

Towards a Fully Disaggregated and Programmable Data Center

Yizhou Shan*

Huawei Cloud
shanyizhou@huawei.com

Zhiyuan Guo

University of California, San Diego
z9guo@eng.ucsd.edu

Will Lin

University of California, San Diego
w5lin@ucsd.edu

Yiying Zhang

University of California, San Diego
yiying@ucsd.edu

Abstract

The data center deployment status quo is challenged by the pressure from the applications it hosts and the hardware it runs on. On the one hand, the applications are becoming more fine-grained, driven by the recent trend of serverless and microservice. On the other hand, the hardware is increasingly faster, heterogeneous, and specialized due to recent network improvements and hardware resource disaggregation trends. Data center designers are at a crossroads facing several challenges when these two trends meet.

In this paper, we envision a futuristic data center consisting of disaggregated devices connected by a highly flexible data center network. We present guidelines and initial solutions for data center designers to navigate design trade-offs. We use a three-layer approach. At the top layer, we embrace two coexisting abstractions and propose a disaggregation-native design methodology. At the bottom layer, we describe the hardware and key features required to build disaggregated devices and the networking infrastructure. We propose an MLIR-based intermediate representation layer between these two layers to decouple the hardware from the software. It transforms software into small executable units to tame heterogeneity and aid runtime migration. We hope our proposal can pave the way for future disaggregated data center deployment.

CCS Concepts

• **Computer systems organization** → **Cloud computing**.

*Work done while at UCSD.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

APSys '22, August 23–24, 2022, Virtual Event, Singapore

© 2022 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-9441-3/22/08.

<https://doi.org/10.1145/3546591.3547527>

Keywords

Resource Disaggregation, Serverless, MLIR, Networking

ACM Reference Format:

Yizhou Shan, Will Lin, Zhiyuan Guo, and Yiying Zhang. 2022. Towards a Fully Disaggregated and Programmable Data Center. In *13th ACM SIGOPS Asia-Pacific Workshop on Systems (APSys '22)*, August 23–24, 2022, Virtual Event, Singapore. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3546591.3547527>

1 Introduction

For decades, data centers are organized as racks of general-purpose servers connected with a static topology of fixed-logic network switches, hosting workloads from various industries and domains. However, this status quo is challenged by the pressure from the applications it hosts and the hardware infrastructure it runs on.

On the one hand, data center applications are becoming more fine-grained, driven by the recent trend of serverless and microservice. With serverless, users can focus on their business logic, avoid deployment chores, and enjoy the pay-as-you-go model. As some pointed out, serverless computing is expected to become the default computing paradigm of the cloud era [54]. Hence, we will see more deployed and mature applications or system software transitioning into smaller, DAG-based serverless counterparts.

Besides, the data center hardware infrastructure is increasingly faster, programmable, and specialized. Several factors contribute to this trend. First of all, cloud vendors are customizing their hardware to avoid energy and computing inefficiencies as a result of the ending of Moore's Law and Dennard scaling [4, 5, 19, 32, 51]. Second, the data center network (DCN) is not only faster with single-digit microsecond latency and more than 100 Gbps bandwidth [6, 22], but also more flexible with programmable control [16, 18] and data planes [9, 16, 18, 19]. The above two trends *together* enabled hardware resource disaggregation, an architecture that breaks servers into network-attached hardware resource pools [27, 40, 55], each of which has disaggregated devices that can be built, scaled, and managed independently. We

make a key observation that most disaggregated devices will have some computation power as a result of co-packaging with SmartNICs-alike components [5, 19, 38, 45, 70, 78]. Consequently, disaggregated devices are programmable and can realize novel paradigms such as near-data processing or advanced memory/storage primitives [3, 10, 52, 58, 60].

All these emerging trends are exciting. However, the data center designers face several key challenges when they try to run emerging systems on new hardware. First, software and hardware are not evolving at the same pace, with hardware having a much longer lead time. This mismatch results in lockstep development and delays time to market [19, 43]. Second, there is a semantic gap between the software and the hardware. Unlocking the full potential of new hardware, especially heterogeneous devices such as FPGA or RMT [9], requires a non-trivial software redesign, *e.g.*, co-designing hypervisors with accelerators [2, 5, 19, 38]. Dynamic offloading frameworks can help but fall short in a distributed setting or dealing with multiple heterogeneous resources [42, 49, 50]. In all, it is not easy for data center designers to exploit the emerging hardware.

The third challenge is that it is unclear how to best build disaggregated devices, knowing what features are required, what are nice to have, and how to best use them. Though disaggregation first appeared more than a decade ago [40], most works build or emulate disaggregated hardware using regular servers [8, 33, 55, 61, 79]. This approach has sub-optimal performance due to software overheads. Little has been explored on building standalone disaggregated devices. Moreover, prior works usually explore a few design choices rather than exploit the optimal one [8, 20, 58, 79]. For instance, some delegate the memory address translation service onto disaggregated memory devices [55, 78]. However, this service could also run on a disaggregated computing device or a programmable switch depending on traffic patterns [36]. This begs the question of whether we have exploited all the trade-offs among programmability, manageability, performance, and cost.

The last challenge concerns DCN, as it is the cornerstone to realizing disaggregation. The disaggregation paradigm works because of the fast network [20, 55]. Furthermore, even though programmable switches and SmartNICs are gradually deployed in production [39], they often lack flexible programming models [41] and commercial multi-tenancy or hot-plug support [17, 67, 71, 72]. Moreover, the long-standing static network topology is hitting its limitations [7, 68, 69]. However, it is unclear how to tailor the DCN for disaggregation in terms of ensuring low latency, enabling flexible programmability and reconfigurability.

To help data center designers tackle these challenges, we propose guidelines and initial solutions to navigate the design trade-offs. The foundation of our vision is a fully disaggregated and programmable data center infrastructure consisting

of disaggregated devices connected by a highly flexible DCN. We adopt a three-layer approach as shown in Figure 1.

The top layer is the abstraction layer through which the users could interact with the underlying hardware. We embrace two coexisting abstractions and propose a design methodology to improve both. To maintain backward compatibility, we expect a traditional virtual machine (VM) abstraction to host legacy workloads, similar to the vNode concept [55]. To best utilize disaggregated hardware, we expect a distributed application execution abstraction expressed as DAGs, similar to [46, 62]. With these two abstractions, a futuristic disaggregated data center could host both legacy and emerging workloads. Nonetheless, this paper focuses on optimizing the second abstraction because the data center workloads are becoming more serverless (*i.e.*, expressed as DAGs), and prior work has shown that this abstraction reduces data movement and achieves better performance [62].

Similar to the cloud native idea, we propose a disaggregation native design methodology, which has a set of non-intrusive but informative ways for developers to consciously make their programs in a way that could benefit from the underlying heterogeneous and disaggregated hardware. We believe that it is best for developers to directly control the placement of their data/code, failure handling policies, and network topologies by explicitly annotating their program.

The abstractions and methodology mitigate the tension, but they cannot solve the fundamental semantic gap between the software and the hardware. Hence, we propose to use an Intermediate Representation (IR) layer. The IR layer decouples the software from the hardware, enabling each to evolve independently, and the IR can bridge the discrepancy. For instance, we can add new heterogeneous hardware without changing the software built on top of our abstractions. Its indirection layer enables software to run on heterogeneous hardware and facilitates runtime migration without any changes. We also plan to use the IR layer to transform high-level disaggregation-native annotations into low-level hardware primitives to reduce developer hazards. Together with the top layer, they address the first two challenges.

In particular, the IR consists of disaggregatable execution units, or *codelets*. We propose to use MLIR [35] to decompose a program into multiple smaller *codelets* and a companion DAG dictating the execution order of the codelets. The codelet is a powerful entity designed to exploit the hardware infrastructure. A codelet encapsulates everything that is needed to run a part of software on multiple devices, *e.g.*, it has data, IR instructions, network topology configurations, and policies such as security requirements. We would compile the codelet for distinct hardware targets using backends such as LLVM or CIRCT [1]. Furthermore, the codelet is designed as the smallest unit of migration among devices.

For the last layer, we envision a fully disaggregated and programmable hardware infrastructure, which has disaggregated devices organized as resource pools connected by a highly flexible DCN. We believe that SmartNICs (or DPUs) are key to building disaggregated devices. In fact, we think every disaggregated device should have a co-packaged SmartNIC that provides network connectivity and computation power for customization [19, 38, 45, 70, 78]. Moreover, we believe that DCN should be programmable and reconfigurable: offering in-network computing [30, 31, 36], while capable of reconfiguring network topologies to best fit workload characteristics [68, 69]. A disaggregated control plane would oversee all the resource pools and the DCN, responsible for resource management, load balancing, etc. We hope this programmable and highly configurable infrastructure could quickly adapt to future workloads' fast-changing requirements.

While navigating the last layer, we propose principled guidelines for building disaggregated devices to address the third challenge. In particular, we identify three key features that any disaggregated devices should have: network connectivity, multi-tenancy and virtualization support, and security defenses, with the last being offered to environments where security is a must. A SmartNIC-empowered disaggregated device could easily offer all.

For the last challenge, we lay out a three-step approach. Similar to prior works [15, 68], we leverage reconfigurable topology to morph network topology to best match workload characteristics. Unlike prior works, we propose enabling the reconfigurability using both circuit switches and packet switches with cut-through forwarding. Second, we expose both network's programmability and reconfigurability to the IR layer, so workloads can leverage both at the same time. Third, we adapt the network disaggregation idea [56] to consolidate networking hardware into a logical computing plane, orchestrated by a global SDN controller. Consequently, the network becomes an ordinary computing resource of which the workloads can take full control.

We now briefly discuss how a distributed database system could benefit from our proposals. First, developers can annotate SQL operations with distinct network, security, and fault-tolerance requirements, e.g., use a device with a trusted execution environment. Then, our IR layer will decompose the database code into multiple codelets and compile each codelet into multiple binaries for heterogeneous devices (e.g., FPGA, RMT, and CPU). During runtime, we will schedule codelets following their embedded policies. If we detect that a codelet is better served on a different device, we will launch a migration. For instance, we can migrate a codelet with SQL aggregations or caching onto a p4 switch to reduce network round-trips. The codelet can also reconfigure network topologies to, say, create virtual racks across multiple physical racks to reduce network latency [68].

With these initial proposals, this vision paper outlines a possible path for data center designers to deal with the fast-changing workloads and hardware landscape. The solutions we proposed are by no means complete. We call for more contribution in this space.

2 Background

This section presents recent networking, hardware, and compiler trends motivating our vision.

2.1 Programmable and Reconfigurable Network

As cloud traffic has doubled roughly every year since 2005 [7], the data center networking infrastructure is constantly changing and is never short of innovations. Around the late 2000s, software-defined networking was proposed to improve management efficiency by decoupling the *control plane* (which decides how to handle the traffic) from the *data plane* (which forwards traffic according to decisions that the control plane makes) and then consolidating the control plane onto a set of commodity servers [16, 18, 34]. Recently, the p4 switch [9] and emerging heterogeneous networking devices [29, 56] further empower the data plane with unprecedented programmability and in-network computing at various link locations. Nowadays, both the data and control planes in data center networks are programmable and able to run customized user computation [18, 31].

Furthermore, emerging switching solutions and fabrics are challenging the status quo on how we interconnect data center hardware. For instance, optical circuit switch can reconfigure the network topologies by changing the physical connections among devices to best fit workload's network traffic pattern [68, 69], rendering a drastic departure from the fixed-topology models [25, 59]. On the other hand, emerging fabrics such as Gen-Z [21] and CXL [13, 23, 37, 44] offer extremely low-latency interconnect solutions for disaggregated devices. In all, the DCN infrastructure is more programmable, modularized, and flexible than ever. It is the cornerstone to realizing a fully programmable and disaggregated data center.

2.2 Hardware Resource Disaggregation

Hardware resource disaggregation is a proposal that breaks regular servers into segregated, network-attached hardware resource pools. It promises to improve a data center's manageability and resource utilization. The disaggregation idea roughly went through four epochs since its inception. Initially, around 2009, it was proposed to address memory capacity issue [40]. Couple of years went back without much progress. As the network continued improving, there was a renewed interest around 2016. The second epoch sees plenty of study, infrastructure, low-level systems targeting disaggregation [20, 26, 55]. The third epoch moves up the stack

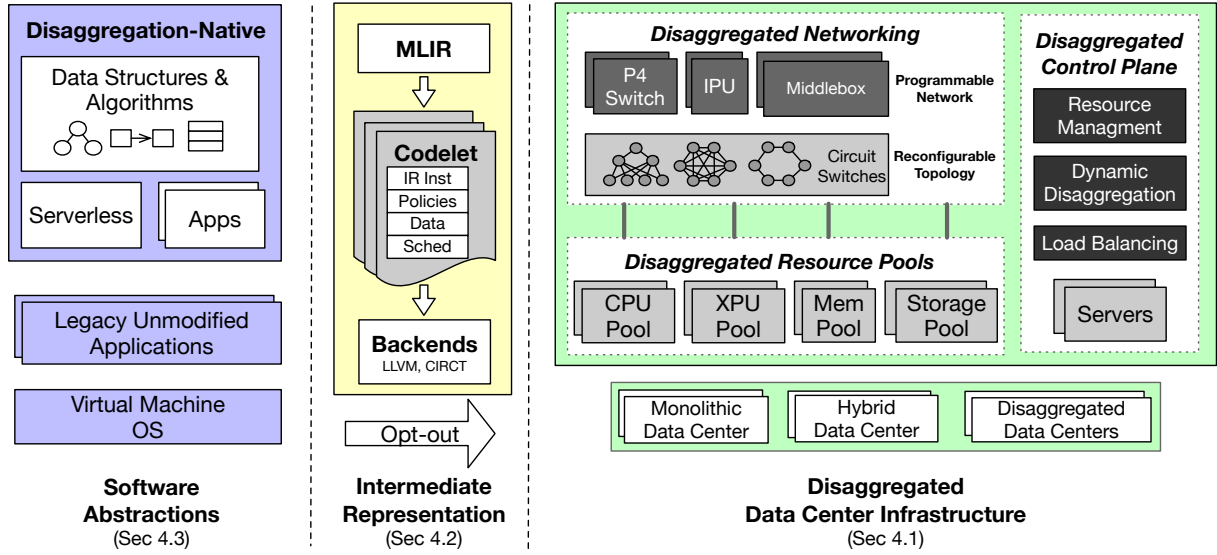


Figure 1: Overview of a Fully Disaggregated and Programmable Data Center. We expect two abstractions, one supports legacy VM workloads and another tailored for emerging serverless workloads. We propose to use an IR layer to help reduce software fragmentation. For disaggregated devices within each pool, we believe that every device is empowered by a SmartNIC-alike device. Networking is the first-class citizen, disaggregated as its own resource pool providing programmability and topology reconfigurability.

with systems co-designing applications or runtime with hardware [53, 64, 65]. The current fourth epoch closely examines how to best match application programming models with the distributed hardware [62] and how to deploy it [22].

Typically, a disaggregated device is a self-contained, standalone board that works without a companion server. In principle, it is directly attached to the DCN fabric. Multiple such devices could be tightly packaged into a standard chassis for a tighter deployment, but they still function and fail independently [37]. Generally, a disaggregated device has two categories of resources: a dominant resource that defines a device’s *personality* and auxiliary resources that make the device functional, *e.g.*, DRAM is a memory device’s personality. The auxiliary resources are offered by co-packaged SmartNICs enabling network connectivity, memory caching, and extra computation power [2, 5, 19, 38, 45, 56, 70, 78].

2.3 Compilers and Programming Frameworks

The confluence of domain-specific abstractions (*e.g.*, graph operators) and accelerators (*e.g.*, GPU and FPGA) calls for the next-generation compilers and programming language (PL) frameworks to reduce software fragmentation and improve compilation for heterogeneous devices. MLIR [35] is a typical example that consolidates high-level abstractions onto a shared framework for common code analysis and optimizations across a set of heterogeneous hardware. We chose to build our IR using MLIR because of this.

3 Overview

In a fully disaggregated data center, all devices and networking hardware are programmable and disaggregated from each other. Moreover, the network topology can be dynamically reconfigured to best match workload requirements. Figure 1 presents the end-to-end deployment and architecture overview of such a data center, which has three main layers: (1) a disaggregated data center infrastructure layer, the key enabling factor of our envisioned model, (2) an IR layer, and (3) an abstraction layer with multiple abstractions and a proposed design methodology. The latter two layers together *enhance* the usability of the infrastructure layer.

Disaggregated Data Center Infrastructure. At a high-level, the infrastructure has three components. (1) The disaggregated resource pools each hosts a type of resource (*e.g.*, compute, memory, and storage) and can be built, scaled, and managed independently. Existing systems built for disaggregation [8, 55] could run top of them without changes. (2) The networking infrastructure consists of circuit switches and programmable networking devices. Network functionalities are first disaggregated from other resource pools and then consolidated into a standalone networking pool, following the network disaggregation idea initially proposed by SuperNIC [56]. Here, we are enlarging its original scope by adding more networking devices into the networking pool, also using circuit switches to enable reconfigurable topologies to adjust topology based on software needs dynamically. (3) Finally, a control plane runs on a small set of regular servers which

is responsible for managing global resources, balancing load, and handling hardware failures.

Software Abstraction. We embrace two abstractions and one methodology to enhance hardware usability. Despite hardware disaggregation’s drastic departure from the traditional computing paradigm, we expect a backward compatible abstraction such that legacy applications and VMs can run without changes. In other words, we need a vNode-like abstraction [55]. But backward compatibility and transparency come at a price. We need another DAG-like abstraction that permits decentralized execution of the application logic, pioneered by FractOS [62]. This abstraction captures the workload serverless-y trend. Besides abstractions, we propose a disaggregation-native, or disaggregation-conscious methodology. The core idea is to expose the heterogeneous hardware nature and co-design software with it. In particular, we advocate for non-transparency and explicit code annotation to expose heterogeneity, reconfigurable network topologies, and hardware failures to the software. Consequently, users can take full control of their data and computation and customize their failure handling to achieve a desired balance among performance, reliability, and cost.

Intermediate Representation (IR). We leverage the IR layer to bridge the gap between the fast changing workloads and the diverse hardware infrastructure. It helps reduce software fragmentation and improve compilation flow for heterogeneous devices. We use MLIR [35] to decompose an application into multiple smaller and disaggregatable *codelets*, together with a companion DAG dictating the execution order of the codelets. The codelet is a rich entity designed to exploit the underlying hardware infrastructure. A codelet consists of data, purposely-defined IR instructions, scheduling primitives, network topology configurations, generated state-management APIs, and a set of policies such as security requirements. We compile the codelet for distinct hardware targets using backends such as LLVM or CIRCT [1]. A codelet is the smallest unit of data and computation migration. The IR is not a one-size-fits-all solution. It would benefit serverless-y workloads expressed as DAGs more than legacy applications running on top of VMs.

4 Deep Dive

We now take a closer look at every layer required to build a disaggregated data center and present our initial proposals. We start with the hardware infrastructure, which is characterized by disaggregated devices, the DCN, and a control plane to tie the two together. We then detail the IR layer and the disaggregation-native methodology.

4.1 Disaggregated Devices

In principle, disaggregated devices are directly attached to the DCN fabric. They are also programmable thanks to co-packaged SmartNIC-like devices with computational resources (§2.2). We identify three core functionalities that any disaggregated devices should consider to incorporate.

The first is **network connectivity**, the must-have feature. Many solutions are available. We categorize them into two types. The first type enables basic (“dummy”) network connectivity. Its sole job is sending and receiving data packets. Ethernet and PCIe fall into this category. The second type enables advanced (“smart”) network connectivity. It could have built-in advanced transport and other higher-level services such as remote memory translation. Typical solutions are RoCE, CXL, and Gen-Z. Although the first type only exposes a bare-metal network interface, it allows flexible customization and co-design with other resources. The second type has a more consistent and easy-to-use interface across devices, but the hardened features discourage modifications and may limit scalability [78]. When deciding on a network connectivity solution, one must consider infrastructure availability. For instance, Ethernet is widely available but emerging fabrics are still under development [22, 37, 44].

The second is **multi-tenancy and virtualization support**, the one provides protection, isolation, and fairness for shared resources. Multi-tenancy is a long-standing subject and has been extensively studied on conventional computing resources (*e.g.*, CPU, GPU, and FPGA), but applying it to disaggregated devices poses new challenges. First, since a disaggregated device could have more than one computing resource, a correct multi-tenancy design therefore requires diligent coordination across resource boundaries. For instance, a recent disaggregated memory device has ASIC, FPGA, and CPU, its authors went to great lengths to ensure proper isolation [78]. The second key challenge arises when novel computing resources are used, especially those with poor multi-tenancy support. For example, RMT, the core computing unit of a programmable switch, currently has no commercial support for multi-tenancy but only has a few research proposals [17, 67, 71, 72].

Last but not least, **security features** could be added for environments that require in-depth security guarantees. We identify that a disaggregated device can provide three levels of security defenses. The first level is ensuring basic data encryption, authentication, and authorization. It secures the device from potentially unwarranted accesses and from using any tempted data. Unfortunately, security mechanisms required for this level are missing in many recent devices [58, 70, 78]. The last two security levels are avoiding covert/side-channels and delivering confidential computing. Both ensure stronger security/privacy guarantees and grant users to off-load highly-sensitive computation and data to disaggregated

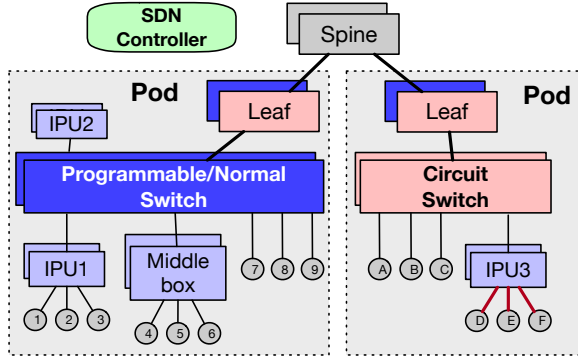


Figure 2: Proposed DCN Infrastructure. A pod with thousands of devices is the unit of deployment, striking a balance between scalability and performance. The gray circles are devices or servers. ToR can be programmable, circuit, or normal switch. Emerging network devices can be directly attached to a ToR switch. Network programmability is enabled by all the blue-colored boxes. Dynamic topology is enabled by red-colored circuit switches and blue-colored devices with cut-through forwarding. Black links are standard Ethernet links. Red links could be novel links such as GNet [22].

devices. Though both have been studied individually for resources like CPU [12], FPGA [76, 77], and GPU [28, 63], they share the same challenges when applied to disaggregated devices. For example, there is no trusted execution environment (TEE) in any programmable switches.

In all, we believe future disaggregated devices *should* invest in all or some of the following features: network connectivity, multi-tenancy and virtualization support, and security defenses. The key technical challenges are dealing with multiple or novel computing resources.

4.2 Data Center Networking

This section discusses issues related to the DCN. Figure 2 shows the envisioned architecture in which the DCN is organized as pods and managed by a logically centralized SDN controller [6, 18, 22].

Topology and Deployment Scale. The DCN is the cornerstone to enabling disaggregation. A tight deployment within a limited scale ensures stable and low communication latency, the key to making disaggregation practical [20, 22, 55]. Hence, we believe future disaggregated data center deployment will be contained within a limited scale due to DCN latency constraint, striking a balance between scalability and performance (e.g., within a 10 us pod with roughly few thousands of devices or servers).

Network Programmability. Network data plane’s programmability is enabled by programmable switch and middlebox [9, 73] as well as emerging networking devices such as IPU [29] and SuperNIC [56]. Together, they make a rich set of resources such as RMT, FPGA, and CPU available to perform in-network computing. We treat the network as

another first-class disaggregated resource. We view both networking switches/devices as regular disaggregated devices, following the same design principles discussed in §4.1. This unification simplifies and provides principled guidelines for future networking device design.

We pose no restriction on how switches, networking devices, and other devices are interconnected. For instance, in Figure 2, IPU2 is connected to a switch while IPU1 is connected to both disaggregated devices and the switch. Besides standard Ethernet or Infiniband, we advocate using novel link layer solutions with features such as end-to-end flow control, lossless transmission, and incast prevention to break the rigorous boundaries in layered protocols [22, 24].

Dynamic Reconfigurable Topology. We propose to enable dynamic reconfigurable topology on top of fixed physical network deployment. Specifically, we want to build conceptual point-to-point connections among selective devices to avoid in-network buffering. At its core, we use circuit-switching to create temporal links and reconfigurable topologies are realized by adjusting those temporal links. As Figure 2 shows, we think both optical circuit switches and packet-based switches with cut-through forwarding can be used. Nonetheless, we believe a key open challenge is to develop an efficient scheduling policy co-designed with the infrastructure along with the running workloads [15, 57, 68, 69].

4.3 Control Plane

The disaggregated control plane oversees all the disaggregated devices and the networking infrastructure. It is responsible for global resource management, health monitoring, load balancing, failure handling, etc. It is similar to the Global Resource Managers proposed by LegoOS [55].

For resource management, we embrace a *two-level approach* in which the control plane only does coarse-grained allocation while the specific devices do finer-grained ones. We would expose a menu of abstractions for failure handling such as best-effort handling and transparent handling. The former exposes failures to applications while the latter hides them. Users could choose the one they see fit.

We observe that the control plane’s most challenging task will be balancing computation, data, and packets among disaggregated devices at high speed during runtime. It involves a complicated choice matrix concerning device capability, data location, network topology, and network bandwidth. Traditional methods may fall short in making such decisions fast enough to produce reasonable outcomes. In response, we identify two non-exclusive approaches to tackle this issue. First, we could utilize reinforcement learning to transform it into a learning problem. Similar methods are proven effective for OS [75], database [48], and data structures [14]. Second, we plan to *onload* this task from the control plane into the

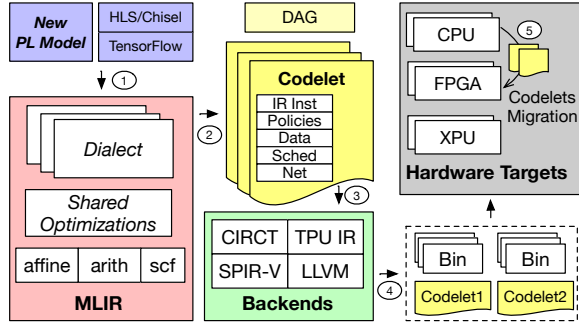


Figure 3: IR and Compilation Flow. *CIRCT is a backend for FPGA. SPIR-V is an IR for GPU. A codelet can be compiled into multiple executable binaries, or bitstreams.*

application layer by introducing explicit data movement and scheduling primitives in the IR layer (§4.4).

4.4 Intermediate Representation

To facilitate the use of heterogeneous devices and runtime migration, we propose to use MLIR to decompose an application into smaller *codelets* and a companion DAG dictating the execution order of codelets. A codelet encapsulates data, computation, policies, network topology configurations, and scheduling primitives. We can compile a codelet into binaries or bitstreams for distinct hardware targets together with generated APIs for codelet’s state management

As Figure 3 shows, ① the MLIR will take existing frameworks or new PL models as inputs. Within MLIR, there could be multiple domain-specific dialects for distinct inputs, *e.g.*, a *p4-dialect* for P4 programs. We propose to use a shared layer to carry out common optimizations and security checks [11]. ② The optimized dialects are then lowered onto a purposely-defined IR. We will package the IR instructions together with other policies and configurations into codelets. ③ Subsequently, the codelets are transformed onto various backends for final compilation. ④ Multiple binaries or bitstreams will be produced for each codelet together with generated APIs that can interact with the states within a codelet binary. A full IR definition is beyond the scope of this paper.

The codelet is a powerful entity designed to fully exploit the underlying hardware infrastructure. First, it can reconfigure network topologies by orchestrating network switches and devices to create temporal link connections. Second, it has explicit intra-codelet data movement and scheduling primitives to overlap computation and communication within a codelet. For inter-codelet scheduling, we plan to add a specially defined DAG-codelet. Third, a codelet serves as the smallest logical unit of migration among devices. To migrate a codelet, we would launch a previously compiled binary of this codelet onto the target hardware and then migrate states using the generated APIs.

4.5 Disaggregation-Native Design Methodology

Disaggregation-native design methodology’s code idea is to co-design software with the heterogeneous infrastructure nature. Many prior works have explored building OS, indexes, trees, and databases for disaggregation [47, 53, 55, 61, 66, 74, 79]. Following their insights, we summarize and propose the following design principles.

First, **enable non-transparency and explicit code annotation**. Software should explicitly annotate their data and code to control placement, co-location, and security. This principle trades ease-of-programming for avoiding unwanted overhead. Our proposed IR can transform annotations into explicit scheduling primitives. Second, **expose failures and enable configurable reliability**. We expect to expose device or network failures to software rather than masking them silently. Consequently, software can customize failure handling logic and choose the right level of reliability (*e.g.*, use erasure-coding rather than three-way replication). Third, **enable network topology-aware designs**. We expect software to be topology-aware and can customize the network topologies to best fit their data access and computation offloading schemes. For example, distributed ML training could morph the topologies to fit its all-reduce pattern. The proposed IR can help generate low-level network reconfiguration primitives. This practice leads to reduced network buffering hence better performance.

5 Conclusion

In this paper, we examine the key challenges due to the recent networking, disaggregation, and serverless trends. We envision a futuristic disaggregated and programmable data center along with a set of design guidelines for data center designers to work towards the north star. In particular, we describe the key features required to build disaggregated devices and discuss DCN’s impact on deployment. We propose an MLIR-based IR layer to reduce software fragmentation and bridge the gap between fast-changing software and diverse hardware. We embrace two abstractions and propose a design methodology to enhance disaggregated data center’s usability. Our vision involves many disciplines and their intersections, hence the solutions we presented are by no means complete. We hope more researchers can get involved and contribute.

Acknowledgment

We would like to thank Yizhou Shan’s thesis committee members, Geoffrey M. Voelker, Stefan Savage, Alex C. Snoeren, and George C. Papan. Their insightful questions during his defense motivate us to write this paper. We thank the anonymous reviewers, Sanidhya Kashyap, Chenxi Wang, Qizhen Zhang, Pengfei Zuo, and Kevin Zhang for their feedback.

References

- [1] CIRCT. <https://circt.llvm.org/>.
- [2] Vmware project monterey snapshot. <https://blogs.vmware.com/vsphere/2022/04/project-monterey-snapshot.html>.
- [3] Marcos K Aguilera, Kimberly Keeton, Stanko Novakovic, and Sharad Singhal. Designing far memory data structures: Think outside the box. In *Proceedings of the Workshop on Hot Topics in Operating Systems*, 2019.
- [4] Amazon. Amazon EC2 Elastic GPUs. <https://aws.amazon.com/ec2/elastic-gpus/>.
- [5] Amazon. AWS Nitro System. <https://aws.amazon.com/ec2/nitro/>.
- [6] Manikandan Arumugam, Deepak Bansal, Navdeep Bhatia, James Boerner, Simon Capper, Changhoon Kim, Sarah McClure, Neeraj Motwani, Ranga Narasimhan, Urvis Panchal, Tommaso Pimpo, Ariff Premji, Pranjal Shrivastava, and Rishabh Tewari. Bluebird: High-performance SDN for bare-metal cloud services. In *19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22)*, 2022.
- [7] Hitesh Ballani, Paolo Costa, Raphael Behrendt, Daniel Cletheroe, Istvan Haller, Krzysztof Jozwik, Fotini Karinou, Sophie Lange, Benn Thomsen, Kai Shi, and Hugh Williams. Sirius: A flat datacenter network with nanosecond optical switching. In *SIGCOMM 2020*. ACM, August 2020.
- [8] Laurent Bindschaedler, Ashvin Goel, and Willy Zwaenepoel. Hailstorm: Disaggregated compute and storage for distributed lsm-based databases. In *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS '20, 2020.
- [9] Pat Bosshart, Glen Gibb, Hun-Seok Kim, George Varghese, Nick McKeown, Martin Izzard, Fernando Mujica, and Mark Horowitz. Forwarding metamorphosis: Fast programmable match-action processing in hardware for sdn. *SIGCOMM Comput. Commun. Rev.*, 2013.
- [10] Matthew Burke, Sowmya Dharanipragada, Shannon Joyner, Adriana Szekeres, Jacob Nelson, Irene Zhang, and Dan R. K. Ports. Prism: Rethinking the rdma interface for distributed systems. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles*, SOSP '21, 2021.
- [11] David Chisnall, Deepak Garg, Catalin Hritcu, and Mathias Payer. Secure Compilation, Report from Dagstuhl Seminar 21481. 2022. <https://nebelwelt.net/files/21DAGSTUHL.pdf>.
- [12] Costan, Victor and Devadas, Srinivas. Intel SGX Explained. <https://eprint.iacr.org/2016/086.pdf>.
- [13] CXL Consortium. <https://www.computeexpresslink.org/>.
- [14] Yifan Dai, Yien Xu, Aishwarya Ganesan, Ramnathan Alagappan, Brian Kroth, Andrea Arpaci-Dusseau, and Remzi Arpaci-Dusseau. From WiscKey to bourbon: A learned index for Log-Structured merge trees. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*.
- [15] Nathan Farrington, George Porter, Sivasankar Radhakrishnan, Hamid Hajabdolali Bazzaz, Vikram Subramanya, Yeshaiah Fainman, George Papen, and Amin Vahdat. Helios: A hybrid electrical/optical switch architecture for modular data centers. *SIGCOMM Comput. Commun. Rev.*, 2010.
- [16] Nick Feamster, Jennifer Rexford, and Ellen Zegura. The road to sdn: An intellectual history of programmable networks. *SIGCOMM Comput. Commun. Rev.*, 2014.
- [17] Yong Feng, Zhikang Chen, Haoyu Song, Wenquan Xu, Jiahao Li, Zijian Zhang, Tong Yun, Ying Wan, and Bin Liu. Enabling in-situ programmability in network data plane: From architecture to language. In *19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22)*, 2022.
- [18] Andrew D. Ferguson, Steve Gribble, Chi-Yao Hong, Charles Killian, Waqar Mohsin, Henrik Muehe, Joon Ong, Leon Poutievski, Arjun Singh, Lorenzo Vicisano, Richard Alimi, Shawn Shuoshuo Chen, Mike Conley, Subhasree Mandal, Karthik Nagaraj, Kondapa Naidu Bollineni, Amr Sabaa, Shidong Zhang, Min Zhu, and Amin Vahdat. Orion: Google's Software-Defined networking control plane. In *18th USENIX Symposium on Networked Systems Design and Implementation (NSDI 21)*, 2021.
- [19] Daniel Firestone, Andrew Putnam, Sambhrama Mundkur, Derek Chiou, Alireza Dabagh, Mike Andrewartha, Hari Angepat, Vivek Bhanu, Adrian Caulfield, Eric Chung, Harish Kumar Chandrappa, Somesh Chaturmohta, Matt Humphrey, Jack Lavier, Norman Lam, Fengfen Liu, Kalin Ovtcharov, Jitu Padhye, Gautham Popuri, Shachar Raindel, Tejas Sapre, Mark Shaw, Gabriel Silva, Madhan Sivakumar, Nisheeth Srivastava, Anshuman Verma, Qasim Zuhair, Deepak Bansal, Doug Burger, Kushagra Vaid, David A. Maltz, and Albert Greenberg. Azure Accelerated Networking: SmartNICs in the Public Cloud. In *15th USENIX Symposium on Networked Systems Design and Implementation (NSDI '18)*.
- [20] Peter X. Gao, Akshay Narayan, Sagar Karandikar, Joao Carreira, Sangjin Han, Rachit Agarwal, Sylvia Ratnasamy, and Scott Shenker. Network Requirements for Resource Disaggregation. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI '16)*.
- [21] GenZ Consortium. <http://genzconsortium.org/>.
- [22] Dan Gibson, Hema Hariharan, Eric Lance, Moray McLaren, Behnam Montazeri, Arjun Singh, Stephen Wang, Hassan M. G. Wassel, Zhehua Wu, Sunghwan Yoo, Raghuraman Balasubramanian, Prashant Chandra, Michael Cutforth, Peter Cuy, David Decotigny, Rakesh Gautam, Alex Iriza, Milo M. K. Martin, Rick Roy, Zuowei Shen, Ming Tan, Ye Tang, Monica Wong-Chan, Joe Zbiciak, and Amin Vahdat. Aquila: A unified, low-latency fabric for datacenter networks. In *19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22)*, 2022.
- [23] Donghyun Gouk, Sangwon Lee, Miryeong Kwon, and Myoungsoo Jung. Direct access, High-Performance memory disaggregation with DirectCXL. In *2022 USENIX Annual Technical Conference (USENIX ATC 22)*, 2022.
- [24] Prateesh Goyal, Preey Shah, Kevin Zhao, Georgios Nikolaidis, Mohammad Alizadeh, and Thomas E. Anderson. Backpressure flow control. In *19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22)*, 2022.
- [25] Albert Greenberg, James R. Hamilton, Navendu Jain, Srikanth Kandula, Changhoon Kim, Parantap Lahiri, David A. Maltz,

- Parveen Patel, and Sudipta Sengupta. VI2: A scalable and flexible data center network. In *Proceedings of the ACM SIGCOMM 2009 Conference on Data Communication*, SIGCOMM '09, 2009.
- [26] Juncheng Gu, Youngmoon Lee, Yiwen Zhang, Mosharaf Chowdhury, and Kang Shin. Efficient Memory Disaggregation with Infiniswap. In *Proceedings of the 14th USENIX Symposium on Networked Systems Design and Implementation (NSDI '17)*.
- [27] Hewlett-Packard. The Machine: A New Kind of Computer. <http://www.hpl.hp.com/research/systems-research/themachine/>.
- [28] Tyler Hunt, Zhipeng Jia, Vance Miller, Ariel Szekely, Yige Hu, Christopher J. Rossbach, and Emmett Witchel. Telekine: Secure computing with cloud GPUs. In *17th USENIX Symposium on Networked Systems Design and Implementation (NSDI 20)*.
- [29] Intel. Intel Unveils Infrastructure Processing Unit. <https://www.intel.com/content/www/us/en/newsroom/news/infrastructure-processing-unit-data-center.html>.
- [30] Xin Jin, Xiaozhou Li, Haoyu Zhang, Nate Foster, Jeongkeun Lee, Robert Soulé, Changhoon Kim, and Ion Stoica. Netchain: Scale-free sub-rtt coordination. In *15th USENIX Symposium on Networked Systems Design and Implementation (NSDI 18)*, 2018.
- [31] Xin Jin, Xiaozhou Li, Haoyu Zhang, Robert Soulé, Jeongkeun Lee, Nate Foster, Changhoon Kim, and Ion Stoica. Nocache: Balancing key-value stores with fast in-network caching. In *Proceedings of the 26th Symposium on Operating Systems Principles*, SOSP '17, pages 121–136, 2017.
- [32] Norman P. Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, Rick Boyle, Pierre Luc Cantin, Clifford Chao, Chris Clark, Jeremy Coriell, Matt Dau Mike Daley, Jeffrey Dean, Ben Gelb, Tara Vazir Ghaemmaghami, Rajendra Gottipati, William Gulland, Robert Hagmann, C. Richard Ho, Doug Hogberg, John Hu, Robert Hundt, Dan Hurt, Julian Ibarz, Aaron Jaffey, Alek Jaworski, Harshit Khaitan Alexander Kaplan, Andy Koch, Naveen Kumar, Steve Lacy, James Laudon, James Law, Diemthu Le, Chris Leary, Zhuyuan Liu, Kyle Lucke, Alan Lundin, Gordon MacKean, Adriana Maggiore, Maire Mahony, Kieran Miller, Rahul Nagarajan, Ravi Narayanaswami, Ray Ni, Kathy Nix, Thomas Norrie, Mark Omernick, Narayana Penukonda, Andy Phelps, Matt Ross Jonathan Ross, Amir Salek, Emad Samadiani, Chris Severn, Gregory Sizikov, Matthew Snelham, Jed Souter, Andy Swing Dan Steinberg, Mercedes Tan, Gregory Thorson, Bo Tian, Horia Toma, Erick Tuttle, Vijay Vasudevan, Richard Walter, Walter Wang, Eric Wilcox, and Doe Hyun Yoon. In-Datacenter Performance Analysis of a Tensor Processing Unit. In *Proceedings of the 44th Annual International Symposium on Computer Architecture (ISCA '17)*.
- [33] Ana Klimovic, Christos Kozyrakis, Eno Thereska, Binu John, and Sanjeev Kumar. Flash storage disaggregation. In *Proceedings of the Eleventh European Conference on Computer Systems*, EuroSys '16, 2016.
- [34] Teemu Koponen, Martin Casado, Natasha Gude, Jeremy Stribling, Leon Poutievski, Min Zhu, Rajiv Ramanathan, Yuichiro Iwata, Hiroaki Inoue, Takayuki Hama, and Scott Shenker. Onix: A distributed control platform for large-scale production networks. In *Proceedings of the 9th USENIX Conference on Operating Systems Design and Implementation*, OSDI'10, 2010.
- [35] Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, and Oleksandr Zinenko. Mlir: A compiler infrastructure for the end of moore's law. *arXiv preprint arXiv:2002.11054*, 2020.
- [36] Seung-seob Lee, Yanpeng Yu, Yupeng Tang, Anurag Khandelwal, Lin Zhong, and Abhishek Bhattacharjee. Mind: In-network memory management for disaggregated data centers. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles*, SOSP '21, 2021.
- [37] Huaicheng Li, Daniel S Berger, Stanko Novakovic, Lisa Hsu, Dan Ernst, Pantea Zardoshti, Monish Shah, Ishwar Agarwal, Mark Hill, Marcus Fontoura, et al. First-generation memory disaggregation for cloud platforms. *arXiv preprint arXiv:2203.00241*, 2022.
- [38] Huaicheng Li, Mingzhe Hao, Stanko Novakovic, Vaibhav Gogte, Sriram Govindan, Dan RK Ports, Irene Zhang, Ricardo Bianchini, Haryadi S Gunawi, and Anirudh Badam. Leapio: Efficient and portable virtual nvme storage on arm socs. In *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems*, 2020.
- [39] Yuliang Li, Rui Miao, Hongqiang Harry Liu, Yan Zhuang, Fei Feng, Lingbo Tang, Zheng Cao, Ming Zhang, Frank Kelly, Mohammad Alizadeh, and Minlan Yu. Hpccl: High precision congestion control. In *Proceedings of the ACM Special Interest Group on Data Communication*, 2019.
- [40] Kevin Lim, Jichuan Chang, Trevor Mudge, Parthasarathy Ranganathan, Steven K. Reinhardt, and Thomas F. Wenisch. Disaggregated Memory for Expansion and Sharing in Blade Servers. In *Proceedings of the 36th Annual International Symposium on Computer Architecture (ISCA '09)*.
- [41] Jiaxin Lin, Kiran Patel, Brent E. Stephens, Anirudh Sivaraman, and Aditya Akella. PANIC: A high-performance programmable NIC for multi-tenant networks. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, 2020.
- [42] Ming Liu, Tianyi Cui, Henry Schuh, Arvind Krishnamurthy, Simon Peter, and Karan Gupta. Offloading distributed applications onto smartnics using ipipe. In *Proceedings of the ACM Special Interest Group on Data Communication*, SIGCOMM '19, 2019.
- [43] Jiacheng Ma, Gefei Zuo, Kevin Loughlin, Haoyang Zhang, Andrew Quinn, and Baris Kasikci. Debugging in the brave new world of reconfigurable hardware. In *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS 2022, 2022.
- [44] Hasan Al Maruf, Hao Wang, Abhishek Dhanotia, Johannes Weiner, Niket Agarwal, Pallab Bhattacharya, Chris Petersen,

- Mosharaf Chowdhury, Shobhit Kanaujia, and Prakash Chauhan. Tpp: Transparent page placement for cxi-enabled tiered memory. *arXiv preprint arXiv:2206.02878*, 2022.
- [45] Jaehong Min, Ming Liu, Tapan Chugh, Chenxingyu Zhao, Andrew Wei, In Hwan Doh, and Arvind Krishnamurthy. Gimbal: Enabling multi-tenant storage disaggregation on smartnic jbofs. In *Proceedings of the 2021 ACM SIGCOMM 2021 Conference*, SIGCOMM '21, 2021.
- [46] Philipp Moritz, Robert Nishihara, Stephanie Wang, Alexey Tumanov, Richard Liaw, Eric Liang, Melih Elibol, Zongheng Yang, William Paul, Michael I. Jordan, and Ion Stoica. Ray: A distributed framework for emerging AI applications. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, 2018.
- [47] Stanko Novakovic, Yizhou Shan, Aasheesh Kolli, Michael Cui, Yiyang Zhang, Haggai Eran, Boris Pismenny, Liran Liss, Michael Wei, Dan Tsafir, and Marcos Aguilera. Storm: A fast transactional dataplane for remote data structures. In *Proceedings of the 12th ACM International Conference on Systems and Storage*, SYSTOR '19, 2019.
- [48] Andrew Pavlo, Gustavo Angulo, Joy Arulraj, Haibin Lin, Jiexi Lin, Lin Ma, Prashanth Menon, Todd C Mowry, Matthew Peron, Ian Quah, et al. Self-driving database management systems. In *CIDR*, volume 4, page 1, 2017.
- [49] Phitchaya Mangpo Phothilimthana, Ming Liu, Antoine Kaufmann, Simon Peter, Rastislav Bodik, and Thomas Anderson. Floem: A programming system for NIC-Accelerated network applications. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, 2018.
- [50] Yiming Qiu, Jiarong Xing, Kuo-Feng Hsu, Qiao Kang, Ming Liu, Srinivas Narayana, and Ang Chen. Automated smartnic offloading insights for network functions. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles*, SOSP '21, 2021.
- [51] Parthasarathy Ranganathan, Daniel Stodolsky, Jeff Calow, Jeremy Dorfman, Marisabel Guevara, Clinton Wills Smullen IV, Aki Kuusela, Raghu Balasubramanian, Sandeep Bhatia, Prakash Chauhan, et al. Warehouse-scale video acceleration: co-design and deployment in the wild. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, 2021.
- [52] Waleed Reda, Marco Canini, Dejan Kostić, and Simon Peter. RDMA is turing complete, we just did not know it yet! In *19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22)*, 2022.
- [53] Zhenyuan Ruan, Malte Schwarzkopf, Marcos K. Aguilera, and Adam Belay. AIFM: High-performance, application-integrated far memory. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, 2020.
- [54] Johann Schleier-Smith, Vikram Sreekanti, Anurag Khandelwal, Joao Carreira, Neeraja J Yadwadkar, Raluca Ada Popa, Joseph E Gonzalez, Ion Stoica, and David A Patterson. What serverless computing is and should become: The next phase of cloud computing. *Communications of the ACM*, 64(5):76–84, 2021.
- [55] Yizhou Shan, Yutong Huang, Yilun Chen, and Yiyang Zhang. Legoos: A disseminated, distributed OS for hardware resource disaggregation. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI '18)*.
- [56] Yizhou Shan, Will Lin, Ryan Kosta, Arvind Krishnamurthy, and Yiyang Zhang. Disaggregating and consolidating network functionalities with supernic. *arXiv preprint arXiv:2109.07744*, 2021.
- [57] Vishal Shrivastav, Asaf Valadarsky, Hitesh Ballani, Paolo Costa, Ki Suh Lee, Han Wang, Rachit Agarwal, and Hakim Weatherspoon. Shoal: A Network Architecture for Disaggregated Racks. In *16th USENIX Symposium on Networked Systems Design and Implementation (NSDI '19)*.
- [58] David Sidler, Zeke Wang, Monica Chiosa, Amit Kulkarni, and Gustavo Alonso. Strom: Smart remote memory. In *Proceedings of the Fifteenth European Conference on Computer Systems*, EuroSys '20, 2020.
- [59] Arjun Singh, Joon Ong, Amit Agarwal, Glen Anderson, Ashby Armistead, Roy Bannon, Seb Boving, Gaurav Desai, Bob Felderman, Paulie Germano, Anand Kanagala, Jeff Provost, Jason Simmons, Eiichi Tanda, Jim Wanderer, Urs Hölzle, Stephen Stuart, and Amin Vahdat. Jupiter rising: A decade of clos topologies and centralized control in google's datacenter network. In *Sigcomm '15*, 2015.
- [60] Arjun Singhvi, Aditya Akella, Dan Gibson, Thomas F. Wenisch, Monica Wong-Chan, Sean Clark, Milo M. K. Martin, Moray McLaren, Prashant Chandra, Rob Cauble, Hassan M. G. Wassel, Behnam Montazeri, Simon L. Sabato, Joel Scherpelz, and Amin Vahdat. Irma: Re-envisioning remote memory access for multi-tenant datacenters. In *Proceedings of the Annual Conference of the ACM Special Interest Group on Data Communication on the Applications, Technologies, Architectures, and Protocols for Computer Communication (SIGCOMM '20)*.
- [61] Shin-Yeh Tsai, Yizhou Shan, and Yiyang Zhang. Disaggregating Persistent Memory and Controlling Them Remotely: An Exploration of Passive Disaggregated Key-Value Stores. In *2020 USENIX Annual Technical Conference (USENIX ATC 20)*, 2020.
- [62] Lluís Vilanova, Lina Maudlej, Shai Bergman, Till Miemietz, Matthias Hille, Nils Asmussen, Michael Roitzsch, Hermann Härtig, and Mark Silberstein. Slashing the disaggregation tax in heterogeneous data centers with fractos. In *Proceedings of the Seventeenth European Conference on Computer Systems*, 2022.
- [63] Stavros Volos, Kapil Vaswani, and Rodrigo Bruno. Graviton: Trusted execution environments on GPUs. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*.
- [64] Chenxi Wang, Haoran Ma, Shi Liu, Yuanqi Li, Zhenyuan Ruan, Khanh Nguyen, Michael D. Bond, Ravi Netravali, Miryung Kim, and Guoqing Harry Xu. Semeru: A memory-disaggregated managed runtime. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, 2020.

- [65] Chenxi Wang, Haoran Ma, Shi Liu, Yifan Qiao, Jonathan Eyolfson, Christian Navasca, Shan Lu, and Guoqing Harry Xu. MemLiner: Lining up tracing and application for a Far-Memory-Friendly runtime. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*, 2022.
- [66] Qing Wang, Youyou Lu, and Jiwu Shu. Sherman: A write-optimized distributed b+ tree index on disaggregated memory. *arXiv preprint arXiv:2112.07320*, 2021.
- [67] Tao Wang, Hang Zhu, Fabian Ruffy, Xin Jin, Anirudh Sivaraman, Dan R. K. Ports, and Aurojit Panda. Multitenancy for fast and programmable networks in the cloud. In *12th USENIX Workshop on Hot Topics in Cloud Computing (HotCloud 20)*, 2020.
- [68] Weitao Wang, Dingming Wu, Sushovan Das, Afsaneh Rahbar, Ang Chen, and T. S. Eugene Ng. RDC: Energy-Efficient data center network congestion relief with topological reconfigurability at the edge. In *19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22)*, 2022.
- [69] Weiyang Wang, Moein Khazraee, Zhizhen Zhong, Zhijiao Jia, Dheevatsa Mudigere, Ying Zhang, Anthony Kewitsch, and Manya Ghobadi. Topoopt: Optimizing the network topology for distributed dnn training. *arXiv preprint arXiv:2202.00433*, 2022.
- [70] Zeke Wang, Hongjing Huang, Jie Zhang, Fei Wu, and Gustavo Alonso. FpgaNIC: An FPGA-based versatile 100gb SmartNIC for GPUs. In *2022 USENIX Annual Technical Conference (USENIX ATC 22)*, 2022.
- [71] Jiarong Xing, Kuo-Feng Hsu, Matty Kadosh, Alan Lo, Yonatan Piasetzky, Arvind Krishnamurthy, and Ang Chen. Runtime programmable switches. In *19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22)*, 2022.
- [72] Jiarong Xing, Yiming Qiu, Kuo-Feng Hsu, Hongyi Liu, Matty Kadosh, Alan Lo, Aditya Akella, Thomas Anderson, Arvind Krishnamurthy, TS Eugene Ng, et al. A vision for runtime programmable networks. In *Proceedings of the Twentieth ACM Workshop on Hot Topics in Networks*, 2021.
- [73] Kaiyuan Zhang, Danyang Zhuo, and Arvind Krishnamurthy. Gallium: Automated software middlebox offloading to programmable switches. In *Proceedings of the Annual Conference of the ACM Special Interest Group on Data Communication on the Applications, Technologies, Architectures, and Protocols for Computer Communication, SIGCOMM '20*, 2020.
- [74] Qizhen Zhang, Yifan Cai, Xinyi Chen, Sebastian Angel, Ang Chen, Vincent Liu, and Boon Thau Loo. Understanding the effect of data center resource disaggregation on production dbms. *Proc. VLDB Endow.*, may 2020.
- [75] Yiying Zhang and Yutong Huang. "learned" operating systems. *ACM SIGOPS Operating Systems Review*, 53(1):40–45, 2019.
- [76] Mark Zhao, Mingyu Gao, and Christos Kozyrakis. Shef: Shielded enclaves for cloud fpgas. *ASPLOS 2022*.
- [77] Mark Zhao and G. Edward Suh. Fpga-based remote power side-channel attacks. In *2018 IEEE Symposium on Security and Privacy (SP)*.
- [78] Zhiyuan Guo and Yizhou Shan and Xuhao Luo and Yutong Huang and Yiying Zhang. Clío: A hardware-software co-designed disaggregated memory system. In *the 27th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '22)*, Lausanne, Switzerland, March 2022.
- [79] Pengfei Zuo, Jiazhaoh Sun, Liu Yang, Shuangwu Zhang, and Yu Hua. One-sided RDMA-Conscious extendible hashing for disaggregated memory. In *2021 USENIX Annual Technical Conference (USENIX ATC 21)*, 2021.