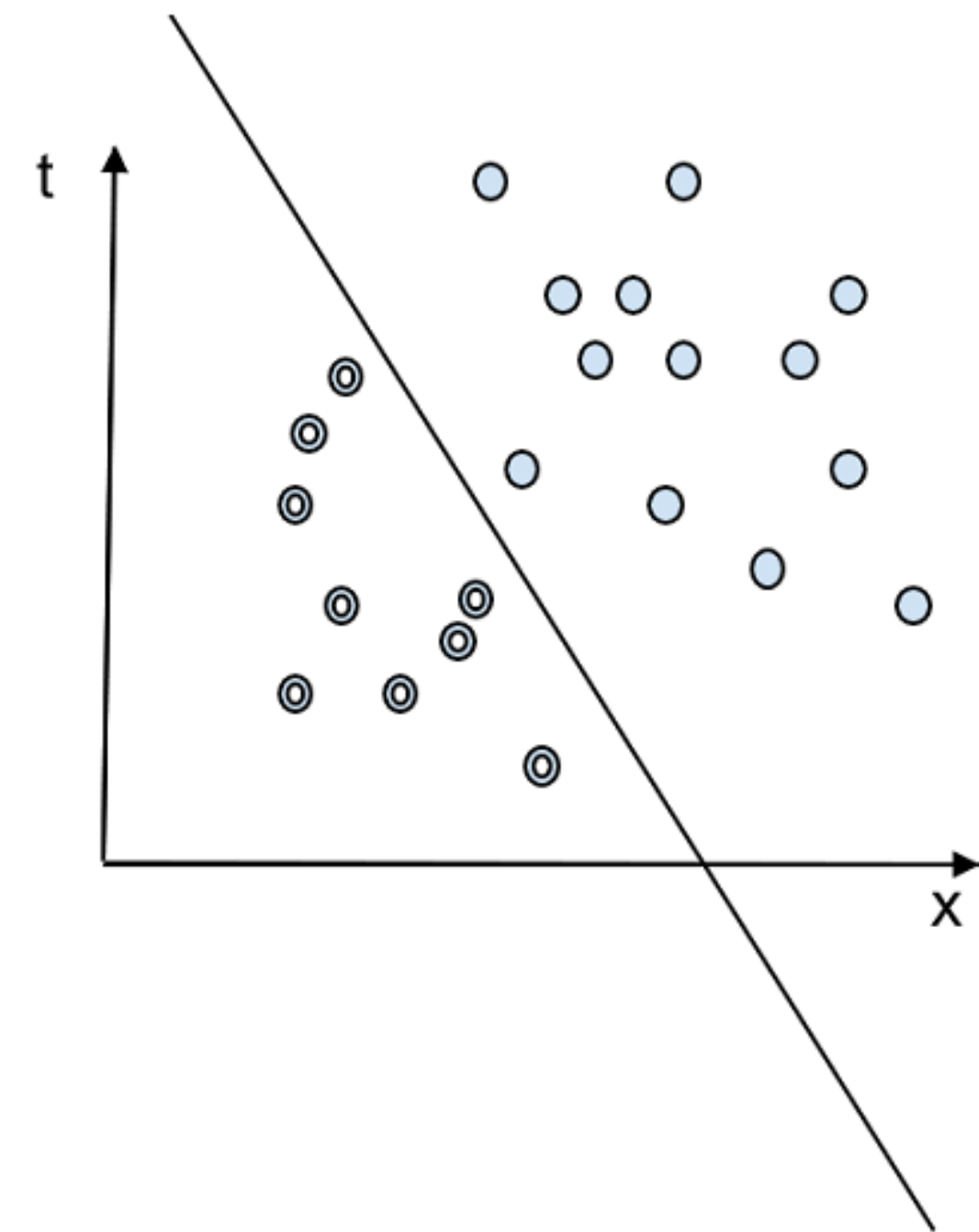
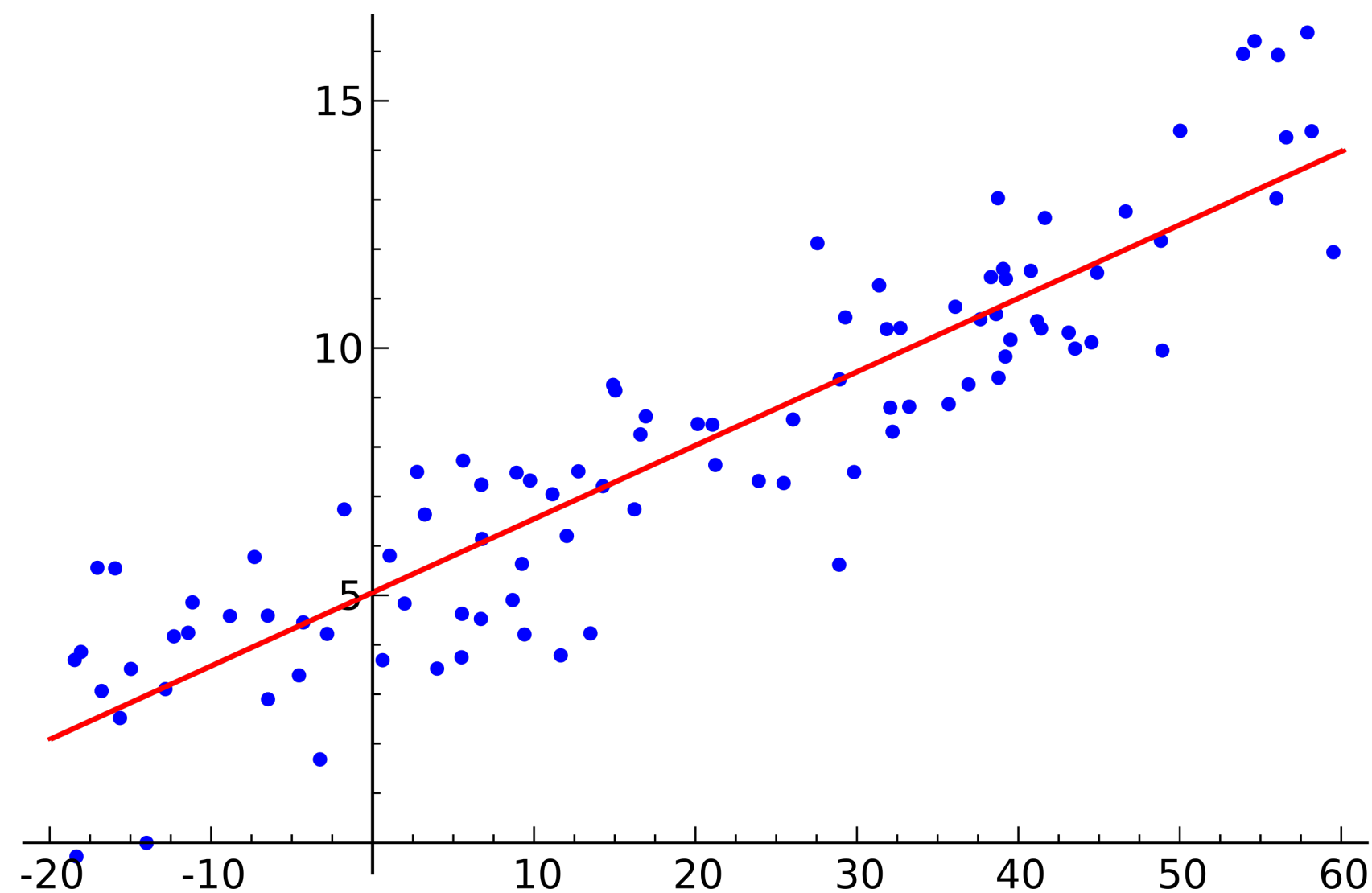


# EECE695D: Efficient ML Systems

**Compute / Memory:  
Deep Neural Networks**

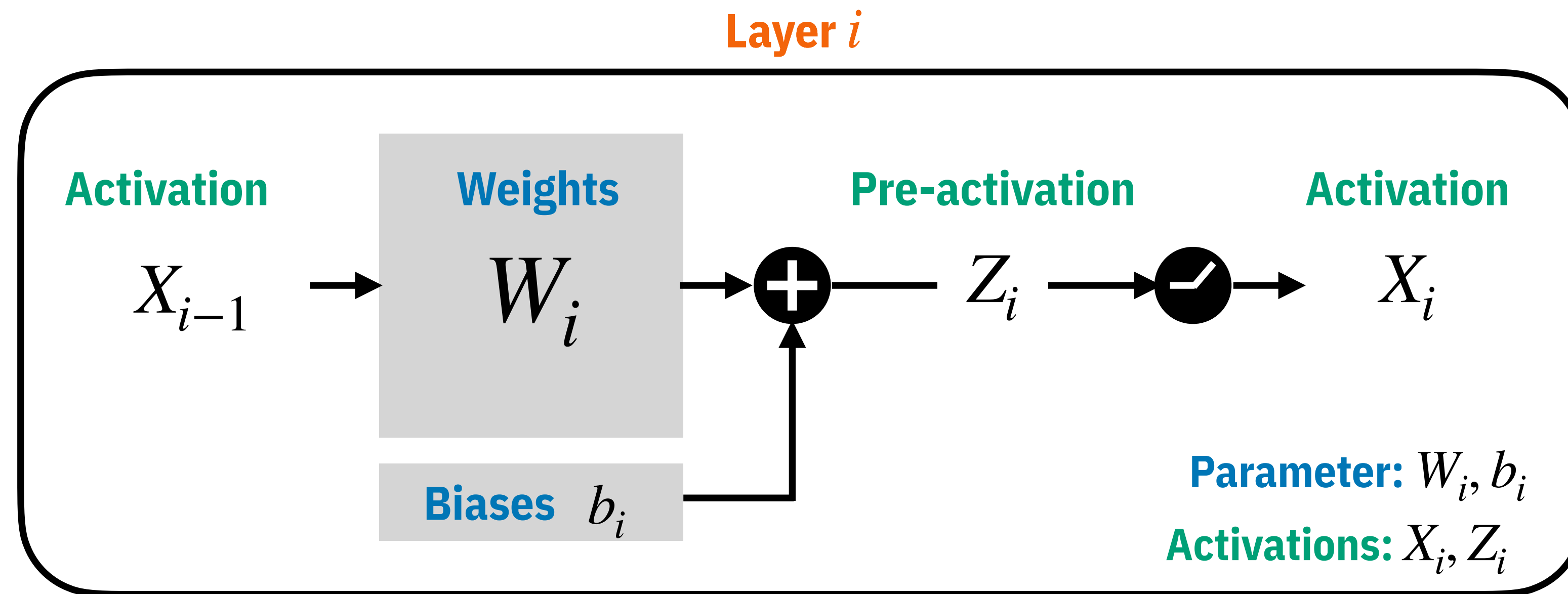
# Recap

- Compute and memory requirements of linear models (linear regression, perceptron)
  - Training vs. Inference
  - Dependency on optimization methods: Exact vs. Indirect



# Deep Neural Networks

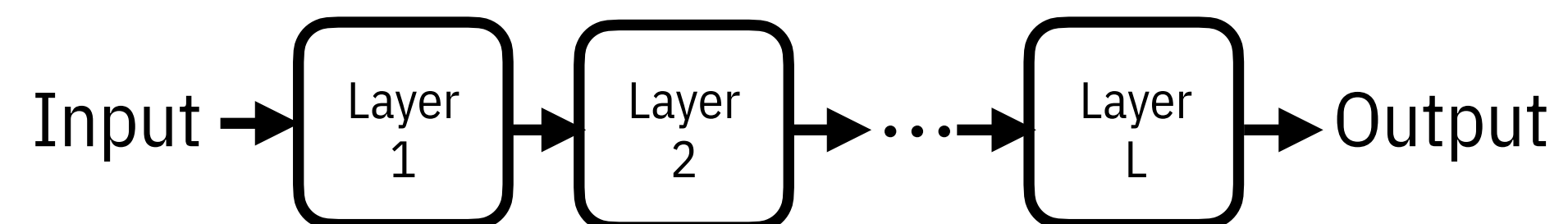
- **Neural networks.** Roughly speaking, a (directed acyclic) graph of **layers**



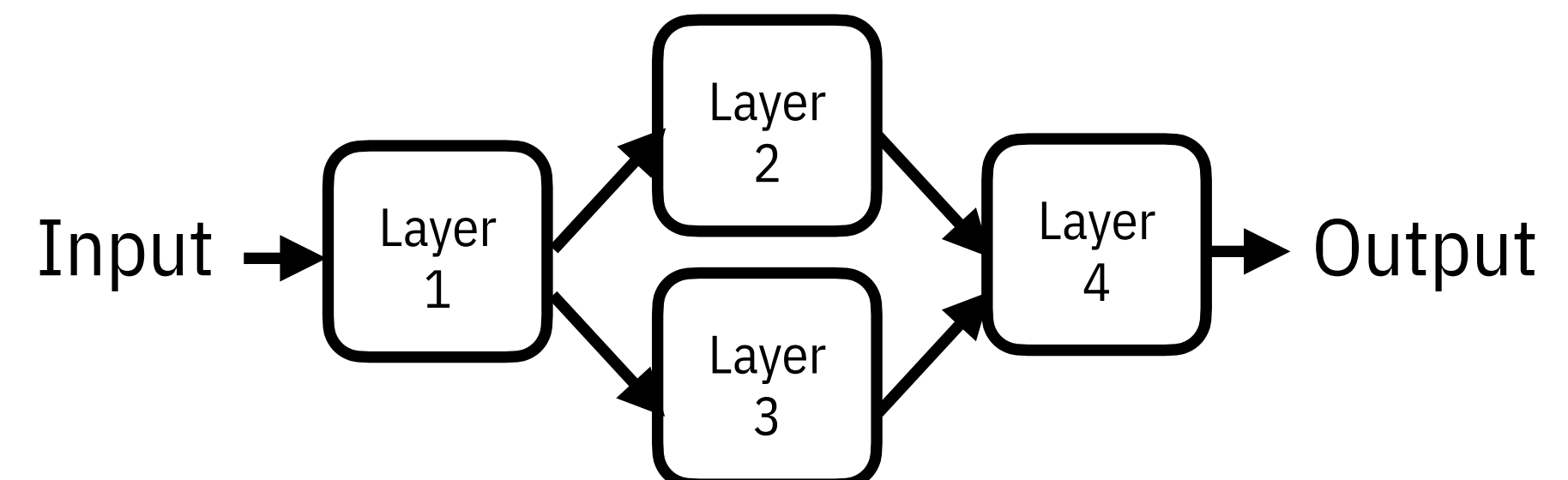
$$Z_i = W_i X_{i-1} + b_i$$

$$X_i = \sigma(Z_i)$$

**Feedforward**

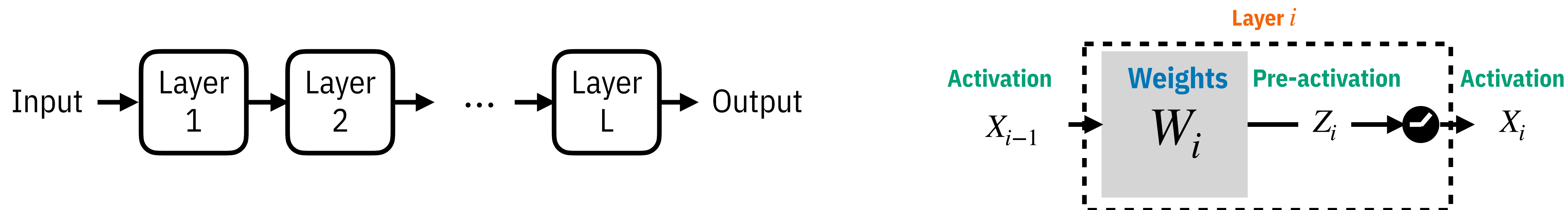


**In general**



# DNN Inference

- Let us consider a simple case: **feedforward net without bias**



- Routine.** The inference procedure can be simplified as

# Currently  $X_{i-1}$  is on cache

1. Load  $W_i$  on cache
2. Compute  $X_i = \sigma(W_i X_{i-1})$
3. Store  $X_i$  on cache #free  $X_{i-1}$
4. Repeat 1–3.

← Memory bandwidth

← Compute

*Note.* Cache (on-chip) memory is very different from the memory that causes “CUDA out of memory error.”  
The latter is an external memory— NVIDIA RTX A6000 has 48GB GDDR6 external memory and 6MB L2 cache.

# DNN Inference

# Currently  $X_{i-1}$  is on cache

1. Load  $W_i$  on cache
2. Compute  $X_i = \sigma(W_i X_{i-1})$
3. Store  $X_i$  on cache #free  $X_{i-1}$
4. Repeat 1–3.

Memory

Load  $W_i$

Load  $W_{i+1}$

Compute

Compute  $X_i$

Time

- **Overlapping.** A common practice is to overlap the loading and processing.

# Currently  $X_{i-1}$  and  $W_i$  on cache

1. Compute  $X_i = \sigma(W_i X_{i-1})$  + Load  $W_{i+1}$  on other cache
2. Store  $X_i$  on cache
3. Repeat 1–2.

Memory

Load  $W_i$

Load  $W_{i+1}$

Load  $W_{i+2}$

Compute

Compute  $X_{i-1}$

Compute  $X_i$

Compute  $X_{i+1}$

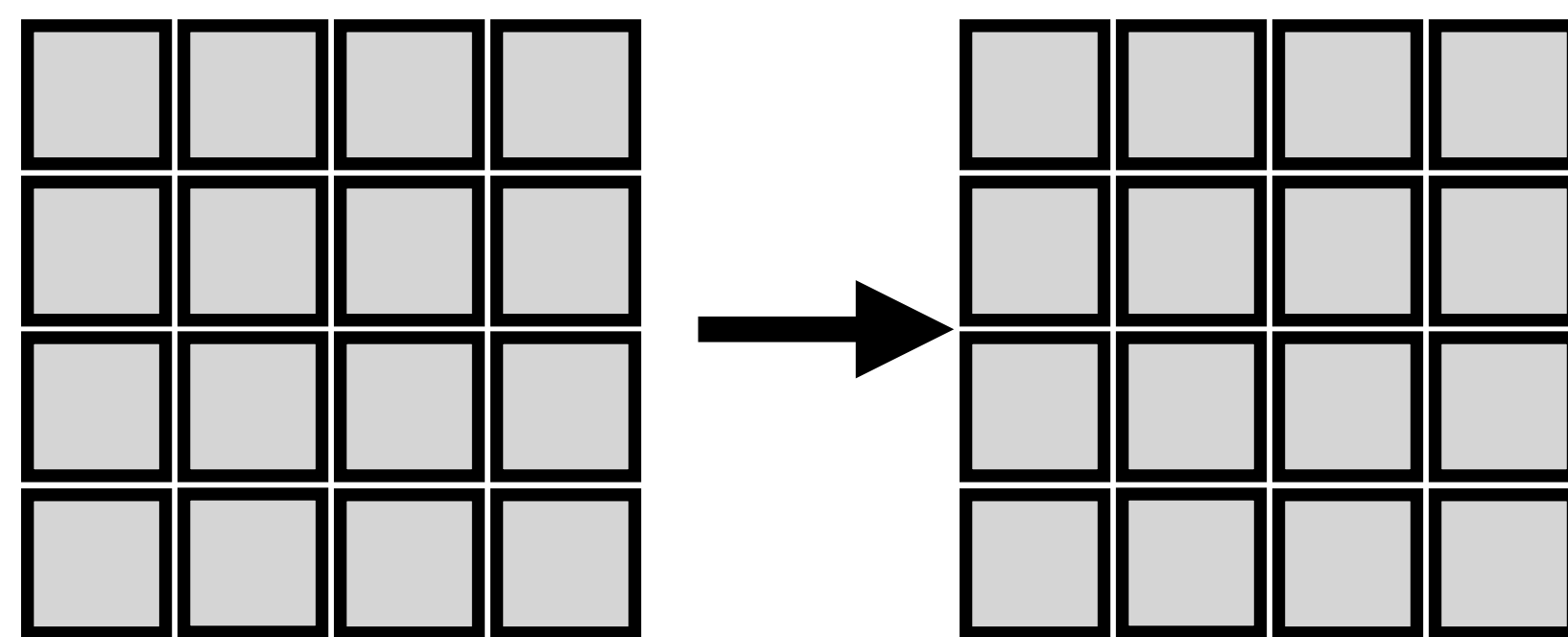
Time

*Note.* An example of compute-bound case

*Note.* Cache should be able to store two weights and two activations—double buffering

# DNN Inference

- **Convolutional Layer.** Parameter sharing—reduced compute, and much more reduced memory.  
(More likely to be compute-bound than fully-connected layers)

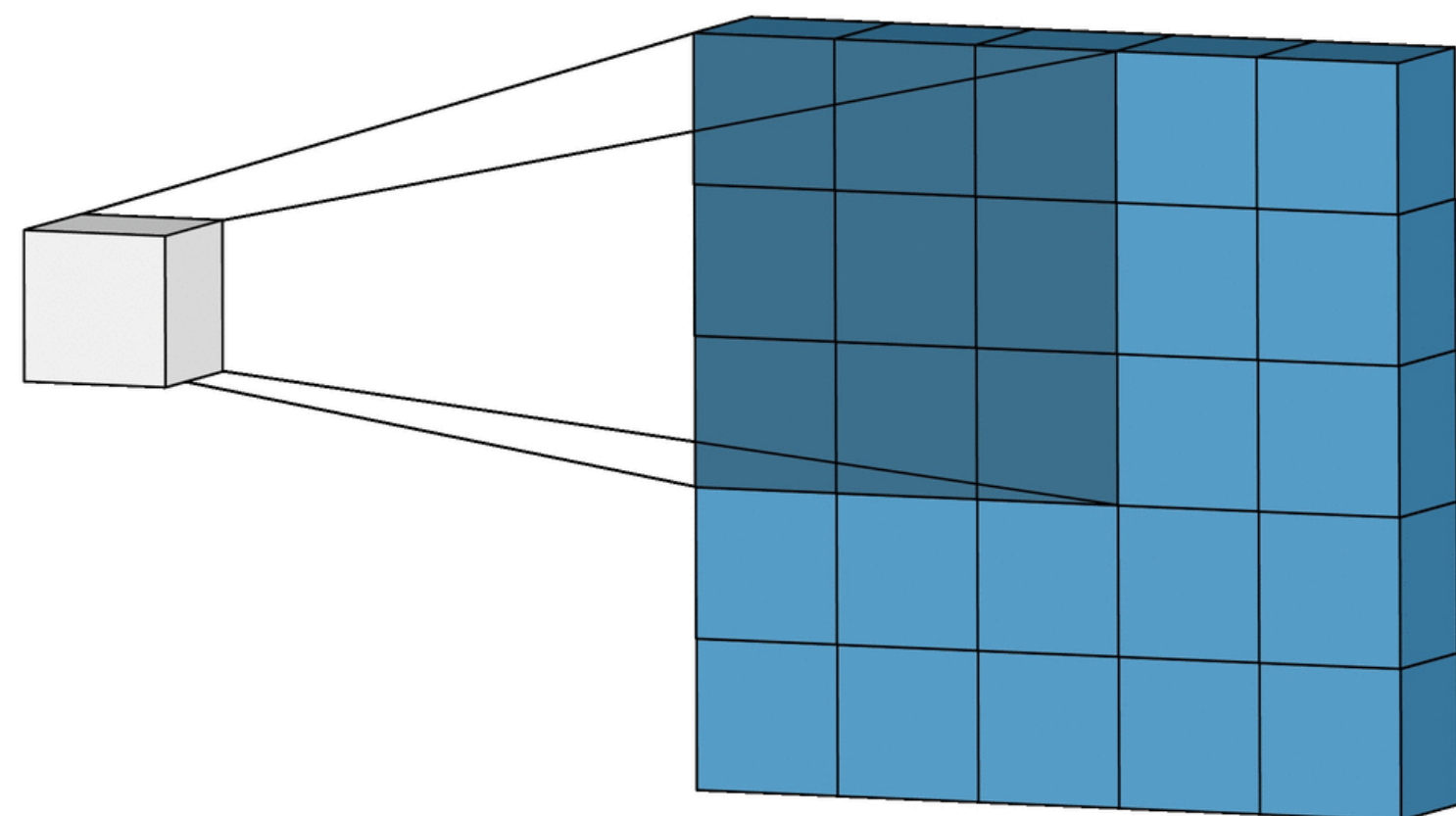


**Activation**

$$X_{i-1} \in \mathbb{R}^{d \times d \times c}$$

**Activation**

$$X_i \in \mathbb{R}^{d \times d \times c}$$



## Option 1) Fully-connected Layer

**Params.** Requires  $c^2 d^4$  parameters.

**Compute.** Requires  $2c^2 d^4$  FLOPs

## Option 2) Convolutional layer ( $3 \times 3$ )

**Params.** Requires  $9c^2$  parameters

**Compute.** Requires  $18c^2 d^2$  FLOPs

*Note 1.* Scaling-up usually done via #channels.

*Note 2.* ConvNets usually uses more layers, too.

*Note 3.* Reduction in memory for parameters, but not for activations!  
(i.e., works well for inference, but for training...?)



# Optimizing DNNs

- **Popular Choice.** We use indirect, iterative methods based on gradient descent (GD)

Mini-batch GD, Stochastic GD, with Momentum, with Nesterov Momentum, AdaGrad, AdaDelta, Adam, AdaMax, AdamP, RAdam, NAdam, RMSProp, Shampoo, Lookahead, SAM, ...

- **Why Not Direct?** (1) Unique direct solution does not really exist

# e.g., ReLU

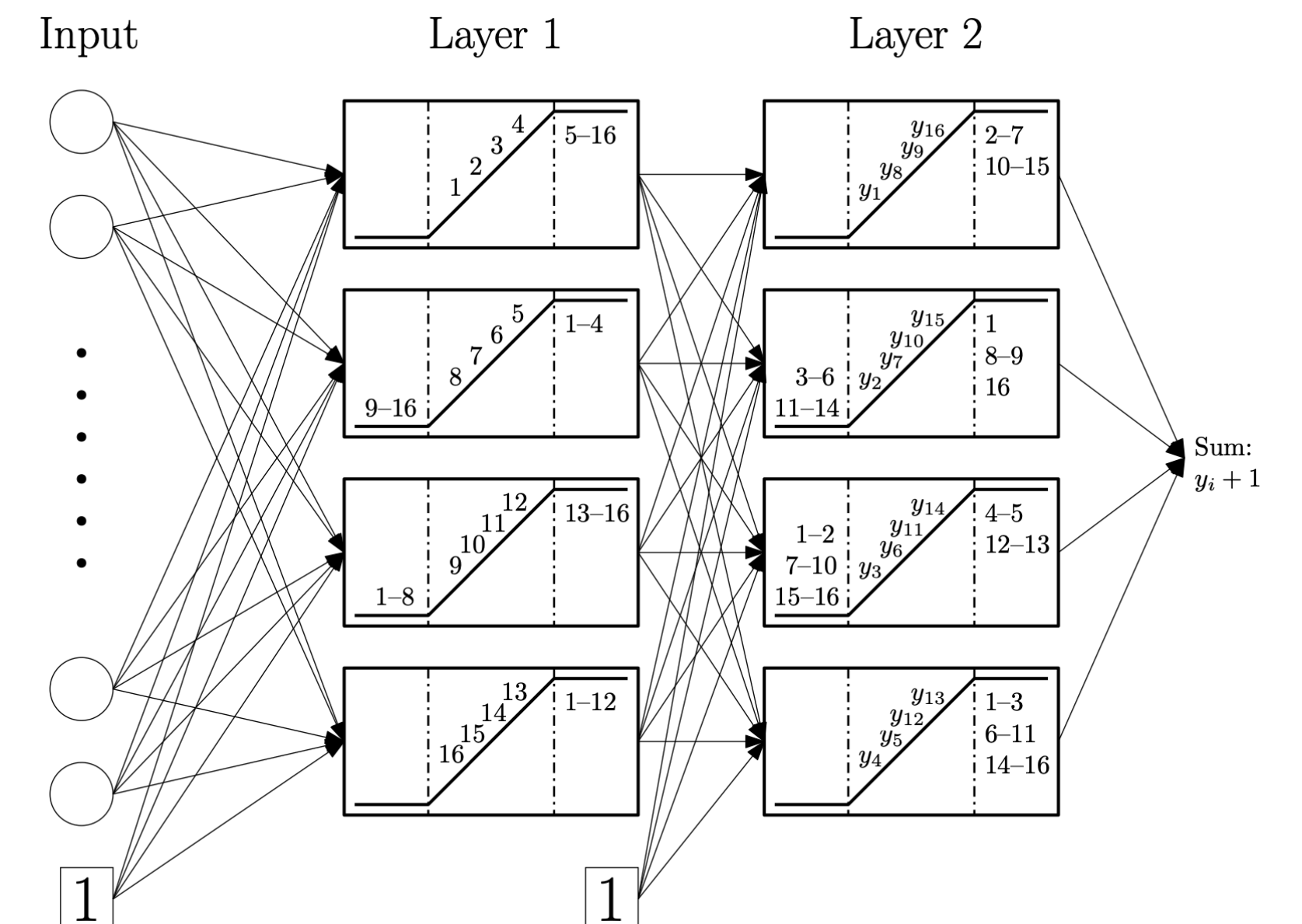
- (2) Even if it exists, expensive to solve

# we've seen  $\mathcal{O}(d^3)$  dependency

- (3) Even when easy to solve, the solution is often wiggly and never generalizes...

$$\mathbf{W}_{j,:}^1 = (-1)^{j-1} \frac{4}{c_{jq} + c_{jq+1} - c_{jq-q} - c_{jq-q+1}} \mathbf{u}^T,$$

$$\mathbf{b}_j^1 = (-1)^j \frac{c_{jq} + c_{jq+1} + c_{jq-q} + c_{jq-q+1}}{c_{jq} + c_{jq+1} - c_{jq-q} - c_{jq-q+1}}.$$



# Optimizing DNNs

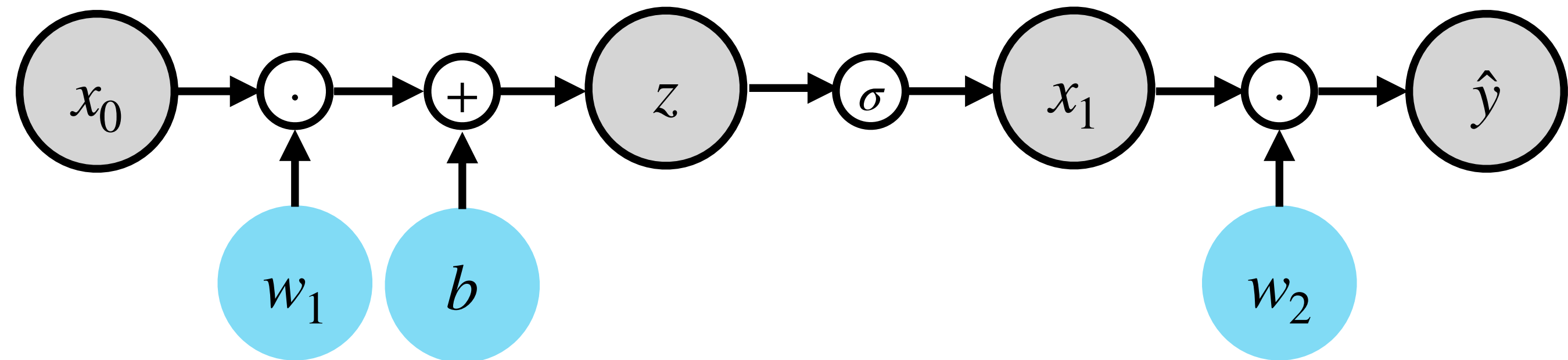
- **Backpropagation.** A “cheap” way to compute gradients (for GD) on multi-layered neural nets.  
Previous approaches required stochastic units & much more compute.
- **Idea.** Computing the gradients can be made less compute-heavy,  
by re-using (1) forward-compute and (2) gradients of succeeding layer activations.
- **Example.** Two-layer ReLU net with 1D parameters

$$z = w_1 \cdot x_0 + b$$

$$x_1 = \sigma(z)$$

$$\hat{y} = w_2 x_1$$

$$L = \frac{1}{2} (y - \hat{y})^2$$



**Q.** Gradients for parameters,  $w_1, w_2, b$ ?  
(evaluated at  $w_1^\circ, w_2^\circ, b^\circ$ )



$$\frac{\partial L}{\partial w_2} \bigg|_{w^\circ, b^\circ} =$$

$$\frac{\partial L}{\partial x_1} \bigg|_{w^\circ, b^\circ} =$$

$$\frac{\partial L}{\partial w_1} \bigg|_{w^\circ, b^\circ} =$$

$$\frac{\partial L}{\partial b} \bigg|_{w^\circ, b^\circ} =$$

$$z = w_1 \cdot x_0 + b$$

$$x_1 = \sigma(z)$$

$$\hat{y} = w_2 x_1$$

$$L = \frac{1}{2} (y - \hat{y})^2$$

# Compare: Forward mode AD

- **Backprop.** An example of reverse mode automatic differentiation.

$$\hat{y} = f(g(h(x)))$$
$$\frac{\partial \hat{y}}{\partial x} = \frac{\partial \hat{y}}{\partial f} \cdot \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial h} \cdot \frac{\partial h}{\partial x}$$

Backward  
Forward

- **Forward.** (Mem. Req.) = Constant # one forward pass for each input, keep one for gradient, one for state  
Runtime  $\propto$  #(Input Variables)
- **Backward.** (Mem. Req.)  $\propto$  #(Intermediate Variables) # single, monolithic backward pass  
Runtime  $\propto$  #(Output Variables)
- **When to use forward?** Ran out of external memory—very common in meta-learning  
Want gradient for very small number of variables—e.g., gradient w.r.t. input vector

# FLOPs of Backprop

- **Rule of Thumb.** (Backward FLOPs)  $\approx 2 \cdot$  (Forward FLOPs)
- **Rough Idea.** Backward FLOPs update both activations and weights, while forward FLOPs are computing activations only.
- **Example.** Three-layer MLP without bias.

$$\hat{y} = W_3 \sigma(W_2 \sigma(W_1 X_0))$$

The diagram illustrates the forward pass of a three-layer MLP. It shows the input  $X_0 \in \mathbb{R}^{d \times N}$  (in a grey box) being multiplied by weight matrices  $W_1 \in \mathbb{R}^{h \times d}$  and  $W_2 \in \mathbb{R}^{h \times h}$  (both in blue boxes), with activation functions  $\sigma$  applied between layers. The final output is  $\hat{y} \in \mathbb{R}^{k \times N}$  (in a green box), which is the result of multiplying the output of the second layer by  $W_3 \in \mathbb{R}^{k \times h}$  (in a blue box).

$$\hat{y} \in \mathbb{R}^{k \times N} = W_3 \in \mathbb{R}^{k \times h} \sigma \left( W_2 \in \mathbb{R}^{h \times h} \sigma \left( W_1 \in \mathbb{R}^{h \times d} X_0 \in \mathbb{R}^{d \times N} \right) \right)$$

$$\hat{y} = W_3 \sigma(W_2 \sigma(W_1 X_0))$$

$$\hat{y} \in \mathbb{R}^{k \times N} = W_3 \in \mathbb{R}^{k \times h} \sigma \left( W_2 \in \mathbb{R}^{h \times h} \sigma \left( W_1 \in \mathbb{R}^{h \times d} X_0 \in \mathbb{R}^{d \times N} \right) \right)$$

### Forward FLOPs

$$\begin{array}{ccc} \text{Layer 3} & \text{Layer 2} & \text{Layer 1} \\ X_3 = W_3 X_2 & \leftarrow Z_2 = W_2 X_1 & \leftarrow Z_1 = W_1 X_0 \\ 2hkN & 2h^2N & 2dhN \end{array}$$

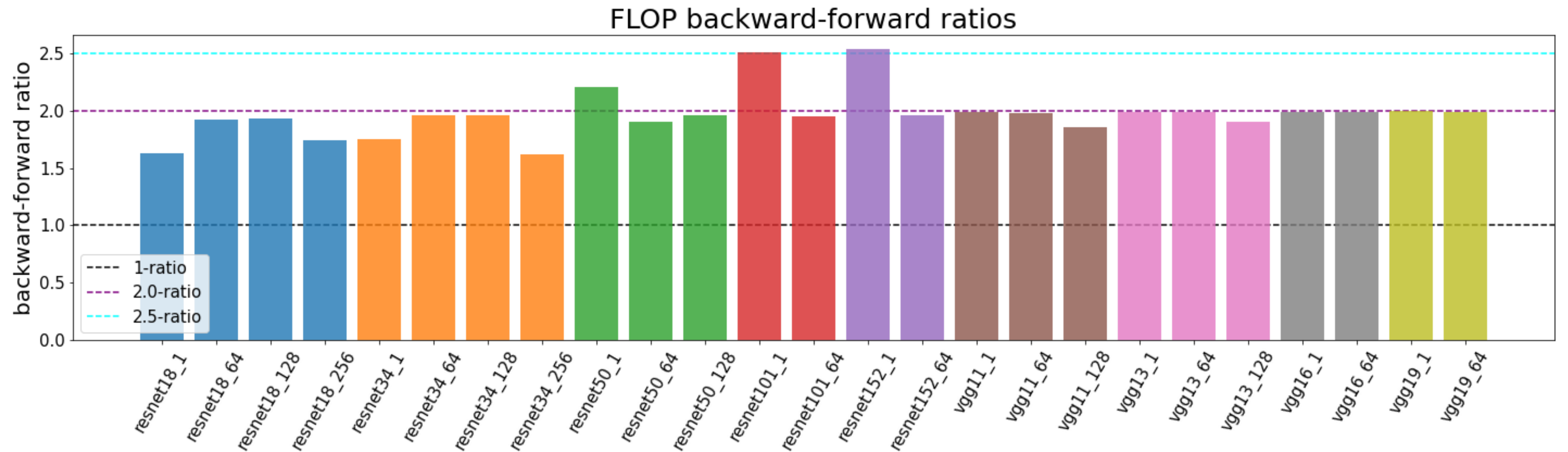
### Backward FLOPs

$$\begin{array}{ccc} \text{Layer 3} & \text{Layer 2} & \text{Layer 1} \\ \frac{\partial L}{\partial W_3} = (\hat{y} - y) X_2^\top & \frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial X_2} X_1^\top & \frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial X_1} X_0^\top \\ 2hkN & 2h^2N & 2dhN \\ \frac{\partial L}{\partial X_2} = W_3^\top (\hat{y} - y) & \frac{\partial L}{\partial X_1} = W_2^\top \frac{\partial L}{\partial X_2} & \end{array}$$

Note. Ignored ReLU!

Note. If memory is large,  
prepare transposed ones!

Note. Don't forget the FLOPs for  
the weight update... how big?



*Note.* Does not mean that optimizing inference == optimizing training!  
(Batch issues, steps until convergence, etc)

# Memory of Backprop

- Revisit the previous example:

$$\hat{y} \in \mathbb{R}^{k \times N} = W_3 \in \mathbb{R}^{k \times h} \sigma \left( W_2 \in \mathbb{R}^{h \times h} \sigma \left( W_1 \in \mathbb{R}^{h \times d} X_0 \in \mathbb{R}^{d \times N} \right) \right)$$

- Forward pass.** Need to have parameters  $W_1$ ,  $W_2$ ,  $W_3$  and data  $X_0$  loaded on memory, and activations  $X_1$ ,  $X_2$ ,  $\hat{y}$  computed and loaded to memory.

- Backward pass.** Additionally have parameter updates for  $\frac{\partial L}{\partial W_1}$ ,  $\frac{\partial L}{\partial W_2}$ ,  $\frac{\partial L}{\partial W_3}$ ,  
using intermediates  $\frac{\partial L}{\partial X_1}$ ,  $\frac{\partial L}{\partial X_2}$ ,  $\frac{\partial L}{\partial \hat{y}}$ .

*Note.* Once you have,  $\left( \frac{\partial L}{\partial X_2}, \frac{\partial L}{\partial X_2} \right)$ , you no longer need  $X_2$ ,  $W_2$  on memory.



# Memory of Backprop

**Forward**

$$Z_1 = W_1 X_0 \longrightarrow Z_2 = W_2 X_1 \longrightarrow X_3 = W_3 X_2$$

**Backward**

$$\begin{array}{l} \frac{\partial L}{\partial W_3} = (\hat{y} - y) X_2^\top \\ \frac{\partial L}{\partial X_2} = W_3^\top (\hat{y} - y) \end{array} \longrightarrow \begin{array}{l} \frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial X_2} X_1^\top \\ \frac{\partial L}{\partial X_1} = W_2^\top \frac{\partial L}{\partial X_2} \end{array} \longrightarrow \begin{array}{l} \frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial X_1} X_0^\top \end{array}$$

**Brain-teaser.** Suppose that all tensors are of uniform size, and that you have enough cache to store four tensors. How many tensor movements (in/out) do you need, in order to perform all computations?

What if you have a cache enough to store five tensors?

# Side Note: Gradient Checkpointing

- In a sense, backpropagation is a way to use extra memory for less compute  
(for storing activations computed during forward)

**Q.** Can we use less extra memory,  
at the cost of slightly increased compute?

**A.** Yes—no need to store all activations.

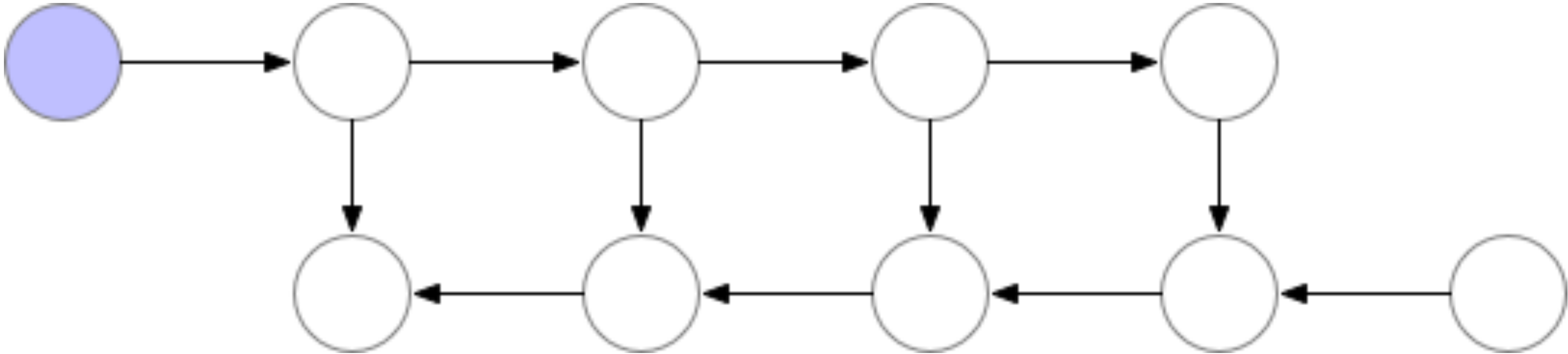
Instead, compute them again with extra  
forward operations!

## Training Deep Nets with Sublinear Memory Cost

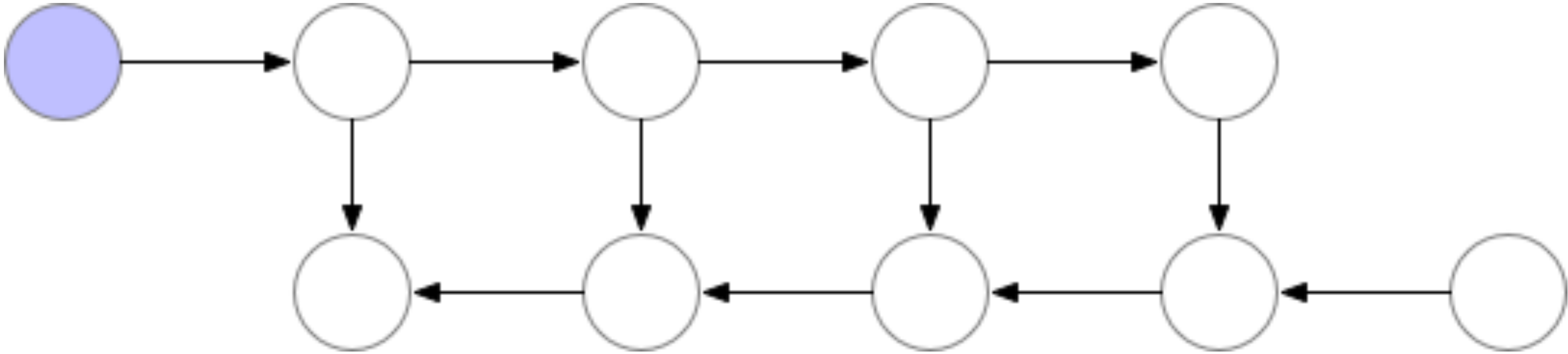
**Tianqi Chen**<sup>1</sup>, **Bing Xu**<sup>2</sup>, **Chiyuan Zhang**<sup>3</sup>, and **Carlos Guestrin**<sup>1</sup>

<sup>1</sup> University of Washington   <sup>2</sup> Dato. Inc   <sup>3</sup> Massachusetts Institute of Technology

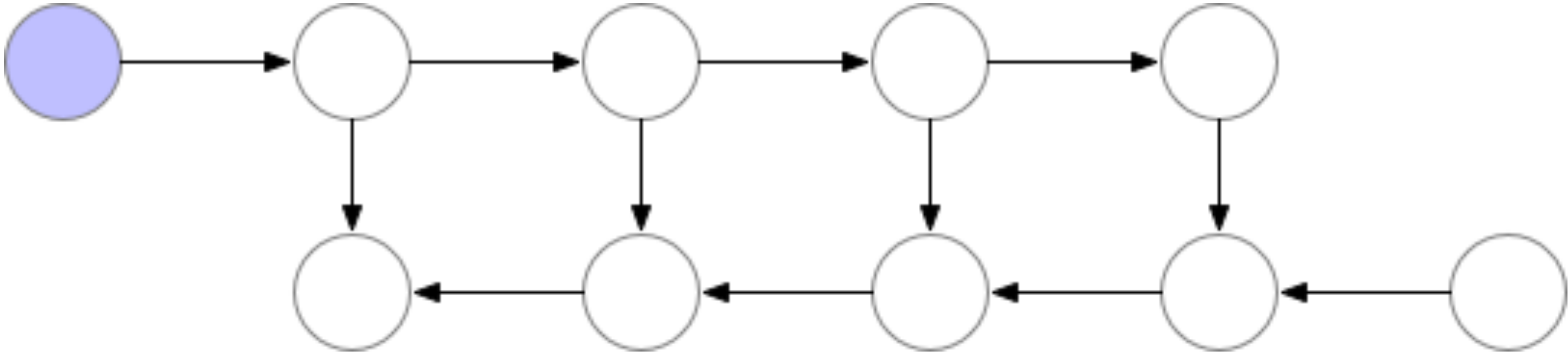
**Vanilla**



**Low-Memory  
(w/o checkpoint)**

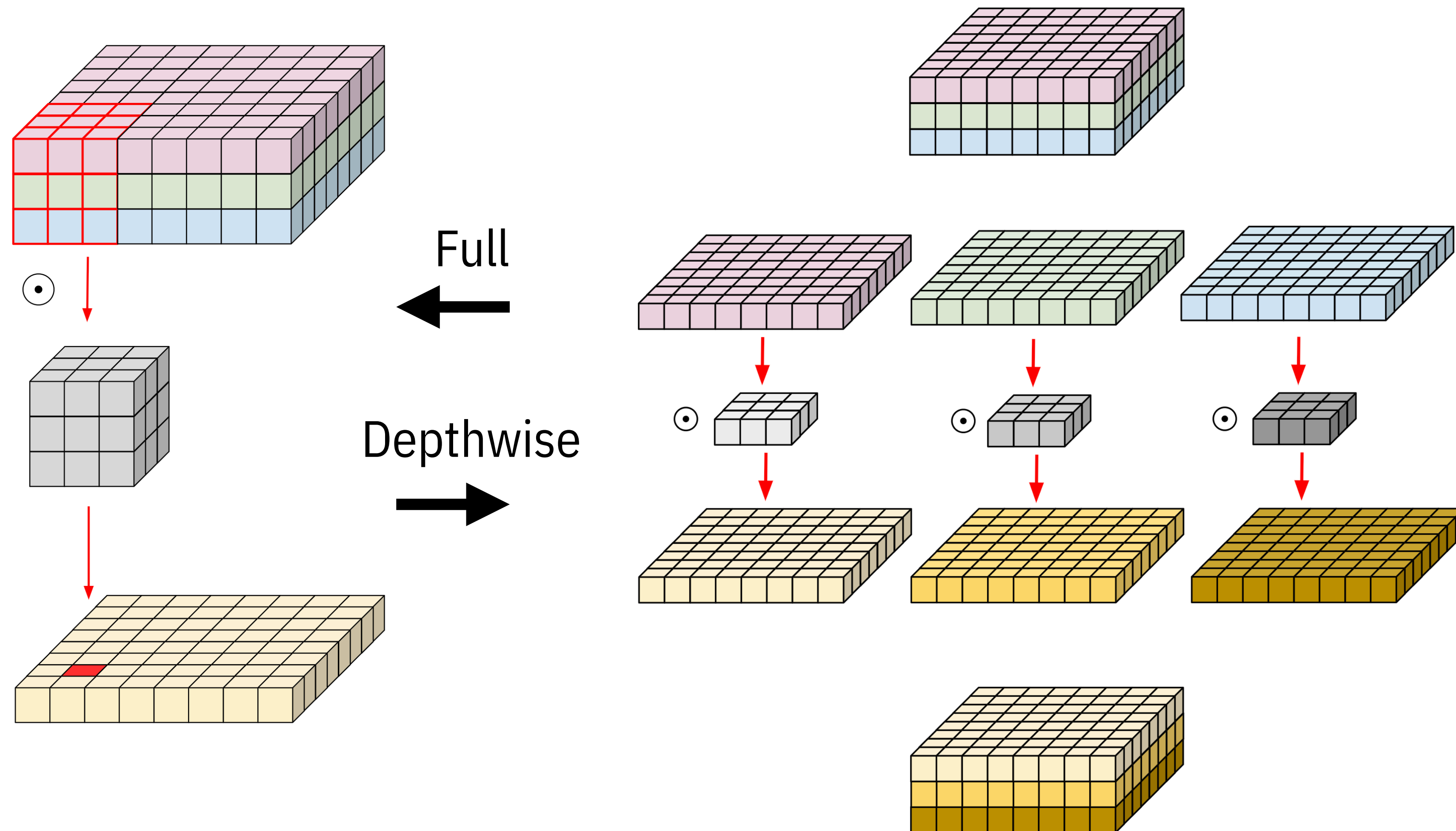


**Checkpointing**



# Side Note: Depthwise Convolution

- In most modern “efficient” neural network architectures (e.g., MobileNet, EfficientNet), people use **depthwise convolutions** instead of full convolutions.



**Params**

*Full.  $c^2k^2$*

*DW.  $ck^2$*

**FLOPs**

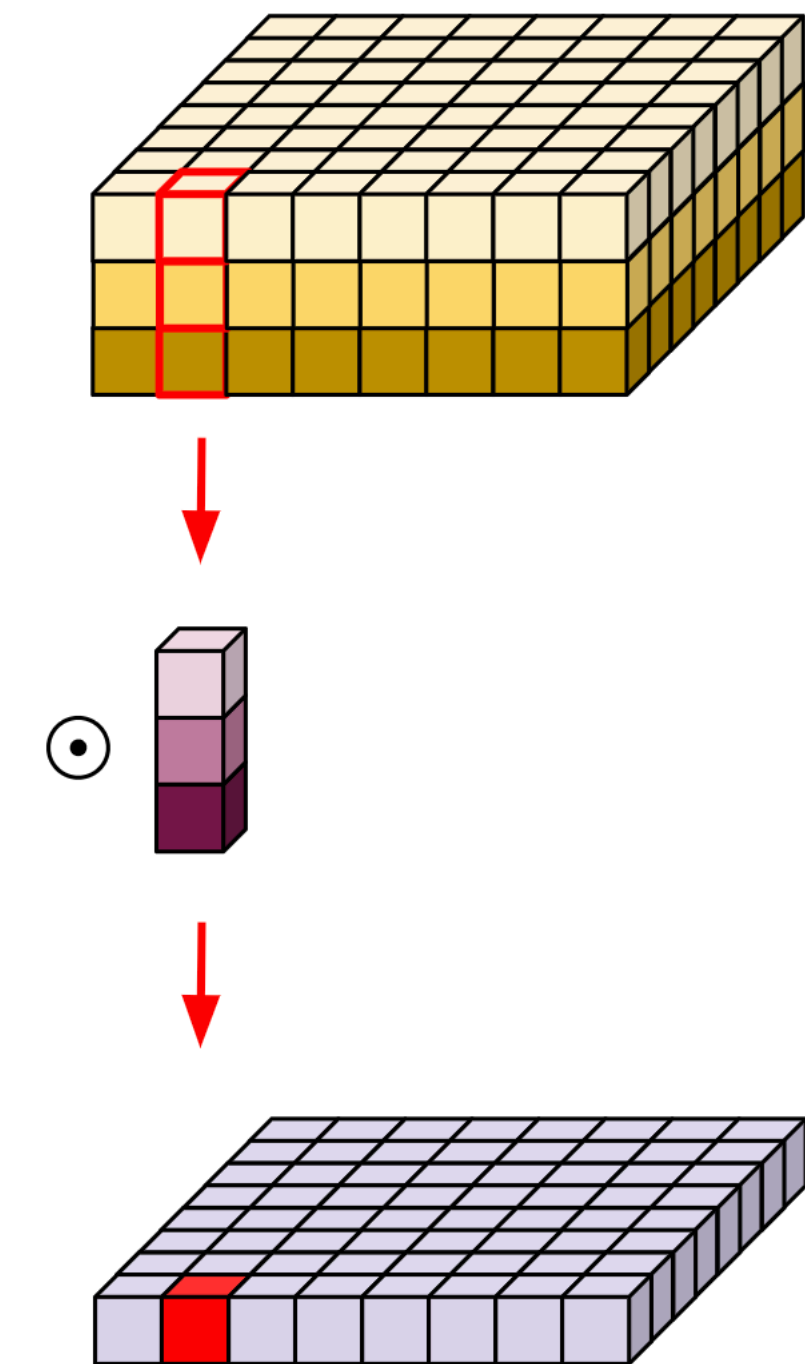
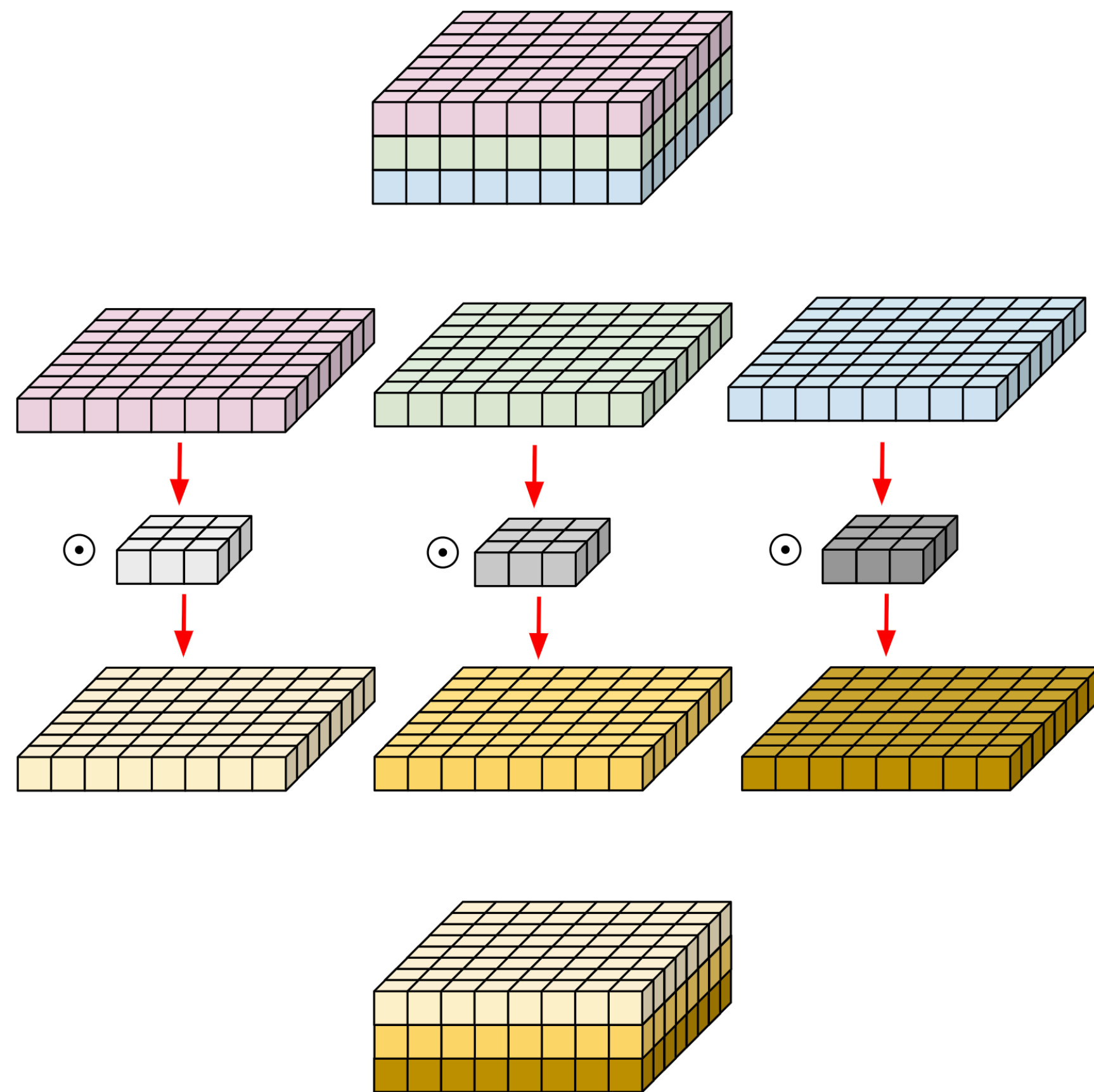
*Full.  $2c^2k^2d^2$*

*DW.  $2ck^2d^2$*

# Side Note: Depthwise Convolution

- **Problem.** Channel information cannot be mixed!  
⇒ Use one-by-one convolution to mix the layers

# e.g., Inverted Bottleneck





# Side Note: Optimization algorithm!

- Some optimizers have extra “state variables” for various reasons.
- For instance, we sometimes utilize “momentum” terms.  
Beware of the memory overhead!!

---

**Algorithm 1:** *Adam*, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation.  $g_t^2$  indicates the elementwise square  $g_t \odot g_t$ . Good default settings for the tested machine learning problems are  $\alpha = 0.001$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\epsilon = 10^{-8}$ . All operations on vectors are element-wise. With  $\beta_1^t$  and  $\beta_2^t$  we denote  $\beta_1$  and  $\beta_2$  to the power  $t$ .

---

**Require:**  $\alpha$ : Stepsize

**Require:**  $\beta_1, \beta_2 \in [0, 1)$ : Exponential decay rates for the moment estimates

**Require:**  $f(\theta)$ : Stochastic objective function with parameters  $\theta$

**Require:**  $\theta_0$ : Initial parameter vector

$m_0 \leftarrow 0$  (Initialize 1<sup>st</sup> moment vector)

$v_0 \leftarrow 0$  (Initialize 2<sup>nd</sup> moment vector)

$t \leftarrow 0$  (Initialize timestep)

**while**  $\theta_t$  not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$  (Get gradients w.r.t. stochastic objective at timestep  $t$ )

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$  (Update biased first moment estimate)

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$  (Update biased second raw moment estimate)

$\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$  (Compute bias-corrected first moment estimate)

$\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$  (Compute bias-corrected second raw moment estimate)

$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$  (Update parameters)

**end while**

**return**  $\theta_t$  (Resulting parameters)

---



# Summary

- Deep networks, as a stack of linear models.
- DNN inference: Compute vs. Memory
  - Parameter sharing — Less bandwidth
  - Overlapping — Compute utilization
- DNN training
  - Backprop — Using extra memory for less compute (smooth tradeoff via checkpointing)
  - Forward FLOPs vs. Backward FLOPs
- **Next.** A bit more about efficiency criteria  
(compute / memory / peak memory / #params)  
Some case studies...