

Machine Learning I

60629A

Parallel computational paradigms for large-scale data processing
— Week #10

Today

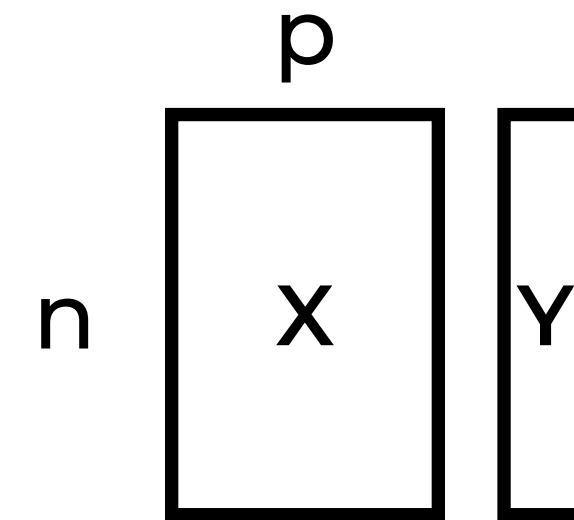
A. Faster computing for machine learning

- Specialized hardware
- Distributed computations
 - Short introduction to MapReduce/Hadoop & Spark

Note: Most lectures so far used stats concepts. Today we'll turn to computer science.

Data & Computation

- We generate massive quantities of data
 1. Google 4K searches/s, Twitter: 6K tweets/s, Amazon: 100s sold products/s
(source: internetlifestats.com)
 2. Banks, insurance companies, etc.
 3. Modestly-sized websites
- Both large n and large p
- In general computation will scale up with the data
- Often fitting an ML models requires one or multiple operations that looks at the whole dataset



e.g., Linear regression $w = (X^T X)^{-1} X^T Y$

Issues with massive datasets

1. Storage
2. Computation

Our World
in Data

This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers.

The chart displays the exponential growth of transistor counts in integrated circuits from 1970 to 2020. The y-axis represents the number of transistors on a logarithmic scale, ranging from 1,000 to 50,000,000,000. The x-axis represents the year. The plot shows a clear upward trend, with transistor counts doubling approximately every two years, as predicted by Moore's Law. Key milestones and labeled chips include:

- 1970s:** Intel 4004, Intel 8008, Intel 8080, Intel 8085, Intel 8086, Intel 8088, Motorola 6800, Motorola 6809, Motorola 68000, TI Explorer's 32-bit Lisp machine chip.
- 1980s:** Intel 80286, Intel 80386, Intel 80486, Intel Pentium, AMD K5, AMD K6, AMD K7, AMD K8, AMD K10, AMD K11, AMD K12, AMD Ryzen, AMD Epyc, Apple A series, Qualcomm Snapdragon, and many others.
- 1990s:** Intel Pentium Pro, Intel Pentium II, Intel Pentium III, AMD K5, AMD K6, AMD K7, AMD K8, AMD K10, AMD K11, AMD K12, AMD Ryzen, AMD Epyc, Apple A series, Qualcomm Snapdragon, and many others.
- 2000s:** Intel Pentium 4, Intel Pentium D, Intel Pentium E, AMD K8, AMD K10, AMD K11, AMD K12, AMD Ryzen, AMD Epyc, Apple A series, Qualcomm Snapdragon, and many others.
- 2010s:** Intel Core i7, Intel Core i5, Intel Core i3, AMD Ryzen, AMD Epyc, Apple A series, Qualcomm Snapdragon, and many others.
- 2020s:** Intel Core i9, Intel Core i7, Intel Core i5, AMD Ryzen, AMD Epyc, Apple A series, Qualcomm Snapdragon, and many others.

Year in which the microchip was first introduced

OurWorldinData.org – Research and data to make progress against the world's largest problems.

Licensed under CC-BY by the authors Hannah Ritchie and Max Roser.

Modern Computation paradigms

- Floating point operations per second (Flop)
- Smartphone ~ 11 TeraFlops
- 1 Tera: 1,000 Giga, 1 Peta: 1,000 Tera, 1 Exa: 1,000 Peta

1. “Single” computers

- Large Computers
- 1,102 petaFlops*

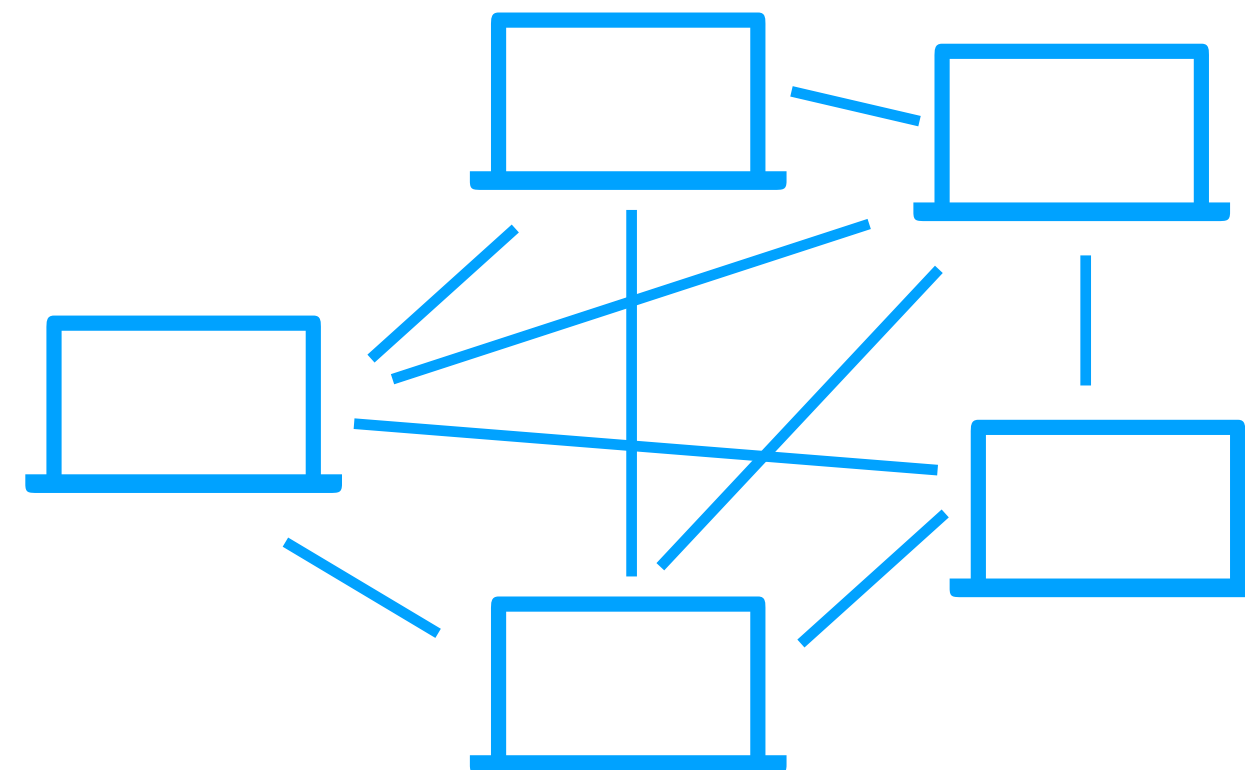
<https://www.top500.org/lists/top500/list/2020/06/>



Photo from Riken

2. Distributed computation

- 2,3 exaFlops*
(Folding@home)



3. Specialized hardware

- Focusses on subset of operations
- Graphical Processing Unit (GPU), Field Programmable Gated Array (FPGA)
- ~1,000 TFlops*



*: these numbers are given as indications and are subject to changes and precisions.

**GPUs & other
specialized hardware**

Hardware

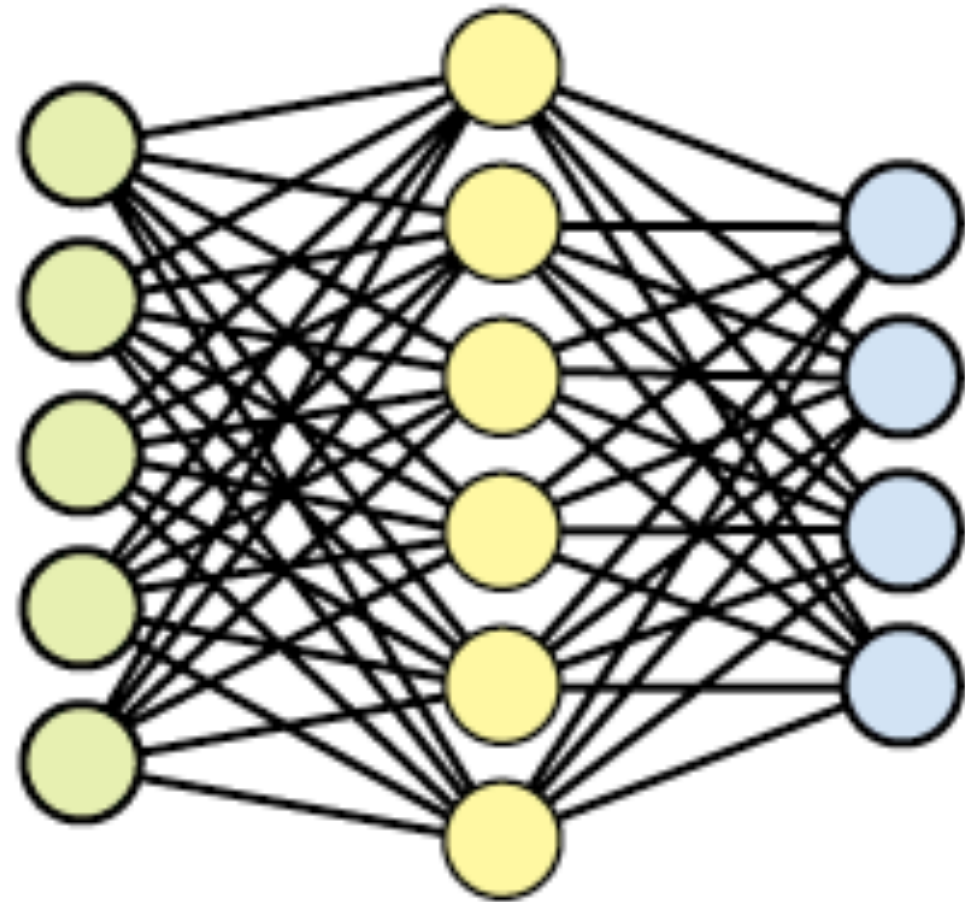


- Central Processing Unit (CPU)
- Computer's cognition
- Executes all the instructions from software
 - Arithmetics, read/writes, logic, etc.

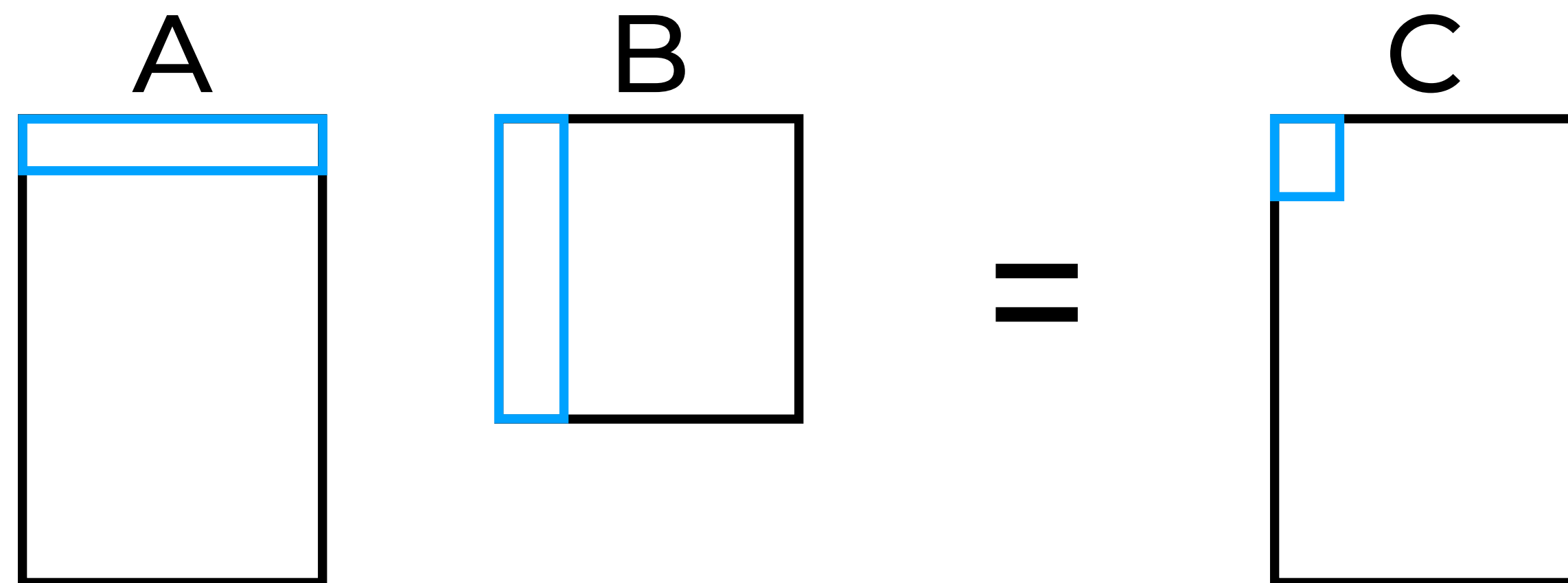
<http://www.personal.psu.edu/users/d/i/dlm99/cpu.html>

Neural Networks

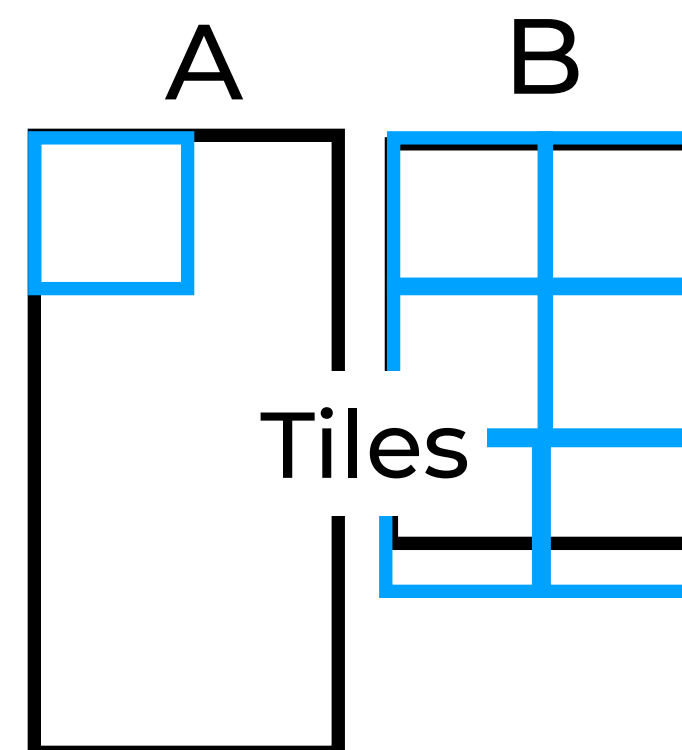
- Linear Algebra:
 - Multiplication *vector X matrix*
 - Neuron activation of multiple neurons
 - Neuron activation for multiple datum
 - Multiplication *matrix X matrix*
 - Activation of multiple neurons for multiple datum
- The exact dimensions of these vector & matrix depend on the data size (mini batch) and the number of neurons



Parallelizing neural network computations



$$C = AB$$



- The values of C can be computed independently from one another
- Matrix multiplication can be parallelized
- When parallelizing we can then obtain very significant gains
- Depend on the size of C
- In practice, divide in tiles
- Note: not all linear algebra operations can be as easily parallelized. Notably, the matrix inverse.

Specialized hardware

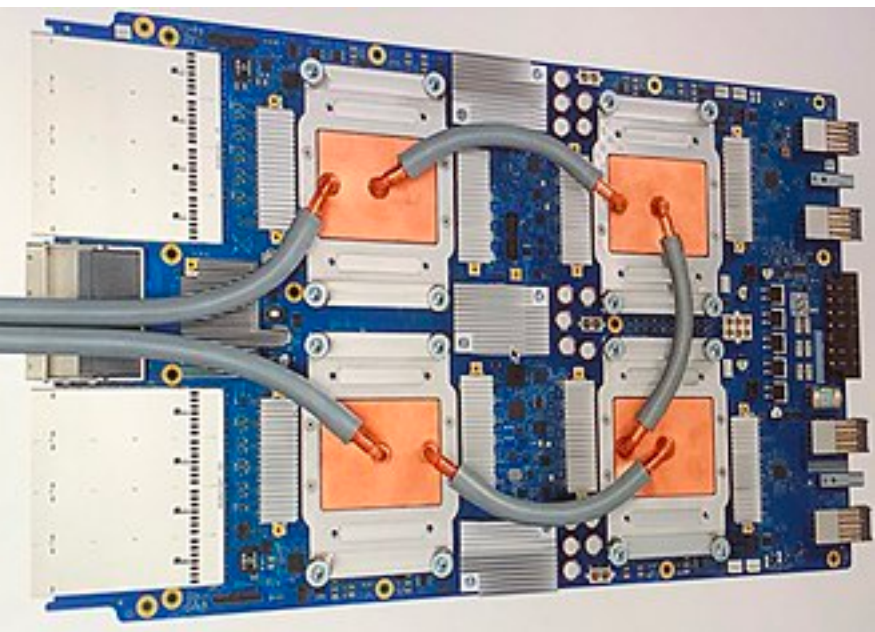


NVidia

- Graphical Processing Unit (GPU)
- Initially designed for 3D games
- Specialized for linear algebra operations
 - Thousands of cores (vs. a few tens for CPUs)
- Fast access to memory
- Multi-GPUs for a single computer



Even more specialized hardware



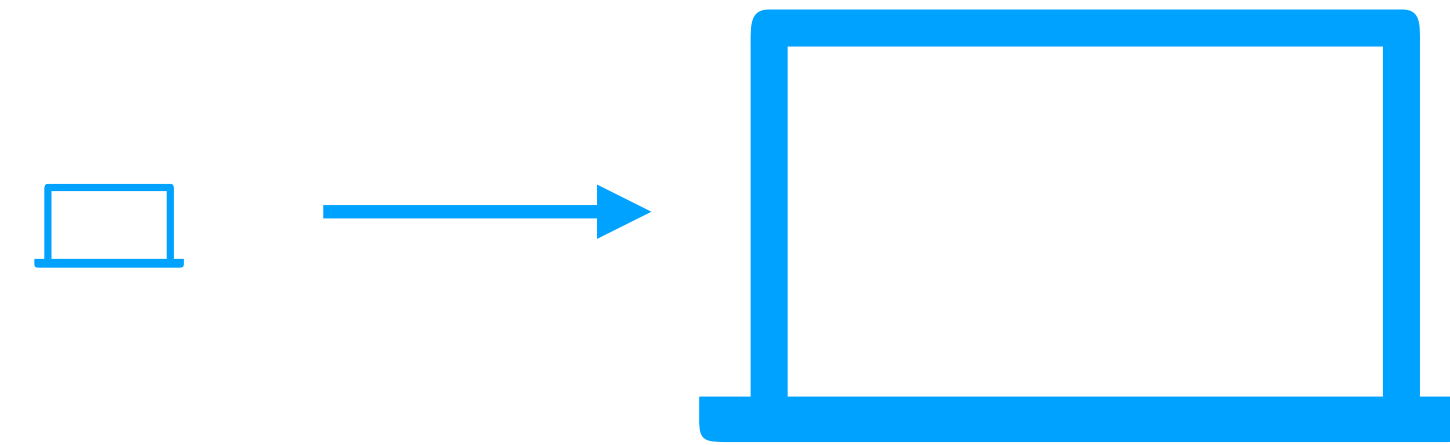
https://en.wikipedia.org/wiki/Tensor_Processing_Unit

- Tensorflow Processing Unit (TPU)
- Developed by Google for neural networks
- Supports matrix multiplication operations
- Training & Test
- Precision of operations is lower compared to GPUs
- Multiple TPUs per “machine”

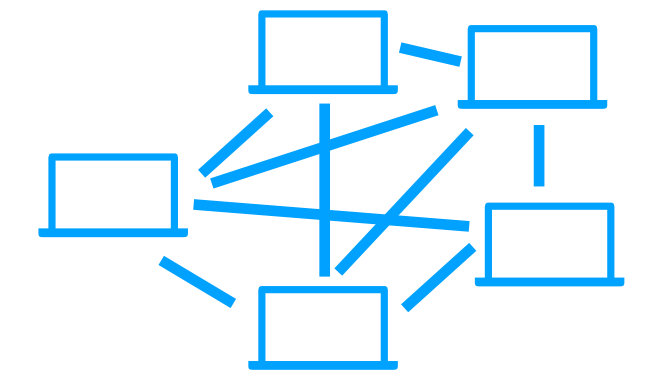
https://colab.research.google.com/github/lcharlin/80-629/blob/master/week10-ParallelComputations/CPU_GPU_TPU.ipynb

- With a few small changes, we can start using a GPU or even a TPU (Google platform only)

Distributed Computing



- Faster computers can help
- What about a large of “slow” computers working together?
- Divide the computation into small problems
 1. All (slow) computers solve a small problem at the same time
 2. Combine the solution of small problems into initial solution

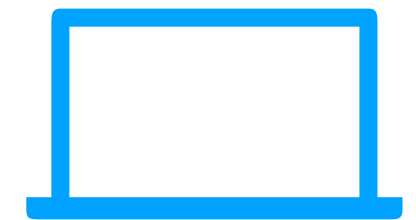


Building our intuition with a simple example

- You are tasked with counting the number of houses in Montreal

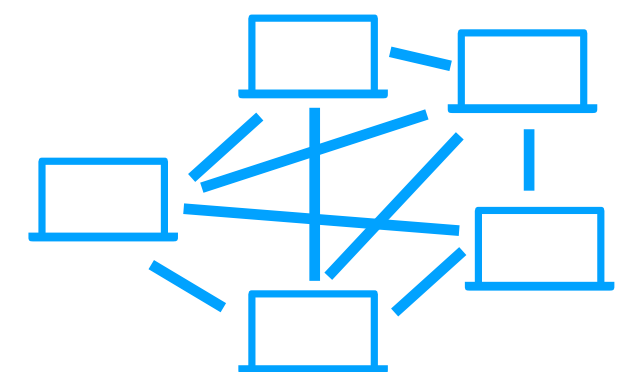
1. Centralized (single computer):

- Ask a marathon runner to jog around the city and count
- Build a system to count houses from satellite imagery



2. Distributed (many computers):

- Ask 1,000 people to each count houses from a small geographical area
- Once they are done they report their result at your HQ



Distributed Computing using MapReduce

Outline

- **MapReduce**
 - **Fundamentals and bag-of-words example**
- **Spark**
 - **Fundamentals & MLlib**

MapReduce

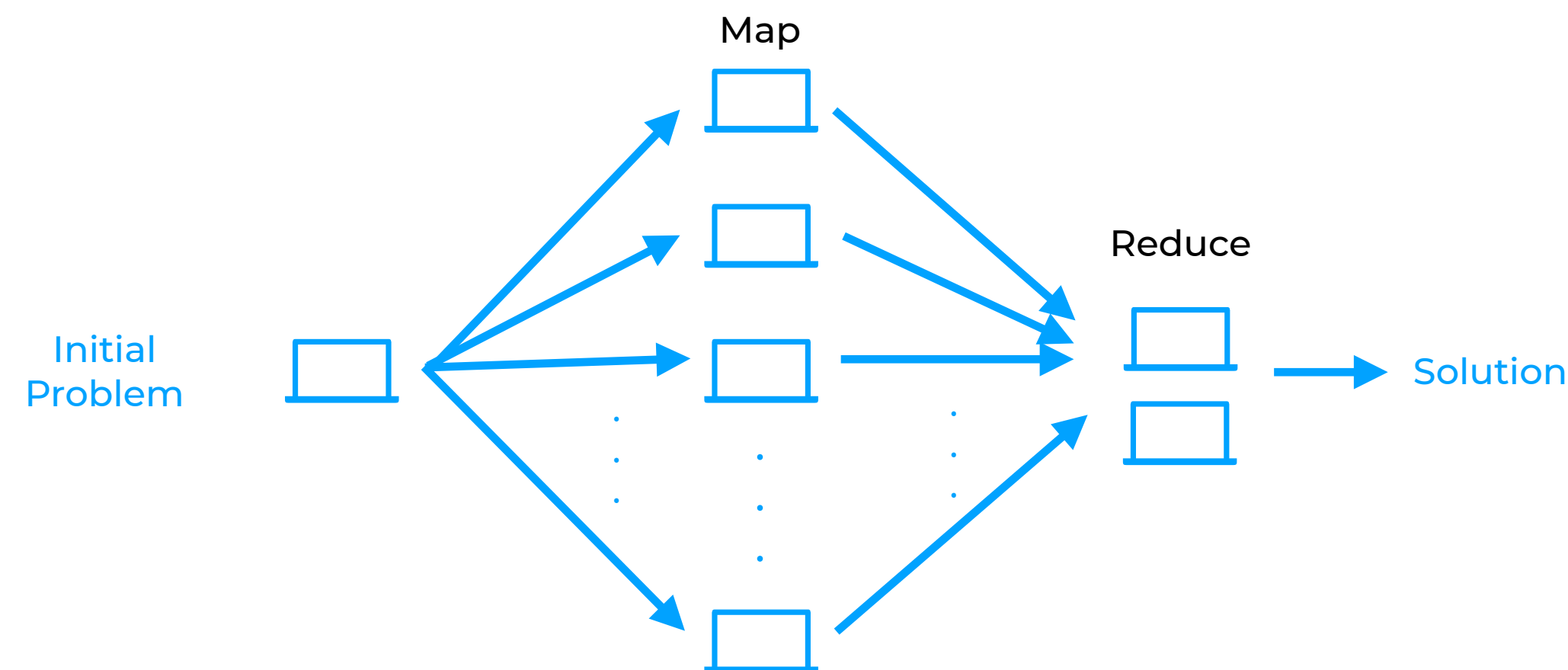
- From Google engineers

**“MapReduce: Simplified Data Processing on Large Clusters”,
Jeffrey Dean and Sanjay Ghemawat, 2004**

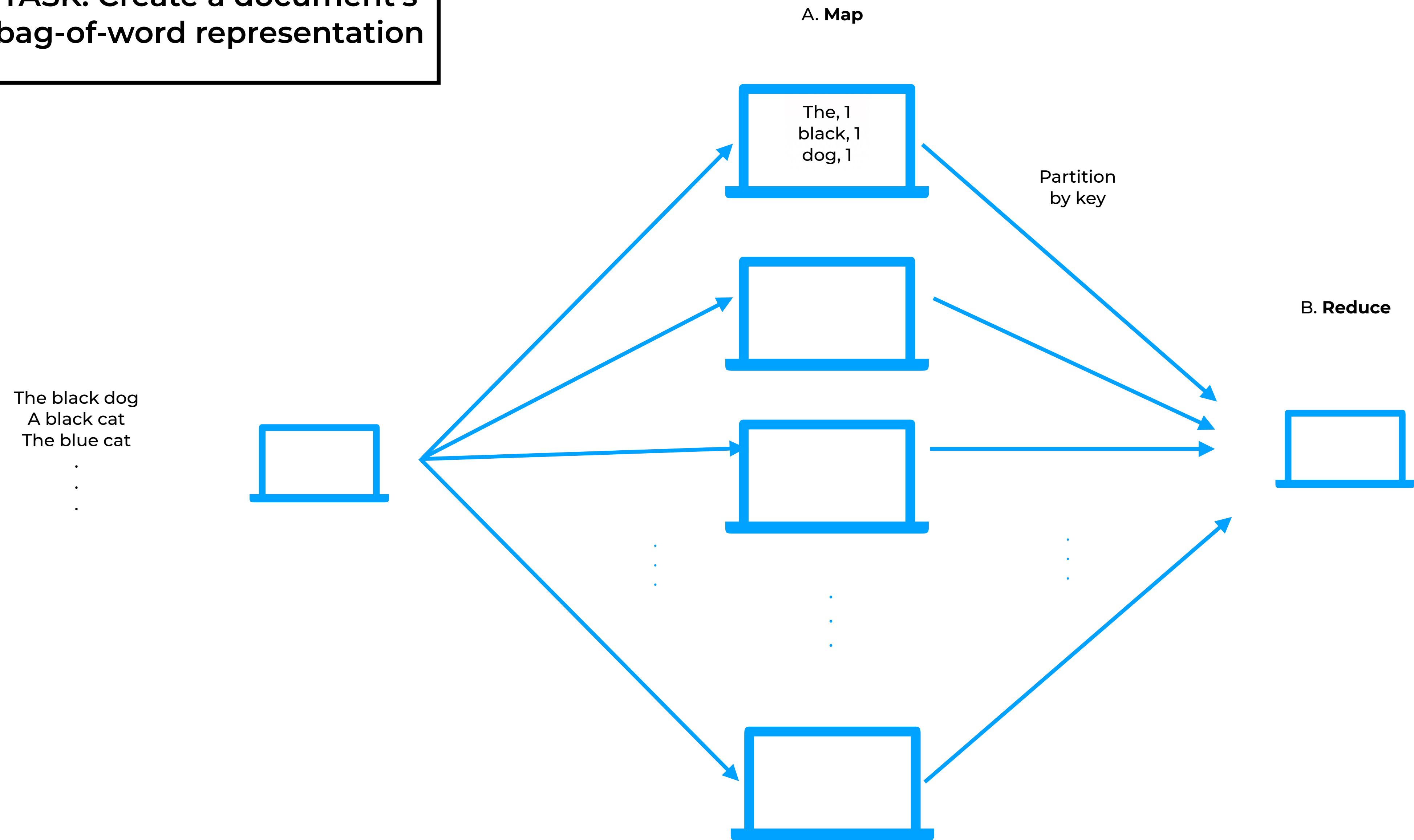
- Now also known as (Apache) Hadoop
- Google built large-scale computation from commodity hardware
- Specific distributed interface
 - Useful for algorithms that can be expressed using this interface

MapReduce

- Two types of tasks:
 - A. Map: Solve a subproblem (filtering operation)
 - Sorting is then performed
 - B. Reduce: Combine the results of map workers (summary operation)



TASK: Create a document's bag-of-word representation



Some details

- Typically the number of subproblems is higher than the number of available machines
 - ~linear speed-up wrt to the number of machines
- If a node crashes, need to recompute its subproblem
- Input/Output
 - Data is read from disk when beginning map/reduce
 - Data is written to disk at the end of map/reduce

MapReduce is quite versatile

- When I was at Google the saying was (roughly):

“If your problem cannot be framed as MapReduce you haven’t thought hard enough about your problem.”

- A few examples of “map-reduceable” problems:
 - Intuition: Your problem needs to be decomposable into map functions and reduce functions
 - Sorting, filtering, distinct values, basic statistics
 - Finding common friends, sql-like queries, sentiment analysis

MapReduce for machine learning

1. Training linear regression

- Reminder: there is a closed-form solution

$$\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

$$\mathbf{w} = \left(\sum_{ij} \mathbf{x}_i^\top \mathbf{x}_j \right)^{-1} \left(\sum_i \mathbf{x}_i^\top \mathbf{Y}_i \right)$$

- Each term in the sums can be computed independently

A. Map

$$\mathbf{x}_0^\top \mathbf{x}_1$$

2. Other models we studied have a closed form solution (e.g., Naive Bayes and LDA)

3. Hyper-parameter search

- A neural network with 2 hidden layers and 5 hidden units per layer and another with 3 hidden layers and 10 hidden units

Shortcomings of MapReduce

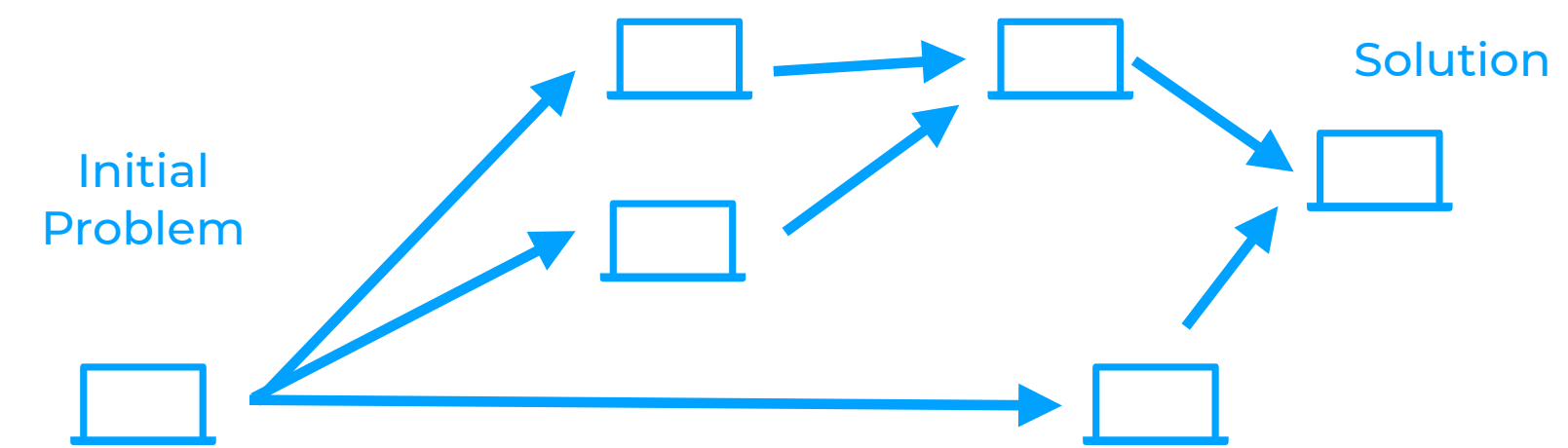
- Many models are fitted with iterative algorithms
 - Gradient descent:
 1. Find the gradient for the current set parameters
 2. Update the parameters with the gradient
- Not ideal for MapReduce
 - Would require several iterations of MapReduce
 - Each time the data is read/written from/to the disk

Distributed computing using Apache Spark

(Apache) Spark

- Advantages over MapReduce

1. Less restrictive computations graph (DAG instead of Map then Reduce)



- Doesn't have to write to disk in-between operations

2. Richer set of transformations

- map, filter, cartesian, union, intersection, distinct, etc.

3. In-memory processing

Spark History

- Started in Berkeley's AMPLab (2009)
- Version 1.0 2014
 - Based on Resilient Distributed Datasets (RDDs)
- Version 2.0 June 2016
 - V2.3 February 2018, V2.4.4 September 2019, V3.5.4 September 2024
- Python + Spark: pySpark
- Good (current) documentation:
 1. Advanced Analytics with Spark, 2nd edition (2017).
 2. Project docs: <https://spark.apache.org/docs/latest/>

DataFrames

- An extra abstraction on top of RDDs (resilient distributed datasets)
 - Encodes rows as a set of columns
 - Each column has a defined type
 - Useful for (pre-processed) machine learning datasets
- Same name as `data.frame` (R) or `pandas.DataFrame`
 - Similar type of abstraction but for distributed datasets
- Two types of operations (for our needs): transformers, estimators.

Spark's “Hello World”

DataFrame → `data = spark.read.format("libsvm").load("hdfs://...")`

`model = LogisticRegression(regParam=0.01).fit(data)`

Estimator

```
graph BT; DF[DataFrame] --> Code1["data = spark.read.format('libsvm').load('hdfs://...')"]; Code1 --> Code2["model = LogisticRegression(regParam=0.01).fit(data)"]; Estimator[Estimator] --> Code2;
```

The diagram illustrates the Spark workflow. A **DataFrame** is used to load data from HDFS into a `data` variable. This `data` variable is then passed to the `fit` method of a `LogisticRegression` estimator. The `Estimator` box points to the `LogisticRegression` class in the code, indicating its role in the model fitting process.

Parallel gradient descent

- Logistic Regression

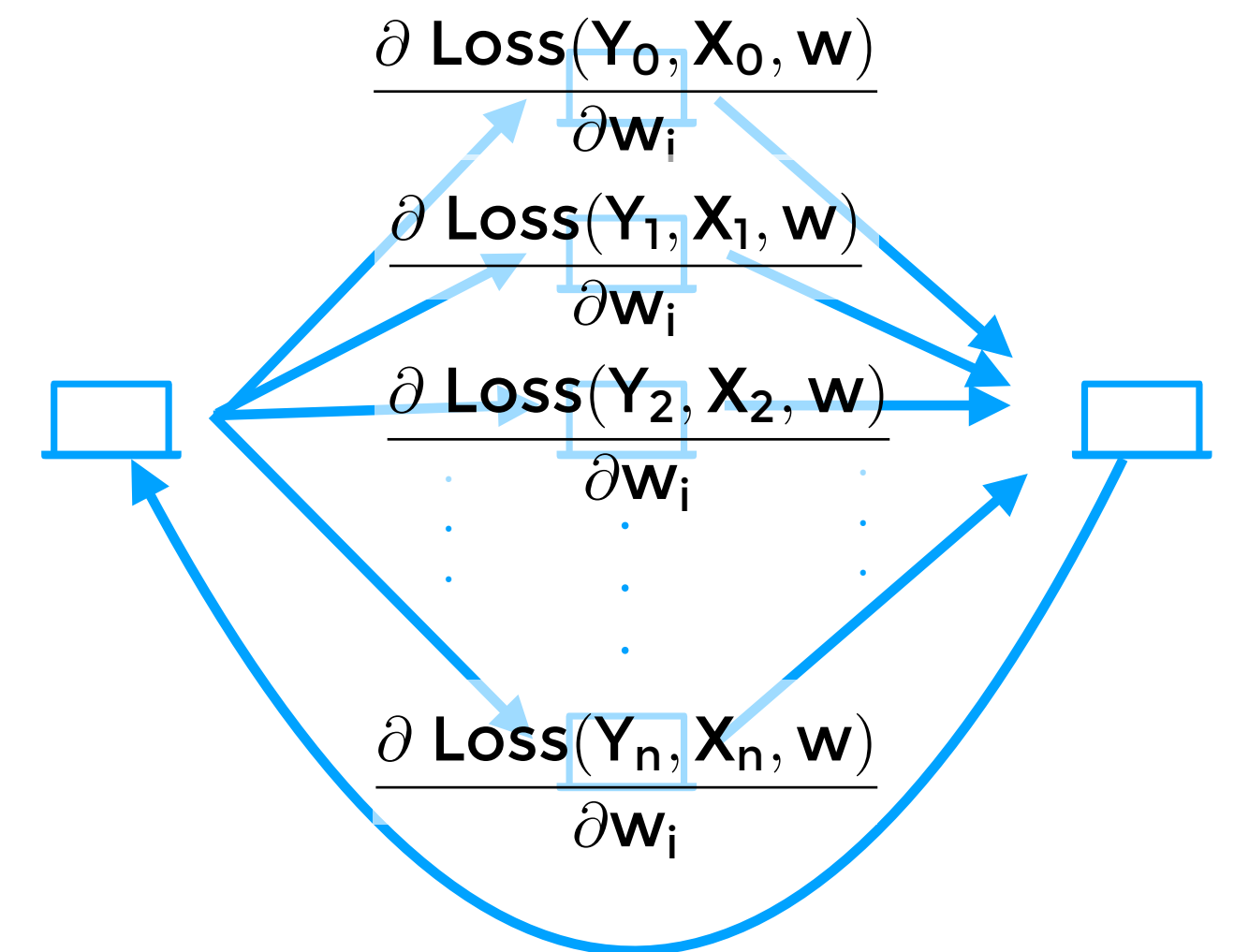
$$y = \frac{1}{1 + \exp(-w_0 - w_1x_1 - w_2x_2 - \dots - w_px_p)}$$

- No closed-form solution, can use gradients

$$\frac{\partial \text{Loss}(Y, X, w)}{\partial w_i}$$

- Loss functions are often decomposable

$$\frac{\partial \sum_j \text{Loss}(Y_j, X_j, w)}{\partial w_i}$$



ML setup

1.

Load your data as a
Spark DataFrame

2.

Machine Learning Library (MLlib) Guide

MLlib is Spark's machine learning (ML) library. Its goal is to make practical machine learning scalable and easy. At a high level, it provides tools such as:

- ML Algorithms: common learning algorithms such as classification, regression, clustering, and collaborative filtering
- Featurization: feature extraction, transformation, dimensionality reduction, and selection
- Pipelines: tools for constructing, evaluating, and tuning ML Pipelines
- Persistence: saving and load algorithms, models, and Pipelines
- Utilities: linear algebra, statistics, data handling, etc.

Classification and Regression - RDD-based API

The `spark.mllib` package supports various methods for [binary classification](#), [multiclass classification](#), and [regression analysis](#). The table below outlines the supported algorithms for each type of problem.

Problem Type	Supported Methods
Binary Classification	linear SVMs, logistic regression, decision trees, random forests, gradient-boosted trees, naive Bayes
Multiclass Classification	logistic regression, decision trees, random forests, naive Bayes
Regression	linear least squares, Lasso, ridge regression, decision trees, random forests, gradient-boosted trees, isotonic regression

<https://spark.apache.org/docs/latest/ml-guide.html>

Takeaways

- Specialized hardware
 - The go-to method for neural nets training
- Distributed computing is useful:
 - for large-scale data (e.g., data that does not fit on a single computer)
 - for faster computing (when you have multiple available computers)
- Current frameworks (e.g., spark) offer easy access to popular ML models + algorithms
 - Useful speedups by decomposing the computation into a number of identical smaller pieces
- Still requires some engineering/coding