# 18-847F: Special Topics in Computer Systems
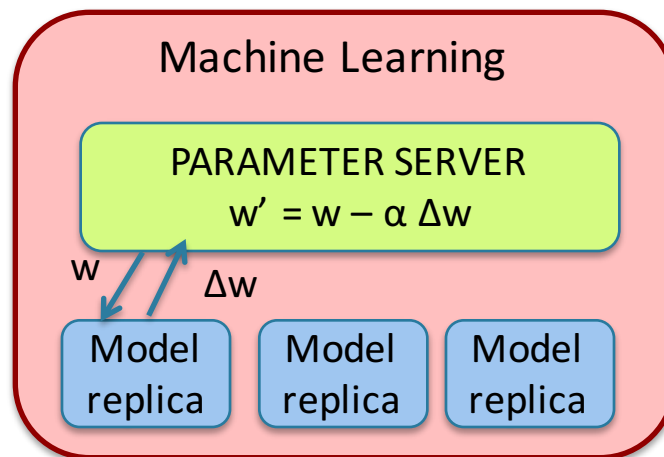
# Foundations of Cloud and Machine Learning Infrastructure
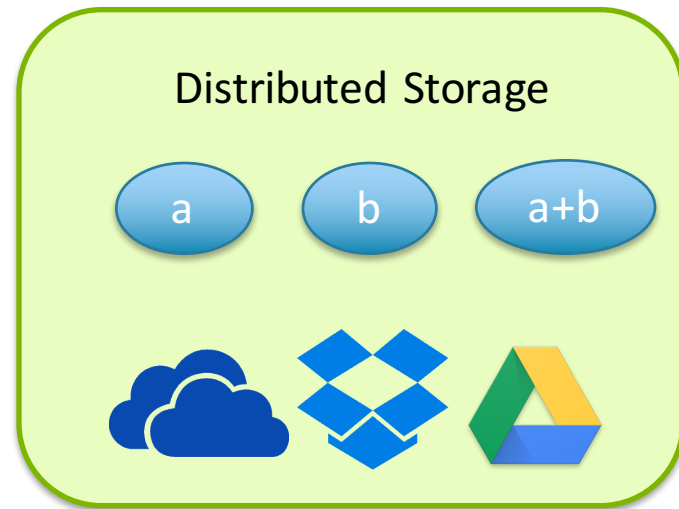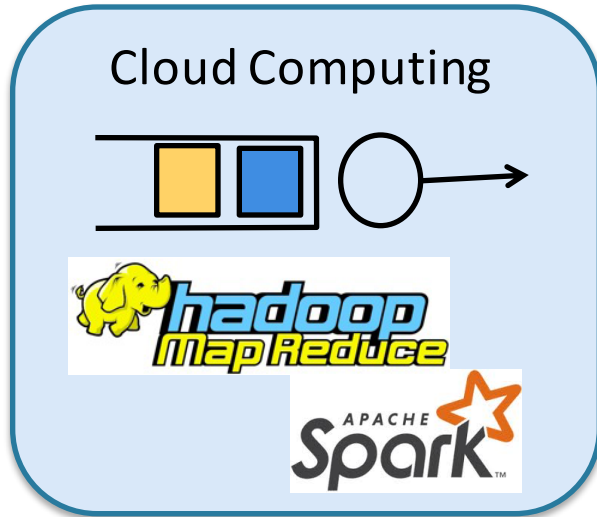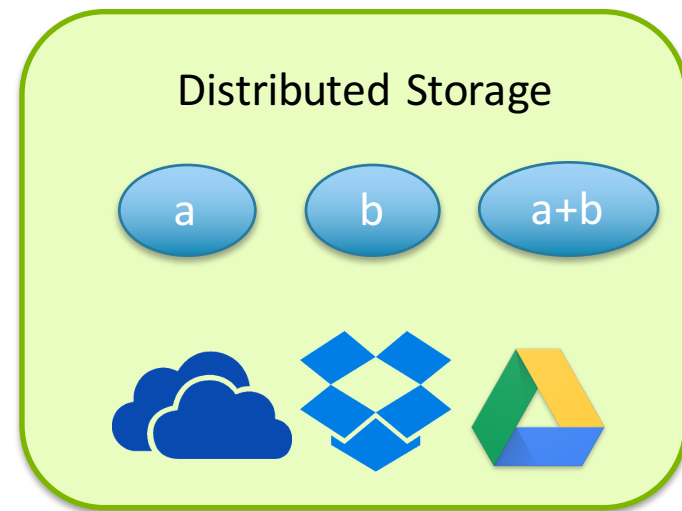
# Lecture 3: Overview of ML Infrastructure

# Foundations of Cloud and Machine Learning Infrastructure

# Topics Covered

## Cloud Computing

## Distributed Storage

a    b    a+b

## Machine Learning

PARAMETER SERVER
$w' = w - \alpha\, \Delta w$

w          $\Delta w$

Model replica    Model replica    Model replica

# Let us recap what we learnt..

**Cloud Computing**

hadoop Map Reduce

APACHE Spark

**Distributed Storage**

a    b    a+b

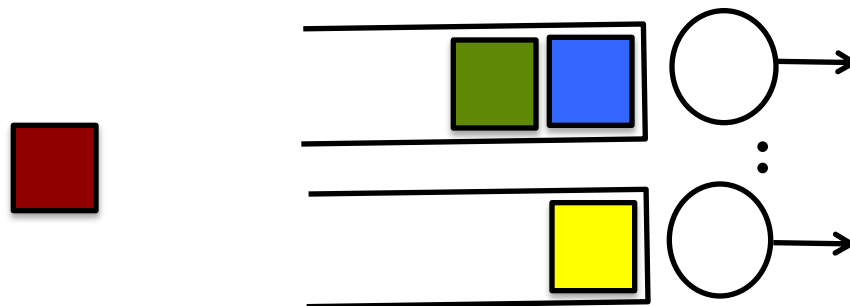# Let us recap what we learnt..

Cloud Computing

- MapReduce, Spark
- Scheduling in Parallel Computing
- Straggler Replication
- Task Replication in Queueing Systems

# Scheduling in Parallel Computing: 1990's

o **Bin-Packing**

    o Need job size estimates

o **Processor Sharing**, i.e. switching b/w threads for different jobs

    o Need processor speed estimates

o **Load-balancing**: Work stealing, Power-of-choice
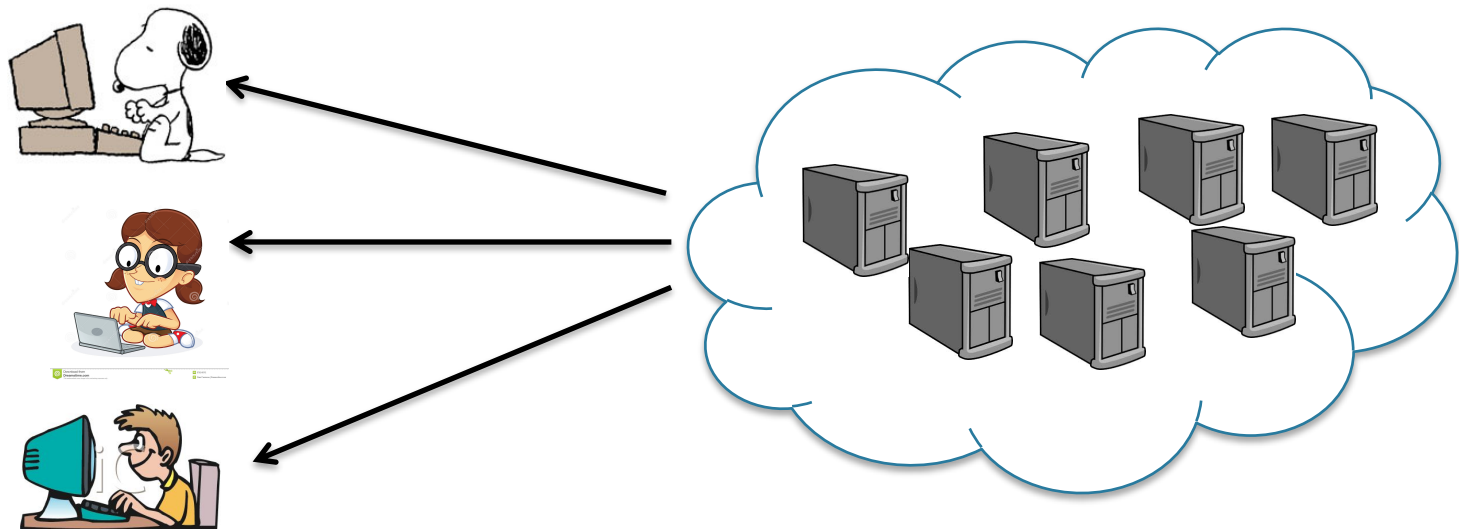
    o Need queue length estimates

# 2000's: The Cloud Computing Era

**KEY ISSUE:** Job sizes, server speeds & queue lengths are unpredictable

**REASON:** Large-scale resource sharing → Variability in service

- Virtualization, server outages etc.
- Norm and not an exception [Dean-Barroso 2013]
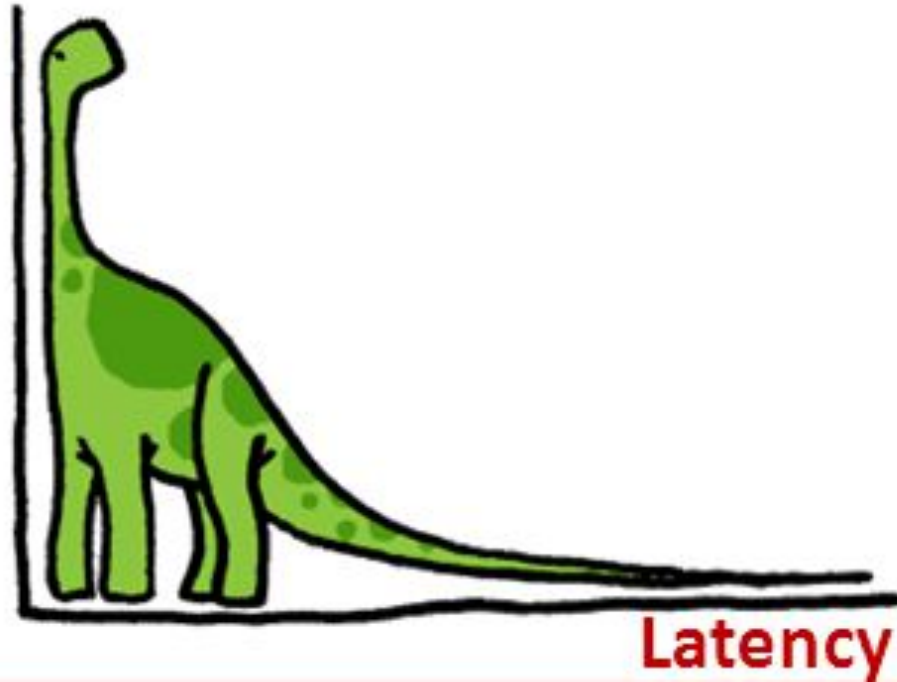
# Cloud Frameworks

MapReduce

Spark: In-memory

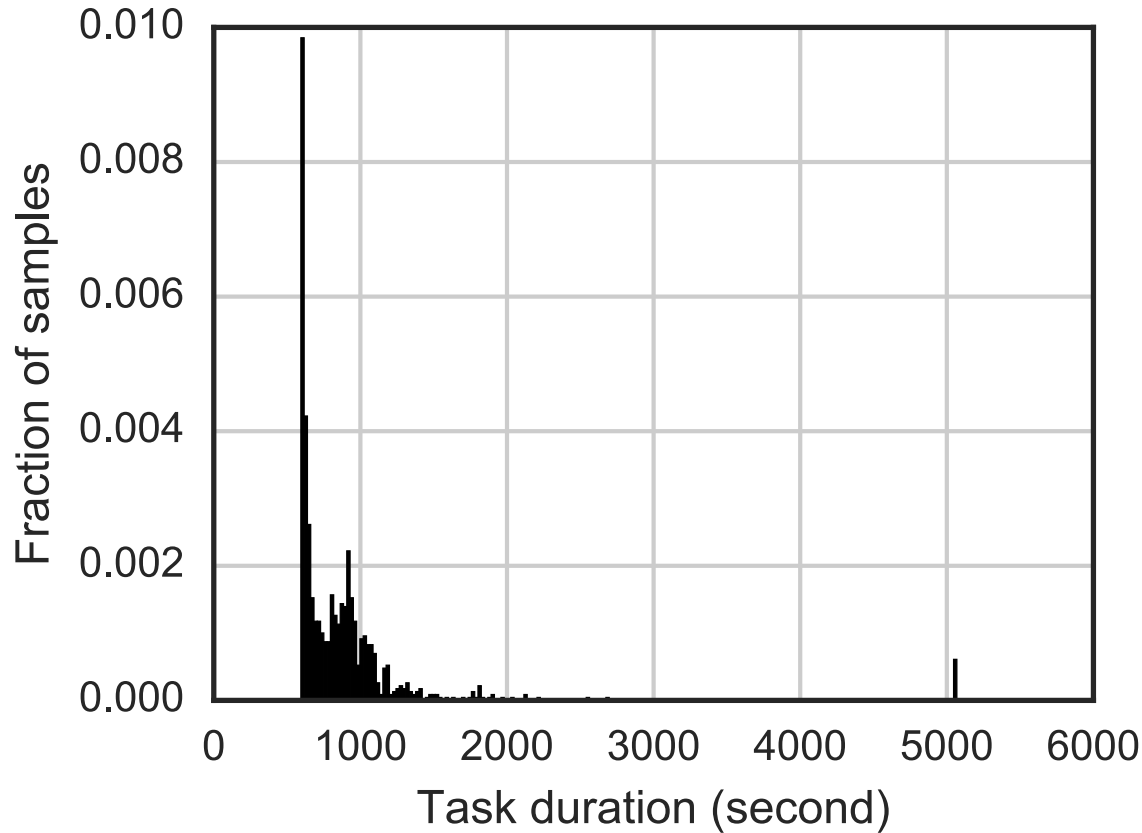Sparrow: Low-Latency Cluster Scheduling

Dolly: Attack of the Clones

# The Tale of Tails



Tail at Scale: 99%ile latency can be much higher than average
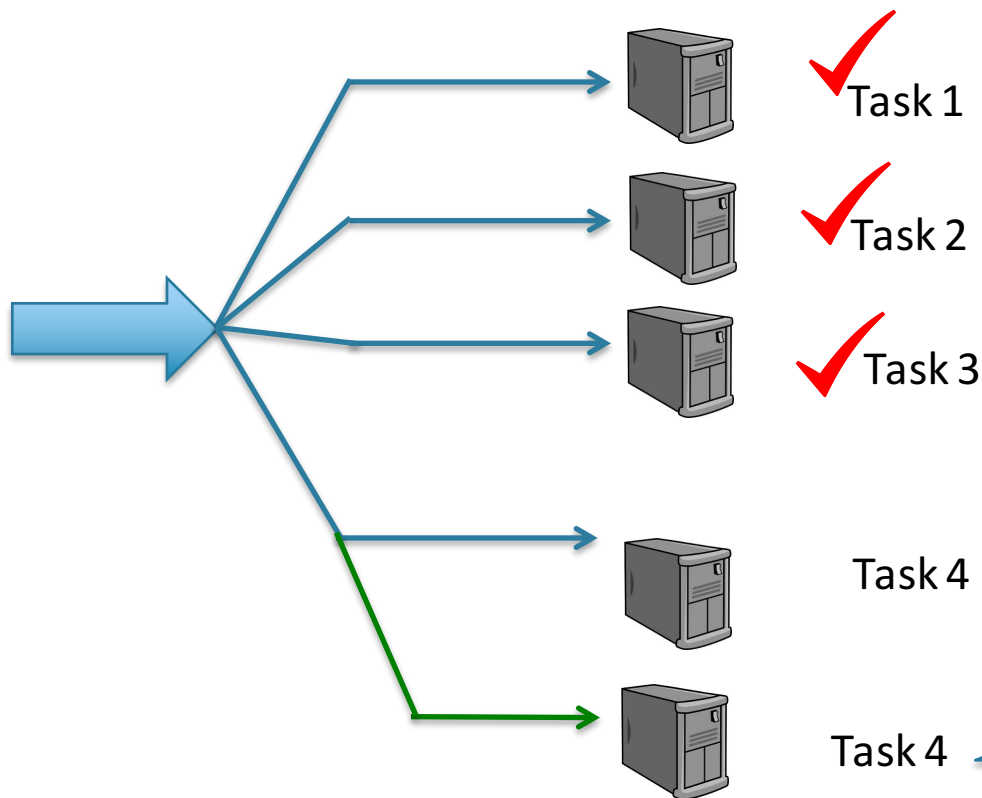
# The Tale of Tails



Tail at Scale: 99%ile latency much higher than average

# Straggler Replication

PROBLEM: Slowest tasks become a bottleneck

SOLUTION: Replicate the stragglers and wait for one copy



**PARAMETERS**

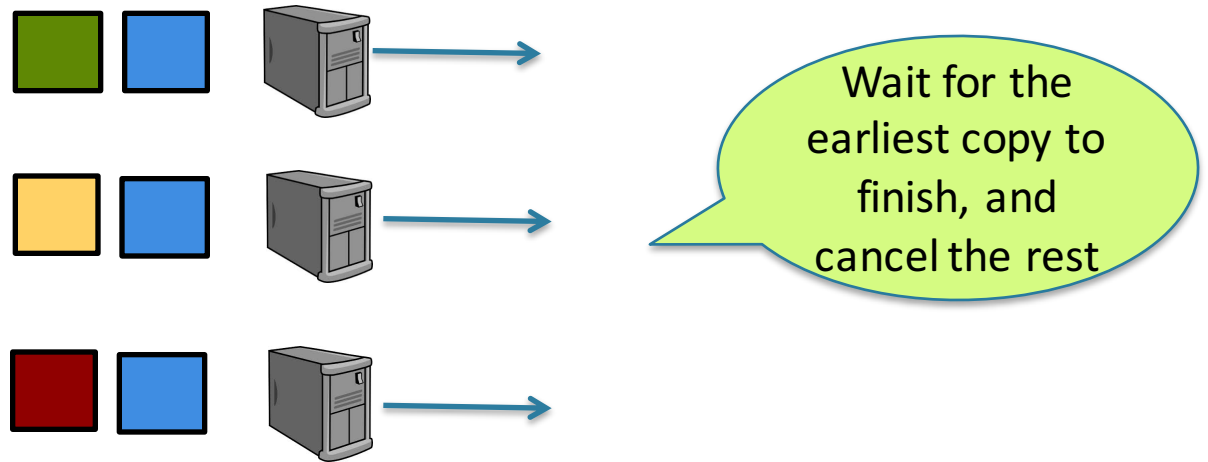p: Frac. of tasks replicated

r: # additional replicas

c: kill/keep original task

Eg. MapReduce, Apache Spark launch 1 replica, keep original copy

# Task Replication in Cloud Computing

**IDEA:** Assign task to multiple servers and wait for earliest copy

Wait for the earliest copy to finish, and cancel the rest

**COST**

o Additional computing time at servers

o Increased queuing delay for other tasks

# Design Questions

o How many replicas to launch?

o Which queues to join?

o When to issue and cancel the replicas?

# Cloud Spot Markets

o  Sell it on the spot market for a lower price!

On Demand price

Price

Spot Instance Price

Morning          Afternoon          Evening          Night

# Guest Lecture: Prof. Carlee Joe-Wong

o Bidding and pricing strategies for spot markets

# Let us recap what we learnt..

## Cloud Computing

## Distributed Storage

a    b    a+b

# Let us recap what we learnt..

o    RAID systems

o    Coding for locality/repair

o    Systems implementation of codes

o    Reducing latency in content

download



Distributed Storage

a    b    a+b

# RAID: Redundant Array of Independent Disks (1987)

o Levels RAID 0, RAID 1, … : design for different goals such as reliability, availability, capacity etc.



o One of the inventors, Garth Gibson is at CMU

# Erasure Coded Storage

o With an (n,k) MDS code, any k out of n chunks are sufficient
  o Facebook, Google, Microsoft use (14,10) or (7,4) codes
  o Currently used for cold data, increasing for hot data



Any k=2 out of n=3 are sufficient

# Codes for Efficient Repair

- Exact repair

- Functional repair

# Guest Lecture: Prof. Rashmi Vinayak

## Hitchhiker Codes and EC-Cache

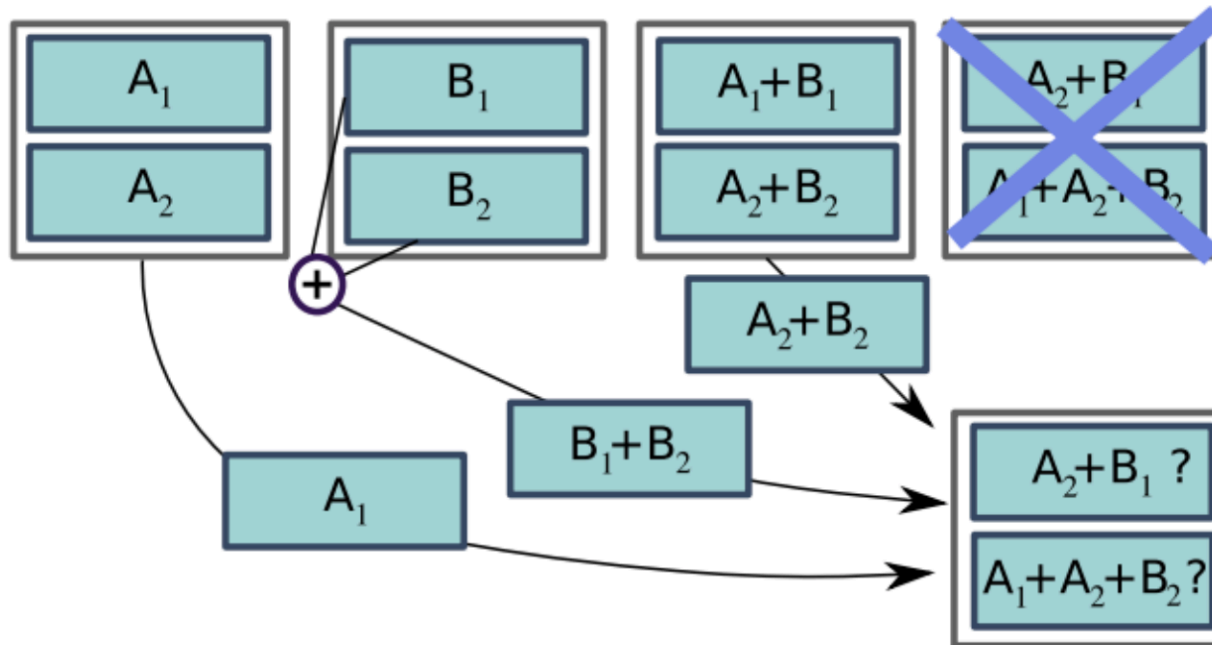| | An MDS Code | | Intermediate Step | | Piggybacked Code | |
|---|---|---|---|---|---|---|
| Node 1 | $a_1$ | $b_1$ | $a_1$ | $b_1$ | $a_1$ | $b_1$ |
| Node 2 | $a_2$ | $b_2$ | $a_2$ | $b_2$ | $a_2$ | $b_2$ |
| Node 3 | $a_3$ | $b_3$ | $a_3$ | $b_3$ | $a_3$ | $b_3$ |
| Node 4 | $a_4$ | $b_4$ | $a_4$ | $b_4$ | $a_4$ | $b_4$ |
| Node 5 | $\sum_{i=1}^{4} a_i$ | $\sum_{i=1}^{4} b_i$ | $\sum_{i=1}^{4} a_i$ | $\sum_{i=1}^{4} b_i$ | $\sum_{i=1}^{4} a_i$ | $\sum_{i=1}^{4} b_i$ |
| Node 6 | $\sum_{i=1}^{4} i a_i$ | $\sum_{i=1}^{4} i b_i$ | $\sum_{i=1}^{4} i a_i$ | $\sum_{i=1}^{4} i b_i + \sum_{i=1}^{2} i a_i$ | $\sum_{i=3}^{4} i a_i - \sum_{i=1}^{4} i b_i$ | $\sum_{i=1}^{4} i b_i + \sum_{i=1}^{2} i a_i$ |
| | (a) | | (b) | | (c) | |

Needs 8 symbols
to repair

Needs 6 symbols
to repair

# The (n,k) fork-join model
## [GJ-Liu-Soljanin 2012,14]

○ Request all n chunks, wait for any k to be downloaded

○ Each chunk takes service time $X \sim F_X$

(3,2) fork-join

$\lambda$
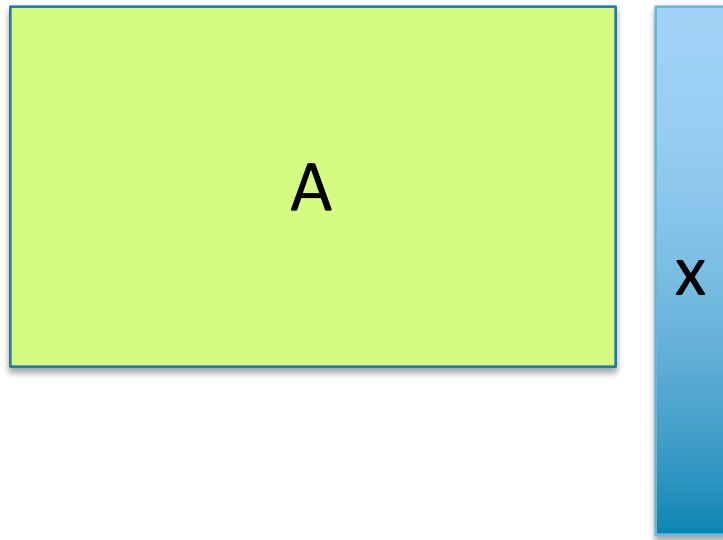
Download
requests

Wait for any 2
out of 3 chunks

k = 1: Replicated Case

k = n: Fork-join system actively studied in 90's

# Coded Computing and ML

o So far: Coding for storage

o Codes can also speed up computing and machine learning

o Example: Matrix-Vector Multiplication

A

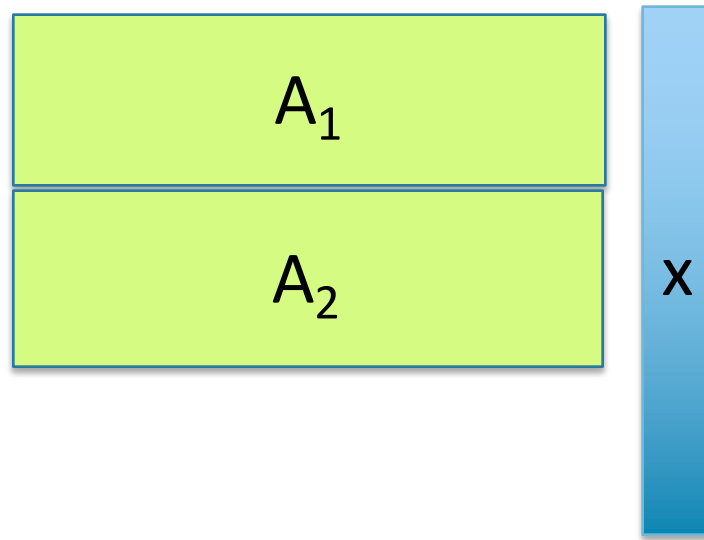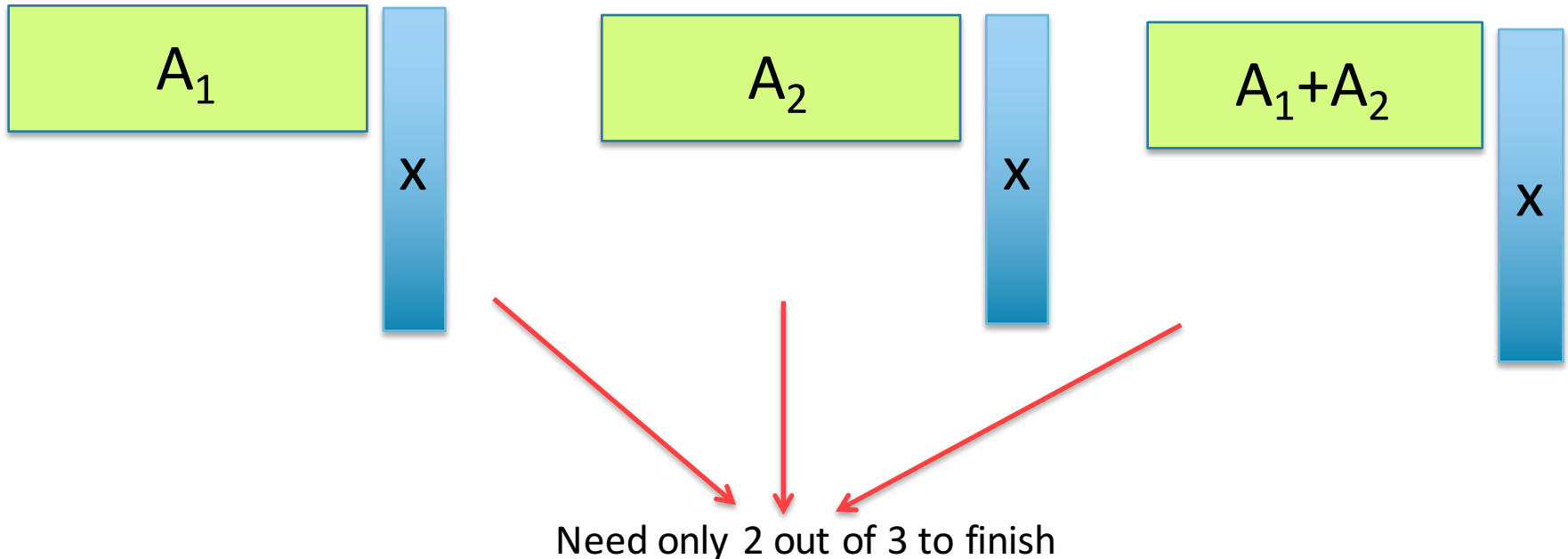x

# Coded Computing and ML

o So far: coding for storage

o Codes can also speed up computing and machine learning!

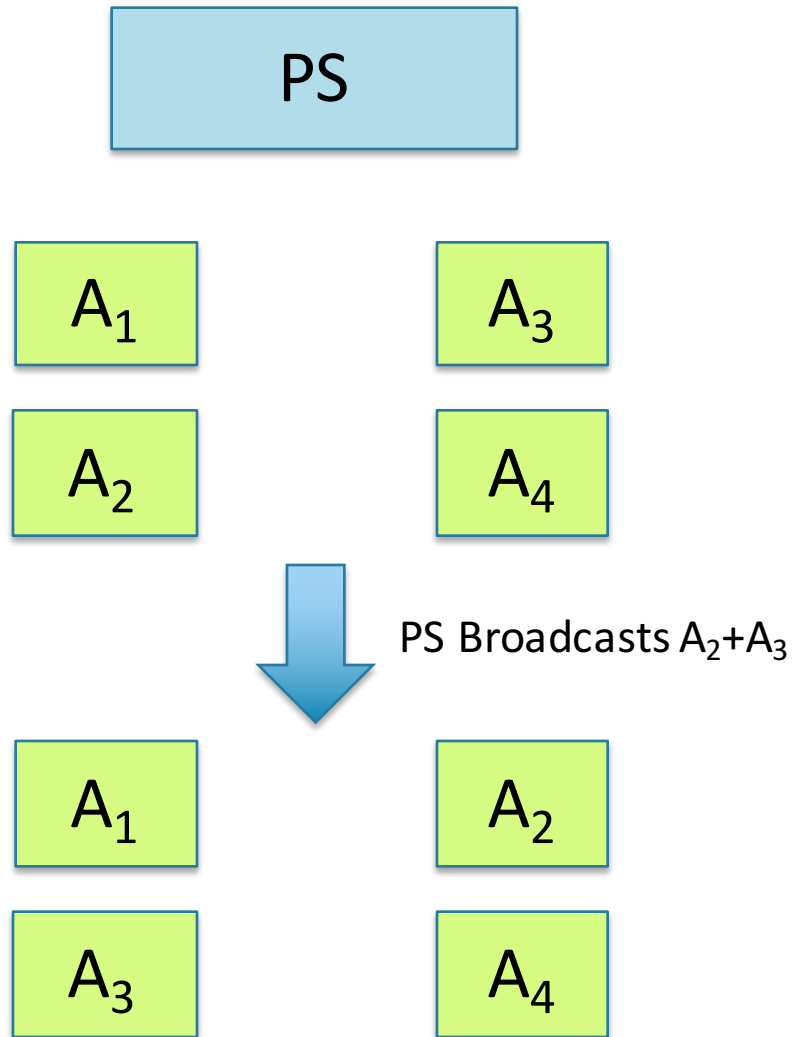o Example: Matrix-Vector Multiplication

# Coded Computing and ML

o So far: coding for storage

o Codes can also speed up computing and machine learning!
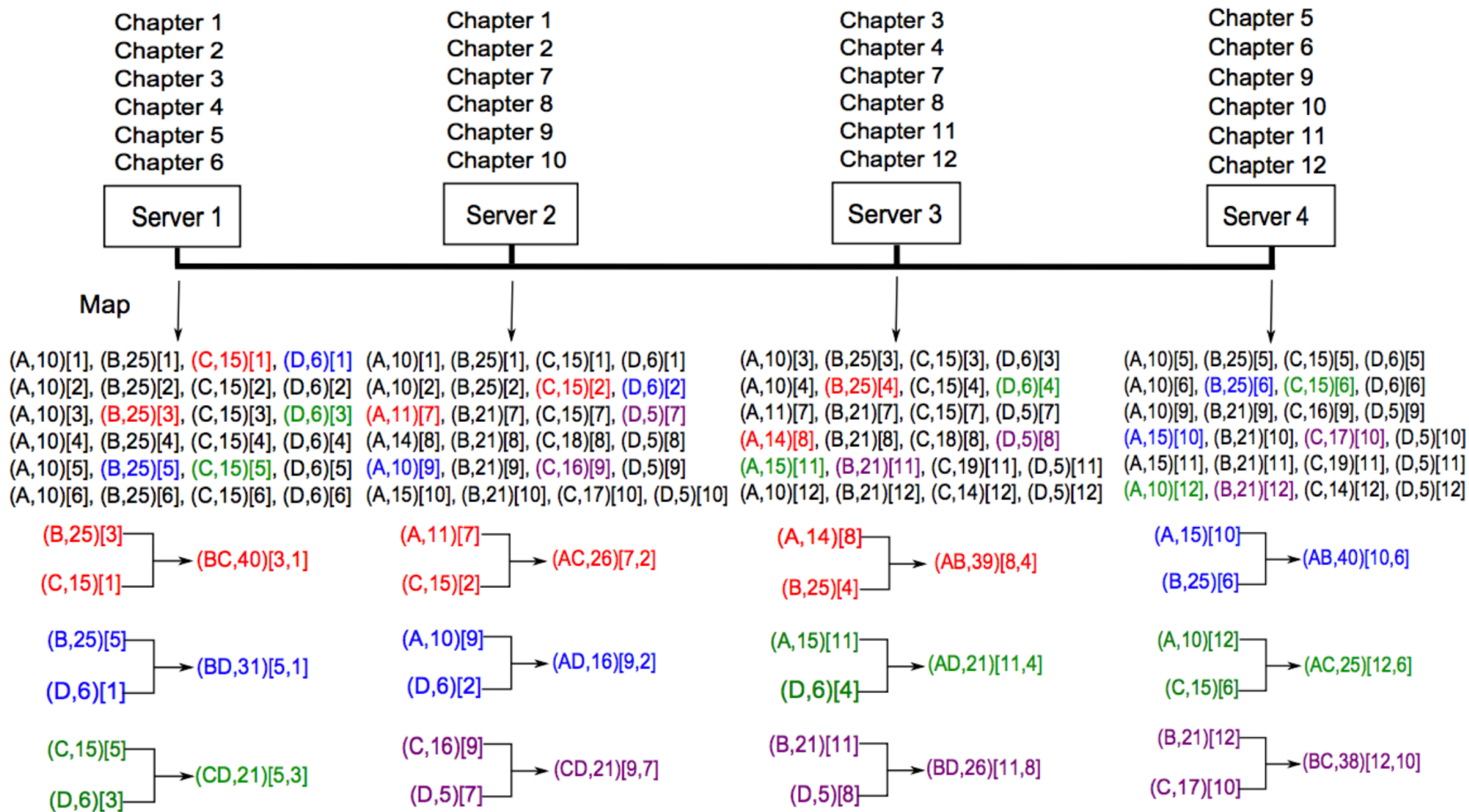
o Example: Matrix-Vector Multiplication

$A_1$    X        $A_2$    X        $A_1+A_2$    X

Need only 2 out of 3 to finish

# Coded Data Shuffling



PS

$A_1$   $A_3$

$A_2$   $A_4$

PS Broadcasts $A_2+A_3$

$A_1$   $A_2$

$A_3$   $A_4$

26

# Coded MapReduce



| Server 1 | Server 2 | Server 3 | Server 4 |
|---|---|---|---|
| Chapter 1 | Chapter 1 | Chapter 3 | Chapter 5 |
| Chapter 2 | Chapter 2 | Chapter 4 | Chapter 6 |
| Chapter 3 | Chapter 7 | Chapter 7 | Chapter 9 |
| Chapter 4 | Chapter 8 | Chapter 8 | Chapter 10 |
| Chapter 5 | Chapter 9 | Chapter 11 | Chapter 11 |
| Chapter 6 | Chapter 10 | Chapter 12 | Chapter 12 |

Map

Server 1:
(A,10)[1], (B,25)[1], (C,15)[1], (D,6)[1]
(A,10)[2], (B,25)[2], (C,15)[2], (D,6)[2]
(A,10)[3], (B,25)[3], (C,15)[3], (D,6)[3]
(A,10)[4], (B,25)[4], (C,15)[4], (D,6)[4]
(A,10)[5], (B,25)[5], (C,15)[5], (D,6)[5]
(A,10)[6], (B,25)[6], (C,15)[6], (D,6)[6]

Server 2:
(A,10)[1], (B,25)[1], (C,15)[1], (D,6)[1]
(A,10)[2], (B,25)[2], (C,15)[2], (D,6)[2]
(A,11)[7], (B,21)[7], (C,15)[7], (D,5)[7]
(A,14)[8], (B,21)[8], (C,18)[8], (D,5)[8]
(A,10)[9], (B,21)[9], (C,16)[9], (D,5)[9]
(A,15)[10], (B,21)[10], (C,17)[10], (D,5)[10]

Server 3:
(A,10)[3], (B,25)[3], (C,15)[3], (D,6)[3]
(A,10)[4], (B,25)[4], (C,15)[4], (D,6)[4]
(A,11)[7], (B,21)[7], (C,15)[7], (D,5)[7]
(A,14)[8], (B,21)[8], (C,18)[8], (D,5)[8]
(A,15)[11], (B,21)[11], (C,19)[11], (D,5)[11]
(A,10)[12], (B,21)[12], (C,14)[12], (D,5)[12]

Server 4:
(A,10)[5], (B,25)[5], (C,15)[5], (D,6)[5]
(A,10)[6], (B,25)[6], (C,15)[6], (D,6)[6]
(A,10)[9], (B,21)[9], (C,16)[9], (D,5)[9]
(A,15)[10], (B,21)[10], (C,17)[10], (D,5)[10]
(A,15)[11], (B,21)[11], (C,19)[11], (D,5)[11]
(A,10)[12], (B,21)[12], (C,14)[12], (D,5)[12]

Server 1:
(B,25)[3] + (C,15)[1] → (BC,40)[3,1]
(B,25)[5] + (D,6)[1] → (BD,31)[5,1]
(C,15)[5] + (D,6)[3] → (CD,21)[5,3]

Server 2:
(A,11)[7] + (C,15)[2] → (AC,26)[7,2]
(A,10)[9] + (D,6)[2] → (AD,16)[9,2]
(C,16)[9] + (D,5)[7] → (CD,21)[9,7]

Server 3:
(A,14)[8] + (B,25)[4] → (AB,39)[8,4]
(A,15)[11] + (D,6)[4] → (AD,21)[11,4]
(B,21)[11] + (D,5)[8] → (BD,26)[11,8]

Server 4:
(A,15)[10] + (B,25)[6] → (AB,40)[10,6]
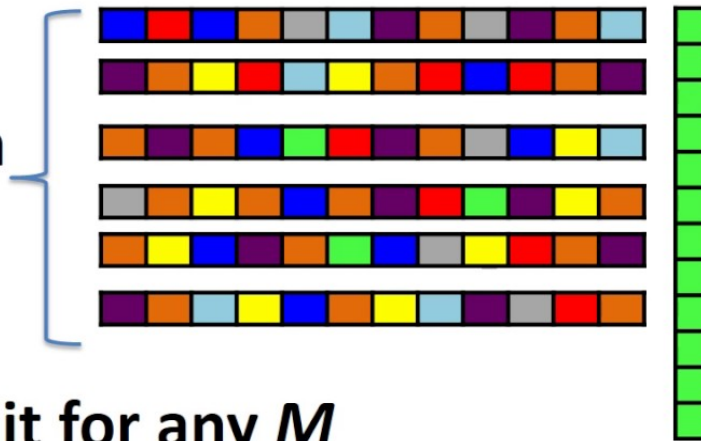(A,10)[12] + (C,15)[6] → (AC,25)[12,6]
(B,21)[12] + (C,17)[10] → (BC,38)[12,10]

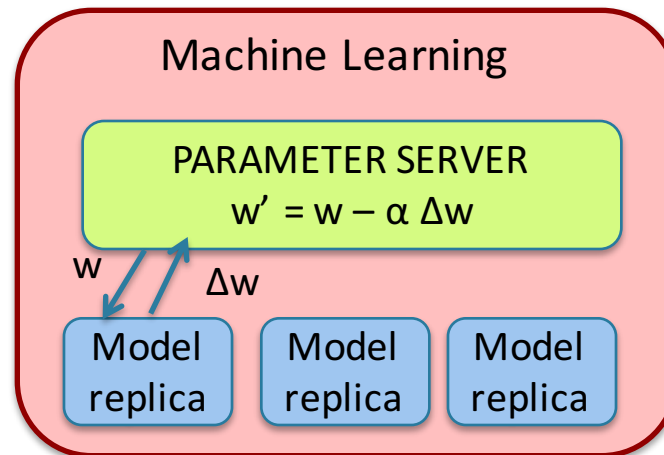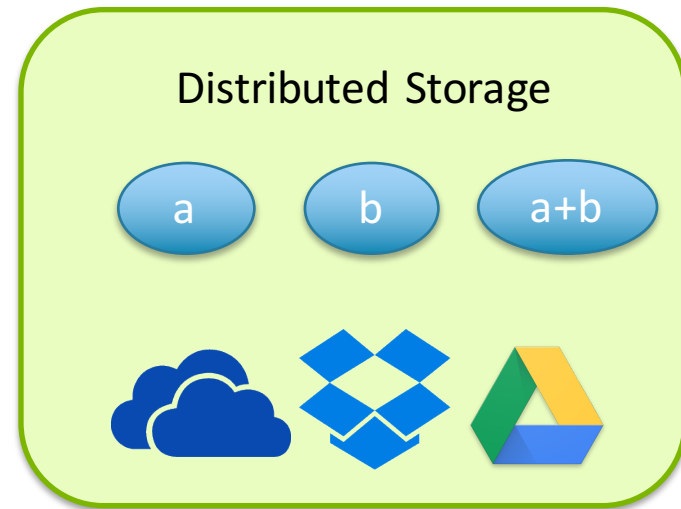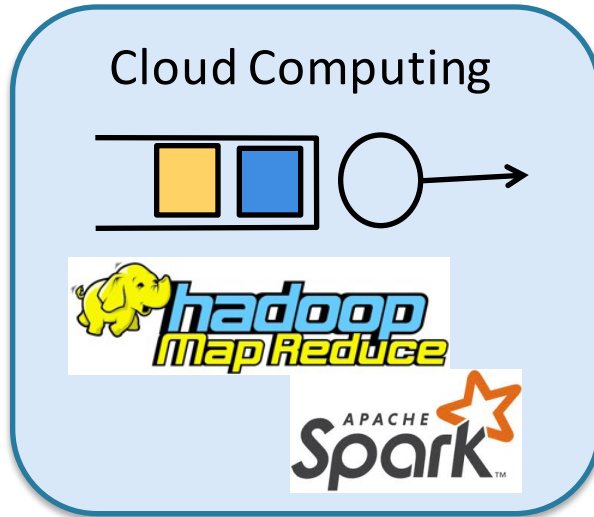# Guest Lecture: Sanghamitra Dutta

Short-dot codes



**(P,M) MDS Code**

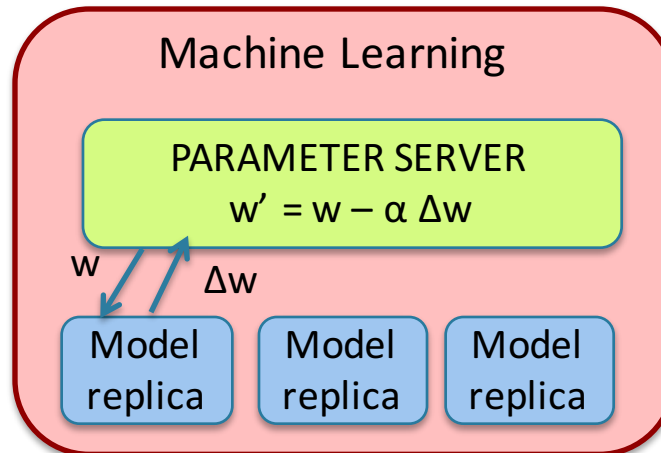*P* dense dot-products of length *N* in *P* parallel processors

Wait for any *M* computations to finish
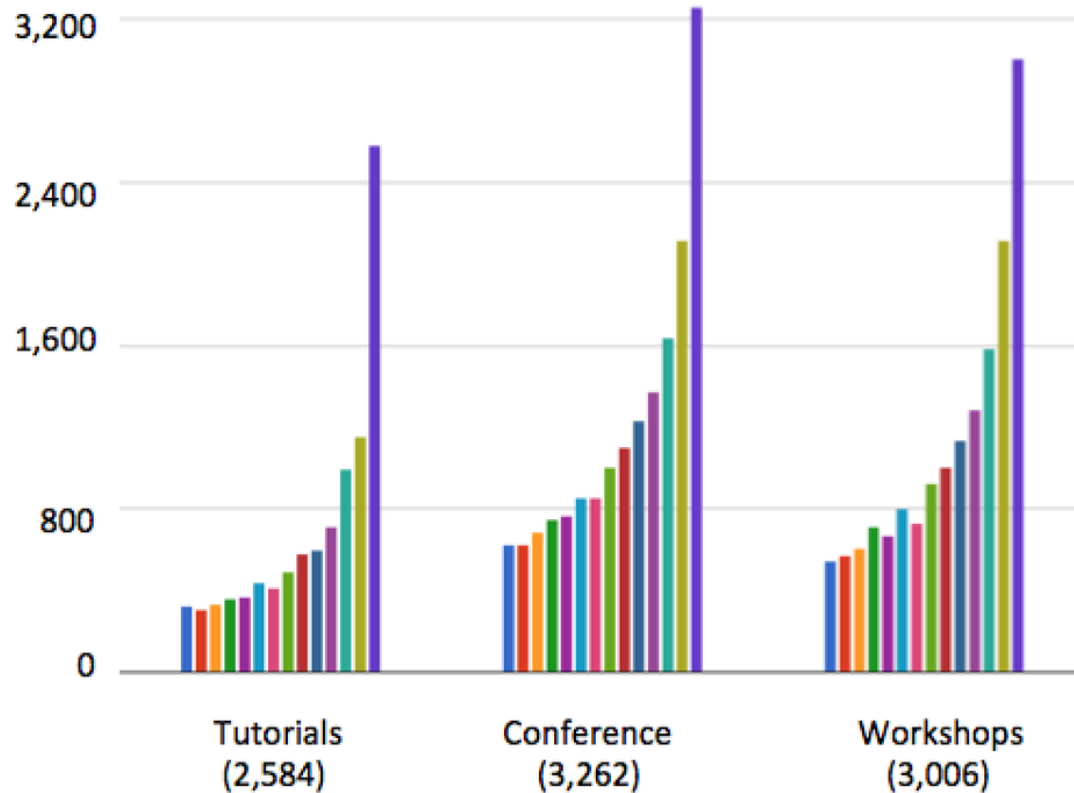
# Last Module: Machine Learning

**Cloud Computing**

**Distributed Storage**

a    b    a+b

**Machine Learning**

PARAMETER SERVER
$$w' = w - \alpha \, \Delta w$$

w     $\Delta w$

Model replica    Model replica    Model replica

# Last Module: Machine Learning

o    SGD Methods, Convergence

o    DistBelief, Alexnet

o    Synchronous, Asynchronous SGD

o    GANs, Reinforcement Learning

## Machine Learning

PARAMETER SERVER
$w' = w - \alpha\, \Delta w$

$w$     $\Delta w$

Model replica    Model replica    Model replica

# The unprecedented ML boom

# The Origins: 1950



Alan Turing

# Neural Networks: Perceptron 1957



Frank Rosenblatt
(1928-1971)

Original Perceptron

(From Perceptrons by M. L Minsky and S. Papert, 1969, Cambridge, MA: MIT Press. Copyright 1969 by MIT Press.)

Simplified model:

23

# Back-propagation Algorithm

Geoff Hinton (U. Toronto, Google)

Journal Home
Current Issue
AOP
Archive

THIS ARTICLE ▾

Download PDF
References

Export citation
Export references

Send to a friend

More articles like this

Table of Contents
< Previous | Next >

### letters to nature

## Learning representations by back-propagating errors

DAVID E. RUMELHART[*], GEOFFREY E. HINTON[†] & RONALD J. WILLIAMS[*]

[*]Institute for Cognitive Science, C-015, University of California, San Diego, La Jolla, California 92093, USA
[†]Department of Computer Science, Carnegie-Mellon University, Pittsburgh, Philadelphia 15213, USA
[†]To whom correspondence should be addressed.

We describe a new learning procedure, back-propagation, for networks of neurone-like units. The procedure repe    ights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the net and the desired output vector. As a result of the weight adjustments, internal 'hidden' units which are not part of the input or output come to represent important features of the task domain, and the regularities in the task are captured by the interactions of these units. The ability to create useful new features distinguishes back-propagation from earlier, simpler methods such as the perceptron-convergence procedure[1].
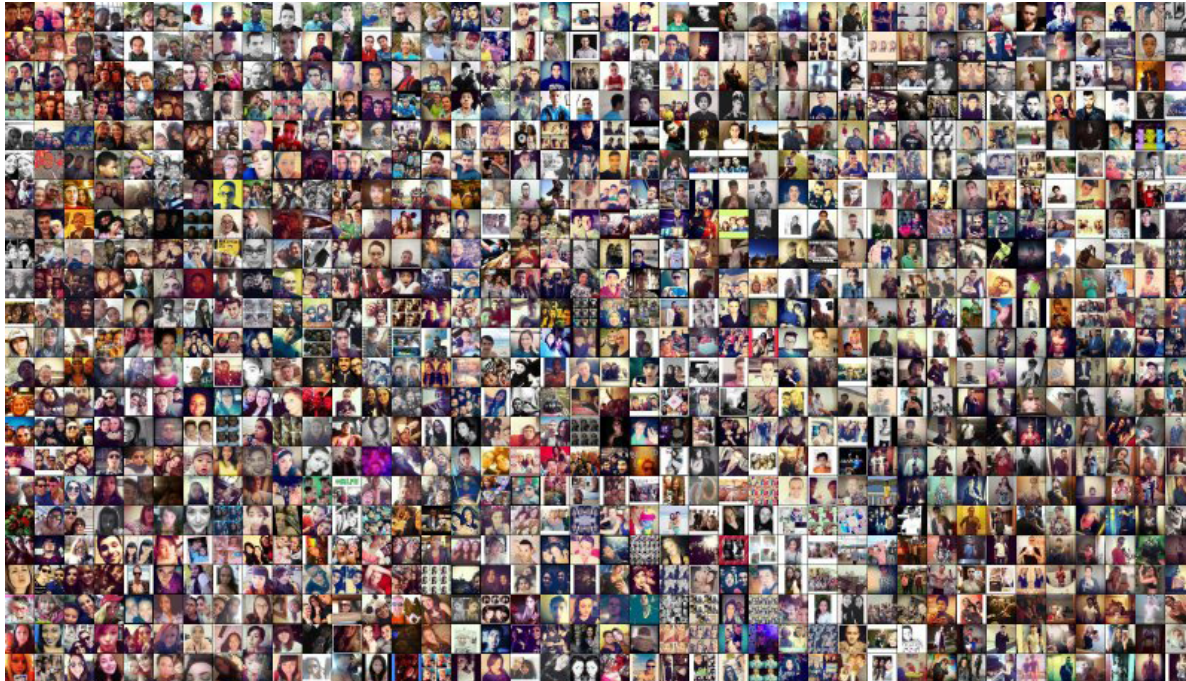
### References

1. Rosenblatt, F. *Principles of Neurodynamics* (Spartan, Washington, DC, 1961).
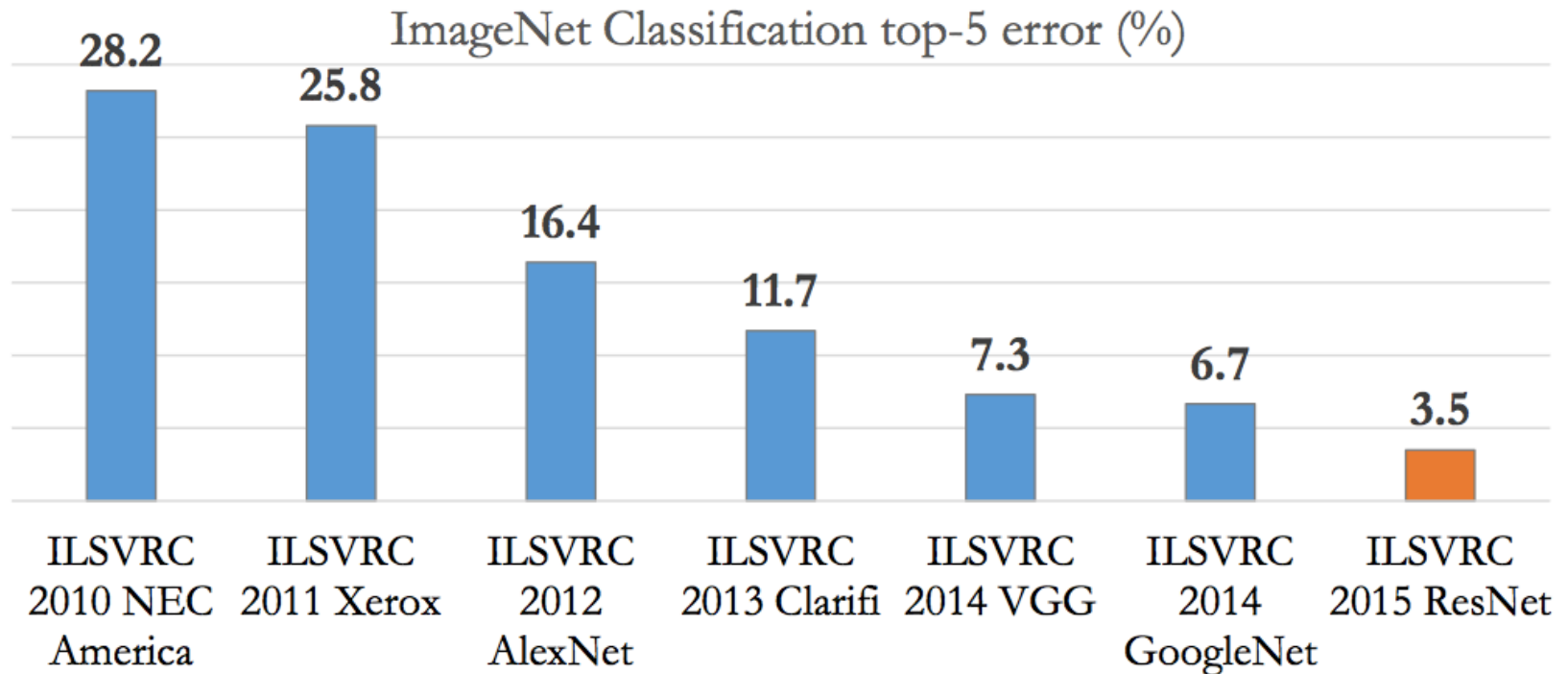
# MNIST (LeCun et al 1998)

# ImageNet and ILSVRC (2012)
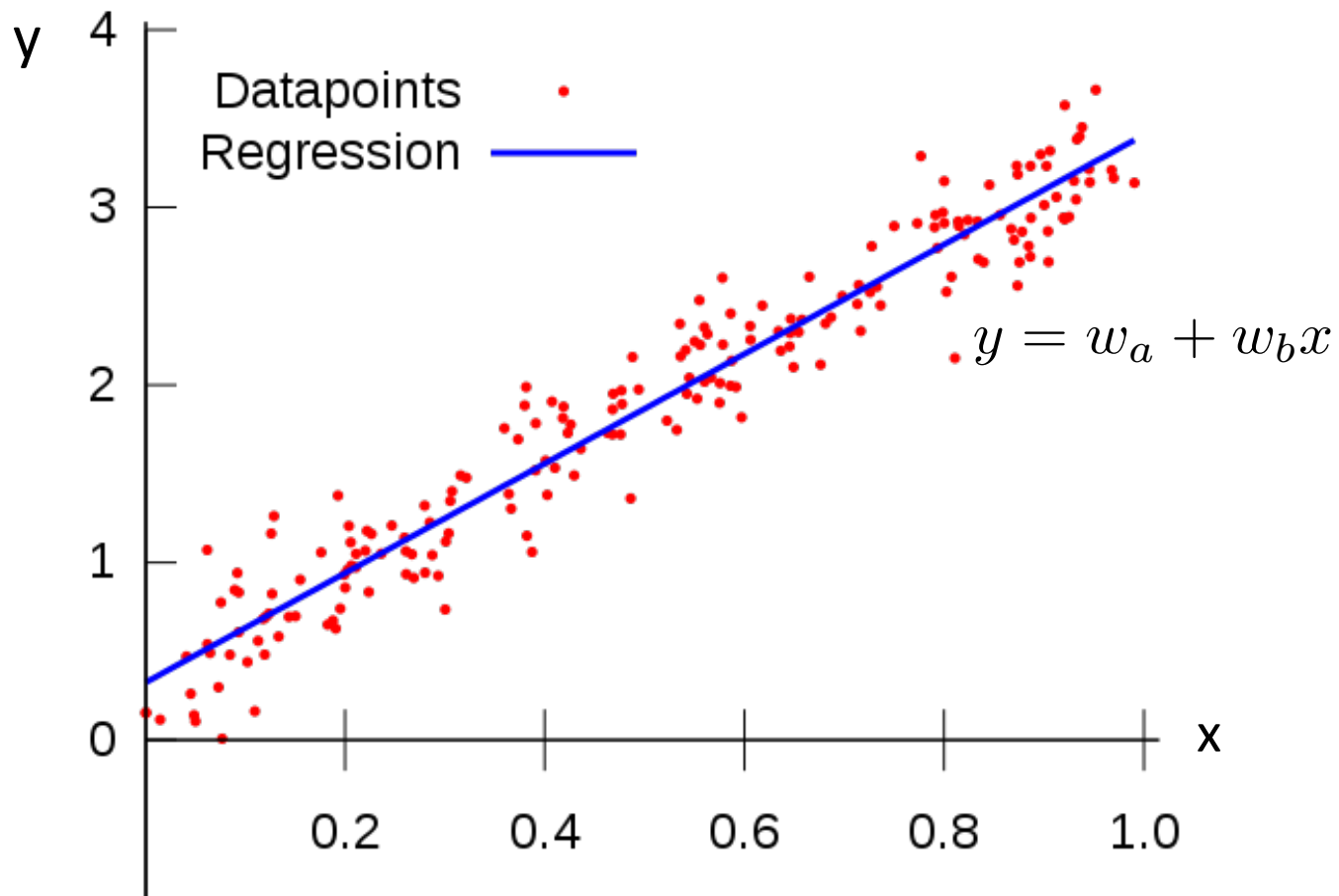


Fei-Fei Li, Stanford

# ImageNet and ILSVRC



ImageNet Classification top-5 error (%)

| Benchmark | Error |
|---|---|
| ILSVRC 2010 NEC America | 28.2 |
| ILSVRC 2011 Xerox | 25.8 |
| ILSVRC 2012 AlexNet | 16.4 |
| ILSVRC 2013 Clarifi | 11.7 |
| ILSVRC 2014 VGG | 7.3 |
| ILSVRC 2014 GoogleNet | 6.7 |
| ILSVRC 2015 ResNet | 3.5 |

# Why the sudden success?

o Availability of massive datasets like Imagenet

o Computing power to train deep neural networks
   o Parallelization
   o GPUs

o Algorithmic advances:
   o Momentum, Adagrad, Adam etc.

# Core of ML: Gradient Descent (GD)

# Simplest ML example: Regression



$$y = w_a + w_b x$$

Given a big dataset of $(\mathbf{x_1}, y_1)$, $(\mathbf{x_2}, y_2)$, $(\mathbf{x_3}, y_3)$, $(\mathbf{x_4}, y_4)$, ….$(\mathbf{x_N}, y_N)$
Find the optimal weights $\mathbf{w}$

# Core of ML: Gradient Descent (GD)

$$\min_{\mathbf{w}} F(\mathbf{w}) = \min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^{N} \nabla (y_i - \mathbf{w}^T \mathbf{x})^2$$

$F(\mathbf{w})$

$F(\mathbf{w}^*)$

# Core of ML: Gradient Descent (GD)

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla F(\mathbf{w_t})$$

# Exercise: Find the update rule for $w_a$ and $w_b$



$$y = w_a + w_b x$$

Given a big dataset of $(\mathbf{x_1}, y_1)$, $(\mathbf{x_2}, y_2)$, $(\mathbf{x_3}, y_3)$, $(\mathbf{x_4}, y_4)$, ….$(\mathbf{x_N}, y_N)$
Find the optimal weights $\mathbf{w} = (w_a, w_b)$

# Gradient Descent (GD)

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \frac{1}{N} \sum_{i=1}^{N} \nabla (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

Too expensive for large datasets

$\mathcal{L}((w_1, w_2))$

# Stochastic Gradient Descent (SGD)

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \boxed{\nabla (y_i - \mathbf{w}^T \mathbf{x}_i)^2}$$

Easy, but possibly too noisy

# Mini-batch SGD

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \frac{1}{m} \sum_{i=1}^{m} \nabla (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

Less noisy, but also computationally tractable

$\mathcal{L}(w_1, w_2)$

# Exercise: How does variance scale with m?

If $Var(\nabla F(\mathbf{w}, \xi_i)) = \sigma^2$

What is the variance of the gradient update in mini-batch SGD?

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \sum_{i=1}^{m} \frac{1}{m} \nabla F(\mathbf{w}_t, \xi_i)$$

# Convergence of SGD

$$\mathbb{E}[F(\mathbf{w}_k) - F_*] \leq \frac{\eta L M}{2c} + (1 - \eta c)^{k-1}\left(F(\mathbf{w}_0) - F_* - \frac{\eta L M}{2c}\right)$$

Decay Rate

Error Floor

How does decay rate and error floor change with
- η (Learning Rate) ?
- M (Second moment of gradient) ?

# Many other variants of SGD

- Momentum SGD

- Nesterov Momentum

- AdaGrad

- Adam

- AdaDelta

- RMS prop
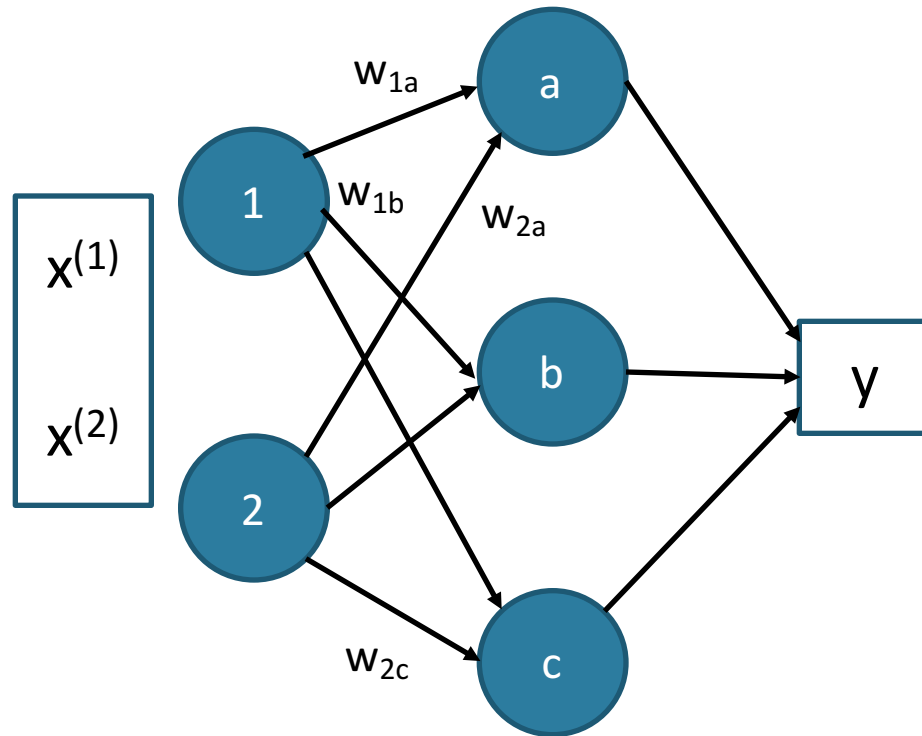
# Many other variants of SGD

# Many other variants of SGD

# SGD and Backpropagation



Given a big dataset of $(\mathbf{x_1}, y_1)$, $(\mathbf{x_2}, y_2)$, $(\mathbf{x_3}, y_3)$, $(\mathbf{x_4}, y_4)$, ....$(\mathbf{x_N}, y_N)$
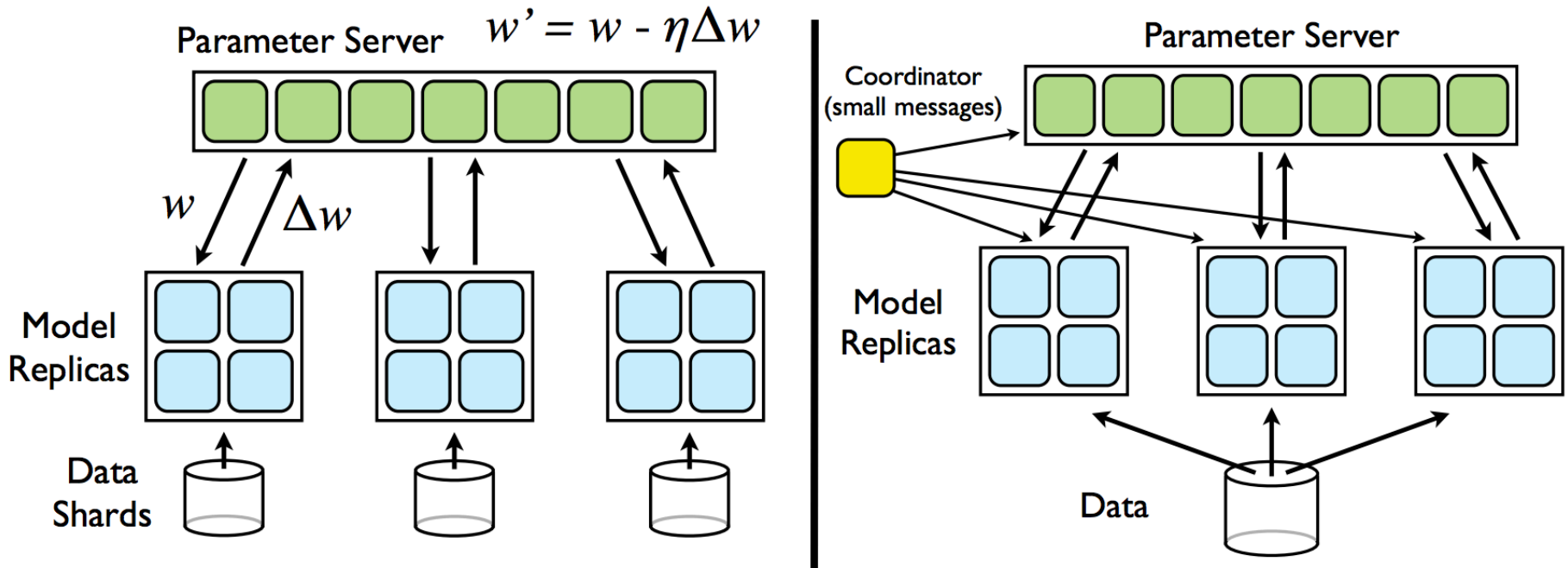Find the optimal weights $\mathbf{w}$

# SGD and Backpropagation



Input to a $= inp_a = w_{1a} x_1 + w_{2a} x_2$

Output of a $= out_a = g(inp_a)$
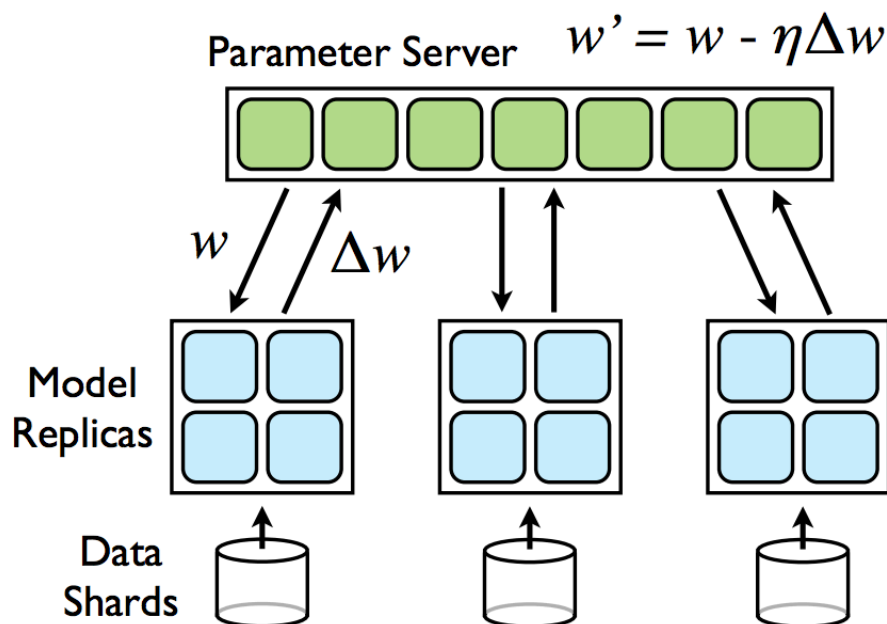
# Distributed Deep Learning

## Data Parallelism

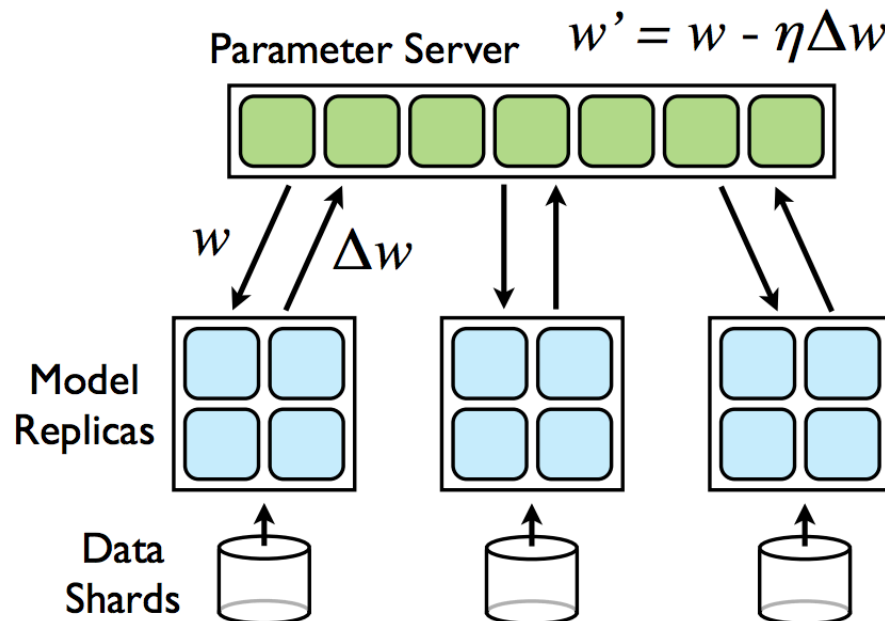# Distributed Deep Learning

Model Parallelism

# Synchronous SGD

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \sum_{k=1}^{K} \frac{1}{K} \nabla F(\mathbf{w}_t, \xi_k)$$
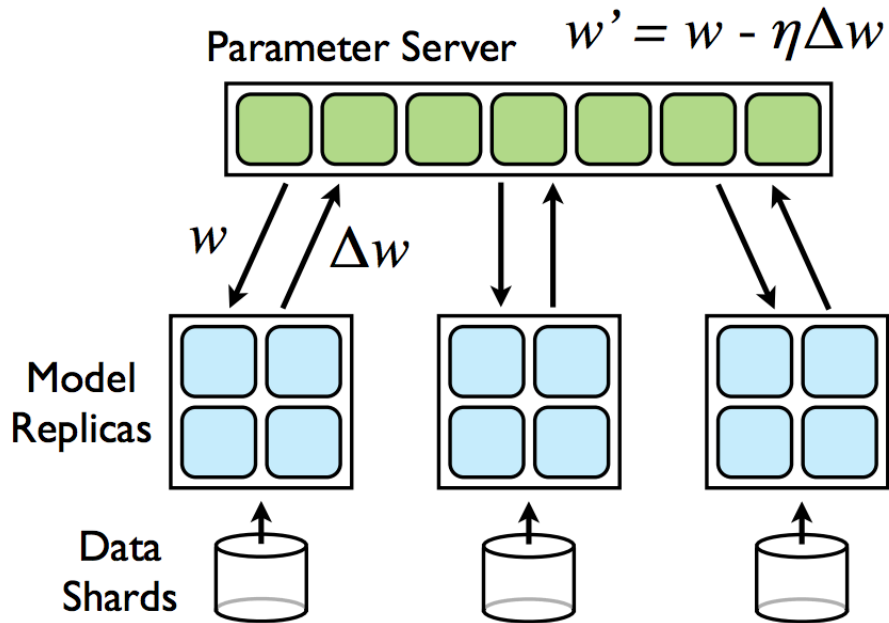
# Q: What is the convergence rate and error floor?

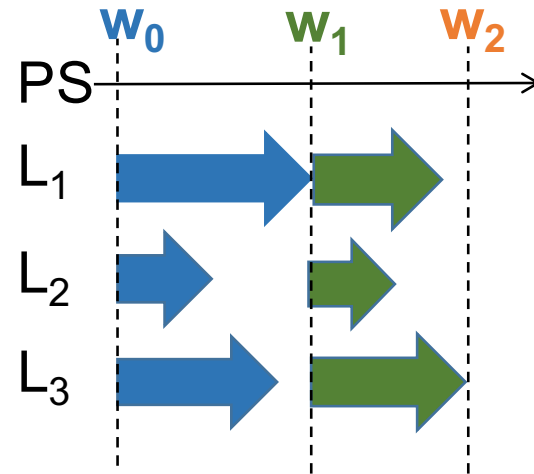$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \sum_{k=1}^{K} \frac{1}{K} \nabla F(\mathbf{w}_t, \xi_k)$$



Parameter Server $w' = w - \eta \Delta w$

$w$ $\Delta w$

Model Replicas

Data Shards

# Q: What is the time to complete each iteration?

$$\mathbb{E}[T] = \mathbb{E}[\max(X_1, X_2, \ldots X_K)]$$

Slowest Learner is the bottleneck
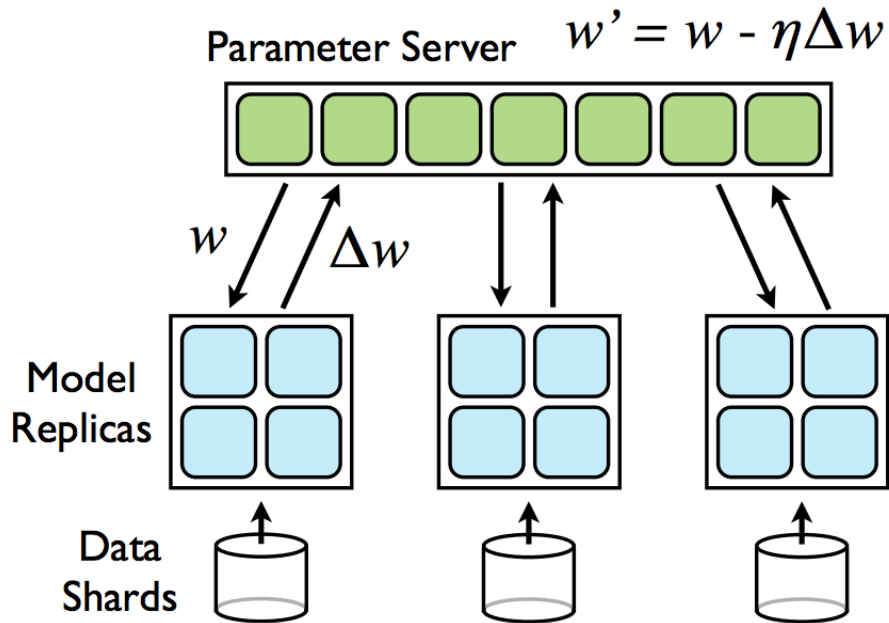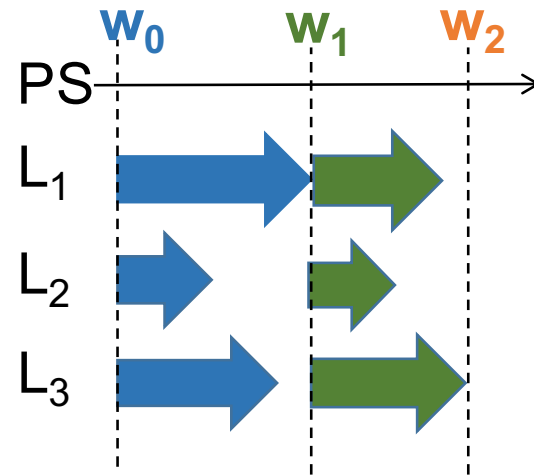
Parameter Server  $w' = w - \eta \Delta w$

Fully Sync-SGD

# Q: How can we reduce it?

$$\mathbb{E}[T] = \mathbb{E}[\max(X_1, X_2, \ldots X_K)]$$

Slowest Learner is the bottleneck

Parameter Server $\quad w' = w - \eta \Delta w$

$w \quad \Delta w$

Model Replicas

Data Shards

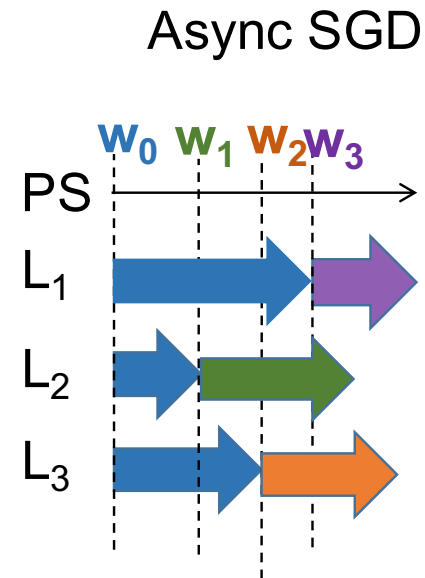Fully Sync-SGD

$W_0 \quad W_1 \quad W_2$

PS

$L_1$

$L_2$

$L_3$
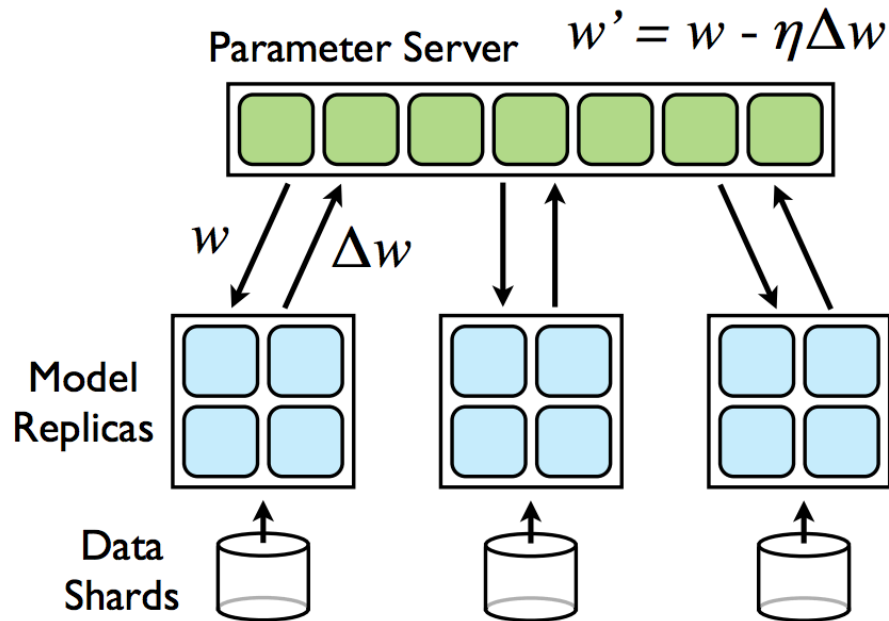
# Asynchronous SGD: Don't wait for all

Asynchronous SGD cuts the latency tail.

But, what effect does it have on the error?
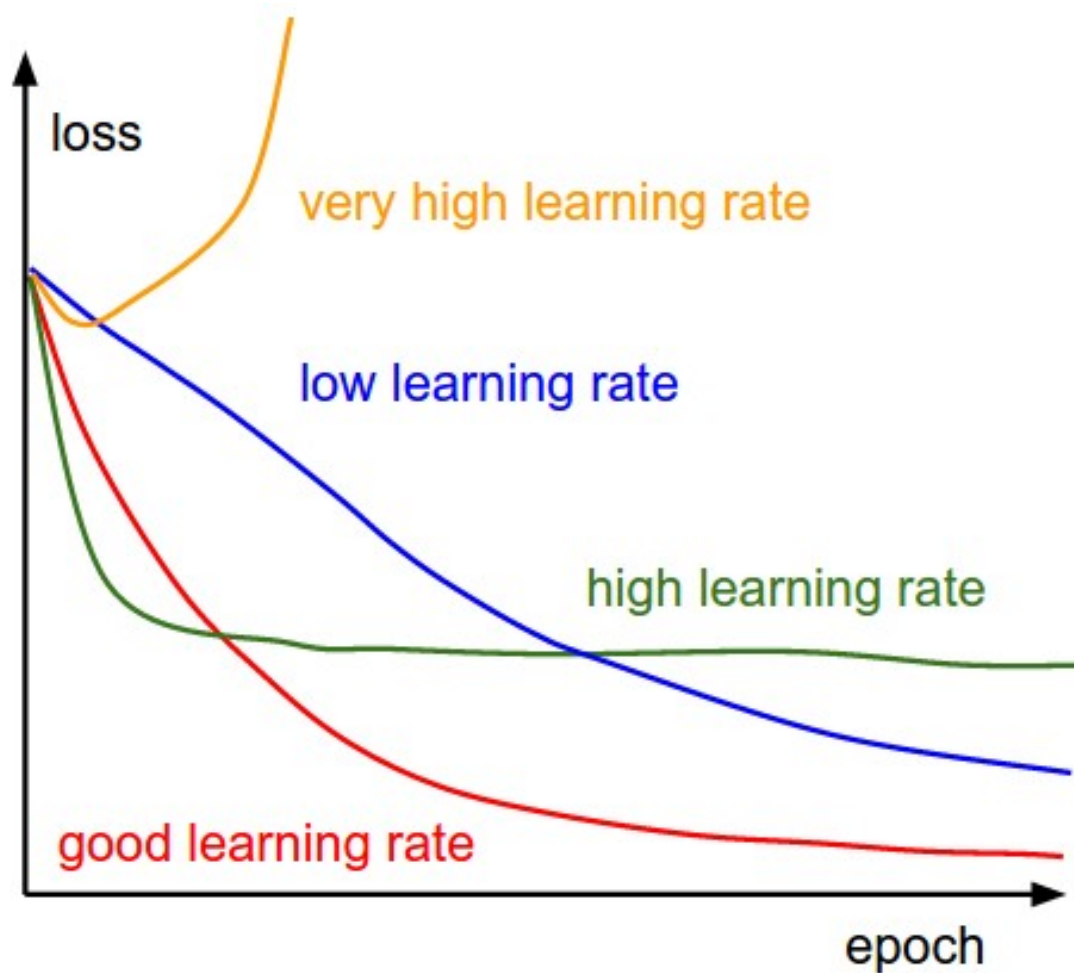
# Variants of Distributed SGD

- Synchronous SGD

- Asynchronous SGD

- HogWild

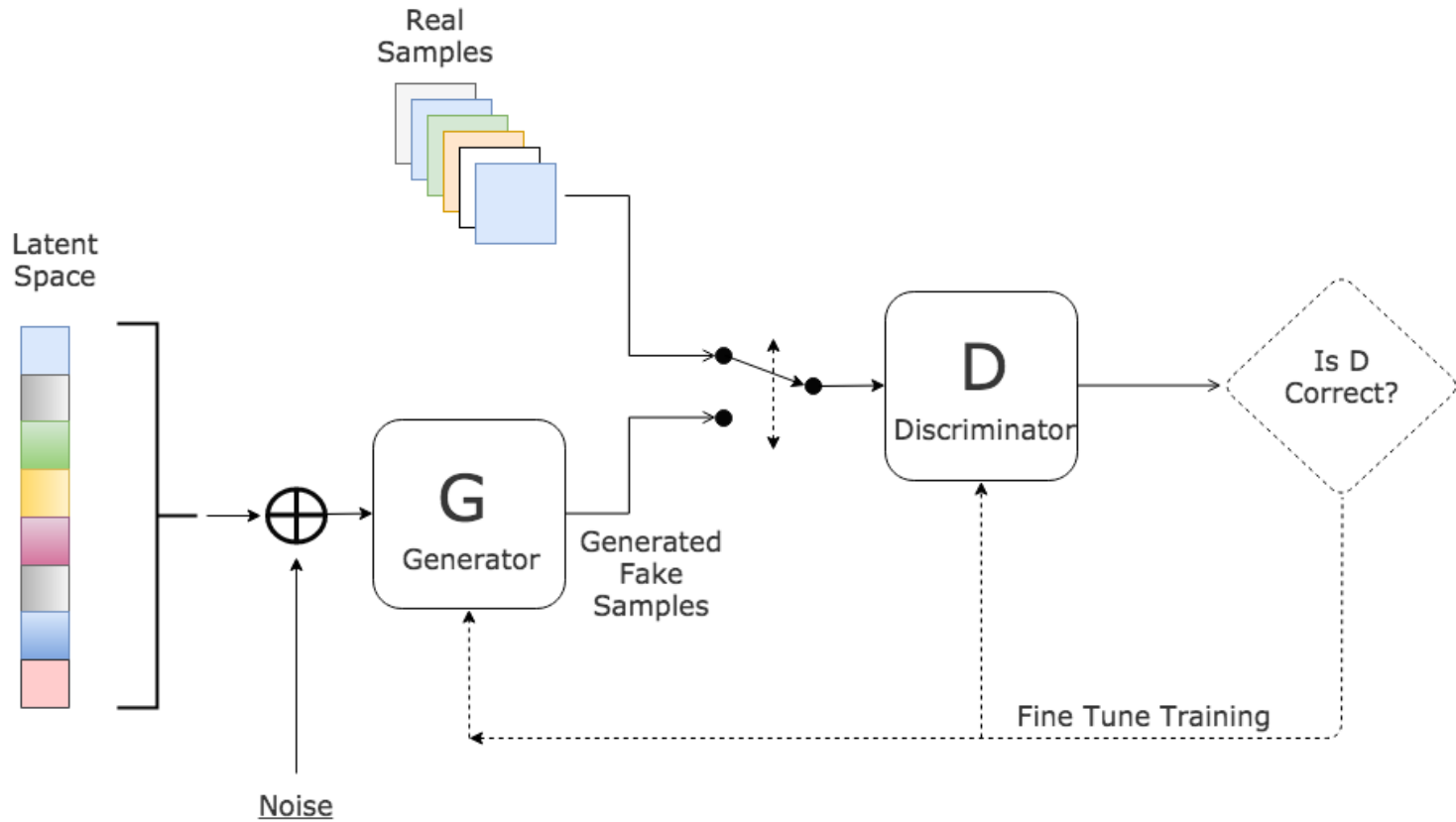- Elastic-Averaging SGD

# Hyper-Parameter Tuning

Need to choose the right

- Learning rate

- Mini-batch size

- Momentum

- Number of layers

- Number of neurons per layer

# Hyper-Parameter Tuning

# Generative Adversarial Networks

# Reinforcement Learning

# Reinforcement Learning