Databases and Big Data

Cloud Computing

Recap

- Big data is everywhere
- Big data carries (potentially) big value
- Big data needs big infrastructure
 - O Who owns it?
 - Who pays for it?

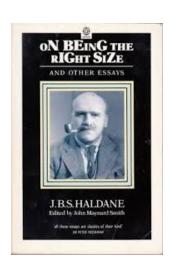
Design Principle

Principle of incommensurate scaling:

- Changing a parameter by a factor of ten usually requires a new design
- "On being the right size"

That means:

 Supporting big data workloads requires different design to RDBMS.



Today

- User's view
- Systems' view
- Google Cloud Example
- AWS Overview

- Informal: computing with large datacenter
- More accurately:
 - Computing as a utility
 - Do your work on someone else's infrastructure

Above the Clouds: A Berkeley View of Cloud Computing

Michael Armbrust, Armando Fox, Rean Griffith, Anthony D. Joseph, Randy Katz, Andy Konwinski, Gunho Lee, David Patterson, Ariel Rabkin, Ion Stoica, and Matei Zaharia (Comments should be addressed to abovetheclouds@cs.berkeley.edu)

UC Berkeley Reliable Adaptive Distributed Systems Laboratory * http://radlab.cs.berkeley.edu/

February 10, 2009

Implication

Computing as utility



Economics (got to be cheap)

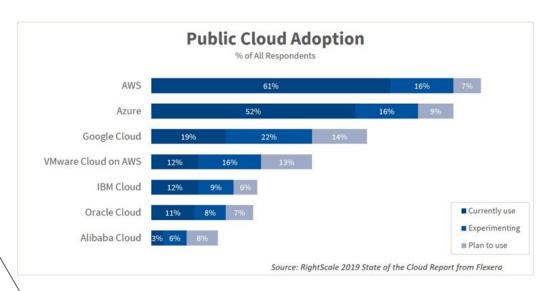
Quality of service (SLA)

Computing on other people's infrastructure



Security, flexibility (your house, your rule)

- Public clouds
 - Anyone can use
 - Arrays of options
- Private clouds
 - Similar interface, but
 - Smaller
 - Privately own









Different design choices

- Public cloud: our focus
 - Infrastructure as a Service (laaS)
 - Platform as a Service (PaaS)
 - Software as a Service (SaaS)

VM, disks, etc.

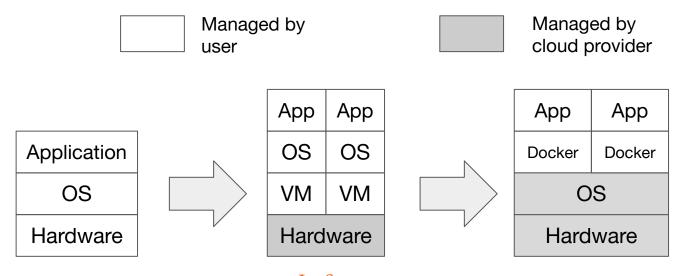
AppEngine, ERM

Github, GoogleDocs, AutoML

(Almost) Everything is a service



• Who manage what?



IaaS

EC2, Azure, etc.





PaaS

Private Cloud

Who manage what?

Let Cloud manage VM/Docker image as well e.g. Lambda Also another PaaS





Managed by cloud provider

App	Арр	
Docker	ker Docker	
OS		
Hardware		



Арр	Арр	
Docker	Docker	
OS		
Hardware		



App	App	
Docker	Docker	
OS		
Hardware		





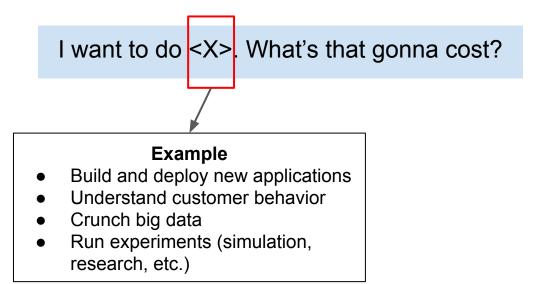




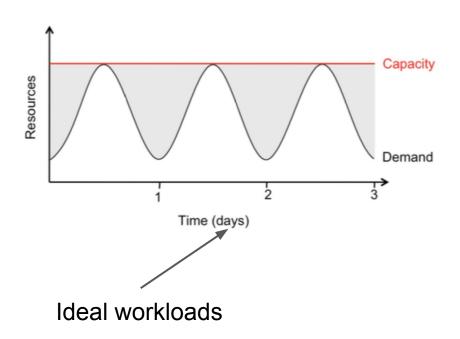




- To buy or to rent?
 - That is the question!



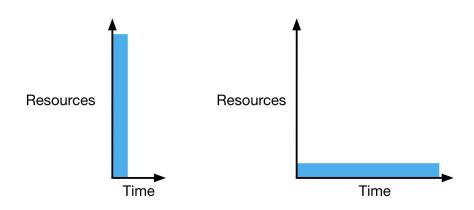
- Pay-as-you-go
 - Charged by service usage
 - VMs, disk, network
 - By minutes, hours, or bytes
 - No upfront cost
 - No deposit



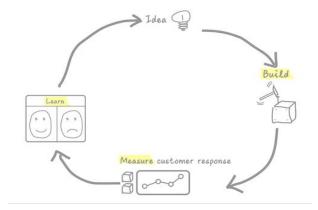
Elasticity

- 1000 VMs for 1 hour
- o vs. 1 VM for 1000 hours
- Same price, 1000x faster results

Get access to 1000s of machines in seconds



- Software as a Service
 - Faster time-to-market
 - Faster response to user
 - Faster access to state-of-the-art technologies



Faster innovation

	Buy	Rent
Utilization (using it all the time?)	100%	Vary
Flexibility (same hardware 5 years later?)	Yes	No
People cost	Yes	No
Opportunity cost (time-to-market)	High	Low
High availability	Costly	Cheap
Security	-	-

People: If you dont mind paying for engineers to build/maintain infrastructure

If you want to ensure High Availability -> Costly

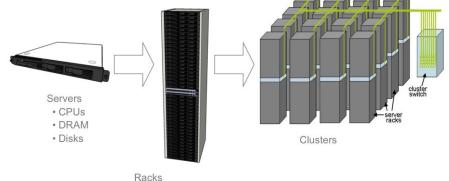
- Possible, because
 - Economies of scale
 - Cloud providers buy, manage hardware in bulk
 - Innovation in software





Cloud Computing - Systems' View

- The hardware
 - Rows and rows of rack-mounted servers
 - 1 data center ~ 100K of servers
 - 100MW (big part is cooling)
 - Few clusters per data center



- 10 00 ----

- 40-80 servers
- Ethernet switch

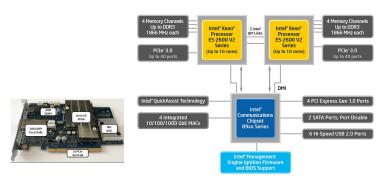
Where your cloud is





- Compute resources
 - o Initially: commodity CPUs
 - Now: GPUs, ASICs

2-socket server







• Storage:

o Basic: HDD, SDD

New: NVM, archival storage

Often separated from Compute:

NAS





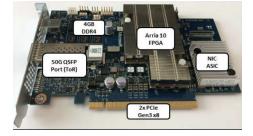


• Network:

- o Basic: 10-100GB, dumb NICs
- New: smart network
 - Programmable switches
 - Network co-processor
 - Smart NICs
 - FPGA







Wide range of choices

Standard	l	III	IV	V	VI
Systems	Web	Database	Hadoop	Haystack	Feed
СРИ	High	High	High	Low	High
	2 x E5-2670	2 x E5-2660	2 x E5-2660	1 x E5-2660	2 x E5-2660
Memory	Low	High	Medium	Med-Hi	High
	16GB	144GB	64GB	96GB	144GB
Disk	Low 250GB	High IOPS 3.2 TB Flash	High 15 x 4TB SATA	High 30 x 4TB SATA	Medium 2TB SATA + 1.6TB Flash
Services	Web, Chat	Database	Hadoop	Photos, Video	Multifeed, Search, Ads

[Facebook servers]

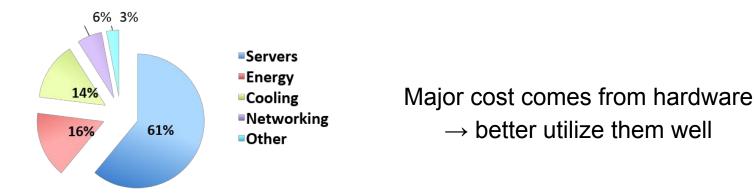
• Useful numbers, again.

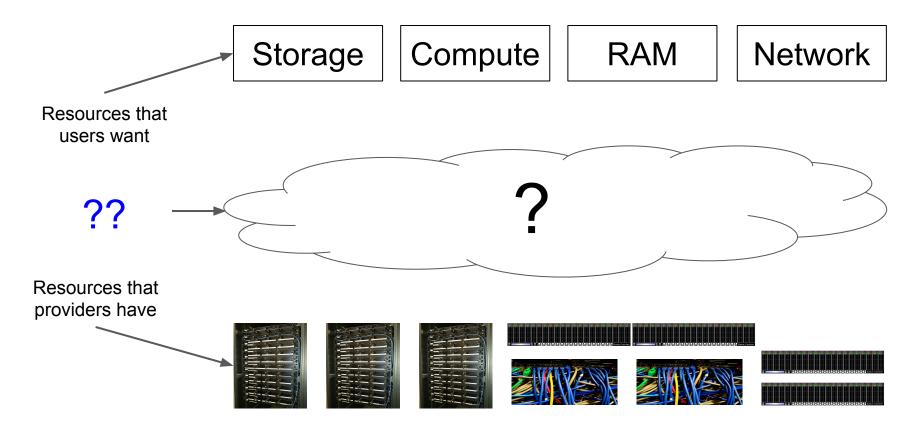
L1 cache reference	0.5 ns
Branch mispredict	5 ns
L3 cache reference	20 ns
Mutex lock/unlock	25 ns
Main memory reference	100 ns
Compress 1K bytes with Snappy	3,000 ns
Send 2K bytes over 10Ge	2,000 ns
Read 1 MB sequentially from memory	100,000 ns
Read 4KB from NVMe Flash	50,000 ns
Round trip within same datacenter	500,000 ns
Disk seek	10,000,000 ns
Read 1 MB sequentially from disk	20,000,000 ns
Send packet CA \rightarrow Europe \rightarrow CA	150,000,000 ns

Total Cost of Ownership (TCO):

[Source: James Hamilton]

- Capital Expense (CapEx): buildings, generators, cooling, hardware, etc.
- Operational Expense (OpEx): electricity, maintenance, people, etc.





- Design goal:
 - Manage data center resources
 - Provide services to users
 - While guaranteeing:
 - High utilization
 - Competitive SLA

Many of them

Different types of services

To justify TCO

To attract more users

Most difficult problem: failures

- ~0.5 overheating (power down most machines in <5 mins, ~1-2 days to recover)
- ~1 PDU failure (~500-1000 machines suddenly disappear, ~6 hrs to come back)
- ~1 rack-move (plenty of warning, ~500-1000 machines powered down, ~6 hrs)
- ~1 network rewiring (rolling ~5% of machines down over 2-day span)
- ~20 rack failures (40-80 machines instantly disappear, 1-6 hours to get back)
- ~5 racks go wonky (40-80 machines see 50% packet loss)
- ~8 network maintenances (4 might cause ~30-minute random connectivity losses)
- ~12 router reloads (takes out DNS and external vIPs for a couple minutes)
- ~3 router failures (have to immediately pull traffic for an hour)
- ~dozens of minor 30-second blips for dns
- ~1000 individual machine failures (2-4% failure rate, machines crash at least twice)
- ~thousands of hard drive failures (1-5% of all disks will die)

Not to mention human errors, software bugs



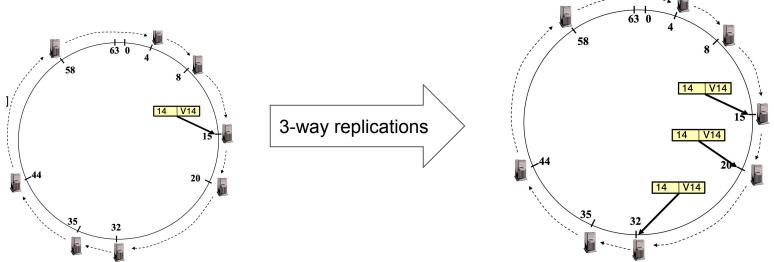
Availability/Reliability must come from software

Question: How to keep functioning under failures?

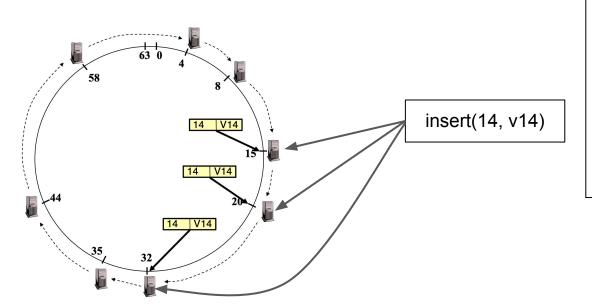
even if node 15 fails, user can go to node 20 and 32 to get it back

Cloud Computing

- How to deal with failures:
 - Redundancy redundancy redundancy
 - Example



Replication



Example

R(k): set of nodes storing k (replicas)

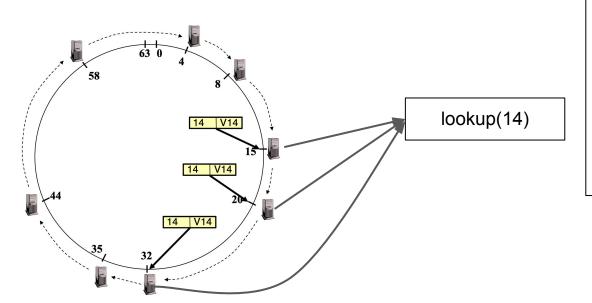
Insert(k,v):

Write to any node in R(k)
Return when any replica returns

Lookup(k):

Read from any node in R(k)
Return the first result

Replication



Example

R(k): set of nodes storing k (replicas)

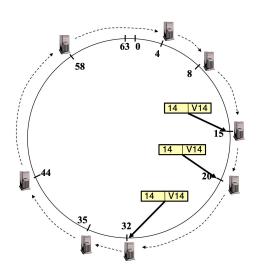
Insert(k,v):

Write to any node in R(k)
Return when any replica returns

Lookup(k):

Read from any node in R(k)
Return the first result

Replication



Eventual Consistency

- -> insert v14 to node 15
- -> then node 15 try to replicate node 20 and 32
- -> but before this replication finish, user read from node 32/20 which is the stale value, not the latest value.

-> read does not return the latest write

Example

R(k): set of nodes storing k (replicas)

Insert(k,v):

Write to any node in R(k)
Return when any replica returns

Lookup(k):

Read from any node in R(k) Return the first result

Remember what problems here?

- CAP theorem: in any distributed system, choose 2
 - Consistency:
 - Read returns the latest write
 - Availability:
 - Every request receives a response
 - Partition Tolerance:
 - Continue to work when network fails



give up on consistency

Question: How to scale?

Question: How to maximize utilization?

Question: How to isolate users?

Question: How to monitor?

. . . .

High-level stack

Storage Layer:
hide complexity of disk
(replicate data even if one disk fails,
still can retrieve data required)

Data Processing

Storage Layer

Resource Management













Hide any failures

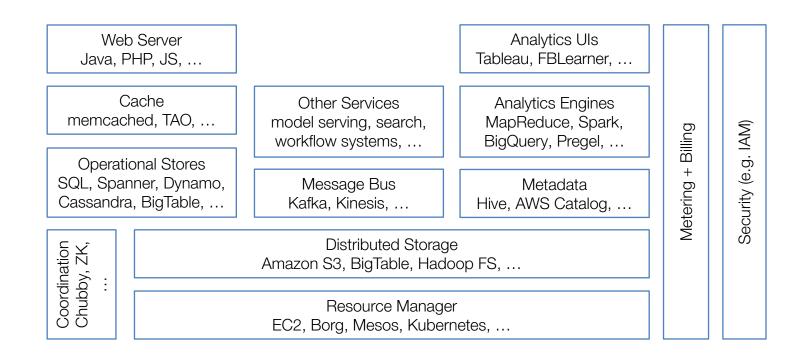
Hide disk failures



Manage hardware

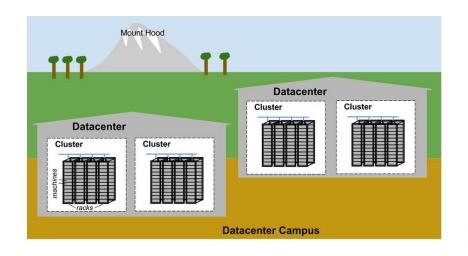
Cloud Computing

Cloud software stack in practice



Google Private Cloud

Google Cloud

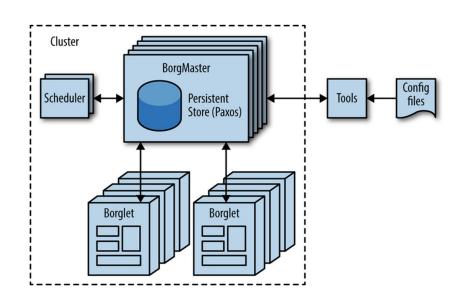




Google Cloud

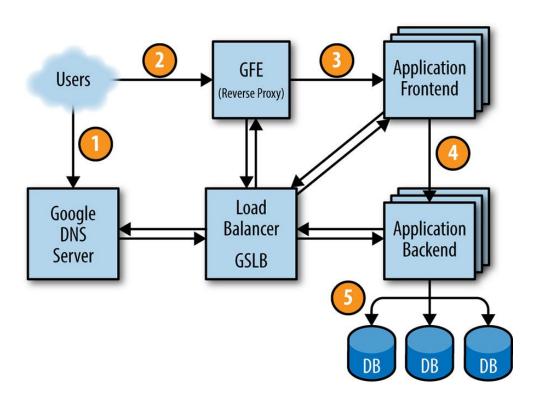
Storage layer Blobstore (small) **SQL** NoSQL Bigtable (big) Spanner Colossus File Local HDD or Flash

Resource management layer



Google Cloud

• Life of a request

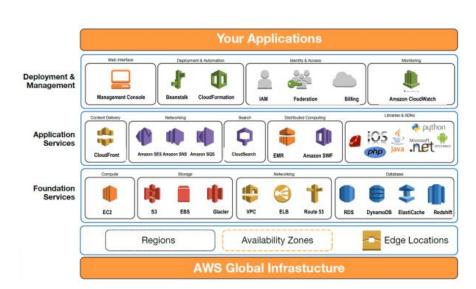


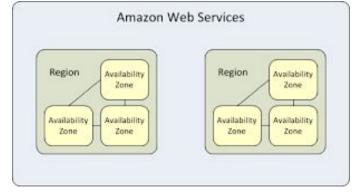
Amazon Web Service

- Most popular
 - The earliest
- Only cover key concepts here
 - You have \$100 to play with

Region

- Where servers are
- Geographical region, e.g. ap-southeast-1
- Multiple availability zones per regions,
 e.g. ap-southeast-1a





- Many services, but this course:
 - Elastic Cloud Computing (EC2)
 - Simple Storage Service (S3)
- We've seen S3
 - Simple bucket-based key-value store (Week 3)
 - Dynamo design (this week)

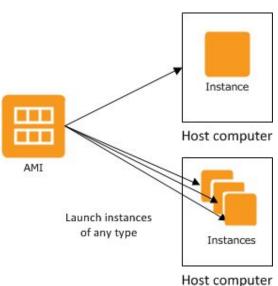




• EC2

- Amazon Machine Image (AIM): like a virtual machine image
- Instance: running image on specific hardware
 - Like a virtual machine
- Instance type: what hardware

Instance	vCPU*	CPU Credits / hour	Mem (GiB)	Storage	Network Performance
t2.nano	1	3	0.5	EBS-Only	Low
t2.micro	1	6	1	EBS-Only	Low to Moderate
t2.small	1	12	2	EBS-Only	Low to Moderate
t2.medium	2	24	4	EBS-Only	Low to Moderate



Free*

 Identity and Access Management (IAM):

Your account has 1 root user (when sign up)

Root user can:

■ Create IAM users

Grant different IAM user access to different

services

to different

IAM user 1

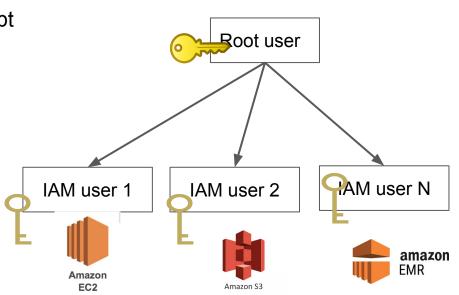
IAM user 2

Amazon
EC2

Amazon S3

Why all these troubles?

- 2 types of credentials
 - Access key:
 - Short, 40-character string
 - To invoke AWS APIs, must be kept secret.
 - SSH key:
 - Standard SSH key
 - To log-in to the EC2 instances



- Security group
 - Firewall: rules for network access
 - Each instance belongs to one security group
 - You can create many groups
 - Default option: publicly available



- How to interact with AWS:
 - AWS GUI console
 - Doesn't scale beyond a handful of instances
 - This course:
 - Use command lines
 - Or boto3 (Python)



Comparison

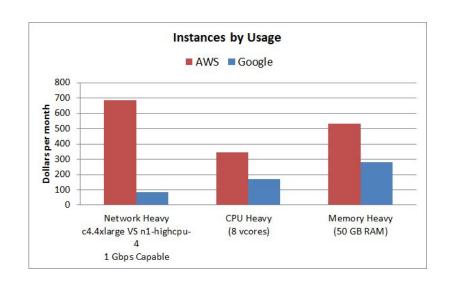
Basic tier: A0, A1, A2, A3, A4
Optimized Compute: D1, D2,
D3, D4, D11, D12, D13
D1v2, D2v2, D3v2, D11v2,...
Latest CPUs: G1, G2, G3, ...
Network Optimized: A8, A9
Compute Intensive: A10, A11,...

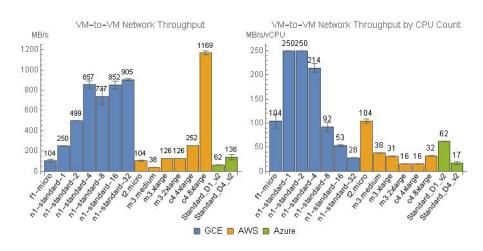
t2.nano, t2.micro, t2.small m4.large, m4.xlarge, m4.2xlarge, m4.4xlarge, m3.medium, c4.large, c4.xlarge, c4.2xlarge, c3.large, c3.xlarge, c3.4xlarge, r3.large, r3.xlarge, r3.4xlarge, i2.2xlarge, i2.4xlarge, d2.xlarge d2.2xlarge, d2.4xlarge,...

n1-standard-1, ns1-standard-2, ns1-standard-4, ns1-standard-8, ns1-standard-16, ns1highmem-2, ns1-highmem-4, ns1-highcpu-2, n1-highcpu-4, n1-highcpu-8, n1-highcpu-16, n1-highcpu-32, f1-micro, g1-small...

Microsoft AZURE Amazon EC2 Google Cloud Engine

Comparison





Summary

- Cloud computing levels the playing field
 - At least for Big Data
- Cloud computing not merely buying hardware in bulk
 - Many technical challenges
- From now on, whenever we speak about big data
 - Assume that it's on the cloud

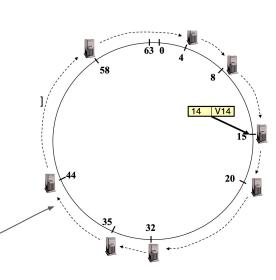
Question 1

- Consider Dynamo
 - Suppose the hashed keys are uniformly distributed in the ring space
 - How do you guarantee that servers have equal shares of the key space?
 - In this figure, they do not

Given a fixed number of machines,

- Master server can assign machine ids to be evenly distributed across the nodes
- Since keys are evenly distributed, all nodes will have **equal** share of key space

If you add a new machine, Need to disrupt the distribution a little bit. And accept the fact that it will not be **evenly** distributed



Question 2

Following up from Q1:

- Suppose the hashed keys are uniformly distributed in the ring space
- Suppose servers have equal shares of the key space
- Does the system guarantee that every server handles the same amount of work (load balanced)?

No,

- Input could have a certain distribution and the hash output follows a certain distribution.
- The read operations may be **heavily skewed**, some keys may be very popular and queried pt -> all go to one server which handles a **disproportionately** amount of work.

- If a few keys are very popular
- replicate them -> send to a few servers
- or redistribute them to some other nodes

