

Warehouse-Scale Computers

1

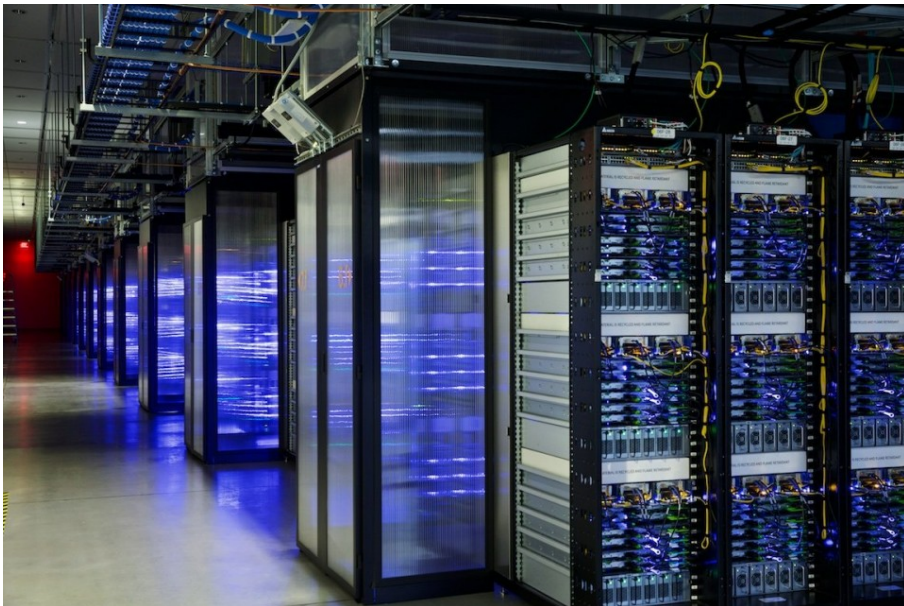
datacenter pictures



Google Council Bluffs, Iowa datacenter

2

datacenter pictures: servers racks



Facebook datacenter, Prineville, Oregon; via OregonLive

3

datacenter pictures: servers

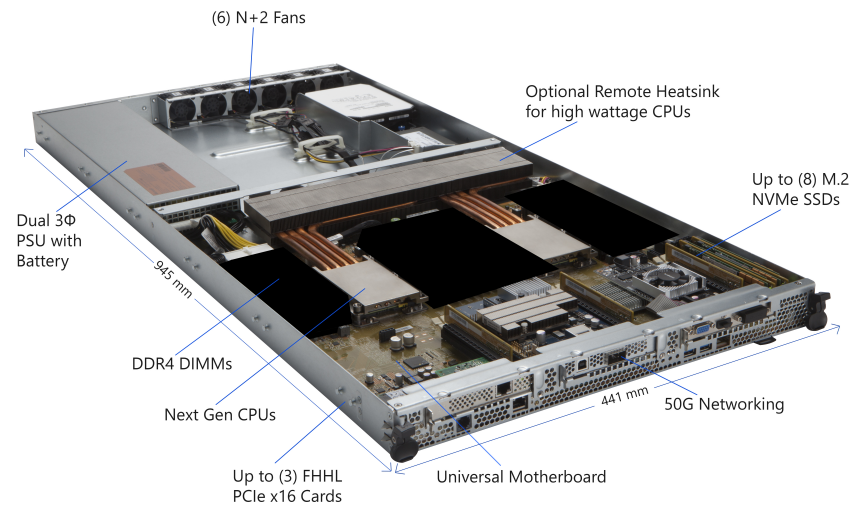
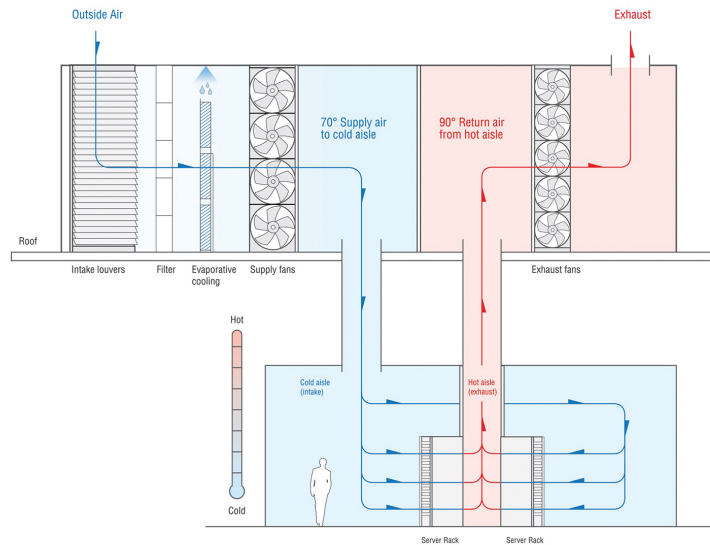


image: Open Compute Project (proposed 2016)

4

datacenter pictures: cooling



Sasser, "A Look at Data Center Cooling Technologies"

5

Mechanical Penthouse



Air mixing section - Return air / out side air / filter corridor



Evaporative cooling / humidification corridor

Facebook/Open Compute Project slide

6

Data Suite

- Hot aisle containment – ductless return
- Cold aisle pressurization – ductless supply



Facebook/Open Compute Project slide

7

datacenter pictures: backup power



Facebook datacenter, Prineville, Oregon; via TechCrunch

8

datacenter pictures: battery room



image: NOAA Center for Weather and Climate Prediction (at University of Maryland)

9

datacenter pictures: battery cabinet



Image: Facebook

10

datacenter pictures: TOR switch



Google "Pluto" Top-of-Rack Switch

11

datacenter sizes

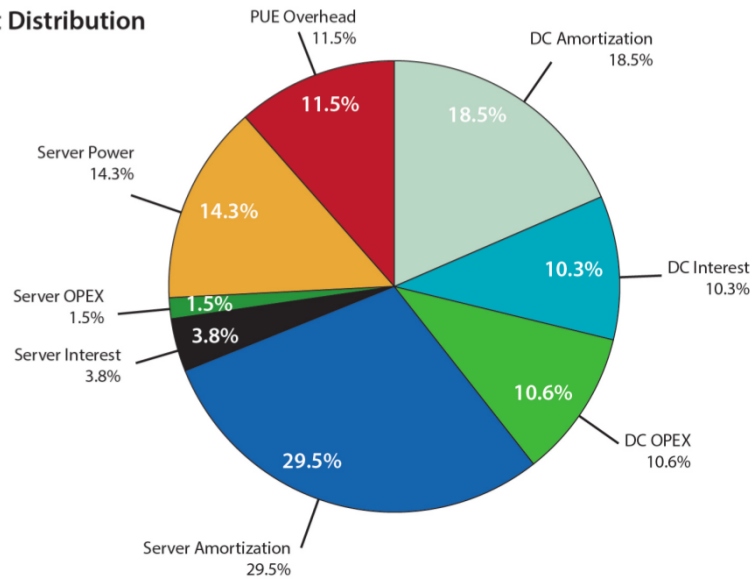
tens to hundreds of megawatts

tens of thousands to hundreds of thousands of servers

12

Money, Money, Money

Case B Cost Distribution



13

Kinds of cost

Operational:

power — e.g. cheap hydroelectric
maintenance — replacement equipment, etc.
people — sysadmins

Capital

buying/renting building + cooling + backup power
buying servers and replacing them when they become outdated

Common metric — cost per Watt

14

Datacenter Applications

“the web”/interactive:

latency matters

reliability matters

“free” parallelism — independent (mostly) requests

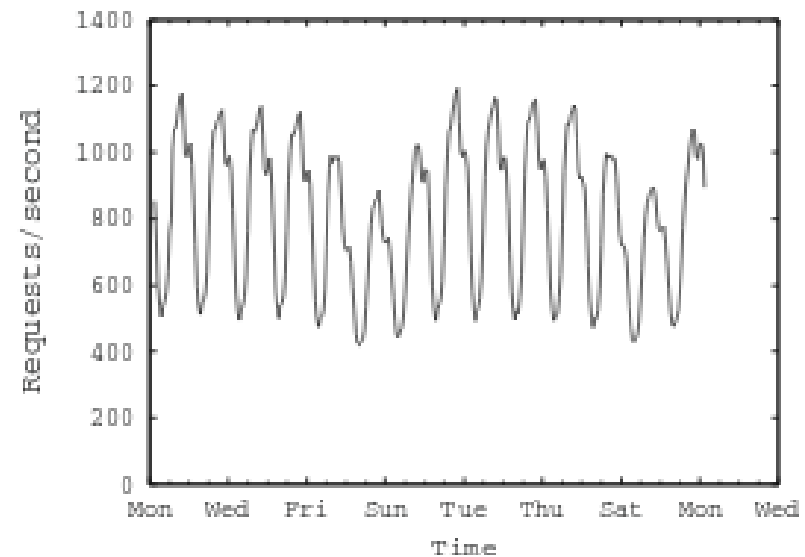
“batch”:

throughput matters

use ‘spare’ capacity from interactive stuff

15

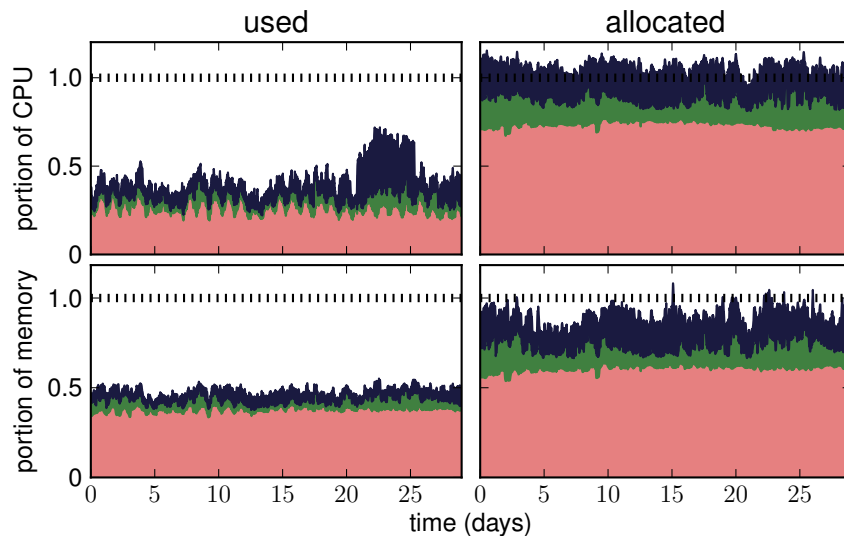
Varying demand



Urdaneta et al, “Wikipedia workload analysis for decentralized hosting”

16

Datacenter applications: consolidation/unpredictability



17

Datacenter versus Supercomputer

both **purpose-built**

different kinds of applications

datacenters tend to be more continuously upgraded

18

DC v SC: Goals

datacenter: focus on **cost-performance**

scale-out: more servers, not bigger machine
bigger individual machines are less efficient per dollar
want to use most mass-produced hardware

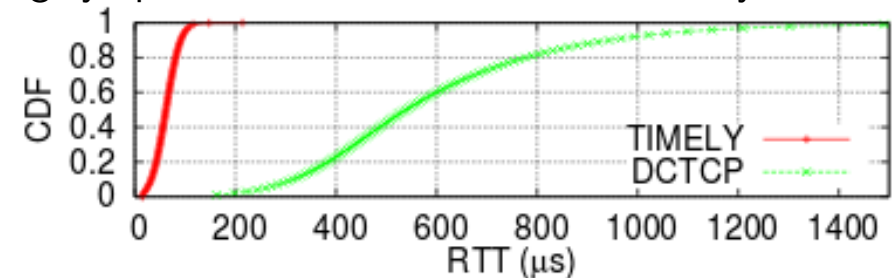
consolidation — run multiple applications together

more software modifications to use worse servers

19

DC v SC: Network

highly optimized datacenter network latency:



supercomputer network latency:
often less than ten microseconds round-trip

20

Datacenter Topology: historical

traditional datacenter topology:

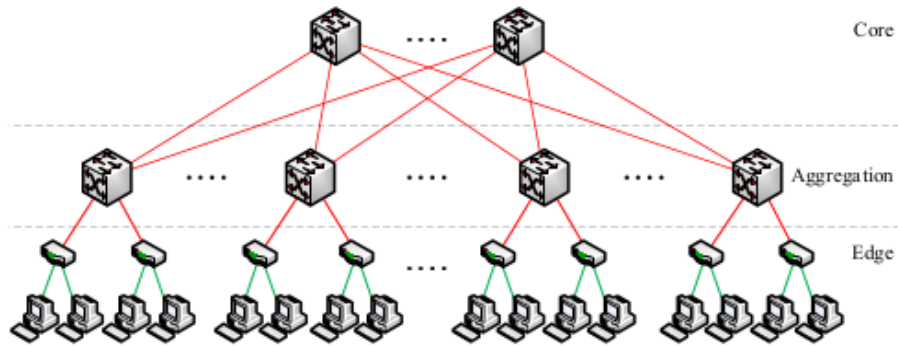


image: Al Fares et al, "A Scalable, Commodity Data Center Network Architecture"

21

Datacenter Topology: four-post

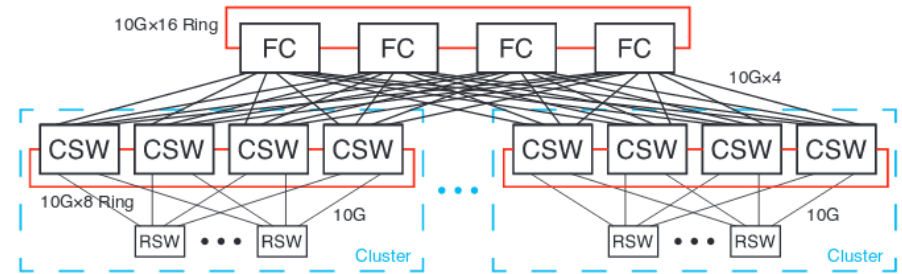


image: Farrington and Andreyev, "Facebook's data center network architecture"

22

Datacenter topology: Clos (1)

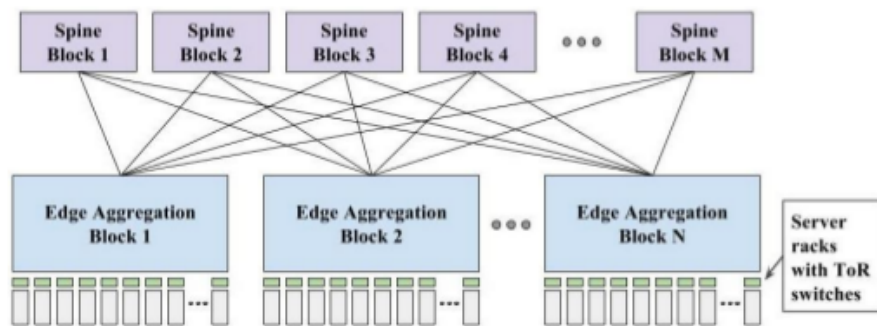


image: Singh et al, "Jupiter Rising: A Decade of Clos Topologies and Centralized Control in Google's Datacenter Network"

23

Datacenter topology: Clos (2)

Aggregation Block (32x10G to 32 spine blocks)

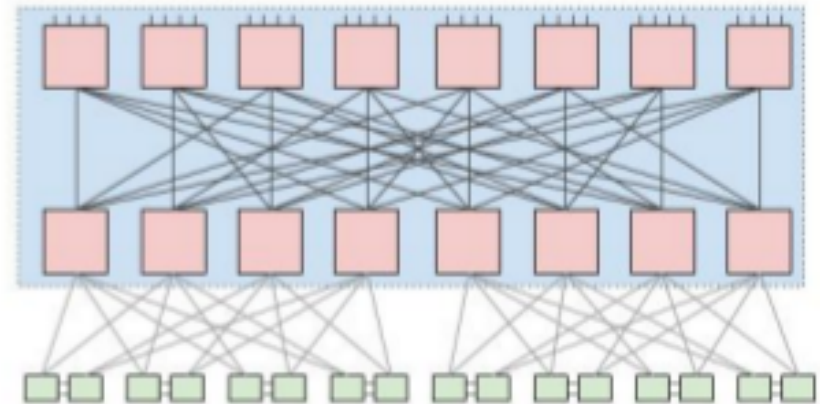


image: Singh et al, "Jupiter Rising: A Decade of Clos Topologies and Centralized Control in Google's Datacenter Network"

24

Datacenter Topology: Clos (3)

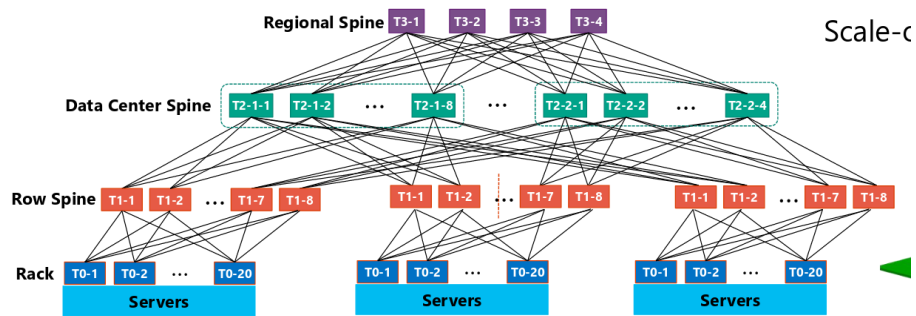


image: Greenberg, "SDN for the Cloud" (Microsoft)

25

DC v SC: Servers

very similar!

mass-produced, usually superscalar processors

usually **high-power** CPUs

... but not the most expensive

26

Server Balance

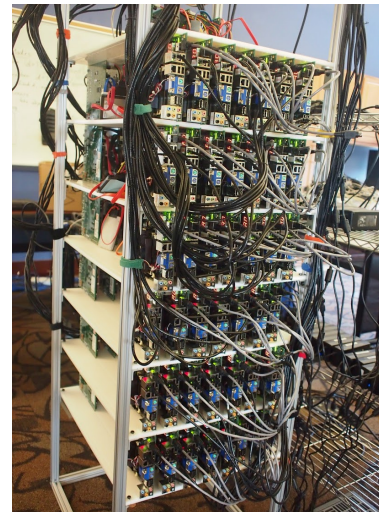
want to maximally use all server resources

balance CPU, memory, storage (disk or SSD)

depends on what applications you run

27

“wimpy” servers



another proposal: cheap, low-power servers at much higher density

28

DC v SC: Storage

storage on normal servers

- less networking required

- computations use local (fast) storage

seperate storage racks

- flat storage hierarchy, more convenient to program

29

DC v SC: Reliability

supercomputer: usually more reliable/expensive components

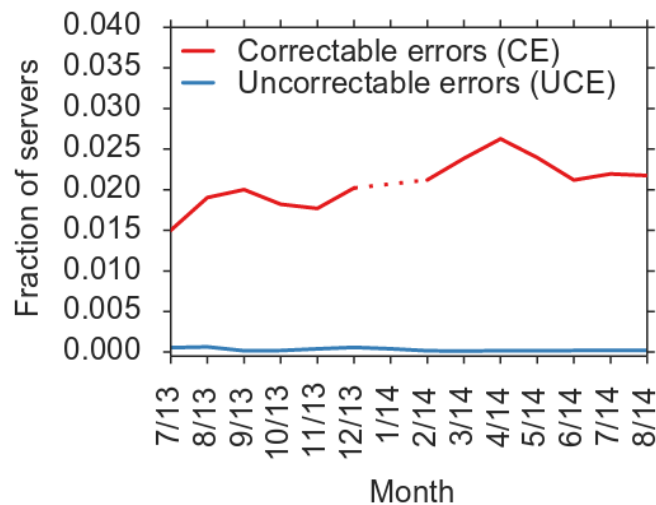
supercomputer: failures — reboot it all

datacenter: **expect failures**

datacenter: failures — work around broken component

30

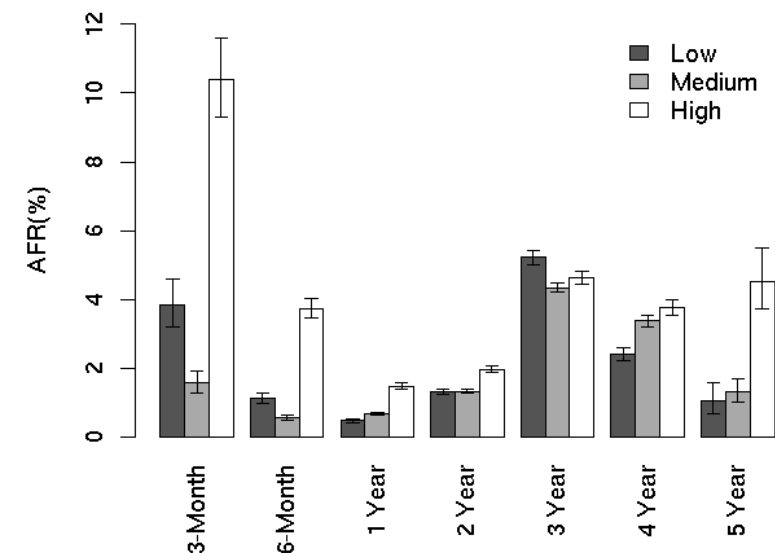
DRAM errors



uncorrectable: approx .03% of servers per month

31

Hard Drive failures



Pinheiro et al, "Failure Trends in a Large Disk Drive Population"

32

trading for software complexity

redundancy — handle failures

means having backup copies of everything

lots of applications per server — scheduling

slower network — compute **close to data**

33

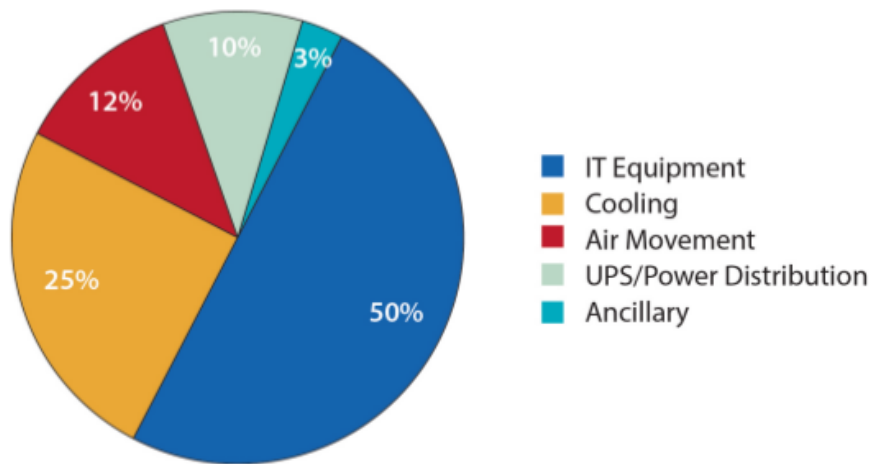
energy efficiency

also a problem for supercomputers, etc.

but optimized much more heavily in the datacenter era

34

old datacenter efficiency



35

PUE

$$\text{PUE} = \frac{\text{total power}}{\text{IT equipment power}}$$

servers and networking equipment

modern large datacenter: < 1.2

before attention to this problem, PUEs of 2 or more were common

36

Achieving high non-IT efficiency

airflow — don't mix hot/cold air

increased ambient temperature

cooling efficiency

- evaporative cooling
- better climates

power: increased electrical efficiency, e.g.:

- avoid AC/DC conversions
- distributed UPS
- get server power supplies that accept utility voltage

37

server efficiency

not especially well studied

similar losses from in-server power supplies, etc.

energy efficiency of components varies a lot

38

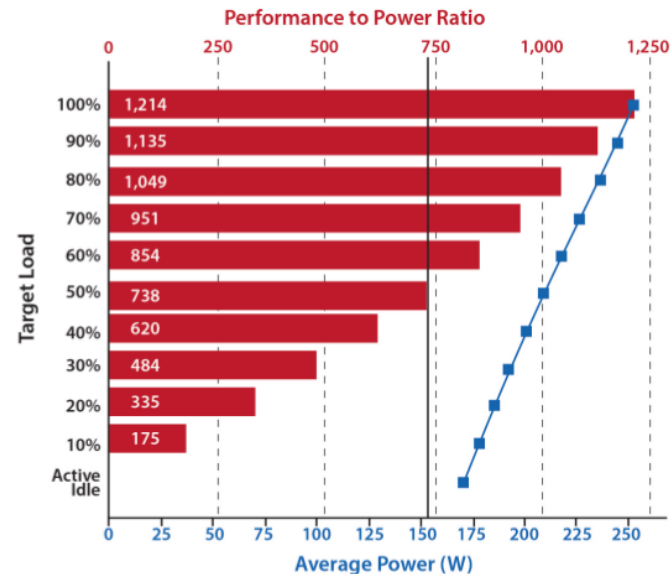
power-capping

underprovision cooling, power distribution, etc.

limit what runs on servers to stay under actual maximum

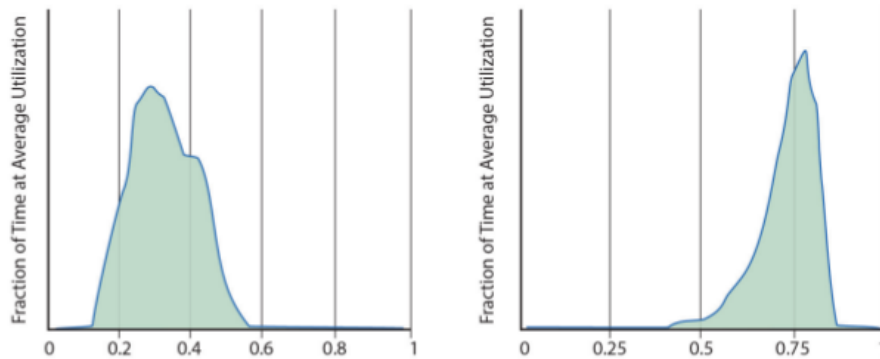
39

power proportionality problem (1)



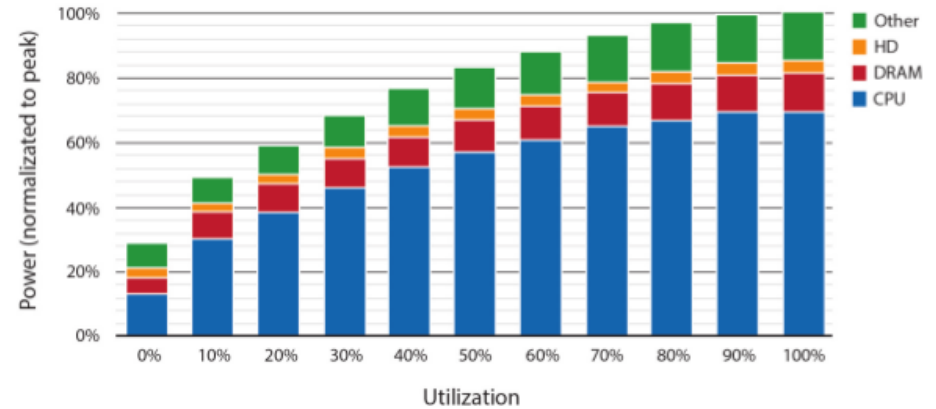
40

power proportionality problem (2)



41

power proportionality problem (3)



42

power-saving modes (1)

what about “sleep” modes?

- save a lot of power

- take milliseconds to seconds to start/end

servers need to be **available** continuously (e.g. stored data)

10% utilized server might be doing some work in every second

not enough time to really save power

43

power-saving modes (2)

processors have lower frequency/voltage modes

problem: doesn't save power in proportion to performance lost

problem: things other than processors use power

44

whack-a-mole in power usage

keep finding things which keep machine from sleeping for long times

keep finding components that use power continuously
tedious engineering problem

45

the datacenter for rent

public clouds — selling datacenter resources

e.g. Amazon Web Services

one way to deal with lower utilization

46

datacenter futures

started with: servers = desktop

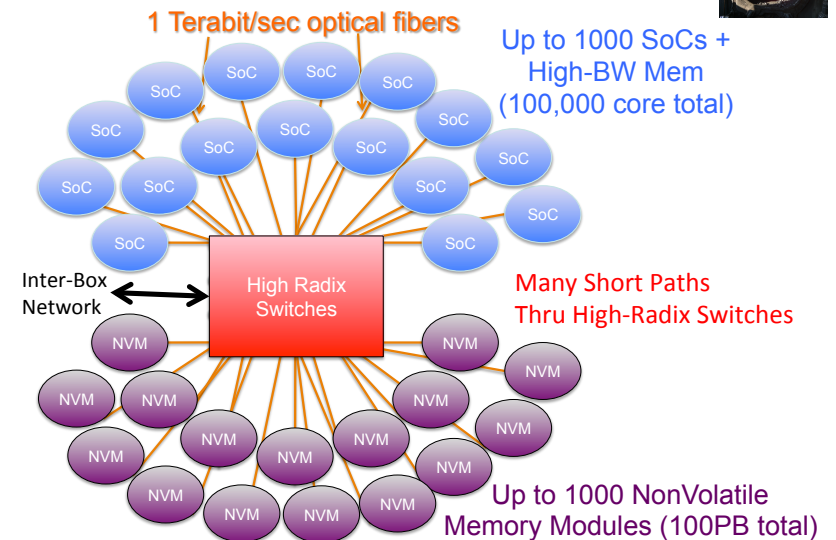
trend now: beefier servers

(revisiting old 'supercomputers'??)

47



FireBox Overview



48

datacenter futures

PCIe as a networking protocol within a rack?

fast, non-volatile RAM-like memories?

customized chips?

GPUs and FPGAs?

ASICs?

49

next time

general areas of HW security:

protect programs from each other — page tables, kernel mode, etc.

protect programs from adversaries — bounds checking, etc.

protect programs from people manipulating the hardware

paper: Smith and Weingart, "Building a high-performance, programmable secure coprocessor"

target audience: e.g. banks want to protect PINs

50

public key cryptography (1)

Smith and Weingart make extensive use of **digital signatures**

digital signatures use a **public/private keypair**

example use case: A wants to email B and have B know A wrote the email

51

public key-cryptography (2)

A generates keypair for communicating with B

public key: given to B; serves as **identity/name**
assumed known by/safe to tell everyone

private key: kept secret by A
assumed **no one else** has private key

52

public key cryptography (3)

two mathematical functions:

$signature = \text{Sign}(A's \text{ private key, message})$

$correct? = \text{Verify}(A's \text{ public key, message, signature})$

Verify will only say correct if private key was used
computationally infeasible to “forge” signature

A uses **Sign** operation, sends message and signature

B uses **Verify** operation; rejects if it says “not correct”

53

certificates

certificates are particular use of digital signature

example: A wants to help B communicate with C

certificate =

Sign(A's private key, “C's public key is XXX”)

certificate “proves” to B what C's public key is
if B trusts A enough

creating a certificate called “certifying”

54