

# CS250B : Modern Computer Systems

## Datacenter Architecture Introduction



Sang-Woo Jun

# Computers at Scale

- ❑ Lots of machines collocated for economy of scale
  - Lower cost of power supply, cooling, management, etc...
  - Cheaper to operate a cluster of 1,000 machines, than to operate 100 clusters of 10 machines!
- ❑ Common terms include clusters, Warehouse-Scale Computers (WSCs)
  - Typically, clusters refer to HPC (High-Performance Computing) clusters for scientific applications
    - Fast processors, tightly coupled via fast network, programmed with MPI, etc
  - Typically, WSCs refer to clusters operated by industry for service
    - Cost-effective components, emphasis on request-level parallelism, flexible resource allocation
- ❑ Physical structure housing these clusters called datacenter

# Some Topics For Datacenter Architecture

- Cost of a datacenter
- Power Supply
- Cooling
- Networking
- Virtualization
- ...

# Scale of Datacenters

## ❑ Google Mayes County Datacenter

- Mayes County, Oklahoma
- Former Gatorade plant (Pepsico)
- 1,400,000 Square feet
  - 24 football fields
  - Largest industrial building in the state
- Total power: 100 Megawatts
  - Total residential power consumption of San Francisco: **168 MW**
  - 1 cent per watt difference is drastic



# Total Cost of Ownership (TCO)

- ❑ Total direct and indirect cost of owning and operating a system
  - Purchasing cost, operating cost, and more
- ❑ Capital Expenditure (CAPEX)
  - Cost to build the datacenter
    - Including cost to buy replacement parts, etc
  - Depreciation (amortized) cost of the datacenter
- ❑ Operational Expenditure (OPEX)
  - Cost to operate the datacenter
  - Power, cooling, human resources, etc

# TCO of a Modern Datacenter

“Power Utilization Effectiveness”

PUE OVERHEAD

1.2%

DC OPEX

17.2%

DC INTEREST

16.7%

DC ARMORTIZATION

29.9%

SERVER POWER

11.6%

SERVER OPEX

1.0%

SERVER INTEREST

2.6%

SERVER AMORTIZATION

19.9%

What we can affect, as engineers

# Factors On Choosing a Location

- ❑ Near customer
- ❑ Land cost, property tax rates
- ❑ Electricity cost
  - Near power plant?
  - Many companies invest in wind farms, etc
- ❑ Proximity to internet backbones
- ❑ Earthquakes, floods, hurricanes...
  - Distributed across the globe just in case

# Example Datacenter Deployment Sites

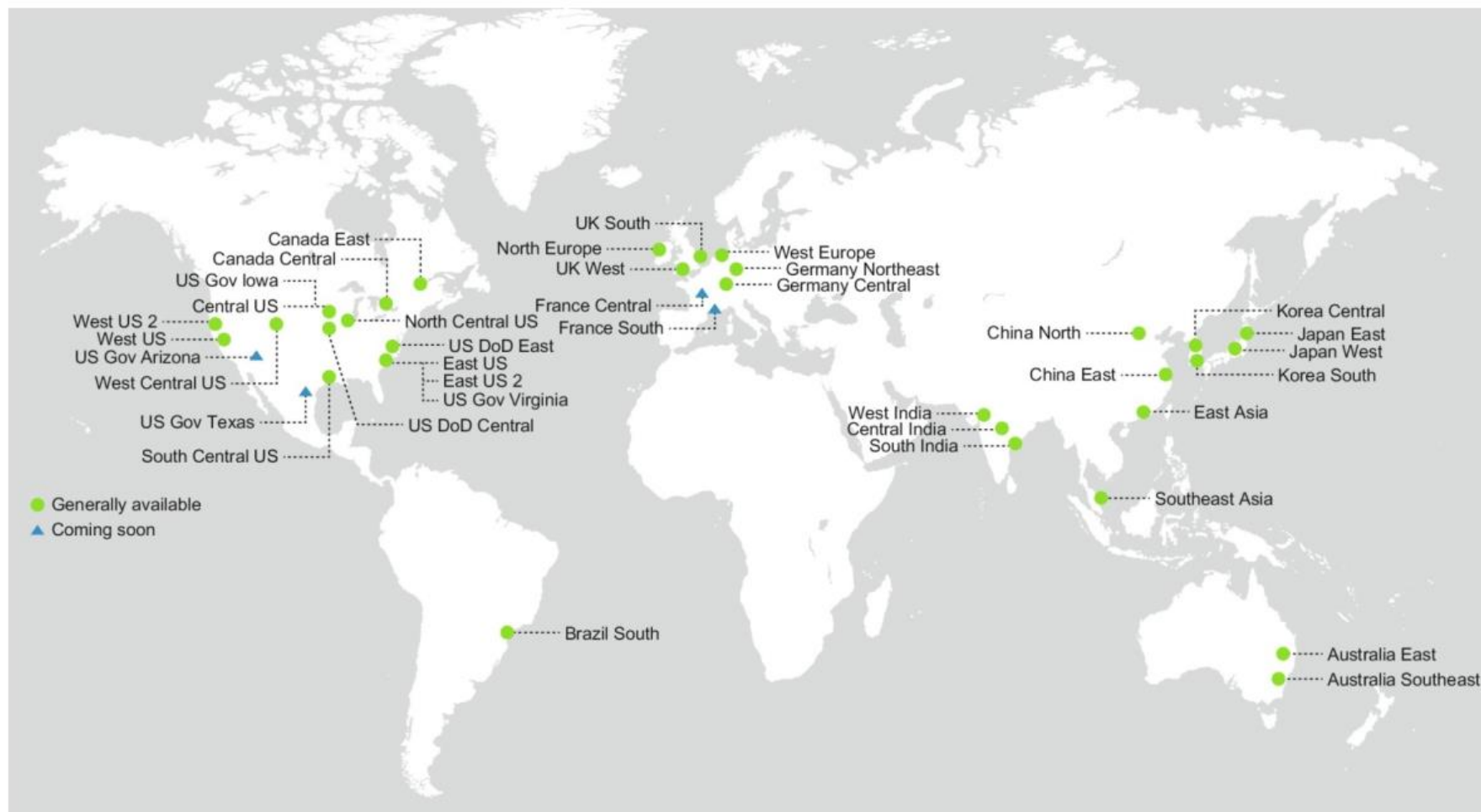
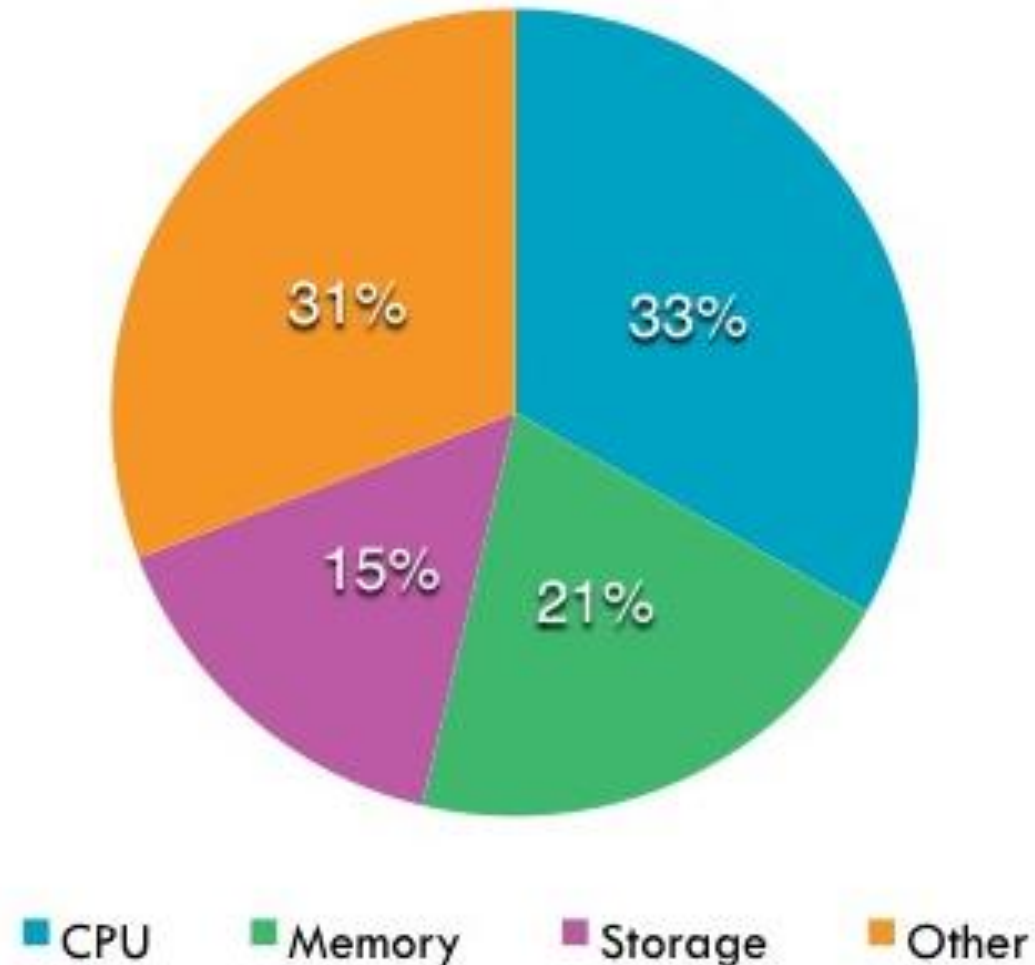


Figure 6.20 In 2017 Microsoft had 34 sites, with four more opening soon. <https://azure.microsoft.com/en-us/regions/>.

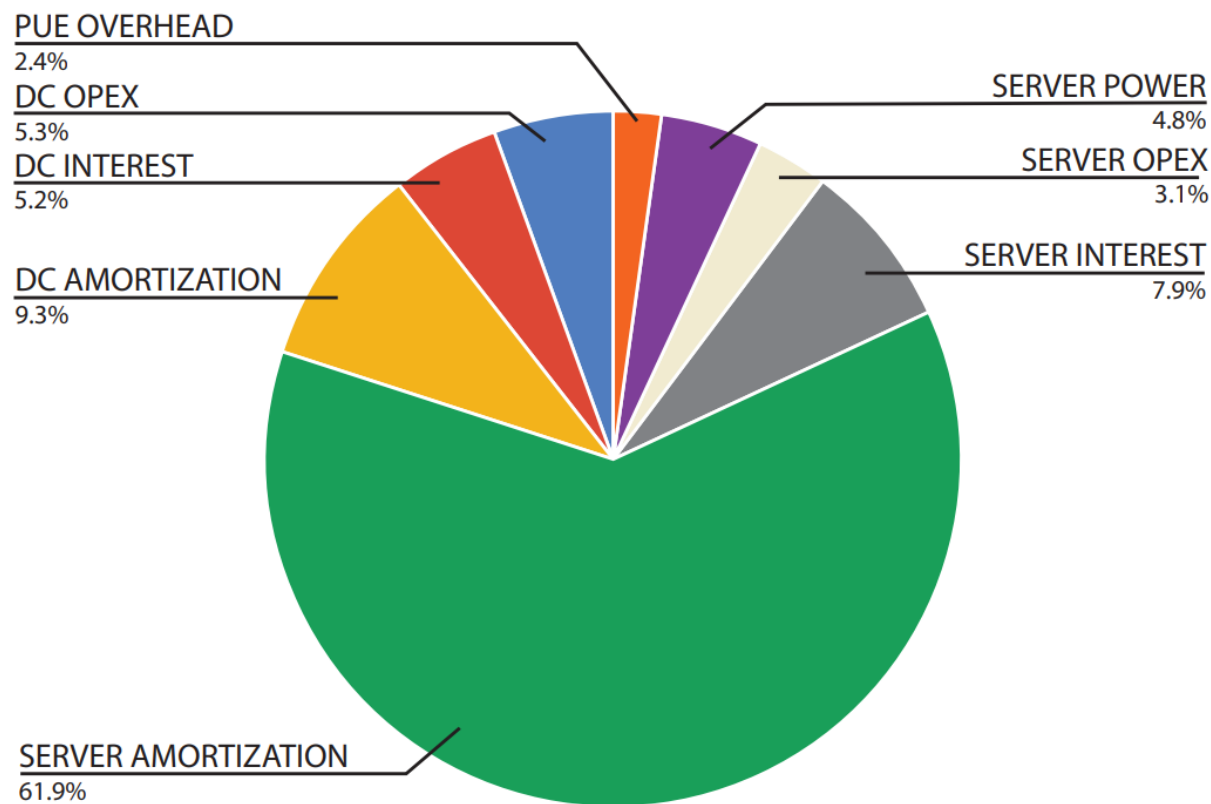


# Cost of a Server Breakdown

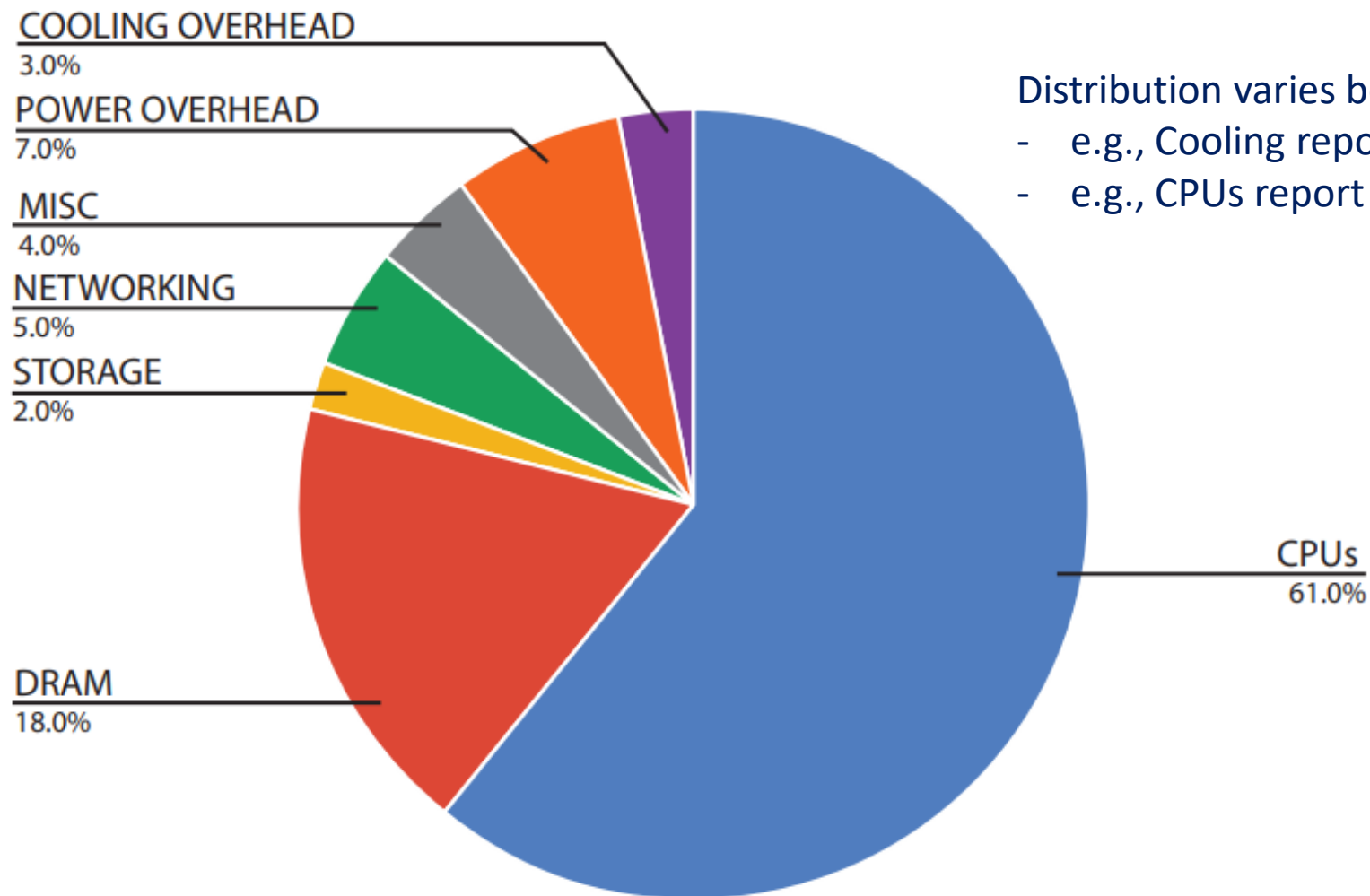


# Aside: TCO of a Classic Datacenter

- ❑ Compared to this, modern use of commodity servers puts less emphasis on server cost



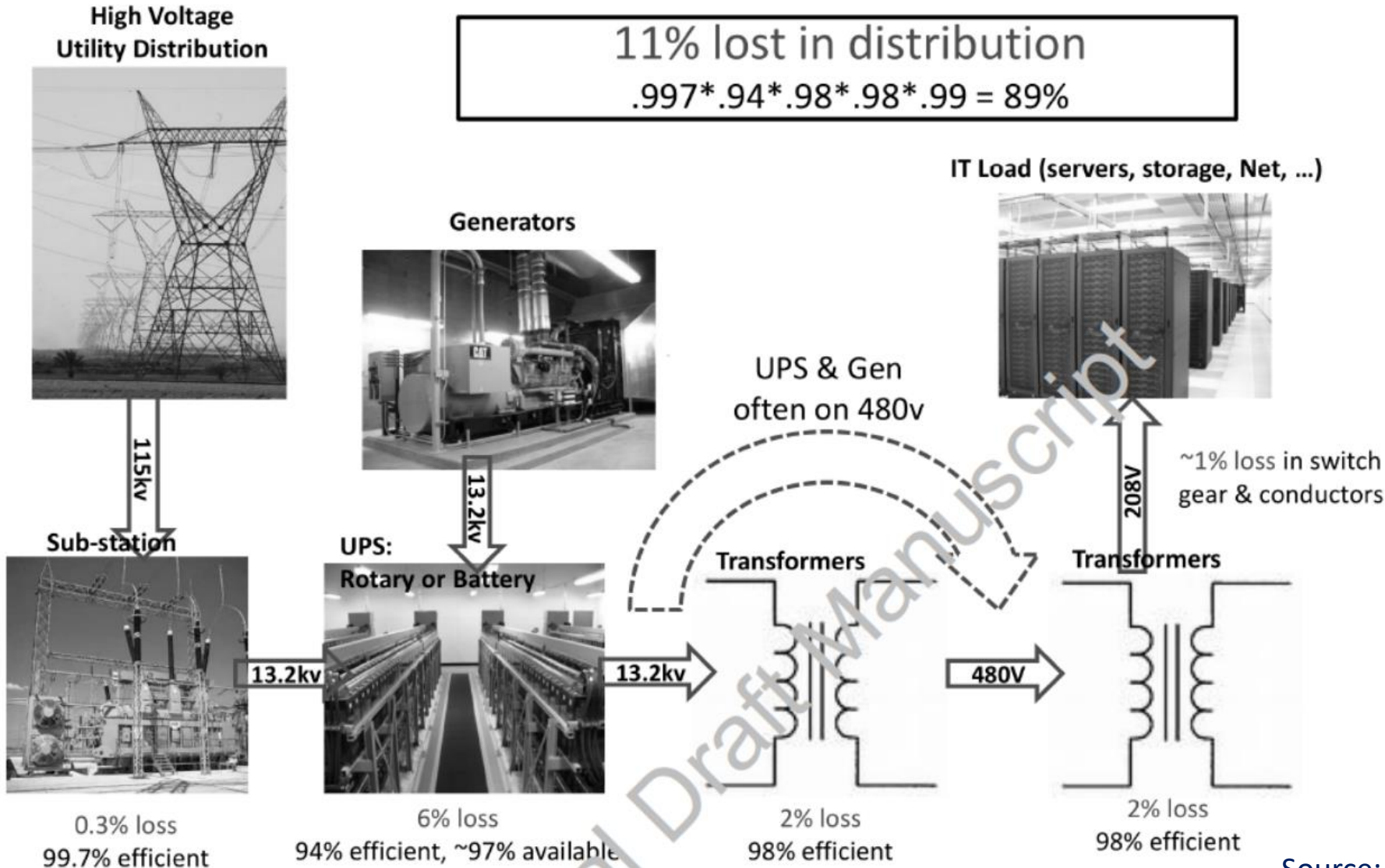
# Power Consumption Breakdown



Distribution varies between instances and workloads

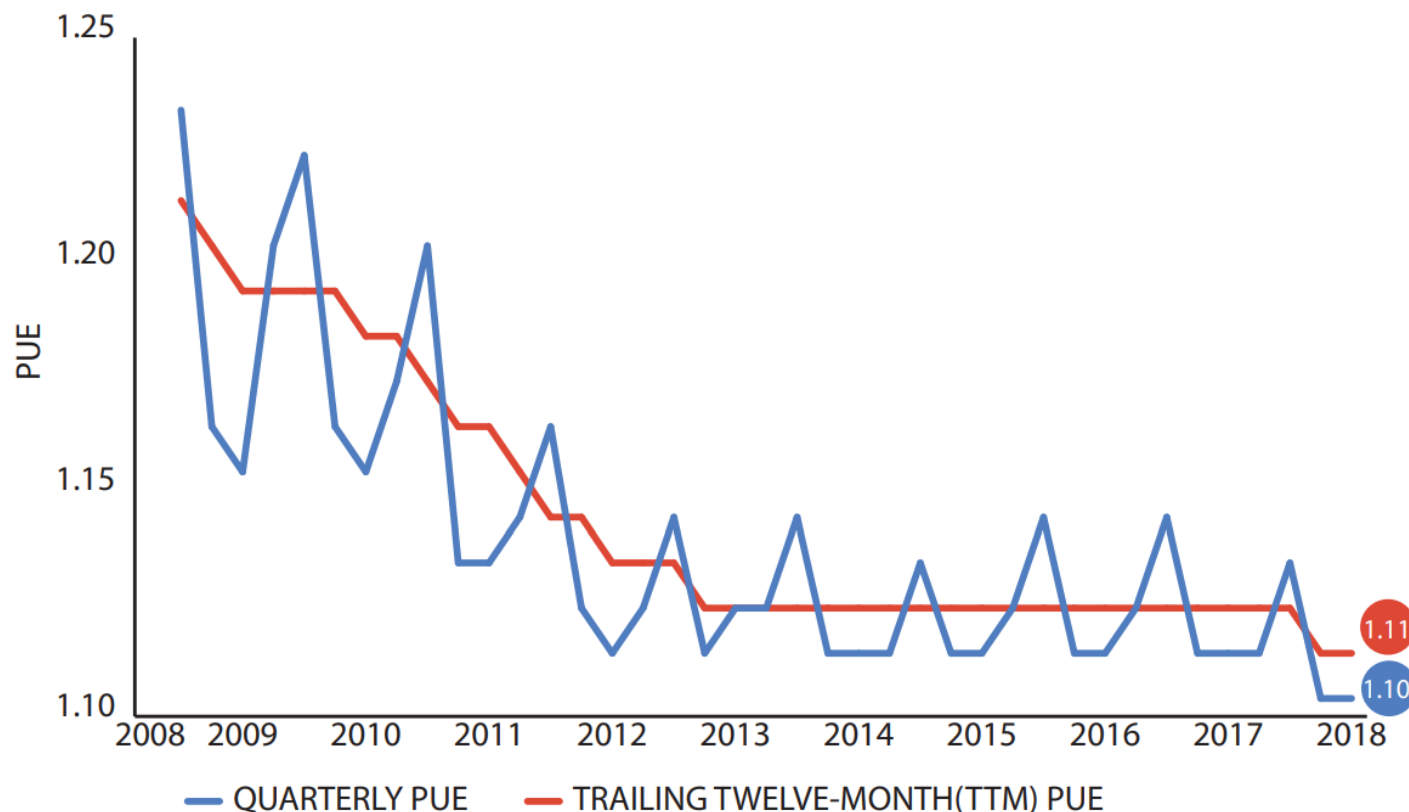
- e.g., Cooling reported 3% - 30% overhead
- e.g., CPUs report 45% - 61% overhead

# Power Distribution in Datacenters



# Power Utilization Effectiveness (PUE)

- Total facility power / IT equipment power
  - Power loss within the facility before reaching servers



# Cooling Infrastructure of a Datacenter

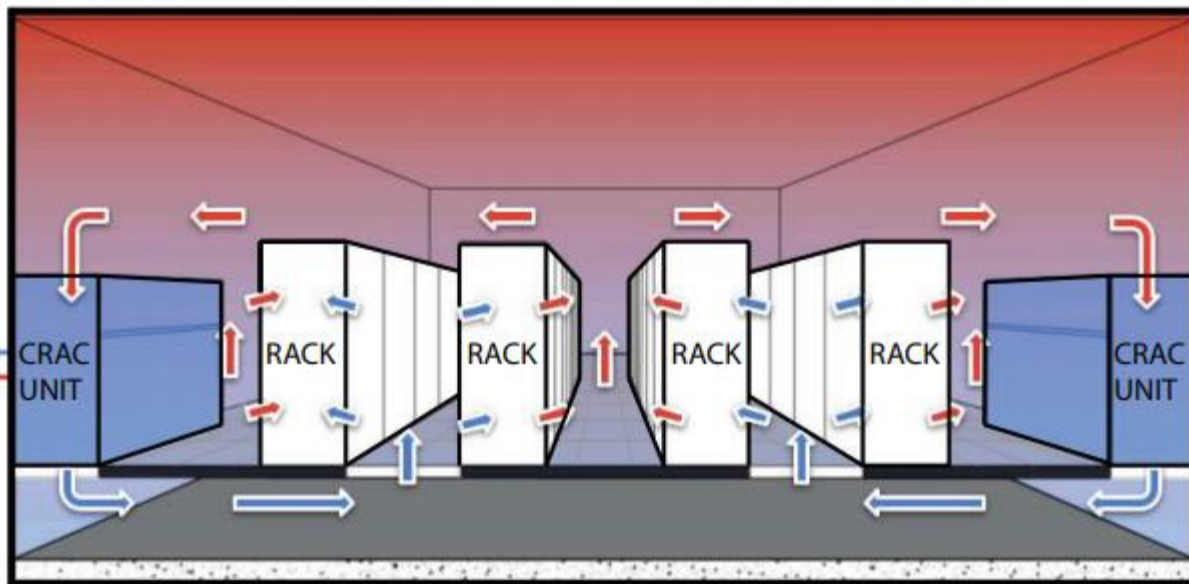
- ❑ Most straightforward: Air-Conditioning based systems
  - Optimal airflow is an intensely studied topic
  - Millions of gallons of water used daily

CEILING

LIQUID SUPPLY

FLOOR TILES

FLOOR SLAB

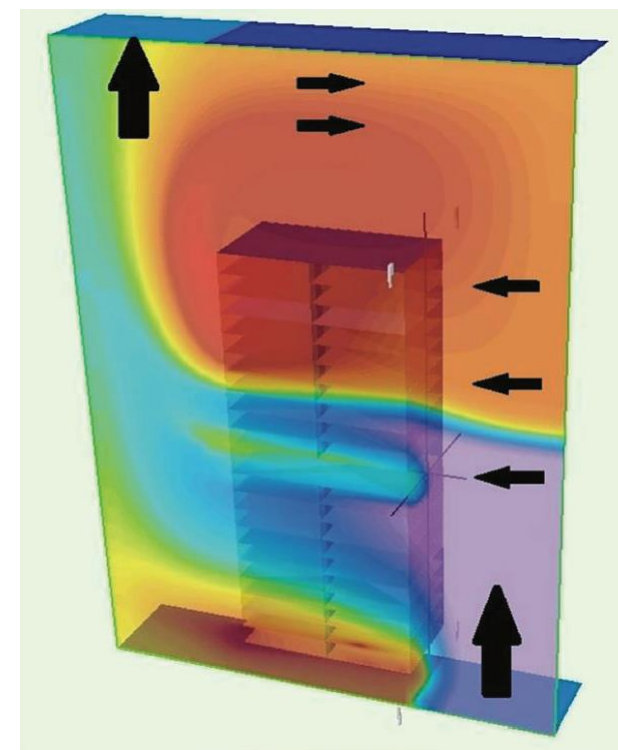


CEILING

LIQUID SUPPLY

FLOOR TILES

FLOOR SLAB



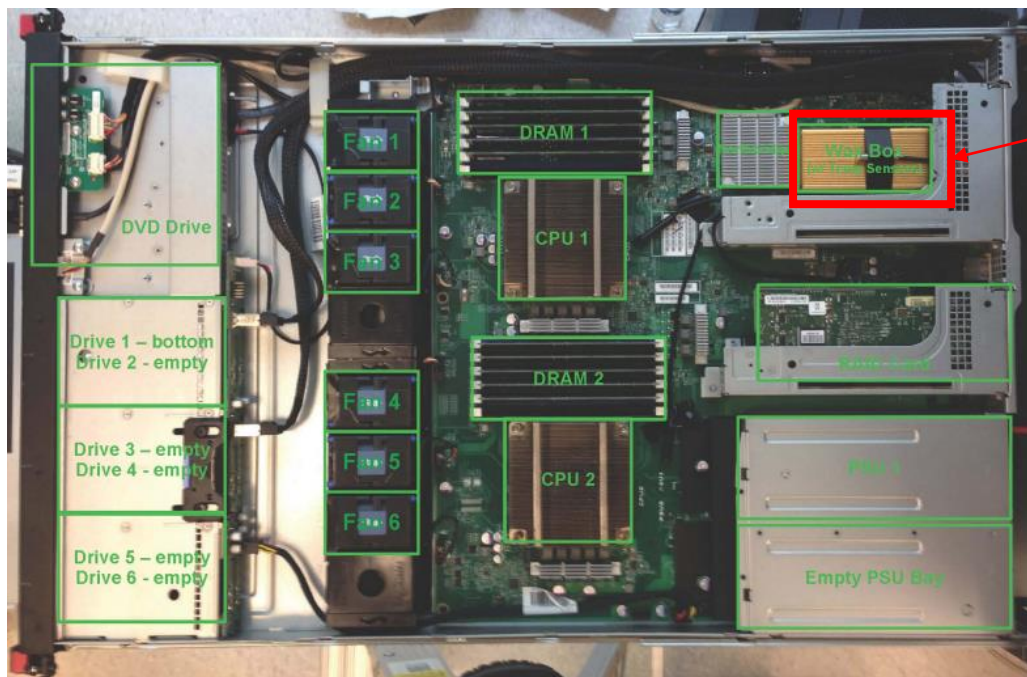
# Unconventional Cooling Techniques

- ❑ Ambient temperature is an important factor in datacenter construction
  - Constructed down mountain valley to take advantage of cold wind
  - Flowing water in walls to effectively remove heat
  - Reportedly no AC requirements except in Summer

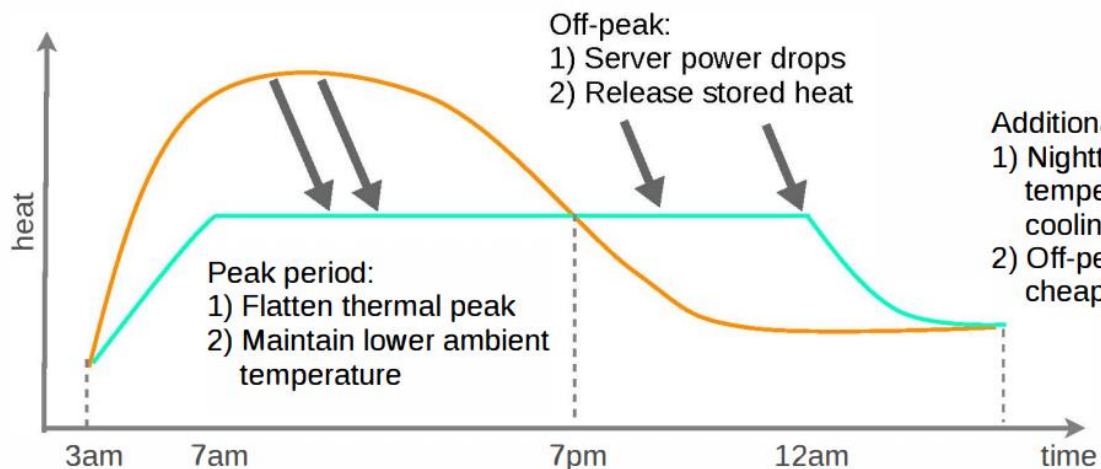


# Unconventional Cooling Techniques

- ❑ “Thermal Time Shifting” using high heat capacity material
  - Flatten the thermal peak by increasing heat capacity of the server
  - 12 percent cooling system reduction!



Box of wax



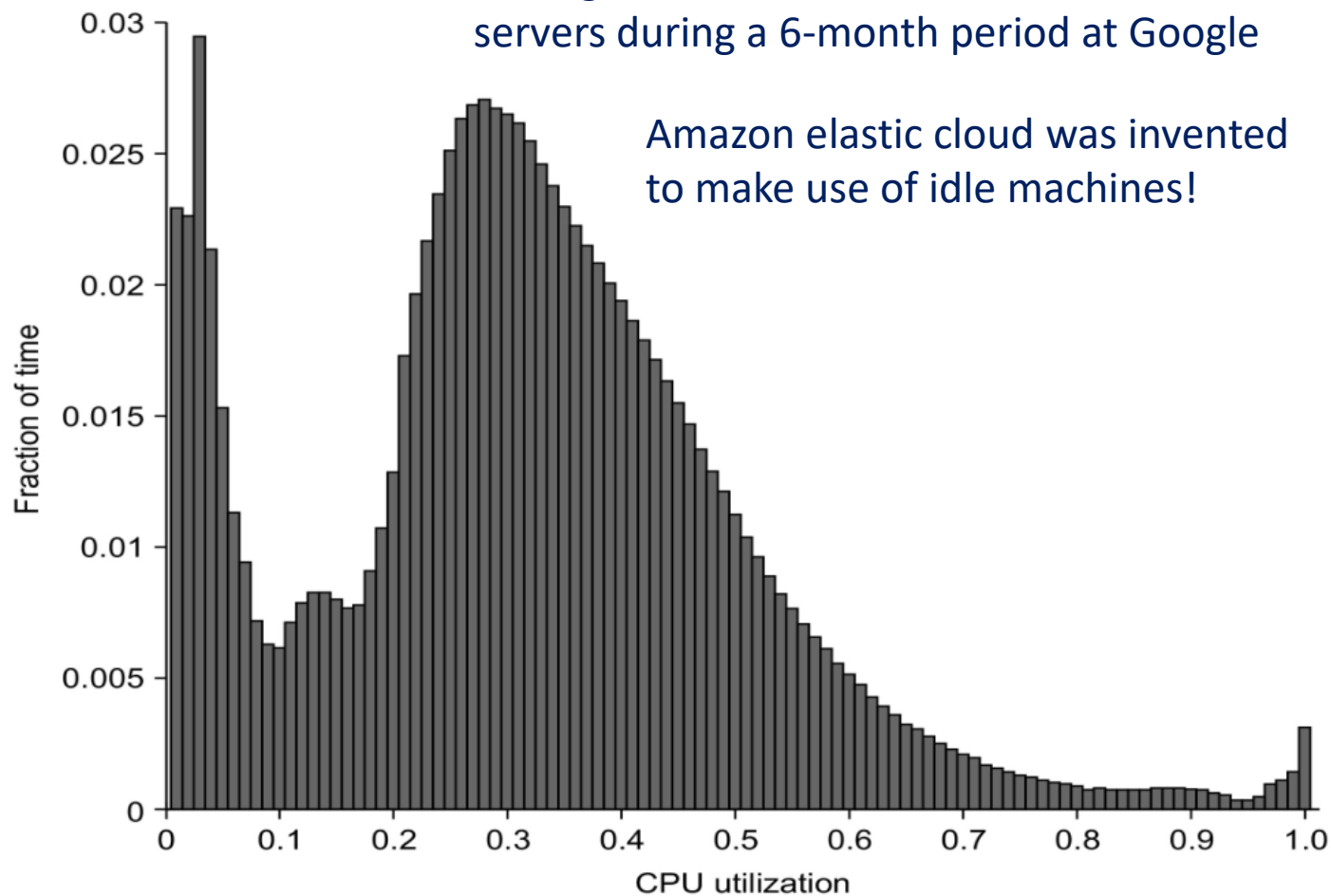
- Additional Advantages:
- 1) Nighttime: lower ambient temperature, more natural cooling opportunities
  - 2) Off-peak time: power is cheaper



# CPU Utilization is Usually Low

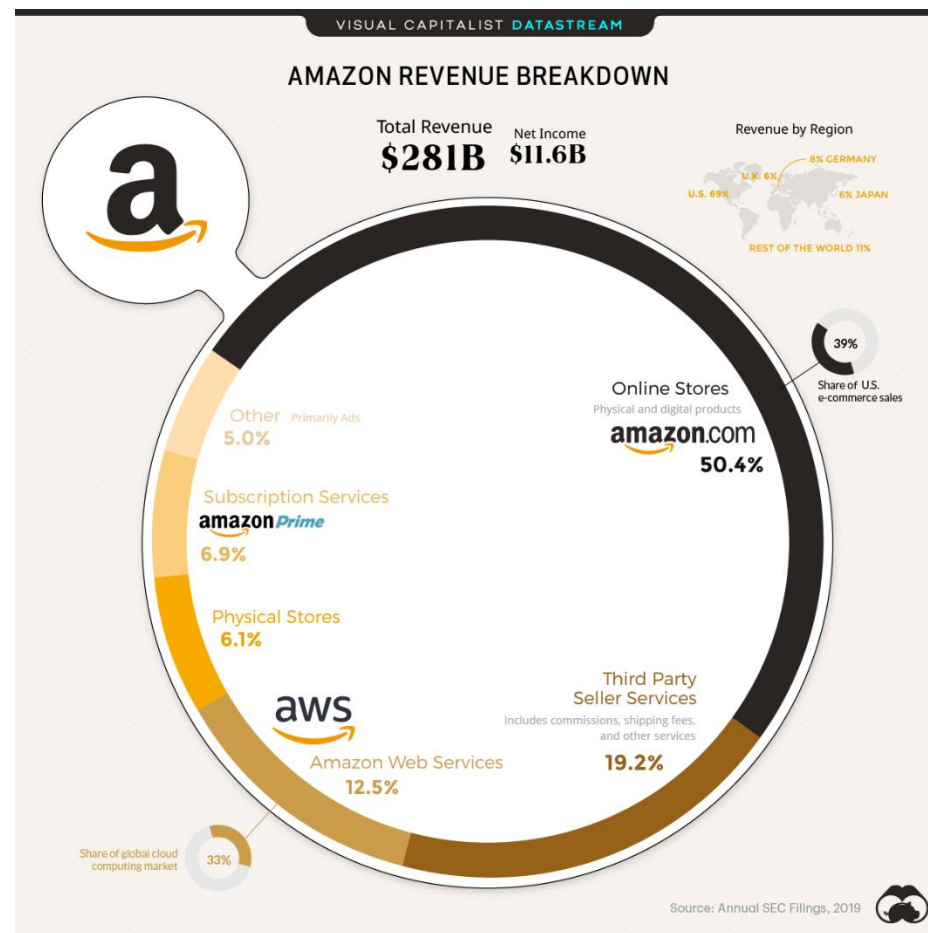
Average CPU utilization of more than 5000 servers during a 6-month period at Google

Amazon elastic cloud was invented to make use of idle machines!



And now...

(<https://www.visualcapitalist.com/how-amazon-makes-its-money/>)

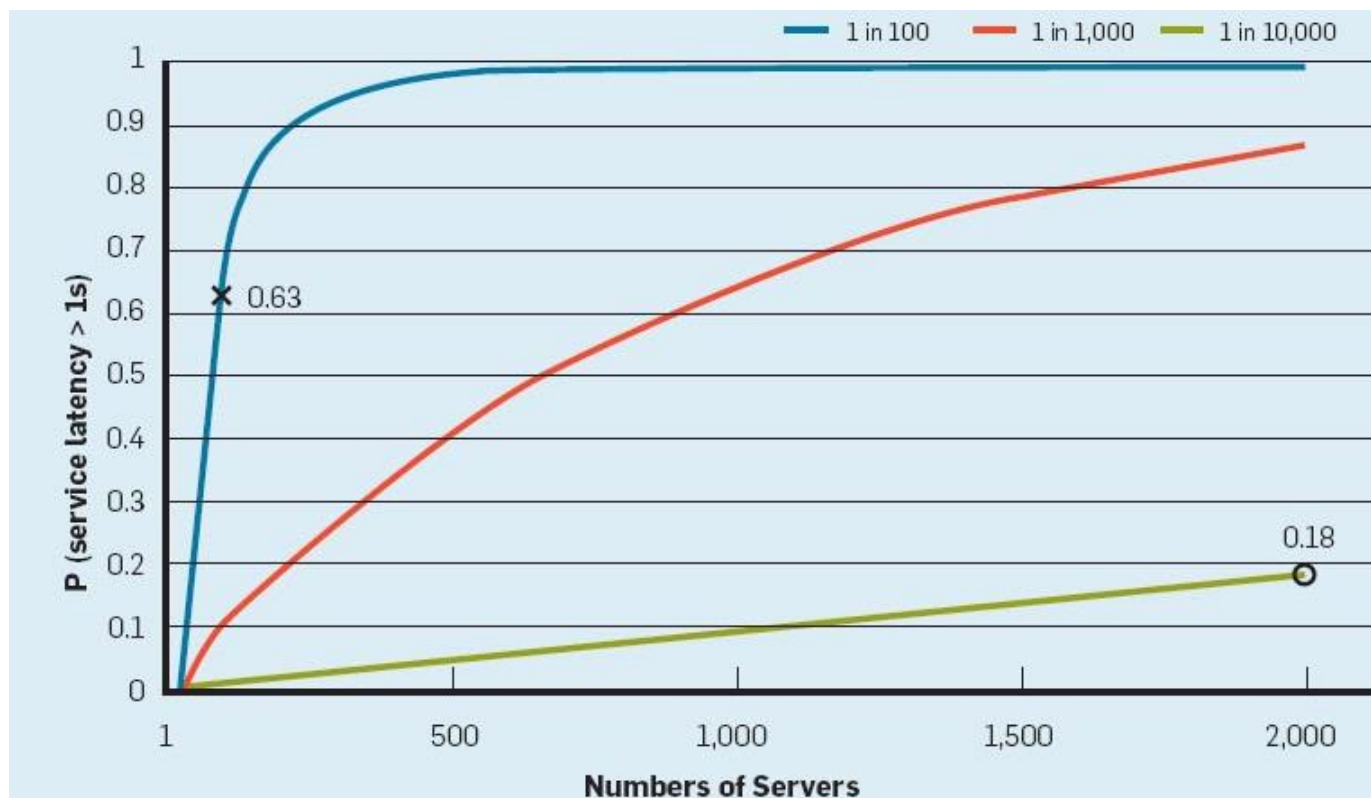


# Latency Issues With Larger Clusters

- ❑ Latency is important because it is felt by users
- ❑ Service Level Objectives (SLOs)/Service Level Agreements (SLAs)
  - e.g., 99% of requests below 100ms
- ❑ Tail latency increases with larger clusters
  - Probability of an outlier, high-latency response increases with larger clusters

# Latency Issues With Larger Clusters

- ❑ Even when server-level outliers are rare, probability piles up with more machines



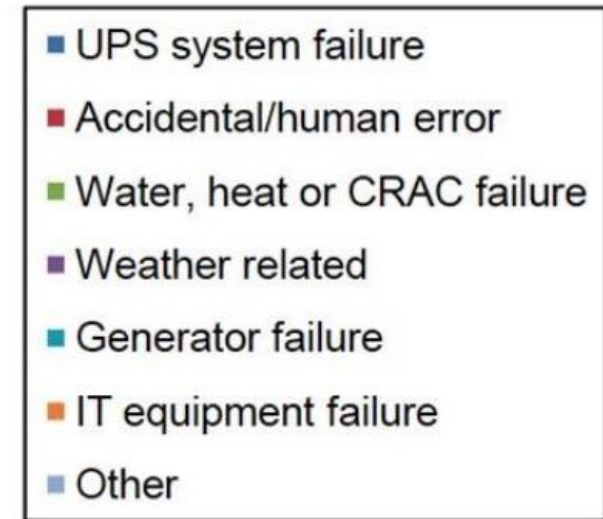
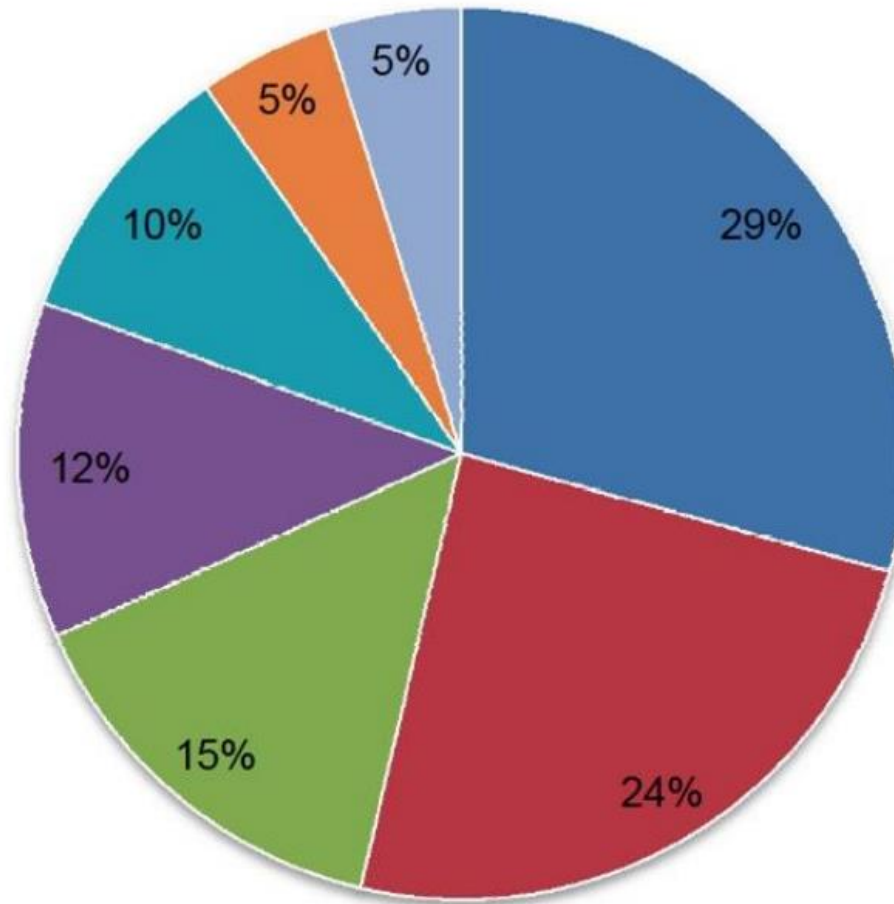
# Outages And Anomalies

Approx. number events in 1st year	Cause	Consequence
1 or 2	Power utility failures	Lose power to whole WSC; doesn't bring down WSC if UPS and generators work (generators work about 99% of time).
4	Cluster upgrades	Planned outage to upgrade infrastructure, many times for evolving networking needs such as recabling, to switch firmware upgrades, and so on. There are about nine planned cluster outages for every unplanned outage.
1000s	Hard-drive failures	2%–10% annual disk failure rate (Pinheiro et al., 2007)
	Slow disks	Still operate, but run 10 × to 20 × more slowly
	Bad memories	One uncorrectable DRAM error per year (Schroeder et al., 2009)
	Misconfigured machines	Configuration led to ~30% of service disruptions (Barroso and HÖlzle, 2009)
	Flaky machines	1% of servers reboot more than once a week (Barroso and HÖlzle, 2009)
5000	Individual server crashes	Machine reboot; typically takes about 5 min (caused by problems in software or hardware).

**Figure 6.1 List of outages and anomalies with the approximate frequencies of occurrences in the first year of a new cluster of 2400 servers.** We label what Google calls a cluster an *array*; see Figure 6.5. Based on Barroso, L.A., 2010. Warehouse Scale Computing [keynote address]. In: Proceedings of ACM SIGMOD, June 8–10, 2010, Indianapolis, IN.

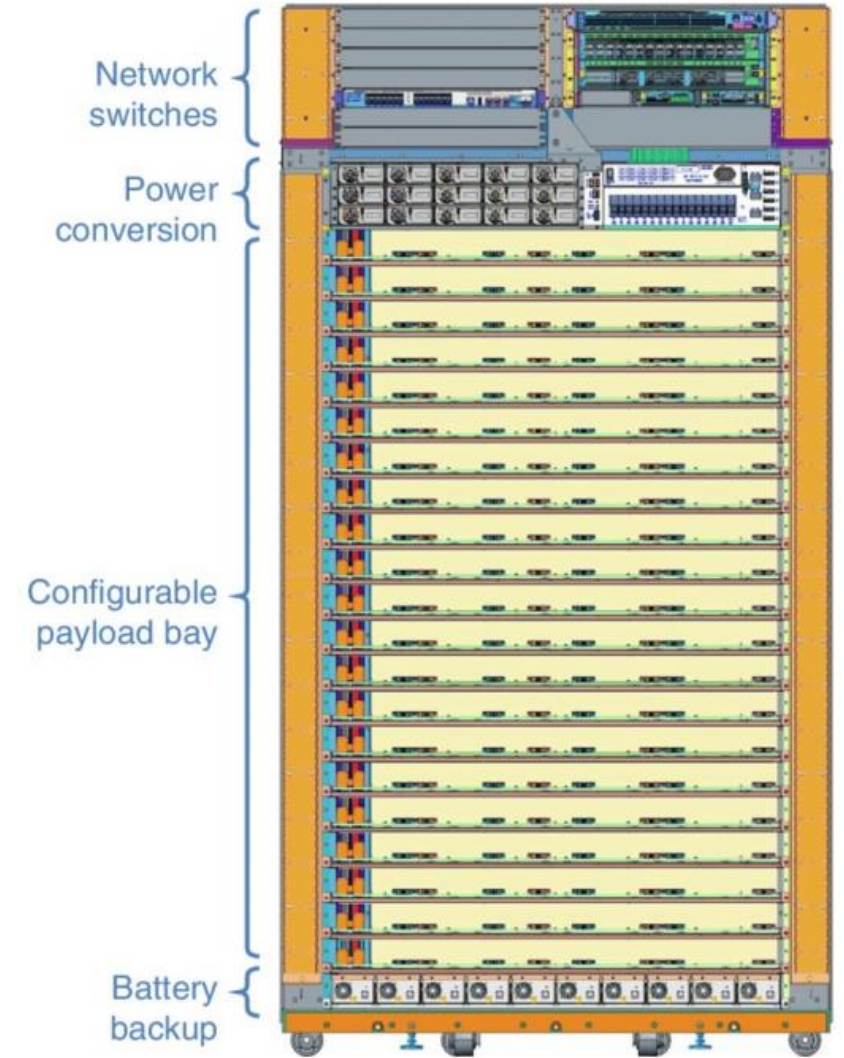
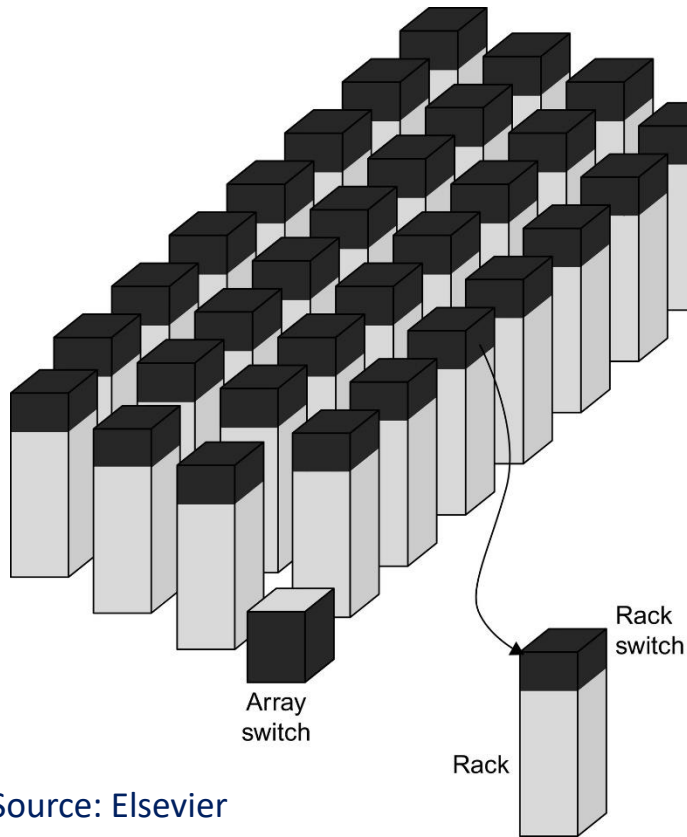
# Sources of Outages

Computed from 41 benchmarked data centers



# Rack-Level Architecture

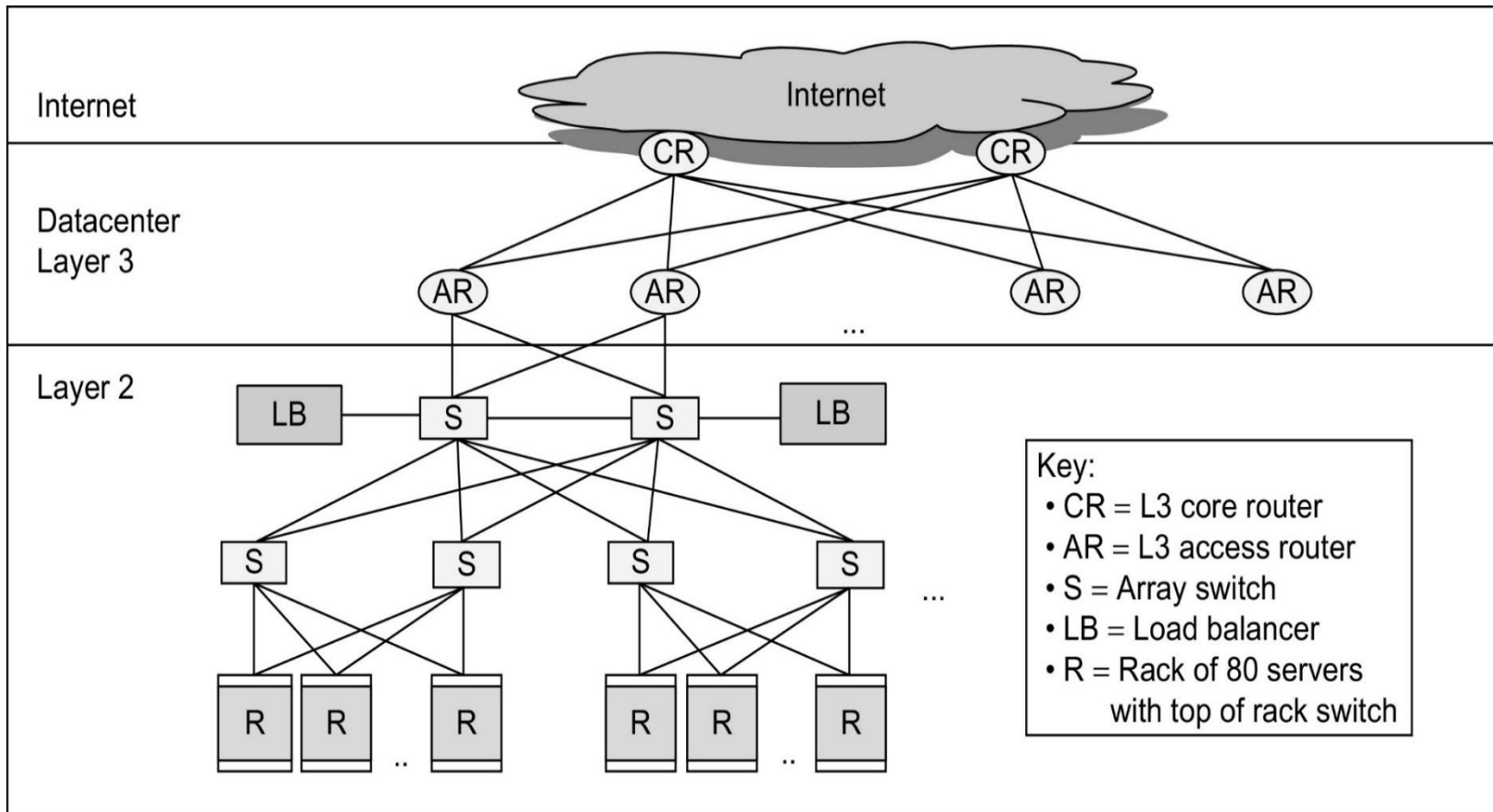
- ❑ Machines organized into racks
- ❑ Racks organized into larger units



**Figure 6.30 A Google rack for its WSC.** Its dimensions are about 7 ft high, 4 ft wide, and 2 ft deep (2 m × 1.2 m × 0.5 m). The Top of Rack switches are indeed at the top of this rack. Next comes the power converter that converts from 240 V AC to 48 V DC for the servers in the rack using a bus bar at the back of the rack. Next is the 20 slots (depending on the height of the server) that can be configured for the various types of servers that can be placed in the rack. Up to four servers can be placed per tray. At the bottom of the rack are high-efficiency distributed modular DC uninterruptible power supply (UPS) batteries.

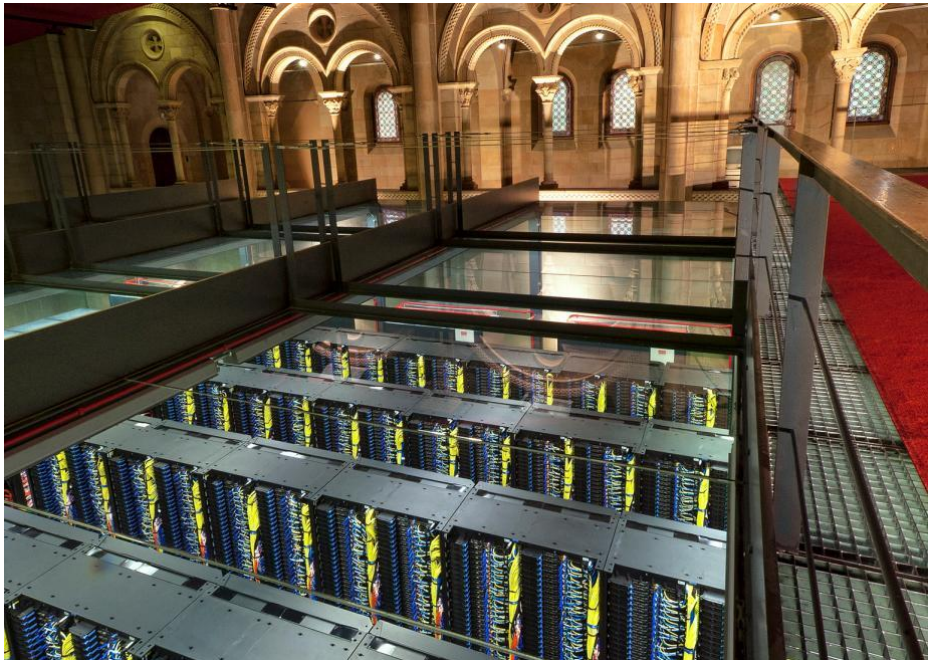
# Older Datacenter Network Architecture

- ❑ Upper-level routers may become bottleneck!



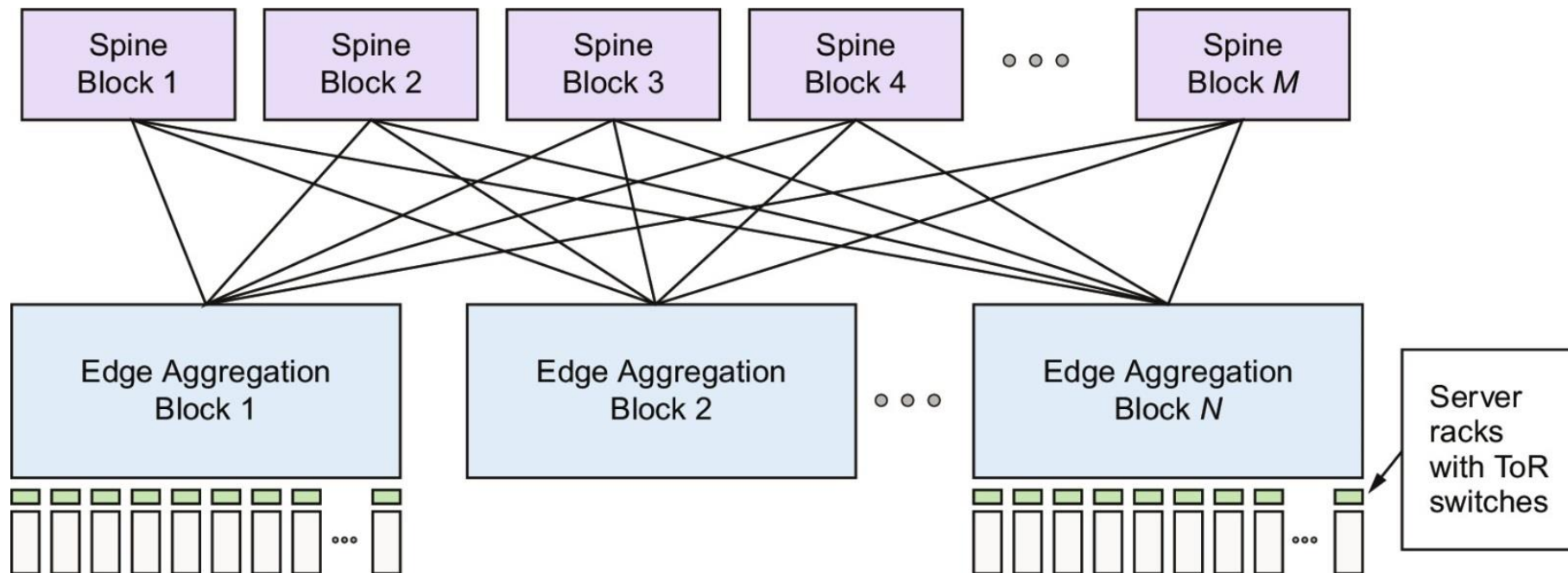
# Array (Aggregation) Switch

- ❑ Switch that connects an array of racks
  - Array switch should have 10 X the bisection bandwidth of rack switch
  - Cost of  $n$ -port switch grows as  $n^2$
  - Often utilize content addressable memory chips and FPGAs





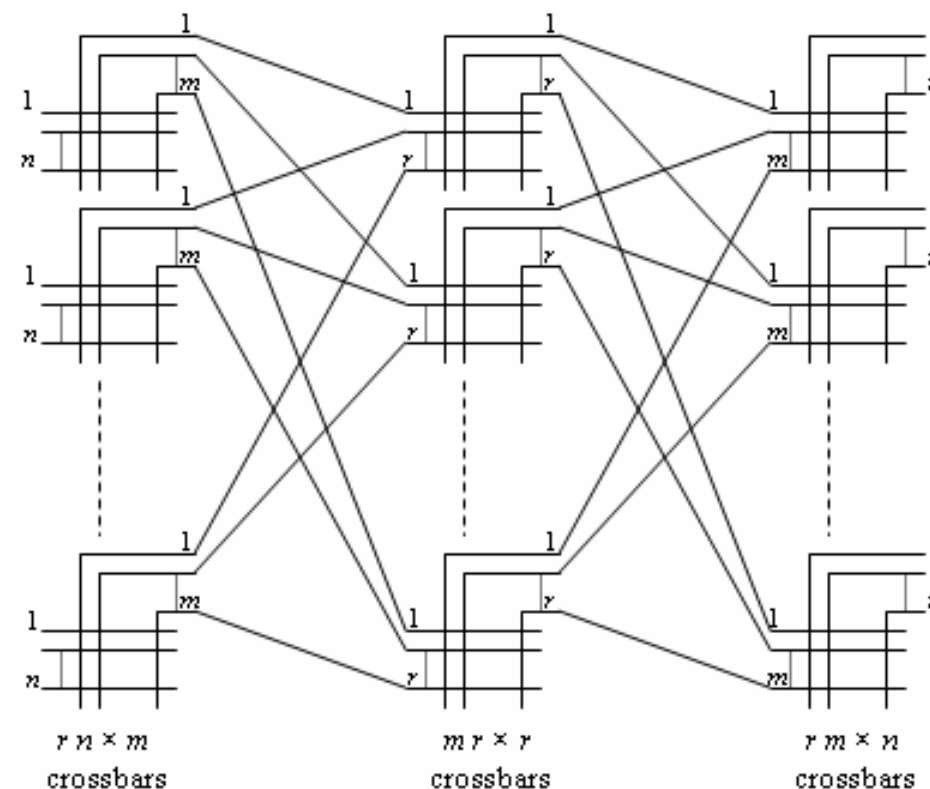
# Newer Solution: Clos Networks



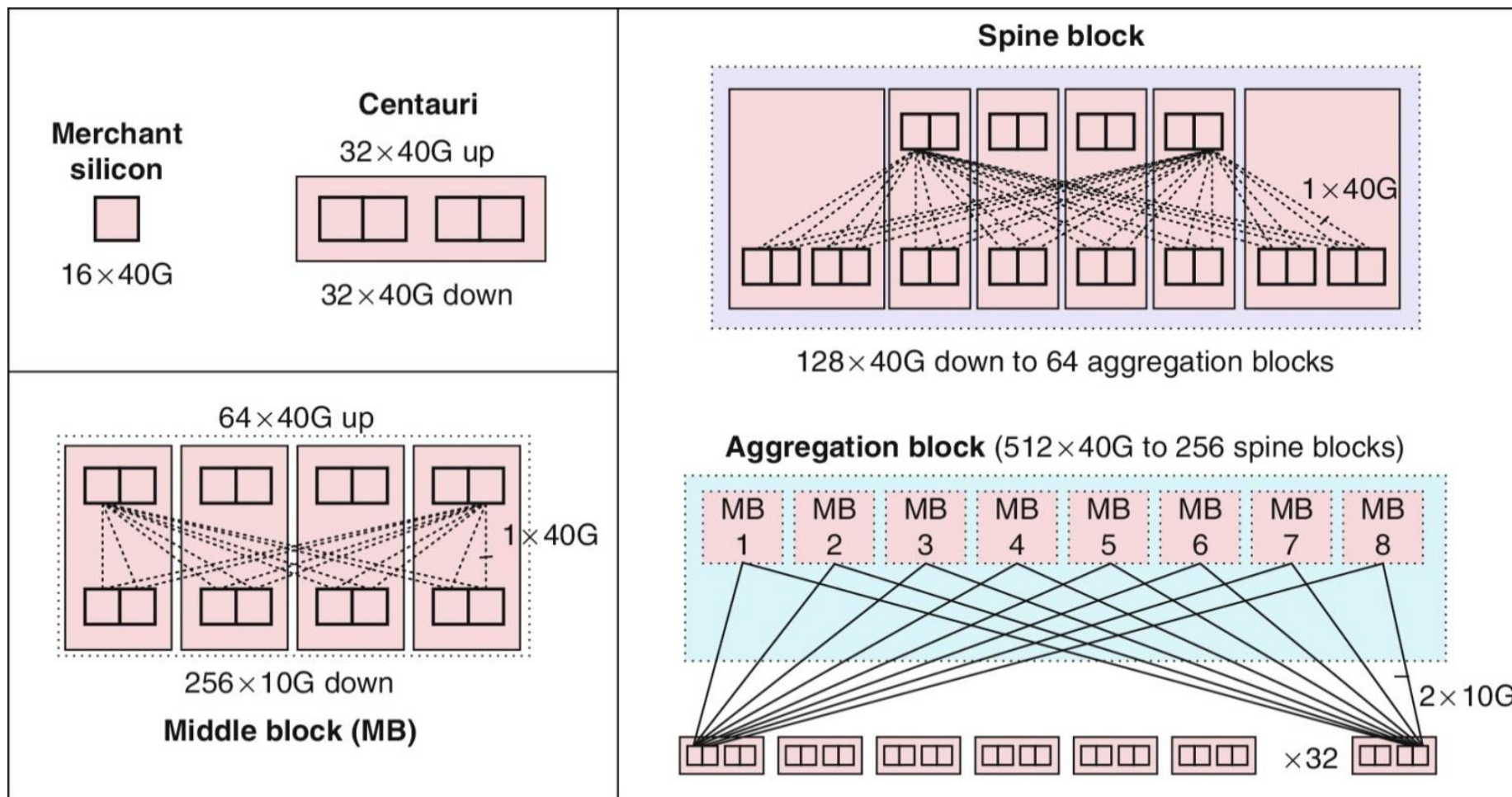
**Figure 6.31** A Clos network has three logical stages containing crossbar switches: ingress, middle, and egress. Each input to the ingress stage can go through any of the middle stages to be routed to any output of the egress stage. In this figure, the middle stages are the  $M$  Spine Blocks, and the ingress and egress stages are in the  $N$  Edge Activation Blocks. Figure 6.22 shows the changes in the Spine Blocks and the Edge Aggregation Blocks over many generations of Clos networks in Google WSCs.

# Newer Solution: Clos Networks

- ❑ Three layers of circuit-switched crossbar switches: ingress, middle, egress
- ❑ Characterized by three parameters
  - “r”: Number of ingress/egress switches
  - “n”: Client nodes per ingress/egress switch
  - “m” Number of middle switches
- ❑ Non-blocking if  $m \geq 2n+1$ 
  - Unused input in ingress can be connected to an unused output in egress without re-arranging existing paths



# Google Jupiter Clos Network



# Infiniband vs. Ethernet

- ❑ Datacenter network technology is typically either Infiniband or Ethernet
  - Both are industry standards
- ❑ IB typically enjoys lower latency due to streamlined protocol
  - Everything up to transport layer implemented in hardware!
- ❑ IB used to enjoy relatively faster bandwidth
  - With 40GbE and 100GbE, not necessarily true
- ❑ IB supports Remote Direct Data Access (RDMA) natively
  - Network transfer does not bother CPU any more, NIC hardware copies data directly to memory – Fast!
  - Currently, RDMA with an Ethernet network infrastructure possible with RDMA over Converged Ethernet (RoCE) protocol, requires special NIC

# Aside: Remote Direct Memory Access (RDMA)

