

# Scalability

David Wihl  
Harvard CS109B, Spring 2018  
Feb 28, 2018

# Outline

- Why Scalability is Important
- Scale up - in-machine scalability
- Scale out - “Data Center as Computer”
- Future Directions

# Preparation for Deep Learning

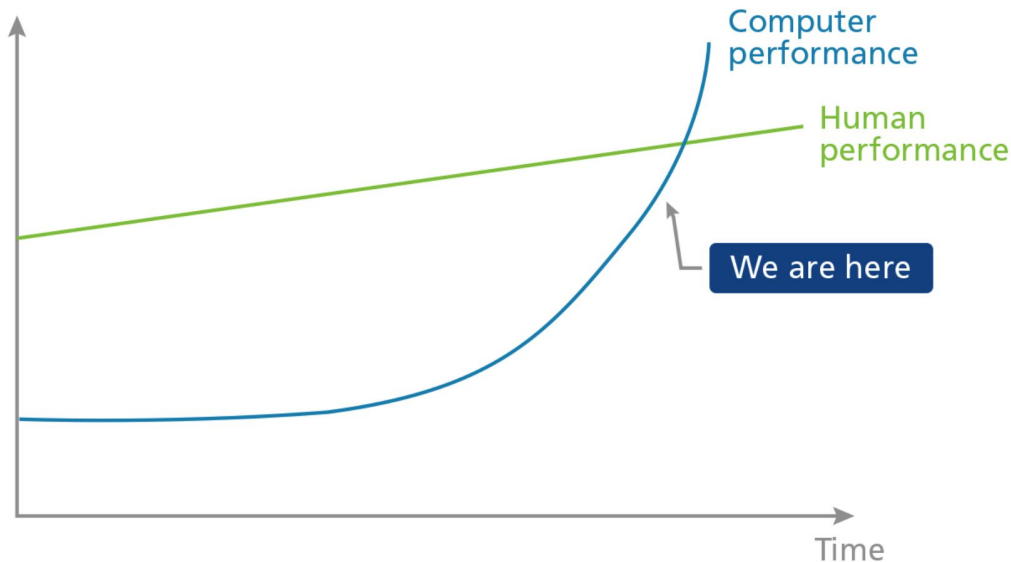
Data Size > RAM memory

Compute Requirements > Hours of CPU time

Current methods are sample inefficient

## Machine learning vs. human learning

Computer performance may outpace human performance



# Interesting Problems in Computer Science

P=NP?

Cryptography / One way functions

Machine Learning, including analysis of large datasets

Non-numeric information: images, music / sounds, text analysis

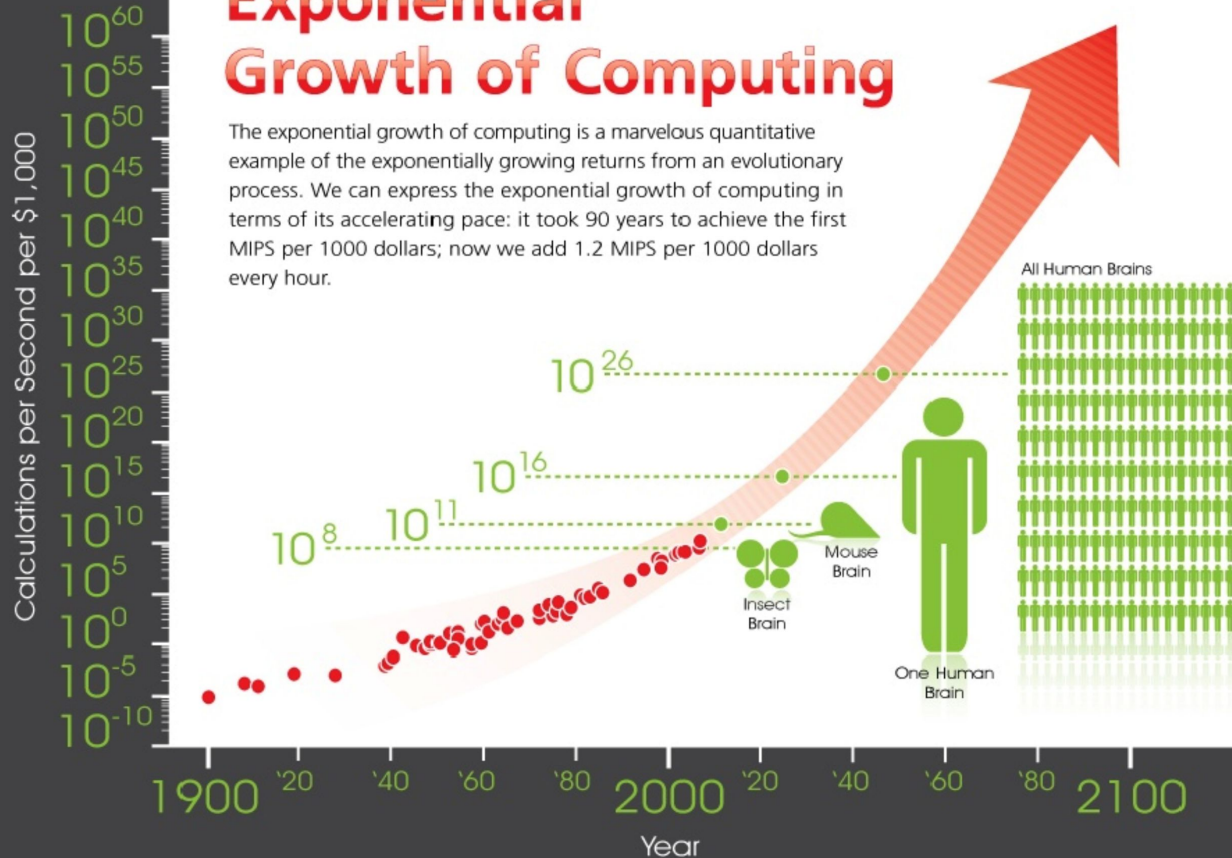
# Global Challenges

Bio / genomics

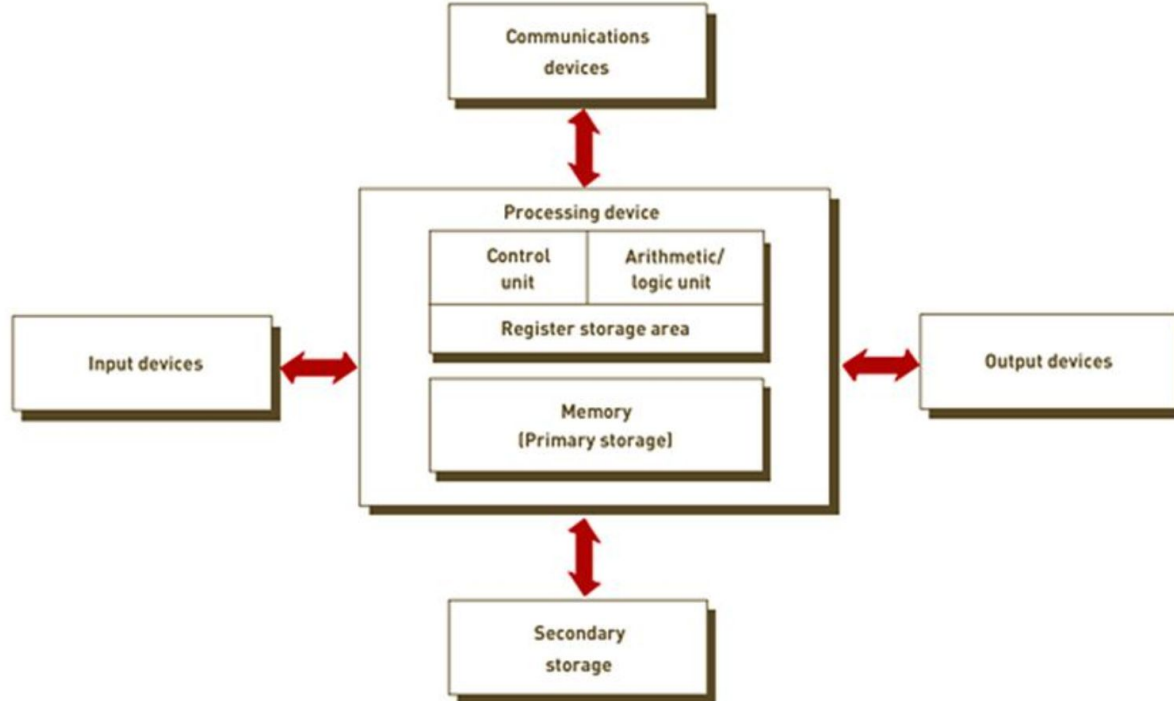
Environmental, including energy

# Exponential Growth of Computing

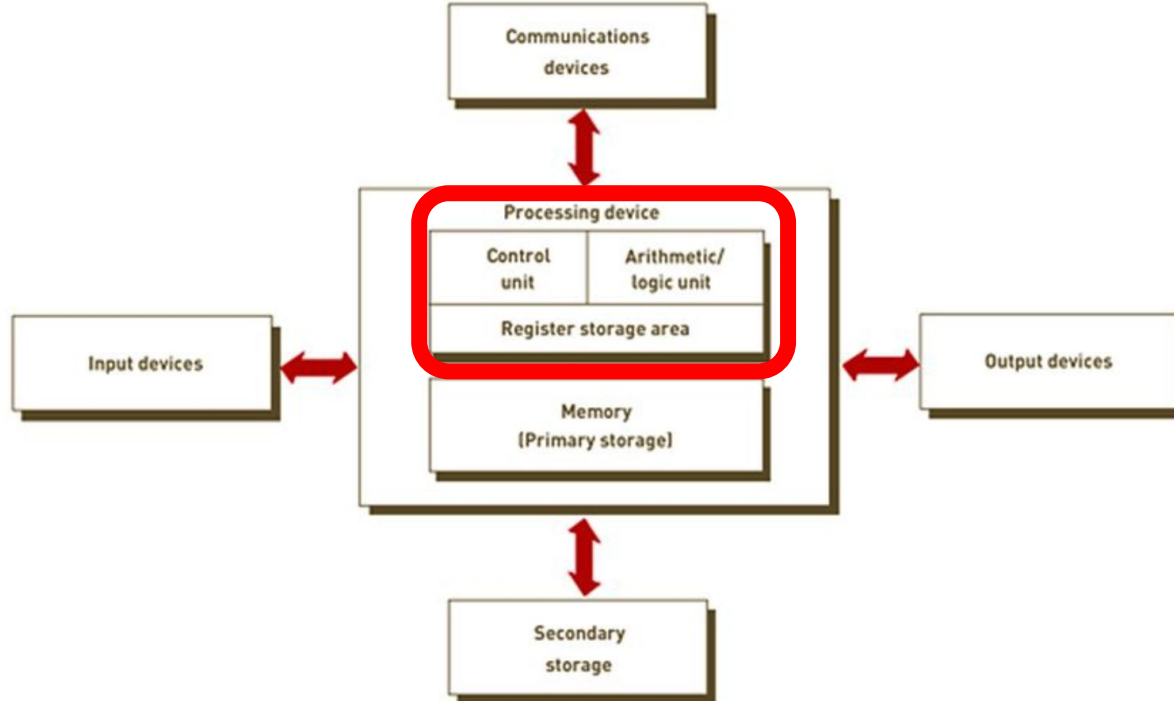
The exponential growth of computing is a marvelous quantitative example of the exponentially growing returns from an evolutionary process. We can express the exponential growth of computing in terms of its accelerating pace: it took 90 years to achieve the first MIPS per 1000 dollars; now we add 1.2 MIPS per 1000 dollars every hour.



# Scale Up - Performance of a Single Node

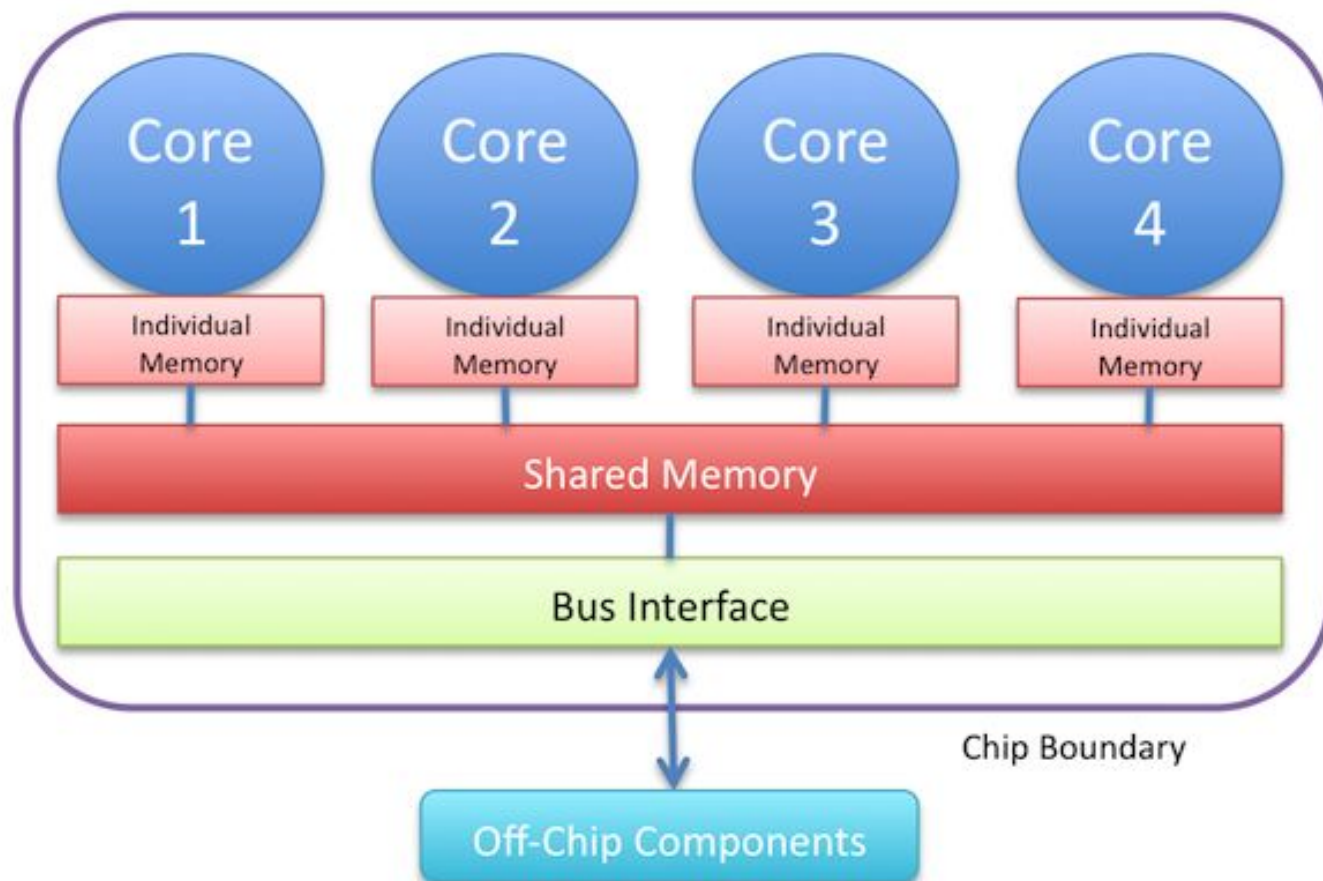


# Scale Up - Performance of a Single Node



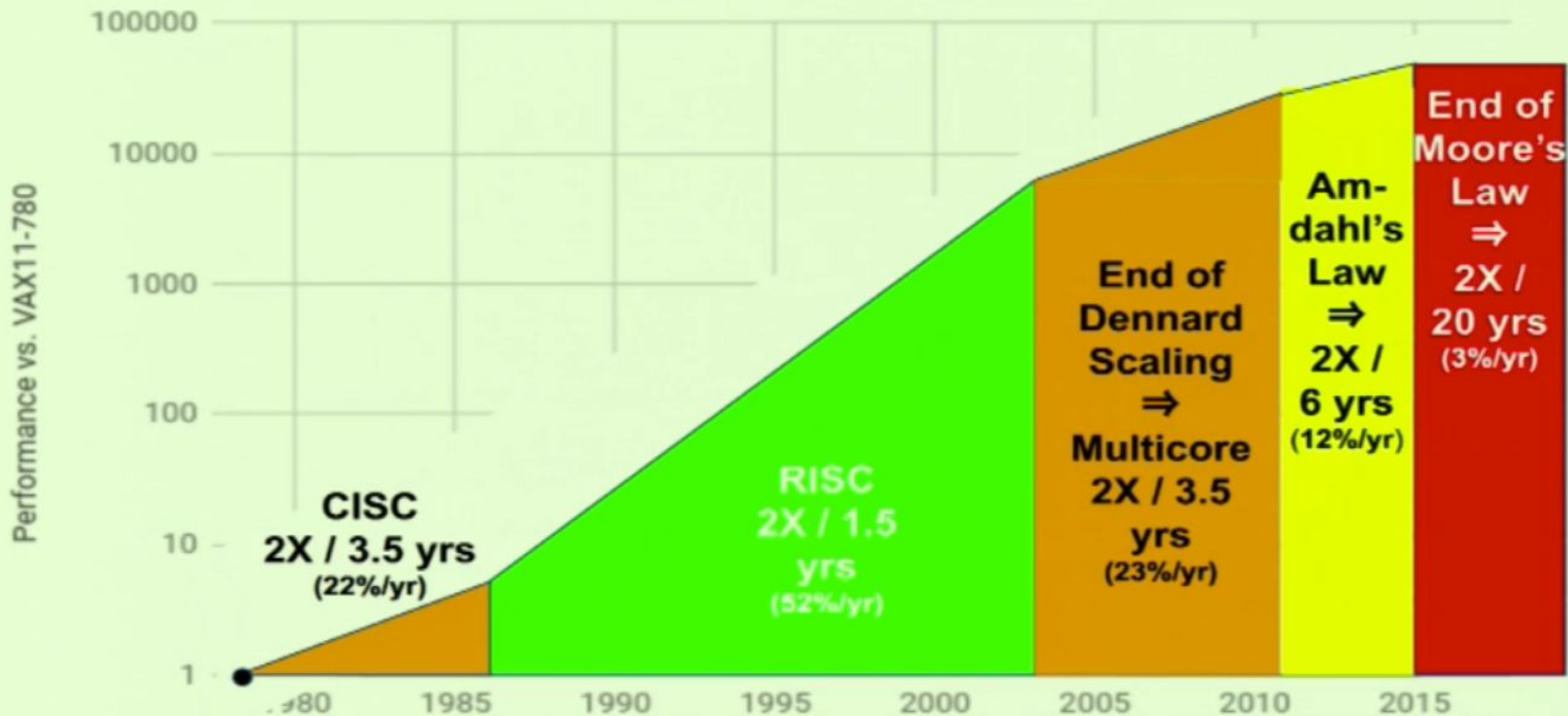


# Multi-core Processor



# End of Growth of Performance?

40 years of Processor Performance



<https://software.intel.com/en-us/blogs/2014/02/19/why-has-cpu-frequency-ceased-to-grow>

# Amdahl's Law: Limits of Parallelization

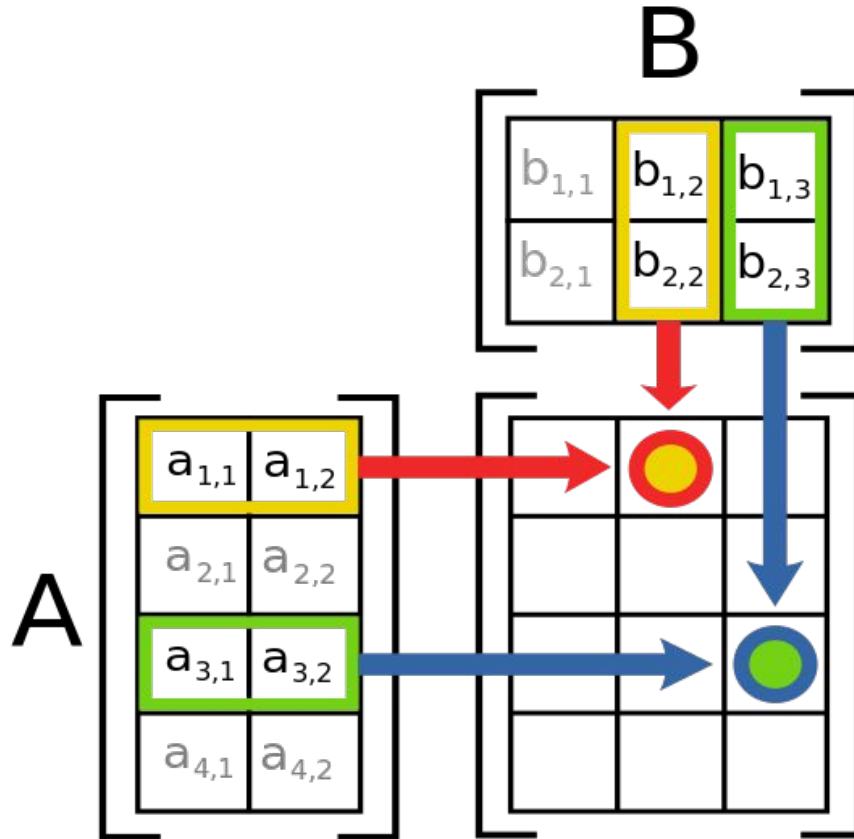
$$S_{\text{latency}}(s) = \frac{1}{(1 - p) + \frac{p}{s}}$$

Example: 30% of execution,  $p=0.3$ , executing on two cores  $s=2$

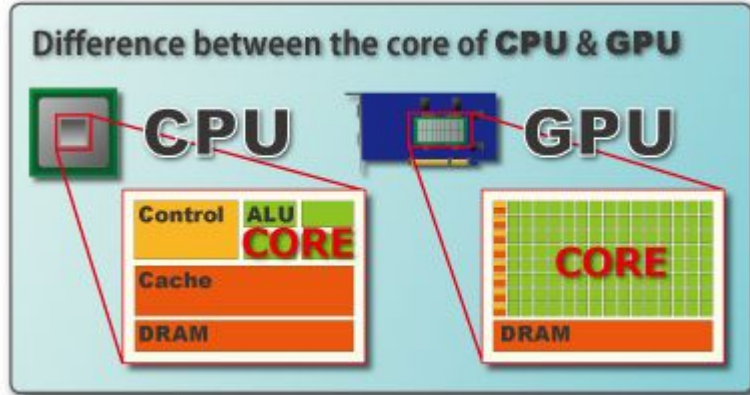
$$= \frac{1}{1 - 0.3 + \frac{0.3}{2}} = 1.18.$$

. . . . .

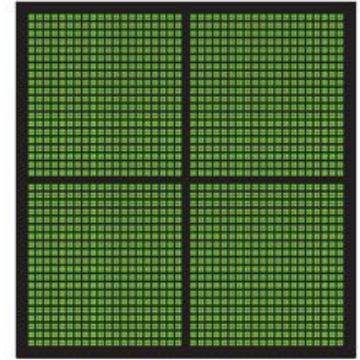
# Parallelizable Floating Point Operations



# Domain Specific Architectures: Enter the VPU, er GPU



CPU  
MULTIPLE CORES



GPU  
THOUSANDS OF CORES

# VALENTINE'S GIFT

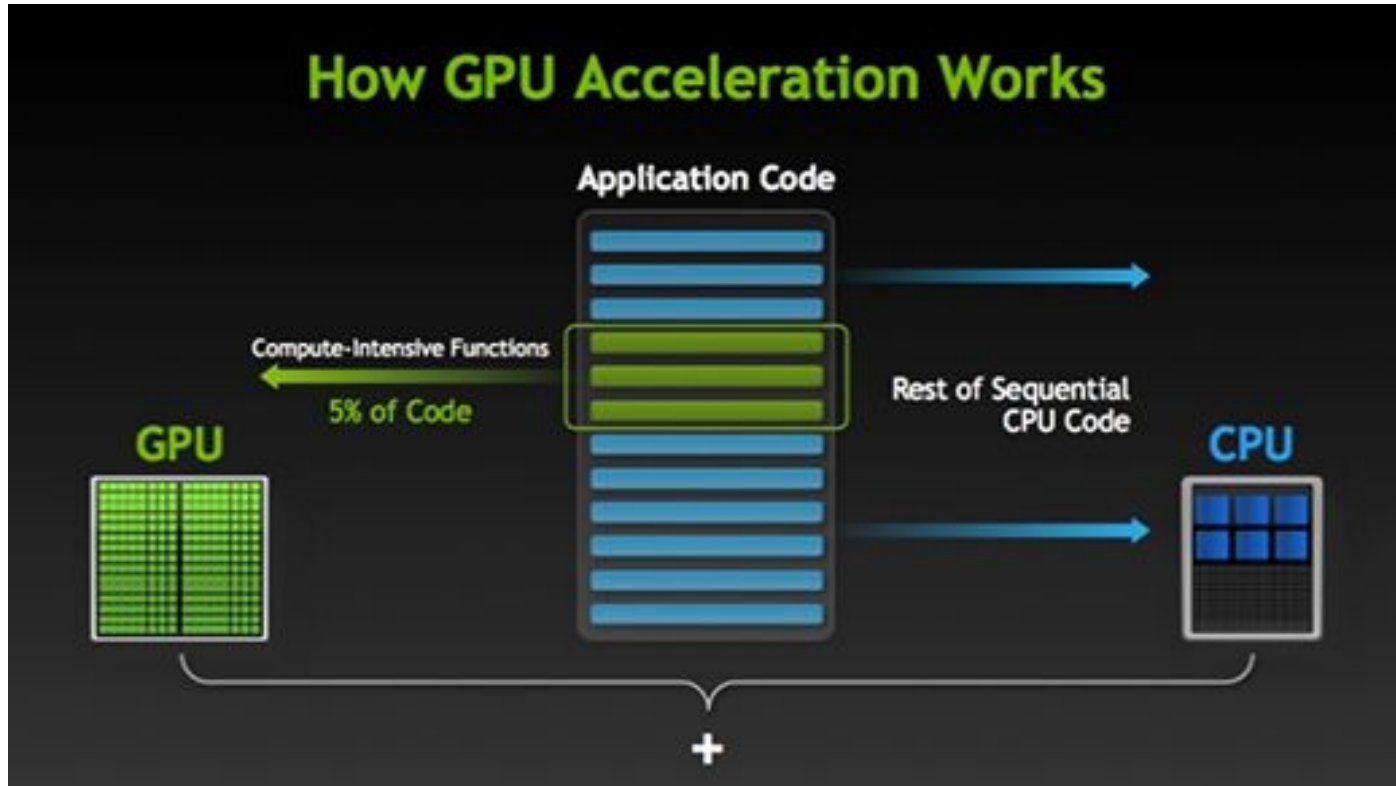
REALITY

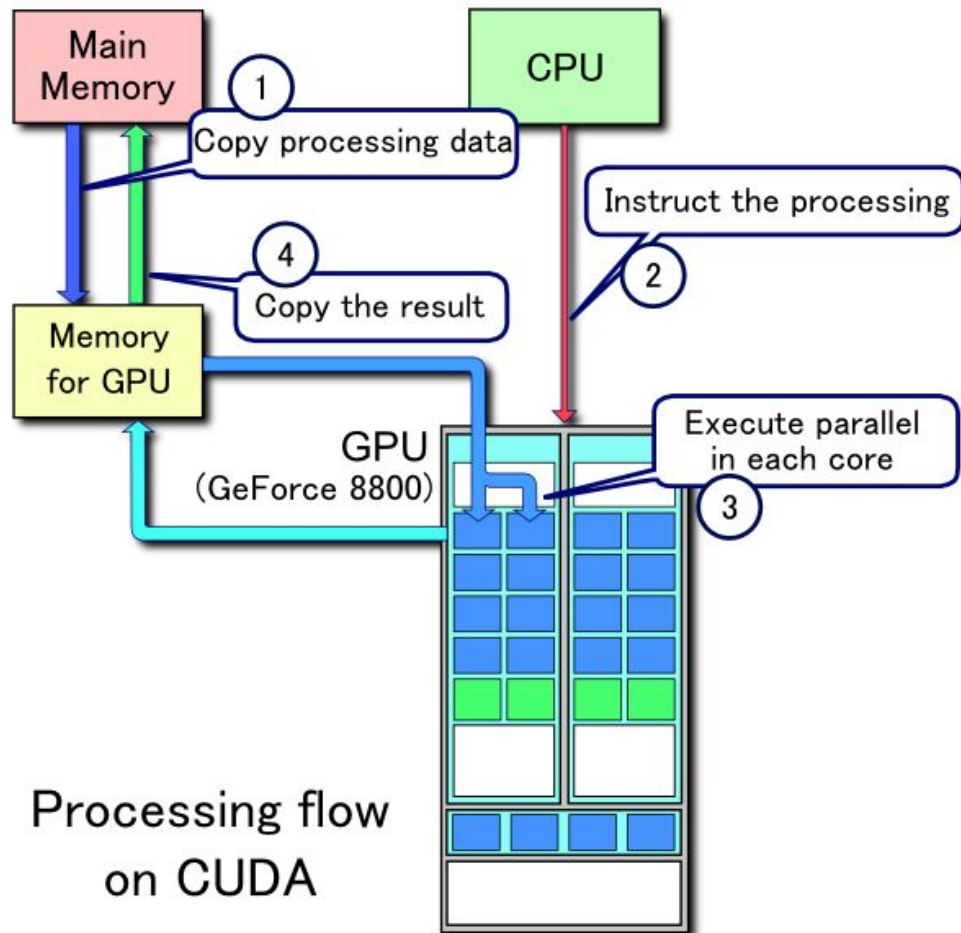


EXPECTATION



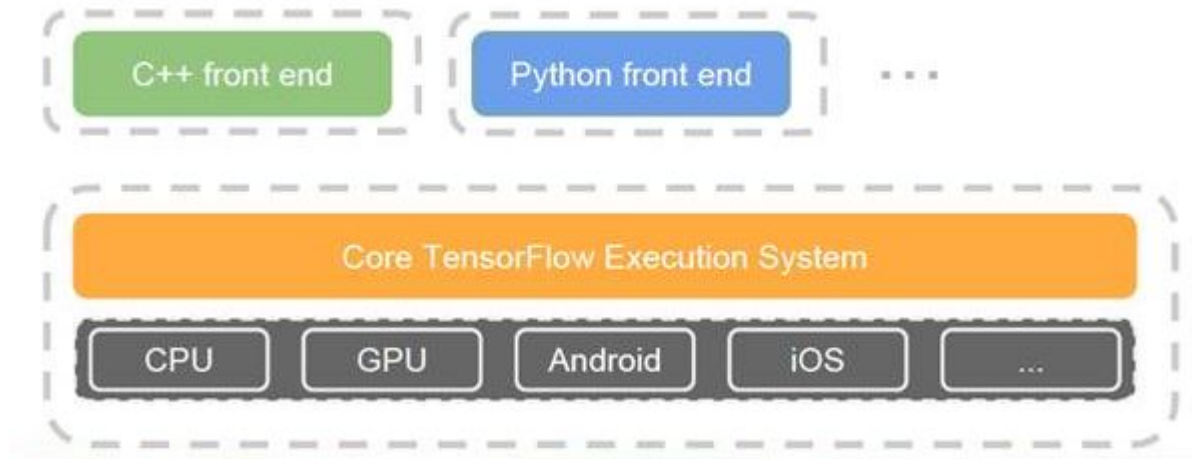
Requires new programming model:







# Abstracted by TensorFlow



# Even More Specialized - The TPU

“A Deep Neural Network Accelerator for the Datacenter”

Reduce **inference** phase by 10X vs. GPUs

Run-time vs. Train-time

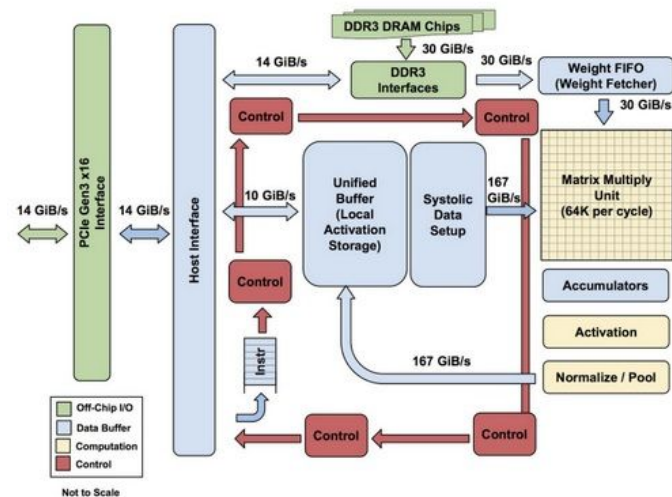
Systolic Execution - Matrices, not 1D vectors

Low precision: 8 bit FP, good enough for ML

Saved building multiple datacenters!

Cliff Young, BA and PhD from Harvard

[https://video.seas.harvard.edu/media/CS+Colloquium+Cliff+Young+2017-04-20/1\\_soxmc716](https://video.seas.harvard.edu/media/CS+Colloquium+Cliff+Young+2017-04-20/1_soxmc716)

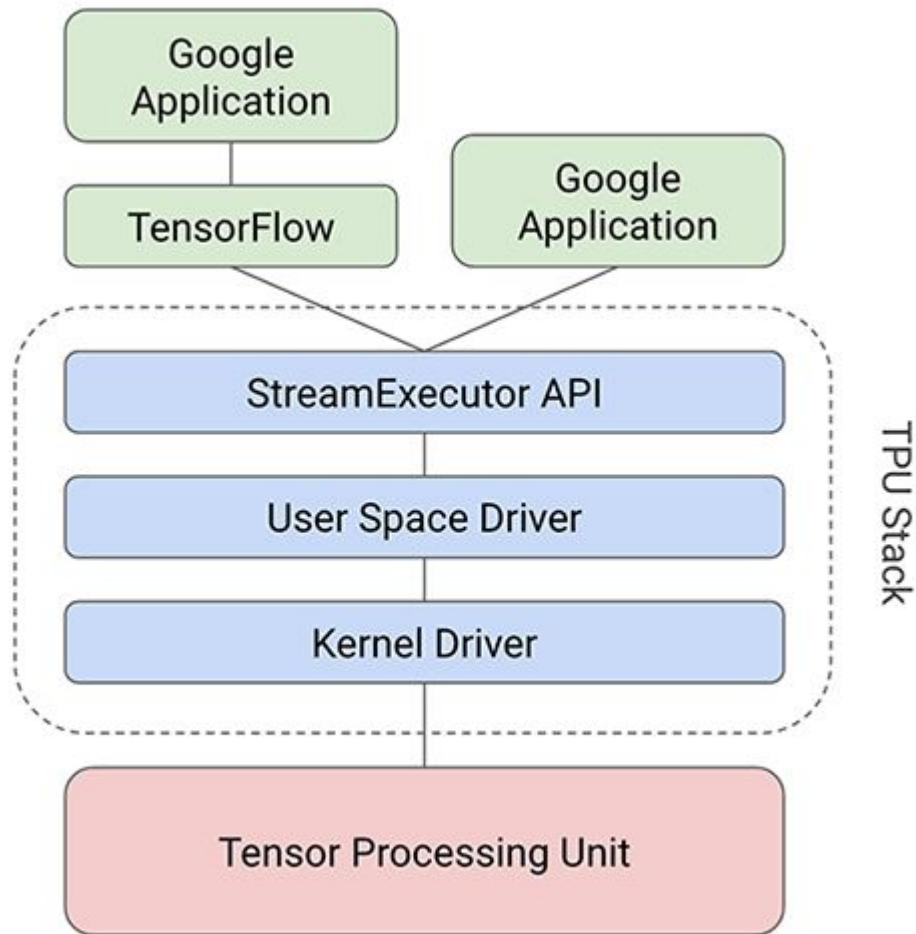


# TPU Limitations

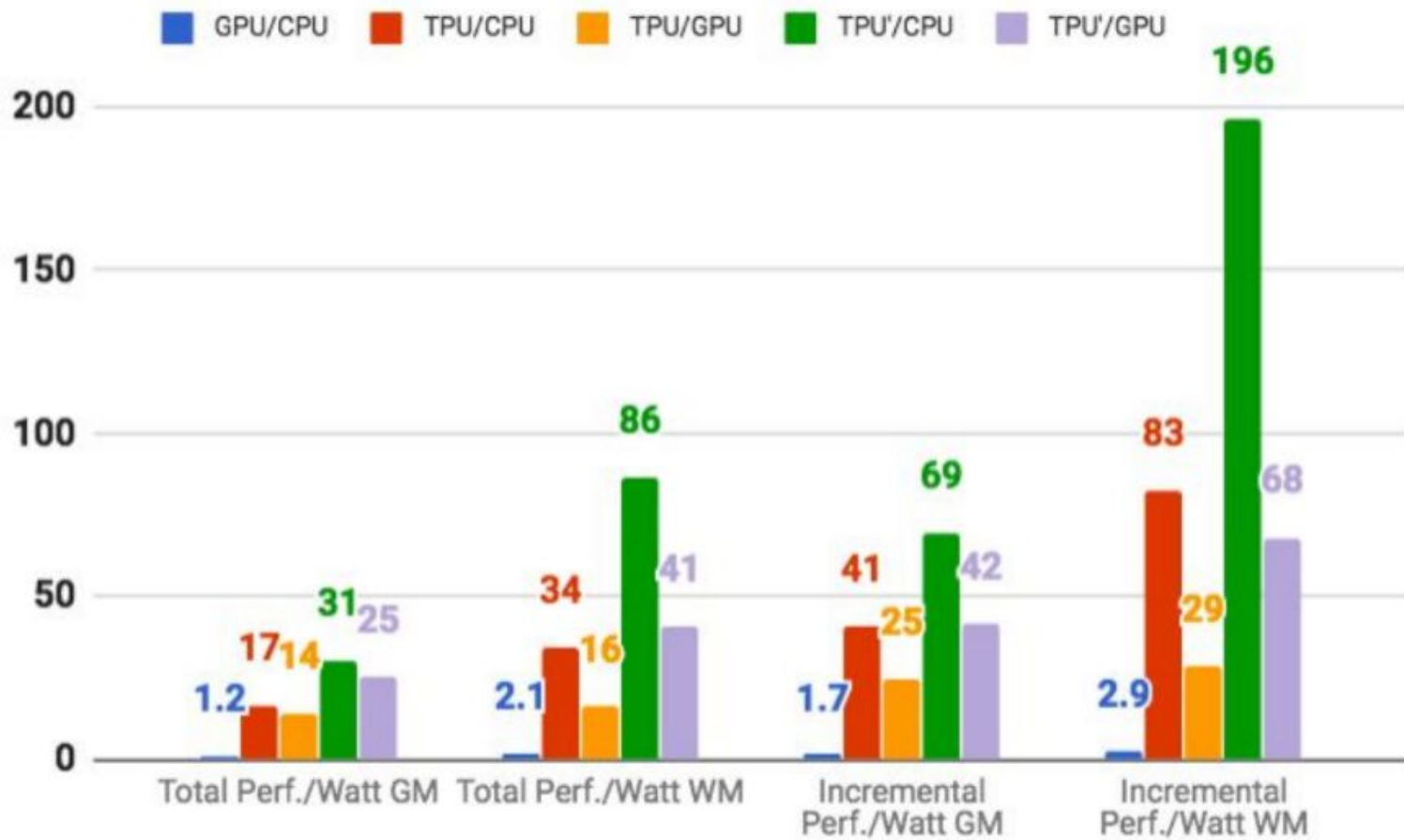
Google Specific Chip

Available to non-Googlers via Google Compute Cloud

Expect to other versions coming from NVidia, Intel, etc.



*Performance/Watt Relative to CPU or GPU*



# Google's TPU form makes me FOMO

Cloud TPU Beta Quota Request X Julien

Sécurisé <https://services.google.com/fb/forms/cloud-tpu-beta-request/>

How large is the largest labeled training dataset you currently wish to process with Cloud TPUs for supervised training? \*

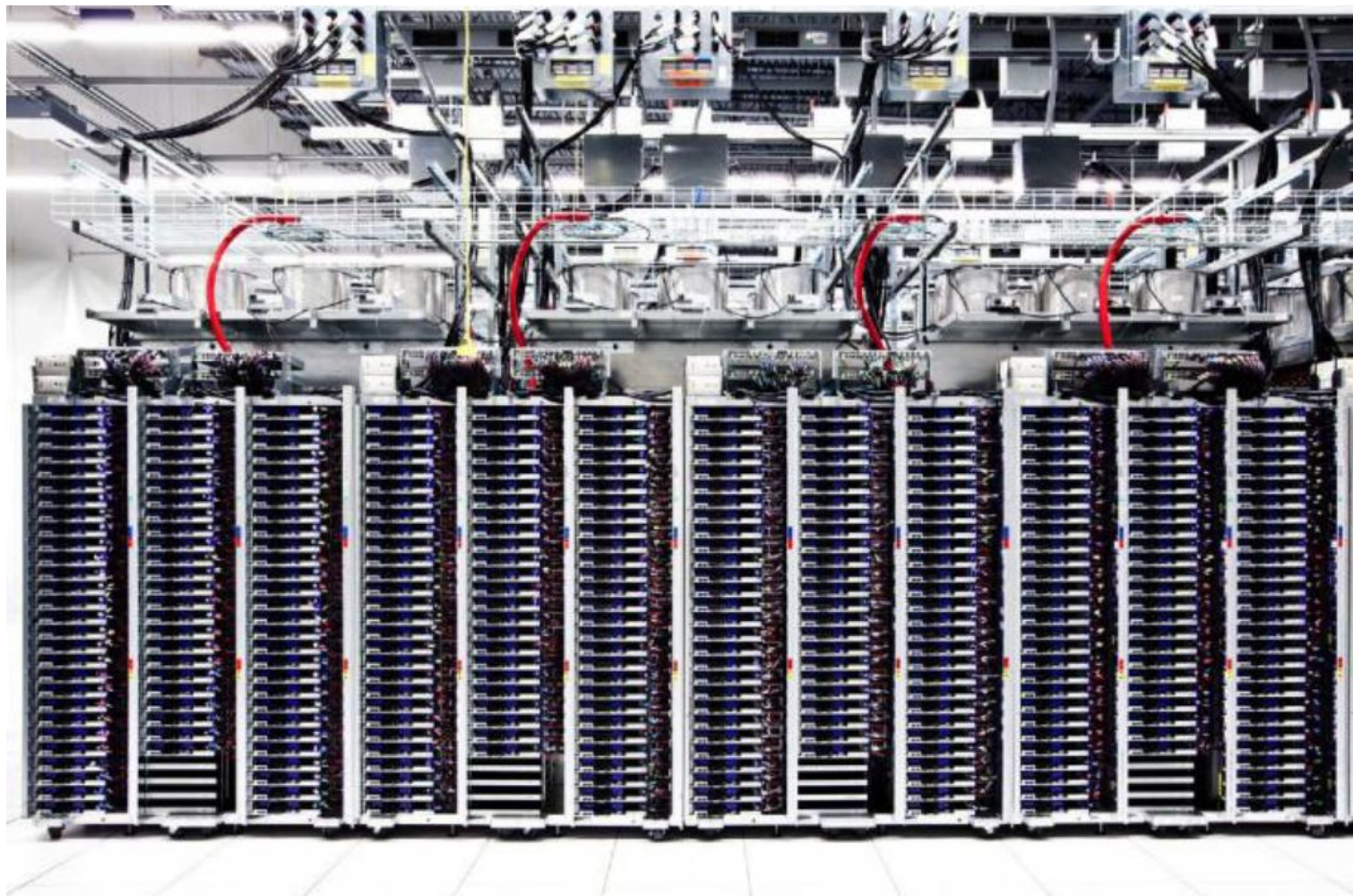
- ☒ < 10 GB
- ☐ 10 - 50 GB
- ☐ 50 - 150 GB
- ☐ 150 - 500 GB
- ☐ 500 GB - 2 TB
- ☐ 2 TB - 10 TB
- ☐ 10 TB - 100 TB
- ☐ 100 TB - 1 PB WTF?? 🤔
- ☐ > 1 PB
- ☐ No storage required; our datasets are generated live

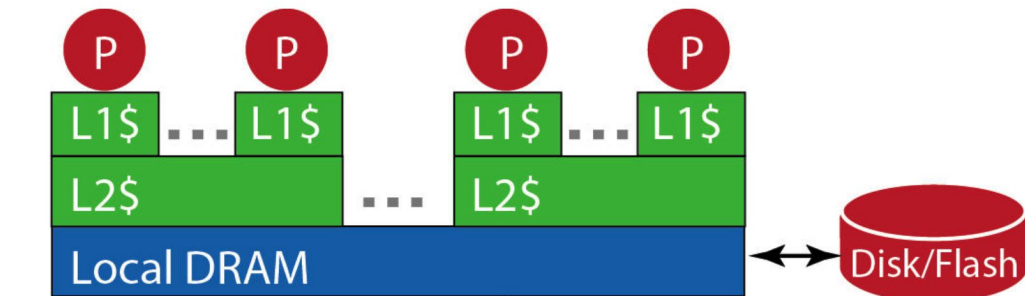
Please estimate the dataset size in GB when fully preprocessed and saved in TFRecord format on GCS, ready for training. If you generate data on the fly, please select "No storage required."

# Scale Out

The Datacenter as a Computer





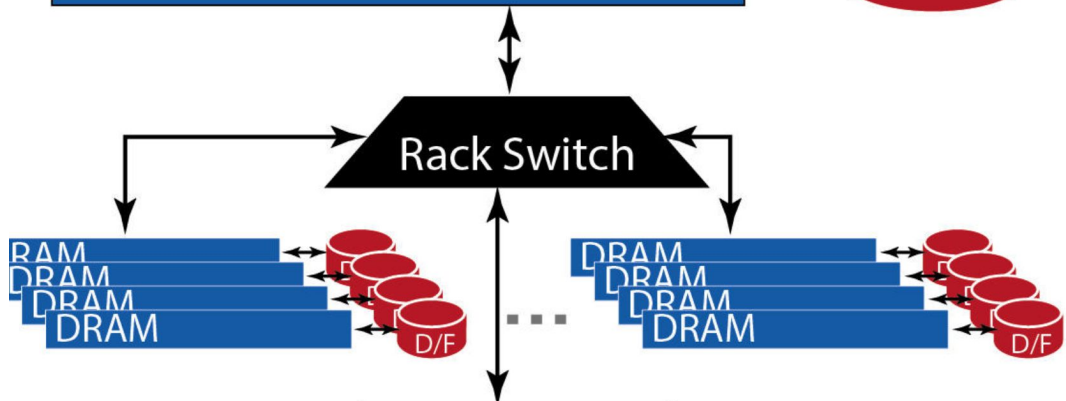


### One Server

DRAM: 16 GB, 100 ns, 20 GB/s

Disk: 2T B, 10 ms, 200 MB/s

Flash: 128 GB, 100 us, 1 GB/s



### Local Rack (80 servers)

DRAM: 1 TB, 300 us, 100 MB/s

Disk: 160 TB, 11 ms, 100 MB/s

Flash: 20 TB, 400 us, 100 MB/s



### Cluster (30 racks)

DRAM: 30 TB, 500 us, 10 MB/s

Disk: 4.80 PB, 12 ms, 10 MB/s

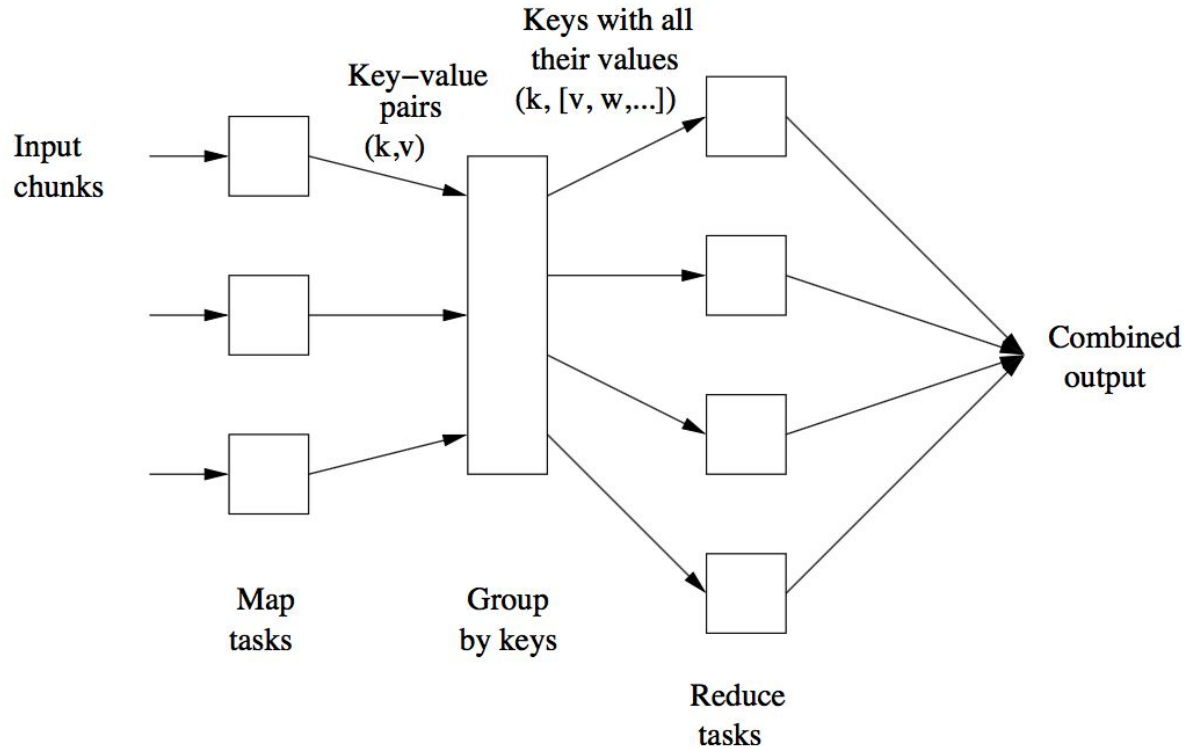
Flash: 600 TB, 600 us, 10 MB/s

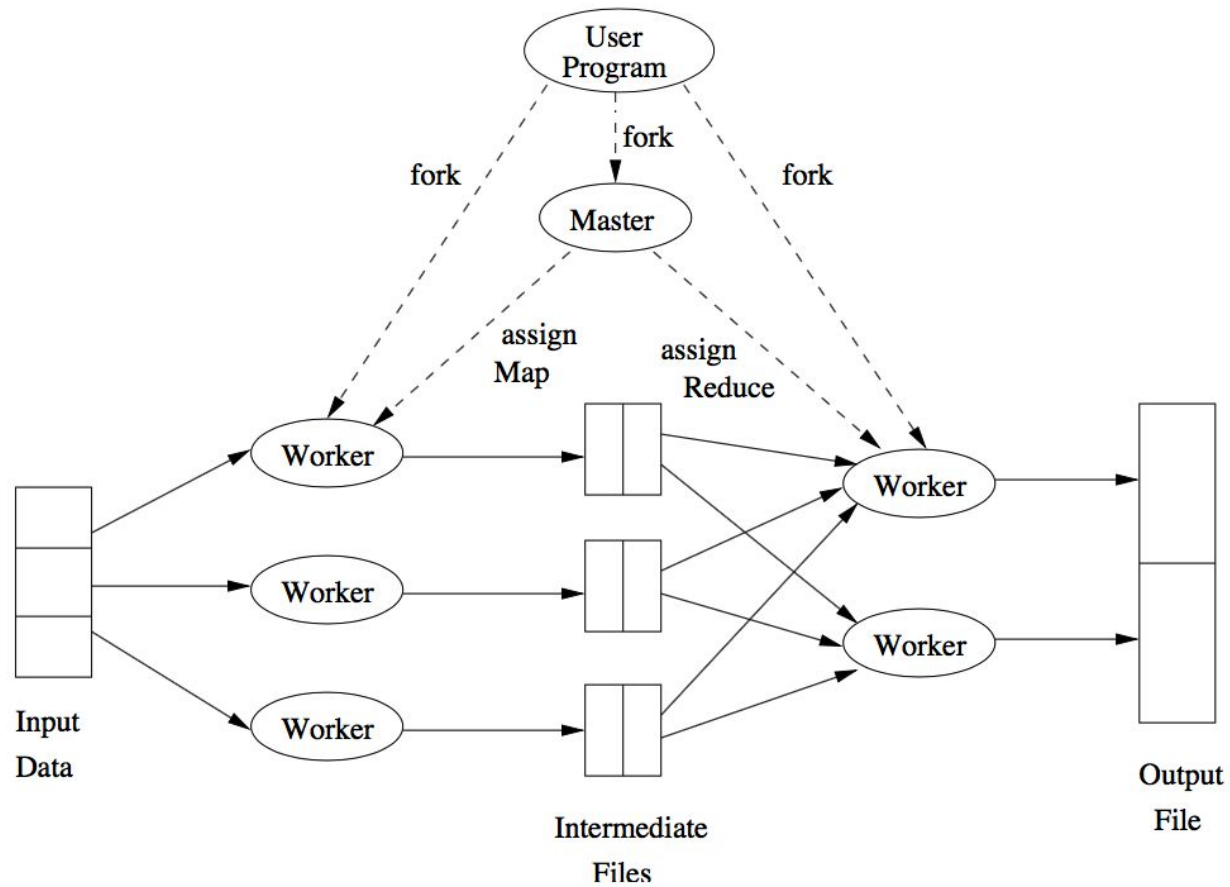


# Datacenter Issues

- Power: internal and external
- Cooling
- Networking, between racks, between datacenters
- Storage fabric
- Security: physical and data
- Server and Failure management

# How to Scale Across Many Nodes: Divide & Conquer





# Scalable Solutions

Matrix operations, including Machine Learning

Finding Similar Items

Mining data

Link analysis

Clustering Algorithms

Recommendation Systems

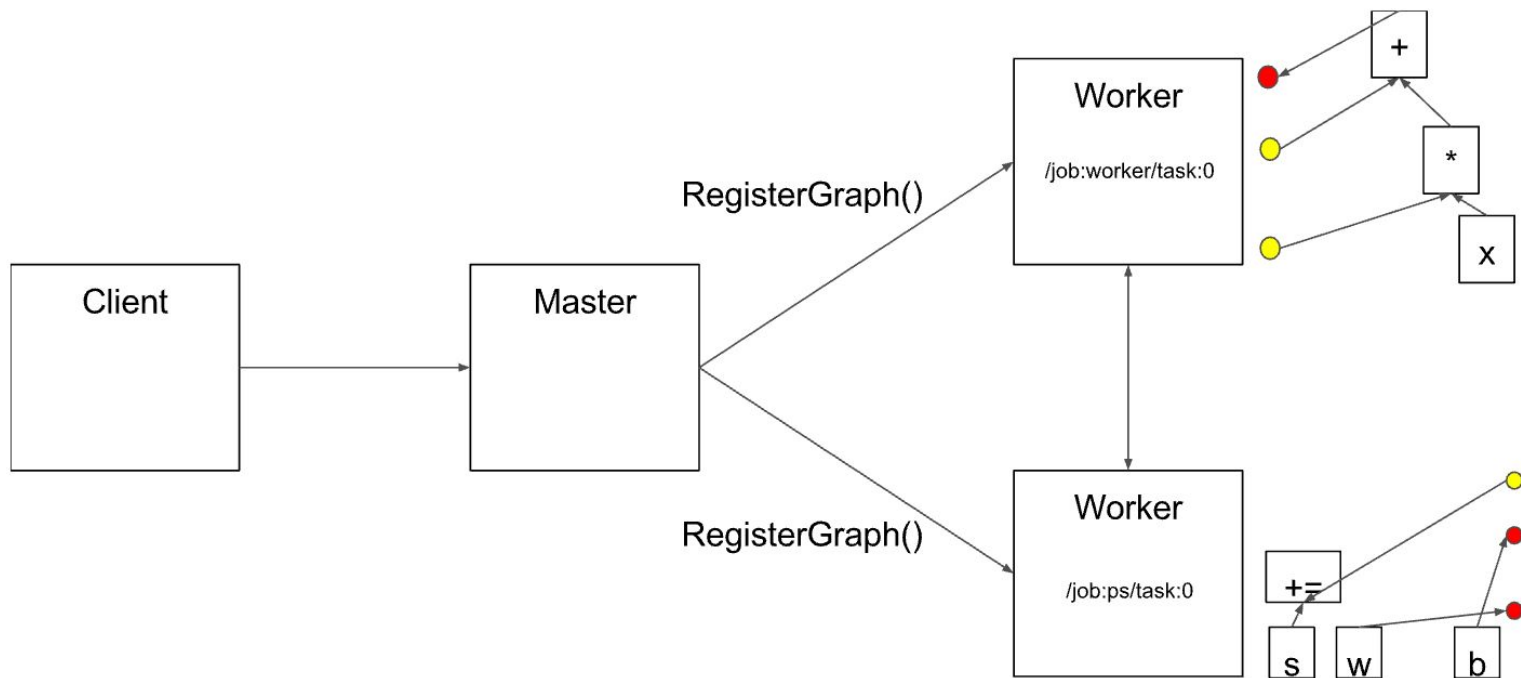
Dimensionality Reduction

# Issues:

- Amdahl's law still applies
- Copy overhead
  - Networking overhead
- Node failures
- Combine / Reconciliation step
- Synchronization
  - shared resources
  - wait time

$$S_{\text{latency}}(s) = \frac{1}{(1 - p) + \frac{p}{s}}$$

# Machine Learning at Datacenter Scale



<https://www.tensorflow.org/extend/architecture> See also **PyTorch Distributed**

# Future Directions

Abstract scalability behind services:

- Google BigTable / BigQuery, AutoML
  - Specialized APIs: NLP, Images, Translation
- AWS: S3 storage, Comprehend, DynamoDB

Datacenters have economies of scale and high capital costs

→ Market consolidation (AWS, Google, Microsoft Azure)

Proprietary APIs and services

→ Lock in

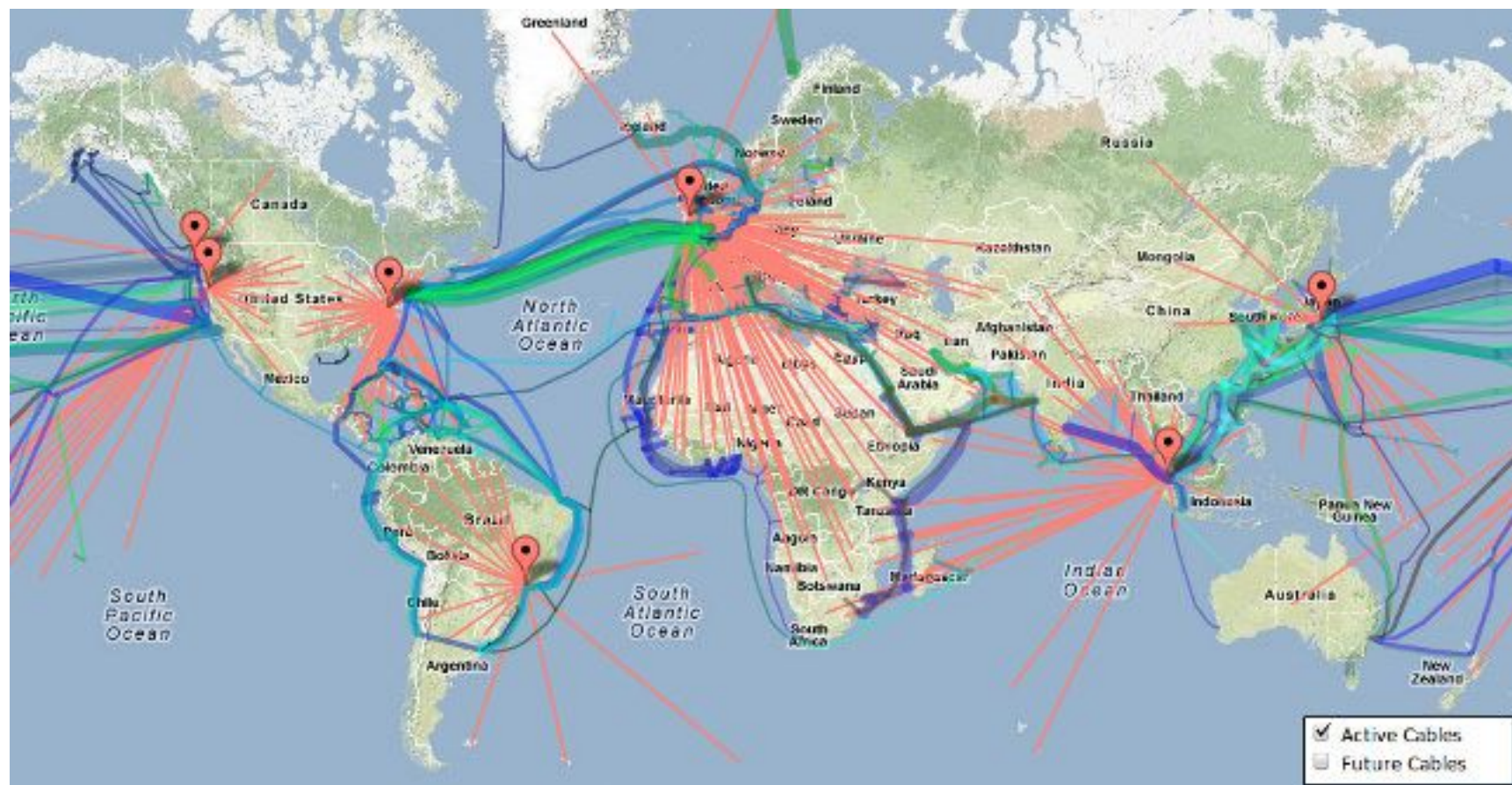
# Planet Scale

## Google worldwide network



Globally connected data centers, fiber network, peering relationships





<https://stackoverflow.com/questions/3907572/where-are-aws-data-center-locations-listed>

# References

Mining massive datasets: <http://www.mmds.org/>

<http://highscalability.com/all-time-favorites/>

<http://www.lecloud.net/tagged/scalability>

The Datacenter as a Computer book <https://research.google.com/pubs/pub41606.html>

TPUs <https://cloud.google.com/blog/big-data/2017/05/an-in-depth-look-at-googles-first-tensor-processing-unit-tpu>

# Quiz Time

