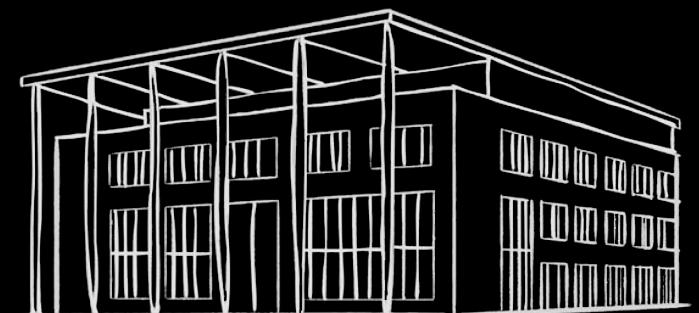


Cloud Computing

Data Science & Analytics Club

10/29/2025

NEAL MAGEE
SCHOOL OF DATA SCIENCE
UNIVERSITY OF VIRGINIA



 UVA Data Science

Introduction & Concepts

Who am I?

Neal Magee

Associate Professor, UVA School of Data Science
Solution Architect, UVA Research Computing
Lead, Internet Operations 2007-2017
.NET/C# Developer 2004-2007
HTML/CGI Developer 1990s

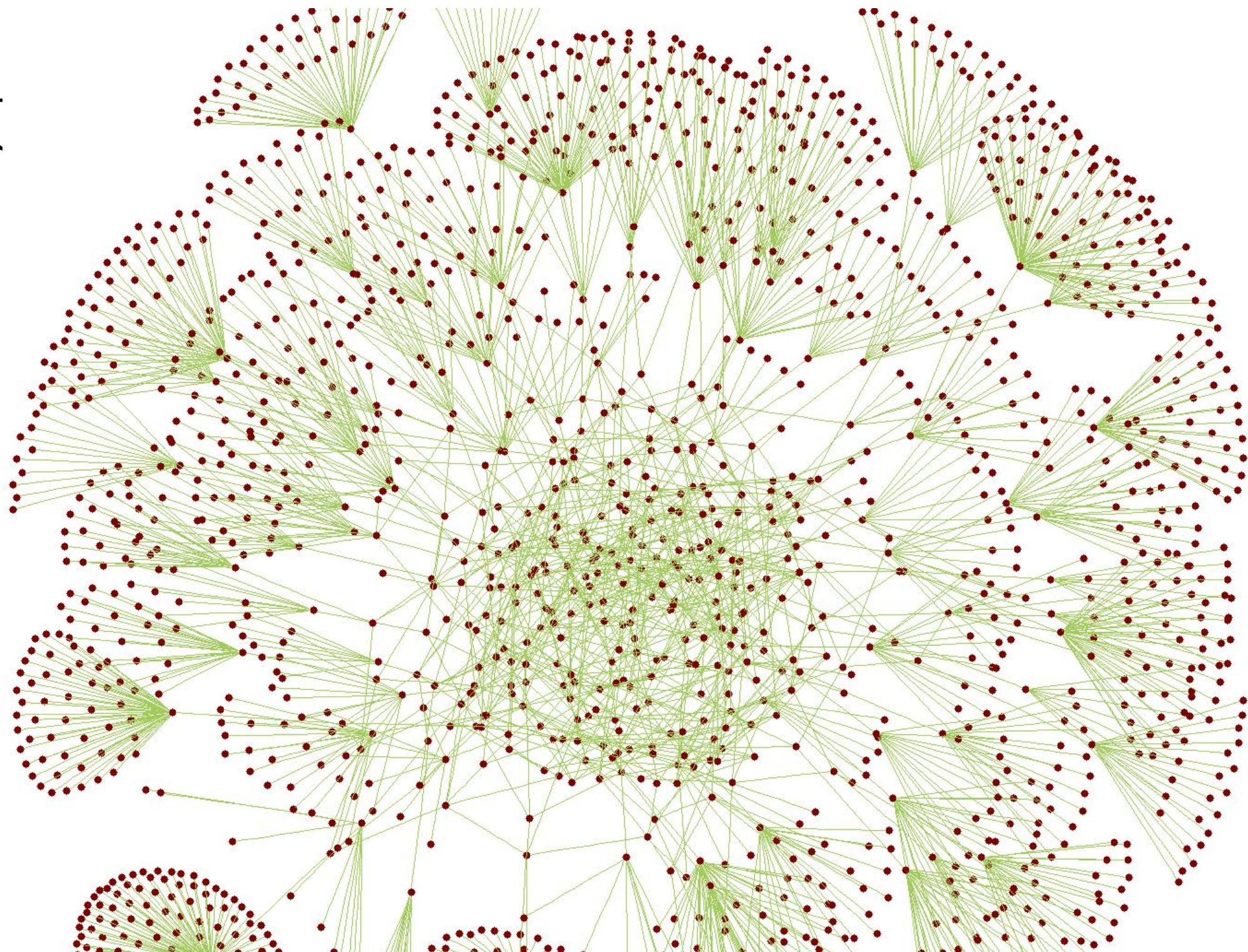


Certifications

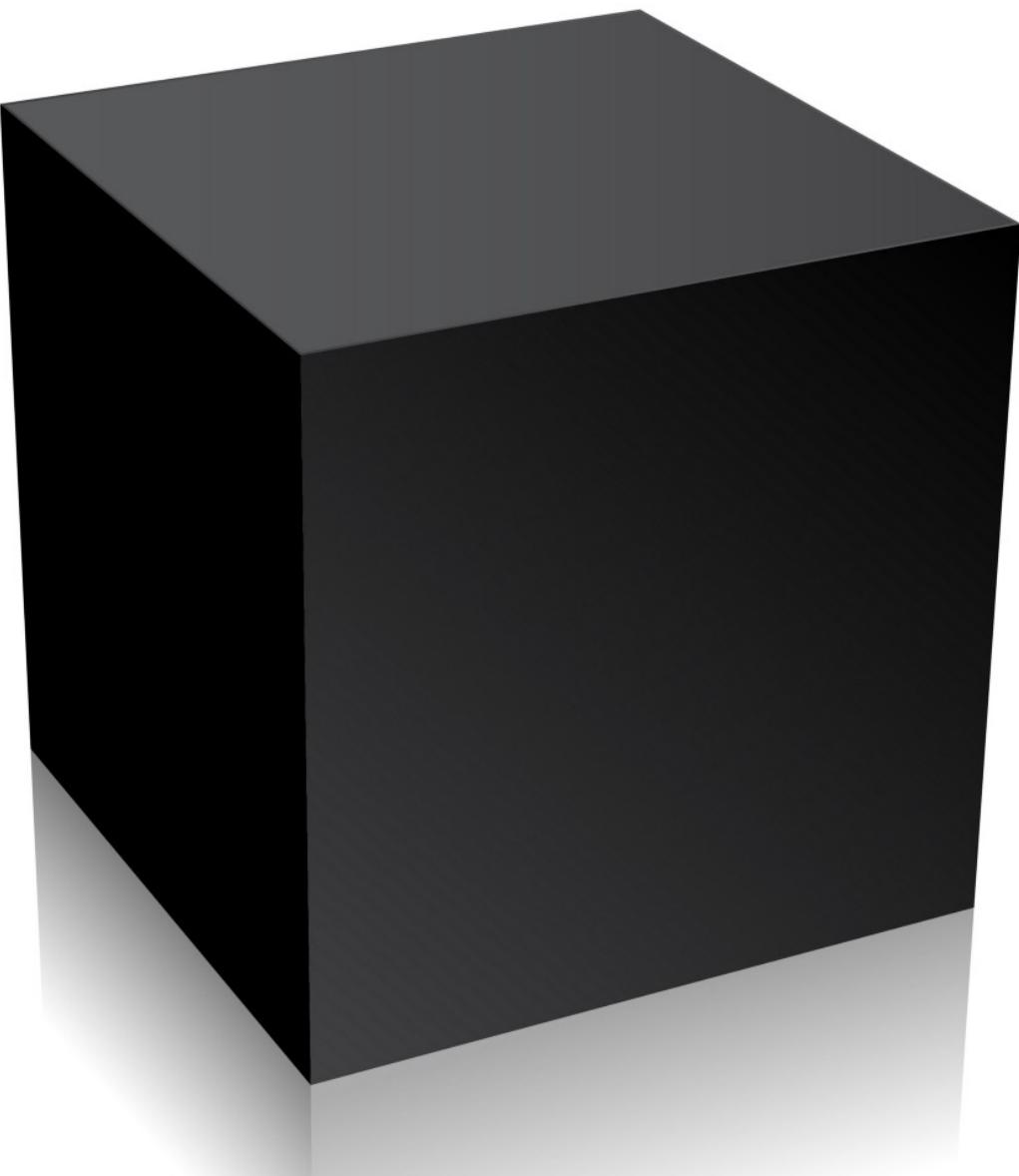
AWS Solution Architect Professional
AWS Solution Architect Associate
AWS Developer Associate

What is the cloud?

Internet



Cloud Computing



INFRASTRUCTURE PLATFORM (IaaS)

OpenStack
vSphere
Azure Stack VMs

AWS EC2
GCE
Azure VMs

CONTAINER PLATFORM (CaaS)

Kubernetes
DC/OS
Docker Datacenter

GKE
ECS
ACS

APPLICATION PLATFORM (PaaS / aPaaS)

CloudFoundry
OpenShift
WaveMaker RAD

Heroku
PCF
Jelastic

FUNCTION PLATFORM (FaaS)

OpenWhisk
Fission
Iron.io

Lambda
GCF
Azure Functions

Public Clouds



Google Cloud Platform



CLOUD INFRASTRUCTURE

Figure 1: Magic Quadrant for Cloud Infrastructure and Platform Services



Source: Gartner (July 2021)

Why the cloud?

1. Unused capacity
2. Capital expenses
3. Scalable
4. Programmable



PARKED 95% OF THE TIME

IDLE 95% OF THE TIME

```

1 [          0.6%] 9 [ ||          2.6%] 17 [          0.0%] 25 [          0.0%
2 [          0.7%] 10 [          0.0%] 18 [          0.0%] 26 [          0.0%
3 [          0.0%] 11 [          0.0%] 19 [          0.0%] 27 [          0.0%
4 [          0.7%] 12 [          1.3%] 20 [          0.6%] 28 [          0.0%
5 [          0.7%] 13 [          0.0%] 21 [          0.0%] 29 [          0.0%
6 [          0.0%] 14 [          0.0%] 22 [          0.0%] 30 [          0.0%
7 [          0.0%] 15 [          0.0%] 23 [          0.0%] 31 [          0.0%
8 [          0.6%] 16 [          0.0%] 24 [          0.0%] 32 [          0.0%
Mem[|||||14.5G/62.8G] Tasks: 116, 606 thr; 1 running
Swp[          431M/63.9G] Load average: 1.20 0.64 0.45
                           Uptime: 300 days(!), 12:21:01

```

PID	USER	PRI	NI	VIRT	RES	SHR	S	CPU%	MEM%	TIME+	Command
1	root	20	0	38204	5988	3908	S	0.0	0.0	41:50.67	/lib/systemd/systemd --system --deserialize 38
186098	root	20	0	270M	4288	3712	S	0.0	0.0	0:00.39	└─ /usr/lib/policykit-1/polkitd --no-debug
186105	root	20	0	270M	4288	3712	S	0.0	0.0	0:00.20	└─ /usr/lib/policykit-1/polkitd --no-debug
186102	root	20	0	270M	4288	3712	S	0.0	0.0	0:00.00	└─ /usr/lib/policykit-1/polkitd --no-debug
162305	clamav	20	0	905M	668M	15892	S	0.0	1.0	9:04.55	└─ /usr/sbin/clamd --foreground=true
169345	clamav	20	0	905M	668M	15892	S	0.0	1.0	0:00.00	└─ /usr/sbin/clamd --foreground=true
161057	clamav	20	0	129M	11820	9320	S	0.0	0.0	1:12.69	└─ /usr/bin/freshclam -d --foreground=true
160514	root	20	0	44148	3268	2740	S	0.0	0.0	0:01.23	└─ /lib/systemd/systemd-udevd
156666	neal	20	0	26052	2912	2304	S	0.0	0.0	0:00.01	└─ SCREEN -R
156667	neal	20	0	21840	5648	3232	S	0.0	0.0	0:00.04	└─ /bin/bash
156694	root	20	0	69884	4344	3708	S	0.0	0.0	0:00.00	└─ sudo su
156695	root	20	0	69332	4120	3524	S	0.0	0.0	0:00.00	└─ su
156696	root	20	0	21308	5124	3244	S	0.0	0.0	0:00.02	└─ bash
156860	root	20	0	11244	3024	2816	S	0.0	0.0	0:00.00	└─ /bin/bash ./update-rivanna.sh
156861	root	20	0	46988	5376	4704	S	0.0	0.0	0:00.79	└─ ssh -tt -i /root/.ssh/nem2p_rivanna nem2p@rivanna.hpc.virginia.edu
144361	root	20	0	21824	2912	2728	S	0.0	0.0	0:00.11	└─ /usr/sbin/incron
140705	syslog	20	0	250M	4128	2792	S	0.0	0.0	2:41.15	└─ /usr/sbin/rsyslogd -n
140711	syslog	20	0	250M	4128	2792	S	0.0	0.0	1:25.87	└─ /usr/sbin/rsyslogd -n
140710	syslog	20	0	250M	4128	2792	S	0.0	0.0	0:00.00	└─ /usr/sbin/rsyslogd -n
140709	syslog	20	0	250M	4128	2792	S	0.0	0.0	1:15.18	└─ /usr/sbin/rsyslogd -n
136291	root	10	-10	5720	3660	2456	S	0.0	0.0	1:24.13	└─ /sbin/iscsid
136290	root	20	0	5220	1724	1596	S	0.0	0.0	0:17.12	└─ /sbin/iscsid
127469	root	20	0	65512	5512	4972	S	0.0	0.0	0:00.06	└─ /usr/sbin/sshd -D
71057	root	20	0	106M	7504	6400	S	0.0	0.0	0:00.00	└─ sshd: neal [priv]
71158	neal	20	0	106M	3316	2216	S	0.0	0.0	0:00.00	└─ sshd: neal@pts/0
71159	neal	20	0	21832	5648	3272	S	0.0	0.0	0:00.08	└─ -bash
71191	neal	20	0	25600	4668	3100	R	2.0	0.0	0:00.88	└─ htop
123780	ntp	20	0	108M	3360	2772	S	0.0	0.0	33:36.02	└─ /usr/sbin/ntpd -p /var/run/ntpd.pid -g -u 115:123
11605	ntp	20	0	108M	3360	2772	S	0.0	0.0	0:00.12	└─ /usr/sbin/ntpd -p /var/run/ntpd.pid -g -u 115:123
87930	root	20	0	55840	11316	3780	S	0.0	0.0	0:30.26	└─ python /usr/sbin/denyhosts --daemon --purge --config=/etc/denyhosts.conf

F1Help F2Setup F3Search F4Filter F5Sorted F6Collapse F7Nice - F8Nice + F9Kill F10Quit



Why the cloud?

1. Unused capacity
2. Capital vs. Operational expenses
3. Scalable
4. Programmable



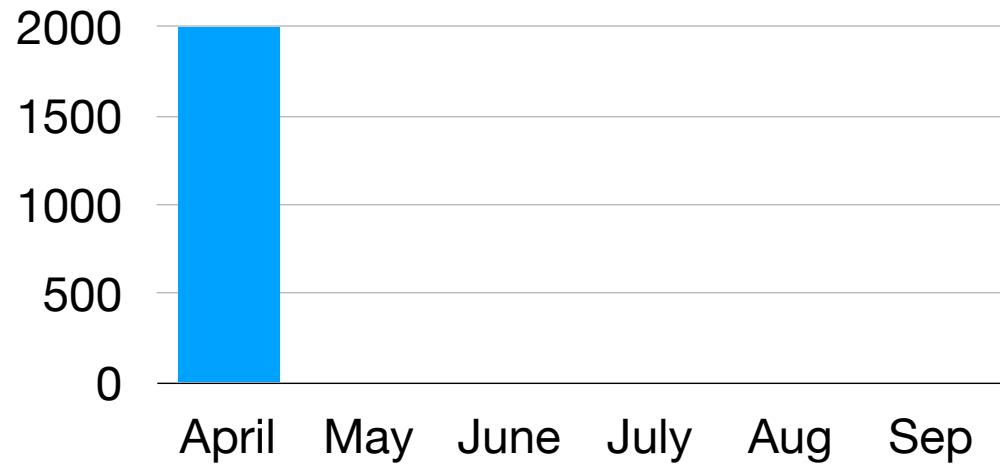


Capital Expenses

Order
Ship
Fixed CPU
Fixed RAM
Fixed Storage

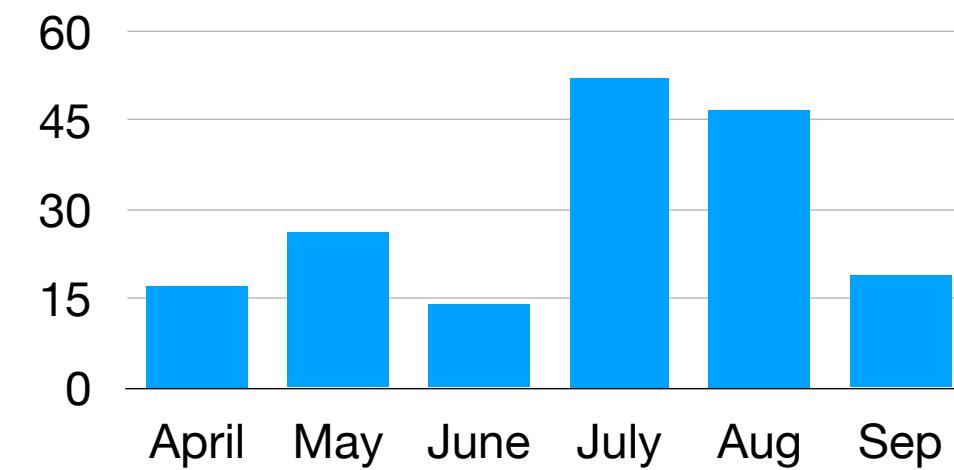
Operational Expenses

Power
Connectivity
Cooling
Engineering



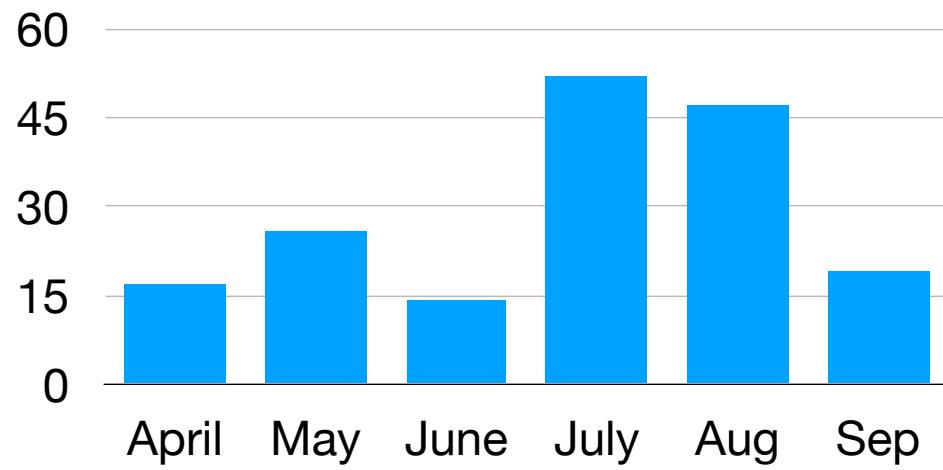
Capital Expenses

Order
Fixed CPU
Fixed RAM
Fixed Storage
Ship



Operational Expenses

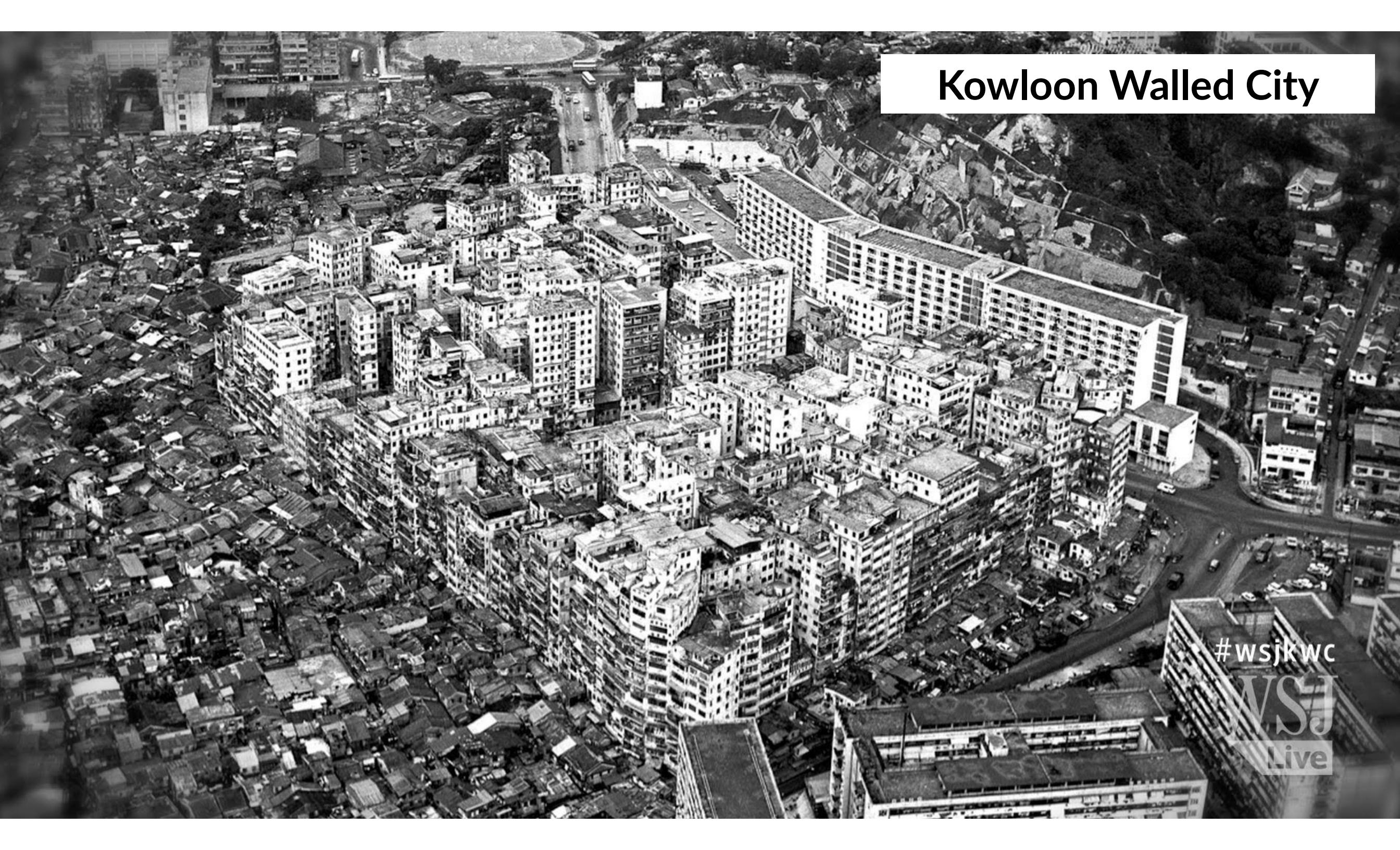
Power
Connectivity
Cooling
Engineering



Cloud = Pay As You Go

Why the cloud?

1. Unused capacity
2. Capital vs. Operational expenses
3. Scalable & Reconfigurable
4. Programmable



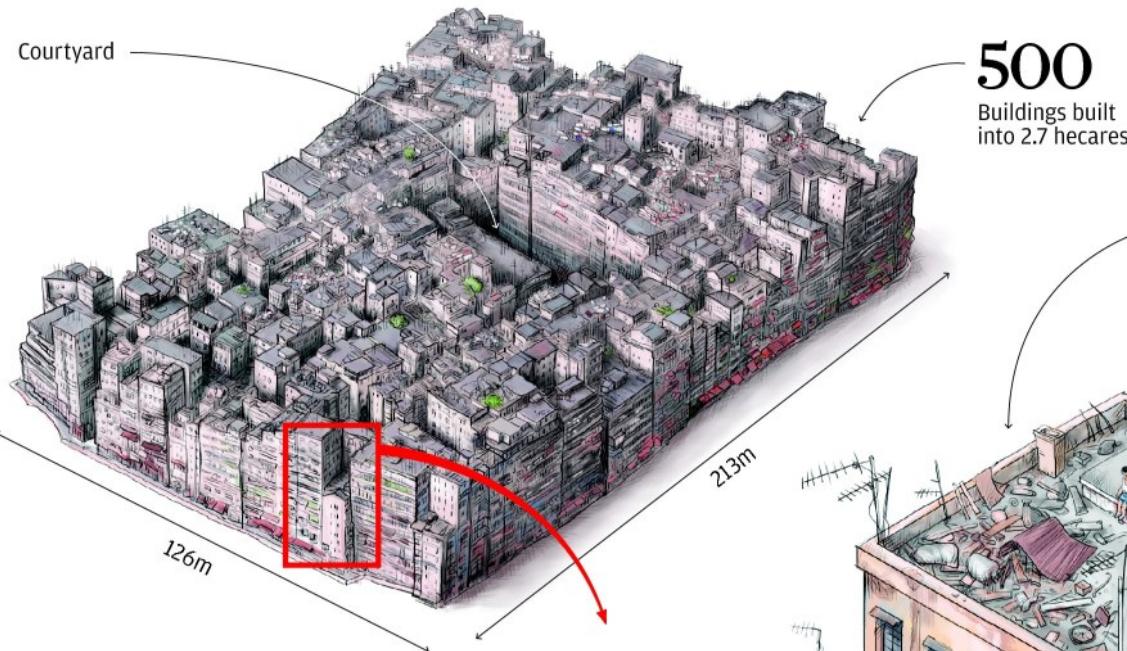
Kowloon Walled City

#wsjkwc
WSJ
Live

City of anarchy

Kowloon Walled City, located not far from the former Kai Tak Airport, was a remarkable high-rise squatter camp that by the 1980s had 50,000 residents. A historical accident of colonial Hong Kong, it existed in a lawless vacuum until it became an embarrassment for Britain.

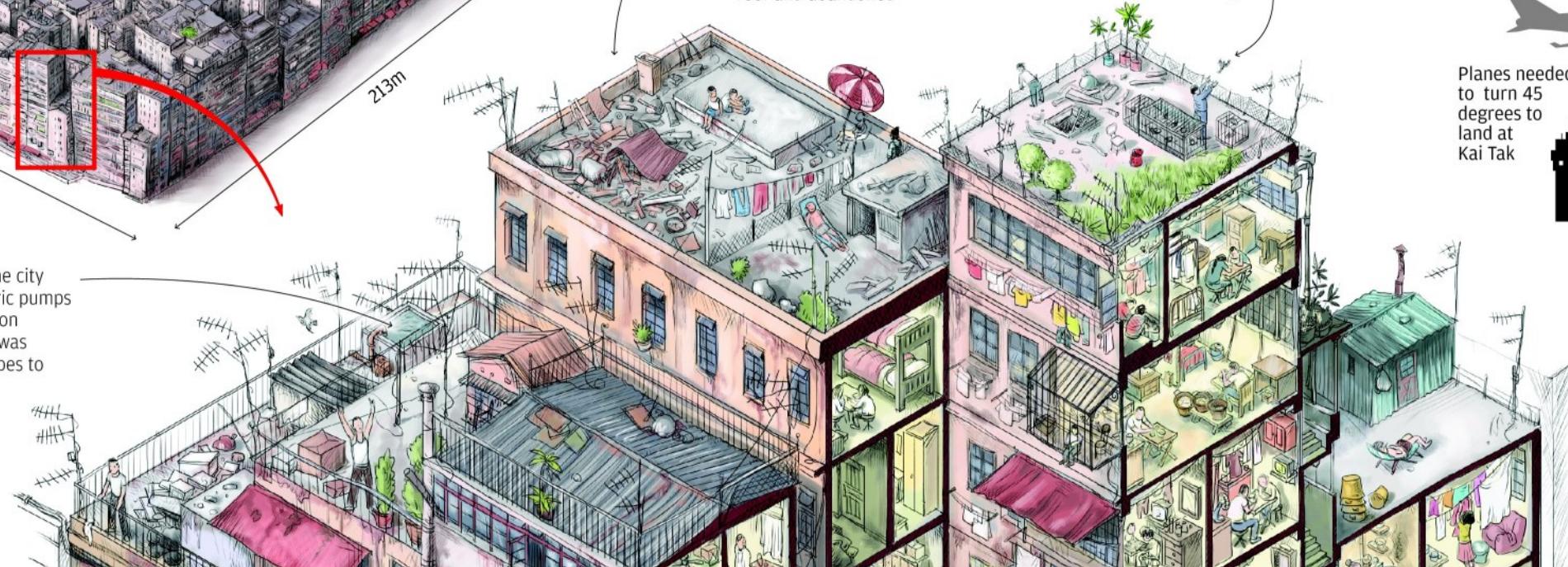
This month marks the 20th anniversary of its demolition.



500

Buildings built
into 2.7 hectares

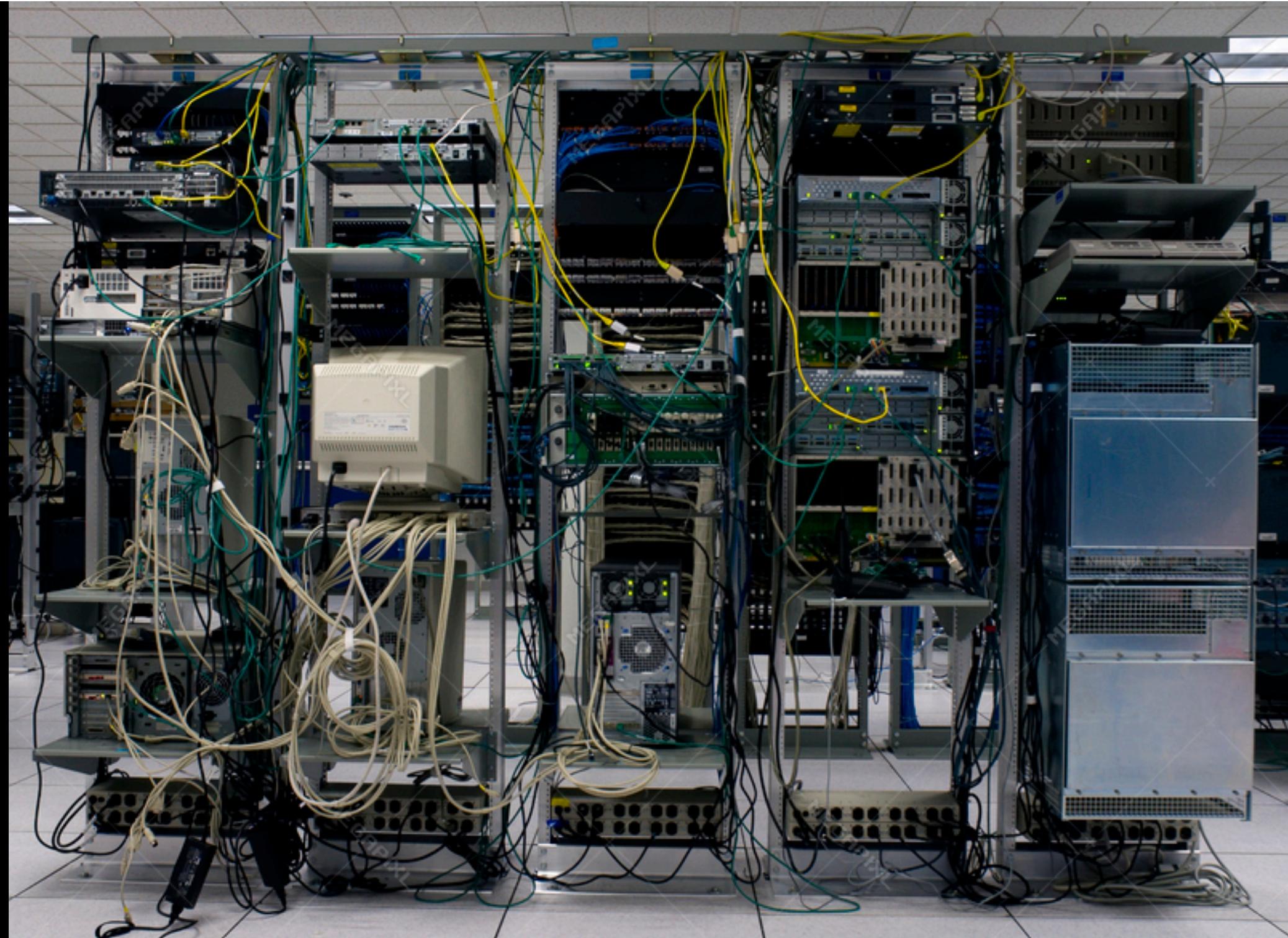
Without municipal services, there was no rubbish collection. Old television sets, broken furniture, discarded mattresses and other bulky items were hauled to the roof and abandoned

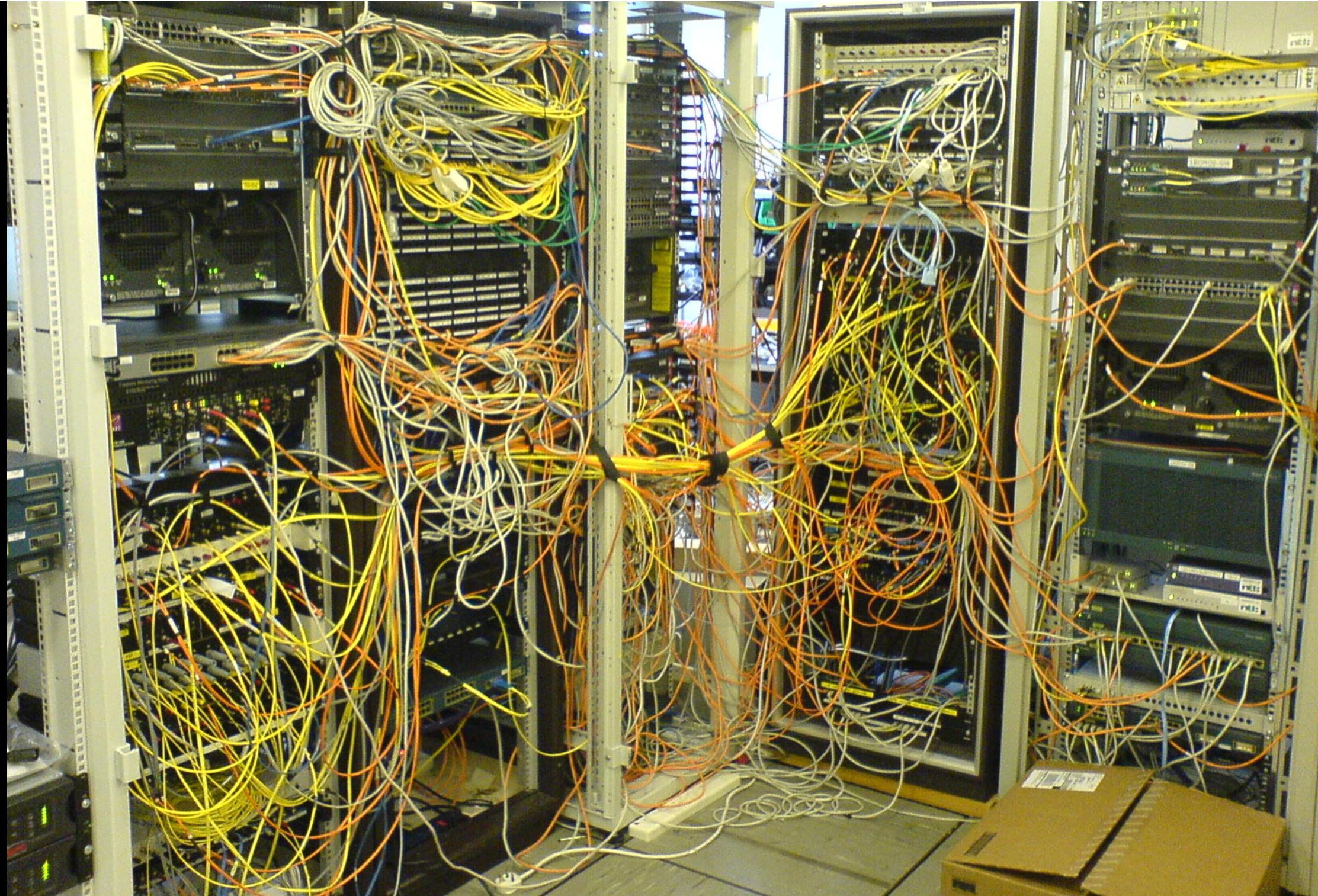




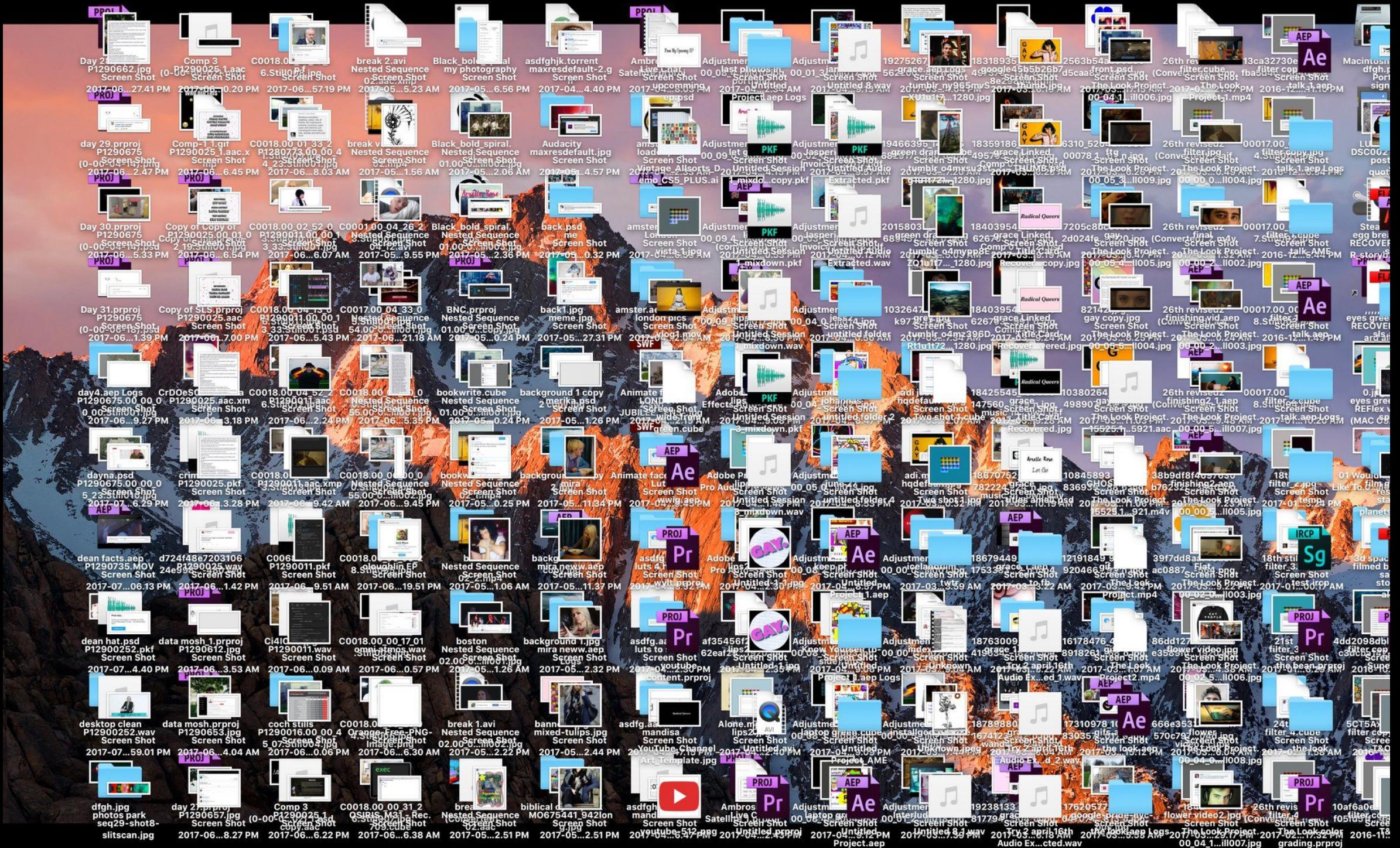
KOWLOON EFFECT

- ACCUMULATE OVER TIME
- TIGHTLY COUPLED
- CANNOT SCALE
- REDUNDANCIES
- HARD TO TROUBLESHOOT









```
1  # Prevent database truncation if the test fails
2
3  abort("The Rails environment is running in production mode!
4      Please use 'rails test' to run tests in the
5      test environment or 'rake db:reset' to
6      destroy the database and start a new one")
7
8
9  require 'spec_helper'
10 require 'rspec/rails'
11
12 require 'capybara/rspec'
13 require 'capybara/rails'
14
15 Capybara.javascript_driver = :webkit
16 Category.delete_all; Category.create
17 Shoulda::Matchers.configure do |config|
18   config.integrate do |with|
19     with.test_framework :rspec
20     with.library :rails
21   end
22 end
23
24 # Add additional requires below this line if you need them
25
26 # Requires supporting ruby files with custom matchers and
27 # helper methods under 'spec/support/' and its subdirectories. By default,
28 # in _spec.rb will both be required on the test page, so you can just
29 # run twice. It is recommended that you do not name files
30 # end with _spec.rb. You can configure this by setting
31 # the matchers_path and helpers_path options in spec_helper.rb.
```

Why the cloud?

1. Unused capacity
2. Capital vs. Operational expenses
3. Scalable & Reconfigurable
4. Programmable



CASE STUDY: NETFLIX

GLOBAL PRESENCE IN 190 COUNTRIES

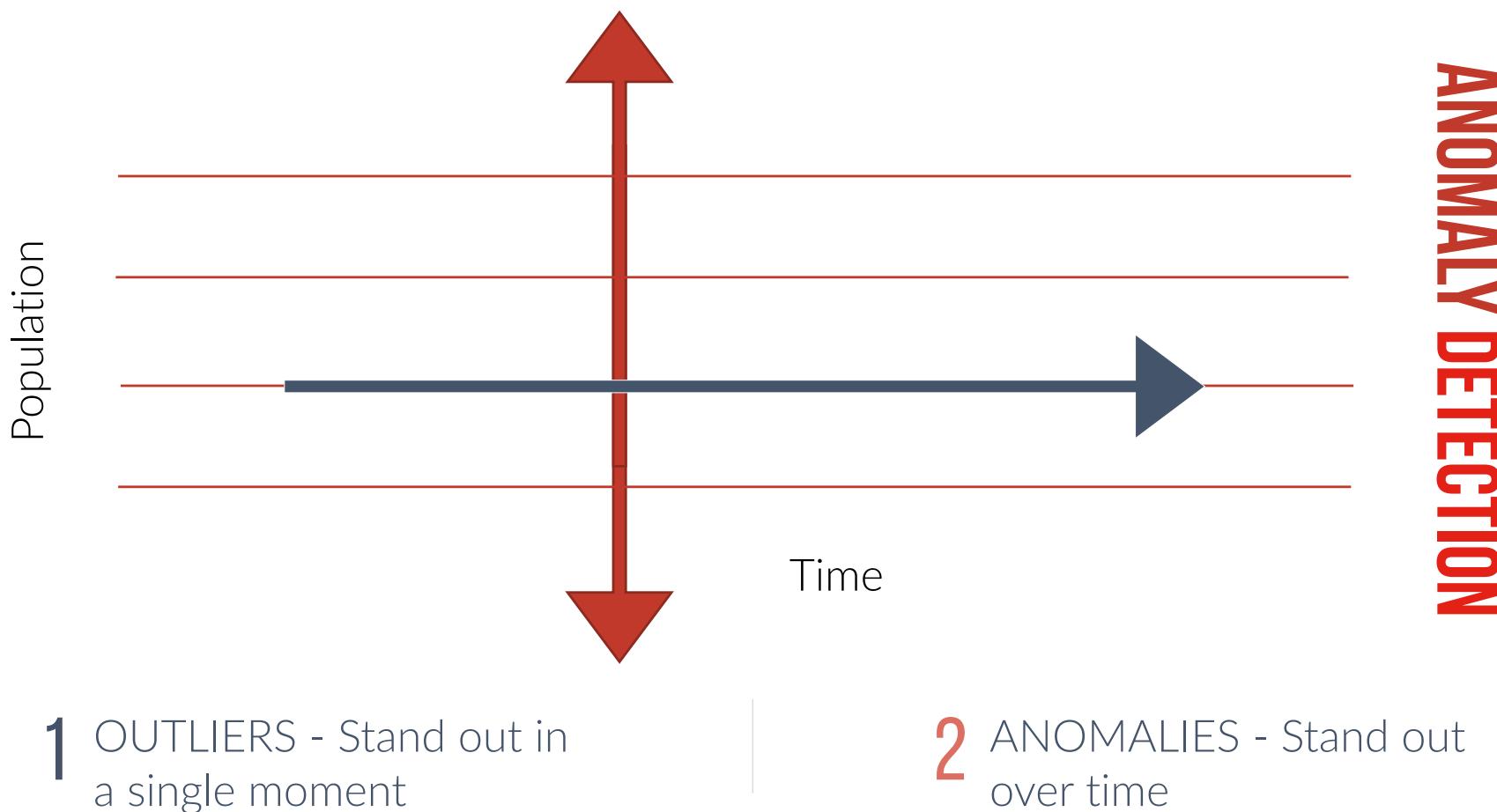


2000+ A/B TESTS ANNUALLY

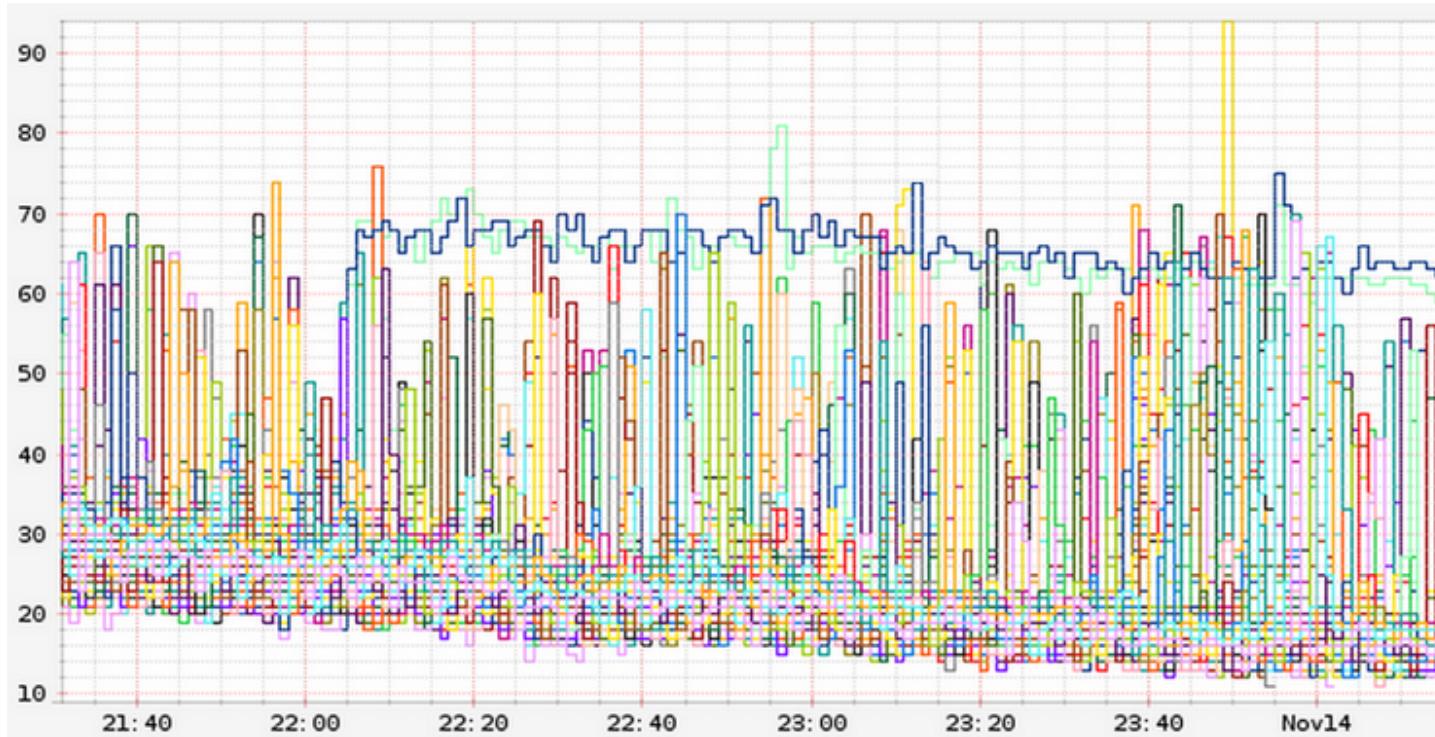
100,000+ INSTANCES PEAK

25-30% OF ALL INTERNET BANDWIDTH

OUTLIER DETECTION

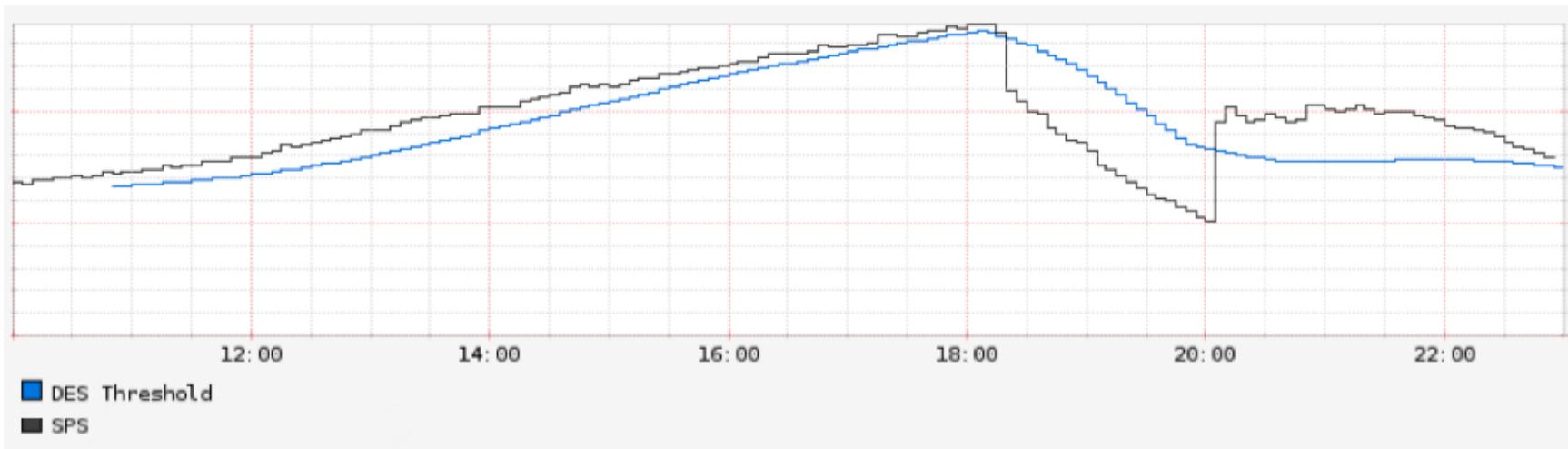


OUTLIER DETECTION



1 OUTLIERS - Stand out in
a single moment

ANOMALY DETECTION



Kalman filter smoothes time series data

Predict the future based on past history

Favor recent history

Threshold-based alerts

2 ANOMALIES - Stand out over time



CHAOS ENGINEERING TOOLS

Provoking distributed systems to understand their behavior under disruption.

NETFLIX SIMIAN ARMY: Fault Injection

- 1 Chaos Monkey
- 2 Janitor Monkey
- 3 Conformity Monkey



AWS FAULT INJECTION SIMULATOR



A managed service for running fault injection experiments to improve an application's performance, observability, and resiliency.

CDN HANDOFF



TIME →



us-east-ATL

us-east-MIA

us-west-POR

us-west-LA

eu-west-DUB

LOAD SHEDDING

NETFLIX

TIME

ACTIVE USER PLAYBACK

High SLA

ACCOUNT MANAGEMENT TASKS

Flexible SLA

TRAILER PREVIEWS

Flexible SLA

Service Disruption

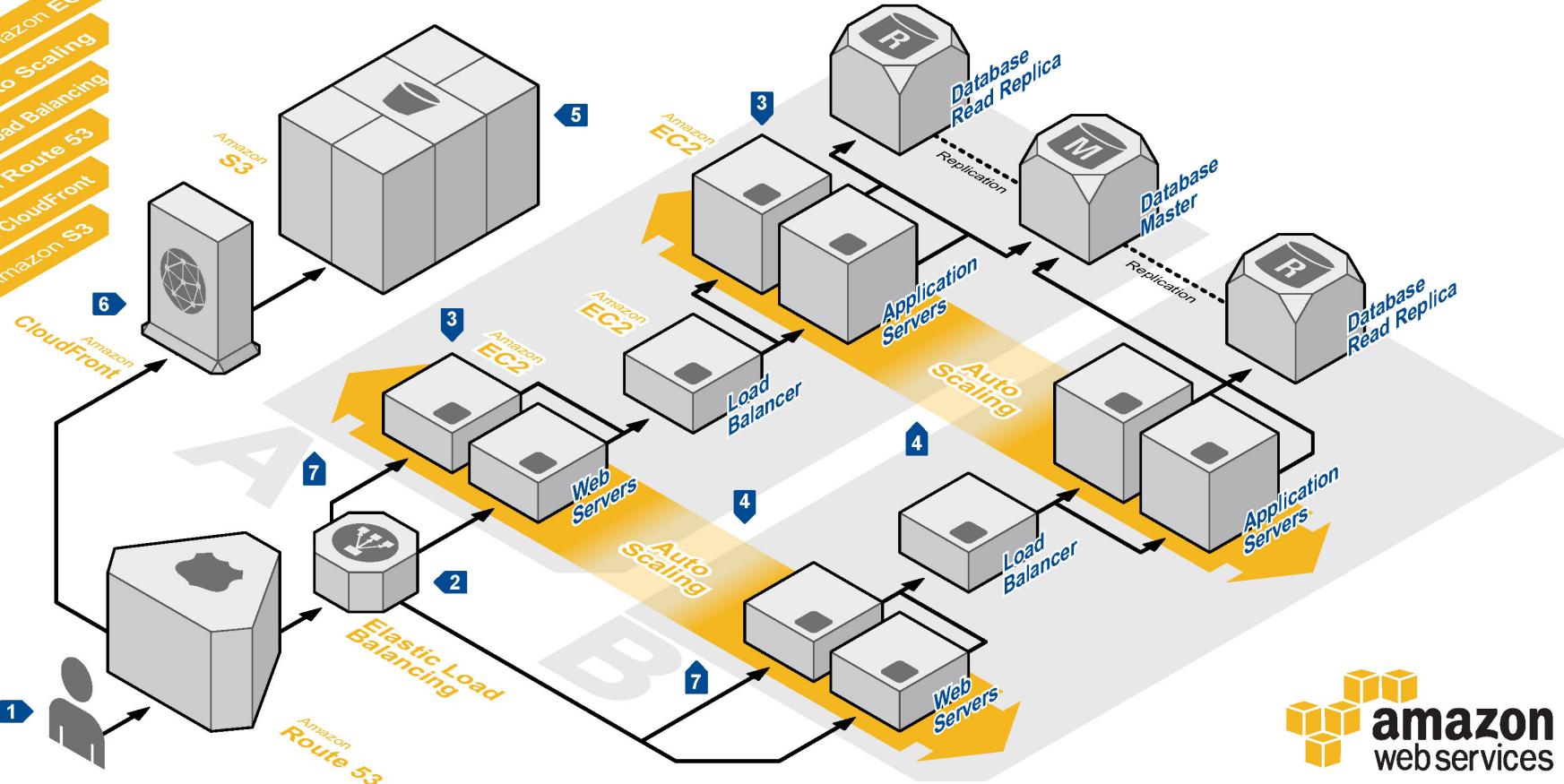


Restoration

Queue/drop lower-priority tasks to sustain priority service layers

WEB APPLICATION HOSTING

Highly available and scalable web hosting can be complex and expensive. Dense peak periods and wild swings in traffic patterns result in low utilization rates of expensive hardware. Amazon Web Services provides the reliable, scalable, secure, and high-performance infrastructure required for web applications while enabling an elastic, scale out and scale down infrastructure to match IT costs in real time as customer traffic fluctuates.



System Overview

1 The user's DNS requests are served by **Amazon Route 53**, a highly available Domain Name System (DNS) service. Network traffic is routed to infrastructure running in Amazon Web Services.

2 HTTP requests are first handled by Elastic Load Balancing, which automatically distributes incoming application traffic across multiple **Amazon Elastic Compute Cloud (EC2)** instances across Availability Zones (AZs). It enables even greater fault tolerance in your applications, seamlessly providing the amount of load balancing capacity needed in response to incoming application traffic.

3 Web servers and application servers are deployed on **Amazon EC2** instances. Most organizations will select an **Amazon Machine Image (AMI)** and then customize it to their needs. This custom AMI will then be used as the starting point for future web development.

4 Web servers and application servers are deployed in an **Auto Scaling** group. Auto Scaling automatically adjusts your capacity up or down according to conditions you define. With Auto Scaling, you can ensure that the number of **Amazon EC2** instances you're using increases seamlessly during demand spikes to maintain performance and decreases automatically during demand lulls to minimize costs.

5 Resources and static content used by the web application are stored on **Amazon Simple Storage Service (S3)**, a highly durable storage infrastructure designed for mission-critical and primary data storage.

6 Static and streaming content is delivered by **Amazon CloudFront**, a global network of edge locations. Requests are automatically routed to the nearest edge location, so content is delivered with the best possible performance.

7 **Availability zones (AZs)** are distinct geographic locations that are engineered to insulate against failures in other AZs. Multiple AZs are combined into a region. Here, the entire web application is deployed in two different AZs for high availability.

- Well-defined
- Loosely coupled
- Scalable
- Reproducible
- Can troubleshoot easily
- Less redundancy

Why the cloud?

1. Unused capacity
2. Capital vs. Operational expenses
3. Scalable
4. Programmable

Infrastructure as Code

Infrastructure as Code

Automatable

Versionable

Repeatable

Distributable

CLOUDFORMATION

```
▼ {  
    "AWSTemplateFormatVersion": "2010-09-09",  
    "Description": "AWS CloudFormation Sample Template LAMP_Multi_AZ: Create  
    demonstrates using the AWS CloudFormation bootstrap scripts to install t  
    creates one or more Amazon EC2 instances, an Elastic Load Balancer and a  
    ▶ "Parameters": { ... }, // 10 items  
    ▶ "Mappings": { ... }, // 3 items  
    ▶ "Conditions": { ... }, // 2 items  
    ▶ "Resources": {  
        ▶ "ElasticLoadBalancer": { ... }, // 2 items  
        ▶ "WebServerGroup": { ... }, // 4 items  
        ▶ "LaunchConfig": { ... }, // 3 items  
        ▶ "WebServerSecurityGroup": { ... }, // 2 items  
        ▶ "DBEC2SecurityGroup": { ... }, // 3 items  
        ▶ "DBSecurityGroup": { ... }, // 3 items  
        ▶ "MySQLDatabase": { ... } // 2 items  
    },  
    ▶ "Outputs": { ... } // 1 item  
}
```

TERRAFORM

```
provider "aws" {
    region = "${var.aws_region}"
}

provider "archive" {}

data "archive_file" "zip" {
    type      = "zip"
    source_file = "hello_lambda.py"
    output_path = "hello_lambda.zip"
}

data "aws_iam_policy_document" "policy" {
    statement {
        sid =
        effect = "Allow"
        principals {
            identifiers = ["lambda.amazonaws.com"]
            type = "Service"
        }
        actions = ["sts:AssumeRole"]
    }
}
```

RESIZE AN INSTANCE

```
aws ec2 stop-instances \  
--instance-ids i-123456789012
```

```
aws ec2 modify-instance-attribute \  
--instance-id i-123456789012 \  
--instance-type p2.16xlarge
```

```
aws ec2 start-instances \  
--instance-ids i-012987654321
```

RESIZING PROGRAMMATICALLY

```
import boto3

instance = "i-123456789012"
ec2 = boto3.client('ec2')

response =
client.modify_instance_attribute(
    InstanceId=instance,
    InstanceType={
        'Value': 'c4.8xlarge',
    }
)
```

Cattle Not Pets



Pets

- * Are given names
- * Are unique & hand-raised
- * When sick you nurse back to health

Cattle

- * Are assigned IDs
- * Are nearly identical
- * When sick you get another

Questions?

Sign Up for AWS Academy

bit.ly/dsac-cloud

GitHub Repo

github.com/nmagee/dsac-cloud