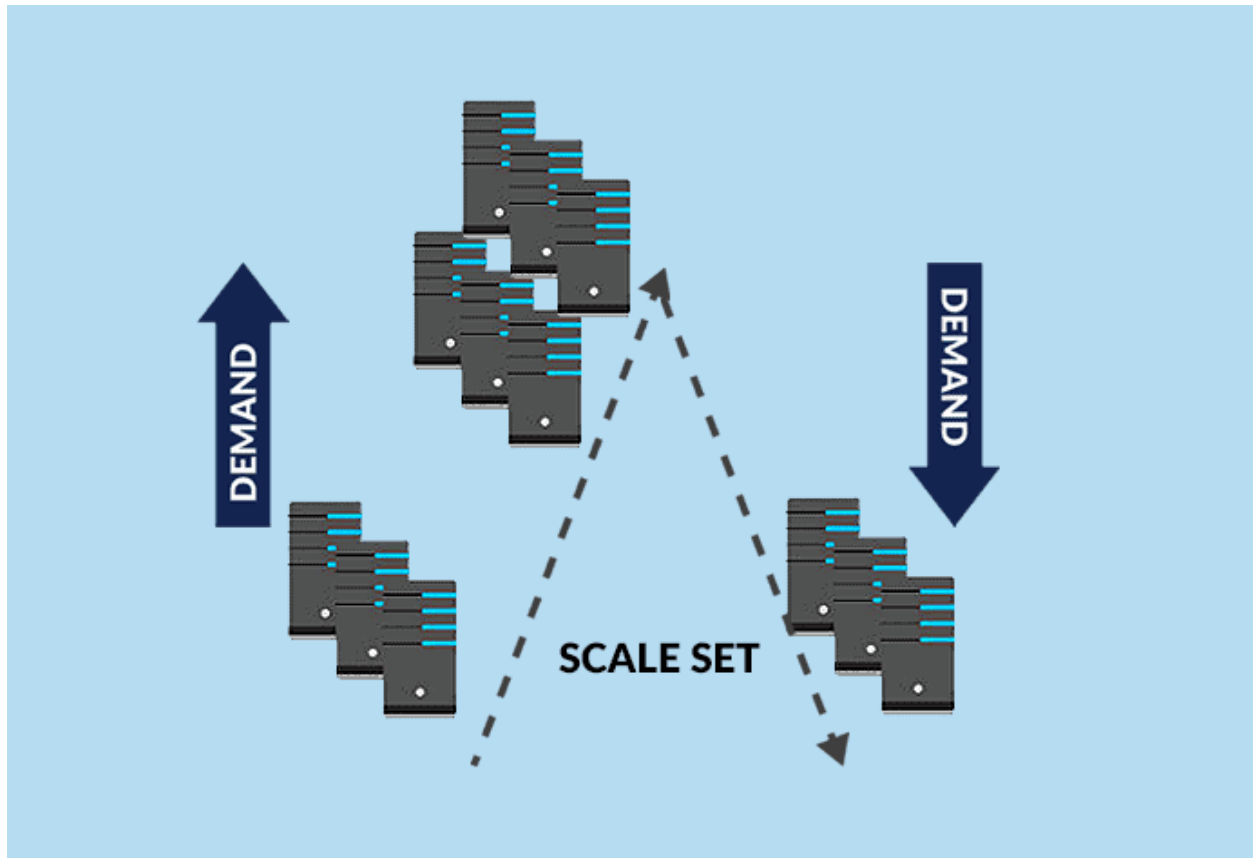


Virtual Machine Scale Set (VMSS)

What are Virtual Machine Scale Sets (VMSS)?

Virtual Machine Scale Sets (VMSS), an interesting service offered by Microsoft Azure, helps to create and manage a set of identical, auto-scaling Virtual Machines (VMs). The number of VM instances can automatically increase or decrease based on scheduled conditions.



VMSS is especially beneficial for applications with variable or unpredictable workloads since it can automatically alter the number of VM instances based on demand. This helps to ensure that your application stays available and responsive to users even during periods of high traffic or increased activity.

You can use VMSS to deploy a group of VMs with identical configurations, such as the operating system, programs, and data. VMSS also provides load balancing, automatic scaling, and connection with other Azure services, making it a powerful and versatile solution for managing large-scale workloads.

Overall, VMSS is an important Microsoft Azure feature that provides a scalable and dependable method for deploying and managing a set of identical VMs.

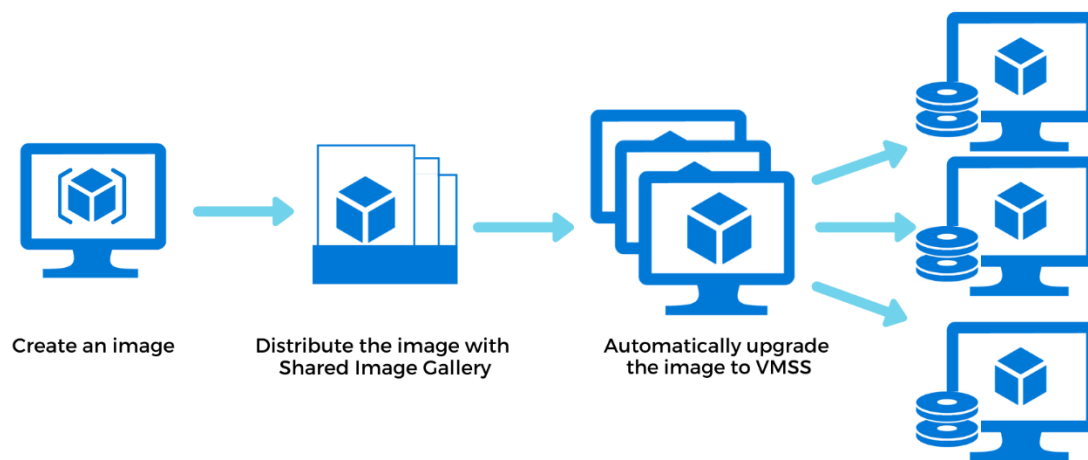
Auto-Scaling

Auto-scaling is a way to automatically scale up or down the number of compute resources that are being allocated to your application based on its needs at any given time.

The key point is that you can now design a scalable architecture that will automatically scale up or scale down to meet your needs over the lifetime of your setup regardless of how fast/slow or big/small your site grows over that time.

Here are the most popular ways of auto-scaling:

- **Horizontal Scaling**
- **Vertical Scaling**



Vertical Scaling:

Vertical Scaling is an attempt to increase or decrease the capacity of a single machine, also called scaling up or down. Here the resources such as processing power, storage, memory, and more are added to an existing work unit.

It is done to increase the capacity of existing hardware or software by adding resources. It can enhance your server without manipulating your code. But it is limited by the fact that you can only get as big as the size of the server.

- **Scaling up** refers to expanding the computational power or capacity of a single virtual machine (VM) instance within the VMSS. Scaling up often entails adding more resources to a single VM, such as increasing the number of CPU cores or RAM, in order to improve performance or accommodate rising workload needs. This can be done manually or automatically depending on predefined rules and metrics.

- **Scale down** is the process of reducing the number of virtual machine (VM) instances in a Virtual Machine Scale Set (VMSS) to match the current demand of the application. Scale down is a significant feature of VMSS since it helps to optimize resource utilization and minimize expenses by deleting unnecessary VM instances. When scaling down in VMSS, VM instances are removed based on the scale-in rules configured for the VMSS. These rules indicate the requirements that must be followed in order to delete VM instances, such as a decrease in CPU consumption, memory usage, or network traffic. When the conditions are met, the VM instances are removed from the VMSS.

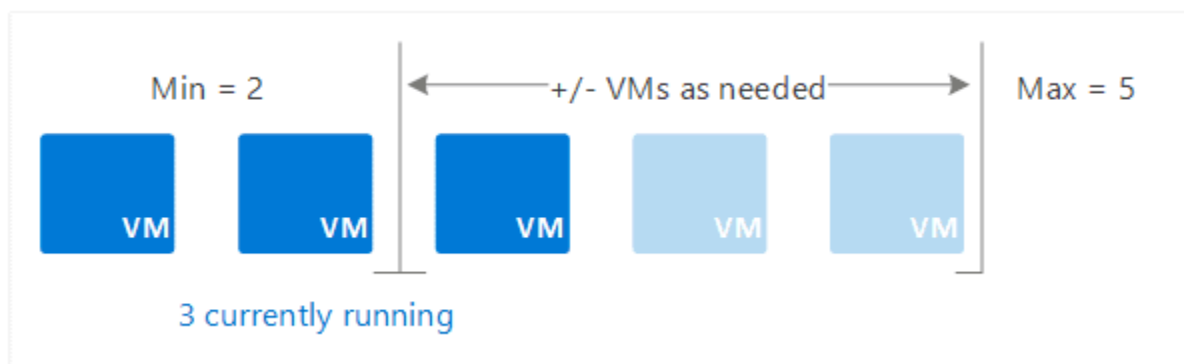
Horizontal Scaling:

Horizontal Scaling is a must-use technology – whenever a high availability of (server) services is required, also called scaling out and in, means adding or removing instances of a resource. Scaling horizontally involves adding more processing units or physical machines to your server or database. It involves growing the number of nodes in the cluster, the application continues running without interruption as new resources are provisioned.

If demand drops, the additional resources can be shut down cleanly and deallocated.

- **Scale-out** is the process of adding extra resources, such as virtual machines (VMs), to a Microsoft Azure Virtual Machine Scale Set (VMSS) to manage an increase in demand for your application. Scaling out a VMSS increases the number of operating VM instances, letting your application to handle additional requests and workload. Scaling out in VMSS is accomplished by adding new instances of the same VM to the set. These instances are provisioned and configured in the same way as the current instances in the set, ensuring that they work seamlessly together.
- **Scale-in** refers to the process of reducing the number of virtual machines (VMs) in the scale set. When a scale-in event occurs, one or more VM instances are removed from the scale set, limiting the system's total capacity. Scale-in is often triggered when resource demand falls and the system no longer requires the same level of capacity. Scaling-in helps to optimize resource utilization and decrease expenses by deleting superfluous VM instances, as the system only uses the resources needed to satisfy current demand.

In Autoscaling you have to specify a minimum and the maximum number of instances to run and add or remove VMs automatically based on a set of rules.



When rule conditions are met, one or more autoscale actions are triggered. You can add and remove VMs, or perform other actions.

Autoscale Settings

An **autoscale setting** is read by the autoscale engine to determine whether to scale up or down. You can create custom autoscaling rules as needed for your situation. Rule types include:

- **Minimum Instance:** The minimum number of instances you want to deploy in your scale set.
- **Maximum Instances:** The maximum number of instances you want to deploy while scaling out. (**Note:** In Azure, you can have a maximum, of 1000 instances)
- **Metric-based** – It measures application load and adds or removes VMs based on that load. For example, add an instance in a scale set when CPU usage is above 50%.
- **Time-based** – For example, trigger an instance every 8 am on Saturday in a given time zone.

Here are a **few points** which are important when we think about going with **Horizontal scaling** or **Vertical Scaling**.

- **Scaling up requires downtime**, in this case, you need to upgrade the server's configuration like RAM, memory, CPU, etc. so while upgrading this configuration your server requires downtime. Once you are done with the update, the restart of the server is required.
- **Scaling up will increase performance but is not available** because it's only one instance and it can go down anytime when it reaches the scaling rules.
- **Scale-Out doesn't require downtime**, in Scale-out it creates new instances of the server it doesn't touch the existing instance, so no downtime is required.
- **Scale-Out, Increase performance and availability as well.** When autoscaling increases no. of instances and handles load/request using a load balancer it increases the performance of the server and also increases availability as well.

Advantages of Virtual Machine Scale set

1. **Scalability:** VMSS allows you to scale your infrastructure up or down automatically according on the demands of your application, ensuring that you have the resources you need when you need them.
2. **High availability:** VMSS distributes virtual machine instances across various fault domains, ensuring that your application remains operational even if one or more virtual machines fail or become unavailable,
3. **Cost-effective:** By dynamically scaling your infrastructure based on demand, VMSS helps you to optimize your resource utilization and decrease your costs, avoiding the need to overprovision resources.
4. **Consistency:** VMSS enables you to deploy and manage a collection of VMs with similar configurations, assuring constant performance and reducing the chance of configuration errors.
5. **Simple management:** VMSS can be simply controlled via Azure Portal, Azure CLI, Azure PowerShell, and Azure Resource Manager templates, making it a flexible and user-friendly option for managing large-scale workloads.
6. **Integration:** VMSS works with other Azure services including Azure Load Balancer, Azure Application Gateway, and Azure Traffic Manager to provide a comprehensive solution for managing and delivering large-scale workloads.

Interview Questions on VMSS

1) What is a Virtual Machine Scale Set (VMSS) in Microsoft Azure?

A Virtual Machine Scale Set (VMSS) is a collection of similar virtual machines (VMs) in Microsoft Azure that may scale up or down automatically dependent on the demand of your application. VMSS offers an automatic scaling solution that may change the number of VM instances based on variables like CPU consumption, memory usage, or network traffic.

VMSS also has built-in load balancing technology that can balance traffic throughout the set's VM instances, ensuring that each instance is used efficiently and effectively. With VMSS, you can deploy a group of identical VM instances with the same configuration, ensuring consistency and simplifying management. VMSS is a robust and adaptable solution for installing and managing virtual machines (VMs) on Microsoft Azure, offering scalability, high availability, cost-effectiveness, load balancing, consistency, and integration with other Azure services.

2) What are the key benefits of using VMSS in Microsoft Azure?

The following are the main advantages of using Virtual Machine Scale Sets (VMSS) in Microsoft Azure:

1. Scalability: VMSS offers automatic scaling of VM instances based on the demand of your application.
2. High availability: VMSS provides fault-tolerant and redundant deployment options to assure your application's high availability.
3. Cost-effectiveness: VMSS allows for more efficient resource usage, which can result in cost savings.
4. Load balancing: VMSS has built-in load balancing capability to optimize traffic distribution among VM instances.
5. Consistency: VMSS enables the consistent deployment and control of identical VM instances with the same configuration.
6. Integration: VMSS works in tandem with other Azure services such as Azure Load Balancer and Azure Application Gateway.

How does VMSS handle automatic scaling and load balancing?

The following mechanisms are used by VMSS to handle automatic scaling and load balancing:

1. Automated scaling: VMSS employs automated scaling rules to adjust the number of VM instances in the set based on the application's demand. These rules are based on data like as CPU utilization, memory usage, or network traffic and can be tailored to the application's specific requirements.
2. Load balancing: VMSS includes load balancing capability that distributes traffic across the set's VM instances. This ensures that each instance is used efficiently and effectively and that the burden is divided equitably across the set.
3. Integration with other Azure services: VMSS connects with other Azure services including Azure Load Balancer, Azure Application Gateway, and Azure Traffic Manager to provide a complete solution for managing and delivering large-scale workloads.

Can I use VMSS with both Windows and Linux operating systems?

Yes, Virtual Machine Scale Sets (VMSS) can be used with both Windows and Linux operating systems. VMSS supports both operating systems and allows you to install and manage groups of identical VMs with the same configuration, regardless of operating system. This simplifies management and ensures consistency across all of your deployments, regardless of platform.

3) What is Azure Managed Disks, and how does it work with VMSS?

Azure Managed Disks is a storage service in Microsoft Azure that facilitates the management and scaling of VM disks. It provides an easy-to-use, scalable, and highly available storage solution for Azure VMs.

When used with Virtual Machine Scale Sets (VMSS), Azure Managed Disks allows you to effortlessly manage the storage of VM disks across a large number of VM instances. VMSS handles the creation and destruction of VM instances automatically, while Azure Managed Disks guarantees that the corresponding disk storage is also generated or removed automatically.

When used with VMSS, Azure Managed Disks provides various advantages, including easier disk management, greater disk reliability, faster disk provisioning, and support for huge disk sizes. Furthermore, you can easily configure your storage settings with Azure Managed Disks to meet your specific needs, such as storage performance and redundancy.

Overall, combining Azure Managed Disks and VMSS simplifies VM disk administration and scalability, resulting in a dependable, scalable, and highly available storage solution for your Azure workloads.

4) How do I create and configure a VMSS in Microsoft Azure?

To construct and configure a Virtual Machine Scale Set (VMSS) in Microsoft Azure, go through the following steps:

1. Create an image: Create an image of the VM you intend to use as the foundation for the VMSS.
2. Create a VMSS: In the Azure portal, create a VMSS and specify the image that you prepared in step 1 as well as other configuration settings such as the number of instances, virtual machine size, and network settings.
3. Configure autoscaling: Configure autoscaling parameters for the VMSS, providing the criteria and circumstances for scaling up or down based on the demand of your application.
4. Configure load balancing: Configure load balancing for the VMSS, specifying the load balancing algorithm and backend pool values.
5. Deploy your application: Deploy your application to the VMSS, either by setting the VM instances individually or by utilizing a custom image that incorporates your application.
6. Monitor and optimize: Monitor the VMSS's performance and usage, and optimize the autoscaling and load balancing parameters as needed to achieve optimal performance and cost-effectiveness.

Overall, constructing and setting a VMSS in Azure entails a few essential processes, including producing an image, creating the VMSS, enabling autoscaling and load balancing, deploying your application, and monitoring and optimizing the VMSS's performance.