

IT Systems Engineering | Universität Potsdam



Memory Disaggregation with Coherent Interconnect Technologies

> <u>Lukas Wenzel</u>, Felix Eberhardt, Andreas Grapentin, Prof. Dr. Andreas Polze <firstname>.<lastname>@hpi.de

Operating Systems and Middleware Group, Hasso Plattner Institute

06.04.2023



Disaggregation

The separation of an aggregate body into its component parts.

from the GNU version of the Collaborative International Dictionary of English

Instead of separate machines with a **fixed**, **tightly coupled** resource configuration

→ different resource pools with **dynamic**, **loosely coupled** associations



Memory Disaggregation Benefits

Operational and provisioning flexibility:

- Cost efficiency
- Energy efficiency (?)





Memory Disaggregation Related Concepts

Non-Uniform Memory Access (NUMA) System

Nodes share a global physical address space, where local and remote memory resources are accessible.

Disaggregated Memory System

Machines have private physical address spaces,

but may map remote memory resources (often exclusively, sometimes shared).

Distributed / Shared-Nothing System

Machines have private physical address spaces, and can only access local memory resources.

Disaggregated

Wenzel (OSM HPI)

Memory







Landscape Early 1990s



Early proposals appeared in the wake of Multicomputer system developments

- Remote memory replaces local disks as swap space through custom page fault handlers or block device drivers
- = Movement of complete pages, often over LAN
- D. <u>Comer</u> and J Griffioen "A New Design for Distributed Systems: The Remote Memory Model." USENIX Summer, 1990.



Disaggregated Memory Wenzel (OSM HPI)

 E. W. <u>Felten</u> and J. Zahorjan "Issues in the Implementation of a Remote Paging System." Technical Report TR 91-03-09, Computer Science Department, University of Washington, 1991.

5

Figure 2: Client and Memory Server Interaction.

Node 0

Client

Nodes

Memory

Remote Physical Memory (local to servers)

Remote

Disks

Node i

Memory

Server

Page Request

Node n

Memory

Server

Get

Client

Page

6

Early proposals appeared in the wake of Multicomputer system developments

- Remote memory replaces local disks as **swap space** through custom page fault handlers or block = device drivers
- Movement of complete pages, often over LAN \equiv
- L. <u>Iftode</u>, et al. "Memory Servers for Multicomputers." In Digest of Papers. Compcon Spring (pp. 538-547). IEEE, 1993.

From Computer History Museum

https://www.computerhistory.org/ revolution/supercomputers/10/74/288



Disaggregated Memory Wenzel (OSM HPI)

/////// Virtual Memory Memory Page Service Server Software Local Physical





Landscape Early 1990s

Landscape Early 1990s



Early proposals appeared in the wake of Multicomputer system developments

- Remote memory replaces local disks as swap space through custom page fault handlers or block device drivers
- = Movement of complete pages, often over LAN
- M. J. <u>Feeley</u>, et al. "Implementing Global Memory Management in a Workstation Cluster." In Proceedings of the 15th ACM symposium on operating systems principles (pp. 201-212), 1995.



Figure 1: Global replacement with hit in the global cache.



Figure 2: Global replacement showing miss in the global cache. The faulted page is read from disk, and the oldest page in the network is either discarded (if clean) or written back to disk. **Disaggregated Memory** Wenzel (OSM HPI)

7

Landscape Late 2000s to early 2010s



Advances in interconnect technology brought new approaches

- = **Backplane** connections via PCIe or coherent socket-to-socket interfaces
- = First steps in **fine grained** remote memory **access**, i.e., cache-line granularity
- K. <u>Lim</u>, et al. "Disaggregated Memory for Expansion and Sharing in Blade Servers." ACM SIGARCH computer architecture news, 37(3), 267-278, 2009.



(a) Memory blade design

Figure 2: Design of the memory blade. (a) The memory blade connects to the compute blades via the enclosure backplane.



(a) Compute blade

Figure 3: Page-swapping remote memory system design. (a) No changes are required to compute servers and networking on existing blade designs. Our solution adds minor modules (shaded block) to the virtualization layer.



This design assumes minor coherence hardware support in every compute blade.

Disaggregated Memory Wenzel (OSM HPI)

8

Landscape Late 2000s to early 2010s



Advances in interconnect technology brought new approaches

- = **Backplane** connections via PCIe or coherent socket-to-socket interfaces
- = First steps in **fine grained** remote memory **access**, i.e., cache-line granularity
- H. <u>Montaner</u>, et al. "A Practical Way to Extend Shared Memory Support Beyond a Motherboard at Low Cost." In Proceedings of the 19th ACM International Symposium on High Performance Distributed Computing (pp. 155-166), 2010.





attached to the motherboard.

Landscape Late 2000s to early 2010s



Advances in interconnect technology brought new approaches

- = **Backplane** connections via PCIe or coherent socket-to-socket interfaces
- = First steps in **fine grained** remote memory **access**, i.e., cache-line granularity
- R. <u>Hou</u>, et al. "Cost Effective Data Center Servers." In 19th International Symposium on High Performance Computer Architecture (HPCA) (pp. 179-187), IEEE, 2013.







Remote Direct Memory Access (RDMA)

- = Protocol and NIC technology allowing **network memory access without processor activity**
- Originally layered on top of InfiniBand, later on Ethernet (RoCE) and even SCTP or TCP transports (iWARP)
- NIC implements RDMA protocol without processor support
- NIC has autonomous DMA access to the memory subsystem
 - > Operating System must allocate buffers and setup mappings and access permissions first





Remote Direct Memory Access (RDMA)

= Protocol and NIC technology allowing **network memory access without processor activity**

Operations (Verbs):

- Message Semantics, two-sided (Send and Receive)
 - Receiver prepares receive buffer and requests Receive
 - Sender fills send buffer and requests Send
 - Sender NIC reads send buffer, transmits to Receiver NIC, which writes receive buffer





Remote Direct Memory Access (RDMA)

= Protocol and NIC technology allowing **network memory access without processor activity**

Operations (Verbs):

- Memory Semantics, one-sided (Read or Write)
 - Reader/Writer sends RDMA request to remote NIC
 - Remote NIC performs memory access on Readers/Writers behalf





RDMA sparked wide range of disaggregated memory proposals

- = **Serving** memory is handled by RDMA NIC with **low overhead**
- = **Using** remote memory requires **explicit initiation** of RDMA transactions
- Variety of approaches place RDMA initiative on different components
 - Libraries
 - Runtime environment
 - Operating system
 - Hypervisor
- Comprehensive survey for further details and references:

M. <u>Ewais</u> and P. Chow "**Disaggregated Memory in the Datacenter: A Survey.**" *IEEE Access* 11: 20688-20712, **2023**.



Classic Coherent Interconnects (like HyperTransport, QPI, UPI, ...) designed for multi-socket machines

- Limited reach beyond board or backplane
- Limited interoperability between vendors and device types

In 2016, several consortia formed to standardize coherent interconnects

- **Gen-Z**: switchable, envisioned flexibility to span multiple racks
- CCIX: point-to-point, announced but not validated support in AMD Epyc 7002 (x86) and Huawei TaiShan 200 (ARM)
- OpenCAPI: point-to-point, evolved from IBM proprietary CAPI interface, available on POWER8 and POWER9

In 2019, **CXL** consortium was formed, previous efforts started to converge (GenZ transferred in 2021, OpenCAPI in 2022)

Landscape CXL

Compute Express Link

- Layered on top of **PCIe**
- Processor (host) implements different **subprotocols** with separate semantics:
 - CXL.io: discovery, register access, interrupts (~regular PCIe)
 - **CXL.cache**: coherent host memory access (device → host, full or io coherency)
 - **CXL.mem**: access to device memory resources (host → device)



Figure 2-1. CXL Device Types

From: Compute Express Link Specification, Revision 3.0, Version 1.0 © 2019-2022 COMPUTE EXPRESS LINK CONSORTIUM, INC. ALL RIGHTS RESERVED.



HPI Hasso Plattner Institut

Coherent interconnects enable new design space for disaggregated memory clients

- = Application **load-store accesses** can be **forwarded transparently** to remote memory server
- K. <u>Katrinis</u>, et al. "Rack-Scale Disaggregated Cloud Data Centers: The dReDBox Project Vision." Design, Automation & Test in Europe Conference & Exhibition (DATE), IEEE, 2016.



 I. <u>Calciu</u>, et al. "Rethinking Software Runtimes for Disaggregated Memory." Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, 2021.

> Simulated System





Figure 5: The Kona Architecture: An FPGA connected to a CPU through a coherent interconnect. Both the CPU and the FPGA have DRAM attached (CMem and FMem, respectively). The FPGA exposes fake physical memory to the CPU (VFMem), backed by remote memory.

Disaggregated Memory

Wenzel (OSM HPI)

17



Coherent interconnects enable new design space for disaggregated memory clients

- = Application **load-store accesses** can be **forwarded transparently** to remote memory server
- V. R. <u>Kommareddy</u>, et al. "DeACT: Architecture-Aware Virtual Memory Support for Fabric Attached Memory Systems." IEEE International Symposium on High-Performance Computer Architecture (HPCA), 2021.





Coherent interconnects enable new design space for disaggregated memory clients

- = Application load-store accesses can be forwarded transparently to remote memory server
- C. <u>Pinto</u>, et al. "ThymesisFlow: A Software-Defined, HW/SW co-Designed Interconnect Stack for Rack-Scale Memory Disaggregation." 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO), 2020.





"RTT latency of this prototype is roughly **950ns**, which includes four crossings of the FPGA stack and six serDES crossings."

Disaggregated Memory

Wenzel (OSM HPI)



Coherent interconnects enable new design space for disaggregated memory clients

- Application **load-store accesses** can be **forwarded transparently** to remote memory server =
- D. <u>Gouk</u>, et al. "Direct Access, High-Performance Memory Disaggregation with DirectCXL." USENIX ATC, 2022.



Disaggregated Memory Wenzel (OSM HPI)

Local

RDMA

DirectC



Coherent interconnects enable new design space for disaggregated memory clients

- = Application **load-store accesses** can be **forwarded transparently** to remote memory server
- C. <u>Wang</u>, et al. "CXL over Ethernet: A Novel FPGA-based Memory Disaggregation Design in Data Centers." arXiv preprint arXiv:2302.08055, 2023.







Coherent interconnects enable new design space for disaggregated memory clients

- = Application **load-store accesses** can be **forwarded transparently** to remote memory server
- C. <u>Wang</u>, et al. "CXL over Ethernet: A Novel FPGA-based Memory Disaggregation Design in Data Centers." arXiv preprint arXiv:2302.08055, 2023.





Coherent interconnects enable new design space for disaggregated memory clients

- = Application load-store accesses can be forwarded transparently to remote memory server
- Y. <u>Sun</u>, et al. "Demystifying CXL Memory with Genuine CXL-Ready Systems and Devices." arXiv preprint arXiv:2303.15375, 2023.



Figure 2: Access latency. Average latency for single AVX512 load (ld), store and write back (st+wb), non-temporal store (nt-st), and sequential pointer chasing in 1GB space(ptr-chase); Sequeutial pointer chasing latency vs. working set size (right). Prefetching at all levels are disabled in both cases





Coherent interconnects enable new design space for disaggregated memory clients

- = Application load-store accesses can be forwarded transparently to remote memory server
- Y. <u>Sun</u>, et al. "Demystifying CXL Memory with Genuine CXL-Ready Systems and Devices." arXiv preprint arXiv:2303.15375, 2023.



Figure 3: Sequential access bandwidth. An experiment showing the maximum possible bandwidth on local-socket DDR5 with 8-channels (a), and CXL memory (b), and remote-socket DDR5 with 1 channel (c). The grey dash line in (b) shows the theoretical max speed of DDR4-2666MT/s



Your comments or experiences?

ThymesisFlow @ HPI Setup

- Hardware
 - □ 2x IBM **IC922** or IBM AC922 or Inspur FP5290G2
 - 2x AlphaData 9V3 FPGAs
 - □ 2x OpenCAPI cables (SlimSAS) 25 Gb/s x8
 - □ 2x 100 Gb/s network cabling







ThymesisFlow @ HPI OpenCAPI 3.0 Protocol

- OpenCAPI 3.0 allows **transparent** remote memory access:
- Accelerator Cards support two orthogonal modes
 - C1 Accelerator acts as another Core and accesses host memory using *virtual addresses* (like CXL.memory)
 - M1 Accelerator acts as another Memory controller and handles memory accesses by the host using *physical addresses* (like CXL.cache)





ThymesisFlow @ HPI Architecture

- Accelerator and Host both post and complete commands:
 - C1-role initiates read or write commands to be handled by host system
 - ThymesisFlow forwards commands from M1 to C1 card via Aurora point-to-point link
 - M1-role handles read or write commands from host system
- Repository: <u>https://github.com/OpenCAPI/ThymesisFlow</u>
- Deep Dive (C. Pinto): <u>https://www.youtube.com/watch?v=4N5G_Am6U-M</u>





ThymesisFlow @ HPI Architecture



- **Memory Server**: C1 mode accelerator uses *virtual address* space of attached process
 - thymesisf-agent application allocates memory to serve and attaches to accelerator
- **Memory Client**: M1 mode accelerator appears at a fixed position in *physical address* space
 - Raw access via mmap()-able memory device at /dev/mishmem-s<n>
 - Registered as a CPU-less NUMA node → application can allocate remote memory through libnuma



Static Offsets

ThymesisFlow @ HPI Demo



andreas.grapentin@ic922-03 ~ \$ lsmod head -n2 Module Size Used by mishmem_s1 5853 0 andreas.grapentin@ic922-03 ~ \$ ls -l /dev/mishmem-s1 crw-rw-rw-1 root root 507, 0 Apr 3 18:55 /dev/mishmem-s1 andreas.grapentin@ic922-03 ~ \$ file /dev/mishmem-s1 /dev/mishmem-s1: character special (507/0) andreas.grapentin@ic922-03 ~ \$ numactl -H available: 4 nodes (0,8,48,64) node 0 cpus: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 node 0 size: 261683 MB	<pre>andreas.grapentin@ic922-04 ~ \$ ls -l /dev/mishmem-s1 -rw-r-rw- 1 root root 4294967296 Apr 3 19:09 /dev/mishmem-s1 andreas.grapentin@ic922-04 ~ \$ /opt/thymesisflow/bin/thymesisf-cli \ > attach-memory \ >afu IBM,RMEM \ >cid 1 \ >cid 1 \ >size 4294967296 \ >port 2 2023-04-03 19:10:25 [INF0] Attaching memory - cid: 1 afu: IBM,RMEM:2 - size: 4294967296 2023-04-03 19:10:26 [DEBUG] circuitid: 1 - type: mresponse_attach - afu IBM,RMEM - size 4294967296 - ea(hex): 100000000 - ea(int):4294967296 - status: 100</pre>
node 0 free: 253409 MB	2023-04-03 19:10:26 [INFO] Successfully allocated memory - EA: 0x100000
node 8 cpus: 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67	000
08 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91	andreas.grapentin@ic922-04 ~ \$
92 93 94 93 node 8 size: 260725 MP	
node 8 free: 252108 MB	
node 48 cpus:	
node 48 size: 0 MB	
node 48 free: 0 MB	
node 64 cpus:	
node 64 size: 0 MB	
node 64 free: 0 MB	
node distances:	
node 0 8 48 64	
0: 10 40 80 80	
8: 40 10 80 80	
48: 80 80 10 80	
64: 80 80 80 10	
<pre>andreas.grapentin@ic922-03 ~ \$ /opt/thymesisflow/bin/thymesisf-cli \</pre>	
> attach-compute \	
>atu IBM,RMEM \	
>C10 1 \	
>5120 429490/290 \	
2023-04-03 19:10:41 [INEO] attaching compute - size: 4204067206 - using	
effective address 4294967296	
Π	
[0] 1:ssh*	"andi@arch-t370:~" 21:10 03-Apr-23

- [STREAM] J. D. <u>McCalpin</u>. **"STREAM: Sustainable Memory Bandwidth in High Performance Computers."** Technical Report, University of Virginia, Charlottesville, **1991**.
- [Imbench] L. <u>McVoy</u>, and C. Staelin. "Imbench: Portable Tools for Performance Analysis." USENIX annual technical conference, **1996**.

Disaggregated Memory Wenzel (OSM HPI)

30



Interested to experiment with or collaborate on extensions to this ThymesisFlow setup?

Message us!

{lukas.wenzel,felix.eberhardt,andreas.grapentin}@hpi.de