# 18-847F: Special Topics in Computer Systems

## Foundations of Cloud and Machine Learning Infrastructure

# Lecture 1: Logistics and Overview

## Foundations of Cloud and Machine Learning Infrastructure

# Graduate Seminar Class

Few Lectures
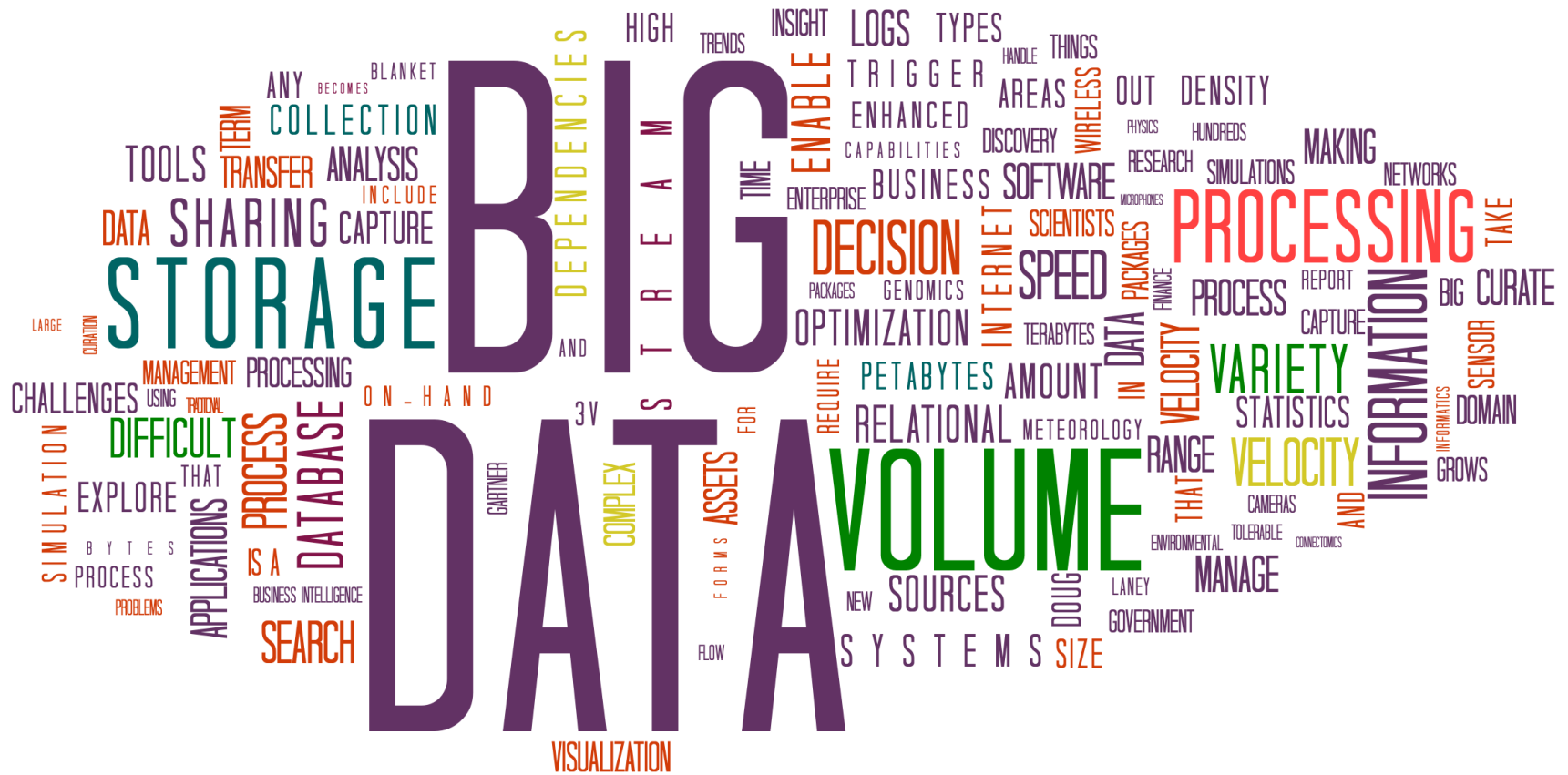
Reading research papers

Student presentations

Class Discussions

Final Research Project (No Exams!)

# Learning Objectives

o Know the state-of-the-art frameworks in cloud and machine learning and their theoretical foundations

o Read and provide constructive criticism of research papers

o Present to an audience, and answer their questions

o Do creative, collaborate research

# Why study Cloud and ML infrastructure?



What are the largest words after 'Big Data'?

# Big Data Gold Rush



Who got rich in the
California gold rush?

# Big Data Gold Rush
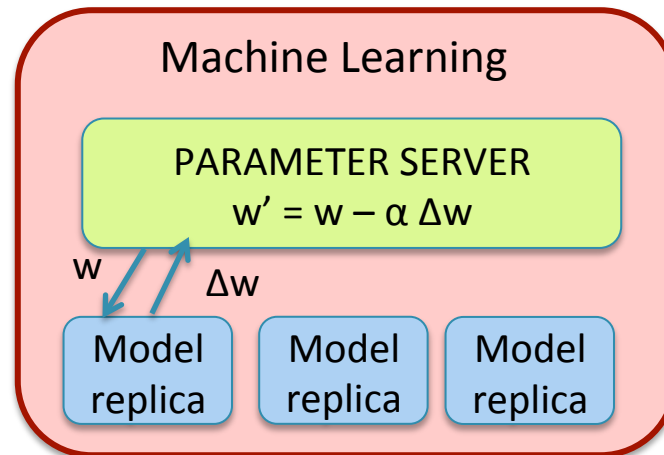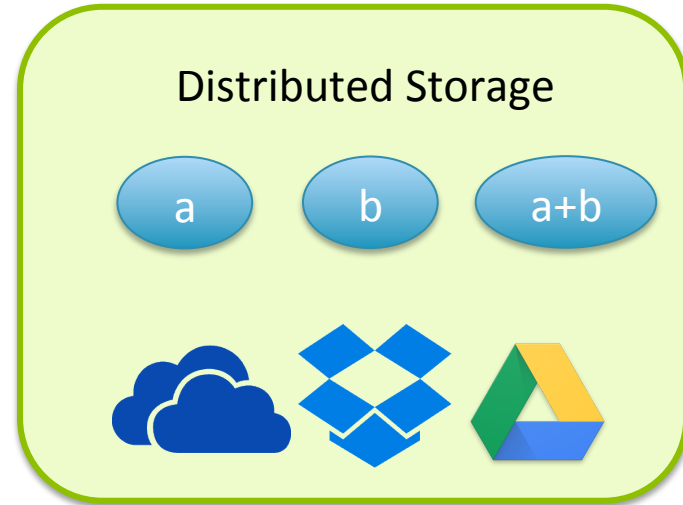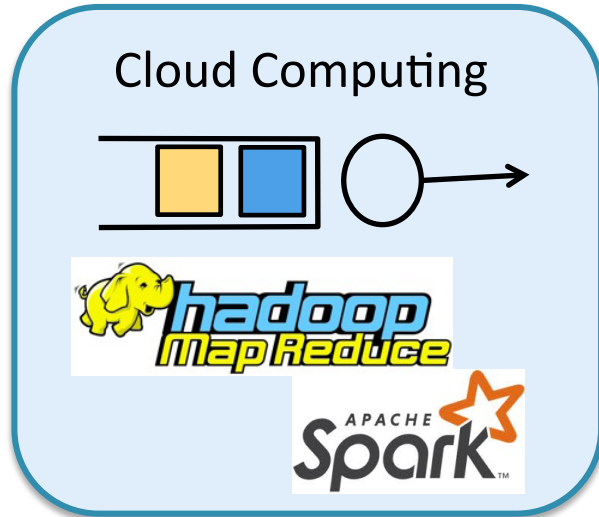


Who got rich in the California gold rush?

In the Big Data rush, it's the infrastructure companies

# Topics Covered

## Cloud Computing

hadoop Map Reduce

APACHE Spark

## Distributed Storage

a     b     a+b

## Machine Learning

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

w     $\Delta w$

| Model replica | Model replica | Model replica |

# Topics Covered

Cloud Computing
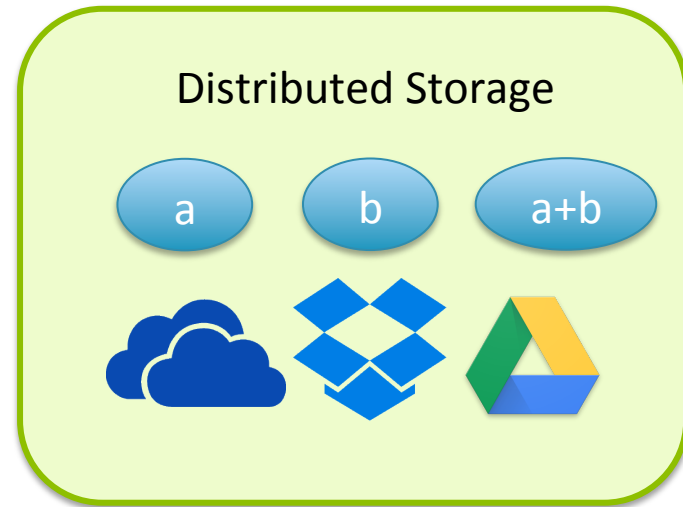
o     Scheduling in Parallel Computing

o     MapReduce, Spark

o     Straggler Replication

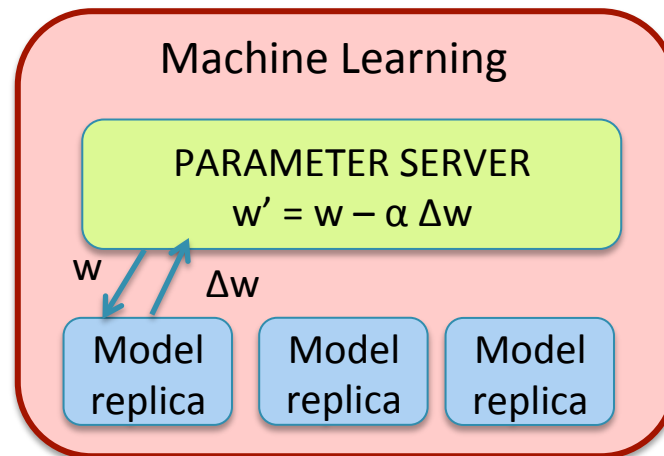o   Task Replication in Queueing Systems

# Topics Covered

o   Coding for locality/repair

o   Systems implementation of codes

o   Reducing latency in content

download

# Topics Covered

o    SGD and its convergence

o    Distributed Deep Learning

o    Hyper-parameter tuning

o    GANs, Deep reinforcement learning

<div style="border:2px solid darkred; background:#f8d0d0; border-radius:20px;">

**Machine Learning**

<div style="background:#d0e860; border:2px solid #8a8a4a; border-radius:10px;">

PARAMETER SERVER

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

</div>

w    $\Delta w$

| Model replica | Model replica | Model replica |

</div>

# Instructor: Gauri Joshi

B.Tech+M.Tech
2005-2010

SM + PhD
2010-2016

Research Staff Member
2016-2017

Assistant Professor
Fall 2017 -

Internships

# Have worked in all these areas

## Cloud Computing

hadoop
Map Reduce

APACHE
Spark™

## Distributed Storage

a    b    a+b

## Machine Learning

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

w        $\Delta w$

Model replica    Model replica    Model replica

# Student Introductions

o Name?

o Department?

o Undergrad/Masters/PhD?

o Previous related classes (if any)?

o What you are looking to learn from this class?

Waiting list will be cleared soon!

# Class Hours and Website(s)

o When: Mon, Wed 4:30-6:00 pm

o Where: Scaife Hall 222

o Class Website (Readings, Schedule):
  https://www.andrew.cmu.edu/user/gaurij/18-847F-Fall-2018.html

o Canvas Site (Readings, Assignments, Projects):
  https://canvas.cmu.edu/

o No prerequisites. Basic knowledge of probability and linear algebra is encouraged.

# Reading Material

Papers will be posted on the class website or on Canvas

o   Book chapters

o   Survey papers

o   Theory papers (Scheduling, Queuing, Coding, Optimization)

o   Systems papers (Cloud, Machine Learning)

Additional reference books listed in the syllabus

# Instructor/TA and Office Hours

Instructor: Prof. Gauri Joshi (gaurij [AT]andrew.cmu.edu)

TA: Jianyu Wang (jianyuw1 [AT]andrew.cmu.edu)

Office Location: CIC 4105

Office Hours: By appointment

# Graduate Seminar Class

A few lectures

Reading research papers

Student presentations

Class Discussions

Final Research Project

# Lectures

o Next week: Deeper Overview of probability and queuing theory

o Guest lectures during the semester by authors of papers relevant to this class

# Graduate Seminar Class

A few lectures

Reading research papers

Student presentations

Class Discussions

Final Research Project

# Homeworks (~50%)

o Submit paper review (due 10:00 am before class)

  o ~Two reviews per week

o Discussion with classmates is okay, but write reviews in your own words.

# Paper Review Format

o Summary of the paper

    o Reflects your understanding of the paper

    o Significance & correctness of results

o Discussion Questions for Class (at least 2)

    o Confusions about the paper, open research directions

o Answers to concept-check questions

# Homework Grading Rubric (Total: 10 pts)

o   Understanding of the paper (4 pts)

o   Discussion Questions (3 pts)

o   Concept-check questions (3 pts)

# Graduate Seminar Class

A few lectures

Reading research papers

Student presentations

Class Discussions

Final Research Project

# Class Presentations (~15%)

o Sign up for presentation at least 1 week in advance

o Each student will present 1-2 times in the semester (depends on # of students registered)

o 20 min presentation, followed by 25 min discussion
- o Motivation and Related work
- o Summary of main results
- o Your views on the paper

# Presentation Grading Rubric (Total: 10 pts)

- o Motivation (1.5 pts)

- o Clarity (1.5 pts)

- o Understanding/Correctness (4 pts)

- o Peer-review Feedback (3 pts)

# Graduate Seminar Class

A few lectures

Reading research papers

Student presentations

Class Discussions

Final Research Project

# Class Participation (~15%)

o The class will be divided into groups of 3-4 students each

o Each group will discuss one of the discussion questions among themselves

o Summarize the discussion to the whole class

# Participation Grading Rubric (Total: 5 pts)

o   Attendance and attention (1.5 pt)

o   Speaking up in class (1.5 pt)

o   Insightful Questions/Comments (2 pt)

# Graduate Seminar Class

A few lectures

Reading research papers

Student presentations

Class Discussions

Final Research Project

# Research Project (~20%)

o Groups of 1-3

o Original research on a topic of your choice
  o Topics aligned with your research allowed and encouraged
  o If you can't think of topics, come talk to Jianyu or me

o Possible Project Types:
  o New theoretical analysis
  o Implementation using one of the frameworks discussed
  o In-depth literature survey of a particular topic

# Timeline

o 1-page proposal due Oct 3

o Publishable quality report (max 5 pg) in ACM format
- o Initial draft due: Nov 21
- o Final report due: Dec 7

o Last week of class: Presentations (~30 min per group)

o Peer-review other presentations

# Project Grading Rubric (Total: 20 pts)

o Originality (1 pts)

o Review of Related Work (1.5 pts)

o Writing and Organization (1.5 pts)

o Technical Results (4 pts)

o Final presentation (10 pts)

o Peer-Review (2 pts)

# In Summary..

o   Paper Reading

o   Submitting Reviews

o   Class Presentations

o   Final Project

Might seem like a lot of work but..

o   You will get fast and efficient at reading papers

o   The project will be a fun, collaborative exercise

o   No exams!

# TO DO

o Fill out the sign-up sheet

o Sign-up for presentations

o Start reading the papers

o Form groups for class projects

o Start thinking about projects

# Topics Covered

## Cloud Computing

hadoop
Map Reduce

APACHE
Spark

## Distributed Storage

a    b    a+b

## Machine Learning

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

w          $\Delta w$

Model replica    Model replica    Model replica

# History and Overview



Cloud Computing



Distributed Storage

# History and Overview

Cloud Computing



o    MapReduce, Spark

o    Scheduling in Parallel Computing

o    Straggler Replication

o    Task Replication in Queueing Systems

# What is the cloud?



A collection of servers that can function as a single computing node, and can be accessed from multiple devices

# 1960's: The Mainframe Era

o    Large, expensive machines
o    Only one per university/institution



IBM 704 (1964)

# 1970's: Virtualization

o   IBM released a VM OS that allowed multiple users to share the mainframe computer



IBM 704 (1964)

# 1980's-1990's: Internet and PCs

o PCs become affordable

o Internet connectivity went on improving

o Virtual Private Networks (VPNs)

o Grid Computing: Connect cheap PCs via the Internet

o On the theory side, queuing theory, traditionally used in operations management rebounded

# 1990's: Scheduling in Parallel Computing

o **Bin-Packing**

J2

J1

J3

J4

Time

Processors

For references see survey
[Weinberg 2008]

43

# 1990's: Scheduling in Parallel Computing

o **Bin-Packing**
  o Need job size estimates



For references see survey
[Weinberg 2008]

# 1990's: Scheduling in Parallel Computing

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

# 1990's: Internet and PCs

- o PCs become affordable
- o Internet connectivity went on improving

- o Virtual Private Networks (VPNs)
- o Grid Computing: Connect cheap PCs via the Internet

- o Many Internet Companies bought their own servers and managed them privately
- o But then the Dotcom bubble burst..

# 2000's: The Cloud Computing Era

o The idea of a flexible, low-cost, scalable, shared computing environment developed



o Computing become a utility, like electricity

# 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]

# The Tale of Tails



Latency

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

# Problem: Stragglers in Parallel Computing

o   A job with hundreds of parallel tasks

o   Machine response time can vary due to virtualization, congestion etc.

o   The slowest tasks are the bottleneck in job completion

Task 1

Task 2

Task n

[Dean "Tail at Scale" 2013]

| Latency | 50%ile | 99%ile |
|---|---|---|
| **1 task finishes** | 1ms | 10ms |
| **All tasks finish** | 40ms | 140 ms |

# Exercise: Tale of Tails

A server finishes a task in 1 sec with probability 0.9, and 10 sec with probability 0.1

o    What is the expected task execution time?

o    If 100 tasks are run in parallel of 100 servers, what is the expected time to complete all of them.

# Exercise: Tale of Tails

A server finishes a task in 1 sec with probability 0.9, and 10 sec with probability 0.1

o What is the expected task execution time?

   1*0.9 + 10*0.1 = 1.9

o If 100 tasks are run in parallel of 100 servers, what is the expected time to complete all of them.

# Exercise: Tale of Tails

A server finishes a task in 1 sec with probability 0.9, and 10 sec with probability 0.1

○ What is the expected task execution time?

$1*0.9 + 10*0.1 = 1.9$

○ If 100 tasks are run in parallel of 100 servers, what is the expected time to complete all of them.

$1*0.9^{100} + 10*( 1 - 0.9^{100}) \sim 10$

# 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

# Straggler Replication Analysis
## [ Wang-GJ-Wornell SIGMETRICS 2014, 15]

**PARAMETERS**

p: Frac. of tasks replicated

r: # additional replicas

c: kill/keep  original task

**METRICS**

E[T] = Time to finish all tasks

E[C] = Total server runtime per task

$$\mathbb{E}[T] = \mathbb{E}[X_{(1-p)n:n}] + \mathbb{E}[Y_{pn:pn}]$$

Y is the residual service time after adding replicas

Central Value Theorem    n -> ∞

Extreme Value Theorem    n -> ∞

$$F_X^{-1}(1-p)$$

Different behavior for Exponential, Light or Heavy tailed Y

# Simulations using Google Cluster Data
## Latency-Cost Trade-off



Increasing fraction of tasks replicated

- MapReduce setting (r=1)
- $r=2$ & keep original copy
- $r=3$ & keep original copy

p=0.0

p=0.5

p=0.5

p=0.5

Expected Latency E[T]

Expected Cost E[C]

Careful choice of replication strategy can be better than the default in MapReduce

# Task Replication in Queueing Systems

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

Task

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

## COST

o Additional computing time at servers

# Task Replication in Queueing Systems

**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

# Analogy: Supermarket Queues


© Getty Images

# Supermarket Queues

# Supermarket Queues



Get a friend to join the other queue!

What if everyone in the supermarket uses this strategy?

# Design Questions

o   How many replicas to launch?

o   Which queues to join?

o   When to issue and cancel the replicas?

# Surprising Insight

In certain regimes, replication could make the
whole system faster, and cheaper!



vs

Effective service rate > Sum of individual servers

# History and Overview

## Cloud Computing

hadoop Map Reduce

APACHE Spark™

## Distributed Storage

a     b     a+b

# History and Overview

o    RAID systems

o    Coding for locality/repair

o    Systems implementation of codes

o    Reducing latency in content
download


Distributed Storage

# 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 was here at CMU

# Coding Theory

o For reliable communication in presence of noise

o Bell Labs was one of the leaders in 1950's

o Key figures: Claude Shannon and Richard Hamming

# Simplest Codes

o **Repetition Code**

  o 0 → 0 0 0  : Rate: 1/3

  o If receive 0 ? ? we can recover from 2 erasures

o **(3,2) code: Data bits: a, b  Parity bit: (a XOR b)**

  o Example:  0 1 1,  1 1 0: Rate 2/3

  o If we receive 0 ? 1 or  ? 1 0 we can correct the failed bit

  o 2 bit symbols:  (0 1)  ?  (1 1)

# (n,k) Reed-Solomon Codes: 1960

o   Data: $d_1, d_2, d_3, \ldots d_k$

o   Polynomial: $d_1 + d_2 x + d_3 x^2 + \ldots d_k x^{k-1}$

o   Parity bits: Evaluate at n-k points:

      x=1:         $d_1 + d_2 + d_3 + d_4$

      x=2:         $d_1 + 2 d_2 + 4 d_3 + 8 d_4$

      x=3 :         ....

      x=4:         ....

      x=n:         ...

o   Can solve for the coefficients from any k coded symbols

# Example: (4,2) Reed-Solomon Code

o   Data: $d_1$, $d_2$ → Polynomial: $d_1 + d_2 x + d_3 x^2 + \ldots d_k x^{k-1}$



o   Can solve for the coefficients from any k coded symbols
o   Microsoft uses (7, 4) code
o   Facebook uses (14,10) code

# Locality and Repair Issues

o Repairing failed nodes is hard with Reed-Solomon Codes..



o If we lose 1 node:

    o Need to contact k other nodes

    o Need to download k times the lost data

# Solution: Locally Repairable Codes

o Codes designed to minimize:
  o Repair Bandwidth
  o Number of nodes contacted

# Replicated Storage

o   Content is replicated on the cloud for reliability



Any 1 out of 3 copies is sufficient

o   Can support more users simultaneously
o   Replicated used for "hot" data, i.e. more frequent accessed

# 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

Q: How many users can we serve, and how fast?

# The (n,k) fork-join model

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

76

# 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

$$A_1$$

$$A_2$$

$$x$$

# 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

$A_1$  x    $A_2$  x

Wait for both to finish

# 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

# Second half: 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

# Second half: 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



NIPS Growth
Total Registrations 3755

# The Origins: 1950

Alan Turing

# Neural Networks: Perceptron 1957



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

Frank Rosenblatt
(1928-1971)

Simplified model:

23

# Back-propagation Algorithm (1986)

Geoff Hinton (U. Toronto, Google)

**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 repeatedly adjusts the weights 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 (%)

| Year | 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: Stochastic Gradient Descent (SGD)

# Simplest ML example: Regression



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 $(x_1, y_1)$, $(x_2, y_2)$, $(x_3, y_3)$, $(x_4, y_4)$, ….$(x_N, y_N)$
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

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

# 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

$\mathcal{L}(\omega_1, \omega_2)$

# 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

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

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?

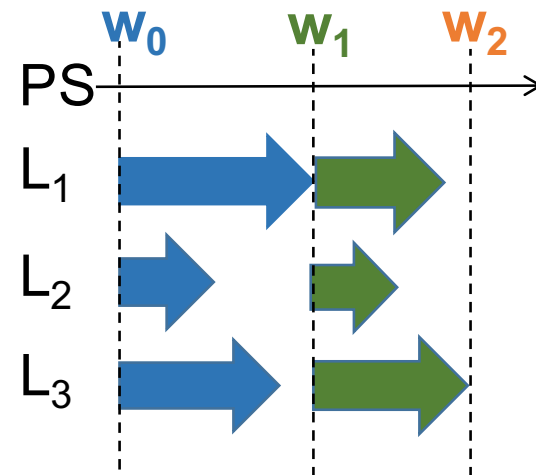$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \sum_{k=1}^{K} \frac{1}{K} \nabla F(\mathbf{w}_t, \xi_k)$$



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

# Q: How can we reduce it?

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

Slowest Learner is the bottleneck

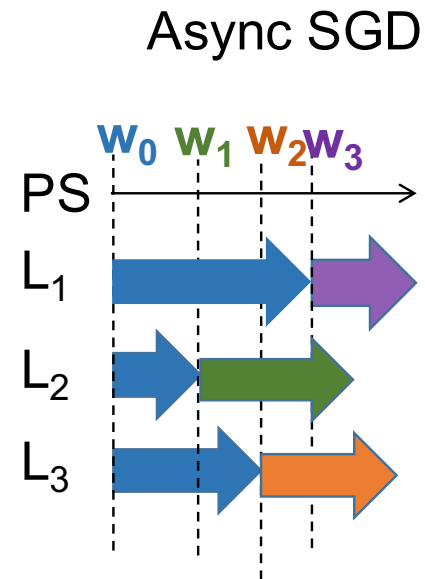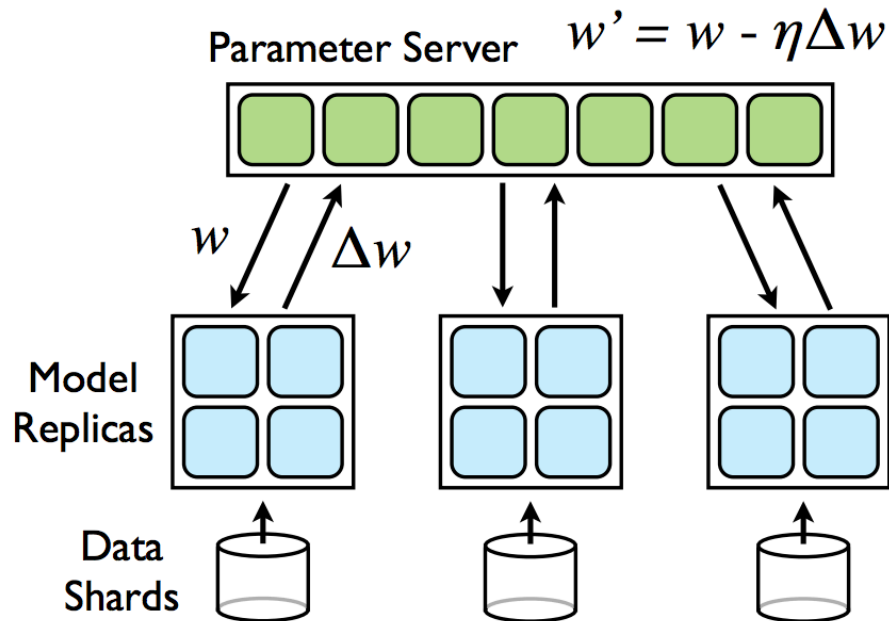Parameter Server  $w' = w - \eta \Delta w$

$w$ / $\Delta w$

Model
Replicas

Data
Shards

Fully Sync-SGD

$w_0$    $w_1$    $w_2$

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?

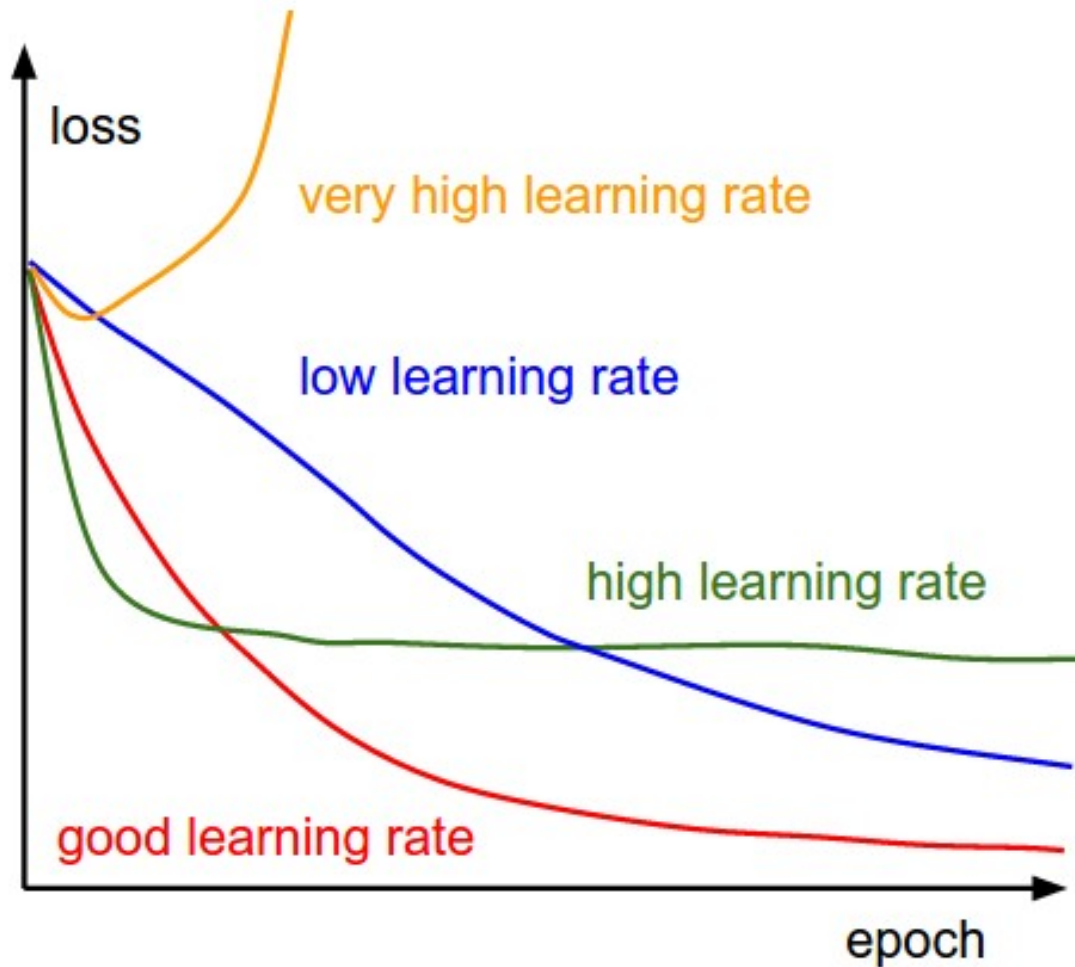# Variants of Distributed SGD

- Synchronous SGD

- Asynchronous SGD

- HogWild

- Elastic-Averaging SGD
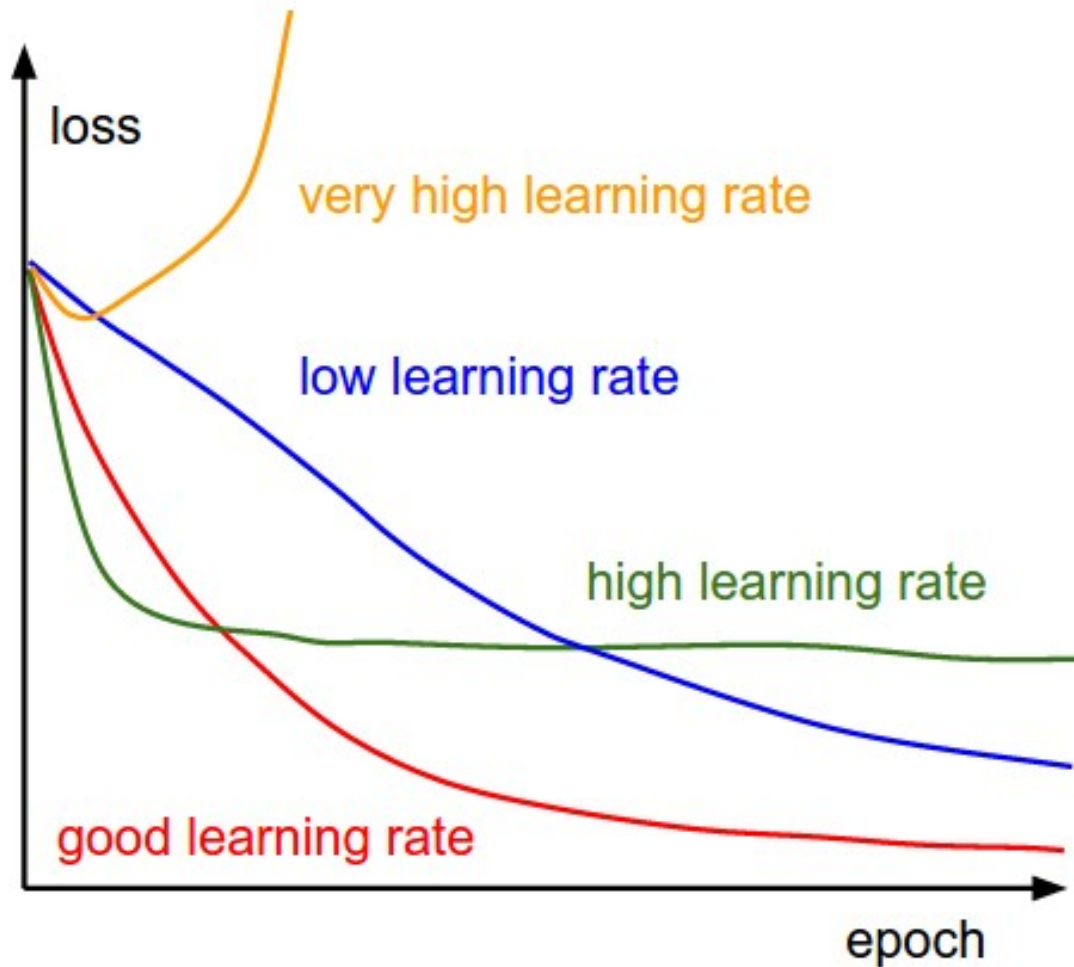
# Hyper-Parameter Tuning

Need to choose the right

- Learning rate

- Mini-batch size

- Momentum

- Number of layers
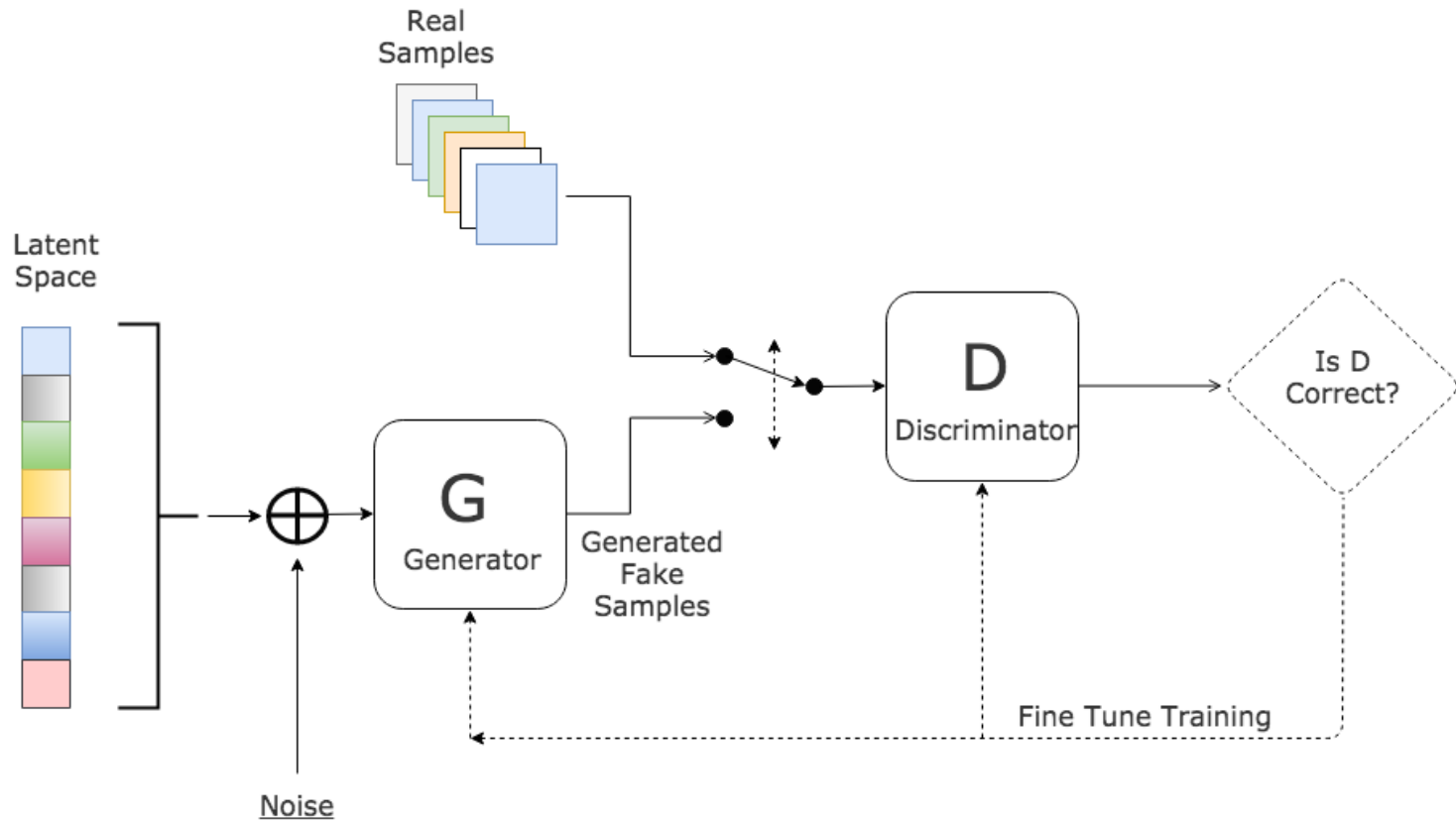
- Number of neurons per layer

# Hyper-Parameter Tuning
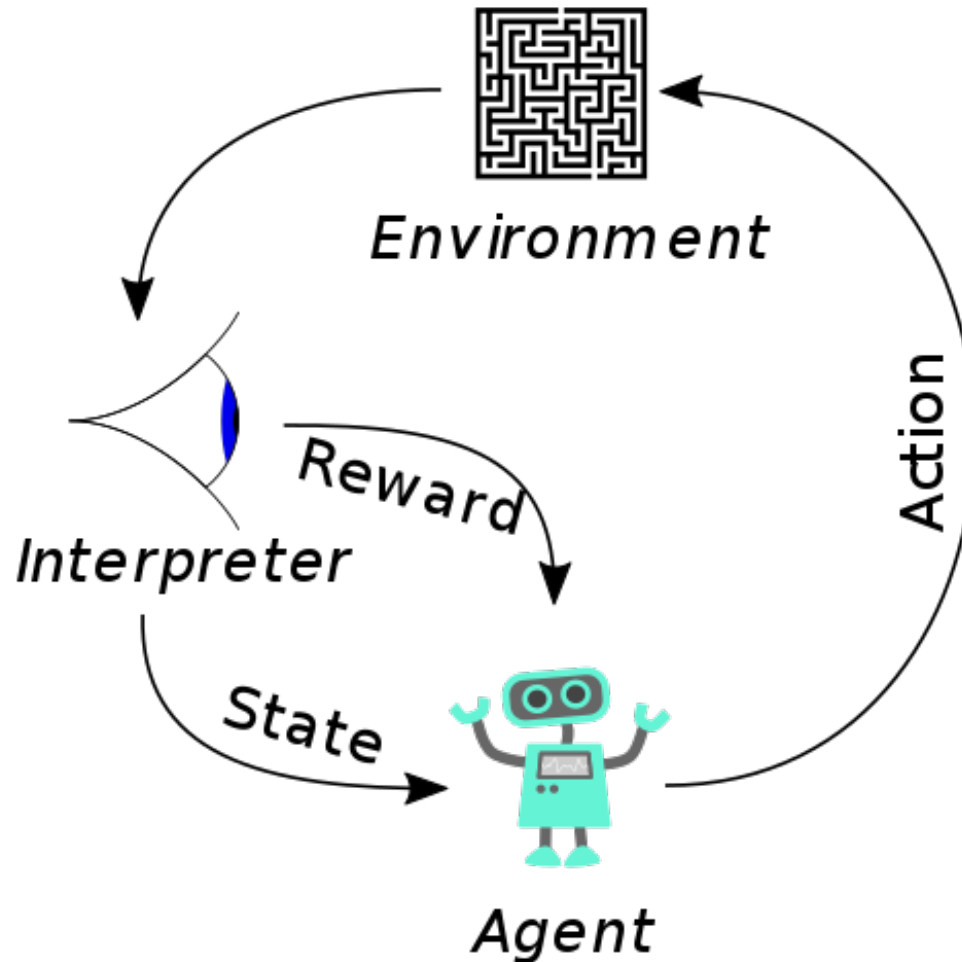
# Multi-armed Bandits and Bayesian Optimization

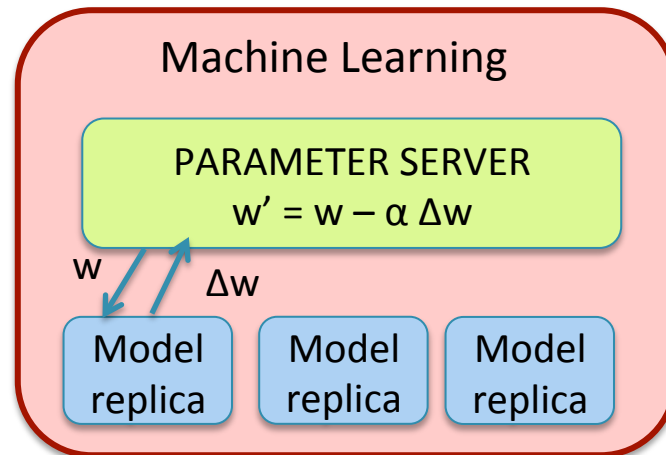# Generative Adversarial Networks

# Reinforcement Learning

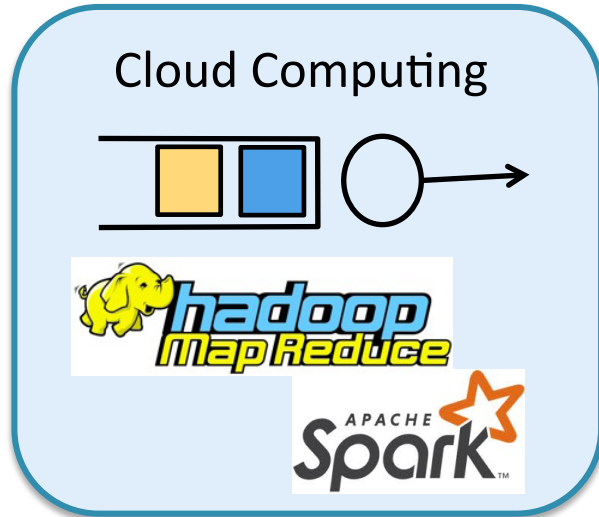# Reinforcement Learning

# 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

# TO DO

- Fill out the sign-up sheet

- Sign-up for presentations

- Start reading the papers

- Form groups for class projects

- Start thinking about projects