

---

# **CS 267**

## **Lecture 17a: Parallel Machine Learning, Part 1**

### **(Supervised Learning)**

**Aydin Buluc**

**<https://sites.google.com/lbl.gov/cs267-spr2019/>**

Slide acks: Amir Gholami, Forrest Landola, Prabhat, and Yang You

## **Outline of the lecture**

---

### **Today: Part 1, Intro and Supervised Learning**

- Machine Learning & Parallelism Intro
- Neural Network Basics
- Scaling Deep Neural Network Training
- Support Vector Machines

### **Thursday: Part 2, Unsupervised Learning**

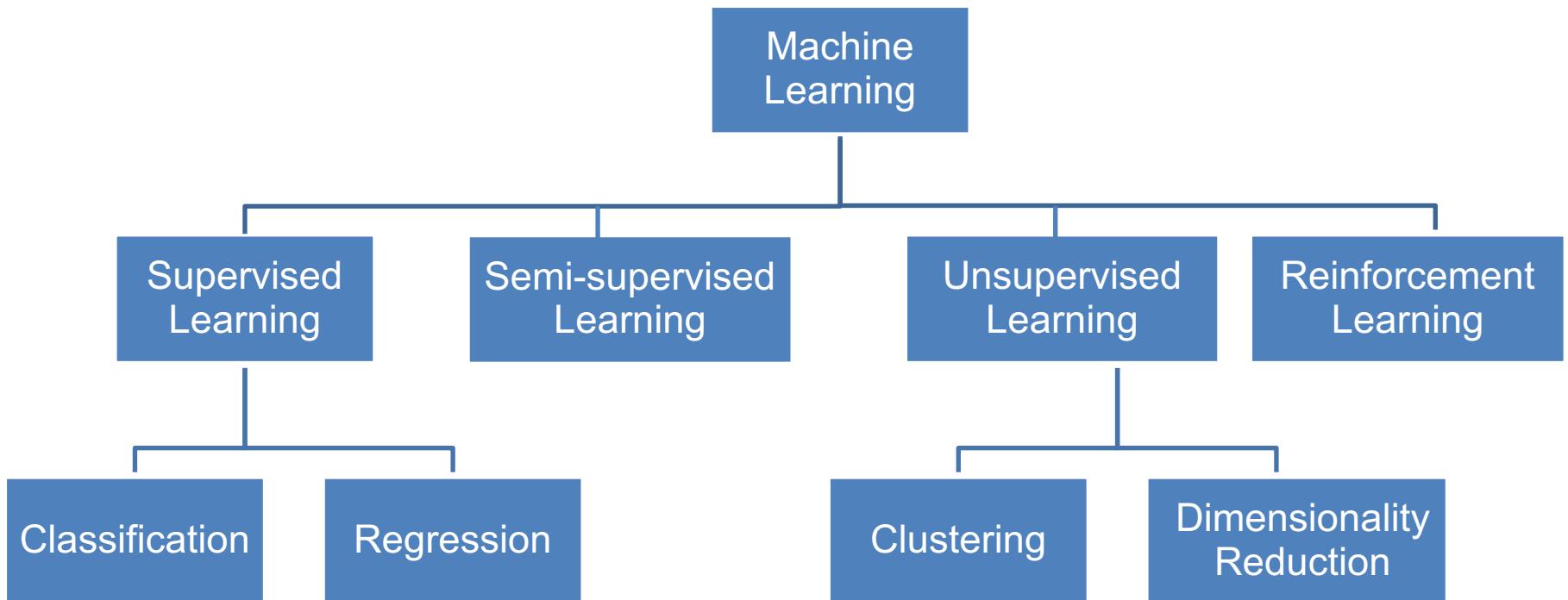
- Non-Negative Matrix Factorization
- Spectral and Markov Clustering
- Sparse Inverse Covariance Matrix Estimation

# Machine Learning Classes and Tasks

---

The central question in Machine Learning:

*“How can we build computer systems that automatically improve with experience, and what are the fundamental laws that govern all learning processes?”*



# Machine Learning Sources of Confusion

---

**Method vs. Task:** A common confusion is between specific learning methods and learning tasks.

- Example #1: Principal Component Analysis is a method for dimensionality reduction
- Example #2: Support Vector Machines are methods used for supervised learning tasks.

Another confusion comes from ***optimization techniques*** vs. **learning methods**.

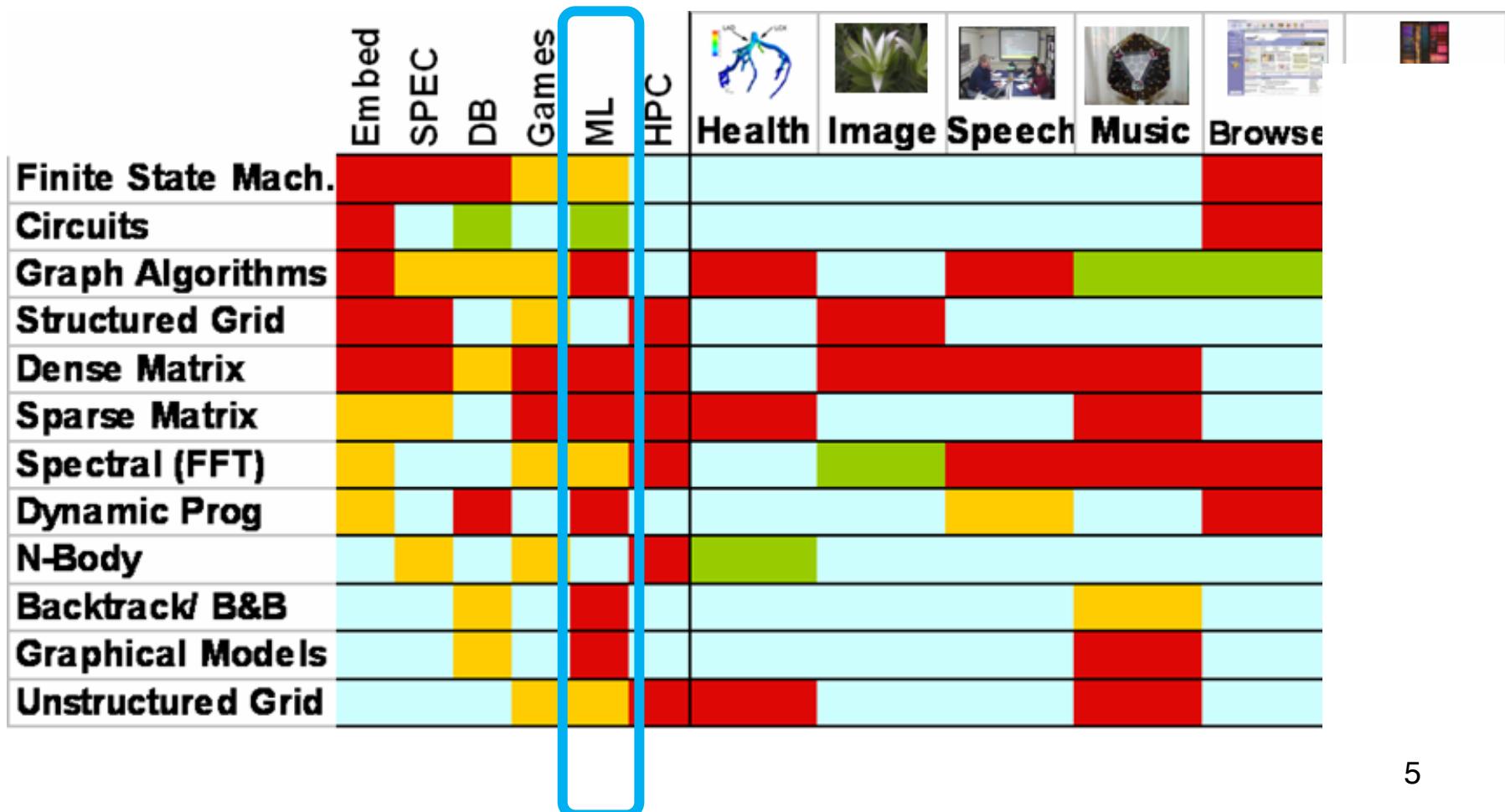
- Example #1: Sequential Minimal Optimization is an optimization technique to train Support Vector Machines
- Example #2: Stochastic Gradient Descent is a popular optimization technique to train Neural Networks.

# Motifs

---

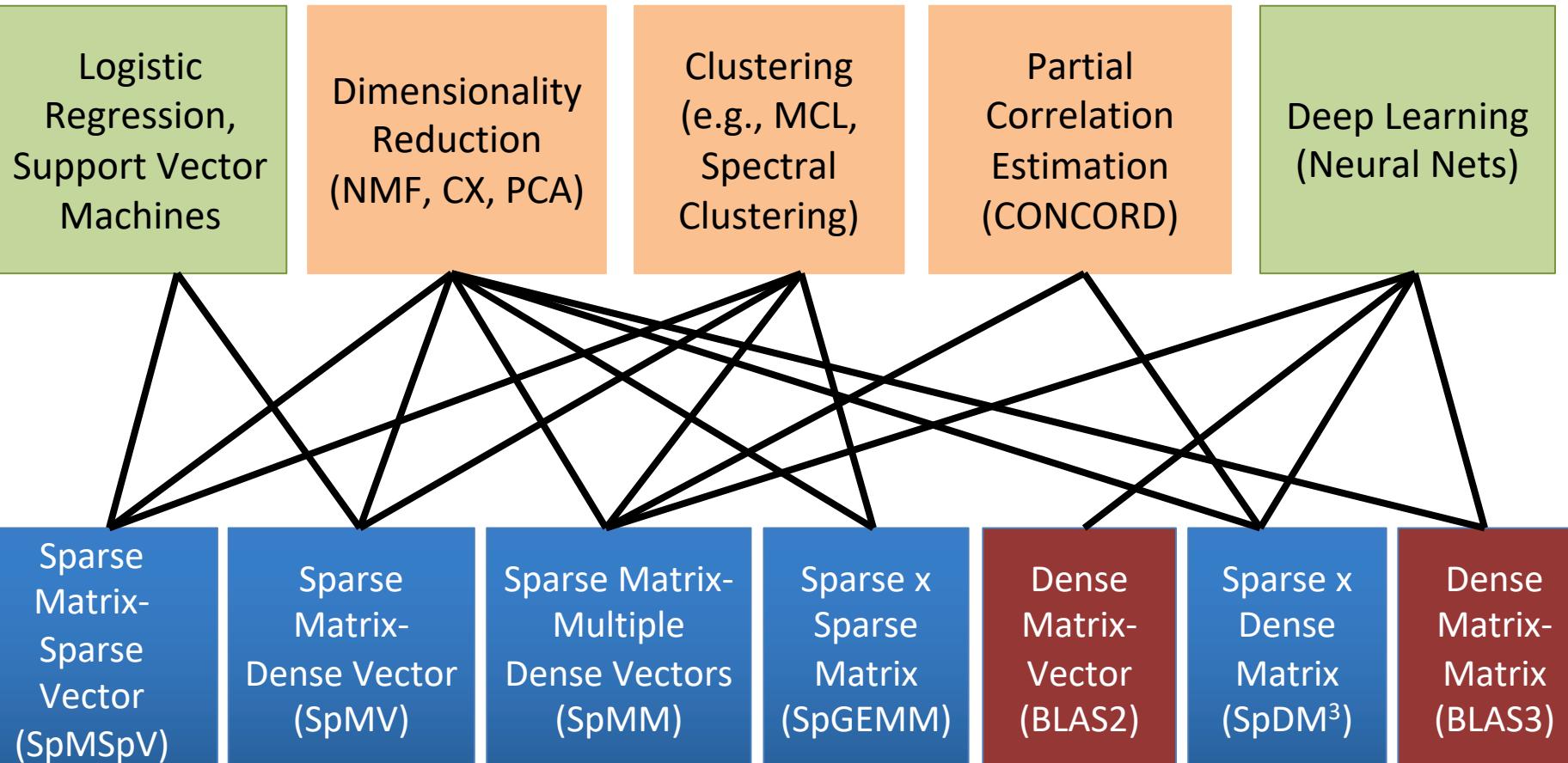
The Motifs (formerly “Dwarfs”) from  
“The Berkeley View” (Asanovic et al.)

**Motifs** form key computational patterns



# Machine Learning relies a lot on Linear Algebra

Higher-level machine learning tasks



Graph/Sparse/Dense BLAS functions (in increasing arithmetic intensity)

# Parallelism in Machine Learning

---

**Implicit Parallelization:** Keep the overall algorithm structure (the sequence of operations) intact and parallelize the individual operations.

Example: parallelizing the BLAS operations in previous figure

- + Often achieves exactly the same accuracy (e.g., model parallelism in DNN training)
- Scalability can be limited if the critical path of the algorithm is long

**Explicit Parallelization:** Modify the algorithm to extract more parallelism, such as working on individual pieces whose results can later be combined

Examples: CA-SVM and data parallelism in DNNs

- + Significantly better scalability can be achieved
- No longer the same algorithmic properties (e.g. HogWild!).

## Outline of the lecture

---

### **Today: Part 1, Intro and Supervised Learning**

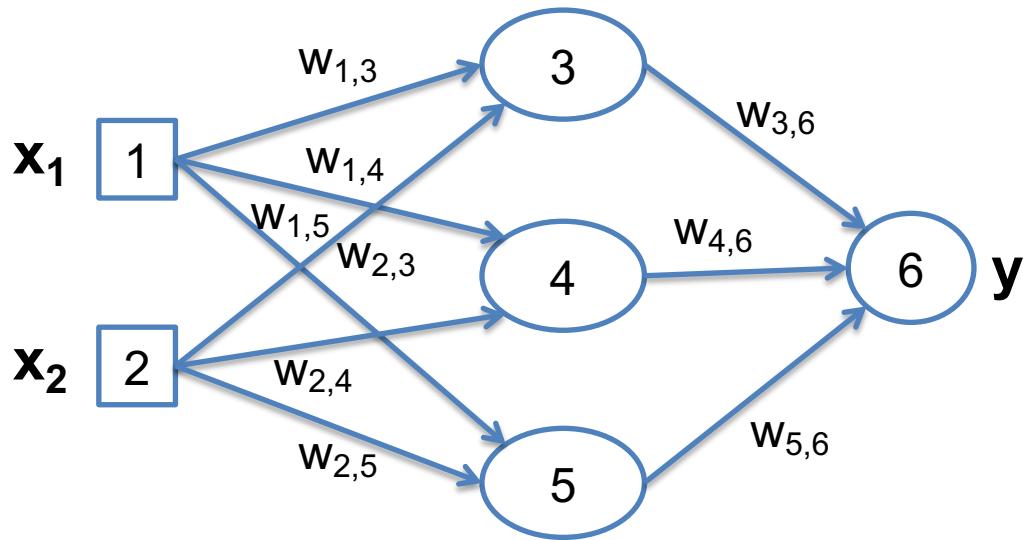
- Machine Learning & Parallelism Intro
- Neural Network Basics
- Scaling Deep Neural Network Training
- Support Vector Machines

### **Thursday: Part 2, Unsupervised Learning**

- Non-Negative Matrix Factorization
- Spectral and Markov Clustering
- Sparse Inverse Covariance Matrix Estimation

# Training Neural Networks

- Training is to **adjust the weights ( $W$ )** in the connections of the neural network, in order to change the function it represents.



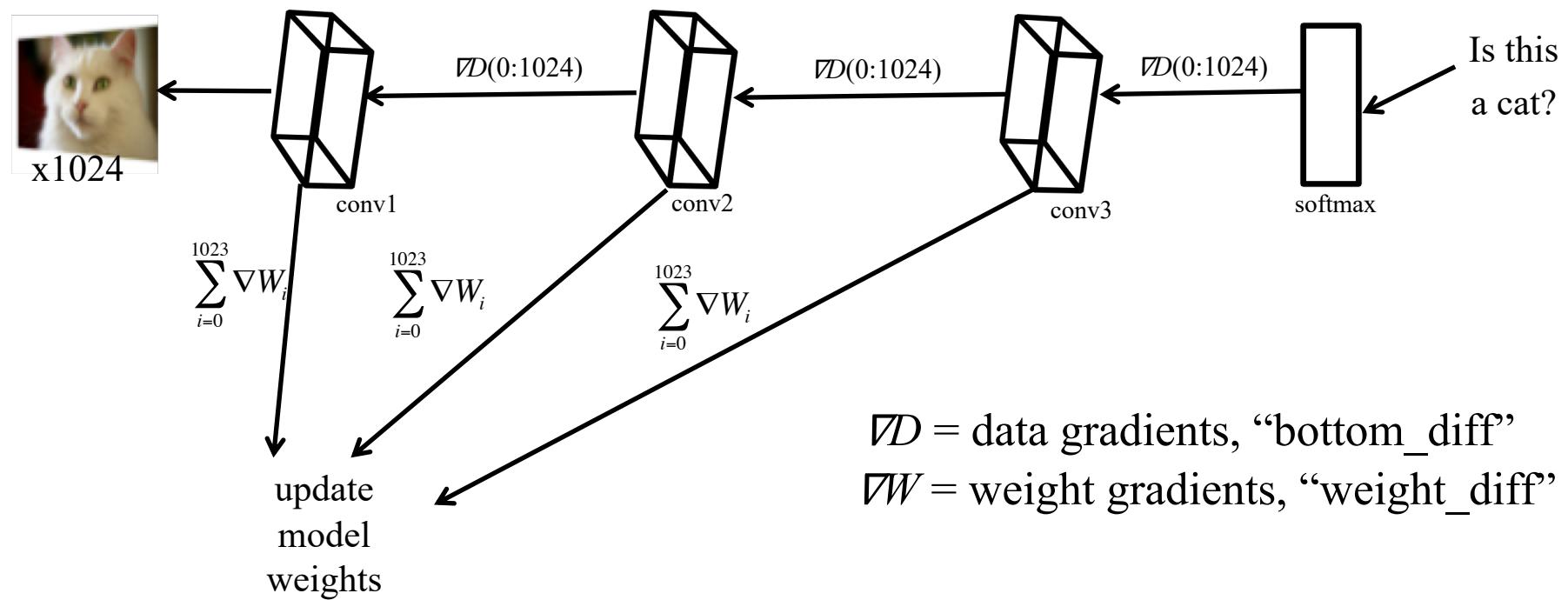
Only parameters are weights for simplicity (i.e. ignore bias parameters)

$W$ : the matrix of weights

A “shallow” neural network with only one hidden layer (nodes 3,4,5), two inputs and one output.

# Deep Neural Network Training

---

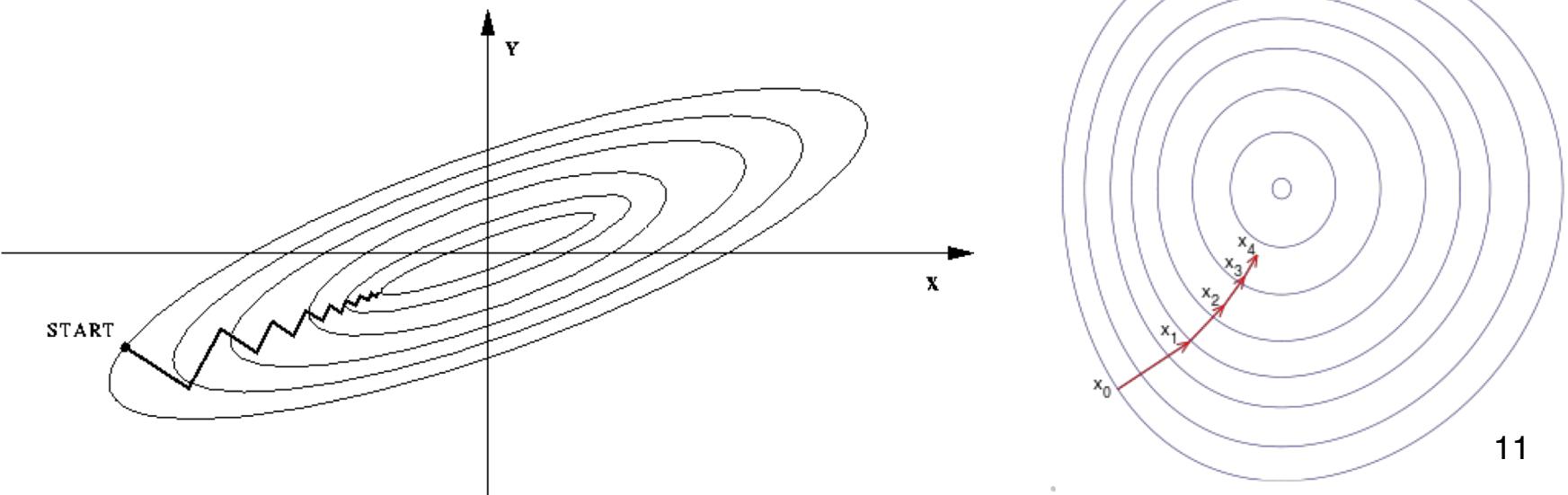


# Gradient Descent

---

$$W^{t+1} \leftarrow W^t - \alpha \cdot \nabla_W f(W^t, x)$$

- Also called the steepest descent algorithm
- In order to minimize a function, move towards the opposite direction of the gradient at a rate of  $\alpha$ .
- $\alpha$  is the step size (also called the learning rate)
- Used as the ***optimization backend*** of many other machine learning methods (example: NMF)



## Stochastic Gradient Descent (SGD)

Assume  $f(W^t, x) = \frac{1}{n} \sum_{i=1}^n f_i(W^t, x)$

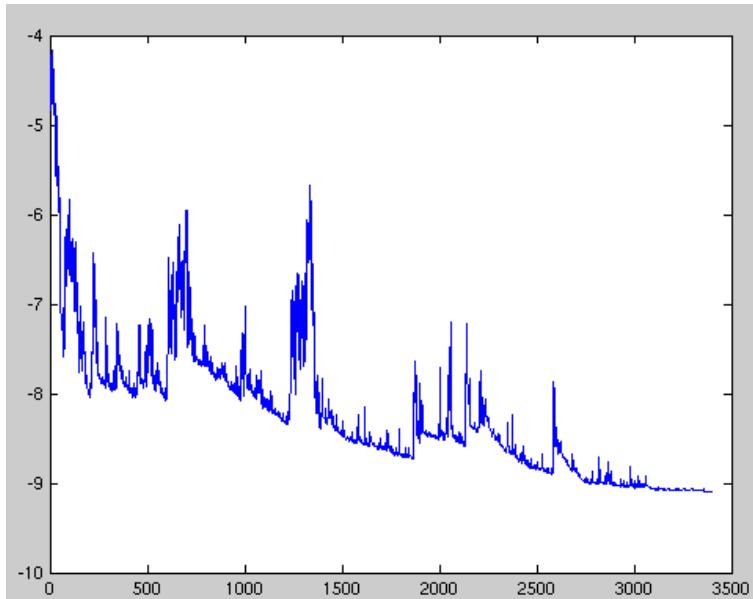
$$W^{t+1} \leftarrow W^t - \alpha \cdot \nabla_W f_i(W^t, x)$$

Pure SGD: compute gradient using 1 sample

$$W^{t+1} \leftarrow W^t - \alpha \cdot \frac{1}{b} \sum_{i=k+1}^{k+b} \nabla_W f_i(W^t, x)$$

Mini-batch: compute gradient using b samples

$f$  is not going down for every iteration



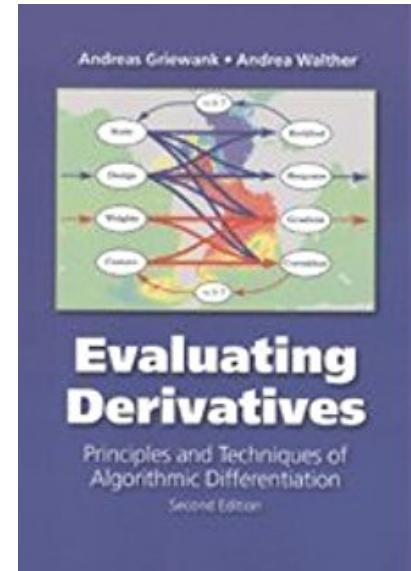
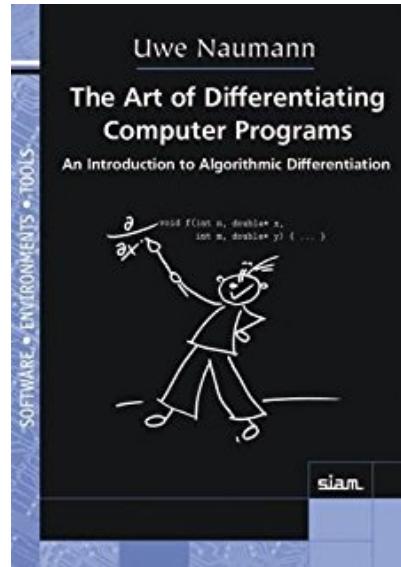
- Actually the name is a misnomer, *this is not a “descent” method*
- But we will stick to it anyway to avoid confusion.
- Performance and parallelism requires batch training
- Larger batch sizes hurt convergence as they get trapped easily
- SGD escapes sharp local minima due to its “noisy” gradients

# Training Neural Networks

---

- Training is performed using an optimization algorithm like SGD
- **SGD needs derivatives.**
- The algorithm **to compute derivatives** on a neural network is called **back-propagation**.
- The back-propagation algorithm is **not a training algorithm**
- **Idea:** Repeated application of the chain rule from calculus

Back-propagation is just  
a special case of the  
**reverse mode**  
**automatic/algorithmic**  
**differentiation**



# Parallelization Opportunities

---

## 1. Data parallelism

Distribute the input (sets of images, text, audio, etc.)

### a) Batch parallelism

- Distribute each full sample to a different processor
- *When people mention data parallelism in literature, this is what they mean 99% of the time*

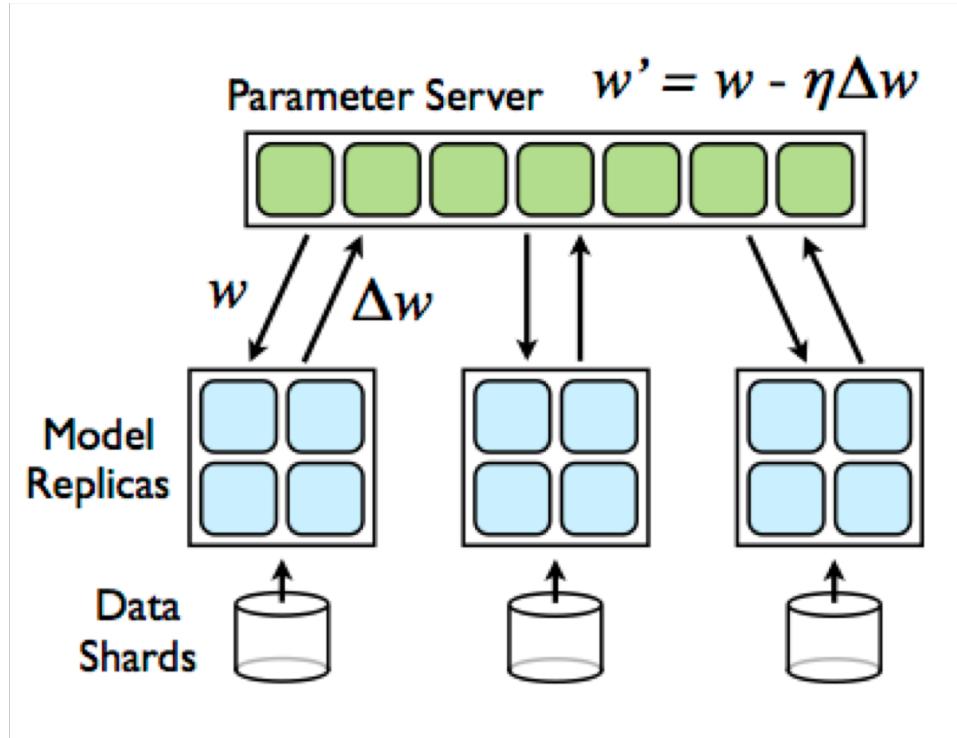
### b) Domain parallelism

- Subdivide samples and distribute parts to processors.
- A previously unexplored approach to parallelism.

## 2. Model parallelism:

Distribute the neural network (i.e. its weights)

## Batch Parallelism #1



Parameter server is some sort of master process

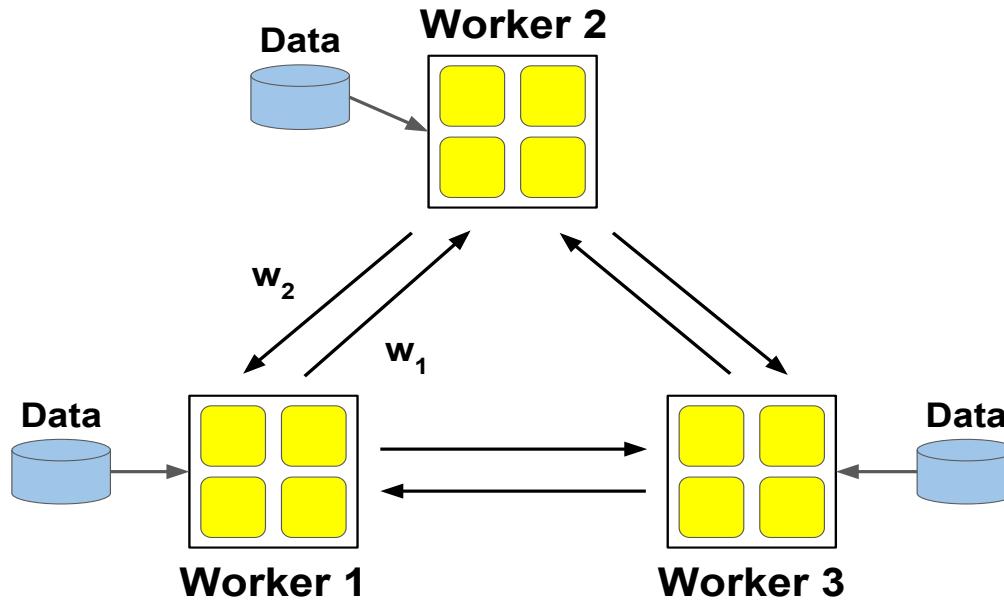
- The fetching and updating of gradients in the parameter server can be done either ***synchronously*** or ***asynchronously***.
- Both has pros and cons. Over-synchronization hurts performance where asynchrony is non-reproducible and might hurt convergence

Dean, Jeffrey, et al. "Large scale distributed deep networks." Advances in neural information processing systems. 2012.

## Batch Parallelism #2

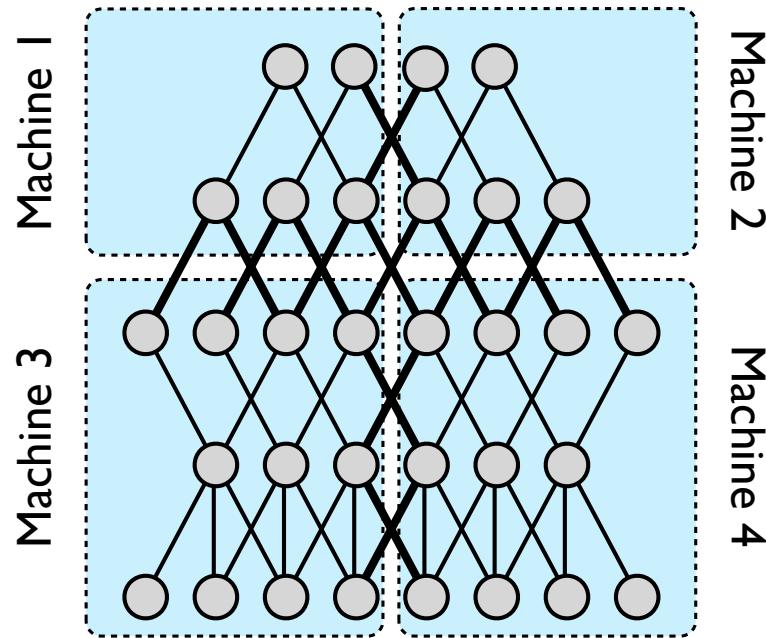
Options to avoid the parameter server bottleneck

1. **For synchronous SGD:** Perform all-reduce over the network to update gradients (good old MPI\_Allreduce)
2. **For asynchronous SGD:** Peer-to-peer gossiping



Peter Jin, Forrest Landola, Kurt Keutzer, "How to scale distributed deep learning?"  
NIPS ML Sys 2016

# Model Parallelism



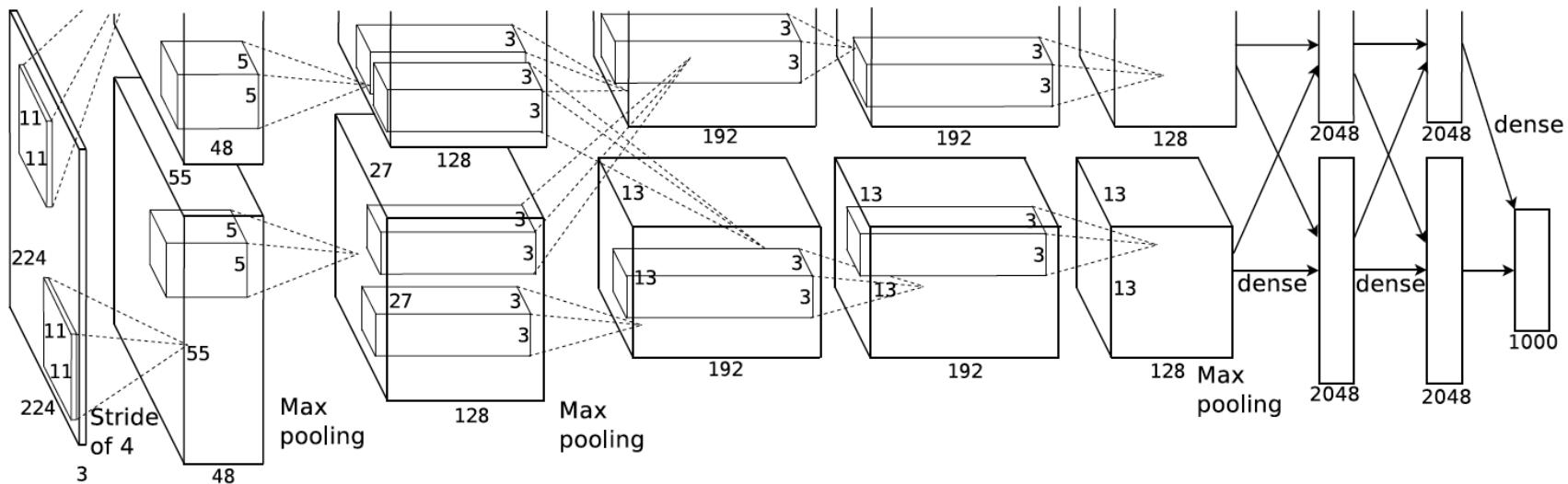
Another misleading (inter layer parallelism is better called pipelining) but nice figure!

- Interpretation #1: Partition your neural network into processors
- Interpretation #2: Perform your matrix operations in parallel

Dean, Jeffrey, et al. "Large scale distributed deep networks." Advances in neural information processing systems. 2012.

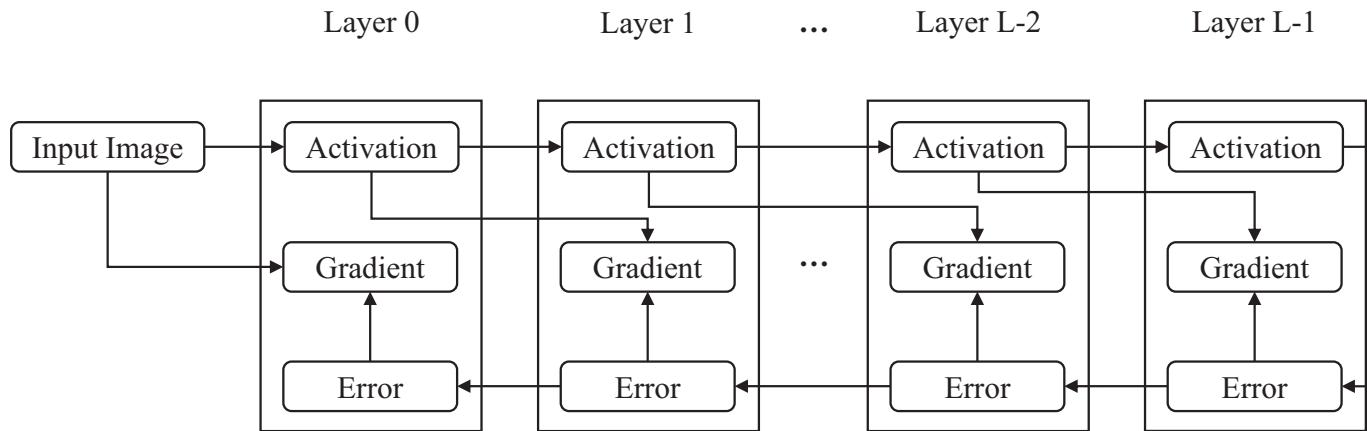
# AlexNet architecture

- ° Reminder: a convolutional deep neural network still often has fully connected layers



## The overlapping opportunities

- In backpropagation, the errors are propagated from the last layer to the first layer and data dependency exists between any two consecutive layers.
- **In contrast, the gradients in different layers are independent of each other.**



Activations are propagated from left to right in forward stage; errors are propagated from right to left in backpropagation stage. Gradients are computed using the activations and errors. Arrows indicated data dependencies.

## SGD training of NNs as matrix operations

---

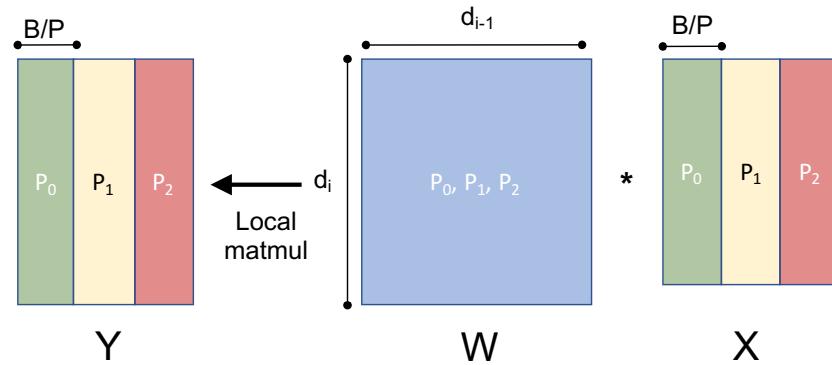
$$N \begin{matrix} M \\ \text{W: weights} \end{matrix} \times \begin{matrix} B \\ X_{in} \end{matrix} = \begin{matrix} B \\ X_{out} \end{matrix} N$$

N = the number of outputs  
M = the number of inputs  
B = the size of the minibatch

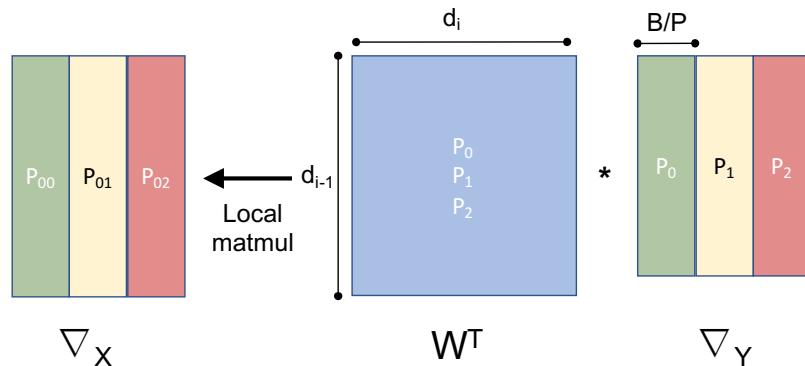
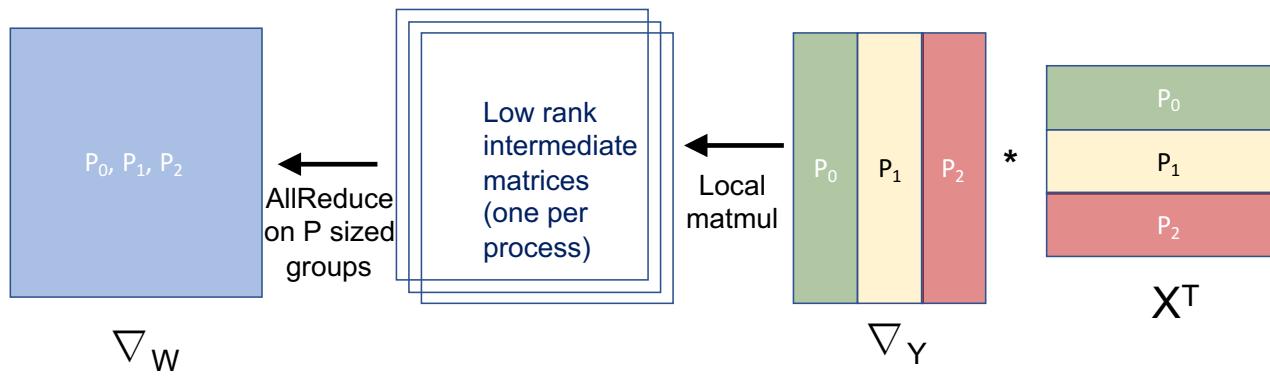
### The impact to parallelism:

- $\mathbf{W}$  is replicated to processor, so it doesn't change
- $\mathbf{X}_{in}$  and  $\mathbf{X}_{out}$  gets skinnier if we only use data parallelism, i.e. distributing  $b=B/p$  mini-batches per processor
- GEMM performance suffers as *matrix dimensions get smaller and more skewed*
- **Result:** Data parallelism can hurt single-node performance

# Data Parallel SGD training of NNs as matrix operations



1. Which matrices are replicated?
2. Where is the communication?
3. Which steps can be overlapped?

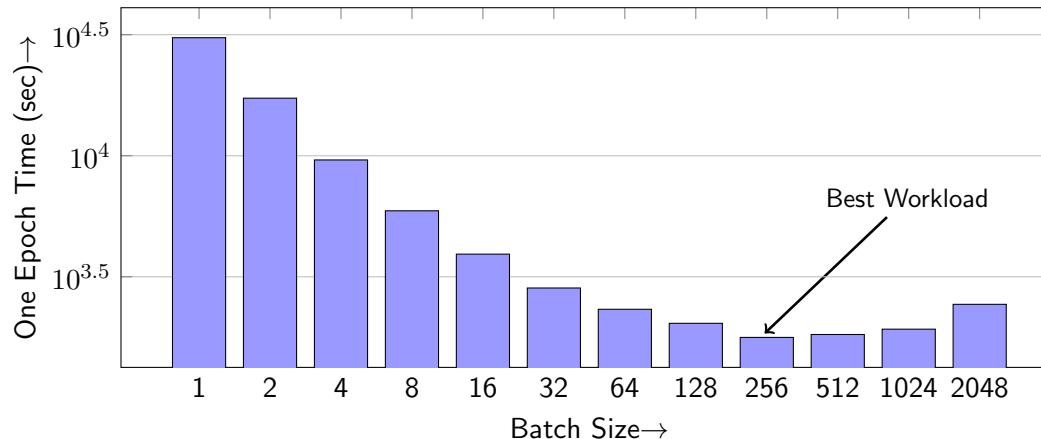


$\nabla_Y = \partial L / \partial Y$  = how did the lost function change as output activations change?  
 $\nabla_X = \partial L / \partial X$   
 $\nabla_W = \partial L / \partial W$

## Batch Parallel Strong Scaling

---

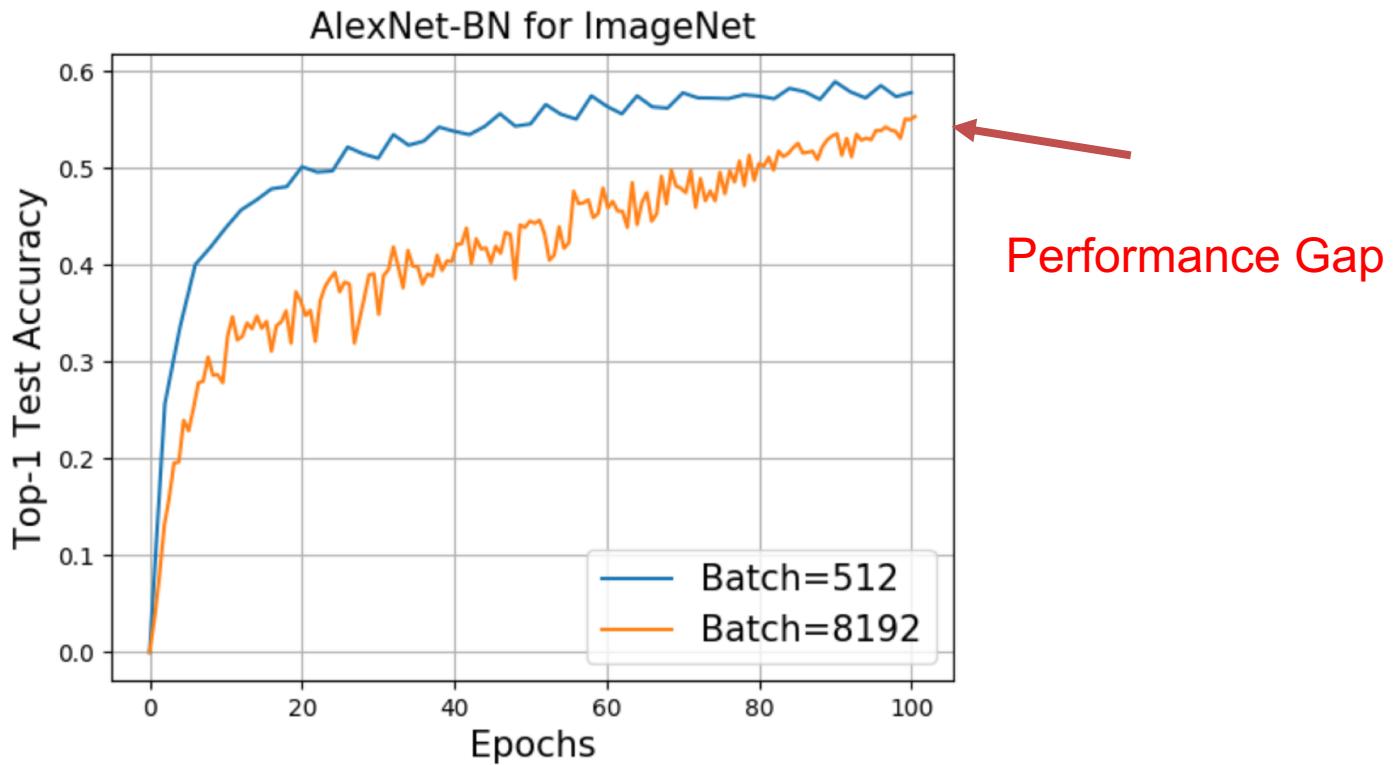
- ° Per-iteration communication cost of batch parallelism is independent of the batch size:
  - larger batch → less communication per epoch (full pass over the data set)
- ° But processor utilization goes down significantly for  $P \gg 1$ 
  - Result: Batch parallel has poor strong scaling



One epoch training time of AlexNet computed on a single KNL

## Problems with Batch Parallelism

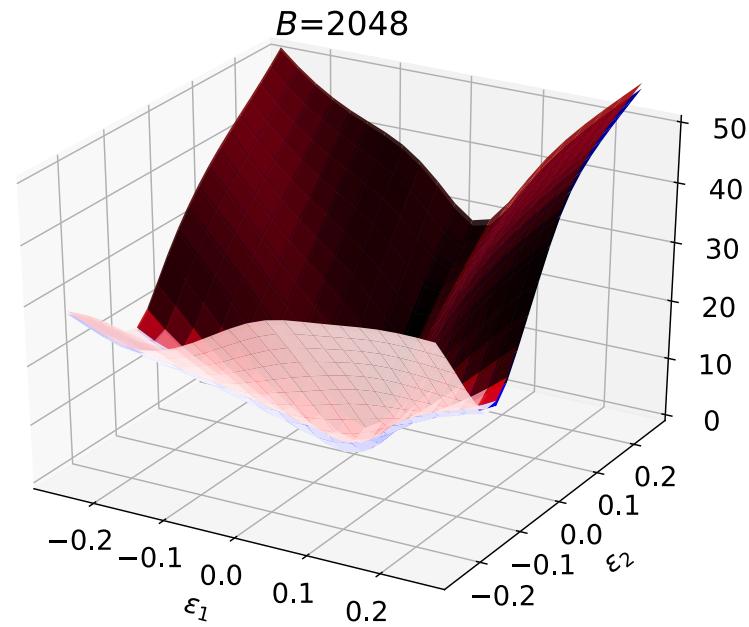
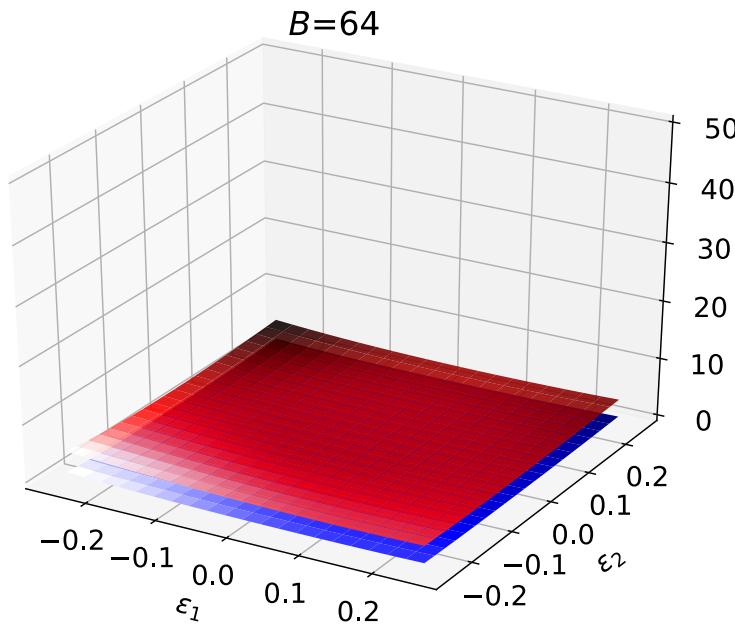
- Batch parallel scaling is limited to B
  - Larger Batch  $\rightarrow$  higher strong scaling efficiency
- But SGD does not perform well for large batch



Ginsburg, Boris, Igor Gitman, and Yang You. "Large Batch Training of Convolutional Networks with Layer-wise Adaptive Rate Scaling." (2018).

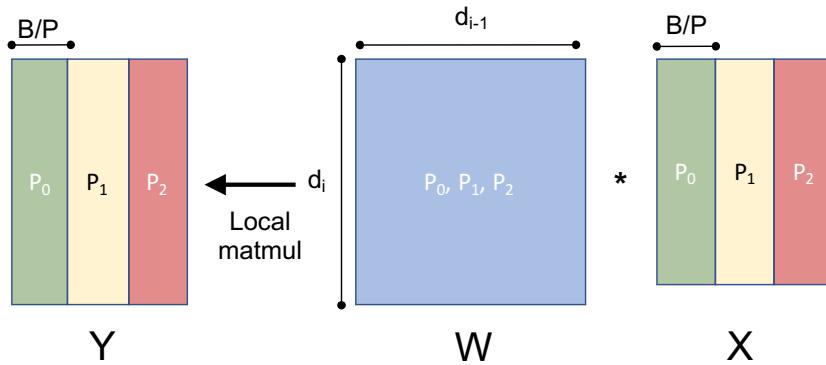
## Problems with Batch Parallel

- The objective function implicitly changes in large batch SGD

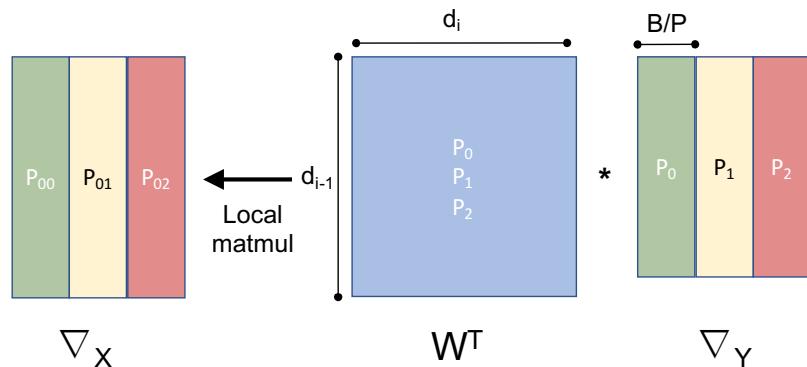
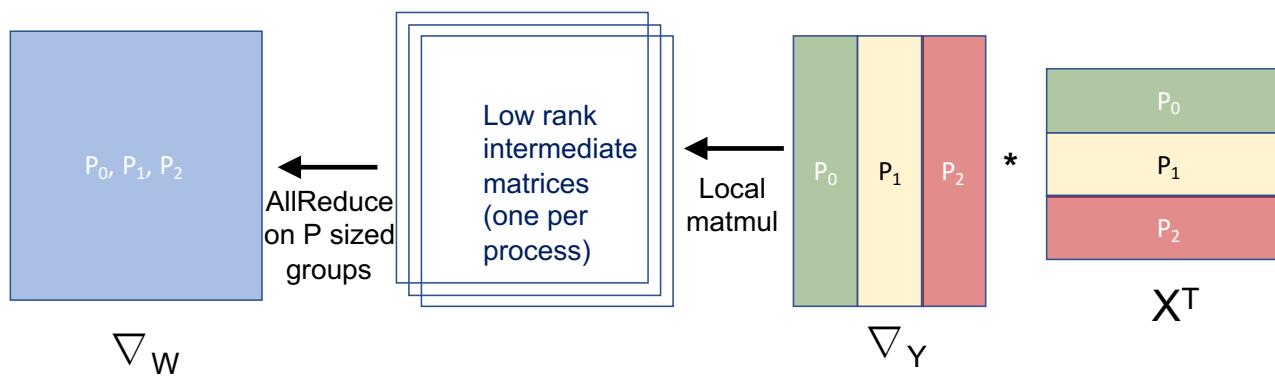


Landscape of loss function for small and large batch at the end of training. Large batch converges to a **suboptimal point**

# Batch Parallel SGD training of NNs as matrix operations

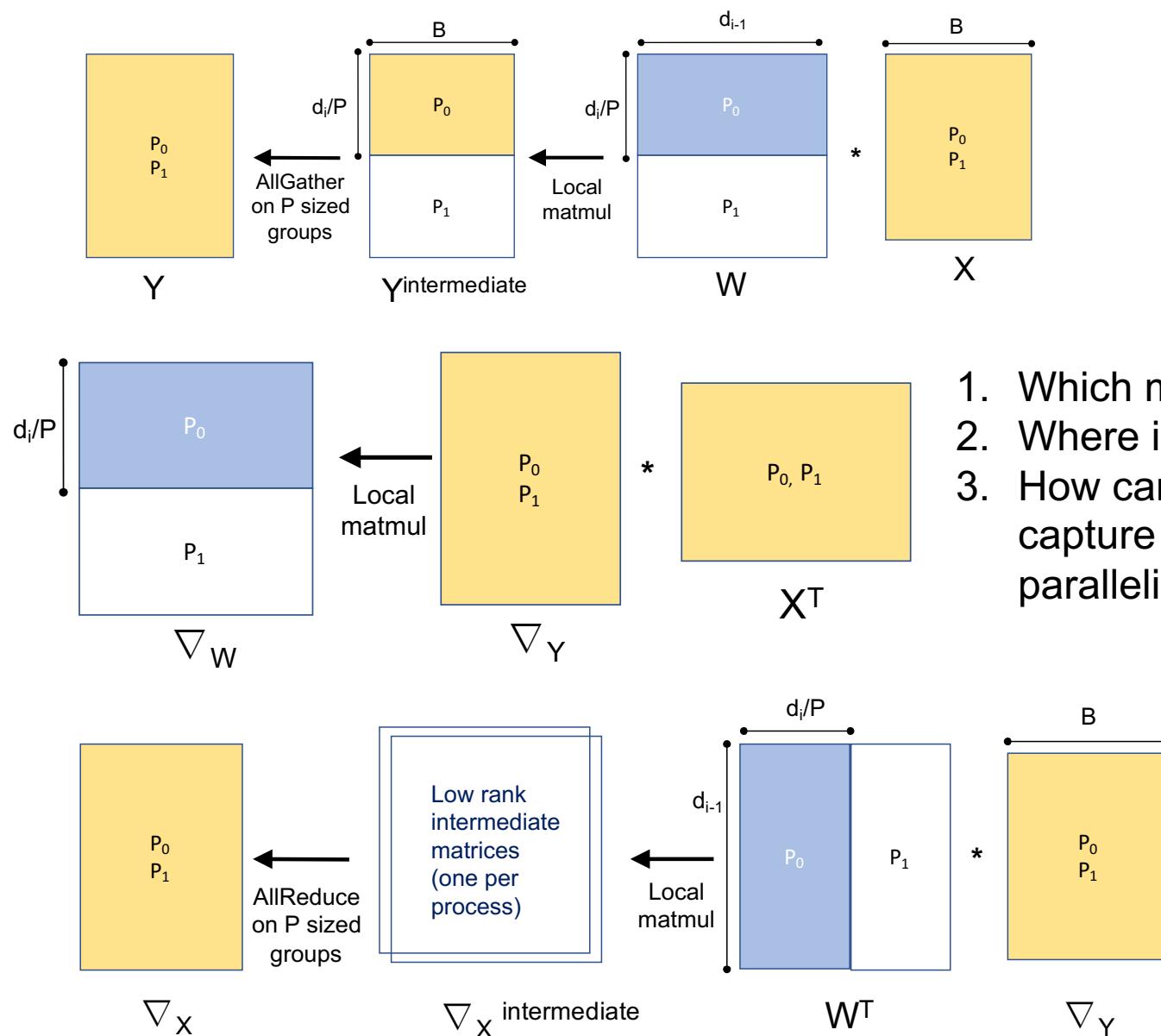


1. Which matrices are replicated?
2. Where is the communication?
3. Which steps can be overlapped?



$\nabla_Y = \partial L / \partial Y$  = how did the lost function change as output activations change?  
 $\nabla_X = \partial L / \partial X$   
 $\nabla_W = \partial L / \partial W$

# Model Parallel SGD training of NNs as matrix operations



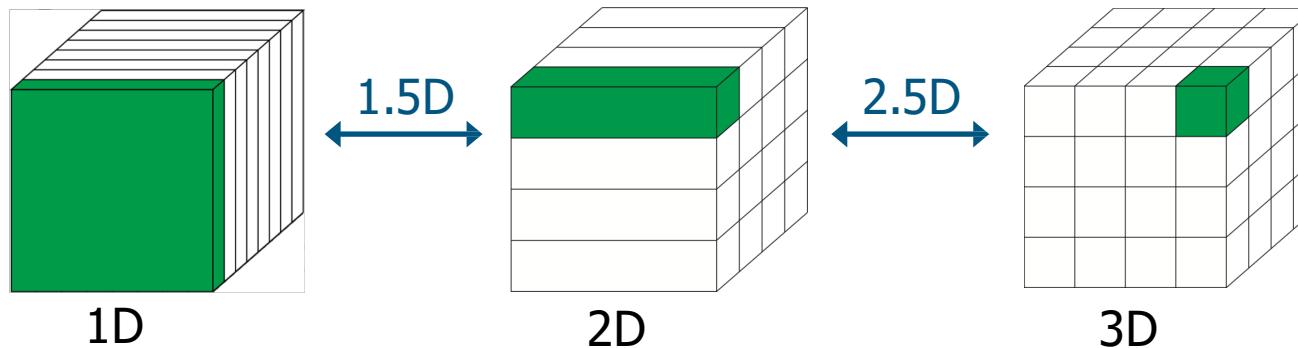
1. Which matrices are replicated?
2. Where is the communication?
3. How can matrix algebra capture both model and data parallelism?

## Combinations of various parallelism opportunities

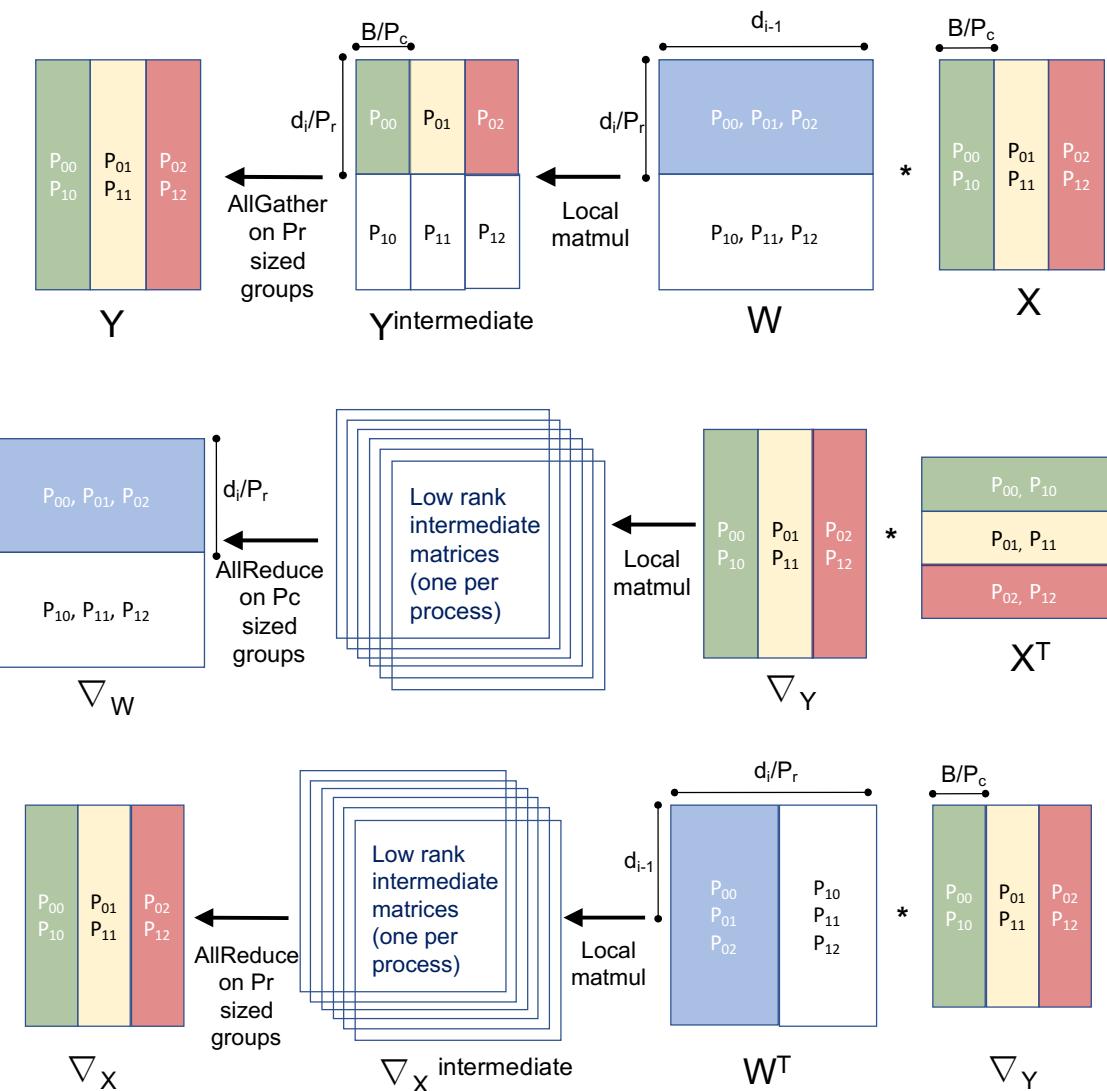
---

- There are various different ways to combine DNN training parallelism opportunities.
  - It helps to think in terms of matrices again.
- We will exploit ***communication-avoiding matrix algorithms***; which trade off some storage (judicious replication) at the expense of reduced communication.
  - Deep Learning community is already OK with data or model replication in many cases

### A succinct classification of parallel matrix multiplication algorithms



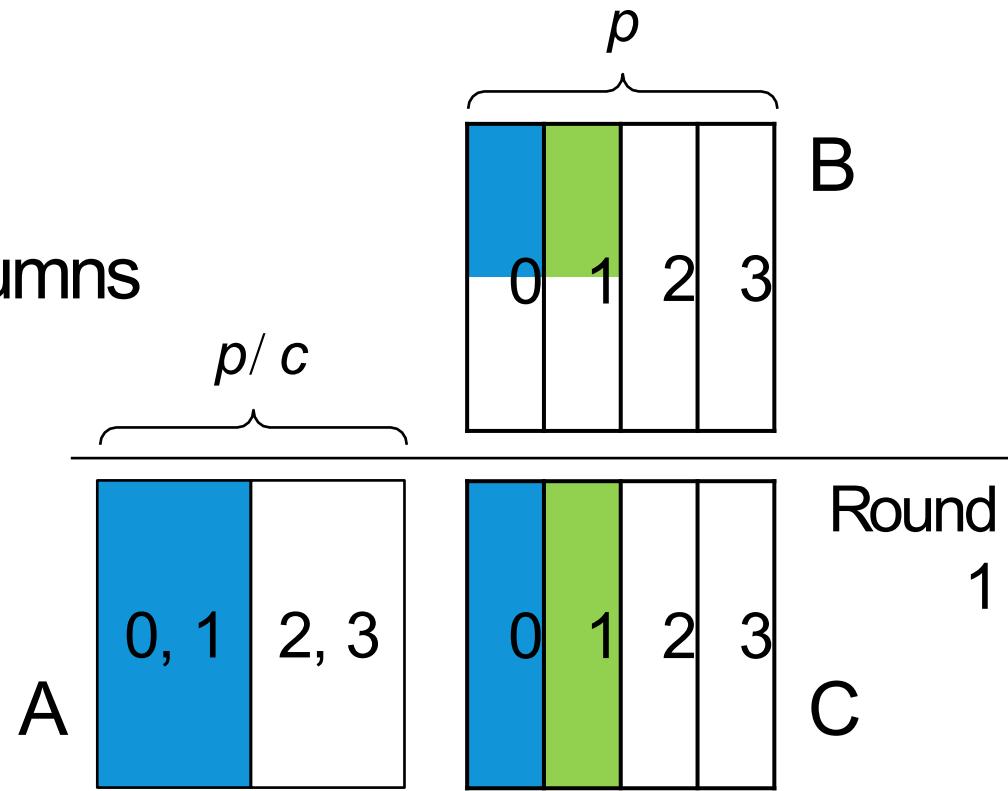
# Batch & Model Parallel SGD training of NNs as matrix operations



Processes are  
2D indexed:  
 $P = P_r \times P_c$

## 1.5D matrix multiplication algorithm

- $p=4, c=2$
- A has  $p/c$  block columns

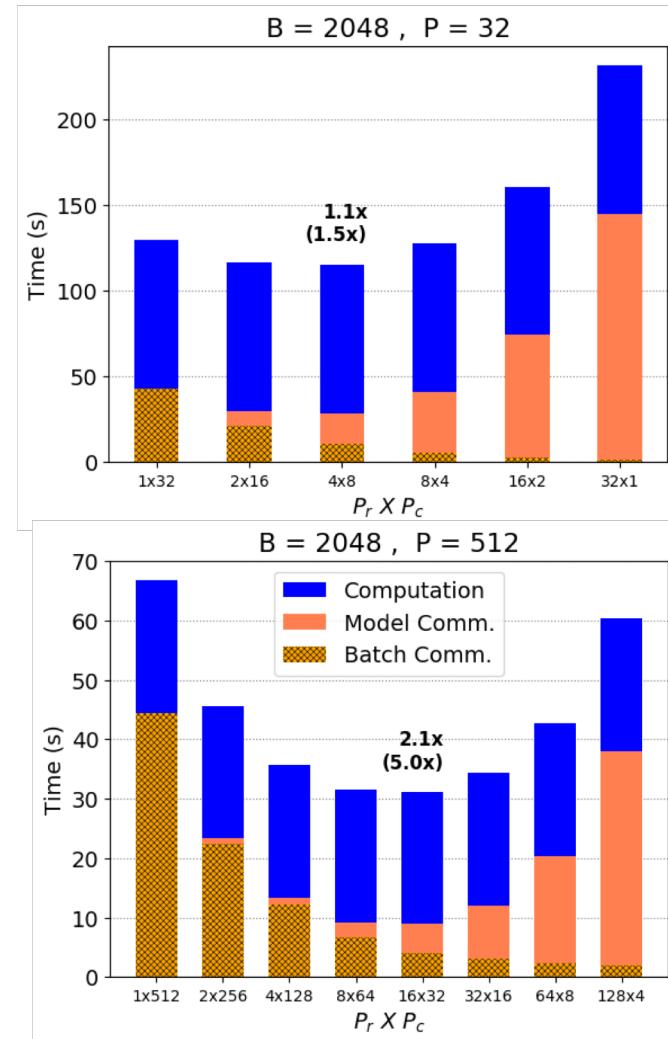
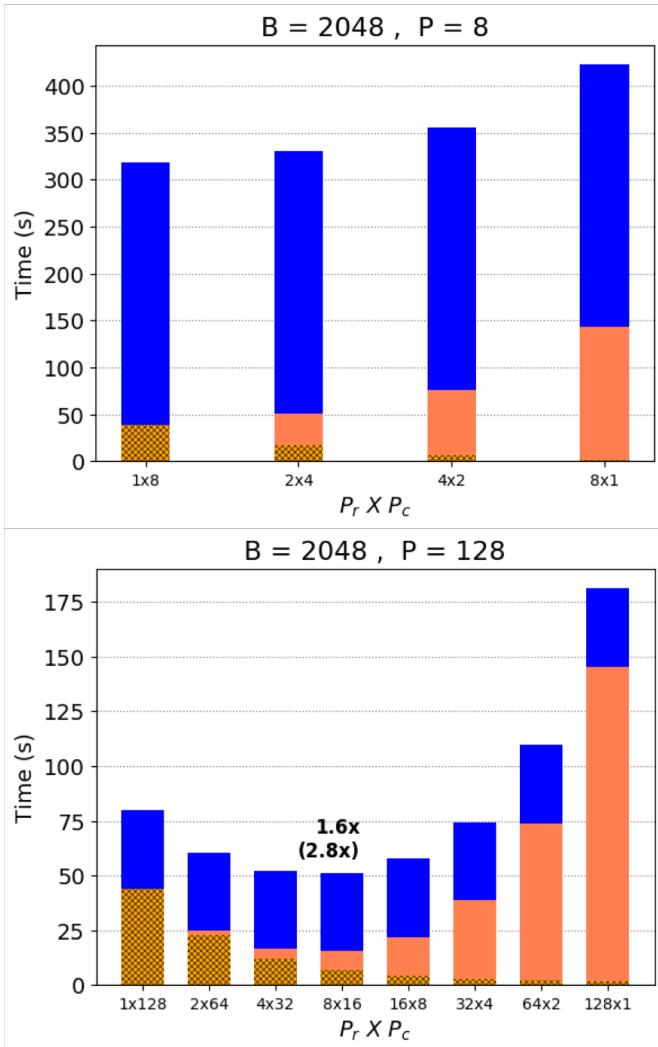


This is the so-called “col A” variant of the algorithm that fits best to our current use case. The cited paper below includes other variants.

Koanantakool, Penporn, et al. "Communication-avoiding parallel sparse-dense matrix-matrix multiplication." IPDPS, 2016

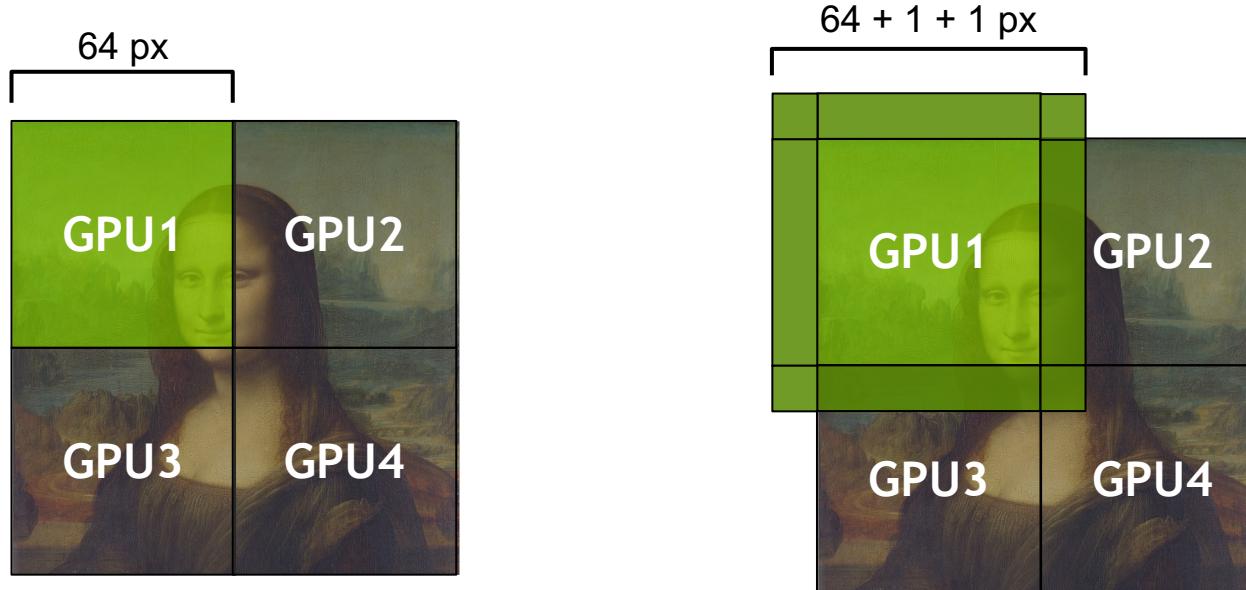
# Integrated Batch + Model Scaling

- For large processes integrated could provide up to 2x speedup



## Domain Parallel

- ° The general idea is the same as *halo regions* or *ghost zones* used to parallelize stencil codes in HPC
  - Before a convolution, exchange local receptive field boundary data

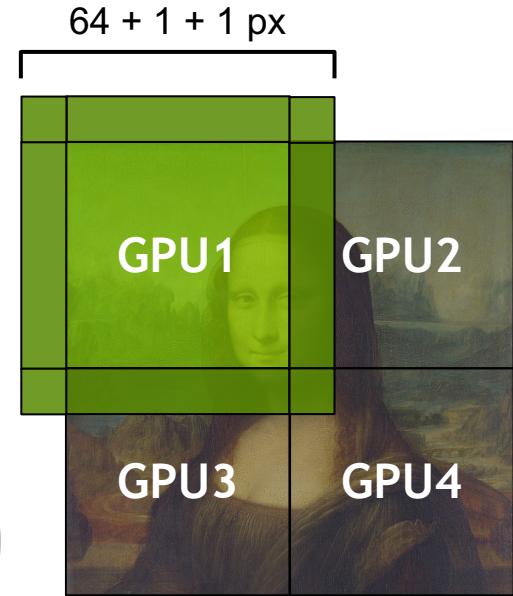


Peter Jin, Boris Ginsburg, and Kurt Keutzer. "Spatially Parallel Convolutions" ICLR Workshop Track, 2018

## Communication Complexity of Domain Parallel

- ° Additional communication for halo exchange during forward and backwards pass
  - Negligible cost for early layers for which activation size is large (i.e. convolutional)

$$T_{comm}(\text{domain}) = \sum_{i=0}^L (\alpha + \beta BX_W^i X_C^i k_h^i / 2) + \sum_{i=0}^L (\alpha + \beta BY_W^i Y_C^i k_w^i / 2) + 2 \sum_{i=0}^L \left( \alpha \log(P) + \beta \frac{P-1}{P} |W_i| \right)$$

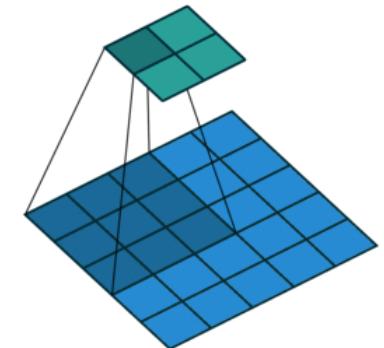


# Domain Parallel Scaling

## ◦ Domain parallel scaling on V100 GPUs

- **B=32, C=64, K=3, D=1, R=1**

| Resolution       | GPUs     | Fwd. wall-clock        | Bwd. wall-clock        |
|------------------|----------|------------------------|------------------------|
| $128 \times 128$ | 1        | 2.56 ms (1.0×)         | 6.63 ms (1.0×)         |
|                  | 2        | 1.52 ms (1.7×)         | 3.50 ms (1.9×)         |
|                  | <b>4</b> | <b>1.23 ms (2.1×)</b>  | <b>2.33 ms (2.8×)</b>  |
| $256 \times 256$ | 1        | 10.02 ms (1.0×)        | 26.81 ms (1.0×)        |
|                  | 2        | 5.34 ms (1.9×)         | 11.79 ms (2.3×)        |
|                  | <b>4</b> | <b>3.11 ms (3.2×)</b>  | <b>6.96 ms (3.9×)</b>  |
| $512 \times 512$ | 1        | 45.15 ms (1.0×)        | 126.11 ms (1.0×)       |
|                  | 2        | 20.18 ms (2.2×)        | 60.15 ms (2.1×)        |
|                  | <b>4</b> | <b>10.65 ms (4.2×)</b> | <b>26.76 ms (4.7×)</b> |



Peter Jin, Boris Ginsburg, and Kurt Keutzer. "Spatially Parallel Convolutions" ICLR Workshop Track, 2018  
Figure: Dumoulin, V., Visin, F.. A guide to convolution arithmetic for deep learning. *arXiv:1603.07285*, 2016.

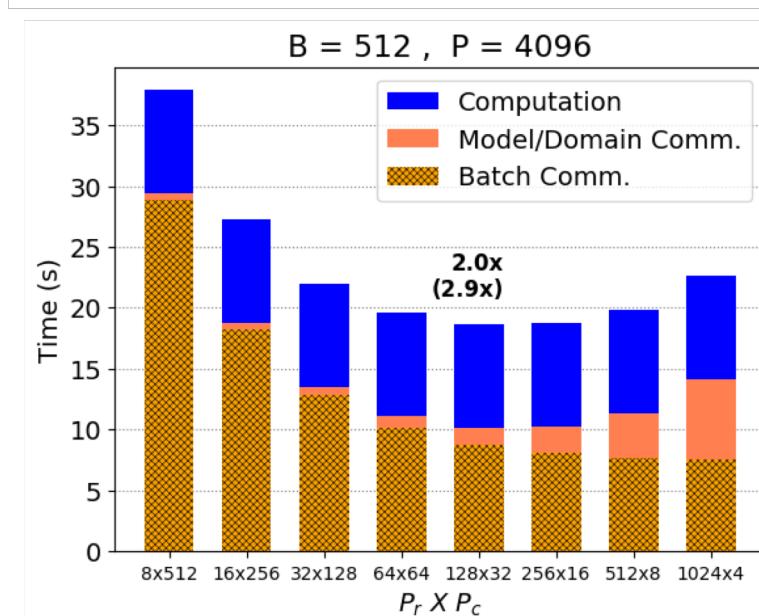
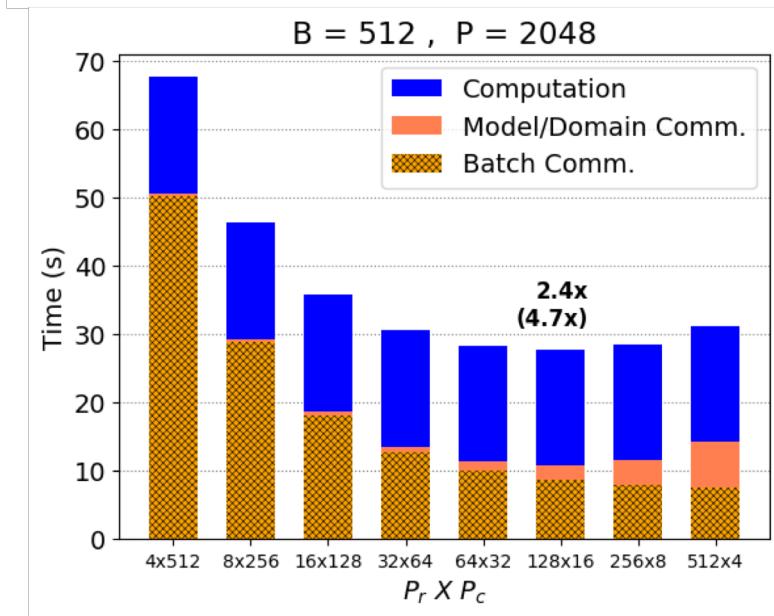
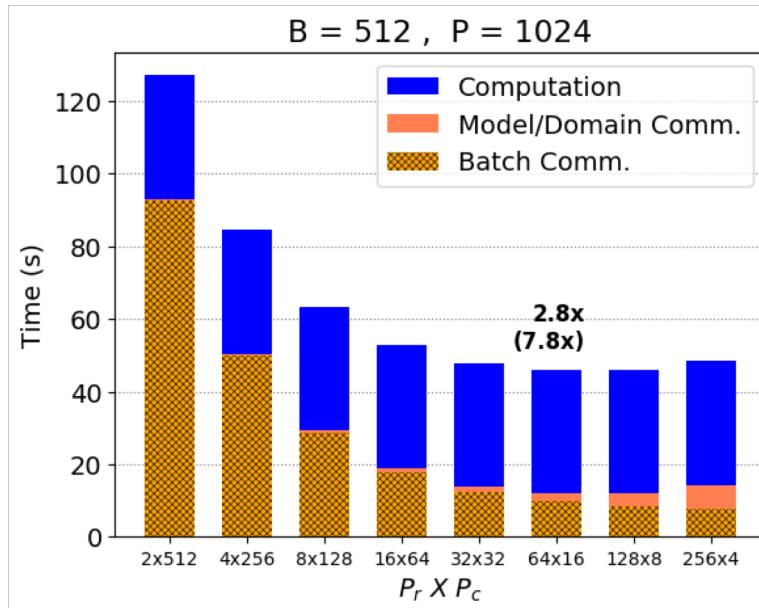
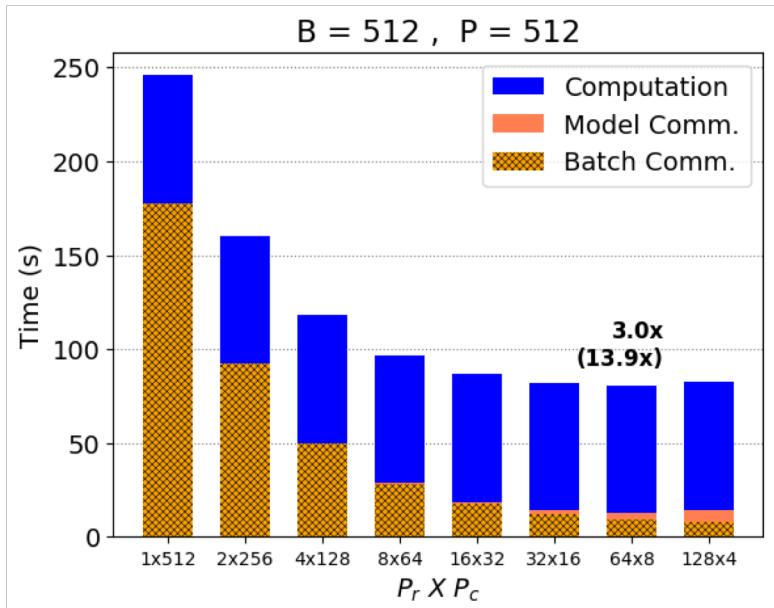
## Integrated Batch, Domain, and Model Parallel

---

- **Batch + Model** for layers with small activation size and large parameters (mostly Fully Connected Layers)
- **Batch + Domain** for layers with large activation sizes (mostly conv layers)

$$T_{comm} = \sum_{i \in L_M} \left( \alpha \log(P_r) + \beta \frac{B}{P_c} \frac{P_r - 1}{P_r} d_i \right) + 2 \sum_{i \in L_M} \left( \alpha \log(P_r) + \beta \frac{B}{P_c} \frac{P_r - 1}{P_r} d_{i-1} \right) \\ + 2 \sum_{i \in L_M} \left( \alpha \log(P_c) + \beta \frac{P_c - 1}{P_c} \frac{|W_i|}{P_r} \right) + \sum_{i \in L_D} \left( \alpha + \beta \frac{B}{P_c} X_W^i X_C^i k_h^i / 2 \right) \\ + \sum_{i \in L_D} \left( \alpha + \beta \frac{B}{P_c} X_W^{i+1} X_C^{i+1} k_w^i / 2 \right) + 2 \sum_{i \in L_D} \left( \alpha \log(P) + \beta \frac{P - 1}{P} |W_i| \right)$$

# Integrated Batch, Domain, and Model Parallel



## Summary of distributed deep learning

---

- Large batch size training often times lead to suboptimal learning
- Integrated parallelism uses communication-avoiding algorithms to extend scaling beyond batch size
- Integrated parallelism optimally combines model and data (batch and domain) parallelism and often performs better than each extreme.
- It is often better to use model parallelism for fully connected layers (large parameters), and data parallelism for convolutional layers (large activations)

Krizhevsky, Alex. "One weird trick for parallelizing convolutional neural networks." *arXiv preprint arXiv:1404.5997* (2014).

## Outline of the lecture

---

### **Today: Part 1, Intro and Supervised Learning**

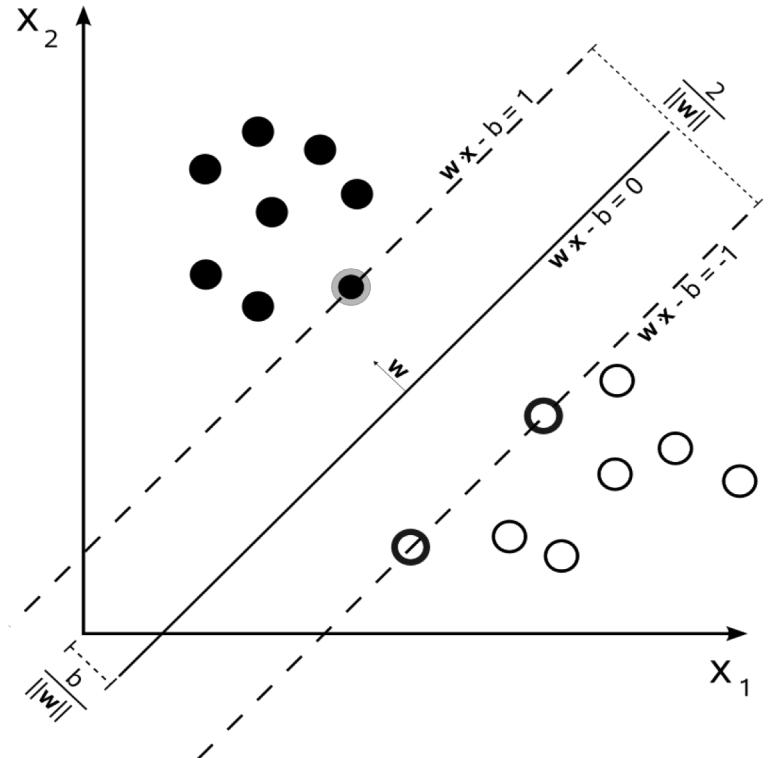
- Machine Learning & Parallelism Intro
- Neural Network Basics
- Scaling Deep Neural Network Training
- Support Vector Machines

### **Thursday: Part 2, Unsupervised Learning**

- Non-Negative Matrix Factorization
- Spectral and Markov Clustering
- Sparse Inverse Covariance Matrix Estimation

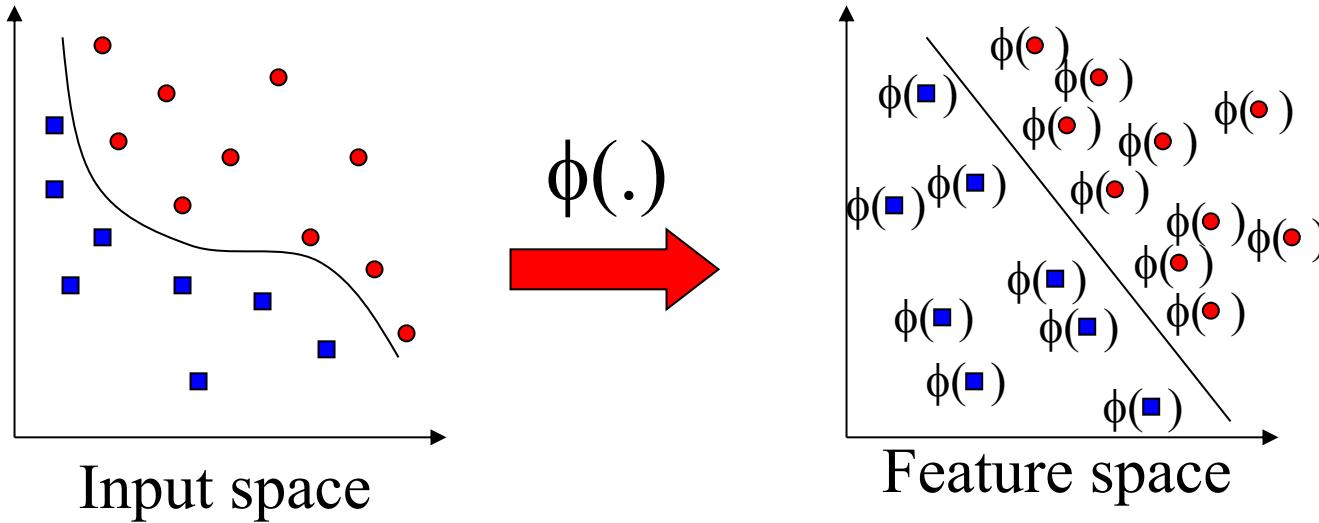
## Support Vector Machine (the linear case)

$$\begin{aligned} \max_{w,b} \quad & \frac{1}{\|w\|} \\ \text{s.t.} \quad & y_i(w^T x_i + b) \geq 1, \quad \forall i \end{aligned}$$



- Only the classification constraints on the support vectors are active
- Naively, leads to a giant quadratic constrained optimization (QP) problem
- Special algorithms, such as Sequential Minimal Optimization (SMO), decompose this giant QP to smaller (in fact *minimal*) QP sub-problems.

# Kernel Support Vector Machine (the non-linear case)



Slide: Eric Xing

Note: feature space is of higher dimension than the input space in practice

- Computation in the feature space can be costly because it is high dimensional
  - The feature space is typically infinite-dimensional!
- ***The kernel trick*** comes to rescue

## Kernel Methods

---

- Kernel methods replace the dot-product similarity of linear models with an arbitrary Kernel function that computes the similarity between  $x$  and  $y$

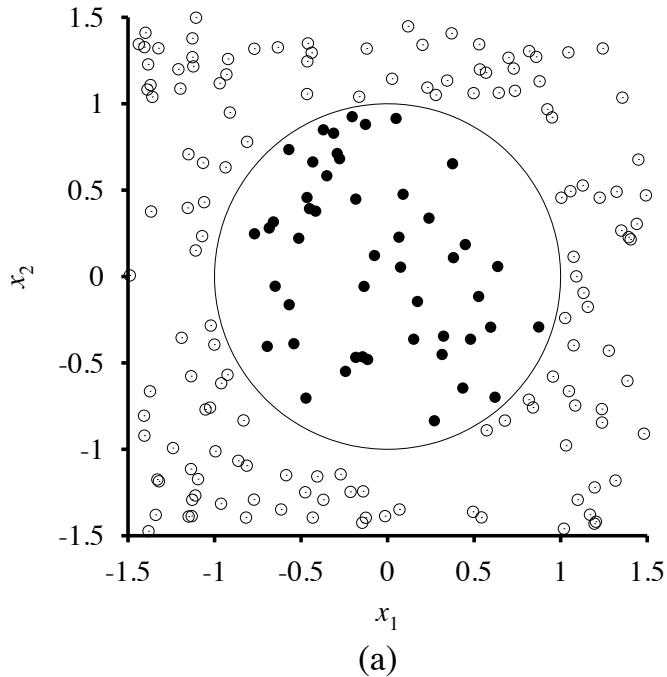
$$K(x, y) : X \times X \rightarrow \mathbb{R}$$

- Kernels need to be symmetric and positive semi-definite.
- Gaussian Kernel (i.e., Radial Basis Function) is the de-facto kernel in ML.

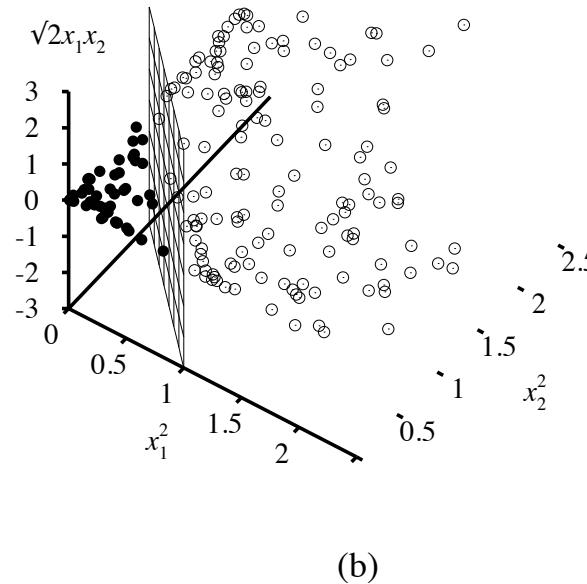
$$K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$$

- We could pre-compute the Kernel (Gram) Matrix for all pairs of samples, but this is too expensive.

# The Kernel Trick



(a)



(b)

Figure source:  
Russell & Norvig

The circular decision boundary in 2D (a) becomes a linear boundary in 3D (b) using the following transformation:  $\phi(x_1, x_2) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$

If we define the **kernel function** as  $\kappa(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i \cdot \mathbf{x}_j)^2$

Then no need to carry out  $\phi(\cdot)$  explicitly as  $\Phi(x) \cdot \Phi(y) = (x \cdot y)^2$

$$(x_1^2, x_2^2, \sqrt{2}x_1x_2) \cdot (y_1^2, y_2^2, \sqrt{2}y_1y_2) = (x_1y_1 + x_2y_2)^2$$

## Major Bottleneck of Kernel SVM

---

- Input dataset: n-by-d matrix ( $n \gg d$ )
  - $X_1, X_2, \dots, X_n$ .  $X_i$  is a vector has d features
- Generate a n-by-n Kernel (Gram) matrix at runtime
  - $K[i][j] = \exp(-r||X_i - X_j||^2)$ , r is positive number
- $O(d^*n^2)$  operations and  $O(n^2)$  memory are huge!
  - a small input generates a large Kernel matrix
  - 357MB input (52K-by-90) = 2000GB Kernel matrix
- Solution: SMO (sequential minimal optimization)
  - using iterative method, avoiding Kernel matrix
  - using ***two rows of Kernel matrix*** each iteration
  - key kernel for sparse inputs: ***sparse matrix times sparse vector*** (as in the BFS case from graph lecture)

# Sequential Minimal Optimization

---

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \boxed{\kappa(\mathbf{x}_i^T \mathbf{x}_j)}$$

The kernel function

$$\text{s.t. } 0 \leq \alpha_i \leq C, \quad i = 1, \dots, m$$

$$\boxed{\sum_{i=1}^m \alpha_i y_i = 0.}$$

The equality constraint

- The smallest possible optimization problem involves “two” Lagrange multipliers at a time
- Because just changing one multiplier would violate the equality constraint

Platt, John. “Fast training of support vector machines using sequential minimal optimization.” Advances in kernel methods. 1999.

# Sequential Minimal Optimization

---

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \boxed{\kappa(\mathbf{x}_i^T \mathbf{x}_j)}$$

The kernel function

$$\text{s.t. } 0 \leq \alpha_i \leq C, \quad i = 1, \dots, m$$

$$\boxed{\sum_{i=1}^m \alpha_i y_i = 0.}$$

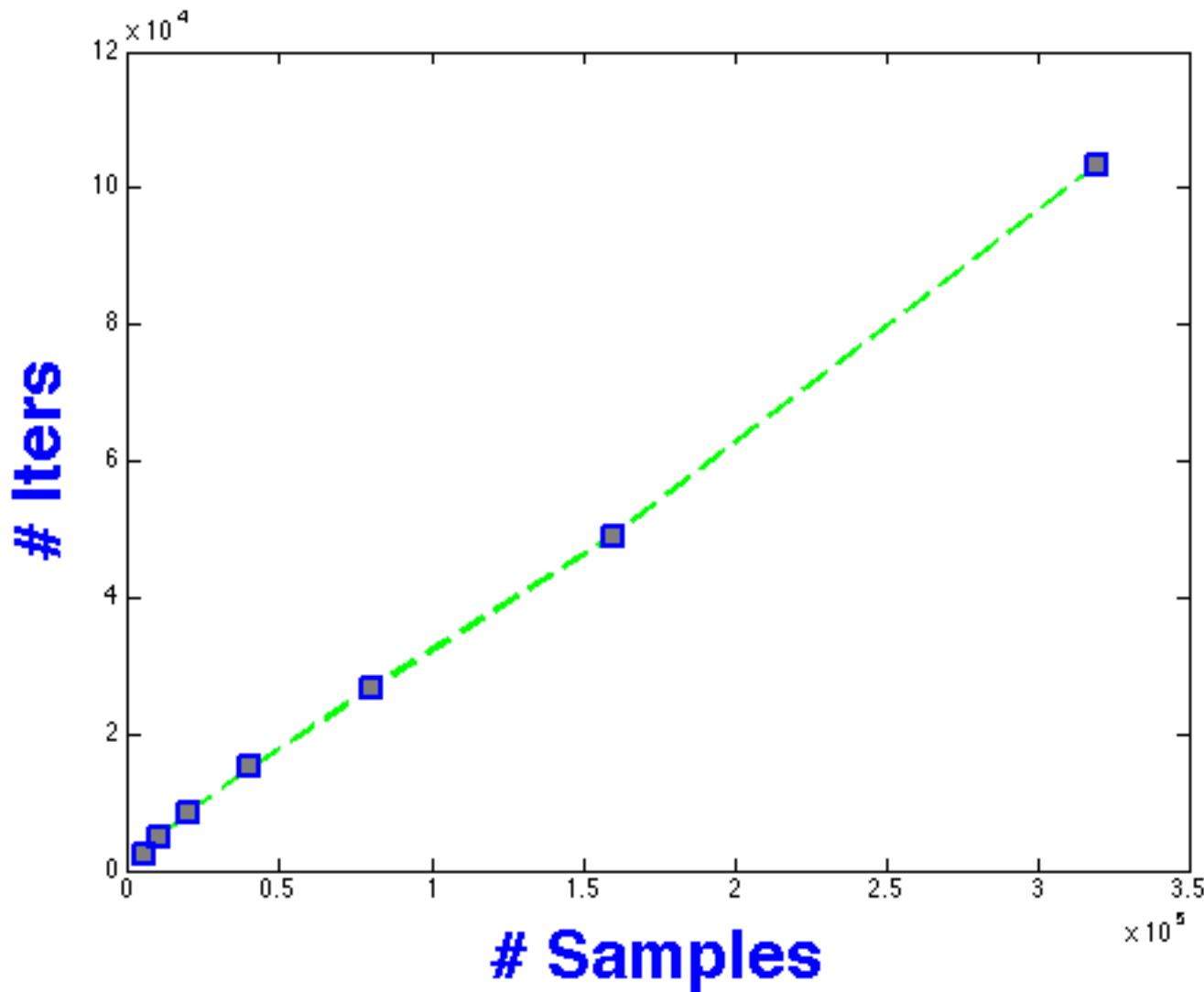
The equality constraint

**Repeat until convergence:**

1. Select some pair  $\alpha_i$  and  $\alpha_j$  to update next (using a heuristic that tries to pick the two that will allow us to make the biggest progress towards the global maximum).
2. Re-optimize  $W(\alpha)$  with respect to  $\alpha_i$  and  $\alpha_j$ , while holding all the other  $\alpha_k$ 's fixed.

#Iterations = O(#samples), bad Weak Scaling!

---

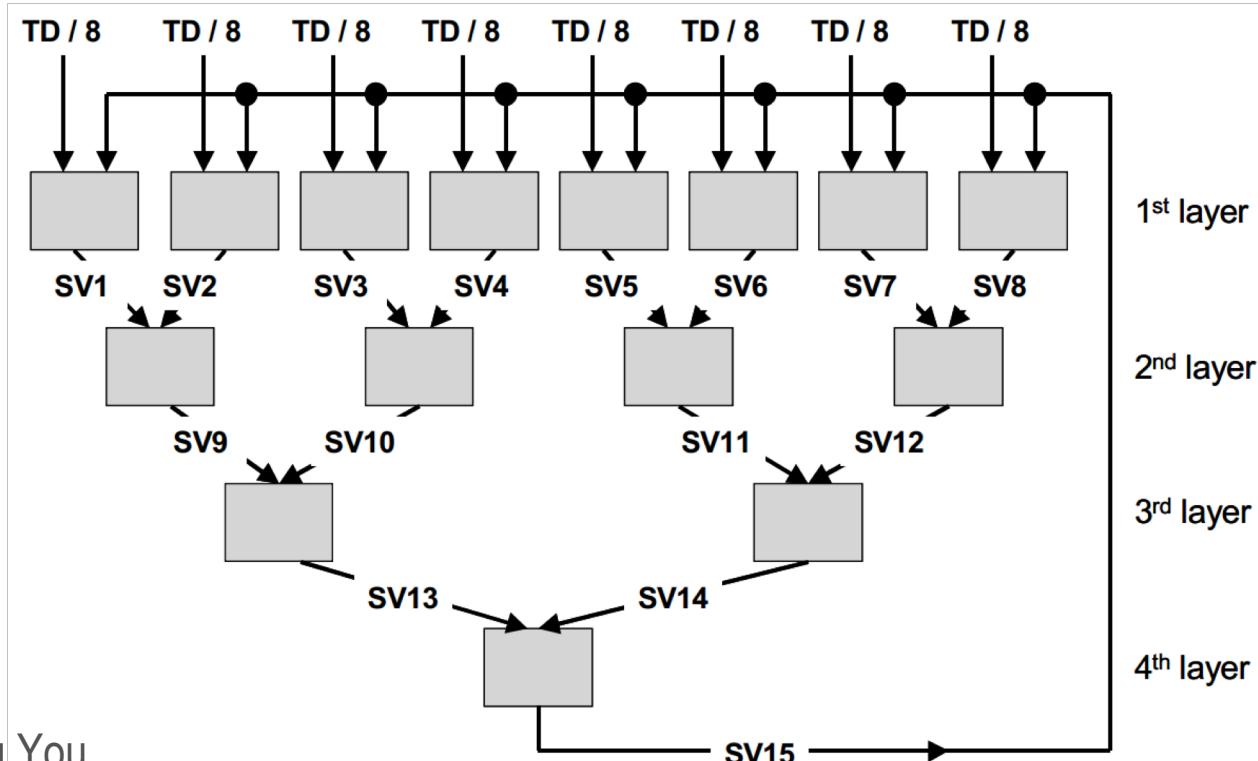


# Cascade SVM (NIPS'04)

Data is partitioned \*evenly\* and processed by multiple SVMs

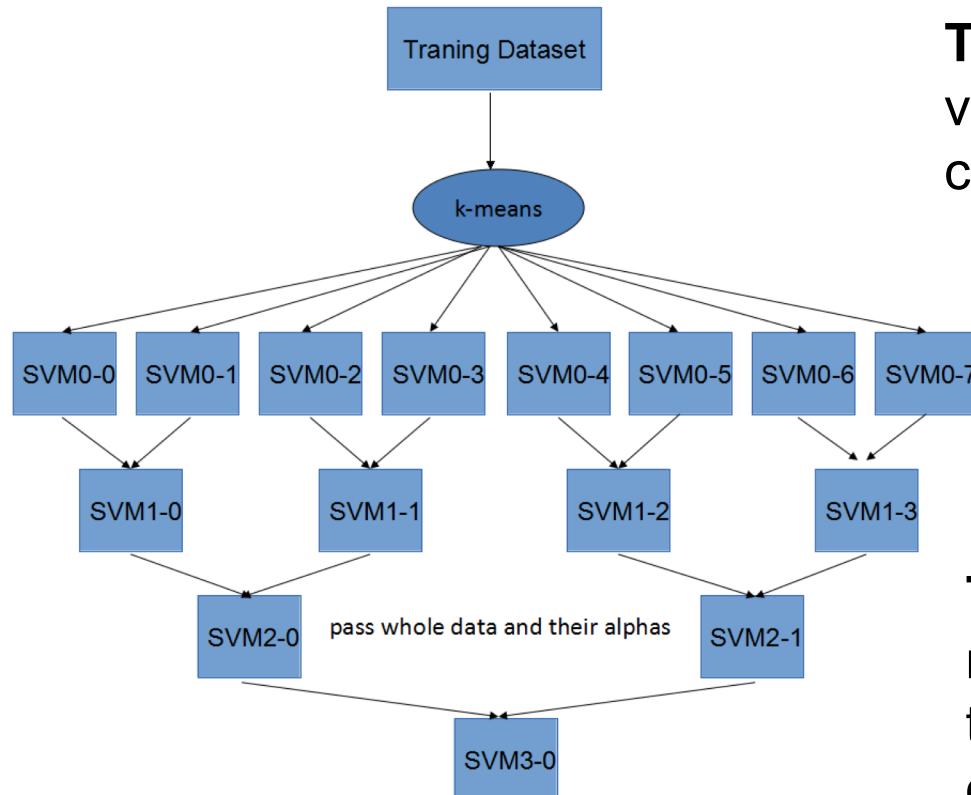
Remove the non support vectors layer-by-layer

- ° data is the support vectors (SV) of previous layer (major reduction)
- ° pass parameters  $\alpha_i$  of SVs to next layer for a better initiation (warm start)



# Divide-and-Conquer SVM: DC-SVM (ICML'14)

- ° Difference between DC-SVM and Cascade
  - ° DC-SVM passes all data layer-by-layer
  - ° DC-SVM uses kernel k-means to partition the dataset

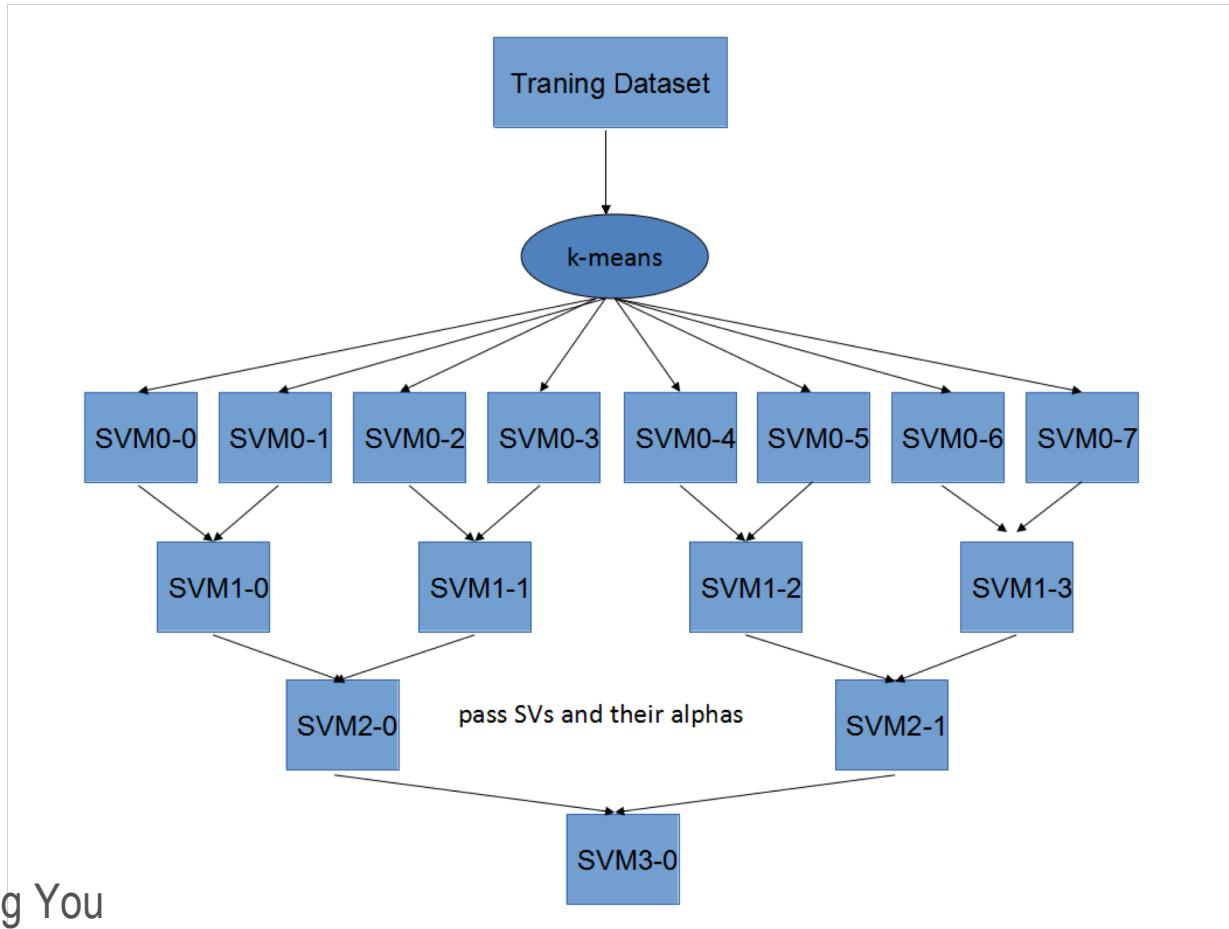


**Theorem #1:** the set of support vectors from the subproblems is close to that of the whole problem

**Theorem #2:** Kernel kmeans minimizes the difference between the solution of subproblems and of the whole problem

## Combine DC-SVM with Cascade: DC-Filter

- ° Only pass support vectors layer-by-layer (reduce workload)
- ° Use kernel k-means to divide the dataset



## Bottleneck of Cascade, DC-SVM, and DC-Filter

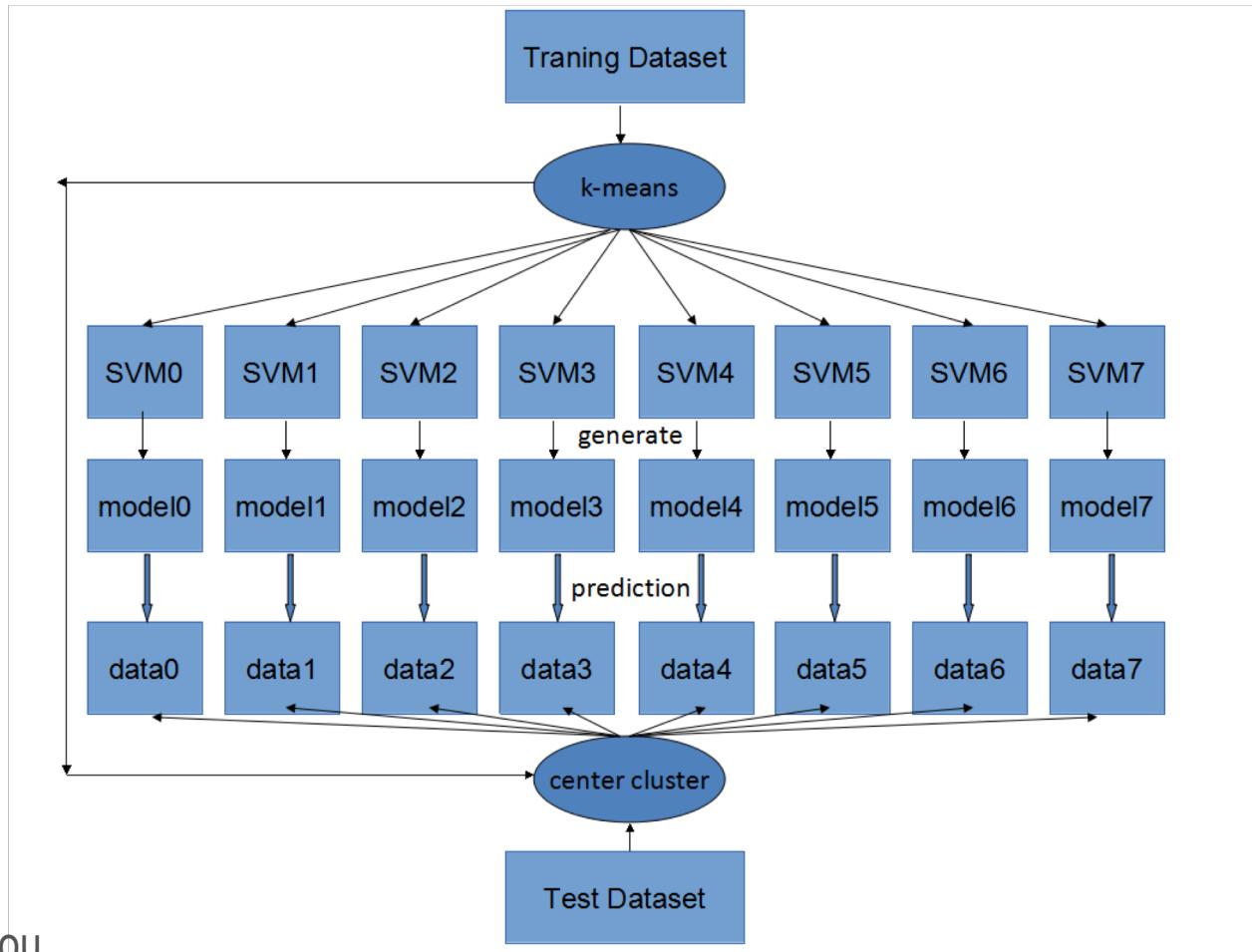
---

- ° Occupy P machines, but bottom level uses 1 machine
- ° Lower levels cost more time than top level

|                             |      |      |       |      |      |      |      |      |
|-----------------------------|------|------|-------|------|------|------|------|------|
| <b>level 1<sup>st</sup></b> | 6000 | 6000 | 6000  | 6000 | 6000 | 6000 | 6000 | 6000 |
| <b>time: 5.49s</b>          | 4.87 | 4.92 | 4.90  | 4.68 | 5.12 | 5.10 | 5.49 | 4.71 |
| <b>iter: 6168</b>           | 5648 | 5712 | 5666  | 5415 | 5936 | 5904 | 6168 | 5453 |
| <b>SVs: 5532</b>            | 746  | 715  | 717   | 718  | 686  | 707  | 721  | 699  |
| <b>level 2<sup>nd</sup></b> | 1461 |      | 1435  |      | 1393 |      | 1420 |      |
| <b>time: 1.58s</b>          | 1.58 |      | 1.50  |      | 1.35 |      | 1.45 |      |
| <b>iter: 7485</b>           | 7485 |      | 7211  |      | 6713 |      | 7035 |      |
| <b>SVs: 5050</b>            | 1292 |      | 1263  |      | 1256 |      | 1239 |      |
| <b>level 3<sup>rd</sup></b> | 2555 |      |       |      | 2495 |      |      |      |
| <b>time: 3.34s</b>          | 3.34 |      |       |      | 3.30 |      |      |      |
| <b>iter: 9081</b>           | 8975 |      |       |      | 9081 |      |      |      |
| <b>SVs: 4699</b>            | 2388 |      |       |      | 2311 |      |      |      |
| <b>level 4<sup>th</sup></b> |      |      | 4699  |      |      |      |      |      |
| <b>time: 9.69s</b>          |      |      | 9.69  |      |      |      |      |      |
| <b>iter: 14052</b>          |      |      | 14052 |      |      |      |      |      |
| <b>SVs: 4475</b>            |      |      | 4475  |      |      |      |      |      |

# CP-SVM: Cluster Partition SVM (1-level divide-and-conquer)

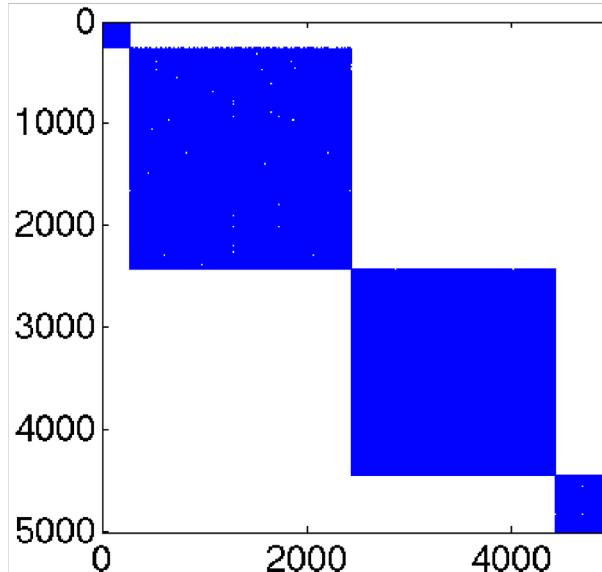
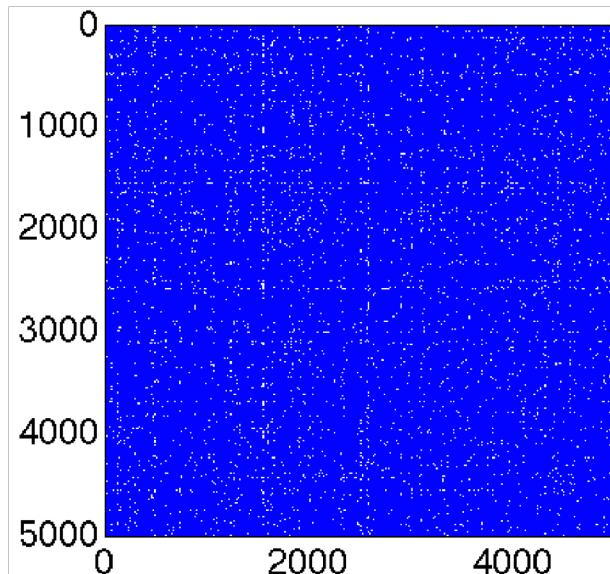
- ° Divide: K-means partitions data into P parts; P models
- ° Conquer: Euclidean distance to select best model



## Why CP-SVM works?

---

- ° When  $\|X_i - X_j\|^2$  is large,  $\exp(-r\|X_i - X_j\|^2)$  is zero
- ° K-means maximize different groups' Euclidean distance
- ° These two matrices have similar F-norm
- ° Analysis assumes the Gaussian kernel: For a given sample, only the support vectors close to it can have an effect on the classification



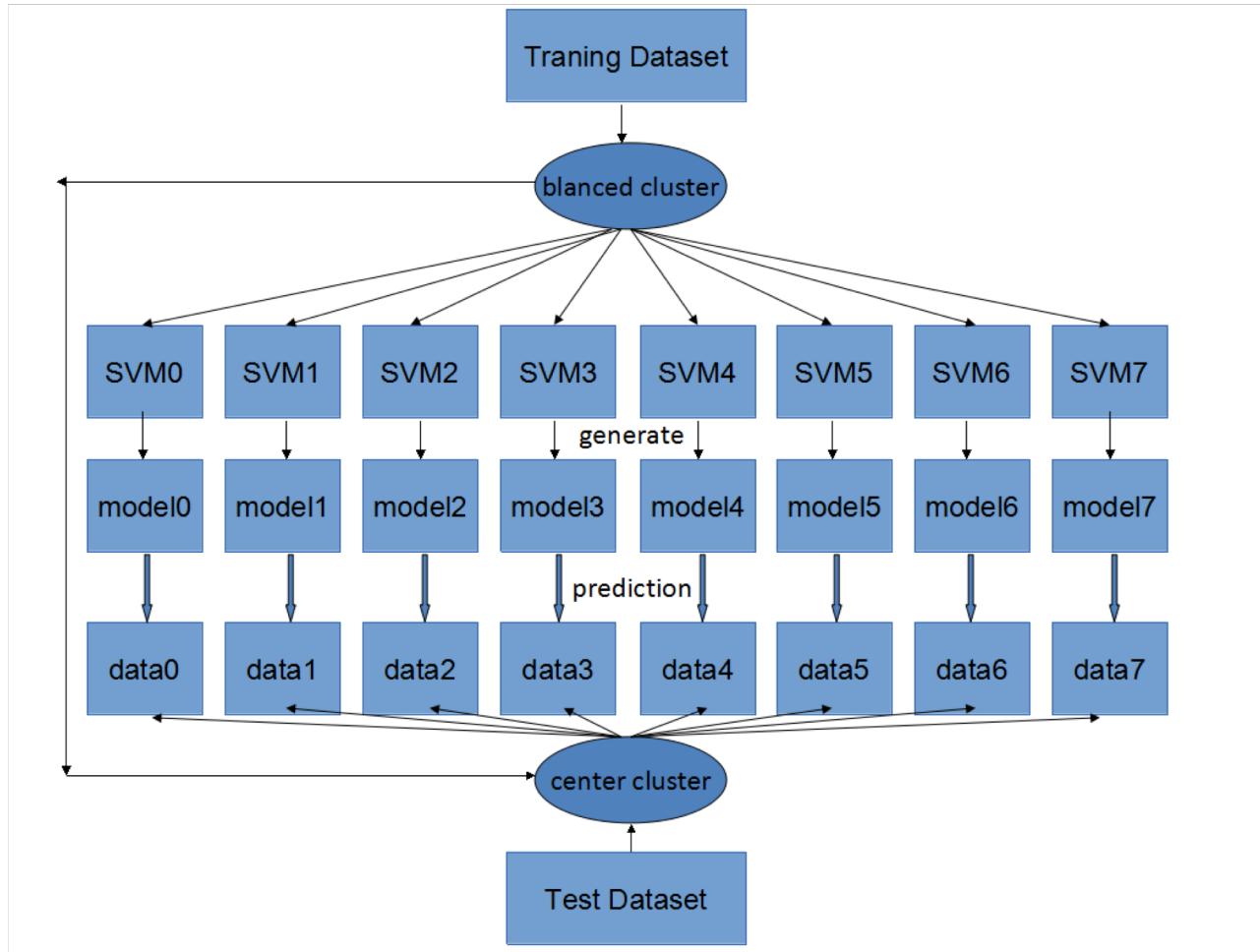
## CP-SVM's bottleneck: load imbalance (because of K-means)

---

- ° Machine 5, time: 0.75s, #iter: 1522, #samples: 4800
- ° Machine 1, time: 0.79s, #iter: 1384, #samples: 4800
- ° Machine 7, time: 0.94s, #iter: 1748, #samples: 4800
- ° Machine 6, time: 1.14s, #iter: 2337, #samples: 4800
- ° Machine 2, time: 1.14s, #iter: 2339, #samples: 4800
- ° Machine 4, time: 1.31s, #iter: 2856, #samples: 4800
- ° Machine 3, time: 5.48s, #iter: 6723, #samples: 9600
- ° Machine 3, time: 6.48s, #iter: 7915, #samples: 9600

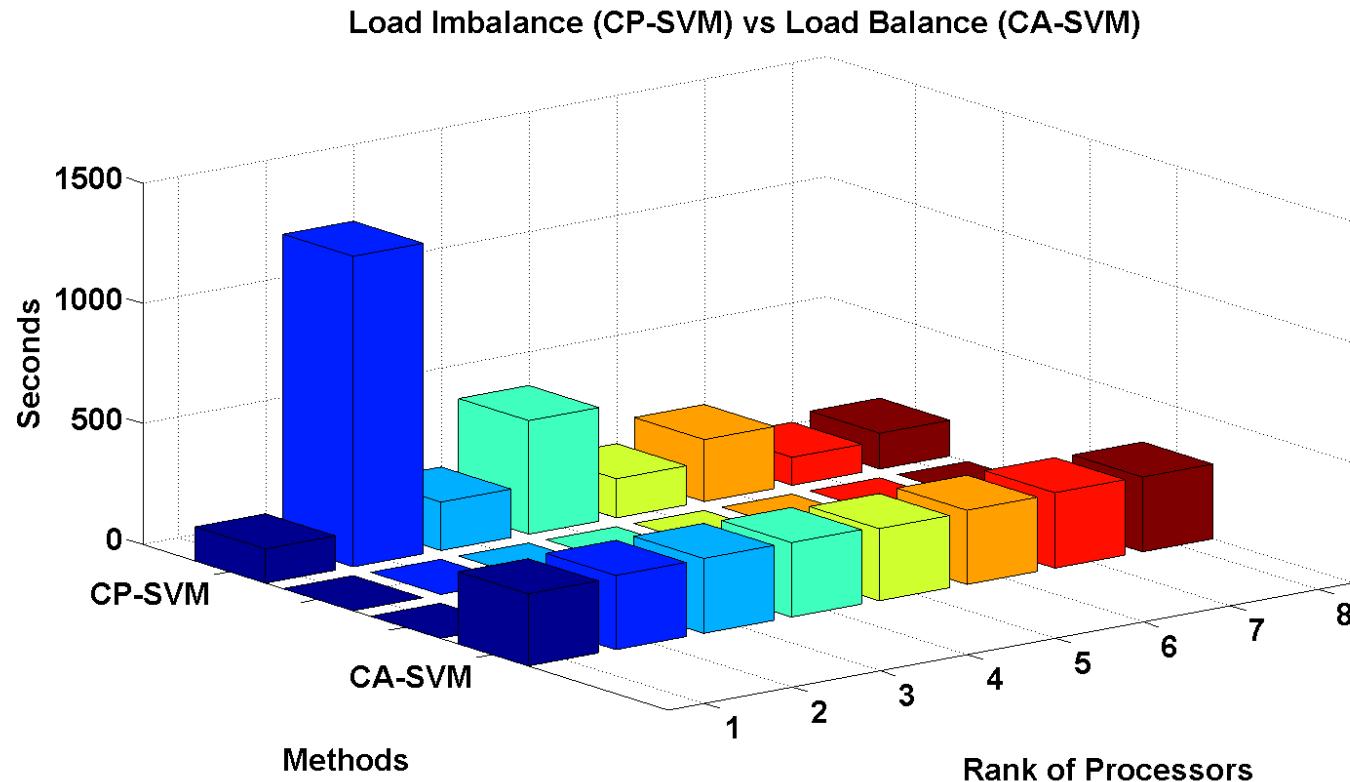
# CA-SVM: Communication Avoiding SVM

- Design a balanced clustering to replace K-means
- Still approximating the distance separation property of k-means as much as possible.



# Load Balance Comparison

- ° Test dataset is epsilon with 128k samples on 8 machines



## 0.2 accuracy loss for 6.6x speedup

---

- ° Overall comparison for IJCNN dataset
- ° All methods use the same number of machines

| Method           | Accuracy | Iterations | Time (Init, Training) |
|------------------|----------|------------|-----------------------|
| <b>Dis-SMO</b>   | 98.7%    | 30,297     | 23.8s (0.008, 23.8)   |
| <b>Cascade</b>   | 95.5%    | 37,789     | 13.5s (0.007, 13.5)   |
| <b>DC-SVM</b>    | 98.3%    | 31,238     | 59.8s (0.04, 59.7)    |
| <b>DC-Filter</b> | 95.8%    | 17,339     | 8.4s (0.04, 8.3)      |
| <b>CP-SVM</b>    | 98.7%    | 7,915      | 6.5s (0.04, 6.4)      |
| <b>CA-SVM</b>    | 98.5%    | 7,450      | 3.6s (0.005, 3.55)    |

You, Yang, et al. "CA-SVM: Communication-Avoiding Support Vector Machines on Distributed Systems", IPDPS, 2015