

Auto Scaling Introduction

Auto Scaling helps you in maintaining correct number of EC2 instances according to application workload. We can create collections of instances called Auto Scaling Groups.

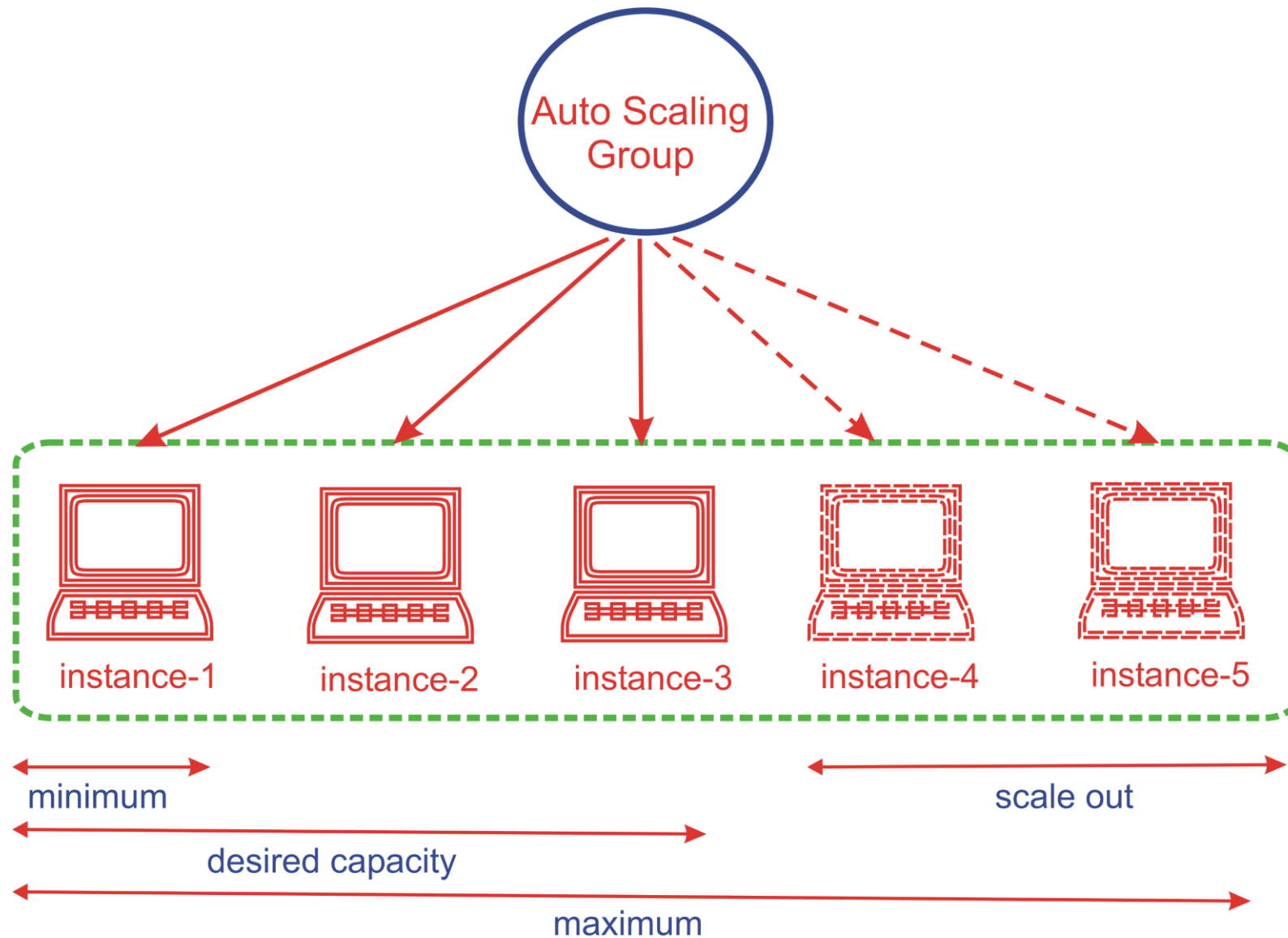
We can specify minimum number of instances in AS Group and AS will ensure that instances do not go below that.

We can specify maximum number of instances in AS Group and AS will ensure that instances do not go above that.

If we specify desired capacity then AS ensures that you always have fixed number of instances.

If we specify scaling policies then AS launches/terminates instances when application load increases/decreases.

Auto Scaling Introduction



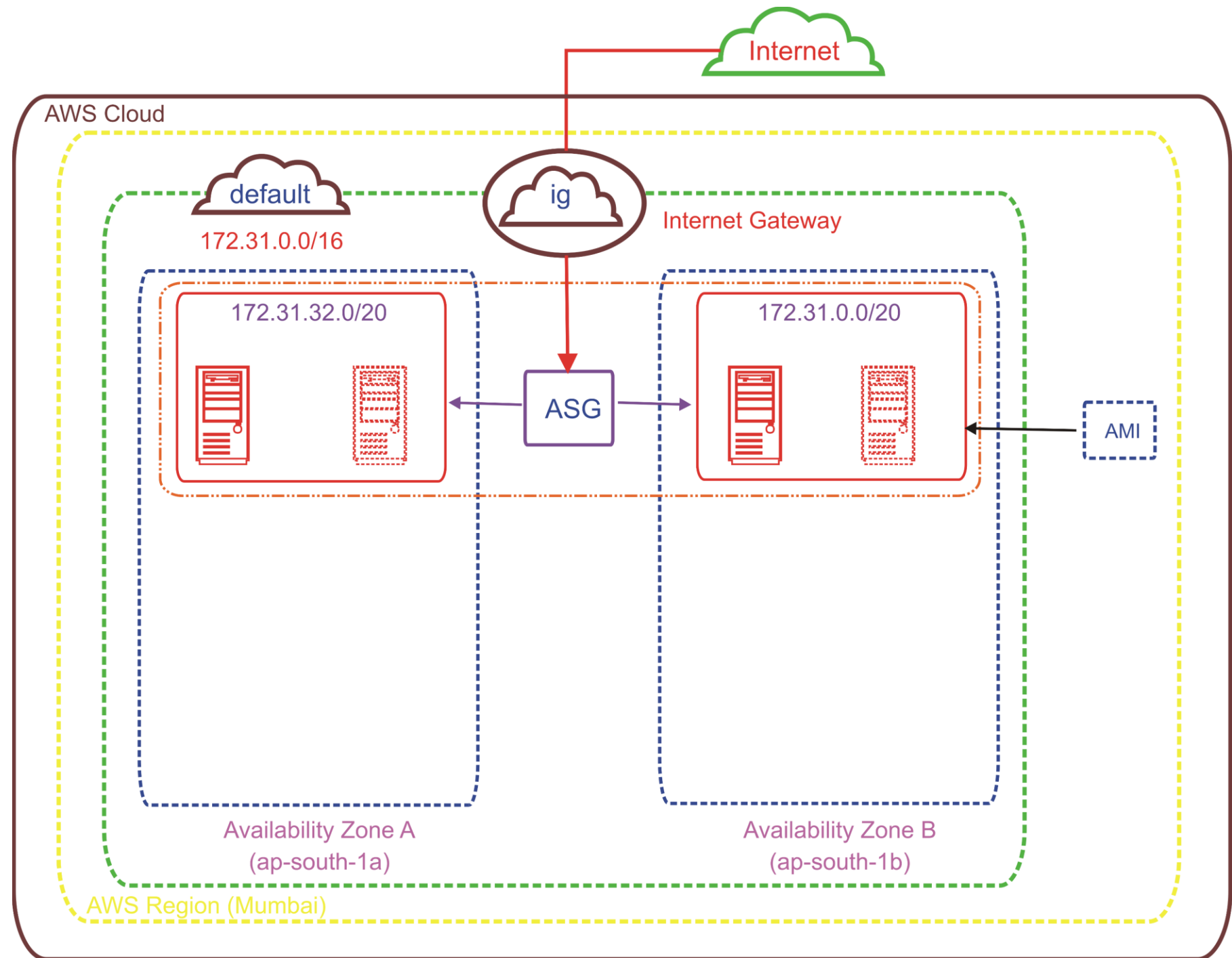
Scaling Plans

- ~~Maintain current number of instances~~
- Manual scaling
- Scale based on demand
- Scale based on schedule



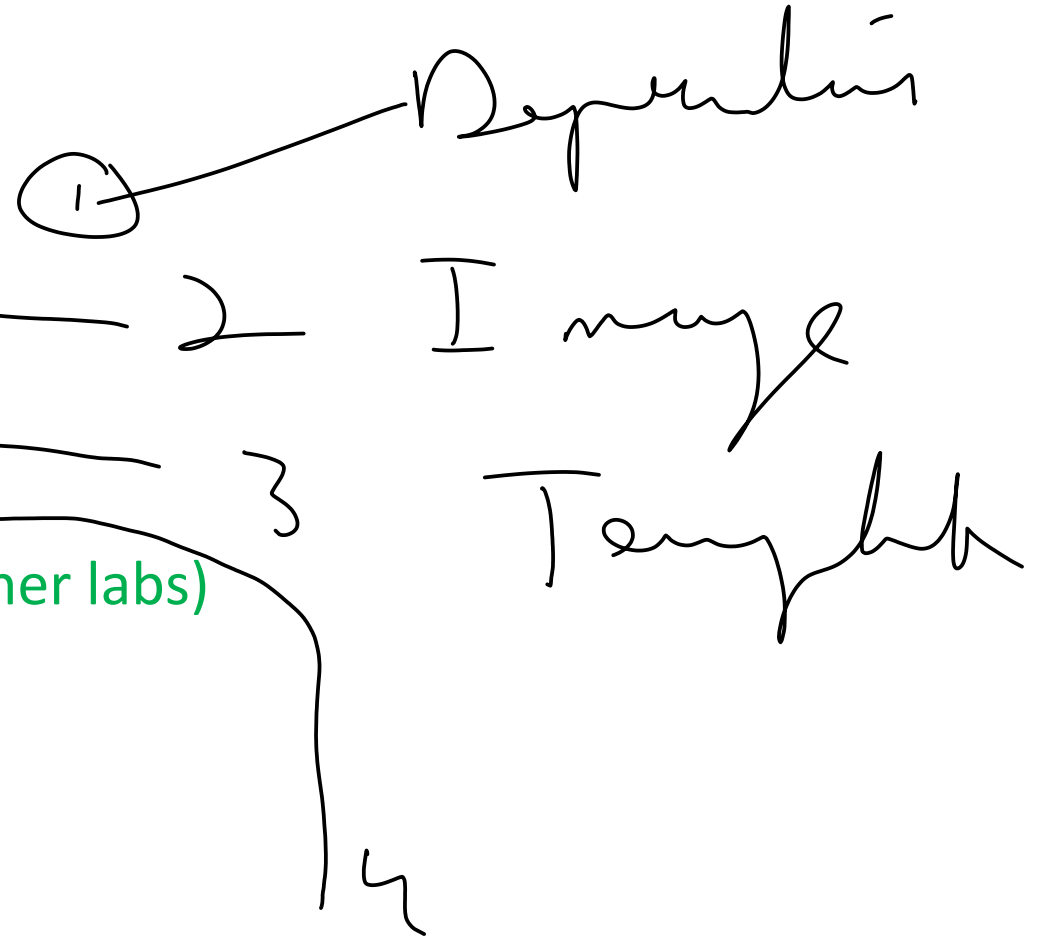
Auto Scaling Lab Setup

Default
VPC



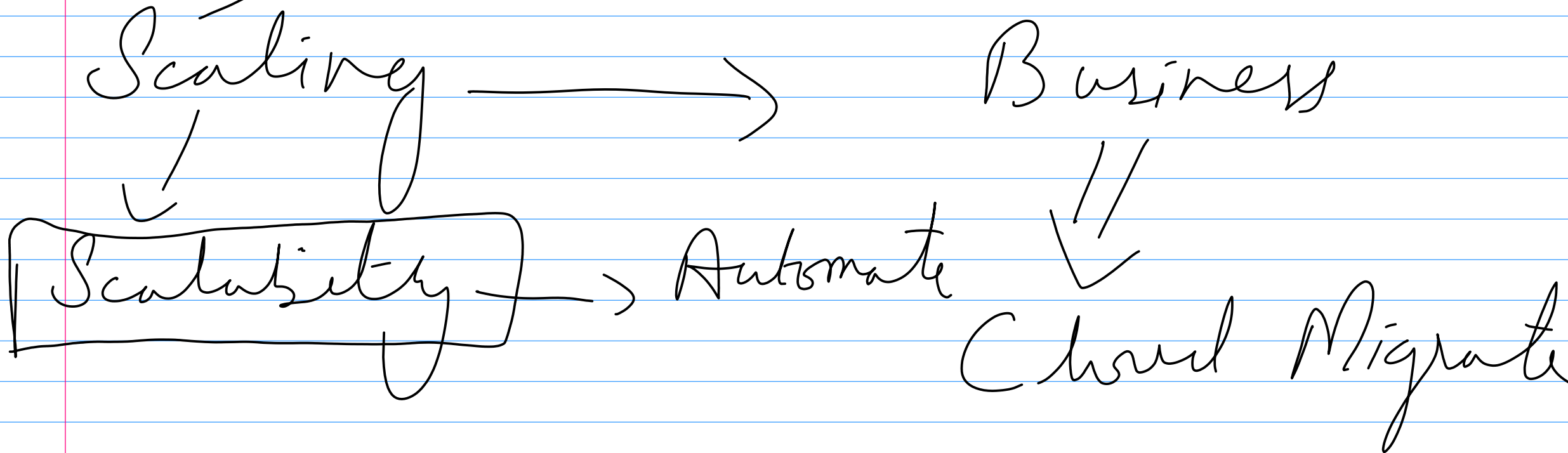
Various Steps

- Create EC2 Webserver Instance with user data
- Create Custom AMI From Instance
- Create Launch Template
- Create Auto Scaling Group
- Test Auto Scaling Group
- Delete Resources (Which are not required in other labs)



Week - 8-a

AutoScaling ∴ (Theory)



Scalability :-

Adding up
Remove

} Resource

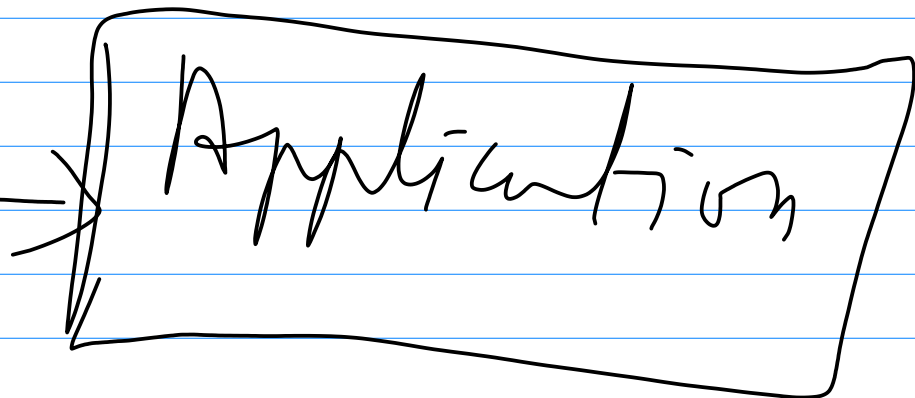
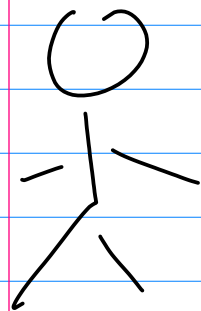
—> as per requirement

Scale up
Scale out

Scale Down
Scale in

EC-2

User



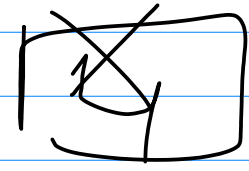
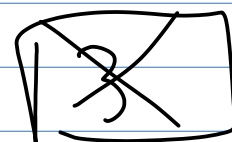
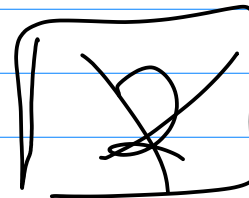
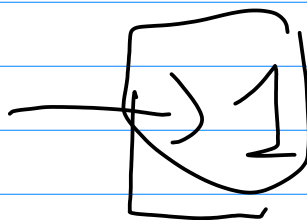
10000 hit/sec

10000000 = ?

1500,

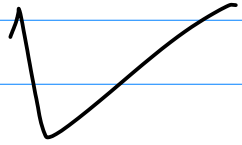
Response

25%

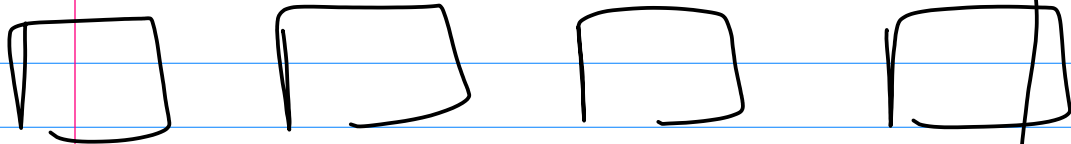


EC2

Scalability

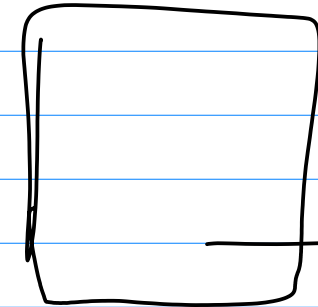


Horizontal



of EC2
instances

Vertical



power

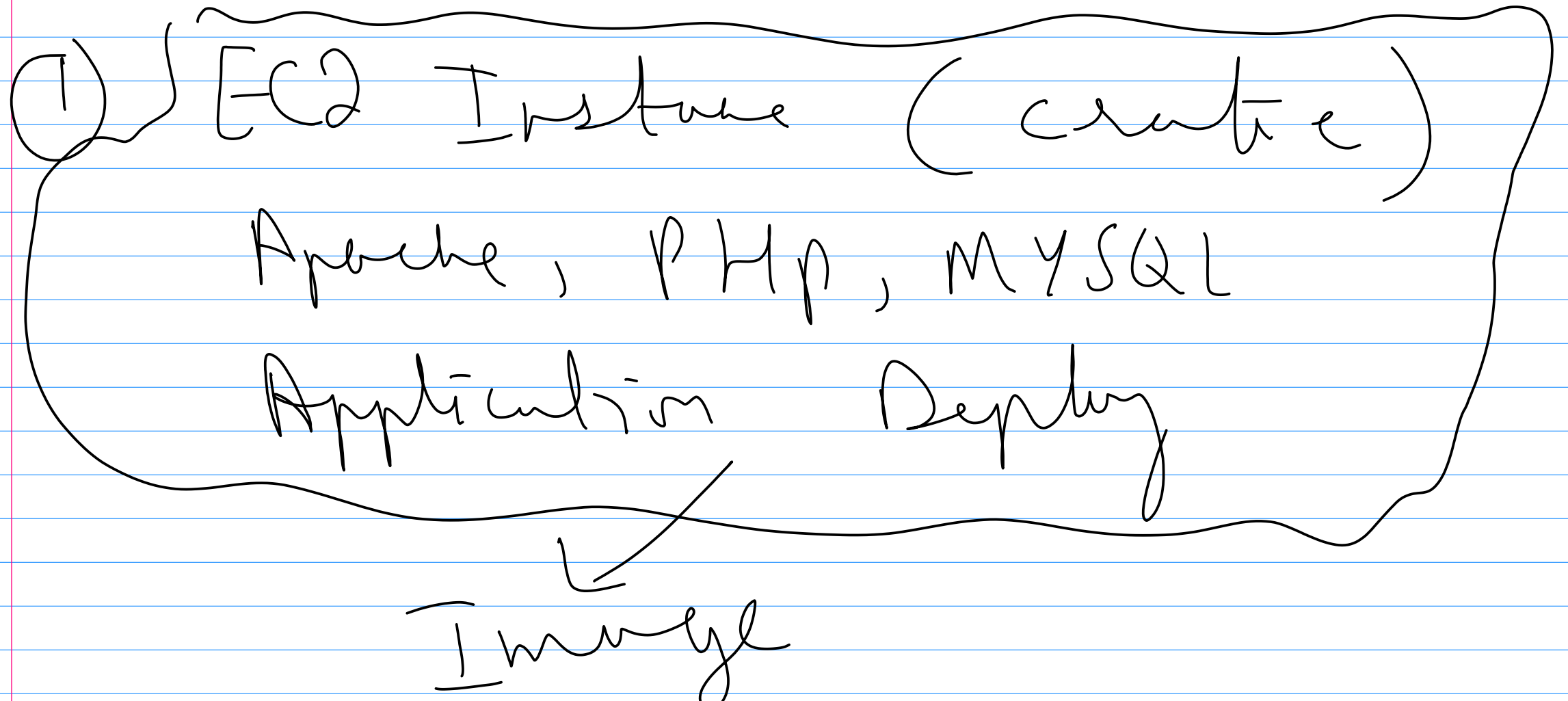
AutoScaling Components

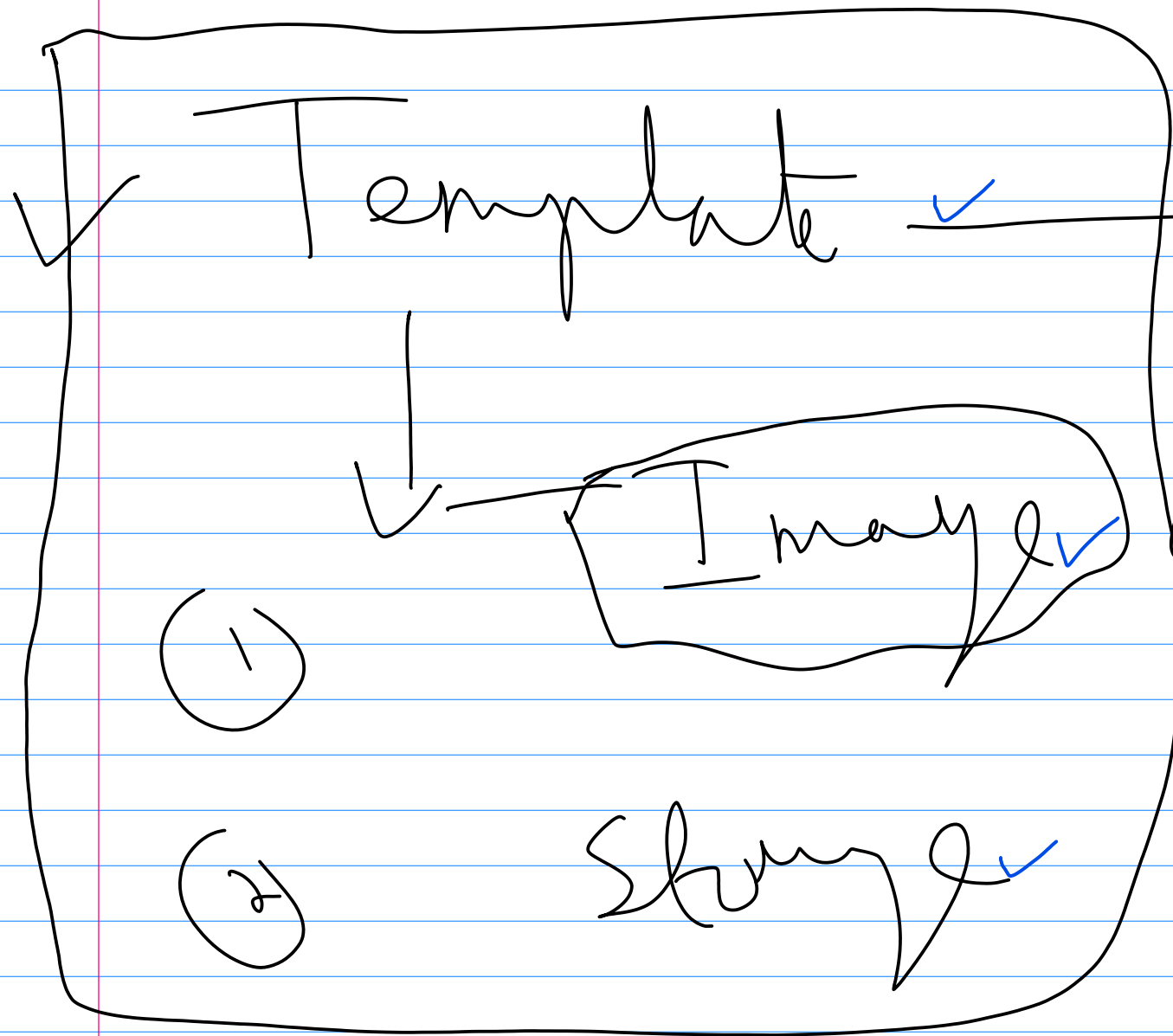
① Launch Configuration
AMI \longrightarrow Custom

② AutoScaling Group

③ Scaling \checkmark Policy \checkmark

AMI →

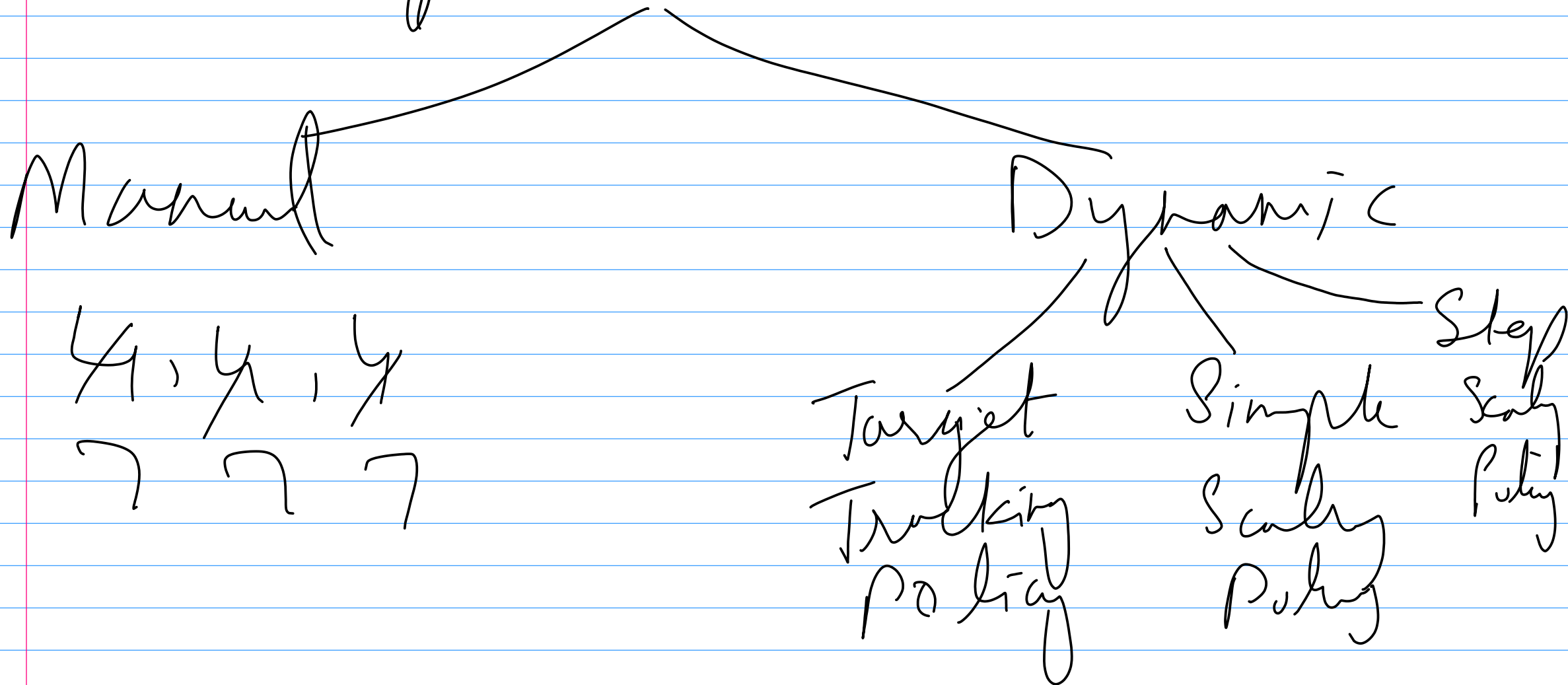




21

Auto Story
Group

3) Scaling Policies



Dynamic Scaling

1) Target-Scaling Policy

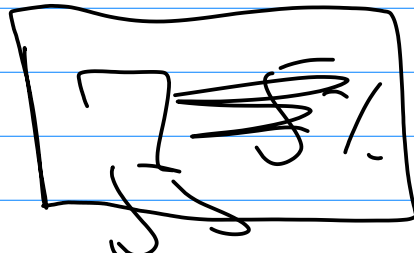
CPU Utilization

→ 80%

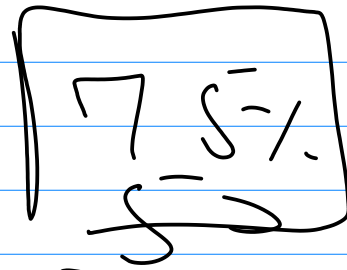
75-1.3



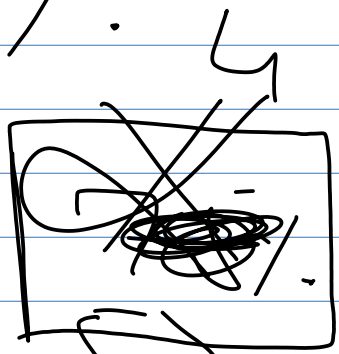
20



20



20



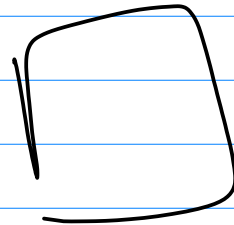
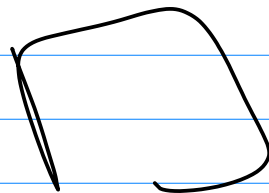
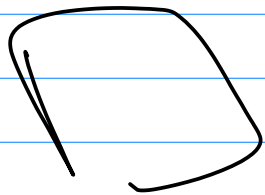
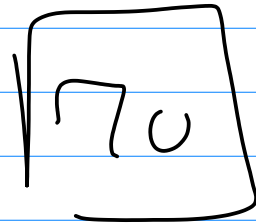
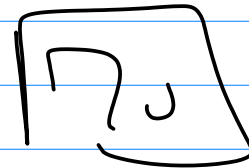
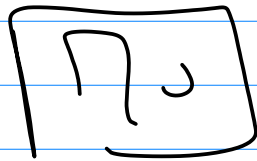
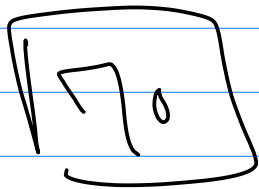
20



2) Simple Scaling

CPU Utilization

70%



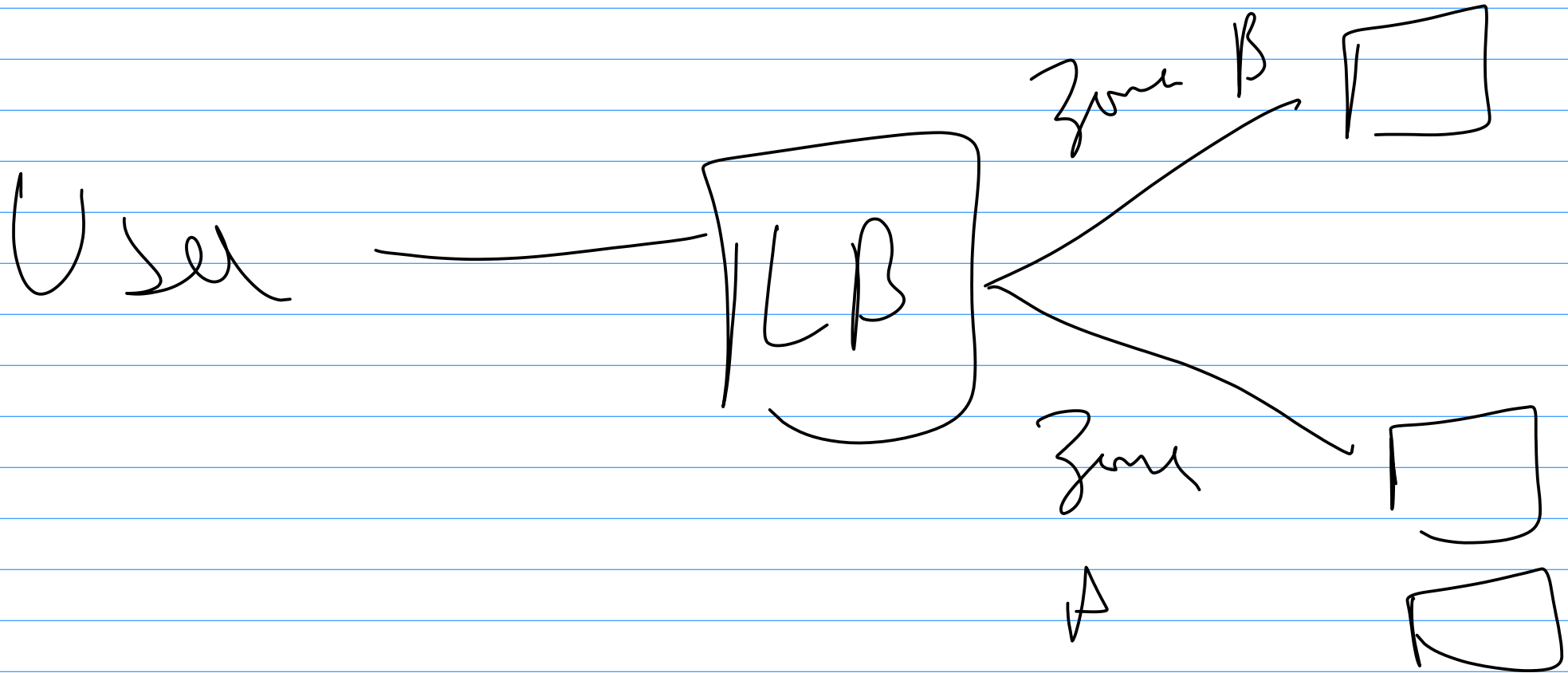
3) Step Scaling Policy ✓

Step 1 — 10 — 20 $\frac{CPU}{\%}$
1 instance

Step — 2 — 30 — 40
2 Instances

5 Instances 60 80

Load Balancing



Warm Up ~~4th~~ Period
Cool Down Period

Warm 300 Sec → 5 Min

Up

Lunch

—

