

Capstone Project – Zhuoyu Zhu

Q7. Handle archiving of results files for free users after 30 minutes

Approach used:

Each time the specified job request is completed and is about to send confirmation email to the user, I first check if the current user is a free user role by checking the data received from SQS. If the user is a free user, I use subprocess method and schedule an archive task in 30 minutes by using sched library. In addition, in case the specified user decides to upgrade to the premium user role within 30 minutes frame, to prevent archiving results file in this scenario, I check the user role again when the scheduled task is about to execute. If the user has upgraded to premium, I don't need to archive this result file any more otherwise we need to archive this result file to Amazon Glacier.

Reason why I used this approach:

Using task scheduler and subprocess make the entire archive process workload more distributed without slowing down user experience on our platform. On the other hand, if you check the time passed within a loop for each job request, you might either need to read data from DynamoDB or SQS. Using these approaches will slow down the annotation platform significantly since the cost of reading data constantly from DynamoDB or SQS within a loop is very expensive. Moreover, the size of SQS is limited, if our platform has hundreds of million of job ready for archive, we have to spend additional efforts to maintain the SQS or distribute some of the workload to another SQS. Using task scheduler and subprocess will solve all of these issues.

Q9. Restore data for users that upgrade to Premium

Approach used:

After the user has upgraded their membership successfully without throwing Stripe API exception, I initiate the archive restoration job request for any result files that were generated while the user was a free user. Then I used an Amazon SQS (zhuoyuzhu_glacier) to accept these request within mpcs_app.py file. After this, within the restore instance, I retrieve these request messages and check status. When the restoration is ready, the python script writes contents onto the result file and upload to S3 result bucket.

Reason why I used this approach:

Since the restoration requests for the given user are relatively lightweight compared to submitting annotation job request, so using Amazon SQS queue should be sufficient for this distributed cloud-based application. The chance about reaching the size limit of SQS is relatively small.

Q11. Add auto-scaling rules to the web server farm

CloudWatch for High-Average-Latency and High-Sum-HTTP-2XXs:

Load balancer: **zhuoyuzhu-elb**

Description

Listeners

Monitoring

Tags

CloudWatch alarms: 2 of 2 in OK Create Alarm

Below are your CloudWatch alarms for the selected resources. Click on an alarm to edit it or click 'view' to see additional options and details in Amazon CloudWatch. [View all CloudWatch alarms](#)

State	Name	More Options
OK	awsapplicationelb-app-zhuoyuzhu-elb-8dcc0e8a688c863a-High-Average-Latency	view
OK	awsapplicationelb-app-zhuoyuzhu-elb-8dcc0e8a688c863a-High-Sum-HTTP-2XXs	view

CloudWatch metrics: Showing data for: Last Hour

Scaling Policy:

Auto Scaling Group: **zhuoyuzhu-auto-scaler**

Details

Activity History

Scaling Policies

Instances

Monitoring

Notifications

Tags

Scheduled Actions

Add policy Actions

add-one-instance-when-exceed-200-HTTP-2XX

Execute policy when:

awsapplicationelb-app-zhuoyuzhu-elb-8dcc0e8a688c863a-High-Sum-HTTP-2XXs breaches the alarm threshold: HTTPCode_Target_2XX_Count > 200 for 60 seconds for the metric dimensions LoadBalancer = app/zhuoyuzhu-elb/8dcc0e8a688c863a

Take the action:

Add 1 instances when 200 <= HTTPCode_Target_2XX_Count < +infinity

Instances need:

120 seconds to warm up after each step

Actions

remove-one-instance-when-latency-below-10ms

Execute policy when:

awsapplicationelb-app-zhuoyuzhu-elb-8dcc0e8a688c863a-High-Average-Latency breaches the alarm threshold: TargetResponseTime < 0.01 for 60 seconds for the metric dimensions LoadBalancer = app/zhuoyuzhu-elb/8dcc0e8a688c863a

Take the action:

Remove 1 instances when 0.01 >= TargetResponseTime > -infinity

Actions

Q12. Add scaling rules to the annotator

Launch configuration for annotator instances:

Launch Configuration: zhuoyuzhu-launch-config-annotator

Details

Copy launch configuration

AMI ID	ami-e3b1c7f5	Instance Type	t2.micro
IAM Instance Profile	instance_profile_zhuoyuzhu	Kernel ID	
Key Name	zhuoyuzhu	Monitoring	false
EBS Optimized	false	Security Groups	sg-7cb8b002
Spot Price		Creation Time	Tue Jun 06 00:51:48 GMT-500 2017
RAM Disk ID		Block Devices	/dev/sda1
User data	View User data	IP Address Type	Only assign a public IP address to instances launched in the default VPC and subnet.

Auto-scaling for annotator instance:

Auto Scaling Group: zhuoyuzhu-auto-scaler-annotator

DetailsActivity HistoryScaling PoliciesInstancesMonitoringNotificationsTagsScheduled Actions

Edit

Launch Configuration	zhuoyuzhu-launch-config-annotator			Availability Zone(s)	us-east-1e
Load Balancers				Subnet(s)	subnet-40dee97a
Target Groups				Default Cooldown	300
Desired	2			Placement Group	
Min	2			Suspended Processes	
Max	10			Enabled Metrics	GroupMaxSize, GroupTotalInstances, GroupDesiredCapacity, GroupInServiceInstances, GroupMinSize, GroupTerminatingInstances, GroupPendingInstances.
Health Check Type	EC2				
Health Check Grace Period	300				
Termination Policies	Default				

CloudWatch for annotator:

Alarm:zhuoyuzhu-sqs-msg-sent-exceed-10-for-5min

DetailsHistory

State Details: State changed to OK at 2017/06/06. Reason: Threshold Crossed: 1 datapoint (0.0) was not greater than the threshold (10.0).

Description:

Threshold: NumberOfMessagesSent > 10 for 5 minutes

Actions: In ALARM: • For group zhuoyuzhu-auto-scaler-annotator use policy add-1-instance-when-msg-sent-exceed-10-for-5min (Add 1 instance)

Namespace: AWS/SQS

Metric Name: NumberOfMessagesSent

Dimensions: QueueName = zhuoyuzhu_job_requests

Statistic: Sum

Period: 5 minutes

Treat missing data missing as:

Percentiles with low samples: evaluate

zhuoyuzhu-sqs-msg-sent-exceed-...
NumberOfMessagesSent > 10

Alarm:zhuoyuzhu-sqs-msg-received-below-5-for-2min

Details

History

State Details: State changed to ALARM at 2017/06/06. Reason: Threshold Crossed: 1 datapoint (0.0) was less than the threshold (5.0).

Description:

Threshold: NumberOfMessagesSent < 5 for 2 minutes

Actions: In ALARM:

- For group **zhuoyuzhu-auto-scaler-annotator** use policy **remove-1-instance-when-msg-received-below-5-for-2min** (Remove 1 instance)

Namespace: AWS/SQS

Metric Name: NumberOfMessagesSent

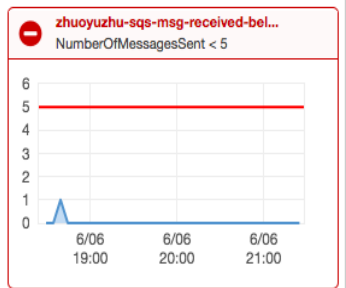
Dimensions: QueueName = zhuoyuzhu_job_requests

Statistic: Sum

Period: 1 minute

Treat missing data missing
as:

Percentiles with evaluate
low samples:



Auto-scaling policy for annotator:

Add policy



add-1-instance-when-msg-sent-exceed-10-for-5min

Actions ▾

Execute policy when: zhuoyuzhu-sqs-msg-sent-exceed-10-for-5min
breaches the alarm threshold: NumberOfMessagesSent > 10 for 300 seconds
for the metric dimensions QueueName = zhuoyuzhu_job_requests

Take the action: Add 1 instances when $10 \leq \text{NumberOfMessagesSent} < +\infty$

Instances need: 120 seconds to warm up after each step

remove-1-instance-when-msg-received-below-5-for-2min

Actions ▾

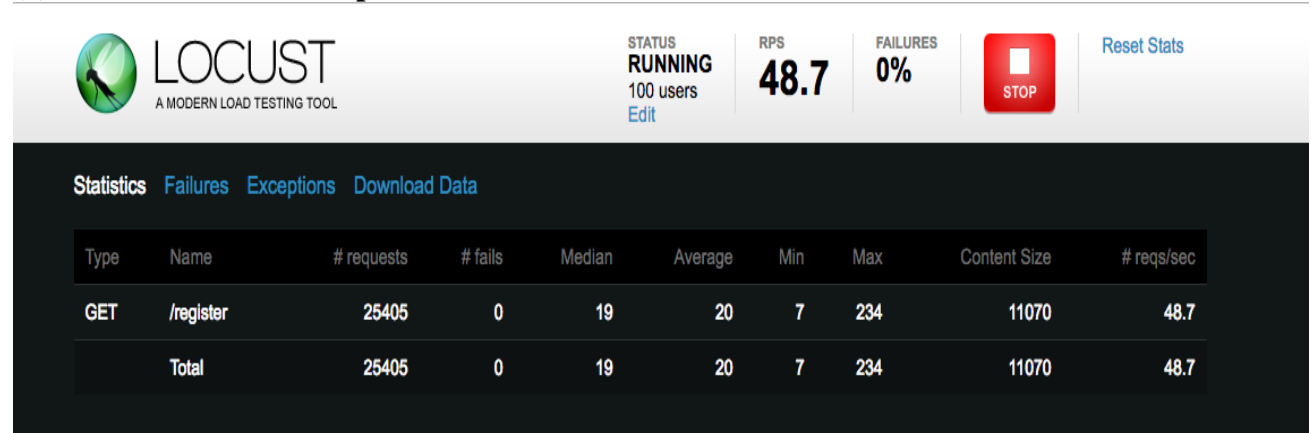
Execute policy when: zhuoyuzhu-sqs-msg-received-below-5-for-2min
breaches the alarm threshold: NumberOfMessagesSent < 5 for 2 consecutive periods of 60 seconds
for the metric dimensions QueueName = zhuoyuzhu_job_requests

Take the action: Remove 1 instances when $5 \geq \text{NumberOfMessagesSent} > -\infty$

Q13. Test under load using Locust (Conclusion are at the bottom of this question)

1. Register template

(1) 100 users and 5 users per second

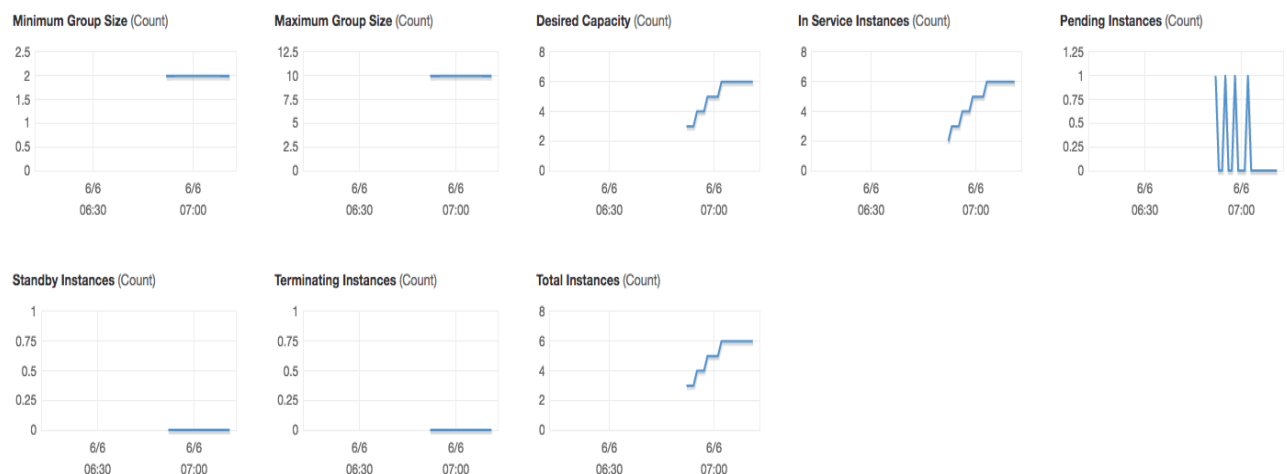


Filter: Any Status

1 to 6 of 6 History Items


Status	Description	Start Time	End Time
Successful	Launching a new EC2 instance: i-03d508c5a00cc9dc1	2017 June 6 02:02:21 UTC-5	2017 June 6 02:05:24 UTC-5
Successful	Launching a new EC2 instance: i-04101310f82a118fc	2017 June 6 01:58:16 UTC-5	2017 June 6 02:00:49 UTC-5
Successful	Launching a new EC2 instance: i-07624d4ae99935881	2017 June 6 01:55:13 UTC-5	2017 June 6 01:57:45 UTC-5
Successful	Launching a new EC2 instance: i-08d86dc3eb0394aff	2017 June 6 01:52:09 UTC-5	2017 June 6 01:54:41 UTC-5
Successful	Launching a new EC2 instance: i-066df995083399c03	2017 June 5 02:15:38 UTC-5	2017 June 5 02:16:11 UTC-5
Successful	Launching a new EC2 instance: i-02fa640c45cc32318	2017 June 5 02:15:38 UTC-5	2017 June 5 02:16:11 UTC-5

zhuoyuzhu-auto-scaler



2. Home Page

(1) 200 users and 20 users per second

**LOCUST**
A MODERN LOAD TESTING TOOL

STATUS
STOPPED
New test

RPS
97.6

FAILURES
0%

[Statistics](#) [Failures](#) [Exceptions](#) [Download Data](#)

Type	Name	# requests	# fails	Median	Average	Min	Max	Content Size	# reqs/sec
GET	/	29458	55	17	45	6	9707	11167	97.6
Total		29458	55	17	45	6	9707	11167	97.6

▼ **Waiting for instance warmup**

Launching a new EC2 instance: i-09569aa883161a5d9

2017 June 6 03:12:09 UTC-5

Description: Launching a new EC2 instance: i-09569aa883161a5d9

Cause: At 2017-06-06T08:11:47Z a monitor alarm awsapplicationelb-app-zhuoyuzhu-elb-8dcc0e8a688c863a-High-Sum-HTTP-2XXs in state ALARM triggered policy zhuoyuzhu-High-Sum-HTTP-2XXs changing the desired capacity from 5 to 6. At 2017-06-06T08:12:07Z an instance was started in response to a difference between desired and actual capacity, increasing the capacity from 5 to 6.

▼ **Successful**

Launching a new EC2 instance: i-0b3547e69934f120c

2017 June 6 03:09:05 UTC-5

2017 June 6 03:11:38 UTC-5

Description: Launching a new EC2 instance: i-0b3547e69934f120c

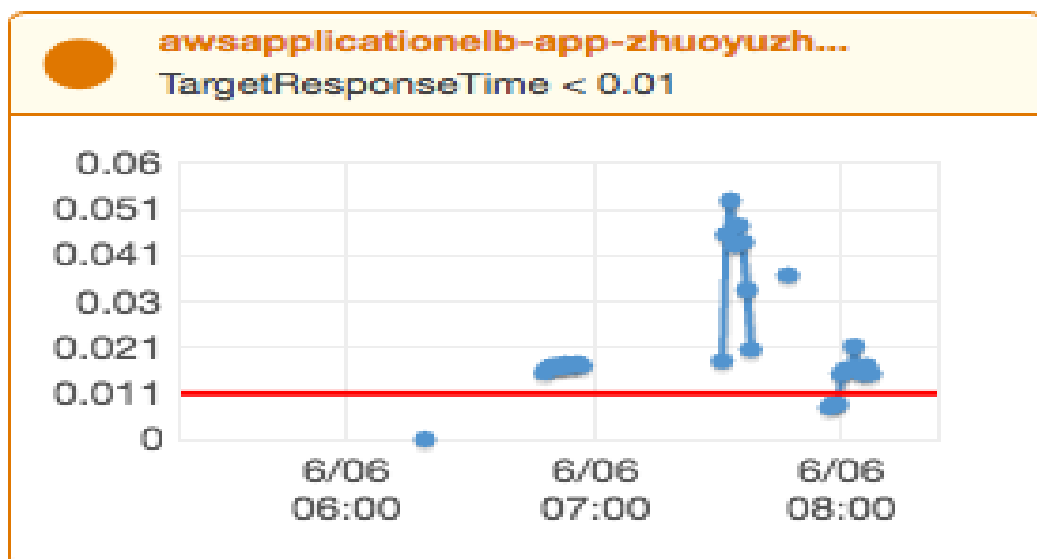
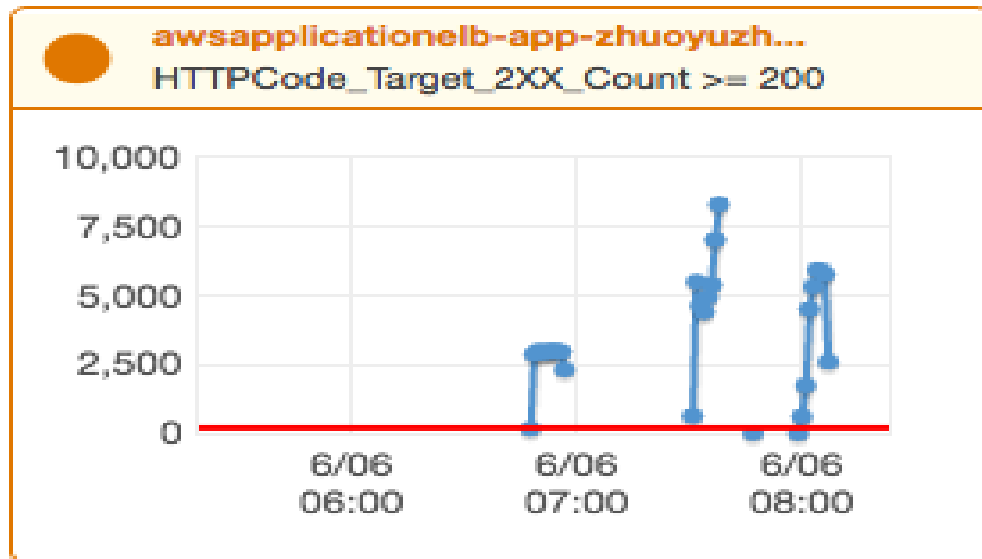
Cause: At 2017-06-06T08:08:47Z a monitor alarm awsapplicationelb-app-zhuoyuzhu-elb-8dcc0e8a688c863a-High-Sum-HTTP-2XXs in state ALARM triggered policy zhuoyuzhu-High-Sum-HTTP-2XXs changing the desired capacity from 4 to 5. At 2017-06-06T08:09:03Z an instance was started in response to a difference between desired and actual capacity, increasing the capacity from 4 to 5.

▼ **Successful**

Launching a new EC2 instance: i-0a8284581d8052882

2017 June 6 03:08:02 UTC-5

2017 June 6 03:08:25 UTC-5



Observation: The majority of the HTTP response code returns as 2XXs which mean our cloud-based annotation application works well when the workload is huge. It keep launch new instances when the alarm is triggered based on the specified timeframe. On the other hand, the other alarm average latency was never triggered. It means when the workload becomes huge, it will takes time to handle these amount of http request.

Q14. Python script for annotator load test

Please refer to the python file, annTest. Run out of time to capture the annotator test screenshots. The approach is using a while loop to keep send job request to the Amazon SNS topic and SQS queue. I used a count variable to record the current number of job requests sent. If the number of job requests exceeds 10, the python script will pause for 120 seconds otherwise it will pause for 25 second after each job request is sent. This will make sure to trigger the scale-out and scale-in policy for annotator to add and remove 1 instance respectively.