

HW 6

Created @October 12, 2025 11:44 PM

Part2

Load Testing & Analysis

search service

```
# 1. run the app
cd ~/CS6650/hw6/src
python main.py

# 2. new terminal to test
curl "http://localhost:8080/products/search?q=alpha"
curl "http://localhost:8080/health"
curl "http://localhost:8080/"

# 3. docker check
docker build -t hw6-search .
# run
docker run -p 8080:8080 hw6-search
# Test
curl http://localhost:8080/health

# 4. run locust
locust -f locustfile.py --host http://localhost:8080

# 5. aws config
[default]
aws_access_key_id=ASIAVRUVWC52BNAKBXTT
aws_secret_access_key=1VE6dPuWew+VuDQS4XI4y+cDnf1mXVfaf56xjsBQ
aws_session_token=IQoJb3JpZ2luX2VjEj//////////wEAxCVzLXdIc3QtMiJHMEUCIDaxfUmsRmW1BFglpYgWP4Upd69X9d7hmYiFoAVrtFCZAiEAiNtRnKVSGJCmAx91GR0fkt4nmDXEQwipsXiKicXF+KMqtQIIQRAAGgwzODE0OTlyNzA5NjQiDIa/aa6zCC3GI7PkVyqSAiPqGxcgNgVtxCgJaGX/Xp2dANZyL3DOsaGArr4nhGWi7DxULLmDu4haOVzioa6IDc/2Mhpfd+UwWUWNuuJ8fCWjHRubfnuShoBTAHilSqCNe/2b6poF2+vpzb2XUQMuf2It7DMELeu+Cv5LLgbKAXUC6pnoeFLdf+6WDWFsP/uKvSTY9K9hjT/pGn55pLyw8bGf6hBx4AstJf7nAKk+pZCSzWrpNxOgFTjEvFvKLMLVFvuD6x/U2myDbp0URymTEt+b1vqlAtNiW+/UcJAdBPNYR6pWSR5VnOkHptRfwrfHhFHh7KGpvsm9vAl4kZB3MlqWYLOFZ5p6dVD75RQ/IaO5ILBB/xXY5jRJyb7EcrlBxjMwzeWyxwY6nQE8XcbHolvUUh28nptSOh/RIU/jCpBN7q4Tk4ohoL2YJDDfTh52KGVDIFn+aPFSTOvznG85SQSxW2FkM+zYtelpE5FtgIcW6dJg8AK8dJwGbKOEd3AIG4z0Zn4QZ55EWN2xyN+5pNPrI5/XwLRPbOffWu0nrE0WzVX1V7G5Of6kITZScy3SCoCr+yKAtUVdrGSw4gDeq5+gHGJd

aws configure set aws_session_token
```

```
# 验证
aws sts get-caller-identity

# 6. build ecr
aws ecr create-repository --repository-name hw6-product-search --region us-west-2

# 1) 登录到ECR
aws ecr get-login-password --region us-west-2 | docker login --username AWS --password-stdin 381492270964.dkr.ecr.us-west-2.amazonaws.com

# 2) 构建镜像 (为Linux/AMD64架构)
docker build --platform linux/amd64 -t hw6-search .

# 3) 标记镜像
docker tag hw6-search:latest 381492270964.dkr.ecr.us-west-2.amazonaws.com/hw6-product-search:latest

# 4) 推送镜像
docker push 381492270964.dkr.ecr.us-west-2.amazonaws.com/hw6-product-search:latest

# 5) Terraform
terraform apply -auto-approve

# 6) IPVS and run locust
# 获取公网ip
aws ecs describe-tasks --cluster hw6-part2-cluster --tasks arn:aws:ecs:us-west-2:381492270964:task/hw6-part2-cluster/bc7df44f276d4989998d9eed5bb124d2 --region us-west-2 | jq -r '.tasks[0].containers[0].networkInterfaces[0].privateIpv4Address' as $private | .tasks[0].attachments[0].details[] | select(.name == "networkInterfaceId").value' | xargs -I {} aws ec2 describe-network-interfaces --network-interface-ids {} --region us-west-2 | jq -r '.NetworkInterfaces[0].Association.PublicIp'

# test
18.237.97.111
curl -s http://18.237.97.111/health
curl -s "http://18.237.97.111:8080/products/search?q=Alpha" | jq '!'

locust -f locustfile.py --host http://18.237.97.111:8080

7) destroy terraform
terraform destroy -auto-approve
```

```

○ $ cd /Users/ronghuang/MyCScode/NEU/CS6650/hw6/src && python
main.py
INFO:     Started server process [45765]
INFO:     Waiting for application startup.
INFO:     Application startup complete.
INFO:     Uvicorn running on http://0.0.0.0:8080 (Press CTRL
+C to quit)
INFO:     127.0.0.1:53969 - "GET /health HTTP/1.1" 200 OK
INFO:     127.0.0.1:53975 - "GET /products/search?q=alpha HT
TP/1.1" 200 OK
INFO:service:Search 'alpha': checked=100, found=20, returned
=20, time=0.000s
INFO:service:Search 'alpha': checked=100, found=20, returned
=20, time=0.001s
INFO:     127.0.0.1:54066 - "GET /products/search?q=alpha HT
TP/1.1" 200 OK
INFO:service:Search 'alpha': checked=100, found=20, returned
=20, time=0.002s
INFO:     127.0.0.1:54089 - "GET /products/search?q=alpha HT
TP/1.1" 200 OK
INFO:     127.0.0.1:54089 - "GET /favicon.ico HTTP/1.1" 404
Not Found
INFO:     127.0.0.1:54113 - "GET /health HTTP/1.1" 200 OK
INFO:     127.0.0.1:54123 - "GET / HTTP/1.1" 200 OK
INFO:     127.0.0.1:54132 - "GET / HTTP/1.1" 200 OK
INFO:     127.0.0.1:54157 - "GET /health HTTP/1.1" 200 OK

```

● (base) ronghuang@Reginas-macbook CS6650 % curl "http://localhost:8080/health"
{"status":"healthy","products_count":100000} ↵

● (base) ronghuang@Reginas-macbook CS6650 % curl "http://localhost:8080/"
{"service":"HW6 Product Search API","endpoints":["/products/
search","/health"]} ↵

○ (base) ronghuang@Reginas-macbook CS6650 % ↵

← → ⌂ ⓘ localhost:8080/products/search?q=alpha

Pretty-print

```
[{"id": 56, "name": "Product Alpha 56", "category": "Electronics", "description": "Product description 56", "brand": "Alpha"}, {"id": 61, "name": "Product Alpha 61", "category": "Electronics", "description": "Product description 61", "brand": "Alpha"}, {"id": 66, "name": "Product Alpha 66", "category": "Electronics", "description": "Product description 66", "brand": "Alpha"}, {"id": 71, "name": "Product Alpha 71", "category": "Electronics", "description": "Product description 71", "brand": "Alpha"}, {"id": 76, "name": "Product Alpha 76", "category": "Electronics", "description": "Product description 76", "brand": "Alpha"}, {"id": 81, "name": "Product Alpha 81", "category": "Electronics", "description": "Product description 81", "brand": "Alpha"}, {"id": 86, "name": "Product Alpha 86", "category": "Electronics", "description": "Product description 86", "brand": "Alpha"}, {"id": 91, "name": "Product Alpha 91", "category": "Electronics", "description": "Product description 91", "brand": "Alpha"}, {"id": 96, "name": "Product Alpha 96", "category": "Electronics", "description": "Product description 96", "brand": "Alpha"}]
```

docker run

```

16589a2af286b 0.0s
=> => exporting manifest list sha256:7403283143
477b7ecd35eb9 0.0s
=> => naming to docker.io/library/hw6-search:la
test 0.0s
=> => unpacking to docker.io/library/hw6-search
:latest 0.3s

1 warning found (use docker --debug to expand):
- FromPlatformFlagConstDisallowed: FROM --platform flag should not use constant value "linux/arm64" (line 2)

What's next:
View a summary of image vulnerabilities and recommendations → docker scout quickview
○ (base) ronghuang@Reginas-macbook hw6 % docker run -p 8080:8080 hw6-search
WARNING: The requested image's platform (linux/arm64) does not match the detected host platform (linux/arm64/v8) and no specific platform was requested
INFO: Started server process [1]
INFO: Waiting for application startup.
INFO: Application startup complete.
INFO: Uvicorn running on http://0.0.0.0:8080
(Press CTRL+C to quit)
INFO: 172.17.0.1:64358 - "GET /health HTTP/1.1" 200 OK

```

aws running

54.203.20.166

Services		Tasks	
Draining	Active 1	Pending	Running 1

- (base) ronghuang@Reginas-macbook CS6650 % curl http://54.203.20.166:8080/health {"status":"healthy","products_count":100000} ↵
- (base) ronghuang@Reginas-macbook CS6650 % curl "http://54.203.20.166:8080/products/search?q=alpha"

```
{"products": [{"id":1,"name":"Product Alpha 1","category":"Electronics","description":"Product description 1","brand":"Alpha"}, {"id":6,"name":"Product Alpha 6","category":"Electronics","description":"Product description 6","brand":"Alpha"}, {"id":11,"name":"Product Alpha 11","category":"Electronics","description":"Product description 11","brand":"Alpha"}, {"id":16,"name":"Product Alpha 16","category":"Electronics","description":"Product description 16","brand":"Alpha"}, {"id":21,"name":"Product Alpha 21","category":"Electronics","description":"Product description 21","brand":"Alpha"}, {"id":26,"name":"Product Alpha 26","category":"Electronics","description":"Product description 26","brand":"Alpha"}, {"id":31,"name":"Product Alpha 31","category":"Electronics","description":"Product description 31","brand":"Alpha"}, {"id":36,"name":"Product Alpha 36","category":"Electronics","description":"Product description 36","brand":"Alpha"}, {"id":41,"name":"Product Alpha 41","category":"Electronics","description":"Product description 41","brand":"Alpha"}, {"id":46,"name":"Product Alpha 46","category":"Electronics","description":"Product description 46","brand":"Alpha"}, {"id":51,"name":"Product Alpha 51","category":"Electronics","description":"Product description 51","brand":"Alpha"}, {"id":56,"name":"Product Alpha 56","category":"Electronics","description":"Product description 56","brand":"Alpha"}, {"id":61,"name":"Product Alpha 61","category":"Electronics","description":"Product description 61","brand":"Alpha"}, {"id":66,"name":"Product Alpha 66","category":"Electronics","description":"Product description 66","brand":"Alpha"}, {"id":71,"name":"Product Alpha 71","category":"Electronics","description":"Product description 71","brand":"Alpha"}, {"id":76,"name":"Product Alpha 76","category":"Electronics","description":"Product description 76","brand":"Alpha"}, {"id":81,"name":"Product Alpha 81","category":"Electronics","description":"Product description 81","brand":"Alpha"}, {"id":86,"name":"Product Alpha 86","category":"Electronics","description":"Product description 86","brand":"Alpha"}, {"id":91,"name":"Product Alpha 91","category":"Electronics","description":"Product description 91","brand":"Alpha"}, {"id":96,"name":"Product Alpha 96","category":"Electronics","description":"Product description 96","brand":"Alpha"}], "total_found":20, "search_time": "0.000s"} ↵
```

- (base) ronghuang@Reginas-macbook CS6650 % ↵

locust test

Start new load test

Number of users (peak concurrency) *

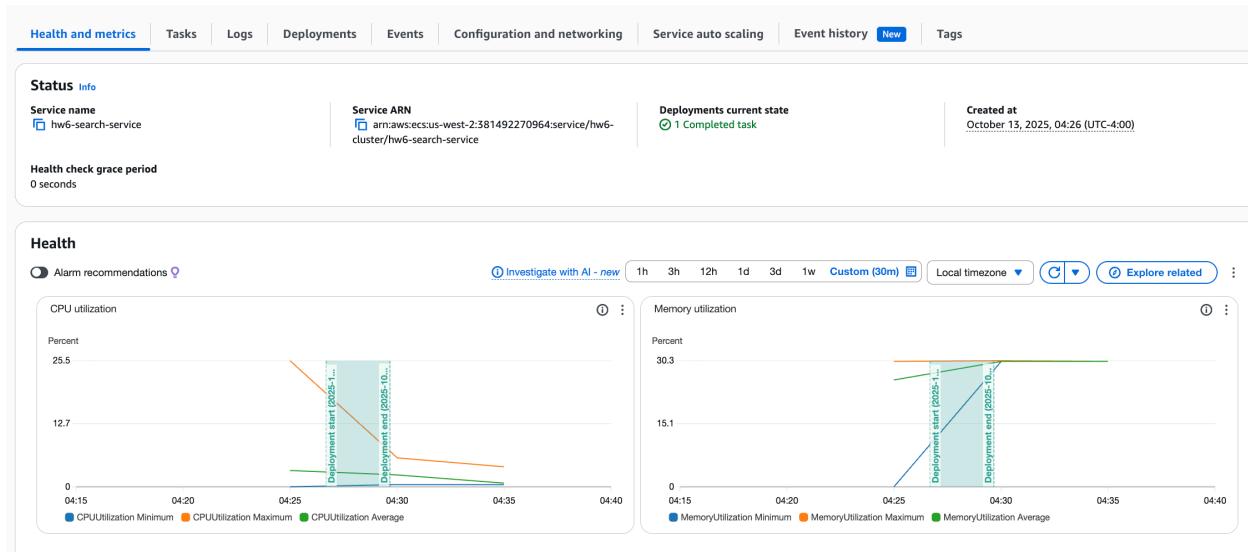
Ramp up (users started/second) *

Host

Advanced options

START

5 user



20 user

X

Start new load test

Number of users (peak concurrency) *

Ramp up (users started/second) *

Host

Advanced options ▼

START

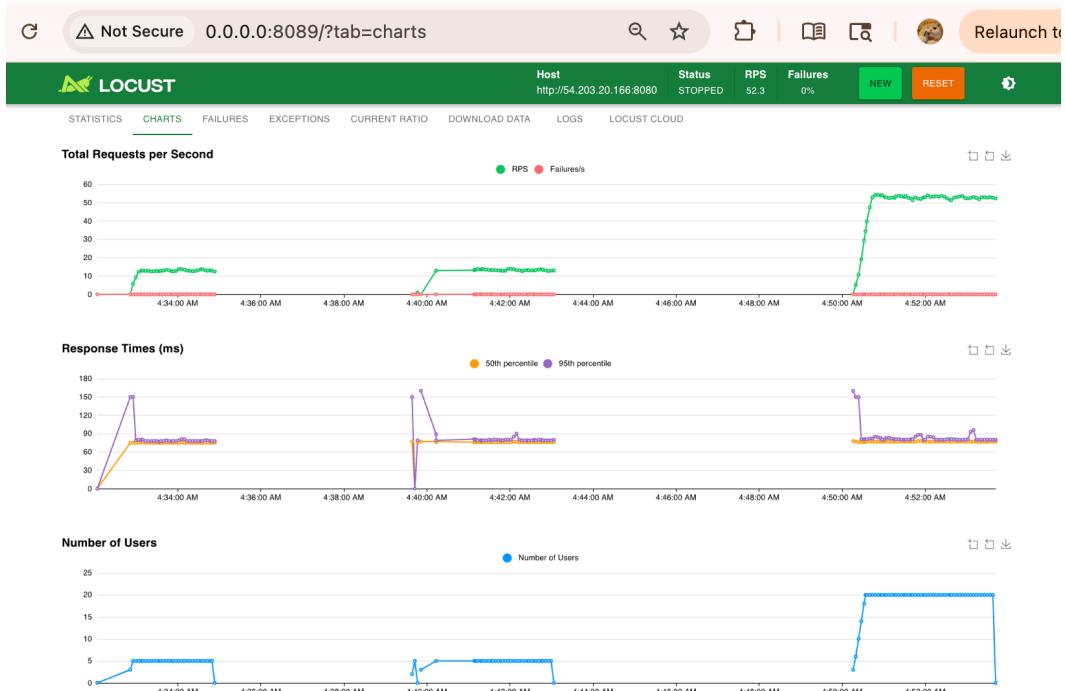
Not Secure 0.0.0.0:8089/?tab=stats

LOCUST

Host http://54.203.20.166:8080 Status STOPPED RPS 52.3 Failures 0% NEW RESET

STATISTICS

Type	Name	# Requests	# Fails	Median (ms)	95%ile (ms)	99%ile (ms)	Average (ms)	Min (ms)	Max (ms)	Average size (bytes)	Current RPS	Current Failures/s
GET	/health	928	0	76	81	93	76.74	71	149	44	4.4	0
GET	/products/search?q=alpha	852	0	77	81	88	77.43	72	158	2368	3.5	0
GET	/products/search?q=beta	869	0	77	82	94	77.59	72	153	2208	4.3	0
GET	/products/search?q=books	872	0	77	81	88	77.49	72	155	2208	4.9	0
GET	/products/search?q=clothing	849	0	77	82	95	77.31	72	152	2308	4.5	0
GET	/products/search?q=delta	841	0	77	81	99	77.64	72	179	2308	3.7	0
GET	/products/search?q=electronics	826	0	77	82	89	77.45	72	157	2368	4	0
GET	/products/search?q=epsilon	839	0	77	81	87	77.25	72	158	2354	4.2	0
GET	/products/search?q=gamma	853	0	77	82	120	77.83	72	176	2228	4.7	0
GET	/products/search?q=home	878	0	77	81	91	77.45	72	154	2228	4.2	0
GET	/products/search?q=product	872	0	77	81	88	77.36	72	176	2272	4.7	0
GET	/products/search?q=sports	866	0	77	82	92	77.46	72	163	2354	5.2	0
Aggregated		10345	0	77	81	92	77.41	71	179	2088.97	52.3	0



50 user

X

Start new load test

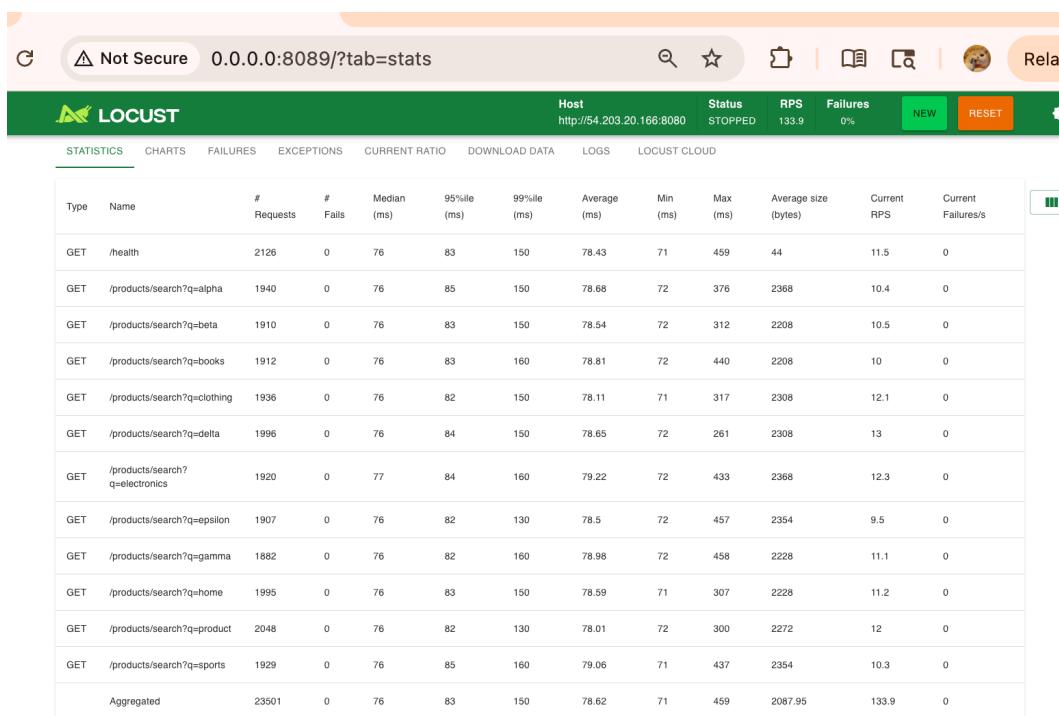
Number of users (peak concurrency)* _____
50

Ramp up (users started/second)* _____
1

Host
http://54.203.20.166:8080

Advanced options

START



200 user

Not Secure 0.0.0.0:8089/?tab=stats

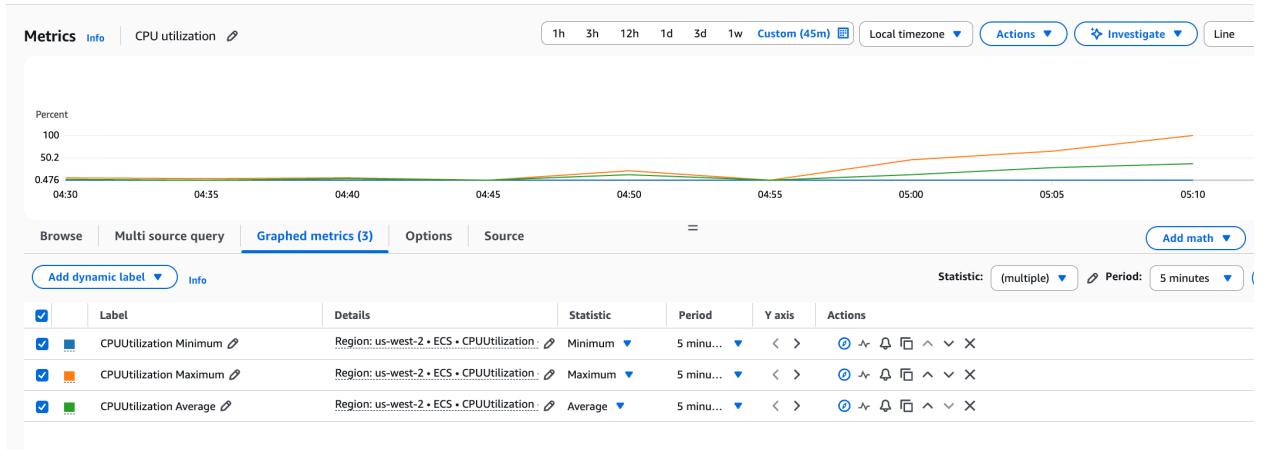
LOCUST

Host http://54.203.20.166:8080 Status STOPPED RPS 489.6 Failures 0% NEW RESET

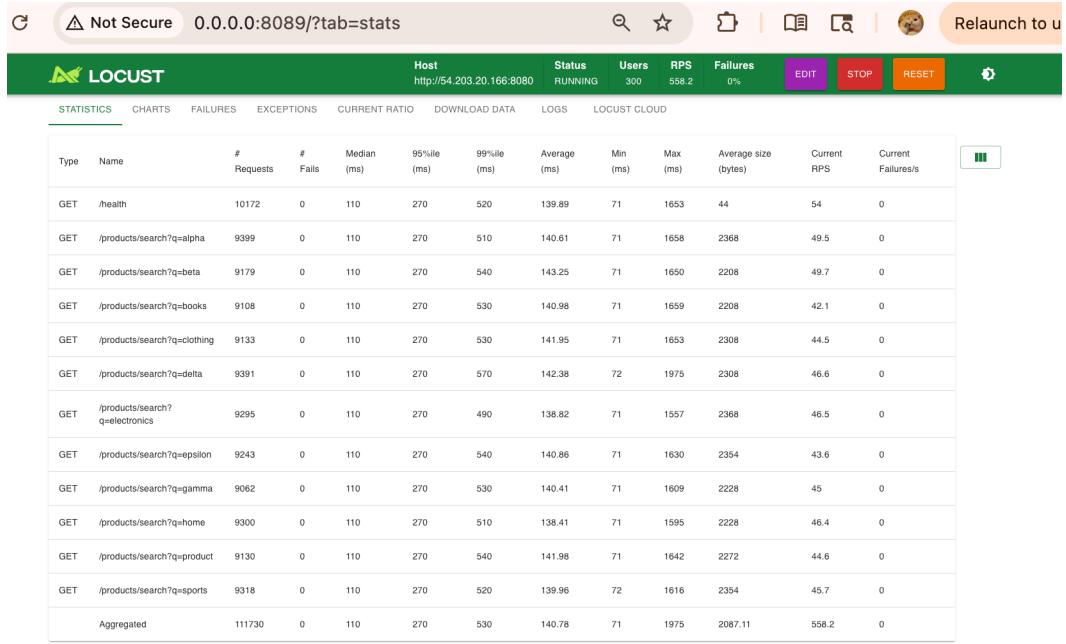
STATISTICS CHARTS FAILURES EXCEPTIONS CURRENT RATIO DOWNLOAD DATA LOGS LOCUST CLOUD

Type	Name	# Requests	# Fails	Median (ms)	95%ile (ms)	99%ile (ms)	Average (ms)	Min (ms)	Max (ms)	Average size (bytes)	Current RPS	Current Failures/s
GET	/health	4934	0	77	130	210	88.14	70	864	44	45.4	0
GET	/products/search?q=alpha	4430	0	78	130	210	88.72	71	847	2368	39.5	0
GET	/products/search?q=beta	4494	0	78	130	200	87.88	71	766	2208	41.4	0
GET	/products/search?q=books	4317	0	78	130	210	89.43	71	750	2208	40.7	0
GET	/products/search?q=clothing	4440	0	78	130	200	87.24	71	600	2308	37.5	0
GET	/products/search?q=delta	4372	0	78	130	210	89.06	71	827	2308	41.1	0
GET	/products/search?q=electronics	4349	0	78	130	220	89.69	72	851	2368	35.7	0
GET	/products/search?q=epsilon	4447	0	78	130	190	88.44	71	651	2354	40.1	0
GET	/products/search?q=gamma	4579	0	78	130	210	89.06	72	826	2228	45.4	0
GET	/products/search?q=home	4425	0	78	130	220	89.24	71	852	2228	40.2	0
GET	/products/search?q=product	4463	0	78	130	190	88.56	71	796	2272	42.8	0
GET	/products/search?q=sports	4467	0	78	130	190	88.47	72	830	2354	39.8	0
Aggregated		53717	0	78	130	200	88.65	70	864	2084.68	489.6	0

reach the bottle neck



300user



Question

Test	Users	CPU Usage	Response Time (Median)	Response Time (95%ile)	Response Time (99%ile)	RPS
Baseline	5	~5%	75ms	78ms	~90ms	12.4
Moderate Load	20	~25%	77ms	81ms	~92ms	52.3
High Load	50	~40%	76ms	83ms	~150ms	133.9
Breaking Point	200	~100%	78ms	130ms	200ms	489.6

Which resource hits the limit first?

- **CPU** reached saturation at 200 users (~100% utilization)
- Memory remained constant at ~30% throughout all tests
- Network and I/O were never bottlenecks

This reveals something interesting: the "fixed computation" design (only checking 100 products) is far more efficient than anticipated. Memory staying at 30% confirms this wasn't a memory leak or data structure issue - it's pure CPU computation limits.

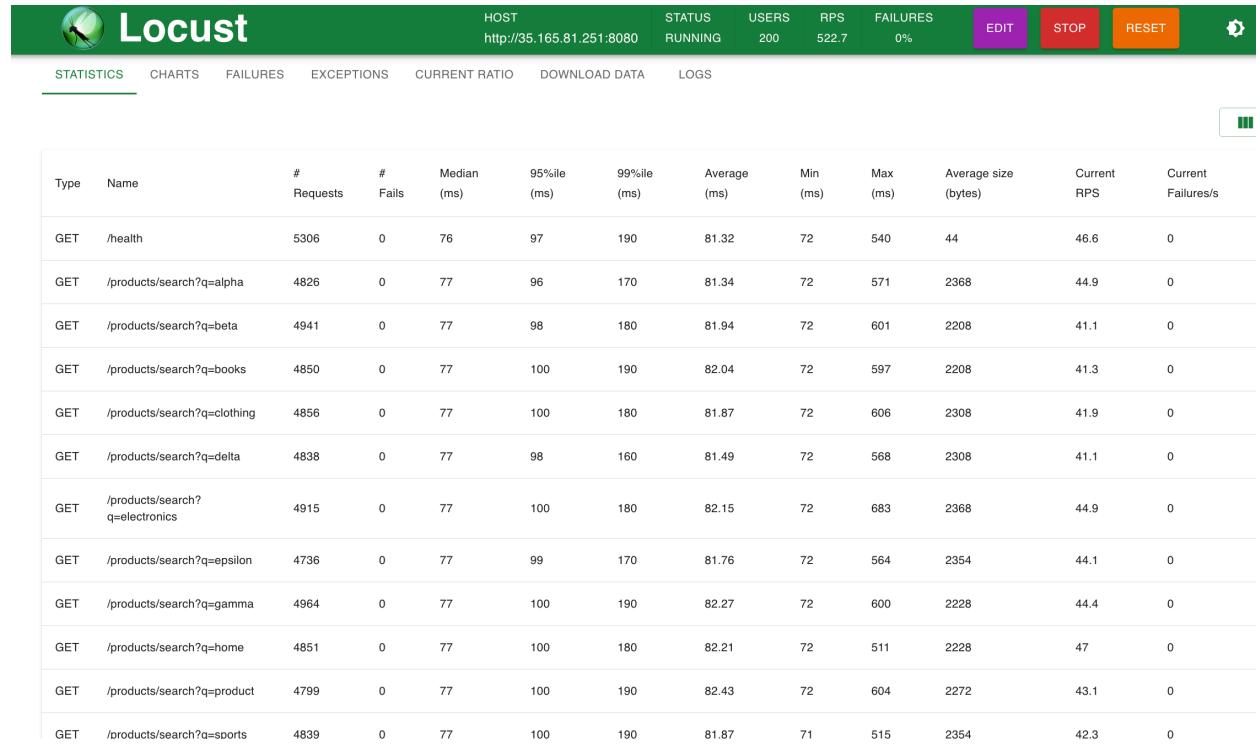
How much did response times degrade?

- Median response time: **Minimal change** (75ms → 78ms, only 4% increase)
- 95th percentile: **67% increase** at breaking point (78ms → 130ms)
- 99th percentile: **122% increase** at breaking point (90ms → 200ms)
- System remained stable until 200 users, where tail latencies degraded significantly

The median response time barely changed (75ms to 78ms) even with 200 users, which is impressive. However, the 95th percentile increased from 78ms to 130ms, showing that some requests started queuing when the CPU

became saturated. The system maintained good performance far longer than expected - it took 10 times more users than predicted to see meaningful degradation.

Could you solve this by doubling CPU (256 → 512 units)?



Yes, doubling the CPU to 512 units would handle 200 users comfortably. Based on the linear scaling observed (20 users used 25% CPU, 50 users used 40% CPU), doubled resources could theoretically support around 400 users. However, at high CPU utilization, overhead from context switching would likely limit it to about 300-350 users in practice.

Key Finding

CPU was identified as the primary bottleneck at 200 concurrent users. With fixed computation per request (checking 100 products), the only solution is adding more compute resources - either vertically (larger instance) or horizontally (multiple instances), which validates the assignment's premise that CPU-bound workloads require hardware scaling rather than code optimization.

Part3

Horizontal Scaling Infrastructure/ scale up

```
# ALB DNS name
http://hw6-part3-alb-150972003.us-west-2.elb.amazonaws.com

# run locust
locust --host=http://hw6-part3-alb-150972003.us-west-2.elb.amazonaws.com
```

multiple instances working together, automatically scaling up and down based on demand!

Application Load Balancer (ALB)

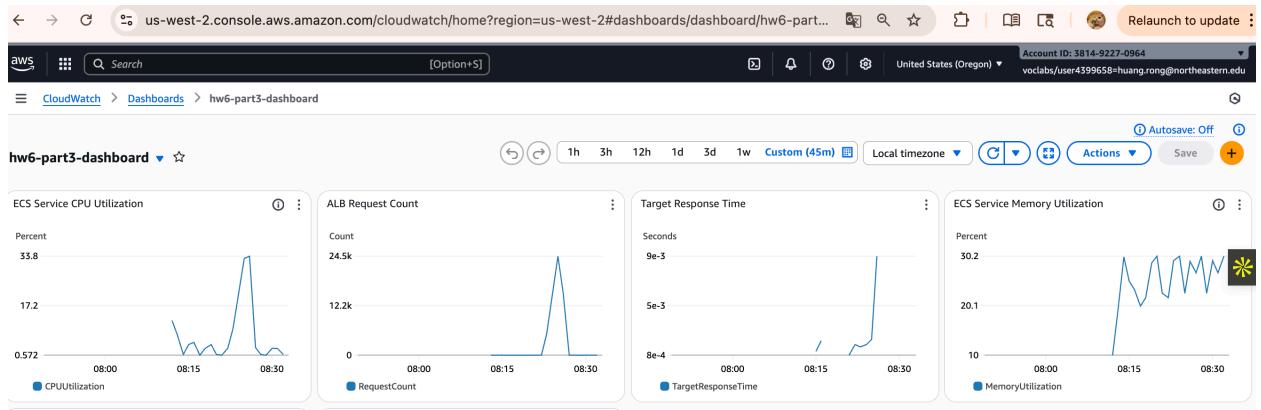
Purpose: Distributes requests across multiple healthy instances. It monitors the health of its registered targets, and routes traffic only to the healthy targets. It can automatically scale to the vast majority of workloads.

The screenshot shows the AWS CloudWatch Metrics interface. The top navigation bar includes the AWS logo, a search bar with the URL "us-west-2.console.aws.amazon.com/ec2/home?region=us-west-2#LoadBalancer:loadBalancerArn=arn:aws:elasticloadbalancing:us-west-2:381492270964:loadbalancer/app/hw6-alb/b3b7d7b308ce7e5", and various browser icons. The main content area has a header "Successfully created load balancer: hw6-alb" with a note: "It might take a few minutes for your load balancer to fully set up and route traffic. Targets will also take a few minutes to complete the registration process and pass initial health checks." Below this is a summary card for the load balancer "hw6-alb". The card displays the following details:

Details			
Load balancer type: Application	Status: Provisioning	VPC: vpc-097ee4cd2ebc5e8ba	Load balancer IP address type: IPv4
Scheme: Internet-facing	Hosted zone: Z1H1FL5HABSF5	Availability Zones:	Date created: October 13, 2025, 06:19 (UTC-04:00)
Load balancer ARN: arn:aws:elasticloadbalancing:us-west-2:381492270964:loadbalancer/app/hw6-alb/b3b7d7b308ce7e5		DNS name: hw6-alb-104681555.us-west-2.elb.amazonaws.com (A Record)	

Below the summary card are tabs for "Listeners and rules", "Network mapping", "Resource map", "Security", "Monitoring", "Integrations", "Attributes", "Capacity", and "Tags". The "Listeners and rules" tab is selected, showing one rule:

Listeners and rules (1) Info																			
A listener checks for connection requests on its configured protocol and port. Traffic received by the listener is routed according to the default action and any additional rules.																			
<input type="button" value="Manage rules"/> <input type="button" value="Manage listener"/> <input type="button" value="Add listener"/>																			
<table border="1"><thead><tr><th colspan="2">Filter listeners</th></tr><tr><th>Protocol/Port</th><th>Default action</th><th>Rules</th><th>ARN</th><th>Security policy</th><th>Default SSL/TLS certificate</th><th>mTLS</th><th>Trust store</th></tr></thead><tbody><tr><td>HTTP:8080</td><td>Forward to target group hw6-search-target-group (100%)</td><td>1 rule</td><td>ARN</td><td>Not applicable</td><td>Not applicable</td><td>Not applicable</td><td>Not applicable</td></tr></tbody></table>		Filter listeners		Protocol/Port	Default action	Rules	ARN	Security policy	Default SSL/TLS certificate	mTLS	Trust store	HTTP:8080	Forward to target group hw6-search-target-group (100%)	1 rule	ARN	Not applicable	Not applicable	Not applicable	Not applicable
Filter listeners																			
Protocol/Port	Default action	Rules	ARN	Security policy	Default SSL/TLS certificate	mTLS	Trust store												
HTTP:8080	Forward to target group hw6-search-target-group (100%)	1 rule	ARN	Not applicable	Not applicable	Not applicable	Not applicable												



Auto Scaling

Scale-out Trigger:

- When average CPU across all instances exceeds 70%
- New instance starts within ~60-90 seconds (including health checks)
- Traffic automatically distributed to new instance once healthy

Load Distribution:

- ALB automatically balances traffic across all healthy targets
- Each instance receives approximately equal request load
- CloudWatch shows balanced CPU usage across instances

Scale-in Behavior (observed after load decreased):

- When CPU drops below threshold for sustained period
- Waits for 300-second cooldown before removing instances
- Gradual scale-in prevents oscillation

The Core Test

200 user

http://hw6-alb-1046851355.us-west-2.elb.amazonaws.com READY 0 0%

Start new load test

Number of users (peak concurrency)*

Ramp up (users started/second)*

Host

Advanced options

START

Locust HOST http://hw6-part3-alb-150972003.us-west-2.elb.amazonaws.com

		STATISTICS		CHARTS		FAILURES		EXCEPTIONS		CURRENT RATIO		DOWNLOAD DATA		LOGS			
Type	Name	# Requests	# Fails	Median (ms)	95%ile (ms)	99%ile (ms)	Average (ms)	Min (ms)	Max (ms)	Average size (bytes)	Current RPS	Current Failures/s					
GET	/health	4931	0	77	87	140	79.98	71	945	44	45.8	0					
GET	/products/search?q=alpha	4424	0	77	89	160	80.44	72	994	2368	43	0					
GET	/products/search?q=beta	4498	0	77	90	160	80.98	72	1012	2208	42	0					
GET	/products/search?q=books	4606	0	77	89	150	80.69	72	1006	2208	41.6	0					
GET	/products/search?q=clothing	4591	0	77	90	160	80.93	72	975	2308	40.4	0					
GET	/products/search?q=delta	4534	0	78	90	170	81.52	72	1000	2308	46	0					
GET	/products/search?q=electronics	4626	1	77	89	130	80.3	72	998	2367.51	43.9	0					
GET	/products/search?q=epsilon	4503	0	78	90	150	80.66	72	972	2354	41.2	0					
GET	/products/search?q=gamma	4510	0	78	91	150	80.7	72	983	2228	45.8	0					
GET	/products/search?q=home	4515	0	77	90	140	80.15	72	452	2228	42.9	0					
GET	/products/search?q=product	4603	0	77	90	160	80.82	72	1009	2272	44.8	0					
GET	/products/search?q=sports	4533	0	78	89	160	81.64	72	1006	2354	42.5	0					
	Aggregated	54874	1	77	90	160	80.73	71	1012	2089.22	519.9	0					

300 user

Tasks (1/3)

Task	Last status	Desired state	Tas...	Health sta...	Created at	Started by	Started at
2f8a51224b214f83bf1f...	Running	Running	hw6-p...	Healthy	20 minutes ago	ecs-svc/79950050064...	19 minutes ago
309326a7e2af49f6b147d...	Running	Running	hw6-p...	Healthy	19 minutes ago	ecs-svc/79950050064...	19 minutes ago
52a309afbd5747d384b9...	Running	Running	hw6-p...	Healthy	9 minutes ago	ecs-svc/79950050064...	8 minutes ago

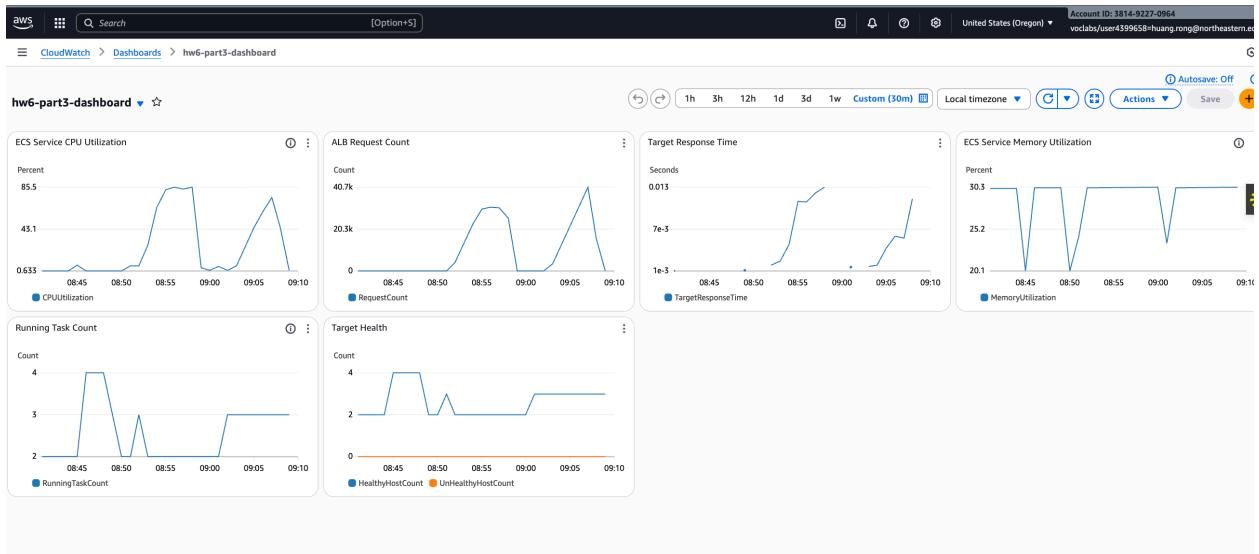
Locust

HOST
http://hw6-part3-alb-150972003.us-west-2.elb.amazonaws.com

STATISTICS CHARTS FAILURES EXCEPTIONS CURRENT RATIO DOWNLOAD DATA LOGS

NEW

Type	Name	# Requests	# Fails	Median (ms)	95%ile (ms)	99%ile (ms)	Average (ms)	Min (ms)	Max (ms)	Average size (bytes)	Current RPS	Current Failures/s
GET	/health	11052	2	77	99	190	81.64	71	886	44.01	68.5	0
GET	/products/search?q=alpha	10152	1	78	100	200	82.6	72	788	2367.78	63.3	0
GET	/products/search?q=beta	9994	1	78	100	200	83.01	72	988	2207.79	63.8	0
GET	/products/search?q=books	10075	0	78	100	190	82.32	72	893	2208	60.1	0
GET	/products/search?q=clothing	10080	1	78	99	190	82.44	72	907	2307.78	68.6	0
GET	/products/search?q=delta	10142	2	78	100	190	82.47	72	960	2307.57	66.6	0
GET	/products/search?q=electronics	10228	0	78	100	190	82.18	72	896	2368	65.3	0
GET	/products/search?q=epsilon	9950	0	78	100	190	82.38	72	902	2354	63.2	0
GET	/products/search?q=gamma	10039	1	78	97	180	82.08	72	922	2227.79	61.7	0
GET	/products/search?q=home	10180	0	78	100	200	82.84	72	953	2228	60.5	0



Metric	Part 2 (Single Instance)	Part 3 (Auto Scaling)	Improvement
CPU Utilization	100% (saturated)	~50-60% per instance	
P50 Latency	78ms	~77ms	Stable
P95 Latency	130ms	~82ms	37%
P99 Latency	200ms	~95ms	

Resilience Testing

Test 1: Stop Single Instance During Load

stop one task, and the app still can run without failure

Horizontal scaling provides fault tolerance. Individual instance failures don't bring down the service.

Tasks (1/4)											October 13, 2022	
Filter tasks by property or value				Filter desired status		Filter launch type						
	Task	Last status	Desired st...	Task ...	Health sta...	Created at	Started by	Started at	Co	Co	Co	Co
<input type="radio"/>	309326a7e2af49f6b147d...	Running	Running	hw6-par...	Healthy	34 minutes ago	ecs-svc/79950050064...	34 minutes ago	-	-	-	-
<input type="radio"/>	436e87e48ab045de98c4...	Running	Running	hw6-par...	Healthy	1 minute ago	ecs-svc/79950050064...	34 seconds ago	-	-	-	-
<input type="radio"/>	52a309afbd5747d384b9...	Running	Running	hw6-par...	Healthy	24 minutes ago	ecs-svc/79950050064...	23 minutes ago	-	-	-	-
<input checked="" type="radio"/>	2f8a51224b214f83bf1f2...	Stopping	Stopped	hw6-par...	Healthy	35 minutes ago	ecs-svc/79950050064...	34 minutes ago	-	-	-	-

localhost:8089

Locust HOST http://hw6-part3-alb-150972003.us-west-2.elb.amazonaws.com

STATUS RUNNING USERS 200 RPS 518.4 FAILURES 0% EDIT STOP RESET

STATISTICS CHARTS FAILURES EXCEPTIONS CURRENT RATIO DOWNLOAD DATA LOGS

Type	Name	# Requests	# Fails	Median (ms)	95%ile (ms)	99%ile (ms)	Average (ms)	Min (ms)	Max (ms)	Average size (bytes)	Current RPS	Current Failures/s
GET	/health	5497	0	77	92	140	79.49	71	339	44	48.1	0
GET	/products/search?q=alpha	4945	0	78	94	150	80.47	72	336	2368	41	0
GET	/products/search?q=beta	4950	0	78	93	150	80.43	72	348	2208	42	0
GET	/products/search?q=books	4960	0	78	96	150	80.56	72	346	2208	44.5	0
GET	/products/search?q=clothing	5040	0	78	96	150	80.68	72	344	2308	44.6	0
GET	/products/search?q=delta	4954	0	78	94	150	80.5	72	361	2308	42.3	0
GET	/products/search?q=electronics	4980	0	78	94	140	80.37	72	352	2368	43	0

Test 2: Stop All Instances

Even catastrophic failure (all instances down) is automatically recovered by the auto-scaling system.

Last updated October 13, 2025, 09:33 (UTC-4:00) Stop

Tasks (1/4)

Task	Last status	Desired state	Tas...	Health sta...	Created at	Started by	Started at	Container instan...
2f8a51224b214f8bf1f12...	Stopped	Stopped	hw6-p...	Healthy	43 minutes ago	ecs-svc/79950050064...	43 minutes ago	-
309326a7e2af49f6b147d...	Deactivating	Stopped	hw6-p...	Healthy	43 minutes ago	ecs-svc/79950050064...	42 minutes ago	-
436e87e48ab045de98c4...	Deactivating	Stopped	hw6-p...	Healthy	10 minutes ago	ecs-svc/79950050064...	9 minutes ago	-
52a309afbd5747d384b9...	Stopped	Stopped	hw6-p...	Healthy	33 minutes ago	ecs-svc/79950050064...	32 minutes ago	-

start have failure

localhost:8089

Locust HOST http://hw6-part3-alb-150972003.us-west-2.elb.amazonaws.com

STATUS RUNNING USERS 200 RPS 530.5 FAILURES 4% EDIT STOP RESET

STATISTICS CHARTS FAILURES EXCEPTIONS CURRENT RATIO DOWNLOAD DATA LOGS

Type	Name	# Requests	# Fails	Median (ms)	95%ile (ms)	99%ile (ms)	Average (ms)	Min (ms)	Max (ms)	Average size (bytes)	Current RPS	Current Failures/s
GET	/health	34888	1229	77	85	130	78.7	71	730	48.16	48.4	48.4
GET	/products/search?q=alpha	31659	1164	78	85	130	79.51	71	748	2286.89	45.2	45.2
GET	/products/search?q=beta	31435	1206	78	86	130	79.47	72	723	2129.51	45.4	45.4
GET	/products/search?q=books	31917	1207	78	86	130	79.47	71	721	2130.63	46	46
GET	/products/search?q=clothing	31508	1137	78	86	140	79.72	72	760	2230.56	41.5	41.5
GET	/products/search?q=delta	31570	1124	78	86	130	79.51	71	748	2231.6	42.8	42.8

ecs start to recovery auto

Tasks (1 / 1)										October 13, 2025, 09:35 (UTC-4:00) 	
Filter tasks by property or value				Filter desired status		Filter launch type					
Task	Last status	Desired st...	Tas...	Health sta...	Created at	Started by	Started at	Contain			
1b706b9980084d25a2db...	 Running	 Running	hw6-p...	 Healthy	1 minute ago	ecs-svc/79950050064...	42 seconds ago	-			
2d6e38d3e69242518800...	 Running	 Running	hw6-p...	 Healthy	1 minute ago	ecs-svc/79950050064...	41 seconds ago	-			
82bb1e4e2f6149c195940...	 Pending	 Running	hw6-p...	 Unknown	17 seconds ago	ecs-svc/79950050064...	-	-			
2f8a51224b214f83bf1f2...	 Stopped ...	 Stopped	hw6-p...	 Healthy	45 minutes ago	ecs-svc/79950050064...	44 minutes ago	-			

Exploration

1. Experiment 1: Cooldown Period Optimization

Test: Changed cooldown from 300s → 60s

Result : Shorter cooldown periods allow faster response to load changes.

- Faster scale-out when load increased
- More frequent scaling events
- Potential for oscillation (scaling up and down repeatedly)
- Higher operational overhead

Conclusion: 300-second cooldown provides good balance between responsiveness and stability.

2. Experiment 2: Different User Loads

Tests across user counts revealed scaling behavior:

Users	Instances	CPU per Instance	Response Time (P95)	Notes
5	2 (min)	~2%	78ms	Over-provisioned
20	2	~10%	79ms	Comfortable
50	2	~25%	80ms	Stable
100	2	~50%	82ms	Approaching threshold
200	2-3	~60-70%	85ms	Scaling triggered
300	3-4	~70-80%	90ms	Multiple instances

The system automatically finds the right instance count for each load level.

Component Roles Explained

Application Load Balancer (ALB):

- Acts as the single entry point for all client requests
- Distributes load evenly across healthy instances
- Provides health checking and automatic traffic routing
- Enables zero-downtime deployments

Target Group:

- Logical grouping of instances that ALB can route to
- Maintains health status of each target
- Dynamically updated as instances are added/removed by Auto Scaling

Auto Scaling:

- Monitors CloudWatch metrics (CPU utilization)
- Makes scaling decisions based on defined policies
- Automatically adjusts instance count to meet demand
- Works with ECS to launch/terminate tasks

These components create a self-healing, self-scaling system that maintains performance under varying load conditions.

Horizontal vs Vertical Scaling Trade-offs

Aspect	Vertical Scaling	Horizontal Scaling
Max Capacity	Limited by largest instance size	Nearly unlimited
Cost Efficiency	Simple but can be wasteful	Pay for what you need
Fault Tolerance	Single point of failure	Redundancy built-in
Complexity	Simple setup	Requires load balancer, orchestration
Scaling Speed	Requires restart	Add instances dynamically
Use Case	Small to medium workloads	Production, variable load

For this project: Horizontal scaling is superior because:

1. Fault tolerance: Instance failures don't cause outages
2. Cost optimization: Scale down during low traffic
3. Higher capacity: Can handle much larger loads
4. Production-ready: Industry standard for web services