

**Lecture 19:**

# **Parallel Deep Network Training**

---

**Parallel Computer Architecture and Programming**

**CMU / 清华大学, Summer 2017**



神威

太湖之光

国家并行计算机工程技术研究中心

# How would you describe this professor?



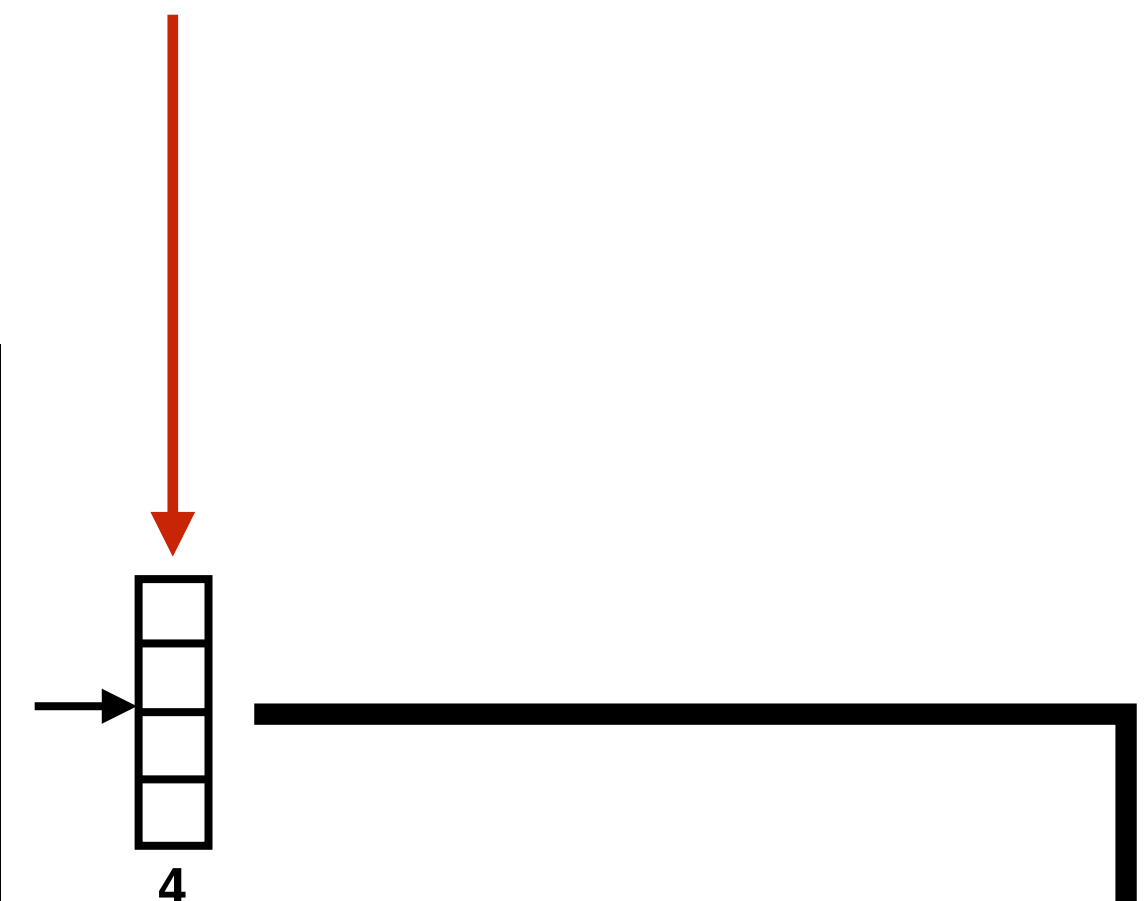
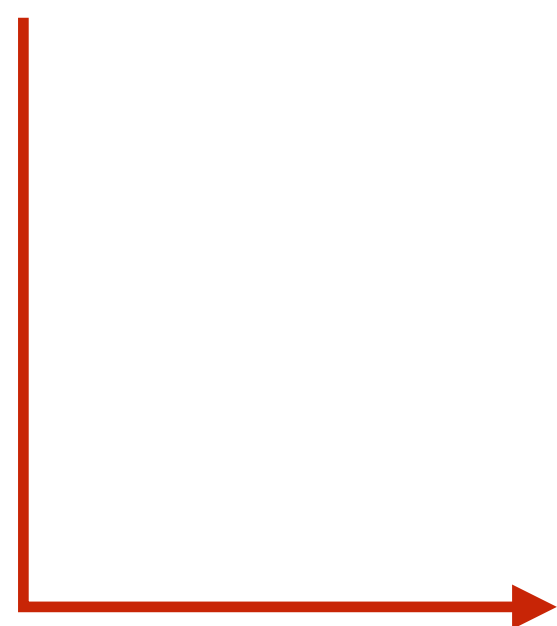
**Easy?**  
**Mean?**  
**Boring?**  
**Nerdy?**

# Professor classification task

Classifies professors as easy, mean, boring, or nerdy based on their appearance.

**Input:**  
image of a professor

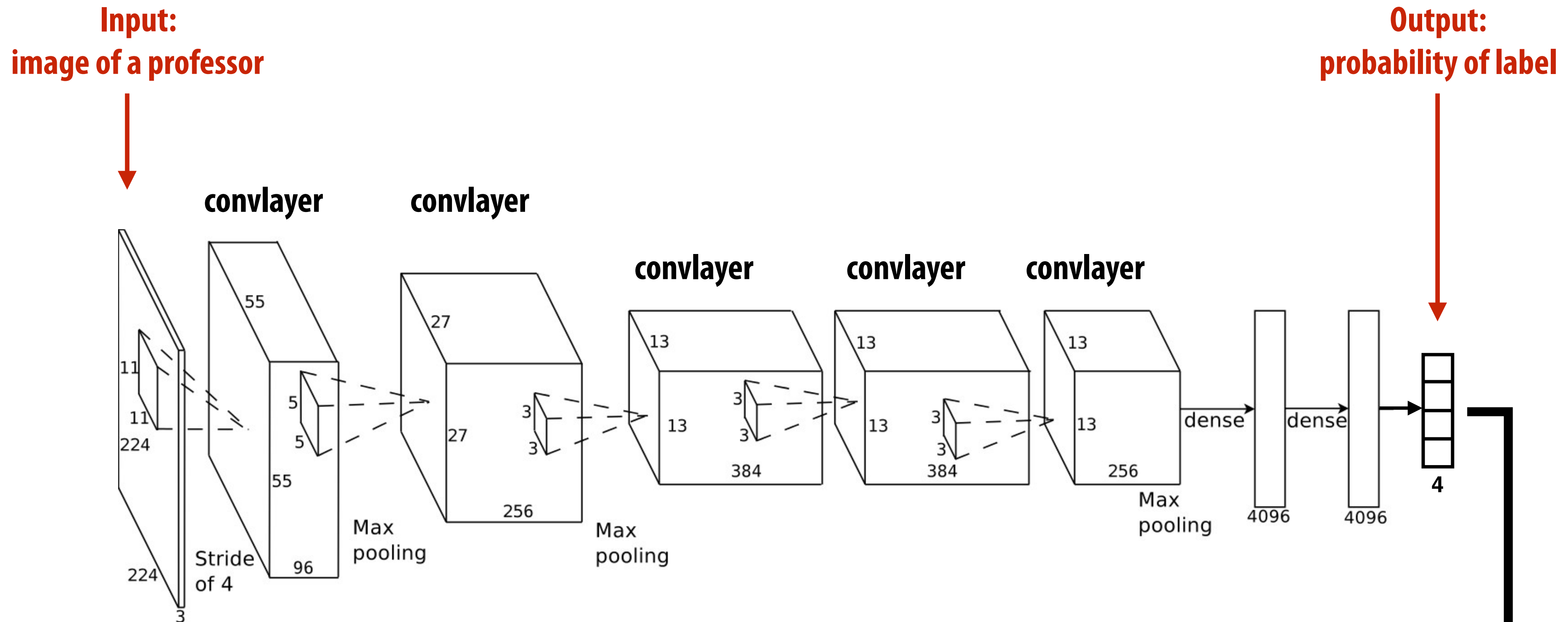
**Output:**  
probability of each of four possible labels



Easy: ??  
Mean: ??  
Boring: ??  
Nerdy: ??

# Professor classification network

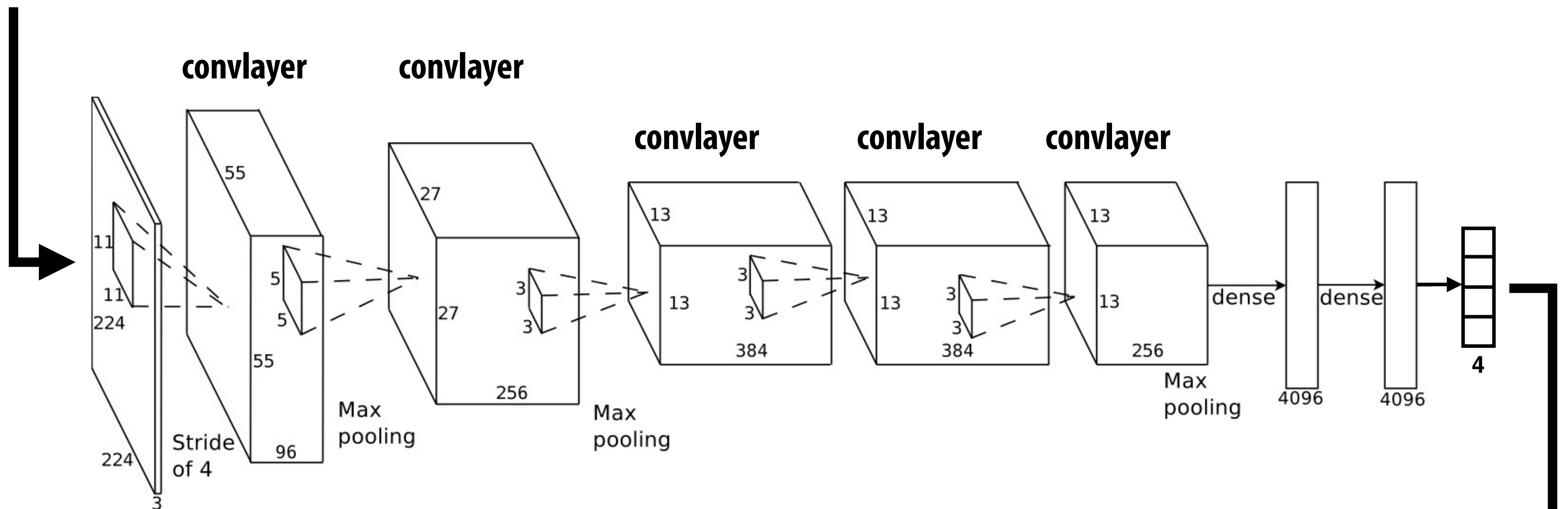
Classifies professors as easy, mean, boring, or nerdy based on their appearance.



Recall: large networks may have  
10's-100's of millions of parameters

Easy: ??  
Mean: ??  
Boring: ??  
Nerdy: ??

# Professor classification network



**Easy: 0.26**  
**Mean: 0.08**  
**Boring: 0.14**  
**Nerdy: 0.52**

# Training data (ground truth answers)



[label omitted]



[label omitted]



[label omitted]



Nerdy



[label omitted]



[label omitted]



[label omitted]



[label omitted]



[label omitted]



Nerdy



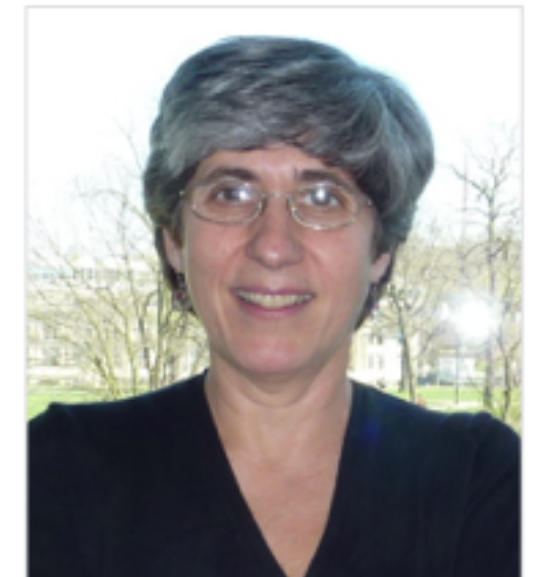
[label omitted]



[label omitted]



Nerdy



[label omitted]



[label omitted]



[label omitted]



Nerdy



[label omitted]



[label omitted]



[label omitted]



Nerdy

# Professor classification network



New image of Kayvon  
(not in training set)

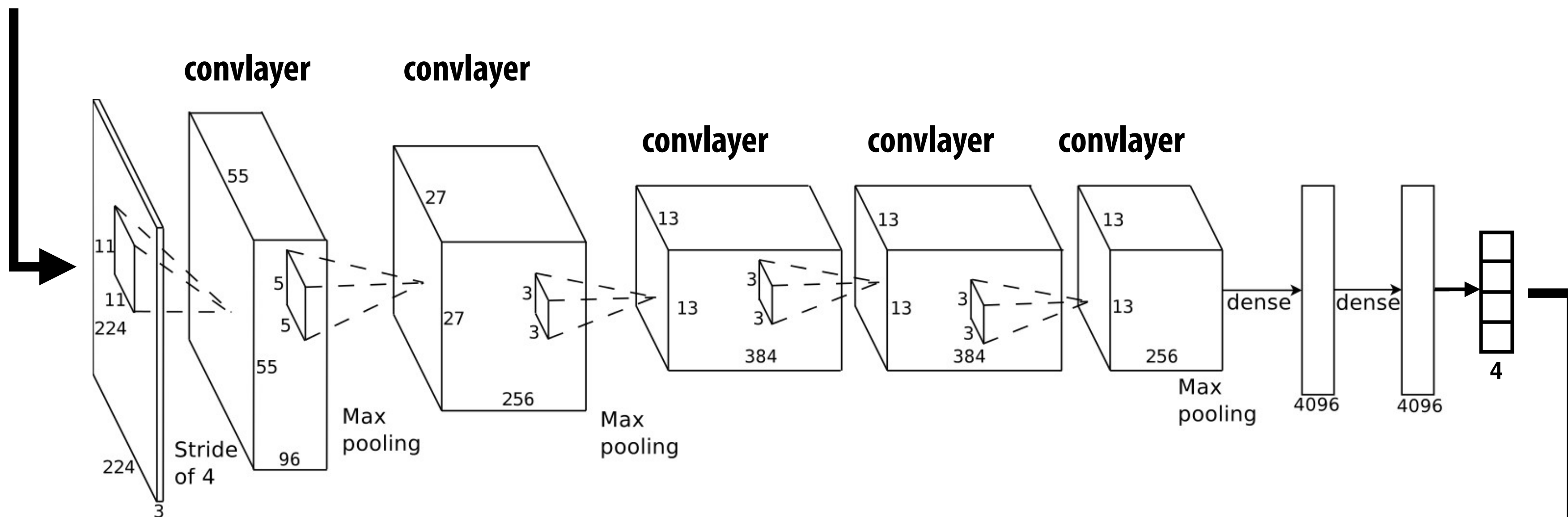
Ground truth  
(what the answer should be)

Easy: 0.0

Mean: 0.0

Boring: 0.0

Nerdy: 1.0



Easy: 0.26

Mean: 0.08

Boring: 0.14

Nerdy: 0.52

Network output



# Error (loss)

**Ground truth:**  
**(what the answer should be)**

**Easy: 0.0**  
**Mean: 0.0**  
**Boring: 0.0**  
**Nerdy: 1.0**

**Network output: \***

**Easy: 0.26**  
**Mean: 0.08**  
**Boring: 0.14**  
**Nerdy: 0.52**

**Common example: softmax loss:**

$$L = -\log \left( \frac{e^{f_c}}{\sum_j e^{f_j}} \right)$$

Output of network for correct category

Output of network for all categories

\* In practice a network using a softmax classifier outputs unnormalized, log probabilities ( $f_j$ ), but I'm showing a probability distribution above for clarity

# Training

**Goal of training: learning good values of network parameters so that the network outputs the correct classification result for any input image**

**Idea: minimize loss for all the training examples (for which the correct answer is known)**

$$L = \sum_i L_i \quad (\text{total loss for entire training set is sum of losses } L_i \text{ for each training example } x_i)$$

**Intuition: if the network gets the answer correct for a wide range of training examples, then hopefully it has learned parameter values that yield the correct answer for future images as well.**

# Intuition: gradient descent

Say you had a function  $f$  that contained hidden parameters  $p_1$  and  $p_2$ :  $f(x_i)$

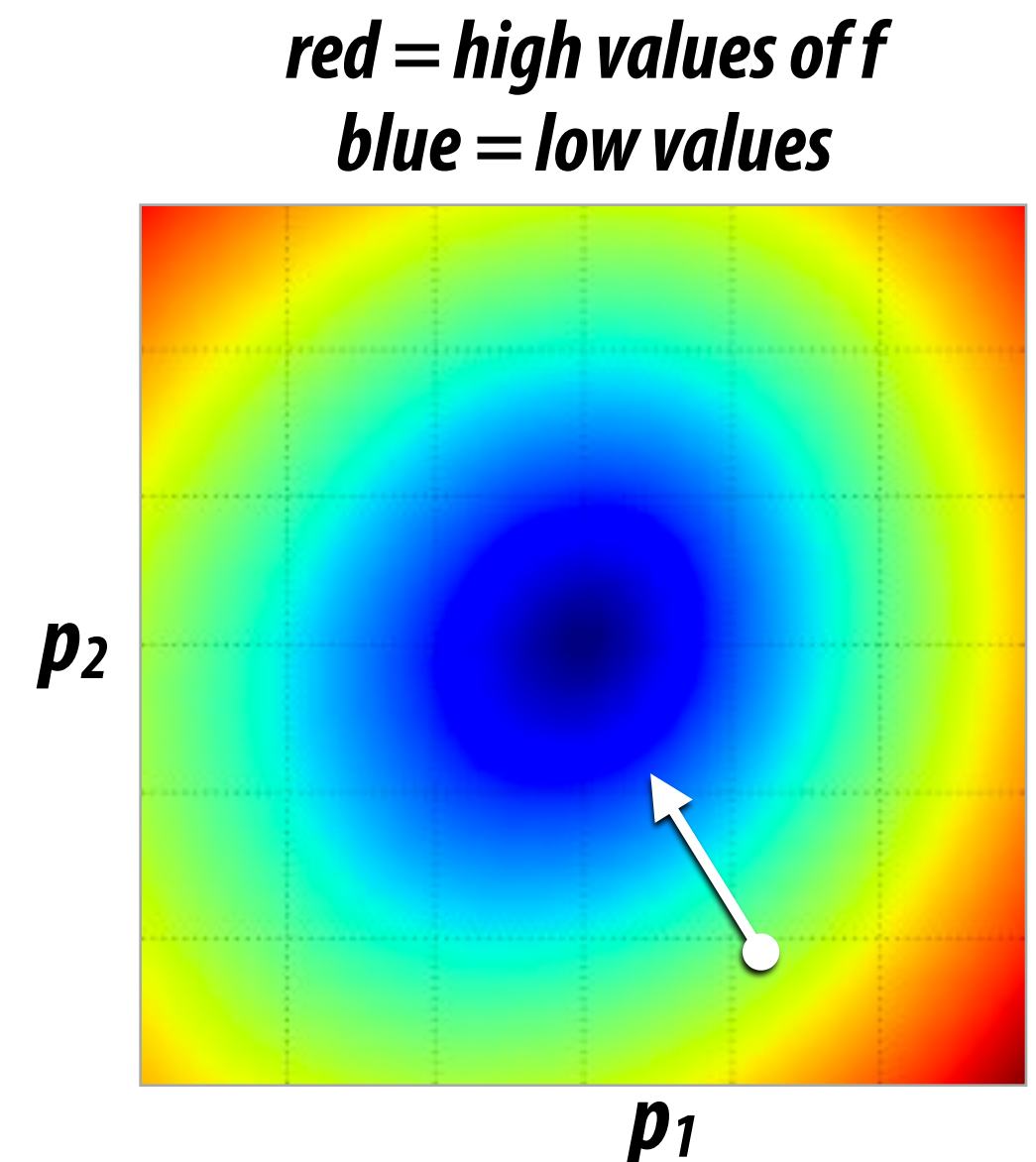
And for some input  $x_i$ , your training data says the function should output 0.

But for the current values of  $p_1$  and  $p_2$ , it currently outputs 10.

$$f(x_i, p_1, p_2) = 10$$

And say I also gave you expressions for the derivative of  $f$  with respect to  $p_1$  and  $p_2$  so you could compute their value at  $x_i$ .

$$\frac{df}{dp_1} = 2 \quad \frac{df}{dp_2} = -5 \quad \nabla f = [2, -5]$$



How might you adjust the values  $p_1$  and  $p_2$  to reduce the error for this training example?

# Basic gradient descent

```
while (loss too high):  
    for each item  $x_i$  in training set:  
        grad += evaluate_loss_gradient(f, params, loss_func,  $x_i$ )  
  
    params += -grad * step_size;
```

**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**

**How do we compute  $d\text{Loss}/dp$  for a deep neural network with millions of parameters?**

# Derivatives using the chain rule

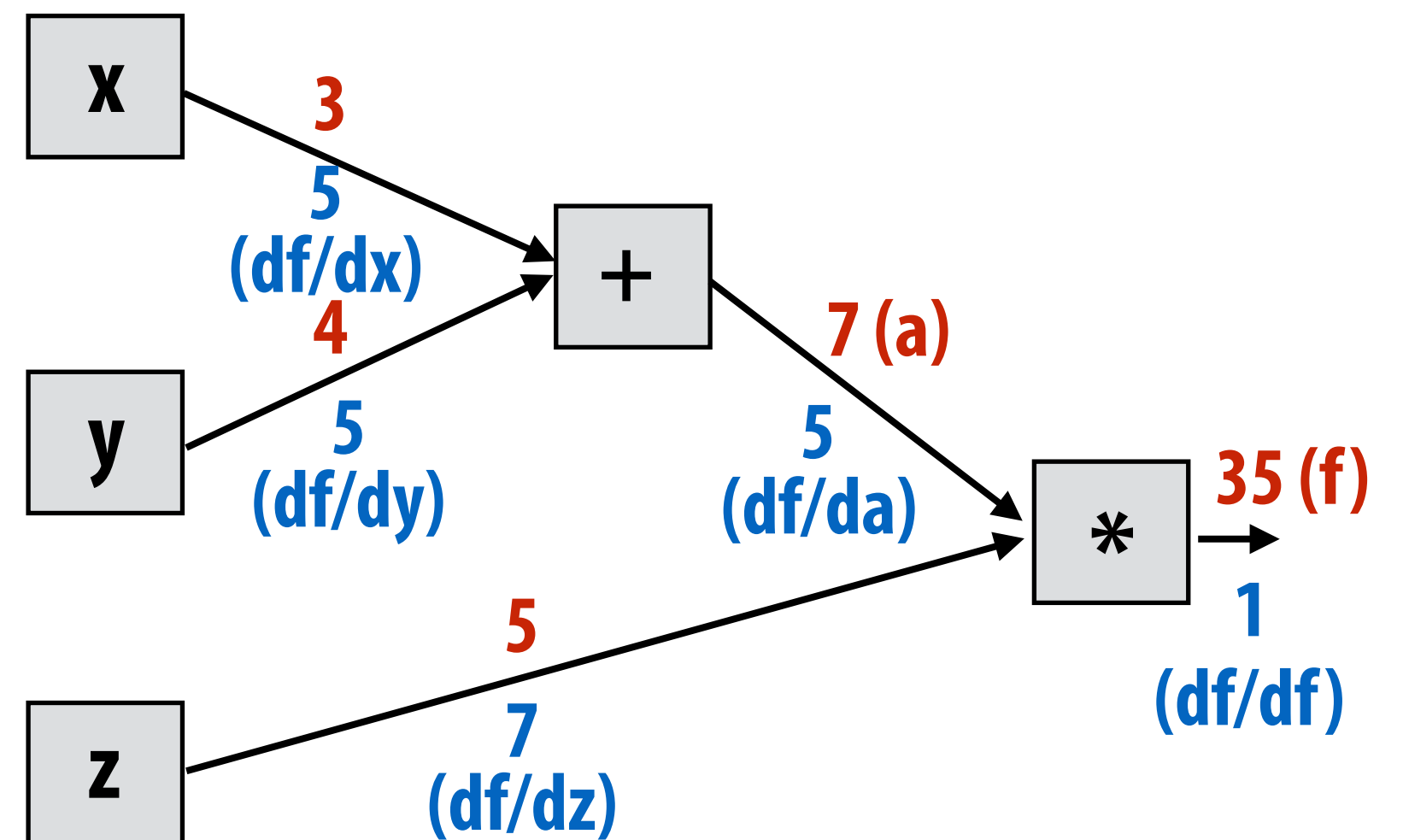
$$f(x, y, z) = (x + y)z = az$$

Where:  $a = x + y$

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

So, by the derivative chain rule:

$$\frac{df}{dx} = \frac{df}{da} \frac{da}{dx} = z$$



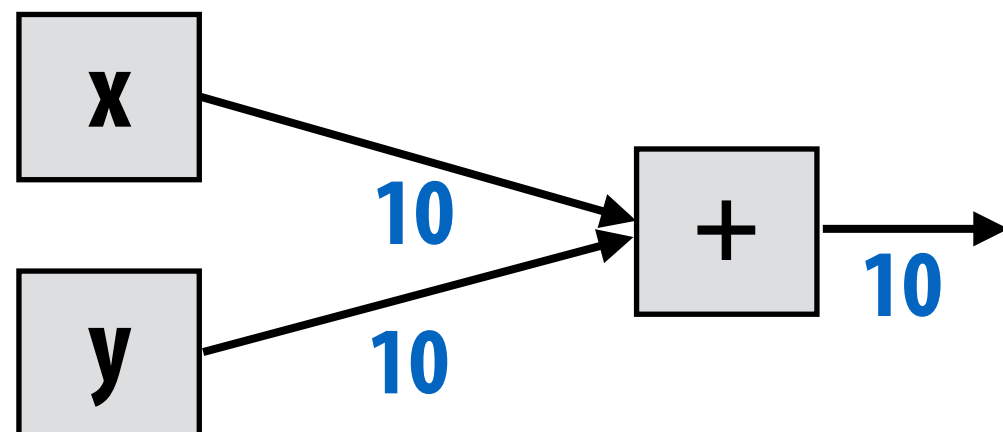
Red = output of node  
Blue =  $df/dnode$

# Backpropagation

Red = output of node

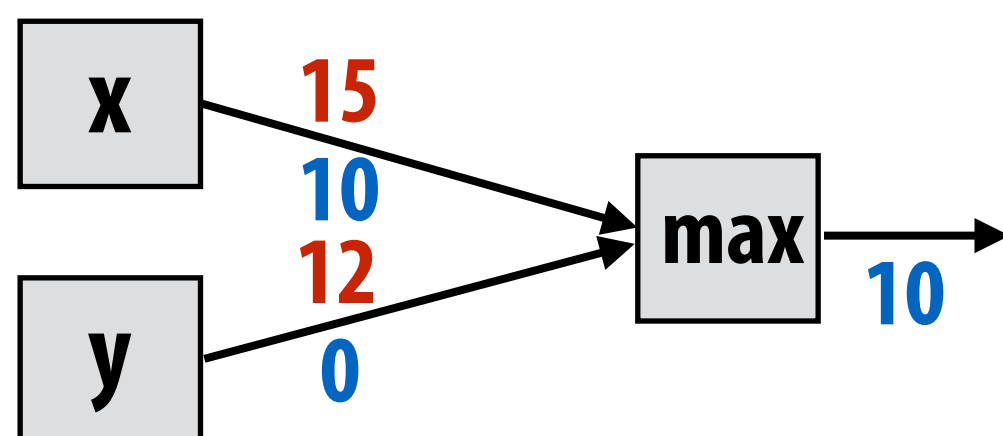
Blue =  $df/dnode$

Recall:  $\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$



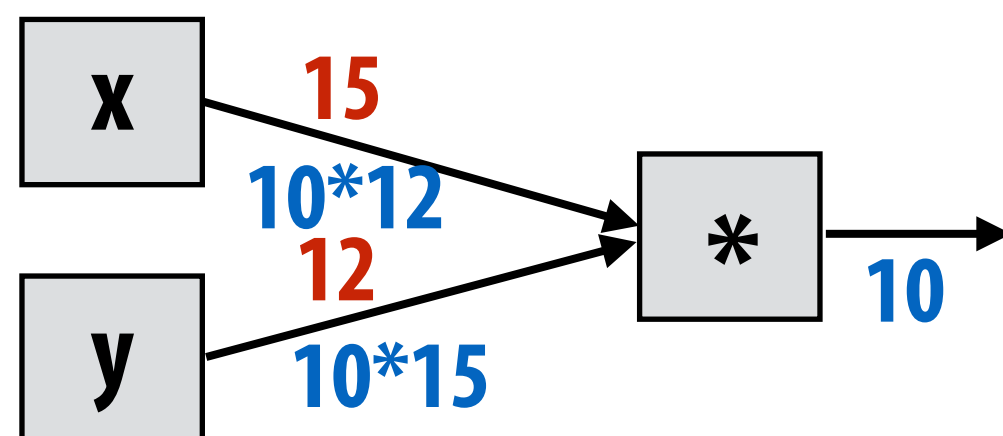
$$g(x, y) = x + y$$

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



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

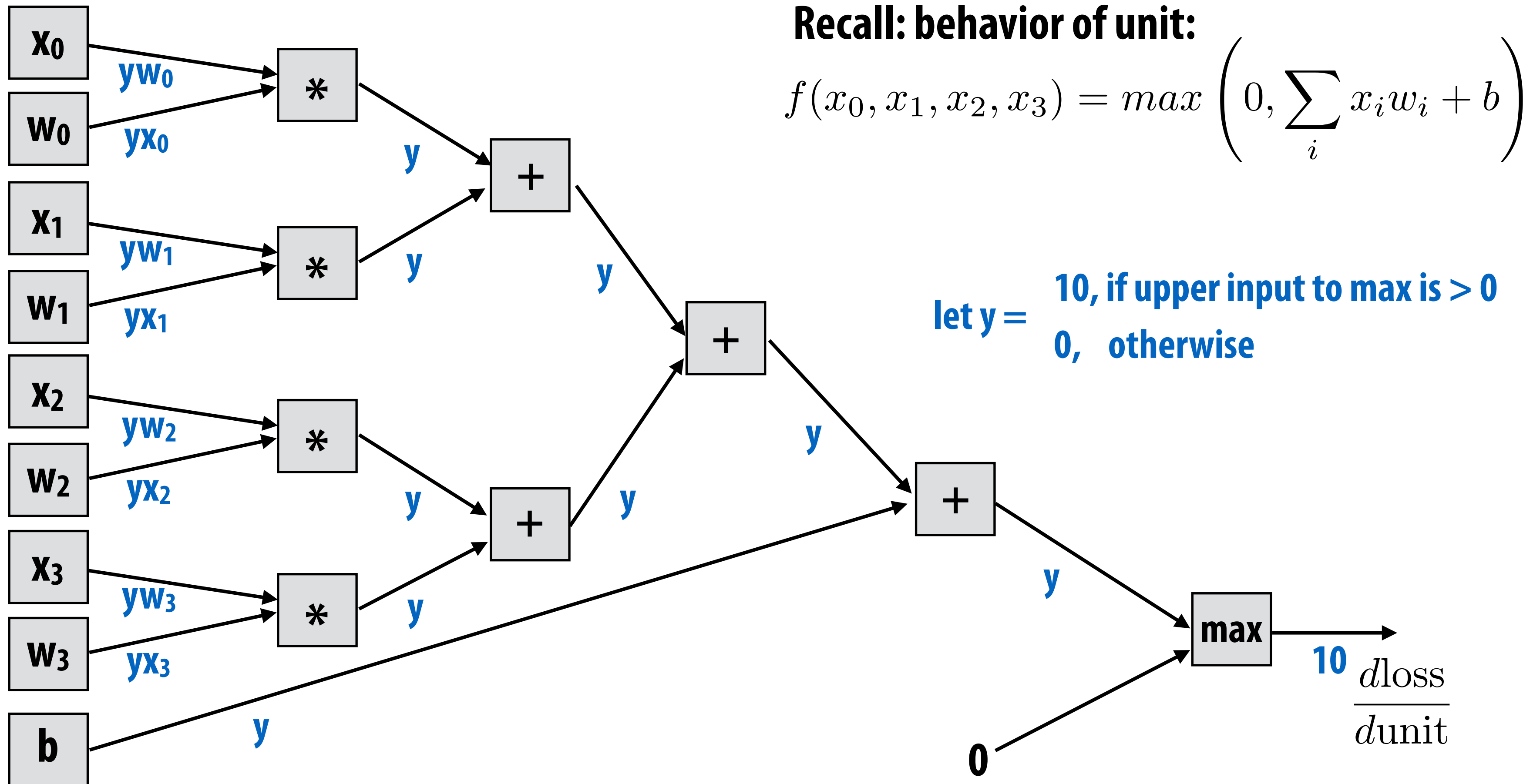
$$\frac{dg}{dx} = \begin{cases} 1, & \text{if } x > y \\ 0, & \text{otherwise} \end{cases}$$



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

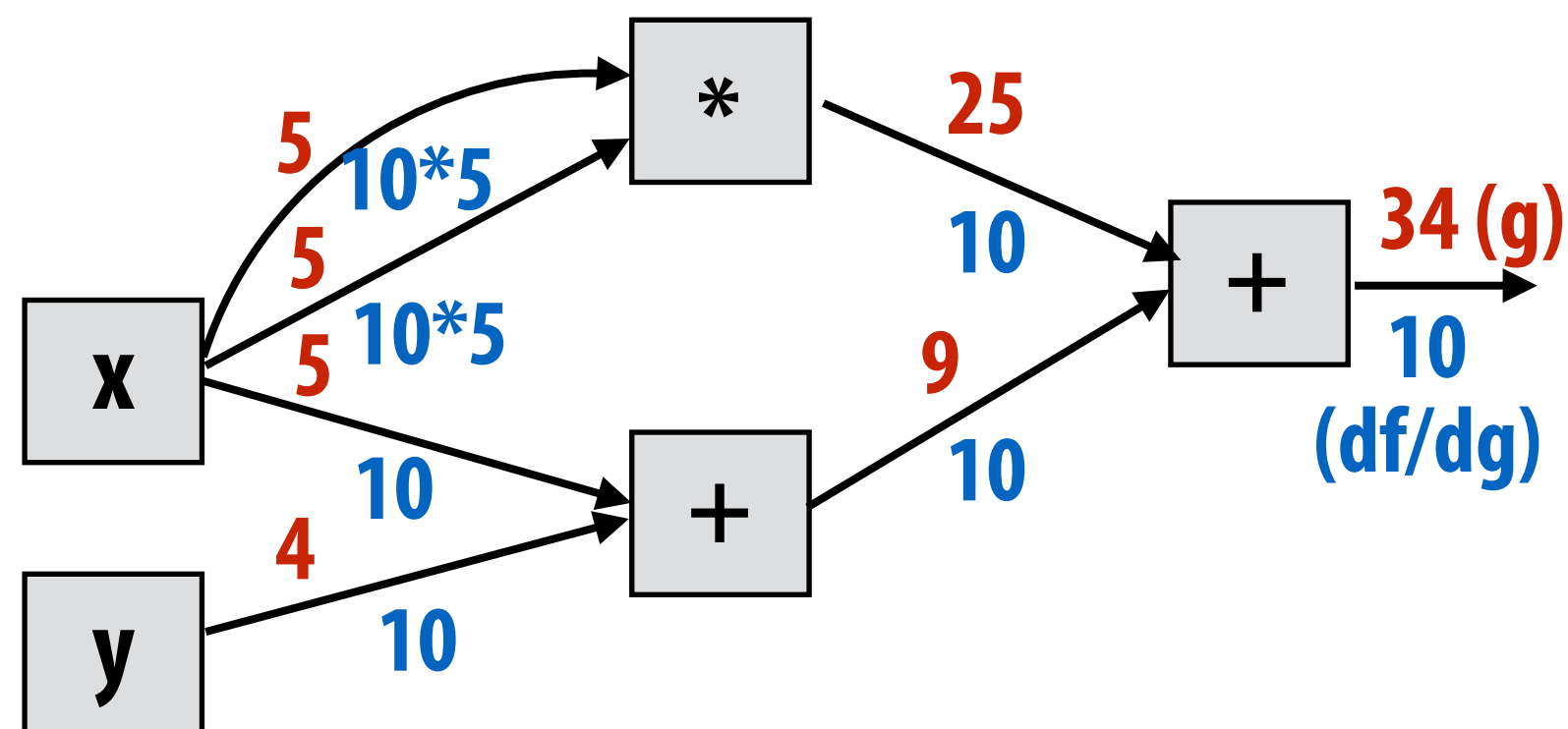
$$\frac{dg}{dx} = y, \frac{dg}{dy} = x$$

# Back-propagating through single unit



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

# Multiple uses of an input variable



Sum gradients from each use of variable:

Here:

$$\begin{aligned} \frac{df}{dx} &= \frac{df}{dg} \frac{dg}{dx} \\ &= 10 \frac{dg}{dx} \\ &= 10(2x + 1) \\ &= 10(10 + 1) = 110 \end{aligned}$$

$$g(x, y) = (x + y) + x * x = a + b$$

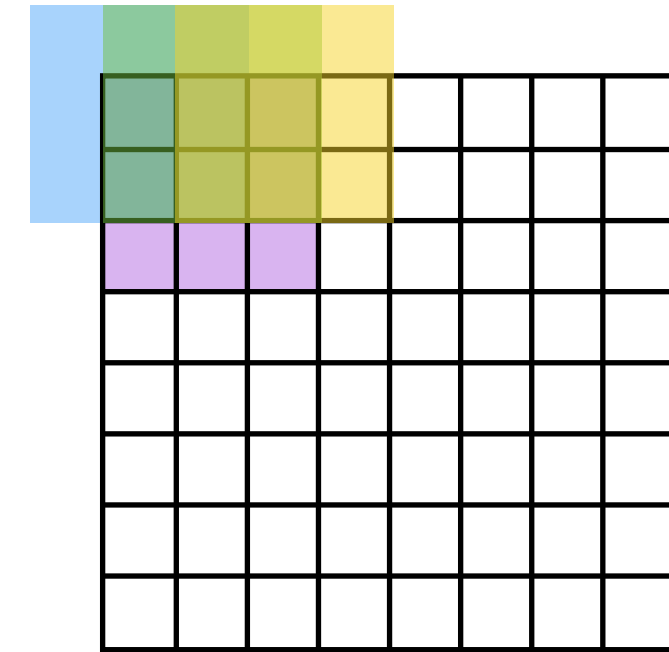
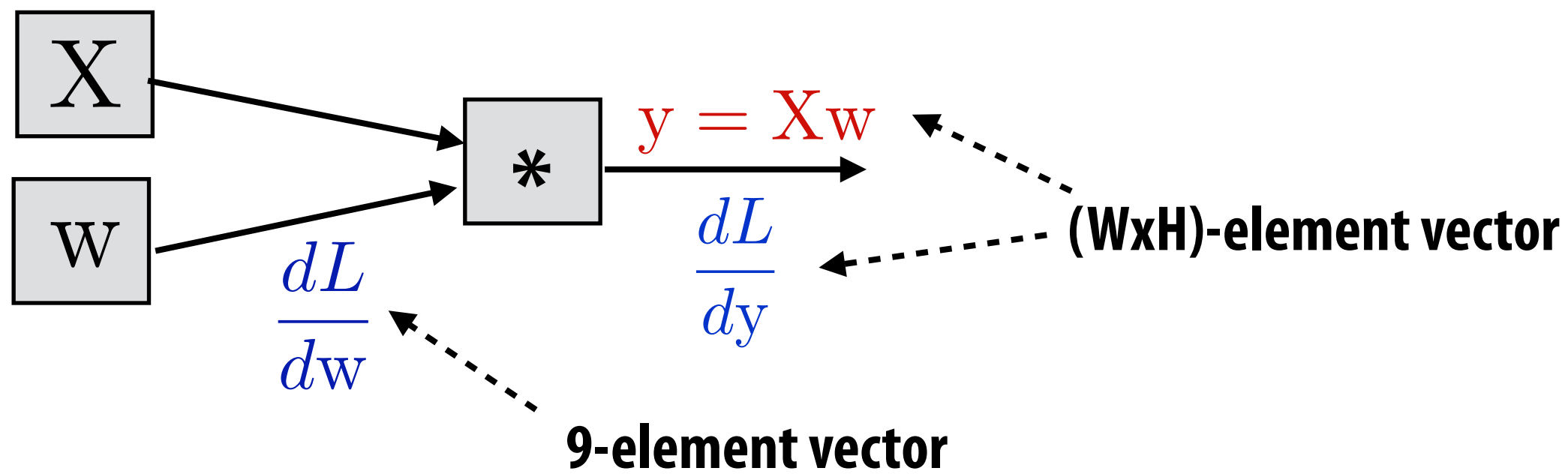
$$\frac{da}{dx} = 1, \quad \frac{db}{dx} = 2x$$

$$\frac{dg}{dx} = \frac{dg}{da} \frac{da}{dx} + \frac{dg}{db} \frac{db}{dx} = 2x + 1$$

**Implication: backpropagation through all units in a convolutional layer adds gradients computed from each unit to the overall gradient for the shared weights**



# Back-propagation: matrix form



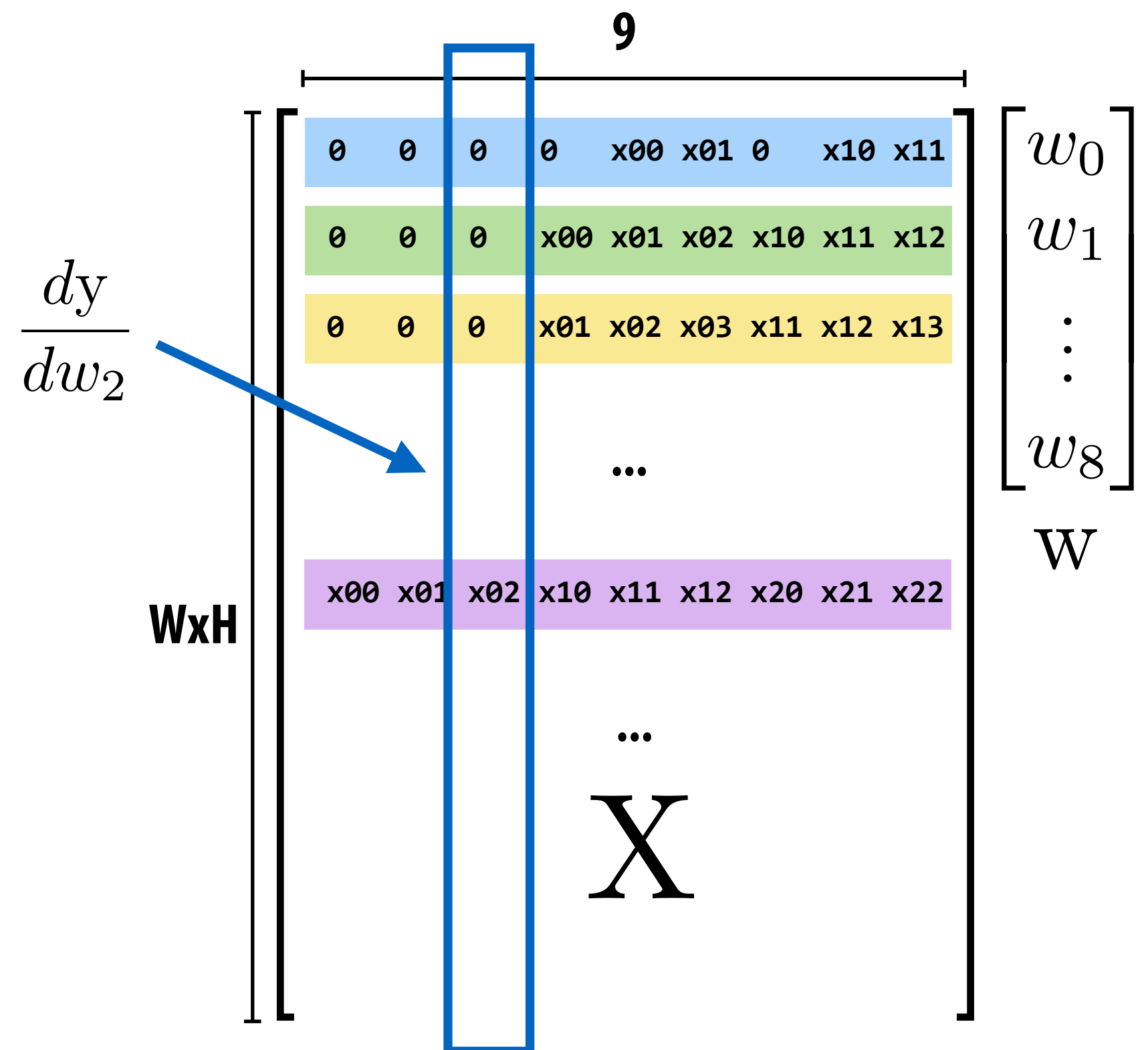
$$\frac{dy_j}{dw_i} = X_{ji}$$

$$\frac{dL}{dw_i} = \sum_j \frac{dL}{dy_j} \frac{dy_j}{dw_i}$$

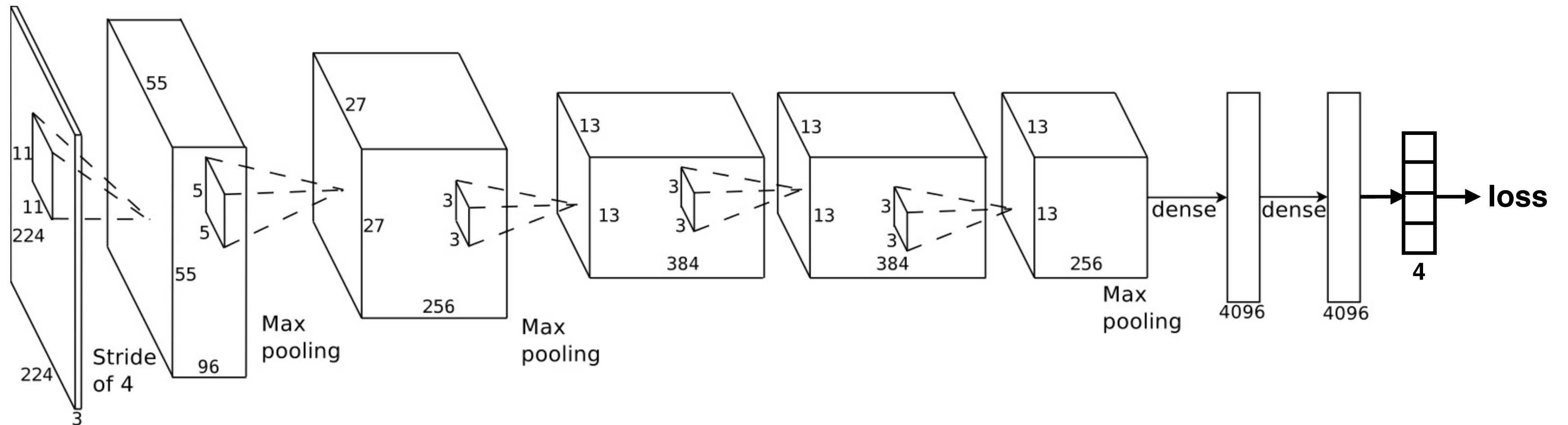
$$= \sum_j \frac{dL}{dy_j} X_{ji}$$

Therefore:

$$\frac{dL}{dw} = X^T \frac{dL}{dy}$$



# Backpropagation through the entire professor classification network



**For each training example  $x_i$  in mini-batch:**

**Perform forward evaluation to compute loss for  $x_i$**

**Compute gradient of loss w.r.t. final layer's outputs**

**Backpropagate gradient to compute gradient of loss w.r.t. all network parameters**

**Accumulate gradients (over all images in batch)**

**Update all parameter values:  $w_{\text{new}} = w_{\text{old}} - \text{step\_size} * \text{grad}$**

# Recall from last class: VGG memory footprint

Calculations assume 32-bit values (image batch size = 1)

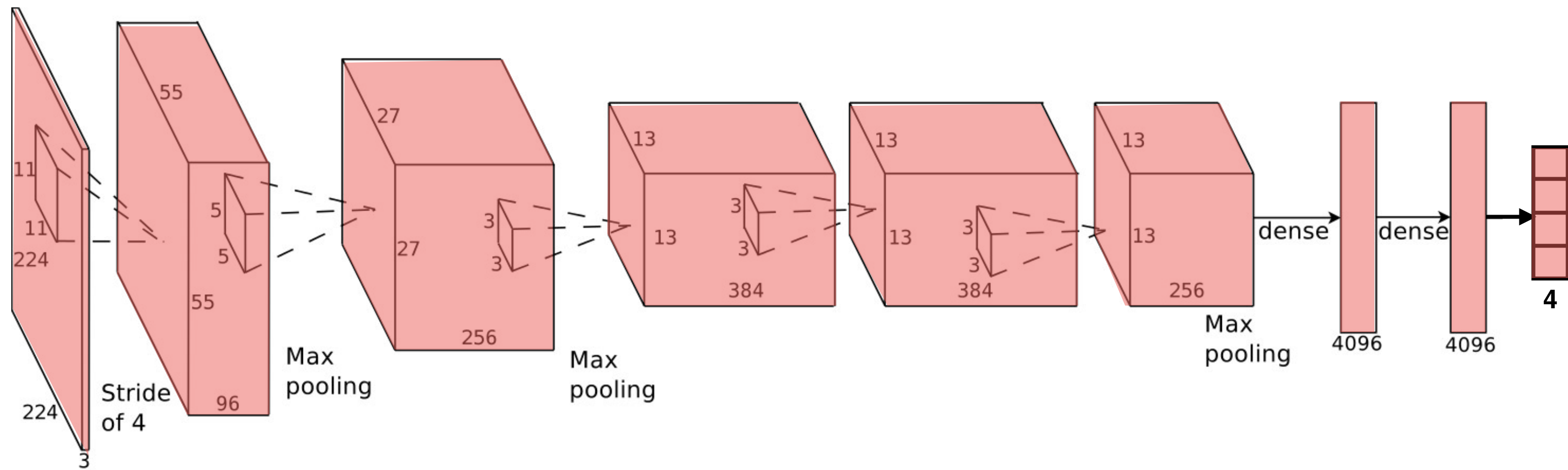
	weights mem:	output size (per image)	(mem)
input: 224 x 224 RGB image	—	224x224x3	150K
conv: (3x3x3) x 64	6.5 KB	224x224x64	12.3 MB
conv: (3x3x64) x 64	144 KB	224x224x64	12.3 MB
maxpool	—	112x112x64	3.1 MB
conv: (3x3x64) x 128	228 KB	112x112x128	6.2 MB
conv: (3x3x128) x 128	576 KB	112x112x128	6.2 MB
maxpool	—	56x56x128	1.5 MB
conv: (3x3x128) x 256	1.1 MB	56x56x256	3.1 MB
conv: (3x3x256) x 256	2.3 MB	56x56x256	3.1 MB
conv: (3x3x256) x 256	2.3 MB	56x56x256	3.1 MB
maxpool	—	28x28x256	766 KB
conv: (3x3x256) x 512	4.5 MB	28x28x512	1.5 MB
conv: (3x3x512) x 512	9 MB	28x28x512	1.5 MB
conv: (3x3x512) x 512	9 MB	28x28x512	1.5 MB
maxpool	—	14x14x512	383 KB
conv: (3x3x512) x 512	9 MB	14x14x512	383 KB
conv: (3x3x512) x 512	9 MB	14x14x512	383 KB
conv: (3x3x512) x 512	9 MB	14x14x512	383 KB
maxpool	—	7x7x512	98 KB
fully-connected 4096	392 MB	4096	16 KB
fully-connected 4096	64 MB	4096	16 KB
fully-connected 1000	15.6 MB	1000	4 KB
soft-max		1000	4 KB

Storing convolution layer outputs (unit “activations”) can get big in early layers with large input size and many filters

Note: multiply these numbers by N for batch size of N images

Many weights in fully-connected players

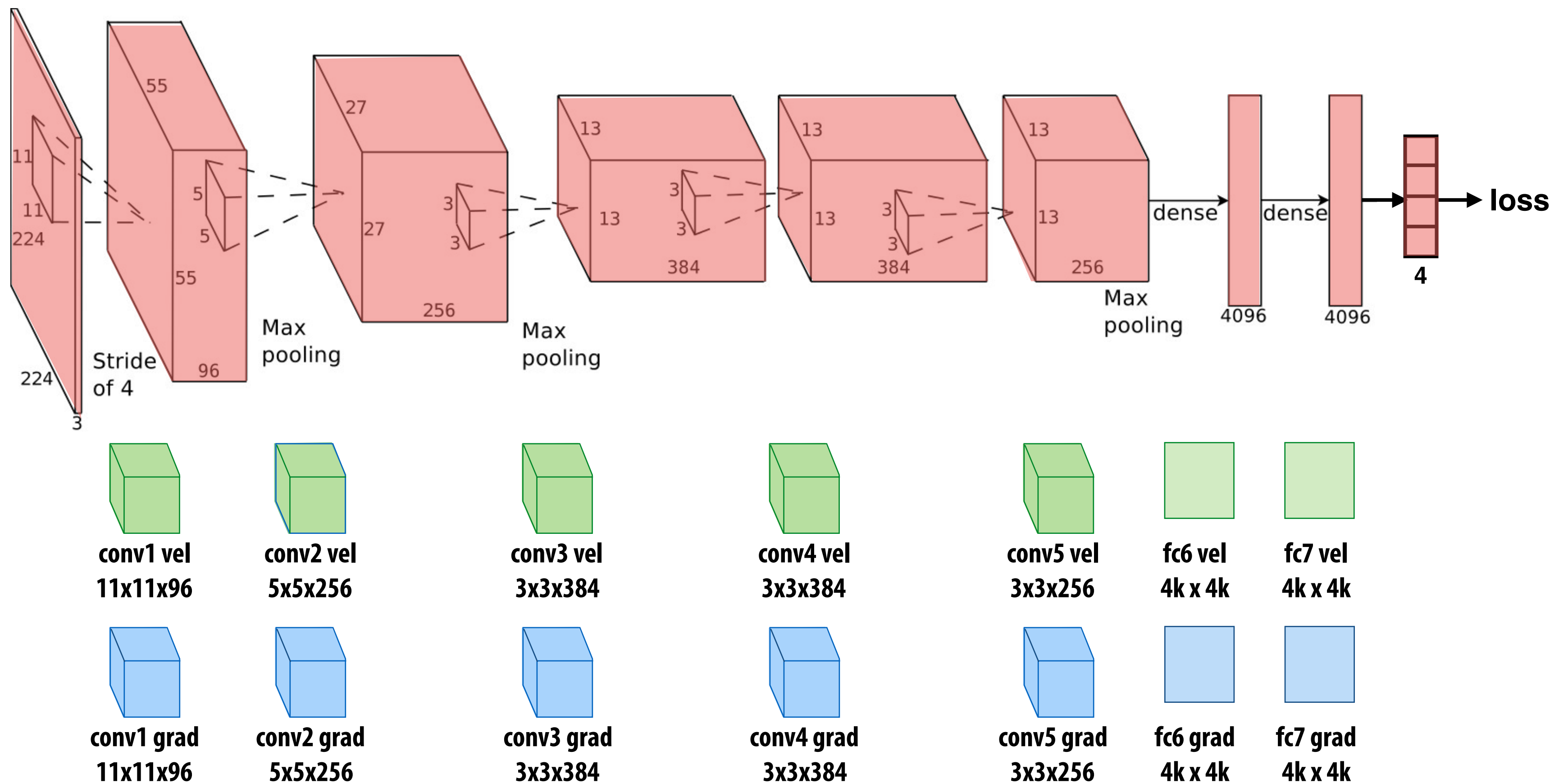
# Data lifetimes during network evaluation



**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

$$vel\_new = mu * vel\_old - step\_size * grad$$

$$w\_new = w\_old + vel\_new$$

# VGG memory footprint

Calculations assume 32-bit values (image batch size = 1)

inputs/outputs get multiplied by mini-batch size

Unlike forward evaluation: cannot immediately free outputs once consumed by next level of network

	weights mem:	output size (per image)	(mem)
input: 224 x 224 RGB image	—	224x224x3	150K
conv: (3x3x3) x 64	6.5 KB	224x224x64	12.3 MB
conv: (3x3x64) x 64	144 KB	224x224x64	12.3 MB
maxpool	—	112x112x64	3.1 MB
conv: (3x3x64) x 128	228 KB	112x112x128	6.2 MB
conv: (3x3x128) x 128	576 KB	112x112x128	6.2 MB
maxpool	—	56x56x128	1.5 MB
conv: (3x3x128) x 256	1.1 MB	56x56x256	3.1 MB
conv: (3x3x256) x 256	2.3 MB	56x56x256	3.1 MB
conv: (3x3x256) x 256	2.3 MB	56x56x256	3.1 MB
maxpool	—	28x28x256	766 KB
conv: (3x3x256) x 512	4.5 MB	28x28x512	1.5 MB
conv: (3x3x512) x 512	9 MB	28x28x512	1.5 MB
conv: (3x3x512) x 512	9 MB	28x28x512	1.5 MB
maxpool	—	14x14x512	383 KB
conv: (3x3x512) x 512	9 MB	14x14x512	383 KB
conv: (3x3x512) x 512	9 MB	14x14x512	383 KB
conv: (3x3x512) x 512	9 MB	14x14x512	383 KB
maxpool	—	7x7x512	98 KB
fully-connected 4096	<b>392 MB</b>	4096	16 KB
fully-connected 4096	<b>64 MB</b>	4096	16 KB
fully-connected 1000	15.6 MB	1000	4 KB
soft-max		1000	4 KB

Must also store per-weight gradients

Many implementations also store gradient "momentum" as well (multiply by 3)

# SGD workload

```
while (loss too high):  
    for each item x_i in mini-batch:  
        grad += evaluate_loss_gradient(f, loss_func, params, x_i)  
    params += -grad * step_size;
```

← **At first glance, this loop is sequential (each step of “walking downhill” depends on previous)**

← **Parallel across images**

↑ **sum reduction**

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

← **trivial data-parallel over parameters**

# DNN training workload

## ■ Huge 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)

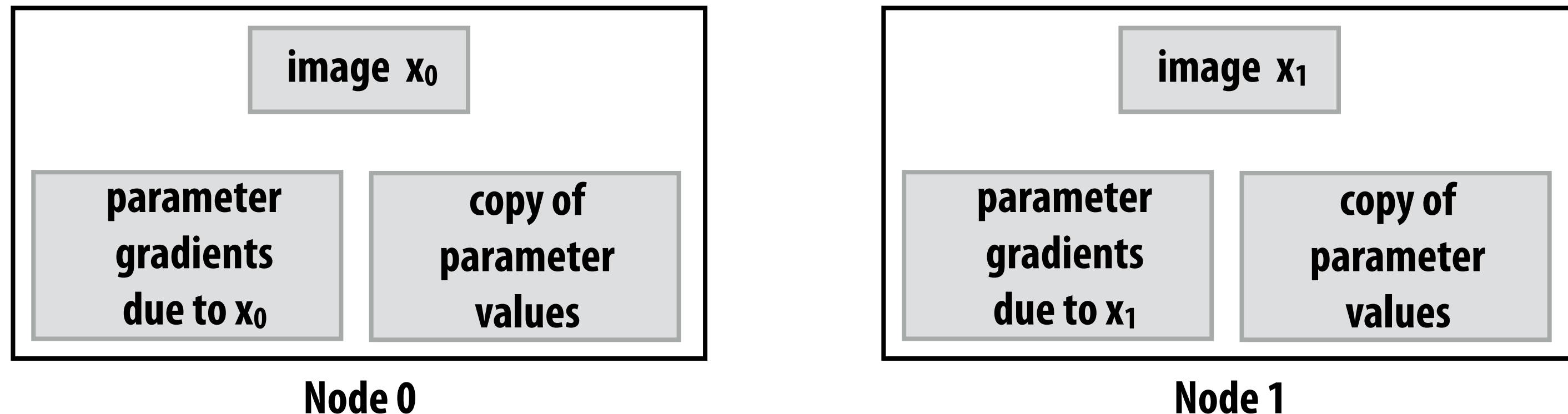
- Each parameter update step depends on previous
- Many units contribute to same parameter gradients (fine-scale reduction)
- Different images in mini batch contribute to same parameter gradients



# 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 * step_size;
```

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



partition mini-batch 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
```

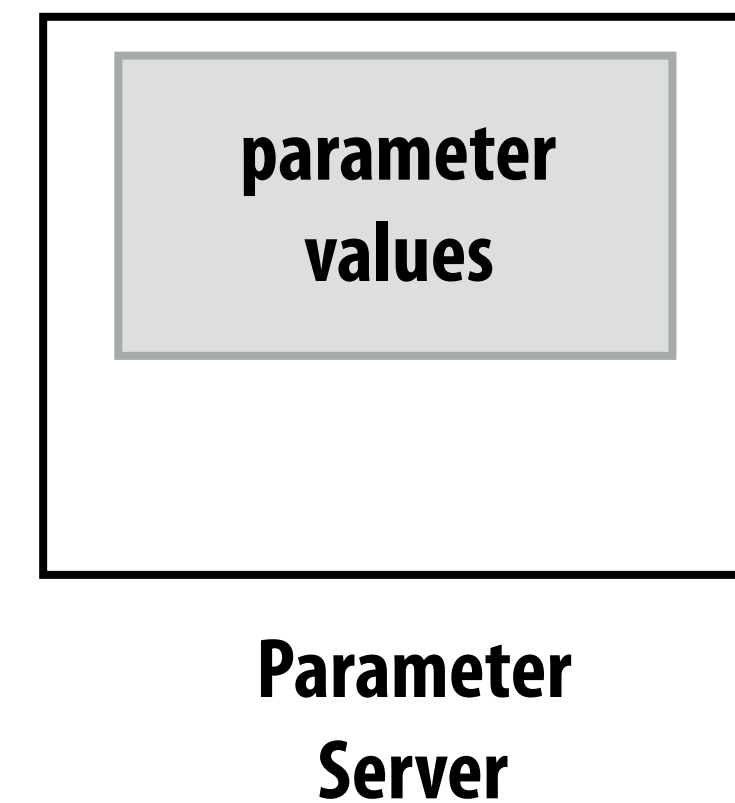
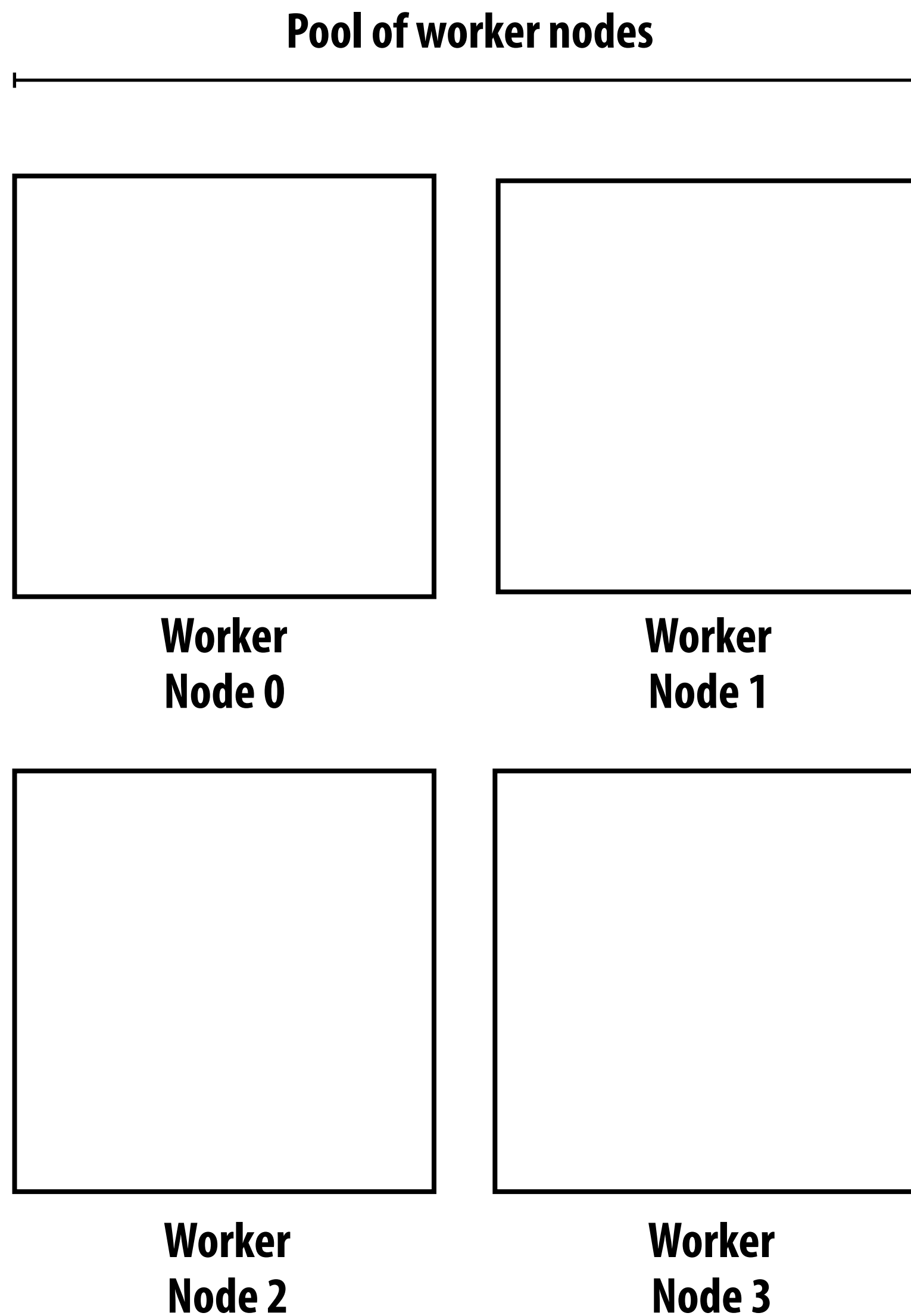
# Challenges of computing at cluster scale

- **Slow communication between nodes**
  - **Commodity clusters do not feature high-performance interconnects (e.g., infiniband) typical of supercomputers**
- **Nodes with different performance (even if machines are the same)**
  - **Workload imbalance at barriers (sync points between nodes)**

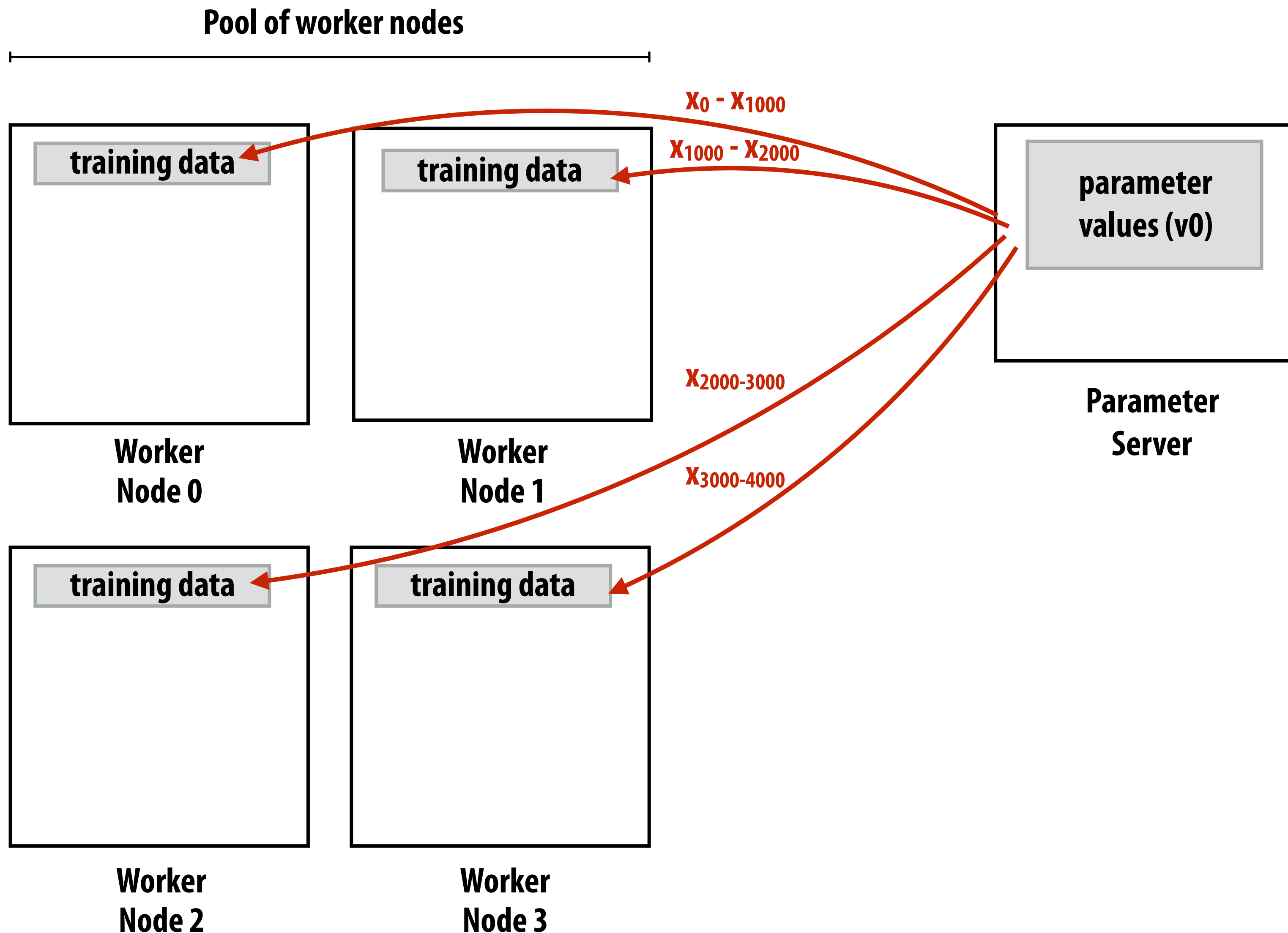
**Modern solution: exploit characteristics of SGD using asynchronous execution!**

# Parameter server design

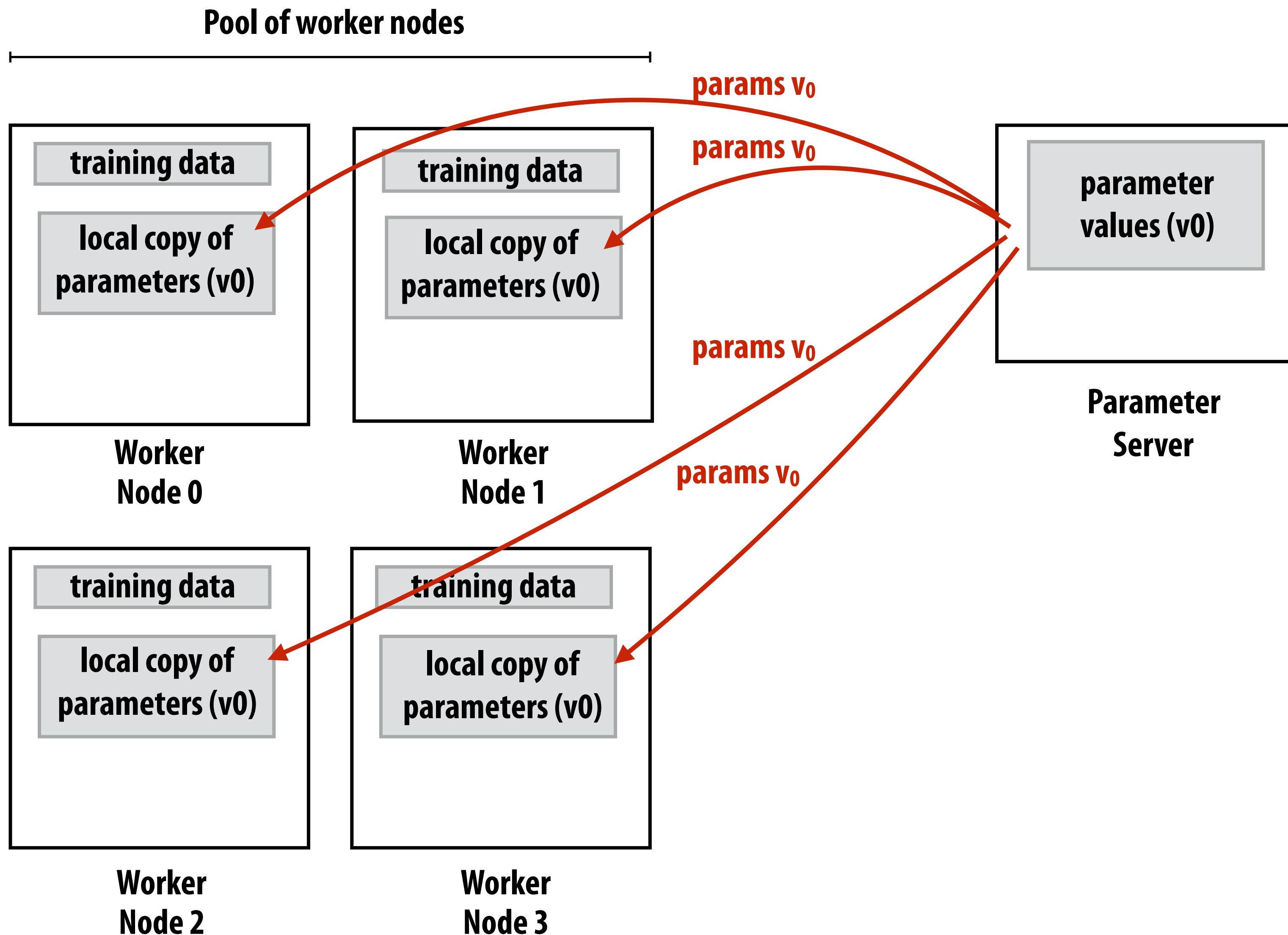
Parameter Server [Li OSDI14]  
Google's DistBelief [Dean NIPS12]  
Microsoft's Project Adam [Chilimbi OSDI14]



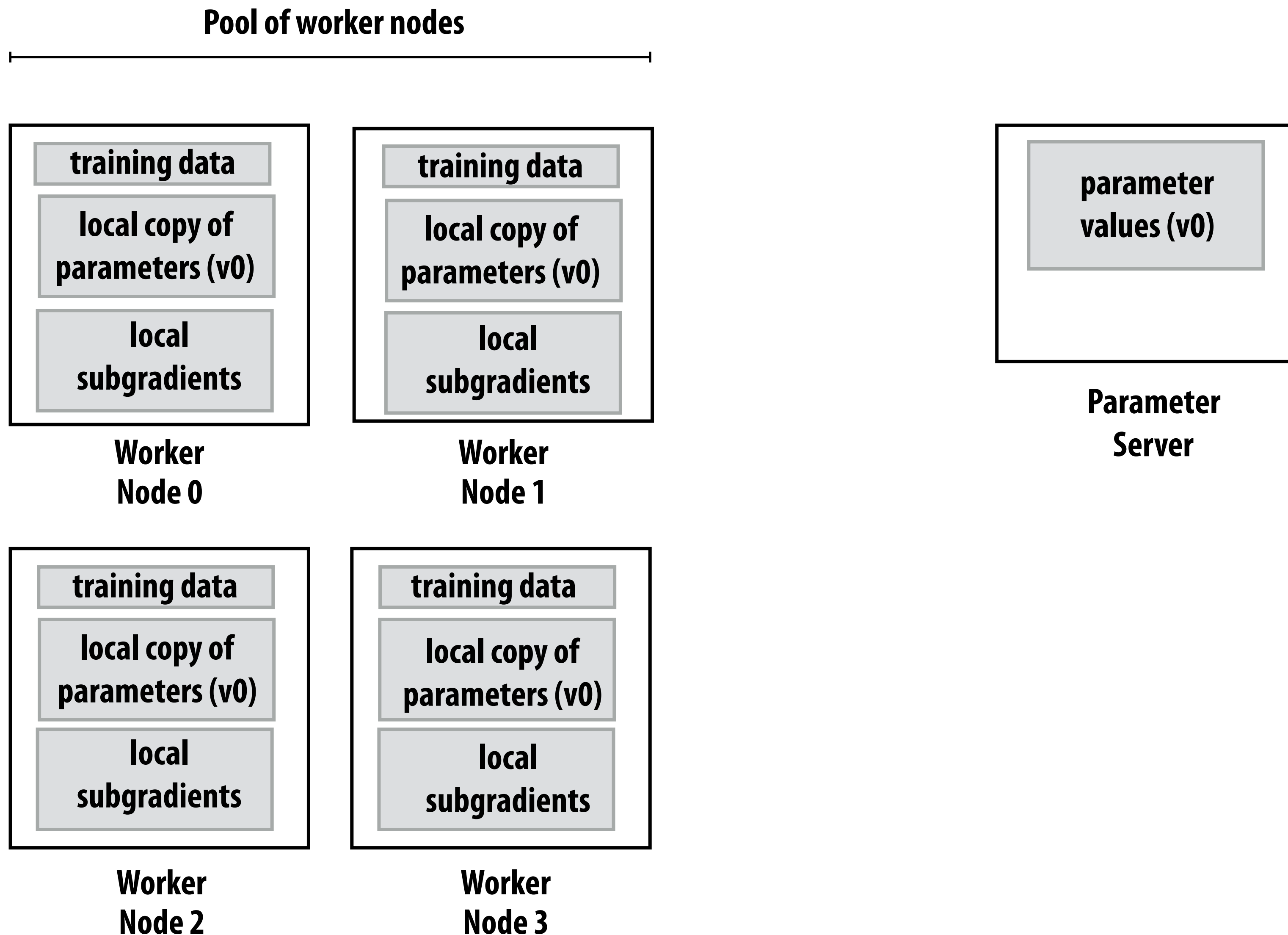
# Training data partitioned among workers



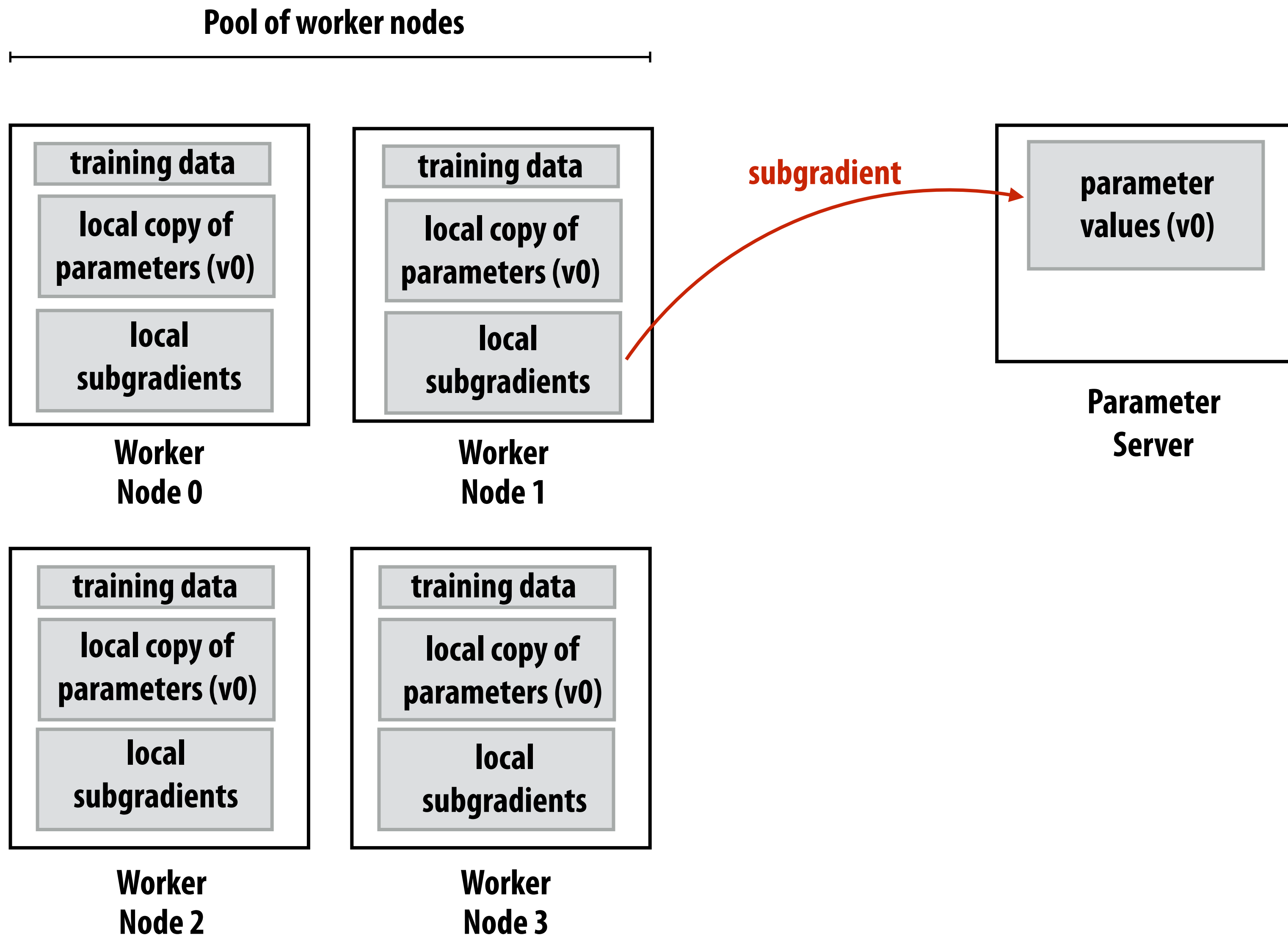
# Copy of parameters sent to workers



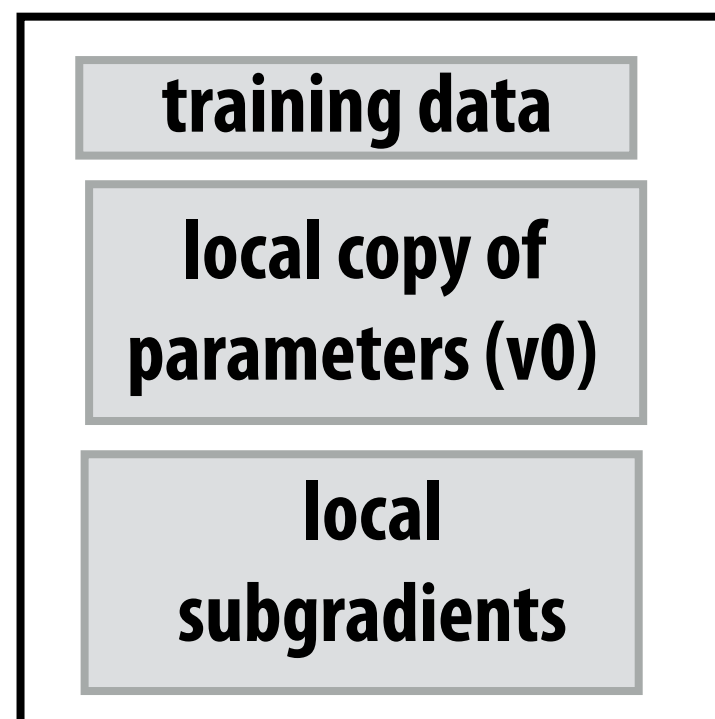
# Workers independently compute local "subgradients"



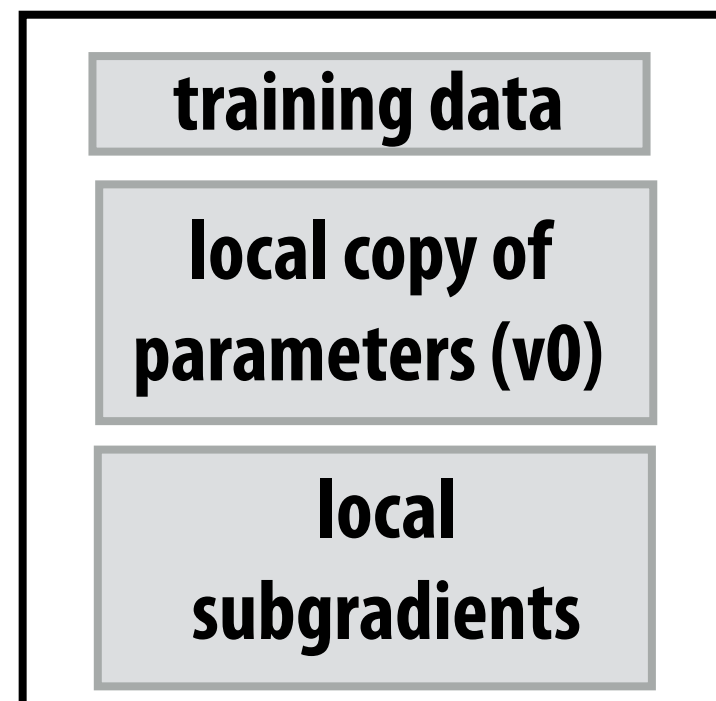
# Worker sends subgradient to parameter server



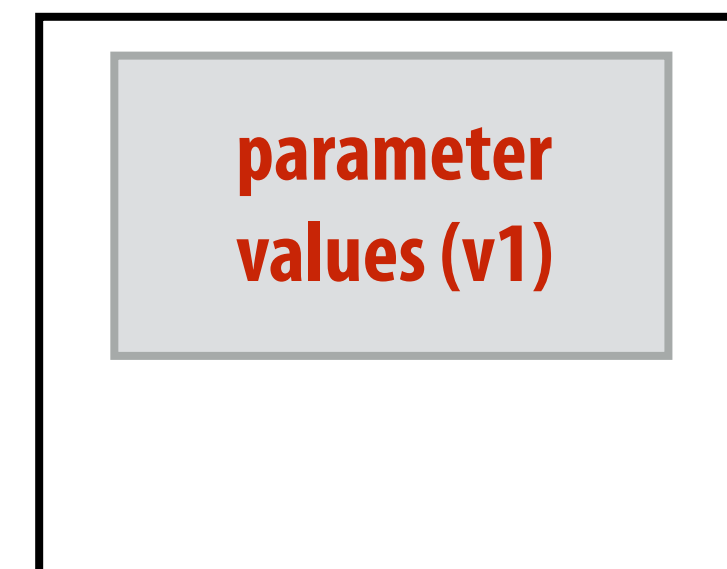
# Server updates global parameter values based on subgradient



Worker Node 0

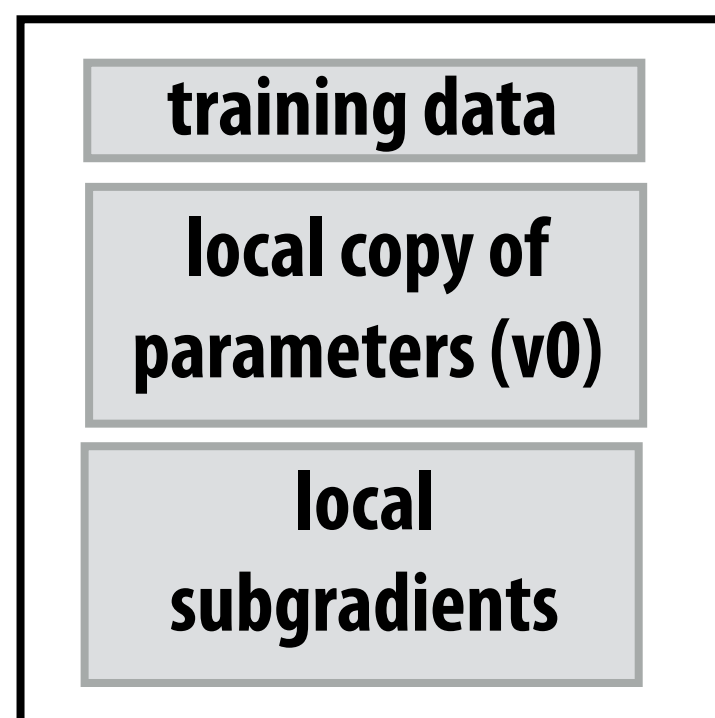


Worker Node 1

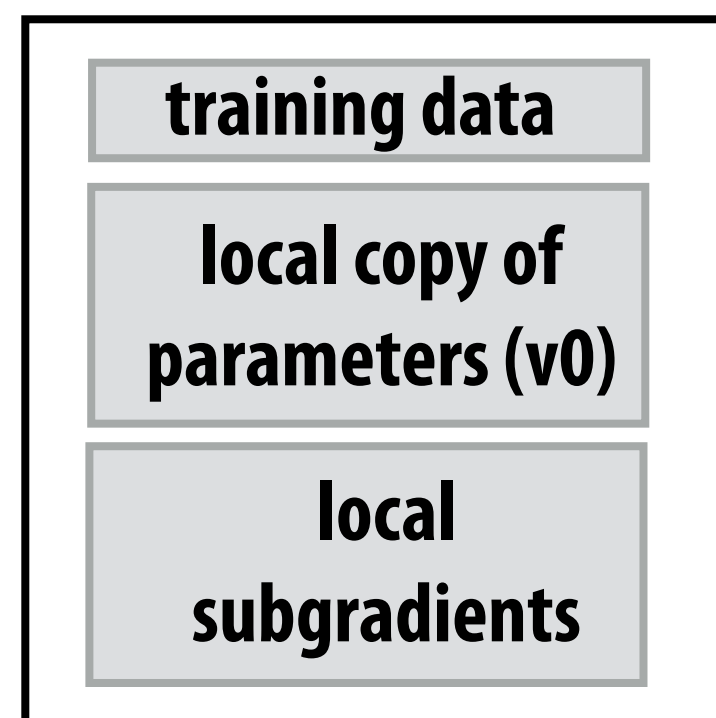


Parameter Server

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



Worker Node 2

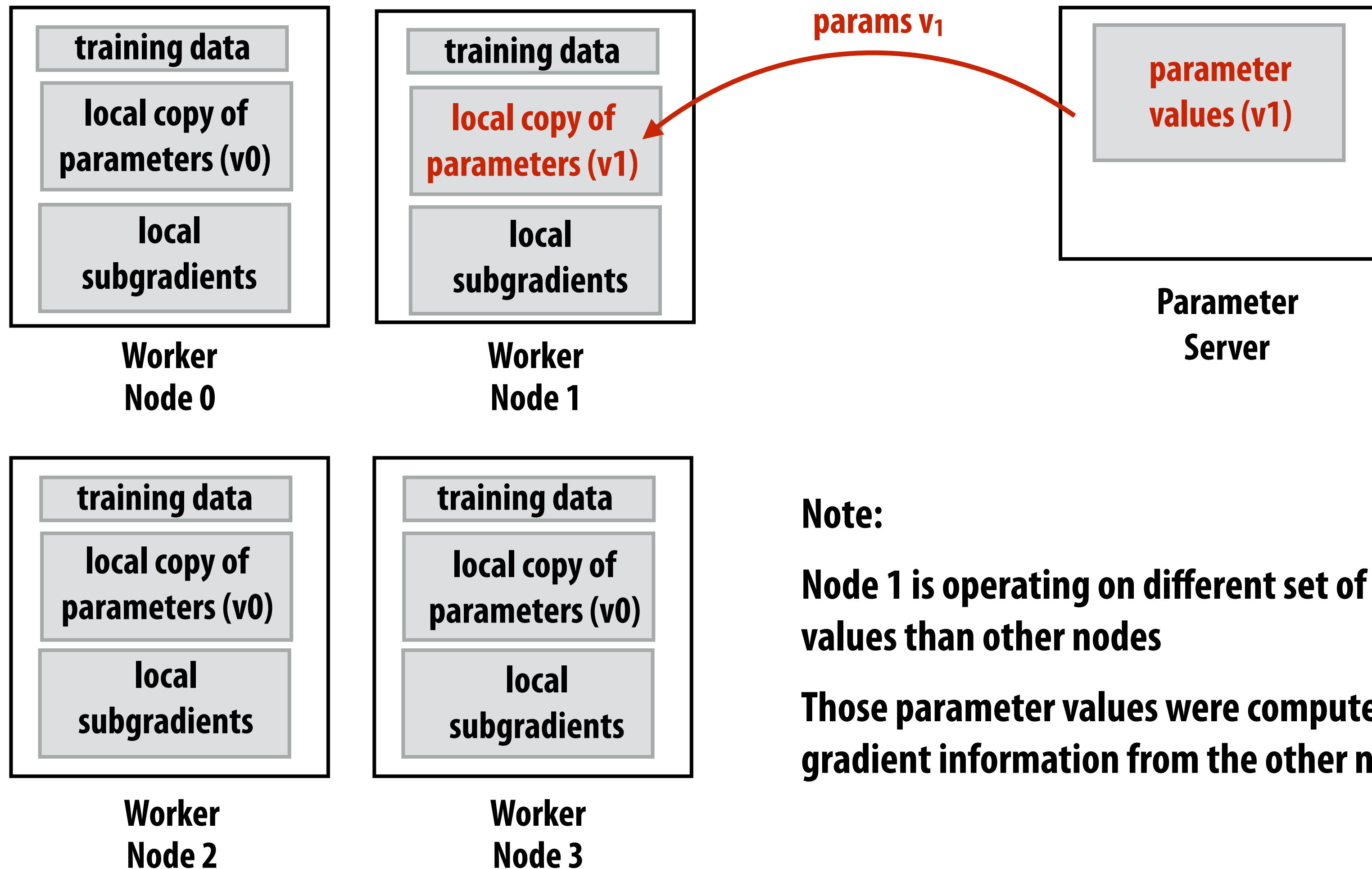


Worker Node 3



# Updated parameters sent to worker

Worker proceeds with another gradient computation step

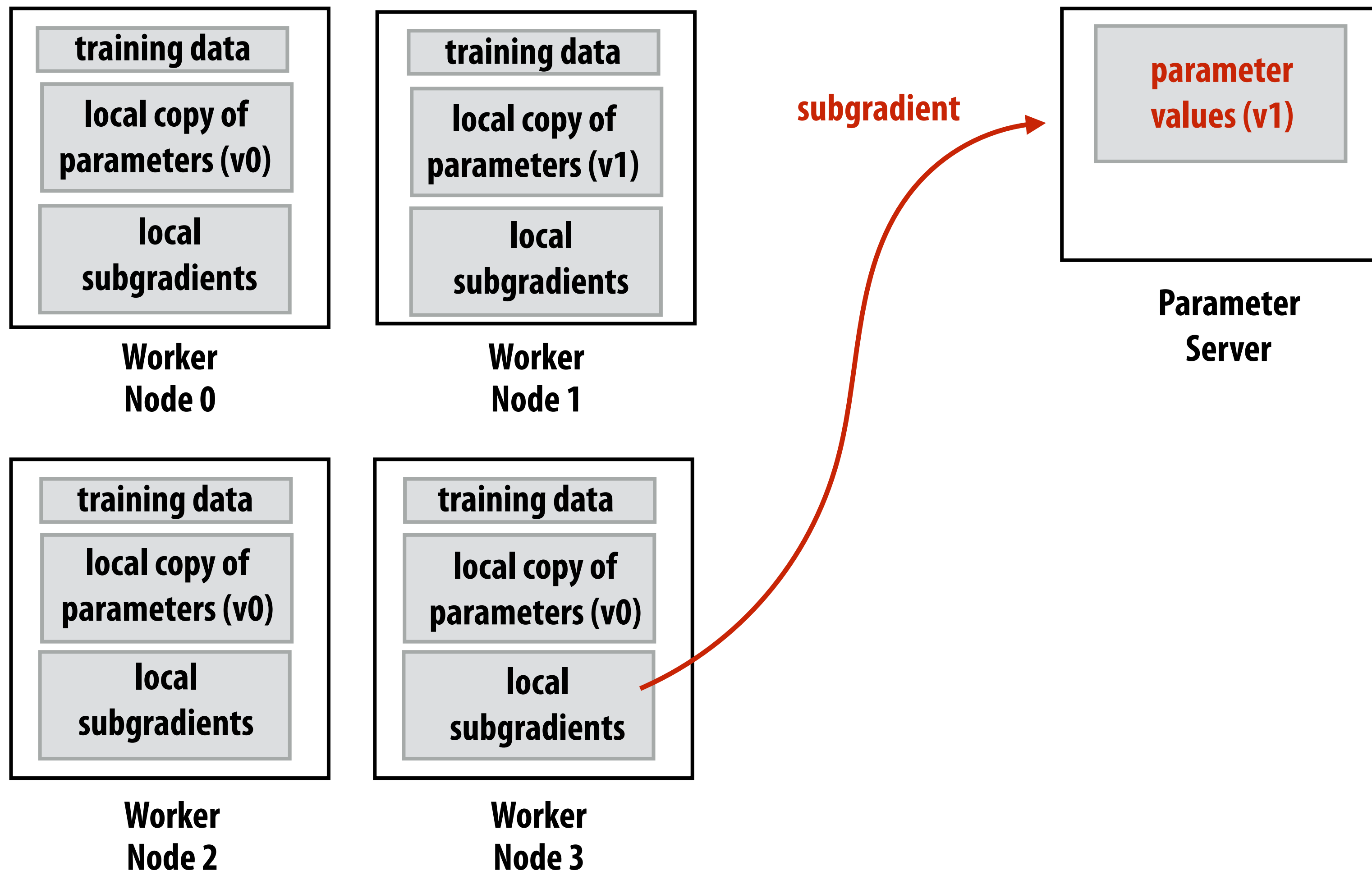


**Note:**

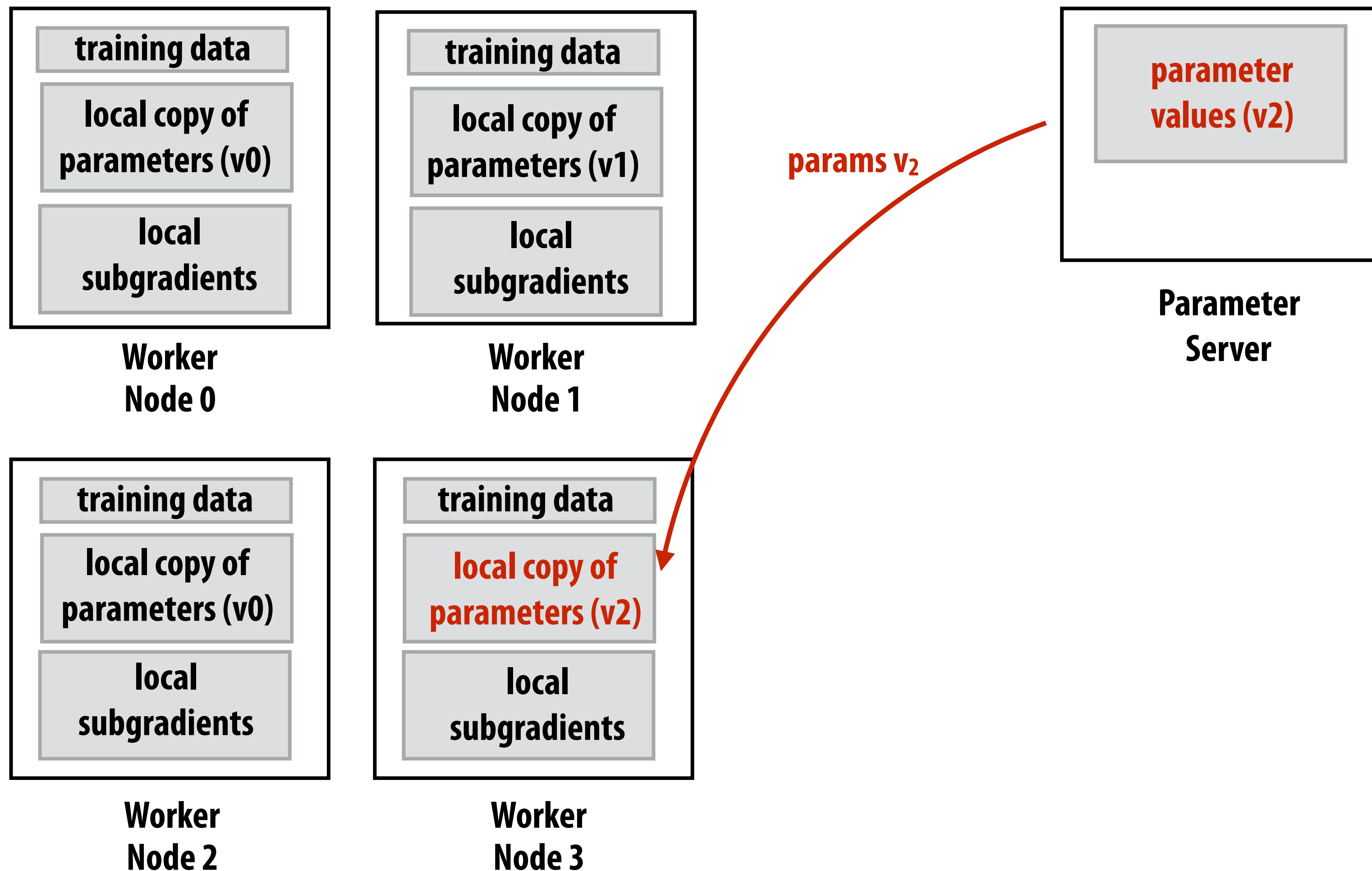
**Node 1 is operating on different set of parameter values than other nodes**

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

# Updated parameters sent to worker (again)



# Worker continues with updated parameters

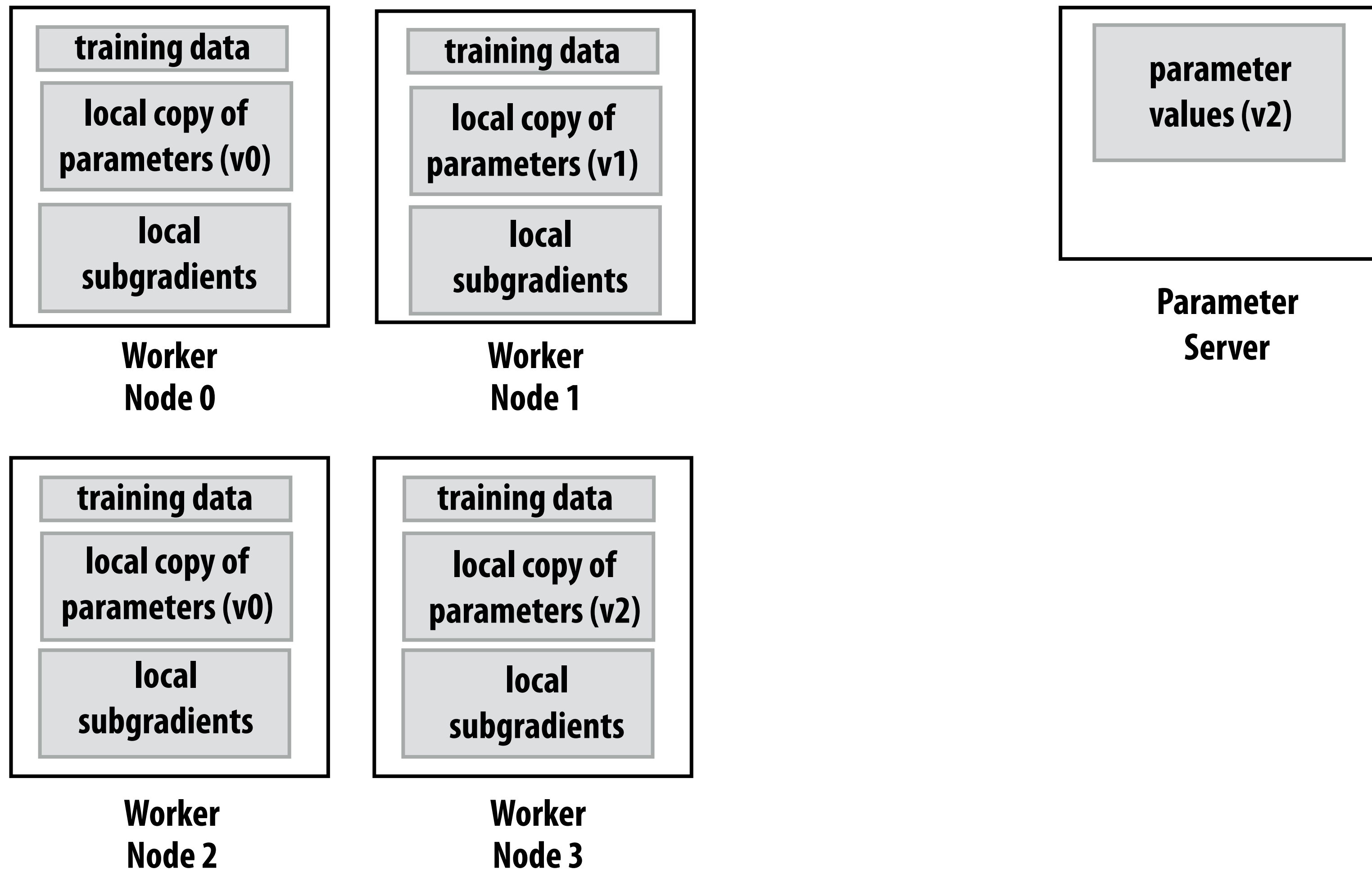


# Summary: asynchronous parameter update

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

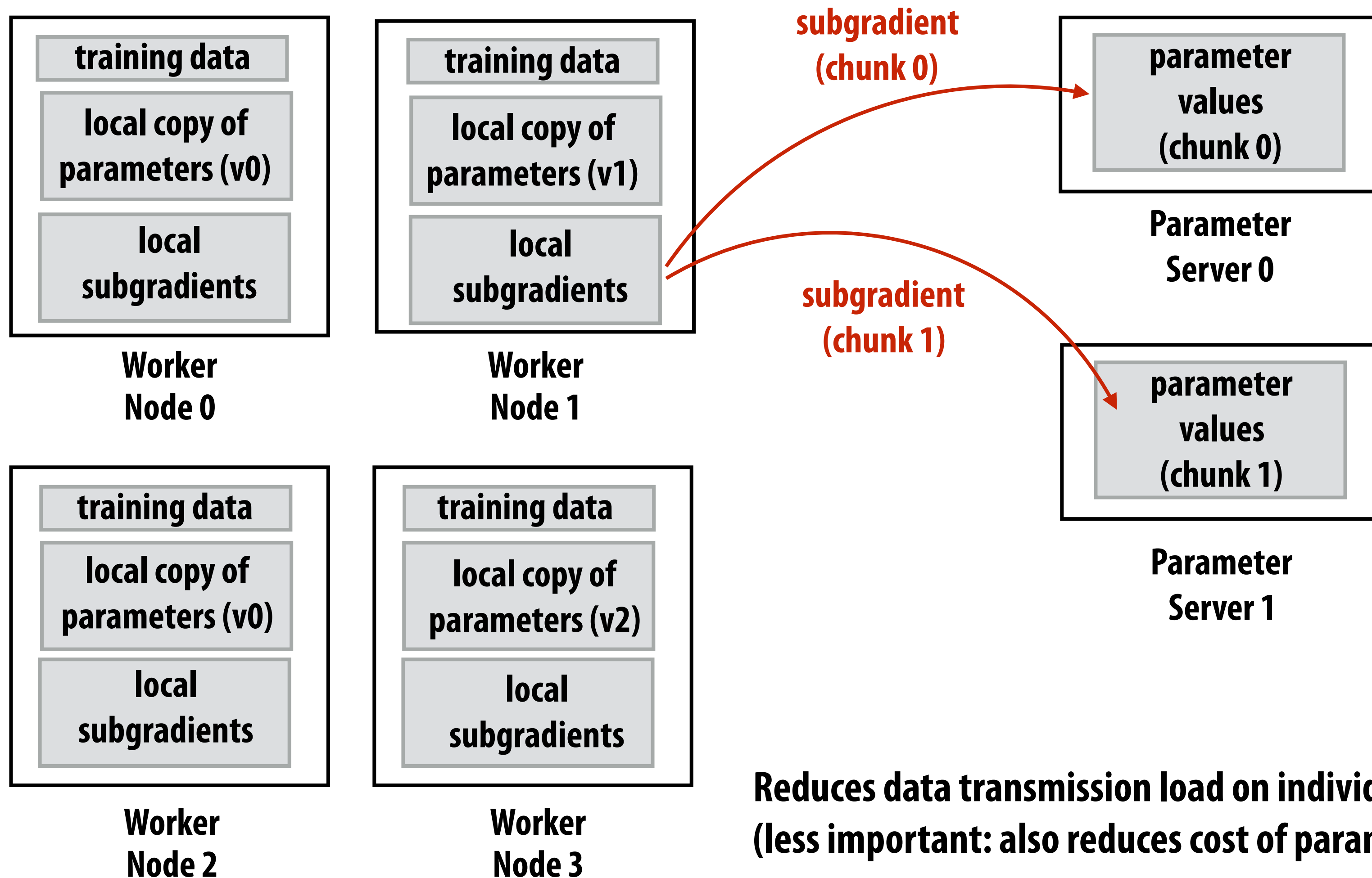
What if there is heavy contention for parameter server?



# Shard the parameter server

Partition parameters across servers

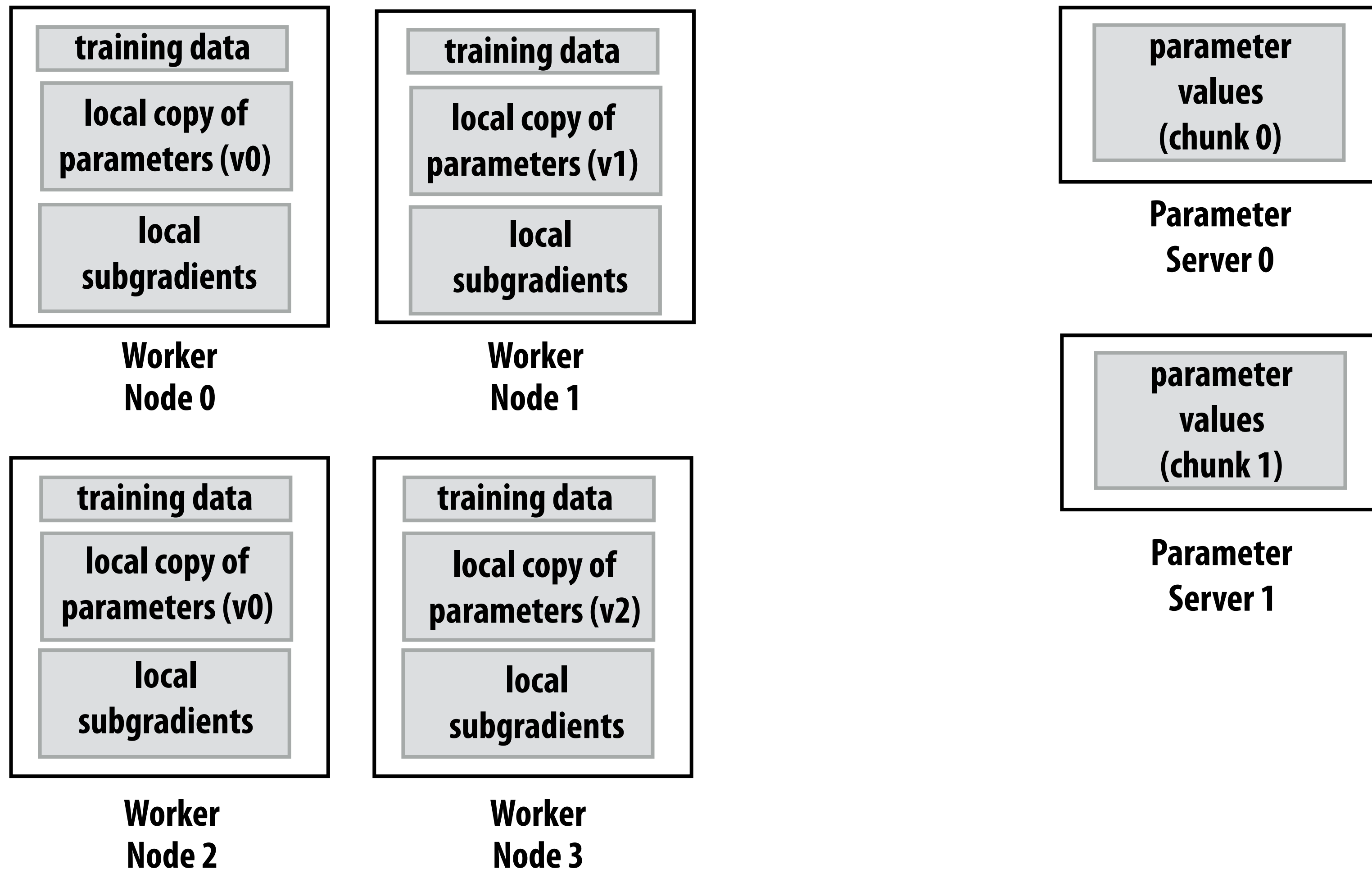
Worker sends chunk of subgradients to owning parameter server



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)

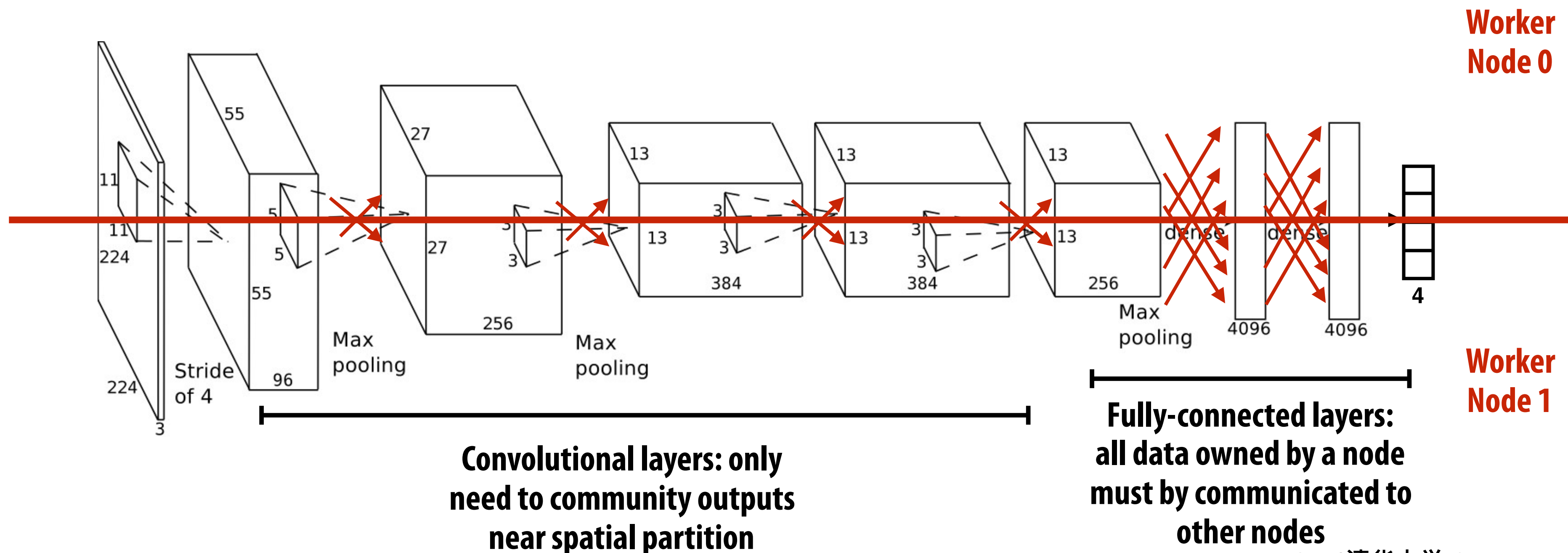


# Model parallelism

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

Reduce internode communication through network design:

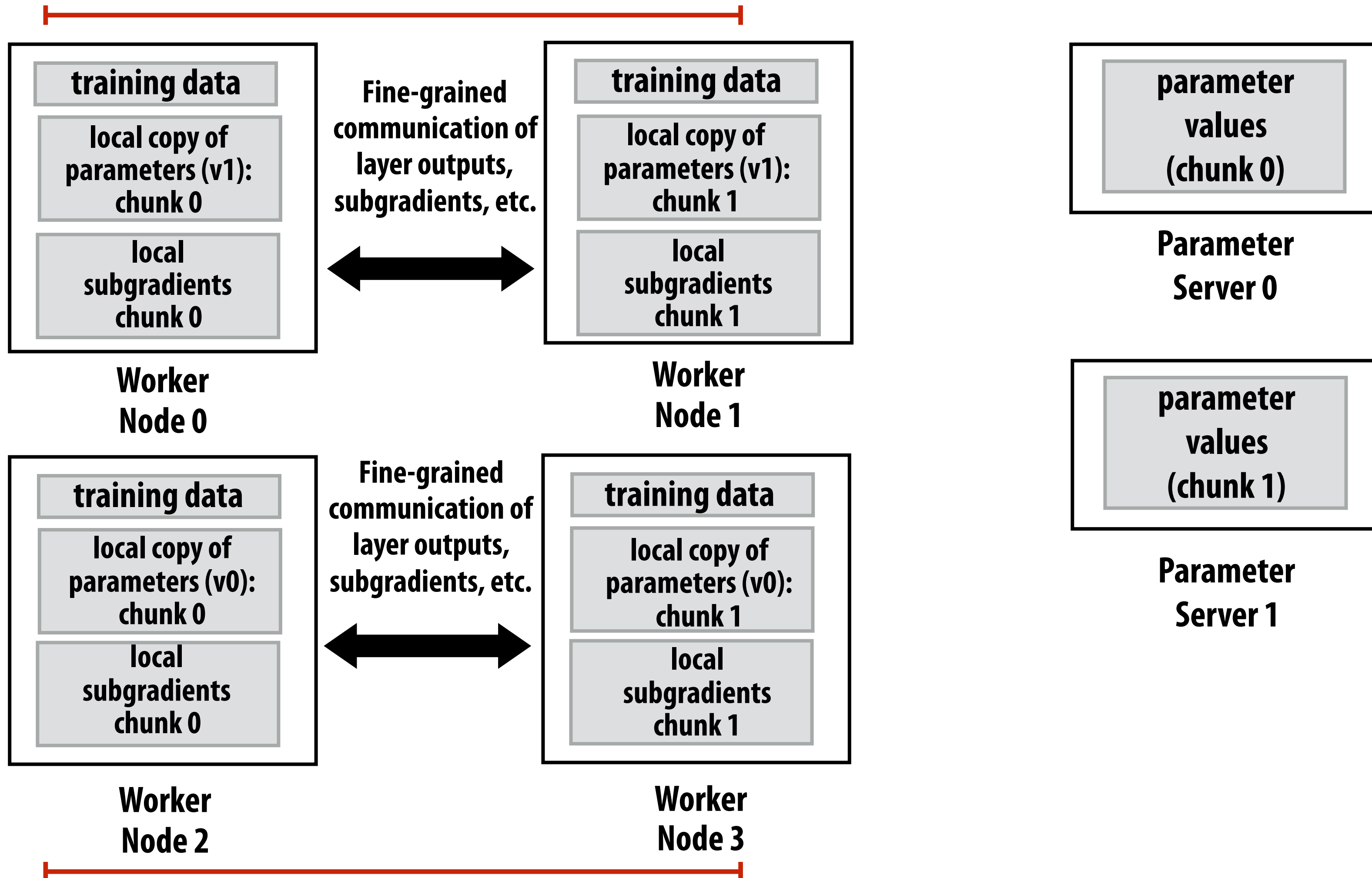
- Use small spatial convolutions (1x1 convolutions)
- Reduce/shrink fully-connected layers





# Training data-parallel and model-parallel execution

Working on subgradient computation  
for a single copy of the model



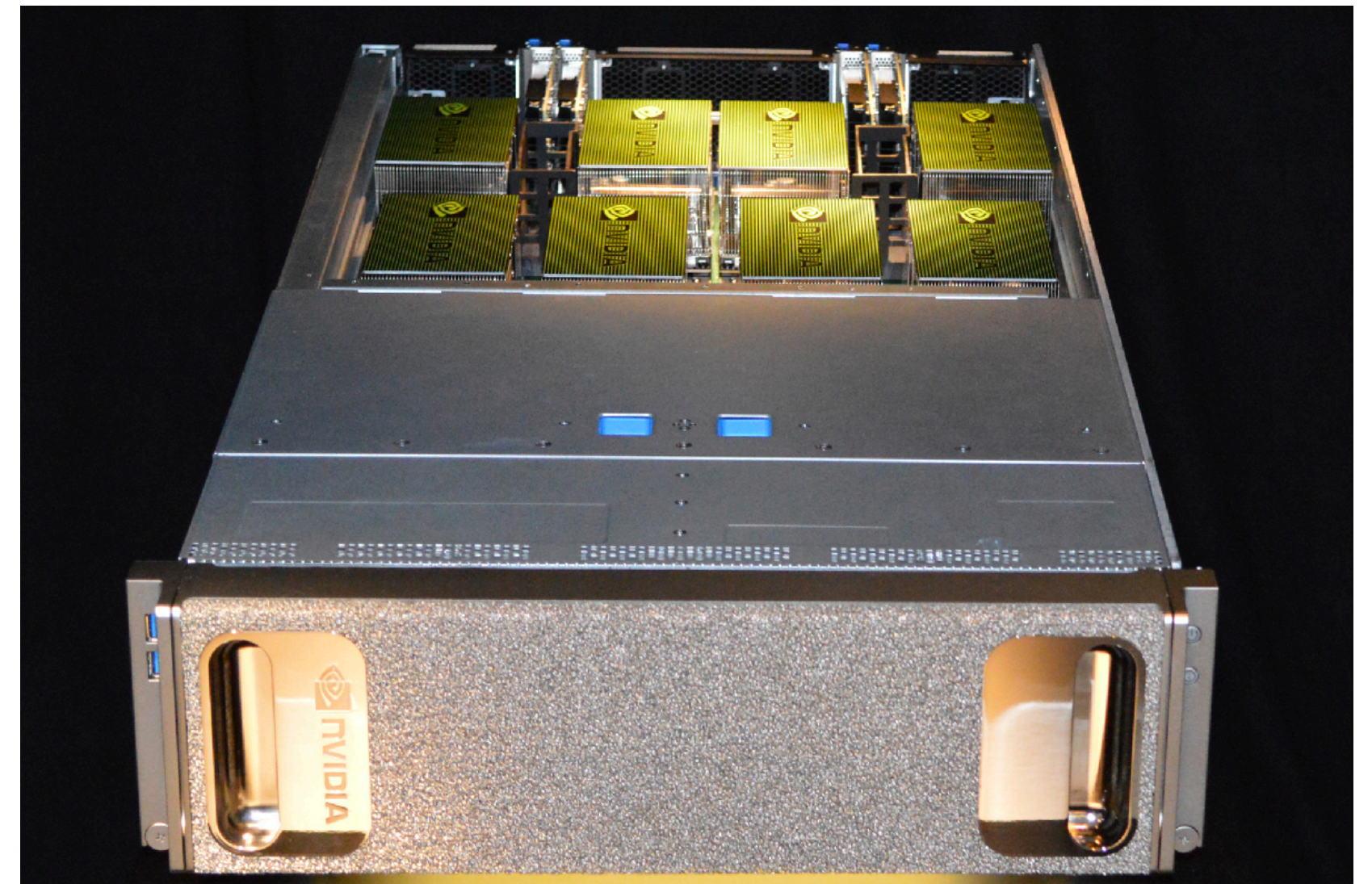
Working on subgradient computation  
for a single copy of the model

# Using supercomputers for training?

- **Fast interconnects critical for model-parallel training**
  - **Fine-grained communication of outputs and gradients**
- **Fast interconnect diminishes need for async training algorithms**
  - **Avoid randomness in training due to computation schedule (there remains randomness due to SGD algorithm)**



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



**NVIDIA DGX-1: 8 Pascal GPUs connected  
via high speed NV-Link interconnect**

# Accelerating data-parallel training

FireCaffe [Iandola 16]

- Use a high-performance Cray Gemini interconnect (Titan supercomputer)
- Use combining tree for accumulating gradients (rather than a single parameter server)
- Use larger batch size (to reduce frequency of communication) and offset by increasing learning rate

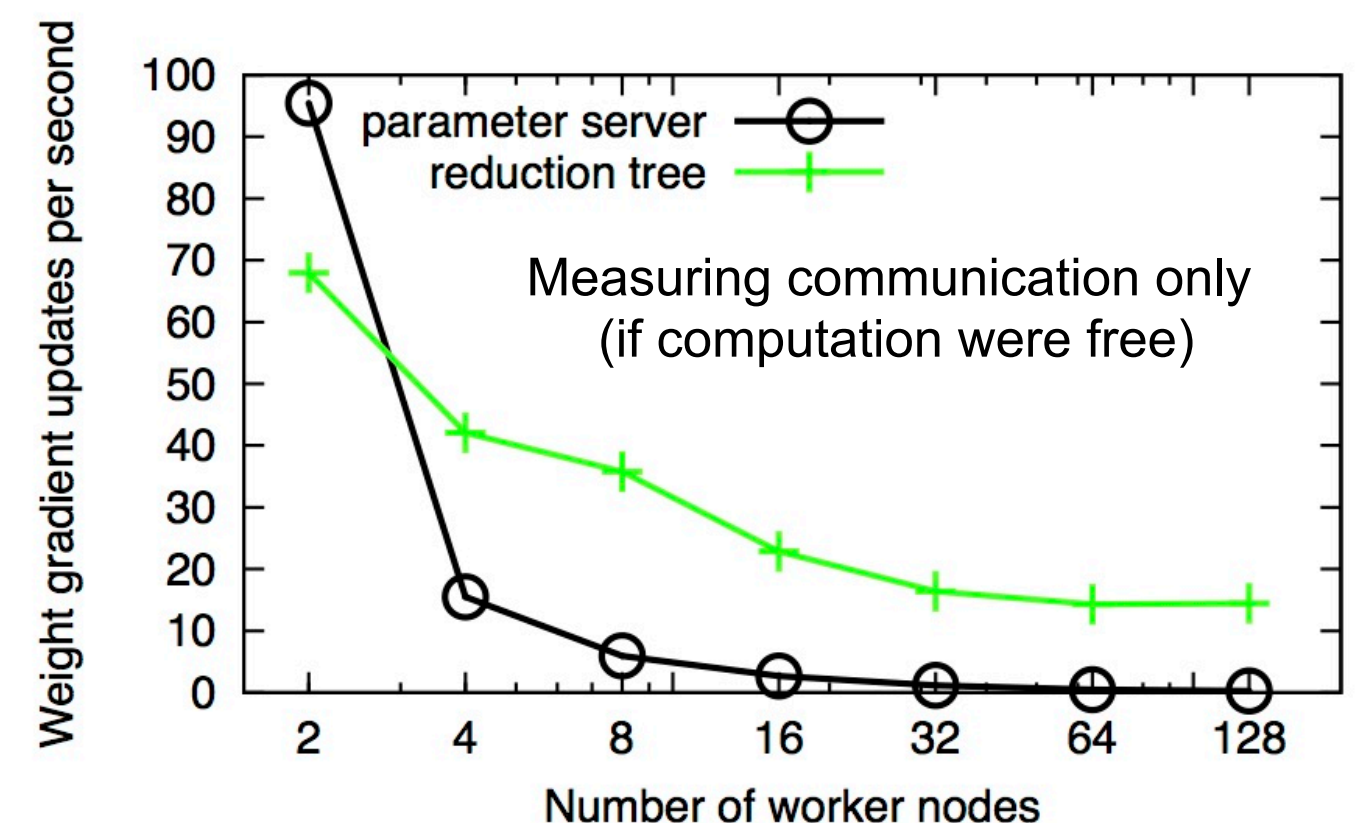
	Hardware	Net	Epochs	Batch size	Initial Learning Rate	Train time	Speedup	Top-1 Accuracy	Top-5 Accuracy
Caffe	1 NVIDIA K20	GoogLeNet [41]	64	32	0.01	21 days	1x	68.3%	88.7%
FireCaffe (ours)	32 NVIDIA K20s (Titan supercomputer)	GoogLeNet	72	1024	0.08	23.4 hours	20x	68.3%	88.7%
FireCaffe (ours)	128 NVIDIA K20s (Titan supercomputer)	GoogLeNet	72	1024	0.08	10.5 hours	47x	68.3%	88.7%

Dataset: ImageNet 1K

**Result: reasonable scalability without asynchronous parameter update: for modern DNNs with fewer weights such as GoogLeNet (due to no fully connected layers)**

Also see: Goyal et al. 2017

“Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour”



# swDNN: DNN library for TaihuLight

<https://github.com/THUHPGC/swDNN>



**Sunway TaihuLight: 40,960 x 256 core processors = 10,485,760 CPU cores**

# Summary: training large networks in parallel

- **Many cluster systems rely on asynchronous update to efficiently use clusters of commodity machines**
  - **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)**
- **Although modern DNN designs (with fewer weights), large mini batch sizes, and efficient use of high performance interconnects (much like any parallel computing problem) enables scalability without asynchronous execution**
- **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! recall the original grid solver example!)**