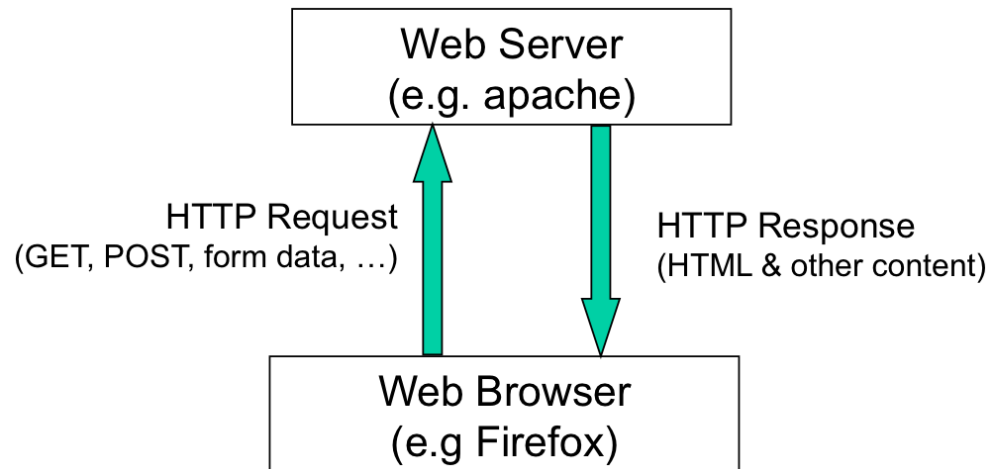# Developer Operations

## Performance & Scaling

# HTTP dialogue

- Communication on the web means browsers talking to web servers over HTTP

```
        ┌─────────────────┐
        │   Web Server    │
        │  (e.g. apache)  │
        └─────────────────┘
           ↑           │
           │           ↓
HTTP Request          HTTP Response
(GET, POST, form data, …)  (HTML & other content)
           │           ↑
        ┌─────────────────┐
        │   Web Browser   │
        │   (e.g Firefox) │
        └─────────────────┘
```

- Within a browsing session the user is (ideally) given an illusion that they are having a one-to-one dialogue with the server

- Part of this illusion is management of **state**

- Another part is achieved by **performance** and **scalability** management

2

# Performance issues

- ## Response time
  - How quickly does a server respond to client requests?
  - Affected by
    - Server load
    - Network load
    - Reliance on third-party services (e.g. DNS)
    - Delays in all software components (server & client OS, process switching, server & client network stack)

# Performance issues

- **Throughput**
  - Rate at which computational work is done
  - Transactions per unit time
  - Affected by
    - Processor speed
    - Storage performance (memory/disk)
    - Some resources overloading badly (e.g. network bottlenecks)
    - Reliance on OS and network services
    - Reliance on third-party services

# Issues related to performance

- **Reliability**
  - Uptime, ideally 99.999...%
  - Consistency

- **Fault tolerance**
  - Achieved through redundancy of components, etc

- **Adaptability**

- **Security**
  - Can add to load on processor, network, …

- **Quality of Service / Quality of Experience**
  - Combination of performance, availability, adaptability, security

# How to improve performance

- Better network speeds

- Faster CPU

- More memory

- More disk space

- Less process switching

- Better design – less bottlenecks, better concurrency (e.g. less resource holding)


- All of the above can help performance of individual servers, but don't **scale** indefinitely due to physical limits
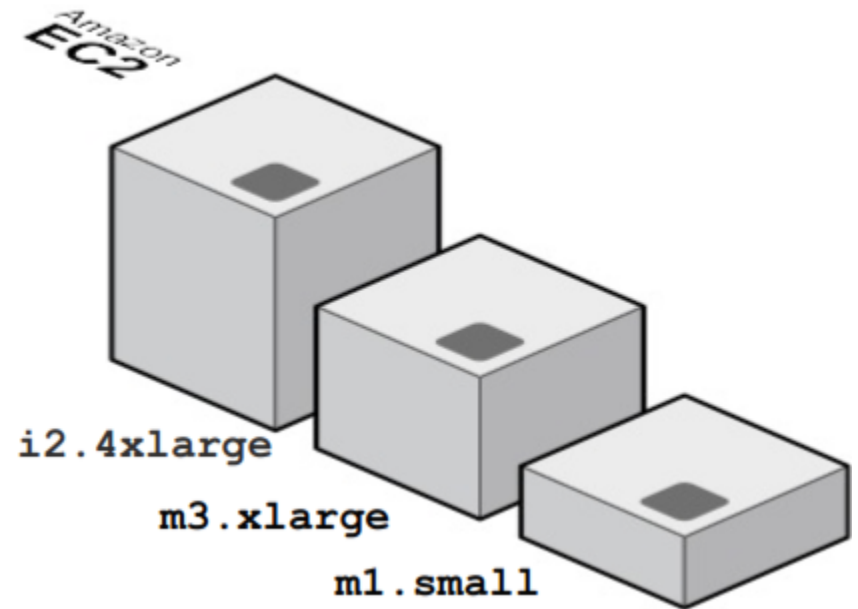  - Also decreasing cost effectiveness

# Scalability

- Ability to grow in size according to need
  - While maintaining performance
- Removal of limits to performance
- Elasticity is one of the main benefits of cloud computing


- One way to remove performance limits is to add resources to each element
  - **Vertical** scaling (scale **up**)
- Another way to remove performance limits is to replicate elements (and balance the load between them)
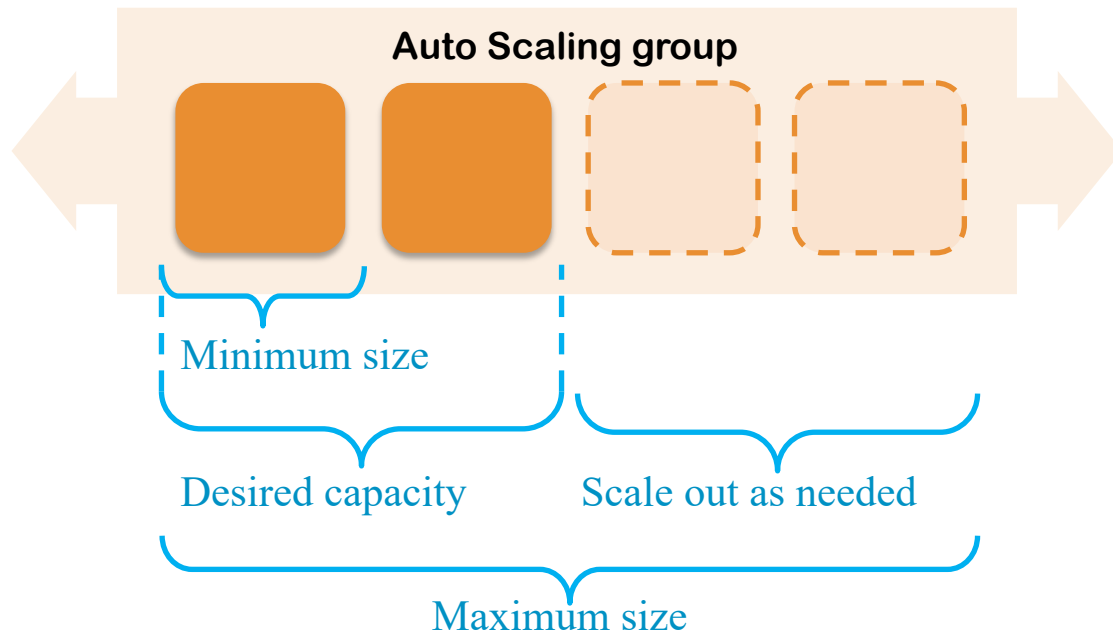  - **Horizontal** scaling (scale **out**)

# Vertical Scaling: Scale up in AWS

- Simple approach.

- High memory/IO/CPU/ Storage.

- Easy to change instance size.

- Will ultimately hit limit.



Amazon EC2

i2.4xlarge

m3.xlarge

m1.small

# Horizontal Scaling: AWS Auto Scaling

**Auto Scaling group**

Minimum size

Desired capacity

Scale out as needed

Maximum size

Auto Scaling groups contain a collection of EC2 instances that share similar characteristics.

Instances in an Auto Scaling group are treated as a logical grouping for the purpose of instance scaling and management.
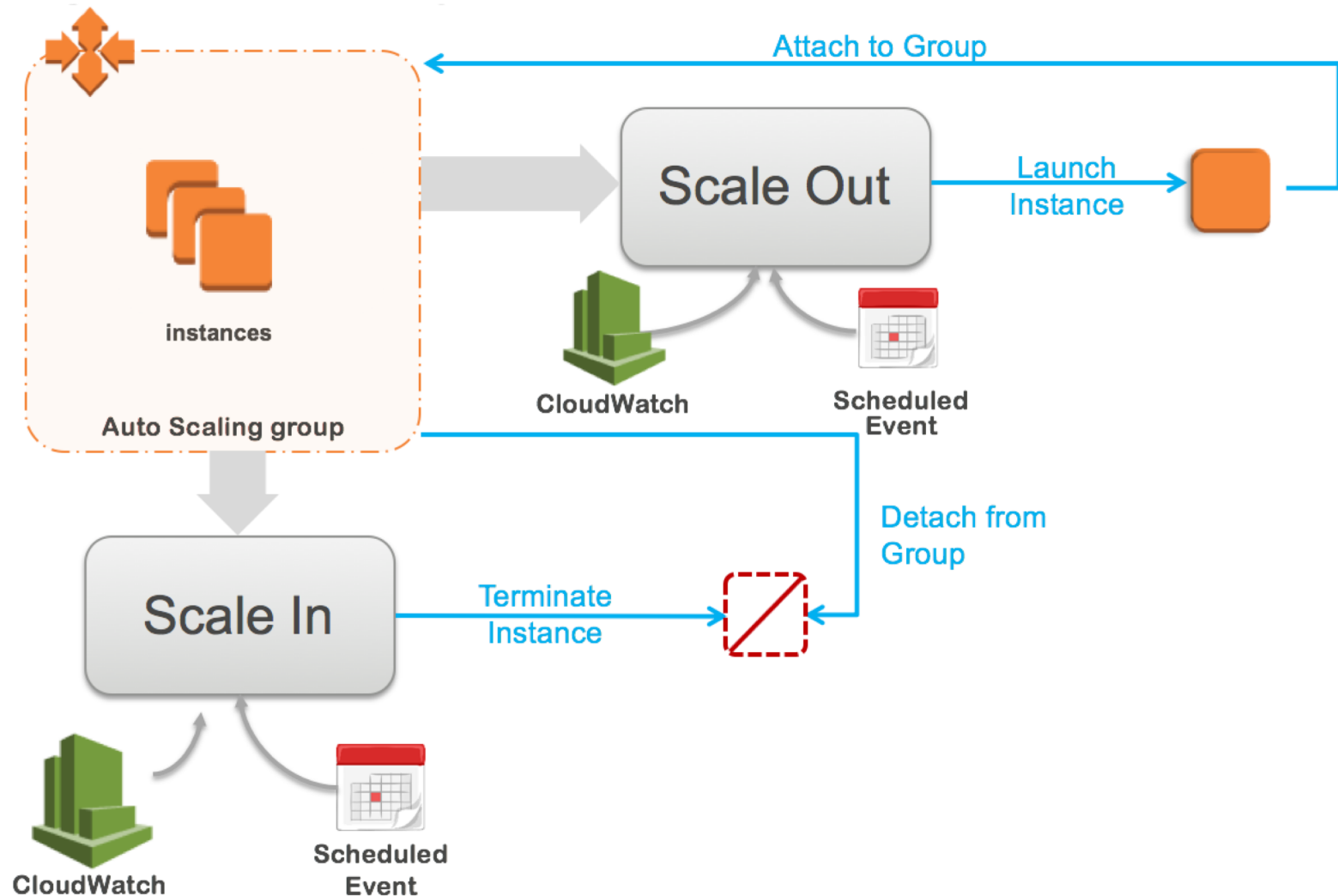
# Dynamic Scaling in AWS

You can create a scaling policy that uses **CloudWatch alarms** to determine when your **Auto Scaling group** should…

scale out                scale in

You can use alarms to monitor:

- Any of the metrics that AWS services send to CloudWatch
- Your own custom metrics

# Auto Scaling Basic Lifecycle

# AWS Launch Configurations

- A **launch configuration** is a template that an Auto Scaling group uses to launch EC2 instances

- When you create a launch configuration, you can specify:
    - AMI
    - Instance type
    - Key pair
    - Security groups
    - Block device mapping
    - User data

- This provides a **blueprint** for each instance that is automatically started by the auto scaler

# Next: Try out AWS auto-scaling

- **Create AMI for launch config**

- **Create Auto Scaling Group**

- **Define Scaling Policies**