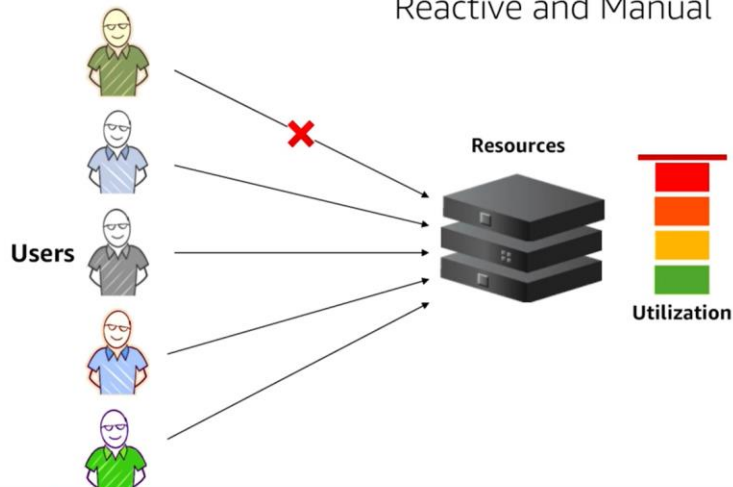


# The Challenge of High Availability

aws training and certification

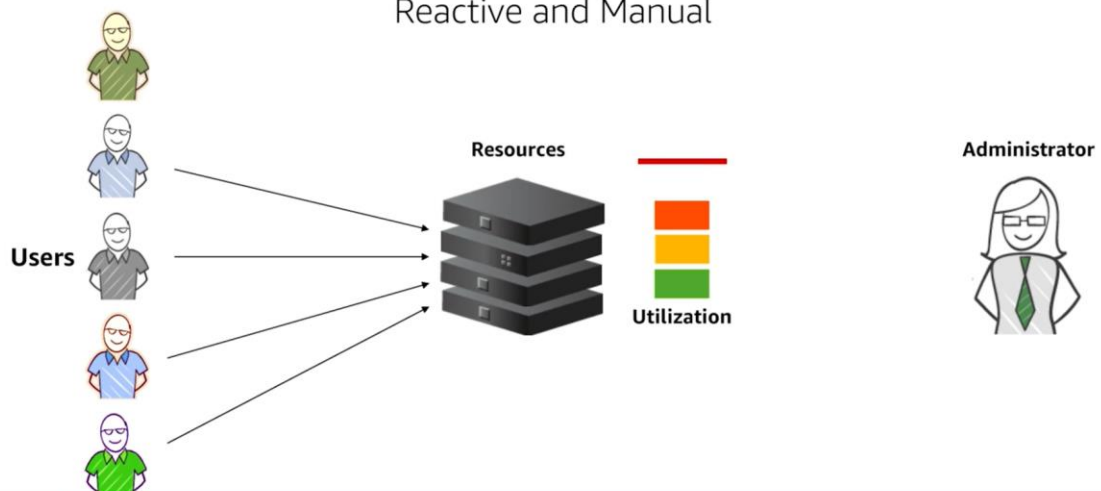
Reactive and Manual



# The Challenge of High Availability

aws training and certification

Reactive and Manual



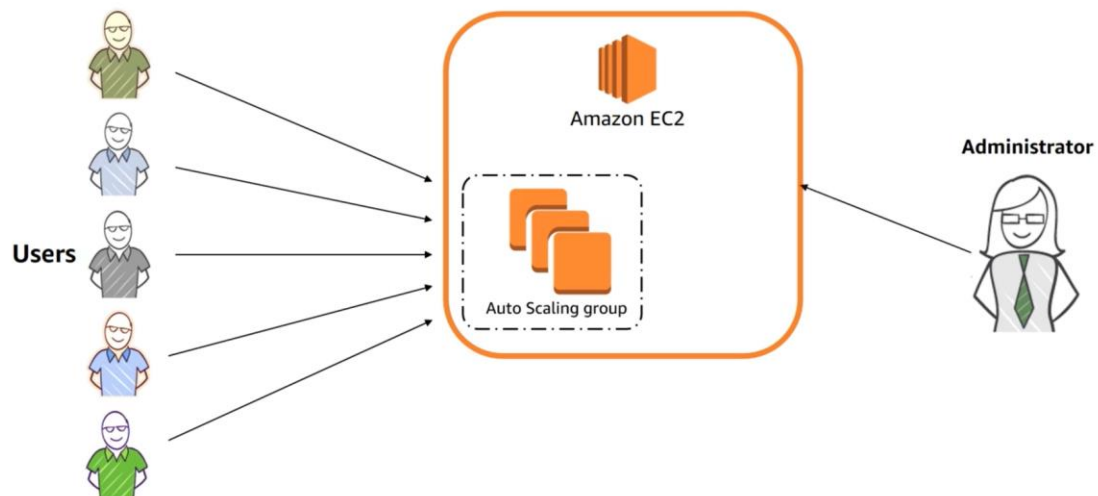
# Automatic Scaling on AWS

aws training and certification



## Auto Scaling AWS Services

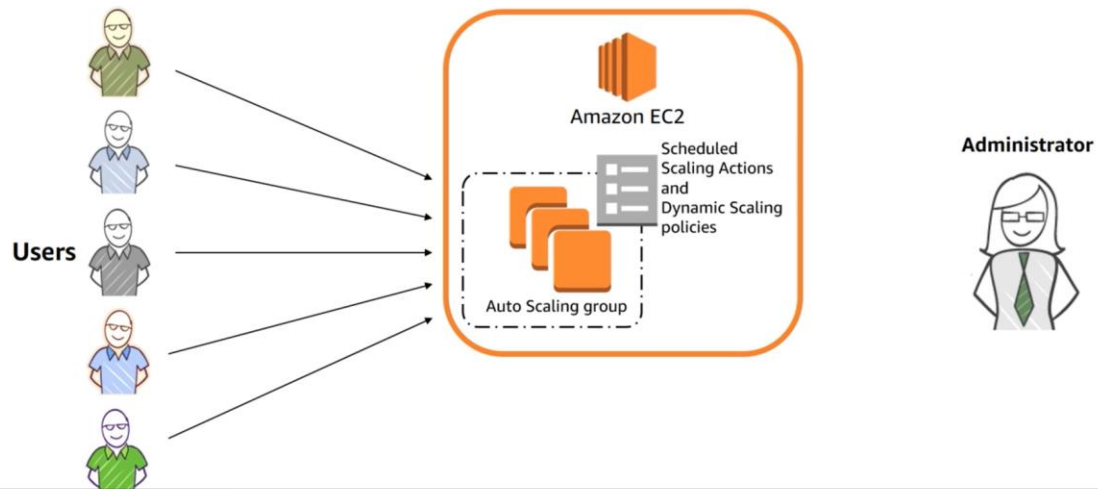
aws training and certification



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

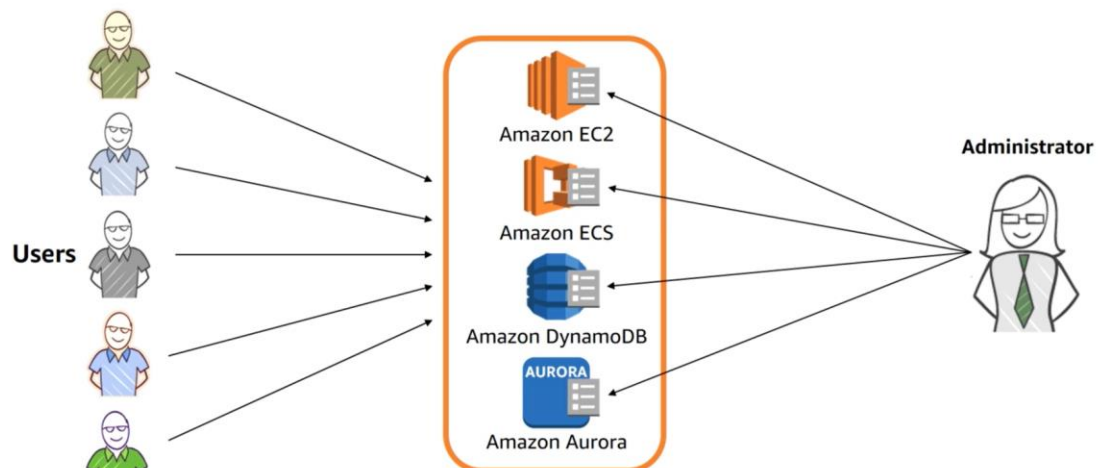
# Auto Scaling AWS Services

aws training and certification

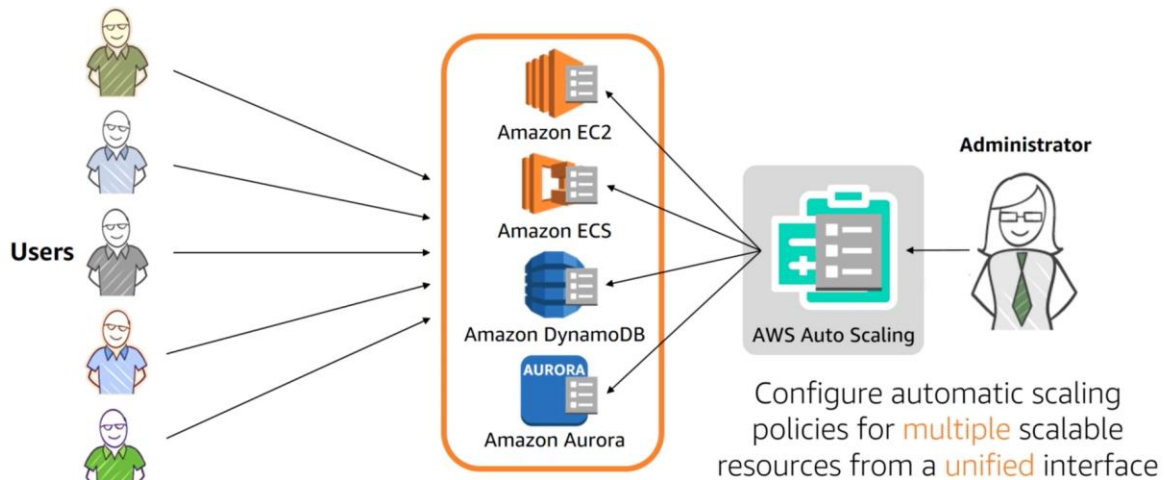


# Automatically Scaling AWS Services

aws training and certification

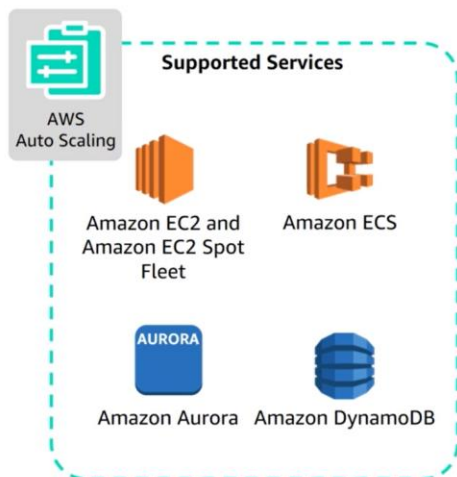


# Introducing AWS Auto Scaling



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

## Overview



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

### AWS Auto Scaling:

- 📦 Leverages existing Amazon EC2 Auto Scaling and Application Auto Scaling services
- 📦 Enables you to select an application as defined by AWS CloudFormation stack or in AWS Elastic Beanstalk

# Configuring AWS Auto Scaling



## STEP 1



**Select  
Application**

## STEP 2



**Scan for  
Scalable Services**

## STEP 3



**Configure  
Scaling Plan**

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

# Scaling Plan Overview



## SCALING STRATEGY

- ☐ AVAILABILITY
- ☒ COST
- ☐ BOTH
- ☐ CUSTOM



**AMAZON EC2  
AUTO SCALING  
GROUP**



SCALING PLAN  
FOR WEB  
APPLICATION 1C

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

# AWS Auto Scaling Benefits



**SCALE EFFICIENTLY  
WITH **ONE INTERFACE****



**SCALE INTELLIGENTLY  
WITH **SCALING PLANS****



**MAINTAIN  
PERFORMANCE WITH  
**CONTINUAL  
MONITORING****



**FREE SERVICE:  
PAY ONLY  
FOR RESOURCES  
YOU need**

# Unified Scaling

aws training and certification

	 Amazon EC2	 Amazon EC2	 Amazon ECS	 Amazon DynamoDB	 Amazon Aurora
<b>Scalable Target</b>	Auto Scaling Group	Spot Fleets	Service	Table or GSI	Cluster
<b>Scalable Dimension</b>	Amazon EC2 instances	Amazon EC2 Spot instances	Tasks	Provisioned Capacity	Aurora Replicas

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

# Predictable Scaling

aws training and certification

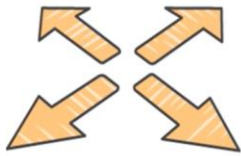


8:00 / 17:13





# Built-In Scaling Strategies



**Optimize for  
Availability**



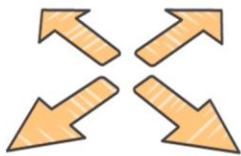
**Balance Availability  
and Cost**



**Optimize for  
Cost**

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

# Built-In Scaling Strategies



**Optimize for  
Availability**

Low resource  
utilization target



**Balance Availability  
and Cost**

Moderate resource  
utilization target



**Optimize for  
Cost**

High resource  
utilization target

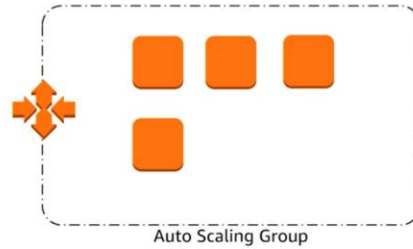


# Target Tracking Scaling Policies

aws training and certification



**METRIC: CPU UTILIZATION**  
**TARGET VALUE: 40%**



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

AWS Auto Scaling

Services Resource Groups

Step 1 Select application

Step 2 Configure scaling plan

Step 3 Review and create

AWS Auto Scaling > Scaling plans > Create scaling plan

### Select application

Select the AWS CloudFormation stack for your application. You cannot use a stack that is in a failed or deleted state.

CloudFormation stacks		
Name	Status	Creation Time
<input type="radio"/> WebApp	UPDATE_COMPLETE	2018-01-23 13:41:28 UTC-0800

Cancel Next

Feedback English (US)

© 2008 - 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved. Privacy Policy Terms of Use

AWS Auto Scaling

Services Resource Groups

Step 1: Select application  
Step 2: Configure scaling plan  
Step 3: Review and create

## Configure scaling plan

Choose how to optimize the scalable resources in your application.

### Scaling plan details

Resources  
WebApp has 3 scalable resources

Name  
  
Must be 1-64 characters long and should not contain the pipe "|" character.

### Auto Scaling groups (1)

☒ Include in scaling plan

Strategy

- ☒ Optimize for availability  
Use a low resource utilization target to provide optimal availability and ensure capacity to absorb spikes in demand.
- ☐ Balance availability and cost  
Use a moderate resource utilization target to provide high availability and reduce costs.
- ☐ Optimize for cost  
Use a higher resource utilization target to ensure lower costs.
- ☐ Custom  
Set your own values

All scalable targets (percent)

100  
80.0  
60.0  
40.0  
20.0  
0

40%

Feedback English (US)

© 2008 - 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved. Privacy Policy Terms of Use

AWS Auto Scaling

Services Resource Groups

Step 1: Select application  
Step 2: Configure scaling plan  
Step 3: Review and create

## Configure scaling plan

Create scaling policies to keep the average CPU utilization of your Auto Scaling groups at 40% and scale from 1 to 10 instances.

### DynamoDB tables and indexes (2)

☒ Include in scaling plan

Strategy

- ☒ Optimize for availability  
Use a low resource utilization target to provide optimal availability and ensure capacity to absorb spikes in demand.
- ☐ Balance availability and cost  
Use a moderate resource utilization target to provide high availability and reduce costs.
- ☐ Optimize for cost  
Use a higher resource utilization target to ensure lower costs.
- ☐ Custom  
Set your own values

All scalable targets (percent)

All scalable targets (percent)

RCU: DataTable

WCU: DataTable

RCU: UserTable

WCU: UserTable

20.0  
0

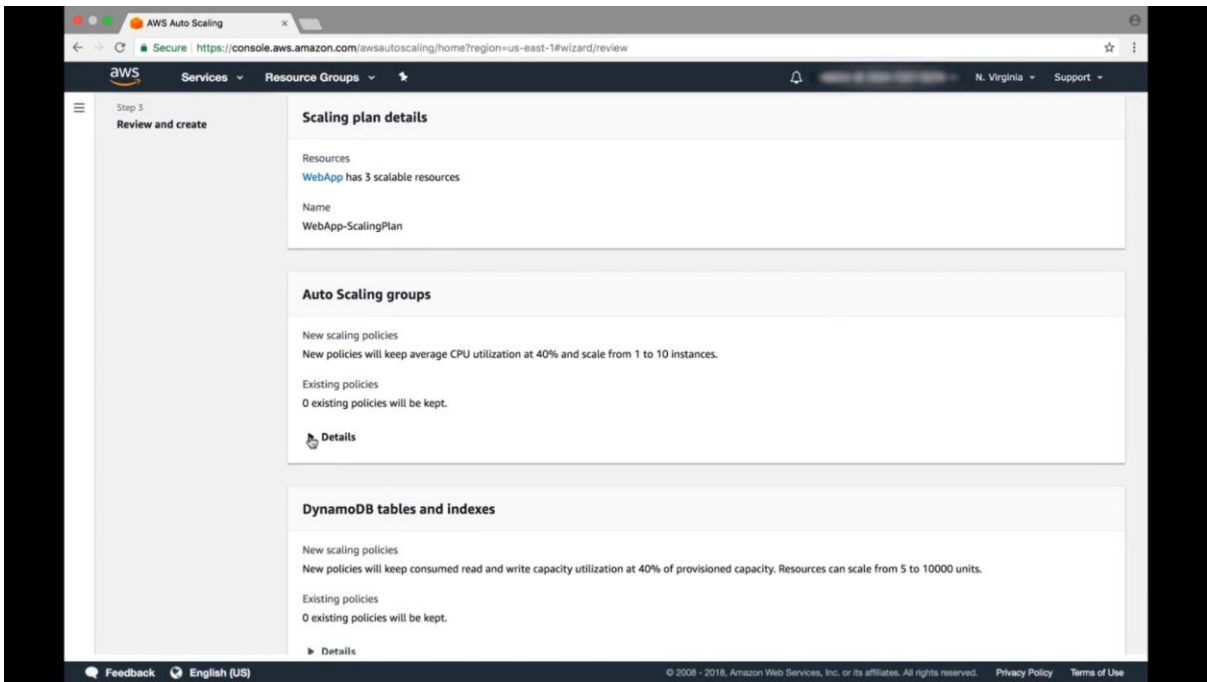
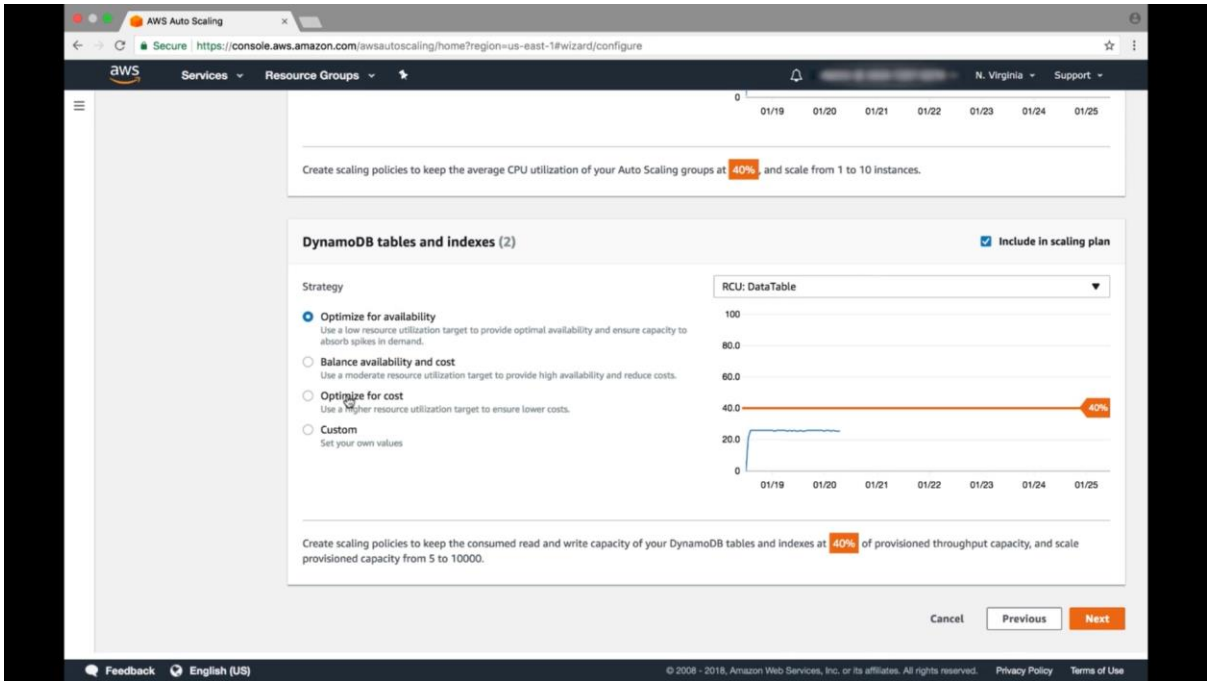
01/19 01/20 01/21 01/22 01/23 01/24 01/25

Create scaling policies to keep the consumed read and write capacity of your DynamoDB tables and indexes at 40% of provisioned throughput capacity, and scale provisioned capacity from 5 to 10000.

Cancel Previous Next

Feedback English (US)

© 2008 - 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved. Privacy Policy Terms of Use



Services

Resource Groups

14:00 / 17:13

Auto Scaling groups

New scaling policies

New policies will keep average CPU utilization at 40% and scale from 1 to 10 instances.

Existing policies

0 existing policies will be kept.

Details

Scaling Policies (1)

Search by resource ID

< 1 >

Group Name	Existing policies	New policy	Action
autoScalingGroup/WebApp-WebServerGroup-13H20HQTSS19C	None	Target Metric: Optimize for availability, Target Value: 40%, Min: 1, Max: 10, Cooldown: Default	Apply new policy

DynamoDB tables and indexes

New scaling policies

New policies will keep consumed read and write capacity utilization at 40% of provisioned capacity. Resources can scale from 5 to 10000 units.

Existing policies

AWS Auto Scaling

Secure | <https://console.aws.amazon.com/autoscaling/home?region=us-east-1#wizard/review>

Services Resource Groups

N. Virginia Support

### DynamoDB tables and indexes

New scaling policies  
New policies will keep consumed read and write capacity utilization at 40% of provisioned capacity. Resources can scale from 5 to 10000 units.

Existing policies  
0 existing policies will be kept.

Details

#### Scaling Policies (4)

Search by resource ID

Resource ID	Scalable metric	Existing policies	New policy	Action
table/DataTable	Read capacity utilization	None	Target Metric: Optimize for availability, Target Value: 40%, Min: 5, Max: 10000, Cooldown: Default	<a href="#">Apply new policy</a>
table/DataTable	Write capacity utilization	None	Target Metric: Optimize for availability, Target Value: 40%, Min: 5, Max: 10000, Cooldown: Default	<a href="#">Apply new policy</a>
table/UserTable	Read capacity utilization	None	Target Metric: Optimize for availability, Target Value: 40%, Min: 5, Max: 10000, Cooldown: Default	<a href="#">Apply new policy</a>
table/UserTable	Write capacity utilization	None	Target Metric: Optimize for availability, Target Value: 40%, Min: 5, Max: 10000, Cooldown: Default	<a href="#">Apply new policy</a>

Cancel Previous **Create scaling plan**

14:15 / 17:13

AWS Auto Scaling

Secure | <https://console.aws.amazon.com/autoscaling/home?region=us-east-1#dashboard>

Services Resource Groups

N. Virginia Support

Scaling plan "WebApp-ScalingPlan" has been created successfully.

AWS Auto Scaling > Scaling plans

Scaling plans (1) Delete **Create scaling plan**

<input type="checkbox"/>	Name	Status	Resource count	Date created
<input type="checkbox"/>	<a href="#">WebApp-ScalingPlan</a>	Active	5	2018-01-25 11:56:13 UTC-0800

14:42 / 17:13

The screenshot displays the AWS Management Console interface for an Auto Scaling Group. The left sidebar shows navigation options like EC2 Dashboard, INSTANCES, IMAGES, ELASTIC BLOCK STORE, NETWORK & SECURITY, and LOAD BALANCING. The main content area shows the 'Create Auto Scaling group' button and a table of existing groups. Below, the 'Activity History' tab is selected, showing a list of scaling events.

Name	Launch Configuration /	Instances	Desired	Min	Max	Availability Zones	Default Cooldown	Health Check Grac
WebApp-Web...	WebApp-LaunchConfig...	4	9	1	10	us-east-1b, us-east-1d	300	0

Status	Description	Start Time	End Time
Not yet in service	Launching a new EC2 instance: i-0622fa25b956a5d5	2018 January 25 11:57:54 UTC-8	
Not yet in service	Launching a new EC2 instance: i-04b9a9e79dc5963	2018 January 25 11:57:54 UTC-8	
Not yet in service	Launching a new EC2 instance: i-03e2171f166f628c	2018 January 25 11:57:54 UTC-8	
Not yet in service	Launching a new EC2 instance: i-08967c88ca42d85d	2018 January 25 11:57:54 UTC-8	
Not yet in service	Launching a new EC2 instance: i-0ef714307132b17a	2018 January 25 11:57:54 UTC-8	
Successful	Terminating EC2 instance: i-04521c6b53371485	2018 January 25 11:43:36 UTC-8	2018 January 25 11:49:21 UTC-8
Successful	Terminating EC2 instance: i-076d016219f0cc59	2018 January 25 11:43:36 UTC-8	2018 January 25 11:49:36 UTC-8
Successful	Terminating EC2 instance: i-02b5bc7e789c6af13	2018 January 25 11:43:36 UTC-8	2018 January 25 11:49:53 UTC-8

Considerations

aws training and certification

## When should you consider using AWS Auto Scaling?

- 📁 If you have applications that use **one or more scalable resources**, and have a **variable load**



## When should you consider using AWS Auto Scaling?

- 📦 If you have applications that use **one or more scalable resources**, and have a **variable load**
- 📦 If you want **more guidance** on defining your application scaling plan



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

## Key Takeaways

- ✓ Monitor your applications continually and automatically adjust capacity.
- ✓ Set up scaling efficiently through automatic resource discovery via a single unified interface.
- ✓ Make smart decisions regarding application high availability with built-in scaling strategies and predictable scaling.
- ✓ Maintain application performance automatically with smart scaling policies.

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.