

CS 4803 / 7643: Deep Learning

Topics:

- (Finish) Computing Gradients
- Backprop in FC+ReLU NNs

- Convolutional Neural Networks

Dhruv Batra
Georgia Tech

Administrativa

- HW1 Reminder
 - Due: 10/02, 11:55pm
- HW0 grades
 - Out week of 10/01

Project

- Goal
 - Chance to try Deep Learning
 - **Combine with other classes / research / credits / anything**
 - You have our blanket permission
 - Extra credit for shooting for a publication
 - Get permission from other instructors; delineate different parts
 - Encouraged to apply to your research (computer vision, NLP, robotics,...)
 - Must be done this semester.
 - Groups of 3-4
- Main categories
 - **Application/Survey**
 - Compare a bunch of existing algorithms on a new application domain of your interest
 - **Formulation/Development**
 - Formulate a new model or algorithm for a new or old problem
 - **Theory**
 - Theoretically analyze an existing algorithm

Computing

- Major bottleneck
 - GPUs
- Options
 - Your own / group / advisor's resources
 - Google Cloud Credits
 - \$50 credits to every registered student courtesy Google
 - <https://colab.research.google.com>
 - Minsky cluster in IC

Administrativa

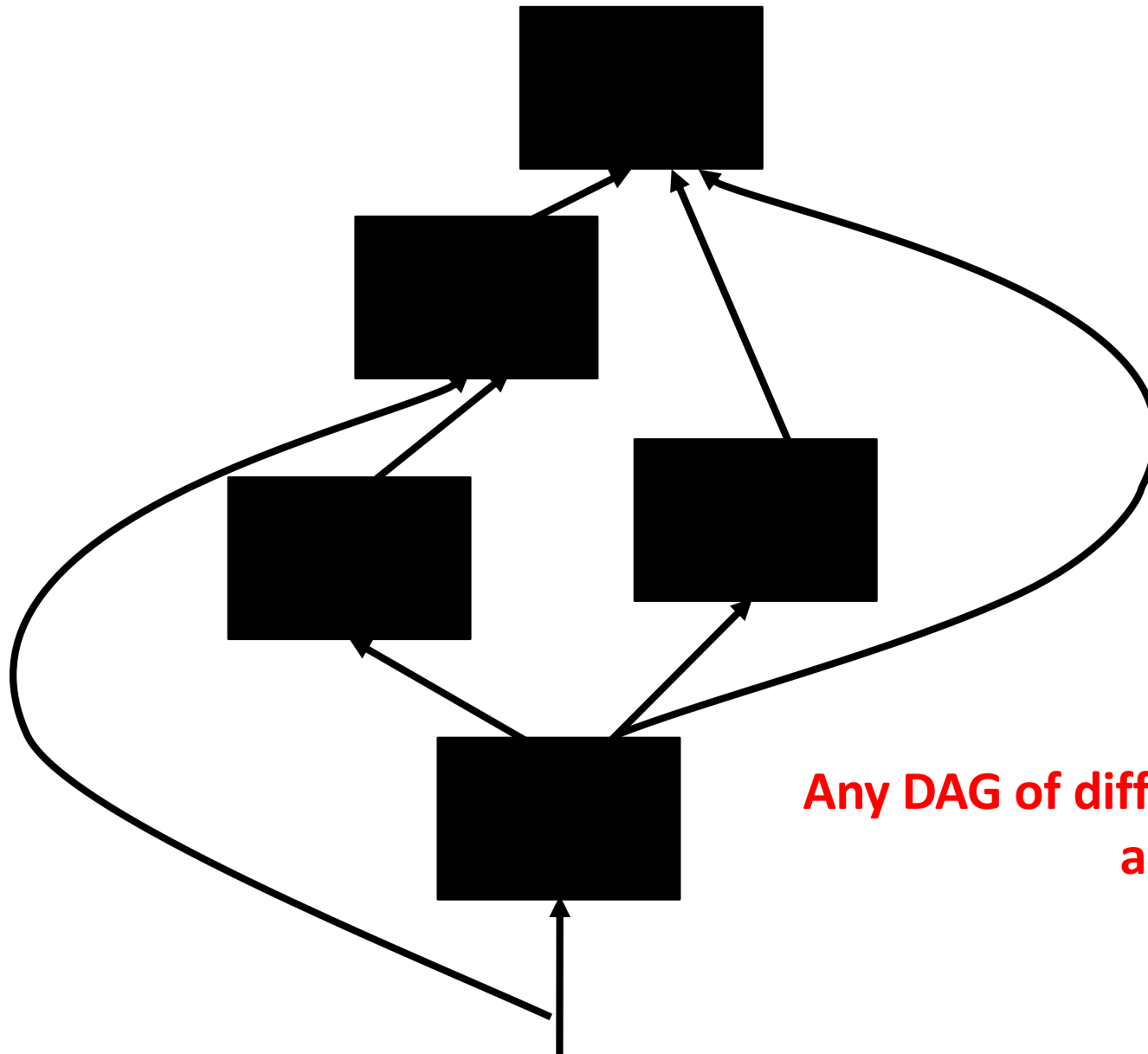
- Project Teams Google Doc
 - https://docs.google.com/spreadsheets/d/1BipWLvvWb7Fu6OSDd-uOCF1Lr_4drKOCRvdhxm_eSHc/edit#gid=0
 - Project Title
 - 1-3 sentence project summary TL;DR
 - Team member names

Recap from last time

How do we compute gradients?

- Analytic or “Manual” Differentiation
- Symbolic Differentiation
- Numerical Differentiation
- Automatic Differentiation
 - Forward mode AD
 - Reverse mode AD
 - aka “backprop”

Computational Graph



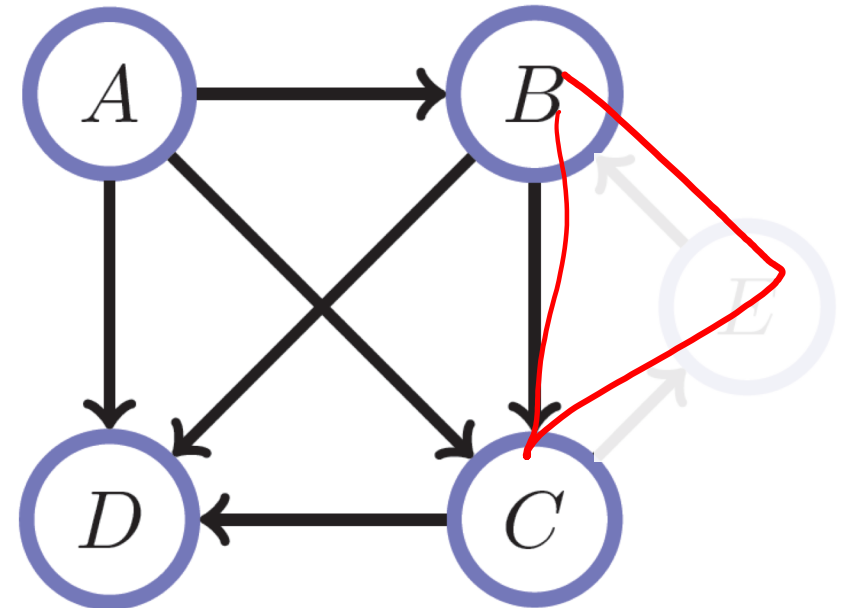
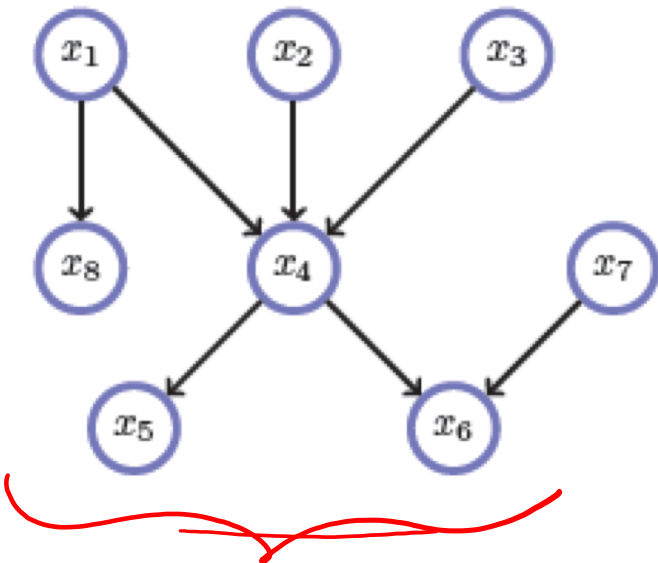
Any DAG of differentiable modules is allowed!

Directed Acyclic Graphs (DAGs)

- Exactly what the name suggests

$$G = (V, E)$$
$$E = \{ (v_i, v_j) \mid v_i, v_j \in V \}$$

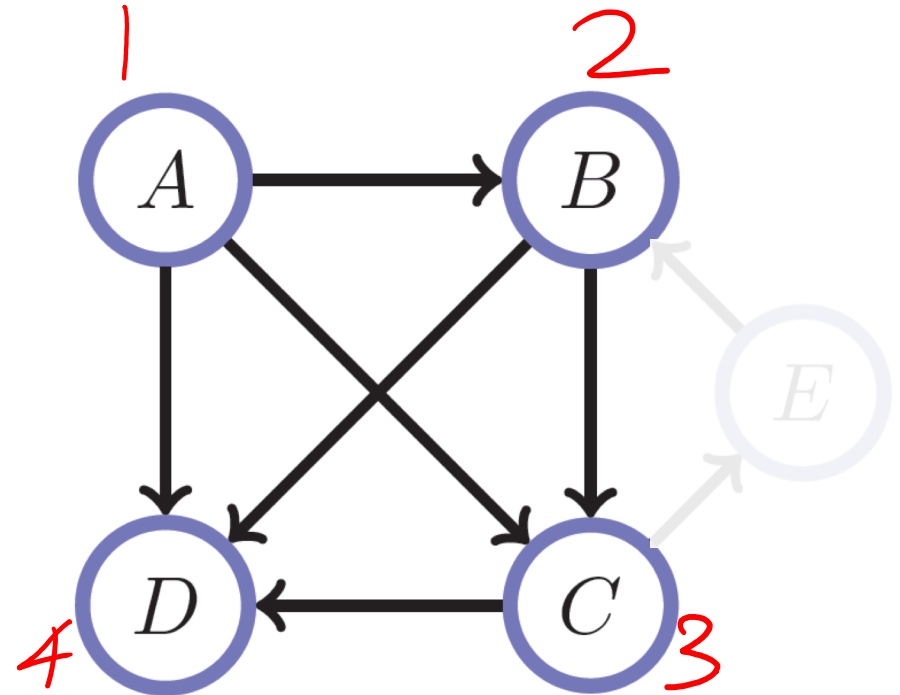
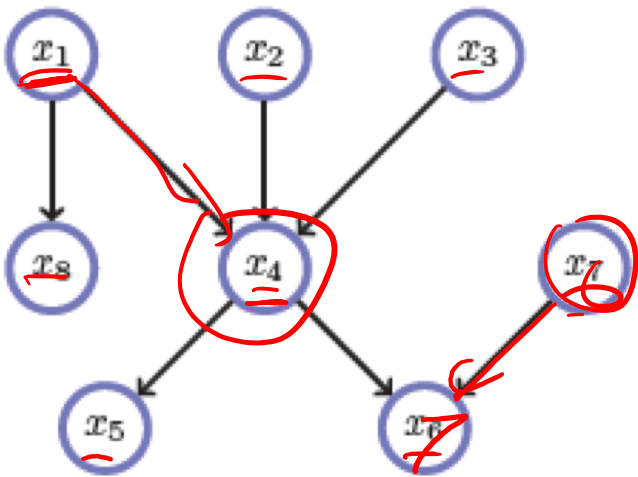
- Directed edges?
- No (directed) cycles
- Underlying undirected cycles okay



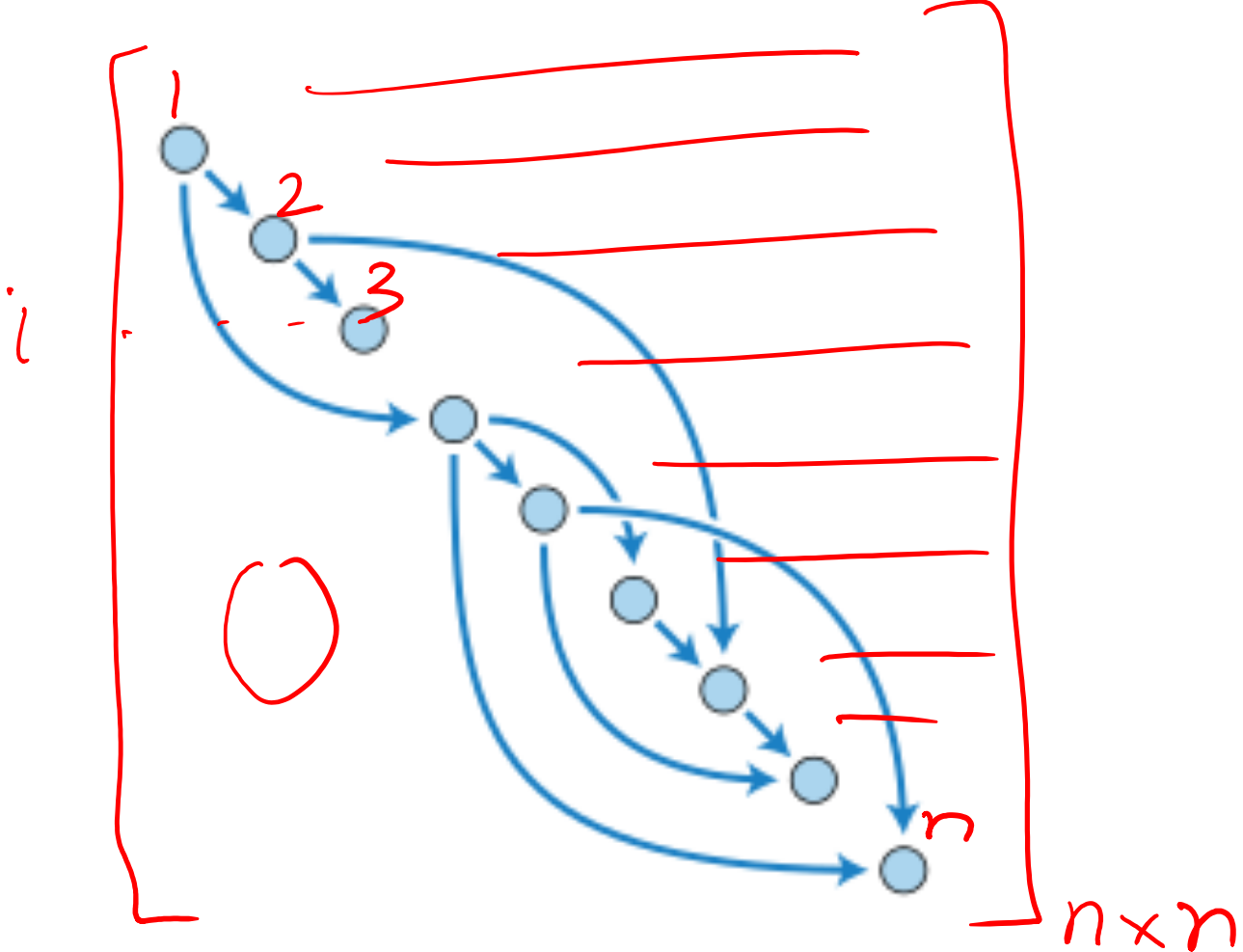
Directed Acyclic Graphs (DAGs)

- Concept
 - Topological Ordering

$$\exists \sigma: V \rightarrow [n] = \{1, \dots, n\}$$
$$\text{s.t. } \forall (v_i, v_j) \in E \quad \sigma(v_i) < \sigma(v_j)$$



Directed Acyclic Graphs (DAGs)

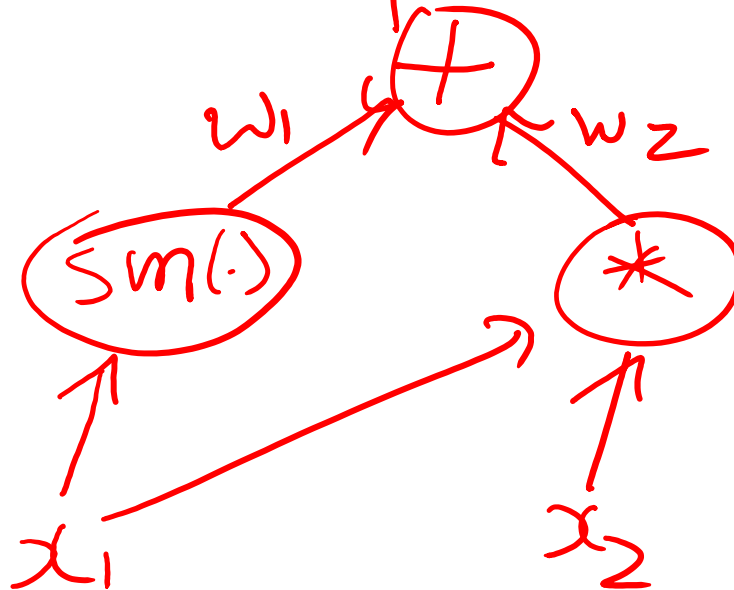


Computational Graphs

- Notation

$$f(x_1, x_2) = \underbrace{x_1 x_2}_{w_2} + \underbrace{\sin(x_1)}_{w_1}$$

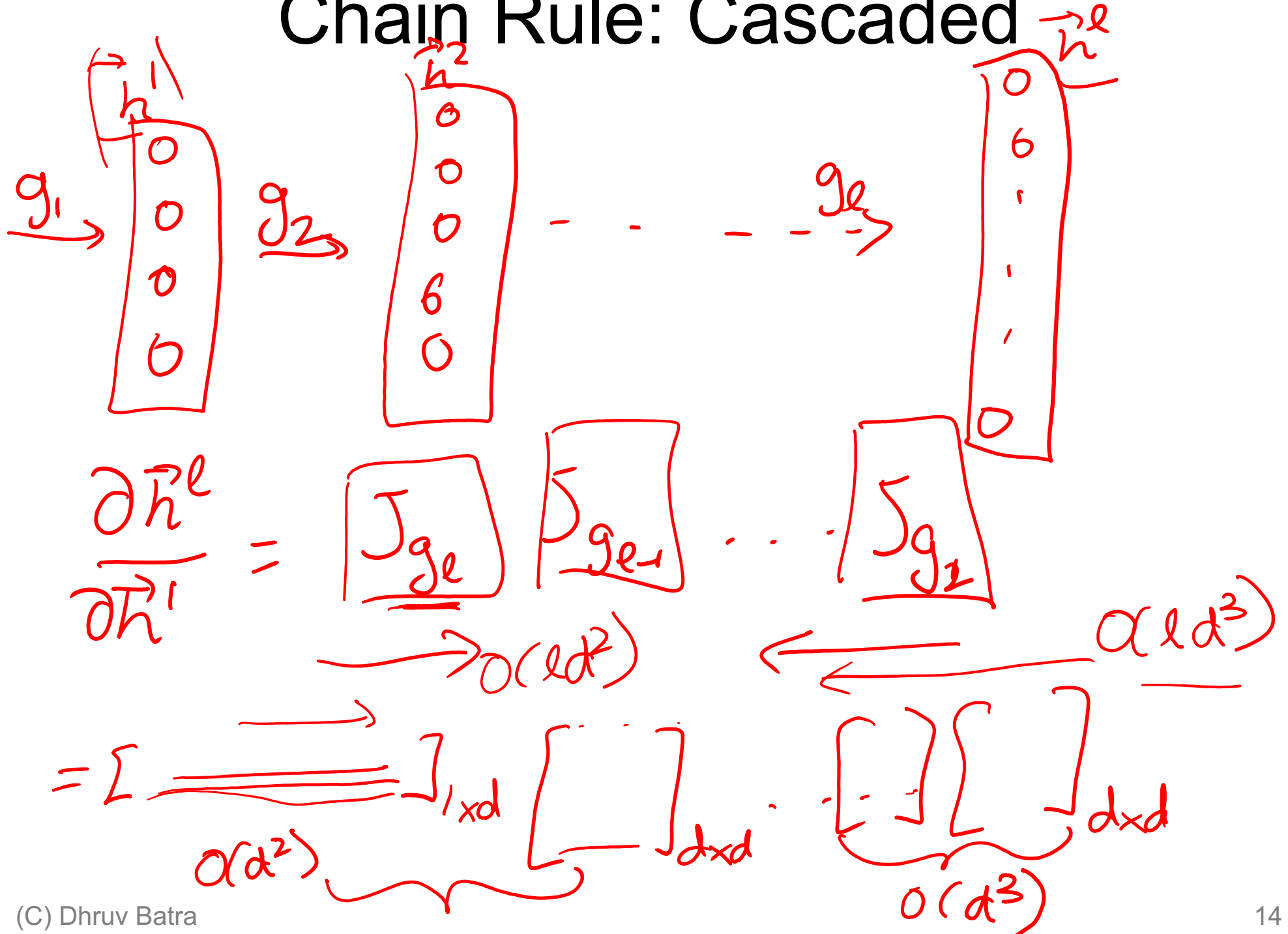
Handwritten notes: $f(x_1, x_2) = w_2 \uparrow$ (under $x_1 x_2$), w_1 (under $\sin(x_1)$)



How do we compute gradients?

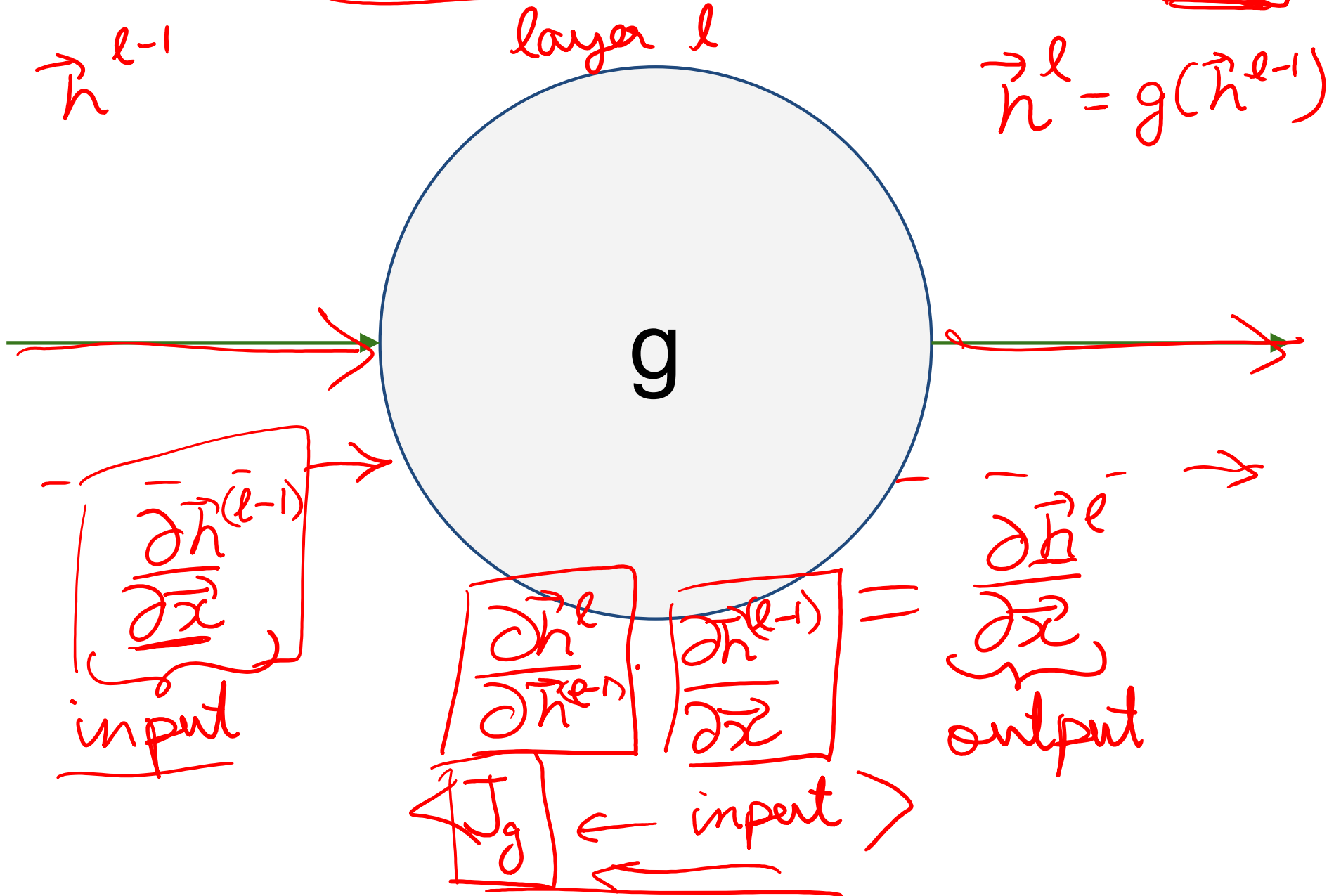
- Analytic or “Manual” Differentiation
- Symbolic Differentiation
- Numerical Differentiation
- Automatic Differentiation
 - Forward mode AD
 - Reverse mode AD
 - aka “backprop”

Chain Rule: Cascaded

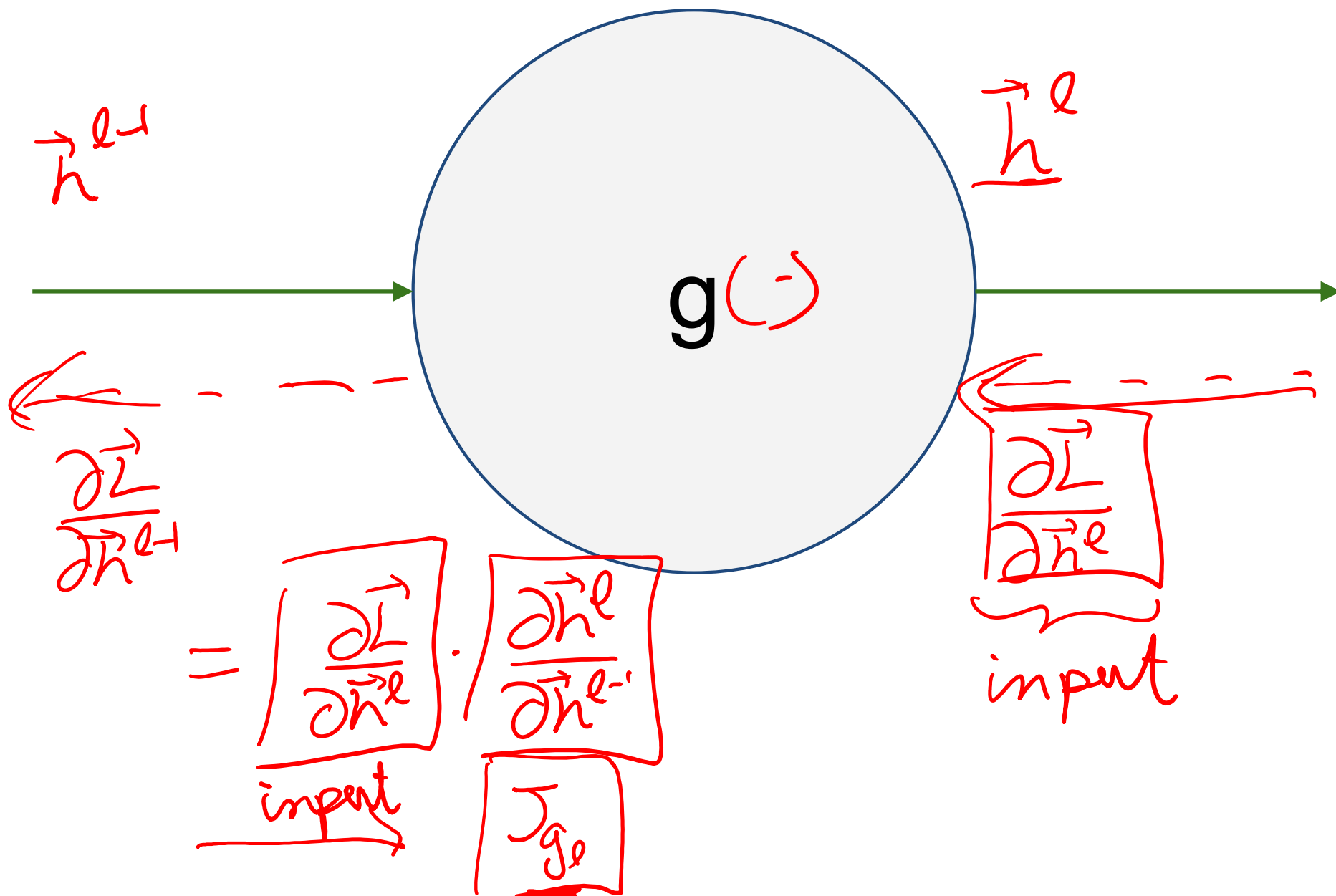


Forward mode AD

$$\frac{\partial \mathcal{L}}{\partial \vec{x}}$$

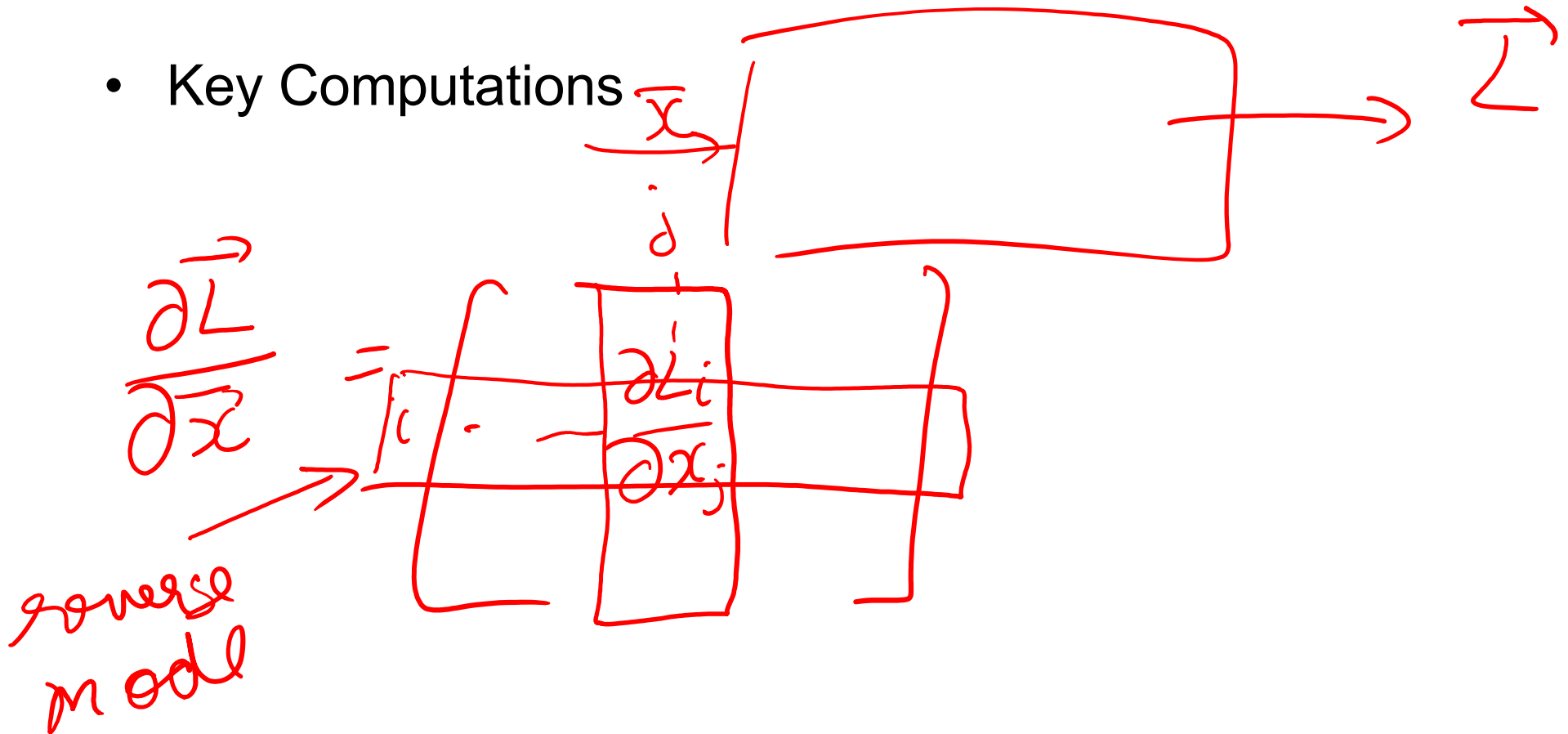


Reverse mode AD



Forward mode vs Reverse Mode

- Key Computations



Example: Forward mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

$$\frac{\partial f}{\partial \vec{x}}$$

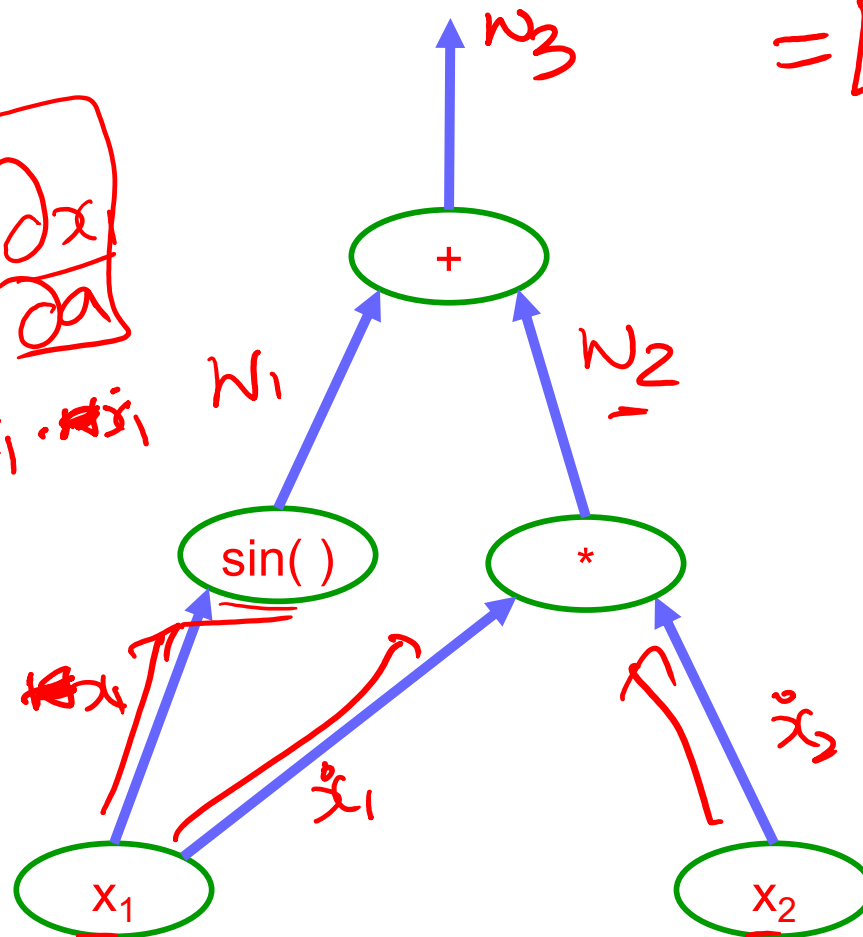
$$= \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right]$$

$$\frac{\partial f}{\partial a} = \dot{f}$$

$$\dot{x}_1 = \frac{\partial x_1}{\partial a}$$

$$\dot{w}_1 = \frac{\partial w_1}{\partial a} = \frac{\partial w_1}{\partial x_1} \cdot \frac{\partial x_1}{\partial a} = \cos(x_1) \cdot \dot{x}_1$$

$$\dot{x}_1 = \frac{\partial x_1}{\partial a}$$



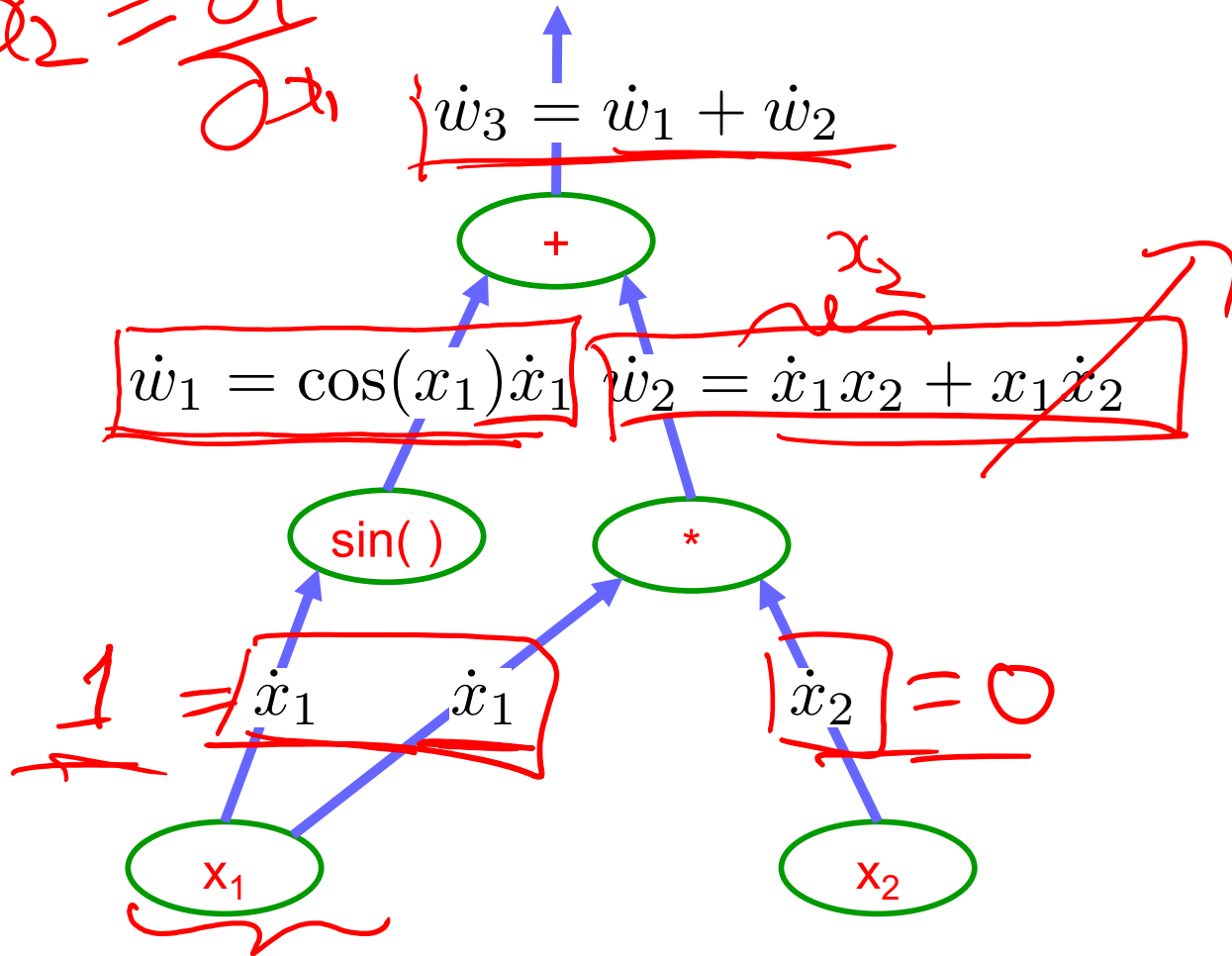
Example: Forward mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

$\frac{\partial f}{\partial a}$

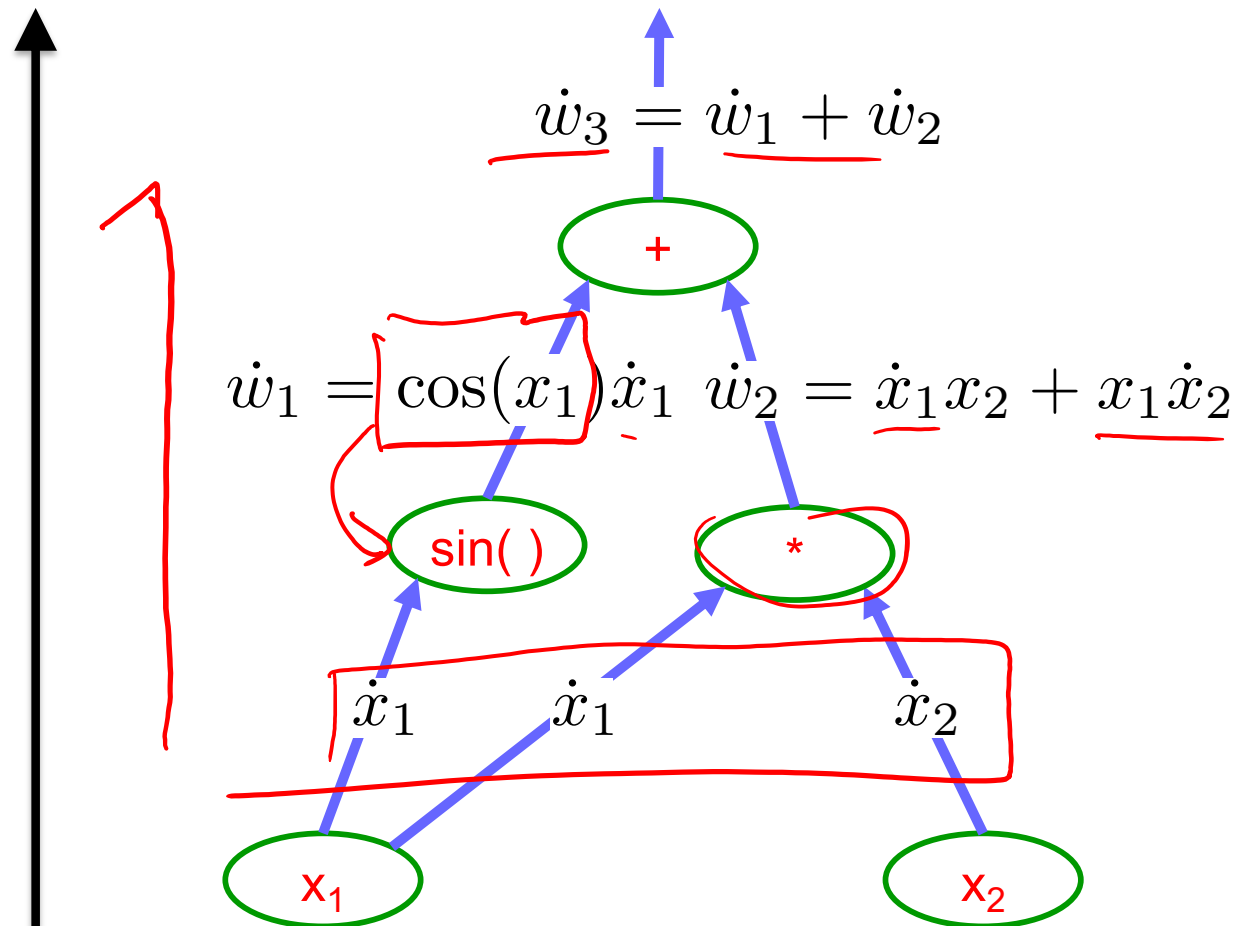
$\cos(x_1) + x_2 = \frac{\partial f}{\partial x_1}$

$\cos(x_1)$



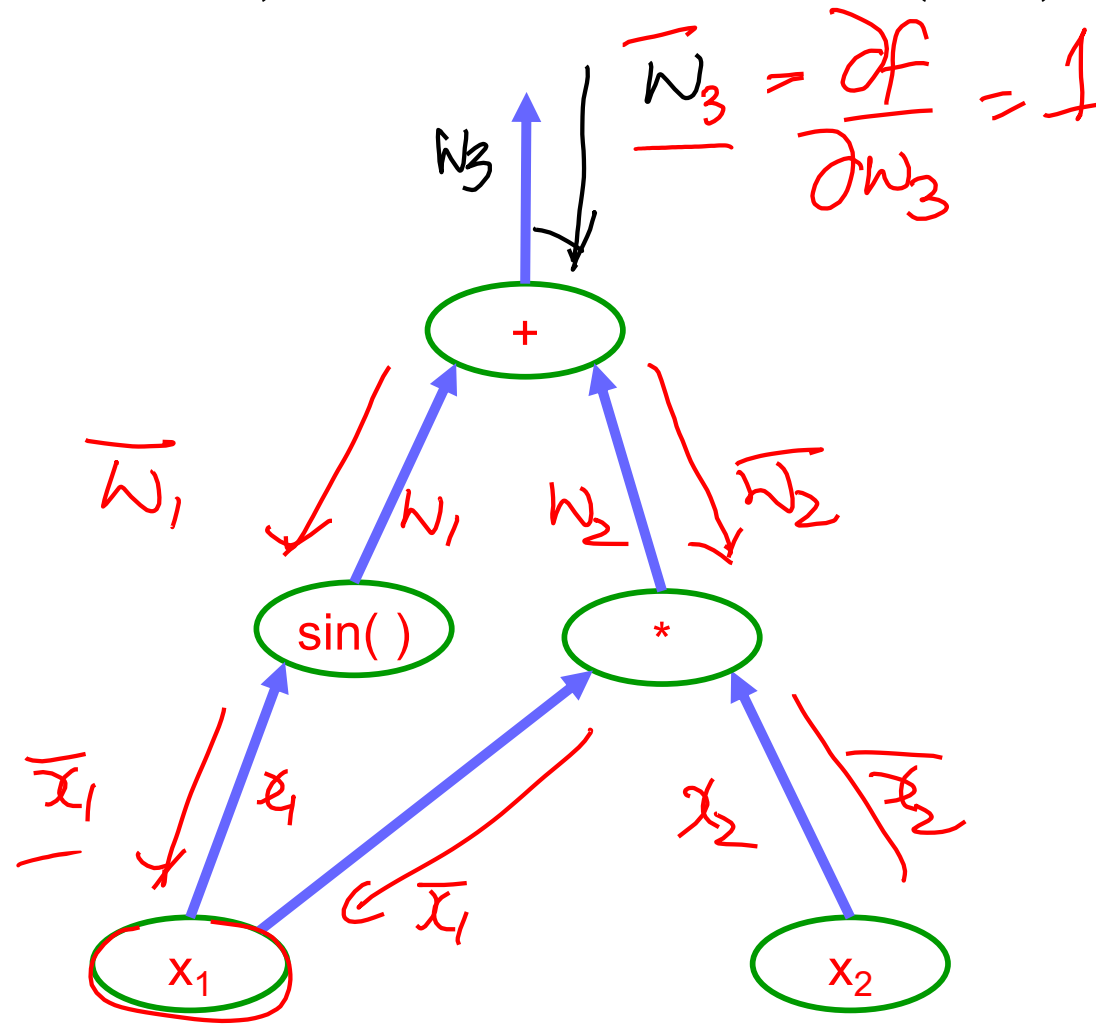
Example: Forward mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$



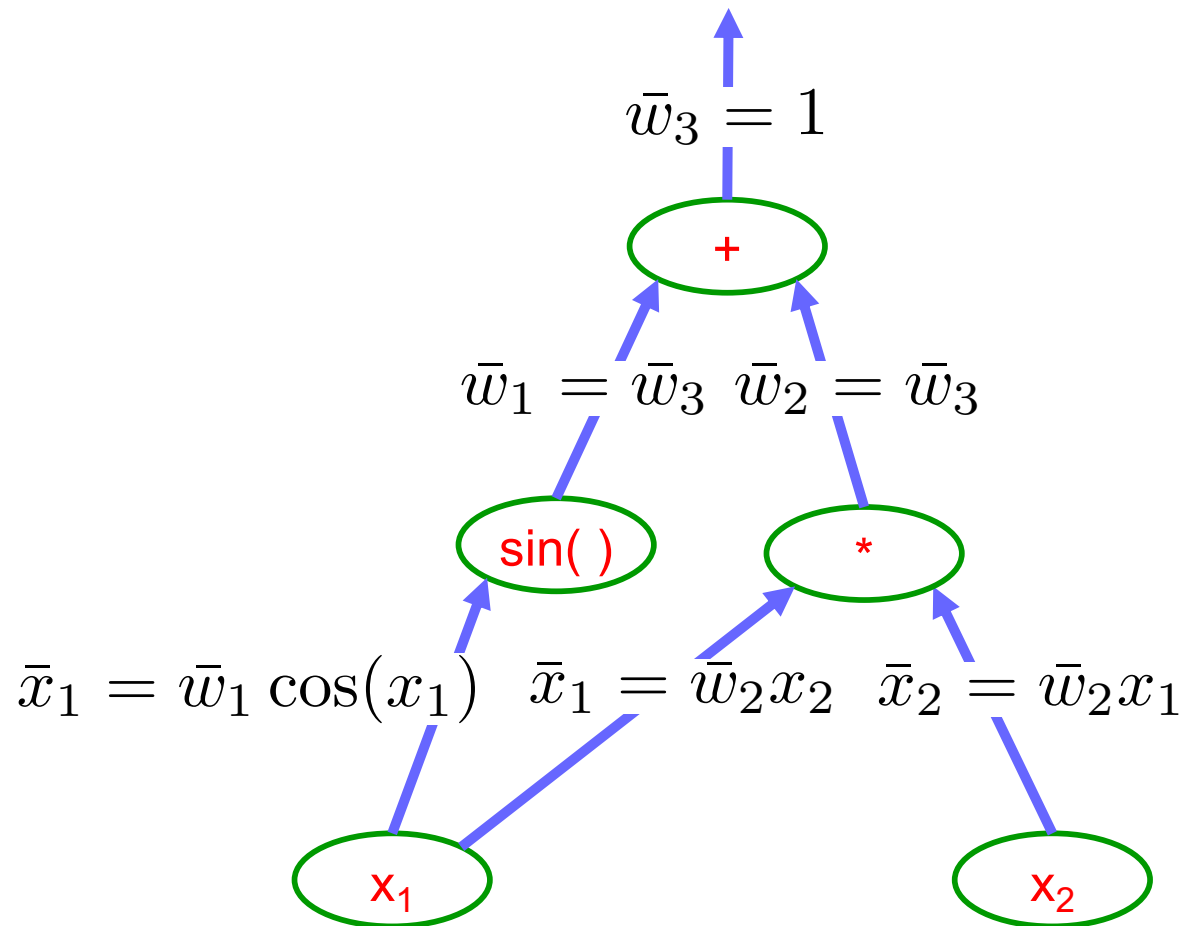
Example: Reverse mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

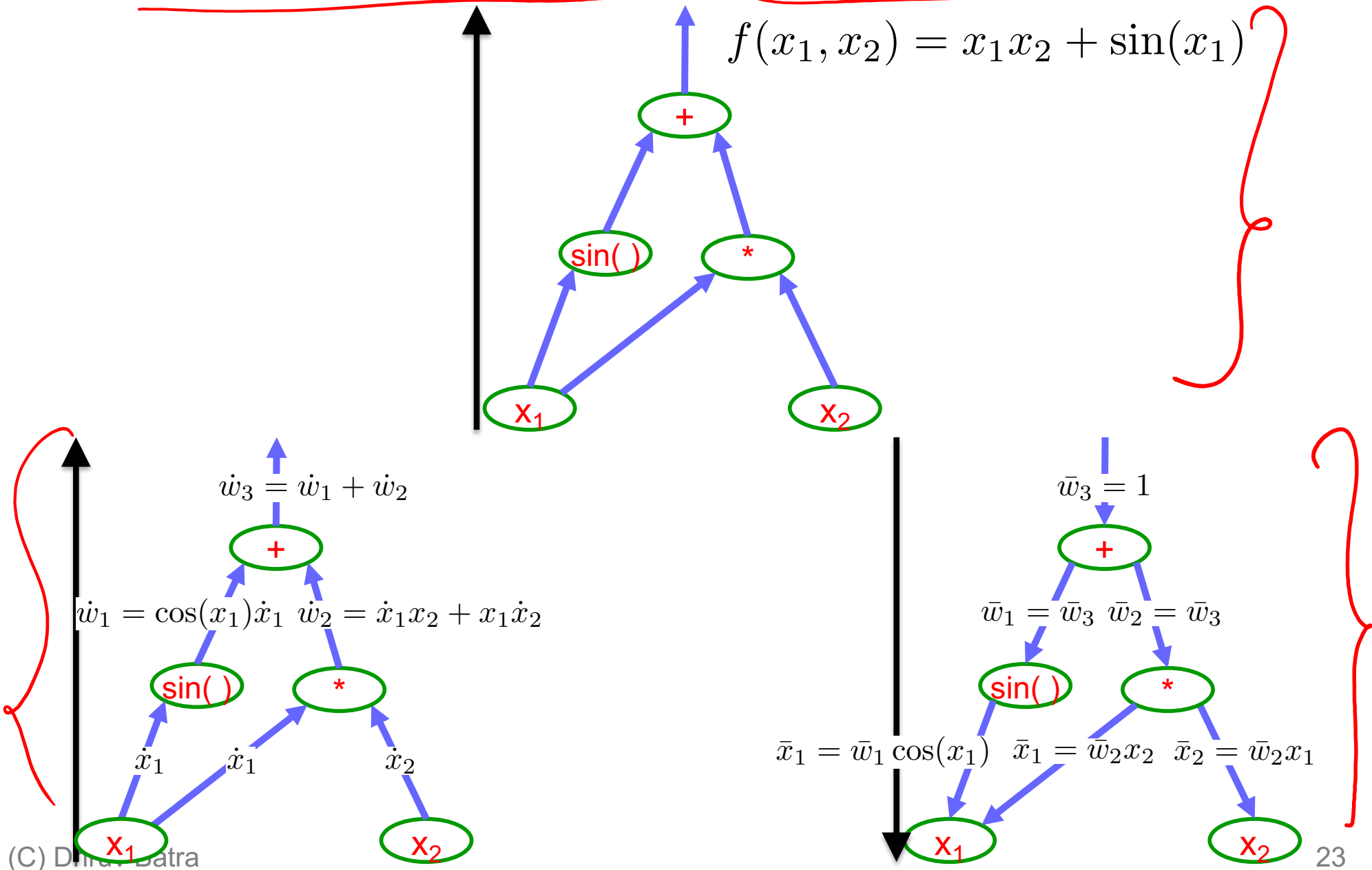


Example: Reverse mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$



Forward Pass vs Forward mode AD vs Reverse Mode AD



Forward mode vs Reverse Mode

- What are the differences?
- Which one is faster to compute?
 - Forward or backward?

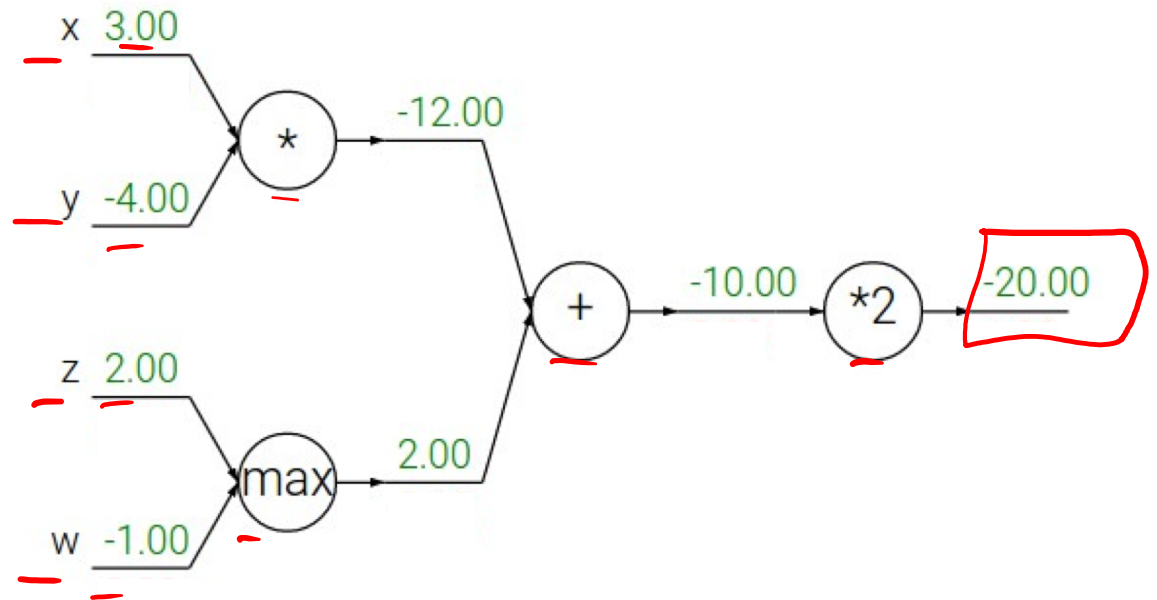
Forward mode vs Reverse Mode

- What are the differences?
- Which one is faster to compute?
 - Forward or backward?
- Which one is more memory efficient (less storage)?
 - Forward or backward?

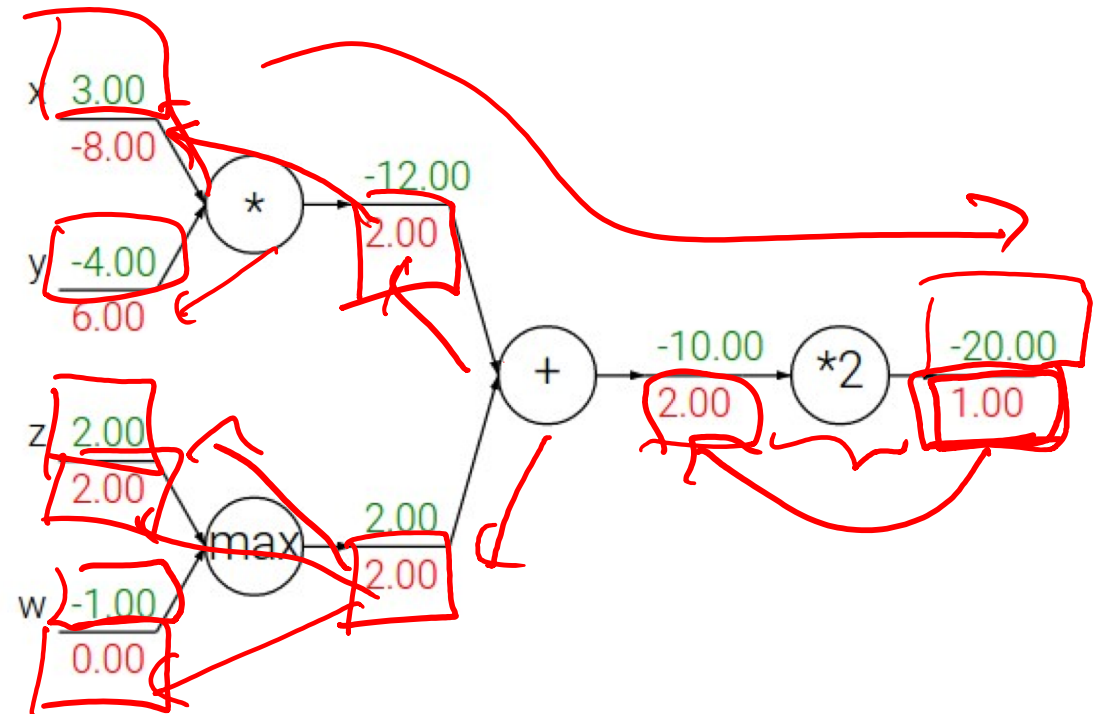
Patterns in backward flow

$f(\dots)$

$$= 2(xy + \max\{z, w\})$$

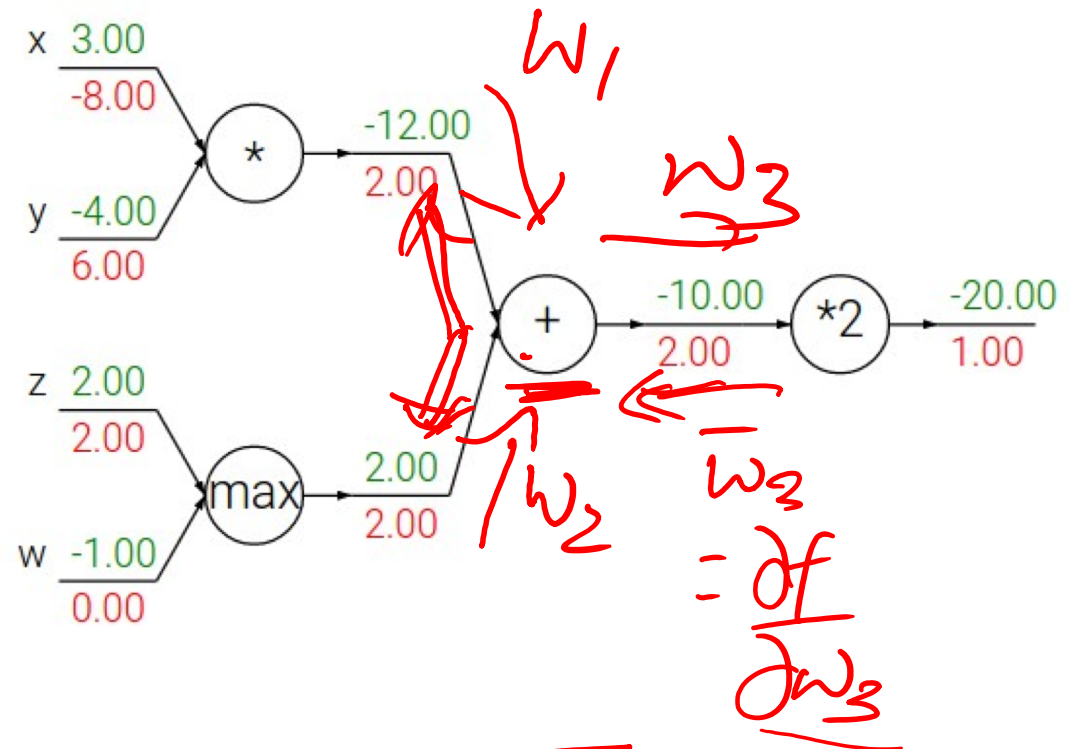


Patterns in backward flow



Patterns in backward flow

add gate gradient distributor



$$w_3 = w_1 + w_2$$

$$\frac{\partial w_3}{\partial w_1} = 1$$

$$\frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial w_3} \cdot \frac{\partial w_3}{\partial w_1}$$

Patterns in backward flow

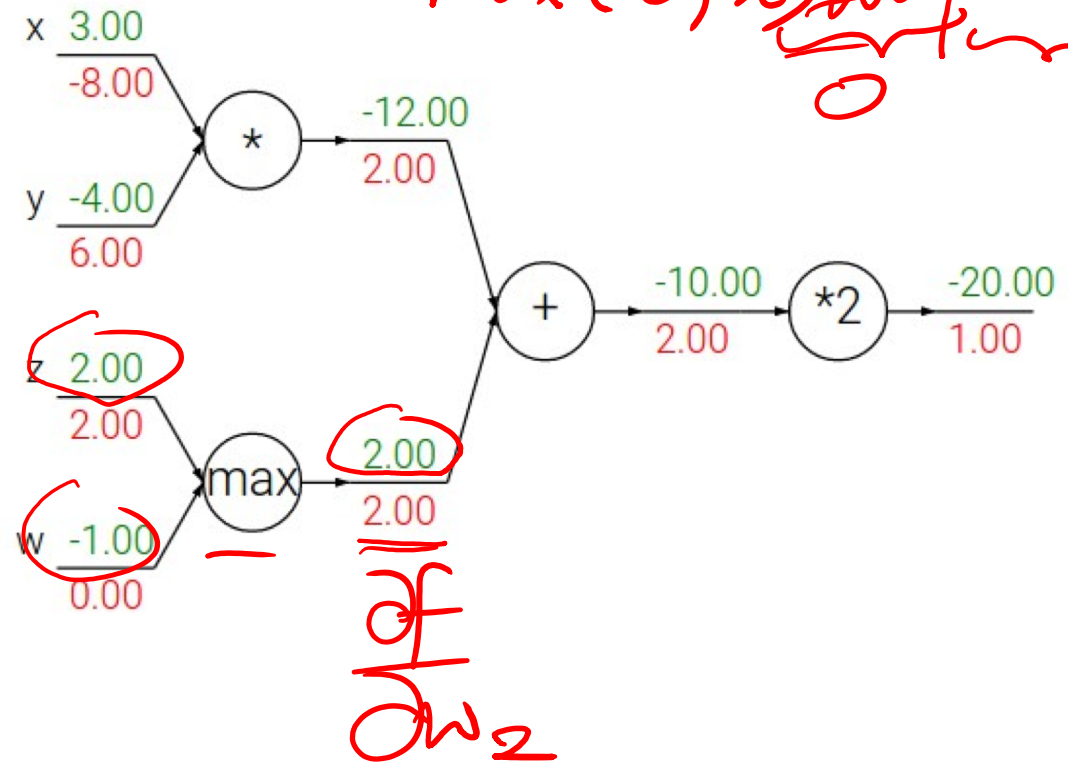
$$\underline{w}_2 = \max\{\underline{z}, \underline{w}\}$$

$$\max(0, x) \begin{cases} +1 \\ 0 \end{cases}$$

add gate: gradient distributor

Q: What is a max gate?

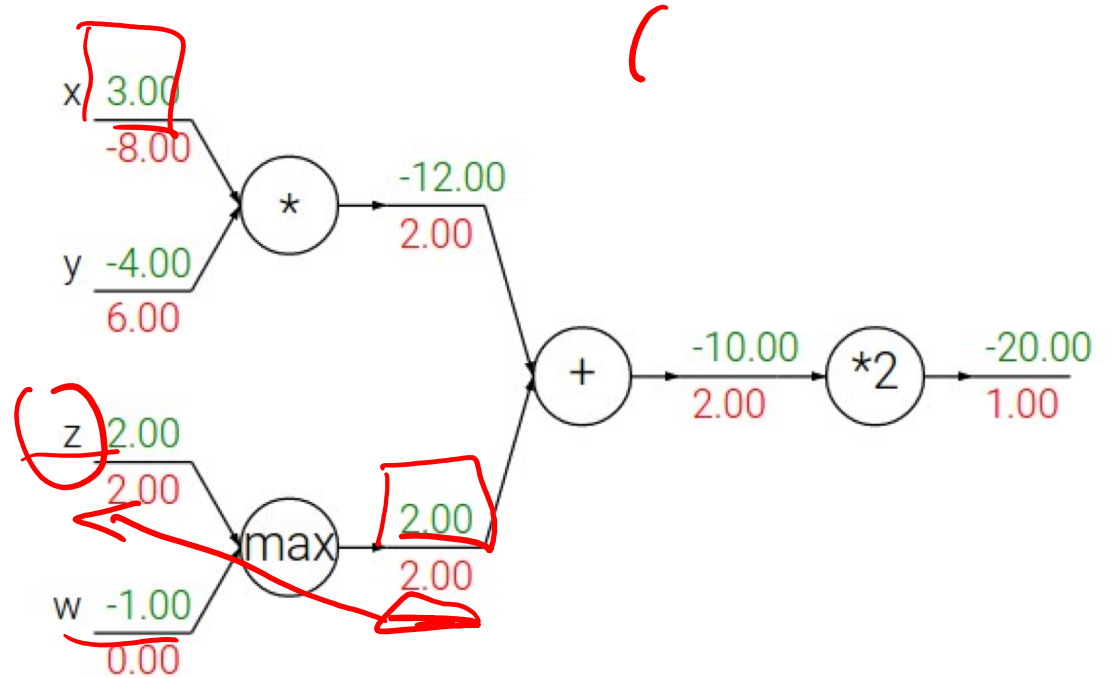
$\frac{df}{dz}$
 $\frac{df}{dw}$



Patterns in backward flow

add gate: gradient distributor

max gate: gradient router

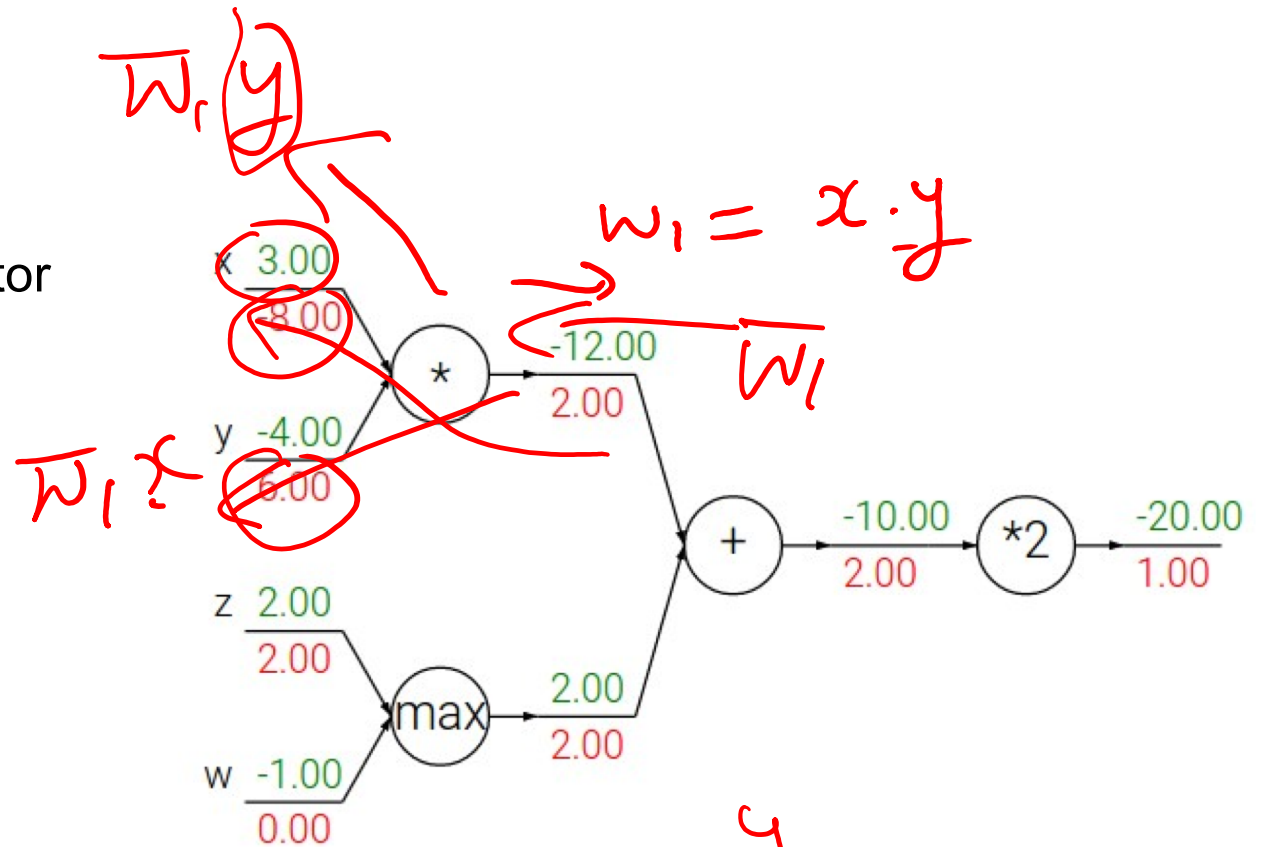


Patterns in backward flow

add gate: gradient distributor

max gate: gradient router

Q: What is a mul gate?



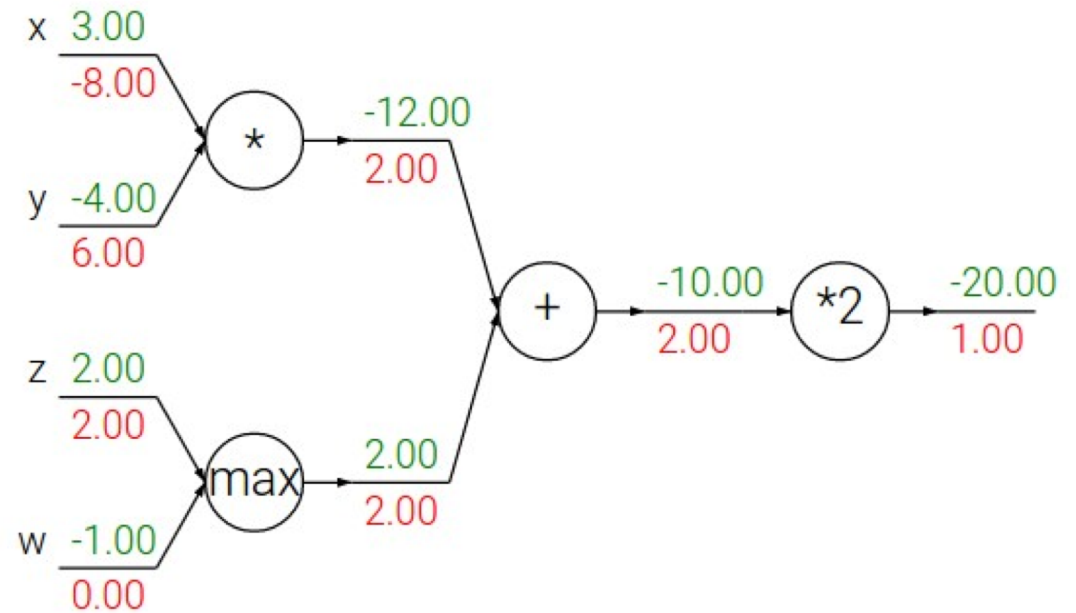
$$\frac{\partial f}{\partial x} = \left[\frac{\partial f}{\partial w_1} \right] \left[\frac{\partial w_1}{\partial x} \right]$$

Patterns in backward flow

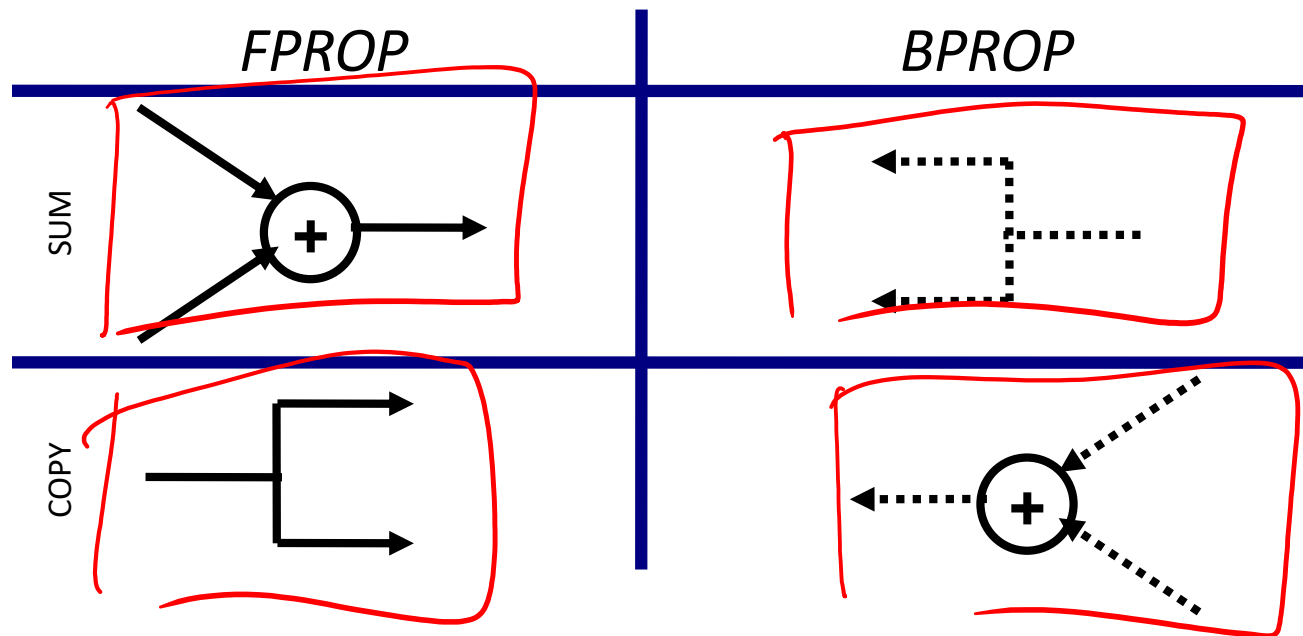
add gate: gradient distributor

max gate: gradient router

mul gate: gradient switcher



Duality in Fprop and Bprop



Plan for Today

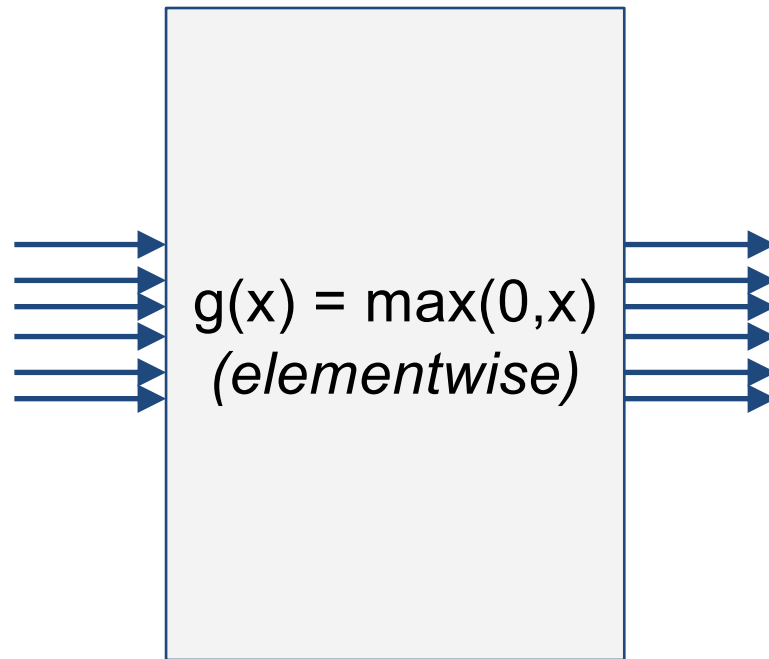
- (Finish) Computing Gradients
 - Backprop in FC+ReLU NNs

- Convolutional Neural Networks

Jacobian of ReLU

$$\vec{h}^{l-1} \in \mathbb{R}^{4096}$$

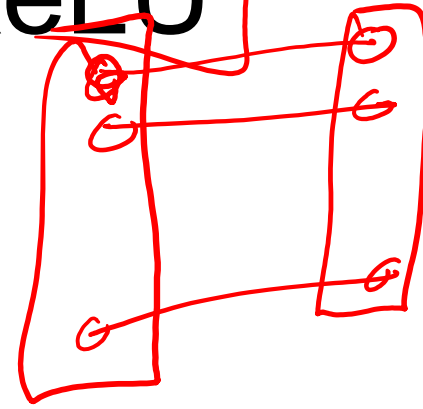
4096-d
input vector



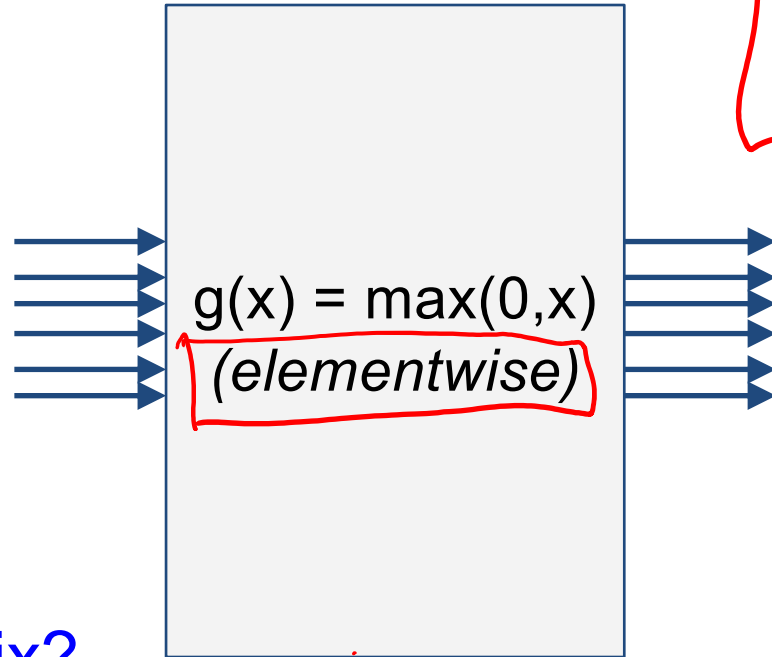
$$\vec{h}^l \in \mathbb{R}^{4096}$$

4096-d
output vector

Jacobian of ReLU



4096-d
input vector



4096-d
output vector

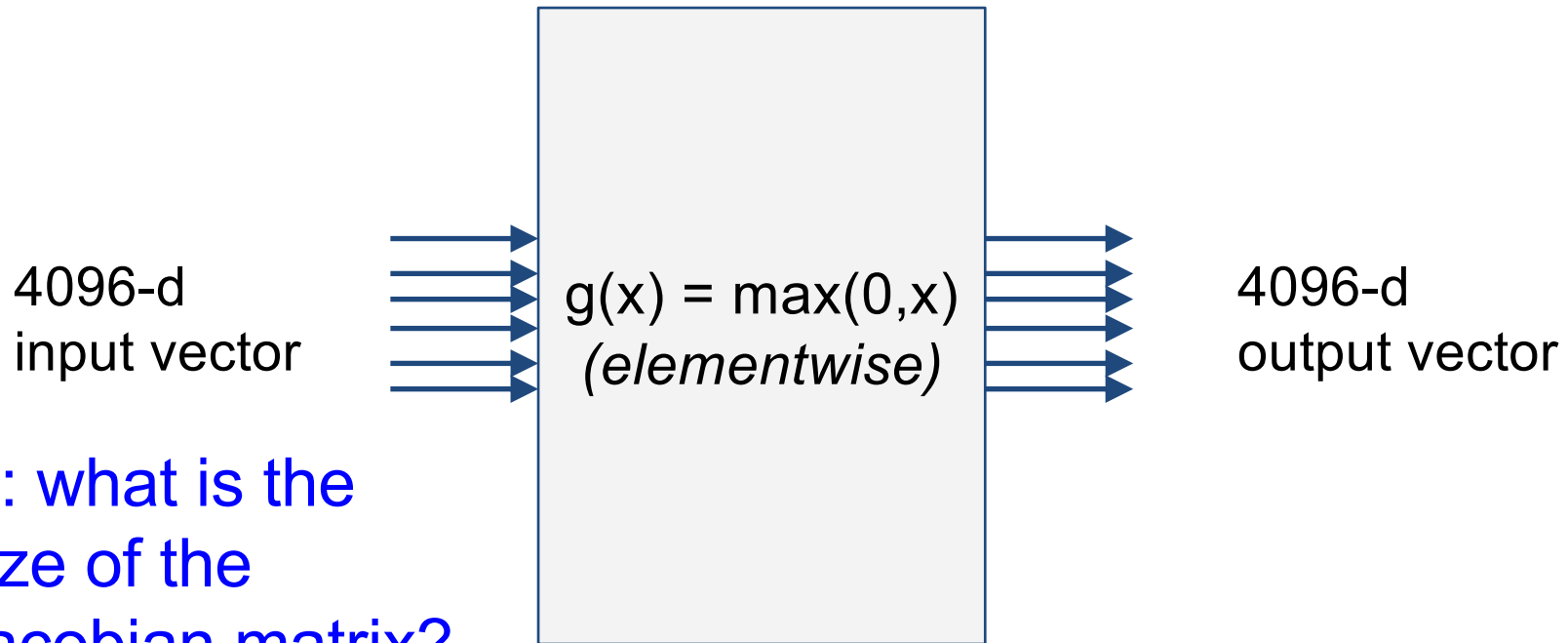
Q: what is the
size of the
Jacobian matrix?

$$\frac{\partial h_i^l}{\partial h_j^{l-1}} = 0 \quad \forall i \neq j$$

$$\frac{\partial \vec{h}^l}{\partial \vec{h}^{l-1}} = \begin{bmatrix} \frac{\partial h_1^l}{\partial h_1^{l-1}} & & \\ & \ddots & \\ & & \frac{\partial h_n^l}{\partial h_n^{l-1}} \end{bmatrix} \rightarrow \begin{bmatrix} \circ & & \\ & \circ & \\ & & \circ \end{bmatrix}$$

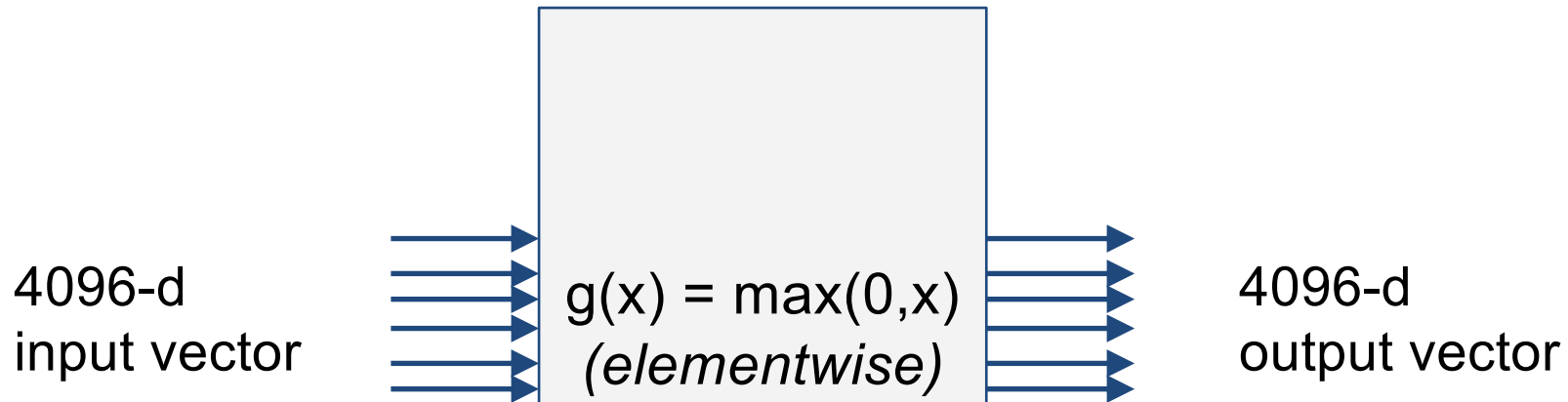
4096 x 4096

Jacobian of ReLU



Q: what is the
size of the
Jacobian matrix?
[4096 x 4096!]

Jacobian of ReLU



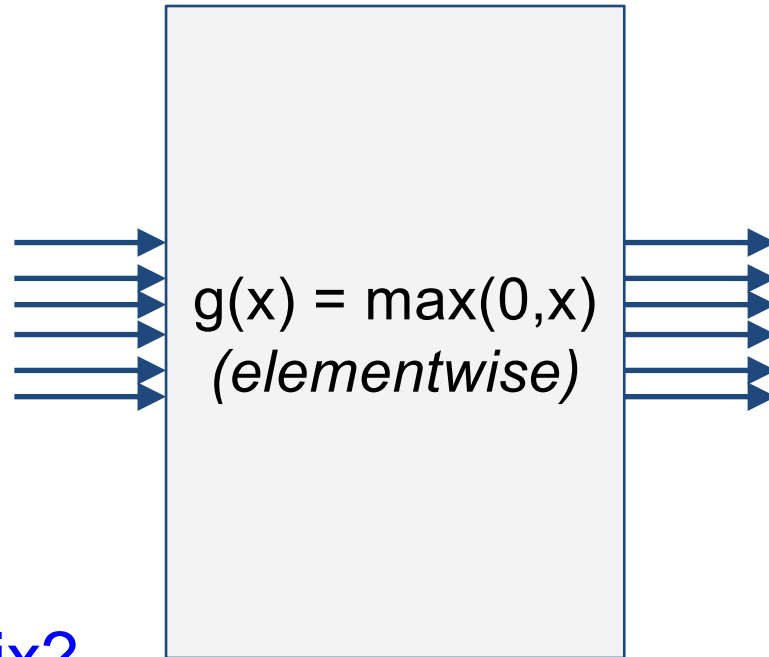
Q: what is the
size of the
Jacobian matrix?
[4096 x 4096!]

in practice we process an
entire minibatch (e.g. 100)
of examples at one time:

i.e. Jacobian would technically be a
[409,600 x 409,600] matrix :\
}

Jacobian of ReLU

4096-d
input vector

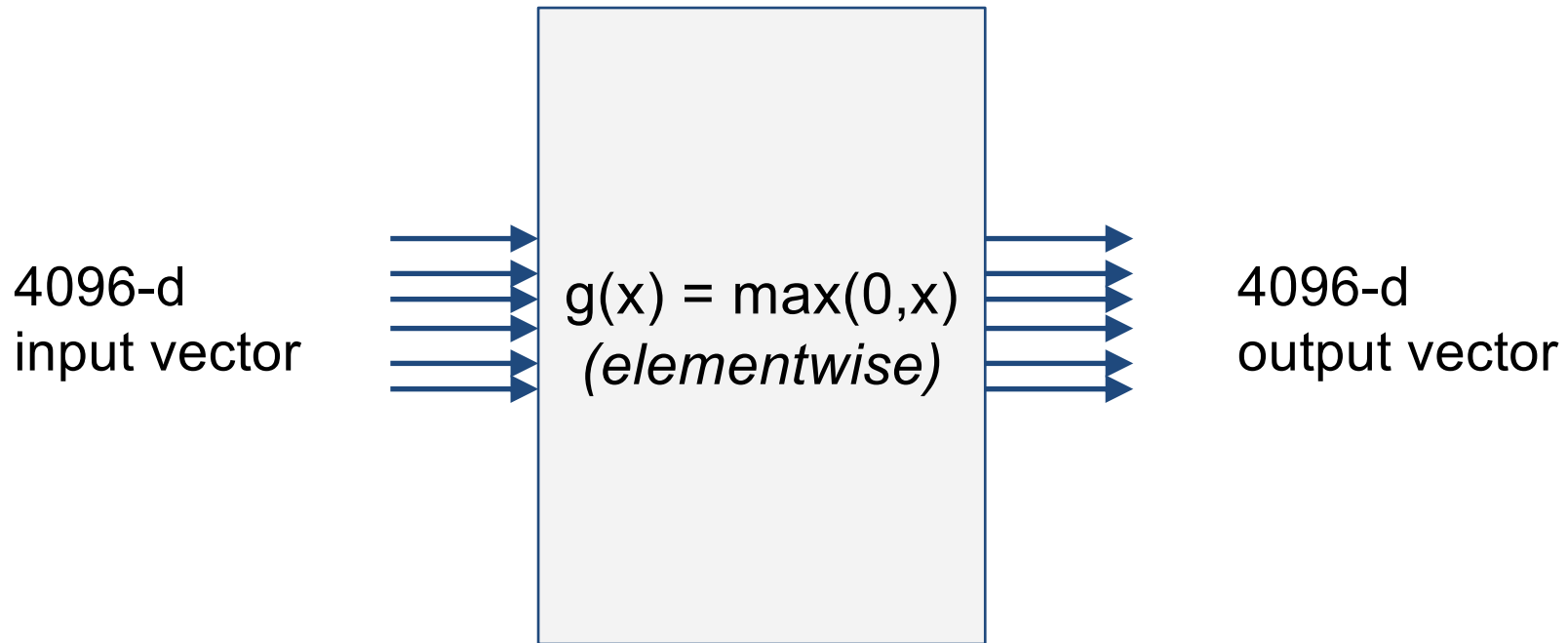


4096-d
output vector

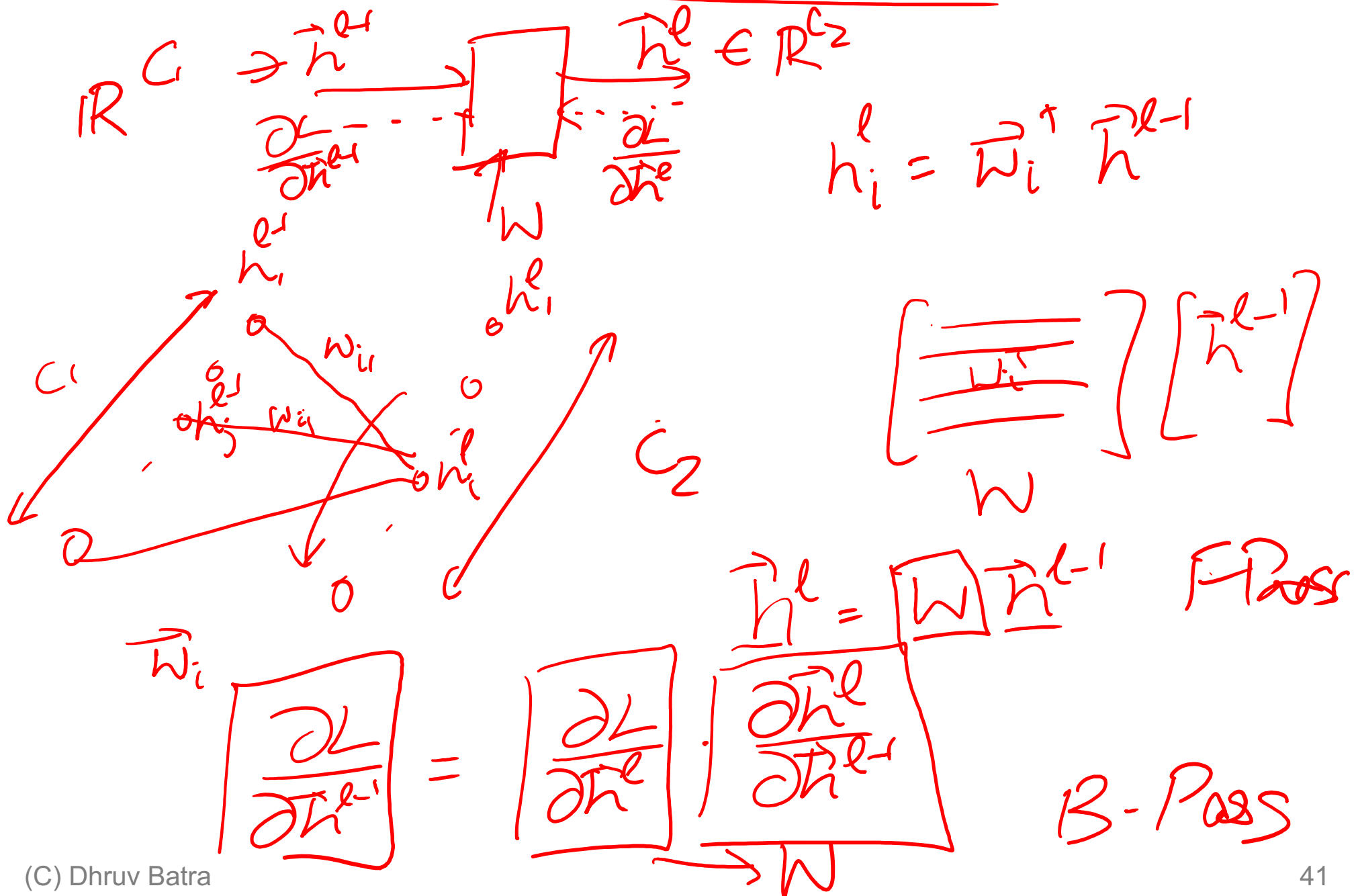
Q: what is the
size of the
Jacobian matrix?
[4096 x 4096!]

Q2: what does it
look like?

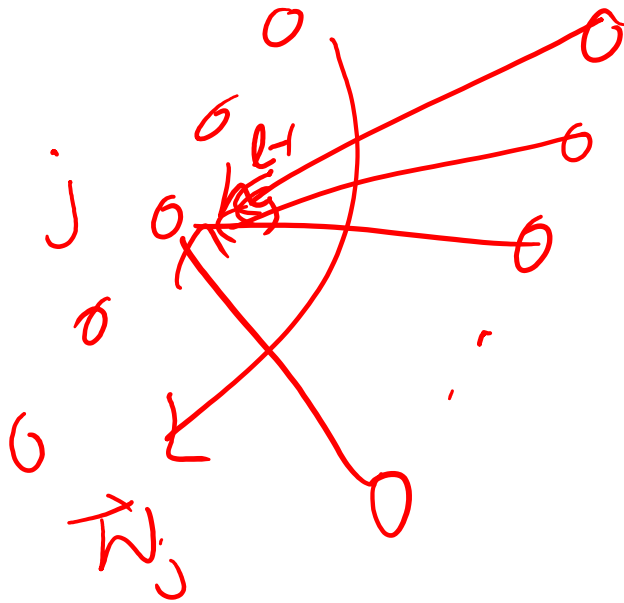
Jacobian of ReLU



Jacobians of FC-Layer



Jacobians of FC-Layer



Jacobians of FC-Layer

Jacobians of FC-Layer



Convolutional Neural Networks

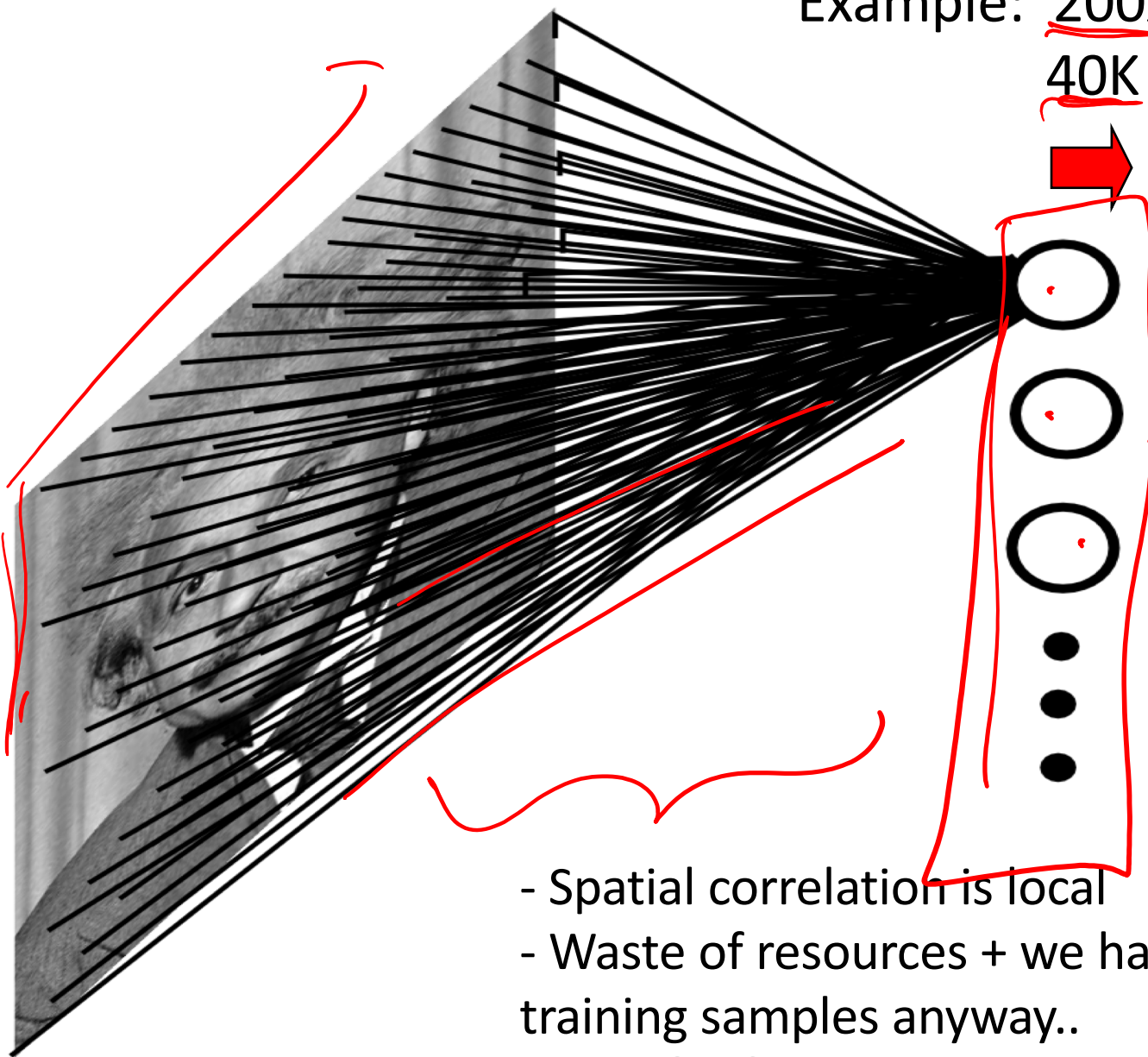
(without the brain stuff)

Fully Connected Layer

Example: 200x200 image
40K hidden units

4×10^4

~2B parameters!!!



4×10^4

4×10^4

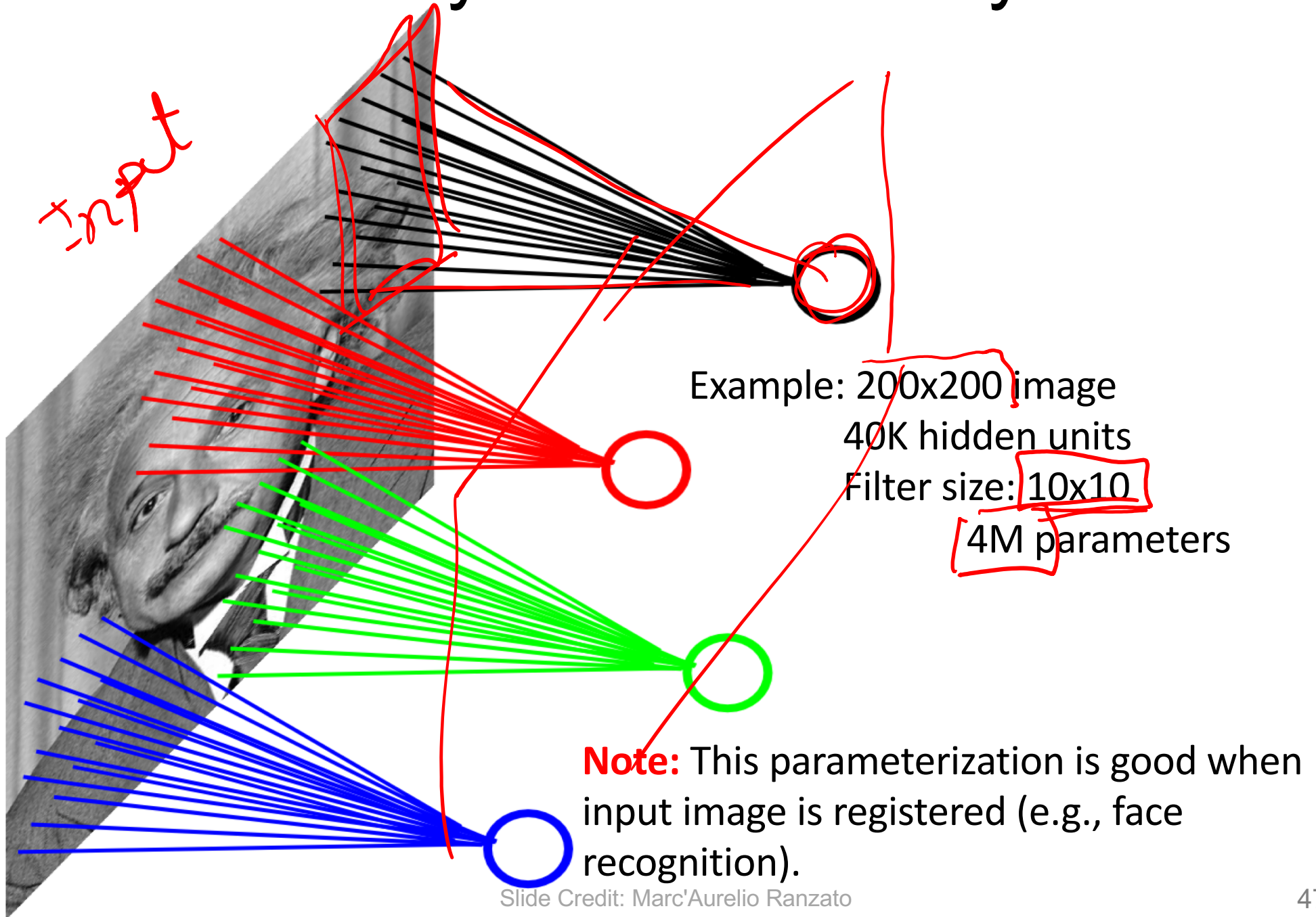
4×10^4

w

1.6×10^9

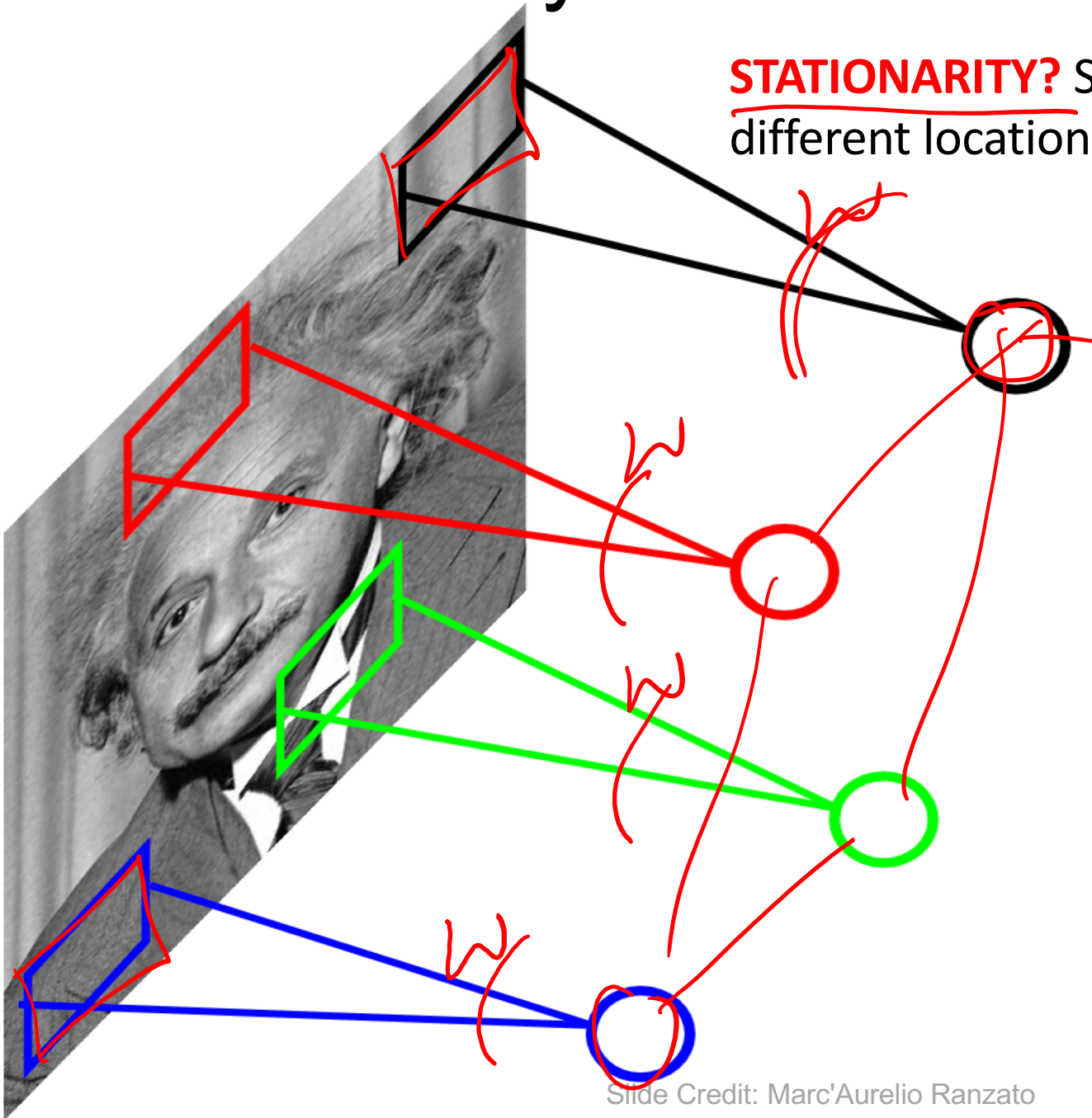
- Spatial correlation is local
- Waste of resources + we have not enough training samples anyway..

Locally Connected Layer

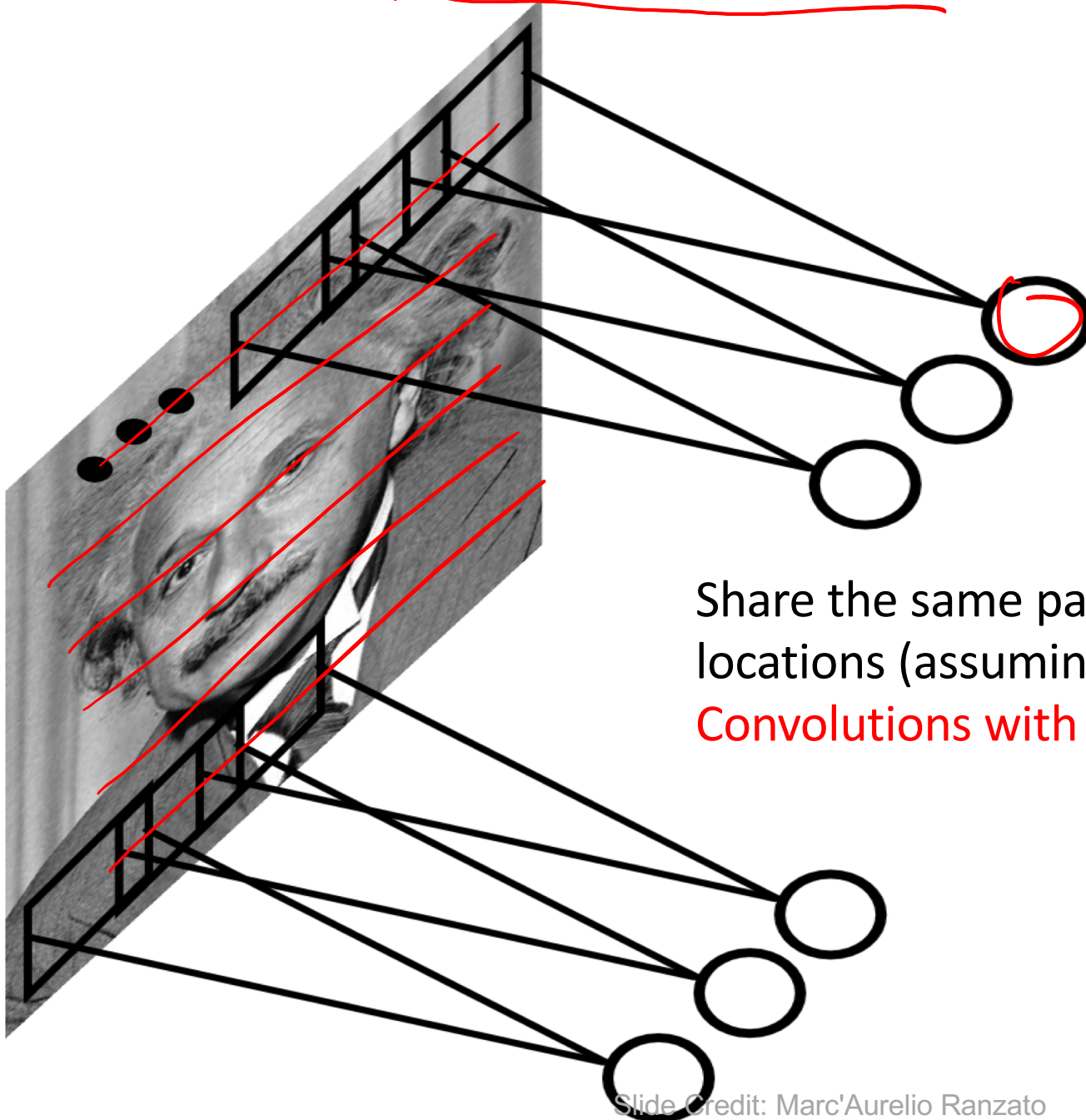


Locally Connected Layer

STATIONARITY? Statistics is similar at different locations



Convolutional Layer



Share the same parameters across different locations (assuming input is stationary):

Convolutions with learned kernels

Convolutions!

math \rightarrow CS \rightarrow programming

Convolutions for mathematicians



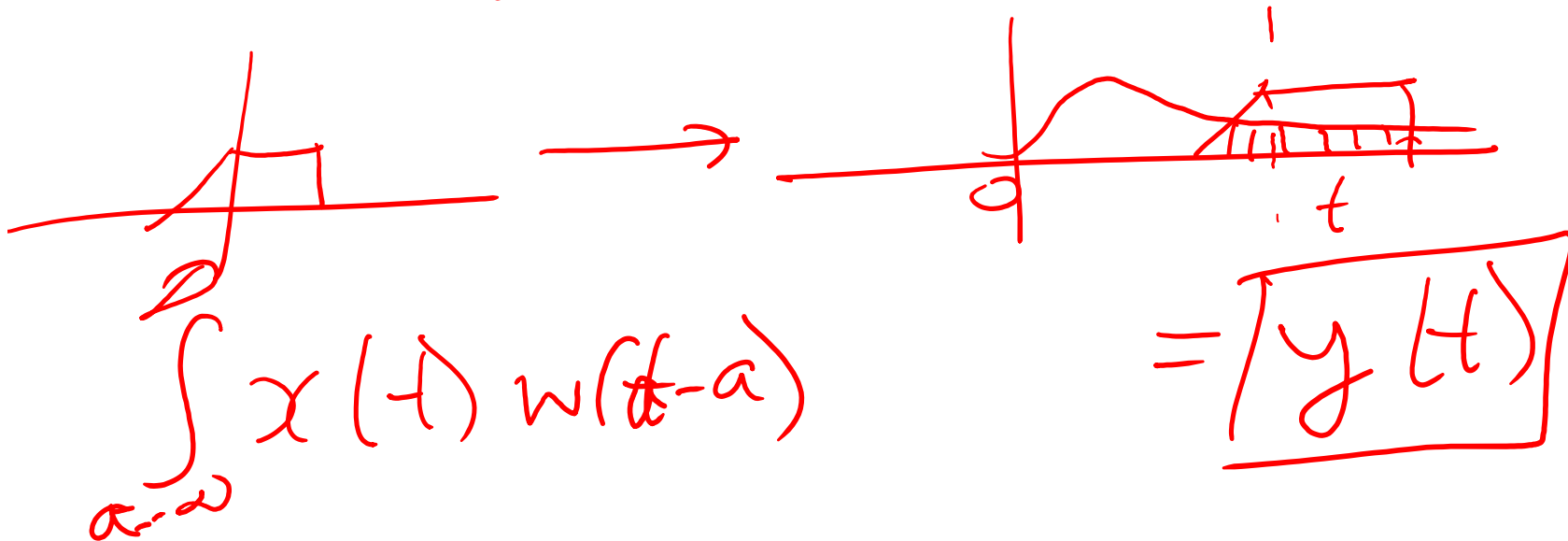
$$\begin{aligned}y(t) &= (x * w)(t) \\ &= (w * x)(t) \\ &= \int_{-\infty}^{\infty} \underline{x(t-a) w(a)} da \\ &= \int_{-\infty}^{\infty} x(a) \underline{w(t-a)} da\end{aligned}$$

Convolutions for mathematicians

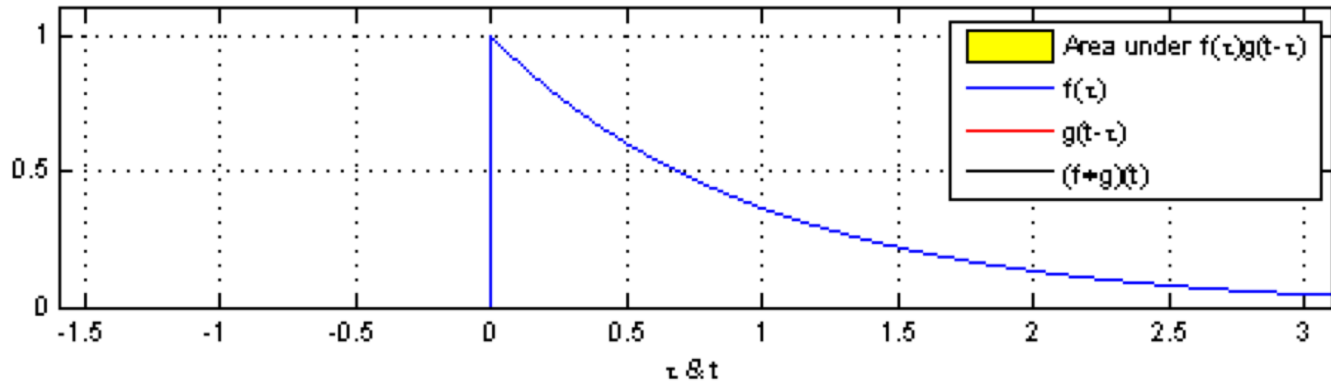
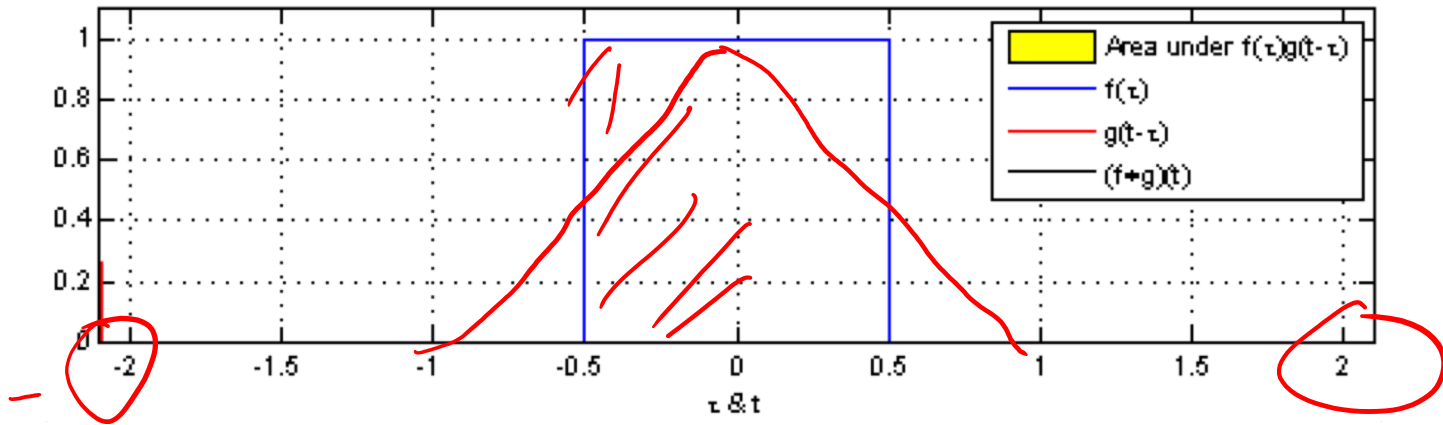
$$w(a) \rightarrow w(-a)$$



$$w(-a) \xrightarrow{\text{-kernel}} w(-(a-t))$$



$$y(t_1, t_2) = \int_{a=-\infty}^{\infty} \int_{b=-\infty}^{\infty} x(t_1 - a, t_2 - b) w(a, b) da db$$



"Convolution of box signal with itself2" by Convolution_of_box_signal_with_itself.gif: Brian Amberg derivative work: Tinos (talk) - Convolution_of_box_signal_with_itself.gif. Licensed under CC BY-SA 3.0 via Commons - https://commons.wikimedia.org/wiki/File:Convolution_of_box_signal_with_itself2.gif#/media/File:Convolution_of_box_signal_with_itself2.gif

Convolutions for computer scientists

$$y[x, c] = \sum_{a=-\infty}^{\infty} \sum_{b=-\infty}^{\infty} x[x-a, c-b] w[a, b]$$

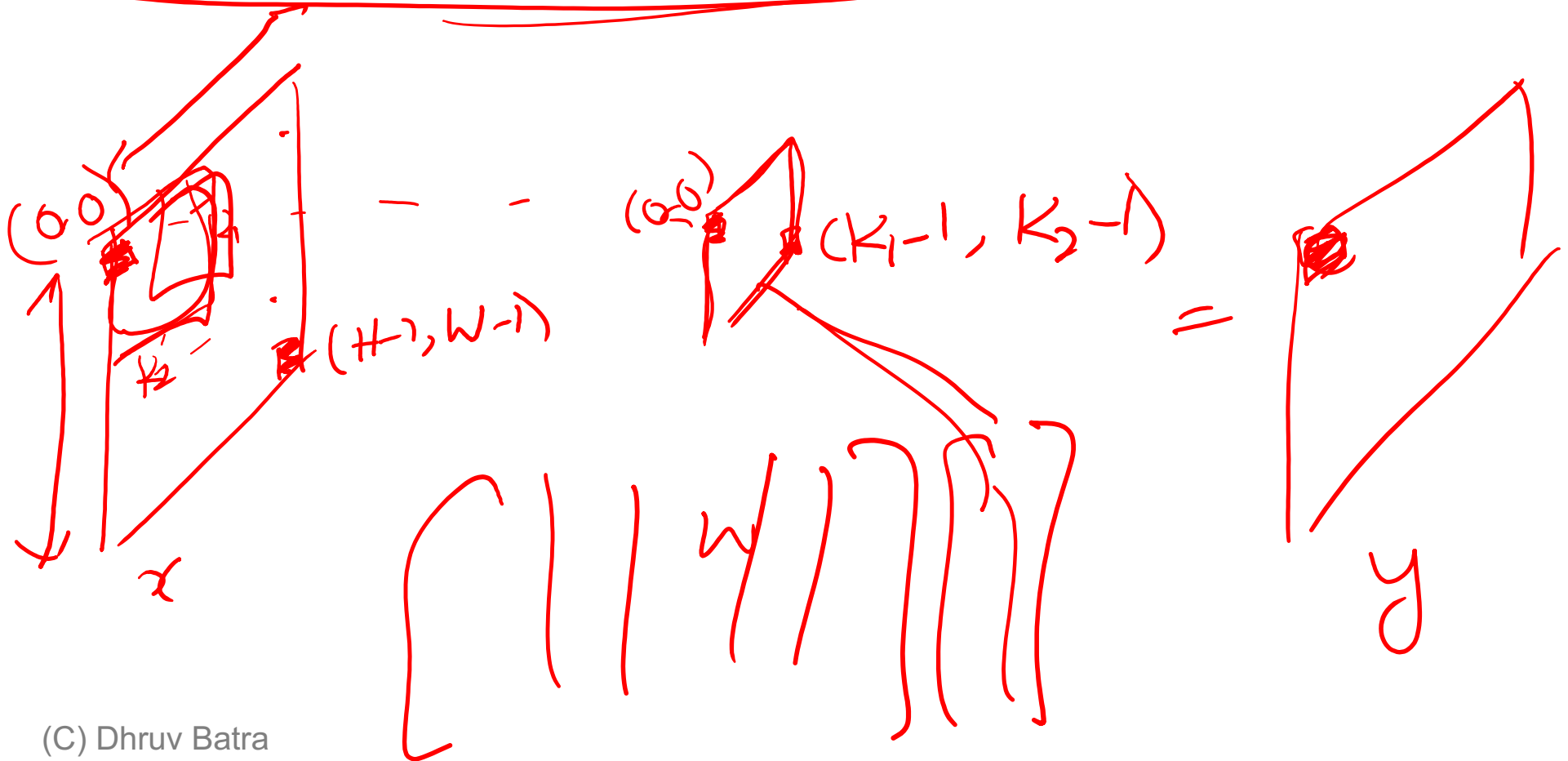


$$y[x, c] = \sum_{a=-\frac{K_1-1}{2}}^{\frac{K_1-1}{2}} \sum_{b=-\frac{K_2-1}{2}}^{\frac{K_2-1}{2}} x[x-a, c-b] w[a, b]$$

Convolutions for computer scientists

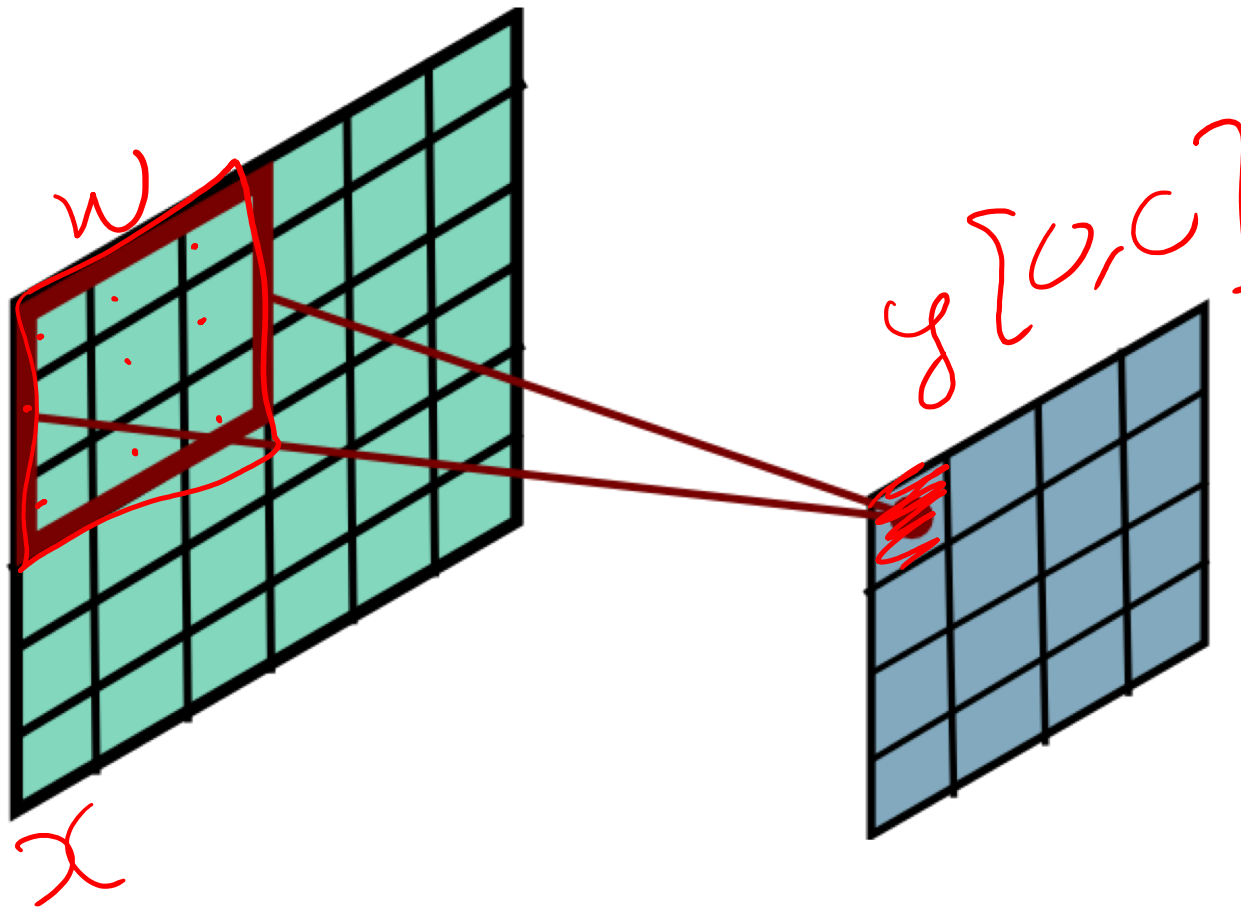
Convolutions for programmers

$$y[x, c] = \sum_{a=0}^{K_1-1} \sum_{b=0}^{K_2-1} x[(x+a), (y+b)] w[a, b]$$

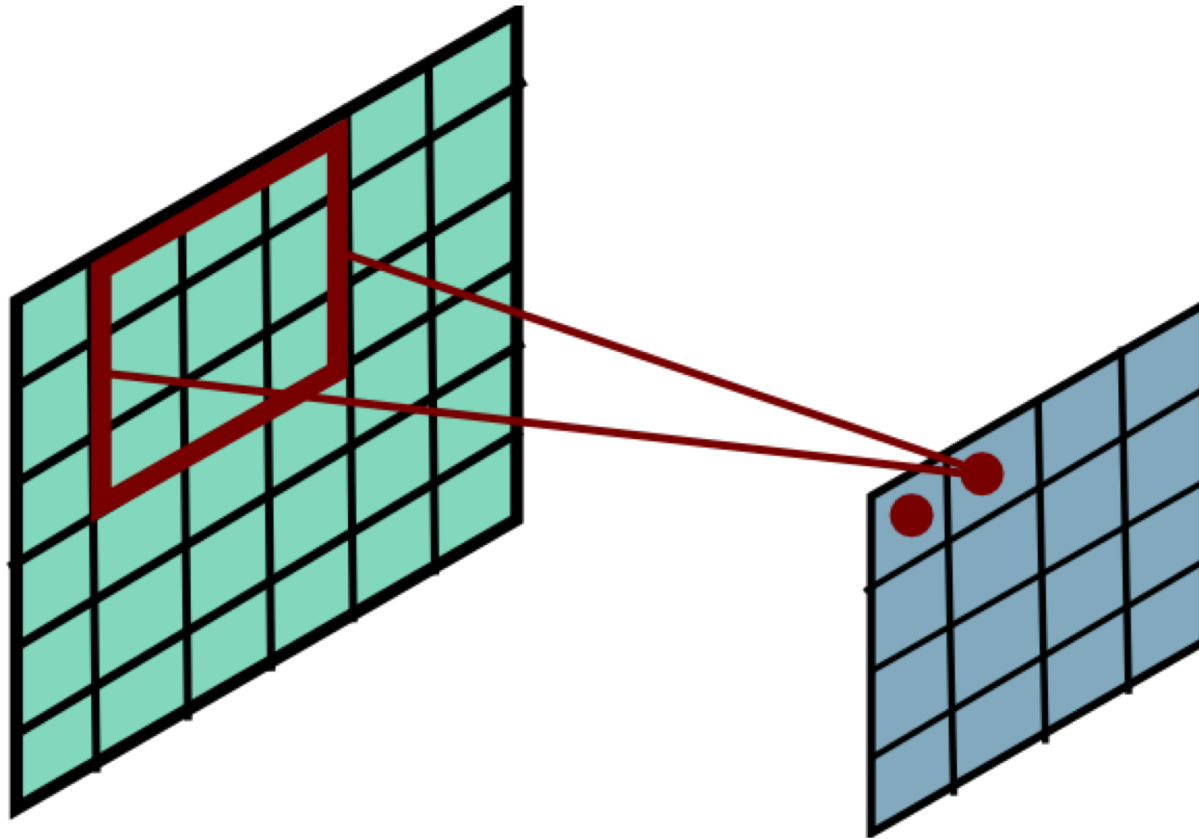


Convolutions for programmers

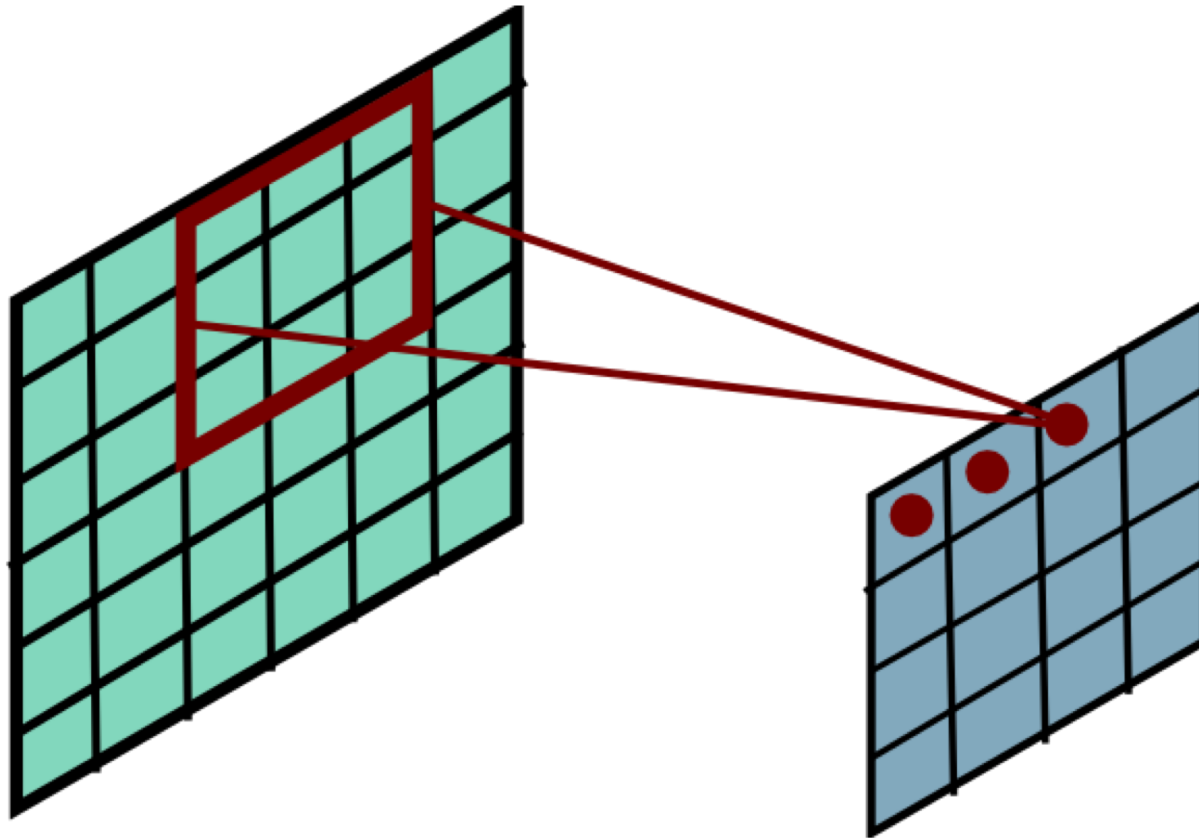
Convolutional Layer



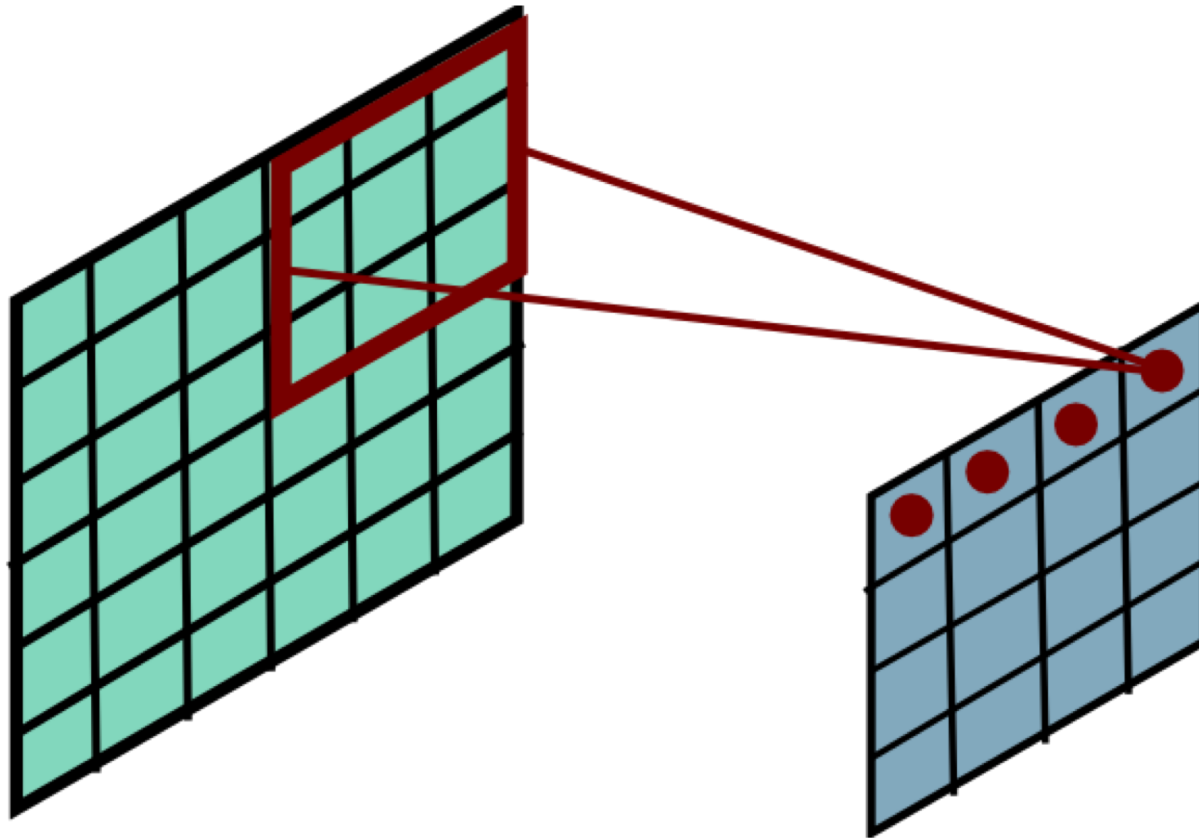
Convolutional Layer



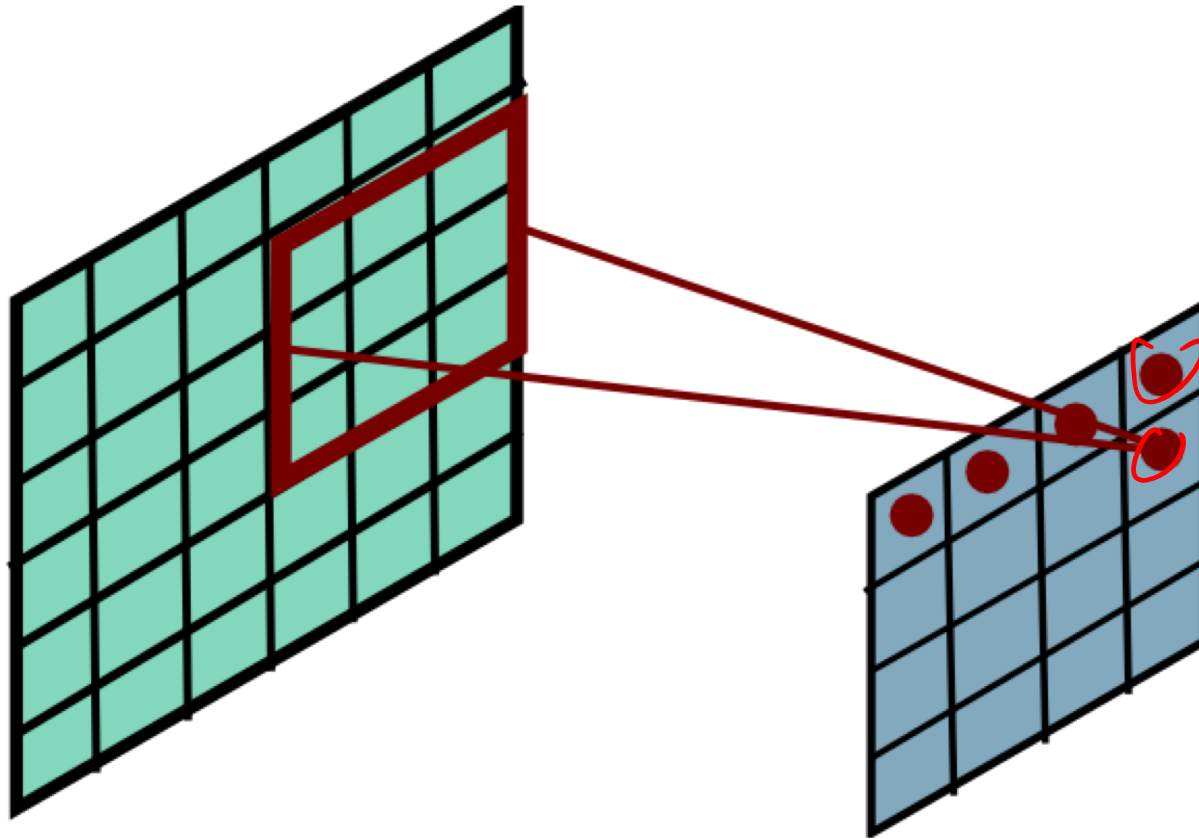
Convolutional Layer



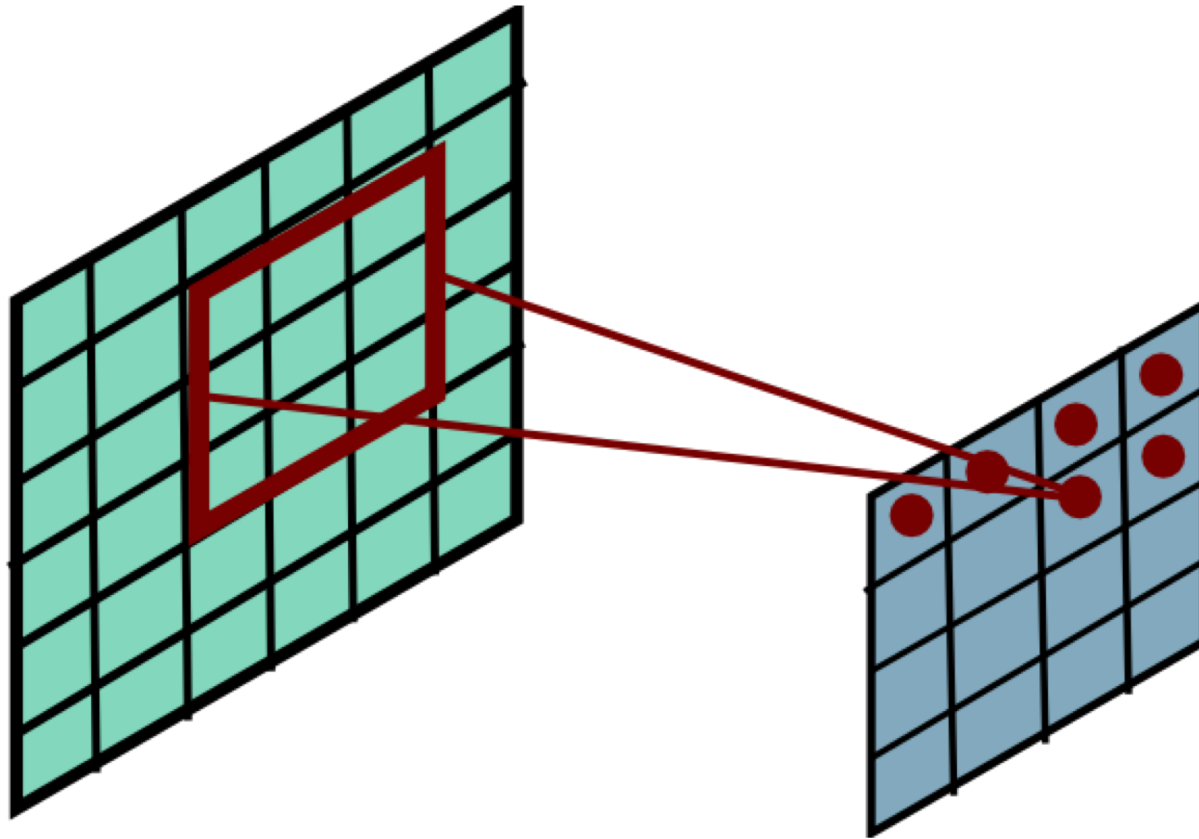
Convolutional Layer



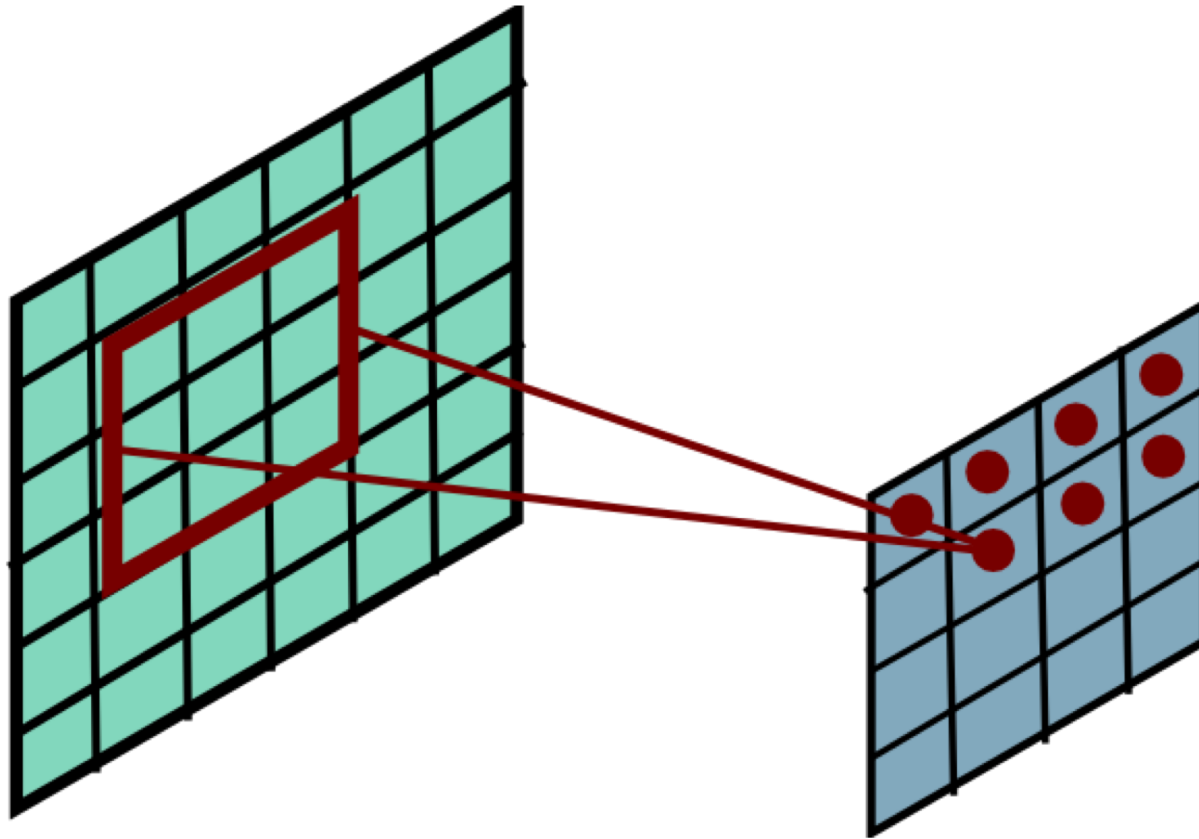
Convolutional Layer



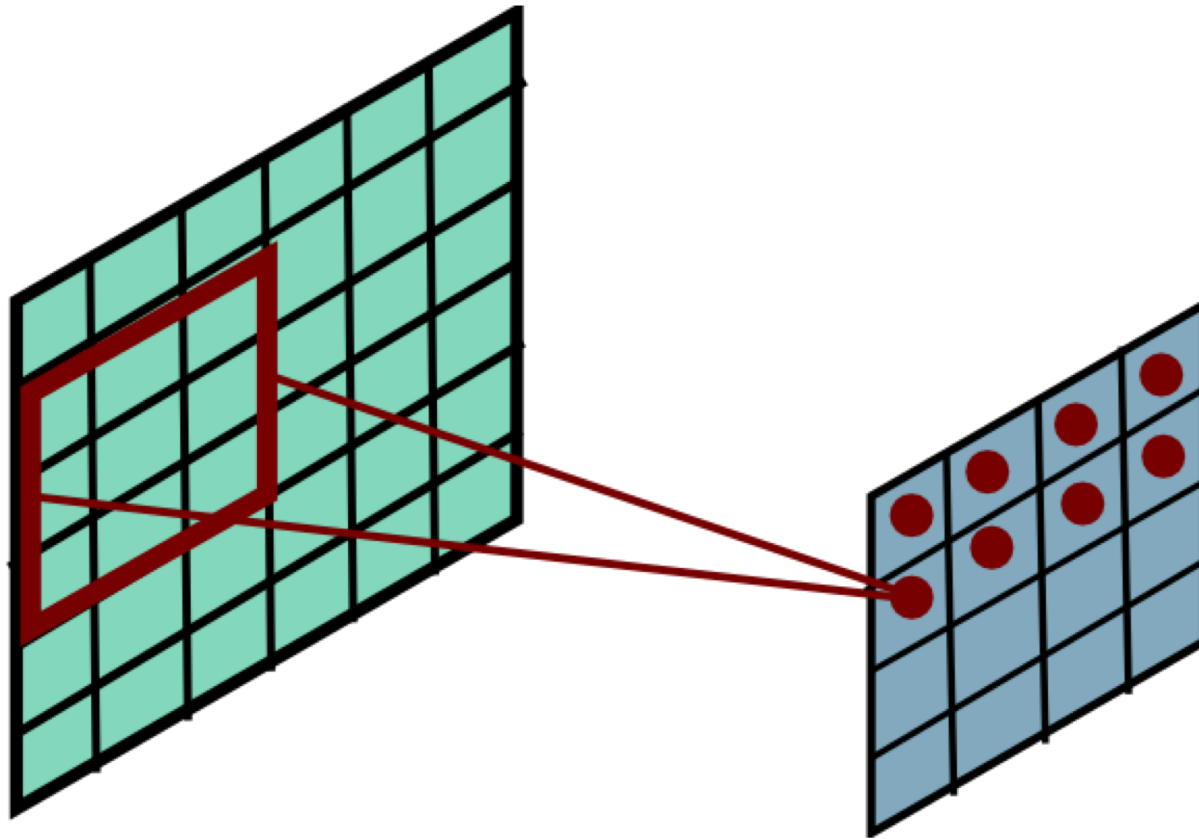
Convolutional Layer



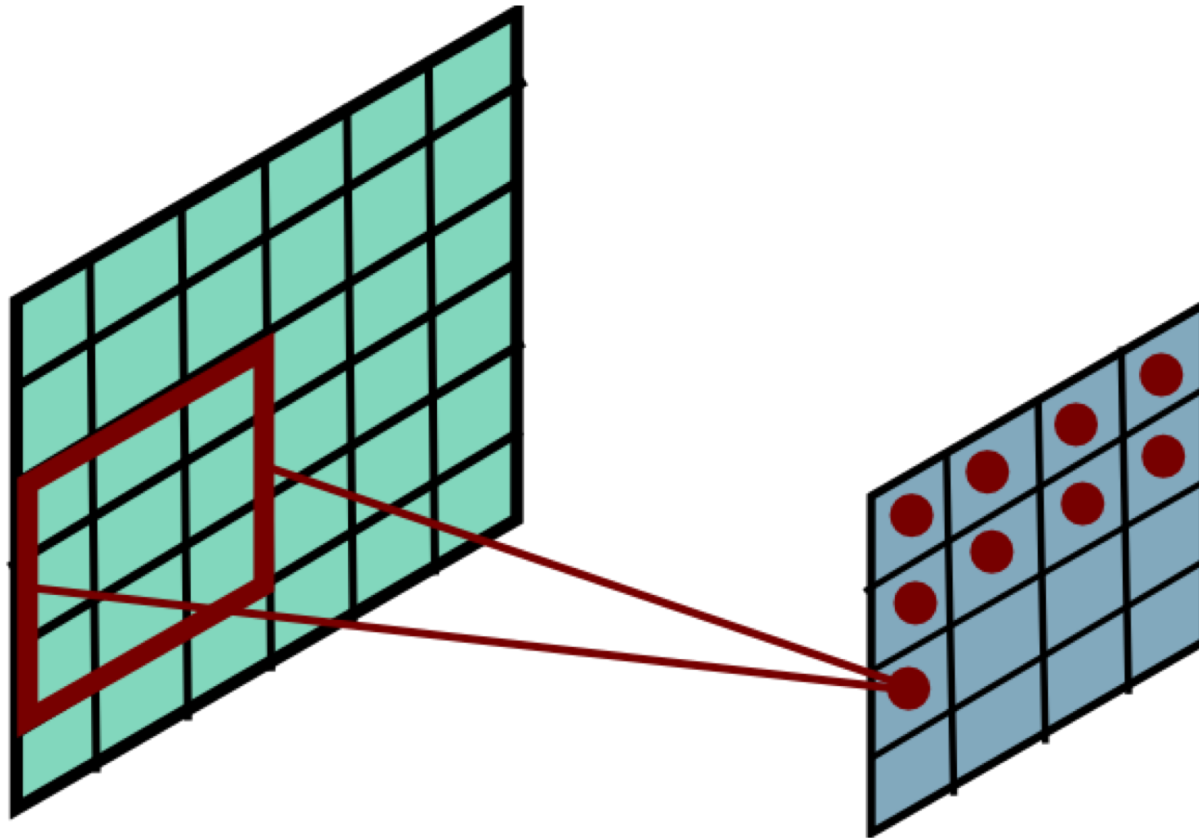
Convolutional Layer



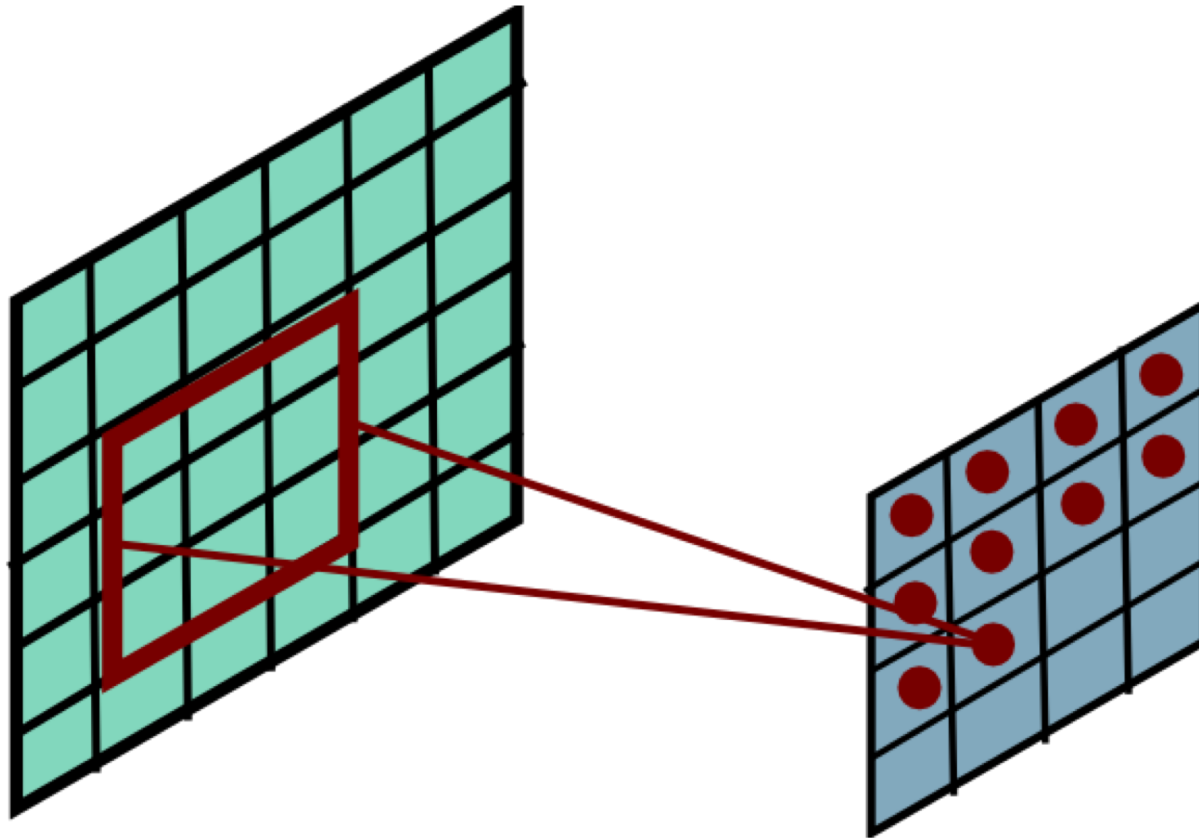
Convolutional Layer



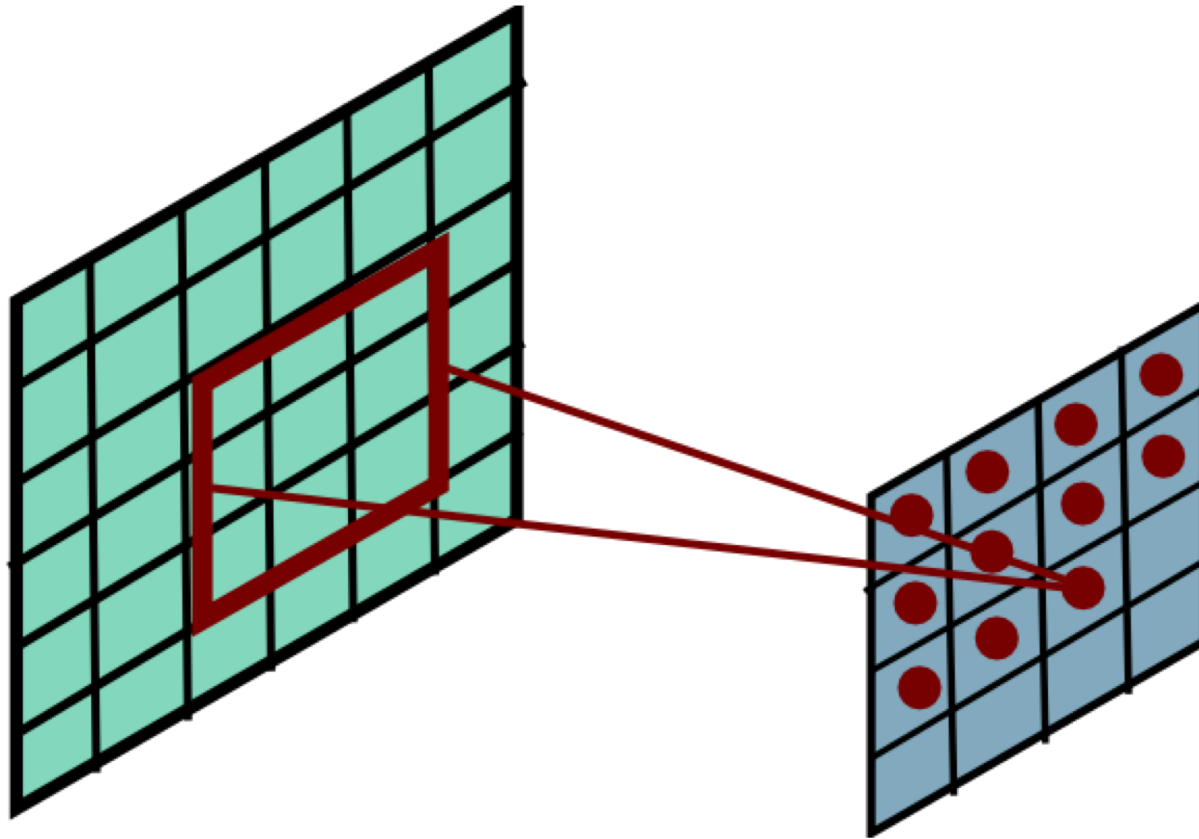
Convolutional Layer



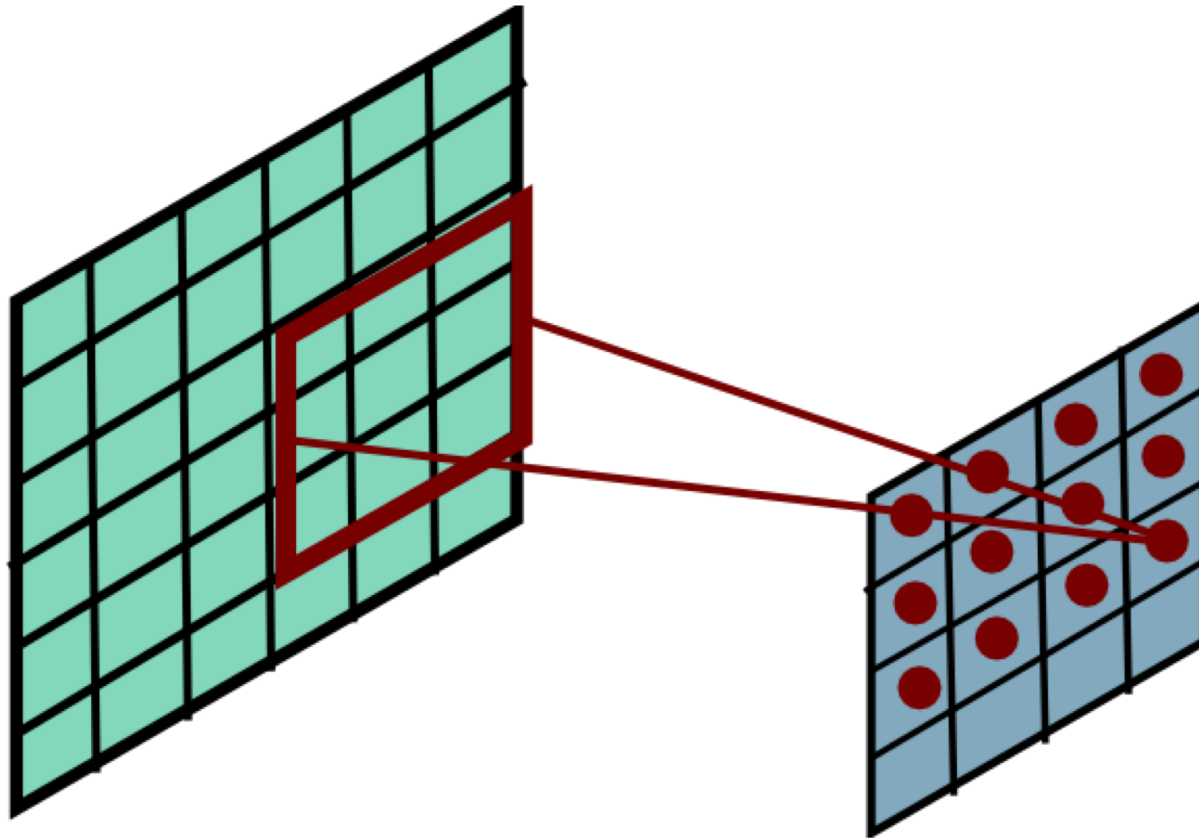
Convolutional Layer



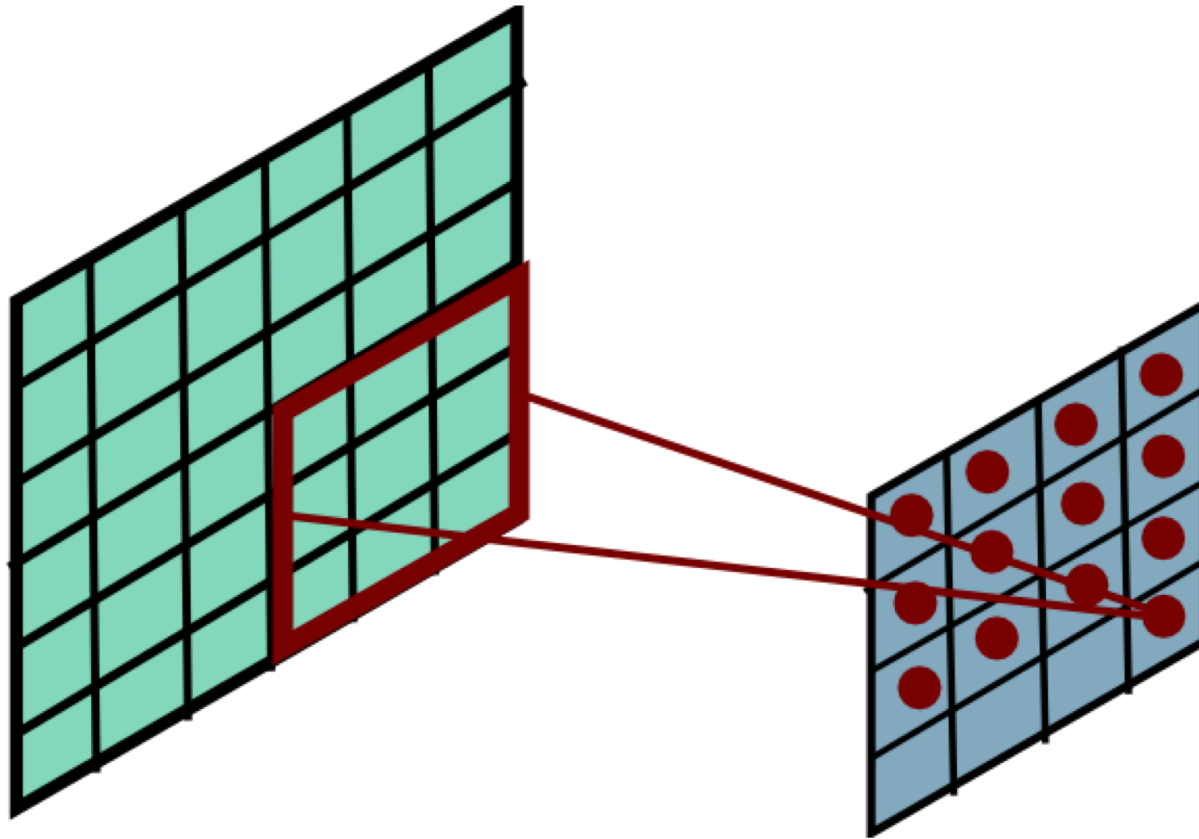
Convolutional Layer



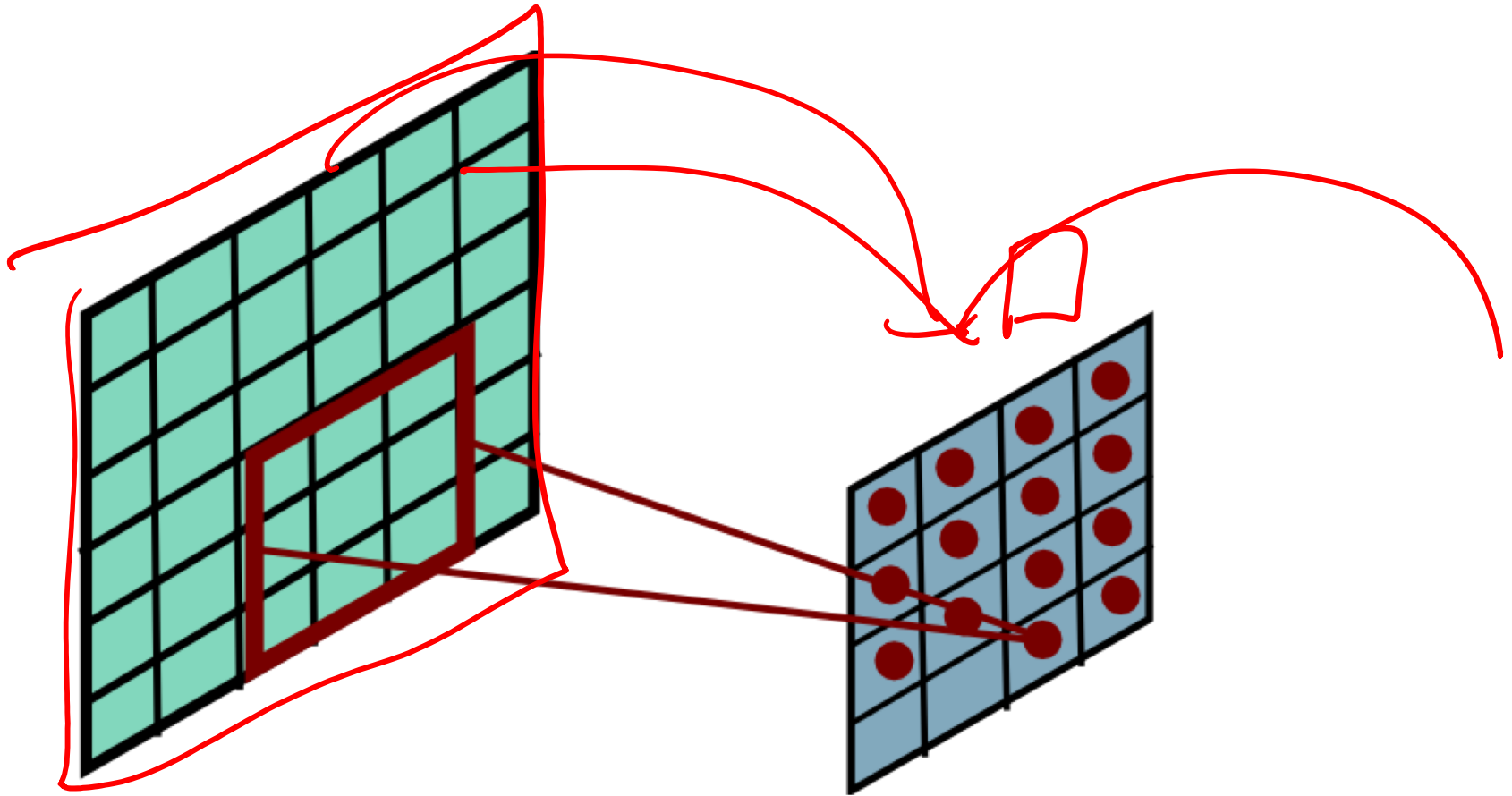
Convolutional Layer



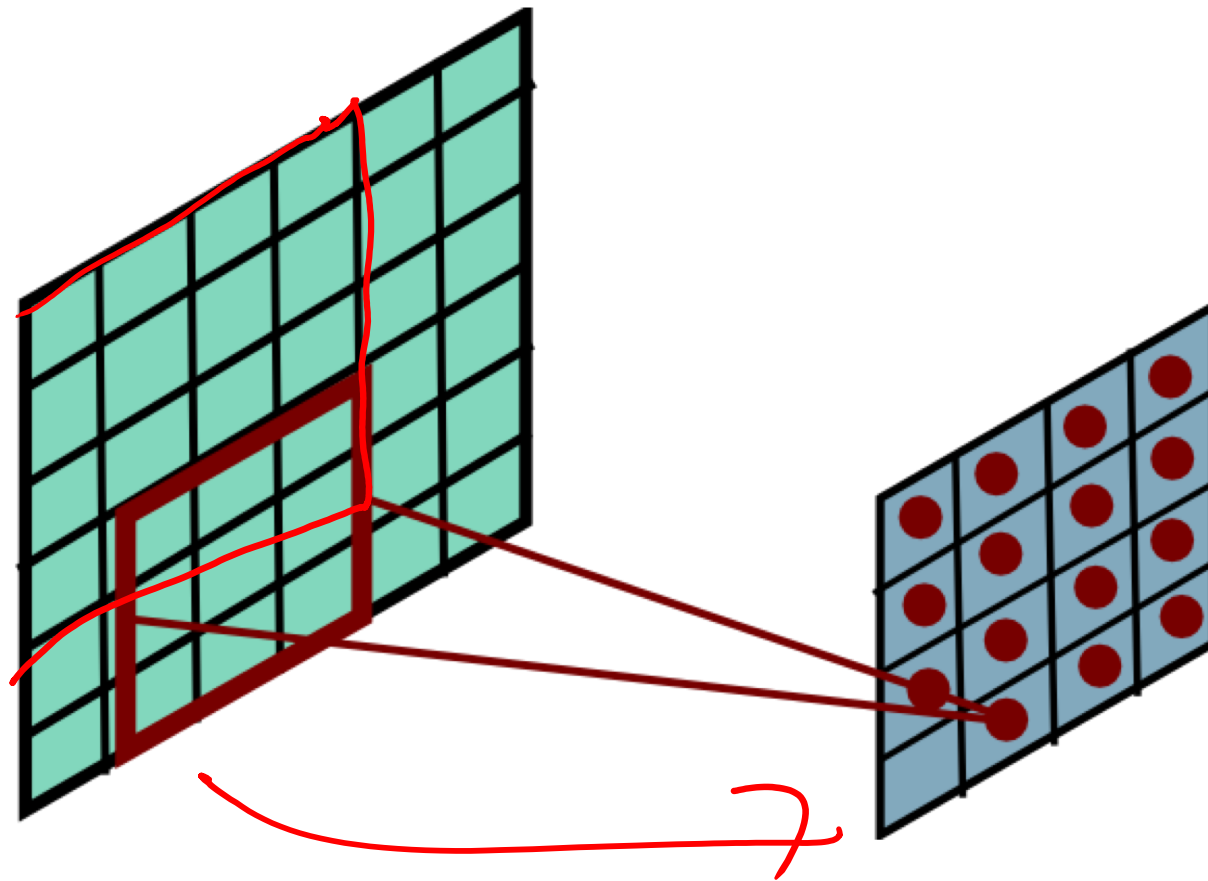
Convolutional Layer



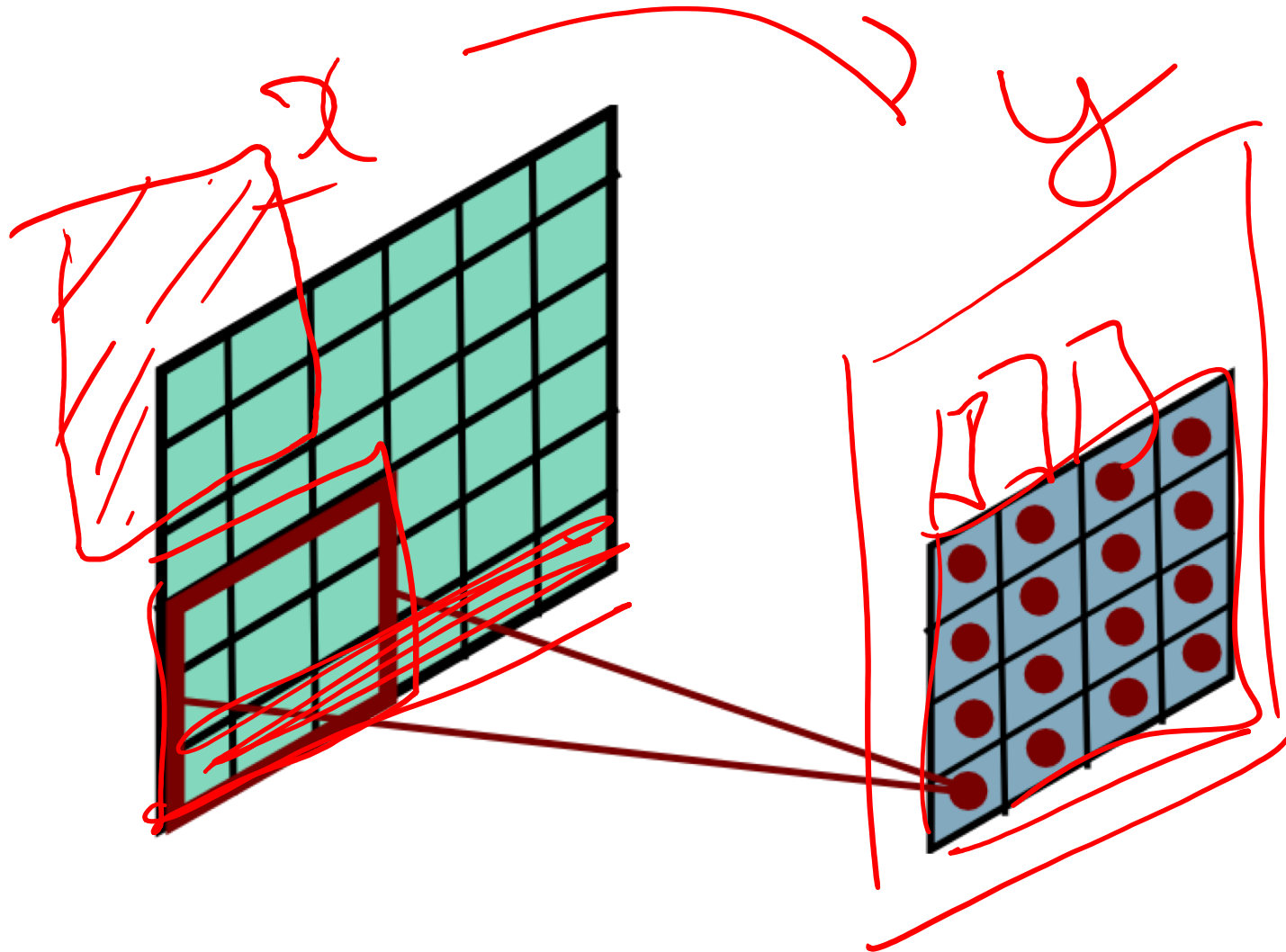
Convolutional Layer



Convolutional Layer



Convolutional Layer

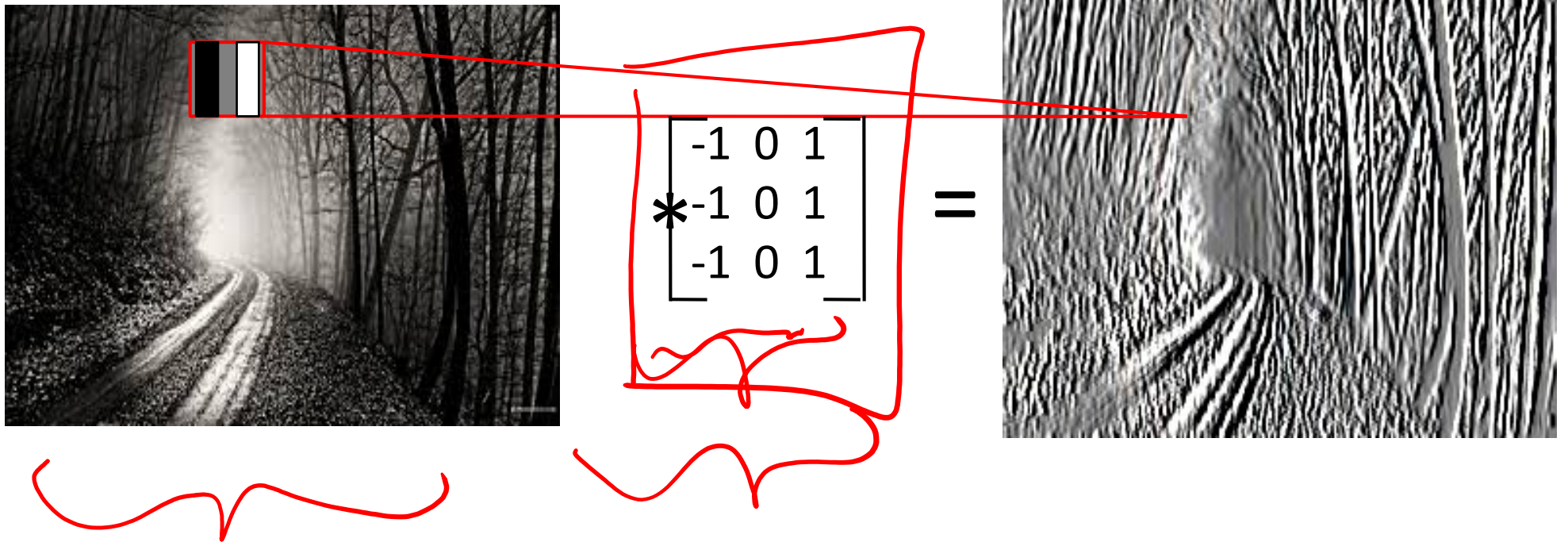


Mathieu et al. "Fast training of CNNs through FFTs" ICLR 2014

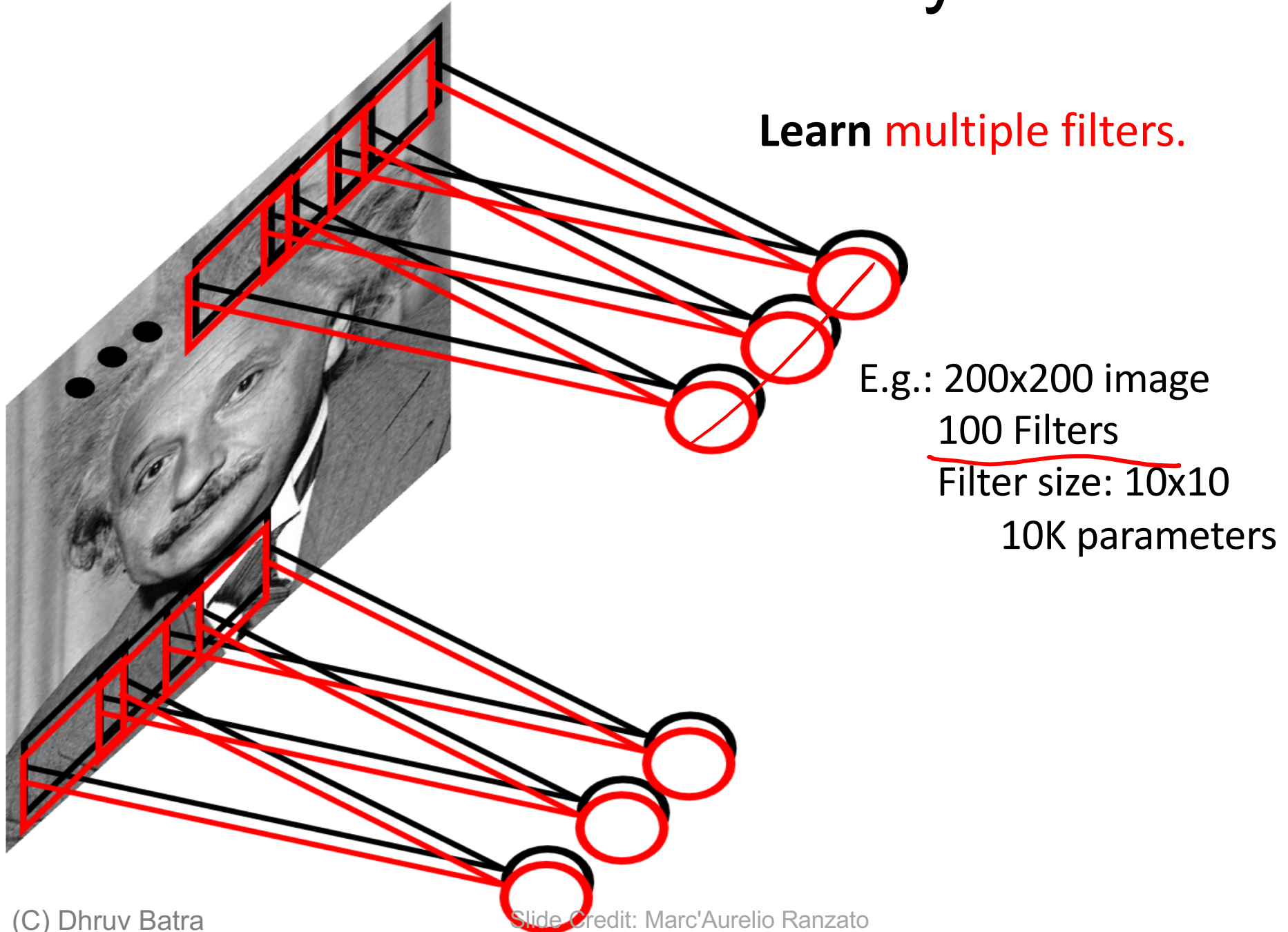
Convolution Explained

- <http://setosa.io/ev/image-kernels/>
- <https://github.com/bruckner/deepViz>

Convolutional Layer



Convolutional Layer

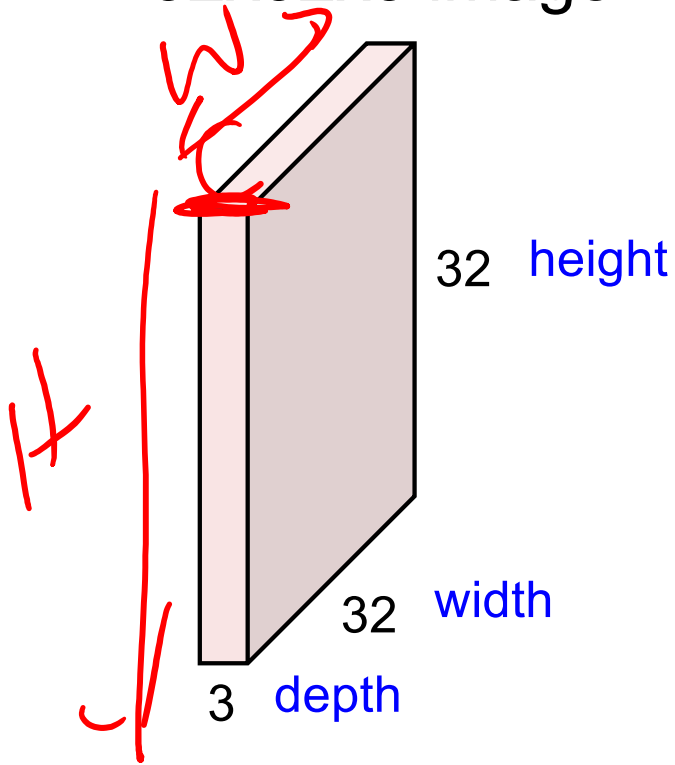


Convolutional Layer

Convolutional Layer

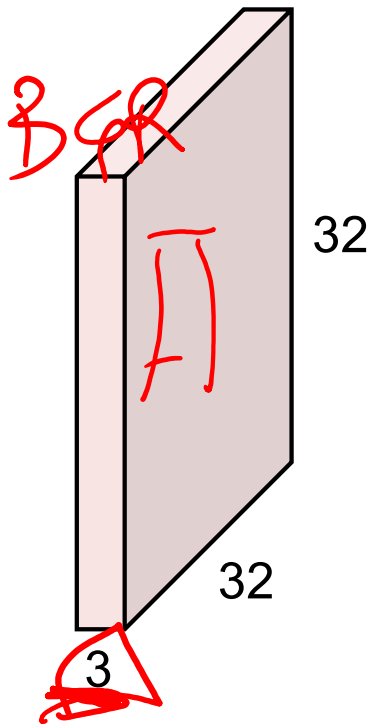
Convolution Layer

32x32x3 image -> preserve spatial structure

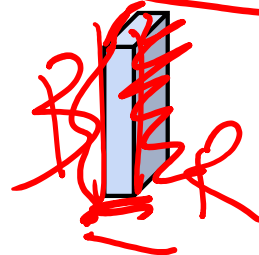


Convolution Layer

32x32x3 image



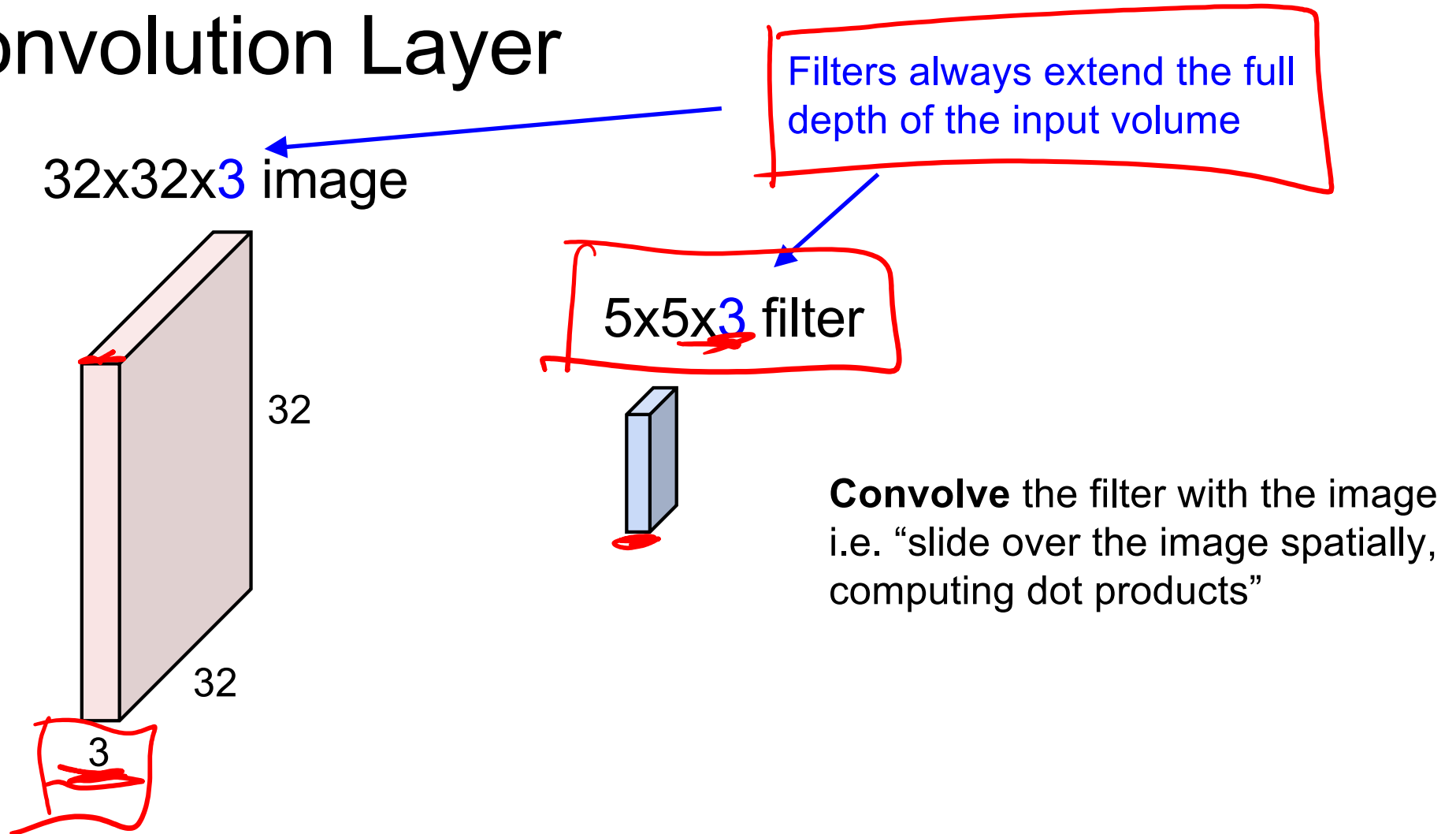
5x5x3 filter



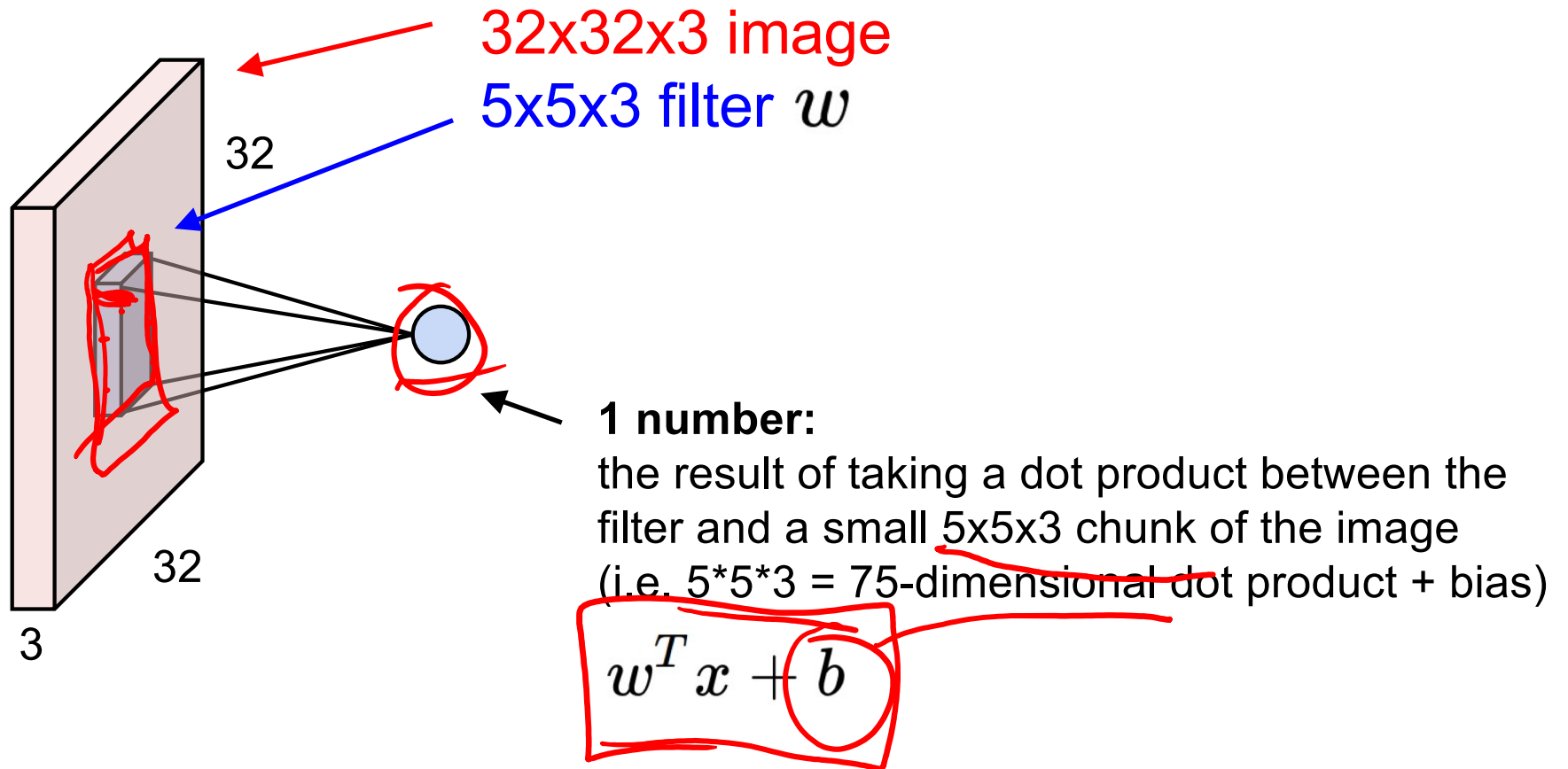
Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

$W \in \mathbb{R}^{k_1 \times k_2 \times 3}$

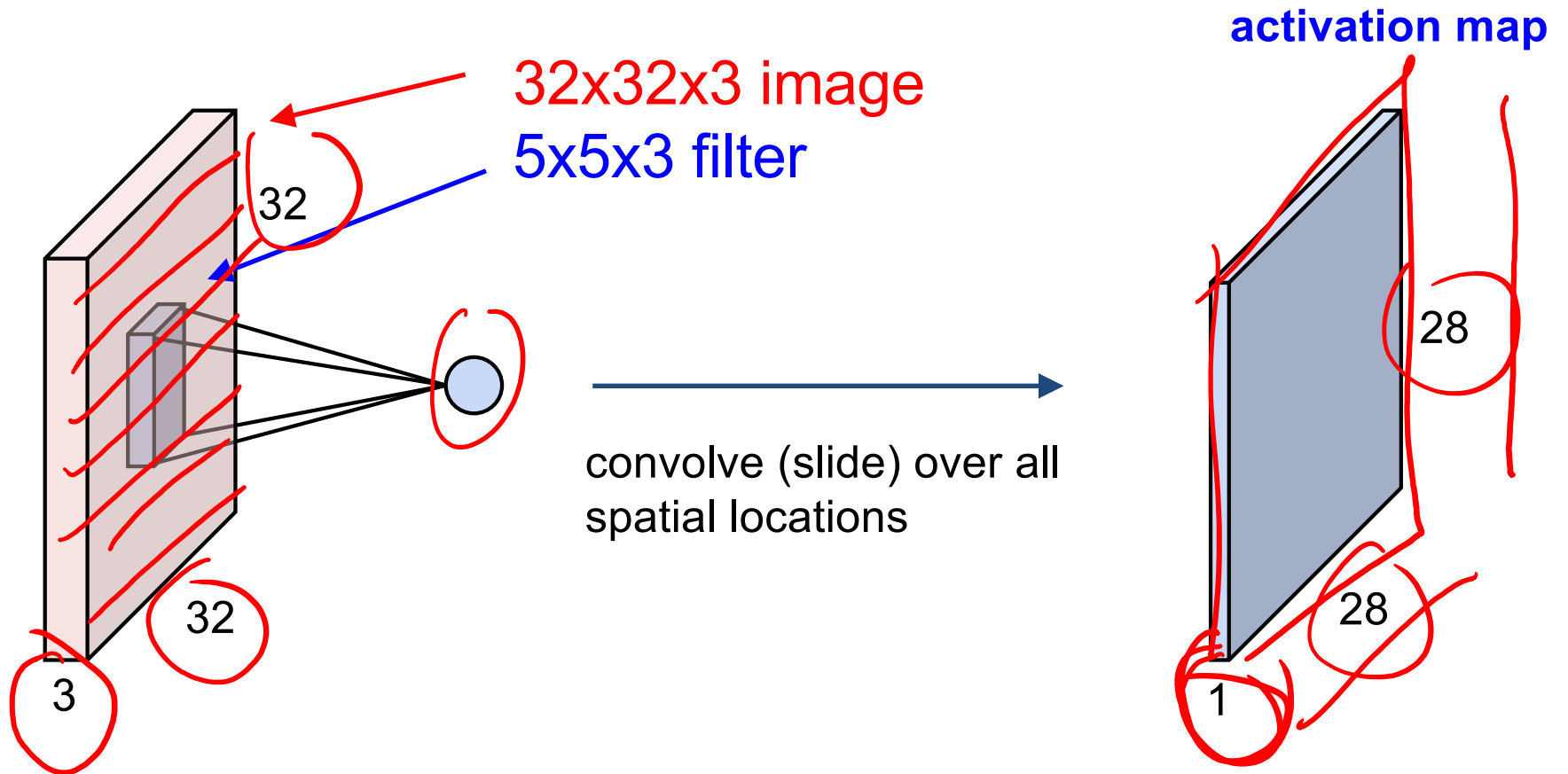
Convolution Layer



Convolution Layer

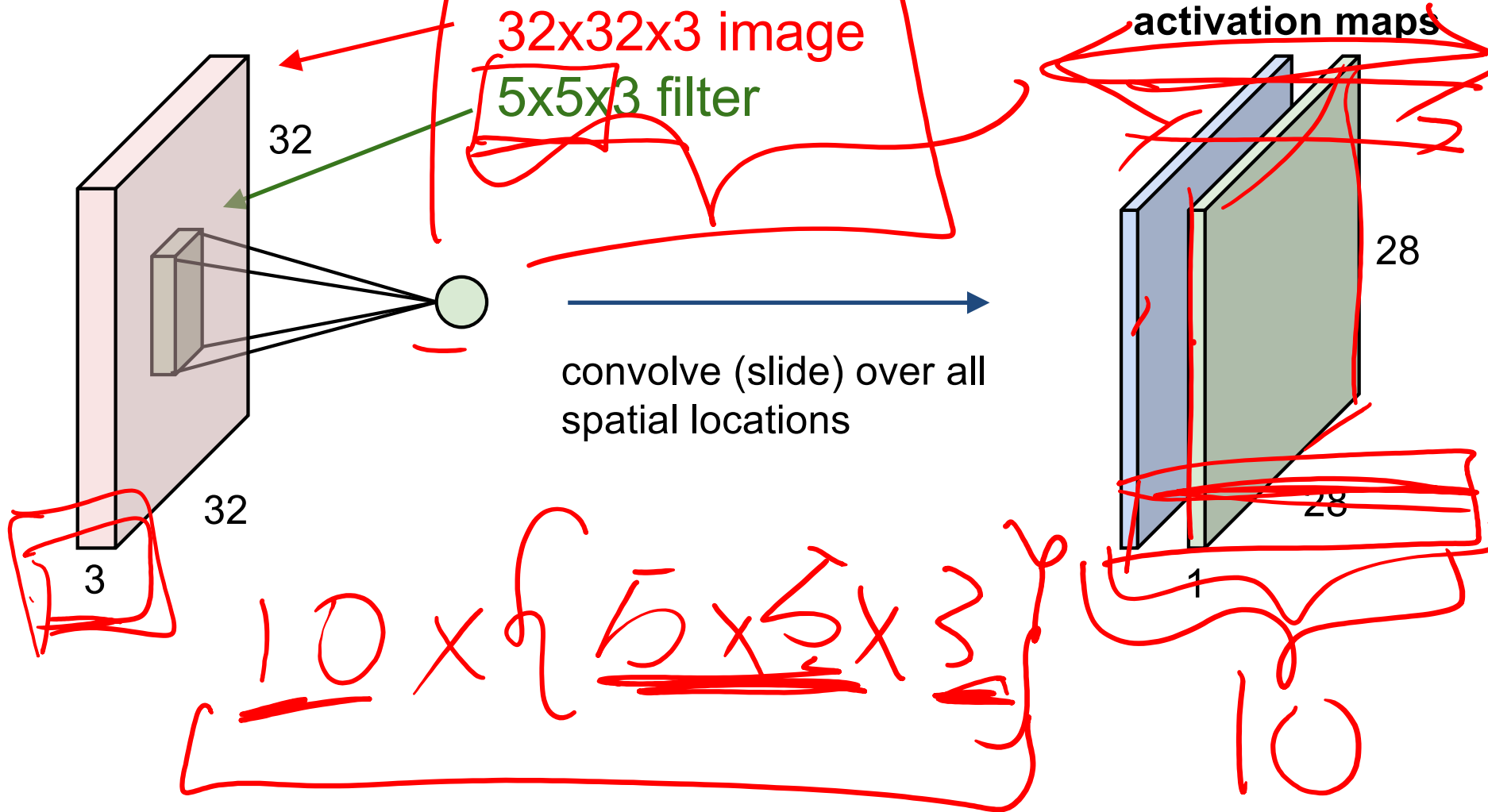


Convolution Layer

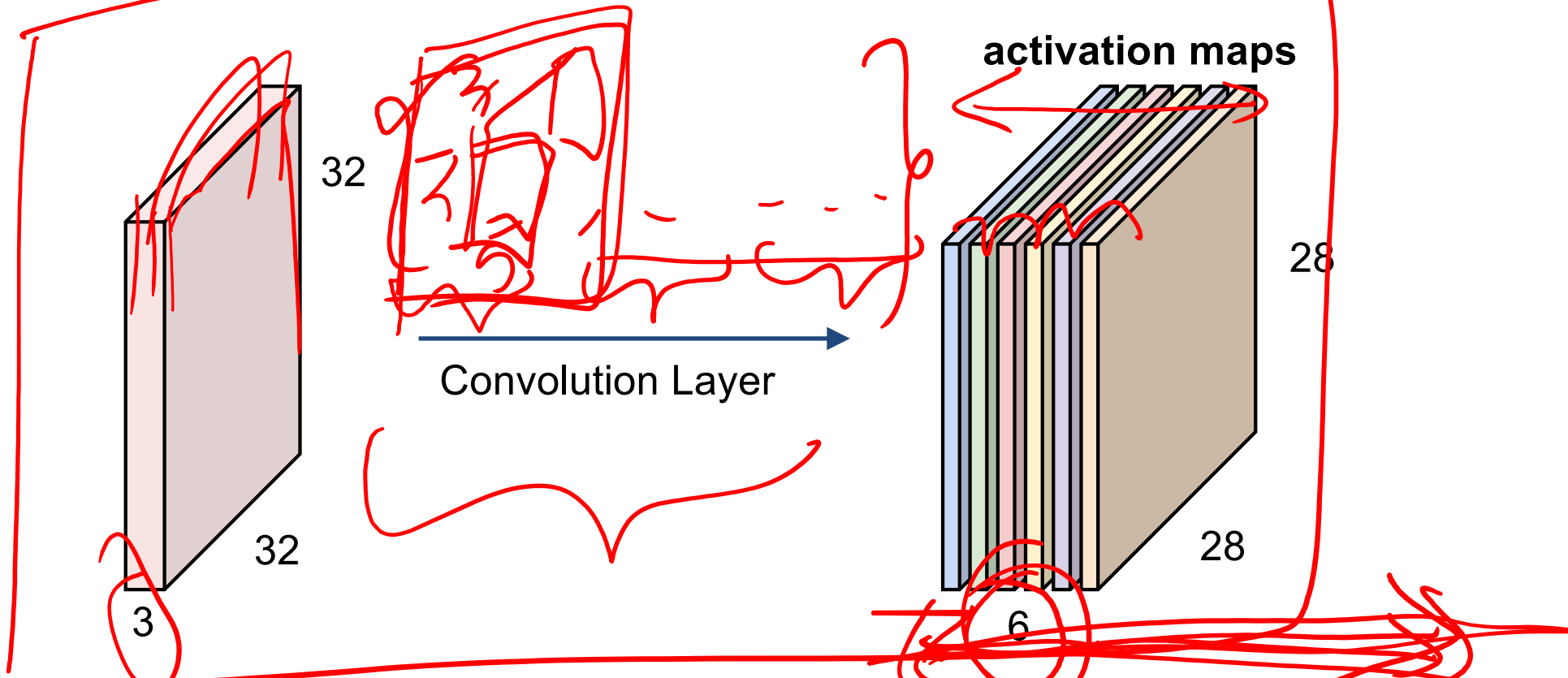


Convolution Layer

consider a second, **green** filter



For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a "new image" of size 28x28x6!