

# Agenda

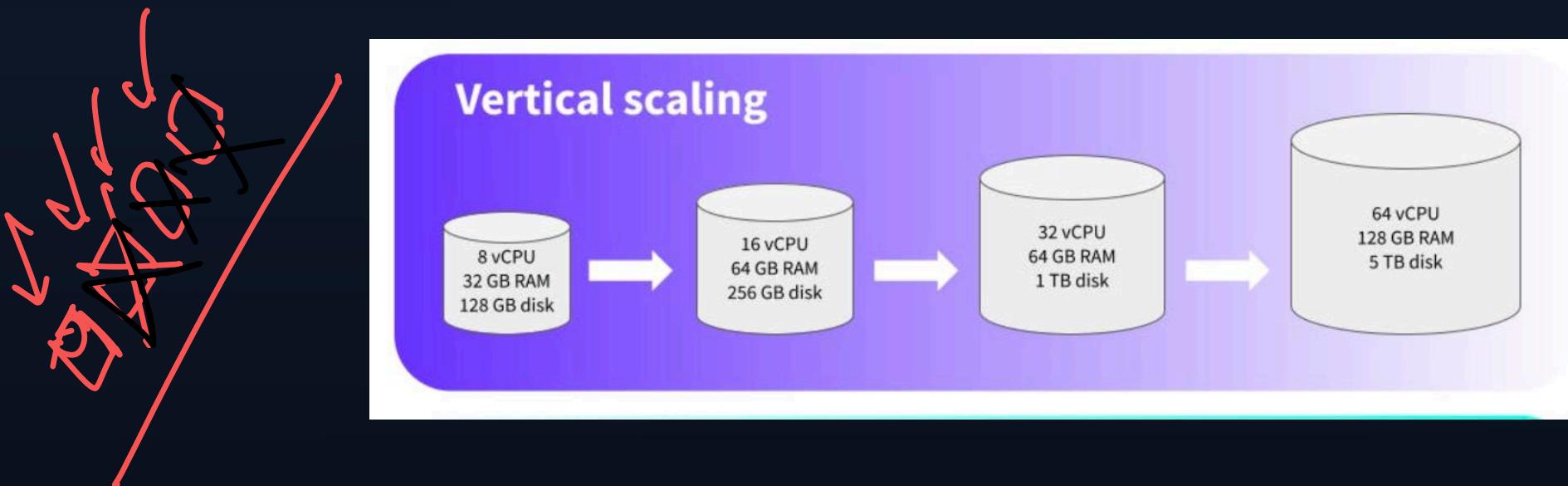
- ① Capacity management
- ② Scaling options
  - ↳ Horizontal/vertical scaling
- ③ Setup HS, VA  
↳ ASG → LT → scheduled Action → lifecycle policy
- ④ Types of LB - ALB, NLB, GLB, CB -

# Quick Recap of EC2 Instance, AMI and Backups

# 1. Vertical Scaling

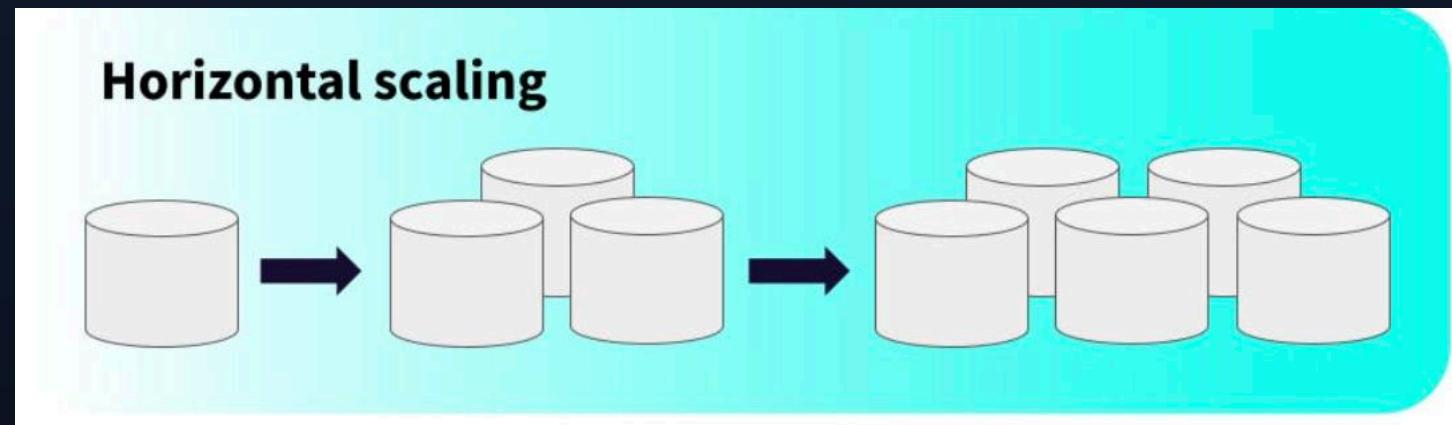
- *Vertical scaling* refers to increasing the capacity of a system by adding capability to the machines it is using (as opposed to increasing the overall number of machines). This is also called *scaling up*.
- Vertical scaling can be upgrading the physical machine that's running your system, or it can be swapping over to a different, more capable machine. As long as the *number* of machines our system uses isn't changing, that's scaling vertically.

*For example, imagine that we have an application with a cloud database that has reached the limits of the server it's running on: a single 8 vCPU GCP instance with 32 GB of RAM.*



## 2. Horizontal Scaling

- *Horizontal scaling* refers to increasing the capacity of a system by adding additional machines (nodes), as opposed to increasing the capability of the existing machines. This is also called *scaling out*.
- *For example, imagine that we again have an application with a cloud database that has reached the limits of the server it's running on: a single 8 vCPU GCP instance with 32 GB of RAM.*

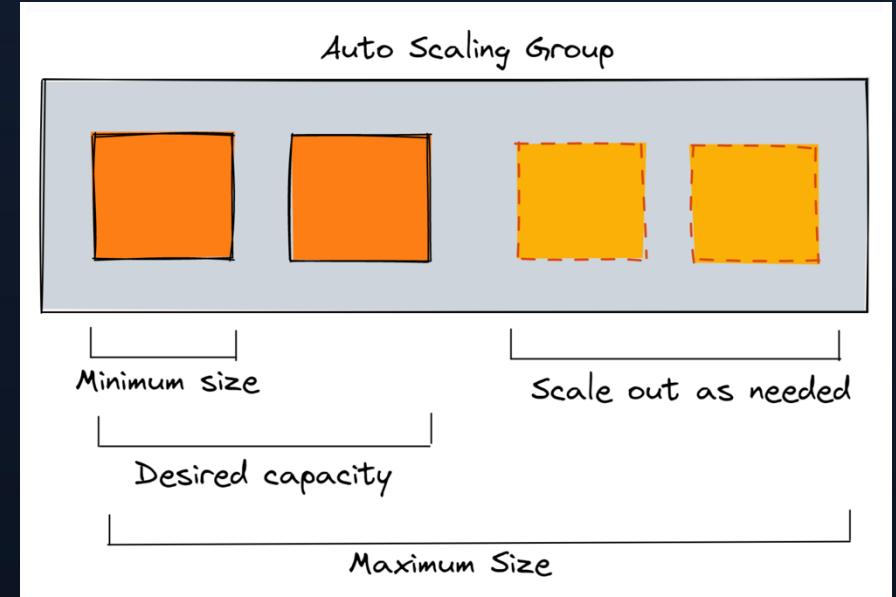


# Introduction to Auto Scaling

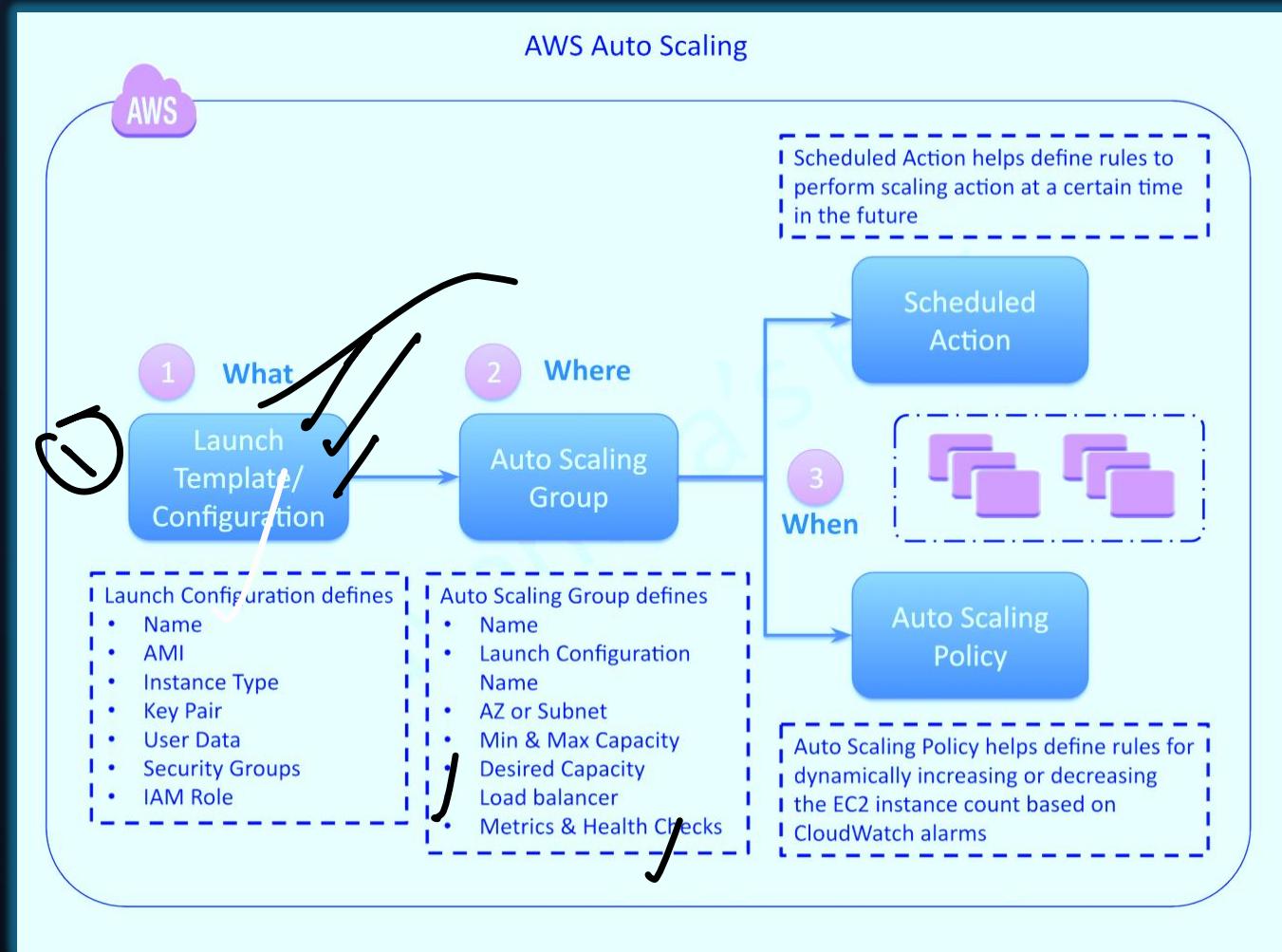


Auto-scaling is a feature of **cloud computing** that helps organizations scale the capacity of cloud services or virtual machines up or down based on traffic levels.

By automatically expanding and lowering fresh instances as demand rises and falls, auto-scaling reduces cost and enables consistent functionality. As a result, amidst changing and often unforeseen requests for services, auto-scaling ensures stability. Auto-scaling also avoids the need to react explicitly to heavy traffic in real-time which would necessitate more tools and instances. Furthermore, auto-scaling allows for the installation, tracking, and deactivation of each unit.



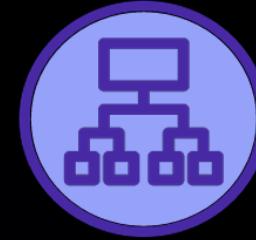
# ASG Components



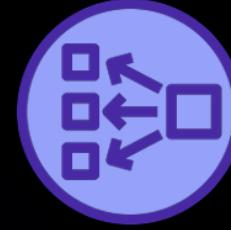
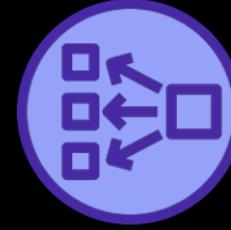
# Types of Load balancer

|                                                                                                                   |                                                                                                               |                                                                                                                  |                                                                                                                 |                                                                                                            |
|-------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------|
|  Application Load Balancer (ALB) |  Network Load Balancer (NLB) |  Gateway Load Balancer (GWLB) |  Classic Load Balancer (CLB) |  AWS Global Accelerator |
| Layer 7                                                                                                           | Layer 4                                                                                                       | Layer 3 gateway/<br>4 load balancer                                                                              | Layer 4/7                                                                                                       | TCP/UDP                                                                                                    |
| <b>Targets</b><br>IP, instances,<br>AWS Lambda,<br>containers                                                     | <b>Targets</b><br>IP, instances, ALB,<br>containers                                                           | <b>Targets</b><br>IP, instances                                                                                  | <b>Targets</b><br>EC2-Classic                                                                                   | <b>Targets</b><br>IP, ALB, NLB                                                                             |
| <b>Protocols</b><br>HTTP, HTTPS, gRPC                                                                             | <b>Protocols</b><br>TCP, UDP, TLS                                                                             | <b>Protocols</b><br>IP                                                                                           | <b>Protocols</b><br>TCP, SSL/TLS,<br>HTTP, HTTPS                                                                | <b>Protocols</b><br>TCP, UDP                                                                               |

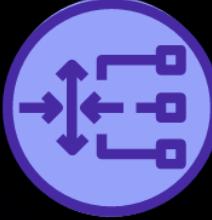
# Which load balancing technology should we use?

| Targets                                                                                      | Requires                                                                                                                                                               |                                                                                     |
|----------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------|
|  Instances  | Layer 7 routing<br>HTTP2/gRPC                                                                                                                                          |  |
|  AWS Lambda | Redirects, web sockets                                                                                                                                                 |                                                                                     |
|  Containers | Fixed response<br>Authentication                                                                                                                                       | <b>Application Load Balancer</b>                                                    |
|  IP        | Web application firewall, AWS Outposts/AWS Local Zones<br>Cookie stickiness, HTTP Desync mitigation<br>Best option for the AWS Load Balancer Controller for containers |                                                                                     |

# Which load balancing technology should we use?

| Targets                                                                                      | Requires                                                                                   |                                                                                                                     |
|----------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------|
|  Instances  | Low latency<br>Zonal isolation                                                             |                                  |
|  ALB        | Long-lived TCP connections                                                                 |                                                                                                                     |
|  Containers | Connection-based<br>Layer 4 load balancing<br>PrivateLink support                          | <br><b>Network Load Balancer</b> |
|  IP        | Elastic IP support<br>Hybrid architecture support<br>AWS Fargate support direct to K8s pod |                                                                                                                     |

# Which load balancing technology should we use?

| Targets                            | Requires                                                                                                            |                                                                                                                     |
|------------------------------------|---------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------|
| <input type="checkbox"/> Instances | Bump in the wire<br>Auto scaling for packet processing devices (firewall, IdP)                                      | <br><b>Gateway Load Balancer</b> |
| ○→ IP                              | Packet preservation for inspection<br>PrivateLink GWLB endpoint<br>Multi-port to same instance<br>Route table entry |                                                                                                                     |

# Which load balancing technology should we use?

| Targets                                                                               | Requires                                                                              |                                                                                     |
|---------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------|
|  NLB | Accelerate latency-sensitive applications                                             |  |
|  ALB | Improve resiliency and availability on a global scale                                 |  |
|  IP  | Simplified global traffic management<br><br>Global set of anycast static IP addresses |  |