

CS 744: DATACENTER AS A COMPUTER

Shivaram Venkataraman Fall 2021

ANNOUNCEMENTS

- Assignments
 - Assignment zero is due!
 - Form groups for Assignment I on Piazza
- Class format
 - Review
 - Lecture
 - Discussion



OUTLINE

- Hardware Trends
- Datacenter design
- WSC workloads
- Discussion

WHY IS ONE MACHINE NOT ENOUGH?

More date - deepit fit in I machine

ML / video encoding La parallelien when you use 71 martine

Single point of failure Applications might need specific Lardware

Compute Berring capacity



WHAT'S IN A MACHINE?

Interconnected compute and storage

Newer Hardware

- GPUs, FPGAs
- RDMA, NVlink



SCALE UP: MAKE MORE POWERFUL MACHINES

Moore's law

- Stated by Intel founder
 Gordon Moore
- Number of transistors on microchip double every 2 years
- Today "closer to 2.5 years"
 Intel CEO Brian Krzanich



DENNARD SCALING IS THE PROBLEM

Suggested that power requirements are proportional to the area for transistors

- Both voltage and current being proportional to length
- Stated in 1974 by
 Robert H. Dennard
 (DRAM inventor)

Broken since 2005



"Adapting to Thrive in a New Economy of Memory Abundance," Bresniker et al

DENNARD SCALING IS THE PROBLEM

Performance per-core is stalled

Number of cores is increasing



"Adapting to Thrive in a New Economy of Memory Abundance," Bresniker et al

MEMORY TRENDS



MEMORY TAKEAWAY



HDD CAPACITY



BACKBLAZE

HDD BANDWIDTH



Figure 4: Maximum sustained bandwidth trend

Disk bandwidth is not growing

SSDS

Performance:

- Reads: 25us latency ~ Out
- Write: 200us latency
- Erase: 1,5 ms

Steady state, when SSD full

- One erase every 64 or 128 reads (depending on page size)

Lifetime: 100,000-1 million writes per page

SSD VS HDD COST



ETHERNET BANDWIDTH

Ethernet

Growing 33-40% per year !





TRENDS SUMMARY

CPU speed per core is flat

Memory bandwidth growing slower than capacity

SSD, NVMe replacing HDDs

Ethernet bandwidth growing

DATACENTER ARCHITECHTURE



STORAGE HIERARCHY (DC AS A COMPUTER V2)



WAREHOUSE-SCALE COMPUTERS

Single organization

Homogeneity (to some extent)

Cost efficiency at scale

- Multiplexing across applications and services
- Rent it out!

Many concerns

- Infrastructure
- Networking
- Storage
- Software

. . .

– Power/Energy

- Failure/Recovery

SOFTWARE IMPLICATIONS

Reliability

Storage Hierarchy

Workload Diversity

Single organization

WORKLOAD: PARTITION-AGGREGATE



WORKLOAD: SCHOLAR SIMILARITY



Reduce Stage





Kroin very borg word MACHINE LEARNING

A intervive each for each warryle

Table 2.1: Six production applications plus ResNet benchmark. The fourth column is the total number of operations (not execution rate) that training takes to converge.

| Type of | Parameters (MiB) | Training | | | Inference |
|---------|---------------------|--------------|-----------|-------------|-------------|
| Neural | | Examples to | ExaOps to | Ops | Ops |
| Network | | Convergence | Conv | per Example | per Example |
| MLP0 | 225 | 1 trillion | 353 | 353 Mops | 118 Mops |
| MLP1 | 40 | 650 billion | 86 | 133 Mops | 44 Mops |
| LSTM0 | 498 | 1.4 billion | 42 | 29 Gops | 9.8 Gops |
| LSTM1 | 800 | 656 million | 82 | 126 Gops | 42 Gops |
| CNN0 | 87 | 1.64 billion | 70 | 44 Gops | 15 Gops |
| CNN1 | 104 | 204 million | 7 | 34 Gops | 11 Gops |
| ResNet | 98 | 114 million | <3 | 23 Gops | 8 Gops |

DISCUSSION



https://forms.gle/nFrMPkSAWTjMcgUp6

Scale-up vs Scale-out

DISCUSSION

100 GB of memory

Scale-up vs Scale-out Scale up [1 big machine] - Workloads that are hard to parallehize - less retwork B/W - less communication overhead b) depends on the workload

Scale - Out [10x 10GB of memory] - handling very large Repaciby / large data - incremental deployments

- Fault tilerance E fail reliability Now down strogglers

DISCUSSION







1 - in - 100 > W1 returns slow response [7: of time

NEXT STEPS

Next class: Storage Systems

Assignment I out Thursday. Submit groups before that!