

Auto Scaling ECS Applications



Justin Menga
FULL STACK TECHNOLOGIST
@jmenga pseudo.co.de

Introduction

Auto Scaling ECS Applications

- ECS Auto Scaling
 - Scales your ECS services
- EC2 Auto Scaling
 - Scales your ECS clusters
- Auto Scaling Solution
 - CloudWatch Events
 - ECS Capacity Manager Lambda Function
 - CloudWatch Metrics and Alarms
 - Auto Scaling Policies
 - CloudFormation support

ECS Auto Scaling Challenges

ECS Auto Scaling

ECS Service Auto Scaling

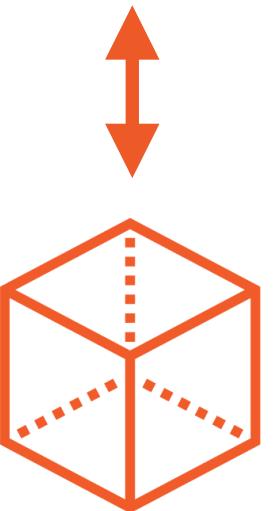


Add bullets here

ECS Cluster Auto Scaling

Add bullets here

Application Load



**ECS Service
Auto Scaling**

ECS Service

Application Load

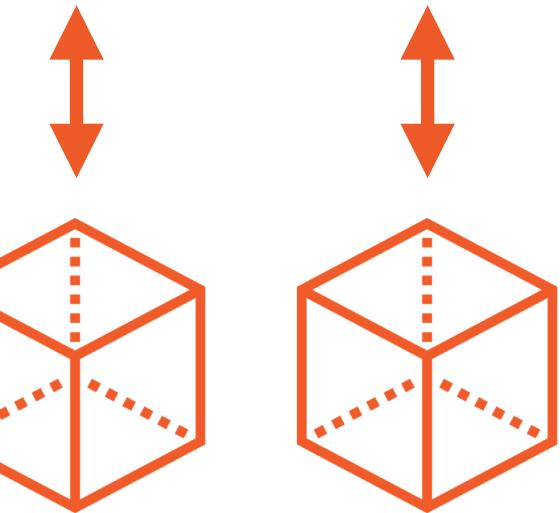
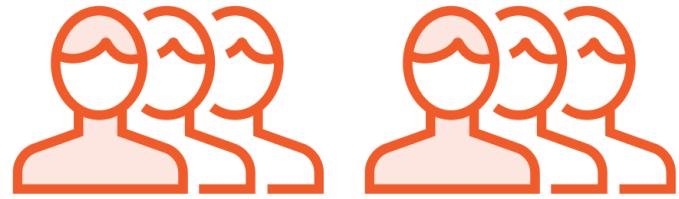


ECS Service

**ECS Cluster
Auto Scaling**

ECS Container Instance

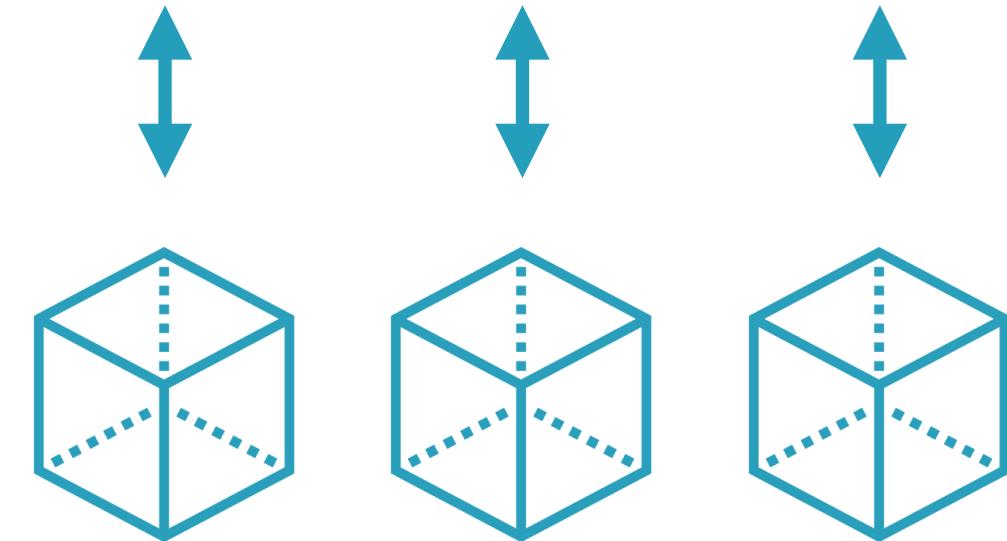
Application Load



**ECS Service
Auto Scaling**

ECS Service

Application Load

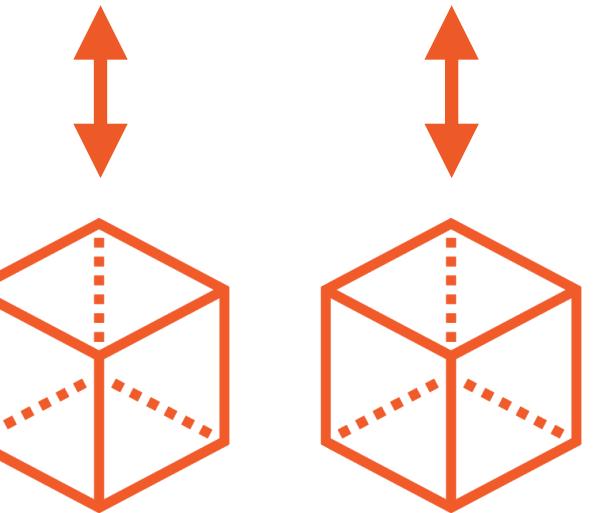
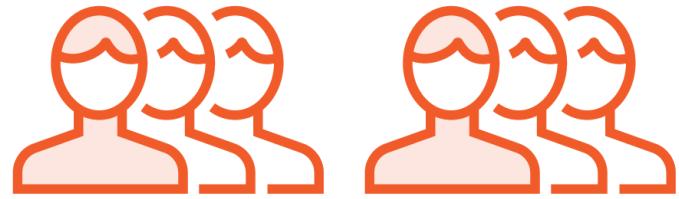


ECS Service

**ECS Cluster
Auto Scaling**

ECS Container Instance

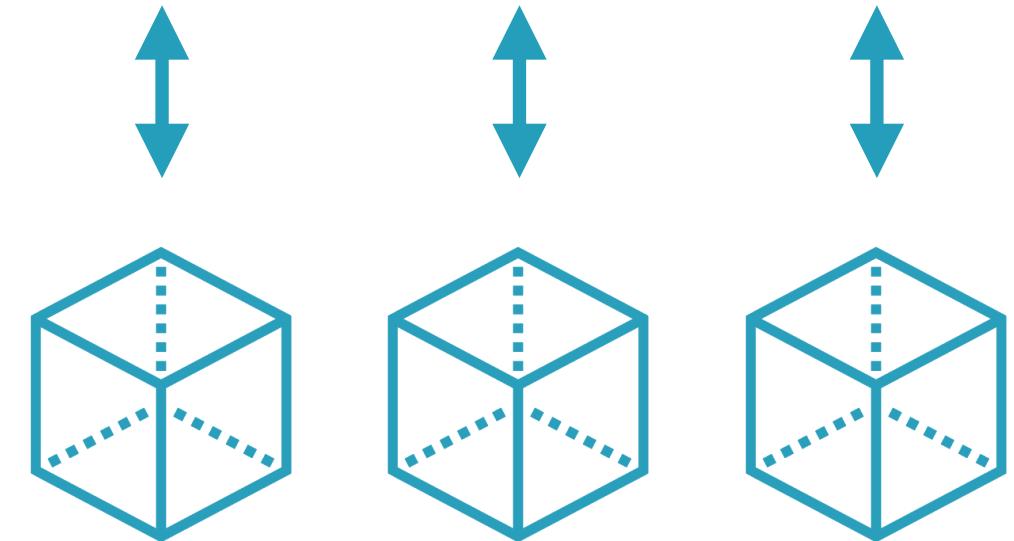
Application Load



**ECS Service
Auto Scaling**

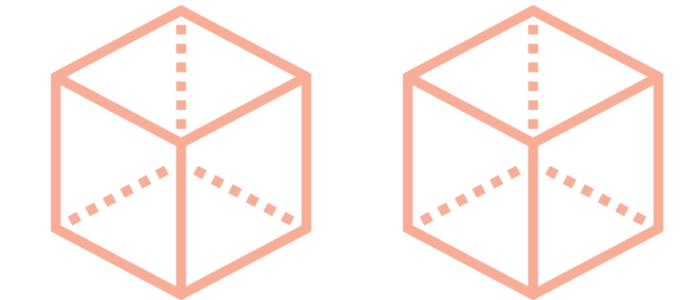
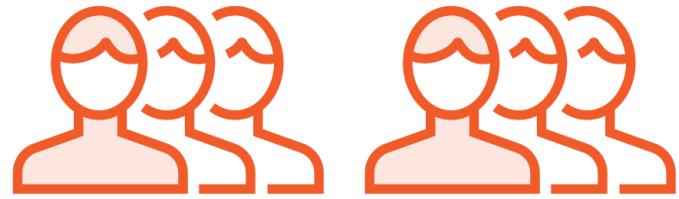
**ECS Cluster
Auto Scaling**

Application Load



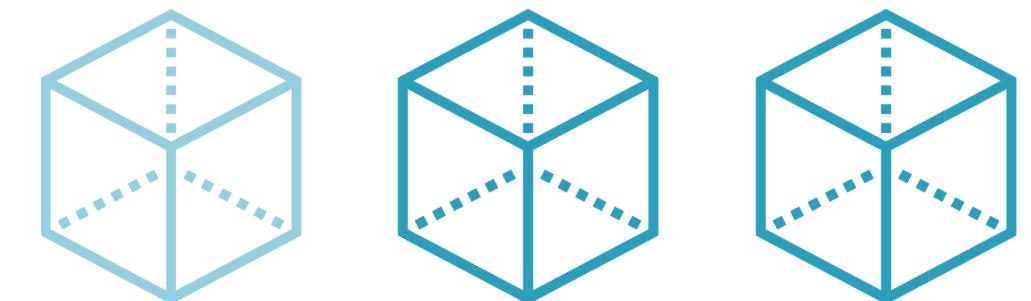
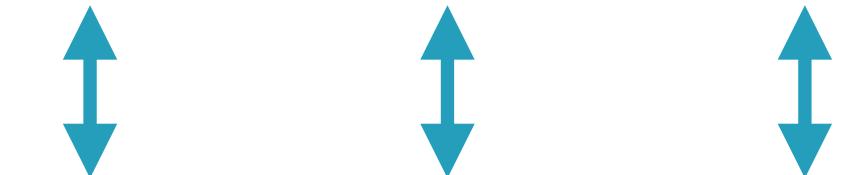
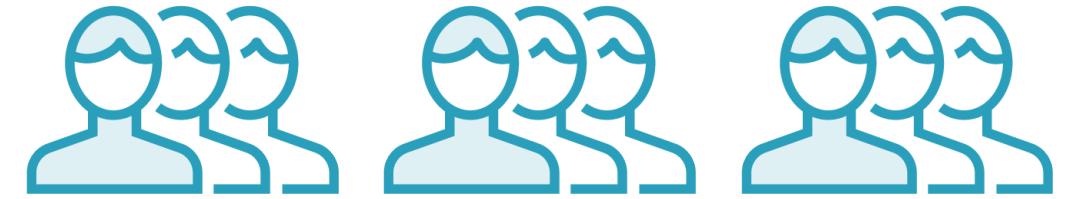
ECS Container Instance

Application Load



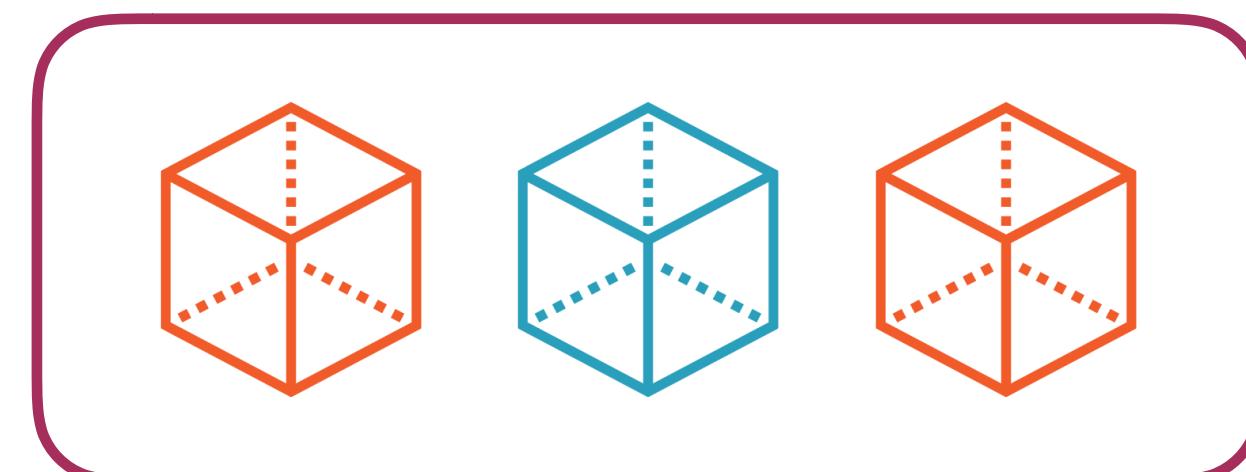
**ECS Service
Auto Scaling**

Application Load



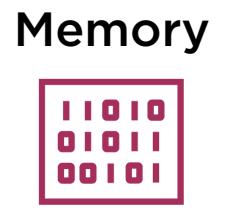
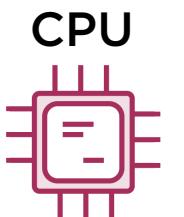
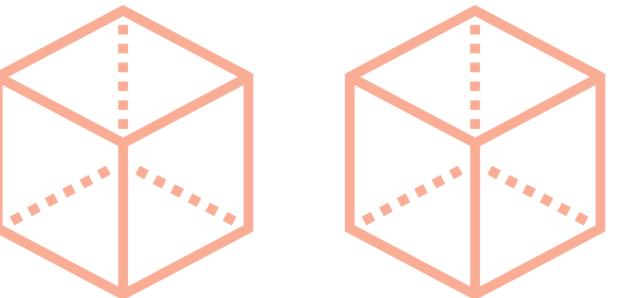
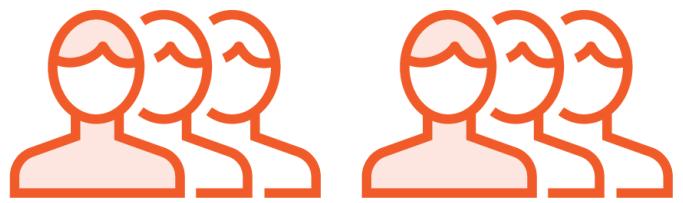
**ECS Cluster
Auto Scaling**

ECS Container Instance



ECS Service Auto Scaling

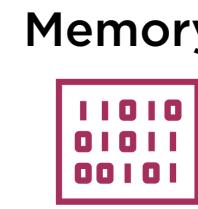
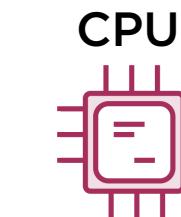
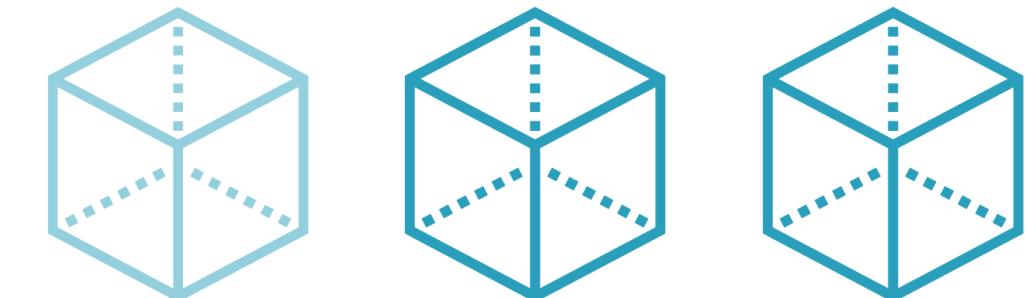
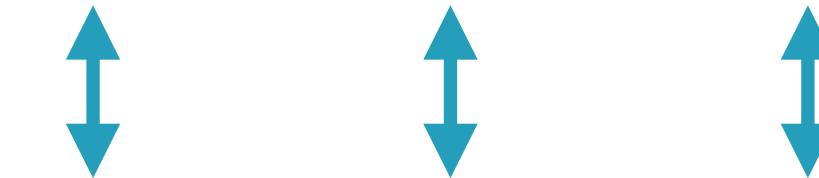
Application Load



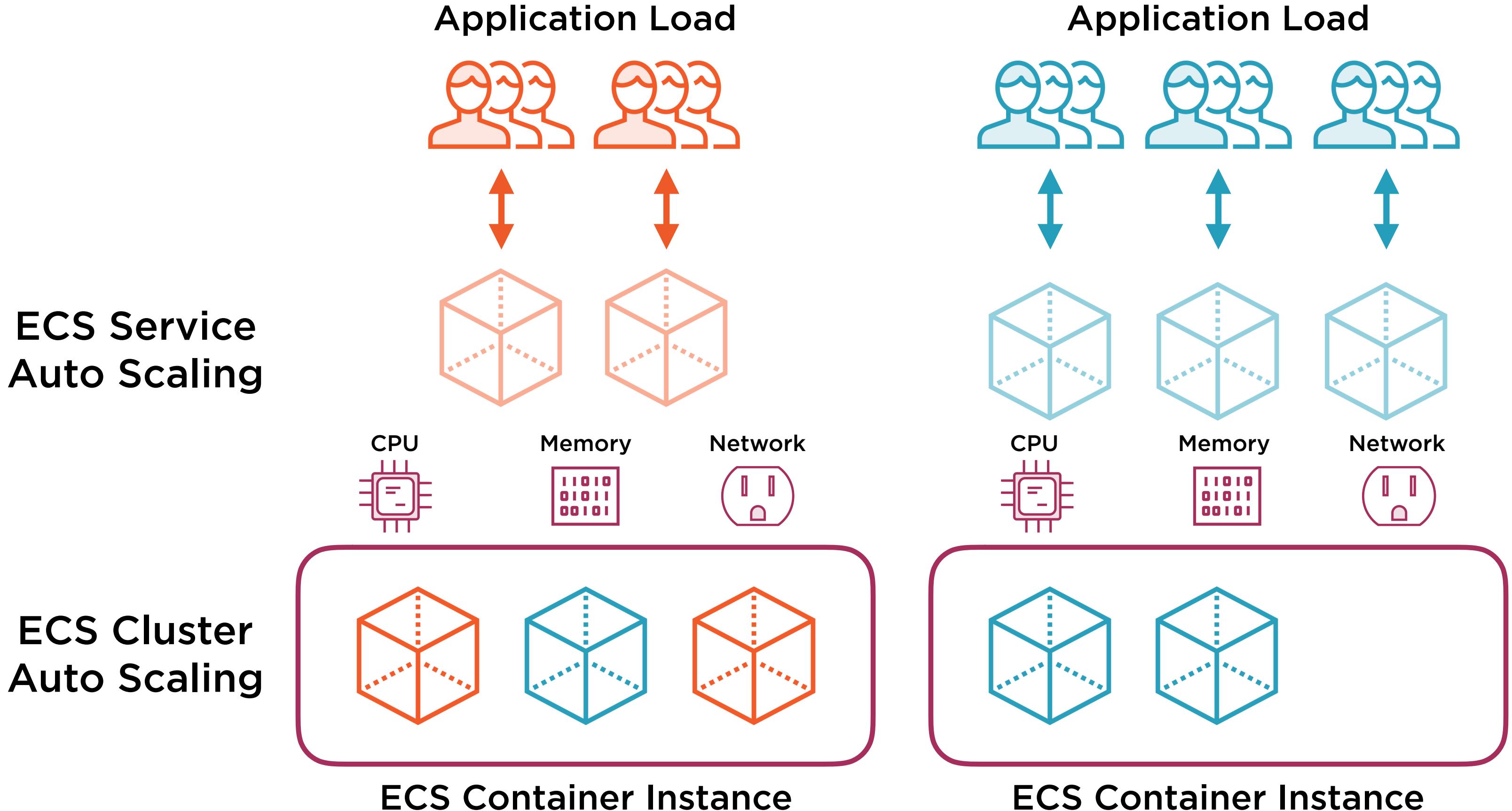
ECS Cluster Auto Scaling

ECS Container Instance

Application Load



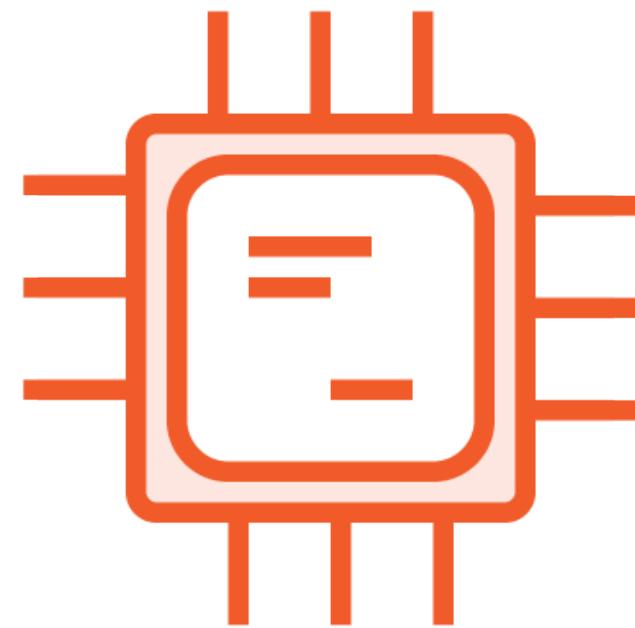
ECS Container Instance



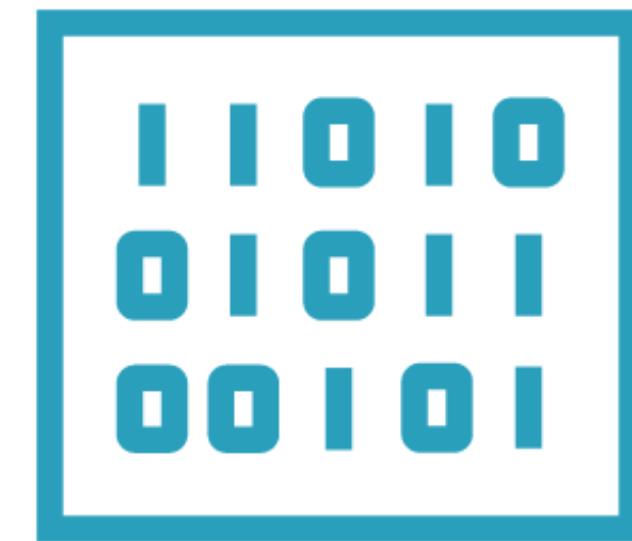
ECS Auto Scaling Challenges

Understanding ECS Resources

ECS Resource Types



CPU



Memory

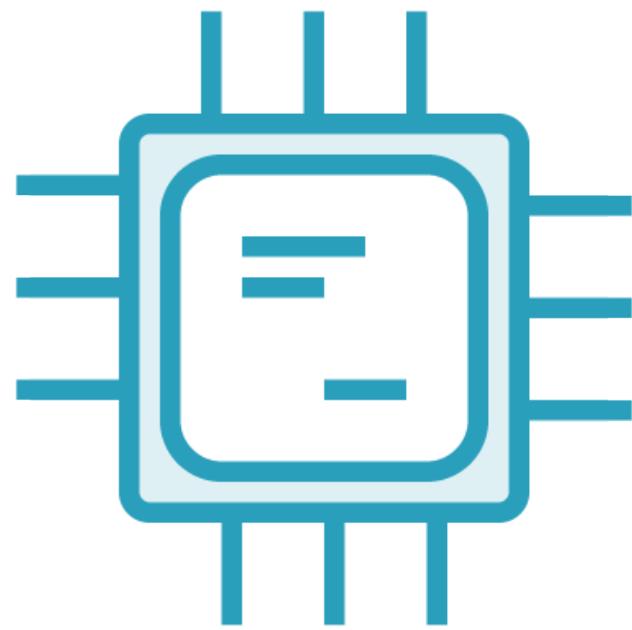


Network Ports

CPU Resources



CPU Reservations
CPU Units



**1024 Units
per Core**

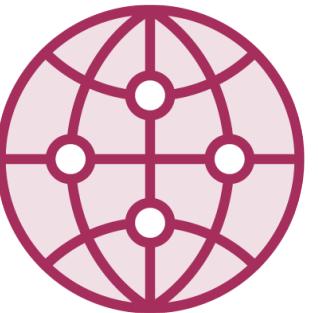


**Used under
CPU Load**

CPU Scheduling based upon CPU Units

Reservations	Container A 256 Units	Container B 256 Units	Container C 512 Units
Unit Ratio	$256 / 1024$	$256 / 1024$	$512 / 1024$
CPU Time	25%	25%	50%

EC2 Container Service

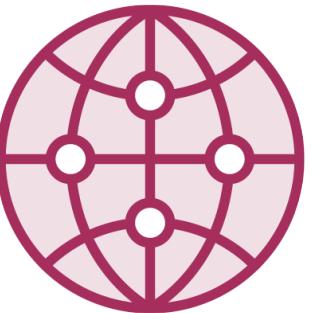


Deploy Container A



Available
CPU
1024

EC2 Container Service

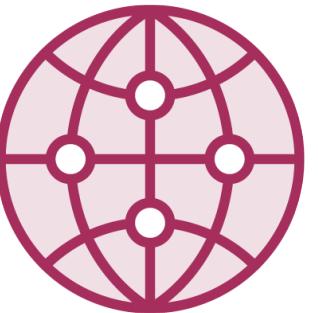


Deploy Container A



Available
CPU
824

EC2 Container Service



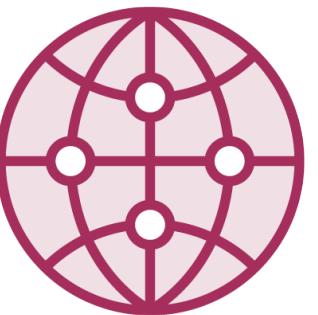
Deploy Container B



ECS Container Instance
1024 CPU Units

Available
CPU
824

EC2 Container Service



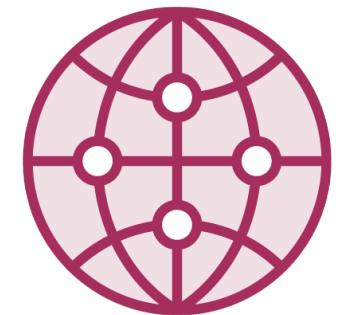
Deploy Container B



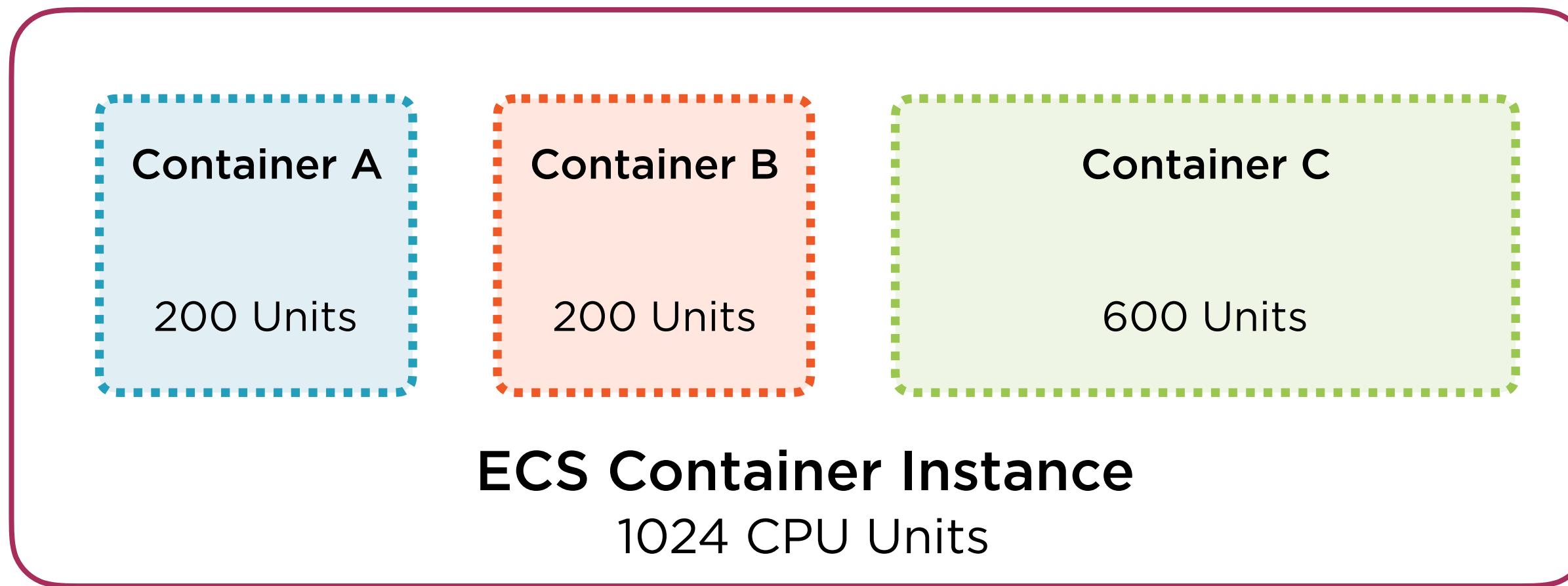
ECS Container Instance
1024 CPU Units

Available
CPU
624

EC2 Container Service

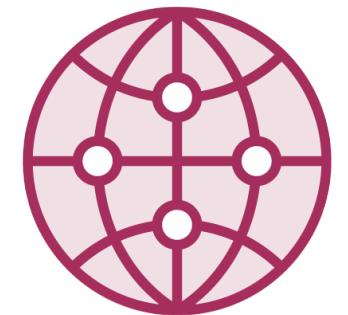


Deploy Container C

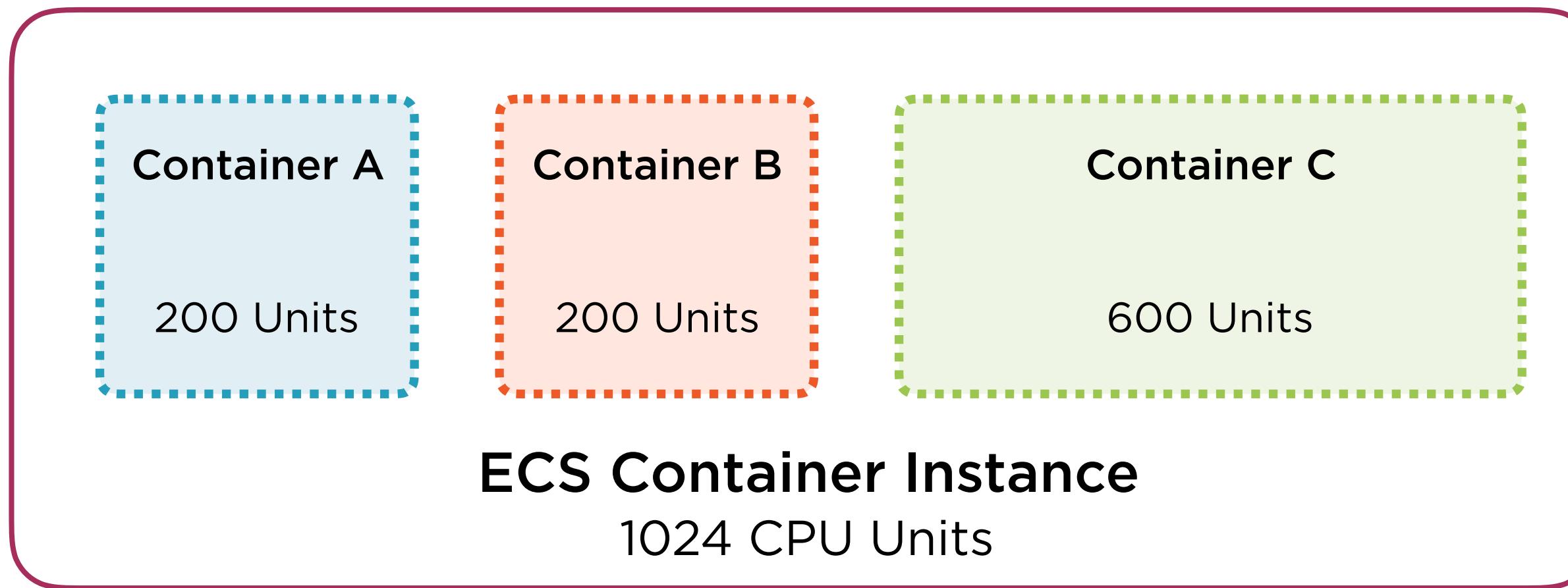


Available
CPU
624

EC2 Container Service

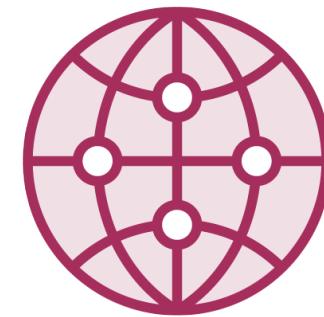


Deploy Container C



Available
CPU
24

EC2 Container Service



Deploy Container D



NO CPU RESOURCES
AVAILABLE

Container A

200 Units

Container B

200 Units

Container C

600 Units

Available
CPU
24

ECS Container Instance
1024 CPU Units

Memory Resources

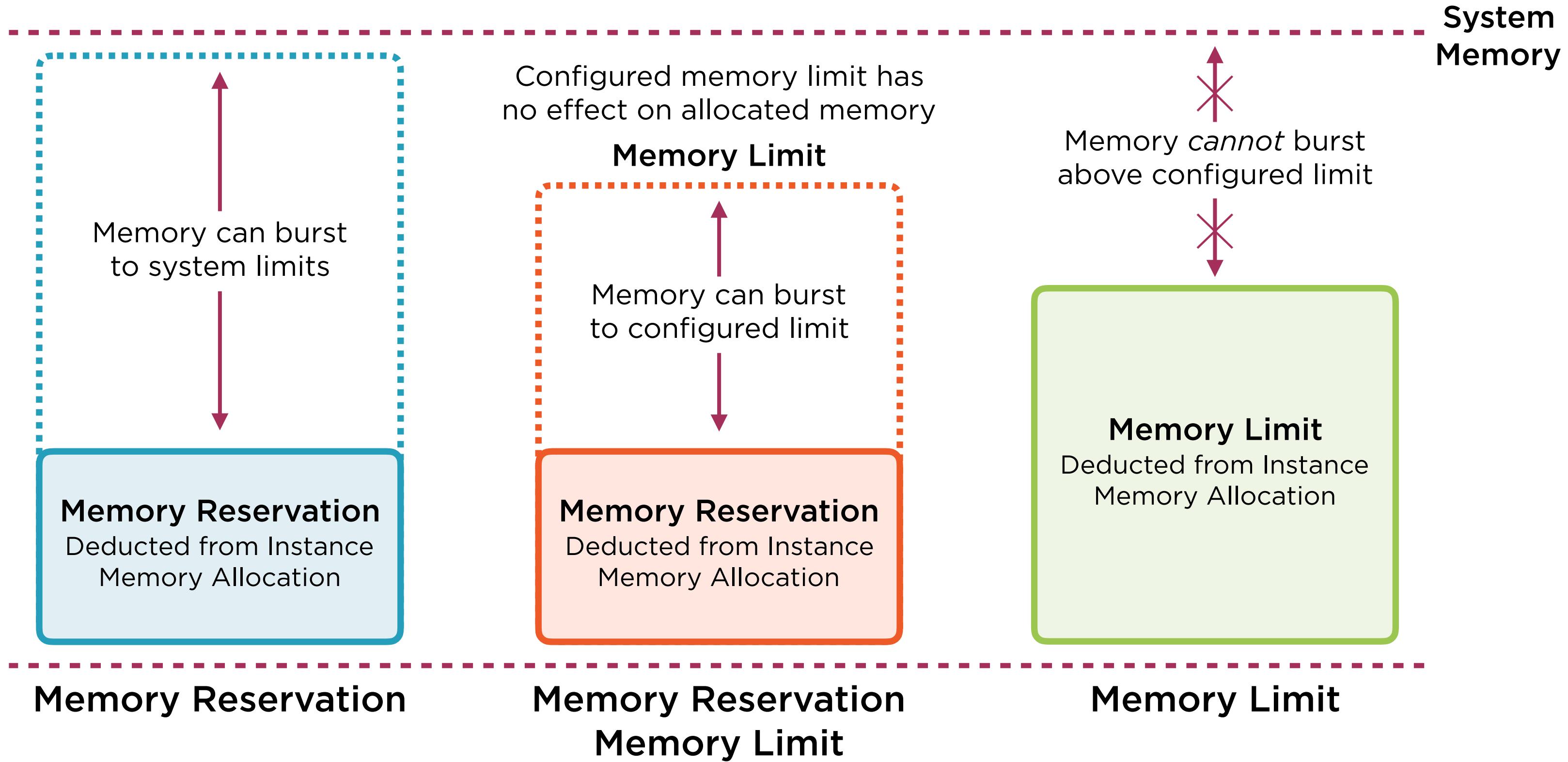


**Memory
Reservations**



**Memory
Limits**

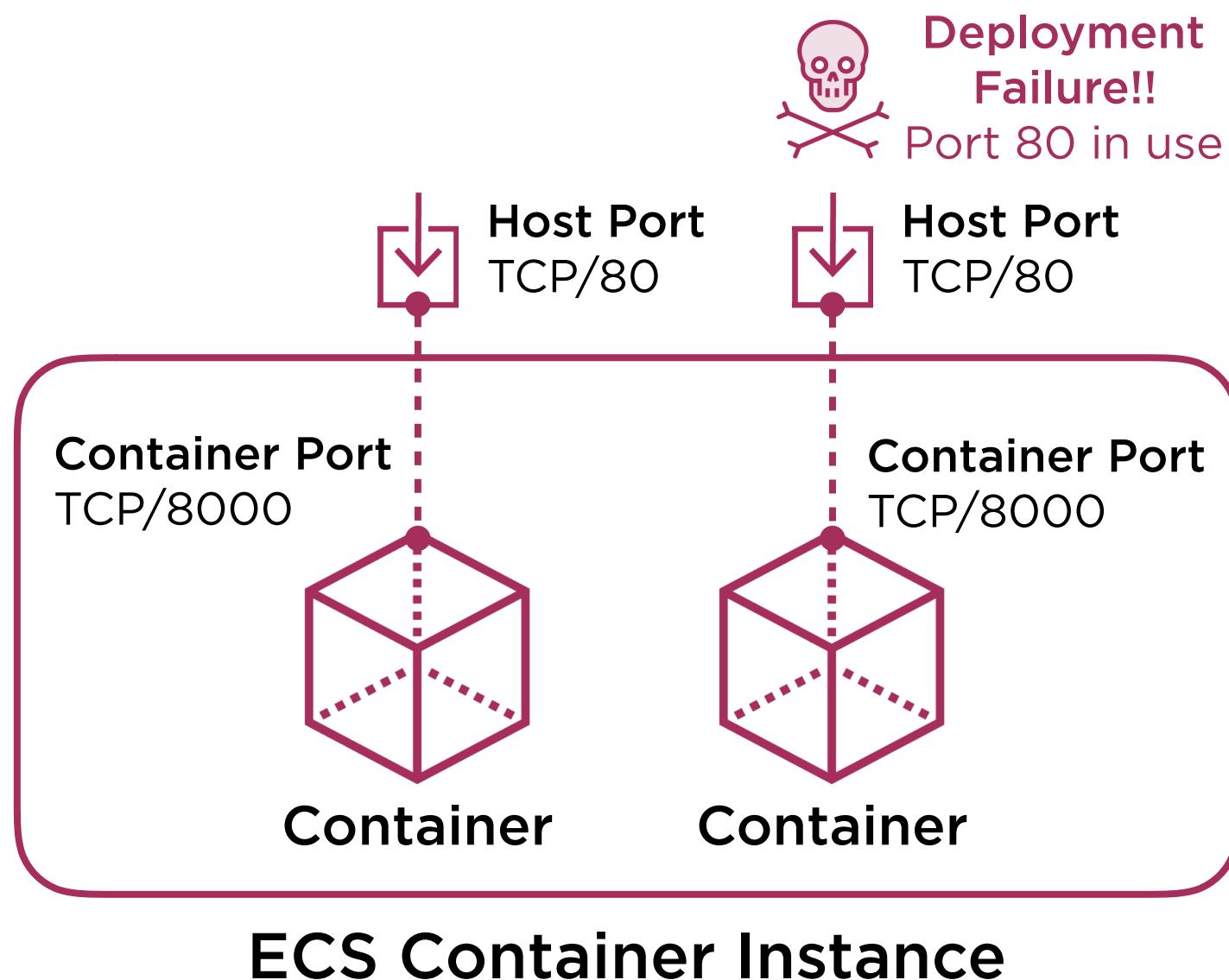
Memory Resources



Network Port Resources

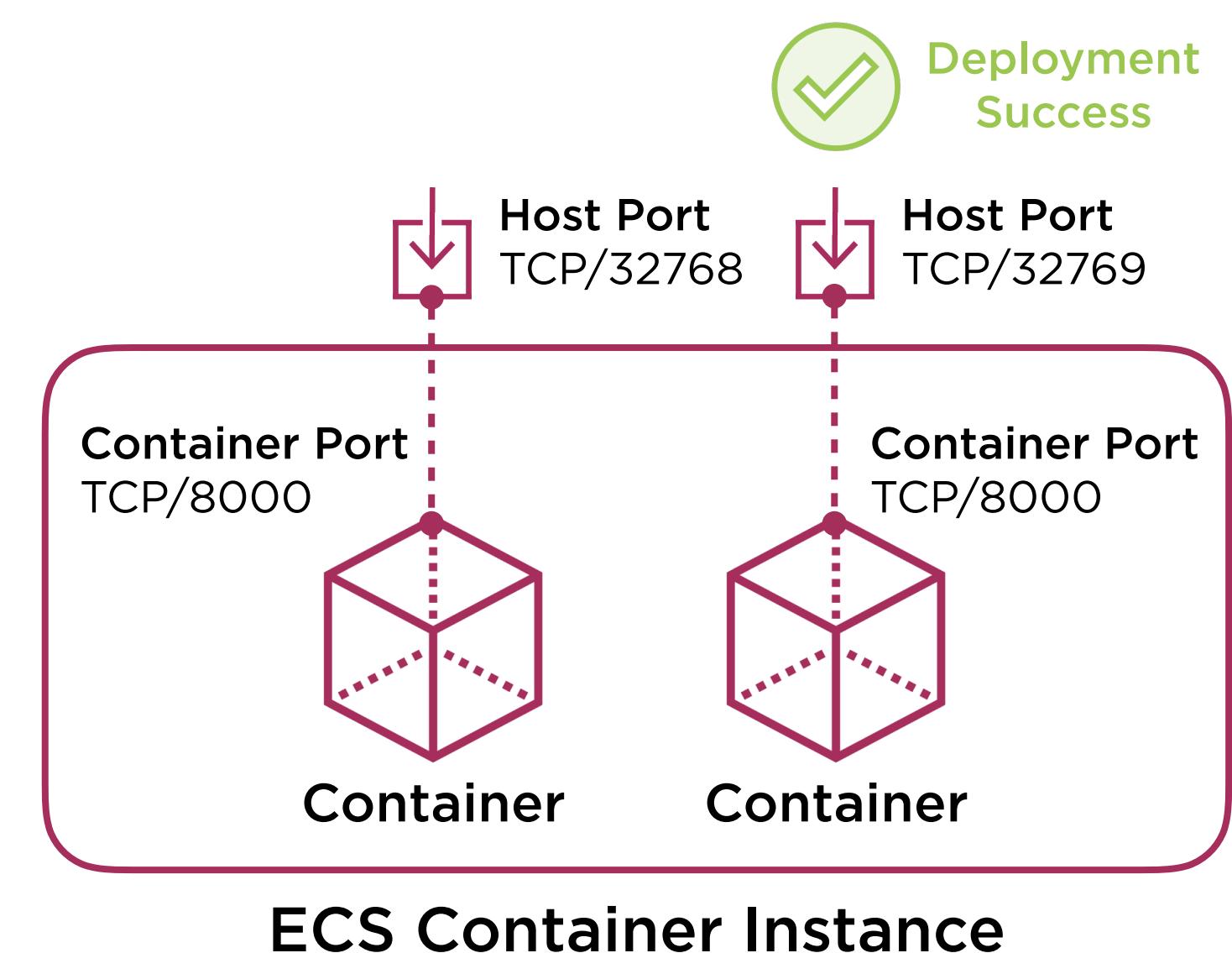
Static Port Mapping

Only one container per port per instance



Dynamic Port Mapping

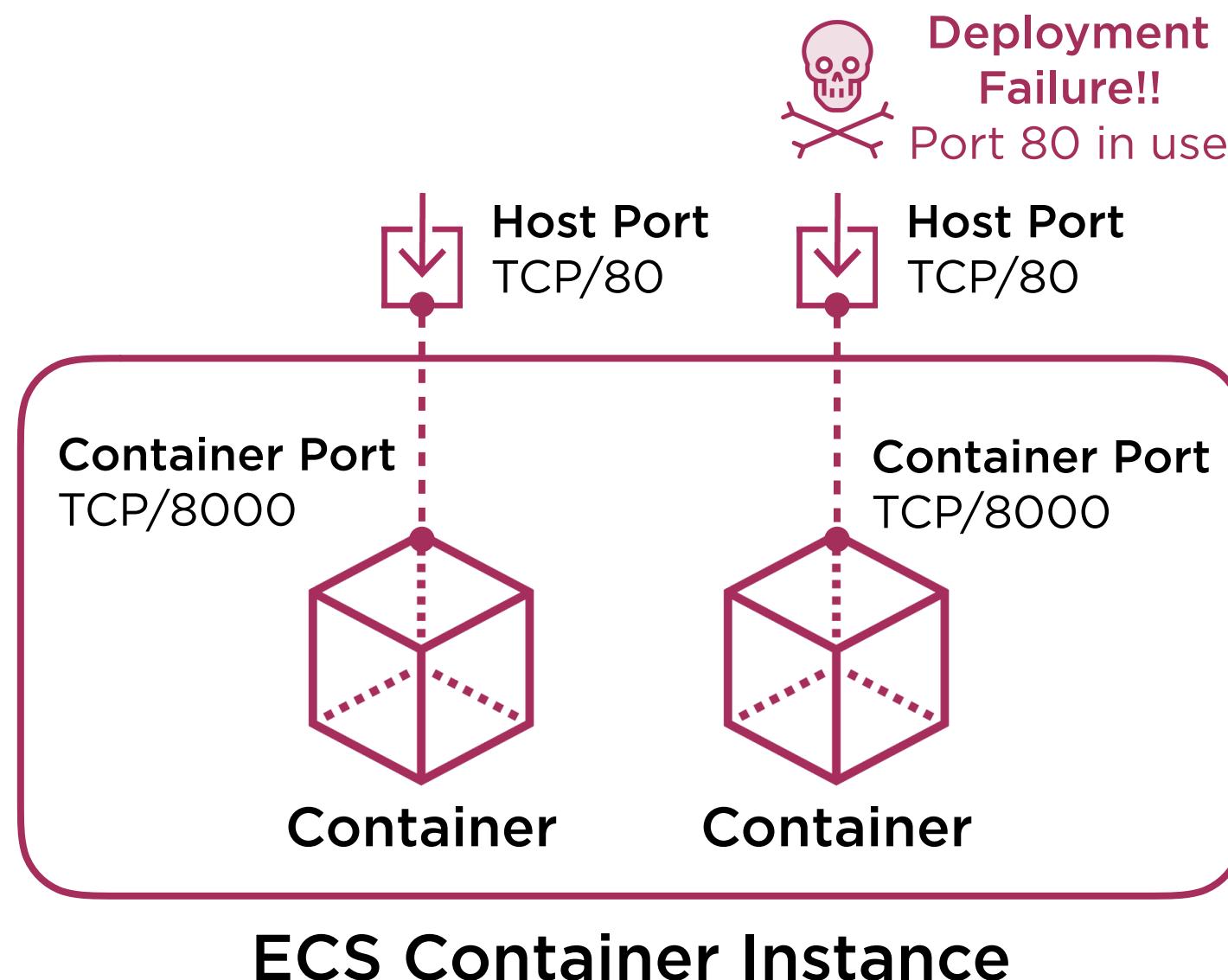
Multiple containers per port per instance



Dynamic Port Mapping

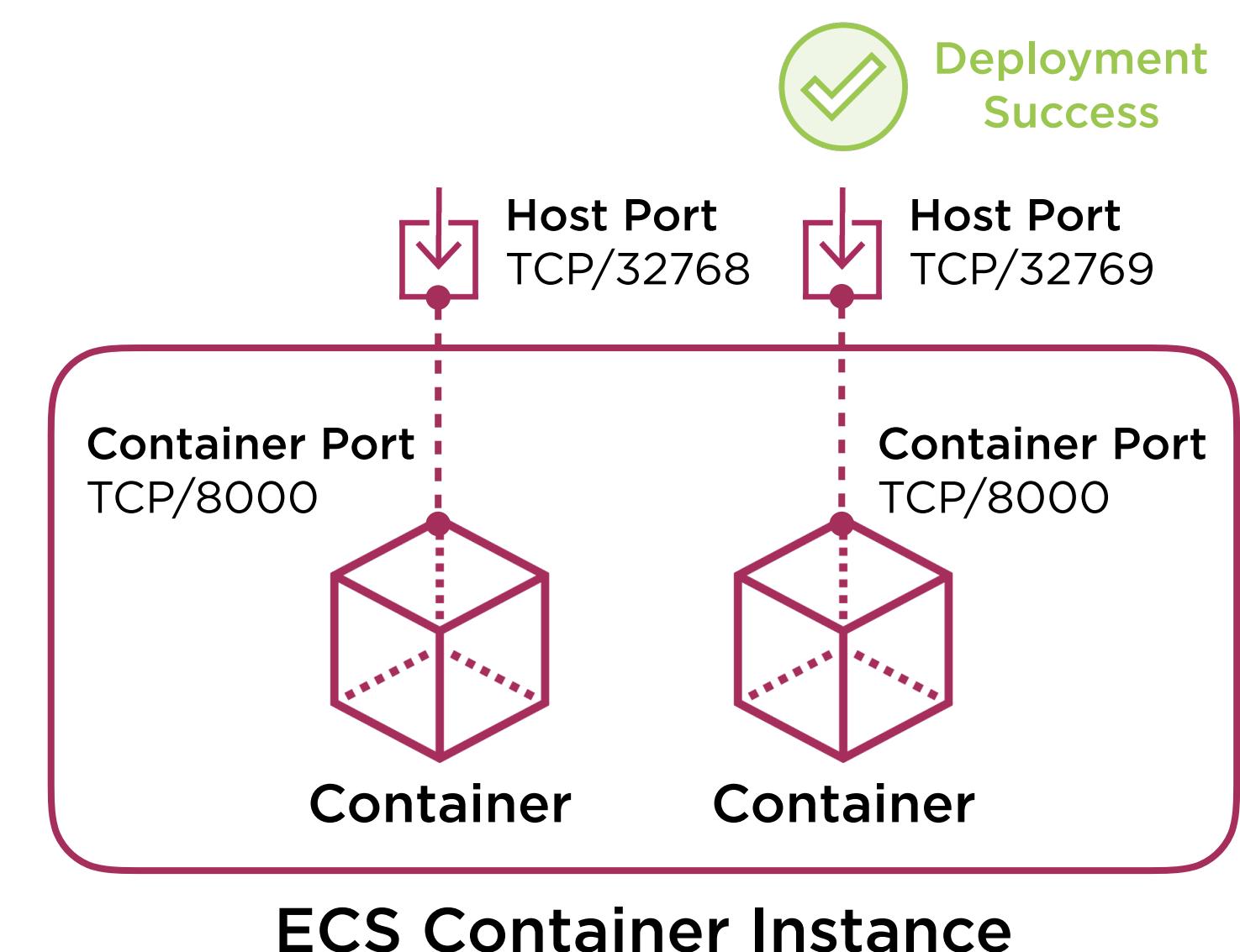
Static Port Mapping

Only one container per port per instance



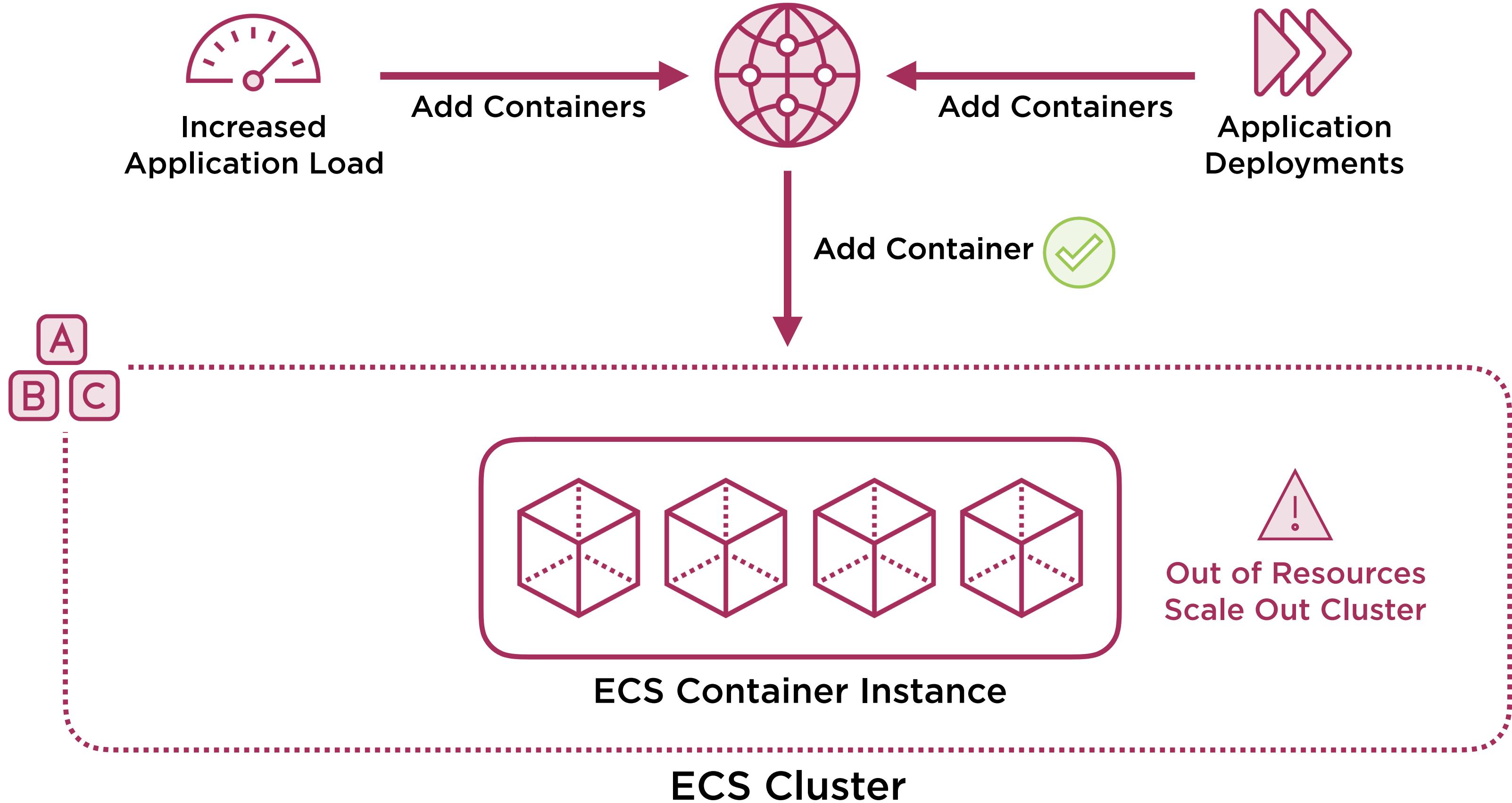
Dynamic Port Mapping

Multiple containers per port per instance

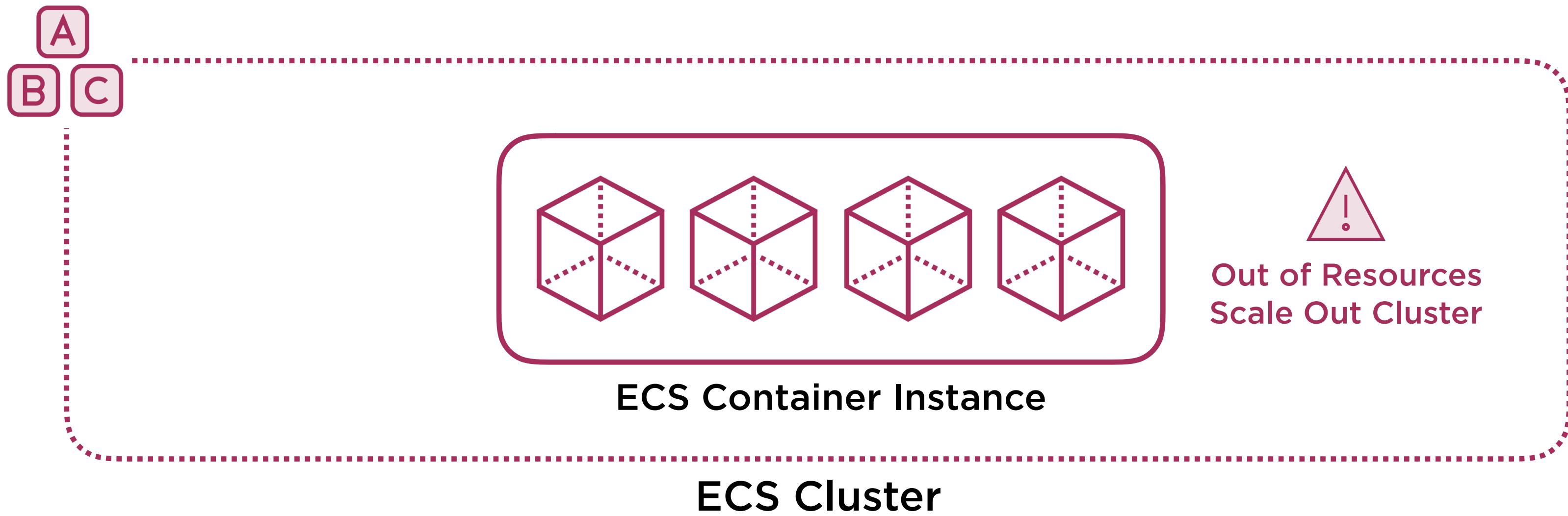
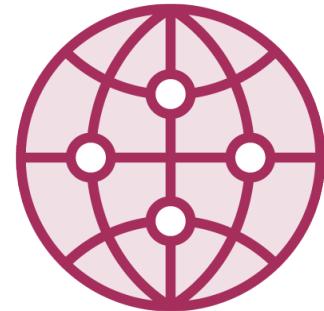


ECS Capacity Management

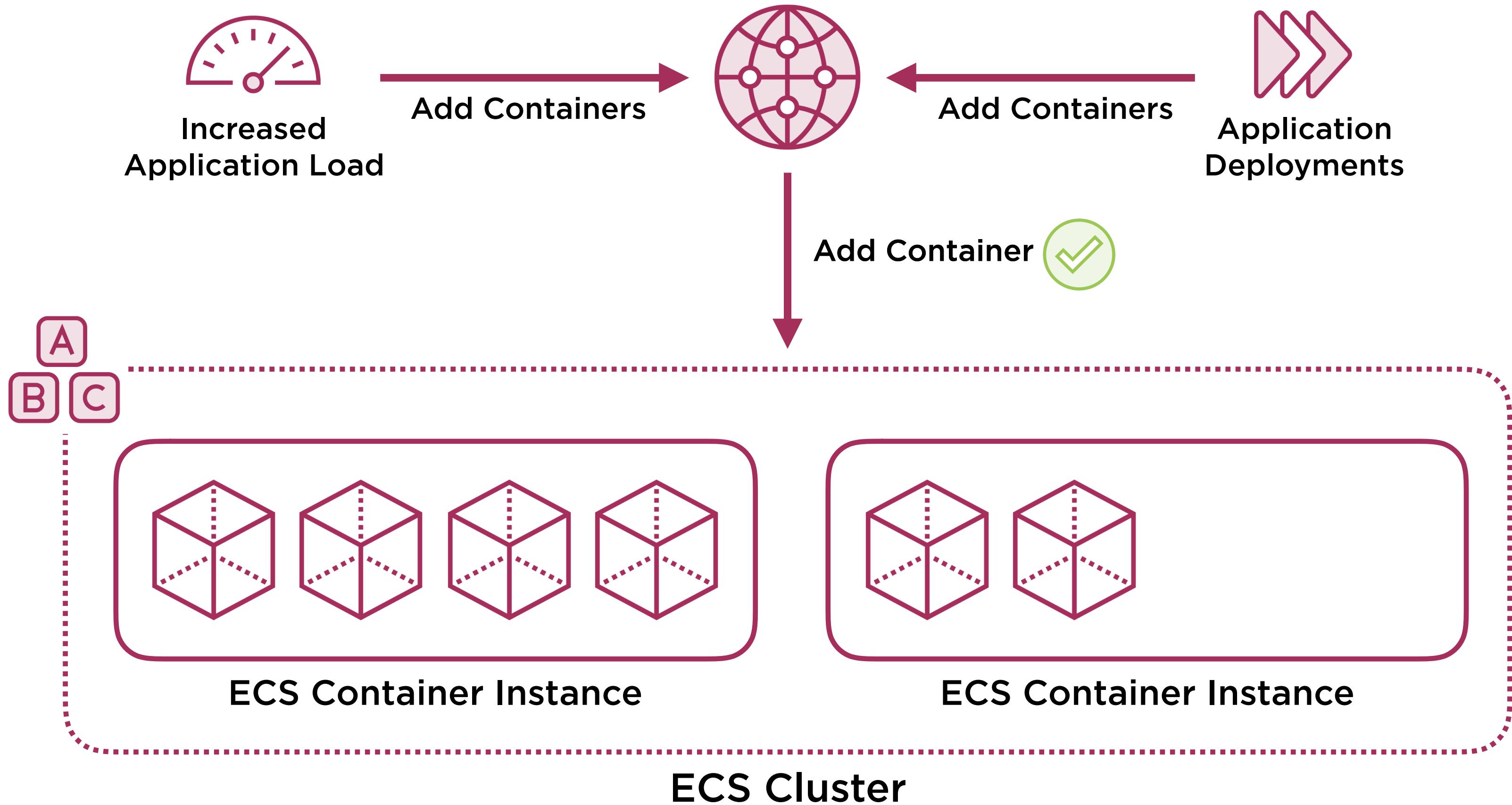
EC2 Container Service



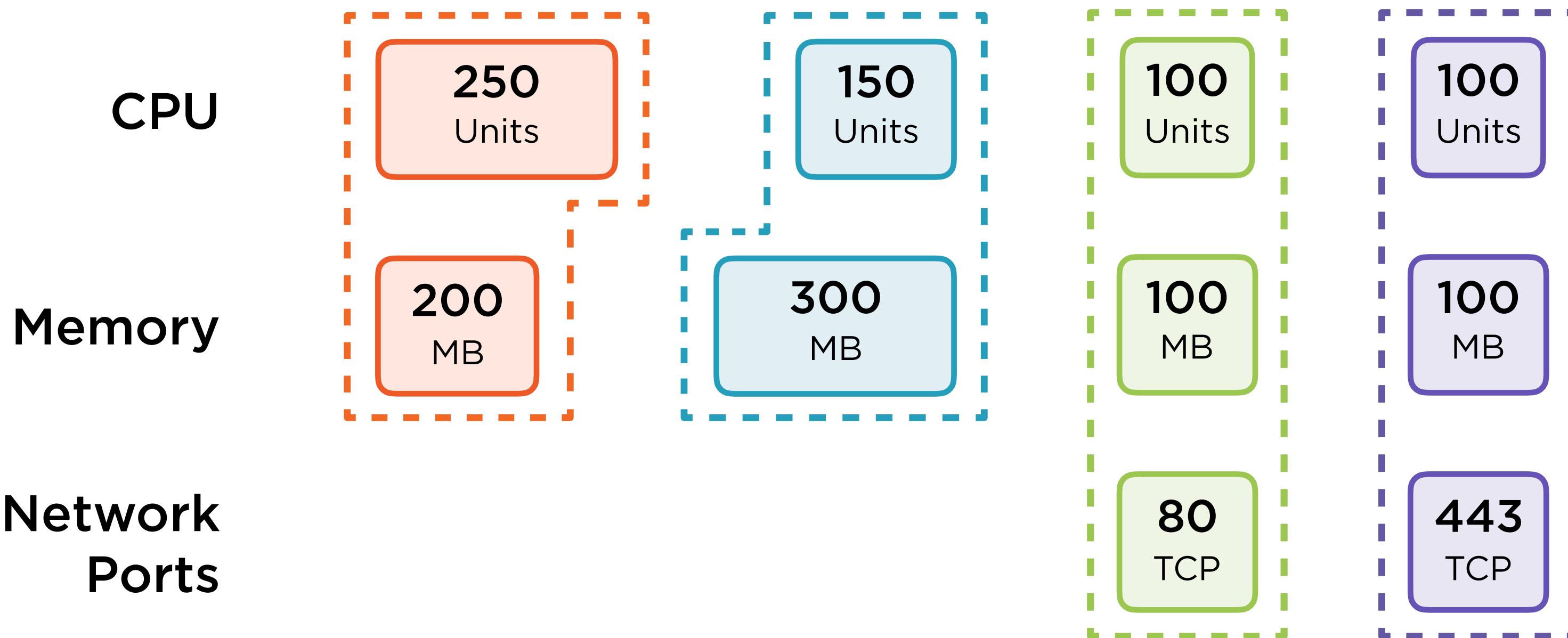
EC2 Container Service



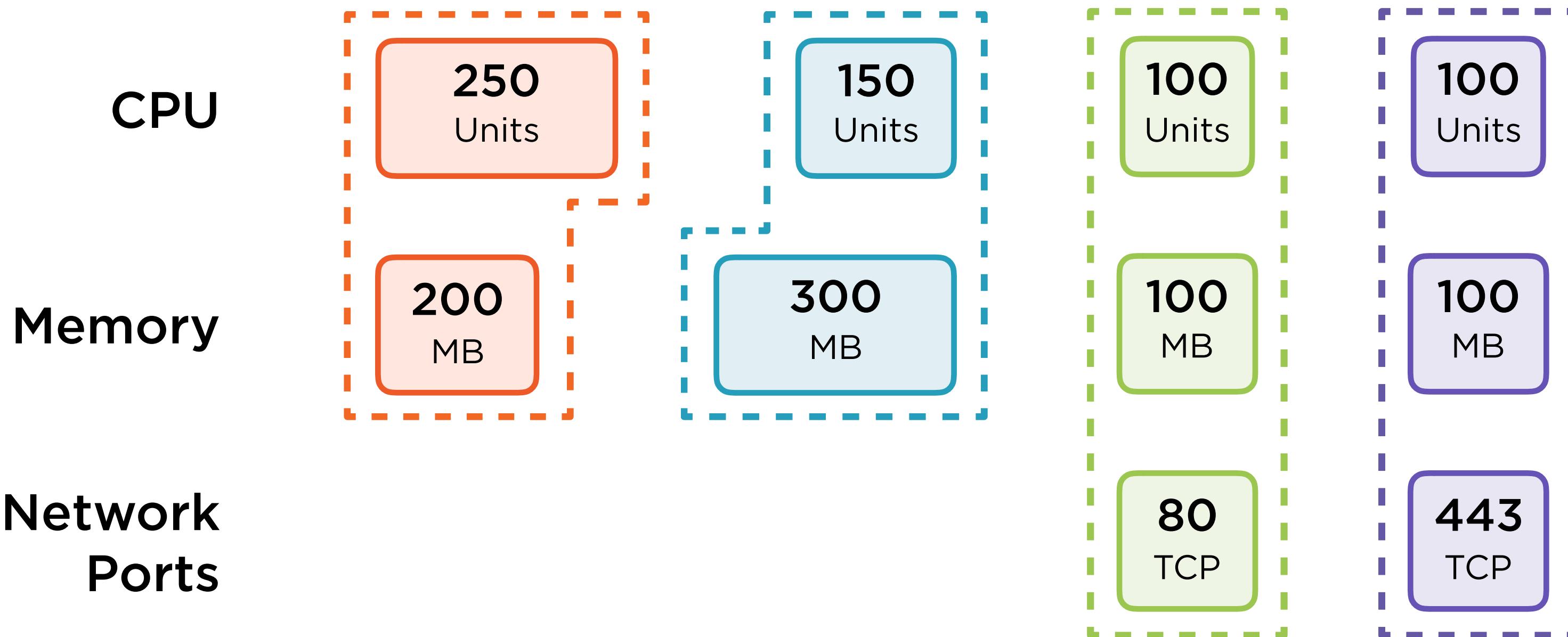
EC2 Container Service

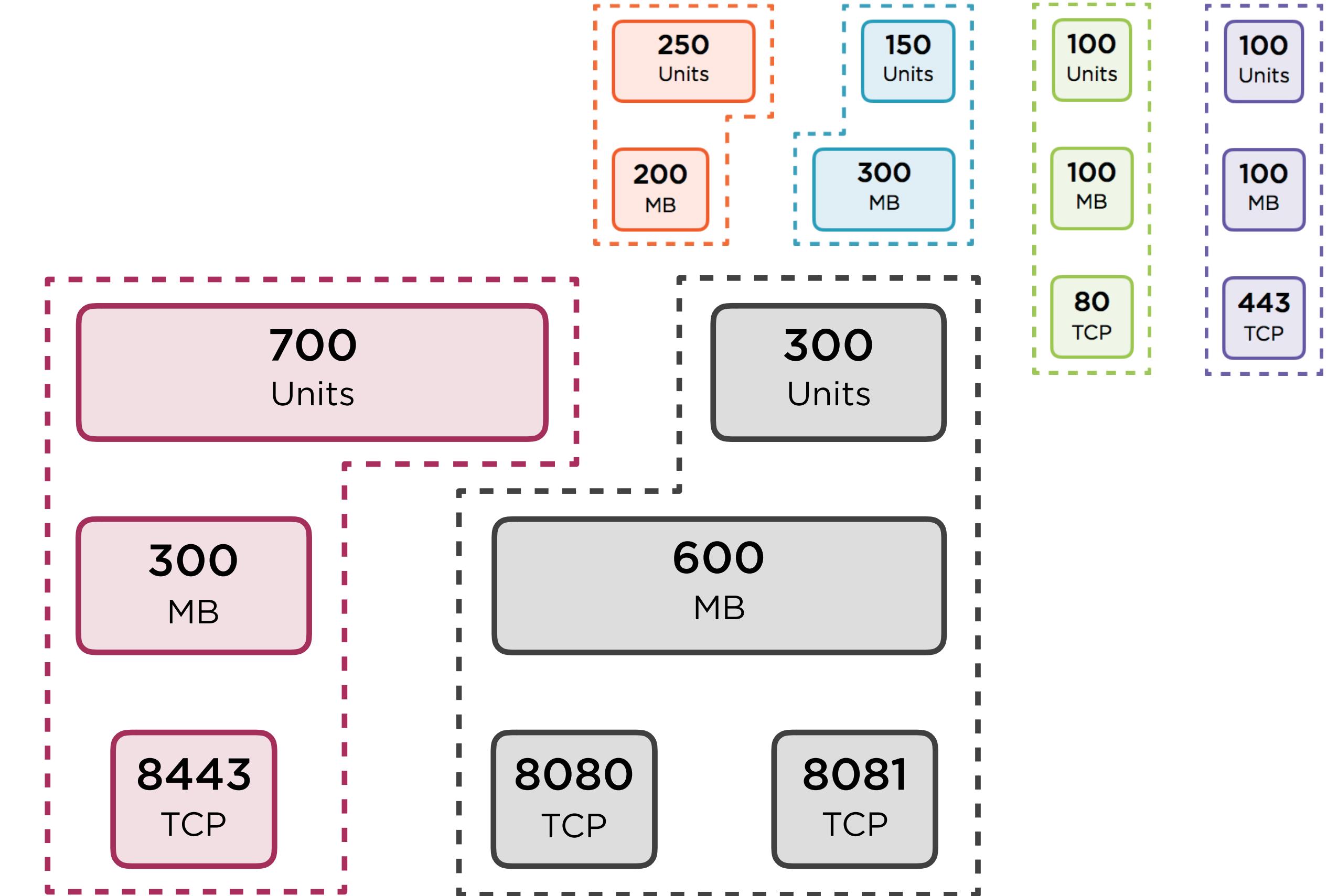


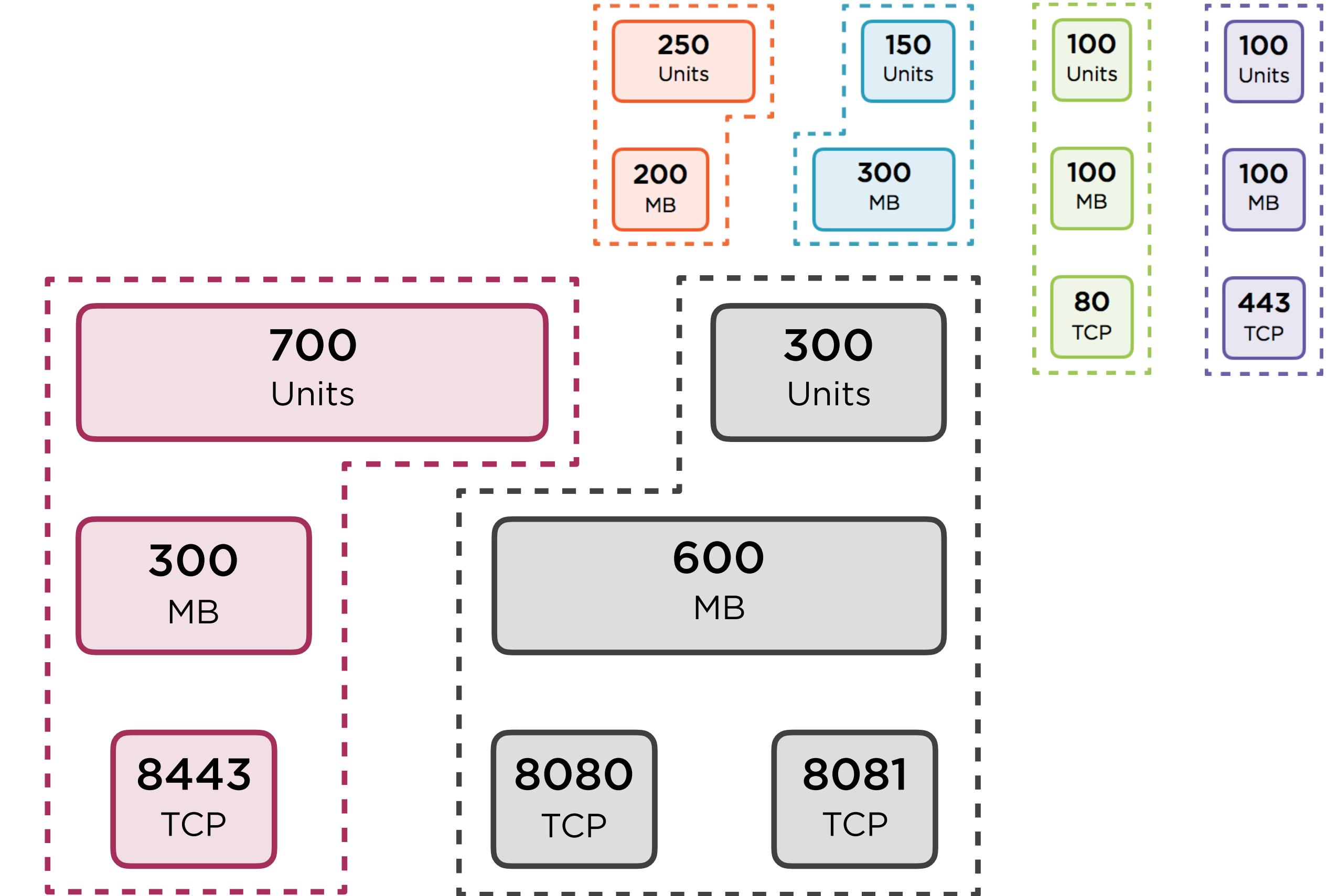
Dealing with Resources for Multiple Containers

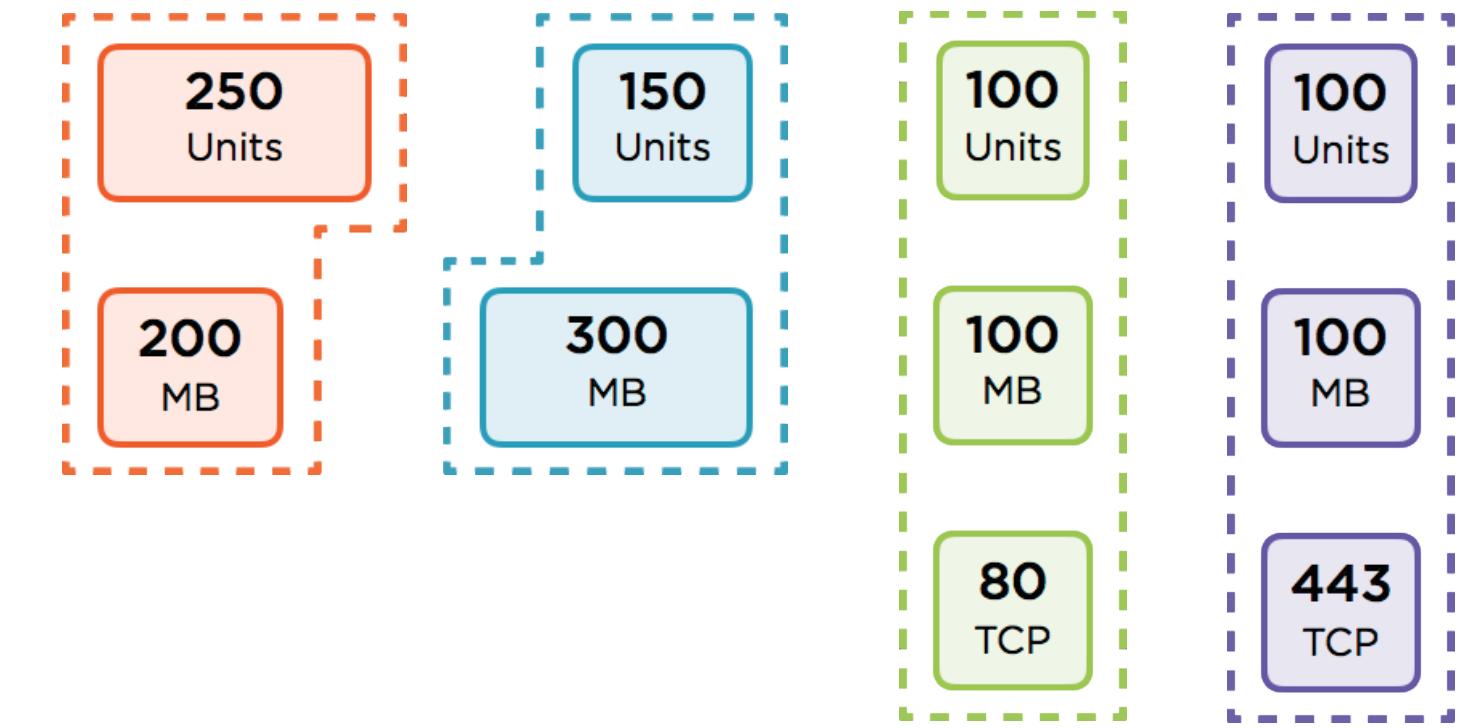
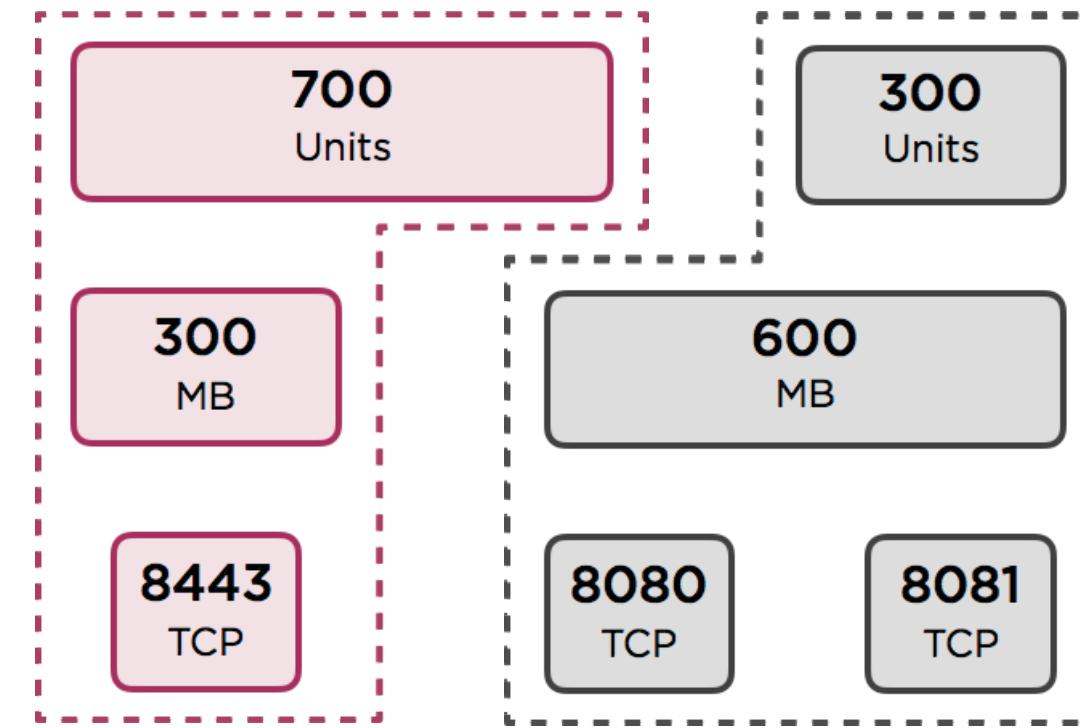


Dealing with Resources for Multiple Containers







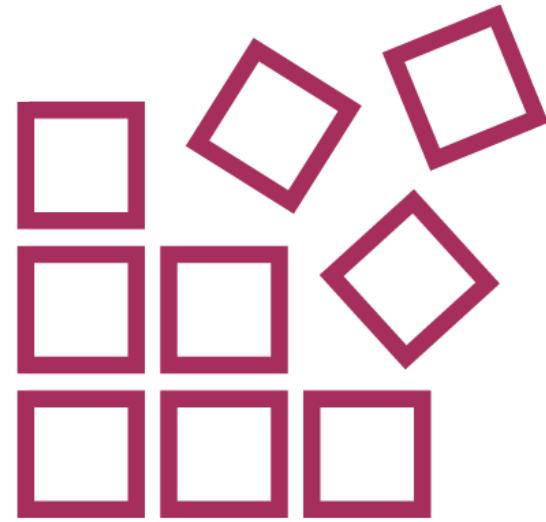


ECS Container Instance



ECS Container Instance

Calculating ECS Capacity



1. Calculate Individual Resource Capacity

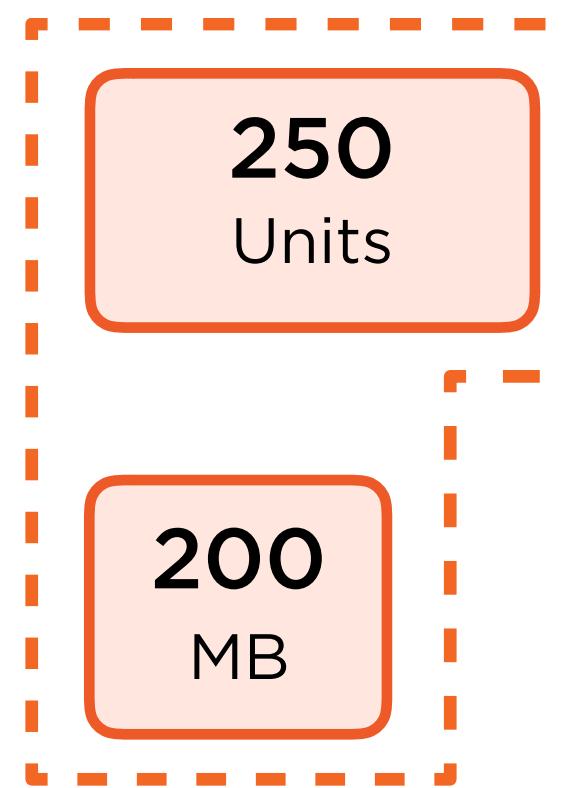
- Determine container that requires largest amount of each resource
- Calculate current capacity in the cluster expressed in terms of containers
- E.g. I can currently schedule two of the containers that consume the most CPU in my cluster
- Repeat for each type of resource

2. Calculate Overall Capacity

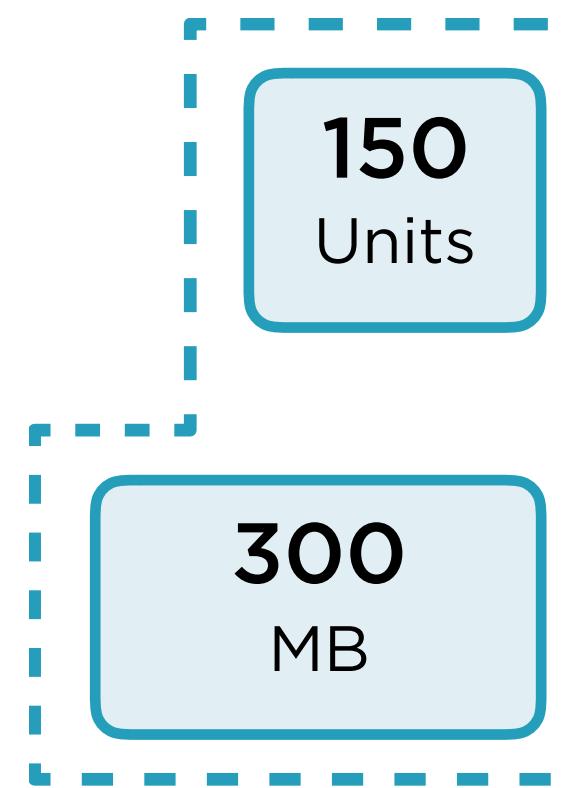
- Take the lowest value of each of the current individual resource capacities

Let's play Container Tetris!

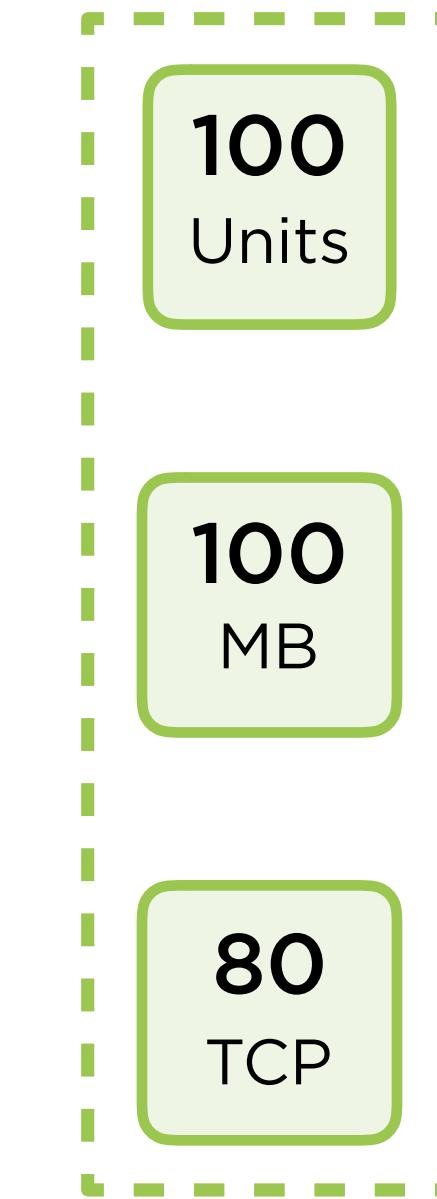
**Worst Case CPU
of 250 Units**



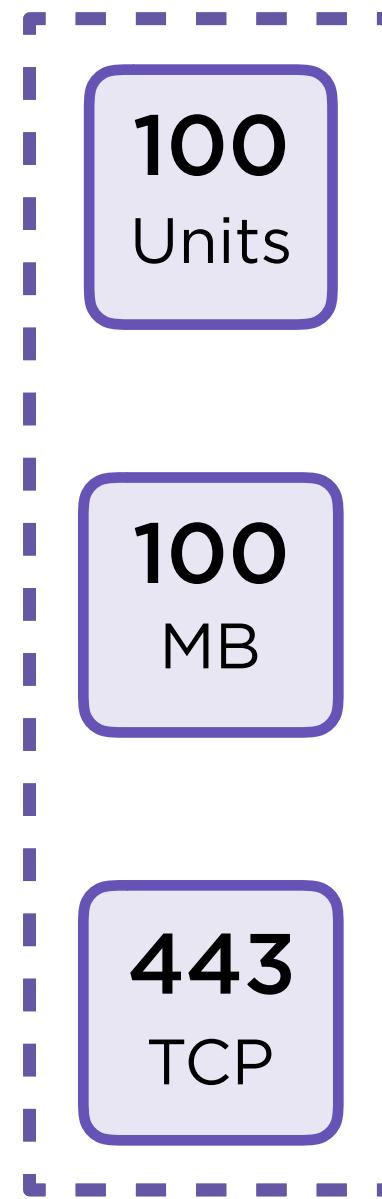
**Worst Case Memory
of 300MB**



**Network Port
TCP/80**

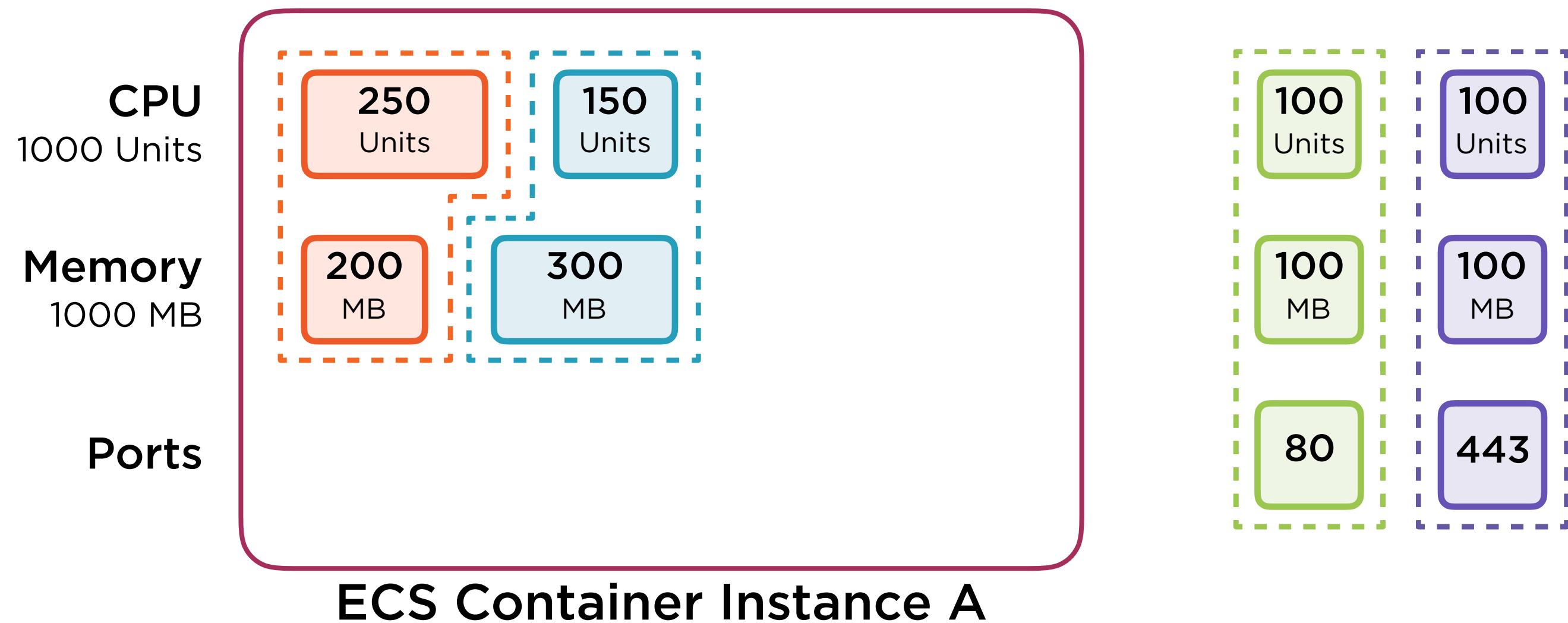


**Network Port
TCP/443**

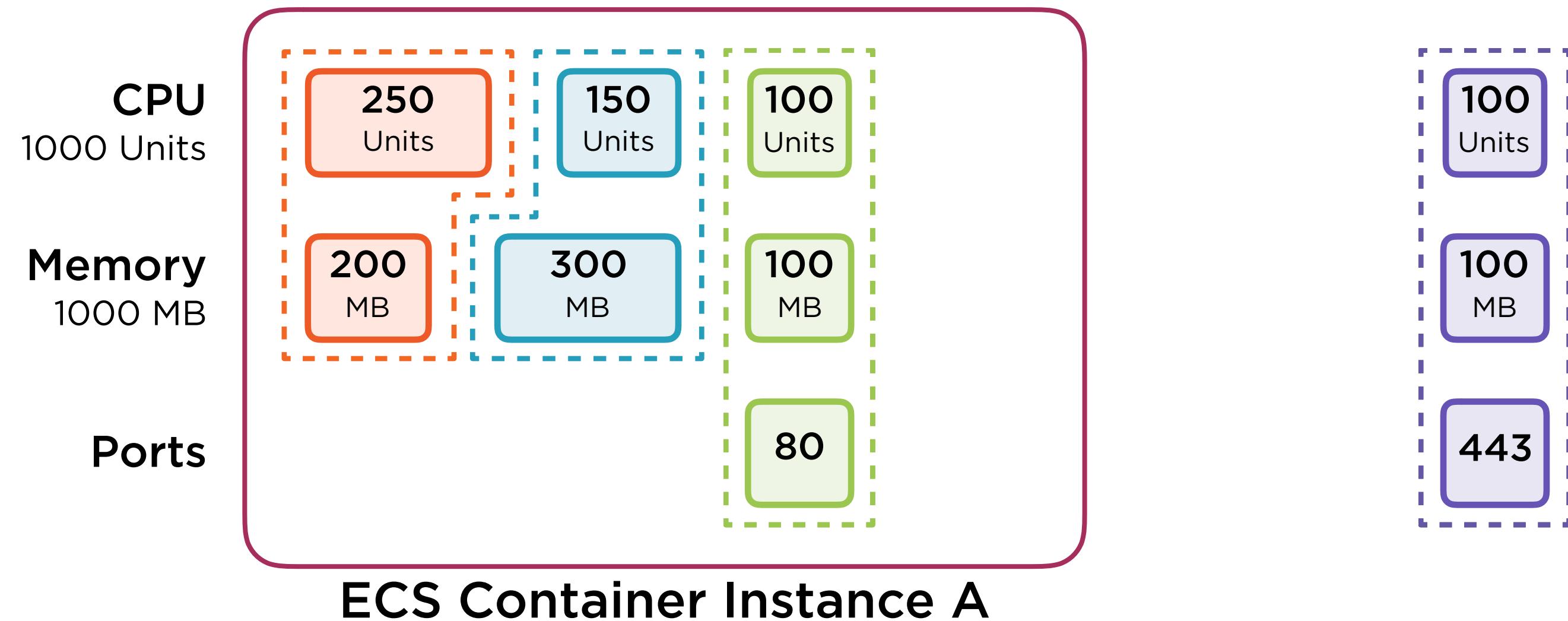


Scaling out ECS Clusters

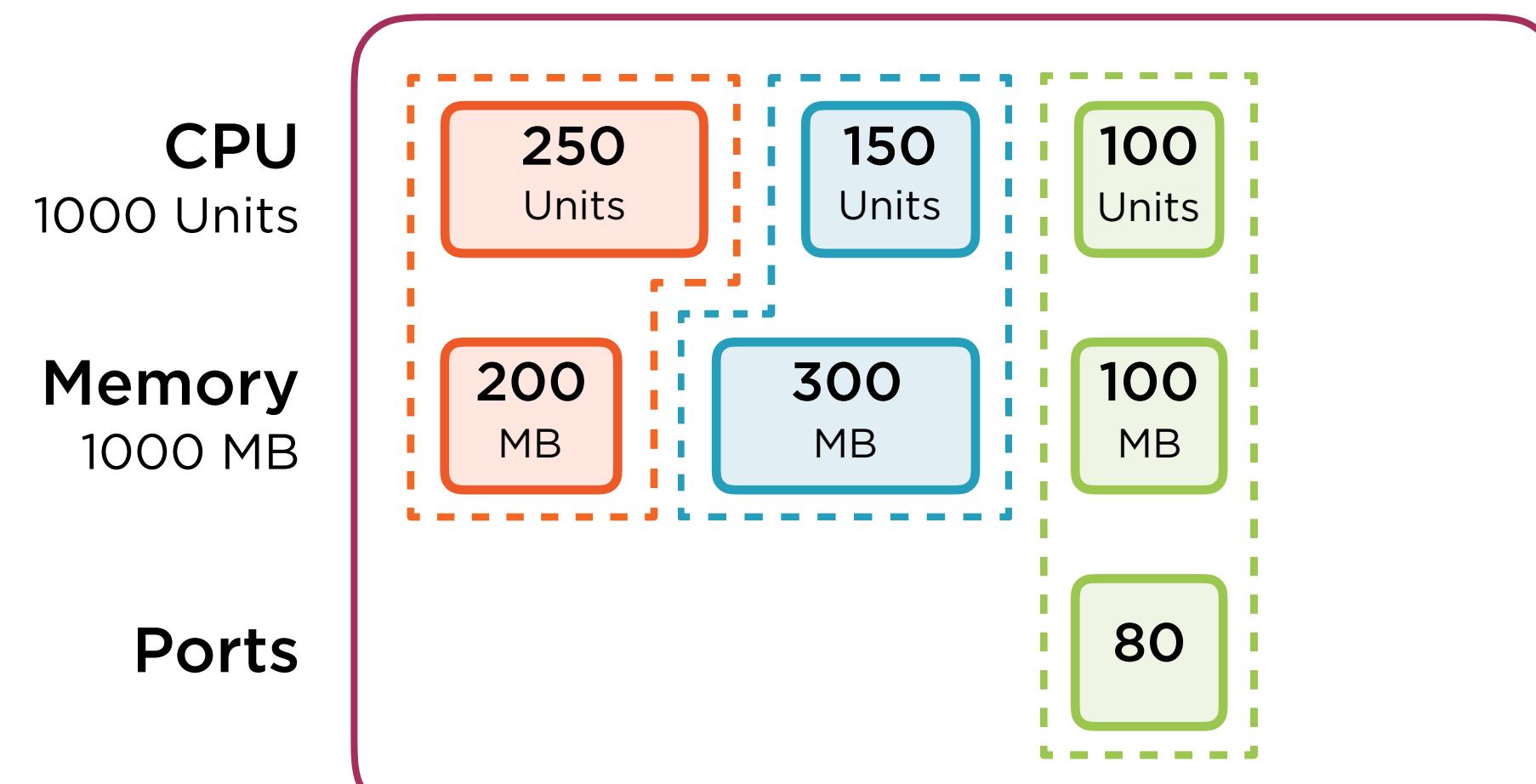
Resource	Current Usage	Remaining Resources	Container Capacity
CPU	400	600	2 x Orange
Memory	500	500	1 x Blue
Port 80	0	1	1 x Green
Port 443	0	1	1 x Purple
Minimum Capacity			1



Resource	Current Usage	Remaining Resources	Container Capacity
CPU	500	500	2 x Orange
Memory	600	400	1 x Blue
Port 80	1	0	0 x Green
Port 443	0	1	1 x Purple
Minimum Capacity			0

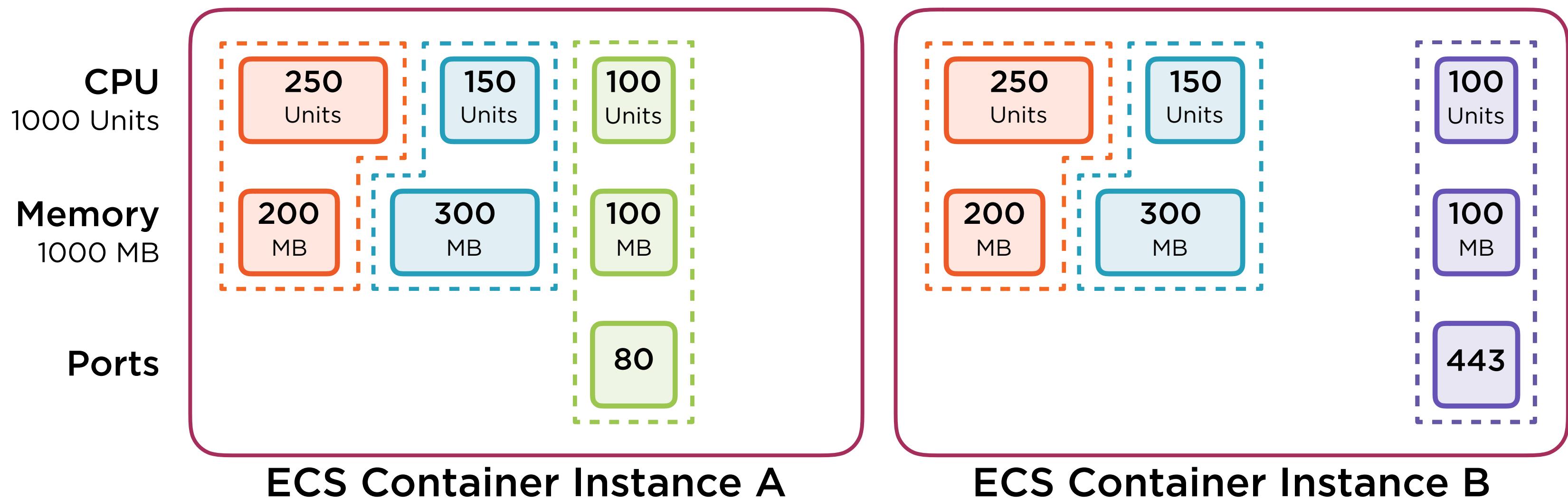


Resource	Current Usage	Remaining Resources	Container Capacity
CPU	500	500	2 x Orange
Memory	600	400	1 x Blue
Port 80	1	0	0 x Green
Port 443	0	1	1 x Purple
Minimum Capacity			0

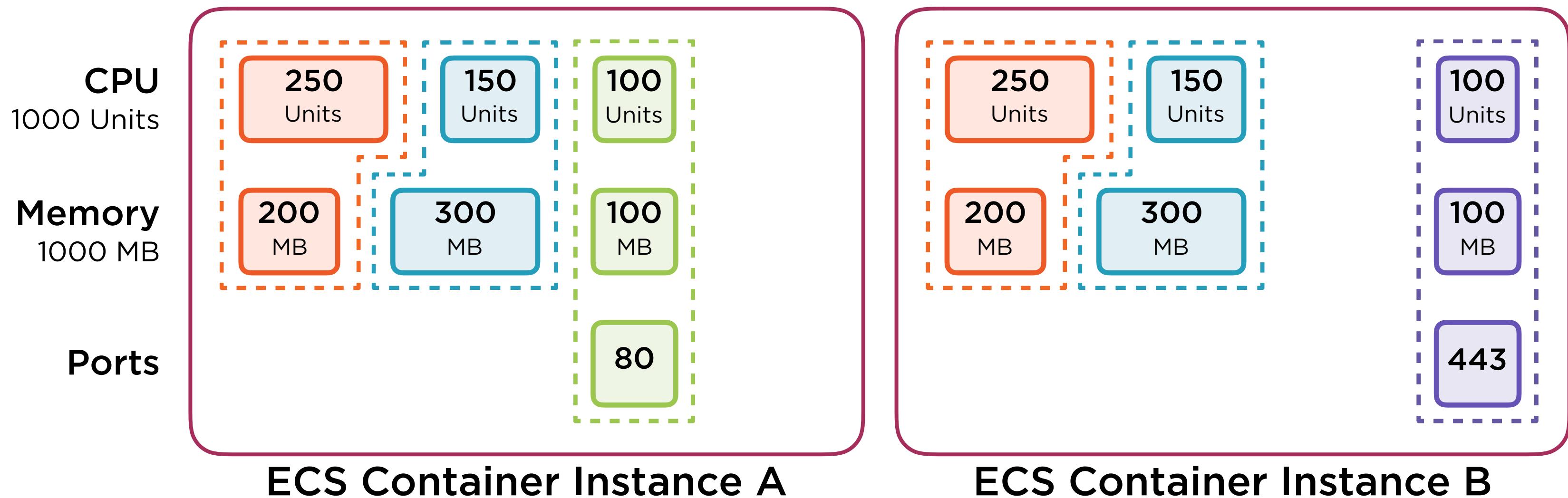


ECS Container Instance A

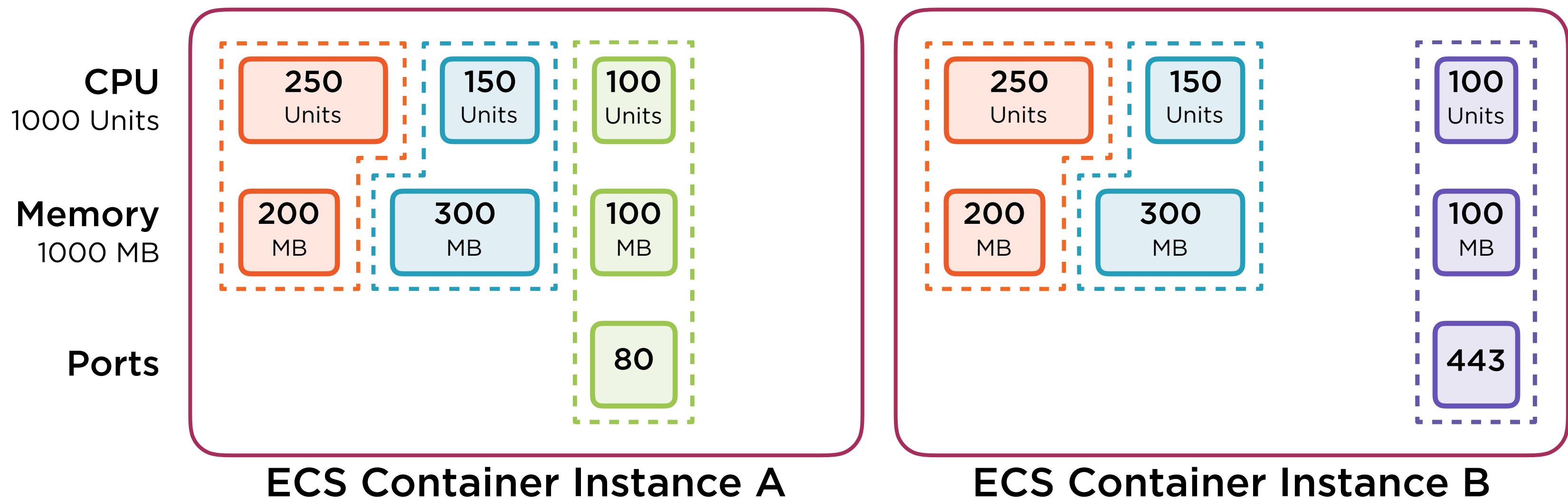
Resource	Current Usage	Remaining Resources	Container Capacity
CPU	500	500	2 x Orange
Memory	600	400	1 x Blue
Port 80	1	0	0 x Green
Port 443	0	1	1 x Purple
Minimum Capacity			0



Resource	Current Usage	Remaining Resources	Container Capacity
CPU	500	500	2 x Orange
Memory	600	400	1 x Blue
Port 80	1	0	0 x Green
Port 443	0	1	1 x Purple
Minimum Capacity			0

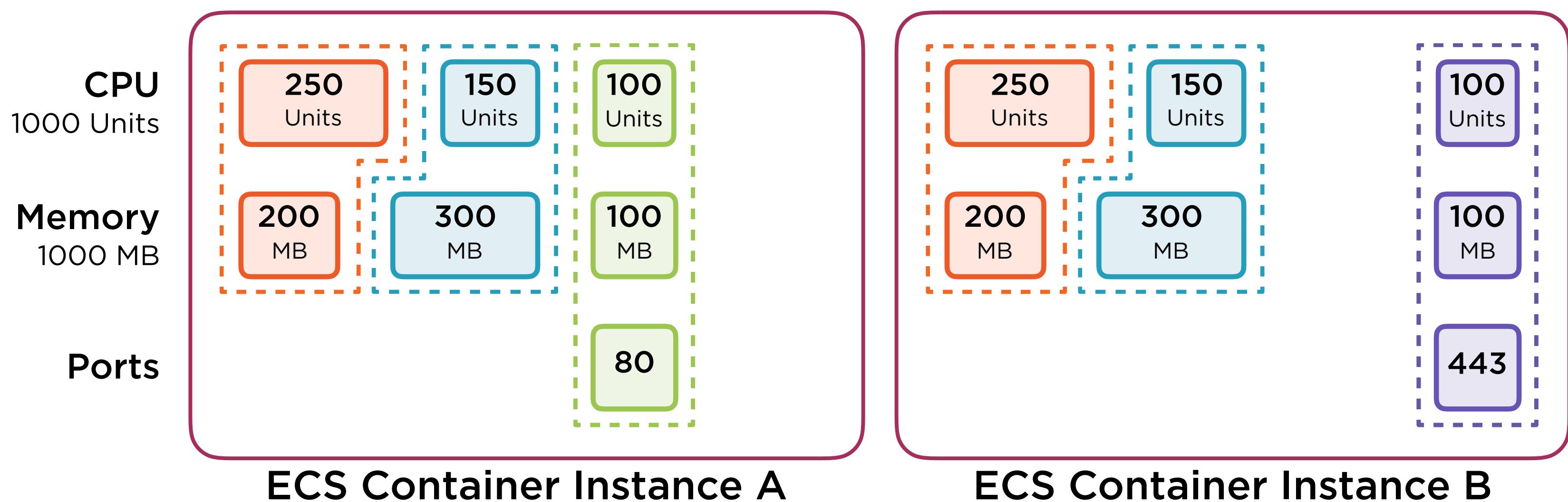


Resource	Current Usage		Remaining Resources	Container Capacity	
CPU	500 (A)	500 (B)	500 (A)	500 (B)	2(A) + 2(B) = 4
Memory	600 (A)	600 (B)	400 (A)	600 (B)	1(A) + 1(B) = 2
Port 80	1 (A)	0 (B)	0 (A)	1 (B)	0(A) + 1(B) = 1
Port 443	0 (A)	1 (B)	1 (A)	0 (B)	1(A) + 0(B) = 1
					Minimum Capacity
					1

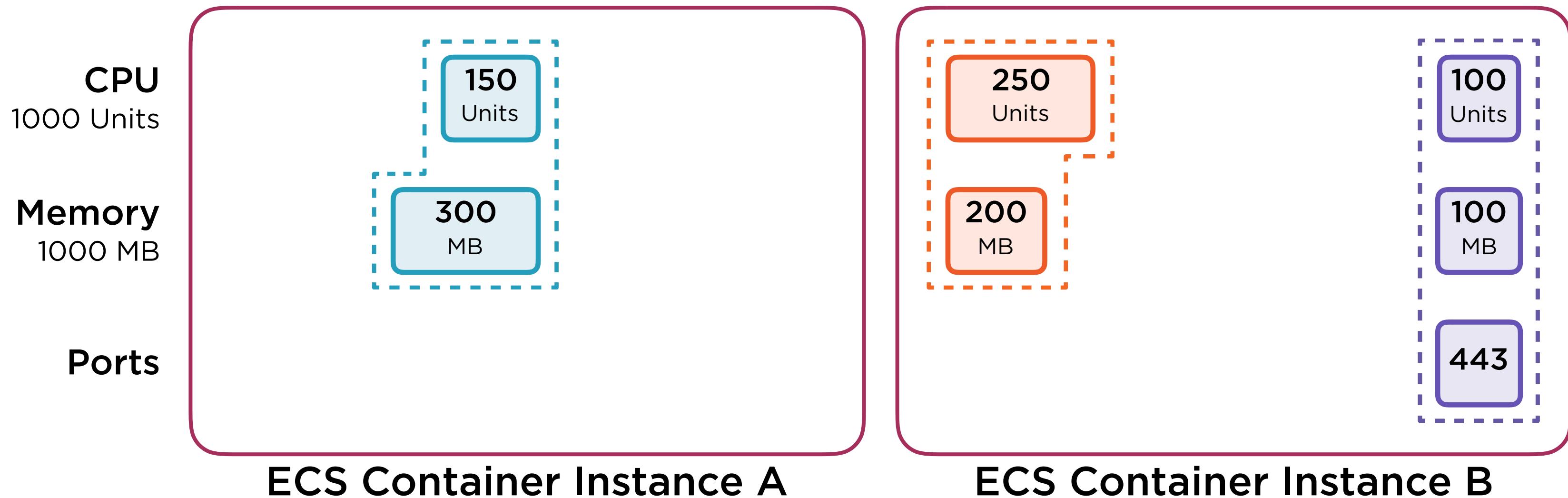


Scaling in ECS Clusters

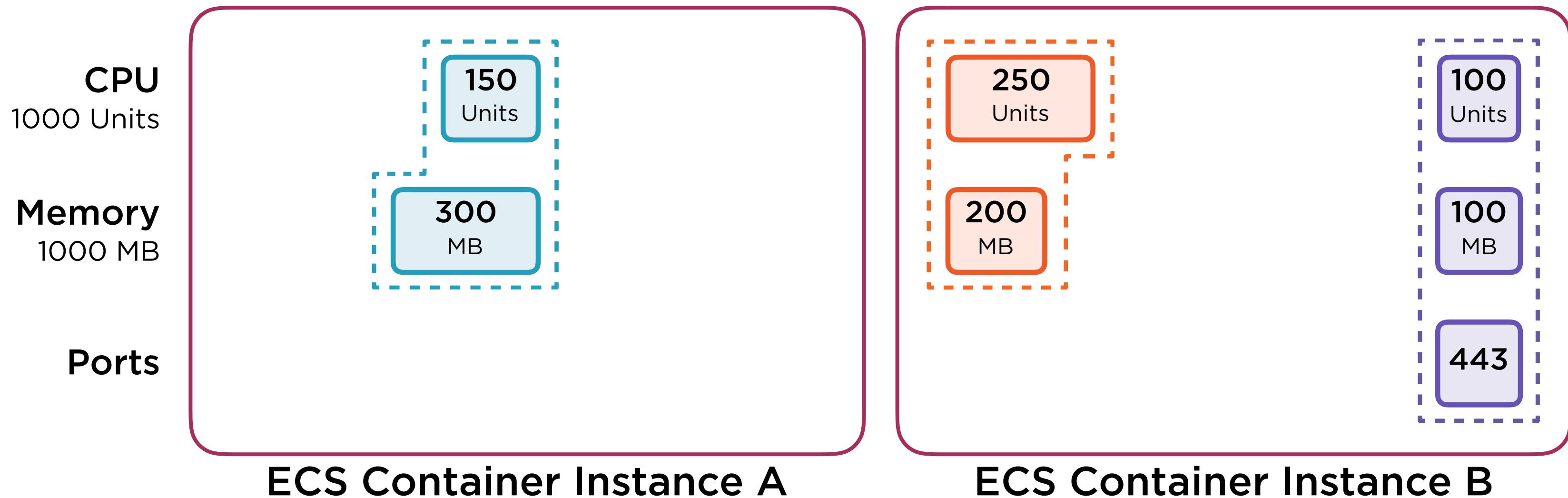
Resource	Current Usage		Remaining Resources		Container Capacity
CPU	500 (A)	500 (B)	500 (A)	500 (B)	$2(A) + 2(B) = 4$
Memory	600 (A)	600 (B)	400 (A)	600 (B)	$1(A) + 1(B) = 2$
Port 80	1 (A)	0 (B)	0 (A)	1 (B)	$0(A) + 1(B) = 1$
Port 443	0 (A)	1 (B)	1 (A)	0 (B)	$1(A) + 0(B) = 1$
					Minimum Capacity 1



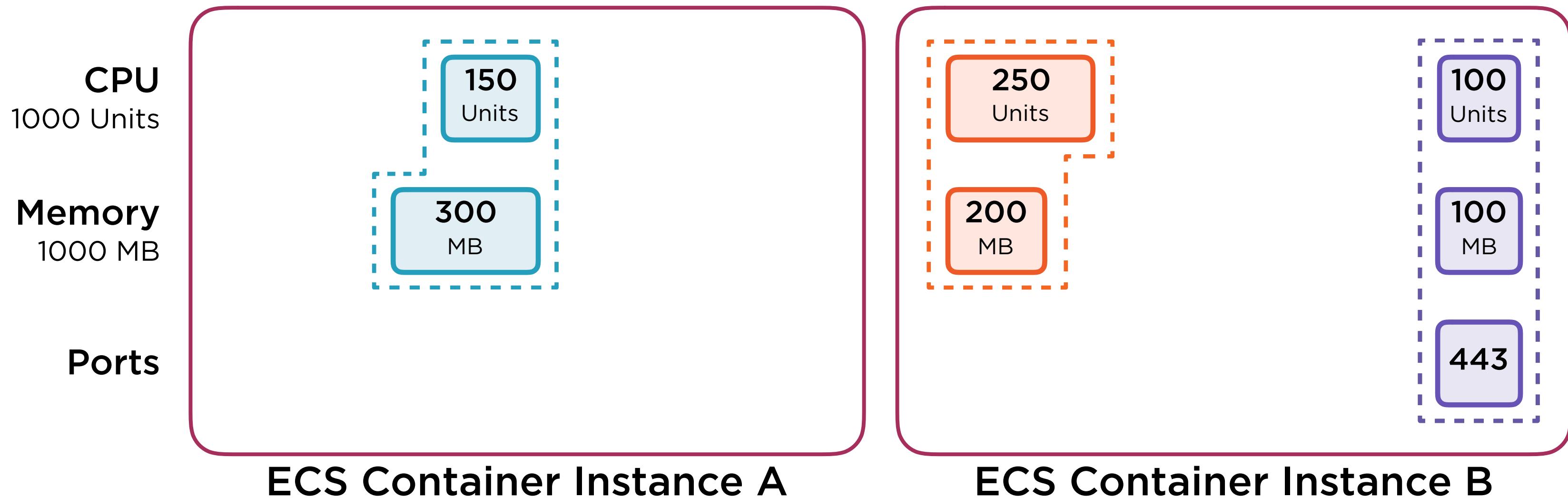
Resource	Current Usage	Remaining Resources	Container Capacity	Instance Capacity	Scale In Capacity	
CPU	150 (A)	350 (B)	850 (A) 650 (B)	$3(A) + 2(B) = 5$	3	$5 - 3 = 2$
Memory	300 (A)	300 (B)	700 (A) 700 (B)	$2(A) + 2(B) = 4$	3	$4 - 3 = 1$
Port 80	0 (A)	0 (B)	1 (A) 1 (B)	$1(A) + 1(B) = 2$	1	$2 - 1 = 1$
Port 443	0 (A)	1 (B)	1 (A) 0 (B)	$1(A) + 0(B) = 1$	1	$1 - 1 = 0$
					Minimum Capacity 0	



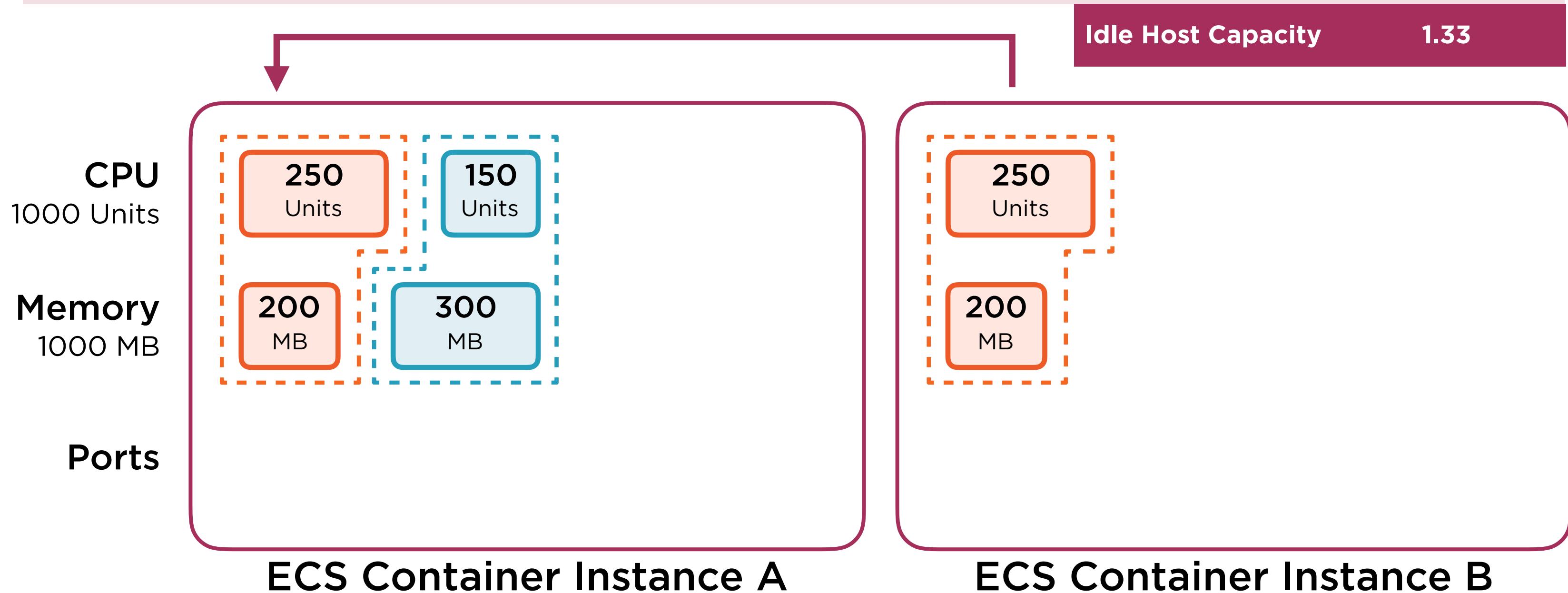
Resource	Current Usage	Remaining Resources	Container Capacity	Instance Capacity	Idle Capacity	
CPU	150 (A)	350 (B)	850 (A) 650 (B)	$3(A) + 2(B) = 5$	3	$5 / 3 = 1.67$
Memory	300 (A)	300 (B)	700 (A) 700 (B)	$2(A) + 2(B) = 4$	3	$4 / 3 = 1.33$
Port 80	0 (A)	0 (B)	1 (A) 1 (B)	$1(A) + 1(B) = 2$	1	$2 / 1 = 2.00$
Port 443	0 (A)	1 (B)	1 (A) 0 (B)	$1(A) + 0(B) = 1$	1	$1 / 1 = 1.00$
					Idle Host Capacity 1.00	



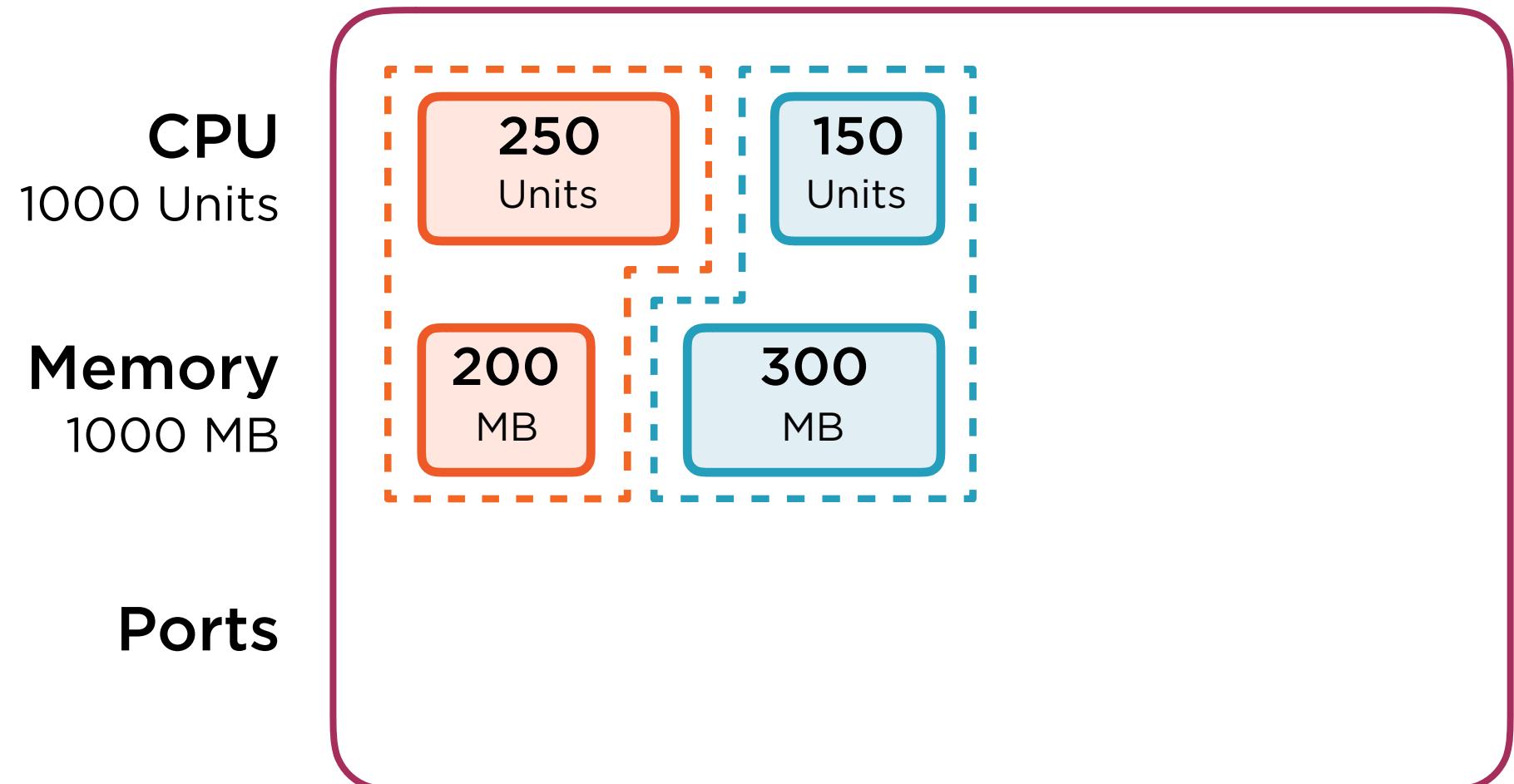
Resource	Current Usage	Remaining Resources	Container Capacity	Instance Capacity	Idle Capacity	
CPU	150 (A)	350 (B)	850 (A) 650 (B)	$3(A) + 2(B) = 5$	3	$5 / 3 = 1.67$
Memory	300 (A)	300 (B)	700 (A) 700 (B)	$2(A) + 2(B) = 4$	3	$4 / 3 = 1.33$
Port 80	0 (A)	0 (B)	1 (A) 1 (B)	$1(A) + 1(B) = 2$	1	$2 / 1 = 2.00$
Port 443	0 (A)	1 (B)	1 (A) 0 (B)	$1(A) + 0(B) = 1$	1	$1 / 1 = 1.00$
					Idle Host Capacity 1.00	



Resource	Current Usage	Remaining Resources	Container Capacity	Instance Capacity	Idle Capacity	
CPU	150 (A)	250 (B)	850 (A) 750 (B)	$3(A) + 3(B) = 6$	3	$6 / 3 = 2.00$
Memory	300 (A)	200 (B)	700 (A) 800 (B)	$2(A) + 2(B) = 4$	3	$4 / 3 = 1.33$
Port 80	0 (A)	0 (B)	1 (A) 1 (B)	$1(A) + 1(B) = 2$	1	$2 / 1 = 2.00$
Port 443	0 (A)	0 (B)	1 (A) 1 (B)	$1(A) + 1(B) = 2$	1	$2 / 1 = 2.00$



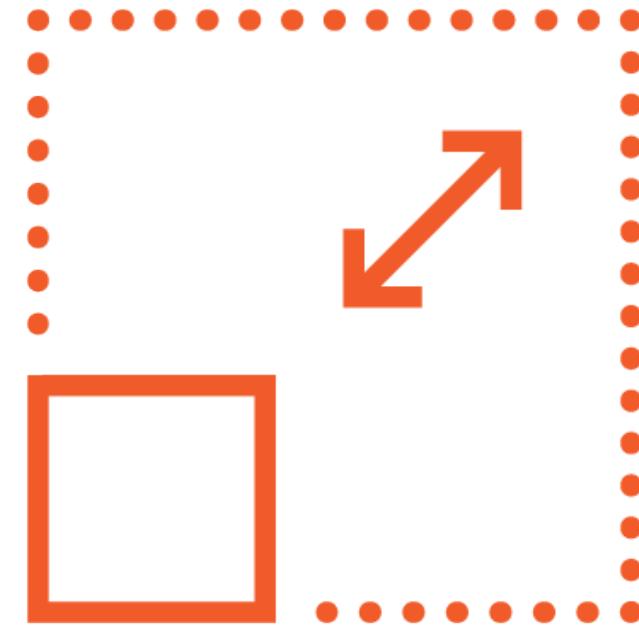
Resource	Current Usage	Remaining Resources		Container Capacity	Instance Capacity	Idle Capacity	
CPU	150 (A)	250 (B)	850 (A)	750 (B)	$3(A) + 3(B) = 6$	3	$6 / 3 = 2.00$
Memory	300 (A)	200 (B)	700 (A)	800 (B)	$2(A) + 2(B) = 4$	3	$4 / 3 = 1.33$
Port 80	0 (A)	0 (B)	1 (A)	1 (B)	$1(A) + 1(B) = 2$	1	$2 / 1 = 2.00$
Port 443	0 (A)	0 (B)	1 (A)	1 (B)	$1(A) + 1(B) = 2$	1	$2 / 1 = 2.00$
						Idle Host Capacity	
						1.33	



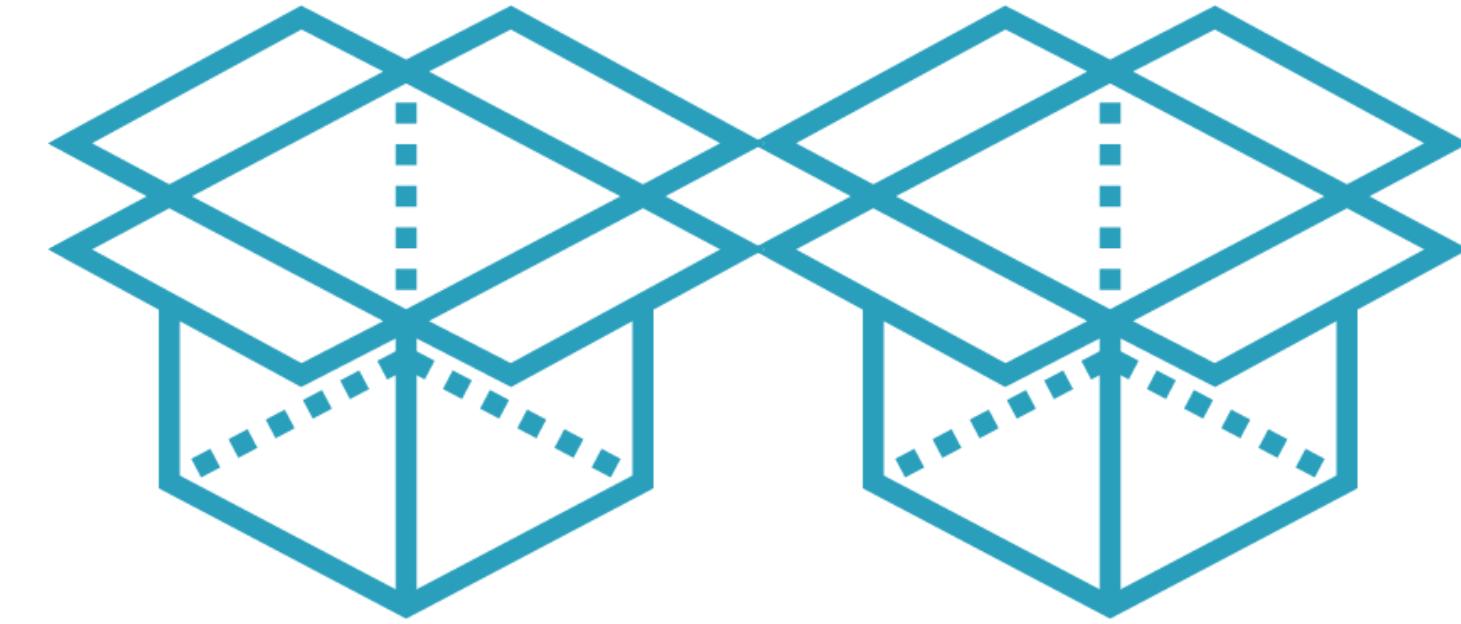
ECS Container Instance A

Auto Scaling Solution Overview

ECS Capacity Summary

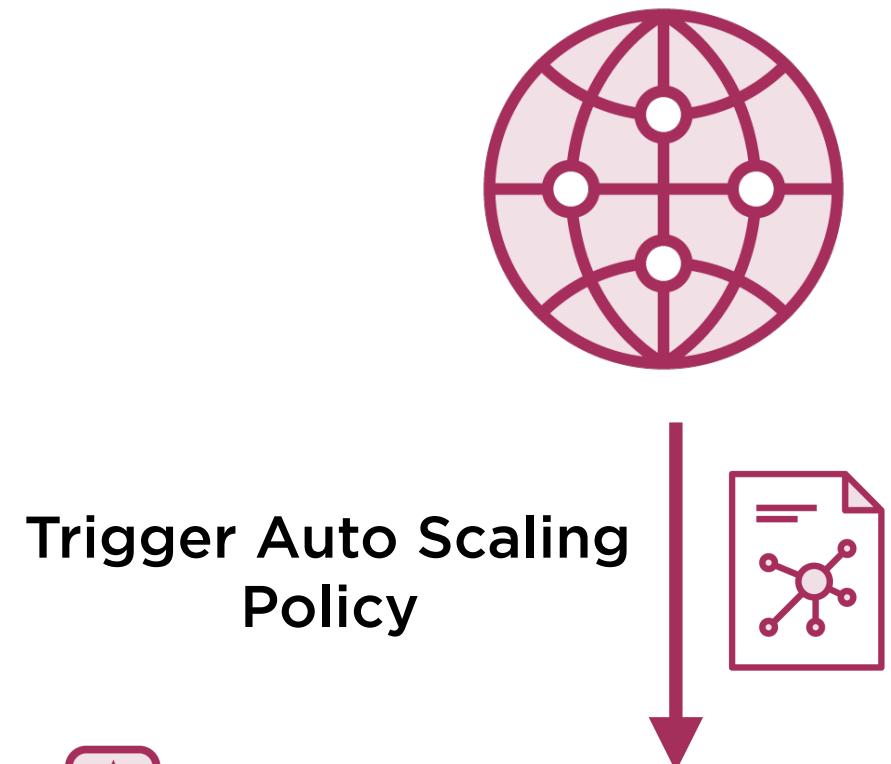


Container Capacity
Scale out when < 1

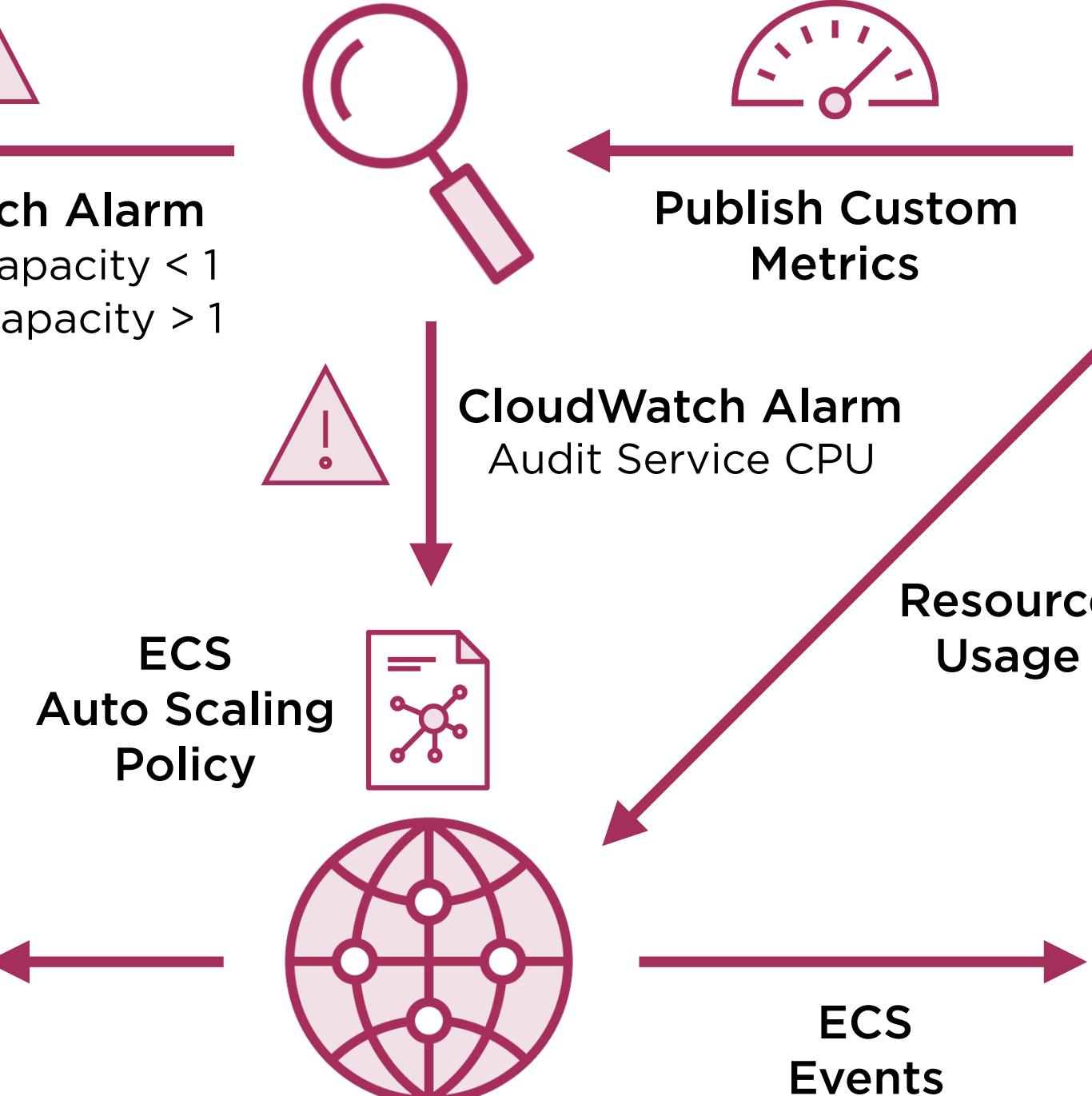


Idle Host Capacity
Scale in when > 1

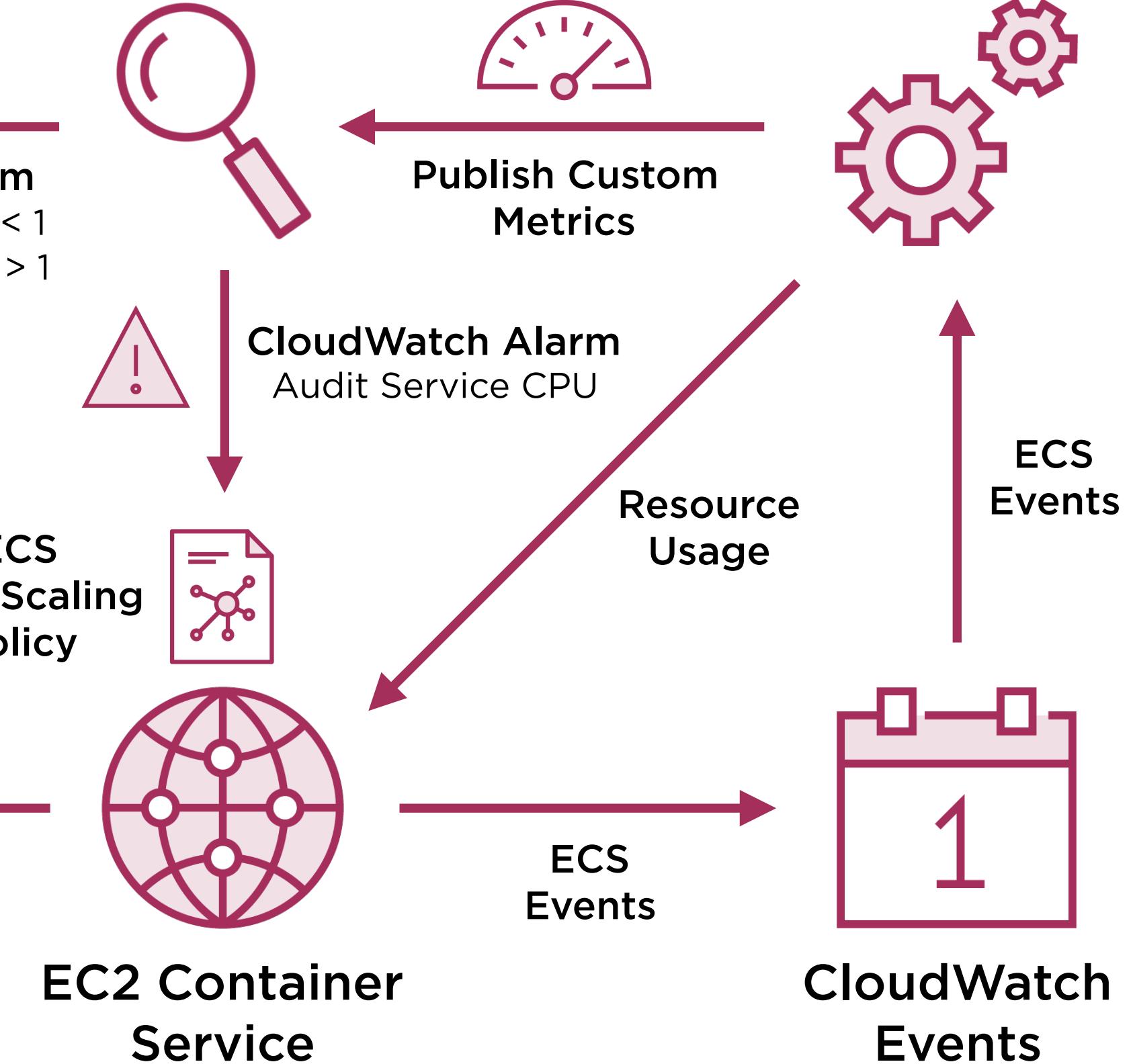
EC2 Auto Scaling Service



CloudWatch Service

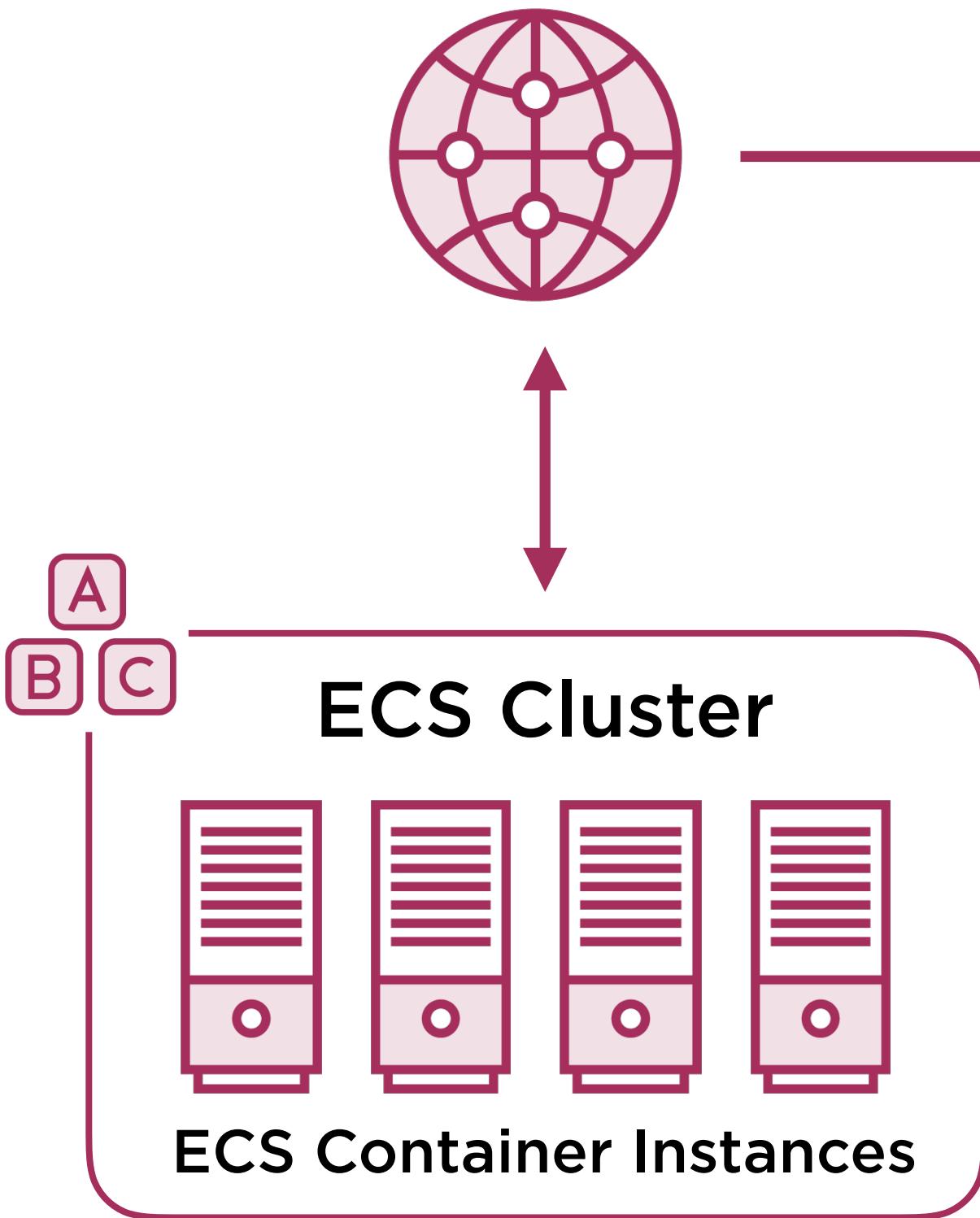


ECS Capacity Lambda Function



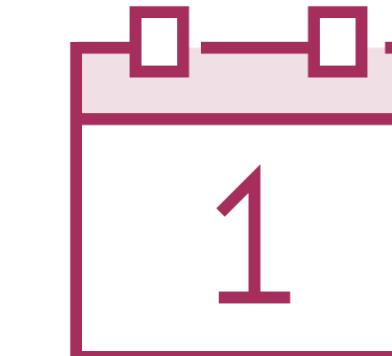
Understanding ECS Events

EC2 Container Service

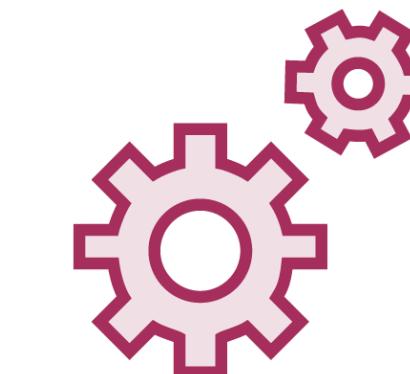


ECS Container Instance
State Change Events
ECS Task
State Change Events

CloudWatch Events



ECS Capacity Lambda Function



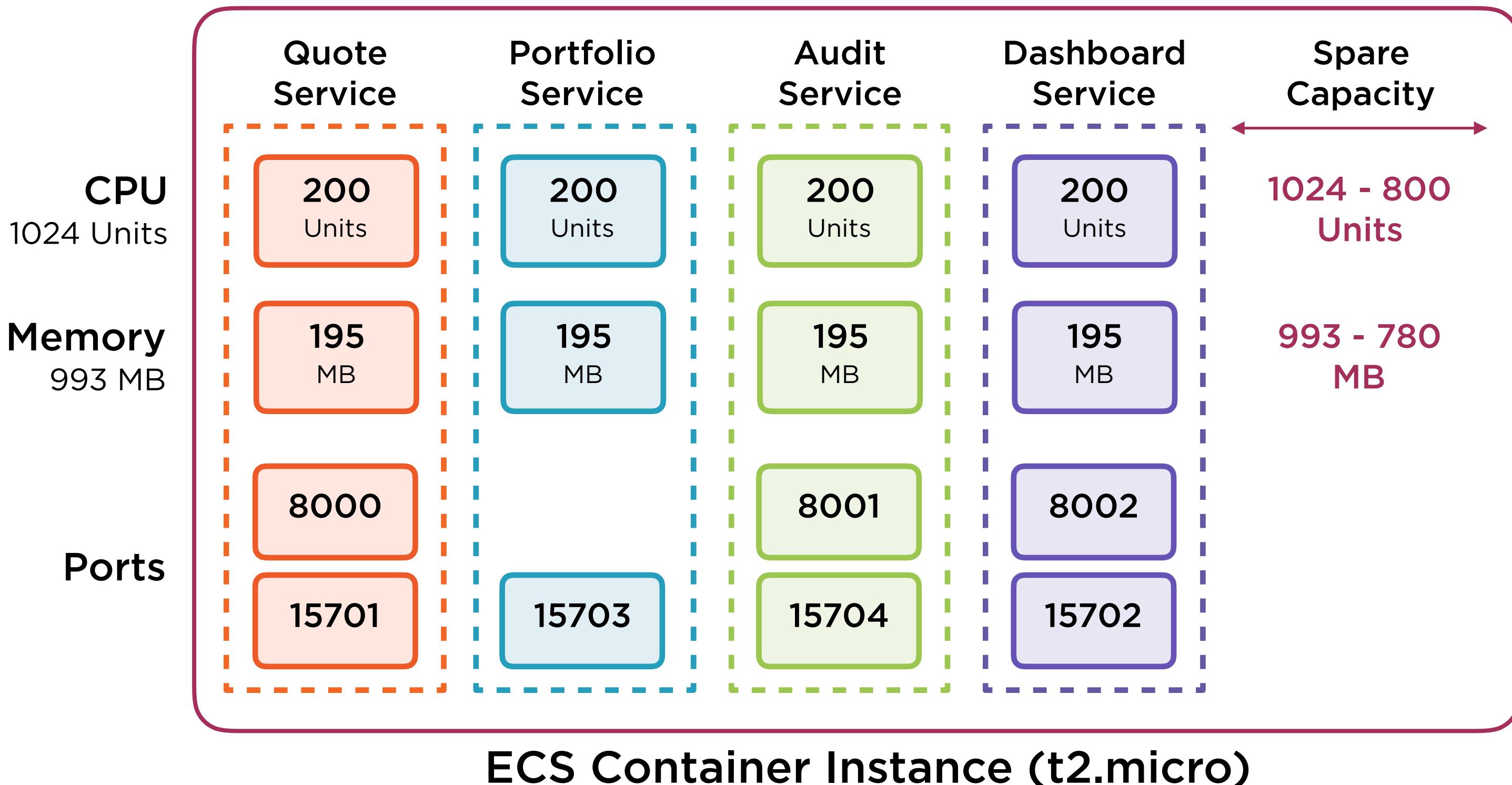
Check Current Capacity

Creating a Capacity Manager Lambda Function

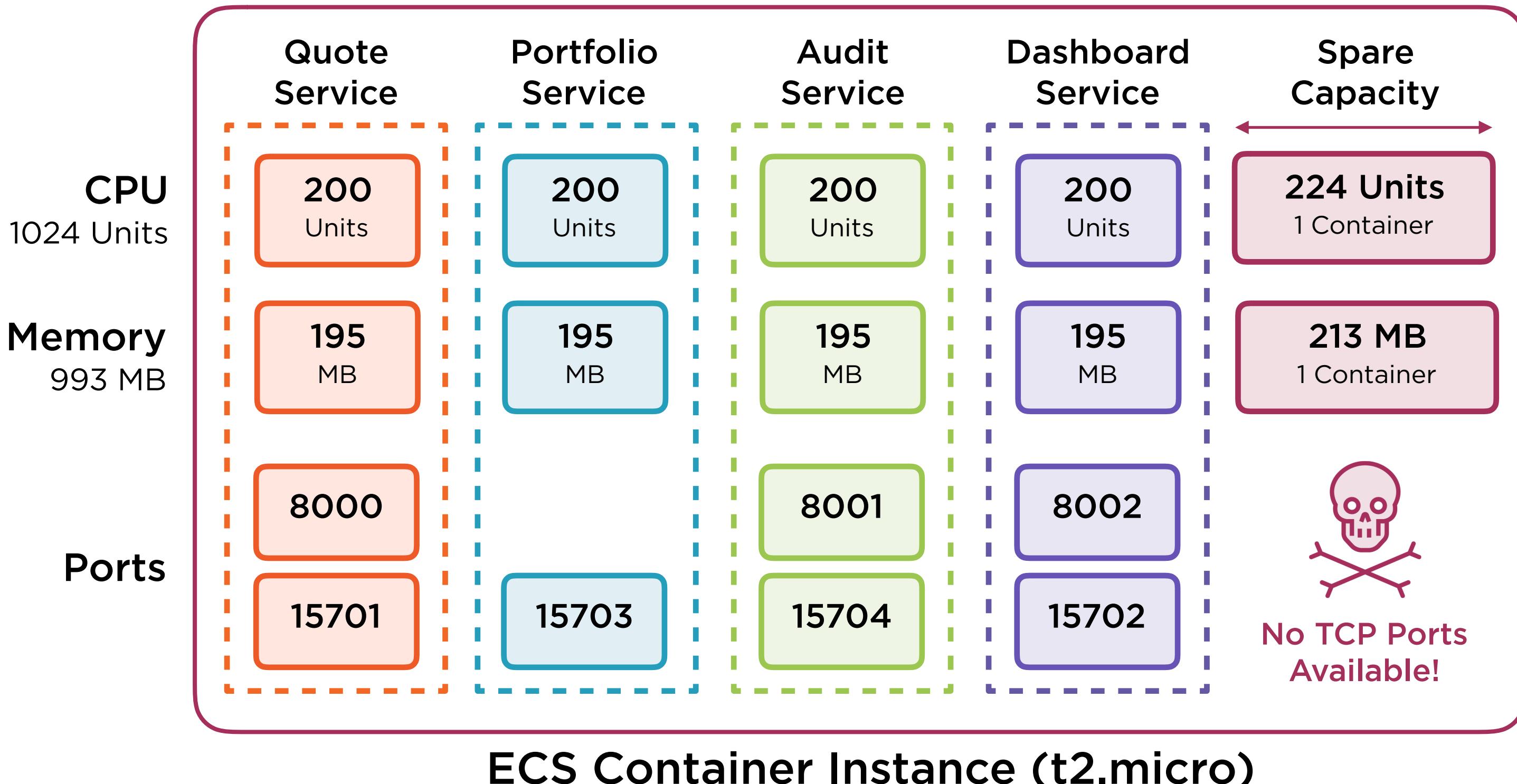
Calculating ECS Container Instance Capacity

Testing the ECS Capacity Manager

Current Microtrader Utilization



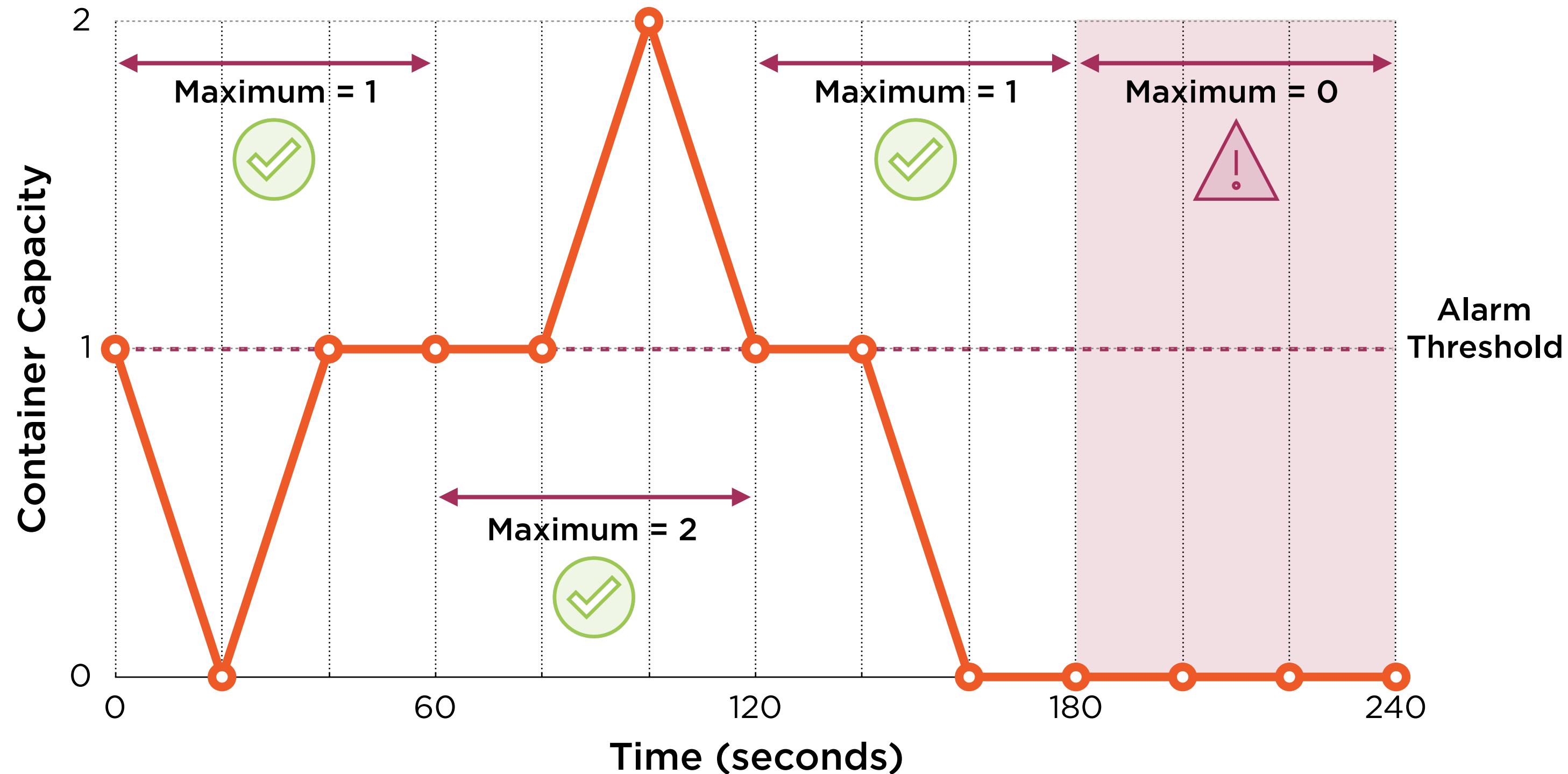
Current Microtrader Utilization



Publishing CloudWatch Custom Metrics

Creating CloudWatch Alarms

Choosing CloudWatch Statistics



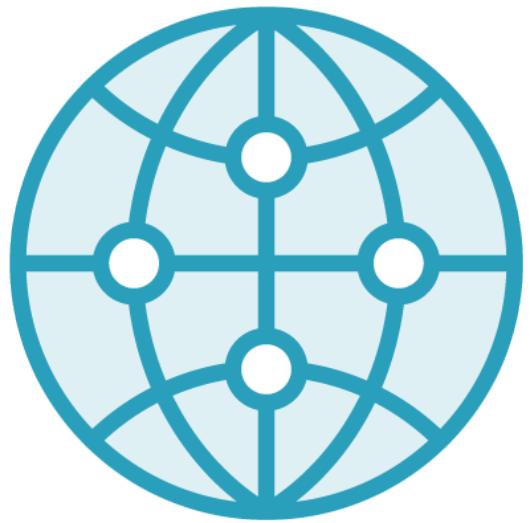
Creating EC2 Auto Scaling Policies

Managing ECS Capacity using CloudFormation

EC2 Auto Scaling using CloudFormation

Configuring ECS Autoscaling

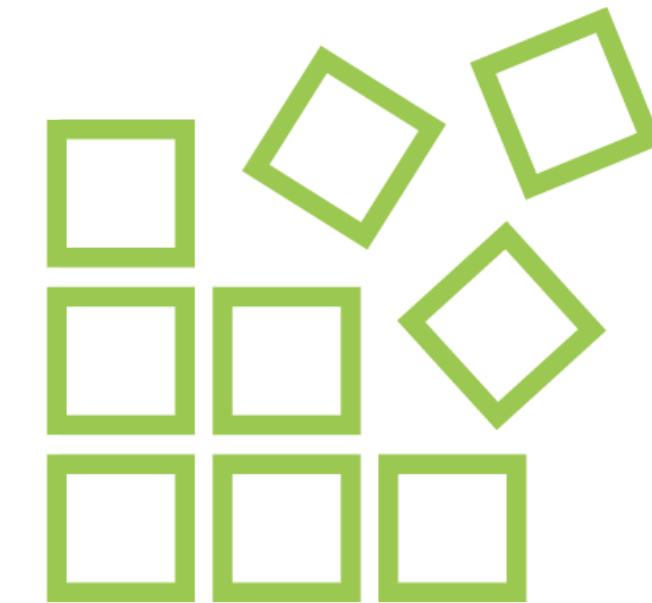
ECS Auto Scaling



Applied at an
ECS Service level

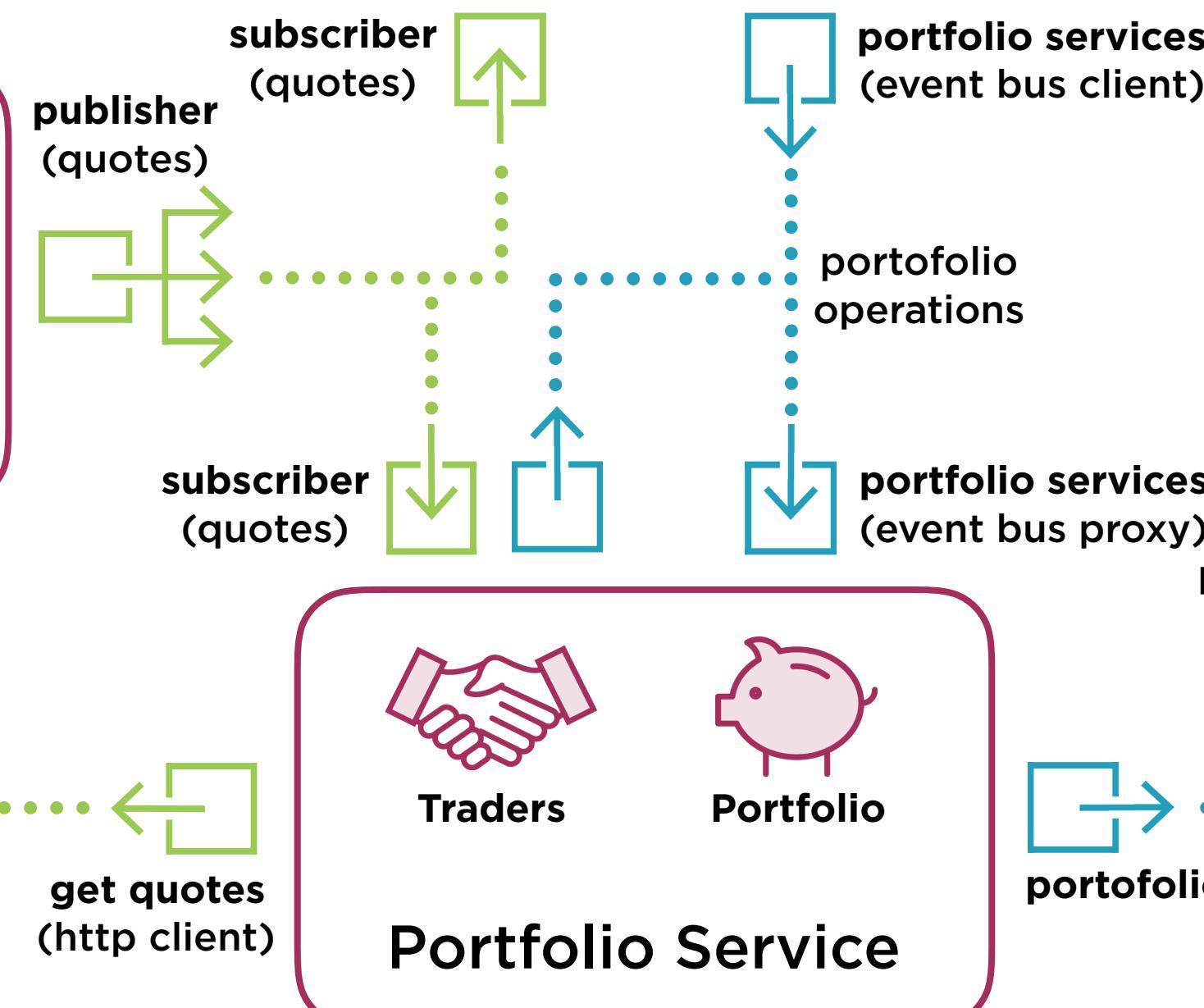
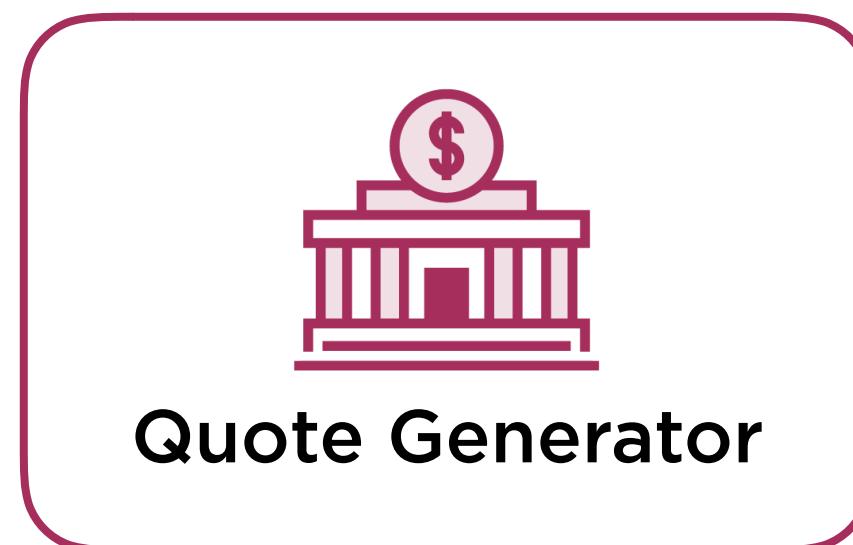


Choose Application
Specific Metrics



Add/Remove
Containers

MARKET_PERIOD
Default = 3000ms



ECS Auto Scaling using CloudFormation

Summary

Auto Scaling ECS Applications

- ECS Auto Scaling
 - Application specific
 - Adds/subtracts containers for your ECS services
- EC2 Auto Scaling
 - Needs to ensure we have spare capacity for additional containers
 - Needs to scale in when we have too much container capacity

Summary

Calculating ECS Capacity

- Consider the worst case capacity for each type of resource
 - Memory
 - CPU
 - Network Ports
- Take the lowest of each resource capacity calculation
- Container Capacity is used to scale out when < 1
- Idle Host Capacity is used to scale in when > 1

Summary

ECS Auto Scaling Solution

- ECS Capacity Lambda function
 - Triggered by ECS container instance state change events
- Custom CloudWatch metrics
 - ContainerCapacity
 - IdleHostCapacity
- CloudWatch alarms trigger EC2 Auto Scaling policies
- ECS Auto Scaling policies triggered by application-specific metrics and increase or decrease container counts

Avoid static network port mappings!