

MULTI-LAYER NETWORKS AND BACK PROPAGATION

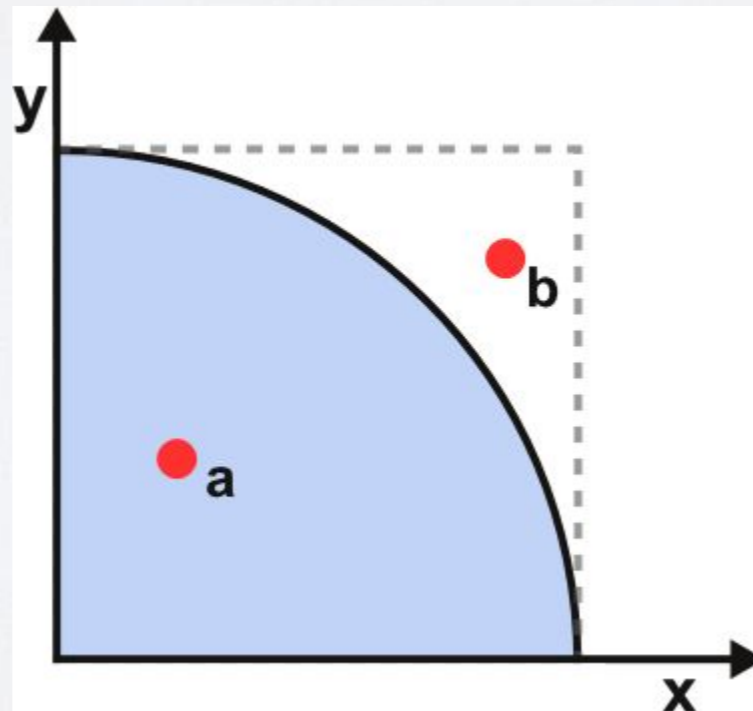
Philipp Krähenbühl

ADMINISTRATIVE

- Prof. Krähenbühl out of town this week
- HW 1
- HW 2

HW 1

- Given $k \{x_i\} \sim U[0, 1]^2$, estimate π



HW 1

- Naive estimator

- $U(x,y) = I(x^2+y^2 < 1)$, using k samples

- $E(U) = \pi/4$, $E(\bar{U}) = \pi/4$

- ```
def forward(self, x):
 return 4*torch.mean((torch.sum(x**2, dim=1)<1).float(), dim=0)
```

# HW 1

- Better estimator
  - $U(x) = \sqrt{1-x^2}$ , using 2k samples
  - Still does monte-carlo estimate, but integrates one dim in close form

- ```
def forward(self, x):  
    return 4*torch.mean(torch.sqrt(1-x**2))
```

HW 1

- Another interpretation
 - $E_{xy}(I\{x^2+y^2<1\}) = E_x E_y(I\{x^2+y^2<1\} | x) = \pi/4$
 - $I\{x^2+y^2<1\} \Leftrightarrow I\{y<\text{sqrt}(1-x^2)\}$
 - $E_y(I\{y<\text{sqrt}(1-x^2)\}|x) = \text{sqrt}(1-x^2)$

HW 1

- Another interpretation
 - $E_x E_y(I_{\{x^2+y^2 < 1\}} | x) = E_x(\text{sqrt}(1-x^2)) = \pi/4$
 - $E(U) = \pi/4$

HW I

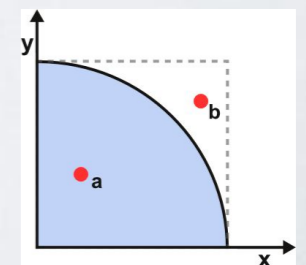
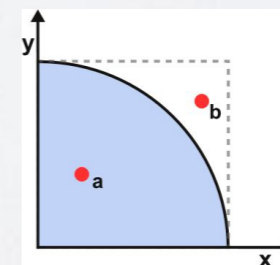
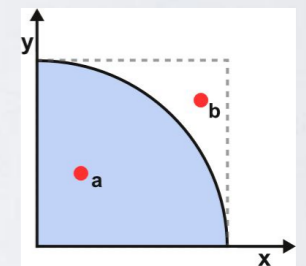
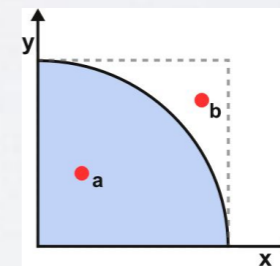
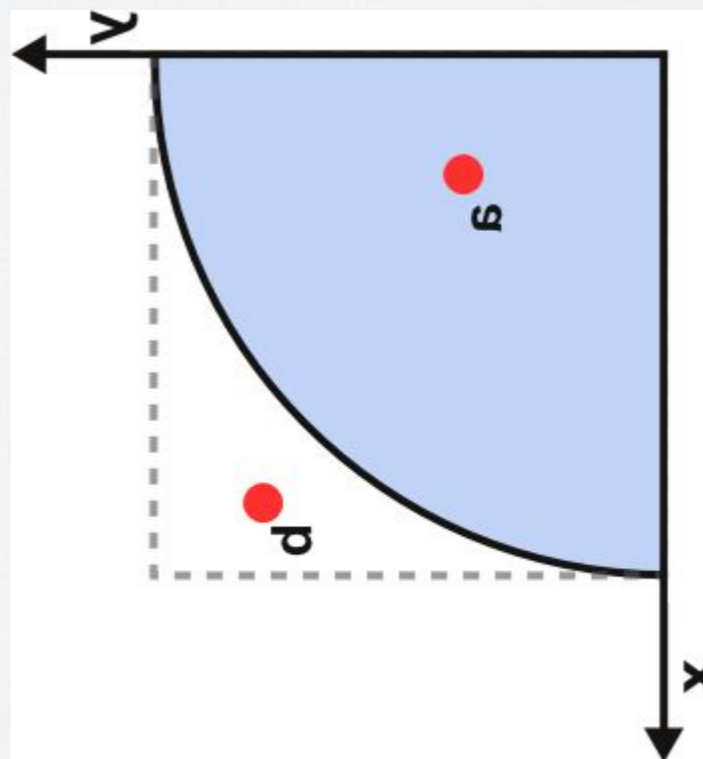
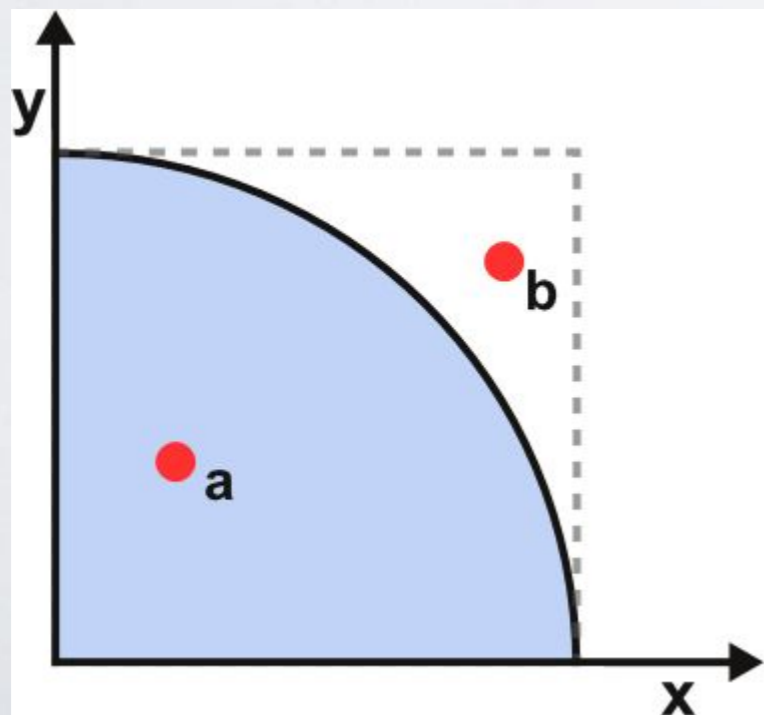
- Better estimator
 - It uses more samples

HW I

- Can we do even better?
 - Squeeze more information from data

- $x = x^*2 - \text{torch.floor}(x^*2)$

- $x = | -x$



RECAP

- Loss function

Last Thu

- How well are fitting the data?

- Gradient computation

Today

- What **local** change to the weights decreases the loss?

- Optimization algorithm

Th

- How do we change the weights to **globally** minimize the loss?

fully connected

ReLU

fully connected

ReLU

fully connected

ReLU

fully connected

ReLU

fully connected

GRADIENT

- Forward pass
 - Compute all latent activations z
 - Compute the loss
- How do we compute the gradient of l w.r.t. W_1 , W_3 , W_5 ?
 - l and all f are differentiable

Loss

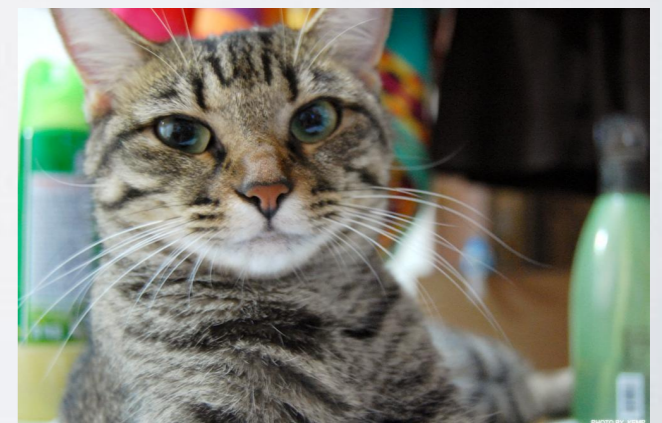
fully connected

ReLU

fully connected

ReLU

fully connected



GRADIENT

- Forward is a gigantic nested function?

$$l(f_5(f_4(f_3(f_2(f_1(x, W_1))), W_3)), W_5)$$

- Compute the gradient using chain-rule

$$l(z_5)$$

$$z_5 = f_5(z_4, W_5)$$

$$z_4 = f_4(z_3)$$

$$z_3 = f_3(z_2, W_3)$$

$$z_2 = f_2(z_1)$$

$$z_1 = f_1(x, W_1)$$

Loss

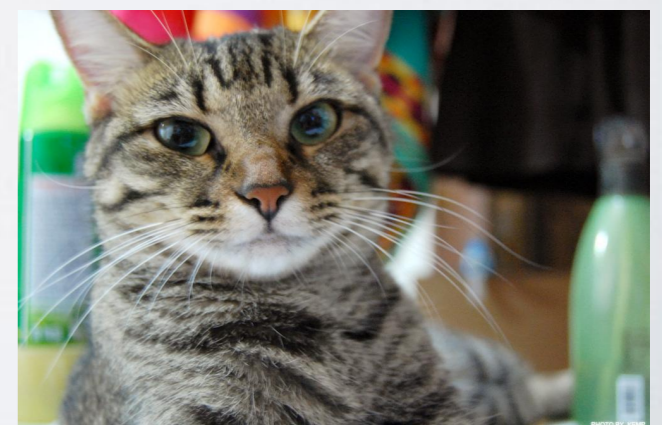
fully connected

ReLU

fully connected

ReLU

fully connected



GRADIENT

$$\frac{d}{dW_5} l(f_5(f_4(f_3(f_2(f_1(x, W_1))), W_3)), W_5) = \frac{d}{dz_5} l(z_5) \frac{d}{dW_5} f_5(z_4, W_5)$$

$$\frac{d}{dW_5} l(z_5) = \frac{d}{dz_5} l(z_5) \frac{d}{dz_4} f_5(z_4, W_5) \frac{d}{dz_3} f_4(z_3) \frac{d}{dW_3} f_3(z_2, W_3)$$

...

$l(z_5)$

$$z_5 = f_5(z_4, W_5)$$

$$z_4 = f_4(z_3)$$

$$z_3 = f_3(z_2, W_3)$$

$$z_2 = f_2(z_1)$$

$$z_1 = f_1(x, W_1)$$

GRADIENT + CHAIN RULE

$$l(z_4)$$

$$z_4 = f_4(z_3)$$

$$z_3 = f_3(z_2, W_3)$$

$$z_2 = f_2(z_1)$$

$$z_1 = f_1(x, W_1)$$

$$\frac{\partial}{\partial W_3} l(z_4)$$

$$= \frac{\partial}{\partial z_4} l(z_4) \frac{\partial}{\partial W_3} z_4$$

$$= \frac{\partial}{\partial z_4} l(z_4) \frac{\partial}{\partial z_3} z_4 \frac{\partial}{\partial W_3} z_3$$

EVALUATION

$$\frac{\partial}{\partial W_3} \ell(z_4) = \frac{\partial}{\partial z_4} \ell(z_4) \frac{\partial}{\partial z_3} z_4 \frac{\partial}{\partial W_3} z_3$$

scalar

Id

EVALUATION

$$\frac{\partial}{\partial W_3} \ell(z_4) = \frac{\partial}{\partial z_4} \ell(z_4) \frac{\partial}{\partial z_3} z_4 \frac{\partial}{\partial W_3} z_3$$

$$1 \times |W_3|$$

EVALUATION

$$\frac{\partial}{\partial W_3} \ell(z_4) = \frac{\partial}{\partial z_4} \ell(z_4) \frac{\partial}{\partial z_3} z_4 \frac{\partial}{\partial W_3} z_3$$

|z3|

EVALUATION

$$\frac{\partial}{\partial W_3} \ell(z_4) = \frac{\partial}{\partial z_4} \ell(z_4) \frac{\partial}{\partial z_3} z_4 \frac{\partial}{\partial W_3} z_3$$

$$|z_3| \times |W_3|$$

EVALUATION

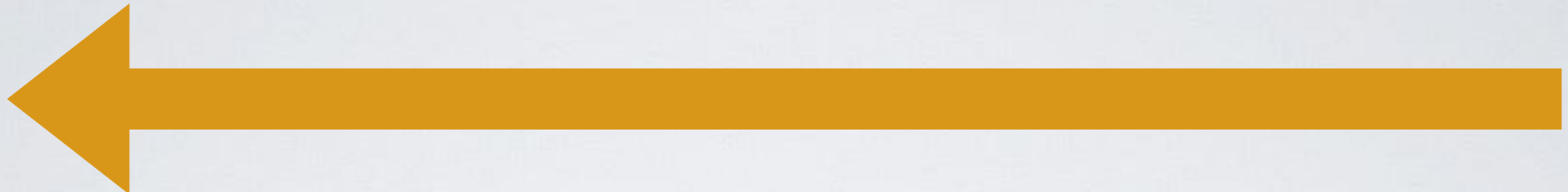
$$\frac{\partial}{\partial W_3} \ell(z_4) = \frac{\partial}{\partial z_4} \ell(z_4) \frac{\partial}{\partial z_3} z_4 \frac{\partial}{\partial W_3} z_3$$

$|z_4| \times |z_3|$
 $|z_3| \times |W_3|$

EVALUATION

$$\frac{\partial}{\partial W_3} \ell(z_4) = \frac{\partial}{\partial z_4} \ell(z_4) \frac{\partial}{\partial z_3} z_4 \frac{\partial}{\partial W_3} z_3$$
$$1 \times |z_4| \quad |z_4| \times |z_3| \quad |z_3| \times |W_3|$$

EVALUATION



$$\frac{\partial}{\partial z_4} \ell(z_4)$$

$$\frac{\partial}{\partial z_3} z_4$$

$$\frac{\partial}{\partial W_3} z_3$$

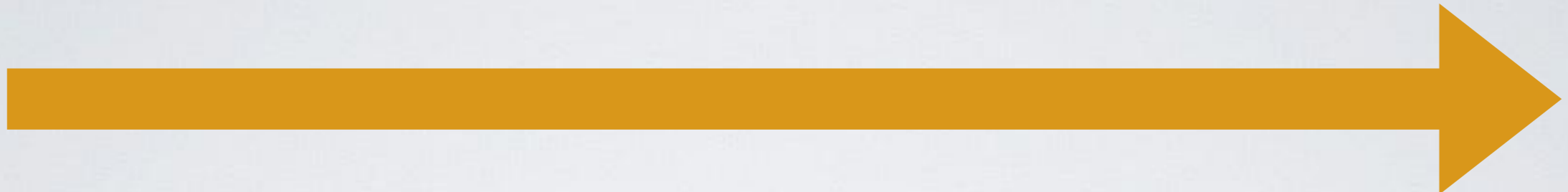
$$1 \times |z_4|$$

$$|z_4| \times |z_3|$$

$$|z_3| \times |W_3|$$

$$O(1|z_4||W_3|) \quad O(|z_4||z_3||W_3|)$$

EVALUATION



$$\frac{\partial}{\partial z_4} \ell(z_4)$$

$$\frac{\partial}{\partial z_3} z_4$$

$$\frac{\partial}{\partial W_3} z_3$$

$$1 \times |z_4|$$

$$|z_4| \times |z_3|$$

$$|z_3| \times |W_3|$$

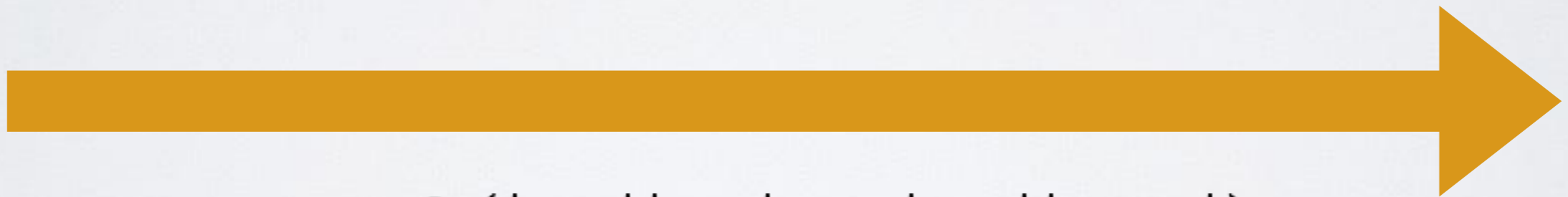
$$O(1|z_4||z_3|)$$

$$O(1|z_3||W_3|)$$

EVALUATION

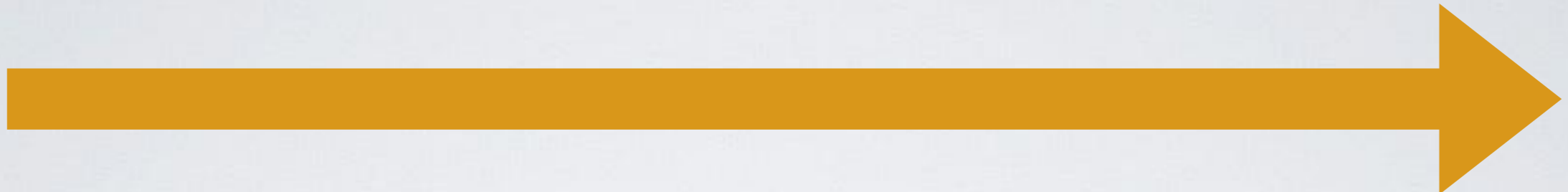


$$O(|z_4||W_3| + |z_4||z_3||W_3|)$$



$$O(|z_4||z_3| + |z_3||W_3|)$$

BACK PROPAGATION



$$\frac{\partial}{\partial z_4} \ell(z_4)$$

$$\frac{\partial}{\partial z_3} z_4$$

$$\frac{\partial}{\partial W_3} z_3$$

DEEP NETWORK

dog vs cat

$$l(z_4)$$

Non-linearity

$$z_4 = f_4(z_3)$$

Linear transformation W

$$z_3 = f_3(z_2, W_3)$$

Non-linearity

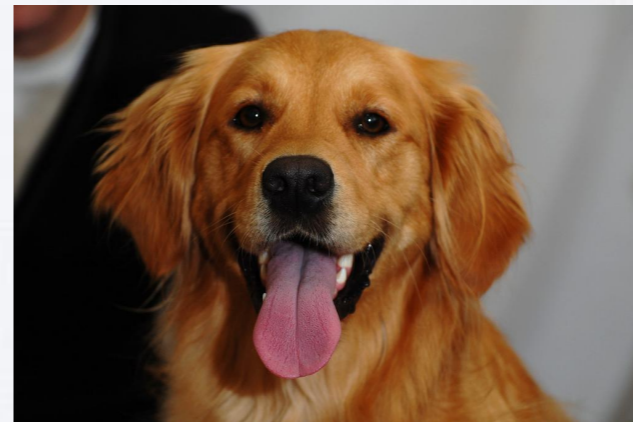
$$z_2 = f_2(z_1)$$

Linear transformation W

$$z_1 = f_1(x, W_1)$$

backward

forward



EFFICIENT COMPUTATION?

$$\frac{\partial}{\partial W_3} \ell(z_4) = \frac{\partial}{\partial z_4} \ell(z_4) \frac{\partial}{\partial z_3} z_4 \frac{\partial}{\partial W_3} z_3$$

$$\frac{\partial}{\partial W_1} \ell(z_4) = \frac{\partial}{\partial z_4} \ell(z_4) \frac{\partial}{\partial z_3} z_4 \frac{\partial}{\partial z_2} z_3 \frac{\partial}{\partial z_1} z_2 \frac{\partial}{\partial W_1} z_1$$

DEEP NETWORK

dog vs cat

$$\ell(z_4)$$

Non-linearity

$$z_4 = f_4(z_3)$$

Linear transformation W

$$z_3 = f_3(z_2, W_3)$$

Non-linearity

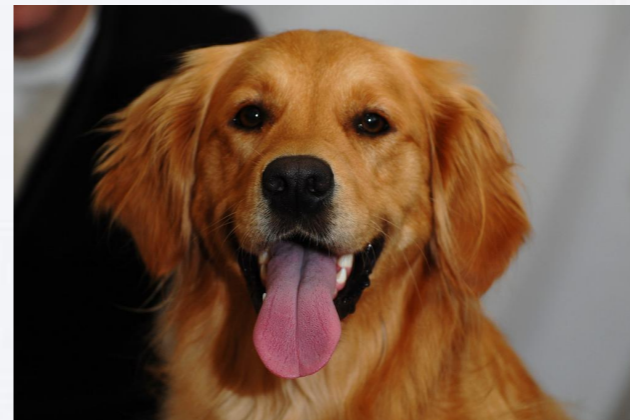
$$z_2 = f_2(z_1)$$

Linear transformation W

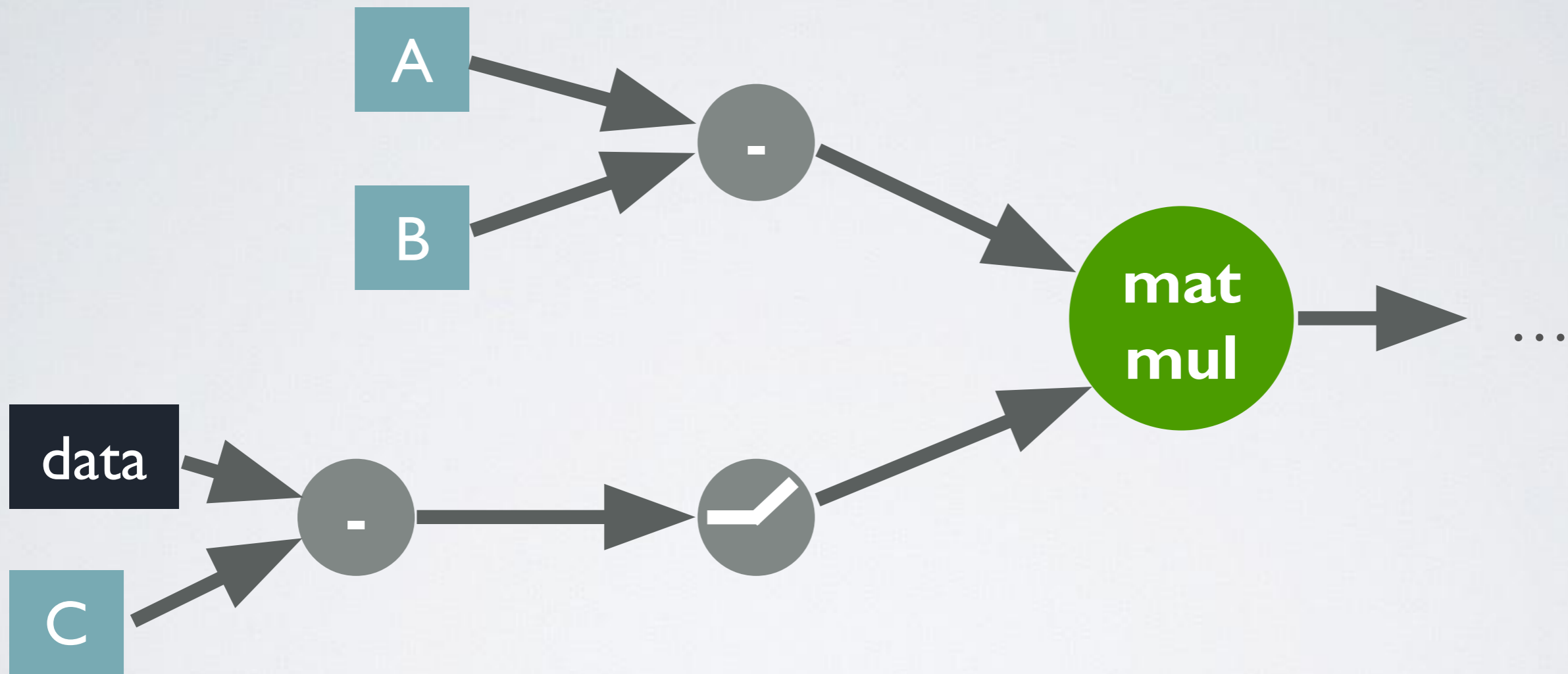
$$z_1 = f_1(x, W_1)$$

backward

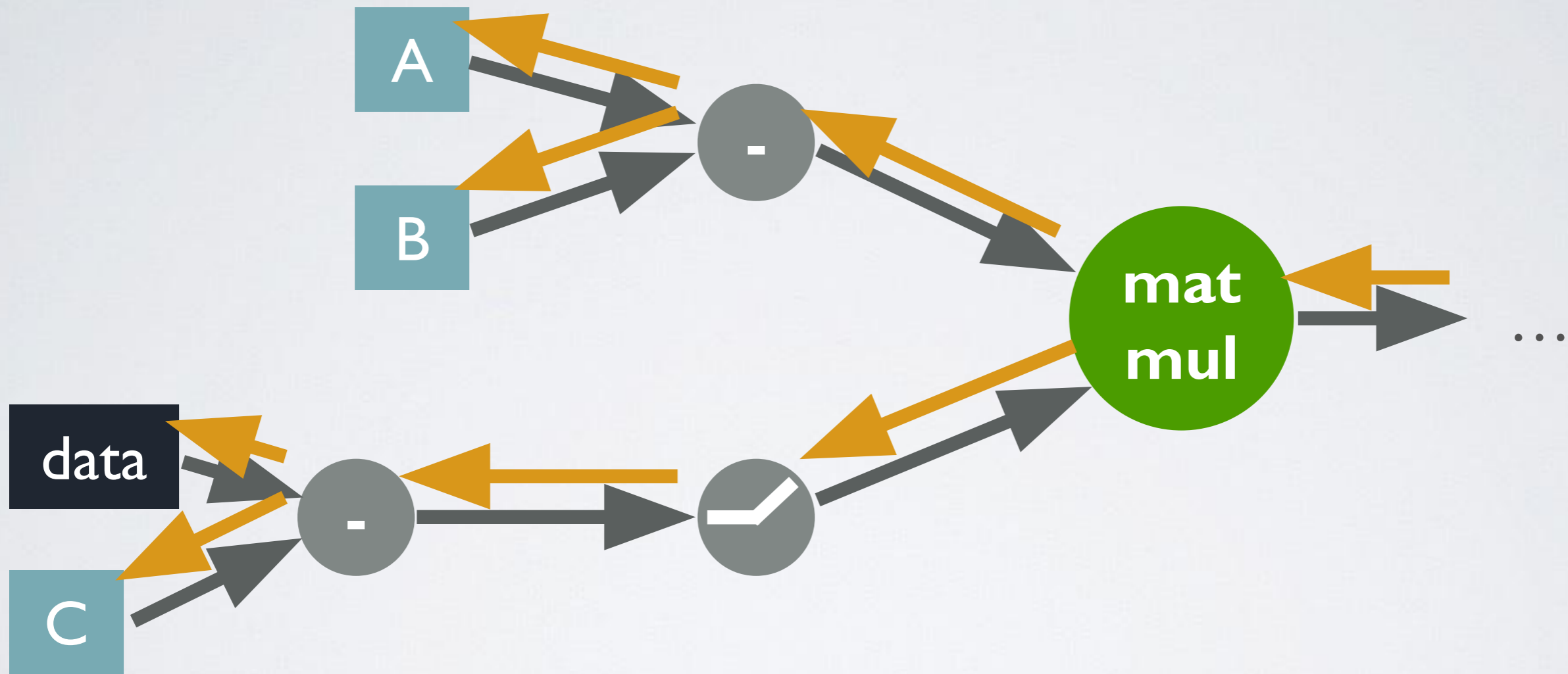
forward

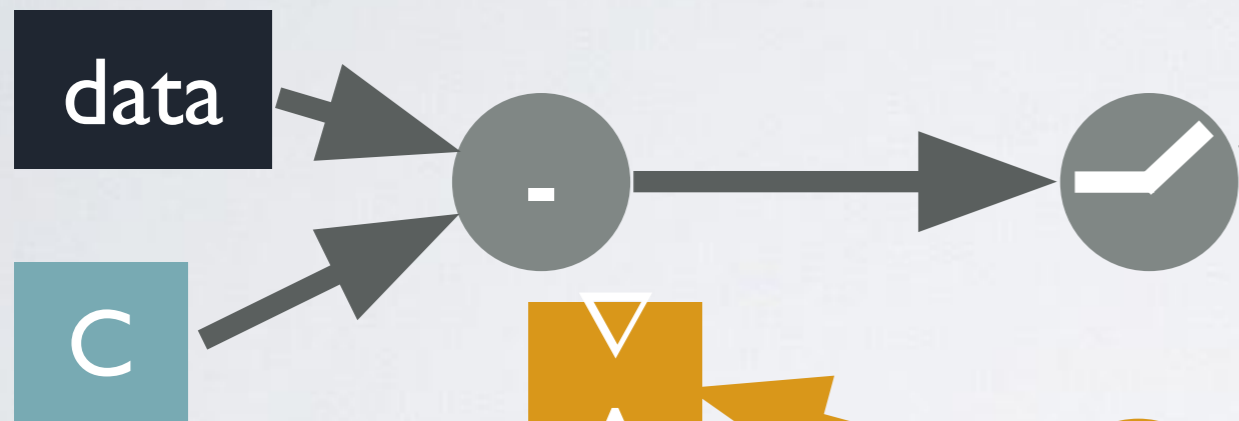
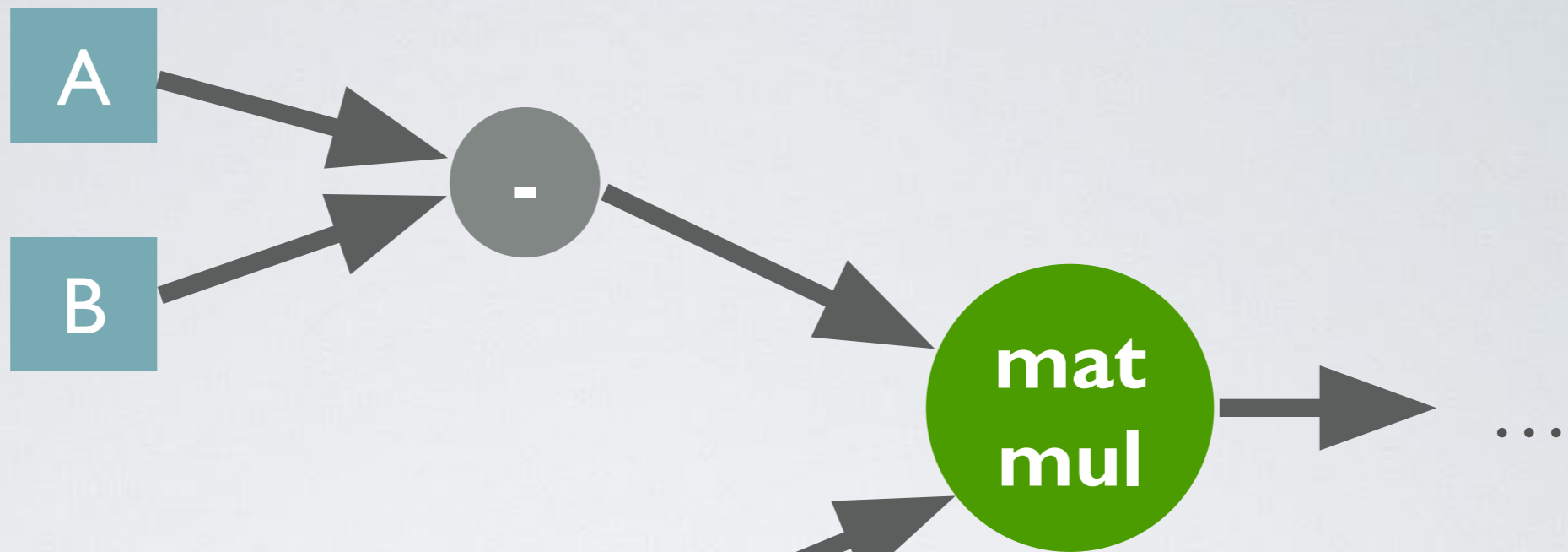


DEEP NETWORK IN GENERAL

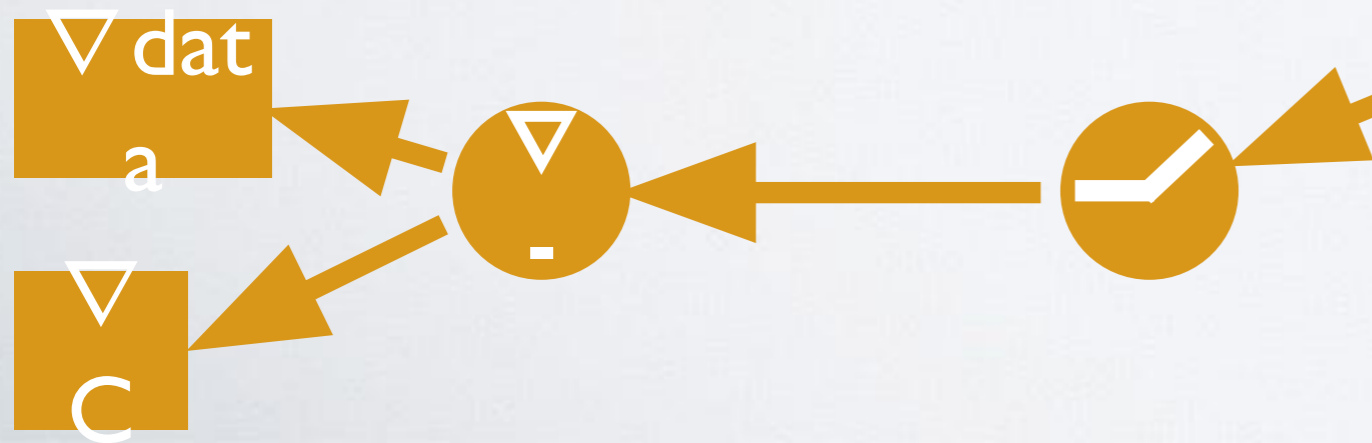
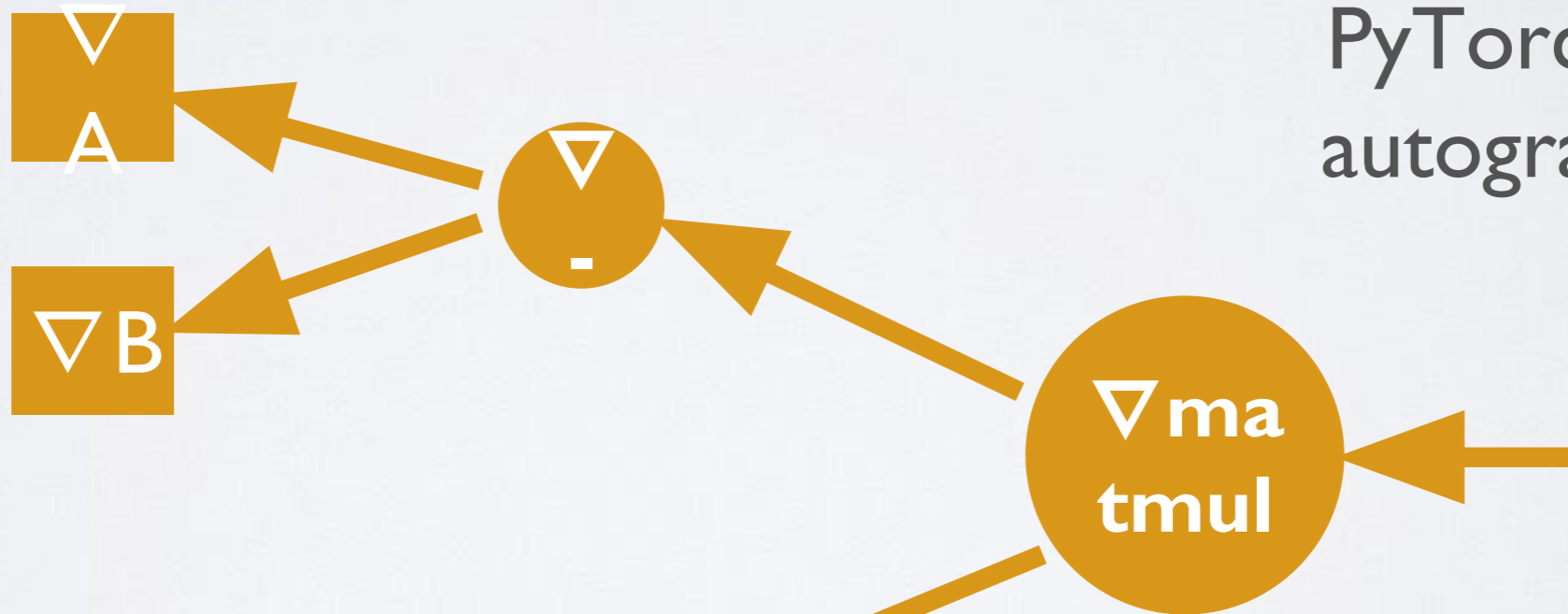


DEEP NETWORK IN GENERAL





in
PyTorch
autograd



Fit to screen

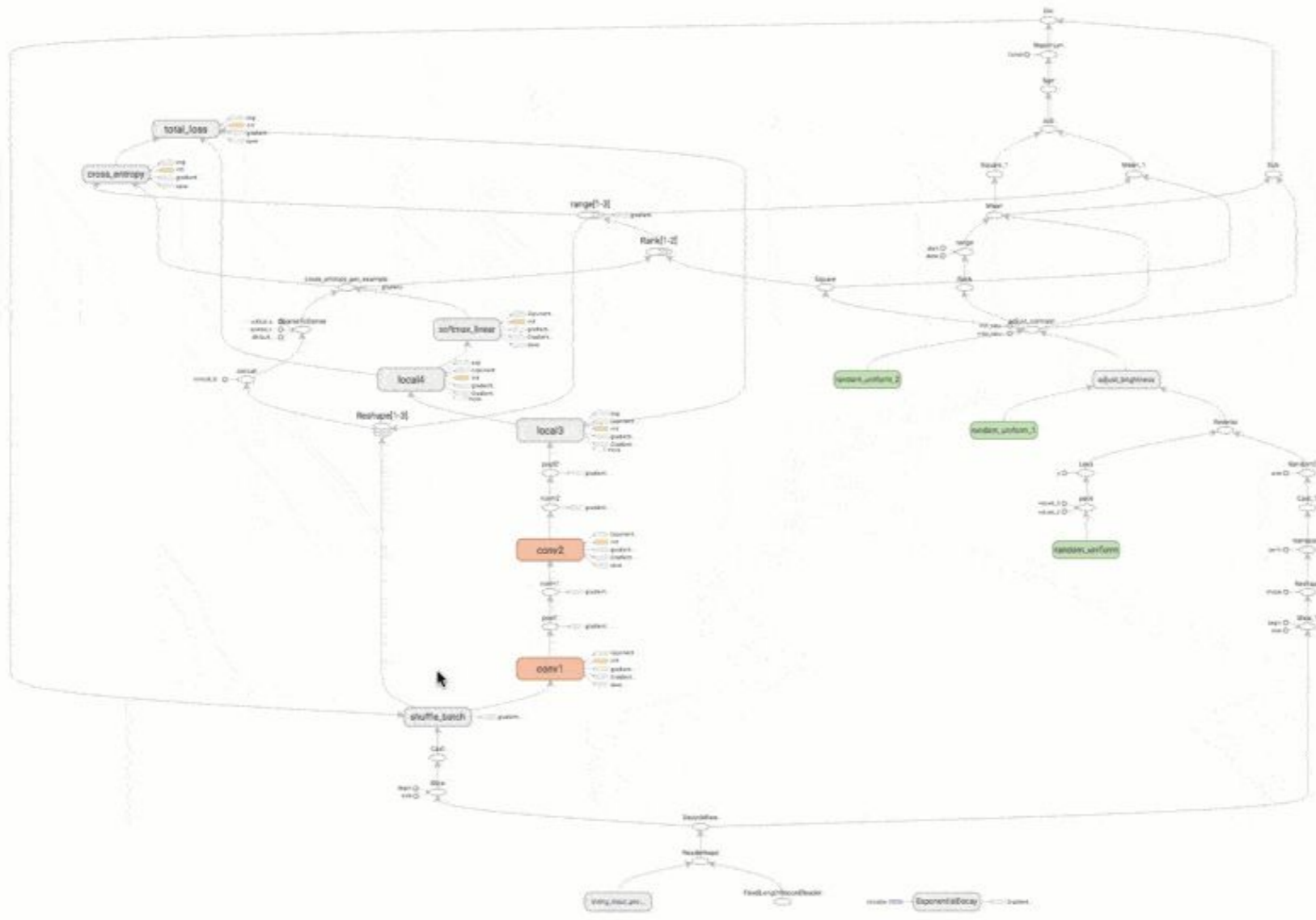
Run `cifar-train`

Upload Choose File

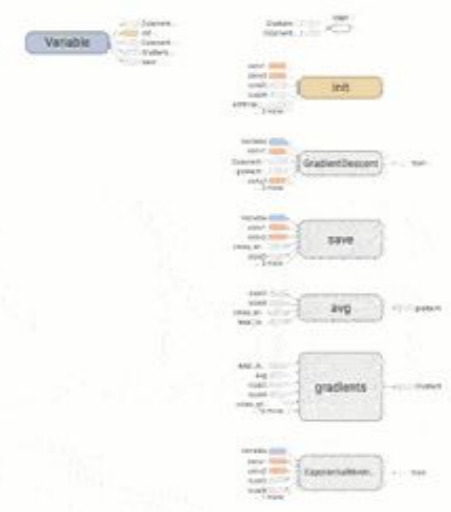
Color Structure
color: same substructure
gray: unique substructure

Graph (* = expandable)
Namespace*
OpNode
Unconnected series*
Connected series*
Constant
Summary
Dataflow edge
Control dependency edge
Reference edge

Main Graph



Auxiliary nodes



In PyTorch, need additional plugins to visualize compute graph

DEEP NETWORK

dog vs cat

$$l(z_4)$$

Non-linearity

$$z_4 = f_4(z_3)$$

Linear transformation W

$$z_3 = f_3(z_2, W_3)$$

Non-linearity

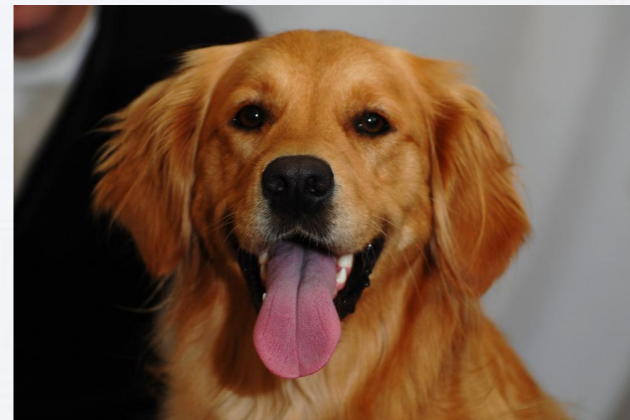
$$z_2 = f_2(z_1)$$

Linear transformation W

$$z_1 = f_1(x, W_1)$$

backward

forward



RECAP

- Loss function

Last Thu

- How well are fitting the data?

- Gradient computation

Today

- What **local** change to the weights decreases the loss?

- Optimization algorithm

Th

- How do we change the weights to **globally** minimize the loss?

fully connected

ReLU

fully connected

ReLU

fully connected

ReLU

fully connected

ReLU

fully connected

NEXT CLASS

dog vs cat

$$\ell(z_4)$$

Non-linearity

$$z_4 = f_4(z_3)$$

Linear transformation W

$$z_3 = f_3(z_2, W_3)$$

Non-linearity

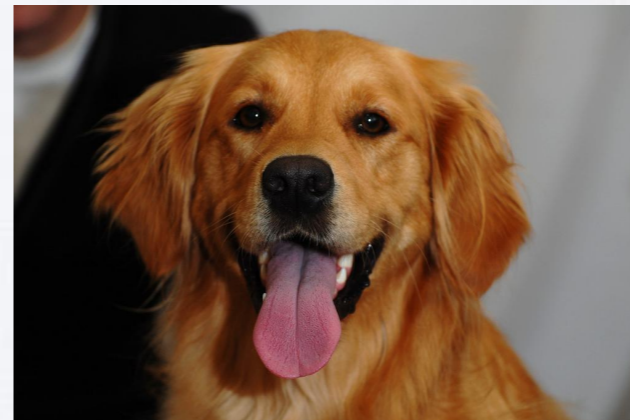
$$z_2 = f_2(z_1)$$

Linear transformation W

$$z_1 = f_1(x, W_1)$$

backward

forward



NEXT WEEK

dog vs cat

$$l(z_4)$$

Non-linearity

$$z_4 = f_4(z_3)$$

Linear transformation W

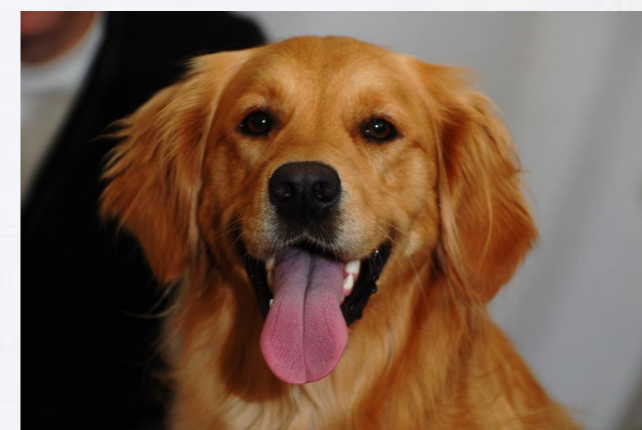
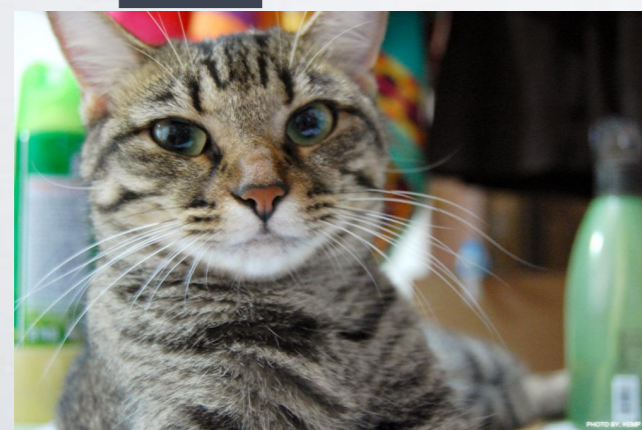
$$z_3 = f_3(z_2, W_3)$$

Non-linearity

$$z_2 = f_2(z_1)$$

Linear transformation W

$$z_1 = f_1(x, W_1)$$



NEXT WEEK

dog vs cat

$$l(z_4)$$

Non-linearity

$$z_4 = f_4(z_3)$$

Linear transformation W

$$z_3 = f_3(z_2, W_3)$$

Non-linearity

$$z_2 = f_2(z_1)$$

Linear transformation W

$$z_1 = f_1(x, W_1)$$

