# CS 7643: Deep Learning

Topics:
- Computational Graphs
  - Notation + example
- Computing Gradients
  - Forward mode vs Reverse mode AD

Dhruv Batra

Georgia Tech

# Administrivia

- HW1 Released
  - Due: 09/22

- PS1 Solutions
  - Coming soon

# 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
  - Encouraged to apply to your research (computer vision, NLP, robotics,…)
  - Must be done this semester.

- 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
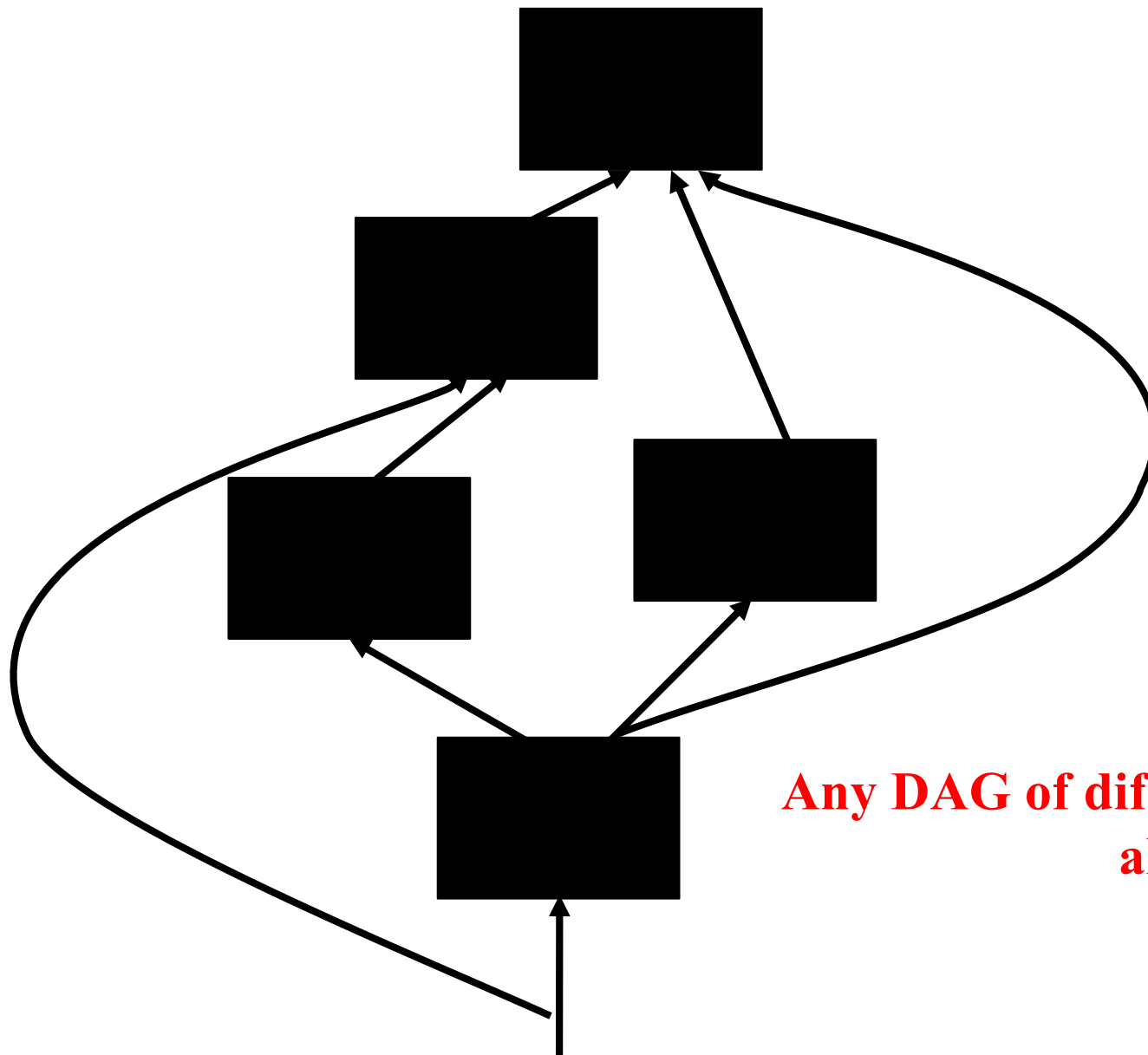    - Theoretically analyze an existing algorithm

# Administrativia

- Project Teams Google Doc
  - https://docs.google.com/spreadsheets/d/1AaXY0JE4IAbHvo DaWlc9zsmfKMyuGS39JAn9dpeXhhQ/edit#gid=0
  - Project Title
  - 1-3 sentence project summary TL;DR
  - Team member names + GT IDs

# Recap of last time

# How do we compute gradients?

- 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!**
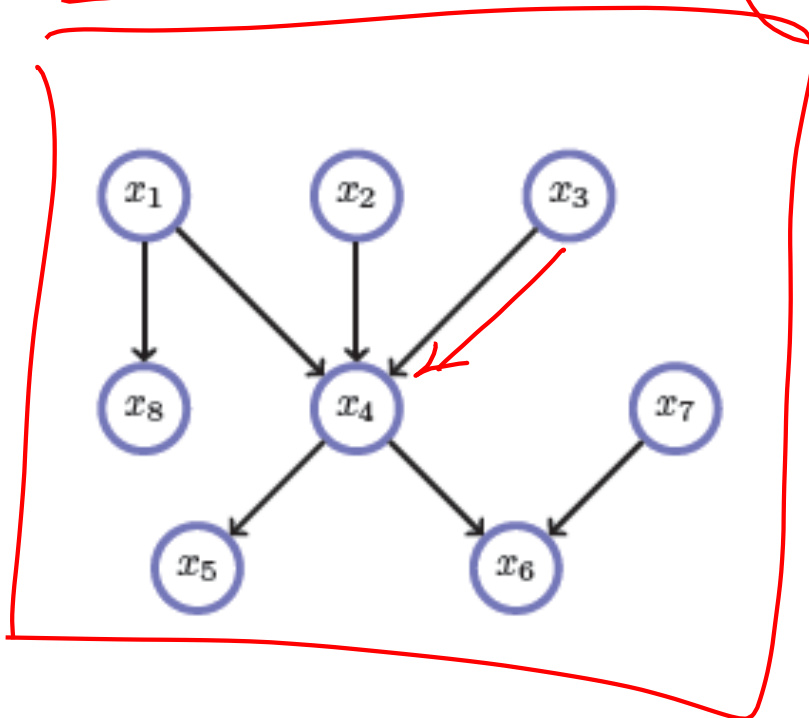
# Directed Acyclic Graphs (DAGs)

- Exactly what the name suggests
  - Directed edges
  - No (directed) cycles
  - Underlying undirected cycles okay

$$E = \{ (v_i, v_j) \}$$

$$\{ \{ v_i, v_j \} \}$$
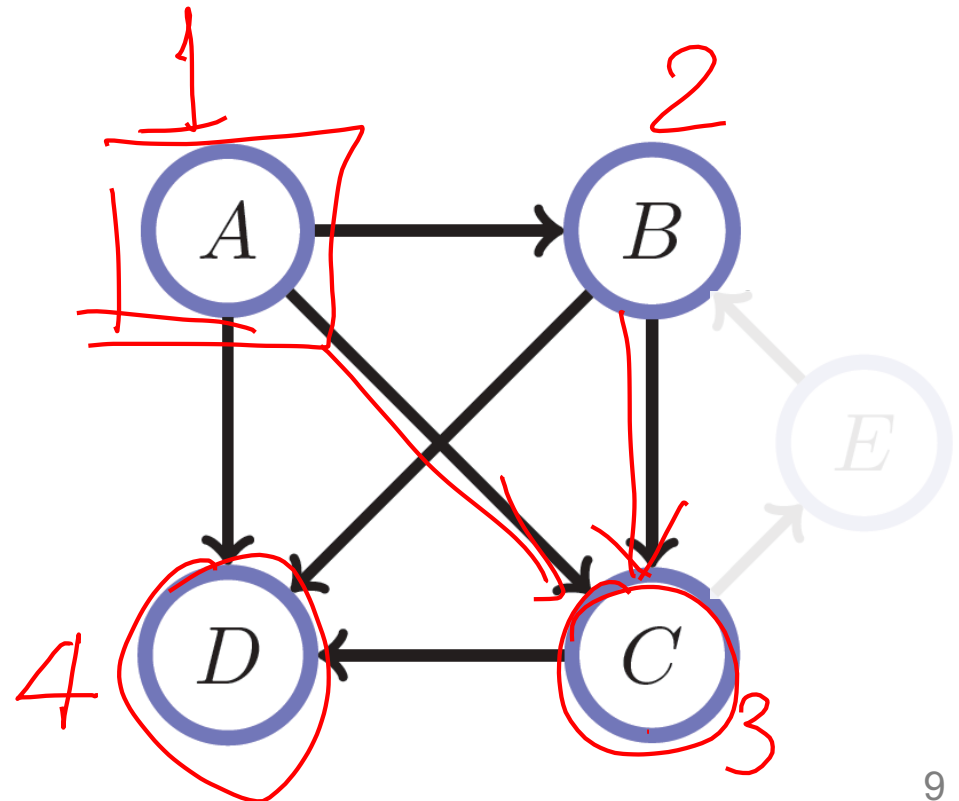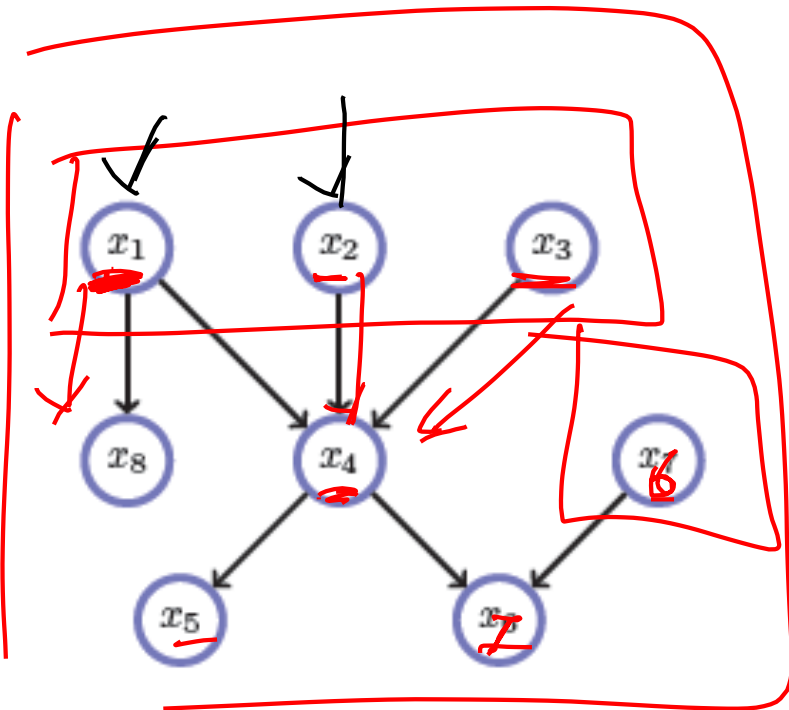
$(v_i, v_j) \ldots v_i)$

$(D, A) \notin E$

# Directed Acyclic Graphs (DAGs)

- Concept
  - Topological Ordering
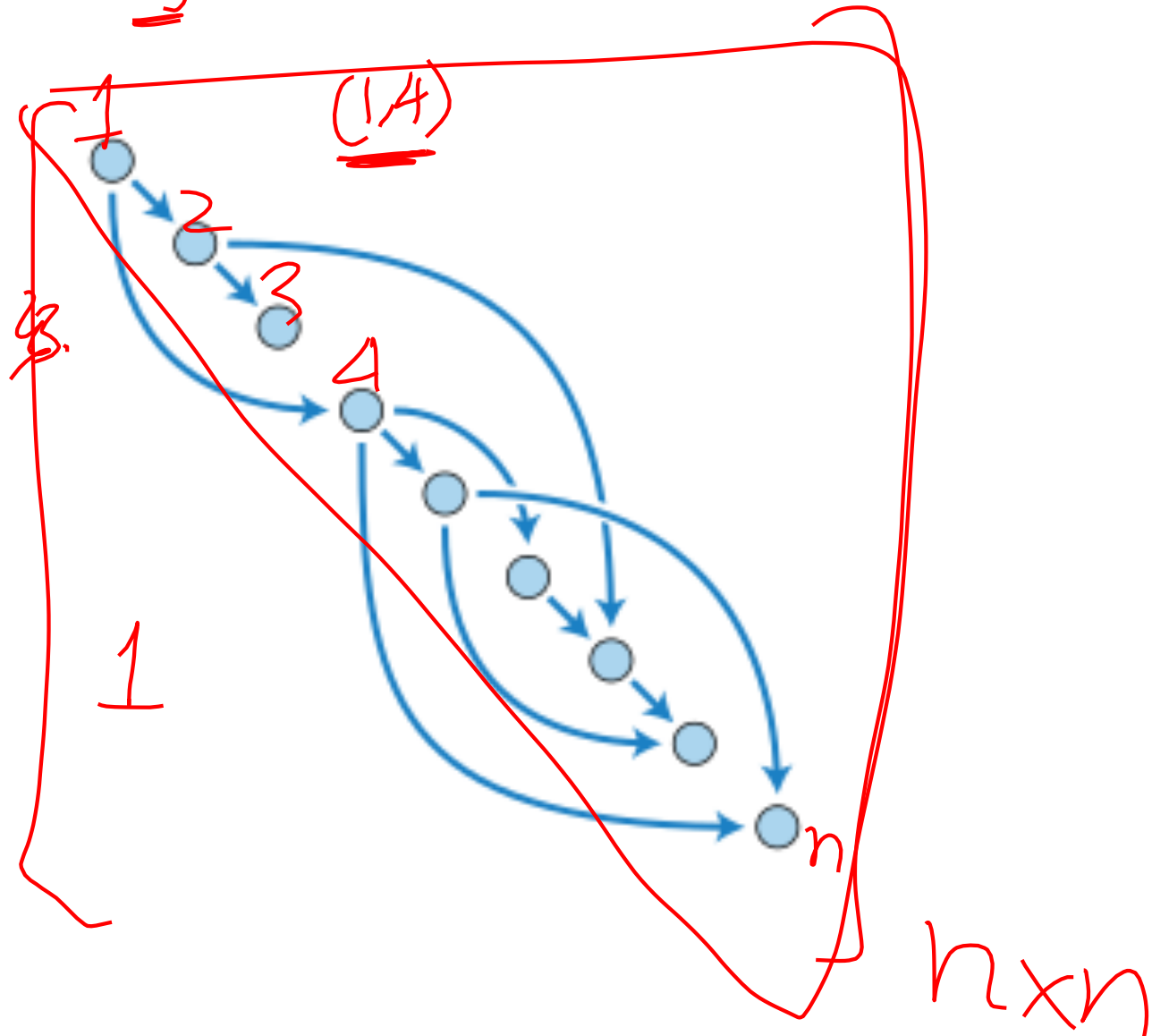
$$\exists\ \sigma : V \to [n] = \{1, \ldots, n\}$$

$$s.t\ (v_i, v_j) \in E \quad \sigma(v_i) < \sigma(v_j)$$

# Directed Acyclic Graphs (DAGs)

$$a_{ij} = 1 \quad (i,j) \in E$$

$$(1,4)$$
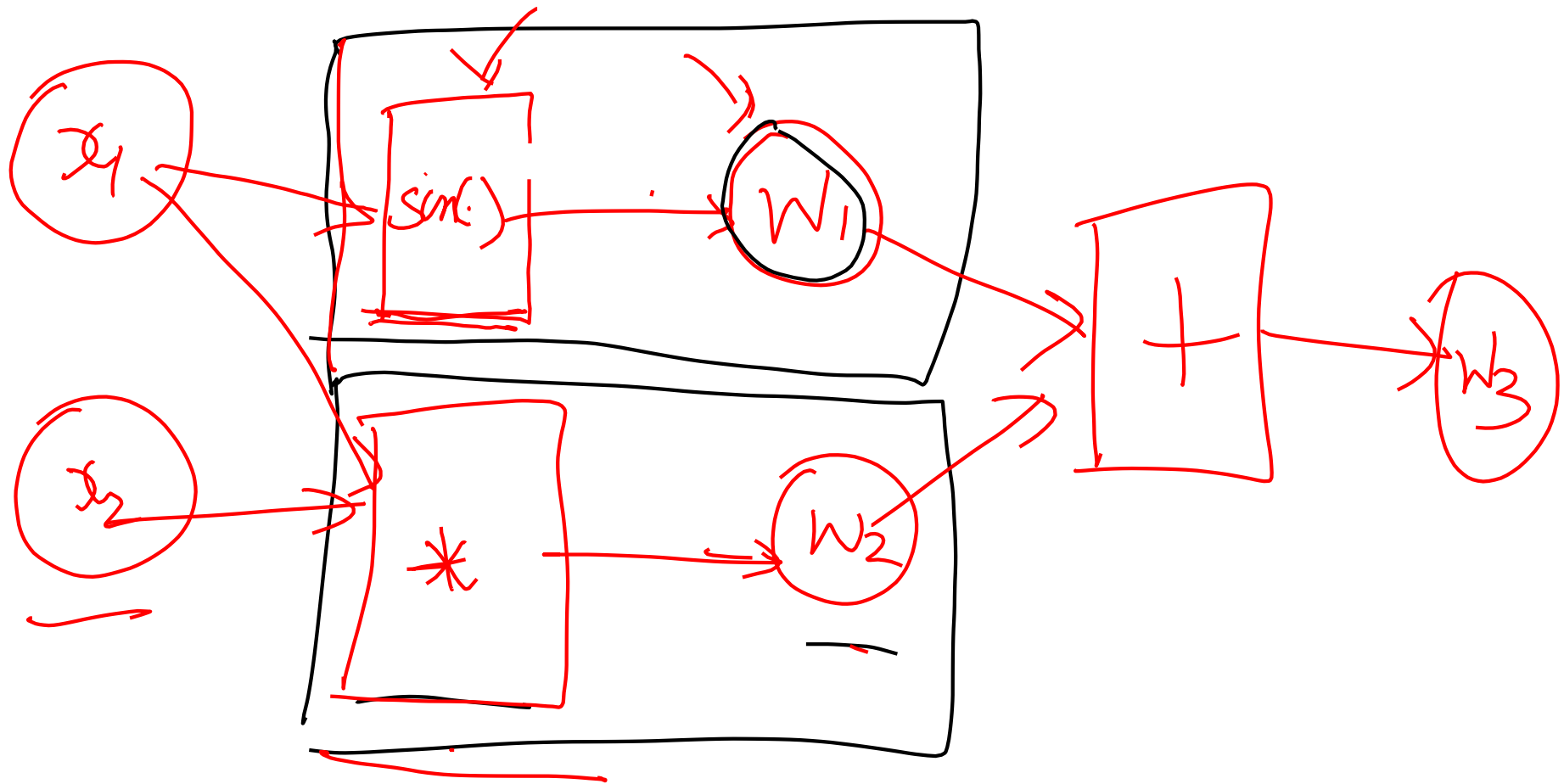
$$A =$$

1

$n \times n$

# Computational Graphs

- Notation #1
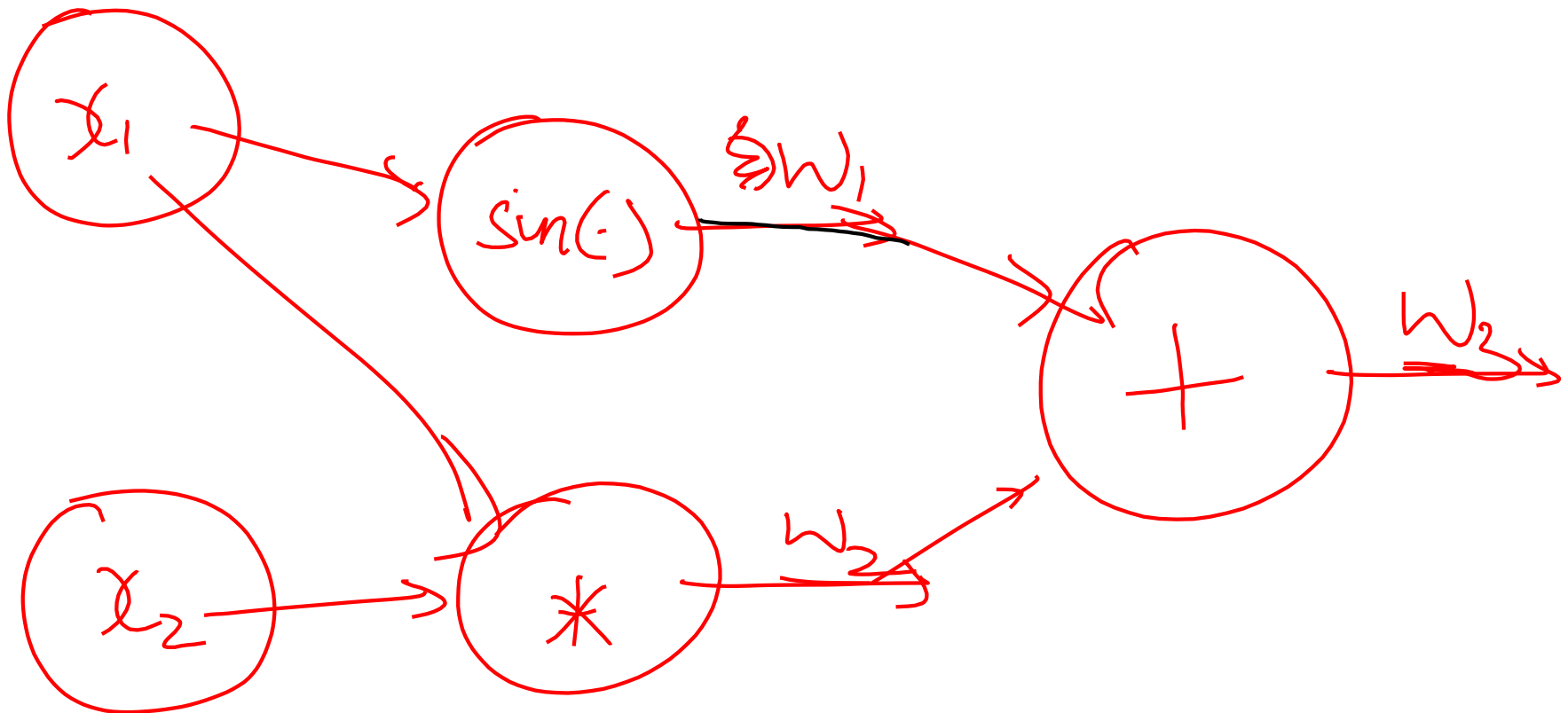
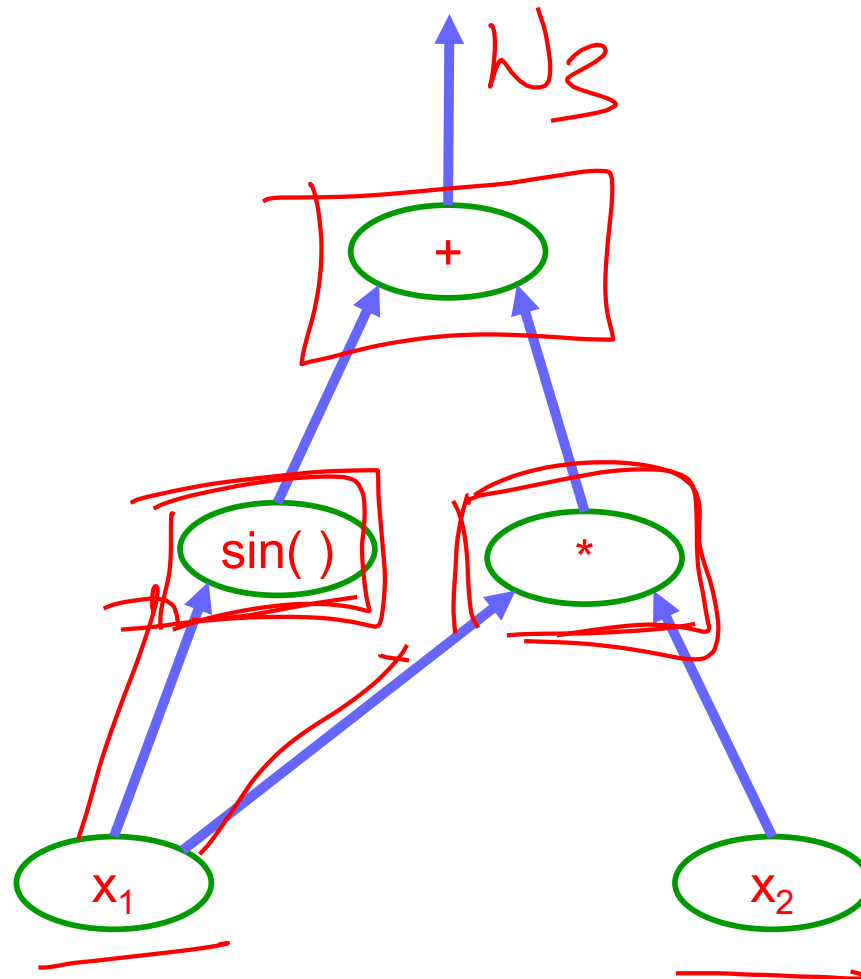$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

# Computational Graphs

- Notation #2

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

# Example

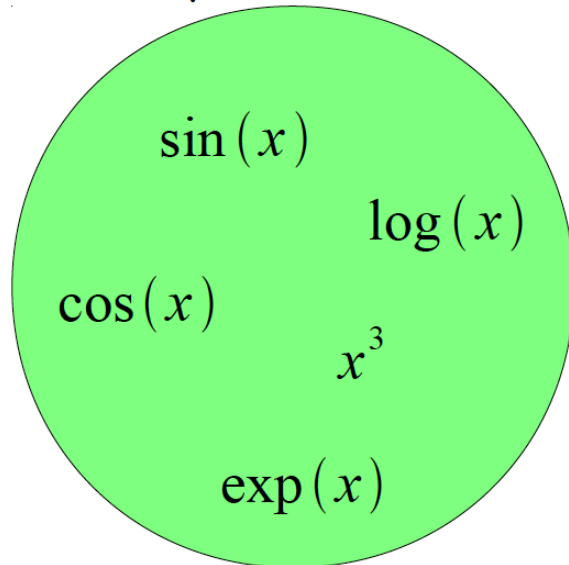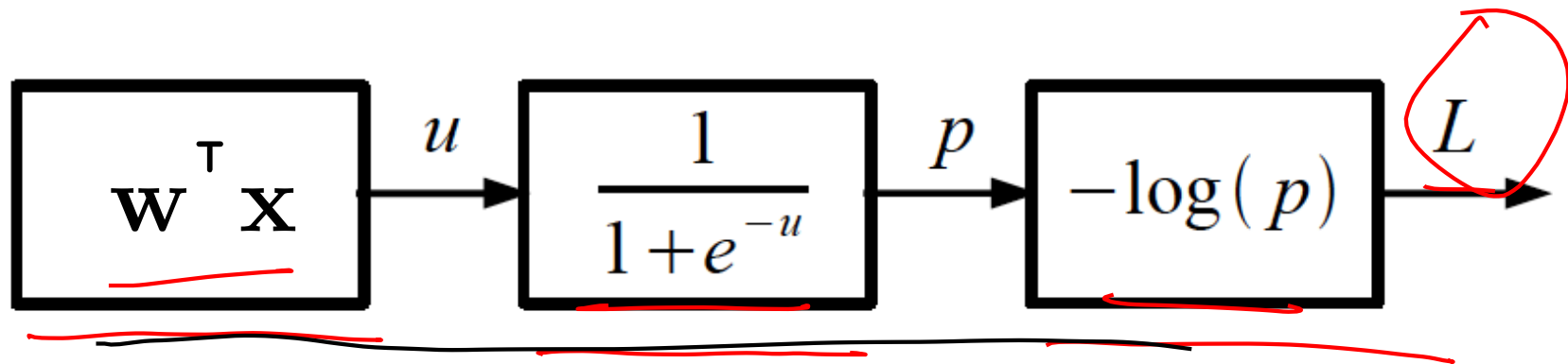$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

# Logistic Regression as a Cascade

Given a library of simple functions

$$\sin(x)$$
$$\log(x)$$
$$\cos(x)$$
$$x^3$$
$$\exp(x)$$

Compose into a

complicate function

$$- \log \left( \frac{1}{1 + e^{-\mathbf{w}^\mathsf{T}\mathbf{x}}} \right)$$

$$\mathbf{w}^\mathsf{T}\mathbf{x} \xrightarrow{u} \frac{1}{1+e^{-u}} \xrightarrow{p} -\log(p) \rightarrow L$$

# Forward mode vs Reverse Mode

- Key Computations

# Forward mode AD

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



layer

$\underline{l}$ $g(\cdot)$

$\vec{h}^{l-1}$

$\frac{\partial \vec{h}^{l-1}}{\partial \vec{x}}$

$\vec{h}^l$

$\frac{\partial \vec{h}^l}{\partial \vec{x}} = \frac{\partial \vec{h}^l}{\partial \vec{h}^{l-1}} \cdot \frac{\partial \vec{h}^{l-1}}{\partial \vec{x}}$

Jacobian

$$\vec{h}^l = g(\vec{h}^{l-1})$$

# Reverse mode AD

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

$$\vec{h}^{\,\ell-1} \qquad \vec{h}^{\,\ell}$$

$$g$$

$$\frac{\partial L}{\partial \vec{h}^{\,\ell-1}} = \frac{\partial L}{\partial \vec{h}^{\,\ell}} \; \frac{\partial \vec{h}^{\,\ell}}{\partial \vec{h}^{\,\ell-1}}$$

input      Jacobian

$$\frac{\partial L}{\partial \vec{h}^{\,\ell}}$$

# Example: Forward mode AD

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

$$\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}$$

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

$a = x_1$

$a = x_2$

$w_3 = f$

$$\frac{\partial w_3}{\partial a} =$$

$\cos(x_1) \cdot \dot{x}_1$

$$\frac{\partial w_1}{\partial a} = \frac{\partial w_1}{\partial x_1} \cdot \frac{\partial x_1}{\partial a}$$

$w_1$

$w_2$

sin( )

\*

$\dot{x}_1$

$$\frac{\partial x_1}{\partial a}$$

$x_1$

$$\frac{\partial x_1}{\partial a}$$

$$\frac{\partial x_2}{\partial a}$$

$x_1$

$x_2$

# Example: Forward mode AD

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

$$\cos(x_1) + x_2 =$$

$$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$$

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

$$\dot{w}_1 = \cos(x_1)\dot{x}_1 \qquad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$$
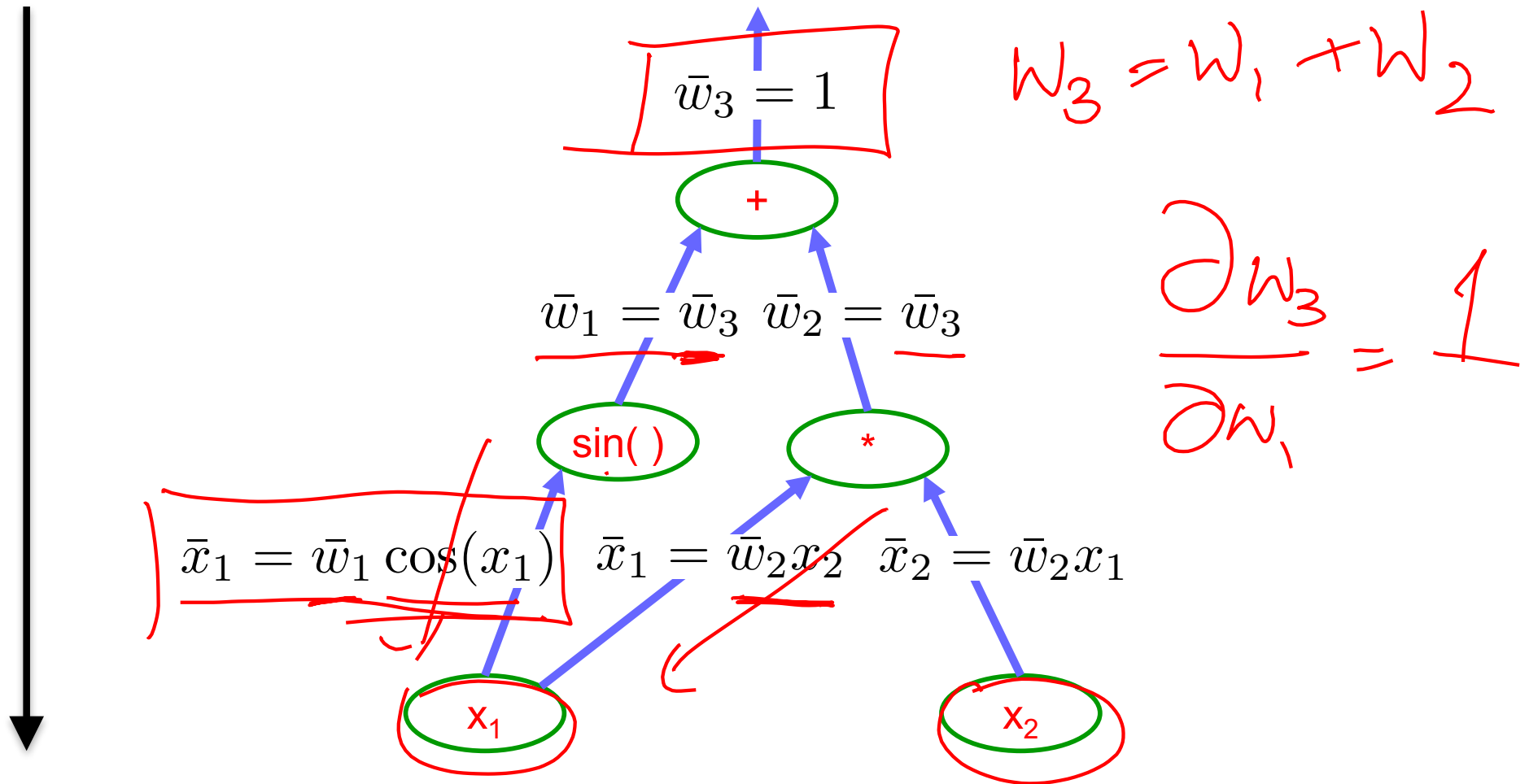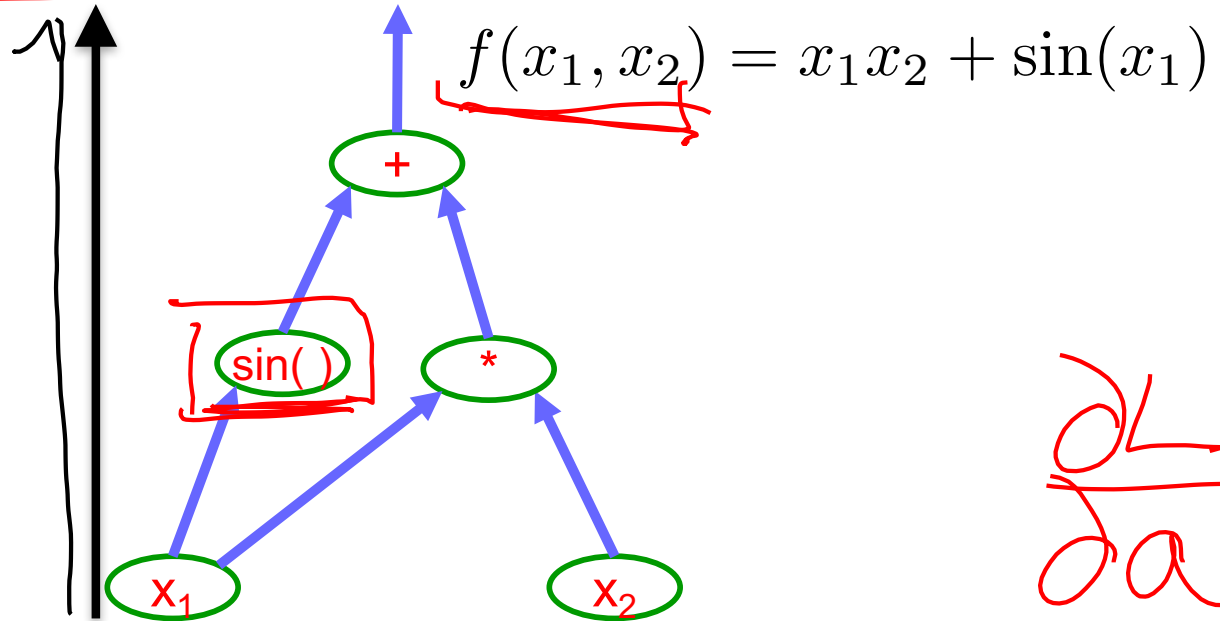
$$1 \qquad \qquad 1$$

sin( )     *

$$1 \qquad 1 \qquad 0$$

$$\dot{x}_1 \qquad \dot{x}_1 \qquad \dot{x}_2$$

$$1 \left| \frac{\partial x_1}{\partial a} \right. \qquad = \qquad = \left| \frac{\partial x_2}{\partial a} \right| 0$$

x_1     x_2

$$a = x_1$$

# Example: Forward mode AD

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



$$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$$

$$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$$

+

sin( )   *

$$\dot{x}_1 \qquad \dot{x}_1 \qquad \dot{x}_2$$

x₁   x₂

# Example: Reverse mode AD

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

# Example: Reverse mode AD

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



$$\bar{w}_3 = 1$$

$$W_3 = W_1 + W_2$$

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

$$\bar{w}_1 = \bar{w}_3 \quad \bar{w}_2 = \bar{w}_3$$

sin( )  *

$$\bar{x}_1 = \bar{w}_1 \cos(x_1) \quad \bar{x}_1 = \bar{w}_2 x_2 \quad \bar{x}_2 = \bar{w}_2 x_1$$

$x_1$  $x_2$

# Forward mode AD vs Reverse Mode AD

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

(center graph)
- + node
- sin( ) node, * node
- $x_1$, $x_2$

(left graph — Forward mode AD)

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

$$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$$

+ node

$$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$$

sin( ) node, * node

$$\dot{x}_1 \quad \dot{x}_1 \quad \dot{x}_2$$

$x_1$, $x_2$

(right graph — Reverse mode AD)

$$\frac{\partial L}{\partial a}$$

$$\bar{w}_3 = 1$$

+ node

$$\bar{w}_1 = \bar{w}_3 \quad \bar{w}_2 = \bar{w}_3$$

sin( ) node, * node

$$\bar{x}_1 = \bar{w}_1 \cos(x_1) \quad \bar{x}_1 = \bar{w}_2 x_2 \quad \bar{x}_2 = \bar{w}_2 x_1$$

$x_1$, $x_2$

# Forward mode vs Reverse Mode

- What are the differences?

- Which one is more memory efficient (less storage)?
  - Forward or backward?

# Forward mode vs Reverse Mode

- What are the differences?

- Which one is more memory efficient (less storage)?
  - Forward or backward?

- Which one is faster to compute?
  - Forward or backward?

# Plan for Today

- (Finish) Computing Gradients
  - Forward mode vs Reverse mode AD
  - Patterns in backprop
  - Backprop in FC+ReLU NNs

- Convolutional Neural Networks

# Patterns in backward flow

$$f(\cdots) = 2\left(xy + \max\{z,w\}\right)$$

# Patterns in backward flow



$$W_3 = N_1 + W_2$$

$$\frac{\partial W_3}{\partial W_1} = 1$$

$$W_1 = x\,y$$

$$\frac{\partial W_1}{\partial x} = y$$

$$f = 2\,W_3$$

$$\frac{\partial f}{\partial W_3} = 2$$

$W_1$

$W_3$

$W_2$

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

x  3.00
-8.00

-12.00
2.00

y  -4.00
6.00

z  2.00
2.00

w  -1.00
0.00

-10.00
2.00

-20.00
1.00

# Patterns in backward flow

$$w_3 = w_1 + w_2$$

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

**add** gate: gradient distributor



x  3.00
   -8.00

y  -4.00
   6.00

-12.00
2.00

z  2.00
   2.00

w  -1.00
   0.00

max  2.00
     2.00

+  -10.00
   2.00

*2  -20.00
    1.00

# Patterns in backward flow

**add** gate: gradient distributor

Q: What is a **max** gate?

# 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?



$W_1 = xy$    $\left[\dfrac{\partial f}{\partial W_1}\right] = 1$

# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router

**mul** gate: gradient switcher

# Gradients add at branches

# Duality in Fprop and Bprop



*FPROP*          *BPROP*

SUM

$W = x + y$

COPY

$x$

$\dfrac{\partial f}{\partial w}$

$\dfrac{\partial f}{\partial x}$

$\dfrac{\partial f}{\partial x}$

# Modularized implementation: forward / backward API

Graph (or Net) object  *(rough psuedo code)*

```python
class ComputationalGraph(object):
    #...
    def forward(inputs):
        # 1. [pass inputs to input gates...]
        # 2. forward the computational graph:
        for gate in self.graph.nodes_topologically_sorted():
            gate.forward()
        return loss # the final gate in the graph outputs the loss
    def backward():
        for gate in reversed(self.graph.nodes_topologically_sorted()):
            gate.backward() # little piece of backprop (chain rule applied)
        return inputs_gradients
```

# Modularized implementation: forward / backward API



x
z
*
y

(x,y,z are scalars)

```
class MultiplyGate(object):
    def forward(x,y):
        z = x*y
        return z
    def backward(dz):
        # dx = ... #todo
        # dy = ... #todo
        return [dx, dy]
```

$$\frac{\partial L}{\partial z}$$

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

# Modularized implementation: forward / backward API



x

z

*

y

(x,y,z are scalars)

```
class MultiplyGate(object):
    def forward(x,y):
        z = x*y
        self.x = x # must keep these around!
        self.y = y
        return z

    def backward(dz):
        dx = self.y * dz # [dz/dx * dL/dz]
        dy = self.x * dz # [dz/dy * dL/dz]
        return [dx, dy]
```

# Example: Caffe layers

# Caffe Sigmoid Layer

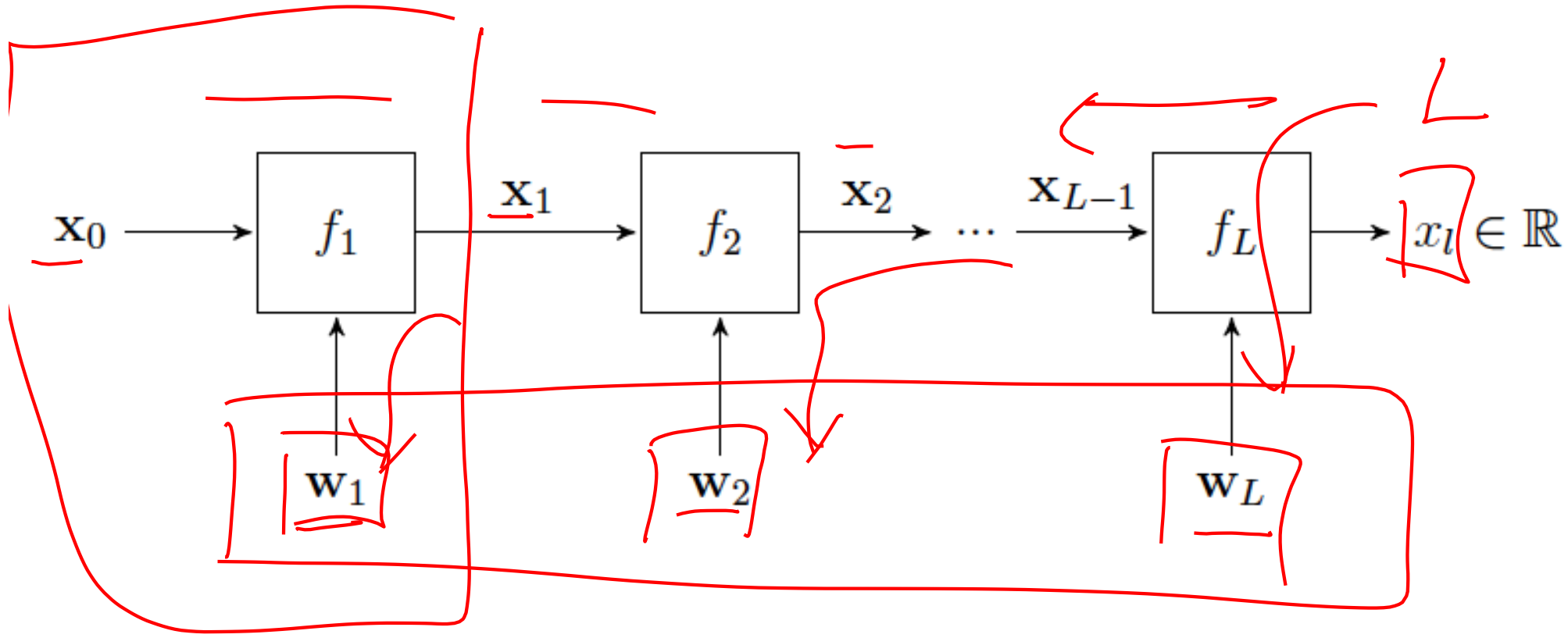```cpp
1    #include <cmath>
2    #include <vector>
3
4    #include "caffe/layers/sigmoid_layer.hpp"
5
6    namespace caffe {
7
8    template <typename Dtype>
9    inline Dtype sigmoid(Dtype x) {
10     return 1. / (1. + exp(-x));
11   }
12
13   template <typename Dtype>
14   void SigmoidLayer<Dtype>::Forward_cpu(const vector<Blob<Dtype>*>& bottom,
15       const vector<Blob<Dtype>*>& top) {
16     const Dtype* bottom_data = bottom[0]->cpu_data();
17     Dtype* top_data = top[0]->mutable_cpu_data();
18     const int count = bottom[0]->count();
19     for (int i = 0; i < count; ++i) {
20       top_data[i] = sigmoid(bottom_data[i]);
21     }
22   }
23
24   template <typename Dtype>
25   void SigmoidLayer<Dtype>::Backward_cpu(const vector<Blob<Dtype>*>& top,
26       const vector<bool>& propagate_down,
27       const vector<Blob<Dtype>*>& bottom) {
28     if (propagate_down[0]) {
29       const Dtype* top_data = top[0]->cpu_data();
30       const Dtype* top_diff = top[0]->cpu_diff();
31       Dtype* bottom_diff = bottom[0]->mutable_cpu_diff();
32       const int count = bottom[0]->count();
33       for (int i = 0; i < count; ++i) {
34         const Dtype sigmoid_x = top_data[i];
35         bottom_diff[i] = top_diff[i] * sigmoid_x * (1. - sigmoid_x);
36       }
37     }
38   }
39
40   #ifdef CPU_ONLY
41   STUB_GPU(SigmoidLayer);
42   #endif
43
44   INSTANTIATE_CLASS(SigmoidLayer);
45
46
47   }  // namespace caffe
```
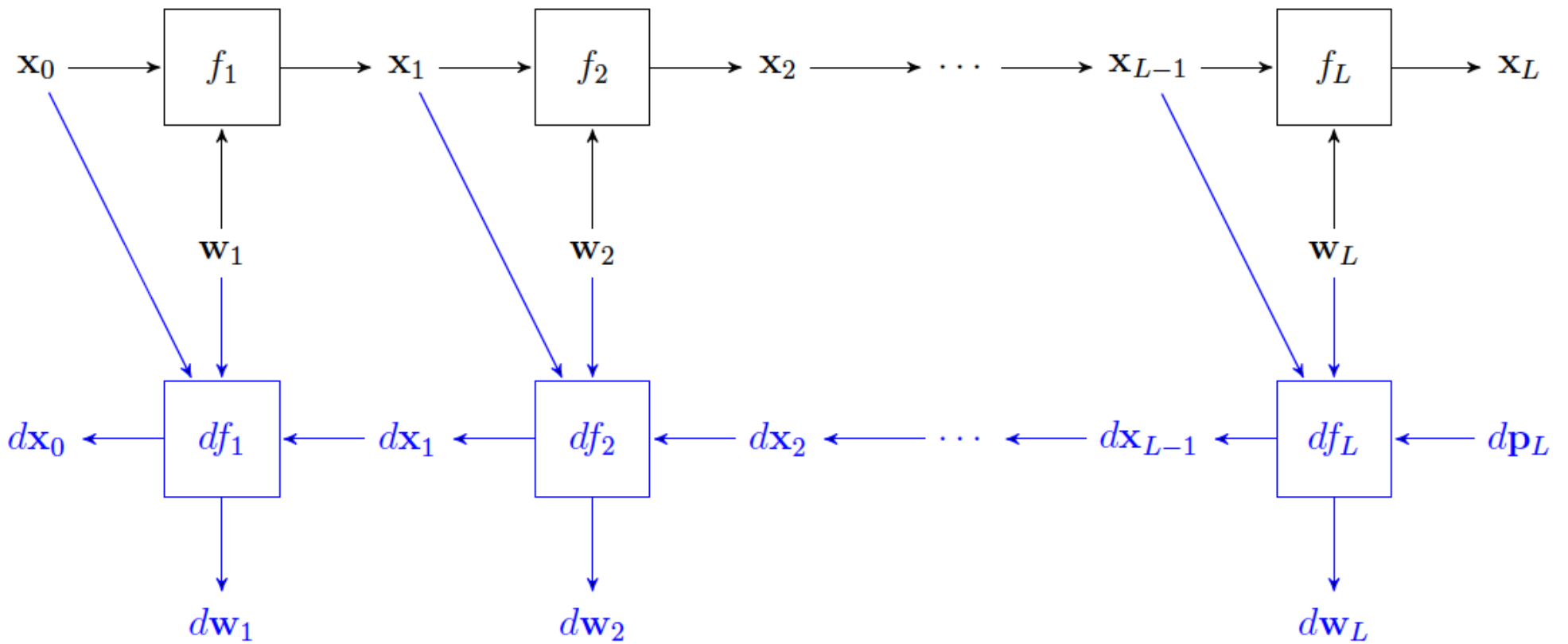
Caffe is licensed under BSD 2-Clause

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$(1 - \sigma(x))\,\sigma(x)$$ * top_diff  (chain rule)

$$\frac{\partial L}{\partial w_l}$$

$$\mathbf{x}_0 \longrightarrow \boxed{f_1} \xrightarrow{\ \mathbf{x}_1\ } \boxed{f_2} \xrightarrow{\ \mathbf{x}_2\ } \cdots \xrightarrow{\ \mathbf{x}_{L-1}\ } \boxed{f_L} \longrightarrow x_l \in \mathbb{R}$$

$\mathbf{w}_1 \qquad \mathbf{w}_2 \qquad \mathbf{w}_L$

# Key Computation in DL: Forward-Prop

Slide Credit: Marc'Aurelio Ranzato, Yann LeCun

# Jacobian of ReLU

$\vec{h}^{\ell}$

$\vec{h}^{(\ell+1)}$

$h_i^{\ell+1} = \max\{0, h_i^{\ell}\}$

4096-d
input vector

f(x) = max(0,x)
*(elementwise)*

4096-d
output vector

$$\frac{\partial h_i^{\ell+1}}{\partial h_j^{\ell}} = 0 \quad i \neq j$$

# Jacobian of ReLU

4096-d
input vector

$f(x) = \max(0, x)$
*(elementwise)*

4096-d
output vector

Q: what is the
size of the
Jacobian matrix?

+1

$\frac{\partial L}{\partial \vec{h}^{\ell+1}}$

4096 × 4096

# Jacobian of ReLU

4096-d
input vector

$f(x) = \max(0,x)$
*(elementwise)*

4096-d
output vector

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

# Jacobian of ReLU



4096-d
input vector

$f(x) = max(0,x)$
*(elementwise)*

4096-d
output vector

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

$f(x) = \max(0,x)$
*(elementwise)*

4096-d output vector

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

Q2: what does it look like?

# Jacobians of FC-Layer

# Jacobians of FC-Layer

$$\vec{h}^{\ell-1} \longrightarrow \boxed{\phantom{xxx}} \longrightarrow \vec{h}^{\ell} \in \mathbb{R}^{C_2}$$

$$\boxed{\frac{\partial \vec{h}^{\ell}}{\partial W}}$$

$$W = C_2 \times C_1$$

$*(j,k)$

$C_2 \times (C_1 C_2)$

$$\frac{\partial \vec{h}_i^{\ell}}{\partial W_{j,k}}$$

$$\frac{\partial h_j^{\ell}}{\partial W}$$

$$\left[\; W_{j,k}^T \;\right]$$

$$h_j^{\ell} = \vec{W}_j^{\,T} \vec{h}^{\ell-1}$$

$$\frac{\partial h_j^{\ell}}{\partial \vec{W}_j} = \left(\vec{h}^{\ell-1}\right)^T$$

# Jacobians of FC-Layer

# Convolutional Neural Networks

## (without the brain stuff)

# Fully Connected Layer $4 \times 10^4$

Example: 200x200 image
40K hidden units

➡ **~2B parameters**!!!



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

# Locally Connected Layer



Example: 200x200 image
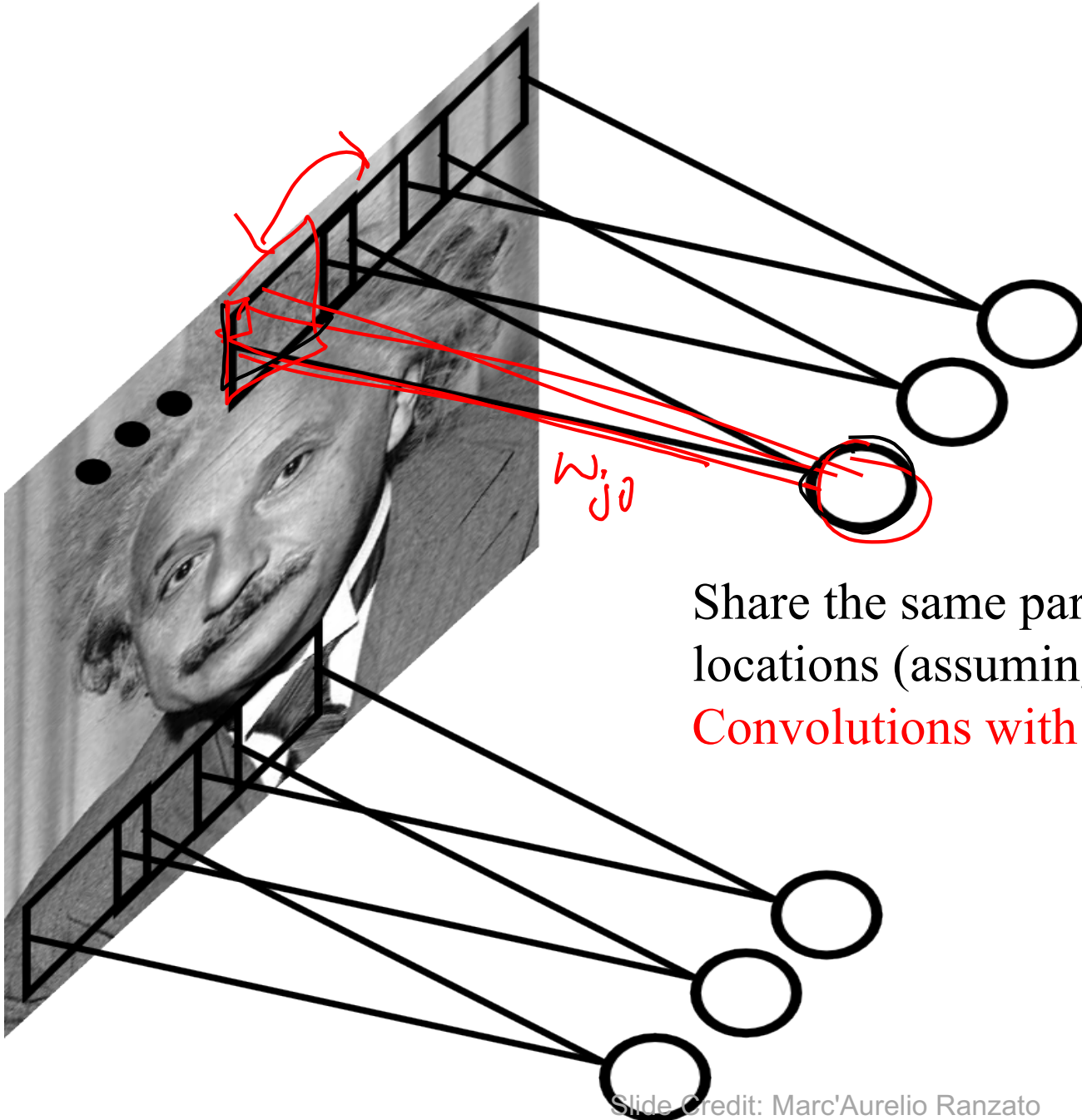40K hidden units
Filter size: 10x10
4M parameters

**Note:** This parameterization is good when input image is registered (e.g., face recognition).

# Locally Connected Layer



**STATIONARITY?** Statistics is similar at different locations

Example: 200x200 image
40K hidden units
Filter size: 10x10
4M parameters

**Note:** This parameterization is good when input image is registered (e.g., face recognition).

# Convolutional Layer



Share the same parameters across different locations (assuming input is stationary):
Convolutions with learned kernels

$w_{j\partial}$

# Convolutions for mathematicians

$$x(t) \quad y(t) \quad w(t)$$

$$y(t) = (x * w)(t) = \int_{a=-\infty}^{\infty} x(t-a) \, w(a) \, da$$

$$= \int x(a) \, w(t-a) \, da$$

$$w(a) \to w(-a)$$

$$w(-a) \to w(-(a-t))$$

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



"Convolution of box signal with itself2" by Convolution_of_box_signal_with_itself.gif: Brian Ambergderivative 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

$$y[r, c] = \sum_{a=-\frac{N_1-1}{2}}^{\frac{N-1}{2}} \sum_{b=-\frac{N-1}{2}}^{+\frac{N-1}{2}} x[r-a, c-b] \, w[a,b]$$

$N_2$

$N_1$

# Convolutions for programmers

$$y[r, c] = \sum_{a=0}^{N_1-1} \sum_{b=0}^{N_2-1} x[r+a, c+b] \, w[a,b]$$

# Convolution Explained

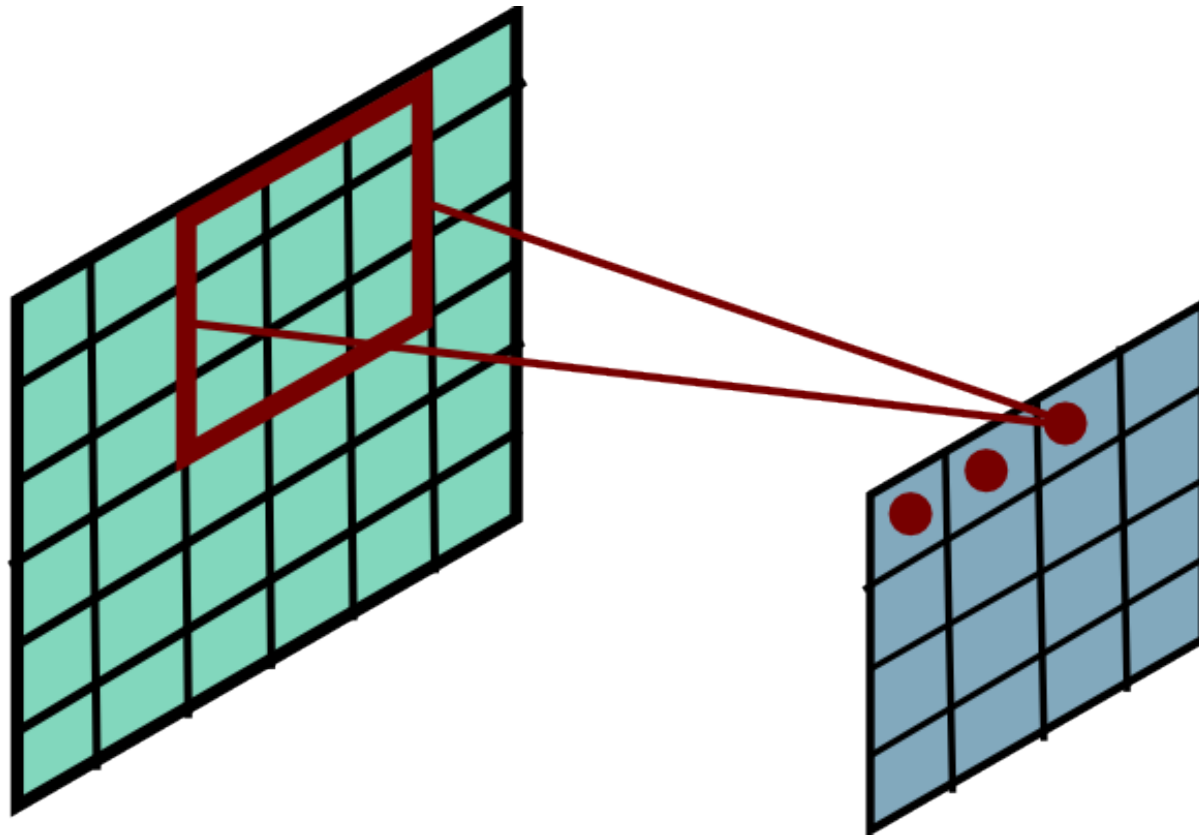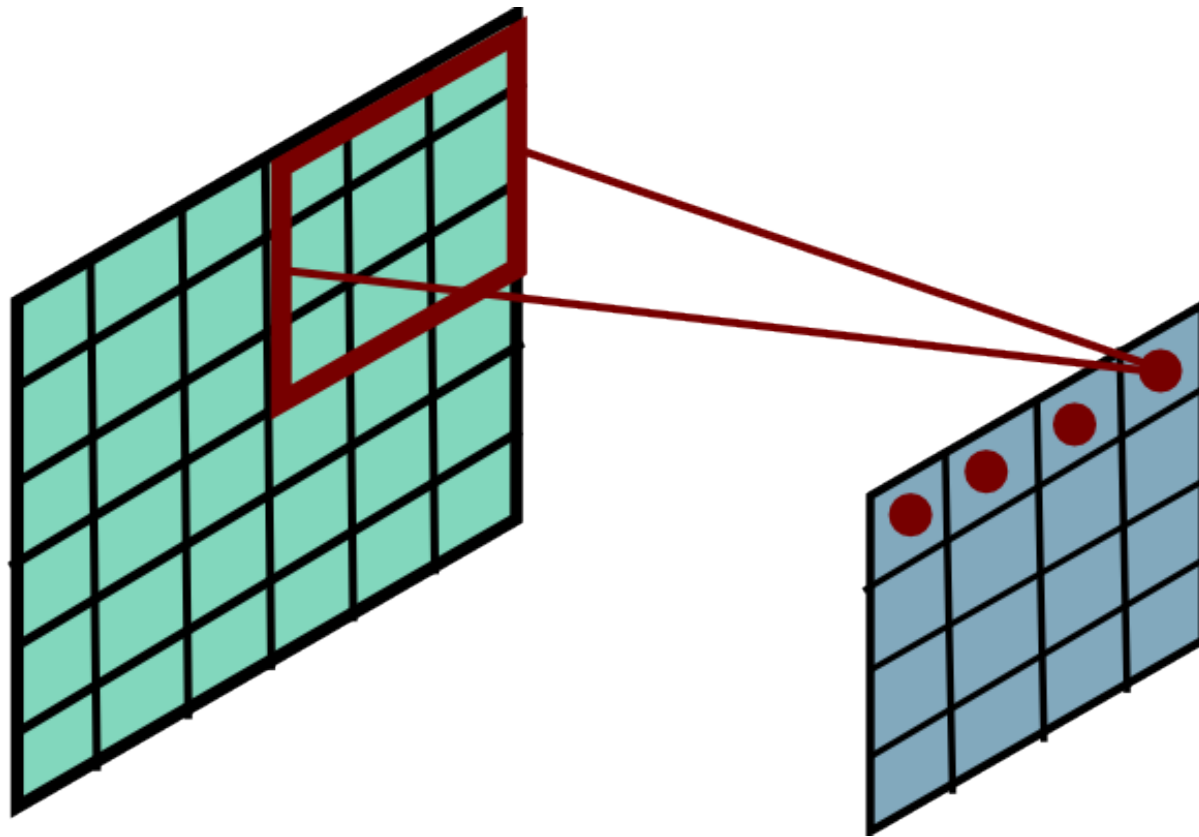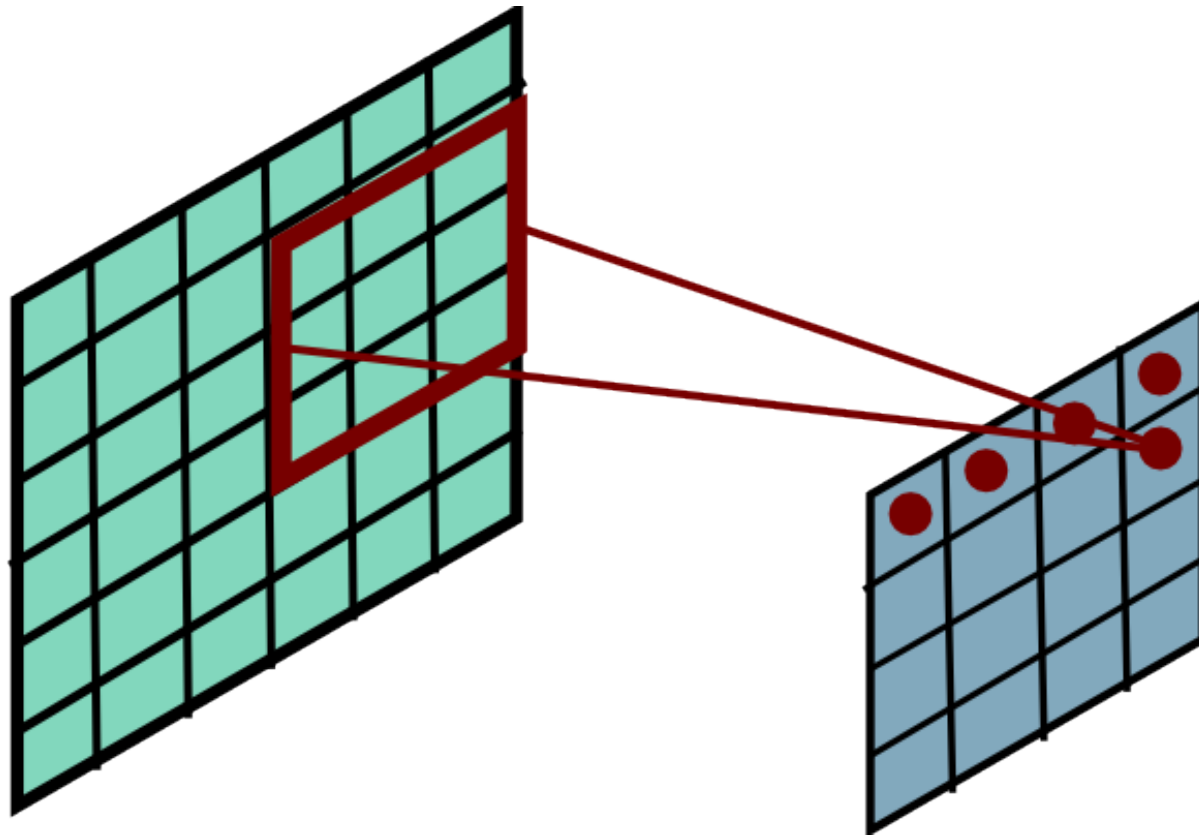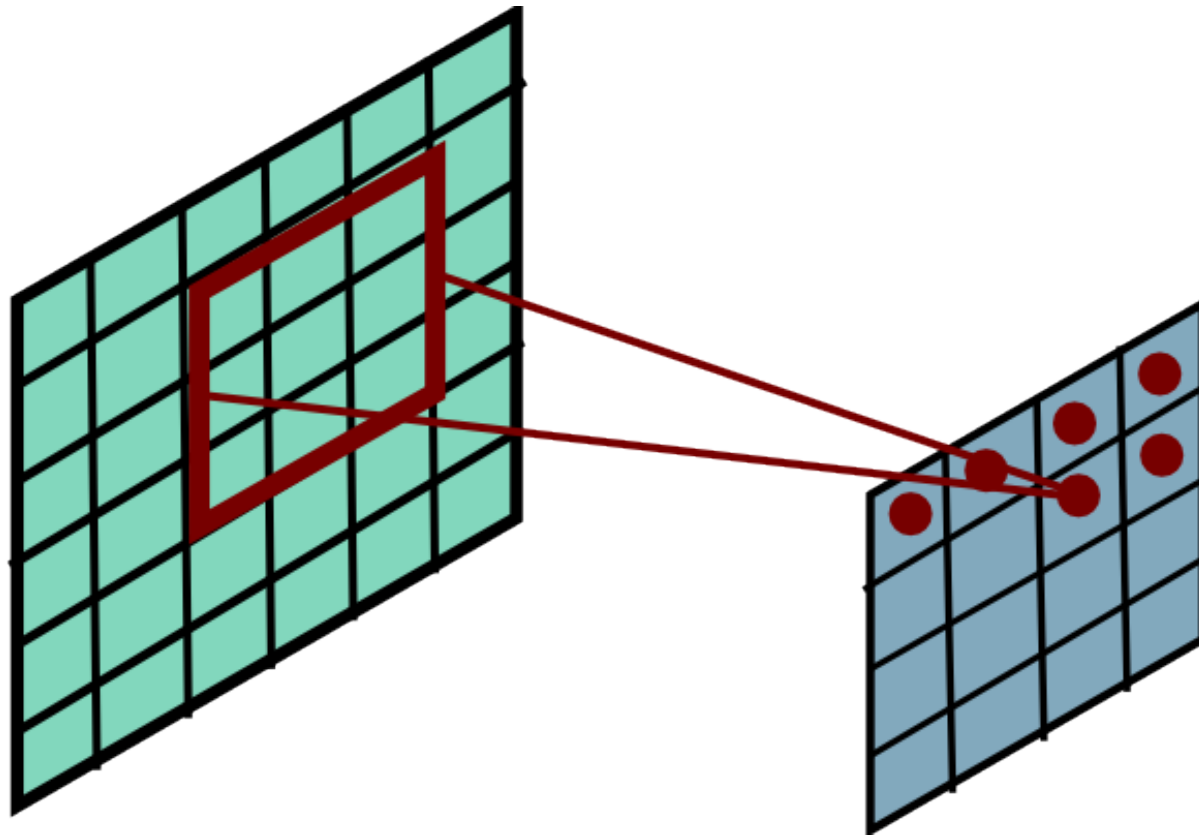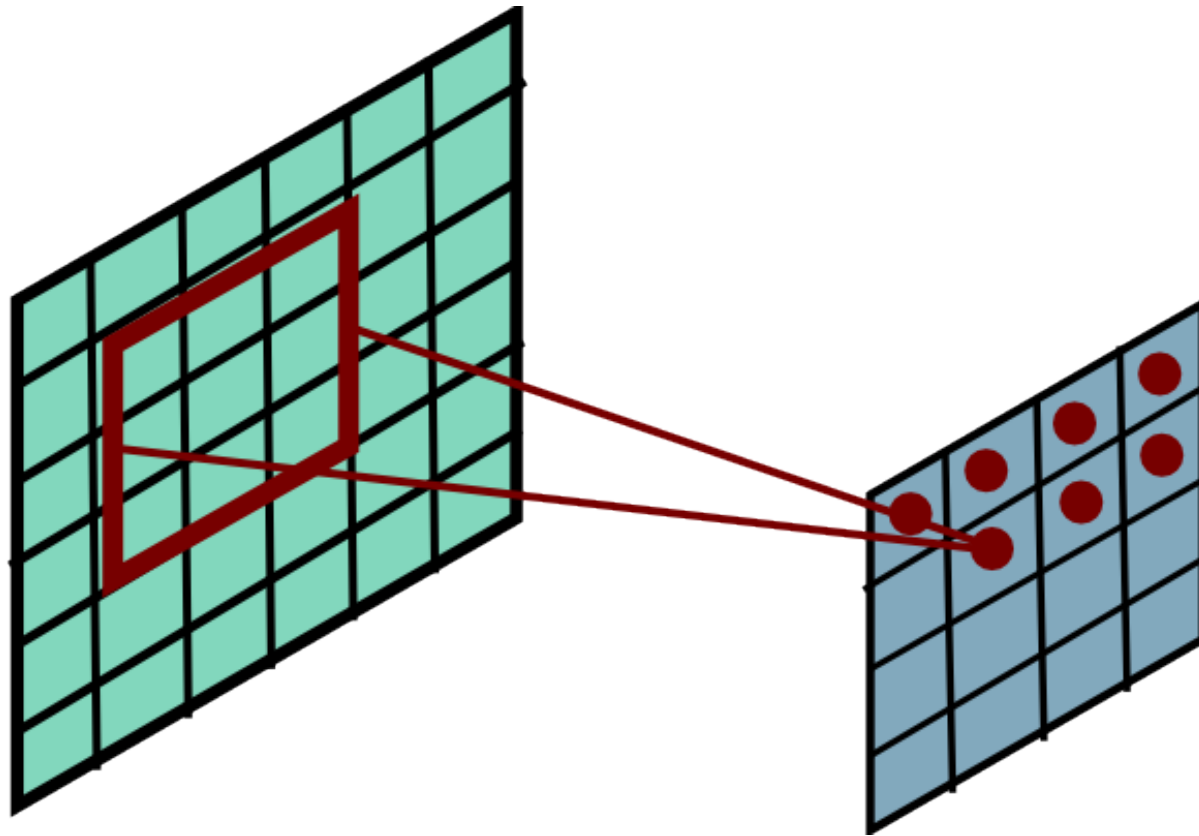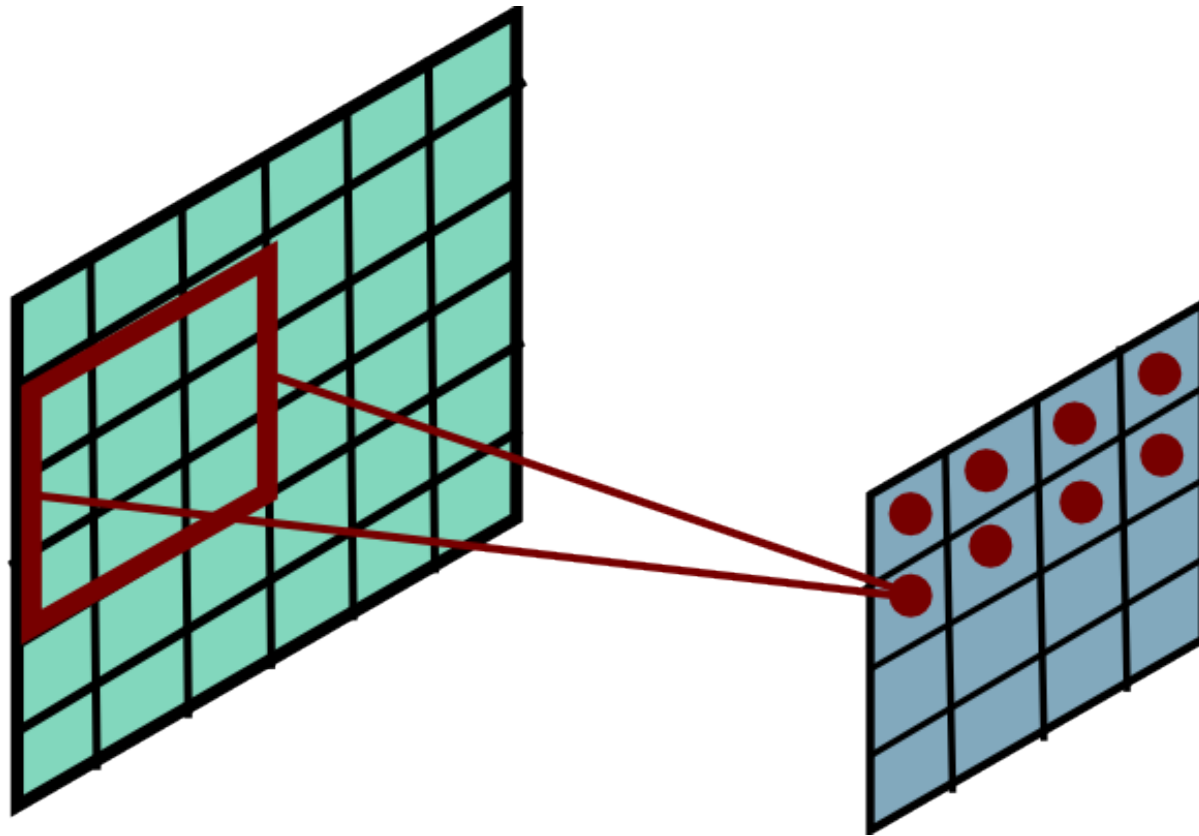- http://setosa.io/ev/image-kernels/
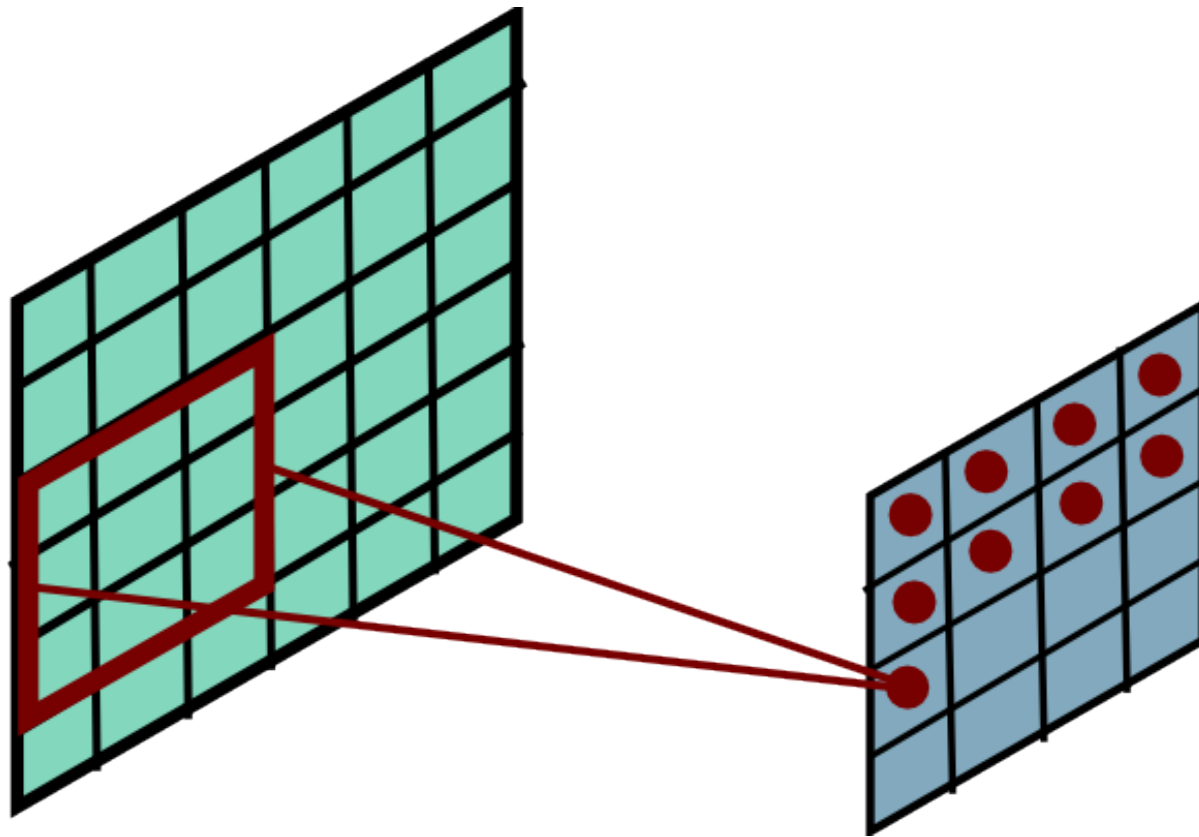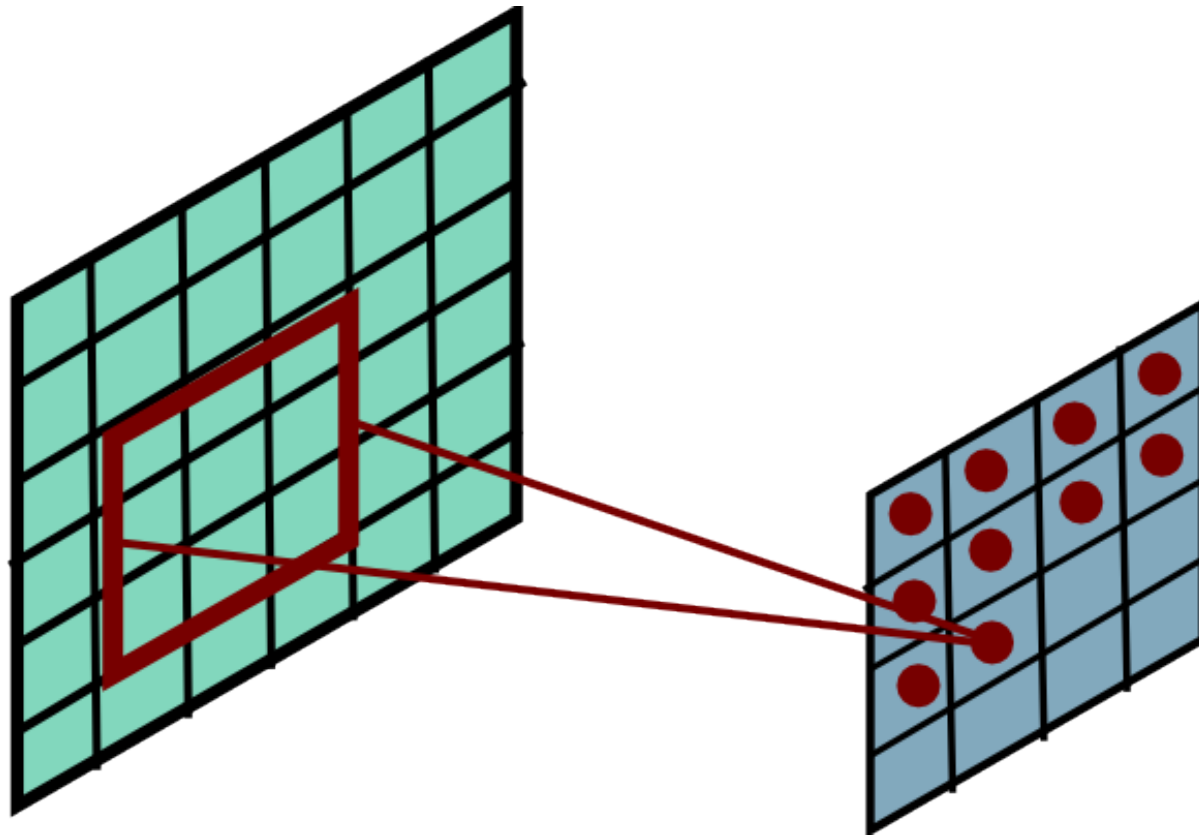
- https://github.com/bruckner/deepViz

# Convolutional Layer

# Convolutional Layer

# Convolutional Layer

# Convolutional Layer

# Convolutional Layer

# Convolutional Layer

# Convolutional Layer

# Convolutional Layer

# Convolutional Layer

Slide Credit: Marc'Aurelio Ranzato

# Convolutional Layer
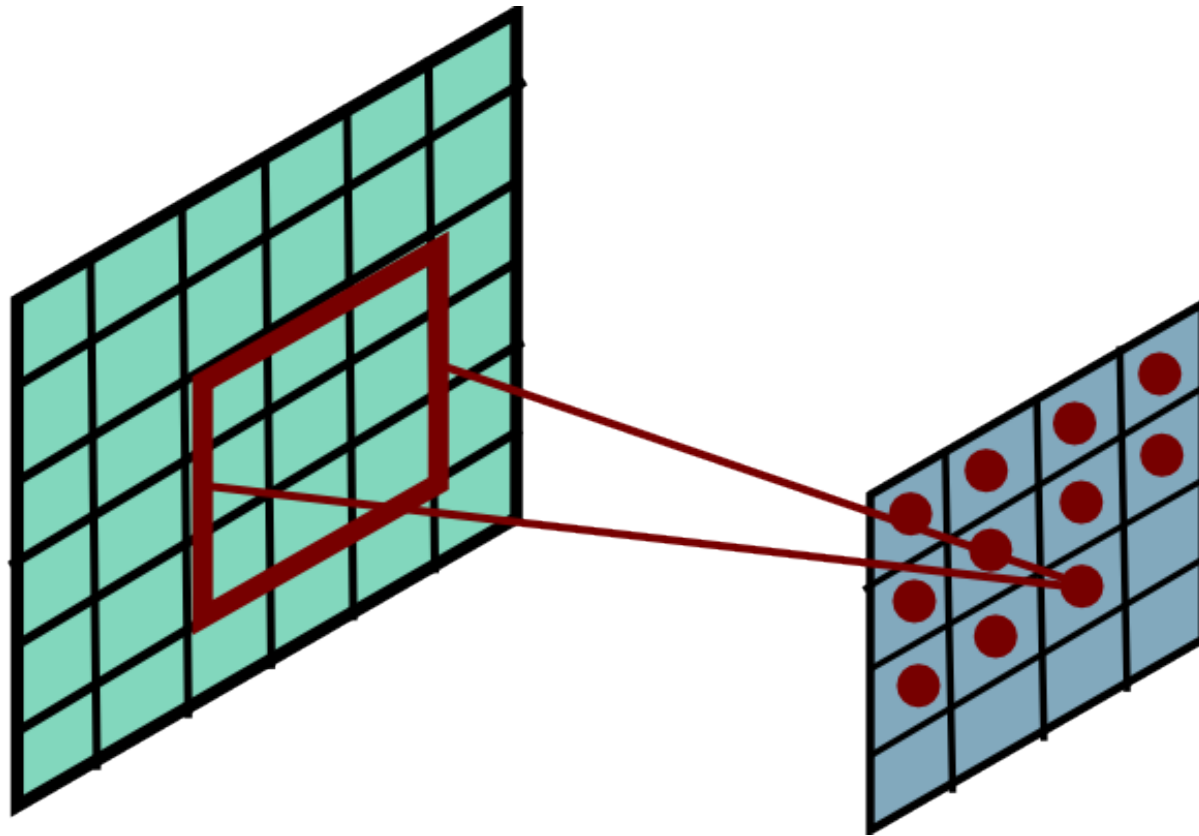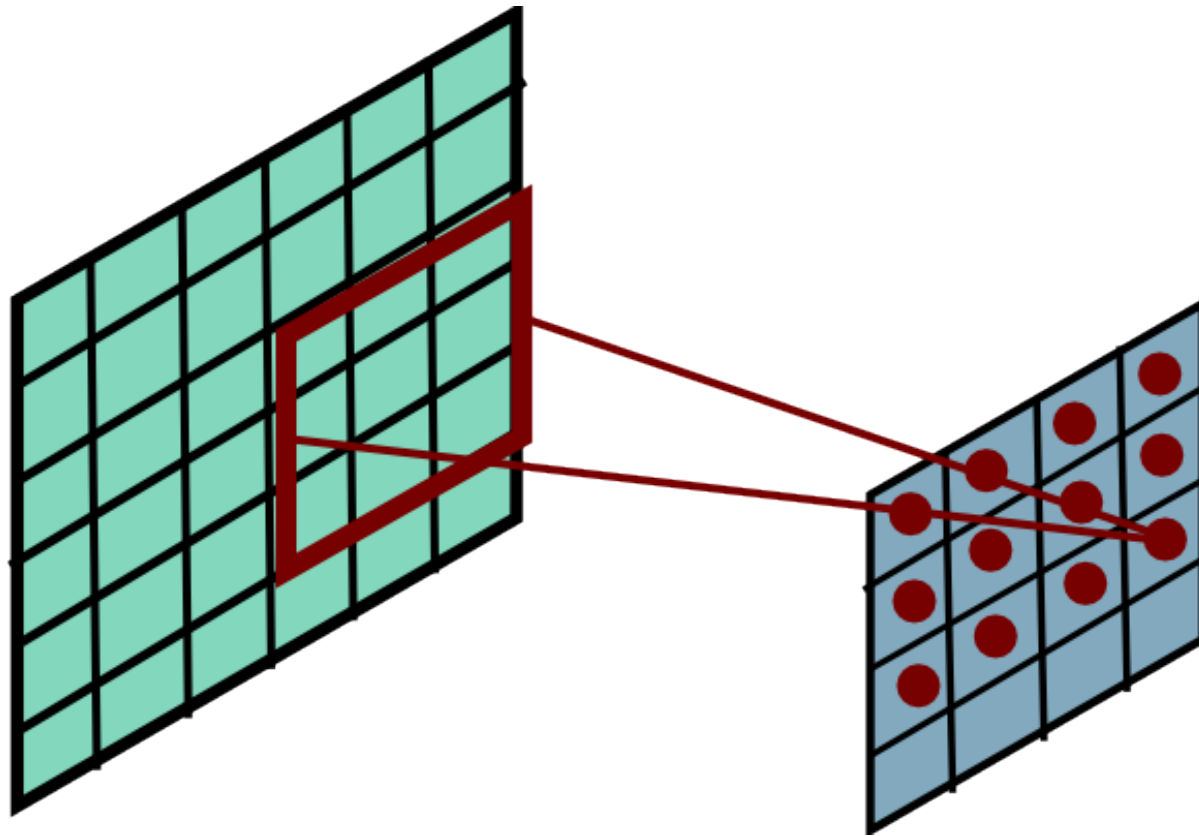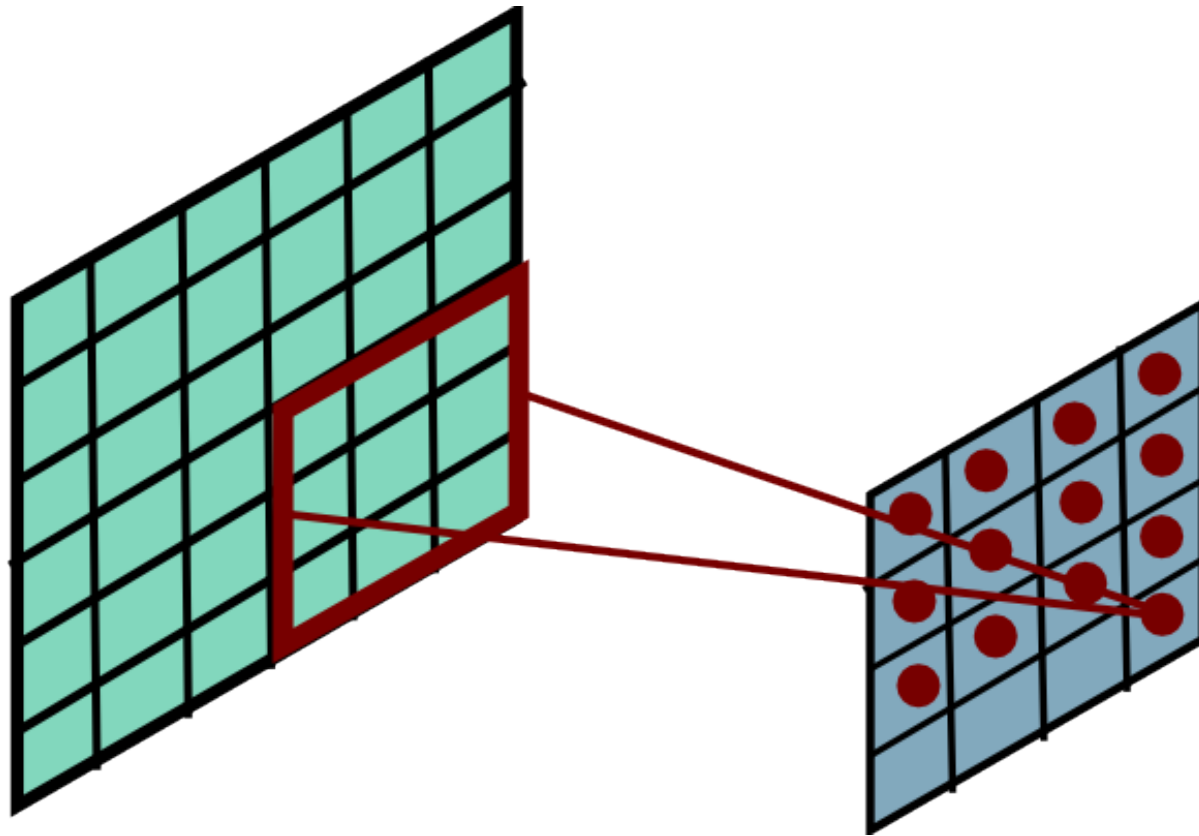
Slide Credit: Marc'Aurelio Ranzato

# Convolutional Layer

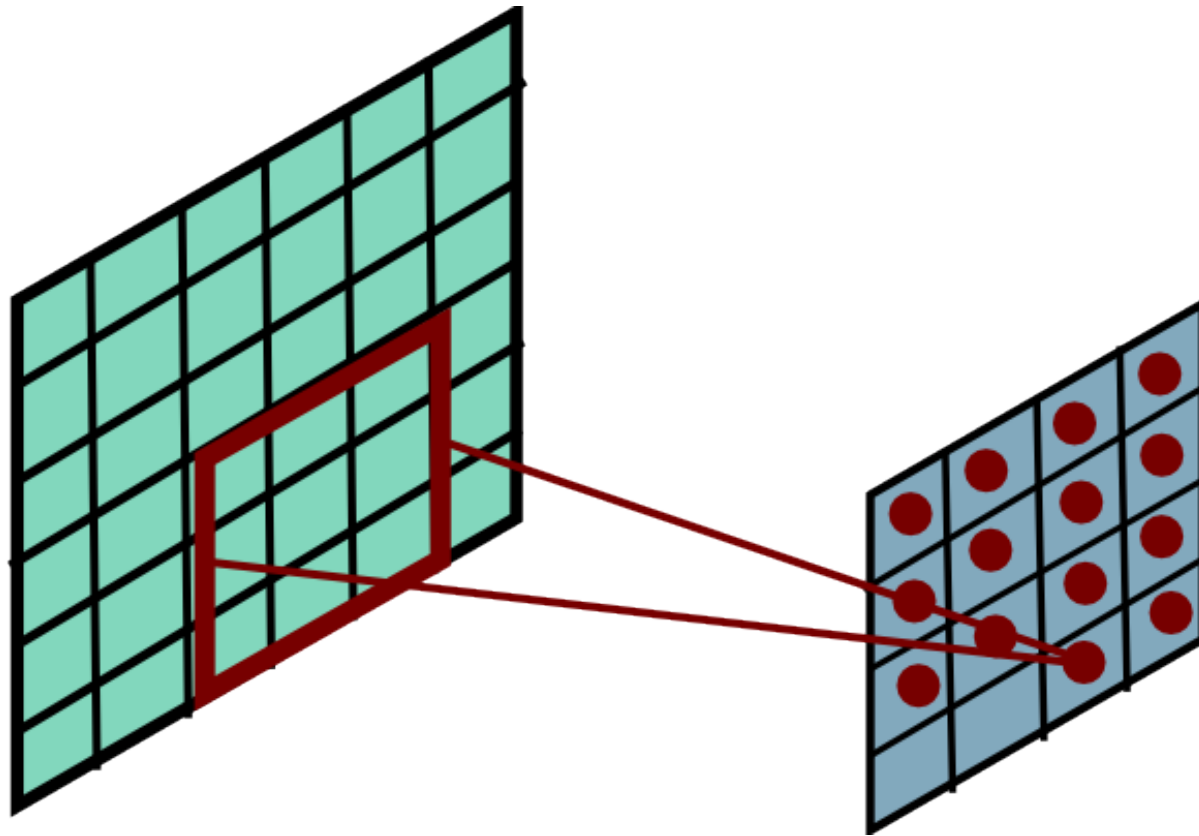# Convolutional Layer

# Convolutional Layer

# Convolutional Layer

# Convolutional Layer
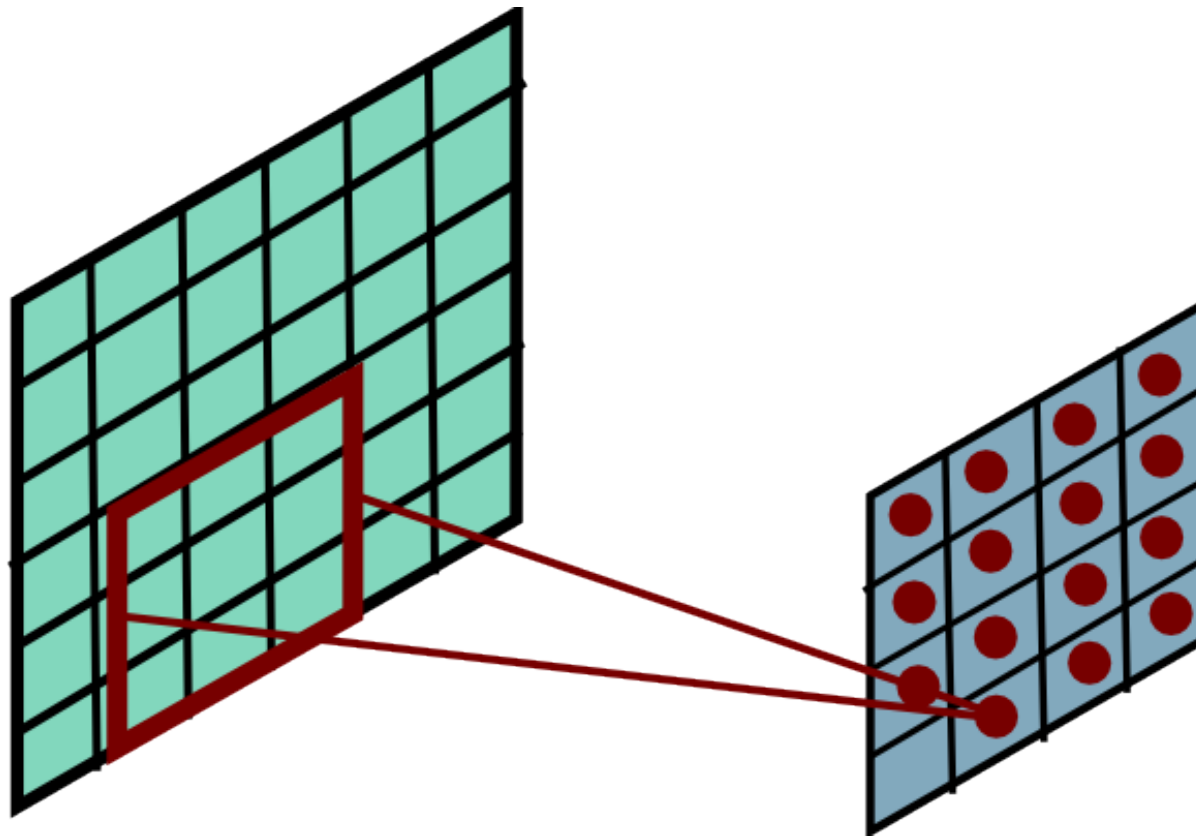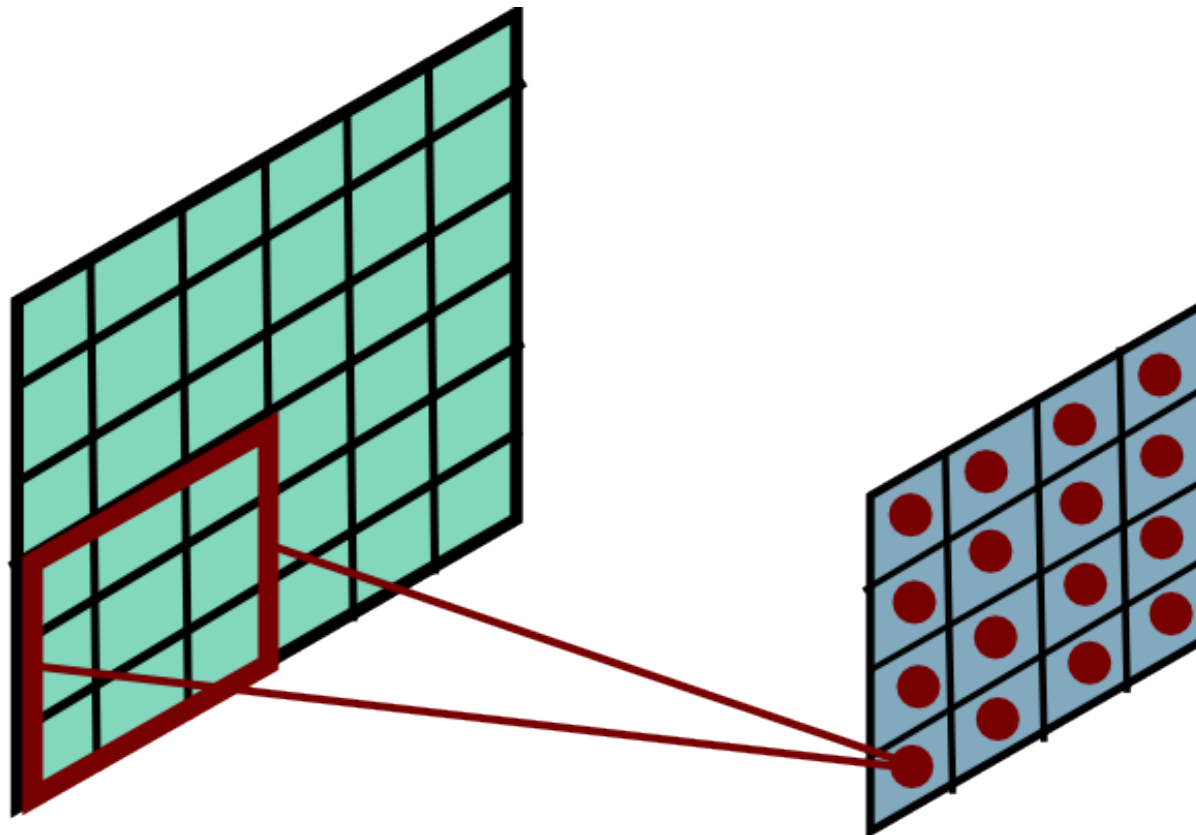


Slide Credit: Marc'Aurelio Ranzato

# Convolutional Layer



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

# Convolutional Layer



$$\ast \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix} =$$

Slide Credit: Marc'Aurelio Ranzato

# Convolutional Layer



**Learn** multiple filters.

E.g.: 200x200 image
100 Filters
Filter size: 10x10
10K parameters

# Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1

**input**

1

3072

$$Wx$$

10 x 3072
weights

**activation**

1

10

# Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1

**input**



1

3072

$Wx$

10 x 3072
weights

**activation**

1

10

**1 number:**
the result of taking a dot product
between a row of W and the input
(a 3072-dimensional dot product)

# Convolutional Layer

# Convolutional Layer

# Convolution Layer

32x32x3 image -> preserve spatial structure



32 **height**

32 **width**

3 **depth**

# Convolution Layer

32x32x3 image

5x5x3 filter

32

32

3

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"

# Convolution Layer

32x32x3 image

Filters always extend the full depth of the input volume

5x5x3 filter

32

32

3

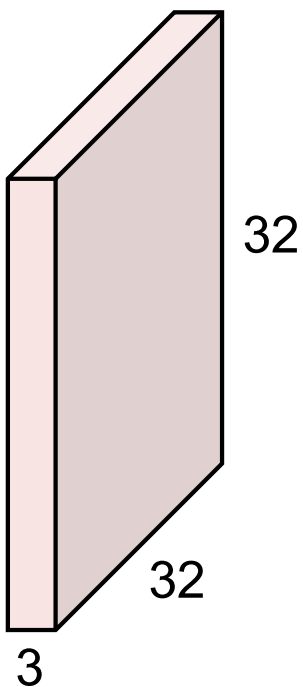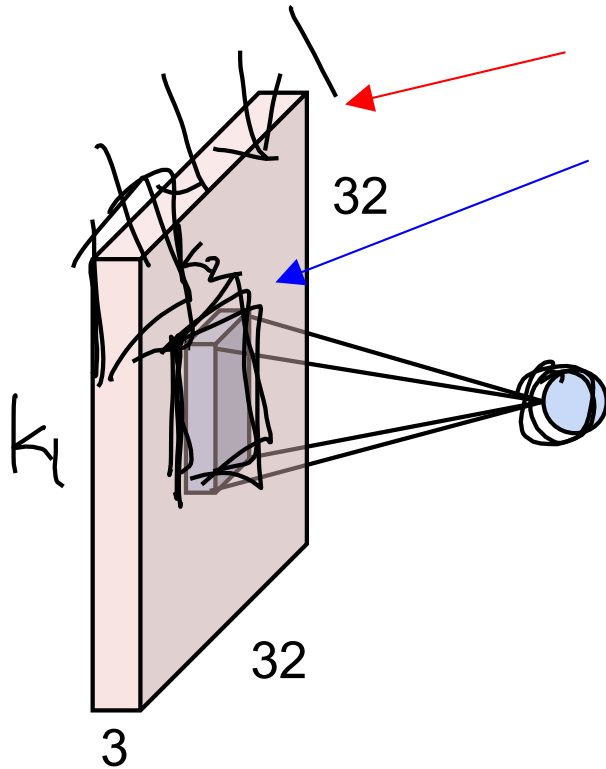**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"

# Convolution Layer

32x32x3 image

5x5x3 filter $w$

$1 \times (k_1 k_2 3)$

$k k 3$

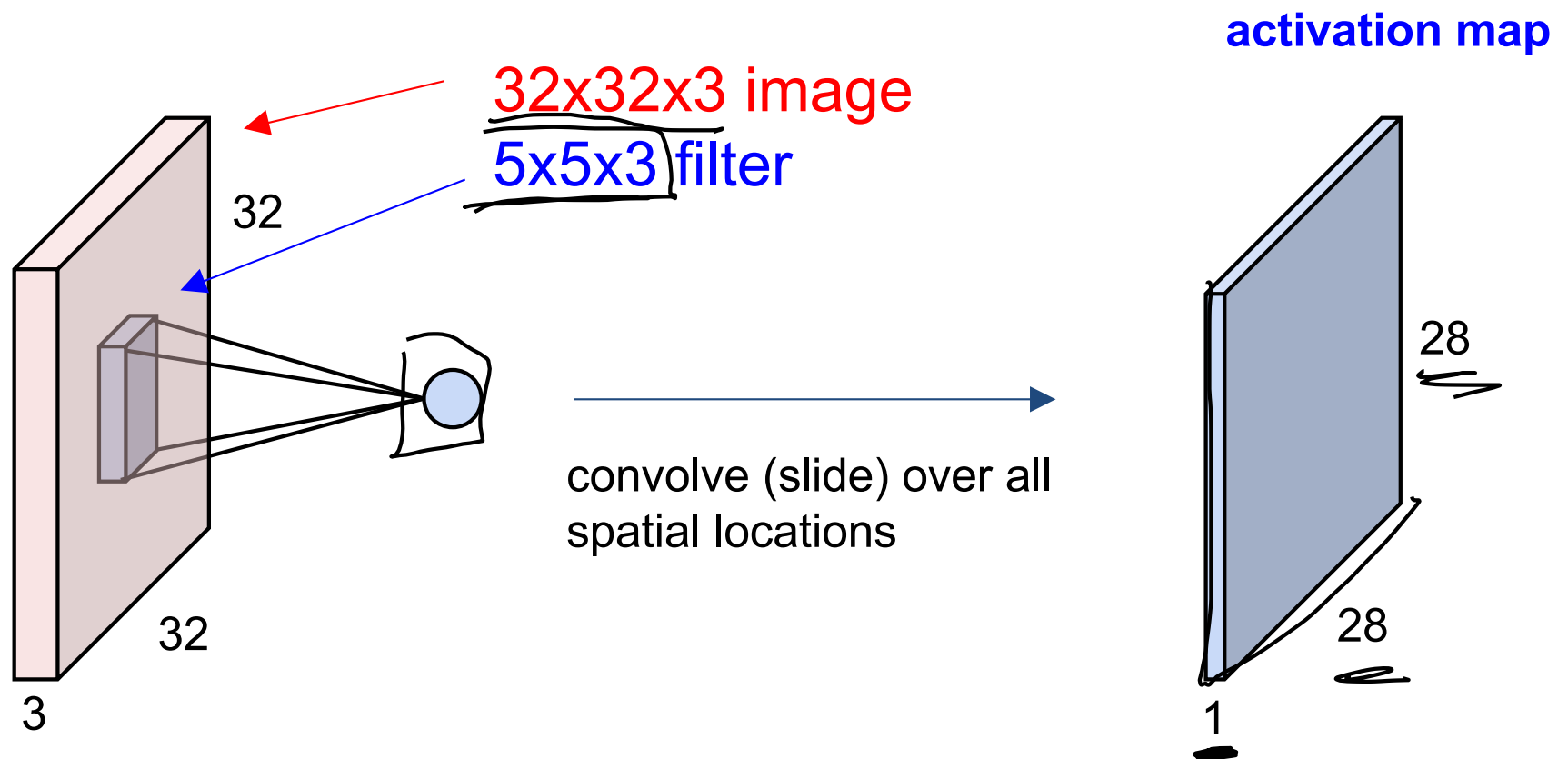**1 number:**
the result of taking a dot product between the filter and a small 5x5x3 chunk of the image
(i.e. 5*5*3 = 75-dimensional dot product + bias)

$$w^T x + b$$

32

32

3

# Convolution Layer



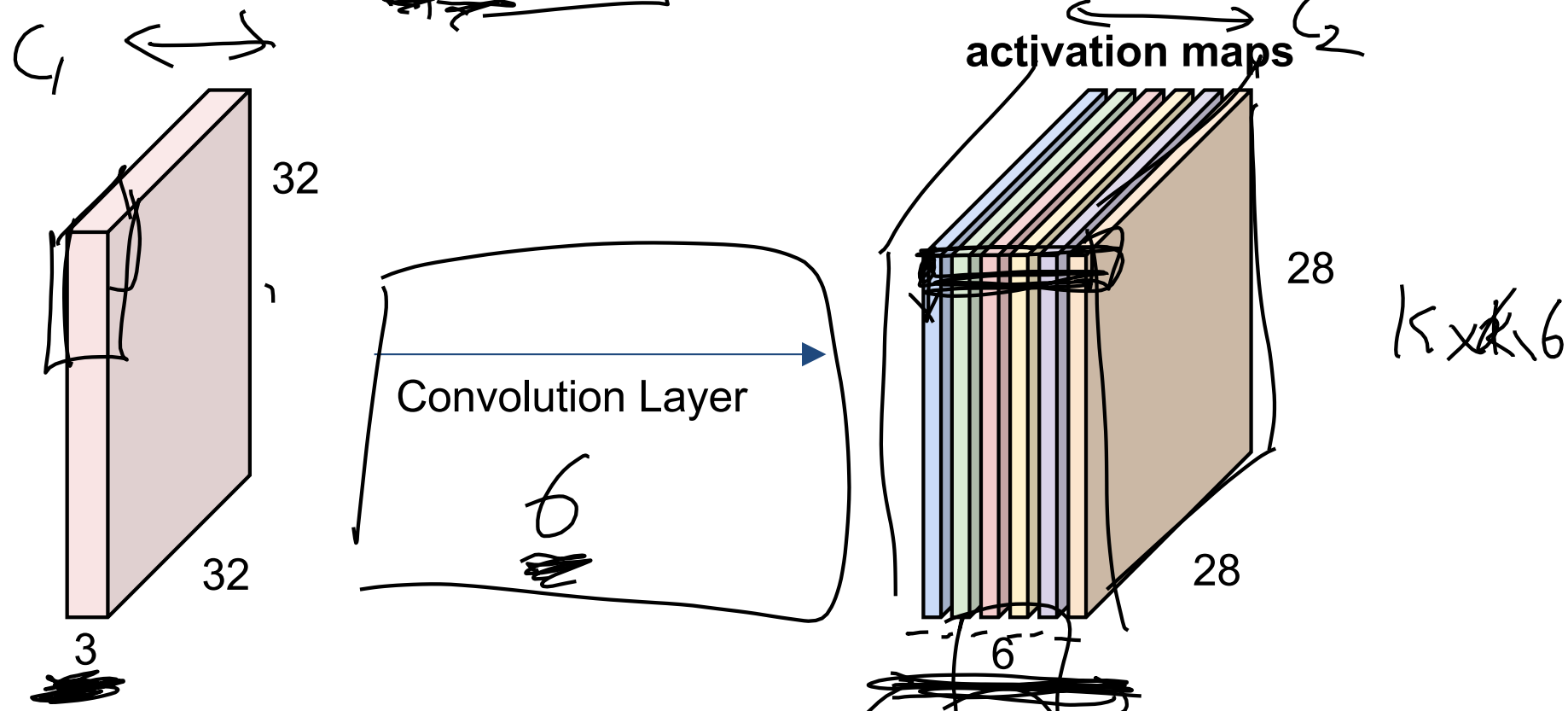**activation map**

32x32x3 image

5x5x3 filter

32

32

3

convolve (slide) over all spatial locations

28

28

1

# Convolution Layer

consider a second, green filter



32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all spatial locations

**activation maps**

28

28

1

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!