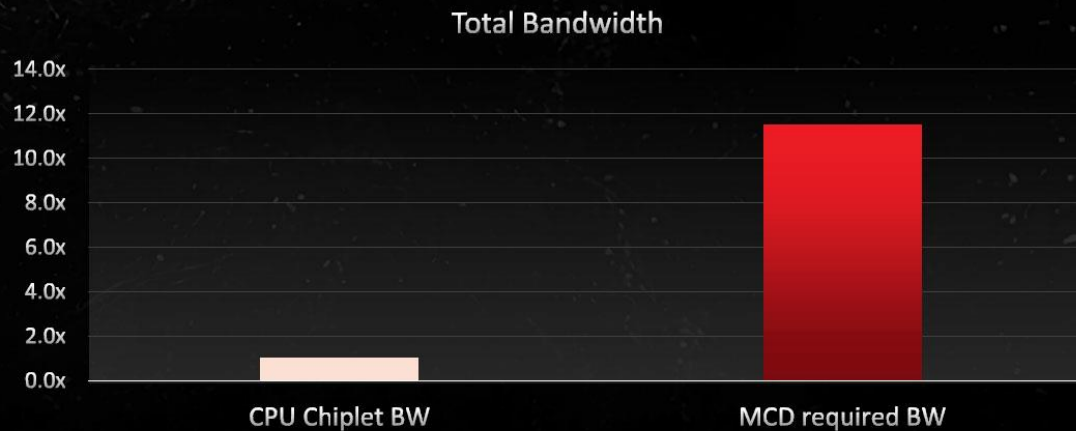
Full N5 performance, better yield for perf/\$ and configurability



AMD RDNA3 Chiplets

CHIPLET TECHNOLOGY HOW TO CONNECT THE CHIPLETS?

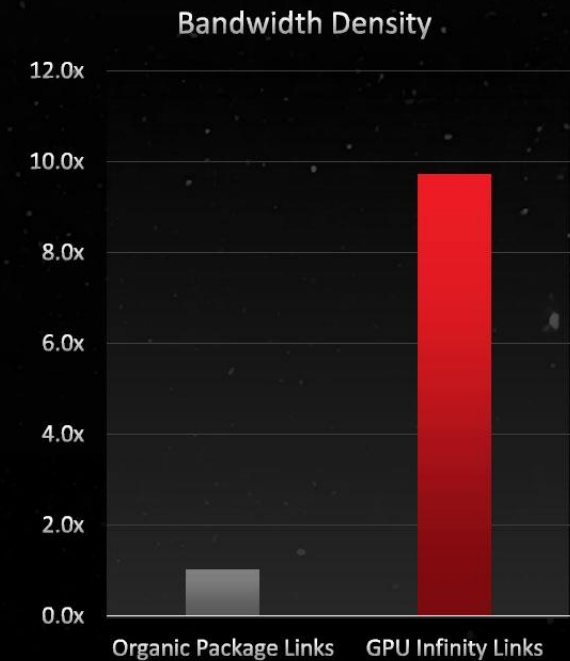
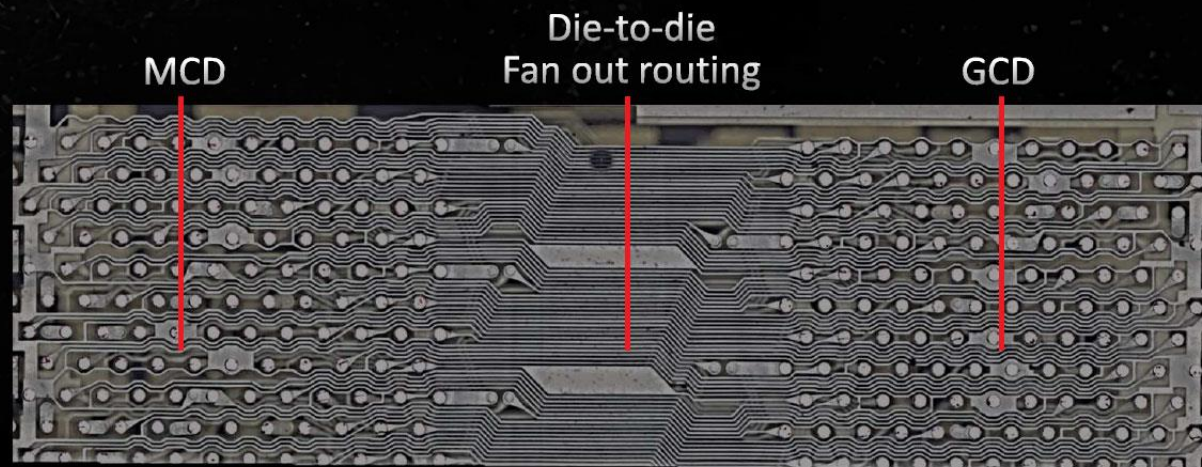
- GCD-MCD partitioning is great, but the bandwidth requirements are still extremely high
- Over 10X what a CPU CCD requires in EPYC
- Breakthrough Advanced packaging and a new interface is required:
 - **High Performance Fanout and Infinity links**



RDNA3 Infinity Links

CHIPLET TECHNOLOGY INFINITY FANOUT LINKS BANDWIDTH DENSITY

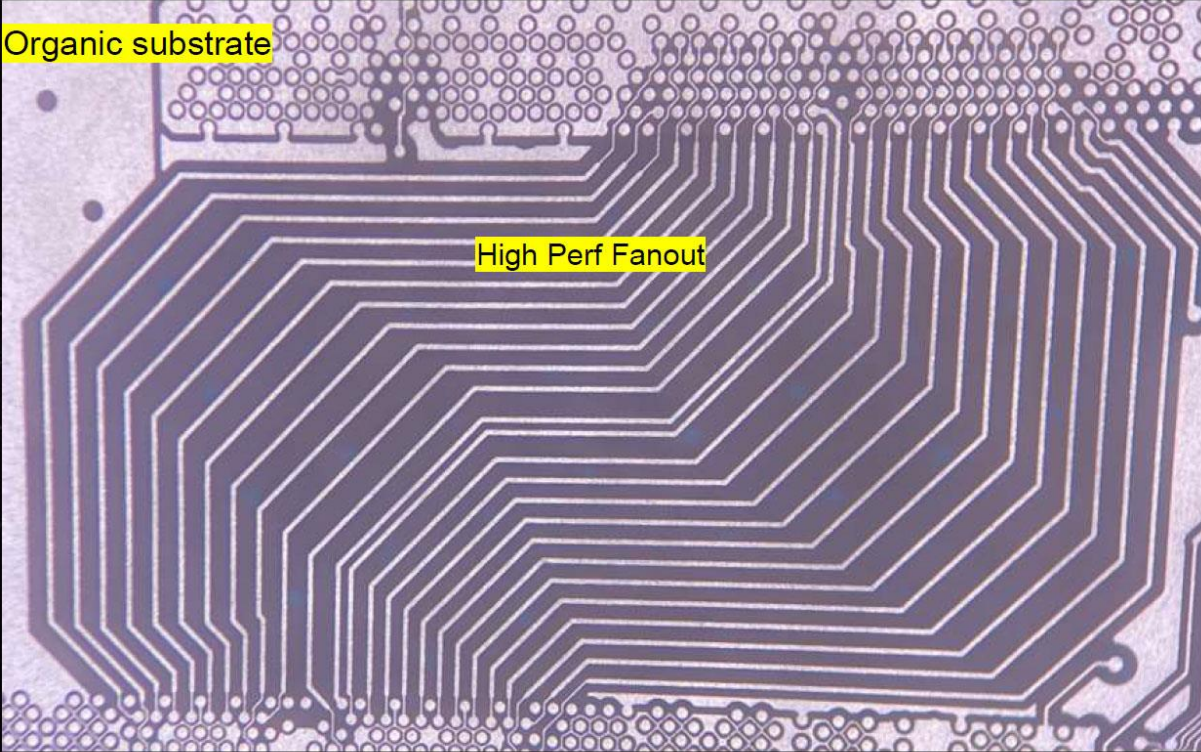
- Infinity Links, operating at 9.2Gb/s with High Performance Fanout provide almost 10X the BW density of the IFOP links used in Ryzen and EPYC
- Enables industry-leading peak bandwidth of 5.3TB/s



RDNA3 Infinity Links

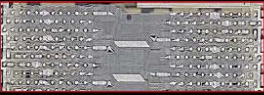
CHIPLET TECHNOLOGY

HIGH PERFORMANCE FANOUT INTERCONNECT



Organic substrate

High Perf Fanout



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

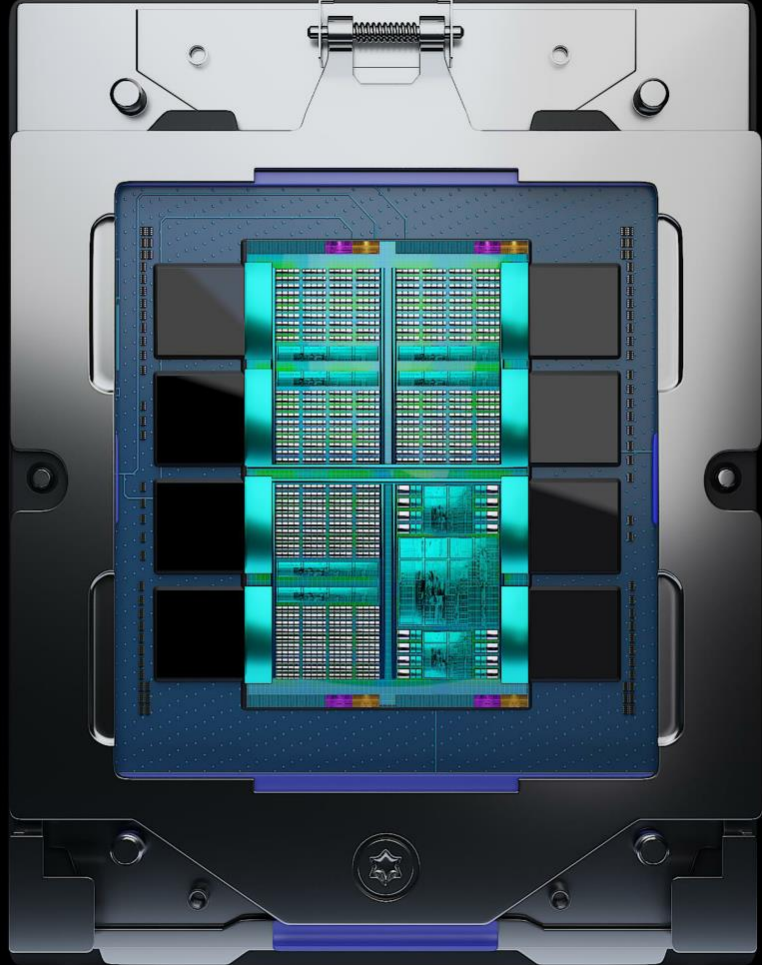
Images approximately to scale

AMD
together we advance_gaming

16



AMD Instinct MI300



The world's first integrated
data center CPU + GPU

AMD INSTINCT™
MI300

Breakthrough architecture to
power the exascale AI era

The image shows a detailed view of the AMD Instinct MI300 accelerator card. It is a large, rectangular card with a dark blue and black color scheme. The central part of the card is a large, square chip with a complex, grid-like pattern of components, illuminated with a bright cyan light. The card is mounted on a metal frame with various connectors and screws. The overall design is sleek and industrial.

AMD Instinct MI300

AMD INSTINCT™ MI300

The world's first data center integrated CPU + GPU


CDNA 3

Next-Gen
Accelerator
Architecture



24 Leadership
Data Center
CPU cores

146B

Transistors

128GB

HBM3

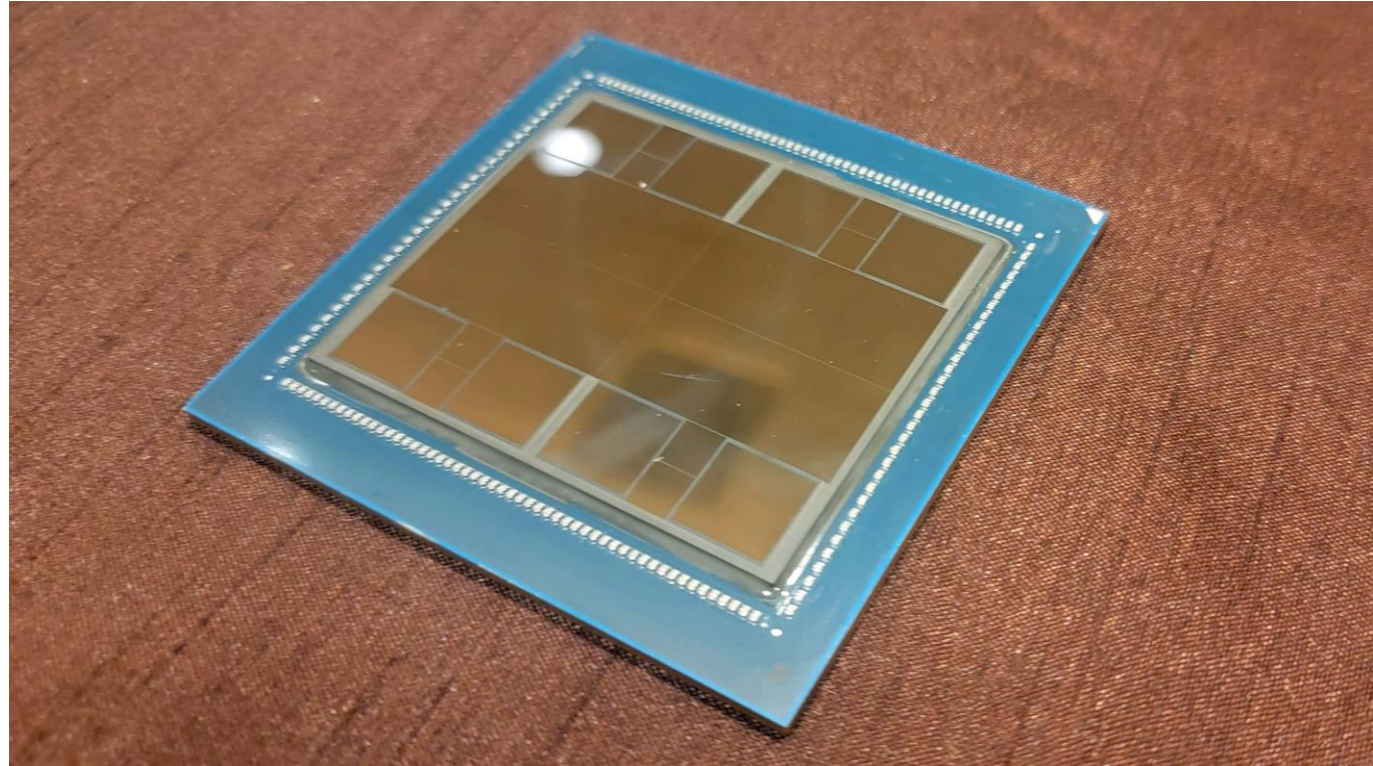
3D

Advanced Chiplet Packaging

AMD Instinct MI300

MI300 – a disaggregated design

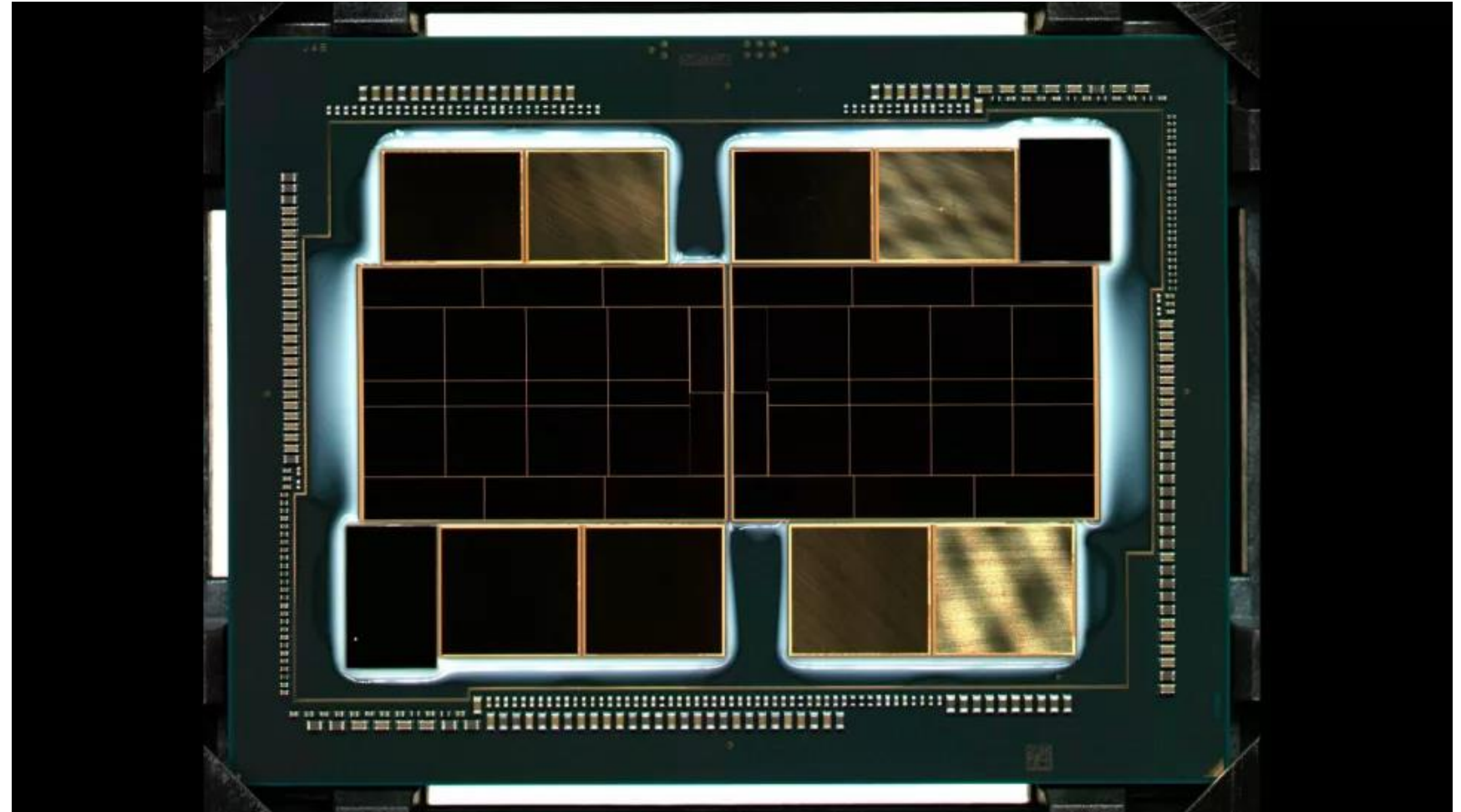
- ❑ Multiple TSMC 5nm chiplets
- ❑ 3D Stacking to place them over a base die
- ❑ On-package High Bandwidth Memory (HBM)



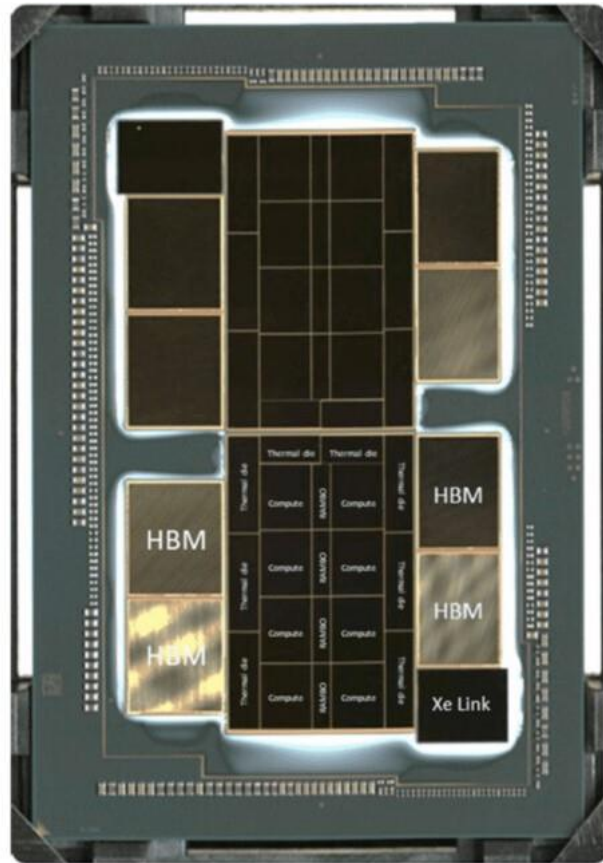
Intel Ponte Vecchio

Intel Ponte Vecchio

- ❑ Intel 10nm process
- ❑ Die Size of 1280mm²



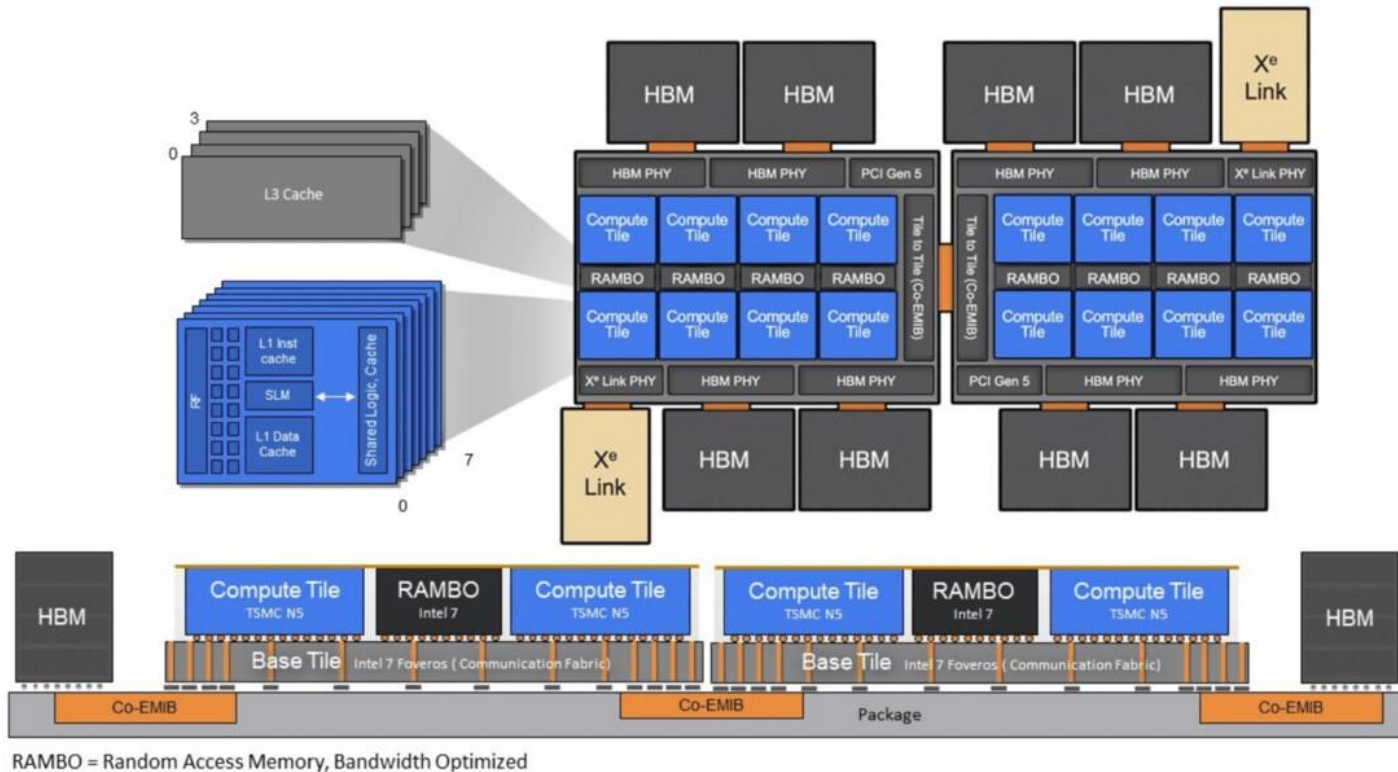
Intel Ponte Vecchio



Integration	Foveros + EMIB
Power Envelope	600W
Transistor count	> 100B
Total Tiles	63 (47 functional + 16 thermal Tiles)
HBM count	8
Package Form factor	77.5 x 62.5 mm (4844 mm ²)
Platforms	3 platforms
IO	4x16 90G SERDES, 1x16 PCIe Gen5
Total Silicon	3100 mm ² Si
Silicon footprint	2330 mm ² Si footprint
Package layers	11-2-11 (24 layers)
2.5D Count	11 2.5D connections
Resistance	0.15 mΩ R _{path} /tile
Package pins	4468 pins
Package Cavity	186 mm ² x4 cavities

St **Figure 2.1.7: Ponte Vecchio chip photographs and key attributes.**

Intel Ponte Vecchio



- ❑ 47 Tiles + 16 thermal tiles
 - ❑ 16 compute tiles, TSMC N5, 2.6TB/s speeds to the chip fabric
 - ❑ 8 tiles for RAMBO cache, Intel 7, 15MB per tile, 1.3 TB/s connection
 - ❑ 2 Foveros base tiles, Intel 7
 - ❑ 2 Xe-Link tiles, TSMC N7
 - ❑ 8 HBM2e tile
 - ❑ 11 Intel's embedded multi-die interconnect bridge (EMIB) tiles
- ❑ The package has 4844mm² with 4468 pins
 - ❑ 2330mm² of silicon for the 47 tiles
- ❑ Fully Integrated Voltage Regulators (FIVR)
- ❑ Compute Express Link (CXL) interface

Figure 2.1.1: 3D and 2D system partitioning with Foveros and EMIB on PVC.

Intel Ponte Vecchio Foveros and EMIB

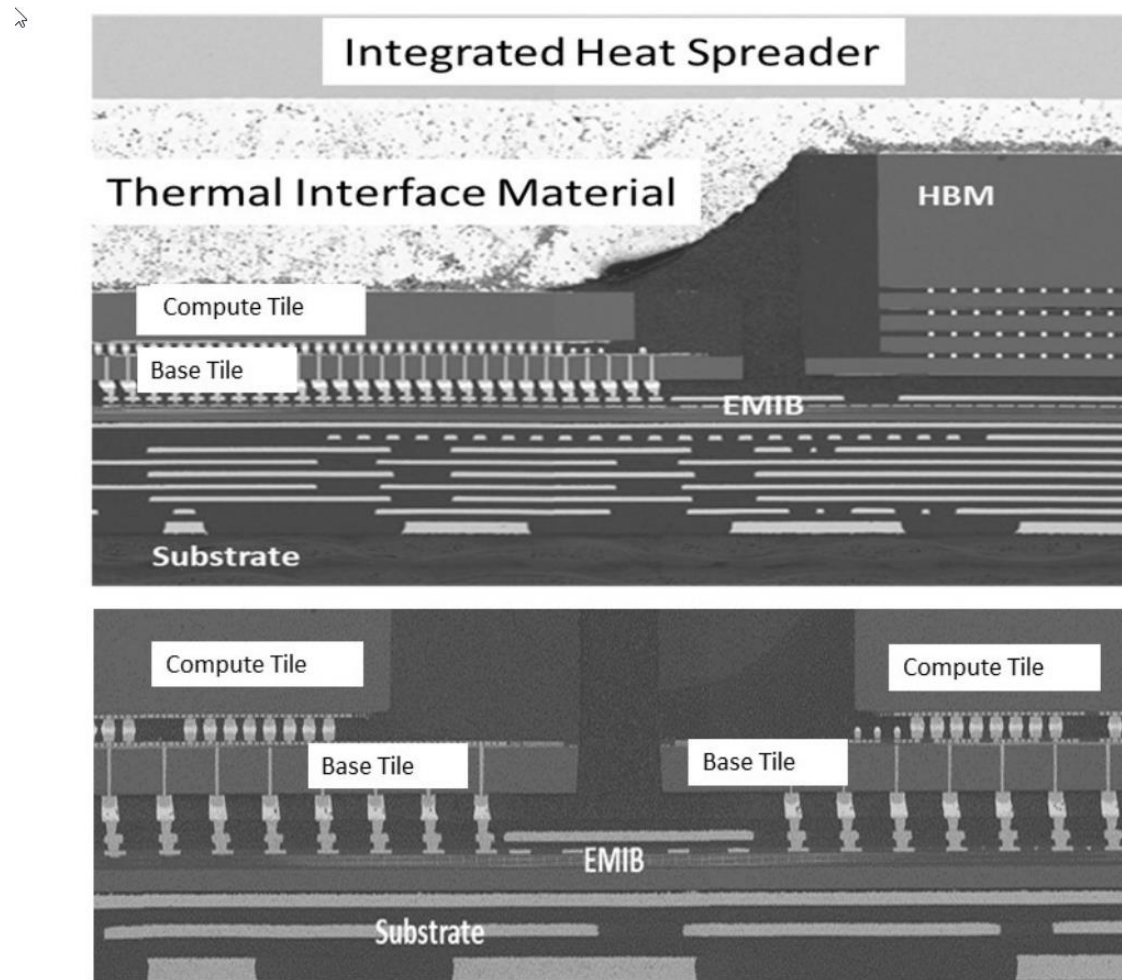
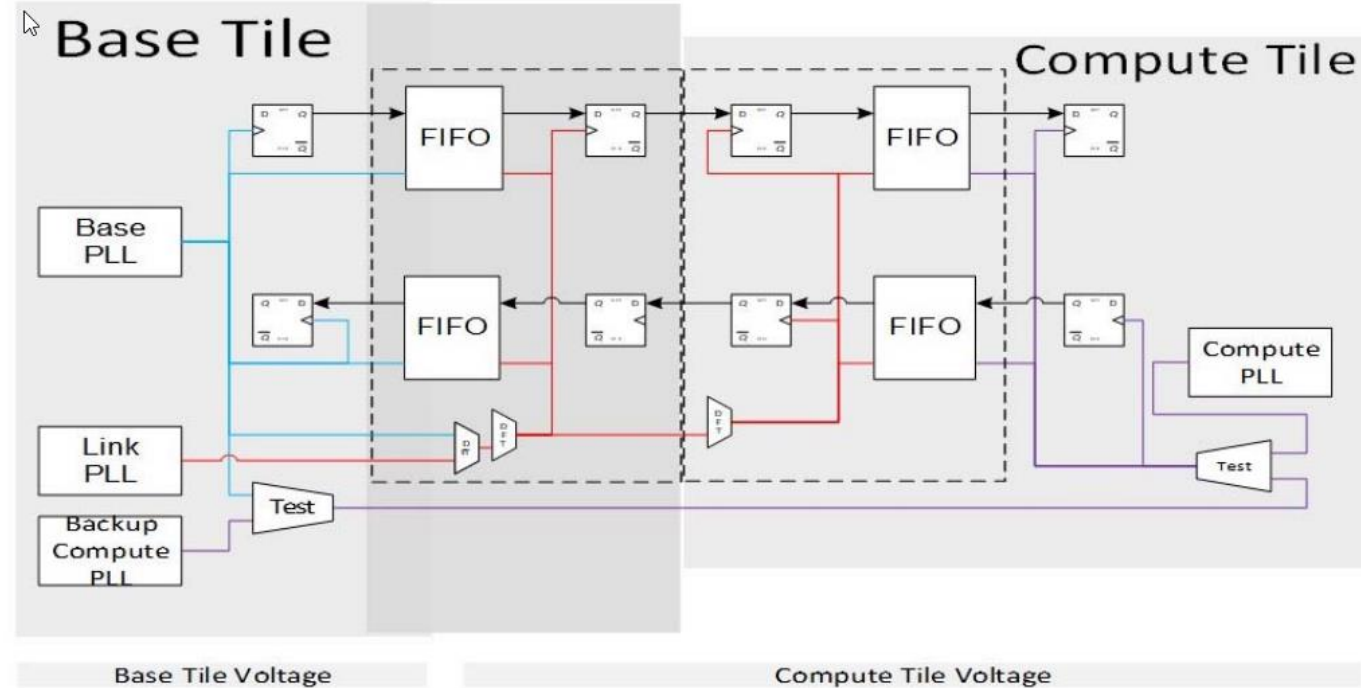


Figure 2.1.2: Process details for Foveros and EMIB.

Intel Ponte Vecchio Die-to-die IO



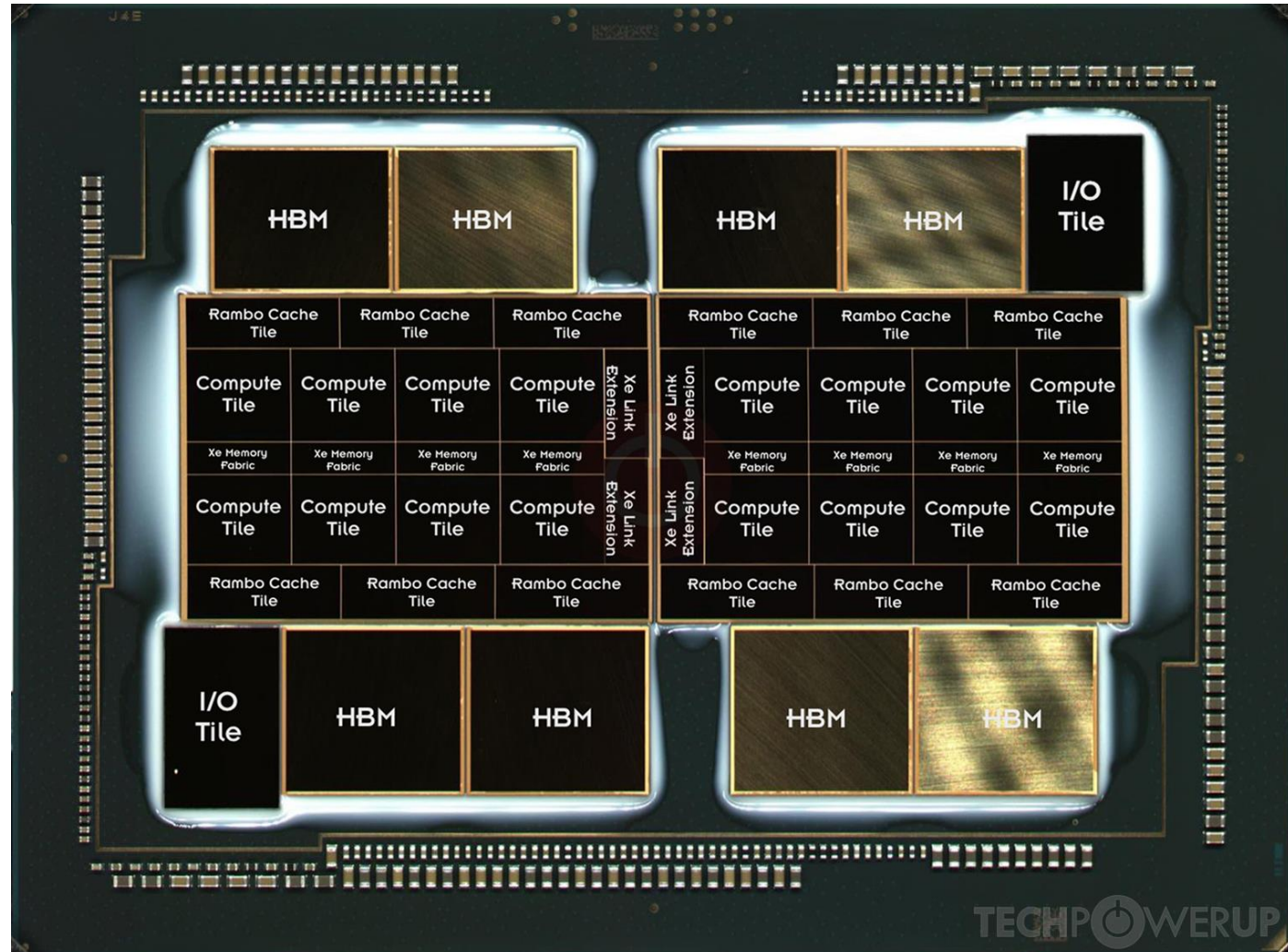
Project	Clock Area	Skew	Freq ratio @ 0.65V	Cdyn/mm ² ratio
Tiger lake Graphics (Intel 7)	1X	1X	1X	1X
PVC 2XCLK Graphics (Intel 7 Foveros)	16X	1X	2.6X	1.3X

Figure 2.1.4: PVC clocking, die-to-die IO and comparisons.

Intel GPU Max 1550 (2023)

Intel GPU Max 1550

- ❑ 600W TDP
- ❑ 128GB of HBM2e, 1024 bit interface
- ❑ Memory bandwidth: 3.27 TB/s



Conclusions

- ❑ Industry has shifted towards chiplets
- ❑ UCle creates a solid foundation for chiplet-based designs
- ❑ Future chips may be monolithic 3D ICs
- ❑ Chiplets solve several issues of current chips, beyond Moore's Law
 - ❑ Reticle limit
 - ❑ Better yields due to smaller silicon area of a single silicon die
 - ❑ Multi-vendor integration
- ❑ Universities need to introduce chiplet-based design in the curricula

Thank you!
Questions?