



How Microsoft Designs its Cloud-Scale Servers



How Microsoft Designs its Cloud-Scale Servers

How is cloud infrastructure server hardware 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 deployment. Cost effective designs are an obvious requirement, but a key difference is the nature of the applications that operate in either environment.

At cloud-scale, applications are engineered to provide the redundancy and resiliency to ensure the 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 its cloud-scale 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 is moving to a model where cloud services are designed to be resilient, enabling hardware designs that do not require high levels of redundancy and availability, and offering significant cost savings. In addition, a component failure does not require immediate triage; the servicing model can move from 24 x 7 to 8:00 am – 5:00 pm, five days a week, 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 scale. We have learned that in order 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 software 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. The MTTR design has a software-based fault domain, meaning that the deployment stamp is determined by the application and that software figures out 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.

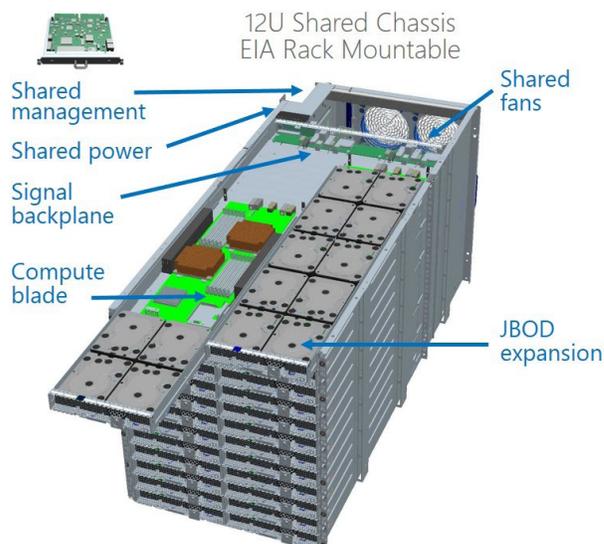
Operational environment for cloud-scale servers

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 high hardware availability of 99.999 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 costs and reducing our environmental impact.

Microsoft’s cloud server architecture

Microsoft has been deploying cloud services since 1995 and we are bringing our experience and learnings to develop a next generation cloud server. The outcome is a new architecture for hardware and software, converging all of our key online services such as Bing, Office 365,



and the Windows Azure platform on a common server framework. Total cost of ownership (TCO) is a key aspect. We strive to keep the acquisition costs low and reduce operational expenses, and we are taking a holistic look at the full server lifecycle – from architecture design through eventual decommissioning.

The result of this effort is 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 greater than 1 million servers across a global footprint of datacenters.

Microsoft's cloud server architecture is based on a modular high-density chassis approach that enables efficient sharing of resources across multiple server nodes. A single 12U chassis can accommodate up to 24 server blades (either compute or storage), where two blades are populated in each 1U slot. A rack can hold three or four chassis depending on the rack height.

The chassis provides common functions to the blades – highly efficient shared power and cooling, integrated and scalable management, and wiring connectivity to the datacenter network spine. By moving these common functions into the chassis, the cost can be reduced for optimal TCO and this approach allows for operating at scale where a typical single deployment might involve tens of thousands of server at a time.

The overall system architecture was designed around the following key principles:

Simplicity – at the scale of 1 million 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 a novel mechanism for blade installation and removal, utilizing a cable-free approach for enabling a plug-and-play infrastructure. This approach allows us to use the same set of core building blocks to easily create different configuration topologies for specific cloud applications and different global datacenter environments.

The result is streamlined manufacturing and assembly, the re-use of common components across deployments for volume economics, simplified maintenance operations during production usage, the ability to target technology refreshes independently across compute, storage, and networking components, and the ability to reuse infrastructure across decommissioning cycles.

The design also supports high density deployments, with up to 96 servers in a 52U rack. This allows better cost amortization for common elements across the server blades, and is a good option for space constrained environments.

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.

The hardware design places emphasis on cost efficiency by the removal of unnecessary features from the perspective of a large-scale cloud operations environment (MTTR based approach). The Field Replaceable Unit (FRU) is the entire blade for both compute and storage, and this design choice reduces the cost structure by minimizing the components and sheet metal required to meet the desired application requirements.

Microsoft cloud server specification *Compute blade*

Processor	
CPU	Dual Intel® Xeon® E5-2400 v2 family
Core QTY	Up to 10 cores / CPU, 20 / Blade
TDP Wattage	Up to 95W
Memory	
Memory Busses and DIMM Slots	3X memory bus / CPU, 6 / Blade
DIMM Type / Speed	6 DIMM slots / CPU, 12 / Blade
	16GB, 2Rx4, 1333MHz, 1.35V
Max Capacity	192 GB / Blade
On-Board Devices	
Storage Controller	Intel® C602 PCH
SATA	4 ports @ 3.0 Gb/s
SATA	2 ports @ 6.0 Gb/s
Server Management	
Chipset	BMC-lite serial thru Chassis Manager
Interface	REST API, CLI thru Chassis Manager
System Firmware	
Version, Vendor	UEFI 2.3.1, AMI
Security	TPM 1.2, Secure Boot
Blade I/O	
PCI-Express Slots	One Gen3 X16 Riser
Networking	Single or Dual 10G Mezzanine Card
SAS	Dual 4X SAS @ 6G Mezzanine Card

Microsoft cloud server specification *JBOD blade*

Storage	
HDDs	LFF 7200 RPM 3.5" SATA
Capacity	20 TB, 10 x 2TB HDD
	30 TB, 10 x 3TB HDD
	40 TB, 10 x 4TB HDD
	50 TB, 10 x 5TB HDD
	60 TB, 10 x 6TB HDD
Connectivity	
Signals	8 Channel SAS @ 6G
Expander	LSI 20-port LSI SAS2X20
Storage Server Configuration Capacities	
Storage Controller : Single Compute with 4 HDDs plus	
+ 1 JBOD	14 HDDs
+ 2 JBOD	24 HDDs
+ 4 JBOD	44 HDDs
+ 6 JBOD	64 HDDs
+ 8 JBOD	84 HDDs

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, 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 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 – secure BIOS 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.

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
- Expected to save 10,000 tons of metal and 1,100 miles of cable per one million servers

Sharing with the industry

Microsoft is contributing these server specifications to the Open Compute Project Foundation in an effort to share its extensive research, development, and experience in operating large cloud-scale datacenters for the benefit of the broader hardware and datacenter community.

Included in the contributions to the Open Compute Project Foundation are:

1. **Hardware and software specifications**
 - Server design, mezzanine card, tray, chassis, and management card
 - Management APIs and protocols (for chassis and server)
2. **Mechanical CAD models**
 - Chassis, server, chassis manager, and mezzanines
3. **Schematics and Gerber files**
 - Chassis manager card, power distribution board, and tray backplane
4. **Source code for Chassis infrastructure**
 - Server management, Fan and Power supply control, Diagnostics and Repair

We believe that openly sharing ideas, specifications and other intellectual property is the key to maximizing innovation and reducing operational complexity in the scalable computing space. The Open Compute Project Foundation provides a structure in which individuals and organizations can contribute and share their intellectual property with members.

Microsoft has extensive experience operating a cloud services' infrastructure since 1995, with a history of innovation, operational excellence and industry leadership. As Microsoft's cloud services portfolio and infrastructure continues to grow, and with new services and applications launching on a rapid basis, we are making thoughtful investments to answer our customer's needs for greater availability, improved performance, increased security, and lower costs.

White paper authors: Kushagra Vaid, Cloud Infrastructure Server Engineering; Mark Shaw, Cloud Server Hardware Development; and Monica Drake, Cloud & Enterprise & Datacenter Product Marketing

For more information, please visit www.microsoft.com/datacenters

A DAY IN THE
MICROSOFT CLOUD 

© 2014 Microsoft Corporation. All rights reserved.

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