

15-319 / 15-619
Cloud Computing

Recitation 4

Sep 16, 2019

Administrative - OH & Piazza

- **Make use of office hours**
 - Make sure that you are able to clearly describe the problem and what you have tried so far to provide the fullest context
 - [Piazza Course Staff](#)
 - [Google calendar in ET](#)
 - [Google calendar in PT](#)
- **Suggestions for using Piazza**
 - Read the Piazza Post Guidelines ([@6](#)) before asking questions
 - Read Piazza questions & answers carefully to avoid duplicates
 - Name the subject properly so that others can find your post
 - Try to ask a **public** question if possible so others can also benefit
 - Don't ask a public question about a quiz question

Administrative - Cloud spending

- **Suggestion on cloud service usage**
 - Monitor AWS expenses regularly
 - Always do the cost calculation before launching services
 - **Keep in mind that there is a 6-hour delay for AWS to update their logs on spending**
 - Terminate your instances when not in use
 - Stopped instances have EBS costs (\$0.1/GB-Month)
 - Make sure spot instances are tagged right after launch

Important Notice

- **DON'T EVER EXPOSE YOUR AWS CREDENTIALS!**
 - Github
 - Bitbucket
 - Anywhere public...
- **DON'T EVER EXPOSE YOUR GCP CREDENTIALS!**
- **DON'T EVER EXPOSE YOUR Azure CREDENTIALS!**
 - ApplicationId, ApplicationKey
 - StorageAccountKey, EndpointUrl

Reflection

- Conceptual content on OLI
 - Modules 3, 4, Quiz 2
- Project theme - **Big data analytics**
 - **Inverted Index:** Implemented an inverted index with MapReduce using TDD
 - **Wiki Data Parallel Processing Analysis:** Use MapReduce to process **36GB** compressed / **128GB** uncompressed wiki data
 - MapReduce application to filter records and calculate aggregate daily pageviews
 - **Data Analytics:** Use Jupyter Notebooks and the pandas library to analyze the data and answer questions

This Week

- **Quiz 3 (OLI Modules 5 & 6)**
 - Due on **Friday**, Sep 20th, 2019, 11:59PM ET
- **Project 2.1 and Reflection**
 - Due on **Sunday**, Sep 22nd, 2019, 11:59PM ET
- **Project 1.2 Discussion**
 - Due on **Sunday**, Sep 22nd, 2019, 11:59PM ET
- **P1.2 Code Review**
 - Due on **Wednesday**, Sep 25th, 2019, 11:59PM ET
- **Primers released this week**
 - P2.2 - Intro to Containers and Docker
 - P2.2 - Kubernetes and Container Orchestration
 - Code Review

Code Review

- Code review is the systematic examination of source code. The goal of code review is to make sure that the code achieves its objective using a sound approach and to improve it if possible.
- We want you to develop good coding habits and skills that will be useful for your careers.
- Please read the “Code Review” primer on TheProject.Zone.
- For Project 1.2, completing code review is worth 5 points, and it will contribute toward the total grade of Project 1.2.

OLI Module 5 - Cloud Management

Cloud Software stack - enables provisioning, monitoring and metering of virtual user “resources” on top of the Cloud Service Provider’s (CSP) infrastructure.

- Cloud middleware
- Provisioning
- Metering
- Orchestration and automation
- Case Study: Openstack - Open-source cloud stack implementation

OLI Module 6 - Cloud Software Deployment Considerations

- Programming on the cloud
- Deploying applications on the cloud
- Build fault-tolerant cloud services
- Load balancing
- Scaling resources
- Dealing with tail latency
- Economics for cloud applications

Project 2 Overview

Scaling and Elasticity with

- VMs
- Containers
- Functions

- **2.1 Scaling Virtual Machines**
 - Horizontal scaling in / out using AWS APIs
 - Load balancing, failure detection, and cost management on AWS
 - Infrastructure as Code (Terraform)
- **2.2 Scaling with Containers**
 - Building your own container-based microservices
 - Docker containers
 - Manage multiple Kubernetes Cluster
 - Multi Cloud deployments
- **2.3 Functions as a Service**
 - Develop event driven cloud functions
 - Deploy multiple functions to build a video processing pipeline

Project 2.1 Learning Objectives

- **Design** solutions and invoke cloud APIs to programmatically provision and deprovision cloud resources based on the current load.
- **Explore** the usability and performance of APIs used in AWS.
- **Configure** and deploy an Elastic Load Balancer along with an Auto Scaling Group on AWS.
- **Develop** elasticity policies to maintain the QoS of a web service that also deals with resource failure.
- **Account** for cost as a constraint when provisioning cloud resources and analyze the performance tradeoffs due to budget restrictions.
- **Experience** using cloud orchestration and automation tools such as Terraform.

Overview of Quality of Service (QoS), Latency and Cloud Elasticity

- Quality of Service (QoS)
- Load patterns for web services
- Vertical scaling (Scale up/down)
- Horizontal scaling (Scale out/in)
- Load balancers
- Autoscaling groups
- Resource monitoring (CloudWatch)

Quality of Service (QoS)

Quantitatively Measure QoS

- **Performance**: Throughput, Latency
(Very helpful in Project 2 & Team Project)
- **Availability**: the probability that a system is operational at a given time *(Project 2)*
- **Reliability**: the probability that a system will produce a correct output up to a given time *(Project 2)*

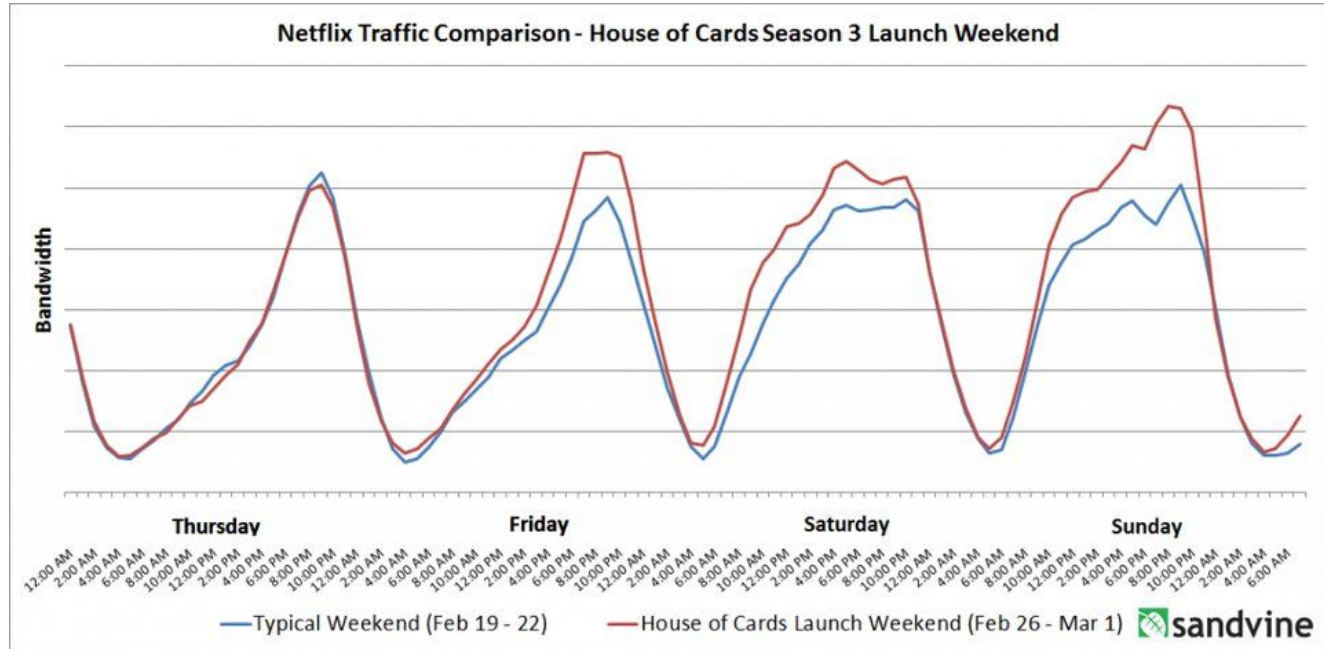
QoS Matters:

- Amazon found every **100ms** of latency cost them **1%** in sales (~\$1B).



Reality, human patterns...

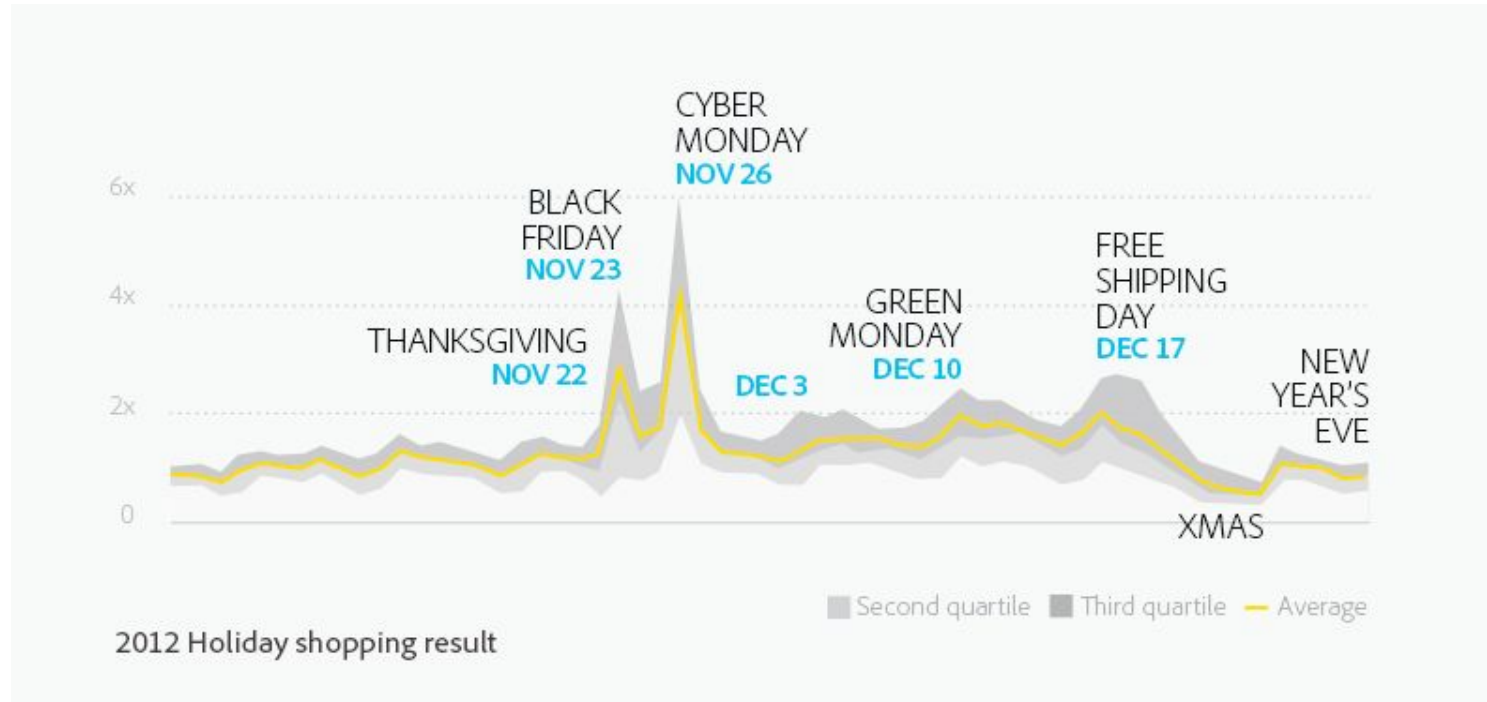
- Daily
- Weekly
- Monthly
- Yearly
- ...



The Ferenstein Wire

Reality, human patterns...

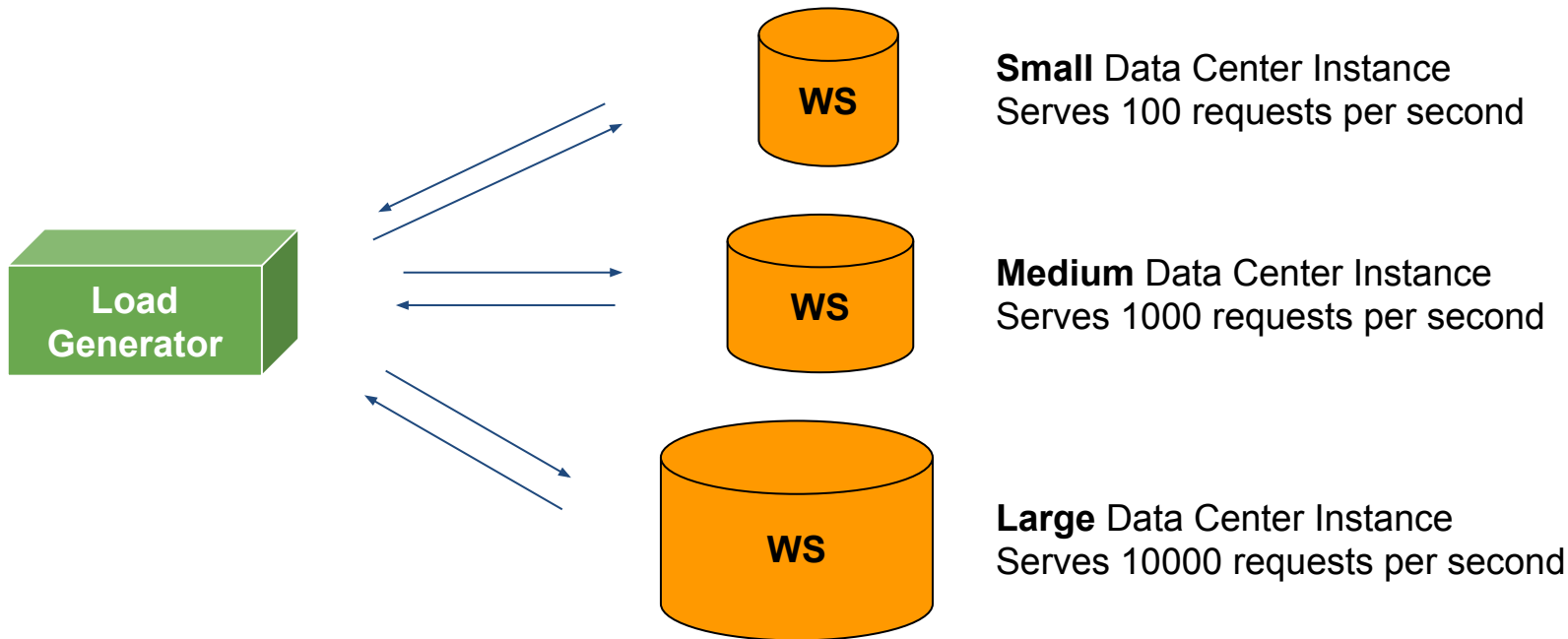
- Daily
- Weekly
- Monthly
- **Yearly**
- ...



Cloud Comes to the Rescue!

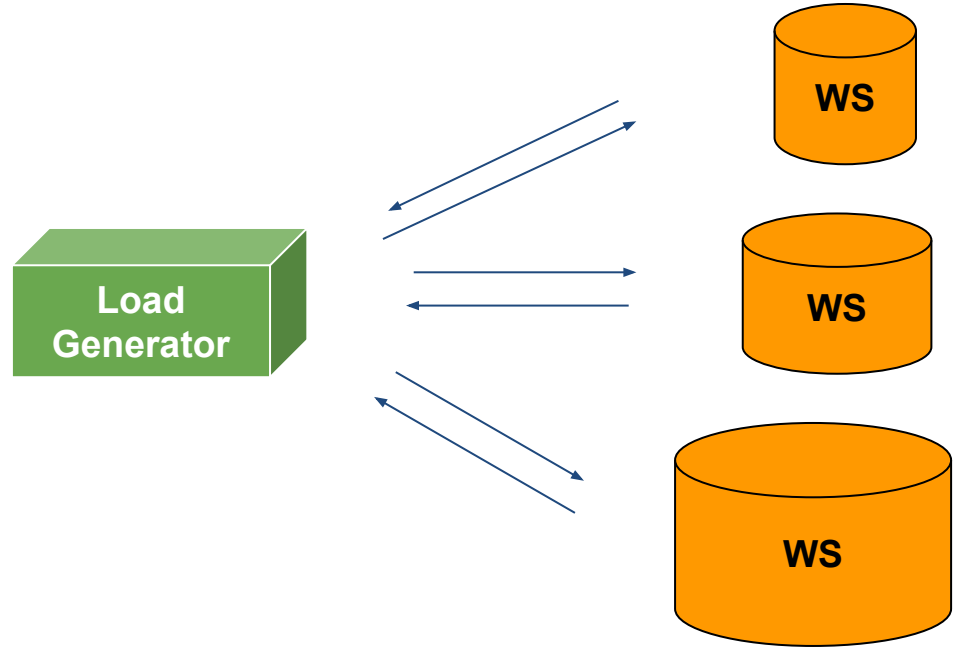
Scaling!

P0: Vertical Scaling

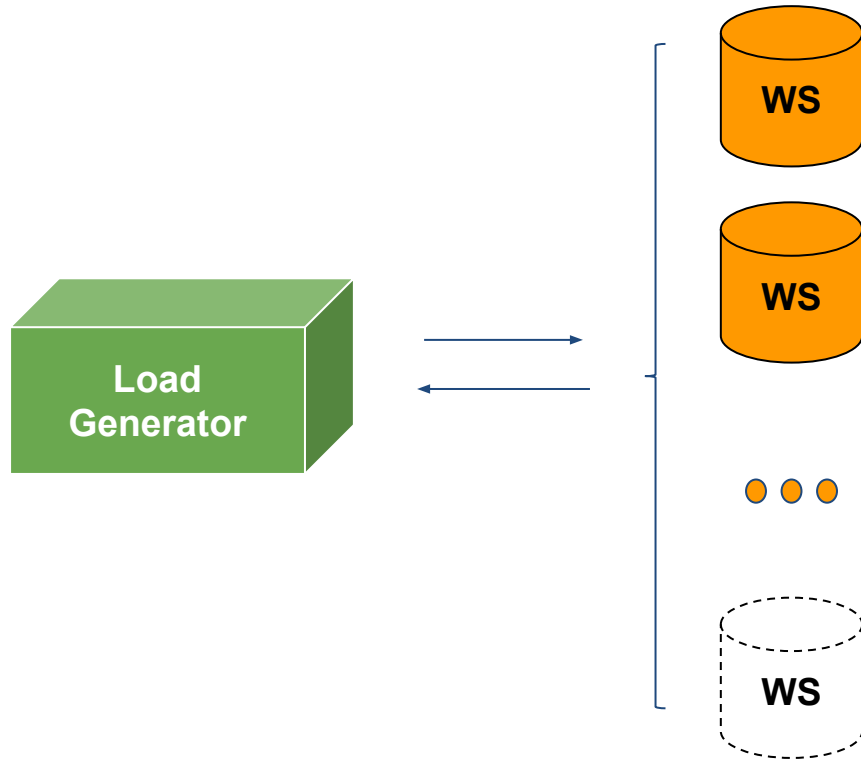


P0: Vertical Scaling Limitation

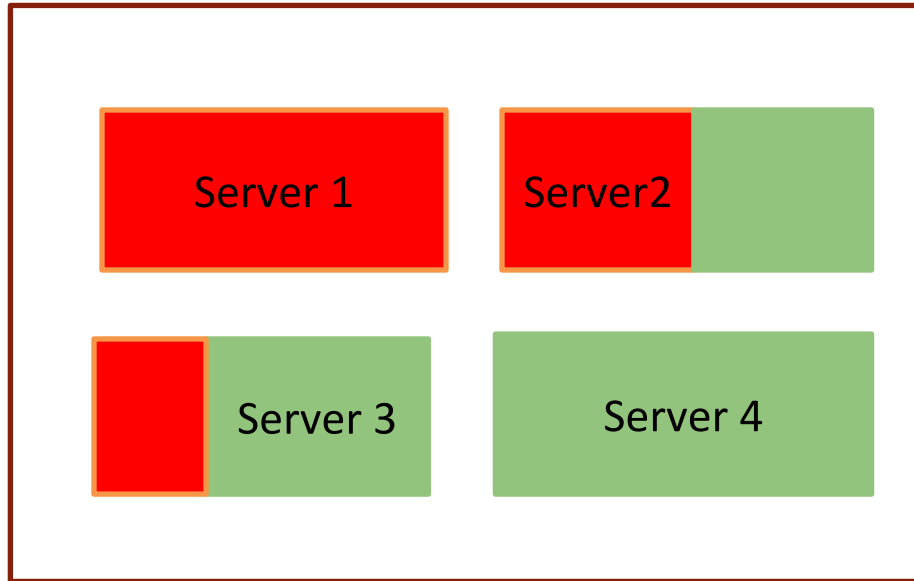
- However, one instance will always have limited resources.
- Reboot/Downtime.



Horizontal Scaling



How do we distribute load?

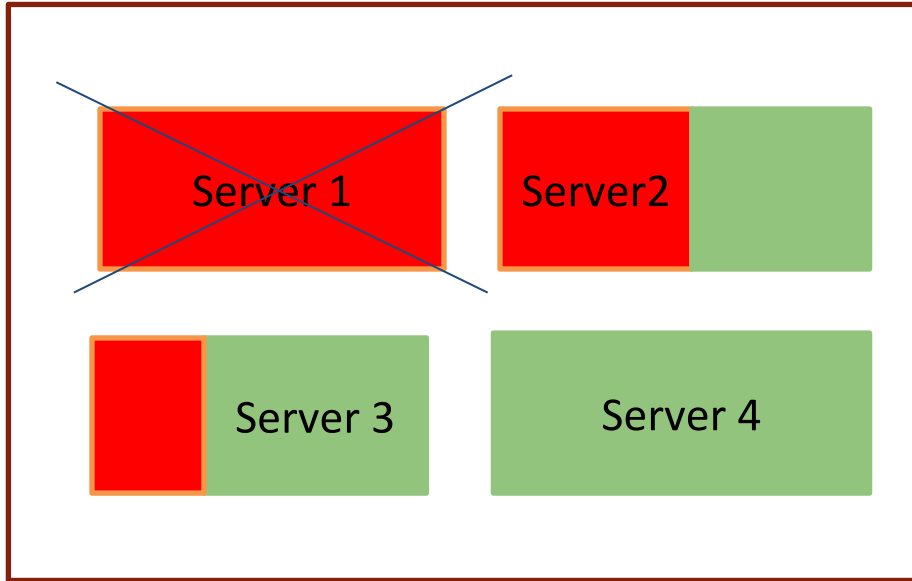


CPU utilization, memory utilization...



Available capacity

Instance Failure?



CPU utilization, memory utilization...

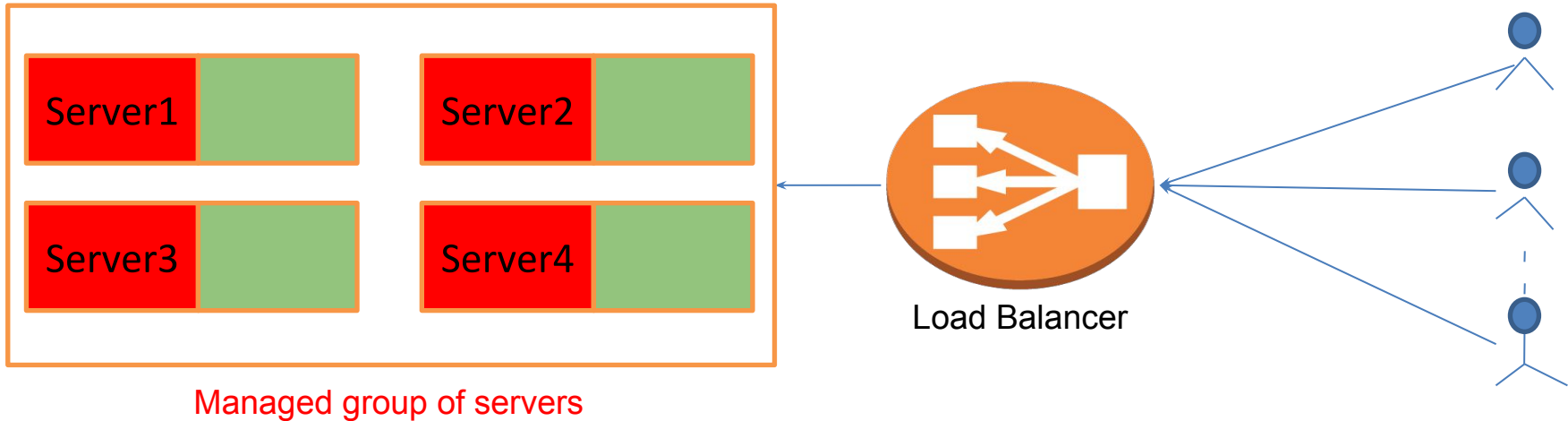


Available capacity

What You Need

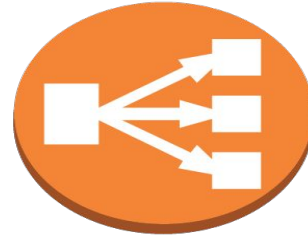
- Make sure that workload is even on each server
- Do not assign load to servers that are down
- Increase/Remove servers according to changing load

How does a cloud service help solve these problems?



Load balancer

- “Evenly” distribute the load
 - A simple distribution strategy
 - Round Robin
 - Health Check
-
- What if the Load Balancer becomes the bottleneck?
 - Elastic Load Balancer (ELB)
 - Could scale up based on load
 - Elastic, but it still takes time
 - Through the warm-up process



Load Balancer

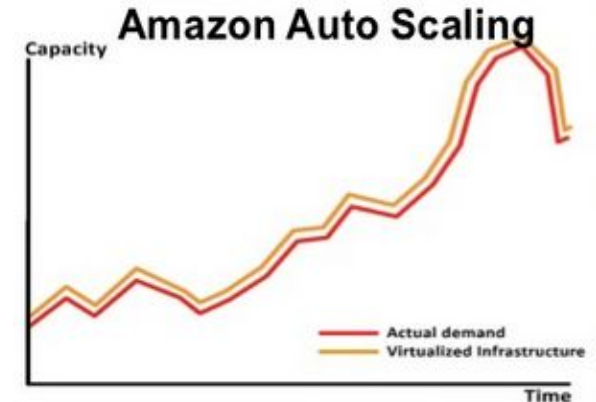
Scaling

Manual Scaling:

- Expensive on manpower
- Over provisioning and low utilization
- Lose customers

Autoscaling:

- Automatically adjust the size based on demand
- Flexible capacity and scaling sets
- Save cost



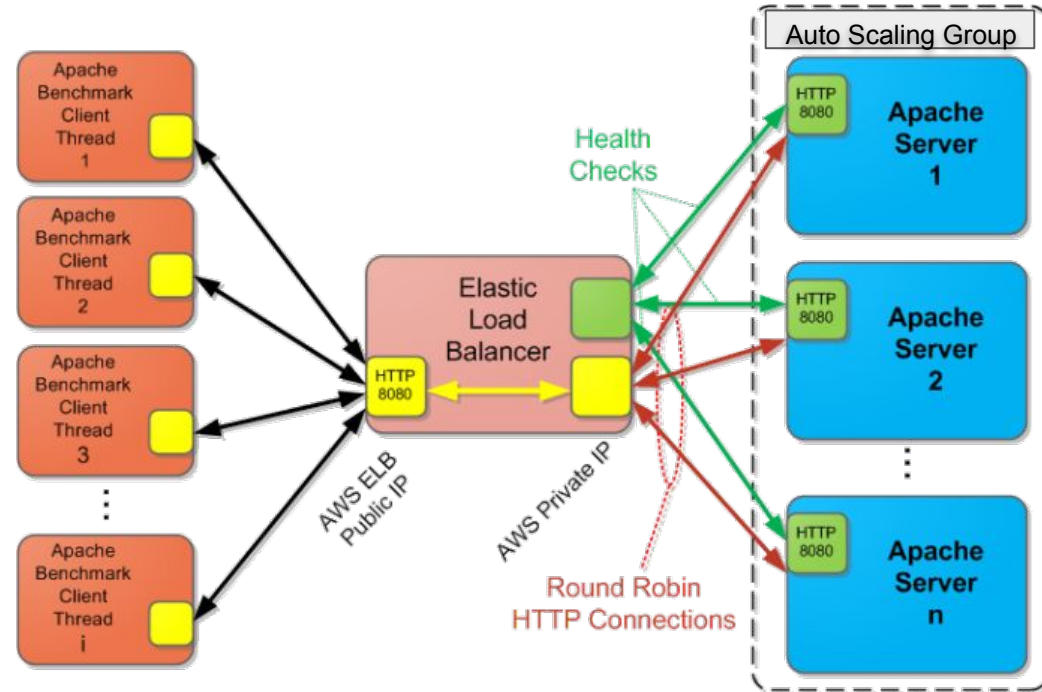
AWS Autoscaling

Auto Scaling on AWS

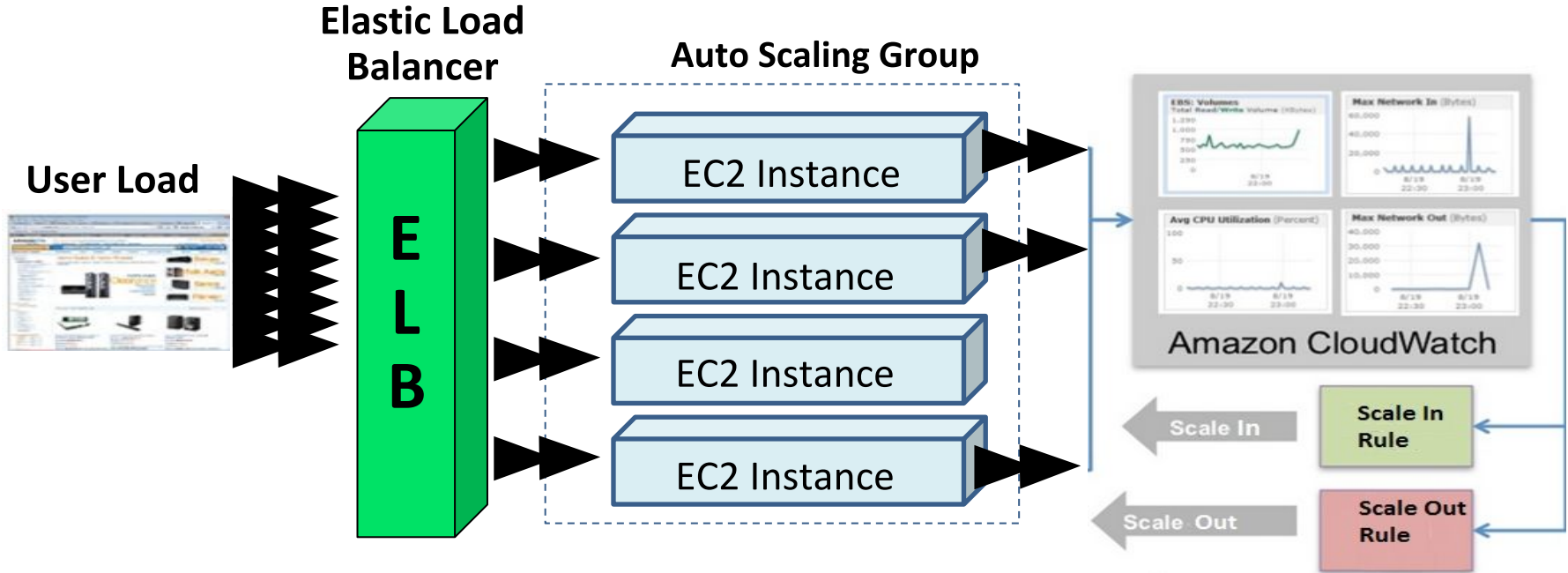
Using the AWS APIs:

- ELB
- Auto Scaling Group
- EC2
- CloudWatch
- Auto Scaling Policy

You can build a load balanced auto-scaled web service.

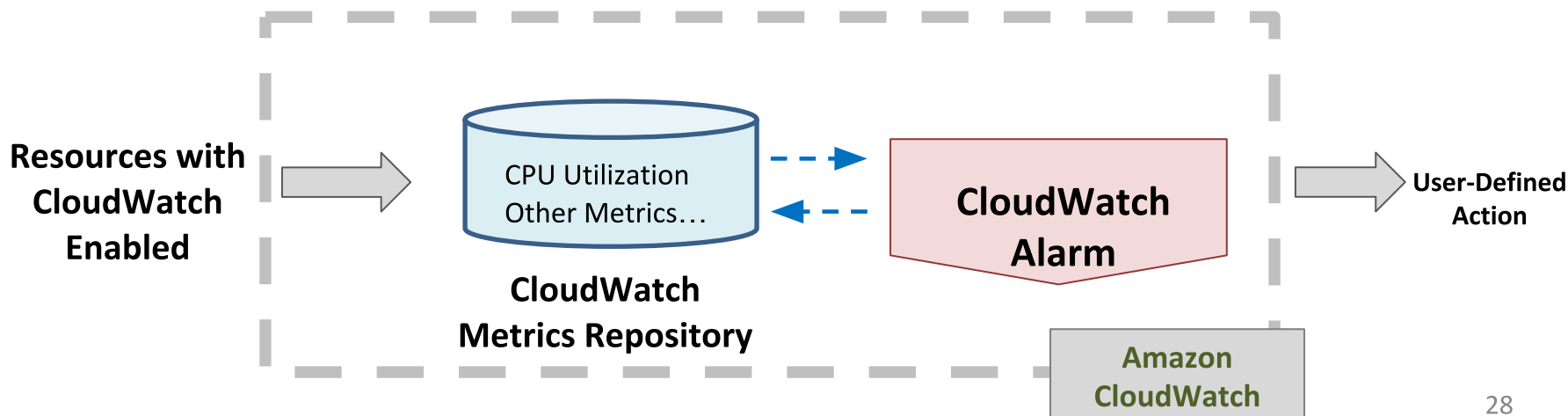


Amazon Auto Scaling Group



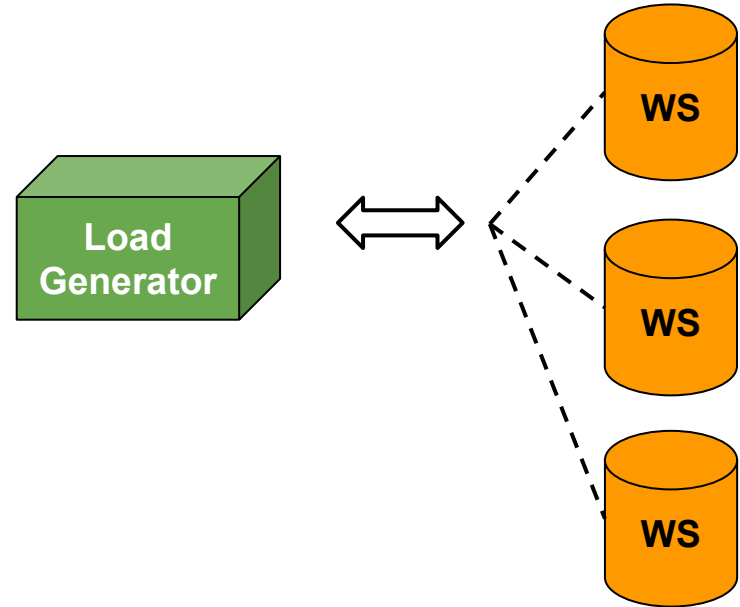
Amazon's CloudWatch Alarm

- Monitor CloudWatch metrics for some specified alarm conditions
- Take automated action when the condition is met



Project 2.1 Scaling on AWS

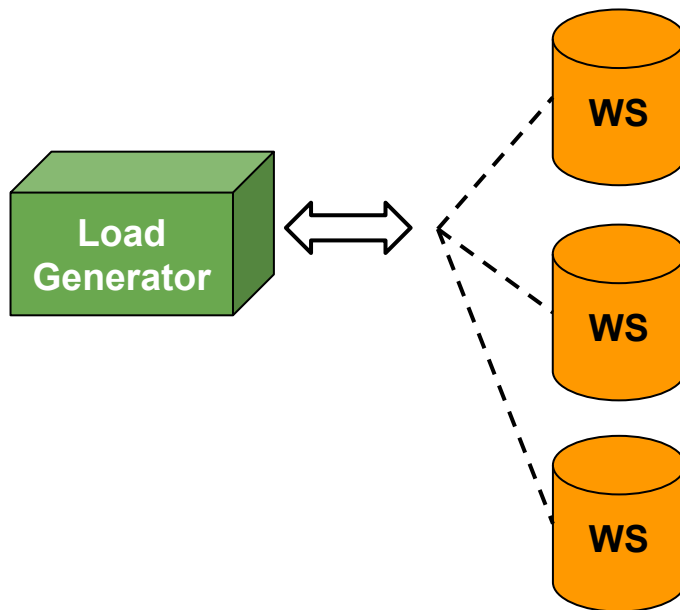
- **Task 1**
 - **AWS Horizontal Scaling**
- **Task 2**
 - AWS Auto Scaling
- **Task 3**
 - AWS Auto Scaling with Terraform



Project 2.1 Scaling on AWS

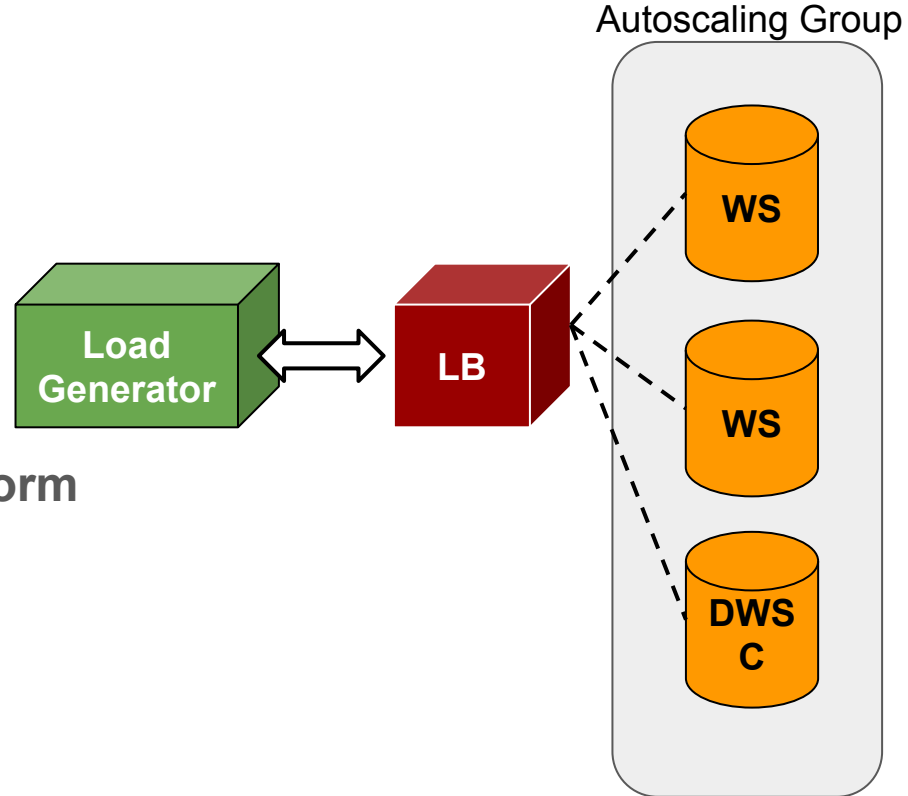
Task 1 - AWS Horizontal Scaling:

- Implement Horizontal Scaling in AWS.
- Write a program that launches the web service instances and ensures that the target total RPS is reached.
- Your program should be fully automated: launch LG → submit password → Launch WS → start test → check log → add more WS...



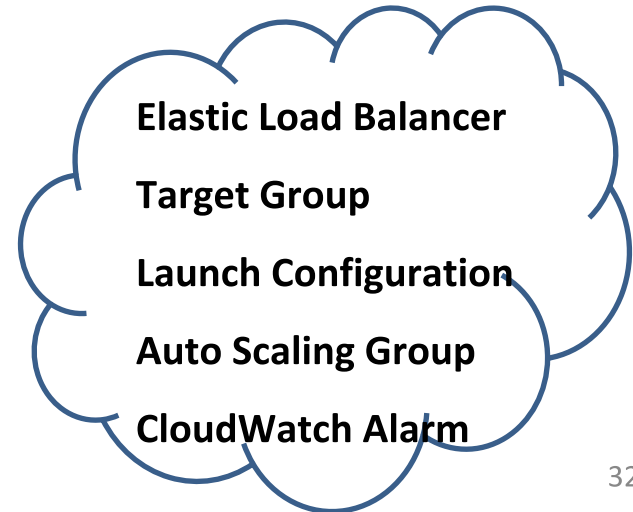
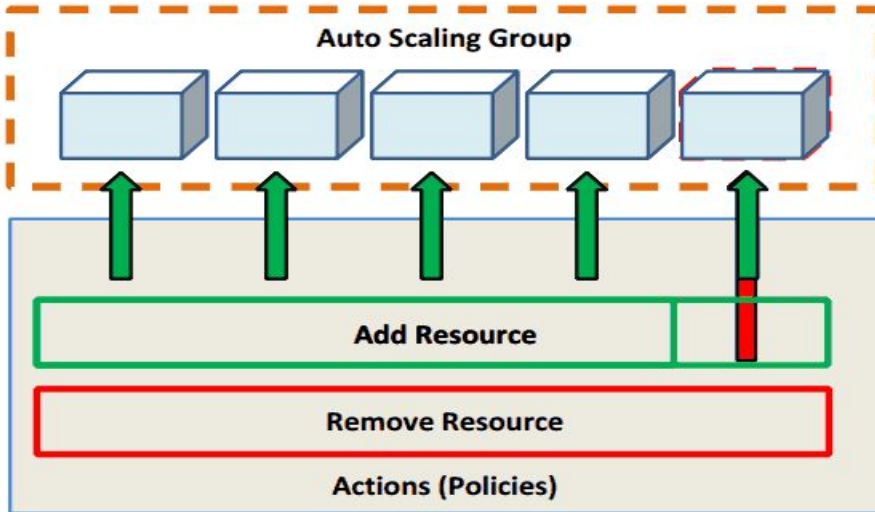
Project 2.1 Scaling on AWS

- Task 1
 - AWS Horizontal Scaling
- Task 2
 - **AWS Auto Scaling**
- Task 3
 - AWS Auto Scaling with Terraform



P2.1 - Task 2

- Programmatically create two security groups, LG, Application Load Balancer (ALB), Auto-Scaling Group (ASG) along with Auto Scaling Policy, launch configuration, and target group.
- Test by submitting a URL request and observe logs from ALB, and CloudWatch.
- Adjust Scale-Out and Scale-In policies if necessary
- Your solution also need to be fault tolerant, health configurations are important



Hints for Project 2.1 AWS Autoscaling

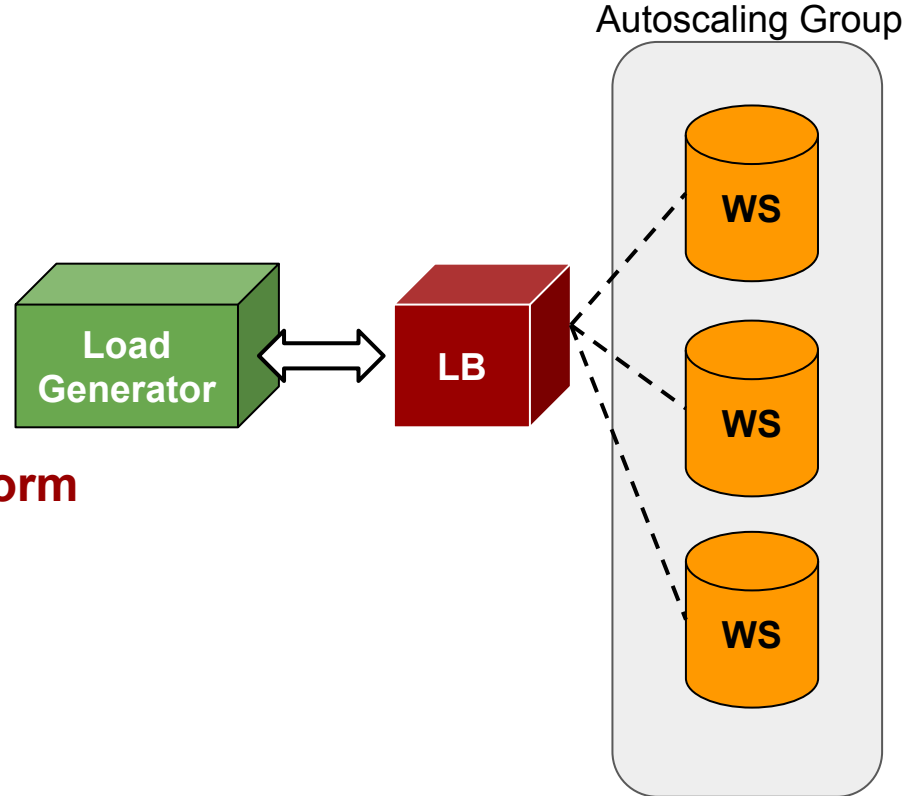


Task 2 - AWS Auto Scaling

- Do a dry run via the console to make sure you understand the workflow completely and mimic that workflow programmatically.
- Autoscaling Test could be very expensive!
 - On-demand, charged by per second, do not blindly launch tests
- Determine if there is a less expensive means to test your solution
- Creating and deleting security groups can be tricky
- CloudWatch and monitoring in ELB is helpful for policy tuning
- Explore ways to check if your instance is ready
- **Understanding the API documents could take time**

Project 2.1 Scaling on AWS

- Task 1
 - AWS Horizontal Scaling
- Task 2
 - AWS Auto Scaling
- **Task 3**
 - **AWS Auto Scaling with Terraform**



Project 2.1 Scaling on AWS



Task 3 - AWS Auto Scaling with Terraform:

- Read the Infrastructure as Code primer to learn about infrastructure automation
- Update your code to create and attach an IAM role to your Load Generator.
- Makes sure that `terraform plan` generates the expected resource

Project 2.1 Code Submission

- At the end of each task, you need to submit your code for that task to the load generator. We will grade your code of each task **separately**.
- You need to package your code into a **.tar.gz** file and name it after your testId. For the Terraform task, you can you use a testID Auto Scaling task.
 - `mkdir testId; tar -zcf <testId>.tar.gz <testId>`
- **Remove all hidden files from your submission(.git / .idea / .DS_store)**
 - `export COPYFILE_DISABLE=true; (Mac Users)`
`tar --exclude='.*' -zcf <testId>.tar.gz <testId>`

Penalties for Project 2.1

Violation	Penalty of the project grade
Spending more than \$20 for this project phase on AWS	-10%
Spending more than \$35 for this project phase on AWS	-100%
Failing to tag all your resources in either parts (EC2 instances, ELB, ASG) for this project with the tag: key=Project, value=2.1	-10%
Submitting your AWS/Andrew credentials in your code for grading	-100%
Using instances other than t3.micro (testing only) or m5.large for Horizontal scaling on AWS	-100%
Using instances other than t3.micro (testing only), m5.large for Autoscaling on AWS	-100%
Submitting executables (.jar, .pyc, etc.) instead of human-readable code (.py, .java, .sh, etc.)	-100%

Penalties for Project 2.1 cont.

Violation	Penalty of the project grade
Attempting to hack/tamper the autograder in any way	-200%
Cheating, plagiarism or unauthorized assistance (please refer to the university policy on academic integrity and our syllabus)	-200%

AWS Cloud APIs



- AWS CLI ([link](#))
- AWS Java SDK ([link](#))
- AWS Python SDK ([link](#))

This Week

- **Quiz 3 (OLI Modules 5 & 6)**
 - Due on **Friday**, Sep 20th, 2019, 11:59PM ET
- **Project 2.1 and Reflection**
 - Due on **Sunday**, Sep 22nd, 2019, 11:59PM ET
- **Project 1.2 Discussion**
 - Due on **Sunday**, Sep 22nd, 2019, 11:59PM ET
- **P1.2 Code Review**
 - Due on **Wednesday**, Sep 25th, 2019, 11:59PM ET
- **Primers released this week**
 - P2.2 - Intro to Containers and Docker
 - P2.2 - Kubernetes and Container Orchestration
 - Code Review

Questions?