

# AMD DELIVERS BREAKTHROUGH MEMORY PERFORMANCE WITH DDR5 DRAM AND COMPUTE EXPRESS LINK<sup>™</sup> (CXL<sup>™</sup>) SUPPORT

AMD EPYC<sup>™</sup> 9004 SERIES PROCESSORS SUPPORT 12-CHANNEL DDR5 MEMORY AND THE CXL OPEN-STANDARD INTERFACE TO DELIVER HIGH PERFORMANCE AND LOW COSTS. Evolving performance, cost, and sustainability requirements are driving the need for new approaches to server memory. CIOs, IT managers, and data center architects face numerous challenges with server memory; these challenges vary by use case.

- Application requirements. Some applications like artificial intelligence (AI) can have an insatiable need for memory. Higher memory capacity and higher memory bandwidth are needed to support AI, database, and analytics workloads. These workloads can help organizations stave off competition and drive better business results.
- Scalability. As server core counts have increased, memory bandwidth has not kept pace. This has been due to memory technology and platform constraints. Figure 1 shows the ratio of CPU floating-point operations per second (FLOPS) versus memory bandwidth (measured in words/second) over several decades. The positive slope of 4.5x/decade demonstrates that the CPU performance is growing at a higher rate than memory bandwidth.<sup>1</sup> This imbalance can create a memory bandwidth bottleneck for applications that require significant data from memory.



Memory stranding. No two workloads need the exact same ratios of compute, storage, memory, and network capacity. For certain workloads, this can lead to a situation where memory might become "stranded" because there are not enough cores to pair with it. Because memory cost represents somewhere between 25 percent and 50 percent of server cost,<sup>2</sup> organizations can end up underutilizing a high-cost resource.

Architects seek denser memory at lower cost. For this, they are turning to persistent or storage-class memory. These solutions can often cost less than traditional local attached memory and may offer lower access latencies with higher throughput, compared to solid-state drives (SSDs). Persistent memory also enables DRAM-like access through system memory to large datasets in near real time.

Persistent memory also can save organizations time because it enables data retention during planned or unplanned restarts. That ability reduces the need for time-consuming data reloads, which can translate into cost savings from increased server utilization.

• Sustainability. Increasingly, organizations are looking to address sustainability goals by reducing their overall carbon footprints.

To address these memory challenges, AMD EPYC<sup>™</sup> 9004 Series processors support 12-channel DDR5 memory, a bandwidth and capacity upgrade from the 8-channel DDR4 memory supported in previous-generation AMD processors and x86 competitive processors. AMD EPYC 9004 Series processors with DDR5 memory deliver 2.25x the memory bandwidth compared to the previous generation, which helps close the memory bandwidth/core gap illustrated in Figure 1.<sup>3</sup> AMD is also pursuing a groundbreaking architectural shift by enabling breakthrough memory expansion with support for Compute Express Link<sup>™</sup> (CXL<sup>™</sup>), an open-standard interface, to enable:

2

- Flexible expansion of both memory bandwidth and capacity
- Innovative memory tiering to fill the current latency gap between direct-attached memory and CXL-attached memory
- Shared memory pools across multiple hosts

AMD is the current leader in CPU core count and holds <u>performance world records</u> across a broad swath of industry-standard benchmarks. With breakthrough memory expansion, 4th Generation AMD EPYC processors are positioned to become the ultimate choice to power data center servers.

This paper provides an overview of how AMD is expanding memory capacity in 4th Generation AMD EPYC processors by supporting both 12-channel DDR5 memory and the CXL interface standard. The paper then outlines how AMD is delivering breakthrough memory expansion and lighting up the ecosystem around it.

# TRADITIONAL MEMORY SYSTEM ARCHITECTURE

Traditional server architectures comprise dedicated processing, memory, and network systems. Memory connects directly to the CPU through DIMM slots. The number of memory channels, the media technology, and the memory speeds set the limit for memory bandwidth and capacity. For example, AMD EPYC 7003 Series processors offer 8 DDR4 memory channels with up to 3,200MHz DIMMs to support 204.8GB/s per-socket theoretical memory bandwidth.

In traditional memory system architectures, there are limits to expanding memory capacity. First, the number of memory channels might be limited by platform constraints. Second, the number of channels is hard-wired so that adding additional channels on-the-go is not feasible. Lastly, adding more capacity by using high-density DIMMs can be expensive from a cost and power perspective. This is where CXL memory expansion and memory disaggregation pave the way for architectural advancements. CXL enables flexibility to add channels "on the go" over the CXL interface and to add a wide variety of heterogeneous memory parts to fit performance or cost criteria.

### INTRODUCTION TO CXL

The CXL open standard is published by the CXL Consortium. The CXL Consortium describes its architecture as "cache coherent." Cache coherent means that changes in data values are propagated throughout a system in a timely fashion. Because CXL is based on resource sharing, the same data will always match even when stored in more than one location.

Because CXL enables the ability to coherently share memory resources between computing devices, it can dramatically increase performance. For example, a server with a CXL-connected accelerator would allow an accelerator to use the CPU's direct-attached memory for workloads that required greater memory capacity. Without CXL, the accelerator would need to access a storage device like an SSD or a hard-disk drive (HDD), which could take more time.

#### **12-CHANNEL DDR5 MEMORY**

AMD EPYC 9004 Series processors support 12-channel DDR5 memory, which offers 4,800 megatransfers per second (MT/s), up from 8-channel DDR4 memory supporting 3,200MT/s on current-generation processors. DDR5 memory is designed to maximize performance for AI, highperformance computing (HPC), and other data-intensive applications and workloads.

DDR5 memory offers key benefits for data center servers:

- Nearly doubles the bandwidth of DDR4 memory
- Powers memory-intensive workloads with 128GB and 256GB DIMMs available today, with support for 512GB DIMMs coming in 2023
- Offers architectural improvements and on-module power-management capabilities



CXL makes use makes use of the PCIe<sup>®</sup> 5.0 ecosystem by using its physical and electrical layers, but it adds an expanded protocol layer. A flexible processor port auto-negotiates between the standard PCIe protocols and the alternate CXL protocols. CXL enables high-bandwidth, low-latency connectivity between a host processor and devices including accelerators, memory buffers, and smart input/output (I/O) devices. Figure 2 shows two hosts sharing CXL-attached memory connected to the CPUs through PCIe 5.0.

The CXL Consortium has created an open ecosystem for data center high-speed enhancements like memory.<sup>4</sup> Its board membership list reads like a list of "who's who" of the semiconductor and data center industries. CXL Board members are from Alibaba Group, AMD, ARM, Cisco, Dell Technologies, Google, Huawei Group, Intel, Meta, Microsoft, NVIDIA, Rambus, and Samsung.

The CXL Consortium introduced the CXL 1.0/1.1 specification in 2019 to address the performance needs of cloud computing, AI, and analytics. Since then, the CXL 2.0 specification has been introduced to define memory pooling, link encryption, and switching. Recently, the CXL 3.0 specification was released to define fabrics and new device types.

CXL defines three classes of devices:

- Type 1: Processing devices without attached memory
- Type 2: Accelerators with attached memory
- Type 3: CXL-attached memory devices

AMD uniquely supports CXL 1.1+, which includes implementing some of the 2.0 specifications, with a focus on supporting type-3 devices. This gives organizations options to expand, tier, or pool memory.

# AMD ENABLES BREAKTHROUGH MEMORY EXPANSION AND A BROAD ECOSYSTEM

AMD enables breakthrough memory capabilities on AMD EPYC 9004 Series processors with both 12-channel DDR5 memory support and CXL interface support with 64 CXL capable lanes, in addition to up to x4 link bifurcation support. Benefits include system flexibility, memory expansion, memory tiering, enhanced security, support for persistent memory, and, finally, support for disaggregated computing. AMD works across a broad ecosystem including ISVs and IHVs like Astera Labs, Marvell, MemVerge, Micron, and others to demonstrate these benefits.

### SYSTEM FLEXIBILITY

Architects gain a previously unimaginable level of system-level flexibility with CXL. CXL interfaces are media-agnostic and offer a high degree of configurability to optimize system lane use. AMD EPYC 9004 Series processors interoperate with all types of memory-expansion controllers and module form factors across various memory generations, capacities, and bandwidths. Architects have flexibility at the link level for both form factor and width. With up to 64 CXL capable lanes and link bifurcation support (up to x4), AMD EPYC 9004 Series processors provide optimum lane usage for system configurations. Currently available engineering sample options include:

- x8 CEM add-in card (AIC)
- x8 E3.s module
- x16 CEM AIC prototype
- x16 CEM AIC prototype

x16 CEM AIC

The CXL interface supports a standard mode of 32GT/s in an x16 lane configuration. To allow for bifurcation, it also supports x8 and x4 lane configurations.<sup>5</sup>



Figure 2. CXL enables pools of shared memory

4

### MEMORY EXPANSION

Memory expansion is a key benefit of CXL. For example, the Micron<sup>®</sup> 2.0 memory-expansion module gives AMD EPYC 9004 Series processors access to up to 4TB of CXL-attached memory and more than 250GB/second memory bandwidth.<sup>6</sup>

CXL architecture is expected to bring considerable performance benefits. The latency for CXL-attached memory is targeted to be approximately one non-uniform memory access (NUMA) hop latency.

To demonstrate the relatively low latency of CXL-attached memory compared to the latency of CPU-attached DRAM, Astera Labs tested its Leo Smart Memory Controllers with AMD EPYC 9004 Series processors. The front-side bandwidth of the CXL x8 PCIe Gen5 port was matched with back-end DDR5 memory bandwidth. Astera Labs found that the memtier benchmark run on 64GB of CXL-attached memory was 96 percent as fast as the same benchmark run on 64GB of CPU-attached memory (see Figure 3).<sup>7</sup>



#### **MEMORY TIERING**

Another benefit of CXL is that it enables memory tiering. Tiering can result in more efficient and cost-effective use of memory. With memory tiering, the fastest, most expensive memory can be used to store hot data. Slower, less expensive memory can be used to store warm or cold data (see Figure 4). AMD EPYC 9004 Series processors can make use of new tiers of memory between DRAM and SSDs. Data tiering can also be used to create a new persistent memory tier.

With memory tiering, AMD EPYC 9004 Series processors support active page migration across local memory and CXL-attached memory. That is, hot pages are migrated into faster DRAM memory, and cold pages are retired to CXL-attached memory.

AMD EPYC 9004 Series processors can use MemVerge<sup>®</sup> Memory Machine<sup>™</sup> software and the AMD Instruction-Based Sampling (IBS) profiler to facilitate tiered memory. MemVerge Memory Machine software presents unified memory to the application. In the background, the software performs tiering services on behalf of the application. It promotes or demotes data from local memory to CXL-attached memory as needed.

#### DATA SECURITY FEATURES

Security features are foundational to any new memory architecture. As threats are getting more sophisticated, data in memory must be protected. AMD EPYC 9004 Series processors are designed with security capabilities in mind from the ground up. In addition, AMD Infinity Guard provides a unique and robust set of security features that help complement industry ecosystem partners at the software and system levels.<sup>8</sup>



Figure 4. Hot data can be stored in memory closer to the CPU

The AMD Infinity Guard protocol secure feature set—with AMD Secure Encrypted Virtualization (SEV) and AMD Secure Nested Paging (SNP) enables the same security capabilities over CXL-attached memory as with local DDR5 memory. This helps ensure that applications have access to proven, production-grade security features and helps enable confidential computing. AMD Infinity Guard works with existing software stacks seamlessly. It is independent of the CXL device implementation, and it works uniformly across CXL device vendors.

#### LOW TCO

Reducing the total cost of ownership (TCO) of memory is a priority for data center managers. Microsoft reports that the cost of DRAM currently accounts for more than 50 percent of the cost of building a server for Microsoft Azure<sup>®</sup>.<sup>9</sup> CXL has the potential to address this issue by:

- Flexibly allowing the use of alternative, lower-cost memory
- Enabling DIMM reuse, which can help lower costs and meet sustainability goals; as systems get upgraded to DDR5, organizations may opt to reuse the DDR4 memory being replaced (DDR4 DIMMs can be reused as CXL-attached memory in the new system)
- Allowing architects to use heterogeneous components to optimize memory for a given workload to help avoid overprovisioning costs

#### PERSISTENT MEMORY SUPPORT

When more persistent memory is needed for a given workload, it can be made available through the CXL interface. Persistent memory retains data once system power is off, which can help save time by eliminating the need for data reloads. This feature helps decrease server downtime and increase efficiency.

Persistent memory also opens opportunities for lowering memory cost. Persistent memory or storage-class memory might be available at a lower cost with "good enough" performance.

#### **DISAGGREGATED COMPUTING**

This concept of resource disaggregation is not new. In 2009, Meta (Facebook) was growing exponentially, and the company realized that it had to rethink its infrastructure to accommodate the huge influx of new people and data, and also control costs and energy consumption. They initiated a project to design the world's most energy efficient data center, one that could handle unprecedented scale at the lowest possible cost. A small team of engineers spent the next two years designing and building one from the ground up: software, servers, racks, power supplies, and cooling. The result now stands in Prineville, Oregon.

It was 38% more energy efficient to build and 24% less expensive to run than the company's previous facilities—and has led to even greater innovation. The Open Compute Project Foundation (OCP) was then formed in 2011 with a mission to apply the benefits of open source and open collaboration to hardware and rapidly increase the pace of innovation in, near and around the data center, and beyond.<sup>10</sup>

The OCP disaggregated components from the traditional server shelf into pools of resources that could be shared across a rack. Pooled storage and pooled memory could be used and flexibly allocated among the server shelves as workloads dictated. By not burdening each server shelf with the maximum amount of memory and storage that might be needed, and instead providing flexible allocation, the overall memory and storage requirements across the rack could be reduced, saving costs. Using pooled resources helps to avoid memory stranding, defined as the scenario when all of the cores in a server have memory allocations and remaining memory capacity can't be used because there are no cores remaining. Additionally, shared power and cooling could be used to improve efficiency and reduce costs.

Now, nearly 10 years later, the computer industry has converged on a similar vision of disaggregation with the CXL open standard published by the CXL Consortium. CXL enables a disaggregated architecture. Pools of shared heterogeneous resources can be efficiently composed to support the needs of specific workloads. This is accomplished through dynamic resource allocation. For example, if memory is not being used, it can be dynamically reassigned to a different host.

AMD EPYC 9004 Series processors use CXL to pool memory to better support workload needs. For example, AMD EPYC 9004 Series processors can work with Marvell<sup>®</sup> memory pooling CXL solutions to connect multiple hosts with a shared memory-pool.<sup>11</sup> Pooling enables dynamic allocation of memory capacity and bandwidth while helping to eliminate stranded memory. When memory capacity is used more efficiently, it can effectively lower TCO.

Cloud service providers (CSPs) have experimented with CXL-enabled memory pooling with positive results. For example, Microsoft found that by pooling memory with CXL across 16 sockets and 32 sockets within a cluster, it could reduce memory demand by 10 percent.<sup>9</sup> This could translate to cutting the cost of its servers by 4 to 5 percent, and it could result in the savings of hundreds of millions of dollars per year.<sup>9</sup>

# AMD LEADERSHIP INNOVATION CONTINUES

Faster AI, database, and analytics workloads are required to drive better business results. To succeed, organizations need their data center servers to deliver high performance while processing large quantities of data. But server memory constraints can stand in the way.

AMD EPYC 9004 Series processors address these memory challenges by supporting both 12-channel DDR5 memory and the CXL open-standard interface to deliver breakthrough memory capabilities. Additionally, AMD is enabling a broad production-ready CXL ecosystem to support adoption of the CXL-enabled AMD EPYC 9004 Series processors.

The AMD EPYC 9004 Series processors help offer customers the benefits of a disaggregated computing architecture: increased flexibility, the ability to expand, tier, and pool memory, support for sustainability goals, and low TCO–all with an overlay of data security capabilities.

LEARN MORE AT www.amd.com/en/processors/epyc-9004-series.

7

- <sup>1</sup> John D. McCalpin, Ph.D. "SC16 Invited Talk: Dr. John D. McCalpin Presents 'Memory Bandwidth and System Balance in HPC Systems." The University of Texas at Austin Texas Advanced
- Computing Center (TACC). October 2016. http://sc16.supercomputing.org/2016/10/07/sc16-invited-talk-spotlight-dr-john-d-mccalpin-presents-memory-bandwidth-system-balance-hpc-systems/index.html.
- <sup>2</sup> The Next Platform. "The Expanding CXL Memory Hierarchy is Inevitable and Good Enough." August 2022.
- www.nextplatform.com/2022/08/22/the-expanding-cxl-memory-hierarchy-is-inevitable-and-good-enough/. <sup>3</sup> AMD claim, EPYC-040: AMD EPYC 9004 CPUs support 12 channels of up to 4800 MHz DDR5 memory, which is 460.8 GB/s of maximum memory throughput per socket. Prior
- generation of AMD EPYC CPUs (7003 series) have a maximum 204.8 GB/s. AMD EPYC 9004 CPUs have 2.25x the memory throughput per CPU. 460.8 + 204.8 = 2.3x (2.25x) the max memory throughput.
- <sup>4</sup> For more information on Compute Express Link (CXL) and the CXL Consortium, visit <u>www.computeexpresslink.org</u>.
- <sup>5</sup> Gary Ruggles. "CXL 2.0 and 3.0 for Storage and Memory Applications." Synopsys. 2022.
- www.synopsys.com/designware-ip/technical-bulletin/cxl2-3-storage-memory-applications.html.
- <sup>6</sup> Micron measured results as of 9/29/2022. Not independently verified by AMD.
- <sup>7</sup> Astera Labs measured results as of 9/29/2022. Test configuration compares the same capacity on CXL-attached memory versus CPU-attached memory. Testing not independently verified by AMD.
- <sup>8</sup> AMD claim, GD-183: AMD Infinity Guard features vary by AMD EPYC<sup>\*\*</sup> processor generations. Infinity Guard security features must be enabled by server OEMs and/or cloud service providers to operate. Check with your OEM or provider to confirm support of these features. Learn more about Infinity Guard at www.amd.com/en/technologies/infinity-guard.
- <sup>9</sup> The Next Platform. "Microsoft Azure Blazes the Disaggregated Memory Trail with Znuma." July 2022.
- www.nextplatform.com/2022/07/11/microsoft-azure-blazes-the-disaggregated-memory-trail-with-znuma/.
- <sup>10</sup> Open Compute Project (OCP). About page. Accessed December 2022. <u>www.opencompute.org/about</u>.
- <sup>10</sup> Marvell, "Marvell Announces Innovative CXL Development Platform for Multi-Host Memory Pooling," November 2022. www.marvell.com/company/newsroom/marvell-announces-innovative-cxl-development-platform-for-multi-host-memory-pooling.html.

©2023 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo, AMD EPYC<sup>™</sup>, and combinations thereof are trademarks of Advanced Micro Devices, Inc. in the United States and/or other jurisdictions. Other product names are for informational purposes only and may be trademarks of their respective companies. Compute Express Link and CXL are trademarks of the CXL Consortium in the US or other countries.