

# Cloud Computing

## Computing as a utility

H2O gas electricity

\$ \$ \$\$

pay based on usage  
meter

def: IT provided as service over internet  
style of comp where scalable

- \* illusion of infinite resources
- \* no up front & fine grained billing ex: hourly
- \* calc./operating system on demand

1969 Leonard Kleinrock ARPANET Project

predicted: comp network → utility

1989 Grid

- solve BIG prob w/ Parallel computing
- by Global Alliance
- Comp resources as a matured service
- failed bc weak internet

1990's

Utility

SaaS

network-based subscription 2 app

2001

NEXT GEN Internet Computing

NEXT GEN Data Centers

why now?

- ① lack of experience w/ large data centers
- ② unknown pricing model
- ③ risk of transferring

④ Tech Factors

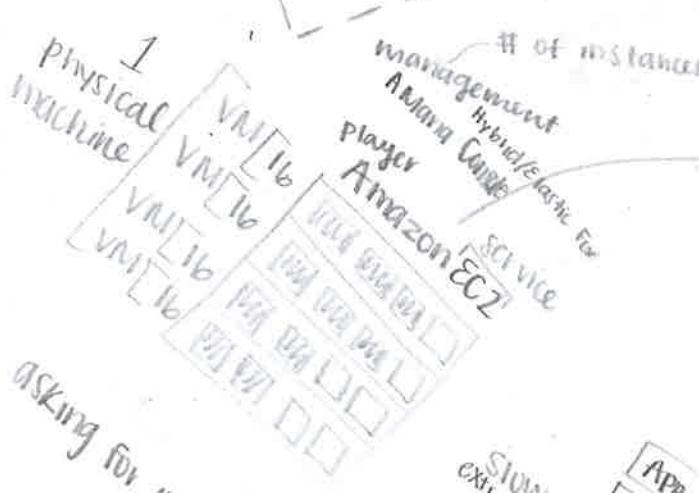
- \* growth of internet
- \* growth of VM
- \* business factors
- \* min capital

⑤ Business Factors

- \* pay-as-you-go
- \* data store on bandwidth

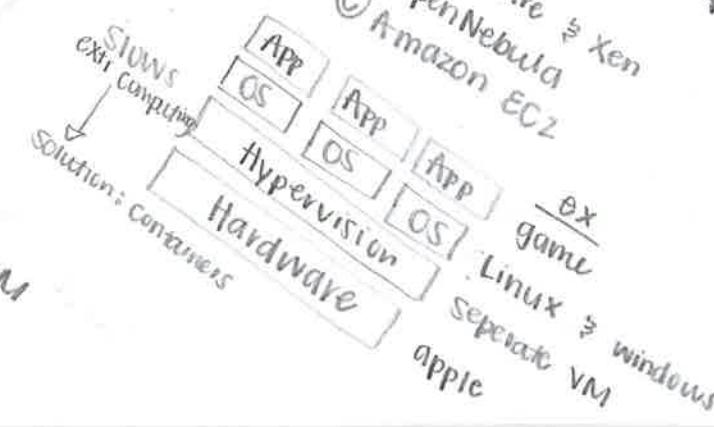
Storage

ex: Amazon S3



**Virtualization**

- basis of cloud comp
- partition hardware → flexible & scalable
- techniques:
  - ① VM ware → Xen
  - ② OpenNebula
  - ③ Amazon EC2

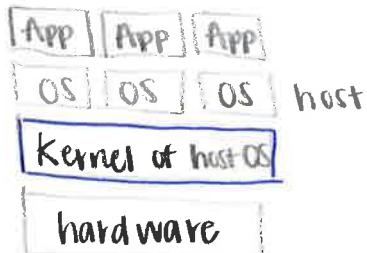


asking for 10 M → get 10 VMs  
not cost effective

# Containers

by Google

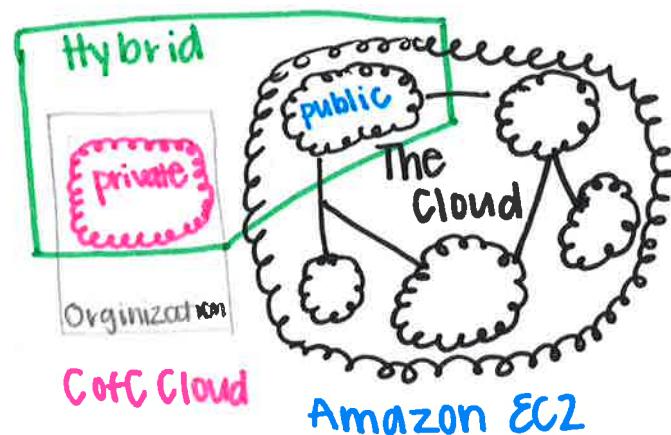
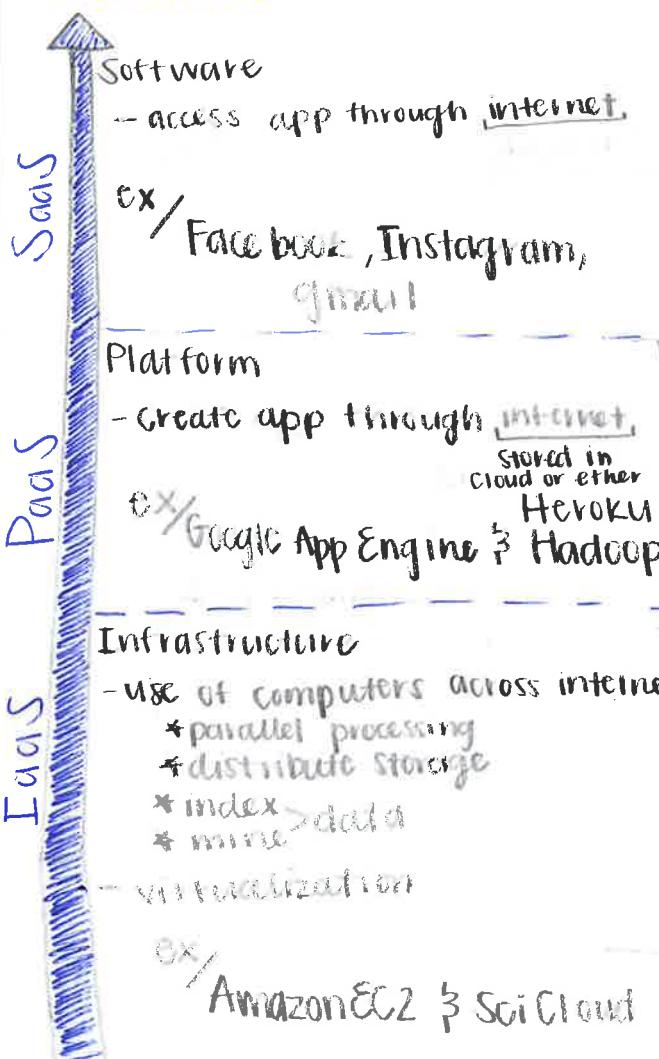
lec 1 | 2



- + light weight (uses only required resources)
- + faster start up than new VM
- + near native performance

ex: Docker, LXC, Linux VServer

# Services



Community  
(multiple universities)

Linux Application MySQL  
Billed on AD Revenue Premium Service  
Software stack

OS on demand

data

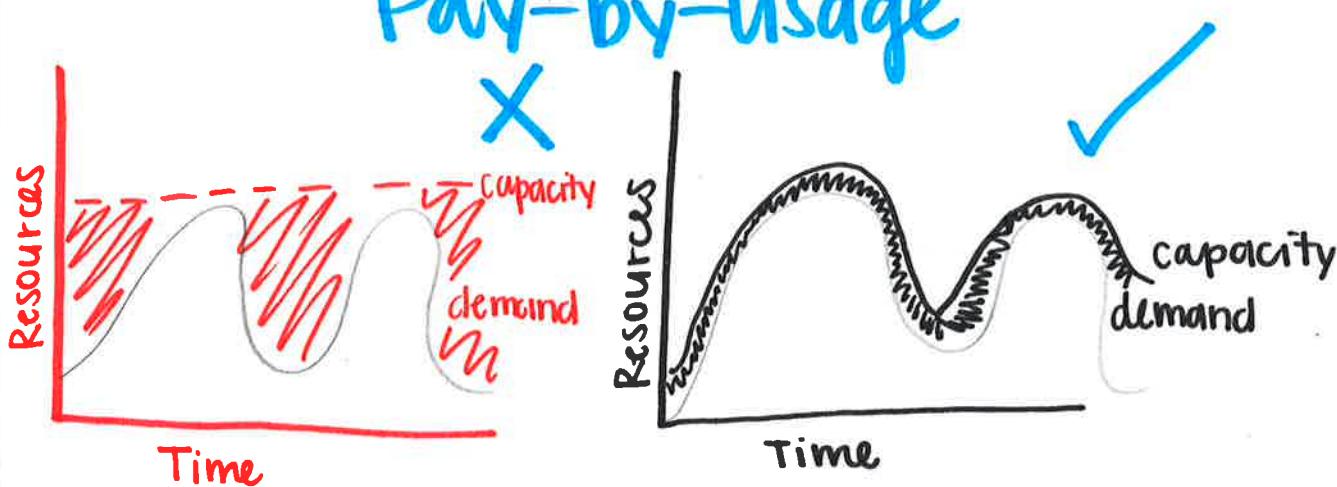
- Store
- runtime
- LB
- Resource Provision

## Short Term Implications

- ① Startup & prototyping
  - ↓ cost of entry
- ② Batch jobs
  - ex banks
    - ↑ push all transactions
    - ↑ recommit @ night
- ③ One-off Task
  - NY Times
    - ↓ digitize all earlier copies
- ④ Cost efficient for scientific comp
- ⑤ Research @ scale
  - \* test which model is best

# Pay-by-Usage

Lec 13



risk of overprocessing

under utilization

## Sunk Cost

### Econ of Cloud Providers

- \* build large datacenter

    \$\$\$ \$100 mil

- \* large companies already own

    ex/google, Amazon

COST of Medium DC > Large DC

    Network  
    Storage  
    Admin

- \* Power

✓ 3.6¢ | Idaho | hydroelectric

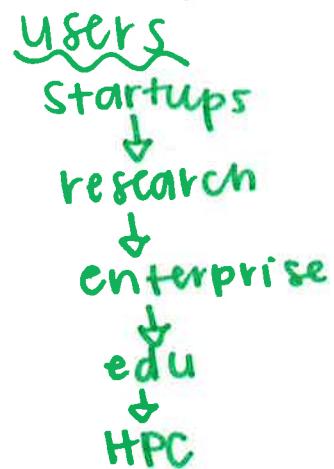
✗ 10¢ | Cali | transmit ele over long distance

✗ 18¢ | Hawaii | ship fuel \$\$\$

- \* Cooling

    dc near rivers

\* world's 5<sup>th</sup> 6<sup>th</sup> most energy consumption



generating revenue

Microsoft : windows is already downloaded to you  
    ↳ they send you access  
    → pay for license

Google : reuse infrastructure

### Policy & Business

Reputation  
Fate Sharing

customers  
SPAM creating content

Cloud providers  
check 4 this

Software  
Licensing

Pay 4 licensing  
ex: windows

high availability servers → commodity servers

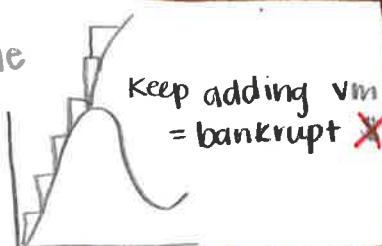
at least 1 failure per day ...

low per server !!

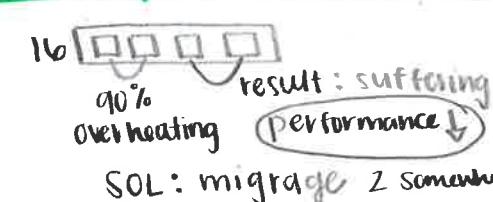
SOLUTION: replicate data & computation

Ex: Mapreduce

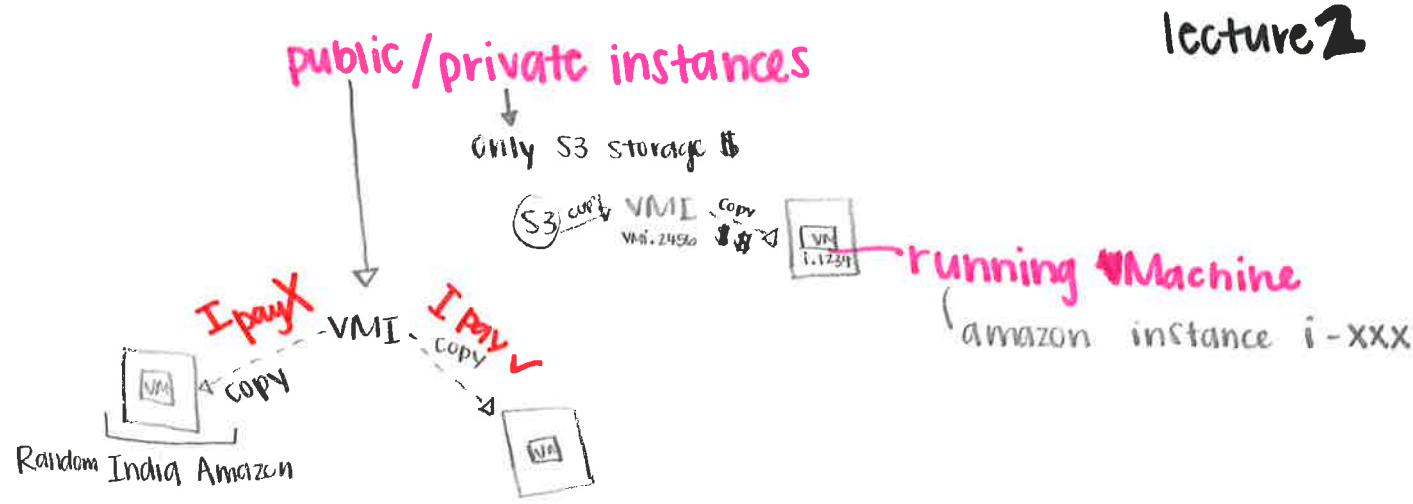
## Adoption

Challenges	Solutions
Availability	<p>what if their servers go down ??</p> <ul style="list-style-type: none"> <li>✓ <b>Multiple Providers</b> instead of 1 google</li> <li>8 amazon 2 google</li> </ul> <p>✓ <b>auto scaling</b> to prevent DDoS attacks</p> 
Data Lock-in	<p>Locked into 1 company ex/ Amazon S3</p> <p>can't migrate from <u>google</u> to <u>amazon</u></p>
Data Confidentiality	<ul style="list-style-type: none"> <li>✓ Encryption</li> <li>✓ Firewall</li> <li>✓ Geographical data store</li> </ul>

## Growth

Data transfer bottleneck	<ul style="list-style-type: none"> <li>✓ Fedex disk</li> <li>✓ Amazon Snowmobile</li> </ul>	
Performance Unpredictability	<ul style="list-style-type: none"> <li>✓ improve VM support</li> <li>✓ Schedule VM</li> <li>✓ Redistribute VM</li> </ul>	
Scalable Solution	Scalable Storage	
in Large Distributed Systems	<ul style="list-style-type: none"> <li>map reduce</li> <li>invent debugger 4 dist VM</li> </ul>	
Scaling FAST	<ul style="list-style-type: none"> <li>✓ auto Scaler for VM</li> <li>✓ snapshots</li> </ul>	





**elastic IP addresses** public IP Amazon provides

no cost = VM + IP address

cost = ~~VM~~ + IP address (not released)

## Pricing Models

### On Demand

Short term unpredictable workload ex/use 50% other 50% cheaper

### Spot Instances

EC2 instance 90% of on demand \$\$ Urgent needs for Large

Amazon will terminate when they need w/out notification

### Reserved Instances

1 to 3 yr deal

### Dedicated Host

Physical Server for me

### Per Sec Billing

manage instances running for irregular time

ex/dev/testing

data processing

batch processing

Serverless computing

## Troubles

- if power on/off = all hard disk data is LOST  
use snapshot
- IP addy is random
- Can NOT turn off public IP addy
- don't forget to terminate instance

# Simple Storage Service (S3)

Ch 10.1

- \* Upload, download, store data
- \* Bucket store data
  - ↓
  - container

100 buckets per account

no limit of # of objects  
5TB >



Stored in location NEVER leave

Can NOT modify / append data

SOL: copy object

make changes

store as new obj

delete old obj

get unique key

\$

.023 per GB/Month for 1st 50TB  
.022 per GB / Month for 450TB

- \* download data anywhere

enable others &amp; charge via amazon dev pay

only through EC2

faster  
FS

## Elastic Block Storage (EBS)

≈ external hard disk

have VM w/ data → kill VM → download data → create new VM → upload data

Solution: VM + volumes

- 1 GB to 16TB
- 20 per account
- stores snapshot
- file storage

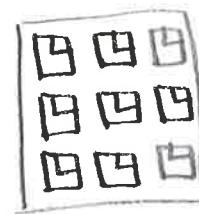
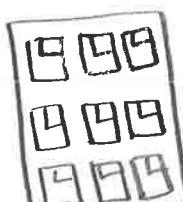
- \* automatically replicates to prevent data loss



Instance / hr

Volume / Month

Network / Month

20 volumes < 16TB  
? snapshots

# Management Provider

## Elastic Fox

- Firefox for managing Amazon EC2
- Hybrid Fox (fork of elastic fox)

## Amazon Management Console

EC2 & EBS

ISV : independent Software Vendors

My own DB

S3

S3 Firefox Organizer

# Private Cloud Provider

for circulizing resources

## Eucalyptus

elastic  
utility  
computing  
computing  
infrastructure  
systems  
security  
logistics  
systems

Web based version of EC2  
function as a software overlay

### purpose

- ✓ tech preview platform for Public Clouds
- ✓ basic dev platform for open source community

everyone  
uses  
openstack

# OpenStack

founded by NASA & Rackspace

EC2 **NOVA** start VM → computational

S3 **SWIFT** store objects

**Glance** store image

EBS **CINDER** block storage

**Compute**  
**Volumes**  
**Network**

glance - img  
Nova - compute  
Swift - obj

# SciCloud

1<sup>st</sup> private cloud @ UT

uses Eucalyptus & OpenStack

# Scaling Applications

## Scaling Information System

Mandatory  
enterprise app deployment

- \* **fault tolerance** - if 1 node fails ALL do NOT
- \* **high availability** - servers are available 99.999% of the time
- \* **scalability** - nodes support a specific load

when Load  $\uparrow$  System should  $\frac{\text{Scale up}}{\text{Scale down}}$

1 Server can support 8-10 req in 100 ms

Large - - - - - 12 req

Med - - - - - 8 req

Small - - - - - 6 req

## 2 Scaling Models

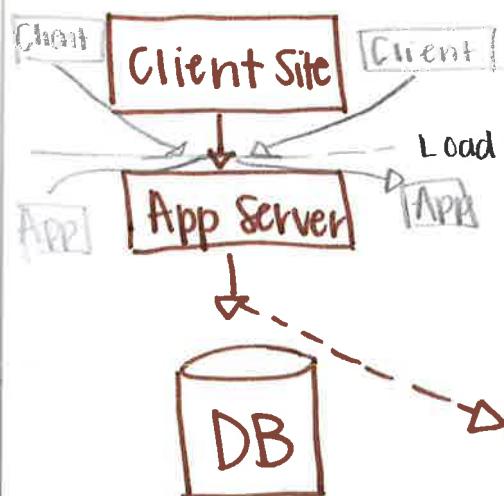
### ① Vertical Scaling $\uparrow$

as load  $\uparrow$  replace current node with more powerful machine  
CON: risk of losing current running job during switch

### ② Horizontal Scaling $\leftrightarrow$

ex / small to medium  
add servers w/ existing servers  
add more nodes to system

ex / MapReduce & Hadoop



Load Balancer - distributes incoming traffic to multiple servers

ex : sites to servers

✓ response time improved

✓ web server farms

✓ hides complexity from site

when pages are dynamically created and clicked multiple times then it gets stored so it doesn't need to be regenerated multiple times

# Load Balancer

## Types

### (A) Network Based

DNS server routes IP  
google.com [US IP - Server  
EU IP - diff server]

### (B) Network Layer

hot flexible

all requests coming from 1 IP get routed to 1 Server

### (C) Transport Layer

all request coming from 1 location → 1 server  
useful if connections are short & established frequently (room 235)

### (D) Application Layer

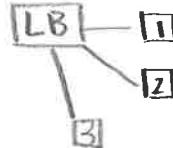
buy something online → Cart



## Classes

### Non adaptive

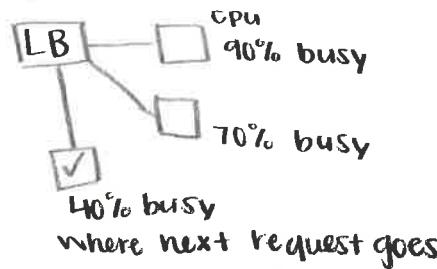
doesn't check current load



ex: Round Robin

Load Distributer

### Adaptive

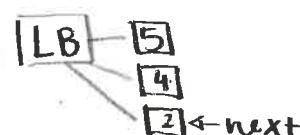


## ALG

Random

Round Robin

Join-Shortest-Queue



Ex

↓  
Nginx in lab  
HA prox  
Pern

# Testing w/ Simulation Load

concurrently

- Benchmark Tools

Tsung, Jmeter

- Multiple Protocols (supported by benchmark tools)

HTTP, XMPP

SSL support

## Bottleneck:



can only support 25 requests (SQL read) \* 1 or 3w\*

SOL: NOSQL or MongoDB

Scalable & elastic

## Auto Scaling

dynamic reaction to defined metrics by scaling

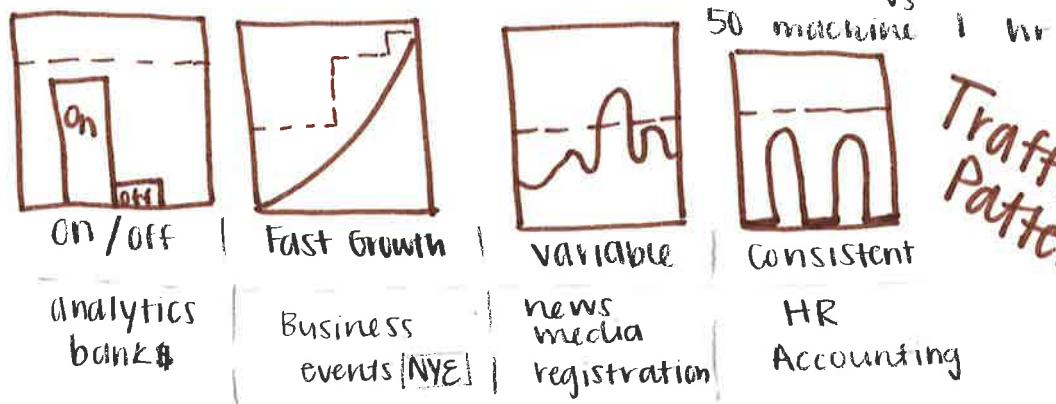
- ✓ ↑ availability
- ✓ Cost Saving
- ✓ Energy Saving

ex

website w/ unknown visitor #  
viral marketing campaigns  
apps that need elasticity  
Scientific app

1 machine 50 hr

vs  
50 machine 1 hr



**Scaling Policy** when to scale (typically threshold-based)

configure optimally is hard

**Resource Provision Policy** How to Scale

ex / Amazon AutoScale

Vendor Neutral

static LB + estimate load on each server

static LB + optimal resource provision policy

# Amazon AutoScaling

lecture 3 | 4

commercial

scale computing resources dynamically  $\Rightarrow$  predictability  
every Friday add 1

A EC2 instances  $\rightarrow$  AutoScaling Groups

define min/max # of instances prevent DDoS attack

## Amazon CloudWatch

monitor EC2 instances  $\rightarrow$  all metrics  $\Rightarrow$  set alarms  
CPU / mem / ...

## Elastic LB

distributes traffic across multiple EC2 instances

\* detects health

diverts traffic from bad ones

## my own autoscaling System

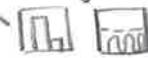
LB

## Measure Performance of Current Setup

- \* Linux shell commands
- \* Collectd - install on every server  $\Rightarrow$  to master
  - no performance impact
- \* Cacti - gives graphs too
  - 20%  $\downarrow$  performance

## Scaling Policy

Time - every Friday @ 1pm



Reactive - threshold: if  $> 75\%$  in last 5 min



Predictive - next = mean of last 5 min



## Resource Provisioning

cloud servers are heterogeneous

Simple - estimation based on Heuristic (if ... then ...)

Optimal - various instance types (diff power & price) make policy about current configuration add servers!

Consider you pay per hour so you don't want hour = START  $\rightarrow$  kill

what is the current machine?

what am I paying for?

ex/ 2-S  $\rightarrow$  12  
1-M  $\rightarrow$  8  
1-L  $\rightarrow$  36  
56  $\downarrow$  add S or kill S

ex

S	6	\$0.25/h
M	12	\$0.4/h

load 6 → use 1 small ✓

load 912 → SOL 1

add 1 small = total 50¢

$$\begin{aligned} \text{COST} &= \text{Cost of } 2S - 10 \text{ min} && \text{but NO Time Lost} \\ &= 50 && - .04 \\ &&& \quad 10 \text{ min profit} \end{aligned}$$

46¢

SOL 2

Kill Small ⚡ add med = total

$$\begin{aligned} \text{COST} &= \text{Cost of } m + 10 \text{ min of } S \\ &= 0.4 && + 0.04 \end{aligned}$$

44¢



## Definitions:

**Region** — 1 task w/ independent characteristics  
own capacity of instances

## Instance

**Type** — each region = MULT instance types

**Time bag** — time → interval where an instance is

**Killing Cost** — \$ LOST when instance is killed before full pay period

**Retain Cost** — \$ of lived duration pay period

## Apply This

① identify scalable components ex: not db → instances

② parameters of ↓ ex: # of small support? 6

③ combine w/ current load

Amazon  
AutoScale

providing the same  
snapshot for all servers  
even new ones added  
during scaling

## Dynamic Deployment Template

ex Cloud ML ↳ TOSCA

# MapReduce

data processing from 1950's

handled @ scale

functional programming

+ distributed file system

Google 2007

1 map & 1 Red

complex

mult map & Red

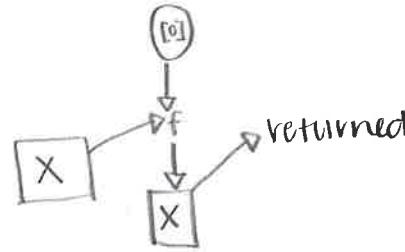
**Map:** • high order function apply function to  $[e1, e2, e3] \rightarrow [m1, m2, m3]$

map + lst

**Fold:** • high order fold (method(), [])  $\rightarrow$  single val

fold f x lst

initial motiv



ex

square  $x = x * x$

map square  $[1, 2, 3] \rightarrow [1, 4, 9]$

ex

fold (+) 0 [1, 2, 3, 4, 5]  $\rightarrow 15$

fold (\*) 1 [1, 2, 3, 4, 5]  $\rightarrow 120$

ex

fold (+) 0 (map square [1, 2, 3])

## Implicit Parallelism

Map & Reduce tasks concurrently  
elements are independent

- partition list between computers

- apply map to each partition

- combine map results w/ fold

associative: order does NOT matter  
 $(a+b)+c = a+(b+c)$

//  
parallelisation

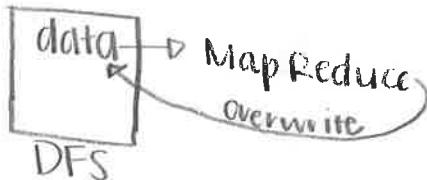
# LARGE Data Probs

- Map
  - \* iterate over Large # of records
  - \* Extract something from each ex: want ID
  - \* shuffle & sort results group ID's w/ "3"
- Reduce
  - \* aggregate results sum of all ID's w/ "3"
  - \* generate Final Output

## MapReduce

Input:  $[(k, v)]$  list  
Line, Line Content

- designed for processing Petabytes
- deployed on Large Cluster of servers
- Scalable (add servers)



**map:**  $(K, V) \xrightarrow{\text{mult } (K', V')} (K', V')$

User defined

function for map  
" " " reduce

**reduce:**  $\underset{\substack{1 \\ \text{all vals} \\ \text{sum}[V']}}{(K', [V'])} \xrightarrow{\text{O}} (K'', V'')^*$

1 val per key

framework takes care of handling map reduce (failed nodes)

dd (1) data distribution

S partition & replication

(2) scheduling

assign workers to map & r

(3) Synchronization

Sorting int data

error hand

(4) ERRORS

detects & restarts

SLOW WORK (5) Slow Workers & re executes

# Partitioner

how map output keys are distributed across reduce

$$\text{Partition}(\text{key}, \text{partition}) = \text{HashFun}(\text{key}) \bmod \text{partitions}$$

$$= 0 \dots 99$$

Partition Group  
 A - 1  
 B - 1  
 E - 3  
 D - 2  
 C - 1

$$\text{partition}(\text{key}, \text{partition}) = \text{Hash}(\text{key}) \% \text{ partitions}$$

# Combiner

optimize network traffic  
 mini reducer stage in mem

$$\text{Combine } (\text{k}', [\text{v}']) \rightarrow (\text{k}'', \text{v}'')^*$$

--- single val

---- (key, single val)  
 passed into reduce

- ① Reduce waits for Map to finish  $\rightarrow$  stored in mem intermediate data
- ② data is ordered by key before reduce

Who implemented MapReduce?

- Google (C++)
- Hadoop
  - opensource by Yahoo, now Apache
- Custom... ex: DB = CouchDB

how to move data 2 worker?

DON'T

- \* limited bandwidth
- \* not enough RAM to hold data in Mem
- \* disk access is slow

move worker 2 data

local disks on the node in cluster

DFS

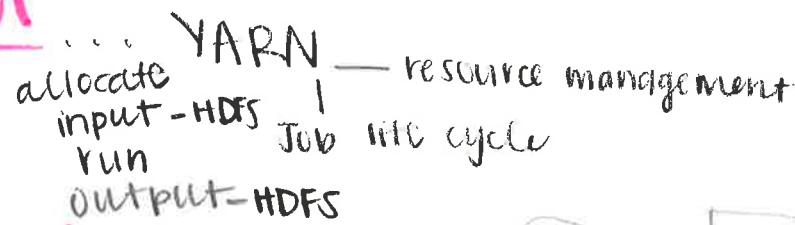
ex: GFS google  
 HDFS hadoop

# DFS

- pick cheaper commodity hardware new scale ↘ ↗
- # of Huge files (multi gb files)
- Once a file is written... must append to it generate new file or overwrite it
- Large reads over RAM
- \* store files as chunks VS GB
  - easier to run Map tasks on (64 mb)
- \* replicate chunk across 3 nodes
- \* 1 MASTER Node handles access to RAM (contains Namenode manager)
  - garbage collection
  - re replication / balancing
- \* no caching
- \* simplified API

## Hadoop Model

① Create cluster



② put data → into blocks  
replicate store in cluster

③ Job  
Run ... copy map 2 nodes

task - run map / reduce on block

④ store results in HDFS

task attempt - if it fails mult node will be abandoned

more nodes = better performance

scalability probs

master node has too many jobs  
resource utilization  
if fails abandons it  
node

Scales to petabytes of data

# MapReduce Alg's

How to get output from input?

Map

$\text{Input } (k, v) \rightarrow \text{Output } (k', v')$

partition( $k, v$ ) → hashfun( $k$ ) % partition  
→ part ID

Combine Reduce

$\text{Input } (k', [v]) \rightarrow \text{Output } (k'', v'')$

## Counting URL Access Frequency

Map

$\text{Input } (\text{Line\#}, \text{Line}) \rightarrow \text{Output } (\text{URL}', 1)$  if appears in Log

Reduce

$\text{Input } (\text{URL}', [1, 1, \dots]) \rightarrow \text{Output } (\text{URL}'', \text{Sum of } [ ])$

\*Sim 2 word count\*

## Distributed Grep

find all rows in .txt files that contain Regular Expression

Map

$\text{Input } (\text{Line\#}, \text{Line}) \rightarrow (\text{Line\#}', \text{Line})$  if has regular expression

w/alt  
Shuffle

Reduce

.... Identity Function....

## Inverted Index

word is in [file1, file2]

Map

$\text{Input } (\text{Line\#}, \text{Line}) \xrightarrow{\text{get each word from line}} (\text{word}, \text{File name})$

map (line#, Line, context){

pageName = Context.get inputsplitfilename()  
for each word in Line:

emit(word, pageName)

Reduce

$\text{Input } (\text{word}, [\text{filename}]) \rightarrow (\text{word}, "[\text{filename}]")$

reduce (word, values){

pageList = []

for each pageName in values:

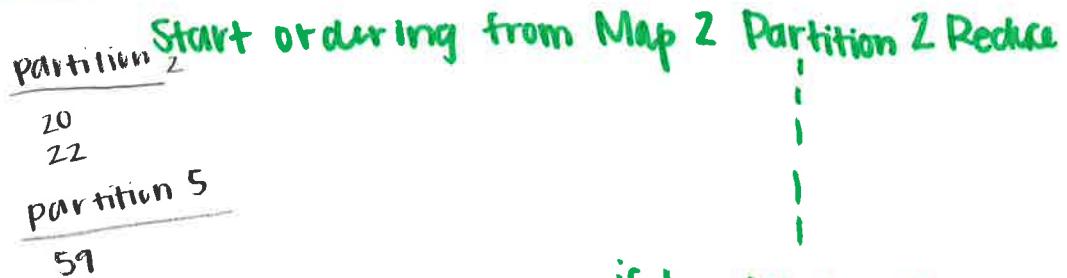
pageList.add(pageName)

emit (word, str (set (pageList)))

# Distributed Global Sort

22
59
64
30
50
20

.txt + S



Map

(Line #, Line) move VAL to key  
 $\rightarrow$  (Line, "")

Reduce

..... Identity function .....

partition  $\rightarrow$  Reduce 0

Reduce

0 < Reduce 1

## Term Co-occurrence Matrix

# i  $\exists$  j in context

terms<sup>2</sup>

keys are large so combiner  
doesn't help

Keys  
Term 1 T2 T3

T1	0	1	5
T2	1	0	1
T3	4	Count of val $T_2 \neq T_3$ in input data	0

"You shall know a word by the company it keeps"

Why? semantic distance

is the sentence real? compute frequency  
of words together

## Map Reduce is inefficient bc synchronization

Separate tasks = need synchronization (hold map data in mem)  
partitioner

X Pairs

("T1, T2", 1)

X too many pairs

("T1, T2",  $\Sigma$  1)

X lots of shuffling

Stripes

$\Sigma$  ("T1, T2", 1)

$\Sigma$  ("T1, T3", 1)  $>$  T1  $\rightarrow$  {T2:1, T3:1}

X storage allocation  
X code = hard

Combine all sentence's T1  $\rightarrow$  {T2:1, T3:5, T4:2}

What is the prob that TermB if A is already in sentence?

$$P(B|A) = \frac{\text{count}(A,B)}{\text{count}(A)}$$

#1 Pairs

- ① gives  $\text{count}(A, B) = ("A, B", 1)$
- ② get count(A) =  $("A", 1)$   
get val 1st

how? Partition all "A"  $\ni "A, _"$  to 1 reduce

COUNT A = Store in mem

might crash

$$\frac{\text{count} "A, _"}{\text{count} A}$$

emit COUNT A

emit COUNT "A, \_"

cheating

#2 Stripes

compute "A, \*"

compute  $P(B|A)$

Trade offs between pairs  $\ni$  stripes

A Creating (K, V) pairs = lots of time

B Size

C Combining ✓ always better to combine values in Map

Complex data types in Hadoop

easy - string

regular expressions

hard - create write comparable

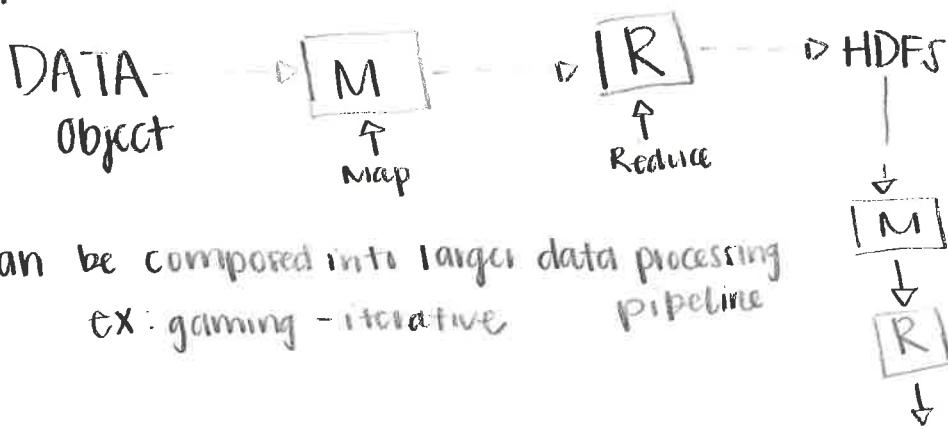
# MapReduce Ex: Lab 5

why would you use Hadoop?

- \* Extract, Transform, Load pipelines  
ex: yelp review dataset → analyze who has the best positive reviews
- \* Reporting & Analysis  
ex: websites when data is too large for sql database
- \* Machine Learning  
ex: recommending products to users  
ex: apply java library to image to abstract face for ppl recognition as key & val
- \* graph algorithm  
ex: trends in social media  
page rank

## Map Reduce Jobs

- \* typically short, code wise  
ex: identity reducer
- \* represents a data flow, rather than tasks



- \* Can be composed into larger data processing pipeline  
ex: gaming - iterative
- \* data MUST be partitioned across many reduce tasks that write out to diff HDFS files ] IF Large
- \* diff input formats

## Designing MapR Alg

Goal: what is the desired output from the input data?

### ① Map Function

- what is Map Input(key, val) ↗ function
- what is Map Output(key, val) ↗ function

### ② Reduce Fun

- what is Input(key, [value]) ↗ function
- what is Output(key, value) ↗ function

## Term co occurrence matrix

# of times word is in sentence

Terms<sup>2</sup> = Matrix size

Keys are large → combiners  
don't help

key	Term <sub>1</sub>	Term <sub>2</sub>	Term <sub>3</sub>
Term <sub>1</sub>	0	2	3
Term <sub>2</sub>	1	0	1
key Term <sub>3</sub>	4	1	0
		Count of val	= depends on input data
		T <sub>2</sub> , T <sub>3</sub>	7

why? how often do terms occur near each other → occurs in —

figure out if this is real eng sentence or not

#1 Pairs approach [ ("T<sub>1</sub>, T<sub>2</sub>", 1)  
key , val ] ✓ easy implementation  
✓ easy understand  
✗ too many pairs

#2 Stripes [ combine (k, v) pairs  
emit a → { b-01 ; c-03 ; ... } ↗ step map ↗ r between  
array or dictionary ↗ less short ↗ shuffled ↗ better combiners  
Reducer: count based on 2nd key ↗ diff to implement

- \* Map & reduce are independent so you process each k, v
- \* tools for sync
- keep int values in memory in hashmap (combine in map)
- sorting functions for keys
- broadcasting... assign value to all values (loop)

conditional Prob(B|A) ...

Reduce → (k, [v]) ↗  
(v<sub>1</sub>, [v<sub>2</sub>]) ↗  
! (A, 1) ↗ complex key

SORT ↗ partition  
only use 1st part

how can we force all A ↗ A<sub>1</sub> — comes to certain map  
compute (a, \*) = # of a then (a, b<sub>1</sub>) ... (a, b<sub>...</sub> ) ↗ reduce

Stripes ... prevent all these prob ↗ (a, \*) . . . . .

# Platform as a Service

Manages... infrastructure

+

Software environment

- \* multi tenant architecture
- \* built in scalability
- \* integrated w/  services ↗ DB
- \* simplify prototyping
- \* \$ for ONLY <sup>cloud</sup> services used      upfront \$ ↓
- \* Vendor Lockin

## Type 1: web.app

- \* built on IaaS 
- \* open-source computing platforms (not tied to IaaS provider)
  - ex: Red Hat OpenShift deployed on diff clouds
- \* social app platforms (upload my own game 2 Facebook)
  - Add ons for Facebook / Google +

ex: google AppEngine, exception bc before IaaS

## Type 2: Function as a Service

platform for hosting microservices / Serverless

"App" callable function

Trigger: upload new image to 

Precondition: file size > 10 MB

Execute: resize\_image(file)

ex: AWS Lambda 



Apache OpenWhisk

# Type 3: Data Processing as a Service

deploying data processing applications

\* provider manages environment & cluster

User → load: app + resource limit

ex: AWS Elastic MapReduce

Google Cloud DataPro

## Google AppEngine

- \* PaaS for dev & hosting apps
- \* easy to build, maintain, & scale
- \* nothing to manage unless you want to optimise

✓ persistent storage → load files

✓ scales servers → autoscale

✓ dynamic traffic

✓ asynchronous task queues

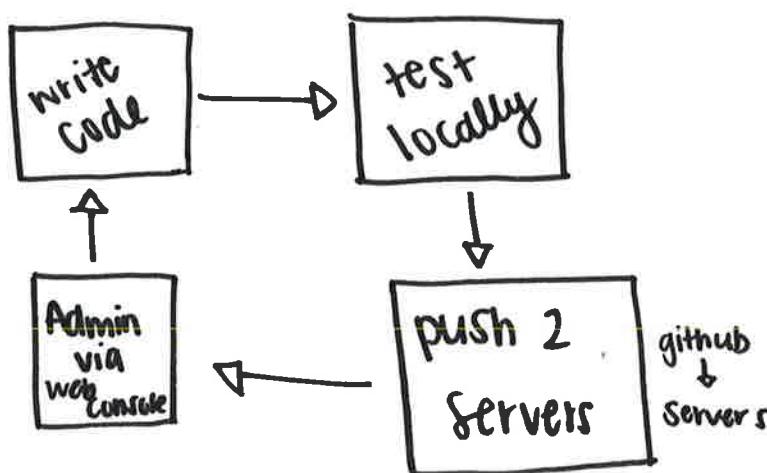
✓ schedule tasks "delete file every night"

✓ integration w/ Google Cloud Services

✓ secure, sandboxed, environment

isolated environment

h/t them login w/ google  
run instances in same server



## Available Cloud Services

Google Cloud SQL fully managed relational db

Datastore schemaless objectstore ... GQL

Blobstore too big for datastore  
load file > get key ... generate links

## Available Data Services

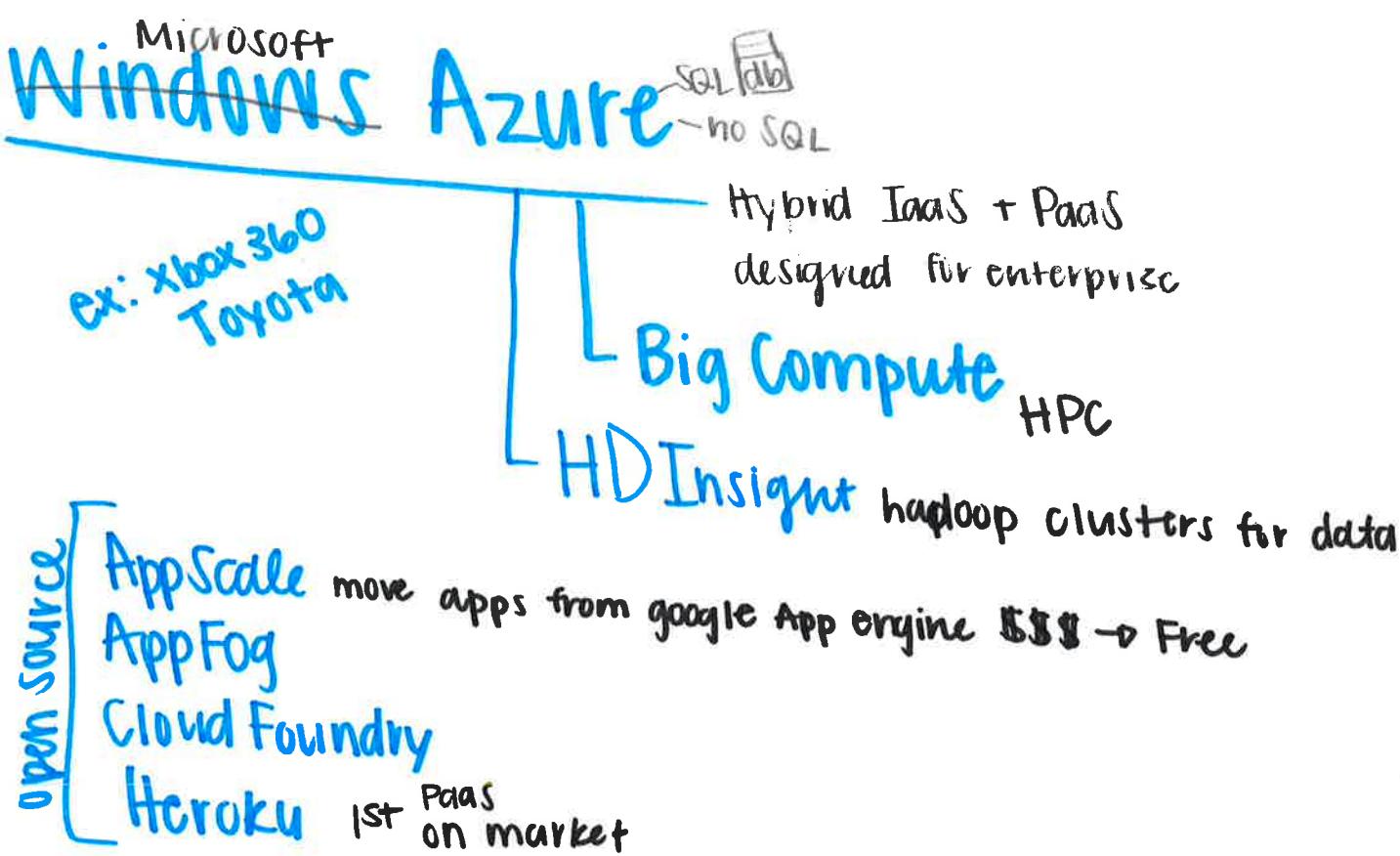
Search through structured data ... html, .txt

Memcache improve performance

Logs monitor performance

Remote lets external apps access Engine services

## AWS Elastic Beanstalk



# AMAZON

## Management Console

90+ services

### IAM: Identity & Access Management

how can I group use 1 credit card?

- \* Manage users create groups
- \* Manage credentials assign security status
- \* Manage permissions create policy per user  
↓ generality

### CloudFormation

template.json

for hardware & software

free (only pay for services)

resource " + " ↓  
in what instance

### Relational



Specify scaling & SQL

done by YOU not Amazon  
easy configuration

### Aurora

specify type of SQL

+ My... & PostgreSQL

they manage

5-10x faster

### NO SQL

Document 

Dynamo - NO SQL 

Neptune - graph 

they manage

### Pipeline

data movement BPM



Moving Data

\* **Snowball** - 50 to 80 TB per Month

" **Edge** - pre processing

- + launch instances (26 - 52 vCPU)
- + AWS Lambda
- + EC2 AMI
- \* GPU
  - |
  - Machine Learning
  - real time video analysis

100 TB

CIVISTER 20 Nodes  
1 cloud server

1 Showmobile - Petabytes to Exabytes

"

1250 Snowball

## Serverless?

Servers are abstracted from developers

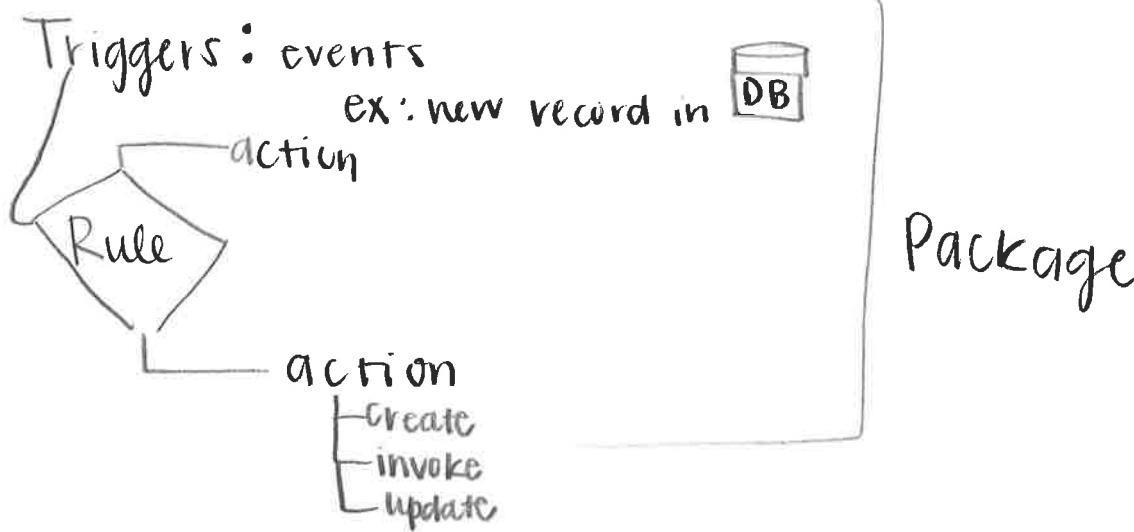
event driven + vs constant running  
if there is a trigger

\$ per second  
Sub

Best fit newer workloads ..... Amazon Alexa hears u  
linked w/ Cloud not standalone

Ex: **Amazon Lambda**  
use multi lang  
Open ....

IBM blawhisk  
based on OpenWhisk  
function is scale  
any lang



- ✓ easy 2 scale
- ✓ rapid Prototyping
- ✓ easy 2 modify functions
- ✓ pay 4 execution
- ✓ applications w/ multi languages
- ✗ hard 2 avoid vendor lock in \$\$\$ 2 migrate
- ✗ tools 4 debugging (hard 2 verify error)
- ✗ architecture is complex
- ✗ slow cold-start  
not kept in mem
- ✗ harder to predict #
- ✗ need stateful computations