



# Microsoft's Open CloudServer

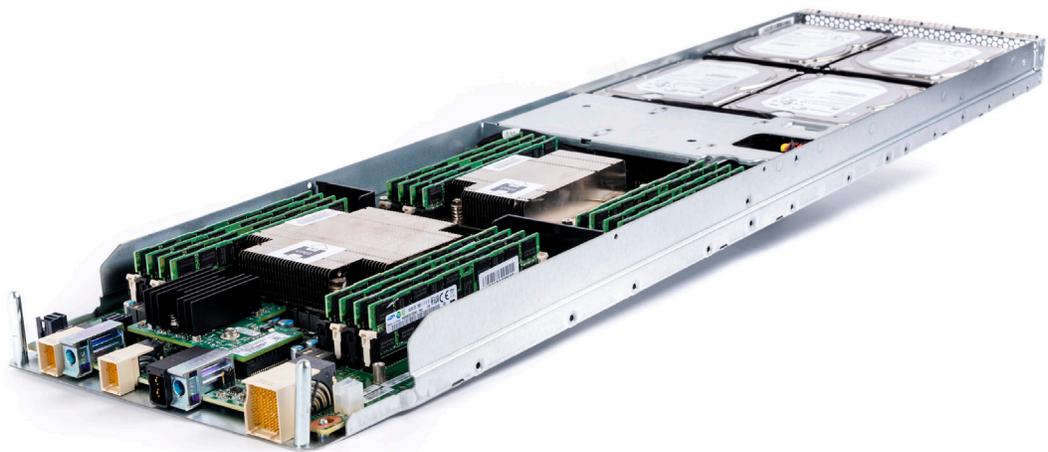


## How is our cloud infrastructure server design different from traditional IT servers?

It begins with scale. From the number of customers that need to be serviced, to the sheer number of servers that must be configured, deployed, managed, and secured, the scale required to support cloud services is several orders of magnitude beyond a traditional IT deployments. Cost effective designs are an obvious requirement, but a key difference is the nature of the applications that operate in either environment.

Today, Microsoft hosts more than 200 cloud services in our global cloud infrastructure of more than 100 datacenters. At this massive, hyper-scale, applications are engineered to provide the redundancy and resiliency needed to ensure these services are available at all times. The implication for hardware design is significant and provides opportunities to reduce physical redundancy, remove unnecessary components, and simplify the operations model.

This strategy brief describes Microsoft's approach to engineering servers and associated hardware for hyper-scale cloud operations.



## From reliable infrastructure to resilient services

In a traditional IT environment, hardware is typically designed for high reliability, requiring several layers of redundancy to ensure the infrastructure is always available – typically 99.999 percent uptime or more. Applications require that the hardware on which they are running is persistently available, may have dedicated redundant servers, and typically lack the capability to shift their workloads to alternate servers in the event of a failure. Failure of any component in the stack results in application downtime. Consequently, each layer of the stack requires multiple levels of redundancy—back-up power feeds and supplies, batteries and generators, back-up storage, back-up cooling, and back-up network connections.

In a cloud infrastructure environment, the sheer scale of operations dictates that at any given time numerous components will be in a failed state. To keep a service up and running in this scenario, the software needs to provide the resiliency. Properly architected applications can instantly shift their workload to alternate hardware – even a different datacenter – in the event of a component failure or configuration error. Hardware availability of 99.9 percent or less is acceptable.

Microsoft has moved to a model where our cloud services are designed to be software resilient, enabling hardware designs that do not require high levels of redundancy and availability, and offering significant efficiencies and cost savings. In addition, a component failure does not require immediate triage; the servicing model can move from a 24 x 7 schedule to an 8:00 am – 5:00 pm, five days a week schedule, thus reducing the cost of support and maintenance operations.

### MTBF vs. MTTR: Implications for server designs

Two critical attributes for any infrastructure design are Mean Time Between Failures (MTBF) and Mean Time To Recovery (MTTR). The combination of these two parameters determines the availability of hardware and the resiliency that must be built into the software.

Microsoft has a long history of providing cloud services at massive scale and we have come to recognize that to expand economically and simplify operations, we need to focus less on the hardware MTBF and instead focus more on the cloud services MTTR. As a result, hardware availability can be compromised from a typical 99.999 percent that is expected in an enterprise IT environment, to availability closer to 99.9 percent. Hardware availability needs only be “good enough,” since the increased software resiliency is providing the mechanism to provide low MTTR.

### Server design options

In server designs optimized for high MTBF, significant levels of redundancy are required and the fault domain will typically be the single server. That forces costly decisions through the design process, since the hardware is the key contributor to service availability.

In an environment optimized for MTTR, server designs follow a different philosophy. Software is intimately familiar with different hardware failures and fault domains such as networking and datacenter power, and the effect these have on the clusters. The software determines the deployment stamp for the application and determines how to keep the service up and running when the inevitable hardware failures occur.

Server stock-keeping unit (SKU) standardizations are also a key attribute that can provide cost advantages through the supply chain and simplify systems management and maintenance. A standardized environment also provides more opportunities for automation, which further reduces the possibility of configuration errors.

## MTBF vs. MTTR: Implications for server designs

Two critical attributes for any infrastructure design are Mean Time Between Failures (MTBF) and Mean Time To Recovery (MTTR). The combination of these two parameters determines the availability of hardware and the resiliency that must be built into the software.

Microsoft has a long history of providing cloud services at massive scale and we have come to recognize that to expand economically and simplify operations, we need to focus less on the hardware MTBF and instead focus more on the cloud services MTTR. As a result, hardware availability can be compromised from a typical 99.999 percent that is expected in an enterprise IT environment, to availability closer to 99.9 percent. Hardware availability needs only be “good enough,” since the increased software resiliency is providing the mechanism to provide low MTTR.

### Server design options

In server designs optimized for high MTBF, significant levels of redundancy are required and the fault domain will typically be the single server. That forces costly decisions through the design process, since the hardware is the key contributor to service availability.

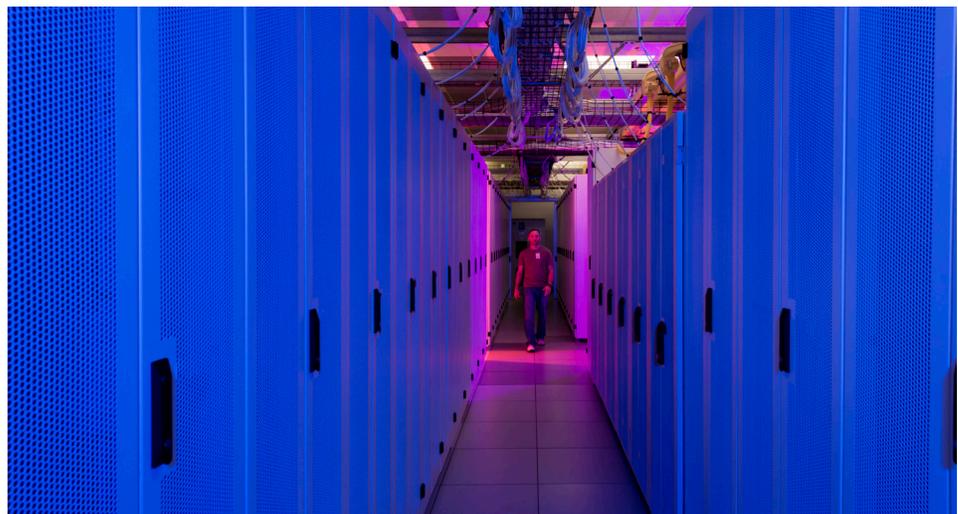
In an environment optimized for MTTR, server designs follow a different philosophy. Software is intimately familiar with different hardware failures and fault domains such as networking and datacenter power, and the effect these have on the clusters. The software determines the deployment stamp for the application and determines how to keep the service up and running when the inevitable hardware failures occur.

Server stock-keeping unit (SKU) standardizations are also a key attribute that can provide cost advantages through the supply chain and simplify systems management and maintenance. A standardized environment also provides more opportunities for automation, which further reduces the possibility of configuration errors.

### Operational environment for cloud servers

With an understanding of MTBF and MTTR and how they impact server design decisions, we can now look broader at the operating environment. A key efficiency driver is total vertical integration – from the processors and memory, solid state drives, and network components to the cloud operating system and management fabric. By standardizing and integrating across the stack, we are able to significantly reduce acquisition cost and lower operating and maintenance expense.

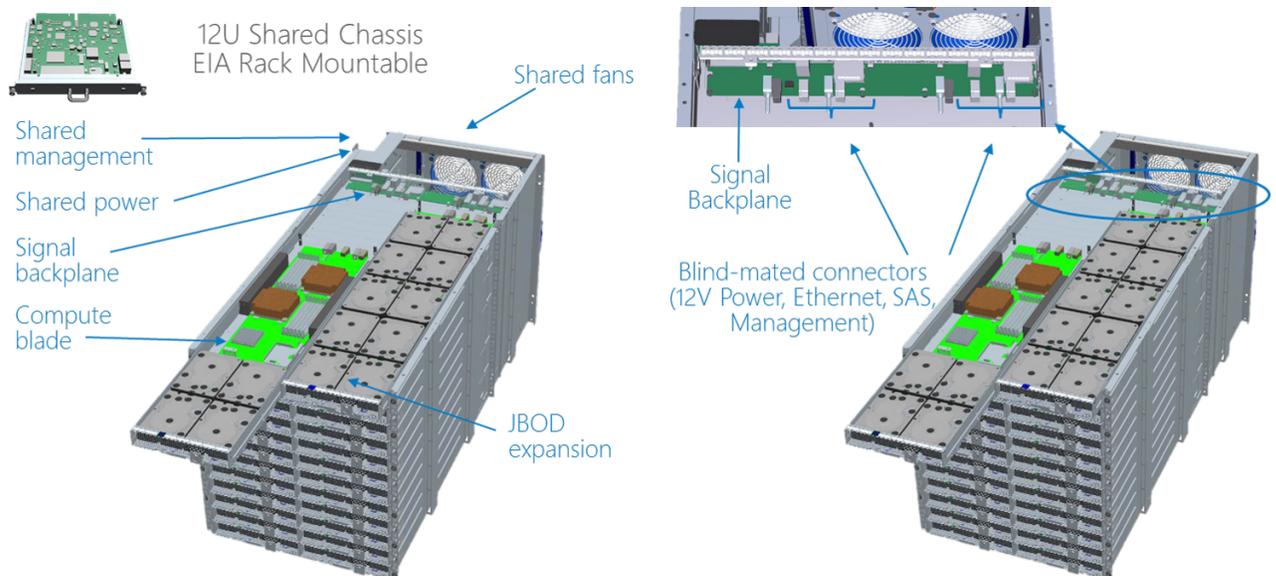
One variable is the application, and each cloud service has specific needs – whether it is specific server capability, memory needs, or network connectivity. Using a common set of building blocks, we can configure the servers to meet those specific needs. In some cases, the service requires higher hardware availability of 99.99 percent or in other cases, a moderate level of availability is acceptable. When a service requires high hardware availability, we can host it in a colocation room in our datacenters. Services engineered for high resiliency don't require the same level of hardware availability, so we host these in our modular datacenters that have relaxed environmental controls and reduced redundancy, saving significant cost and reducing our environmental impact.



## Microsoft's CloudServer

Microsoft has been deploying online services for over two decades and we brought our experience and learnings to develop the Open CloudServer V1, which we shared with the Open Compute Foundation and industry in early 2014. This new architecture for hardware and software converged all of our key cloud services, such as Bing, Office 365, and the Microsoft Azure platform on a common platform framework. Reducing total cost of ownership (TCO) was a key aspect. We strived to keep the acquisition costs low and reduce operational expenses, and we took a holistic look at the full server lifecycle— from architecture design through eventual decommissioning.

The result of this effort was a fully integrated design from the silicon, to the rack, and all the way to the datacenter level. It incorporates the blade, storage, network, systems management, and power mechanicals. And it comes together in a highly efficient single modular design. This cloud server design has been optimized for managing and operating an installed base of more than one million servers across a global footprint of Microsoft's datacenters.



### Performance results

The servers built against this design are currently in production in Microsoft datacenters and are yielding significant advantages over the traditional enterprise servers they replace:

- Up to **40% cost savings** and **15% power efficiency** benefits vs. traditional enterprise servers
- Up to **50% improvement** in deployment and service times
- Up to **75% improvement** in operational agility vs. traditional enterprise servers
- Is expected to save **10,000 tons of metal** and **1,100 miles of cable** for a deployment of one million servers

## Open CloudServer v2

In a sign of commitment to the Open Compute project, we released the V2 specifications in October 2014. As with our original design specification, we donated the hardware specifications for the chassis, the PCBA board, gerber files, and the mechanical CAD models for the metal chassis. In addition, we open sourced the management software and the tools we use during deployments, service, and repair. Vendors will be able to build blades that interoperate within the chassis and users will be able to deploy, manage, and service complete cloud systems.

The heart of the Open CloudServer V2 upgrade is a new compute blade. This blade supports the latest Intel processor, enabling 28 cores per blade. More cores enable more virtual machines, leading to a need to rebalance the system. The new design added capacity increases in all of the other subsystems on the blade.

The V2 blade design supports the transition from 10G networking to 40G networking. The new design utilizes RDMA over Converged Ethernet, or ROCE v2, to improve network efficiencies moving data to and from storage. It increases the SSD flash capacity by 4 times by transitioning away from SATA-based SSD flash storage to the PCI-Express based M.2 form factor. The M.2 flash cards are used in most laptops and tablets today, and are transitioning to the high performance NVMe interface. The M.2's small "stick" form factors enables thermal efficiency improvements, resulting in lower power consumption for cooling fans.

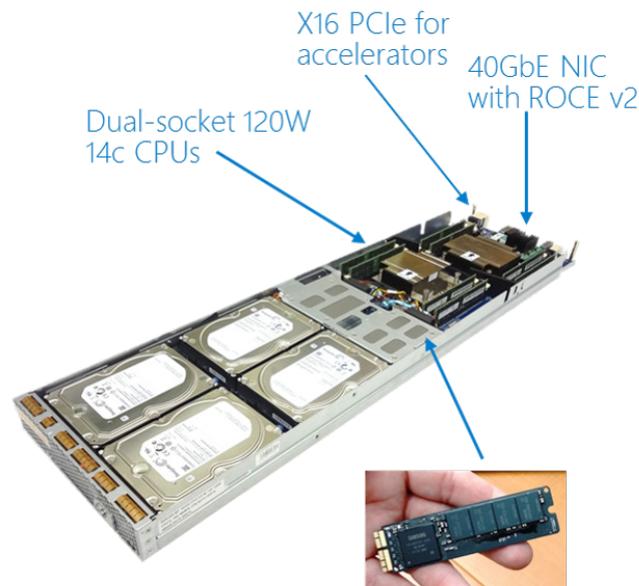
To improve the efficiencies of the processors, the design enables an expansion card that can accommodate items such as GPU and FPGA accelerator cards.

In the chassis, the capacity of the power supplies was increased and optimized for our hyper-scale cloud datacenter configurations. Utilizing a standard 19" EIA rack was very important in creating a flexible platform. To increase overall efficiency, the management, cooling, and power is pooled and shared across 24 blades. The blades can be any combination of compute or storage to meet flexible requirements.

Management of the blades and chassis is through the Chassis Manager, essentially an x86 PC. The Chassis Manager is responsible for setting fan speeds, monitoring the health of power supplies and fans, gathering event logs, and for monitoring the minimum set of out of band features required for running the servers.

There are six shared power supplies that allow for balanced three-phase power at the datacenter, allowing full utilization. Six large fans are used to reduce power consumption and enable fault redundancy. The assembly is designed for low-cost manufacturing and pre-rack assembly before arriving at the datacenter.

The overall result is that a full 52U rack with 96 servers can support 2688 cores, 48TB memory, 2.3 PB disk storage, 3/4ths PB flash storage, and 3.8 Tbps network bandwidth.



*Microsoft Open CloudServer V2 enhancements*

M.2 Flash solid-state storage

## Microsoft Open CloudServer v2 in datacenter deployments

The overall system architecture delivers the following advantages in hyper-scale cloud deployments:

**Simplicity** – at the scale of one million-plus deployed servers, simplicity of design is essential, as even the smallest issue can get magnified and potentially cause unexpected downtime and SLA violations for the services running on the infrastructure.

**Modularity** – the hardware system provides standardized interoperability at the blade, chassis, and rack level by modularizing the interfaces between these components. The hardware system utilizes blind-mate interface for blade installation and removal, enabling a cable-free approach for enabling a plug-and-play infrastructure—and reduces a common source of server failures.

**Efficiency** – shared power and cooling at the chassis level provides efficiency for both power consumption and datacenter power provisioning. The overall design also utilizes minimal materials for cost and mass reduction. Wherever possible, the design leverages existing commodity industry standard components and refrains from creating custom designs unless there is a tangible TCO benefit.

**Operational agility** – the hardware design enables streamlined supply chain operations by enabling pre-assembly of the modular components, and validation of the deployment unit during the factory

integration process. This minimizes deployment errors such as mismatched wiring, incorrect bill of materials (BoM), and incorrect software configurations. The result is faster deployment and a more rapid turnover to operations, enabling efficient use of procured assets and capacity provisioning. The service model during production operations is low touch given the cable-free system design. This reduces human errors relating to blade servicing and maintenance, improving overall TCO.

**Security at hyper-scale** – managing servers, storage, and networking devices requires mechanisms to access these endpoints in a secure manner for command execution. The hardware design provides a dedicated chassis manager which has multiple levels of security built in – a Trusted Platform Module on every compute node, secure BIOS/UEFI boot, authenticated roles (user/admin) and SSL certificate-based REST API command protocol. The combination of these security mechanisms ensures that when operating at scale, the datacenter assets are easily manageable without any security compromises that can impact operations or data confidentiality.

**Safety and compliance** – at hyper-scale, it is imperative that Microsoft's systems are safe to be used in any environment worldwide. The Open CloudServer system meets all safety and electromagnetic certifications to be operated in any datacenter around the globe.

## Sharing with the industry

Microsoft has contributed the V2 Open CloudServer specifications to the Open Compute Project in an effort to share its extensive research, development, and experience in operating hyper-scale cloud datacenters for the benefit of the broader hardware community.

Included in the contributions to the Open Compute Project are:

- **Hardware and software specifications**
  - Server design, mezzanine card, tray, chassis, and management card
  - Management APIs and protocols (for chassis and server)
- **Mechanical CAD models**
  - Chassis, server, chassis manager, and mezzanines
- **Gerber files**
  - Chassis manager card, power distribution board, and tray backplane
- **Source code for Chassis infrastructure**

Server management, fan and power supply control, diagnostics and repair

With more than 200 cloud services meeting the needs of more than one billion customers in 90 global marketplaces, Microsoft is one of only three companies operating a hyper-scale cloud infrastructure. By openly sharing our best practices, server specifications, and other intellectual property with the industry, we believe we can help maximize innovation and reduce operational complexity in the scalable computing space for everyone.

Microsoft has extensive experience operating a cloud services infrastructure since 1995. As Microsoft's cloud services portfolio and infrastructure continues to grow we are making thoughtful investments to answer customer needs for greater availability, improved performance, increased security, and lower costs.

Contributors:

**Monica Drake**

**Mark Shaw**

**Kushagra Vaid**

**For more information, please visit [www.microsoft.com/datacenters](http://www.microsoft.com/datacenters)**

© 2015 Microsoft Corporation. All rights reserved.

This document is for informational purposes only. MICROSOFT MAKES NO WARRANTIES, EXPRESS OR IMPLIED, IN THIS SUMMARY.