



Rack-Scale Composable Infrastructure

Primer by GigaIO Networks

CONTENTS

Introduction.....	3
The importance of composable infrastructure for data center scalability – up, down or out ...	3
The Future is Disaggregated	4
Understanding what the disaggregated composable infrastructure is (and is not).....	5
Introducing GigaIO: innovator of a new disaggregated composable infrastructure	6
GigaIO Rack-Scale Composable Infrastructure offers unique and compelling advantages.....	7
Rack-Scale Composable Infrastructure with FabreX offers flexible data center scale	8
About GigaIO	10
Figure 1– Data Center Roadmap.	4
Figure 2– Driving out Latency	7
Figure 3 – TCO Example	9

This document is for informational purposes only and may contain typographical errors and technical inaccuracies. The content is provided as is, without express or implied warranties of any kind.

© GigaIO Networks, Inc. All rights reserved. GigaIO and its affiliates cannot be responsible for errors or omissions in typography or photography. GigaIO, the GigaIO logo, and FabreX are trademarks of GigaIO Networks Inc. Other trademarks and trade names may be used in this document to refer to either the entities claiming the marks and names or their products. GigaIO disclaims proprietary interest in the marks and names of others. December 2020 | Rev 1.0



Introduction

Do you ever attend meetings where the conversation turns to an unfamiliar technology? Suddenly, many names and acronyms are being thrown around and you have no idea what they mean. Confused, you scan the room and notice that everyone else is nodding as if they know exactly what's being discussed.

We've all been there before. How about CDI, HCI, CI, SDI...?? After you read this primer, you'll be in the know about CDI (Composable Disaggregated Infrastructure) and how it fits in with hyperconverged infrastructure (HCI) and converged infrastructure (CI). You will understand how GigalO's Rack-Scale Composable Infrastructure solution that takes a novel approach to disaggregating compute, storage, and networking into shared resource pools, can significantly reduce your total cost of ownership (TCO) – just to throw in another acronym - at data center scale.

The importance of composable infrastructure for data center scalability – up, down or out

Vast volumes of data are proliferating by the minute – and that calls for a fundamental change in how data centers are architected. The challenge of collecting, organizing, analyzing, moving, and storing so much information can be daunting.

At the same time, advances in data analytics and artificial intelligence (AI) domains such as machine learning (ML) and deep learning (DL) are driving the need for ever-greater processing performance – yet another imperative for the data center to evolve and adopt new approaches with bigger, faster, more secure hardware – accelerators, proprietary ASICs, and storage arrays. These highly complex systems are prone to experience bottlenecks and interconnect issues that drag down responsiveness and utilization.

As more applications introduce support for accelerators (GPUs and FPGAs), which can reduce time to result from weeks to literally minutes, users are clamoring for more of these expensive resources. Yet, industry data shows they are only utilized 15% to 20% of the time, stranded behind the traditional data center's rigid architecture.

Legacy data center infrastructures were not designed for today's workflow requirements. The scalable modern data center needs a solution that can integrate compute, storage and other communication I/O into a single-system cluster fabric, scaling resources up and out across the cluster as needed. This solution should free resources from their silos to be shared with other network users who draw from these resource pools through a *disaggregated composable infrastructure (DCI)*, an emerging category of infrastructure designed to maximize IT resource usage and improve business agility.



The Future is Disaggregated

“Infrastructure disaggregation is the fundamental underpinning of digital transformation” according to John Abbott of 451 Research. The CEO of no less than NVIDIA, Jensen Huang, has also heralded the entire data center as a single composable computer. Why are these industry luminaries hopping onto the disaggregation wagon? There are a few key reasons which are compelling IT managers to reevaluate the very architecture of their data centers.

1. **BUSINESS AGILITY** – As noted above, the legacy IT set up is becoming a bottleneck, and the high cost of making changes, combined with the long lead time to make those changes, prevent companies from entering new market and pursuing new opportunities;
2. **FINANCIAL EFFICACY** – The new consumption models - like subscription and pay-as-you-go - popularized by the private and public cloud make it possible to balance CapEx and OpEx spending like never before;
3. **CUSTOMER VALUE** - Artificial Intelligence and advanced analytics have increased expectations to use all available data to provide more context to decision-making.

As a result of the shifting ground described above, data centers are looking to deploy highly distributed infrastructure which self-optimizes, based on an underlying interconnected fabric where data becomes fully available, and processing can be done wherever and whenever it is needed.

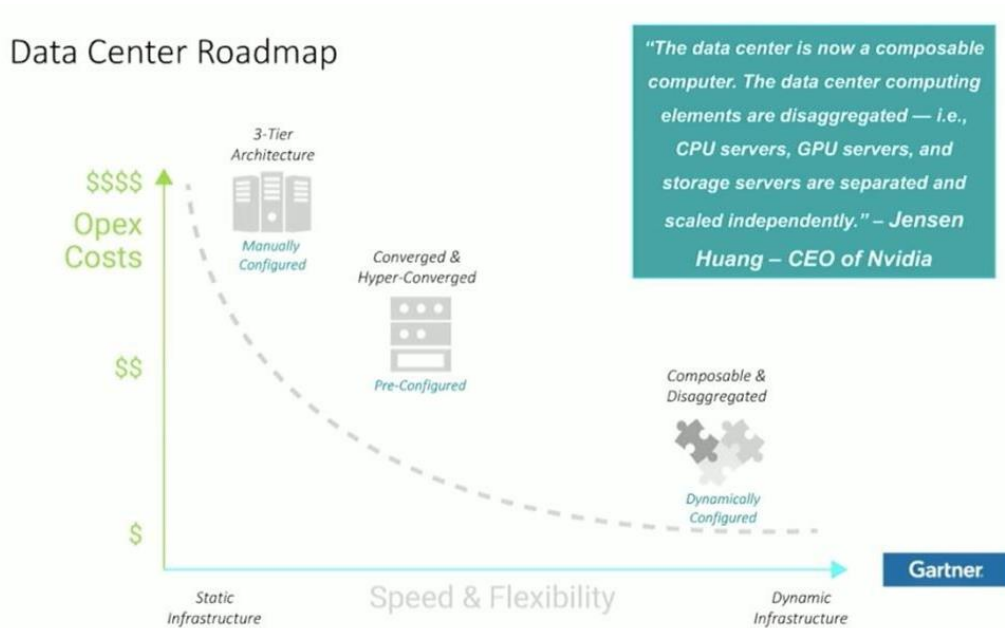


FIGURE 1 – DATA CENTER ROADMAP.



Until now, composable infrastructure choices have been limited to a few proprietary, single vendor solutions that are incompatible with legacy assets, forcing companies to restructure their data center to take advantage of true composability. They have also largely been limited to composing primarily or exclusively storage resources, and doing so over Ethernet as the physical interconnect. While Ethernet might be fine for legacy storage, the latest additions to the data center rack (accelerators and 3D XPoint SSDs, to name a few) cannot be composed efficiently with the inadequate latency and bandwidth of legacy networks. That's all changed now with a game-changing new Rack-Scale Composable Infrastructure from GigaIO.



Understanding what the disaggregated composable infrastructure is (and is not)

In a disaggregated composable infrastructure, server elements such as compute, storage, and network devices are treated as resource pools to be provisioned at will, depending on workflow requirements for optimum performance. Traditional server ingredients (storage, acceleration and so forth) are disaggregated into separate, easily shareable elements. Ideally, the solution should derive maximum utilization of these elements from reducing CapEx and OpEx. But, to disaggregate storage and accelerators, infrastructure interconnects need to support the lowest latency and highest bandwidth possible.

The key phrase is "*disaggregated composable*" – the ability to compose according to a specific workflow need, virtually drawing resources from a server and sharing them as needed. The industry is still learning about disaggregated composable infrastructure, and that unfamiliarity leads to several misconceptions:

- **MISCONCEPTION #1: disaggregation reduces performance** – In fact, the disaggregation and recomposing of optimal resources for each segment of a workflow can mean an actual gain in overall performance.
- **MISCONCEPTION #2: composability can happen on its own** – Unfortunately, it's not as easy as mix and match: before elements can be pooled, they must be disaggregated, which is the process of freeing and grouping the resources into pools while making them available to other resource users. Workflows can then share these now-pooled



resources. Composer software controls how much of a given disaggregated resource is allocated to each server, based on workflow.

- **MISCONCEPTION #3: composed servers need to be close** – Some IT professionals mistakenly believe physical location is important in the same way a GPU should be positioned next to the processor. The disaggregated composable infrastructure allows server resources to be widely dispersed within and outside the rack without losing performance.
- **MISCONCEPTION #4: DCI versus HCI** – A disaggregated composable infrastructure is entirely different from a hyper-converged infrastructure (HCI), a software-defined system in which all traditional data center elements – compute, networking, storage and management – are pooled into an integrated, unified solution, but trapped and usable only within a single system.
- **MISCONCEPTION #5: Composable infrastructures are always proprietary** – With the release of GigaIO's Rack-Scale Composable Infrastructure, based on open standards such as Redfish® APIs and PCIe, that's no longer so. Although other composable infrastructure vendors require their proprietary hardware and software – essentially creating new resource silos of their own – GigaIO can deliver the advantages of composability on virtually all disaggregated server systems and associated resources.
- **MISCONCEPTION #6: Disaggregated composability is a less secure infrastructure** – Not true. The ability to have secure interoperability at the PCIe layer end-to-end means a higher level of trust can be gained from this infrastructure class, which reduces intrusion points.

Introducing GigaIO: innovator of a new disaggregated composable infrastructure

GigaIO's goal is to advance the modern data center by creating the world's only enterprise-class, composable infrastructure – an open standards solution with boundless flexibility, cloudlike agility, future-readiness, security and affordability.

The GigaIO FabreX™ hyper-performance network is the culmination of that effort. FabreX eliminates conversion layers while enabling the data center to run at full efficiency with remarkably low latency and high throughput via a vastly improved interconnect fabric, eliminating the I/O network bottleneck.

Although other composable infrastructure solutions present themselves as "open," in fact, they require such limiting factors as proprietary hardware, software licenses, additional management panes of glass, and so forth. In contrast, GigaIO allows customers to select not only their preferred servers and components but also their composition software from a list of off-the-shelf options; besides, GigaIO is an active supporter of industry open-standards initiatives.

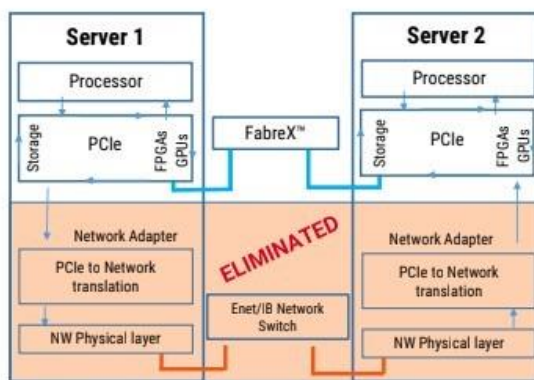


GigalO Rack-Scale Composable Infrastructure offers unique and compelling advantages

Flexibility on the only open platform – Other data center solutions might lock users into proprietary hardware and software resources that are incompatible with legacy systems, requiring costly forklift upgrades and limiting growth options. In contrast, GigalO's Rack-Scale Composable Infrastructure operates over the Redfish API, which is entirely agnostic, enabling users to run their choice of new or existing servers, storage, accelerators, and so forth.

Speed and simplicity of PCIe – The GigalO solution is 100% PCI Express-compliant, instantly and automatically working with all other PCIe-based resources, eliminating the need for additional, multi-step communication processes.

FabreX Drives out Latency to Deliver Disaggregation



Sub-Microsecond Latency

The Only Pure PCIe Network Throughout the Rack to Connect Both Resources and Servers

FIGURE 2– DRIVING OUT LATENCY

The industry's lowest latency and highest bandwidth – GigalO has the numbers to back up its claim:

- **Hardware latency of 350 ns end-to-end**, server to server, versus 1300 ns (and higher) for the alternatives. This industry-leading latency drives superior PCIe performance across the cluster.
- **Bandwidth of 256 Gb/sec up to 512 Gb/sec**, delivered by its third-generation FabreX implementation and fourth-gen solution.

These results validate GigalO's positioning as a young but proven leader in composable infrastructure solutions.

Improved Security – Unlike networks that rely on Ethernet switching to other networks, which creates the risk of a potential security injection point, GigalO is end-to-end PCIe with Trusted



Platform Module (TPM) on the switch, virtually eliminating the potential for undetected attacks. What's more, GigaIO works on servers with a Silicon Root of Trust embedded in hardware to resist firmware attacks, automatically detecting intrusion and rolling back to a last known safe state.

Reduce TCO by up to 50% – Users can easily share costly resources across multiple servers to achieve maximum utilization and prevent overprovisioning. This capability helps lower the CapEx that results when each compute or storage resource is on its upgrade cycle. What's more, reduced power consumption and automated management produce lower OpEx.

Rack-Scale Composable Infrastructure with FabreX offers flexible data center scale

Organizations worldwide are in continuous need of breakthrough data center solutions that help them acquire, analyze and glean meaningful insights and handle highly compute-intensive workflows quickly and securely.

GigaIO with its Rack-Scale Composable Infrastructure powered by the FabreX hyperperformance network, is a fundamentally new architecture that offers a new paradigm enabling organizations to scale data center resources up and out as needed – easily, rapidly and reliably.

Best of all, it allows for maximum resource utilization for lower TCO.

[Calculate your savings](#) with our TCO tool where you input your data and we provide you with results like these:



FabreX Composable Infrastructure

INFRASTRUCTURE INPUTS		
Number of GPU servers	0	8 GPUs/server
Number of App servers	8	CPU-only app servers
Number of JBOGs	2	8 GPUs/JBOG
Number of GPUs	16	
GPU model	V100-32	
Network b/w per port (Gb/s)	128	
Cluster Network type	Fabrex	
WORKLOAD INPUTS		
Number of jobs	150	
% of jobs that need GPUs	50%	
OUTPUTS		
Elapsed Run Time	151	hours

vs.

InfiniBand Converged Infrastructure

INFRASTRUCTURE INPUTS		
Number of GPU servers	4	8 GPUs/server
Number of App servers	4	CPU-only app servers
Number of JBOGs	0	8 GPUs/JBOG
Number of GPUs	32	
GPU model	V100-32	
Network b/w per port (Gb/s)	100	
Cluster Network type	InfiniBand	
WORKLOAD INPUTS		
Number of jobs	150	
% of jobs that need GPUs	50%	
OUTPUTS		
Elapsed Run Time	156	hours

System Purchase Price: \$290K

Power Consumption: 11.7 kW

Total Operation Expense per yr: **\$18,308**

System Purchase Price: \$412K

Power Consumption: 15.9 kW

Total Operation Expense per yr: **\$24,866**

FIGURE 3 – TCO EXAMPLE

Or, [schedule a demonstration](#). We look forward to exploring with you how we can help increase the efficiency and agility of your data center, while keeping your costs down.



About GigaIO

Headquartered in Carlsbad, California, GigaIO democratizes AI and HPC architectures by delivering the elasticity of the cloud at a fraction of the TCO (Total Cost of Ownership). With its universal dynamic infrastructure fabric, FabreX™, and its innovative open architecture using industry-standard PCI Express/soon CXL technology, GigaIO breaks the constraints of the server box, liberating resources to shorten time to results. Data centers can scale up or scale out the performance of their systems, enabling their existing investment to flex as workloads and business change over time. For more information, contact info@gigaio.com or visit www.gigaio.com. Follow GigaIO on [Twitter](#) and [LinkedIn](#).

