

# EC/CS 528: Cloud Computing

## Overview of Cloud Computing

Instructor: Alan Liu

# Announcements

- Piazza Discussion:
  - [piazza.com/bu/fall2022/eccs528](https://piazza.com/bu/fall2022/eccs528)
  - Access code: cloudcomputing
- User account and skill survey:
  - Group and project matchings.
  - Prepare multiple project choices.
  - Q&A.
- New project update:
  - MLOps with Databricks in Public Clouds

# How to start the project?

- Meet weekly with your mentor
  - Schedule a weekly meeting time
  - Record the meeting; mentors talk fast, being able to replay what they said can be super valuable
- Each person should say:
  - What have accomplished since last meeting?
  - What are you going to accomplish by next one?
  - Are they blocked?
- Don't be blocked until weekly meeting:
  - Set up mechanism to ask quick questions to each other, and to mentor, e.g., slack
- Remember you are a team

# Building the “cloud” from scratch - spec and buy

The screenshot shows a Dell product configuration page. The browser address bar is 'dell.com'. The page has tabs for 'Tech Specs & Customization', 'Product Details', 'Reviews', and 'Drivers, Manuals & Support'. Under 'Product Details', there are radio button options for 'Performance Optimized' (selected and 'Included in price'), 'Memory Mirroring', 'Fault Resilient Memory-Vmware', 'Multi Rank Sparing Memory Mode', and 'Single Rank Sparing Memory Mode'. A yellow warning box states: 'Your current hardware configuration exceeds this power supply unit's wattage. Please either upgrade the power supply unit or downgrade the CPU, memory, processor or network adapter.' Below this are two 'Help Me Choose' sections. The first section shows memory options: '32GB RDIMM, 3200MT/s, Dual Rank' (selected, 'Included in price', 'Support & Services price has changed'), '64GB RDIMM, 3200MT/s, Dual Rank' (\$1,336.78 /ea.), '16GB RDIMM, 3200MT/s, Dual Rank' (\$360.00 /ea.), and '8GB RDIMM, 3200MT/s, Single Rank' (\$236.01 /ea.). The second section shows RAID options: 'C1, No RAID for HDDs/SSDs (Mixed Drive Types Allowed)' (selected, 'Included in price') and 'C2, RAID 0 for HDDs or SSDs (Matching Type/Speed/Capacity)' (\$0.00). On the right, a 'Special Offers' section shows a 42% discount with code 'SERVER42'. A pricing summary shows a List Price of \$43,205.22, Total Savings of \$16,013.24, Shipping is Free, and a Dell Price of \$27,191.98. A green 'Add to Cart' button is at the bottom right.

memory

AID



Then receive and assemble...



then you have to run it...

# Issues?

[What do you think?]

- People and skills
  - N areas of expertise =  $O(N)$  people
- Scaling?

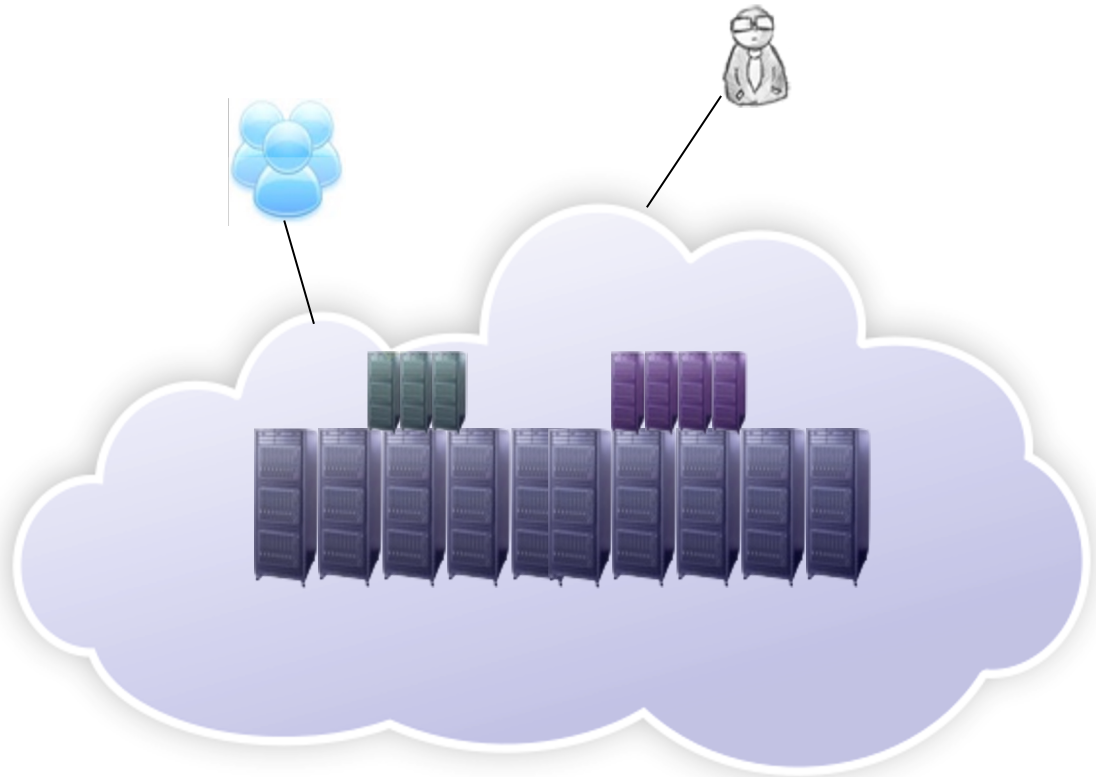
# Why is Cloud Computing transformative?

- Major change in computation is managed and used:
  - Economics of central utility: Price of computers, Operational efficiency, Location (e.g., cheap power, distribution), Co-location other customers, Utilization shared capacity, shared services (e.g., DR)
  - “As with the factory-owned generators that dominated electricity production a century ago, today's private IT plants will be supplanted by large-scale, centralized utilities.” -- Nicholas Carr
- Availability of massive capacity on demand; elastically scale up and down:
  - Startups don't need to be acquired by Google or MS: a startup won't get money today to buy HW.
  - What happens when massive HPC becomes available to everyone?
- Gets rid of key impediments for developing & distributing SW
  - Avoids need for broad HCL, OS support, ... many highly specialized software products...

# Cloud in a nutshell

- On-demand access
- Economies of scale

**All computing will  
move to the cloud**





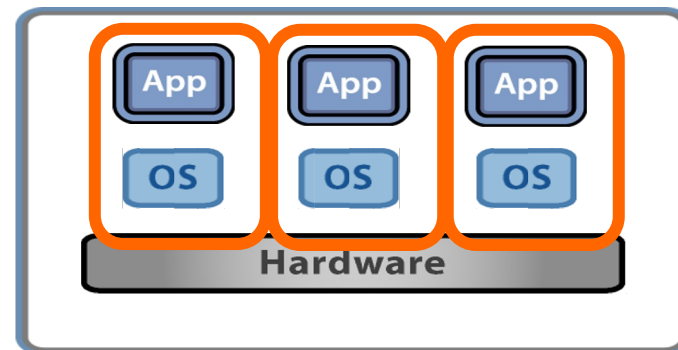
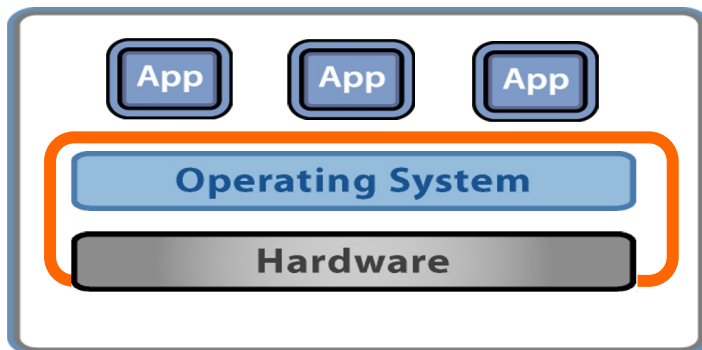
# This is really nothing new...

## Original vision of Utility/grid computing:

*"If computers of the kind I have advocated become the computers of the future, then computing may someday be organized as a public utility just as the telephone system is a public utility... The computer utility could become the basis of a new and important industry."*

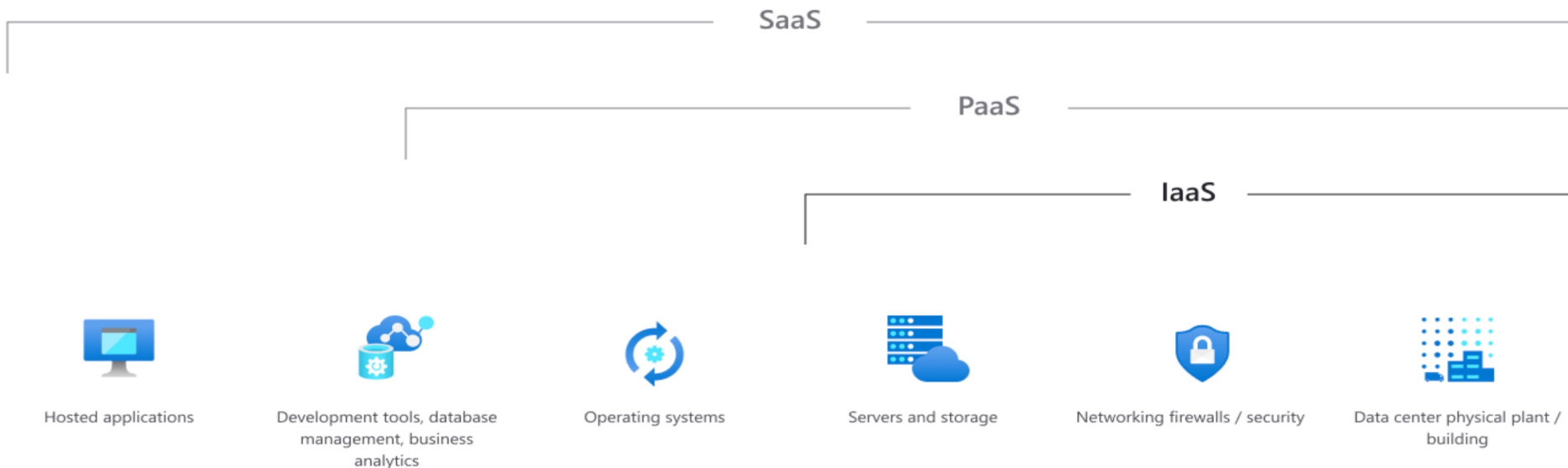
When was this statement from?

Why now?



# Layers of Cloud

- Infrastructure as a Service (IaaS): AWS, Azure, OpenStack, MOC...
- Platform as a Service (PaaS): Salesforce's Force.com, Google App engine, AWS, MSFT Azure
- Software as a Service (SaaS): Hosted applications: Gmail, Facebook, Google docs, eBay



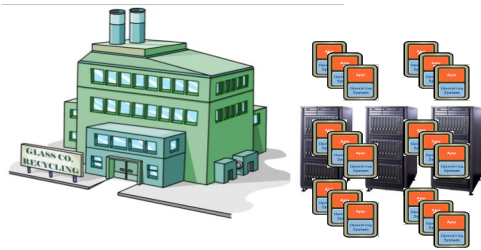
# Motivation for using cloud

- Cloud is not inexpensive today
  - 2-20x more expensive than local
- Administrators do not come in fractional units; if you are small cheaper
- Offers elasticity: can deal with massive fluctuations on demand
- Offers huge variety of services:
  - cloud provider can afford to amortize cost over a huge number of customers

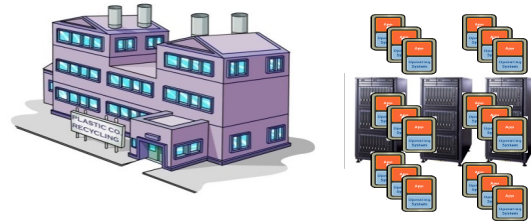
# Examples

- Microsoft's [Azure](#)
- Amazon's [AWS](#)
- Google's [Cloud Services](#)

# Remember this?



Host it R us.

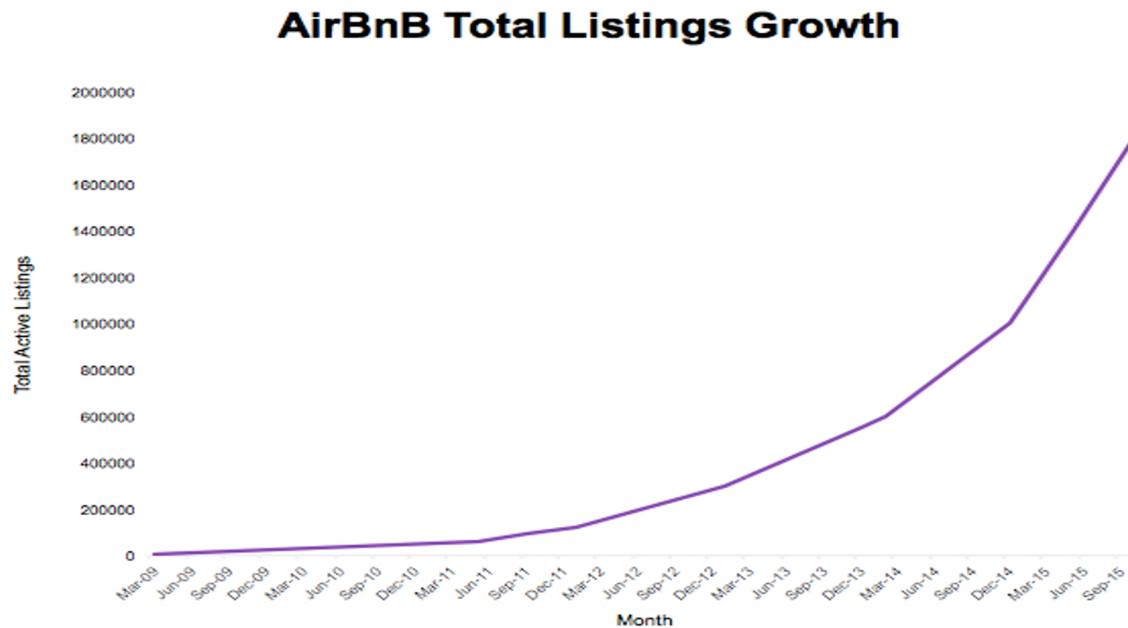


Host 4 Less



# AirBnB Example

- Success of market depends on network of renters and landlords;
  - starts really small



# AirBnB

<https://aws.amazon.com/solutions/case-studies/airbnb/>

- 2010 – 24 EC2 instances, 300 GB of data
- 2015 – 1000 EC2 instances, 50 TBytes data
  
- Grew up entirely on AWS, no data center, no capital purchases, no racking/stacking, no acquisition networking...
  - 5-person operations team
  - Piggyback on AWS for external network, availability zones
  
- Rapid growth easily accommodated.

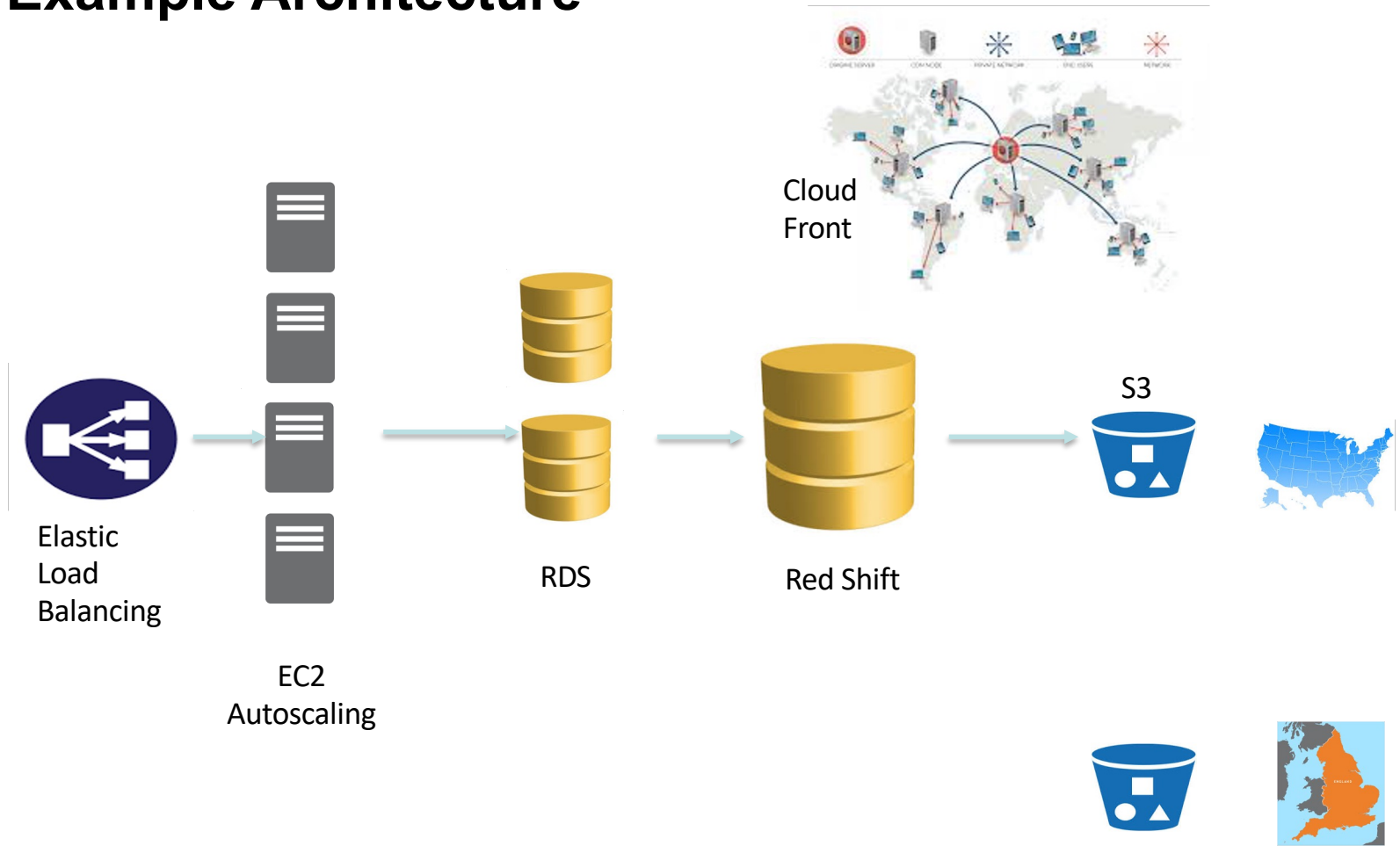
## Coursera

<https://aws.amazon.com/solutions/case-studies/coursera/>

- Massive on-line courses from Stanford, Duke...
- Went from 0 to 3.2 million users in first year
- Accessed from around the world
- Spikes common, e.g., 75% increase in load in 5 minutes



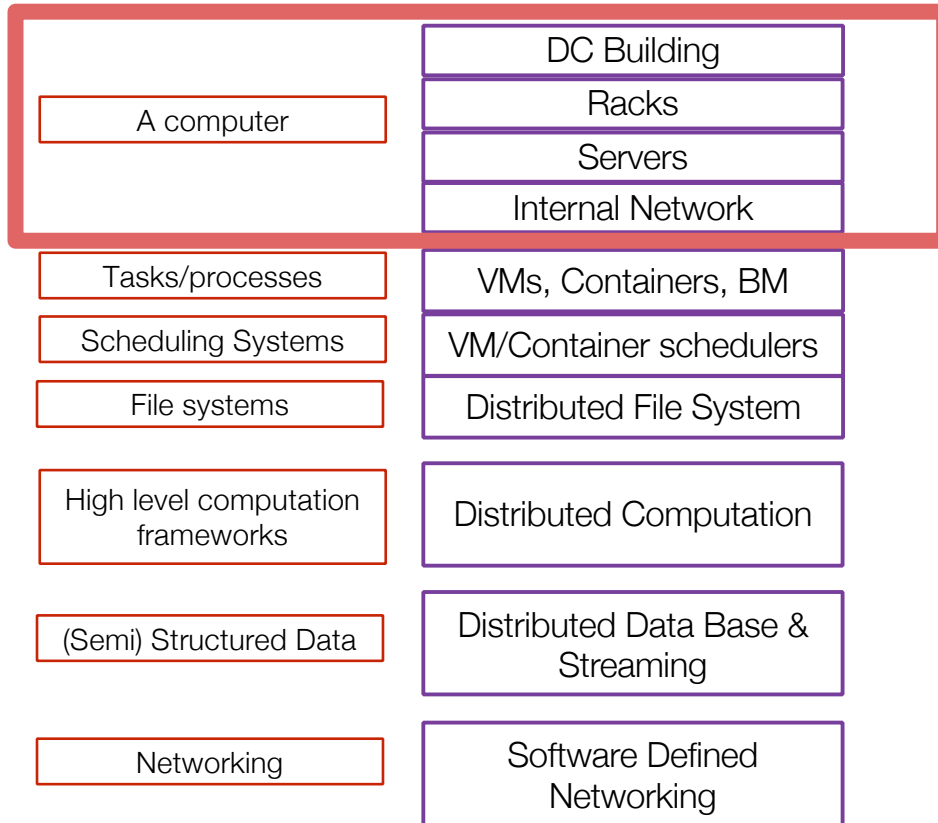
# Example Architecture



## Technology discussed

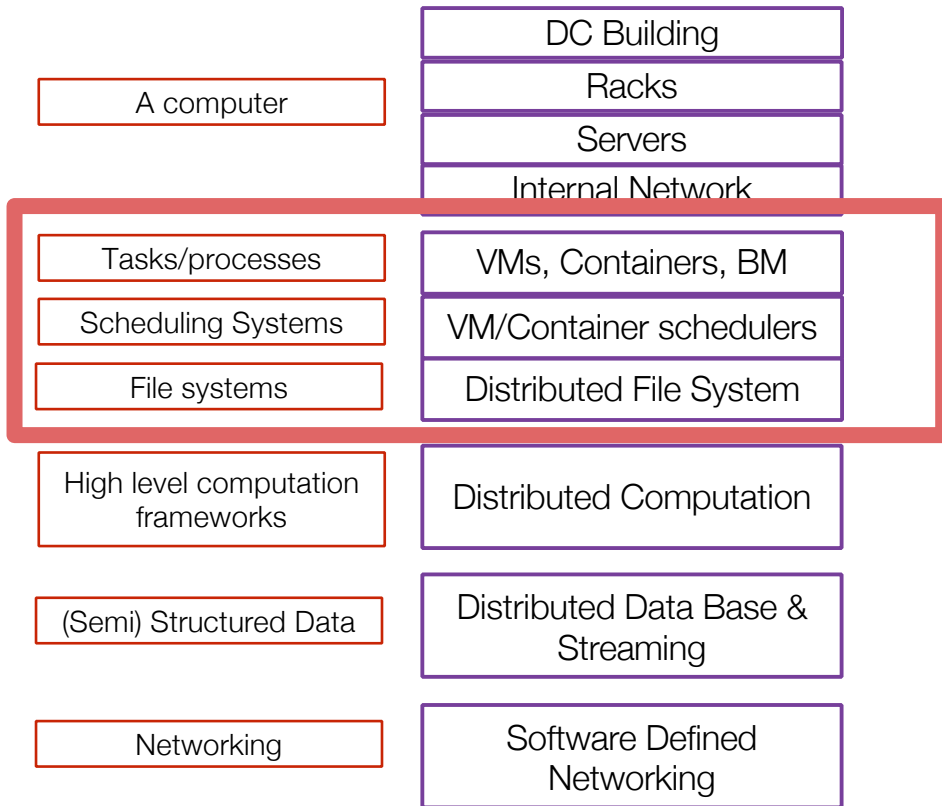
- EC2 & Elastic Load Balancing & EC2 Autoscaling – increase/decrease number of servers as needed.
- Relational Database Service (RDS) – managed service set up DB, patching, read-only replicas, across regions, backups automatically, snapshots
- Cloud Front – CDN, moved from 500 msec to 50msec average latency
- Red Shift – Data warehouse

# Layers of data center



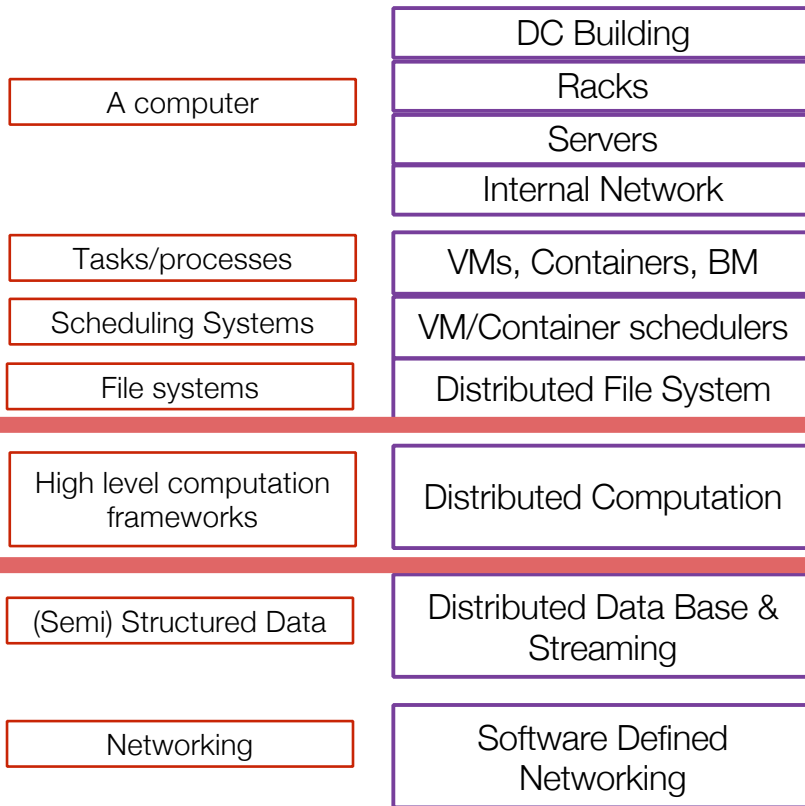
Hardware level:  
How do you build cloud-scale systems?

# Layers of data centers



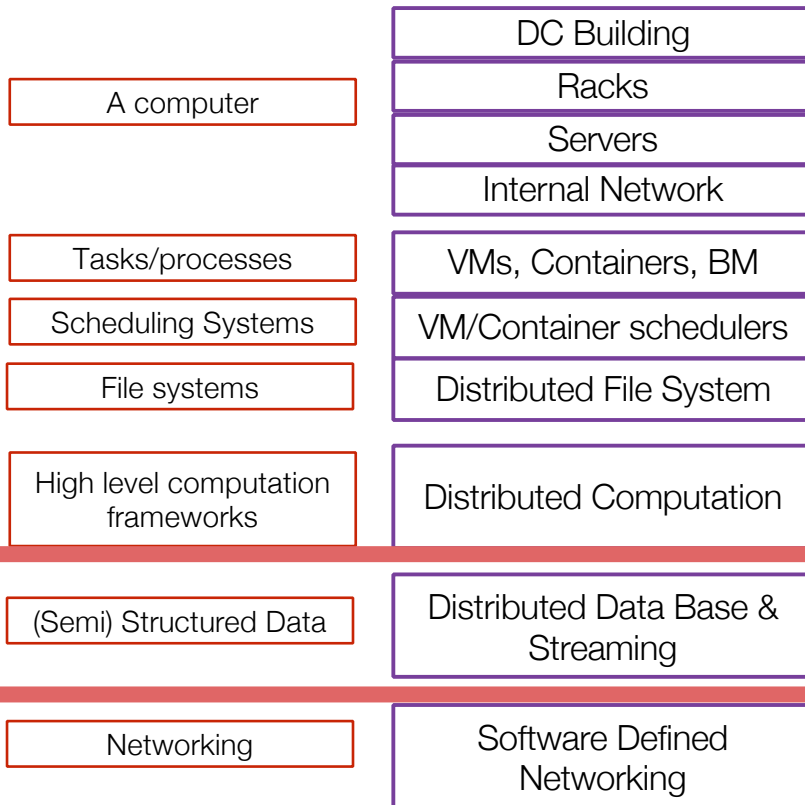
“Operating system”:  
How do you manage and run cloud applications? What about file systems?

# Layers



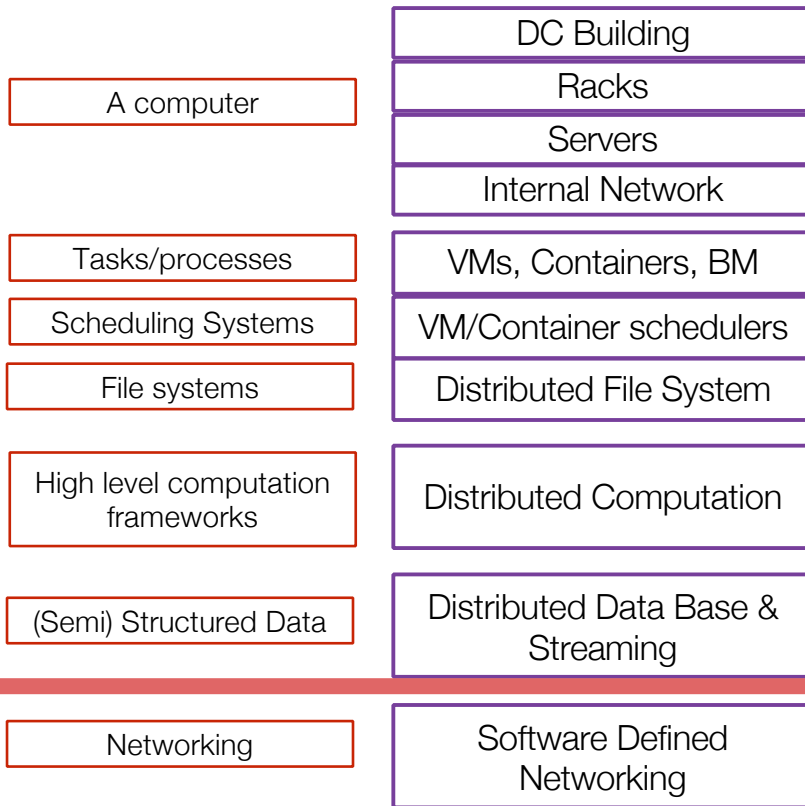
Frameworks:  
How do you write a distributed application?

# Layers



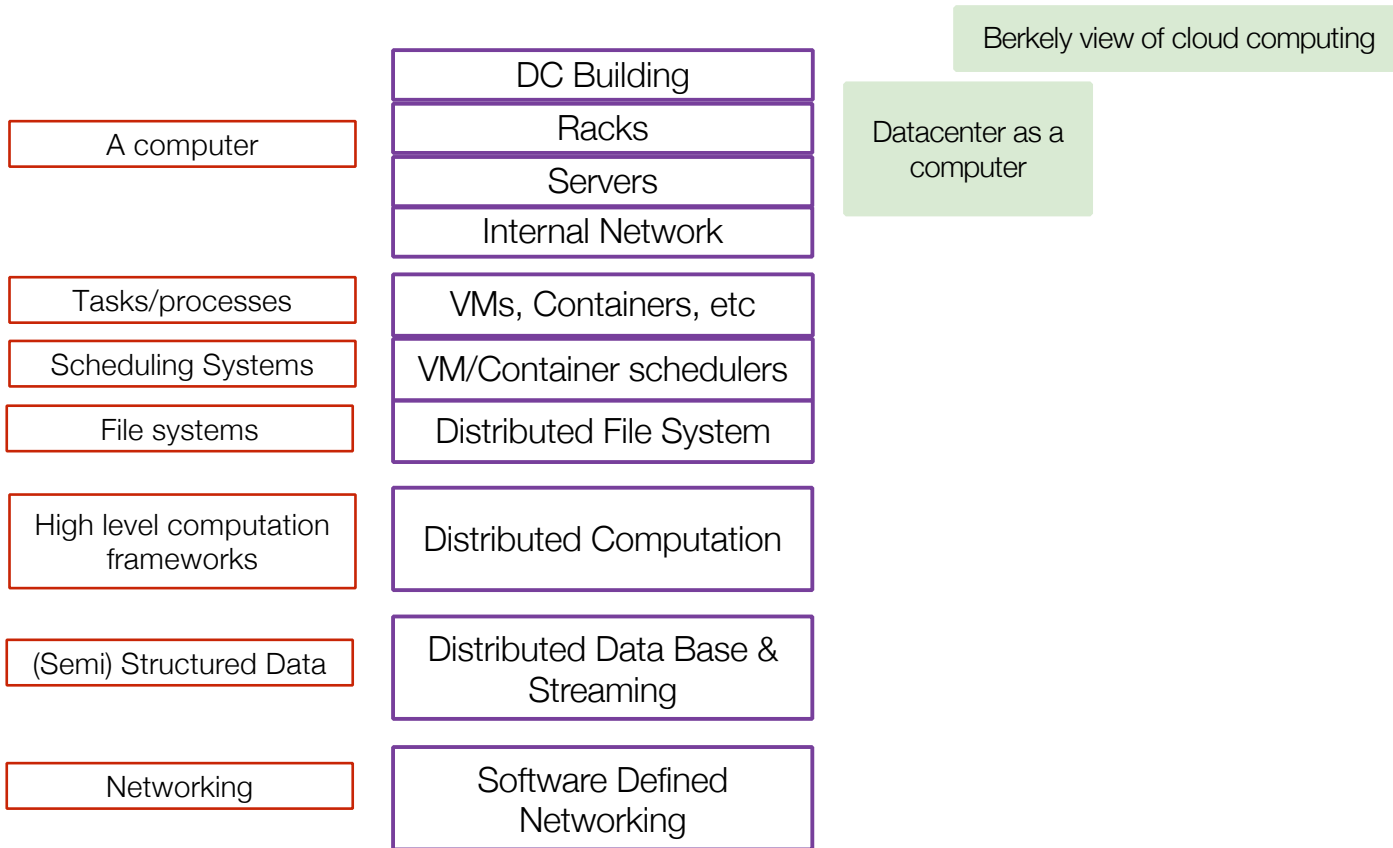
Above the file system:  
How do you manage and work with  
structured data?

# Layers



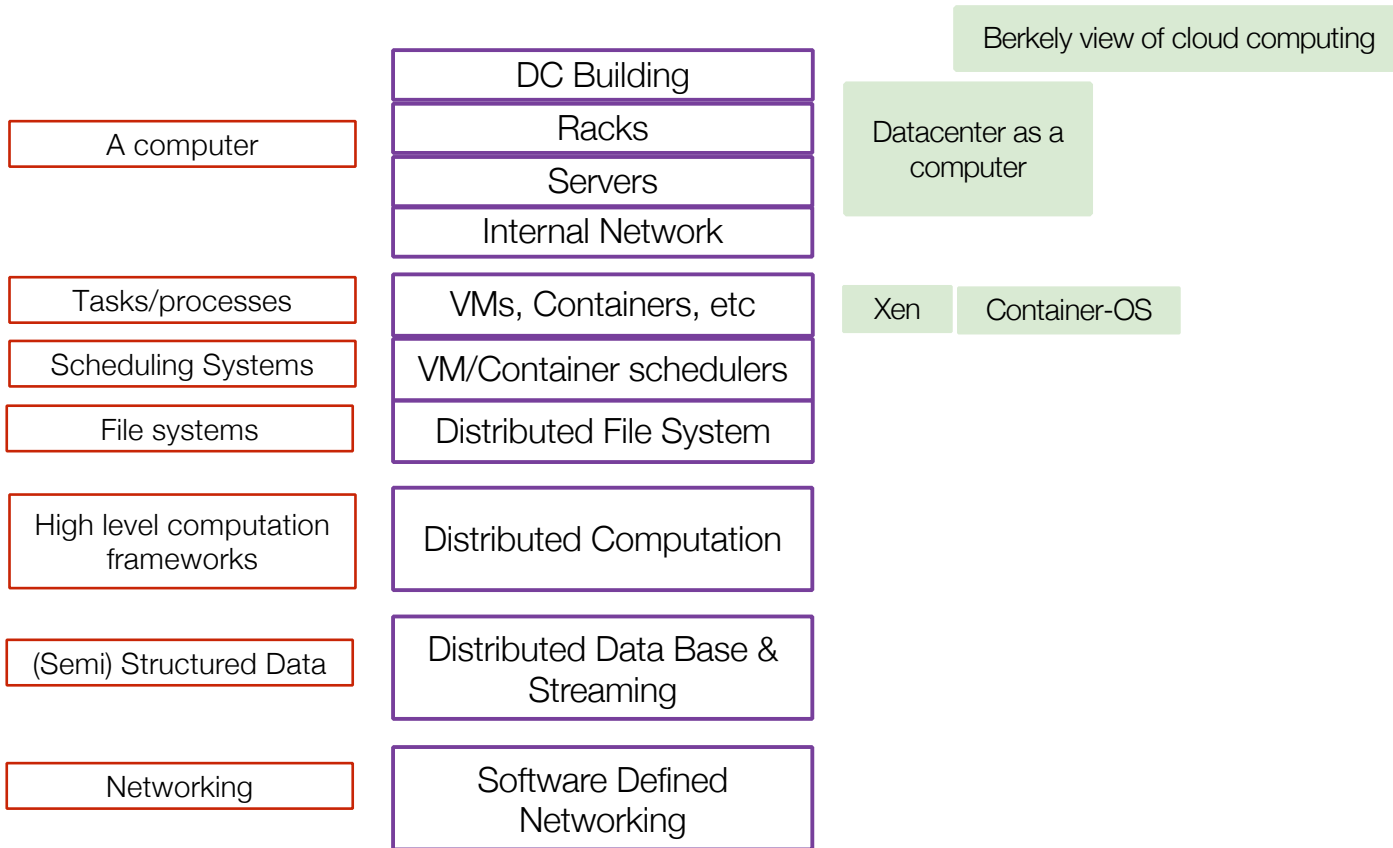
**Networking:**  
How do the parts of a cloud-scale system talk to each other?

# Top-down view of the course

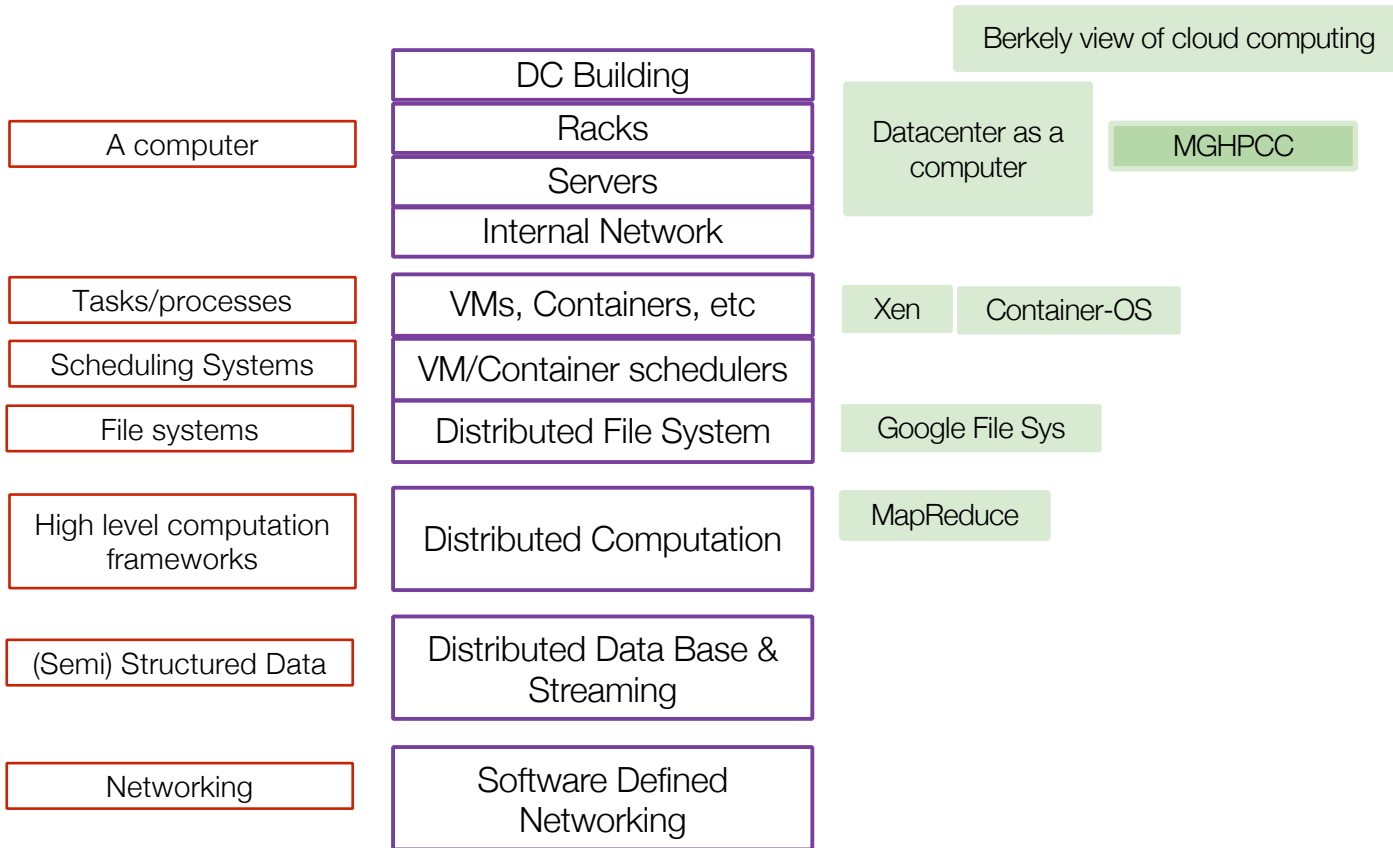




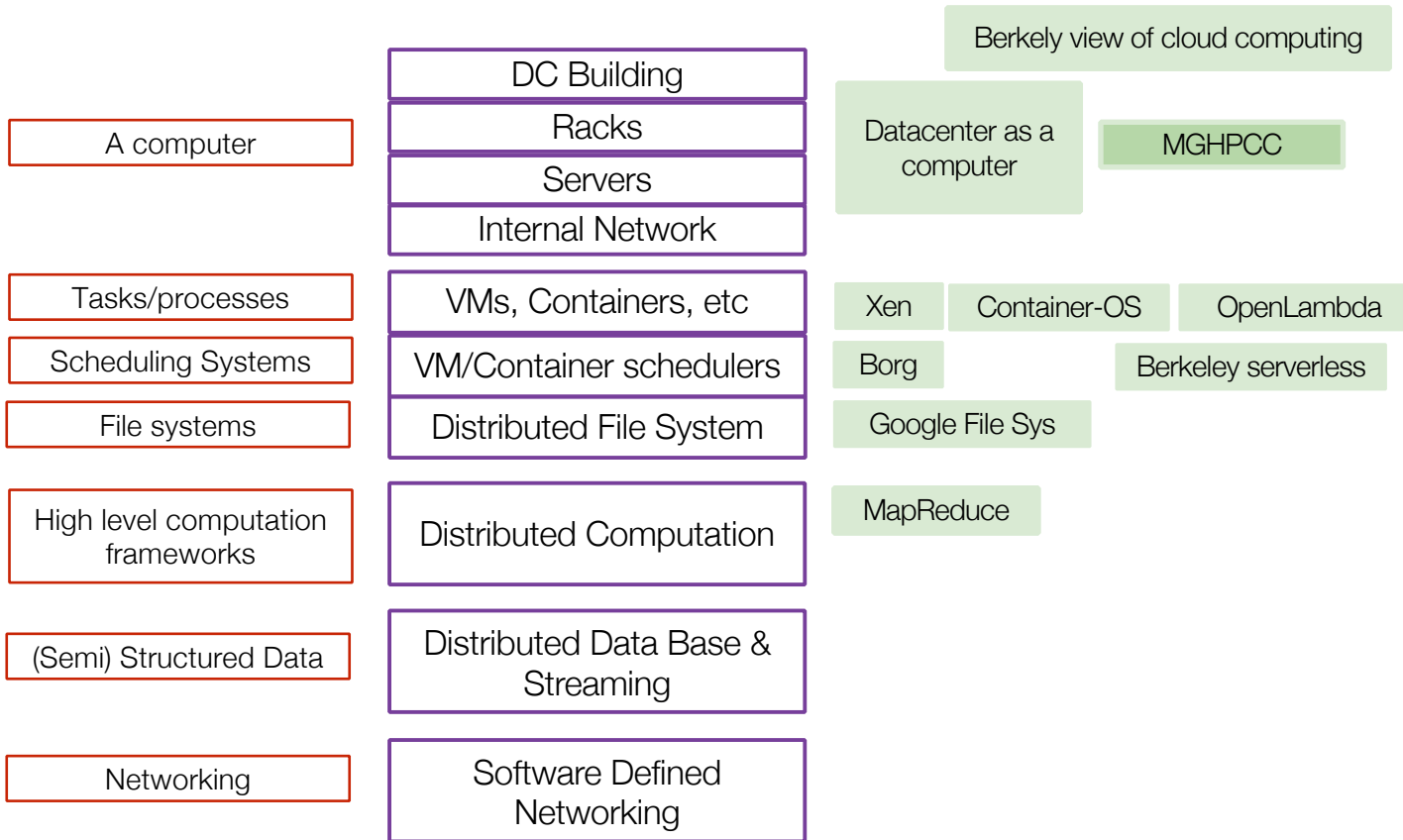
# Top-down view of the course



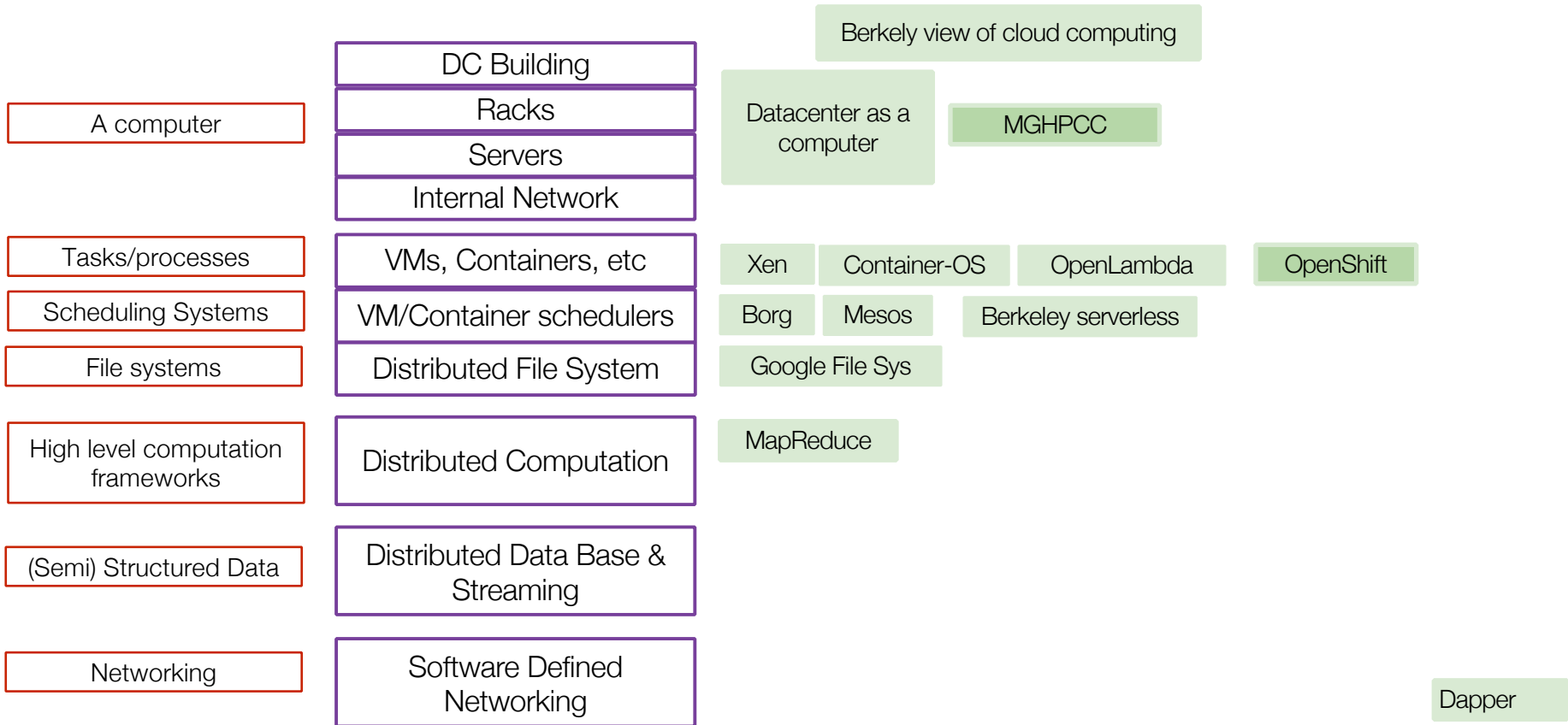
# Top-down view of the course



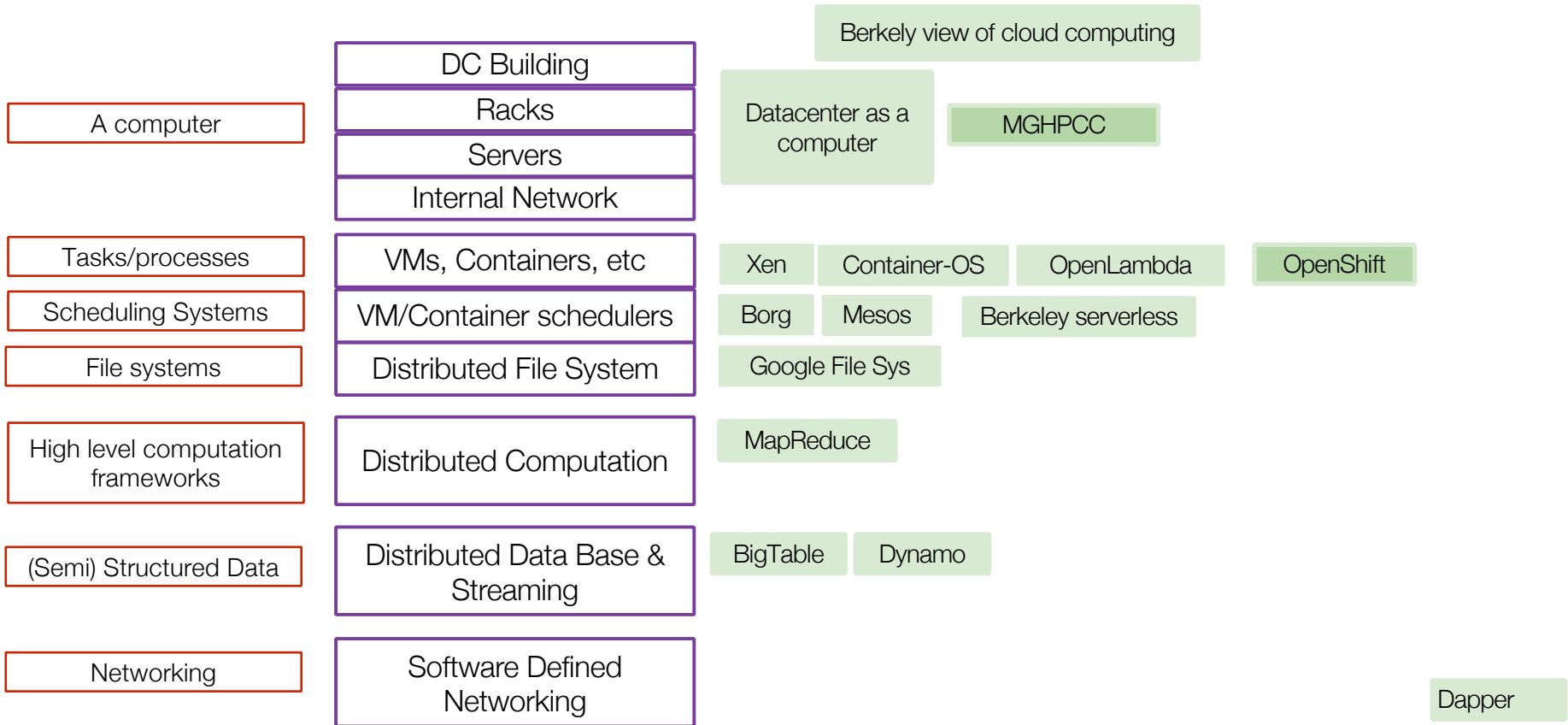
# Top-down view of the course



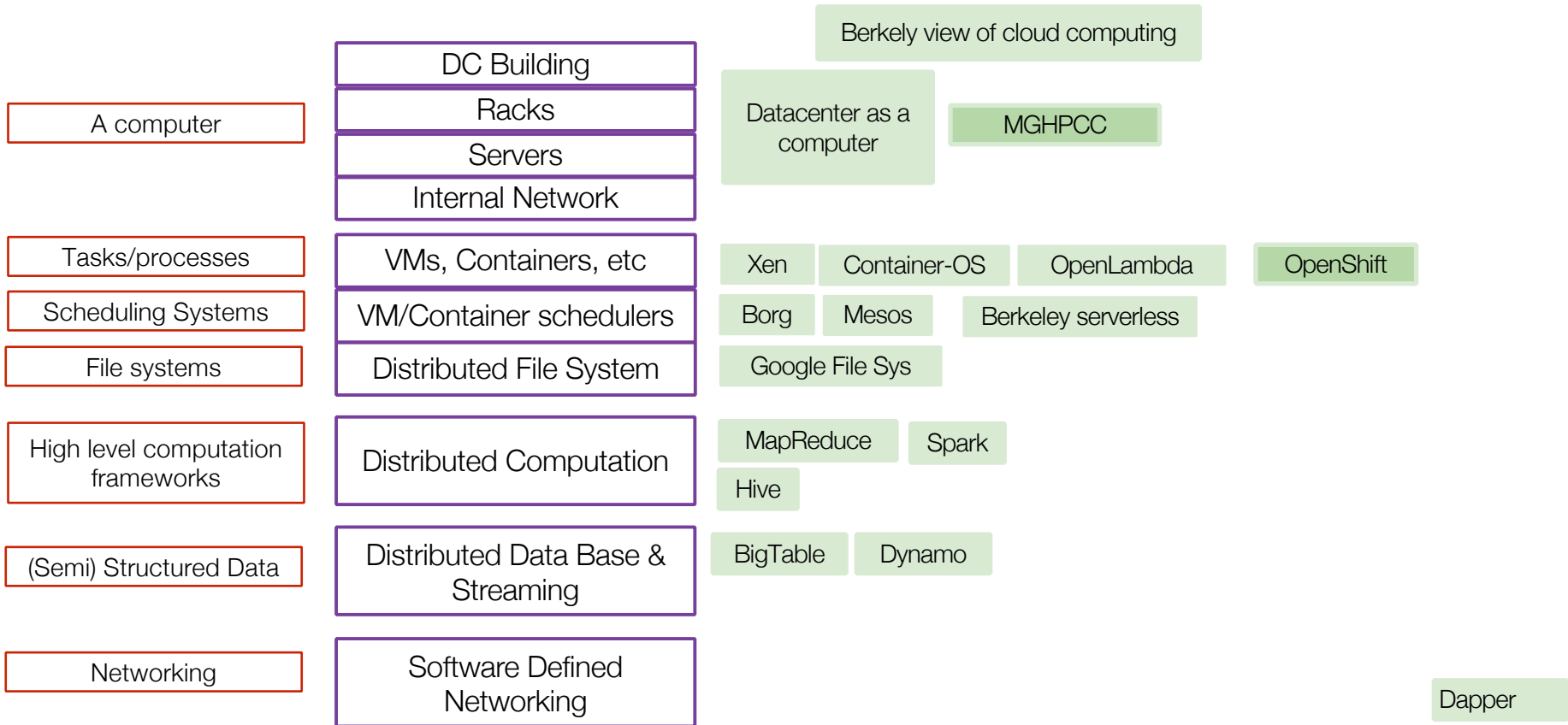
# Top-down view of the course



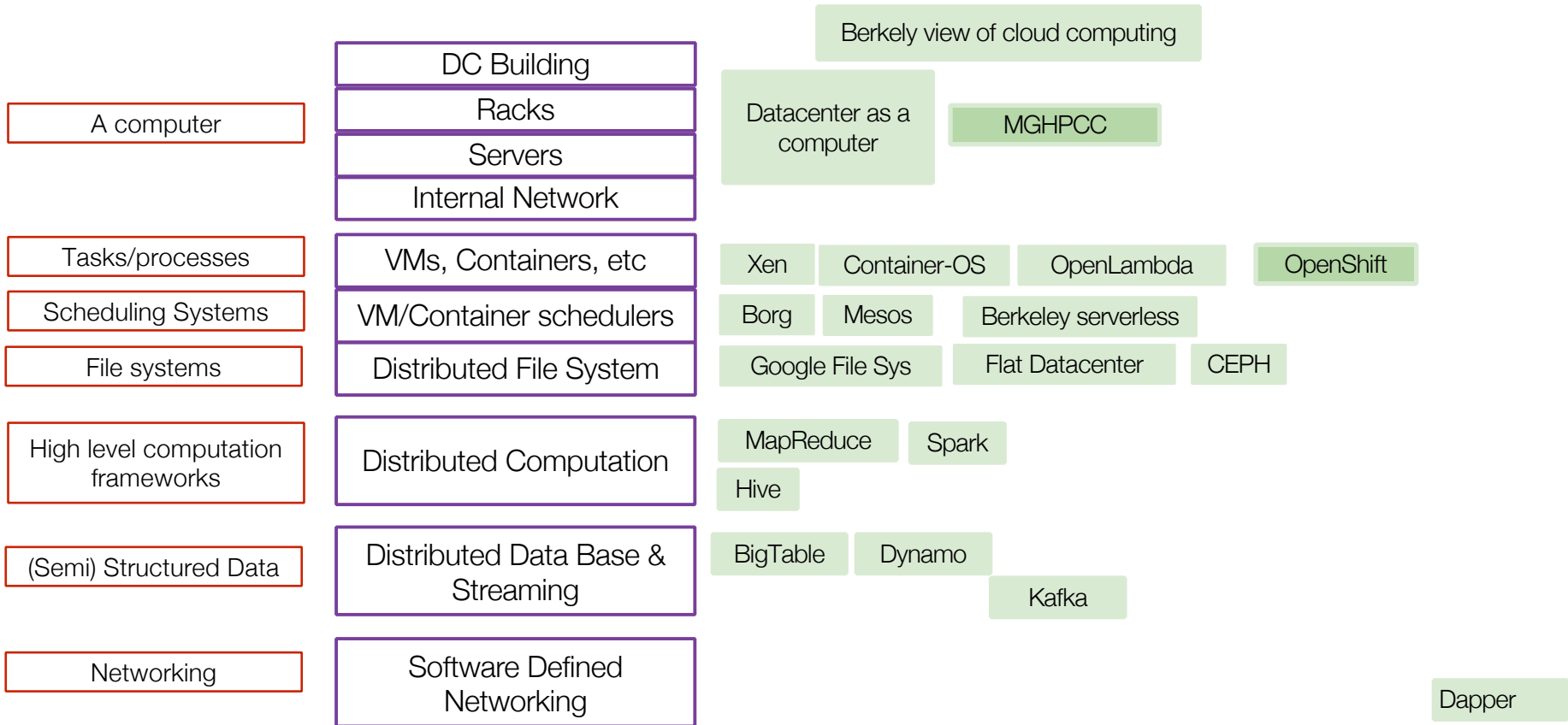
# Top-down view of the course



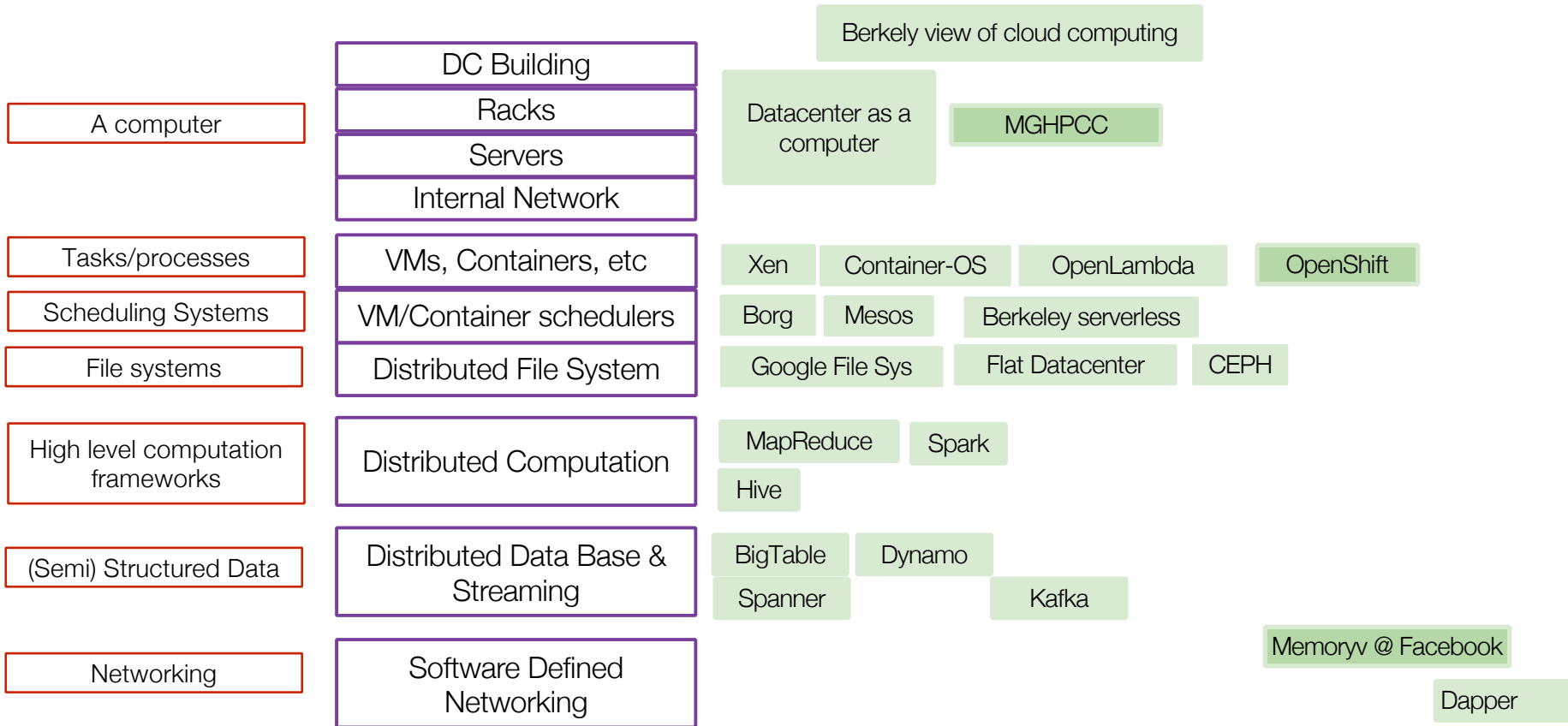
# Top-down view of the course



# Top-down view of the course

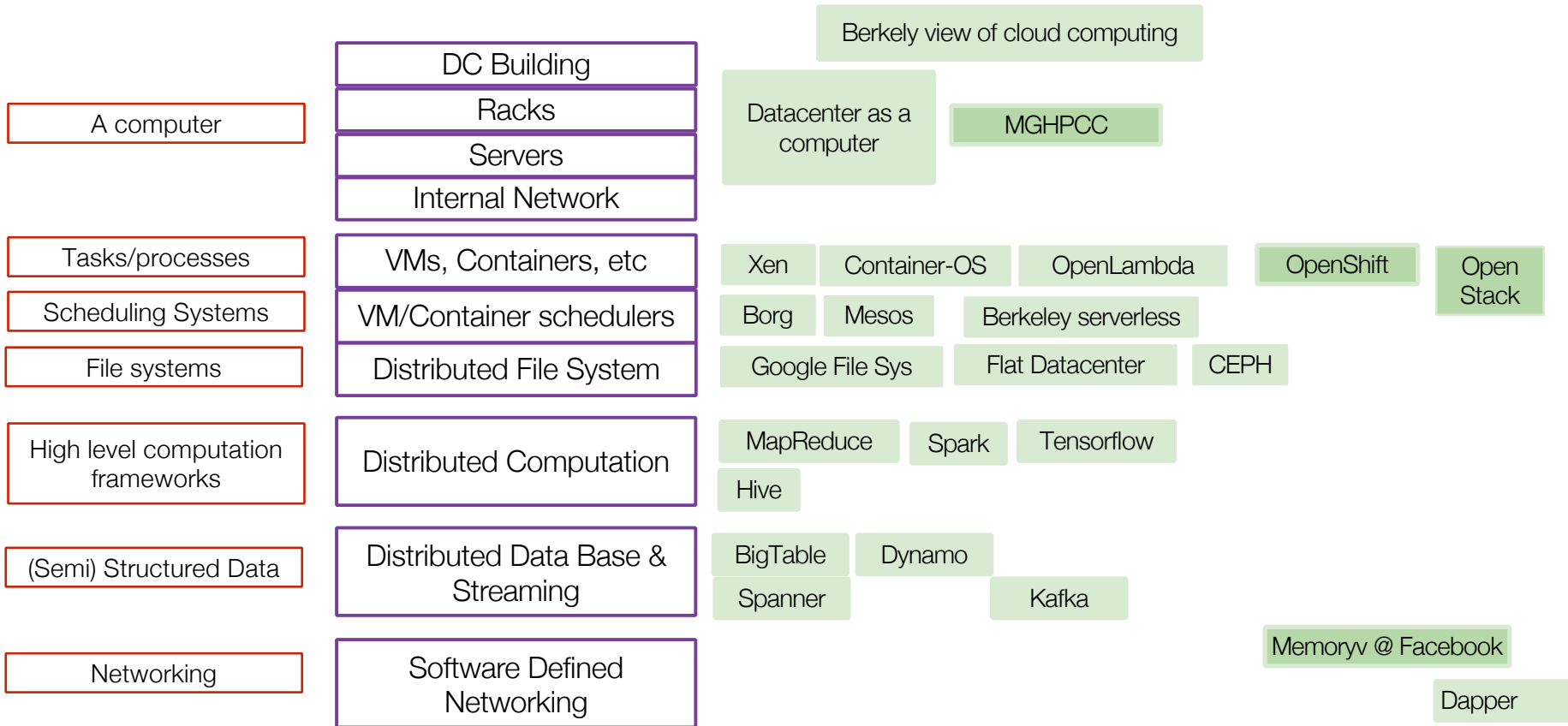


# Top-down view of the course

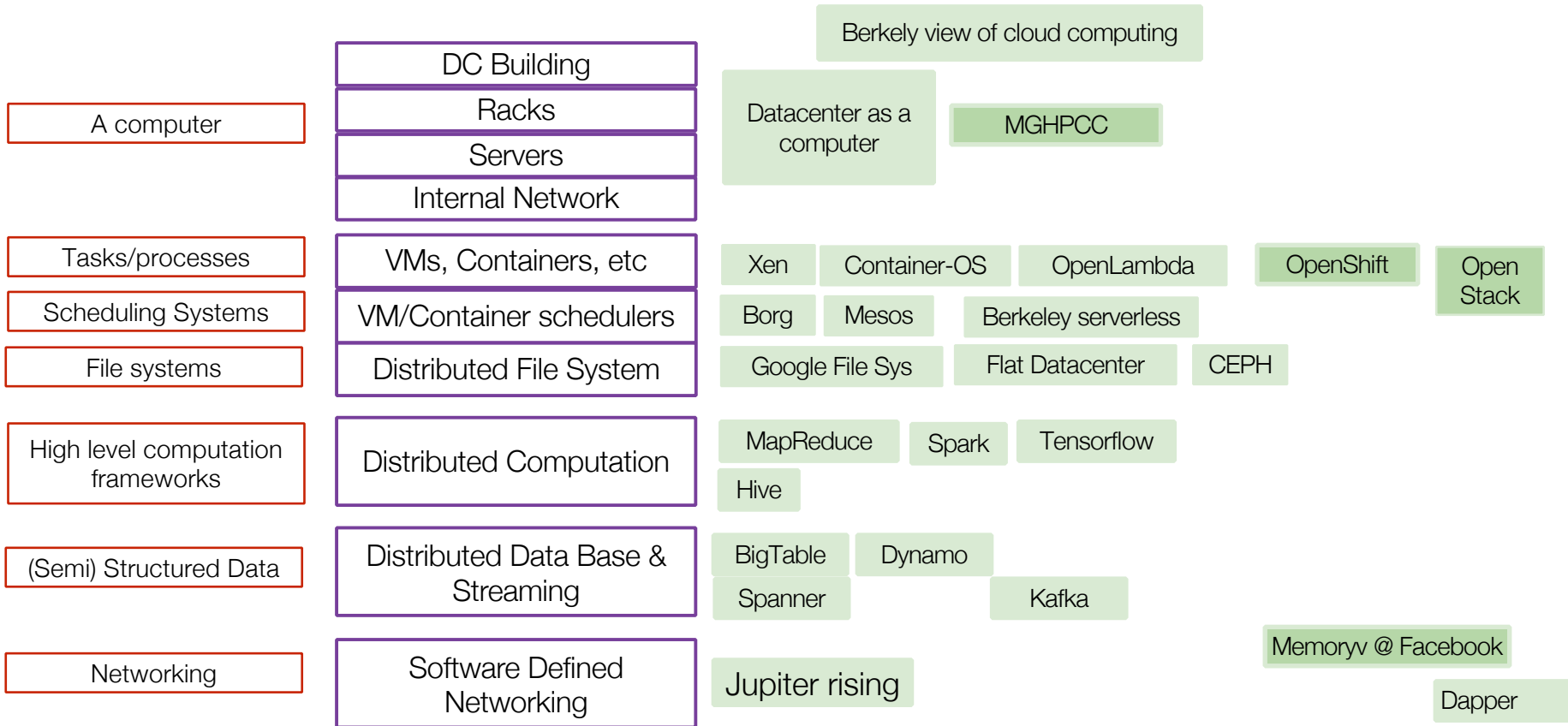




# Top-down view of the course



# Top-down view of the course



# Transformation

- Transformed how SW is developed:
  - continuous deployment; changes tested with real customers
  - example Facebook failure last year
  - massive advantage over waterfall
- It's all about distributed applications
  - change from pets to cattle
  - care about 99th% tail latency
  - stateless servers
  - huge set of higher level services: Containers as a Service, Functions as a Service, Analytics as a Service...

# The challenges

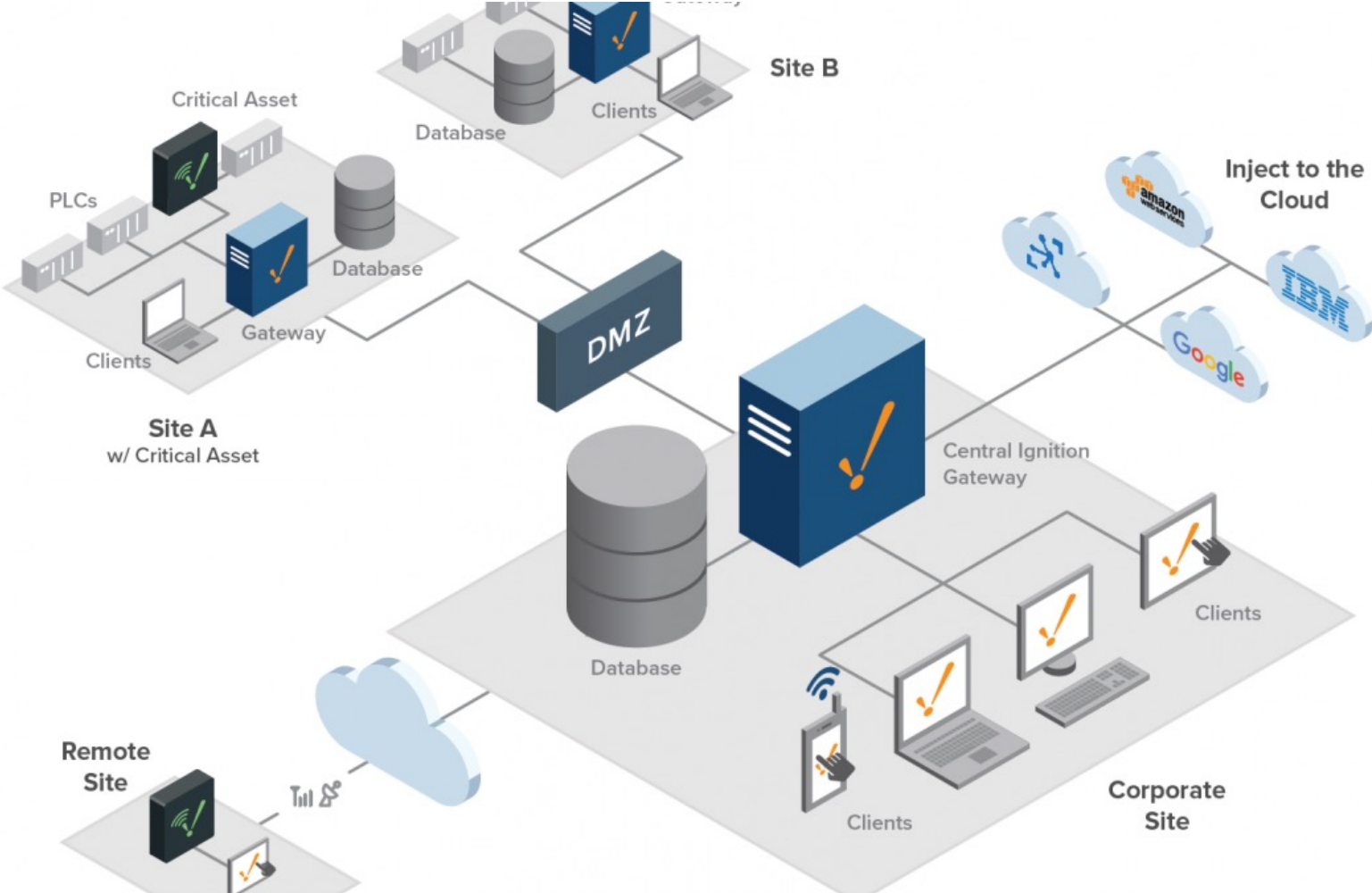
- Monoculture from security perspective
- Emerging oligopoly:
  - Lack of competition limits sources innovation
  - Price is outrageously expensive
- Effort to lock in users: e.g., networking
- Big brother..., or perhaps just Giants whose incentives are not aligned with privacy and marketplace; Consider Facebook

# The Datacenter as a Computer

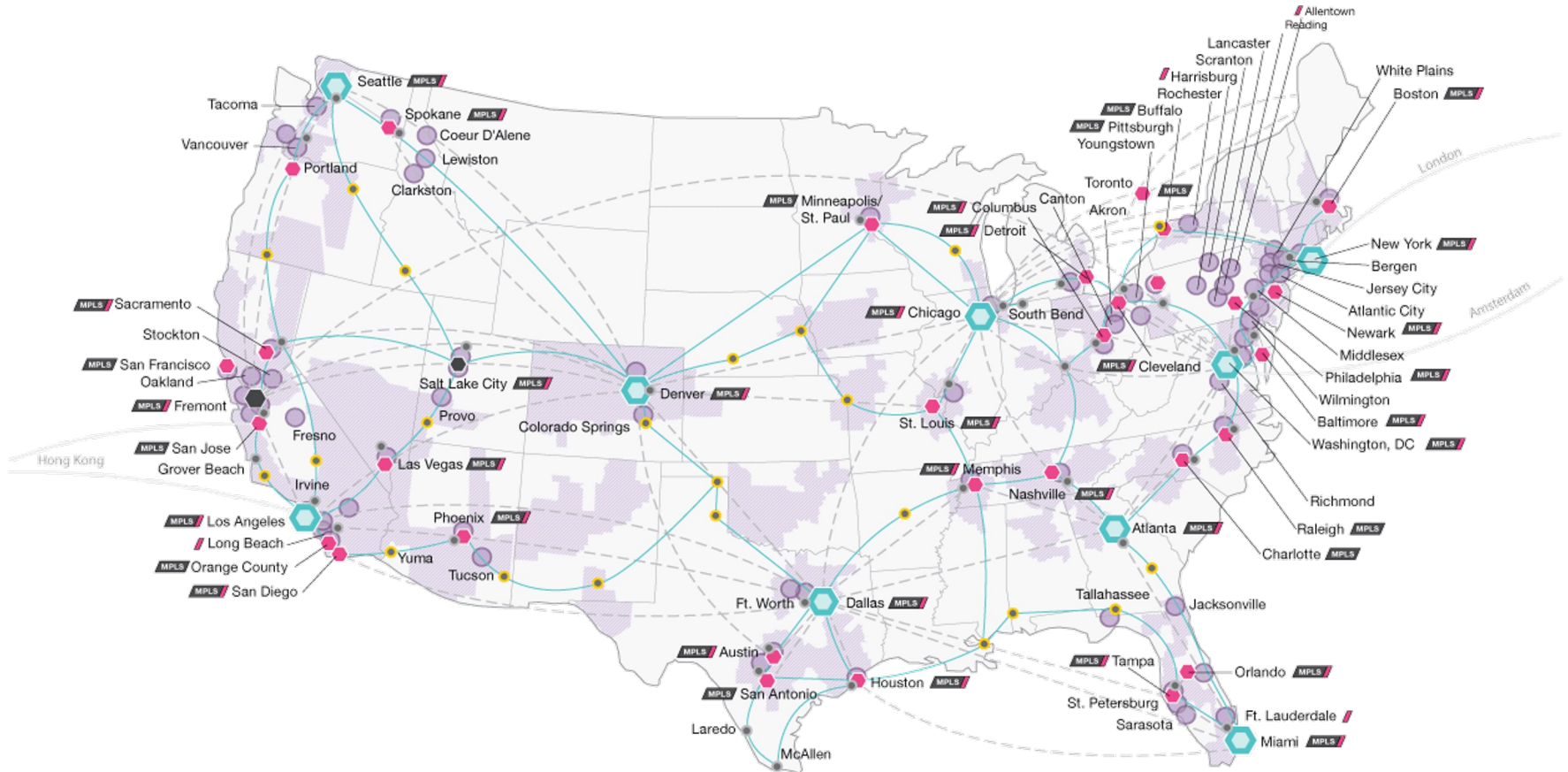
*An Introduction to the Design of Warehouse-Scale Machines – 2<sup>nd</sup> Edition*

Luiz André Barroso, Jimmy Clidaras, Urs Hölzle

# Enterprise IT



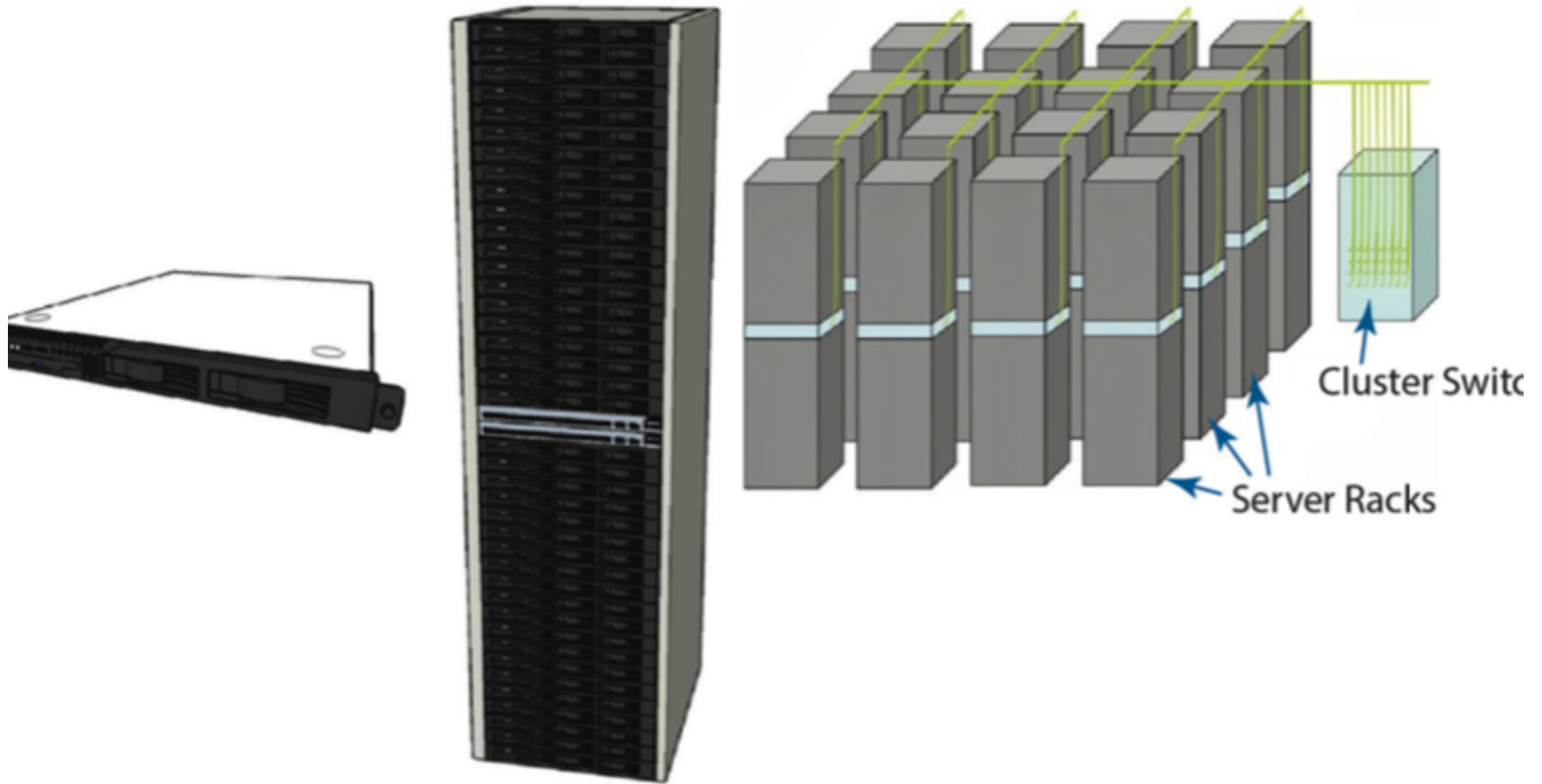
# Traditional “wide-area” networks



<b>LEGEND</b>			
● Core IP Node	/// Media Gateway	MPLS MPLS IP-VPN PoP	▨ Broadband Wireless Spectrum
● Metro IP Node	● Long Haul Termination (All Bandwidths)	--- Nx10 and Nx100 Gigabit Ethernet	○ XO Market
○ Core IP Node w/ Peering	● Long Haul Termination (OC-48 & Above Only)	— 10G/100G Inter-City Long Haul Network	



# Elements of data center

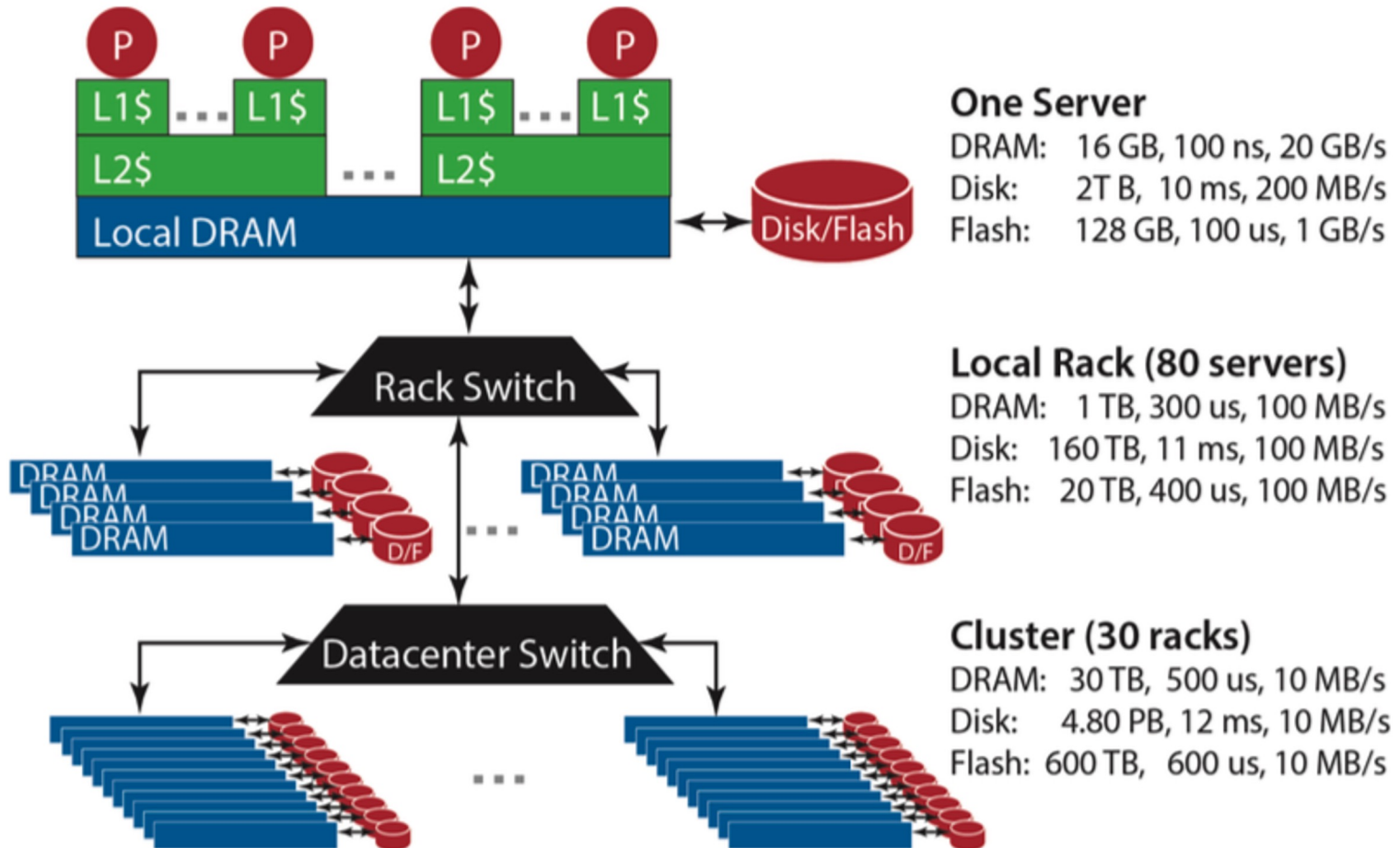




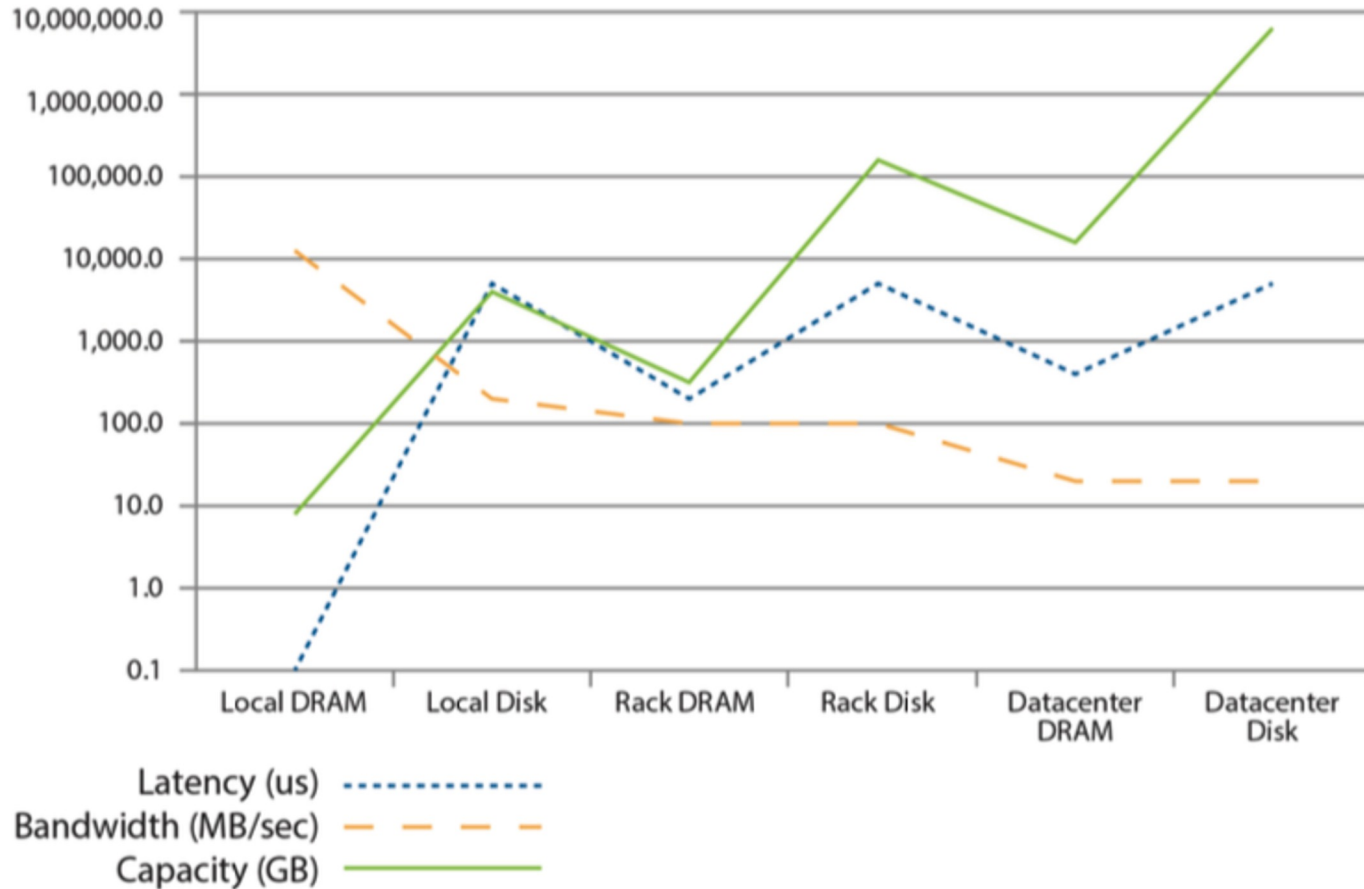
## Storage assumptions

- Storage distributed across all machines
- Software like GFS distributes, versus NAS appliance
  - Redundancy even if rack level failure
  - Multiplex server resources (NIC/enclosure/power)
  - Exploits cheap desktop disks
- Typically network oversubscribed
  - E.g., 32 \* 40Gig links nodes, 4 \*100Gig Links up

# Storage Hierarchy



# Latency, bandwidth & capacity



# Inside a data center



# ***Self-introduction***

**Q&A**