**Lecture 9:**

# Parallel Deep Network Training

**Visual Computing Systems**
**Stanford CS348K, Fall 2018**

# Basic gradient descent

```
while (loss too high):
   for each epoch: // a pass through the training dataset
      for each item x_i in training set:
         grad = evaluate_loss_gradient(f, params, loss_func, x_i)
         params += -grad * learning_rate;
```

**Mini-batch stochastic gradient descent (mini-batch SGD):**

**choose a random (small) subset of the training examples to use to compute the gradient in each iteration of the while loop**

```
while (loss too high):
   for each epoch: // a pass through the training dataset
      for all mini batches in training set:
         grad = 0;
         for each item x_i in minibatch:
            grad += evaluate_loss_gradient(f, params, loss_func, x_i)
         params += -grad * learning_rate;
```

**How do we compute dLoss/dp for a deep neural network with millions of parameters?**

# Quick review of back-propagation

# Derivatives using the chain rule

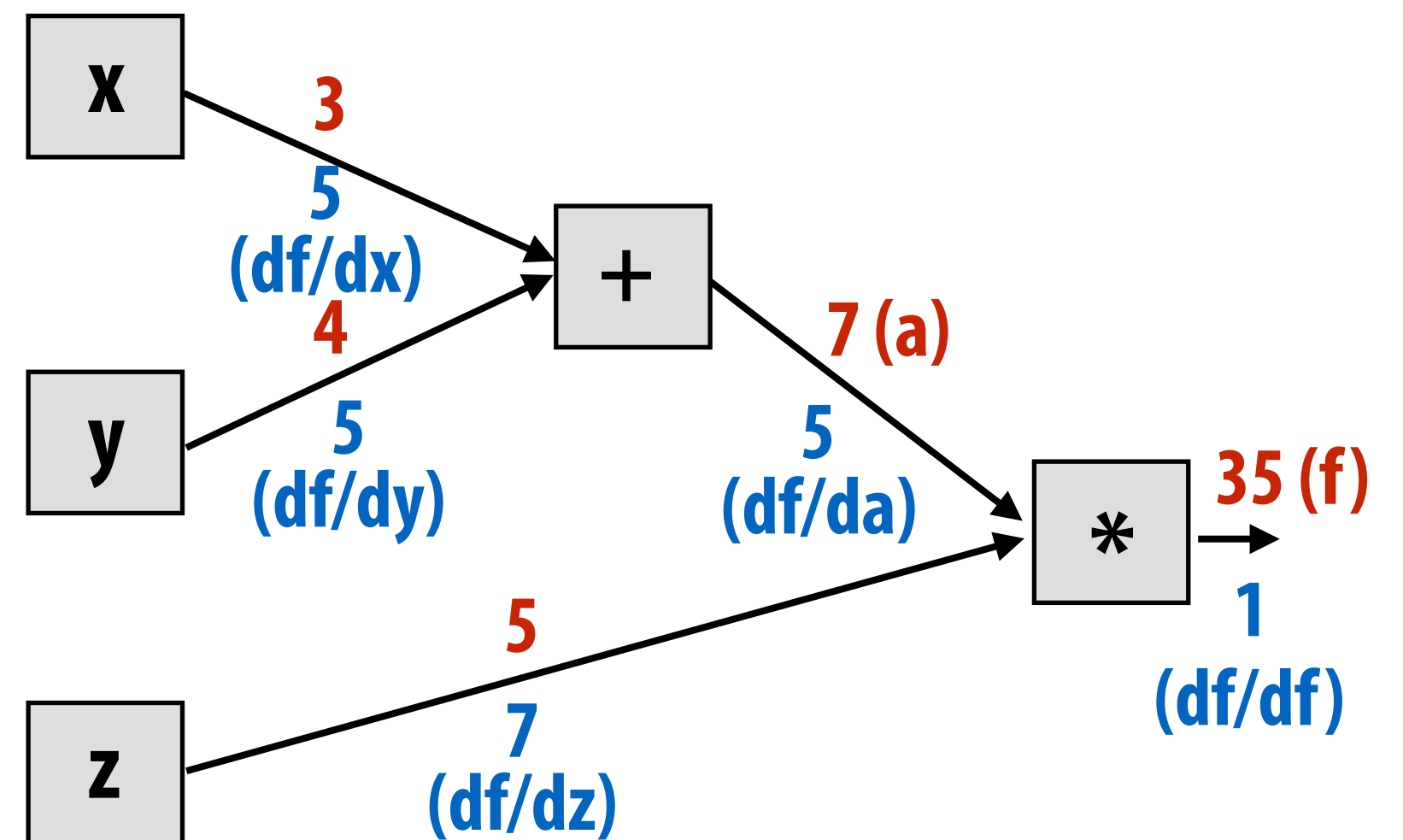$$f(x, y, z) = (x + y)z = az$$

**Where:** $a = x + y$

$$\frac{df}{da} = z \qquad \frac{da}{dx} = 1 \qquad \frac{da}{dy} = 1$$

**So, by the derivative chain rule:**
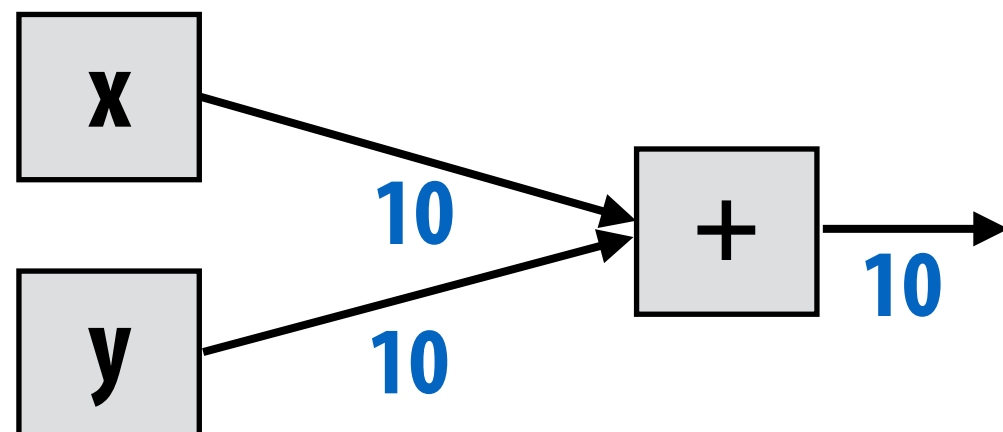
$$\frac{df}{dx} = \frac{df}{da}\frac{da}{dx} = z$$



Red = output of node
Blue = df/dnode

# Backpropagation

**Recall:** $\dfrac{df}{dx} = \dfrac{df}{dg}\dfrac{dg}{dx}$

x

**10**

**+** **10**

y

**10**

$g(x, y) = x + y$

$\dfrac{dg}{dx} = 1 \, , \, \dfrac{dg}{dy} = 1$

x **15**
**10**
**12** **max** **10**
y **0**

$g(x, y) = \max(x, y)$

$\dfrac{dg}{dx} = \begin{matrix} \textbf{1, if x > y} \\ \textbf{0, otherwise} \end{matrix}$

x **15**
**10*12**
**12** **\*** **10**
y **10*15**

$g(x, y) = xy$

$\dfrac{dg}{dx} = y \, , \, \dfrac{dg}{dy} = x$

# Back-propagating through single unit



**Recall: behavior of unit:**

$$f(x_0, x_1, x_2, x_3) = max\left(0, \sum_i x_i w_i + b\right)$$

**let y =** 10, if upper input to max is > 0
0, otherwise

$$\frac{d\text{loss}}{d\text{unit}}$$

**Observe: output of prior layer must be retained in order to compute weight gradients for this unit during backprop.**

# Data lifetimes during network evaluation



**Weights (read-only) reside in memory**

**After evaluating layer i, can free outputs from layer i-1**

# Data lifetimes during training



- Must retain outputs for all layers because they are needed to compute gradients during back-prop
- Parallel back-prop will require storage for per-weight gradients (more about this in a second)
- In practice: may also store per-weight gradient velocity (if using SGD with "momentum") or step size cache in adaptive step size schemes like Adagrad

```
vel_new = mu * vel_old - step_size * grad
w_new = w_old + vel_new
```

# SGD workload

```
while (loss too high):
```
At first glance, this loop is sequential (each step of "walking downhill" depends on previous)

```
for each item x_i in mini-batch:
    grad += evaluate_loss_gradient(f, loss_func, params, x_i)
```
Parallel across images

sum reduction

large computation with its own parallelism
(but working set may not fit on single machine)

```
params += -grad * step_size;
```
trivial data-parallel over parameters

# DNN training workload

- **Large computational expense**

  - Must evaluate the network (forward and backward) for millions of training images
  - Must iterate for many iterations of gradient descent (100's of thousands)
  - Training modern networks on big datasets takes days

- **Large memory footprint**

  - Must maintain network layer outputs from forward pass
  - Additional memory to store gradients/gradient velocity for each parameter
  - Recall parameters for popular VGG-16 network require ~500 MB of memory (training requires GBs of memory for academic networks)
  - Scaling to larger networks requires partitioning DNN across nodes to keep DNN + intermediates in memory

- **Dependencies /synchronization (not embarrassingly parallel)**

  - Each parameter update step depends on previous
  - Many units contribute to same parameter gradients (fine-scale reduction)
  - Different images in mini batch contribute to same parameter gradients

# Synchronous data-parallel training (across images)

```
for each item x_i in mini-batch:
    grad += evaluate_loss_gradient(f, loss_func, params, x_i)
params += -grad * learning_rate;
```

## Consider parallelization of the outer for loop across machines in a cluster



| image $x_0$ | | image $x_1$ |
|---|---|---|
| parameter gradients due to $x_0$ | copy of parameter values | parameter gradients due to $x_1$ |

**Node 0**                                          **Node 1**

```
partition dataset across nodes
for each item x_i in mini-batch assigned to local node:
    // just like single node training
    grad += evaluate_loss_gradient(f, loss_func, params, x_i)
barrier();
sum reduce gradients, communicate results to all nodes
barrier();
update copy of parameter values
```

# Synchronous training

- **All nodes cooperate to compute gradients for a mini-batch \***

- **Gradients are summed (across the entire machine)**
  - **All-to-all communication**
  - **Good implementations will sum gradients for layer _i_ when computing backprop for _i_+1 (overlap communication and computation).**

- **Update model parameters**
  - **Typically done without wide parallelism (e.g. each machine computes its own update)**

- **All nodes proceed to work on next mini-batch given new model parameters**

**\* If curious about batch norm in a parallel training setting. In practice each of _k_ nodes works on a set of _n_ images, with batch norm statistics computed independently for each set of n (mini-batch size is _kn_).**

# Challenges of scaling out (many nodes)

- **Slow communication between nodes**

  - **Commodity clusters do not feature high-performance interconnects (e.g., infiniband) typical of supercomputers**

  - **Synchronous SGD involves all to all communication after each minibatch**

- **Nodes with different performance (even if machines are the same)**
  - **Workload imbalance at barriers (sync points between nodes)**

**Alternative solution: exploit properties of SGD by using asynchronous execution**

# Parameter server design

**Pool of worker nodes**

**Worker
Node 0**

**Worker
Node 1**

**Worker
Node 2**

**Worker
Node 3**

**parameter
values**

**Parameter
Server**

# Training data partitioned among workers



Pool of worker nodes

training data

Worker
Node 0

training data

Worker
Node 1

$x_0 - x_{1000}$

$x_{1000} - x_{2000}$

$x_{2000-3000}$

$x_{3000-4000}$

parameter
values (v0)

Parameter
Server

training data

Worker
Node 2

training data

Worker
Node 3

# Copy of parameters sent to workers

Pool of worker nodes



params v0

params v0

params v0

params v0

training data

local copy of parameters (v0)

Worker Node 0

training data

local copy of parameters (v0)

Worker Node 1

training data

local copy of parameters (v0)

Worker Node 2

training data

local copy of parameters (v0)

Worker Node 3

parameter values (v0)

Parameter Server

# Data parallelism: workers independently compute local "subgradients" on different pieces of data

**Pool of worker nodes**

| training data |
| local copy of parameters (v0) |
| local subgradients |

**Worker Node 0**

| training data |
| local copy of parameters (v0) |
| local subgradients |

**Worker Node 1**

| parameter values (v0) |

**Parameter Server**

| training data |
| local copy of parameters (v0) |
| local subgradients |

**Worker Node 2**

| training data |
| local copy of parameters (v0) |
| local subgradients |

**Worker Node 3**

# Worker sends subgradient to parameter server

**Pool of worker nodes**

**Worker Node 0**
- training data
- local copy of parameters (v0)
- local subgradients

**Worker Node 1**
- training data
- local copy of parameters (v0)
- local subgradients

**subgradient**

**Parameter Server**
- parameter values (v0)

**Worker Node 2**
- training data
- local copy of parameters (v0)
- local subgradients

**Worker Node 3**
- training data
- local copy of parameters (v0)
- local subgradients

# Server updates global parameter values based on subgradient

| Worker Node 0 |
|---|
| training data |
| local copy of parameters (v0) |
| local subgradients |

| Worker Node 1 |
|---|
| training data |
| local copy of parameters (v0) |
| local subgradients |

| Worker Node 2 |
|---|
| training data |
| local copy of parameters (v0) |
| local subgradients |

| Worker Node 3 |
|---|
| training data |
| local copy of parameters (v0) |
| local subgradients |

**Parameter Server**

parameter values (v1)

```
params += -subgrad * step_size;
```

# Updated parameters sent to worker
## Then worker proceeds with another gradient computation step

**params $v_1$**

| Worker Node 0 |
| --- |
| training data |
| local copy of parameters (v0) |
| local subgradients |

| Worker Node 1 |
| --- |
| training data |
| local copy of parameters (v1) |
| local subgradients |

| Parameter Server |
| --- |
| parameter values (v1) |

| Worker Node 2 |
| --- |
| training data |
| local copy of parameters (v0) |
| local subgradients |

| Worker Node 3 |
| --- |
| training data |
| local copy of parameters (v0) |
| local subgradients |

**Notice:**

**Node 1 is operating on different set of parameter values than other nodes**

**Those parameter values were computed without gradient information from the other nodes**

# Updated parameters sent to worker (again)



**Worker Node 0**
- training data
- local copy of parameters (v0)
- local subgradients

**Worker Node 1**
- training data
- local copy of parameters (v1)
- local subgradients

**Worker Node 2**
- training data
- local copy of parameters (v0)
- local subgradients

**Worker Node 3**
- training data
- local copy of parameters (v0)
- local subgradients

**subgradient**

**Parameter Server**
- parameter values (v1)

# Worker continues with updated parameters



Worker Node 0
- training data
- local copy of parameters (v0)
- local subgradients

Worker Node 1
- training data
- local copy of parameters (v1)
- local subgradients

Worker Node 2
- training data
- local copy of parameters (v0)
- local subgradients

Worker Node 3
- training data
- local copy of parameters (v2)
- local subgradients

params $v_2$

Parameter Server
- parameter values (v2)

# Summary: asynchronous parameter update

- **Idea: avoid global synchronization on all parameter updates between each SGD iteration**

  - Algorithm design reflects realities of cluster computing:
    - Slow interconnects
    - Unpredictable machine performance

- **Solution: asynchronous (and partial) subgradient updates**

- **Will impact convergence of SGD**

  - Node N working on iteration $i$ may not have parameter values that result the results of the $i$-1 prior SGD iterations

# Bottleneck?
## What if there is heavy contention for parameter server?

**Worker Node 0**

| training data |
| local copy of parameters (v0) |
| local subgradients |

**Worker Node 1**

| training data |
| local copy of parameters (v1) |
| local subgradients |

**Parameter Server**

| parameter values (v2) |

**Worker Node 2**

| training data |
| local copy of parameters (v0) |
| local subgradients |

**Worker Node 3**

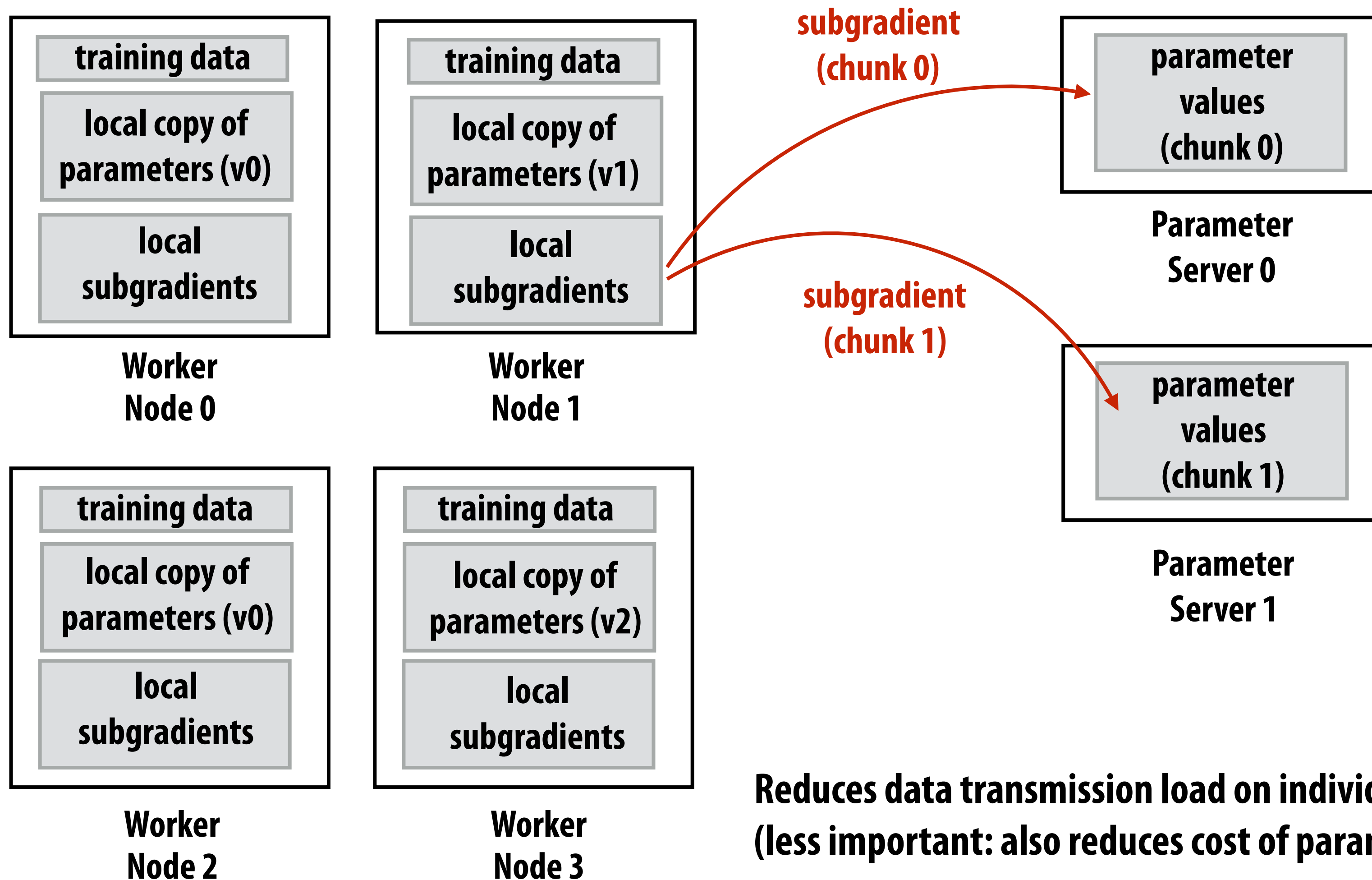| training data |
| local copy of parameters (v2) |
| local subgradients |

# Shard the parameter server

**Partition parameters across servers**

**Worker sends chunk of subgradients to owning parameter server**

| Worker Node 0 | Worker Node 1 |
|---|---|
| training data | training data |
| local copy of parameters (v0) | local copy of parameters (v1) |
| local subgradients | local subgradients |

**subgradient (chunk 0)**

**subgradient (chunk 1)**

parameter values (chunk 0)

**Parameter Server 0**

parameter values (chunk 1)

**Parameter Server 1**

| Worker Node 2 | Worker Node 3 |
|---|---|
| training data | training data |
| local copy of parameters (v0) | local copy of parameters (v2) |
| local subgradients | local subgradients |

**Reduces data transmission load on individual servers (less important: also reduces cost of parameter update)**

# What if model parameters do not fit on one worker?

**Recall high footprint of training large networks
(particularly with large mini-batch sizes)**

| Worker Node 0 | Worker Node 1 | Parameter Server 0 |
|---|---|---|
| training data | training data | parameter values (chunk 0) |
| local copy of parameters (v0) | local copy of parameters (v1) | |
| local subgradients | local subgradients | |

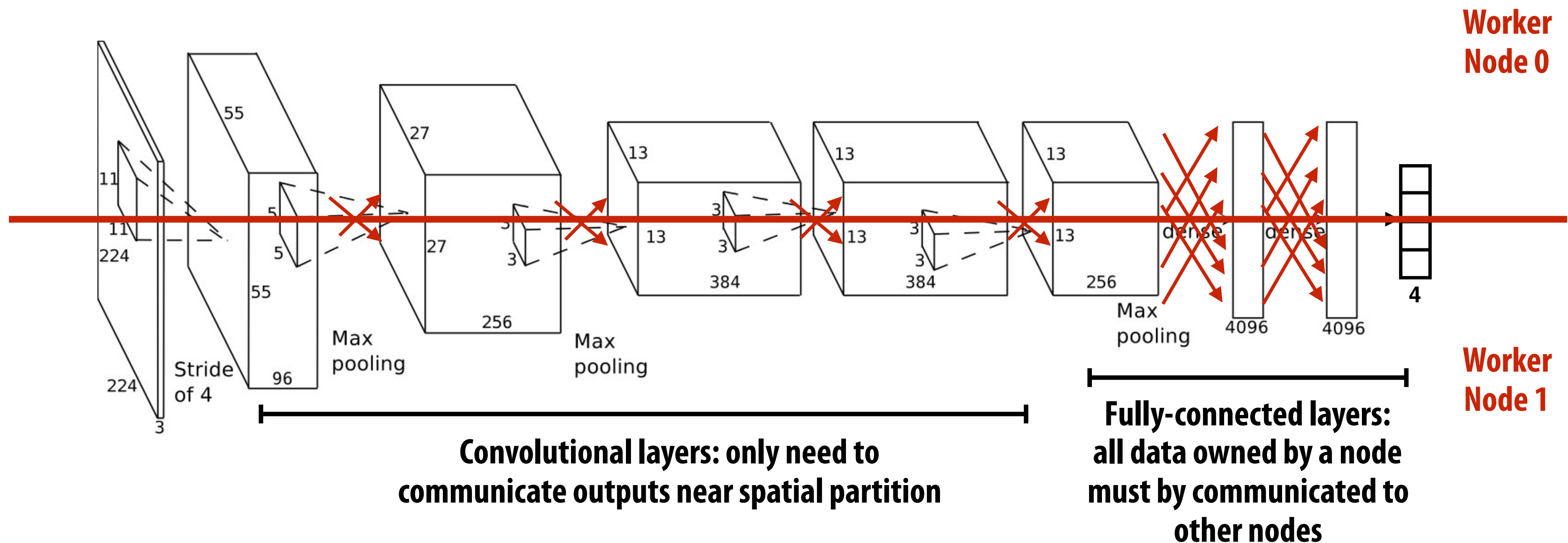| Worker Node 2 | Worker Node 3 | Parameter Server 1 |
|---|---|---|
| training data | training data | parameter values (chunk 1) |
| local copy of parameters (v0) | local copy of parameters (v2) | |
| local subgradients | local subgradients | |

# Model parallelism

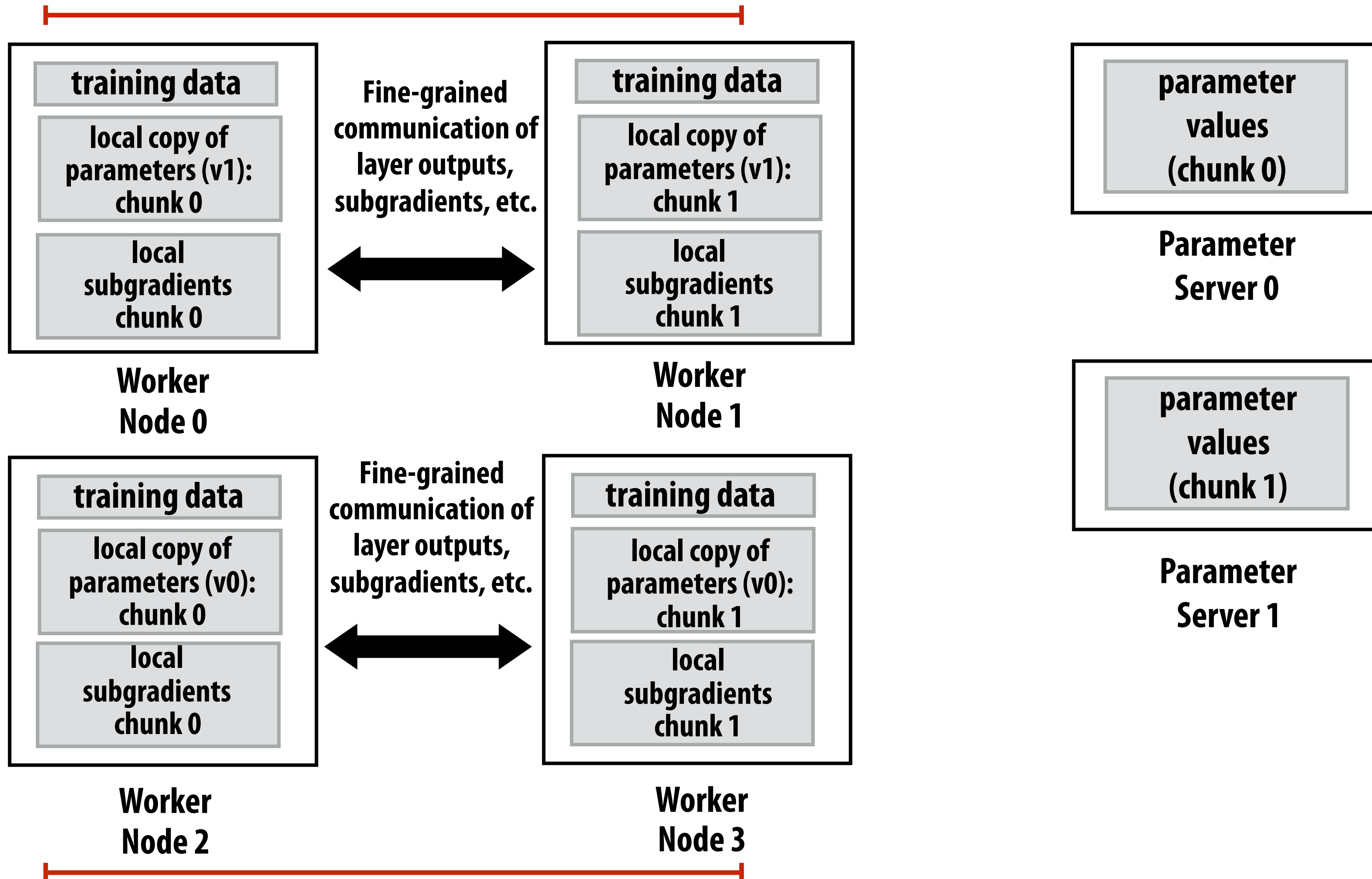## Partition network parameters across nodes (spatial partitioning to reduce communication)

## Reduce internode communication through network design:

- Use small spatial convolutions (1x1 convolutions)
- Reduce/shrink fully-connected layers



**Worker Node 0**

**Worker Node 1**

Convolutional layers: only need to communicate outputs near spatial partition

Fully-connected layers: all data owned by a node must by communicated to other nodes

# Data-parallel and model-parallel execution

**Working on subgradient computation for a single copy of the model**

| | | |
|---|---|---|
| **training data** | | **parameter values (chunk 0)** |
| **local copy of parameters (v1): chunk 0** | **Fine-grained communication of layer outputs, subgradients, etc.** | |
| **local subgradients chunk 0** | | |

**training data**

**local copy of parameters (v1): chunk 1**

**local subgradients chunk 1**

**Parameter Server 0**

**Worker Node 0**

**Worker Node 1**

**training data**

**local copy of parameters (v0): chunk 0**

**local subgradients chunk 0**

**Fine-grained communication of layer outputs, subgradients, etc.**

**training data**

**local copy of parameters (v0): chunk 1**

**local subgradients chunk 1**

**parameter values (chunk 1)**

**Parameter Server 1**

**Worker Node 2**

**Worker Node 3**

**Working on subgradient computation for a single copy of the model**
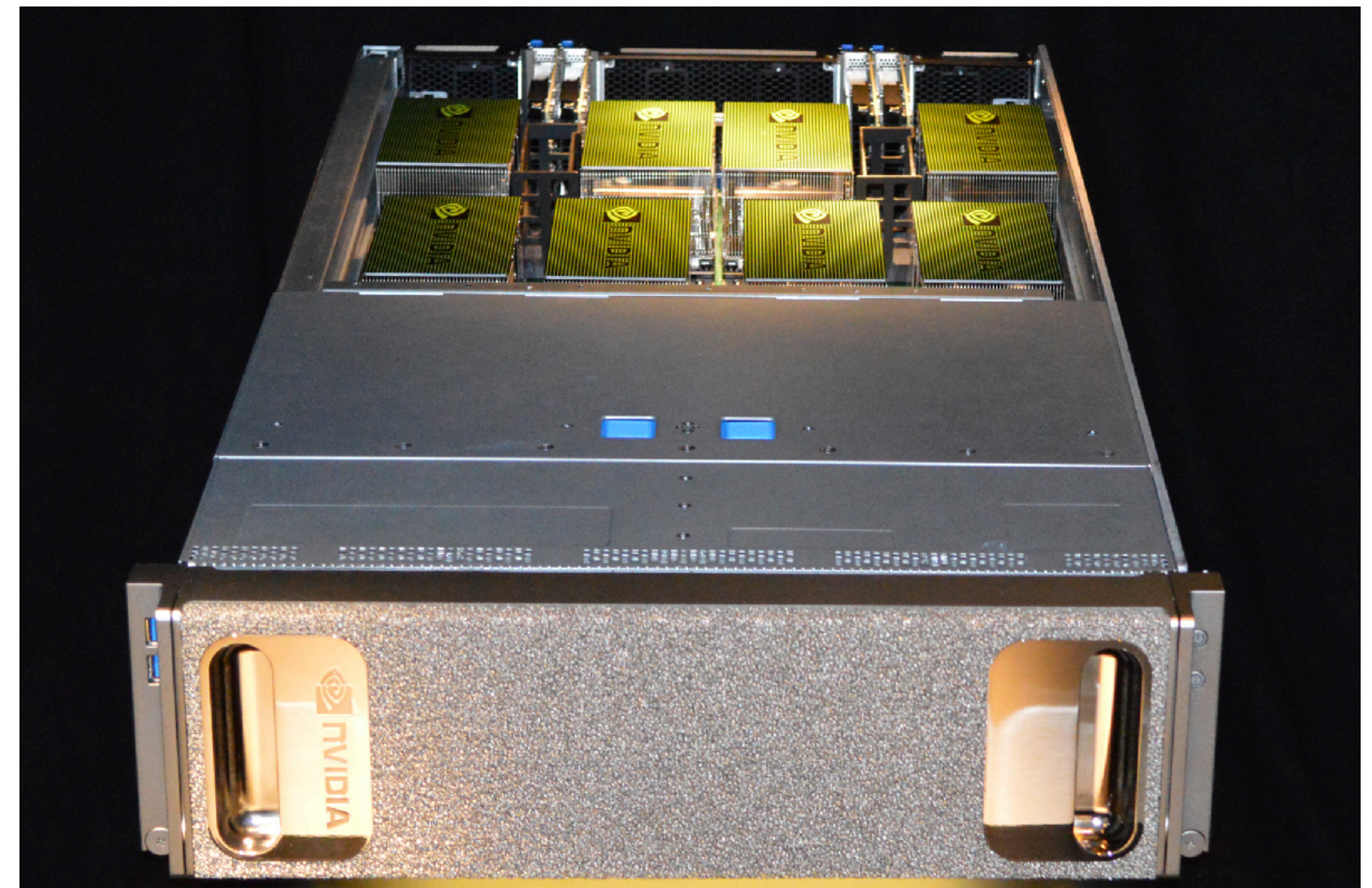
# Asynchronous vs. synchronous debate

- **Asynchronous training: significant distributed system complexity incurred to combat bandwidth/latency constraints of modern cluster computing**

- **Interest in ways to improve scalability of synchronous training**
  - **Better hardware**
  - **Better algorithms for existing hardware**

# Better hardware: using supercomputers for training

- **Fast interconnects critical for model-parallel training**
  - Fine-grained communication of outputs and gradients

- **Fast interconnects diminish need for async training algorithms**
  - Avoid randomness in training due to schedule of computation (there remains randomness due to stochastic part of SGD algorithm)



**OakRidge Titan Supercomputer
(low-latency interconnect)**



**NVIDIA DGX-1: 8 GPUs connected via
high speed NV-Link interconnect
($150,000 in 2018)**

# Modified algorithmic techniques (again): improving scalability of synchronous training…

- **Larger mini-batches increase compute-to-communication ratio: communicate gradients summed over B training inputs**

```
for each item x in mini-batch on this node:
    grad += evaluate_loss_gradient(f, loss_func, params, x)
barrier();
sum reduce gradients across all nodes, communicate results to all nodes
barrier();
update copy of local parameter values
```

- **But large mini-batches (if used naively) reduce accuracy of model trained**

# Increasing learning rate with mini-batch size: linear scaling rule

**Recall: minibatch SGD parameter update**

$$w_{t+1} = w_t - \eta \frac{1}{n} \sum_{x \in \mathcal{B}} \nabla l(x, w_t)$$

**Consider processing of k minibatches (k steps of gradient descent)**

$$w_{t+k} = w_t - \eta \frac{1}{n} \sum_{j<k} \sum_{x \in \mathcal{B}_j} \nabla l(x, w_{t+j})$$

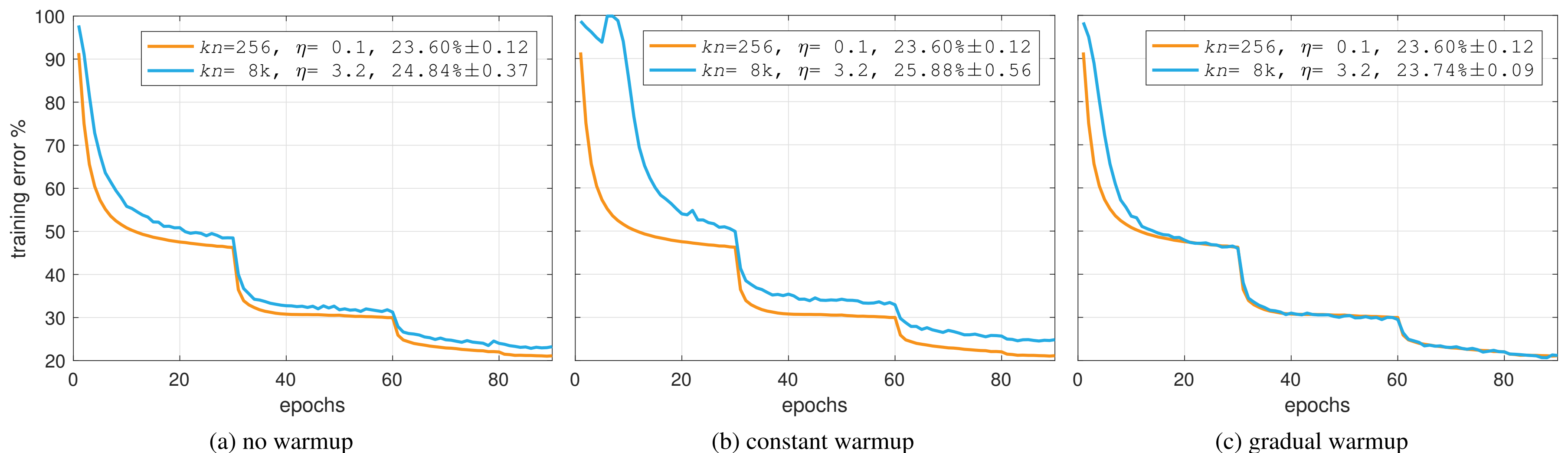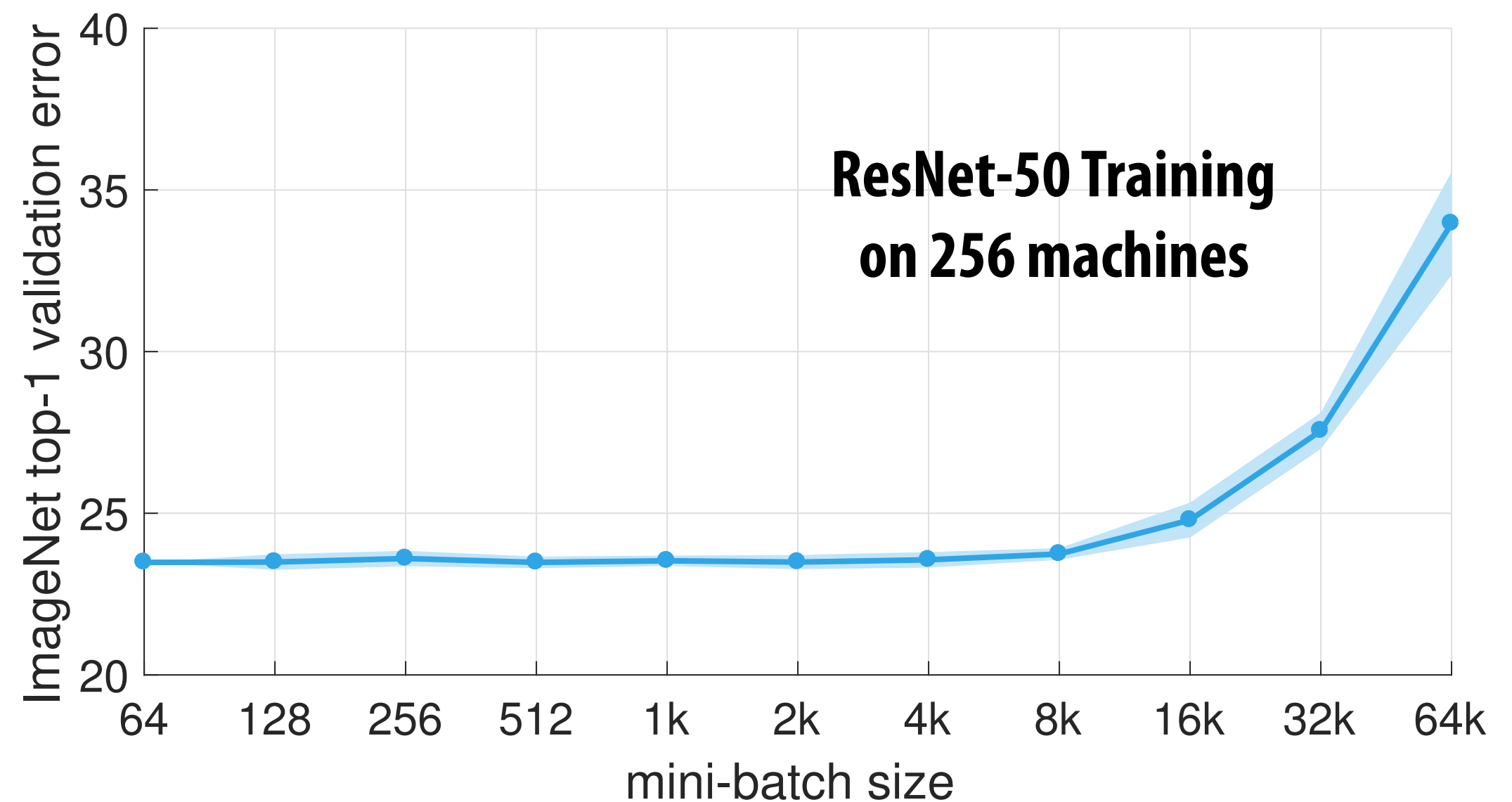**Consider processing one minibatch that is of size kn (one step of gradient descent)**

$$\hat{w}_{t+1} = w_t - \hat{\eta} \frac{1}{kn} \sum_{j<k} \sum_{x \in \mathcal{B}_j} \nabla l(x, w_t)$$

**Suggests that if** $\nabla l(x, w_t) \approx \nabla l(x, w_{t+j})$ **for *j* < *k* then minibatch SGD with size *n* and learning rate** $\eta$ **can be approximated by large mini batch SGD with size *kn* <span style="color:red">if the learning rate is also scaled to</span>** $k\eta$

# When does $\nabla l(x, w_t) \approx \nabla l(x, w_{t+j})$ not hold?

1. **At beginning of training**
   - **Suggests starting training with smaller learning rate (learning rate "warmup")**

2. **When minibatch size begins to get too large (there is a limit to scaling minibatch size)**



ResNet-50 Training on 256 machines



(a) no warmup

| | |
|---|---|
| $kn$=256, $\eta$= 0.1, 23.60%±0.12 | |
| $kn$= 8k, $\eta$= 3.2, 24.84%±0.37 | |

(b) constant warmup

| | |
|---|---|
| $kn$=256, $\eta$= 0.1, 23.60%±0.12 | |
| $kn$= 8k, $\eta$= 3.2, 25.88%±0.56 | |

(c) gradual warmup

| | |
|---|---|
| $kn$=256, $\eta$= 0.1, 23.60%±0.12 | |
| $kn$= 8k, $\eta$= 3.2, 23.74%±0.09 | |

**Minibatch size = 256 (orange) vs. 8192 (blue)**

# Gradient compression

- **Each node computes gradients for minibatch, but only sends gradients with magnitude above a threshold**

- **Locally accumulate gradients below threshold over multiple SGD steps (then send when exceed threshold)**

$$G_0^k = 0$$

**for all iterations _t_:**

$$G_t^k = G_{t-1}^k + \eta \frac{1}{Nb} \sum_{k=1}^{N} \sum_{x \in B_k} \nabla f(x; w_t)$$

**Compress and send ONLY the elements of** $G_t^k$ **greater than threshold.**
**(then locally zero out sent elements)**

**SGD update on each note only uses the sent weights.**

# Handling momentum

**Consider basic momentum in SGD:**

$$u_t = mu_{t-1} + \sum_{k=1}^{N} \left( \nabla_{k,t} \right), \quad w_{t+1} = w_t - \eta u_t \qquad \boxed{\nabla_{k,t} = \frac{1}{Nb} \sum_{x \in \mathcal{B}_{k,t}} \nabla f(x, w_t)}$$

**Consider weight update with momentum after T iterations of SGD**

$$w_{t+T}^{(i)} = w_t^{(i)} - \eta \left[ \cdots + \left( \sum_{\tau=0}^{T-2} m^\tau \right) \nabla_{k,t+1}^{(i)} + \left( \sum_{\tau=0}^{T-1} m^\tau \right) \nabla_{k,t}^{(i)} \right]$$

**Basic sparse update: (what's the problem?)**

$$v_{k,t} = v_{k,t-1} + \nabla_{k,t}, \quad u_t = mu_{t-1} + \sum_{k=1}^{N} sparse\left(v_{k,t}\right), \quad w_{t+1} = w_t - \eta u_t$$

**Problem: momentum discount not applied correctly after sparse update interval T (assume sparse gradients propagated after T iterations of SGD)**

$$w_{t+T}^{(i)} = w_t^{(i)} - \eta \left( \cdots + \nabla_{k,t+1}^{(i)} + \nabla_{k,t}^{(i)} \right)$$

**Fix: locally accumulate and communicate gradient velocities, not gradients:**

$$u_{k,t} = mu_{k,t-1} + \nabla_{k,t}, \quad v_{k,t} = v_{k,t-1} + u_{k,t}, \quad w_{t+1} = w_t - \eta \sum_{k=1}^{N} sparse\left(v_{k,t}\right)$$

# Summary: training large networks in parallel

- **Data-parallel training with asynchronous update to efficiently use clusters of commodity machines with low speed interconnect**
  - Modification of SGD algorithm to meet constraints of modern parallel systems
  - Effects on convergence are problem dependent and not particularly well understood
  - Efficient use of fast interconnects may provide alternative to these methods (facilitate tightly orchestrated solutions much like supercomputing applications)

- **Modern DNN designs, large minibatch sizes, careful learning rate schedules enable scalability without asynchronous execution on commodity clusters**

- **High-performance training of deep networks is an interesting example of constant iteration of algorithm design and parallelization strategy (a key theme of this course!)**