



1. Explain Autoscaling concept along with Task Scaling and EC2 scaling in ECS ?

5 Simple Steps to Set It Up:

1. Turn on Auto Scaling for your ECS Service.
2. Set a rule (e.g., "Keep CPU usage around 60%").
3. Pick a metric to watch, like CPU or custom CloudWatch metric.
4. Attach a scaling policy (target tracking or step scaling).
5. ECS will now automatically add or remove tasks based on the load!

The Core Concept: Load Distribution

Think of it like a restaurant kitchen:

- 1 chef handling 100 orders = very stressed chef (high CPU)
- 4 chefs handling 100 orders = 25 orders per chef = relaxed chefs (low CPU)

ECS Auto Scaling works the same way - it adds more "chefs" (tasks) to reduce the workload per "chef" (task).

E-commerce Website

Initial Setup

- **Service:** Online shopping website
- **Current Tasks:** 2 containers
- **Target CPU:** 60% (our comfort zone)
- **Min Tasks:** 2, **Max Tasks:** 8

Scenario 1: Black Friday Traffic Spike

Before Scaling:

Task 1: Processing 500 requests/min → 85% CPU

Task 2: Processing 500 requests/min → 85% CPU

Average CPU: 85% (ABOVE 60% target)

ECS Auto Scaling Decision:

- "CPU is too high! Need more tasks to share the load"
- **Action:** Scale out from 2 → 4 tasks

After Scaling:

Task 1: Processing 250 /min → 45% CPU

Task 2: Processing 250 requests/min → 45% CPU

Task 3: Processing 250 requests/min → 45% CPU

Task 4: Processing 250 requests requests/min → 45% CPU

Average CPU: 45% (BELOW 60% target ✅)

Key Insight: Same 1000 requests/min, but now distributed across 4 tasks instead of 2!

Scenario 2: Late Night Low Traffic 🌙

Before Scaling:

Task 1: Processing 75 requests/min → 25% CPU

Task 2: Processing 75 requests/min → 25% CPU

Task 3: Processing 75 requests/min → 25% CPU

Task 4: Processing 75 requests/min → 25% CPU

Average CPU: 25% (BELOW 40% target)

ECS Auto Scaling Decision:

- "CPU is too low! We have too many tasks for this workload"
- **Action:** Scale in from 4 → 2 tasks (but wait 5 minutes to be sure)

After Scaling:

Task 1: Processing 150 requests/min → 50% CPU

Task 2: Processing 150 requests/min → 50% CPU

Average CPU: 50% (BELOW 60% target ✅)

🔧 Step-by-Step Implementation

For Fargate (Recommended for Beginners)

yaml

Example Configuration

Service: my-web-app

Launch Type: Fargate

Min Tasks: 2

Max Tasks: 10

Target CPU: 60%

Scale-out cooldown: 60 seconds

Scale-in cooldown: 300 seconds

Console Steps:

1. ECS Console → Select your service
2. Auto Scaling tab → "Configure Service Auto Scaling"
3. Set min/max task count
4. Choose "Target Tracking" scaling policy
5. Select "CPUUtilization" metric
6. Set target value: 60

For EC2 Launch Type (Advanced)

Additional Complexity: You need TWO layers of scaling:

1. **ECS Service Auto Scaling** (scales tasks)
2. **EC2 Auto Scaling Group** (scales EC2 instances)

Problem Scenario:

Current: 2 EC2 instances, each running 2 tasks

High traffic hits → ECS wants to scale to 6 tasks

But: Only capacity for 4 tasks (2 per instance)

Result: 2 tasks remain "PENDING" (can't be placed)

Solution: EC2 ASG detects pending tasks and launches more EC2 instances.

Complete Scaling Timeline Example

Starting Point: 2 tasks, normal traffic

Time: 9:00 AM | Tasks: 2 | CPU: 40% | Status: Stable

Traffic Spike Begins:

Time: 9:05 AM | Tasks: 2 | CPU: 75% | Status: Above target!

Time: 9:06 AM | Tasks: 4 | CPU: 45% | Status: Scaled out ✅

Even More Traffic:

Time: 9:15 AM | Tasks: 4 | CPU: 80% | Status: Above target again!

Time: 9:16 AM | Tasks: 6 | CPU: 55% | Status: Scaled out again ✅

Traffic Normalizes:

Time: 11:00 AM | Tasks: 6 | CPU: 30% | Status: Below target

Time: 11:05 AM | Tasks: 6 | CPU: 30% | Status: Waiting (cooldown period)

Time: 11:10 AM | Tasks: 4 | CPU: 45% | Status: Scaled in ✅

Time: 11:15 AM | Tasks: 4 | CPU: 45% | Status: Stable

⚠️ Common Questions & Answers

Q: "How does adding tasks reduce CPU if the same work needs to be done?" A: The total work stays the same, but it gets divided among more workers. Like having more cashiers at a store - the same number of customers, but each cashier serves fewer people.

Q: "Why the cooldown periods?" A: To prevent "thrashing" - rapid scaling up and down. It's like waiting to see if the lunch rush is really over before sending staff home.

Q: "What happens if I hit the maximum task limit?" A: ECS stops scaling out. Your CPU might stay above target, but you won't get more tasks. You need to either increase the limit or optimize your application.

Key Takeaways

1. **Auto Scaling = Load Distribution:** More tasks share the same workload
2. **Metrics are Averages:** ECS looks at average CPU across all tasks
3. **Application Design Matters:** Your app must support multiple instances
4. **Fargate is Simpler:** AWS manages infrastructure for you
5. **EC2 Needs Two-Layer Scaling:** Tasks + Instances
6. **Cooldowns Prevent Chaos:** Wait periods ensure stable scaling decisions

This is fundamentally about **horizontal scaling** - adding more copies of your application to handle load, rather than making each copy more powerful (vertical scaling).

2What is ECS Cluster Auto Scaling?

ECS Cluster Auto Scaling (CAS) is used **only for EC2 launch type**.

It automatically **adds/removes EC2 instances** in your ECS cluster **based on task demand**.

How it works:

- Works with **Capacity Providers** linked to **Auto Scaling Groups (ASG)**.
- If ECS needs to place tasks but there's no space:
 - CAS increases EC2 count in the ASG.
- If EC2 instances are underutilized:
 - CAS scales in the ASG.

- This ensures **task placement never fails** due to lack of resources.

EC2 Scaling vs Task Scaling vs Fargate Scaling

1. EC2 Scaling (Infrastructure Level)

This is about adding/removing EC2 instances (physical servers).

Example:

- You have 2 EC2 instances (servers)
- Each server can run multiple containers
- When CPU usage is high across all servers, Auto Scaling Group adds a 3rd EC2 instance
- Now you have 3 servers available to run containers

Before scaling: [EC2-1] [EC2-2]

After scaling: [EC2-1] [EC2-2] [EC2-3]

2. ECS Task Scaling on EC2 (Application Level)

This is about adding/removing containers (tasks) on your existing EC2 instances.

Example:

- You have 2 EC2 instances
- Currently running 4 tasks total (2 per server)
- High demand triggers task scaling → ECS adds 2 more tasks
- Now you have 6 tasks distributed across the same 2 servers

Before: EC2-1[Task1, Task2] EC2-2[Task3, Task4]

After: EC2-1[Task1, Task2, Task5] EC2-2[Task3, Task4, Task6]

3. Fargate Scaling (Serverless)

With Fargate, you don't manage EC2 instances. AWS handles the infrastructure.

Example:

- You just specify: "I want 5 tasks running"
- Fargate automatically provisions the underlying compute
- Need more capacity? Scale to 8 tasks
- Fargate handles all the server management behind the scenes

Before: [Task1] [Task2] [Task3] [Task4] [Task5]

After: [Task1] [Task2] [Task3] [Task4] [Task5] [Task6] [Task7] [Task8]

Real-World Scenario

Your Image Conversion Service:

With ECS on EC2:

1. **Normal load:** 2 EC2 instances, 10 tasks total
2. **Traffic increases:** Task scaling adds 20 more tasks on existing servers
3. **Servers getting overloaded:** EC2 scaling adds more instances
4. **More tasks needed:** Task scaling distributes new tasks across all instances

With ECS on Fargate:

1. **Normal load:** 10 tasks running
2. **Traffic increases:** Scale to 30 tasks
3. **No server management:** Fargate handles everything automatically

Key Differences

Aspect	EC2 Scaling	Task Scaling (EC2)	Fargate Scaling
What scales	Physical servers	Containers/Apps	Containers only

Aspect	EC2 Scaling	Task Scaling (EC2)	Fargate Scaling
You manage	Server capacity	App instances	Nothing
Speed	Slower (boot time)	Faster	Fastest
Cost	Pay for servers	Pay for servers	Pay per task
Complexity	High	Medium	Low

Mental Model:

- **EC2 Scaling** = Adding more restaurants (buildings)
- **Task Scaling** = Adding more chefs in existing restaurants
- **Fargate Scaling** = Adding more chefs, AWS manages the restaurants

2. What are the different ECS scheduling and placement strategies?

Amazon ECS Scheduling Strategies

ECS offers two main **schedulers** to manage task placement and execution:

1 Service Scheduler (For long-running apps)

Used for **services that need to stay running**, such as web apps or APIs.

◆ Service Types

- **REPLICA** (default):
 - Runs and maintains a **specified number of tasks**.
 - Automatically replaces failed/stopped tasks.
 - Best for **stateless apps, APIs, and web servers**.
- **DAEMON**:
 - Runs **one task per EC2 instance** that matches placement constraints.
 - Ideal for **cluster-wide agents** (e.g., logging, monitoring).

⚙️ Task Placement Strategies (EC2 only)

- **binpack**:
 - Packs tasks on instances using **least available CPU/memory**.
 - Optimizes for **cost-efficiency**.
- **spread**:
 - Spreads tasks evenly by **AZ, instance ID, or custom attribute**.
 - Ensures **high availability** and **fault tolerance**.
- **random**:
 - Places tasks **randomly**.
 - Default for EC2 when no strategy is specified.

2 Standalone Task Scheduler (For one-time/batch jobs)

Used for **short-lived, run-once tasks** like batch processing or scripts.

◆ Ways to Run Standalone Tasks

- **RunTask (manual)**:
 - Use API or CLI to **start a task once**.
 - Task **runs and stops** on its own.
- **EventBridge Scheduler (automated)**:
 - Schedule **cron-like recurring tasks** (e.g., daily at 2 AM).

- Triggers **RunTask** automatically.

1. Which of the following is the easiest way to configure ECS auto scaling?

- A. Step scaling with EC2 instance metrics
- B. Manual task definition update
- C. Target tracking policy using CloudWatch metrics
- D. Scheduled Lambda triggers

C. Target tracking policy using CloudWatch metrics ✓

2. Which launch type supports ECS Cluster Auto Scaling?

- A. Fargate
- B. EC2 only
- C. Both EC2 and Fargate
- D. Lambda

B. EC2 only ✓

3. What's the difference between ECS Task Scaling and EC2 Scaling?

- A. Task scaling changes the VPC
- B. Task scaling changes the number of containers; EC2 scaling changes the number of servers
- C. EC2 scaling affects task definitions
- D. There is no difference

B. Task scaling changes the number of containers; EC2 scaling changes the number of servers ✓

4. Which ECS scheduling strategy ensures one task runs per EC2 instance?

- A. REPLICA
- B. DAEMON
- C. RUNONCE
- D. CRON

B. DAEMON ✓
