**Lecture 9:**

# Parallel Deep Network Training

**Visual Computing Systems**
**Stanford CS348V, Winter 2018**

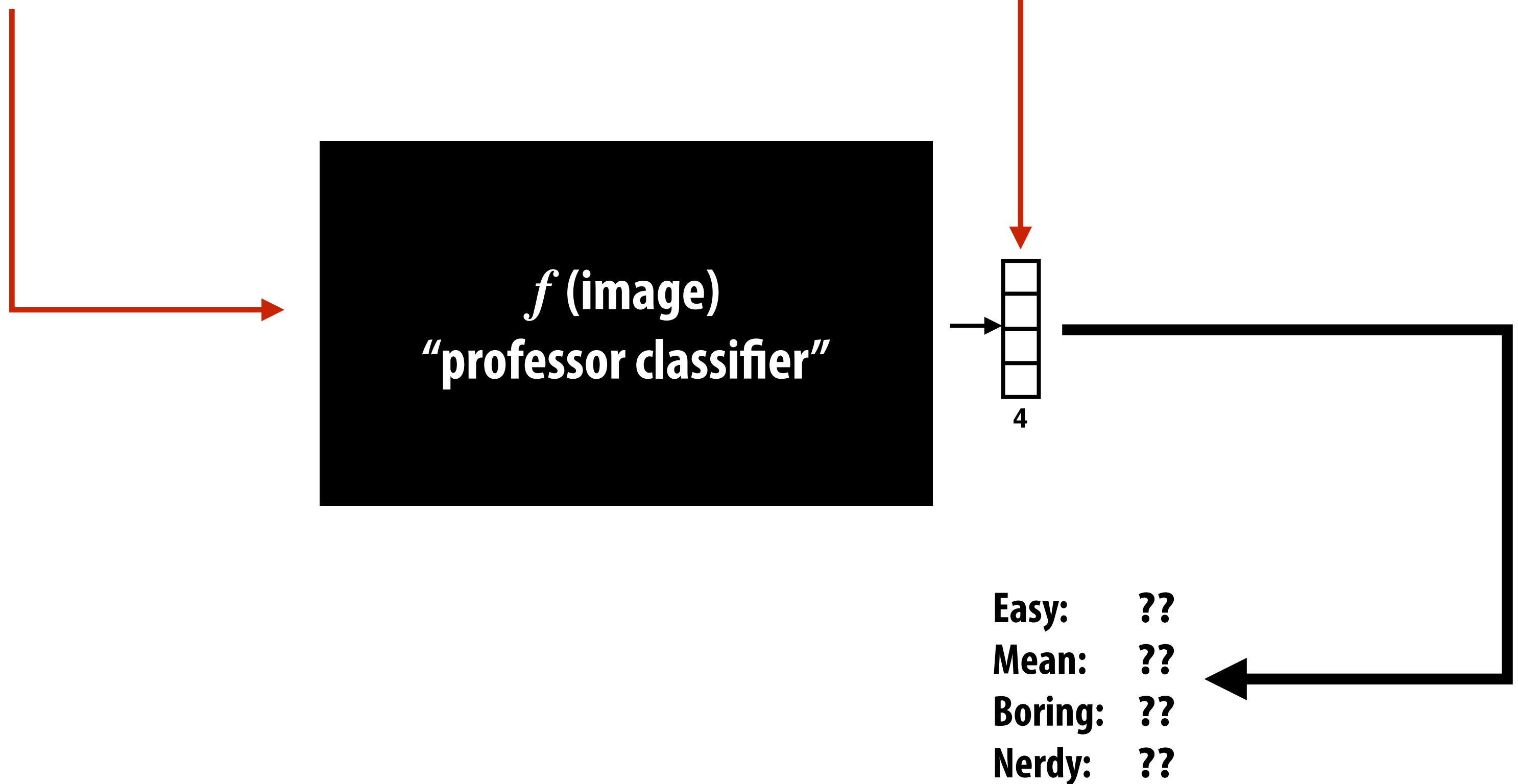# How would you describe this professor?



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**

$f$ (image)
"professor classifier"

4

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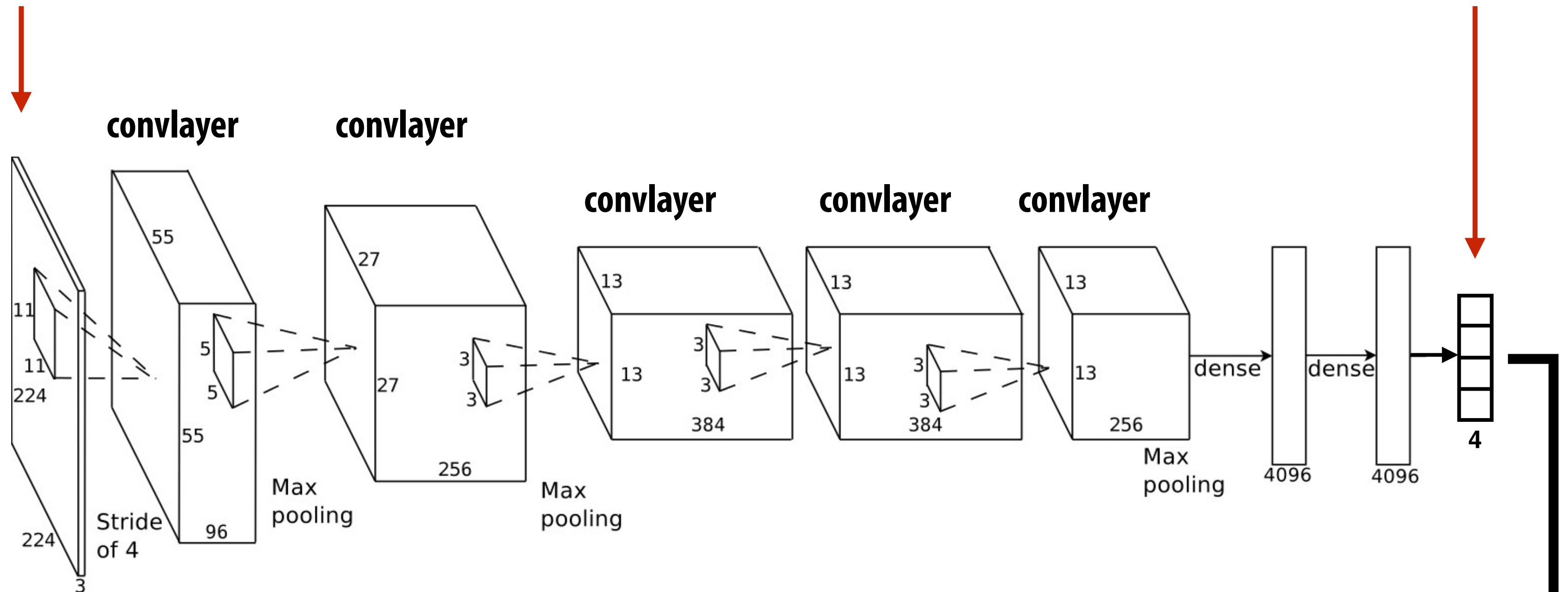Boring:  **??**
Nerdy:  **??**

# Professor classification network

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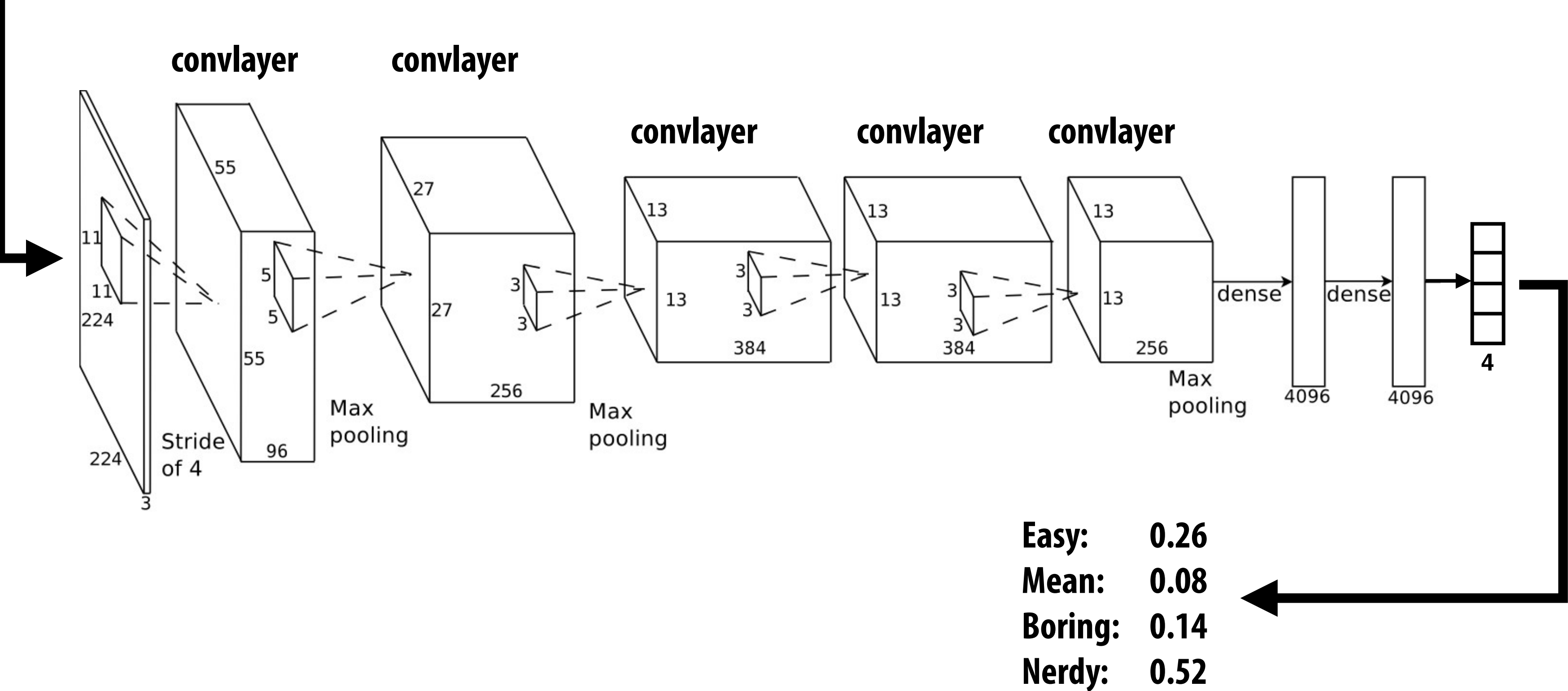


Input:
image of a professor

Output:
probability of label

convlayer  convlayer  convlayer  convlayer  convlayer

dense  dense

4

55  27  13  13  13

11  5  3  3  3

11  5  3  3  3

224  27  13  13  13

224  55  256  384  384  256

3  96  Max pooling  Max pooling  Max pooling  4096  4096

Stride of 4

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

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

# Professor classification network



convlayer     convlayer     convlayer     convlayer     convlayer

Easy:     0.26
Mean:     0.08
Boring:     0.14
Nerdy:     0.52
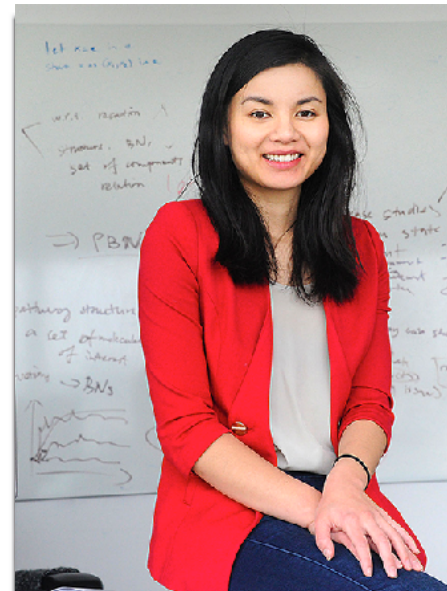
# Training data (ground truth answers)



[label omitted]    [label omitted]    [label omitted]    **Nerdy**    [label omitted]    [label omitted]    [label omitted]
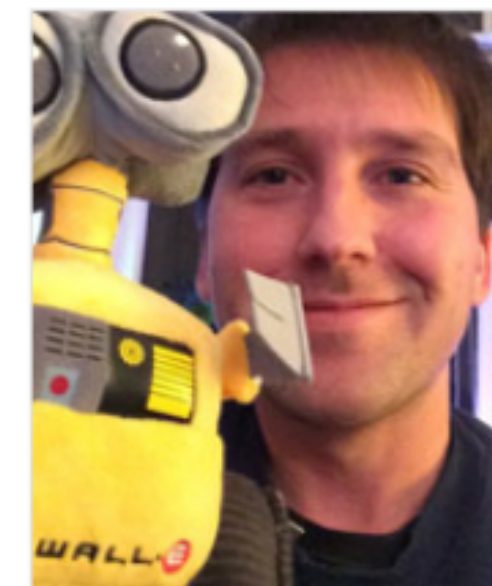
[label omitted]    [label omitted]    **Nerdy**    [label omitted]    [label omitted]    **Nerdy**    [label omitted]

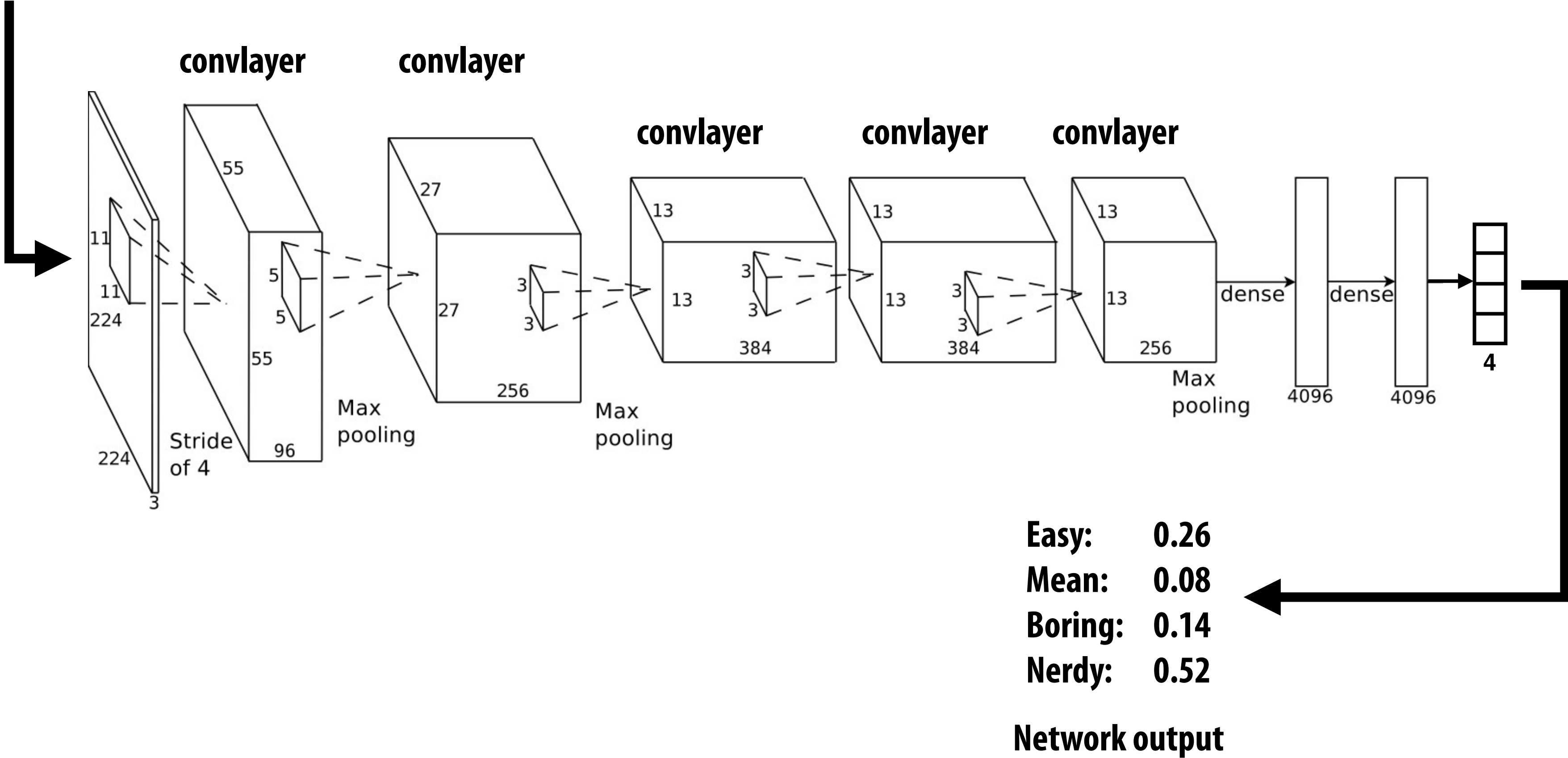[label omitted]    **Nerdy**    **Nerdy**    [label omitted]    [label omitted]    [label omitted]    **Nerdy**

# Professor classification network



New image of Kayvon
(not in training set)

convlayer  convlayer  convlayer  convlayer  convlayer

55  27  13  13  13
11  5   3   3   3
11  5   3   3   3
224 27  13  13  13
256 384 384 256

dense  dense  4

Stride of 4  Max pooling  Max pooling  Max pooling  4096  4096

224  96  3

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

Network output

Stanford CS348V, Winter 2018

# 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** |

**Output of network for correct category**

**Common example: softmax loss:**

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

**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)**

$$L = \sum_i L_i$$   **(total loss for entire training set is sum of losses $L_i$ 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 *p₁* and *p₂*:   $f(x_i)$
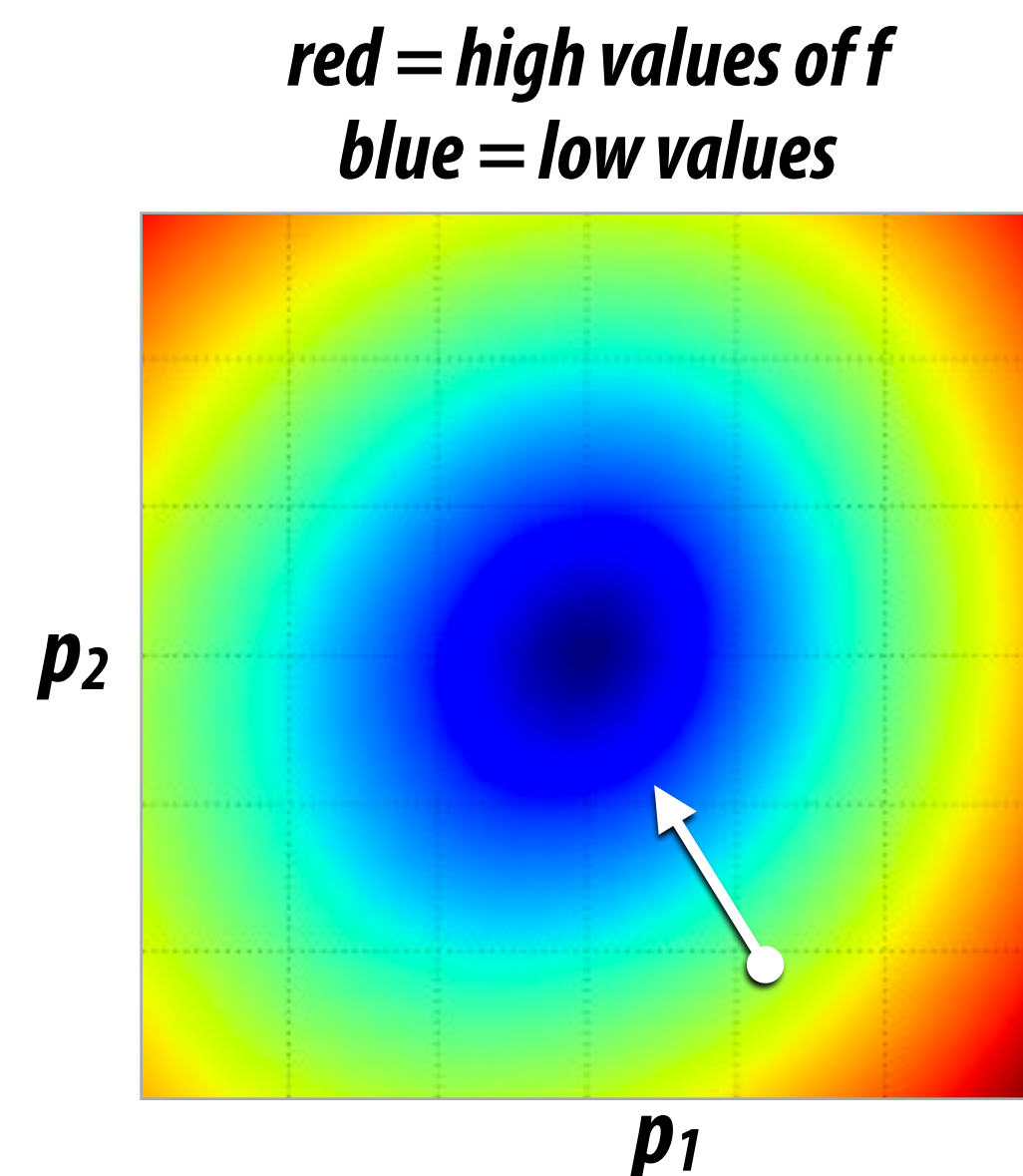
And for some input *xᵢ*, your training data says the function should output 0.

But for the current values of *p₁* and *p₂*, 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 *p₁* and *p₂* so you could compute their value at *xᵢ*.

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

*red = high values of f*
*blue = low values*



*p₂*

*p₁*

How might you adjust the values *p₁* and *p₂* to reduce the error for this training example?

# 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
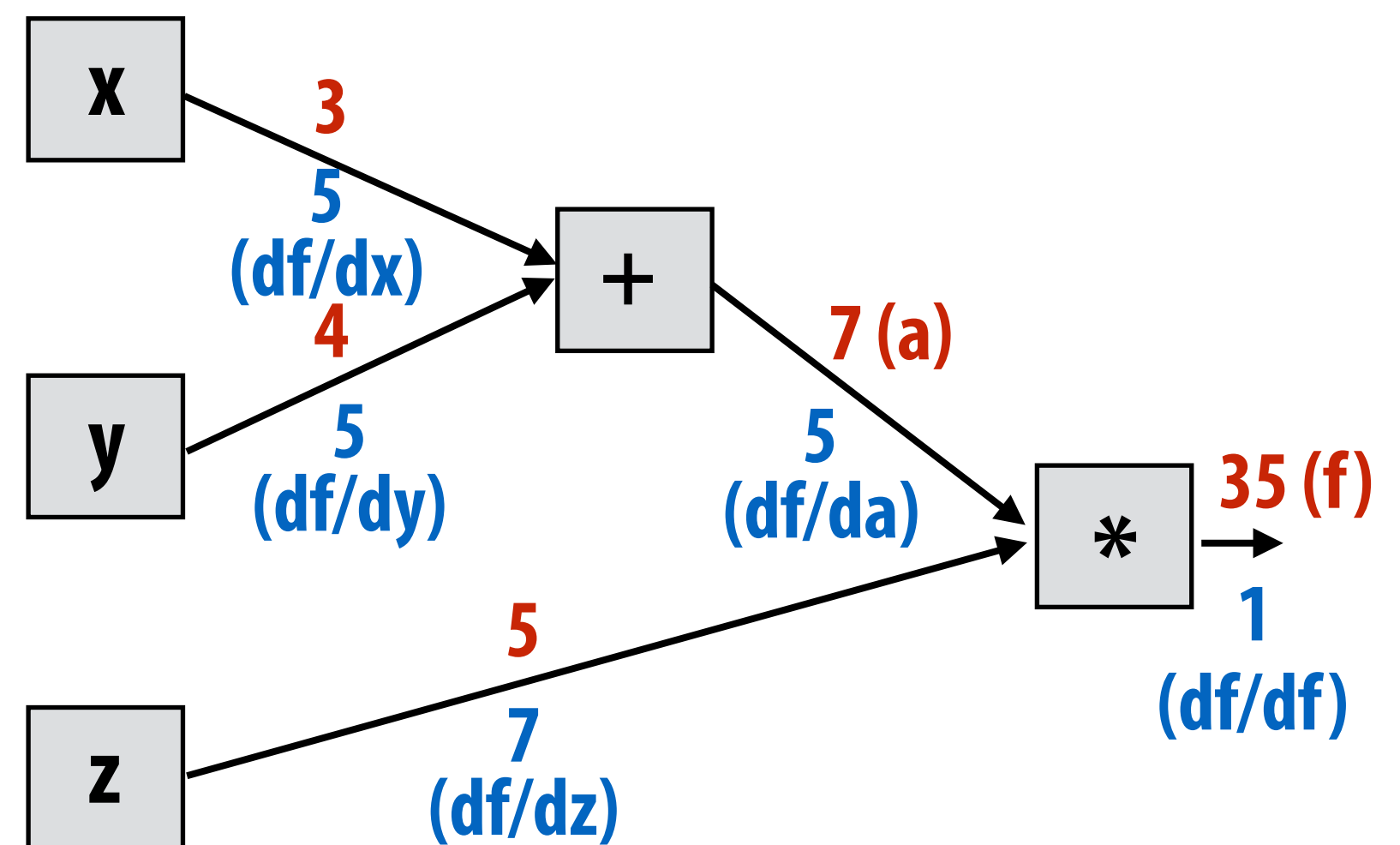
# Derivatives using the chain rule

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

**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

y

**10**

**10**

**+**

**10**

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

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

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

**max**

**10**

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

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

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

**\***

**10**

$g(x, y) = xy$

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

# Back-propagating through single unit

$x_0$

$w_0$

$yw_0$

$yx_0$

*

$x_1$

$w_1$

$yw_1$

$yx_1$

*

$y$

+

$y$

$x_2$

$w_2$

$yw_2$

$yx_2$

*

$y$

$x_3$

$w_3$

$yw_3$

$yx_3$

*

$y$

+

$y$

+

$y$

$y$

+

$y$

**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**

max

**10** $\dfrac{d\text{loss}}{d\text{unit}}$

**0**

b

$y$

**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:**

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

$$= 10\frac{dg}{dx}$$

$$= 10(2x+1)$$

$$= 10(10+1) = 110$$

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

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

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

**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



X

w

$*$

$y = Xw$

$\dfrac{dL}{dy}$

$\dfrac{dL}{dw}$

**(WxH)-element vector**

**9-element vector**

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

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

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

**Therefore:**

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

**9**

| 0 | 0 | 0 | 0 | x00 | x01 | 0 | x10 | x11 |
|---|---|---|---|-----|-----|---|-----|-----|
| 0 | 0 | 0 | x00 | x01 | x02 | x10 | x11 | x12 |
| 0 | 0 | 0 | x01 | x02 | x03 | x11 | x12 | x13 |

...

| x00 | x01 | x02 | x10 | x11 | x12 | x20 | x21 | x22 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|

...

$\dfrac{dy}{dw_2}$

**WxH**

X

$\begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_8 \end{bmatrix}$

w

# 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 minibatch)**

**Update all parameter values:** `w_new = w_old − learning_rate * grad`

# Recall from last class: VGG memory footprint

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

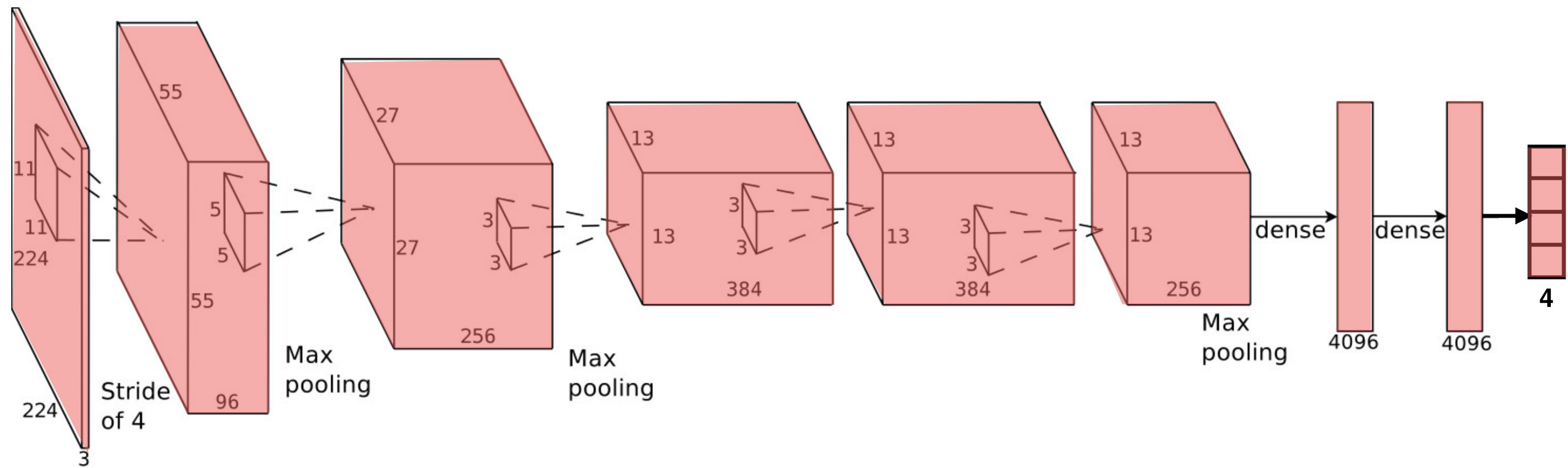| | weights mem: | output size (per image) | (mem) | |
|---|---|---|---|---|
| input: 224 x 224 RGB image | — | 224x224x3 | 150K | **Storing convolution layer outputs (unit "activations") can get big in early layers with large input size and many filters** |
| conv: (3x3x3) x 64 | 6.5 KB | 224x224x64 | 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 | **Note: multiply these numbers by N for batch size of N images** |
| 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 | **Many weights in fully-connected players** |
| 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 | |

# 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



conv1 vel
11x11x96

conv2 vel
5x5x256

conv3 vel
3x3x384

conv4 vel
3x3x384

conv5 vel
3x3x256

fc6 vel
4k x 4k

fc7 vel
4k x 4k

conv1 grad
11x11x96

conv2 grad
5x5x256

conv3 grad
3x3x384

conv4 grad
3x3x384

conv5 grad
3x3x256

fc6 grad
4k x 4k

fc7 grad
4k x 4k

- 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 | 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 |

Must also store per-weight gradients

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

# 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**



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 (e.g., MPI_Allreduce)**
- **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

training data

**Worker Node 0**

**Worker Node 1**

$x_0 - x_{1000}$

$x_{1000} - x_{2000}$

$x_{2000-3000}$

$x_{3000-4000}$

parameter values (v0)

**Parameter Server**

training data

training data

**Worker Node 2**

**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

# Workers independently compute local "subgradients"

**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

**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 (v0)

# 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

**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 (v0)

# Server updates global parameter values based on subgradient

| | |
|---|---|
| **Worker Node 0** | |

**training data**

**local copy of parameters (v0)**

**local subgradients**

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

**parameter values (v1)**

Parameter
Server

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

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



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

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

**params v1**

**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

**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 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



training data

local copy of parameters (v0)

local subgradients

**Worker Node 0**

training data

local copy of parameters (v1)

local subgradients

**Worker Node 1**

training data

local copy of parameters (v0)

local subgradients

**Worker Node 2**

training data

local copy of parameters (v2)

local subgradients

**Worker Node 3**

parameter values (v2)

**Parameter Server**

params $v_2$

# Summary: asynchronous parameter update

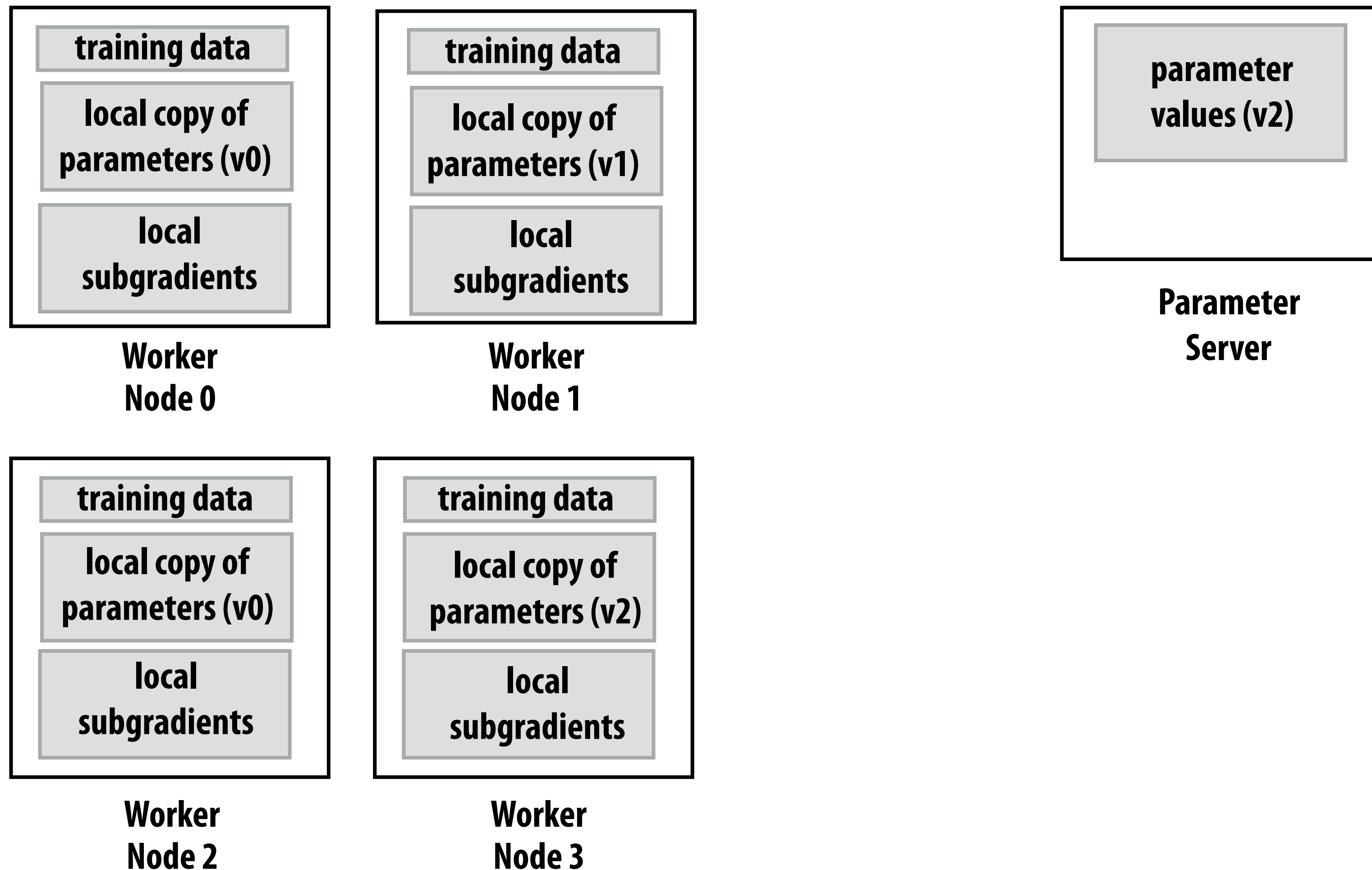- **Idea: avoid global synchronization on all parameter updates between each SGD iteration**
  - Design reflects realities of cluster computing:
    - Slow interconnects
    - Unpredictable machine performance

- **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?

**training data**

**local copy of parameters (v0)**

**local subgradients**

Worker Node 0

**training data**

**local copy of parameters (v1)**

**local subgradients**

Worker Node 1

**parameter values (v2)**

Parameter Server

**training data**

**local copy of parameters (v0)**

**local subgradients**

Worker Node 2

**training data**

**local copy of parameters (v2)**

**local subgradients**

Worker Node 3

# Shard the parameter server

**Partition parameters across servers**
**Worker sends chunk of subgradients to owning parameter server**

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

**Worker Node 0**

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

**Worker Node 1**

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

**Worker Node 2**

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

**Worker Node 3**

**subgradient (chunk 0)**

| parameter values (chunk 0) |

**Parameter Server 0**

**subgradient (chunk 1)**

| parameter values (chunk 1) |

**Parameter Server 1**

**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 / local copy of parameters (v0) / local subgradients | training data / local copy of parameters (v1) / local subgradients | parameter values (chunk 0) |

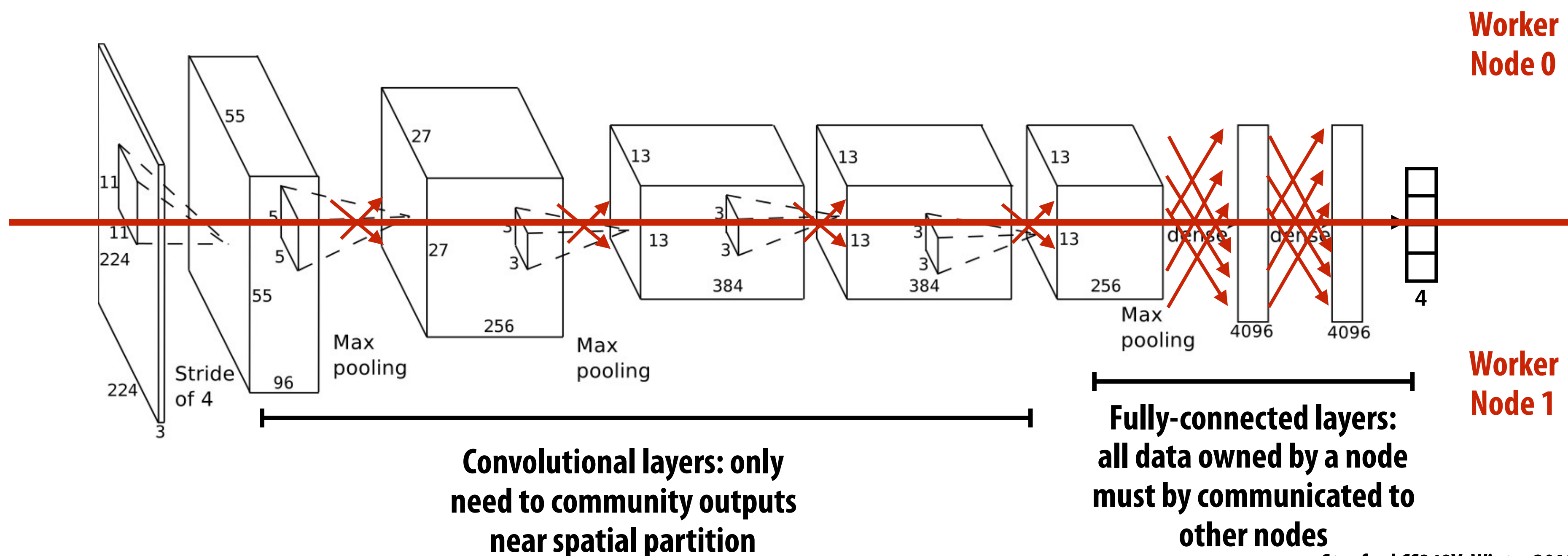| Worker Node 2 | Worker Node 3 | Parameter Server 1 |
|---|---|---|
| training data / local copy of parameters (v0) / local subgradients | training data / local copy of parameters (v2) / local subgradients | parameter values (chunk 1) |

# Model parallelism

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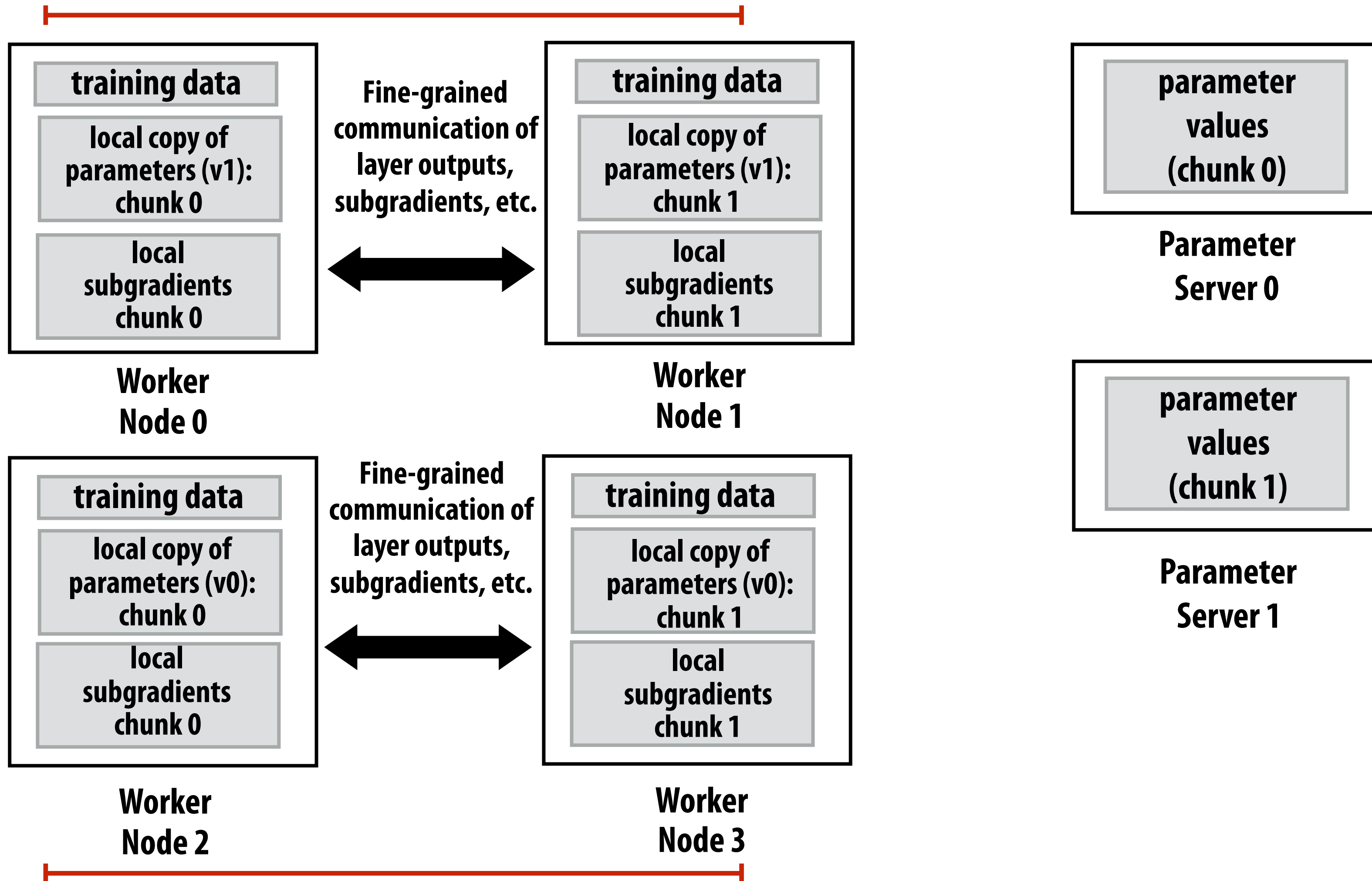(spatial partitioning to reduce communication)**

**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 community outputs near spatial partition

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

# Training data-parallel and model-parallel execution

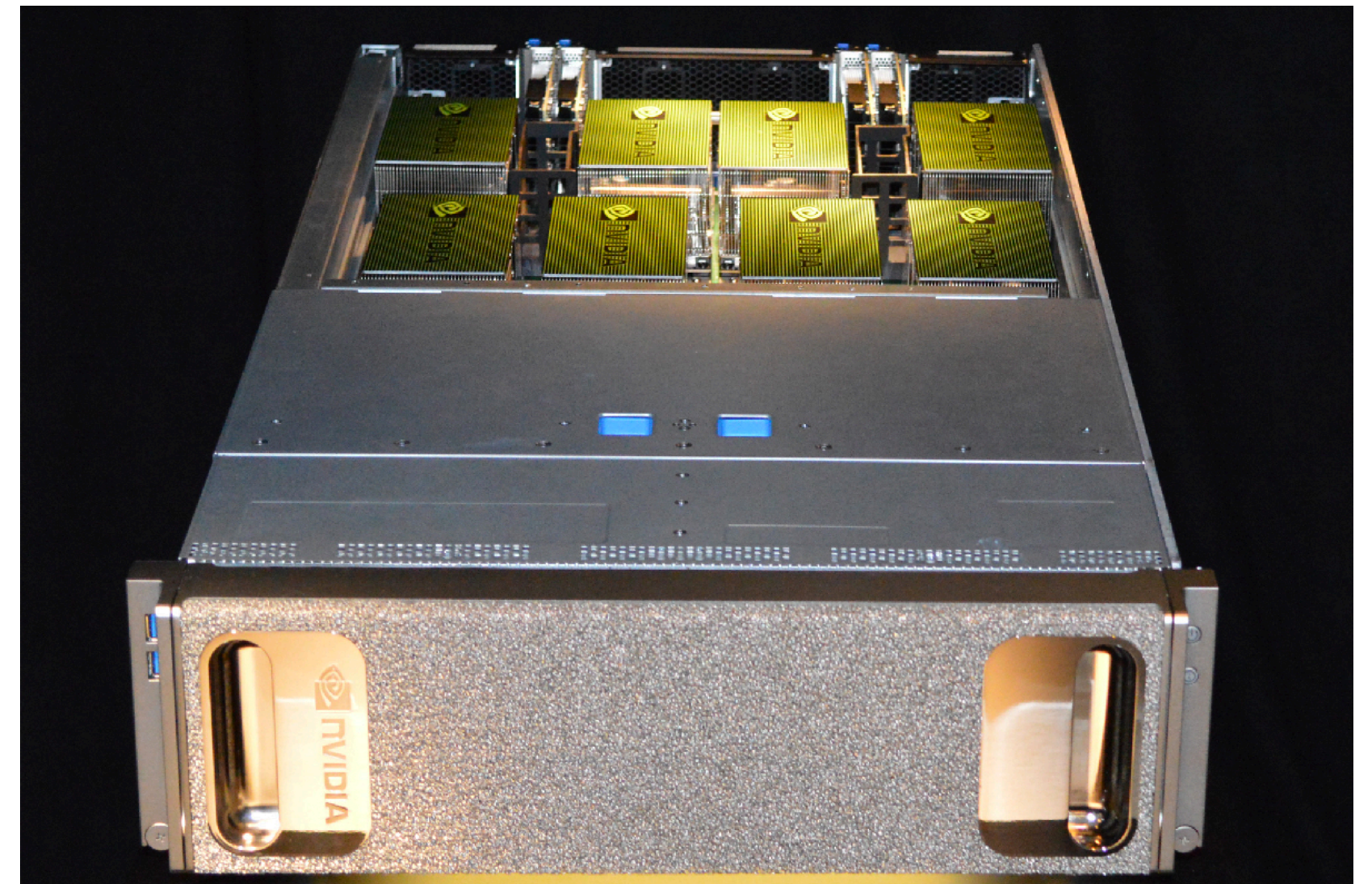**Working on subgradient computation
for a single copy of the model**

**Worker
Node 0**

- training data
- local copy of parameters (v1): chunk 0
- local subgradients chunk 0

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

**Worker
Node 1**

- training data
- local copy of parameters (v1): chunk 1
- local subgradients chunk 1

**Worker
Node 2**

- training data
- local copy of parameters (v0): chunk 0
- local subgradients chunk 0

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

**Worker
Node 3**

- training data
- local copy of parameters (v0): chunk 1
- local subgradients chunk 1

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

**Parameter
Server 0**

- parameter values (chunk 0)

**Parameter
Server 1**

- parameter values (chunk 1)

# 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 (yes, 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**

# Better 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**

# Linear scaling rule

**Recall: minibatch SGD parameter update**

size of mini batch = n
SGD learning rate = $\eta$

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

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

$$w_{t+k} = w_t - \eta \frac{1}{n} \sum_{j<k} \sum_{x \in \mathcal{B}_j} \nabla l(x, w_{t+j})$$
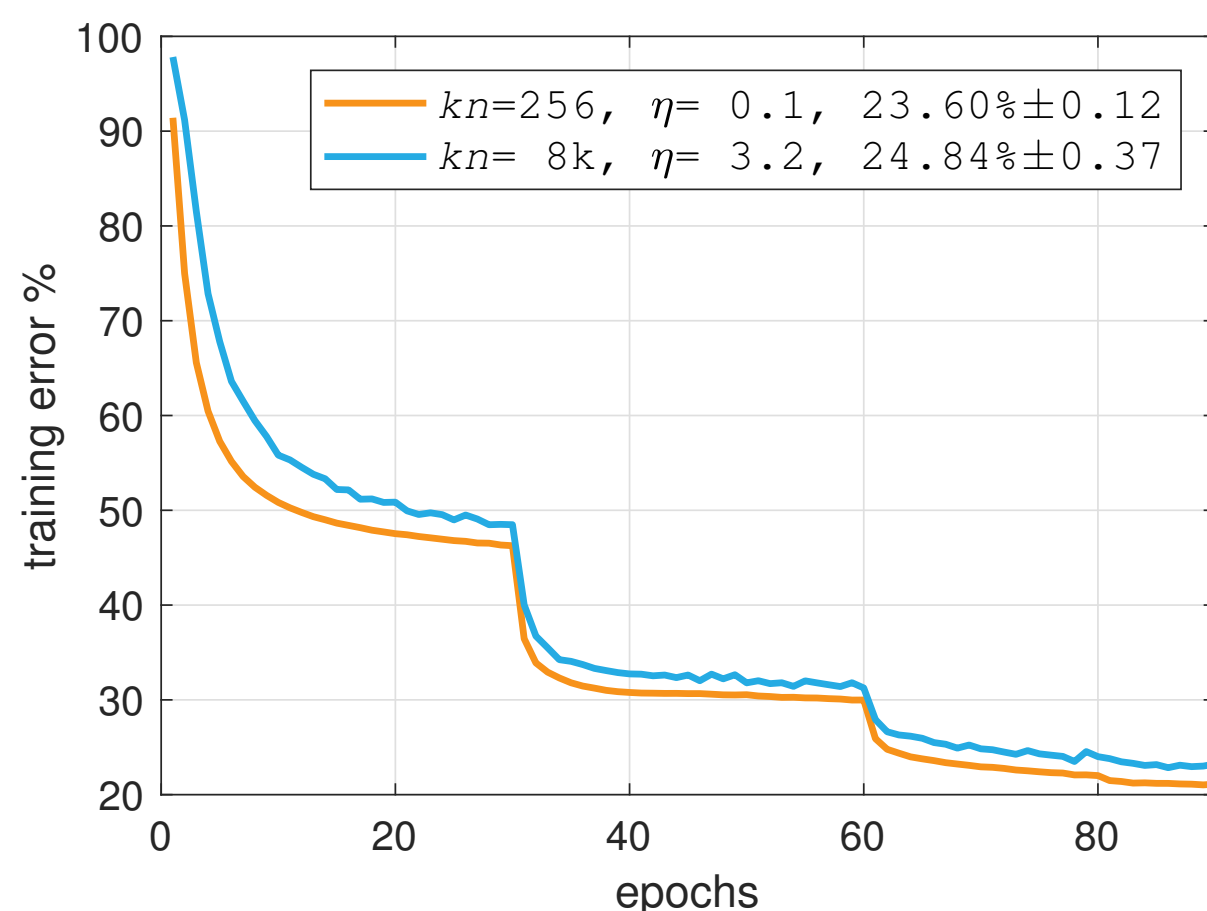
**Consider processing one minibatch that is of size kn (one step of gradient descent)**

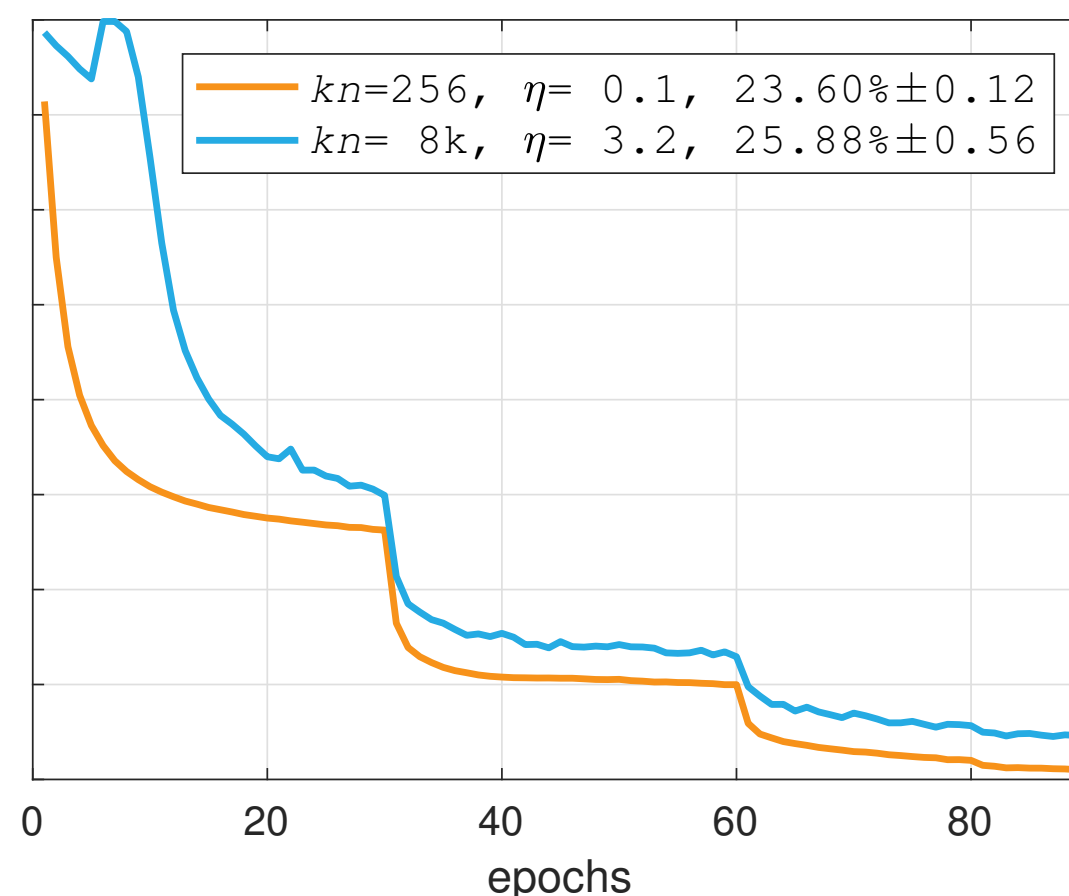$$\hat{w}_{t+1} = w_t - \hat{\eta} \frac{1}{kn} \sum_{j<k} \sum_{x \in \mathcal{B}_j} \nabla l(x, w_t)$$

**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> $k\eta$

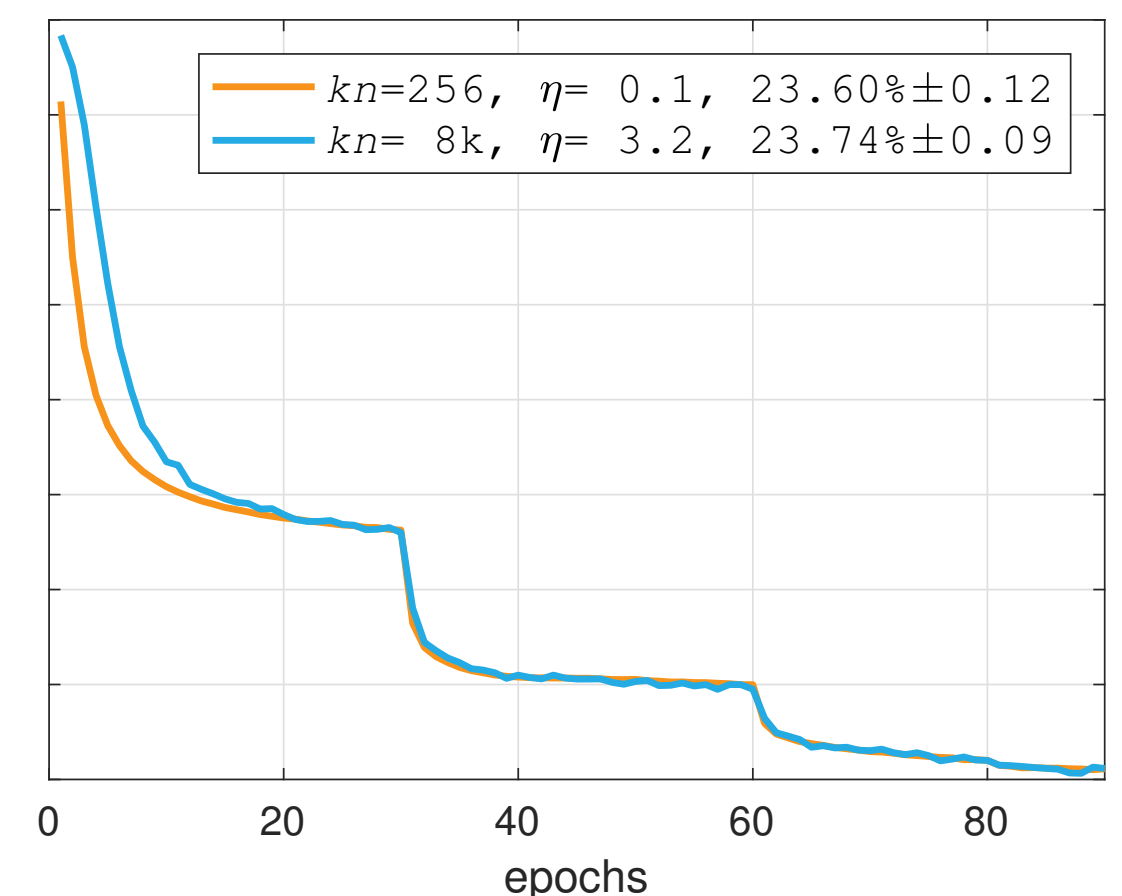# When does $\nabla l(x, w_t) \approx \nabla l(x, w_{t+j})$ not hold?

■ **At beginning of training**

- **Suggests starting training with smaller learning rate (learning rate "warmup")**

■ **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

(b) constant warmup

(c) gradual warmup

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

# Summary: training large networks in parallel

- **Many cluster systems rely on data-parallel training with 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)

- **Modern DNN designs (with fewer weights), 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!)**