

# CSC 561: Neural Networks and Deep Learning

## Learning (part 3)

Marco Alvarez

Department of Computer Science and Statistics  
University of Rhode Island

Spring 2025



# Neural Networks

## Learning the network: Part 3

11-785, Spring 2025

Lecture 5

1  
2

## Training neural nets through Empirical Risk Minimization: Problem Setup

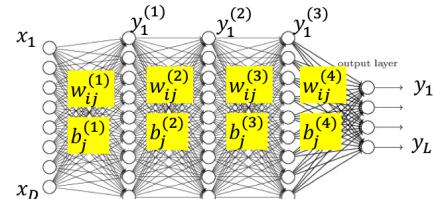
- Given a training set of input-output pairs  $(X_1, d_1), (X_2, d_2), \dots, (X_T, d_T)$
- The divergence on the  $i^{\text{th}}$  instance is  $\text{div}(Y_i, d_i)$ 
  - $- Y_i = f(X_i; W)$
- The loss (empirical risk)

$$\text{Loss}(W) = \frac{1}{T} \sum_i \text{div}(Y_i, d_i)$$

- Minimize  $\text{Loss}$  w.r.t  $\{w_{ij}^{(k)}, b_j^{(k)}\}$  using gradient descent

2

## Notation



- The input layer is the  $0^{\text{th}}$  layer
- We will represent the output of the  $i$ -th perceptron of the  $k^{\text{th}}$  layer as  $y_i^{(k)}$ 
  - Input to network:**  $y_i^{(0)} = x_i$
  - Output of network:**  $y_i = y_i^{(N)}$
- We will represent the weight of the connection between the  $i$ -th unit of the  $k-1^{\text{th}}$  layer and the  $j$ -th unit of the  $k^{\text{th}}$  layer as  $w_{ij}^{(k)}$ 
  - The bias to the  $j$ -th unit of the  $k^{\text{th}}$  layer is  $b_j^{(k)}$

3

4

## Training Neural Nets through Gradient Descent

Total training Loss:

$$\text{Loss} = \frac{1}{T} \sum_t \text{Div}(\mathbf{Y}_t, \mathbf{d}_t)$$

- Gradient descent algorithm:
- Initialize all weights and biases  $\{w_{ij}^{(k)}\}$  Assuming the bias is also represented as a weight
- Using the extended notation: the bias is also a weight
- Do:
  - For every layer  $k$  for all  $i, j$ , update:
    - $w_{i,j}^{(k)} = w_{i,j}^{(k)} - \eta \frac{d\text{Loss}}{dw_{i,j}^{(k)}}$
- Until  $\text{Loss}$  has converged

7

## The derivative

Total training Loss:

$$\text{Loss} = \frac{1}{T} \sum_t \text{Div}(\mathbf{Y}_t, \mathbf{d}_t)$$

- Computing the derivative

Total derivative:

$$\frac{d\text{Loss}}{dw_{i,j}^{(k)}} = \frac{1}{T} \sum_t \frac{d\text{Div}(\mathbf{Y}_t, \mathbf{d}_t)}{dw_{i,j}^{(k)}}$$

- So we must first figure out how to compute the derivative of divergences of individual training inputs

9

## Calculus Refresher: Basic rules of calculus

For any differentiable function

$$y = f(x)$$

with derivative

$$\frac{dy}{dx}$$

the following must hold for sufficiently small  $\Delta x$  →  $\Delta y \approx \frac{dy}{dx} \Delta x$

10

## Calculus Refresher: Basic rules of calculus

For any differentiable function

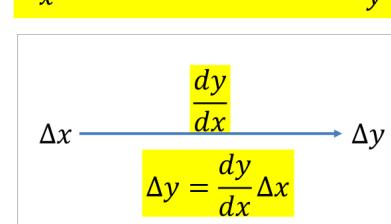
$$y = f(x)$$

with derivative

$$\frac{dy}{dx}$$

the following must hold for sufficiently small  $\Delta x$  →  $\Delta y \approx \frac{dy}{dx} \Delta x$

Introducing the "influence" diagram:  
x influences y



The derivative graph:  
The edge carries the derivative.

Node and edge weights multiply

12

## Calculus Refresher: Basic rules of calculus

For any differentiable function

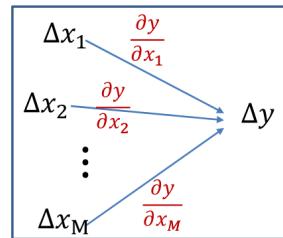
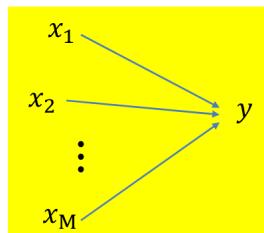
$$y = f(x_1, x_2, \dots, x_M)$$

with partial derivatives

$$\frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \dots, \frac{\partial y}{\partial x_M}$$

the following must hold for sufficiently small  $\Delta x_1, \Delta x_2, \dots, \Delta x_M$

$$\Delta y \approx \frac{\partial y}{\partial x_1} \Delta x_1 + \frac{\partial y}{\partial x_2} \Delta x_2 + \dots + \frac{\partial y}{\partial x_M} \Delta x_M$$

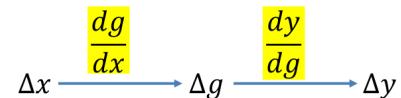


16

## Calculus Refresher: Chain rule

For any nested function  $y = f(g(x))$

$$\frac{dy}{dx} = \frac{dy}{dg(x)} \frac{dg(x)}{dx}$$

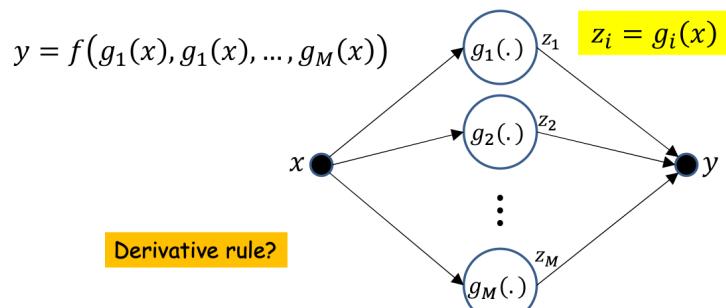


$$\Delta y = \frac{dy}{dg(x)} \frac{dg(x)}{dx} \Delta x$$

18

0

## Distributed Chain Rule: Influence Diagram

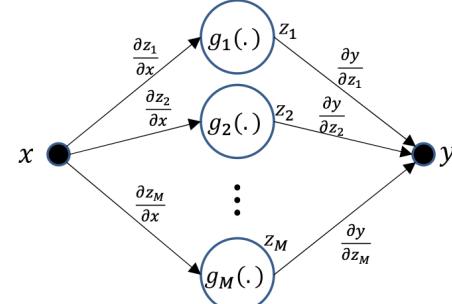


- $x$  affects  $y$  through each of  $g_1 \dots g_M$

20

## Distributed Chain Rule: Influence Diagram

$$y = f(g_1(x), g_2(x), \dots, g_M(x))$$



21

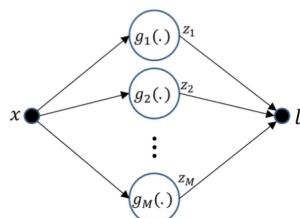
2

## Calculus Refresher: Chain rule summary

For any nested function  $l = f(y)$  where  $y = g(z)$

$$\frac{dl}{dz} = \frac{dl}{dy} \frac{dy}{dz}$$

For  $l = f(z_1, z_2, \dots, z_M)$   
where  $z_i = g_i(x)$



$$\frac{dl}{dx} = \frac{\partial l}{\partial z_1} \frac{dz_1}{dx} + \frac{\partial l}{\partial z_2} \frac{dz_2}{dx} + \dots + \frac{\partial l}{\partial z_M} \frac{dz_M}{dx}$$

23

3

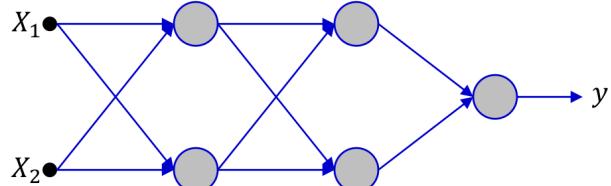
## Our problem for today

- How to compute  $\frac{d\text{Div}(Y,d)}{dw_{i,j}^{(k)}}$  for a single data instance

24

4

## A first closer look at the network

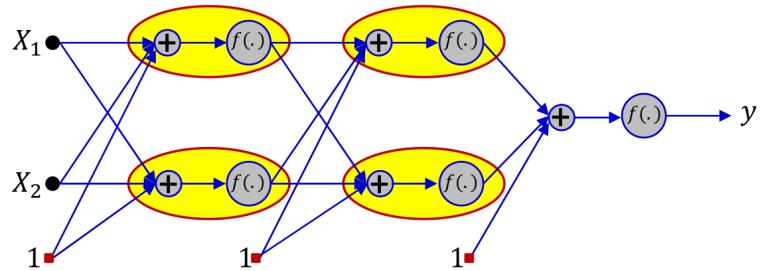


- Showing a tiny 2-input network for illustration
  - Actual network would have many more neurons and inputs

27

5

## A first closer look at the network

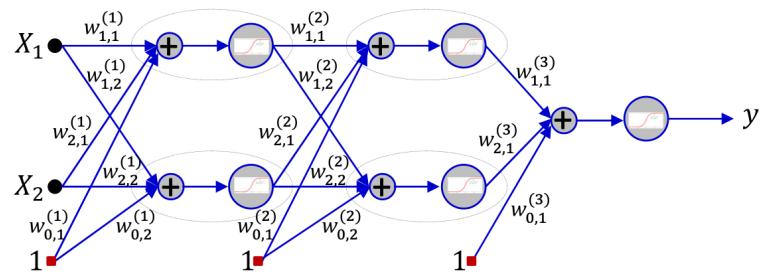


- Showing a tiny 2-input network for illustration
  - Actual network would have many more neurons and inputs
- Explicitly separating the affine function of inputs from the activation

28

6

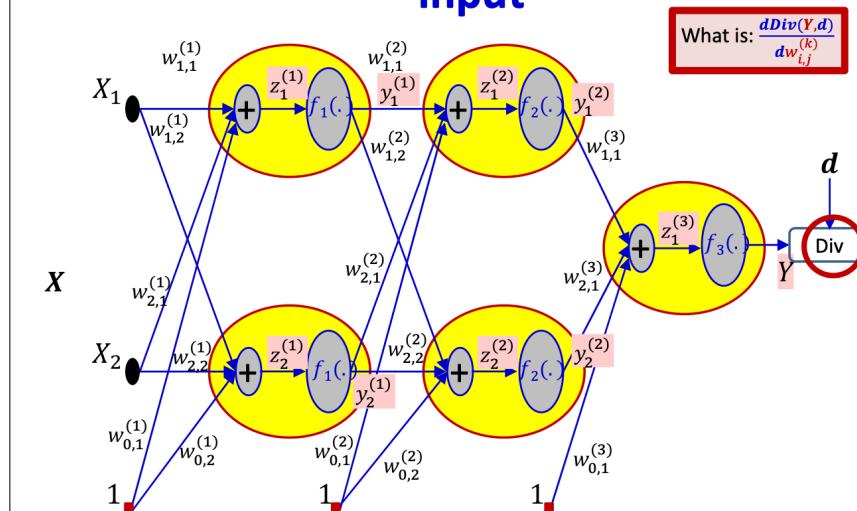
## A first closer look at the network



- Showing a tiny 2-input network for illustration
  - Actual network would have many more neurons and inputs
- Expanded **with all weights shown**
- Let's label the other variables too...

29

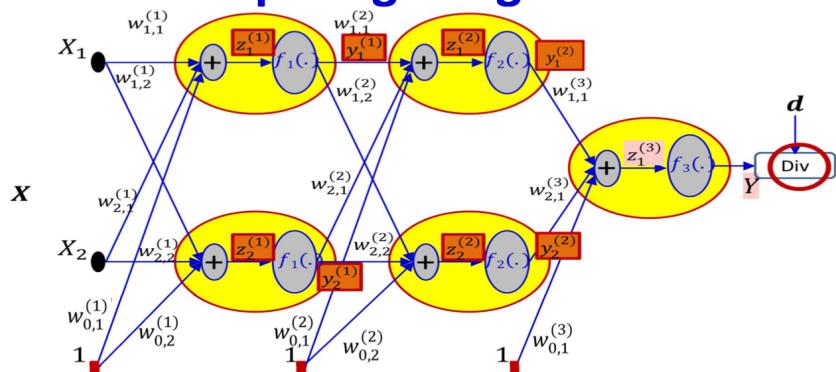
## Computing the derivative for a single input



31

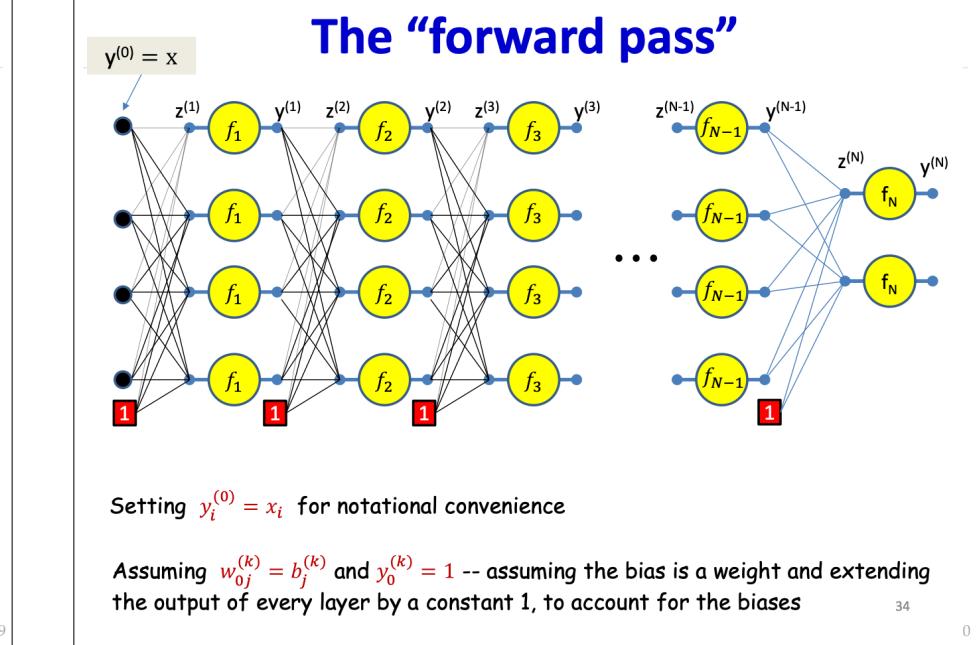
8

## Computing the gradient



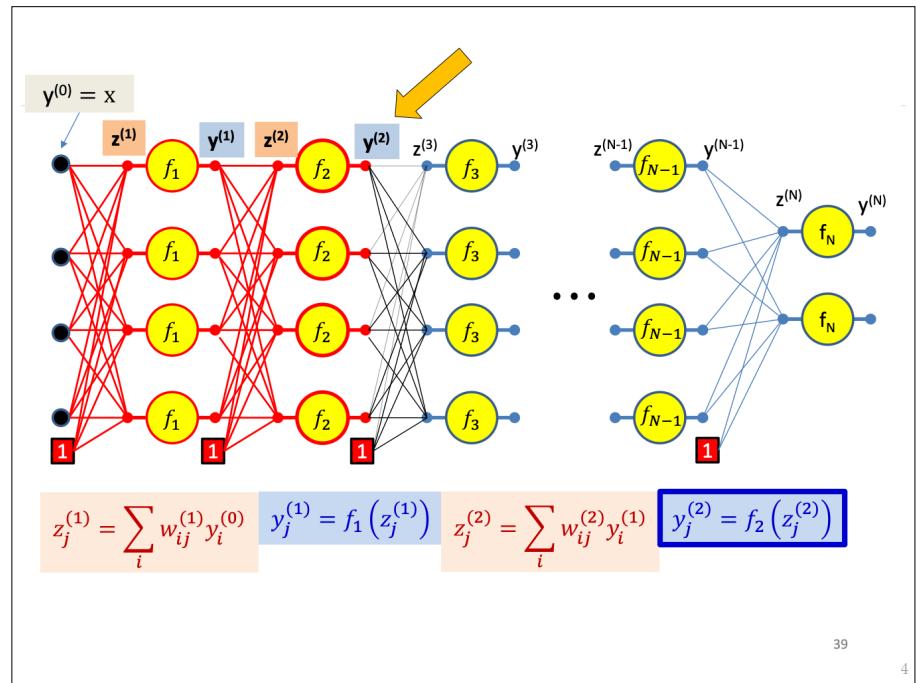
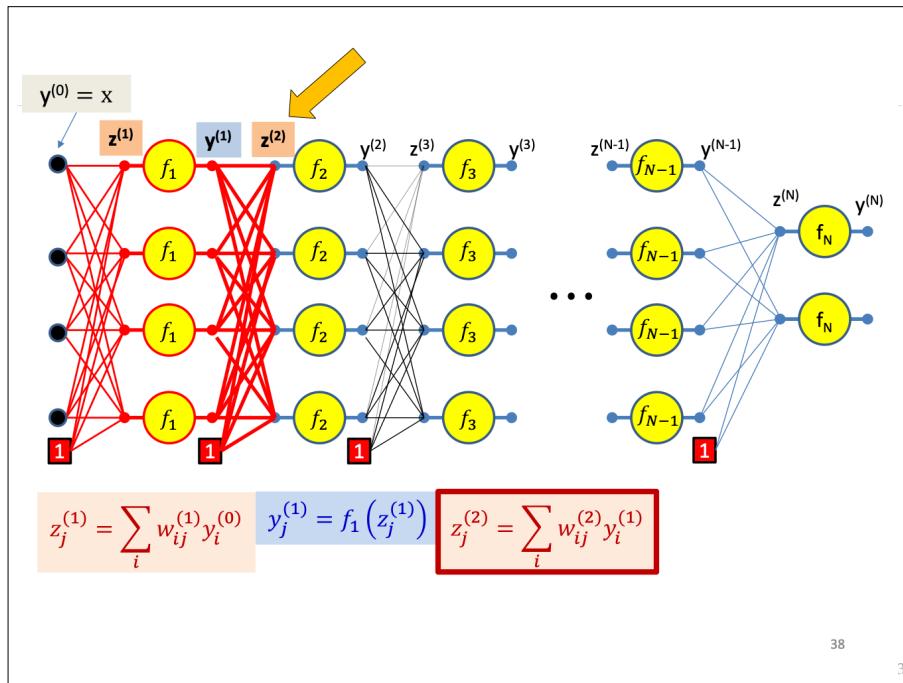
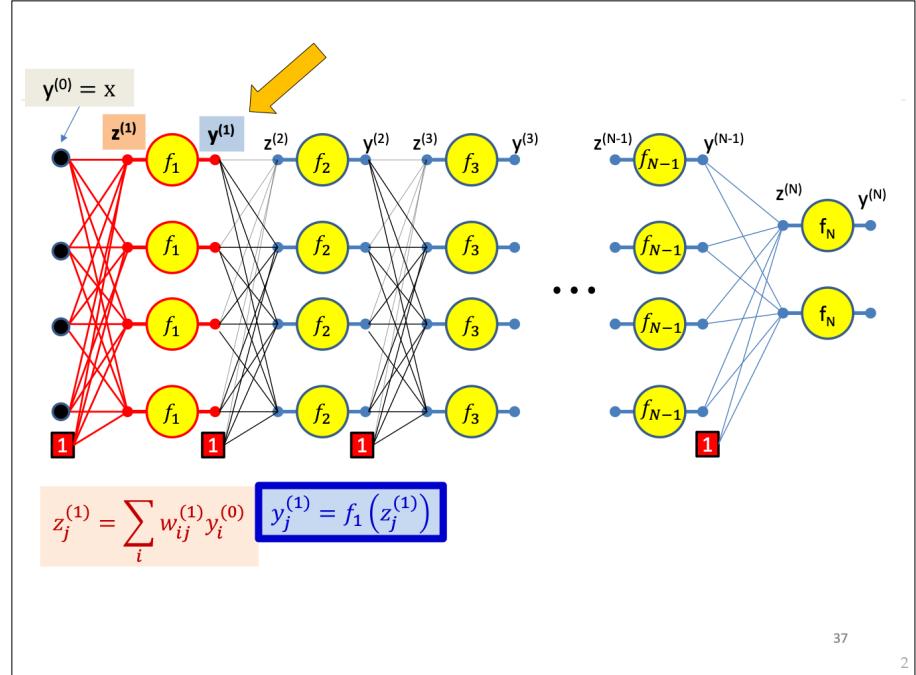
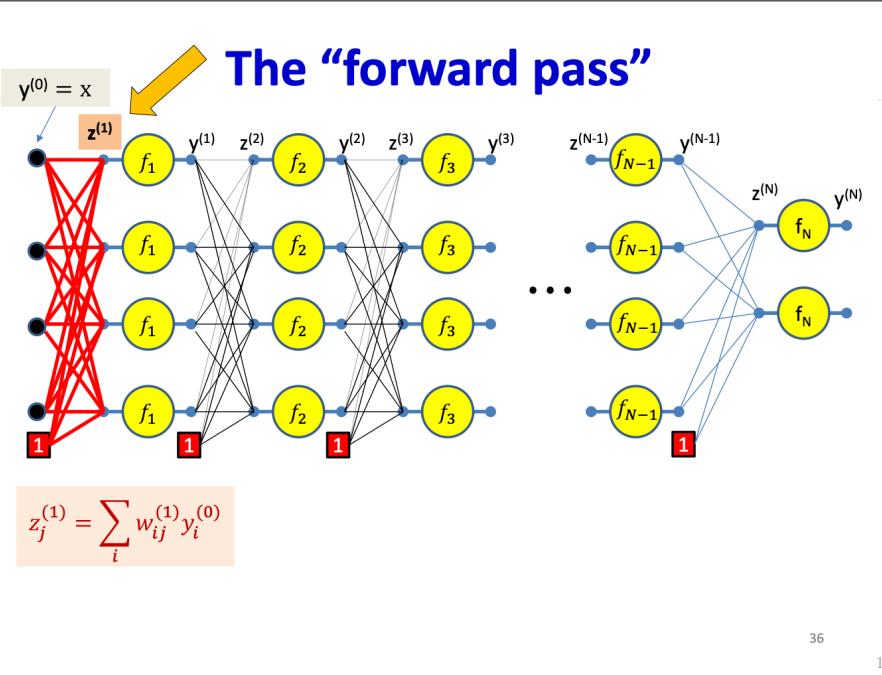
- Note: computation of the derivative  $\frac{d\text{Div}(Y,d)}{dw_{i,j}^{(k)}}$  requires intermediate and final output values of the network in response to the input

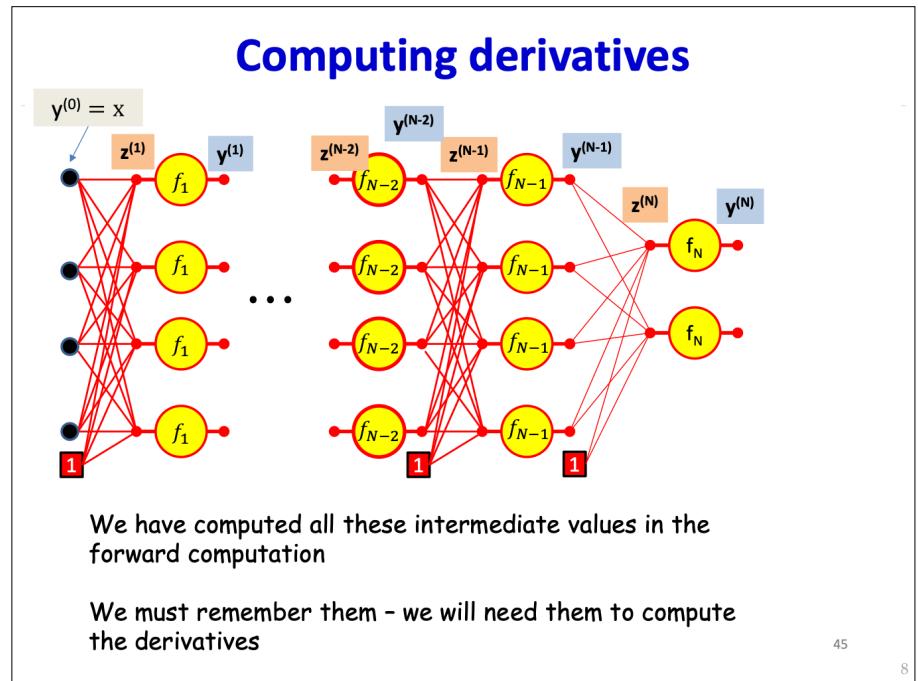
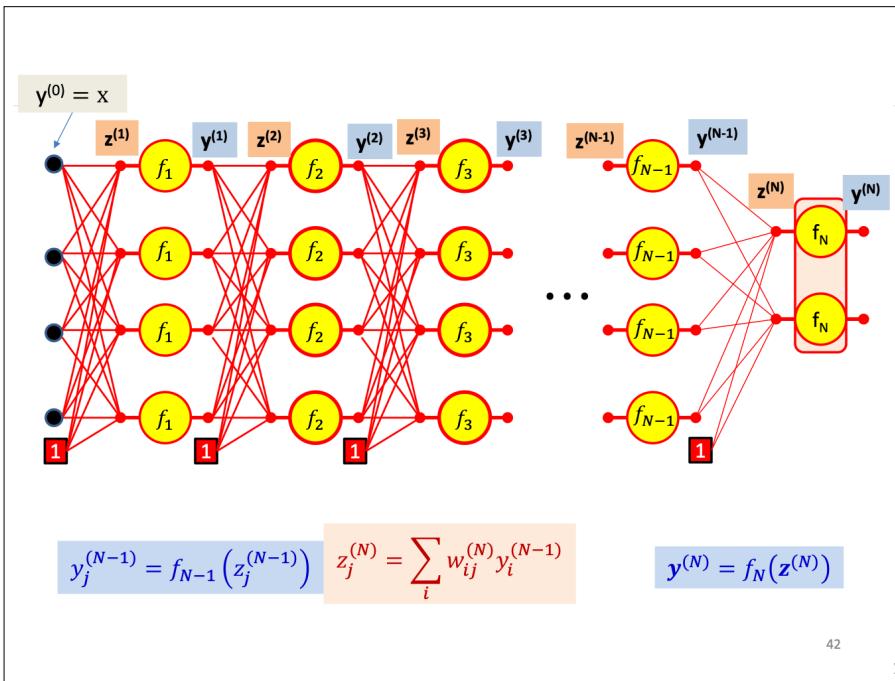
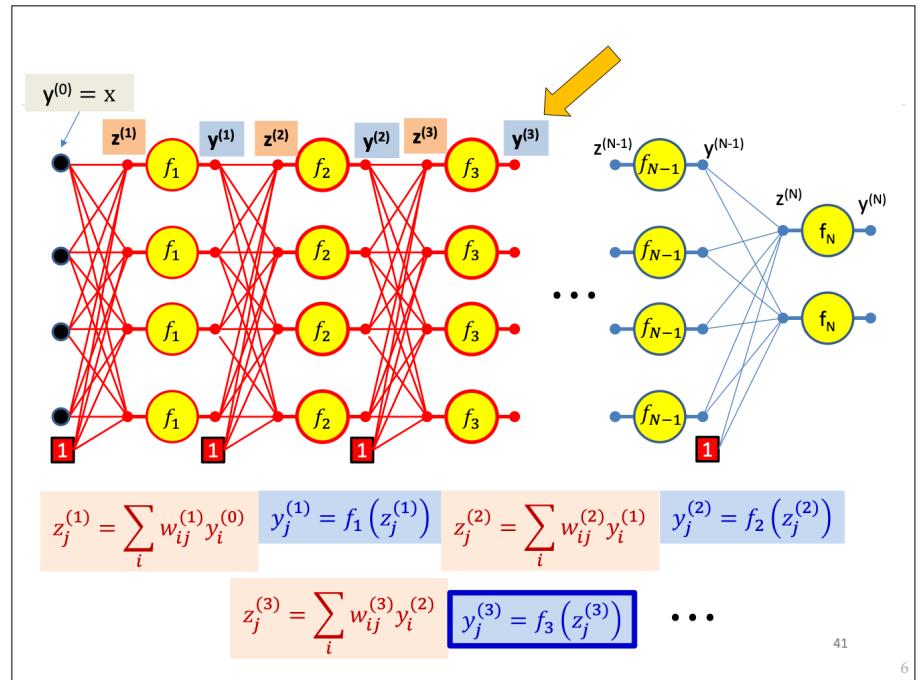
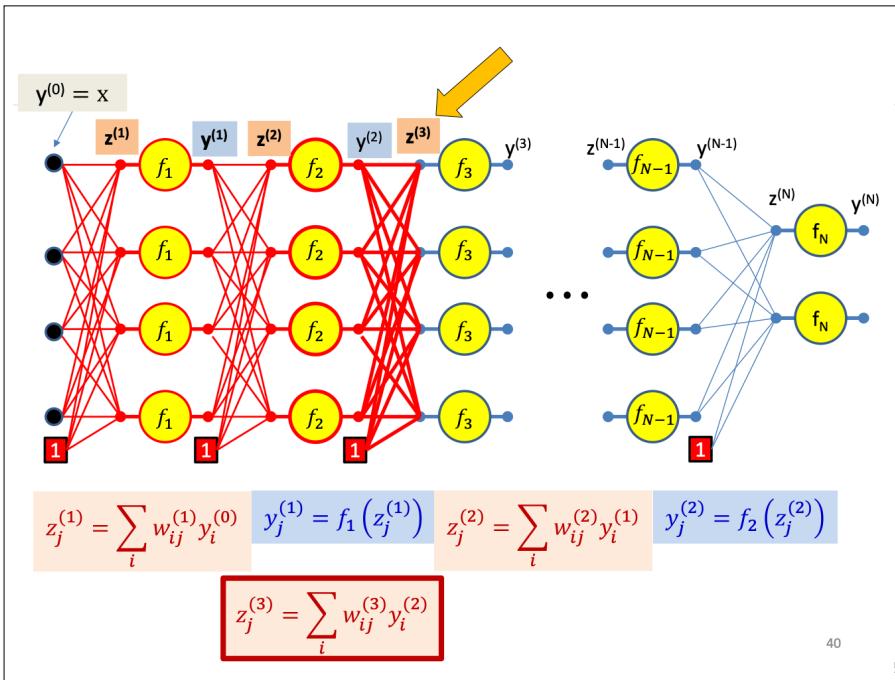
32



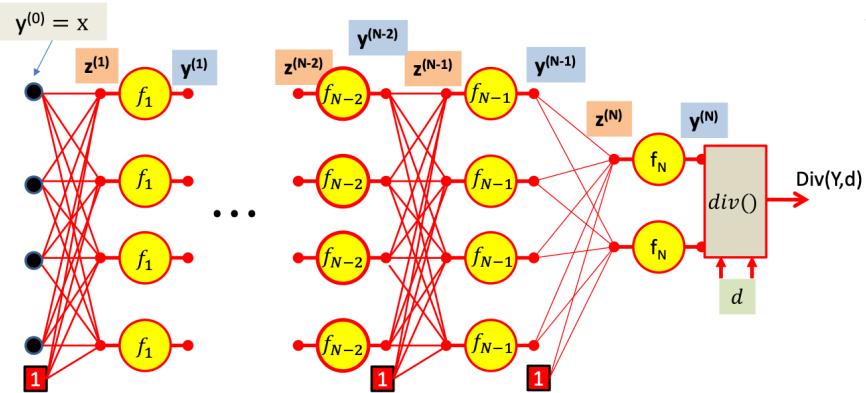
34

0





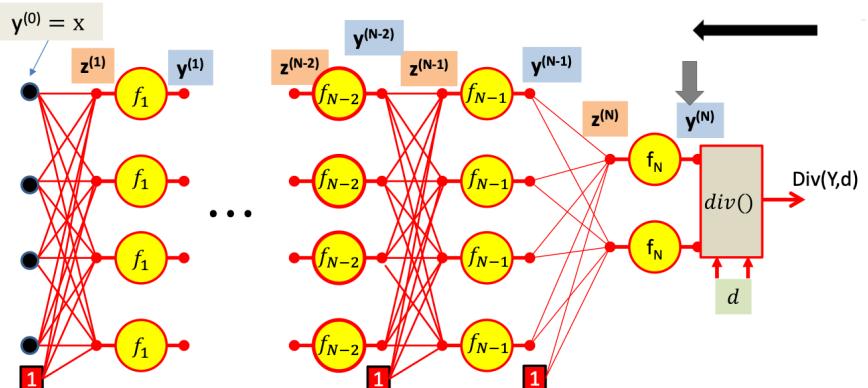
## Computing derivatives



First, we compute the divergence between the output of the net  $y = y^{(N)}$  and the desired output  $d$

46

## Computing derivatives

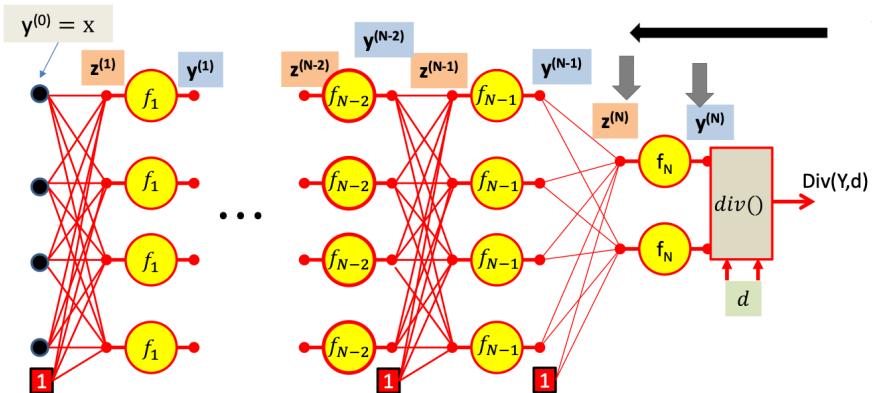


We then compute  $\nabla_{Y^{(N)}} \text{div}(\cdot)$  the derivative of the divergence w.r.t. the final output of the network  $y^{(N)}$

47

0

## Computing derivatives

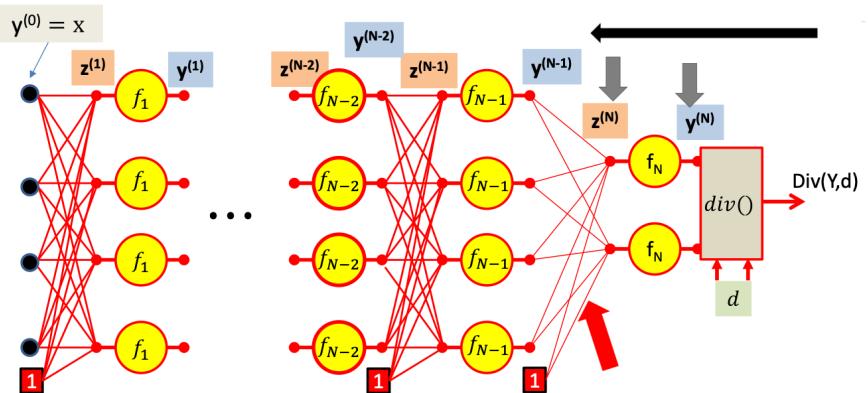


We then compute  $\nabla_{Y^{(N)}} \text{div}(\cdot)$  the derivative of the divergence w.r.t. the final output of the network  $y^{(N)}$

We then compute  $\nabla_{z^{(N)}} \text{div}(\cdot)$  the derivative of the divergence w.r.t. the pre-activation affine combination  $z^{(N)}$  using the chain rule

48

## Computing derivatives

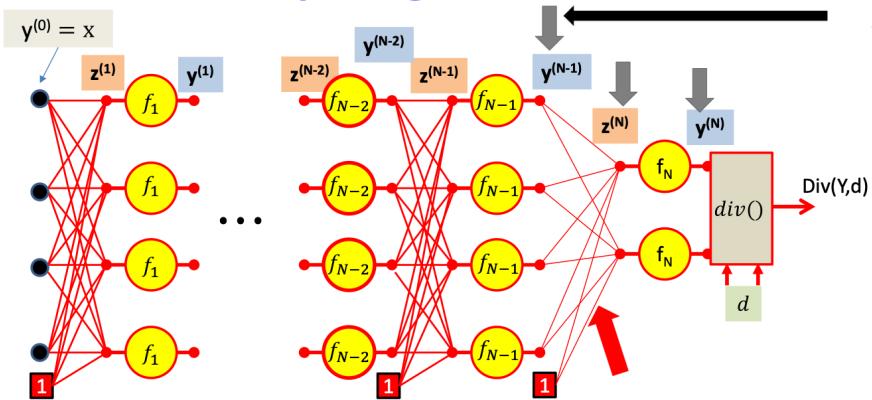


Continuing on, we will compute  $\nabla_{W^{(N)}} \text{div}(\cdot)$  the derivative of the divergence with respect to the weights of the connections to the output layer

49

2

## Computing derivatives

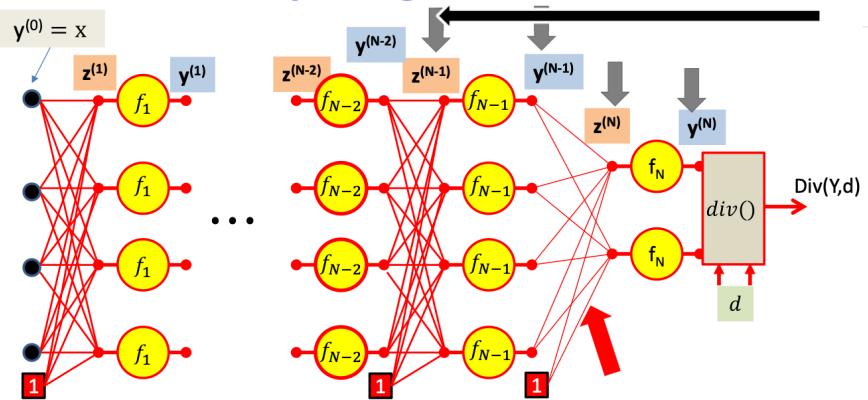


Continuing on, we will compute  $\nabla_{W^{(N)}} \text{div}(\cdot)$  the derivative of the divergence with respect to the weights of the connections to the output layer

Then continue with the chain rule to compute  $\nabla_{Y^{(N-1)}} \text{div}(\cdot)$  the derivative of the divergence w.r.t. the output of the N-th layer

50

## Computing derivatives

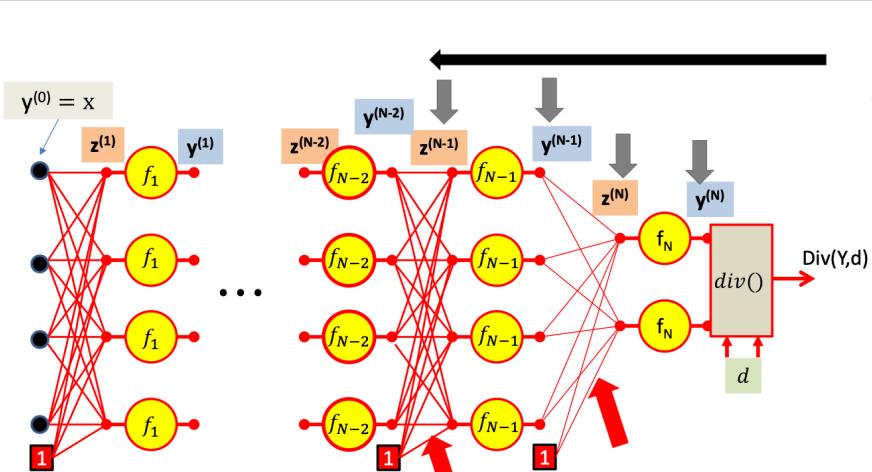


We continue our way backwards in the order shown

$$\nabla_{Z^{(N-1)}} \text{div}(\cdot)$$

51

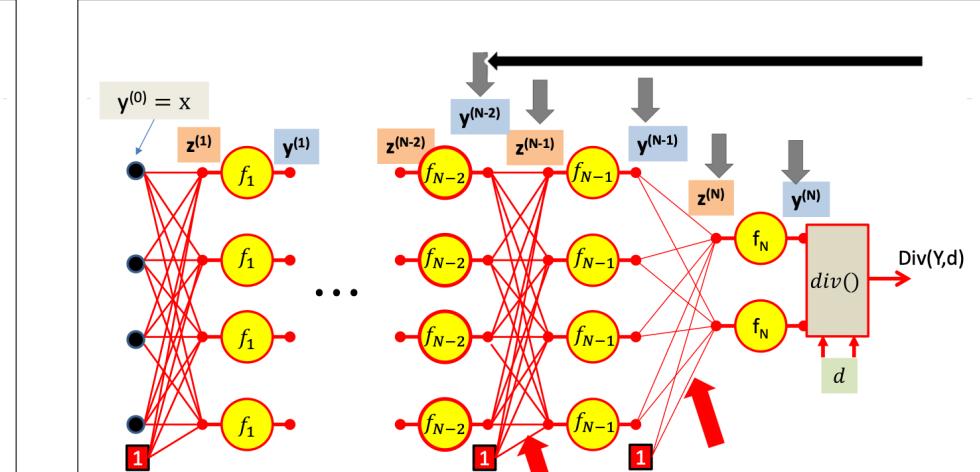
4



We continue our way backwards in the order shown

$$\nabla_{W^{(N-1)}} \text{div}(\cdot)$$

52

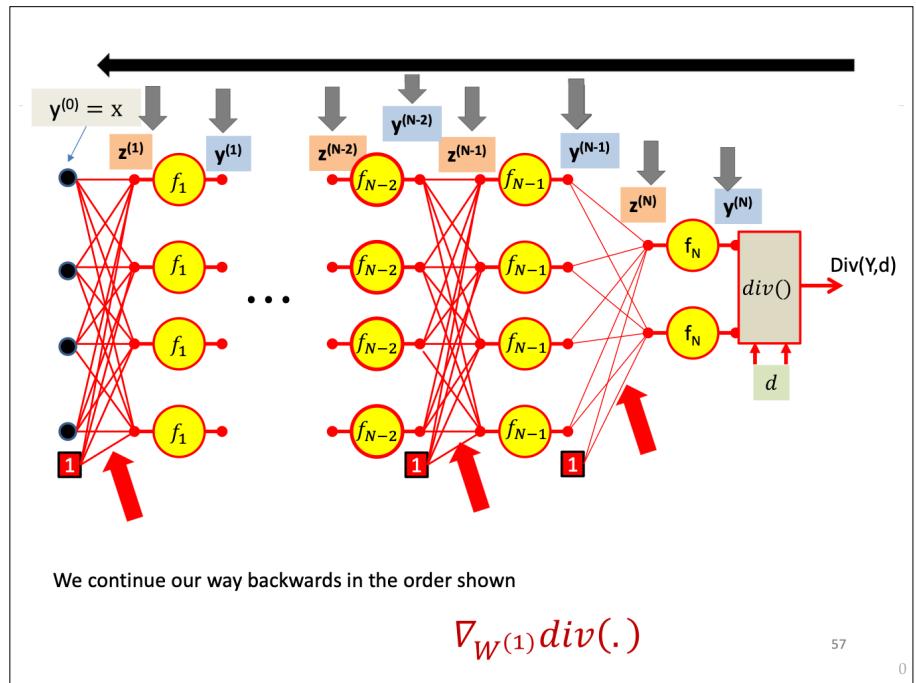
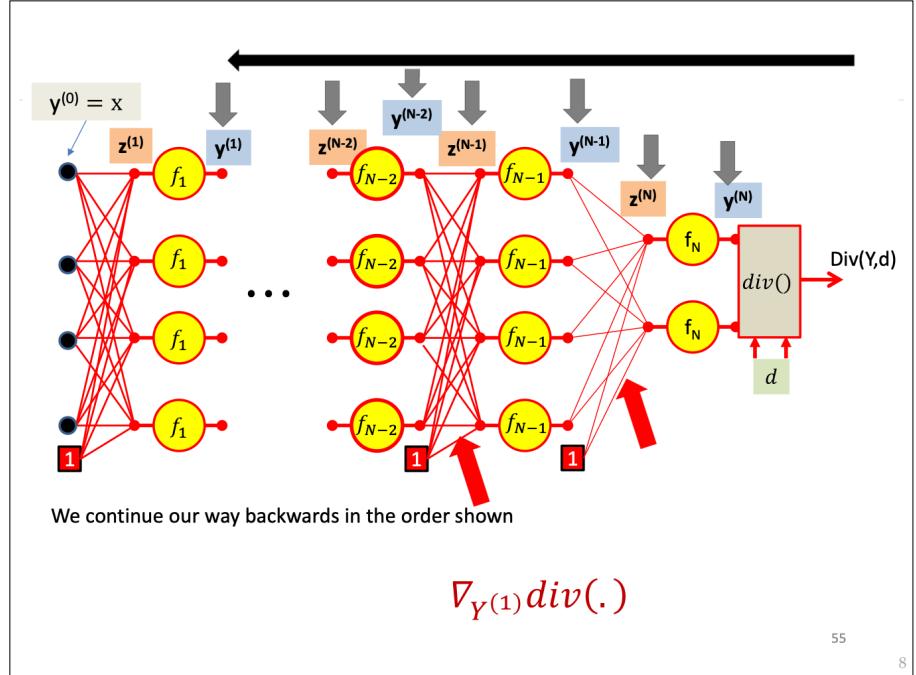
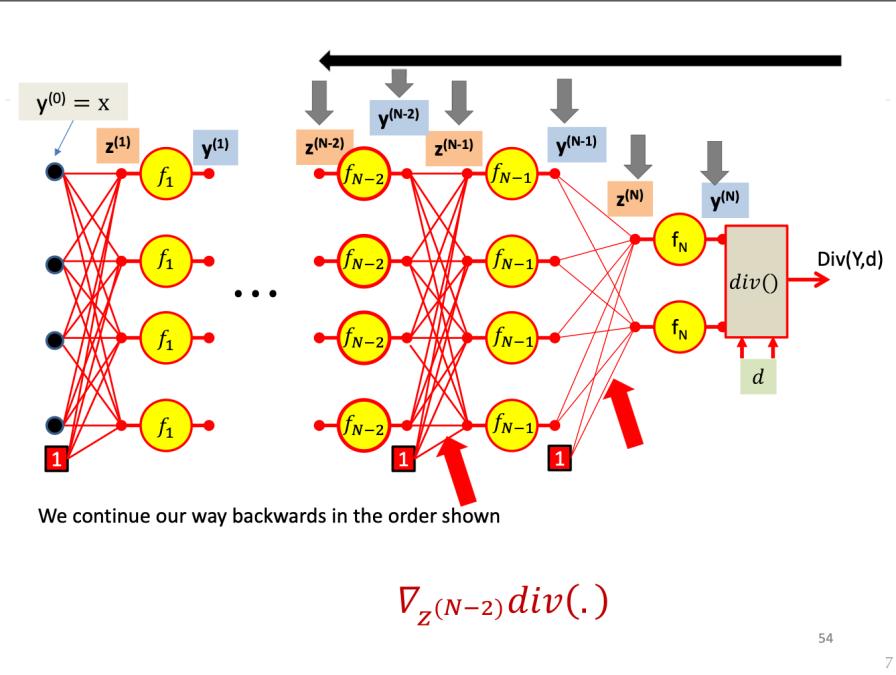


We continue our way backwards in the order shown

$$\nabla_{Y^{(N-2)}} \text{div}(\cdot)$$

53

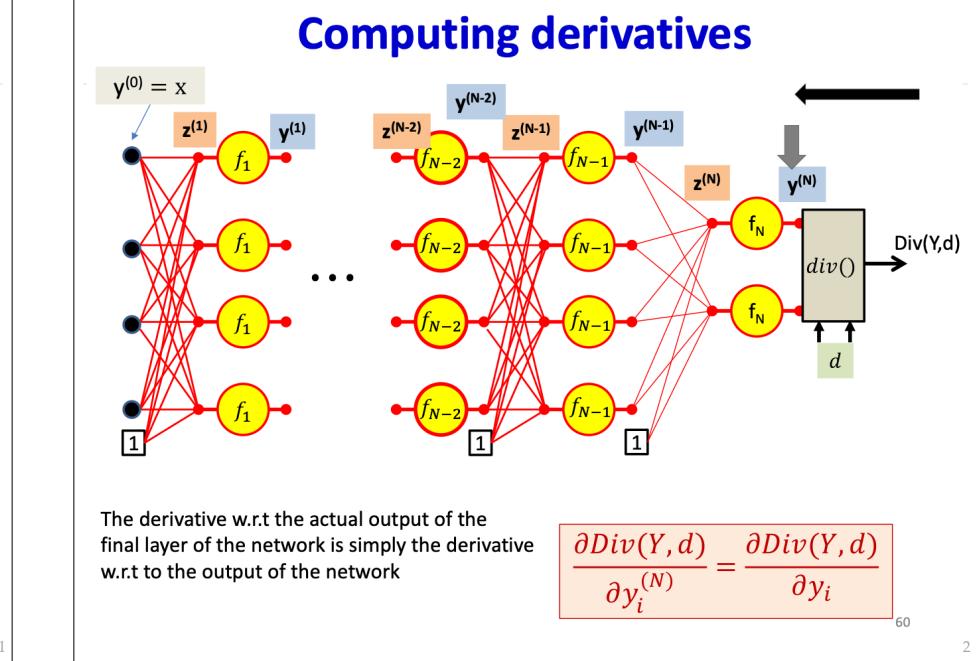
6



## Backward Gradient Computation

- Let's actually see the math..

58



## Calculus Refresher: Chain rule

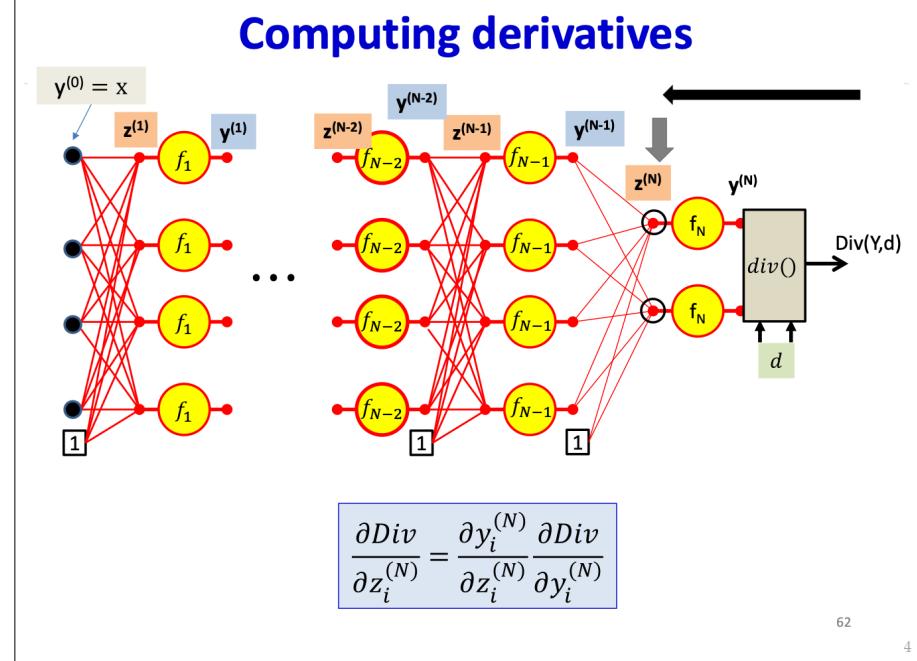
For any nested function  $l = f(y)$  where  $y = g(z)$

$$\frac{dl}{dz} = \frac{dl}{dy} \frac{dy}{dz}$$

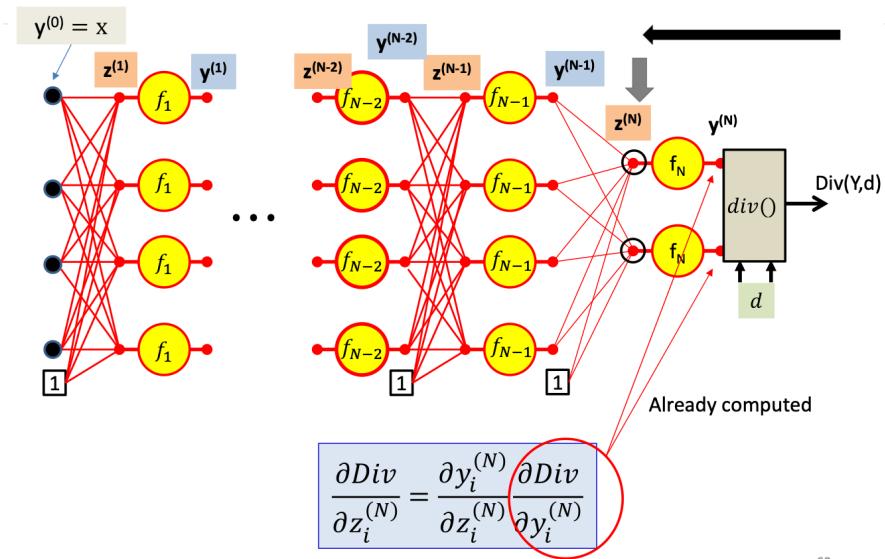
$$\frac{dl}{dz} = \frac{dl}{dy} \frac{dy}{dz}$$

$z \longrightarrow y \longrightarrow l$

3

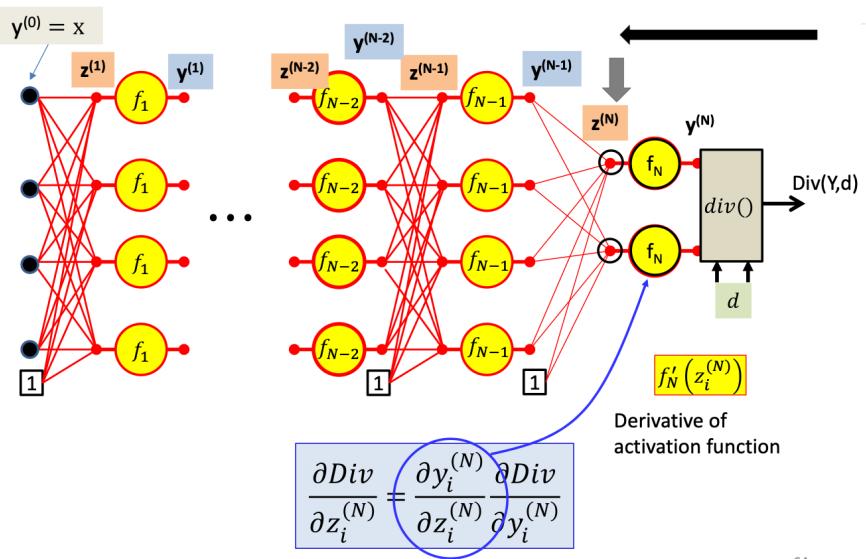


## Computing derivatives



63

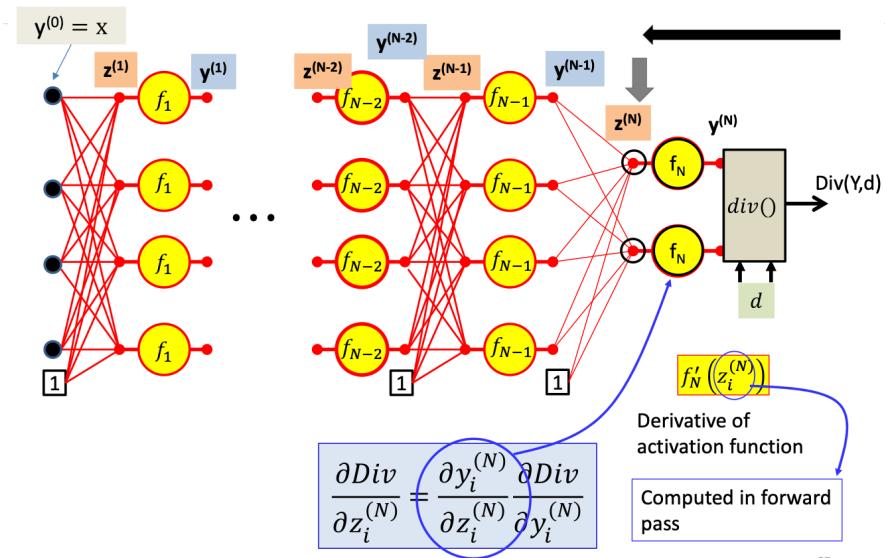
## Computing derivatives



64

6

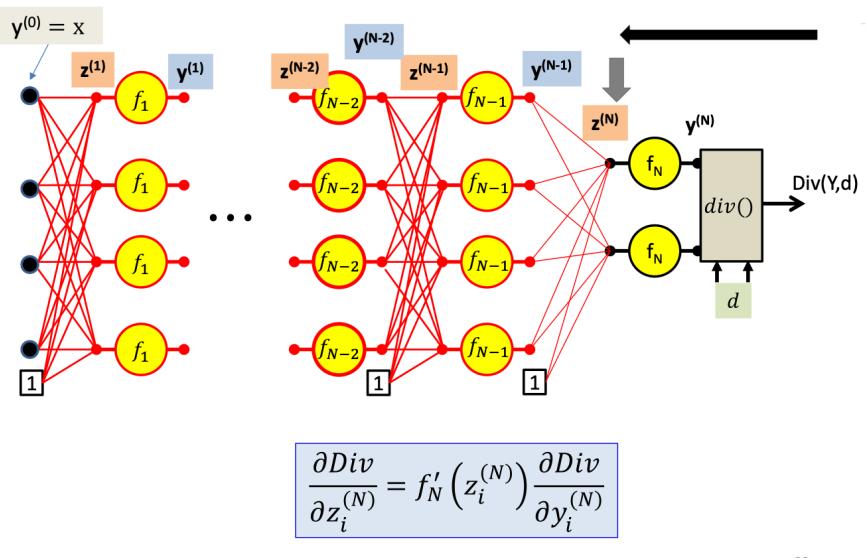
## Computing derivatives



65

7

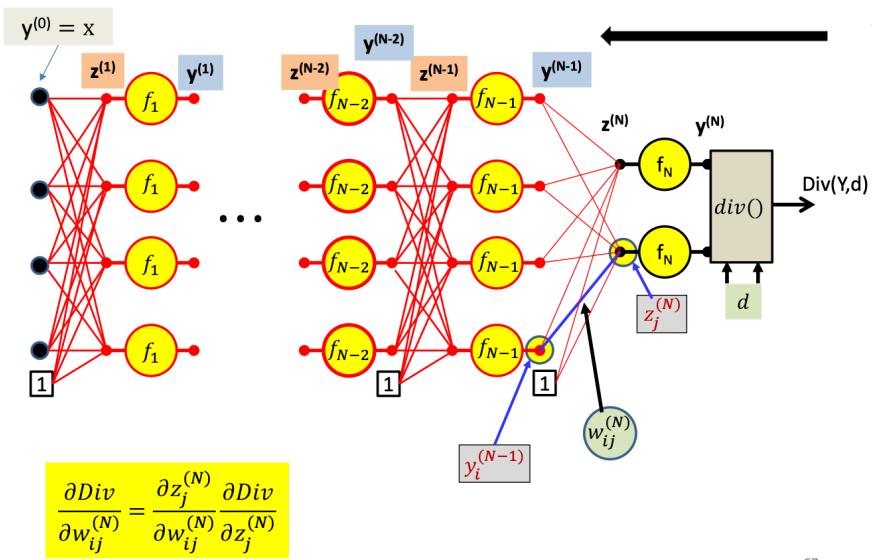
## Computing derivatives



66

8

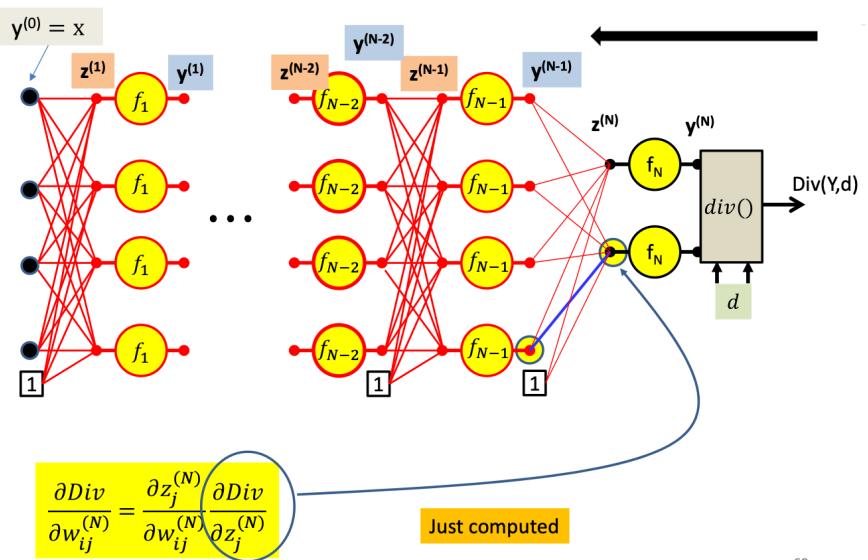
## Computing derivatives



67

9

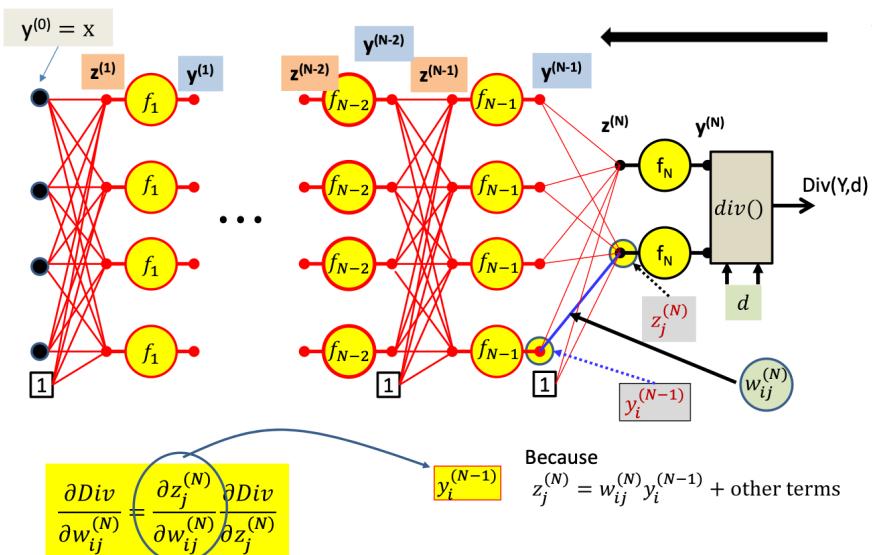
## Computing derivatives



68

0

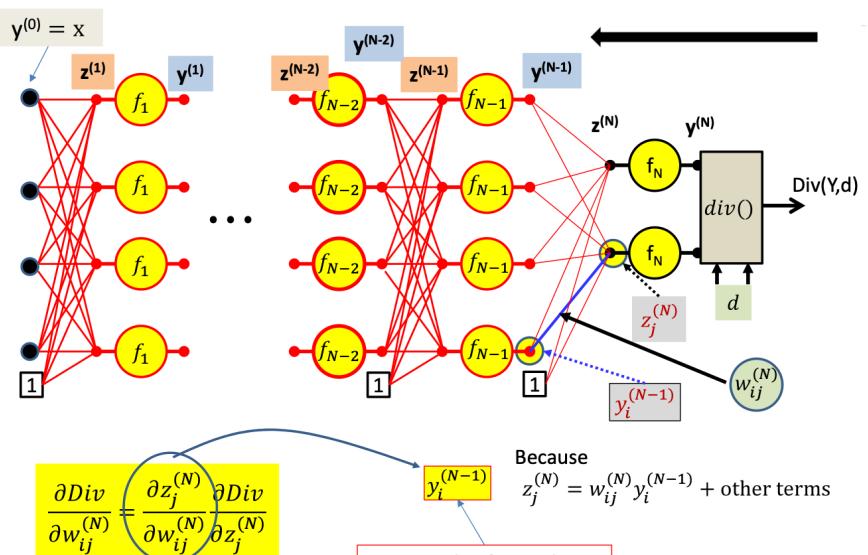
## Computing derivatives



69

1

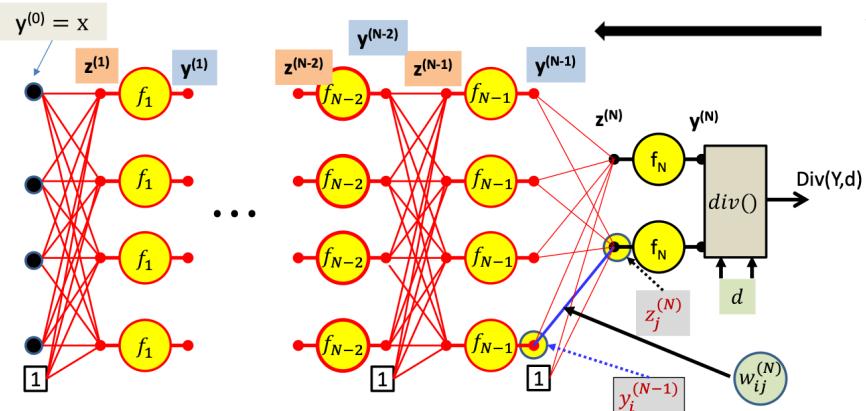
## Computing derivatives



70

2

# Computing derivatives

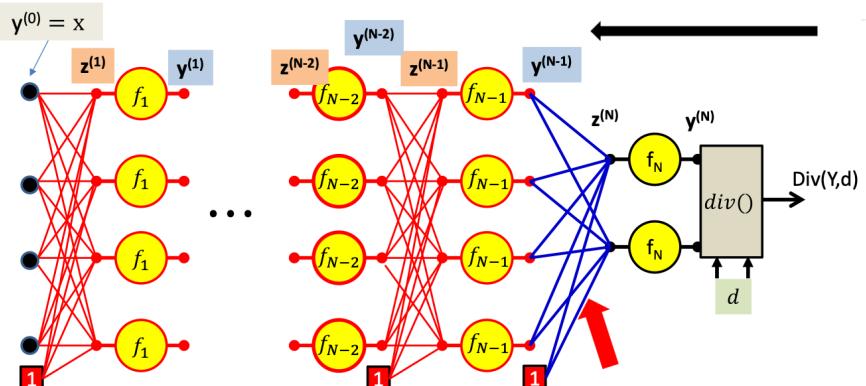


$$\frac{\partial \text{Div}}{\partial w_{ij}^{(N)}} = y_i^{(N-1)} \frac{\partial \text{Div}}{\partial z_j^{(N)}}$$

71

3

# Computing derivatives



$$\frac{\partial Div}{\partial w_{ij}^{(N)}} = y_i^{(N-1)} \frac{\partial Div}{\partial z_j^{(N)}}$$

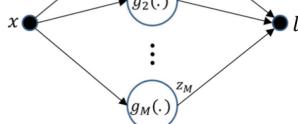
For the bias term  $y_0^{(N-1)} = 1$

72

11

# Calculus Refresher: Chain rule

**For**  $l = f(z_1, z_2, \dots, z_M)$   
**where**  $z_i = g_i(x)$

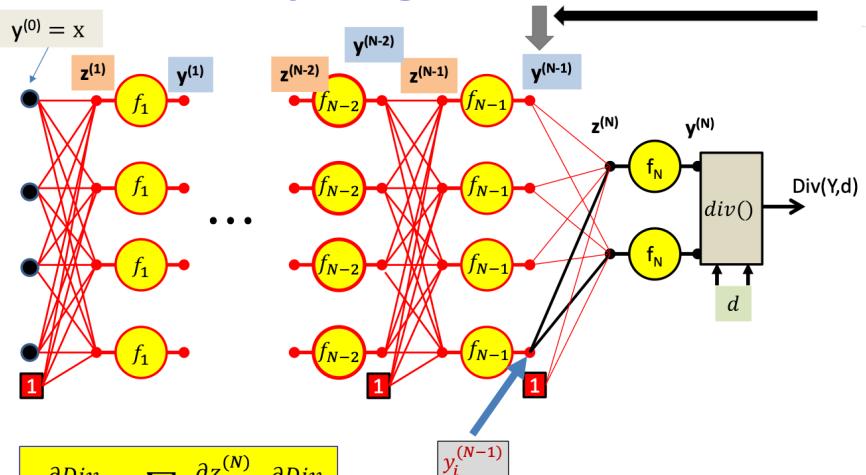


$$\frac{dl}{dx} = \frac{\partial l}{\partial z_1} \frac{dz_1}{dx} + \frac{\partial l}{\partial z_2} \frac{dz_2}{dx} + \dots + \frac{\partial l}{\partial z_M} \frac{dz_M}{dx}$$

3

15

# Computing derivatives

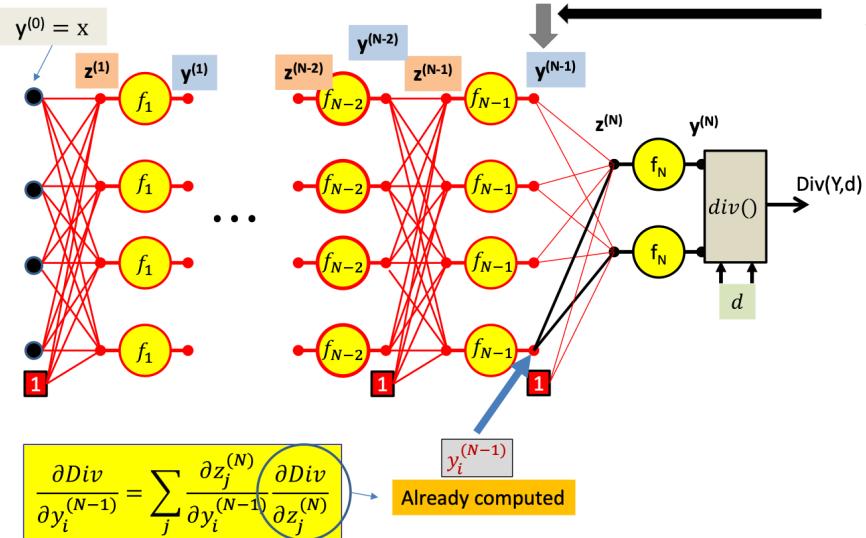


$$\frac{\partial Div}{\partial y_i^{(N-1)}} = \sum_j \frac{\partial z_j^{(N)}}{\partial y_i^{(N-1)}} \frac{\partial Div}{\partial z_j^{(N)}}$$

74

5

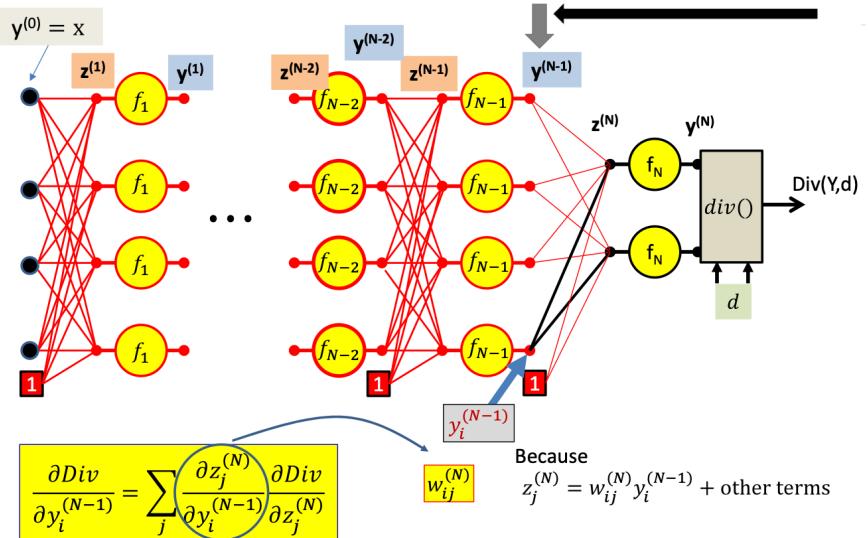
## Computing derivatives



75

7

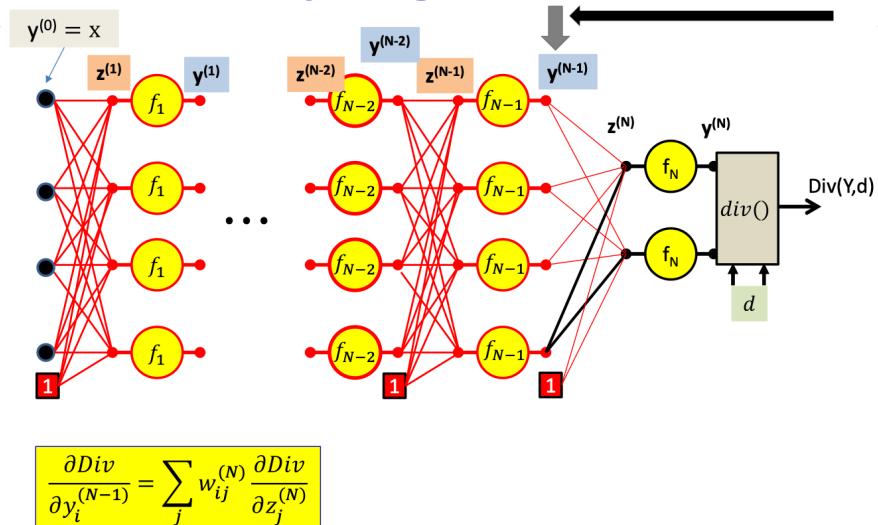
## Computing derivatives



76

8

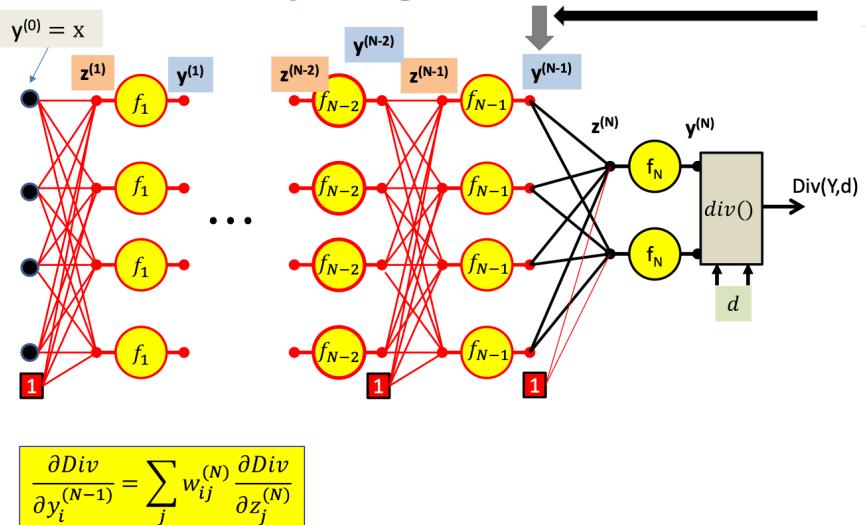
## Computing derivatives



77

9

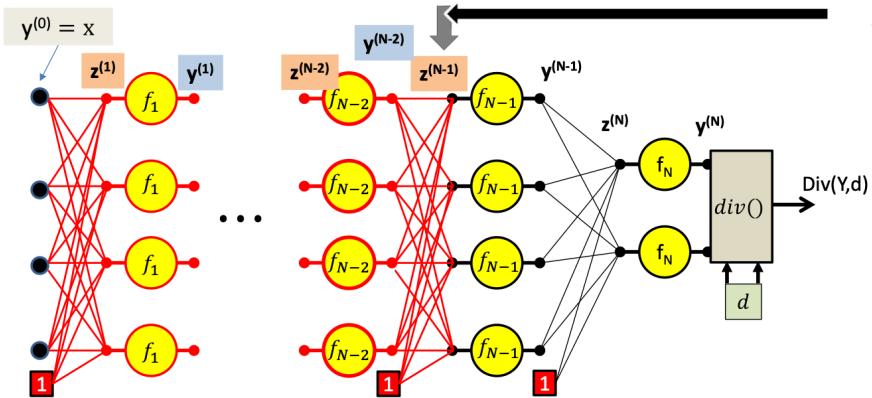
## Computing derivatives



78

0

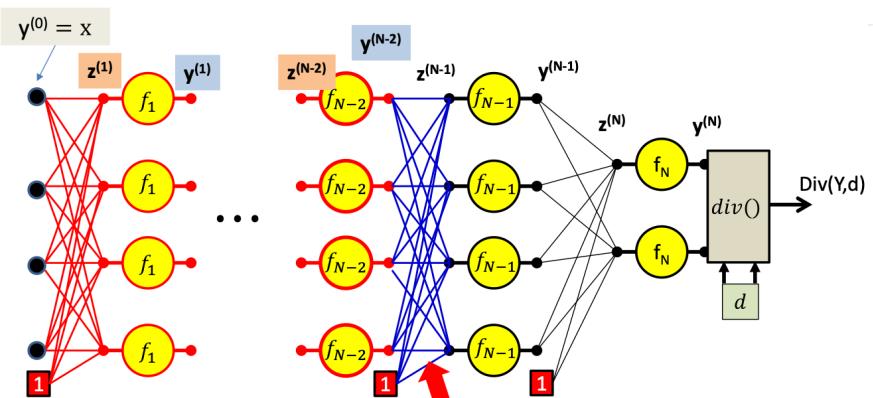
## Computing derivatives



We continue our way backwards in the order shown

$$\frac{\partial \text{Div}}{\partial z_i^{(N-1)}} = f'_{N-1}(z_i^{(N-1)}) \frac{\partial \text{Div}}{\partial y_i^{(N-1)}}$$

79

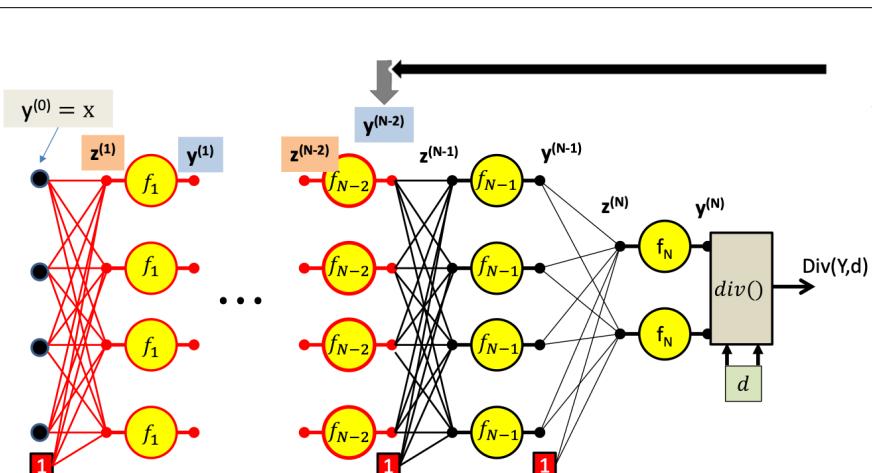


We continue our way backwards in the order shown

$$\frac{\partial \text{Div}}{\partial w_{ij}^{(N-1)}} = y_i^{(N-2)} \frac{\partial \text{Div}}{\partial z_j^{(N-1)}} \quad \text{For the bias term } y_0^{(N-2)} = 1$$

80

2

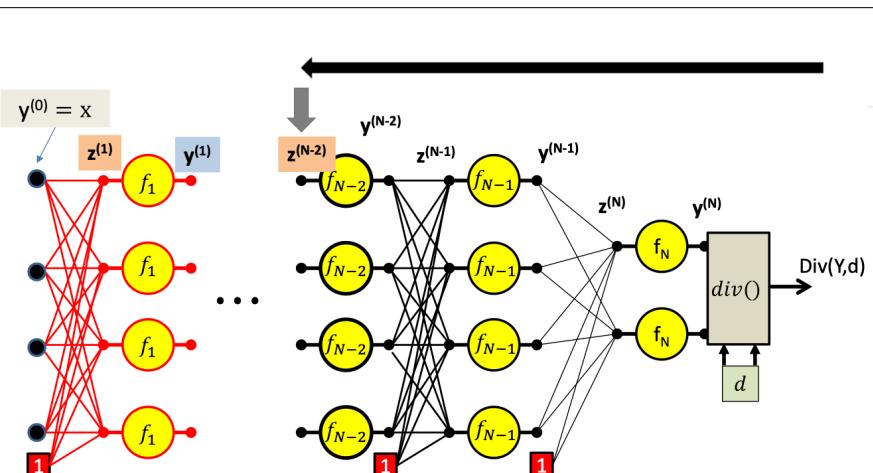


We continue our way backwards in the order shown

$$\frac{\partial \text{Div}}{\partial y_i^{(N-2)}} = \sum_j w_{ij}^{(N-1)} \frac{\partial \text{Div}}{\partial z_j^{(N-1)}}$$

81

3

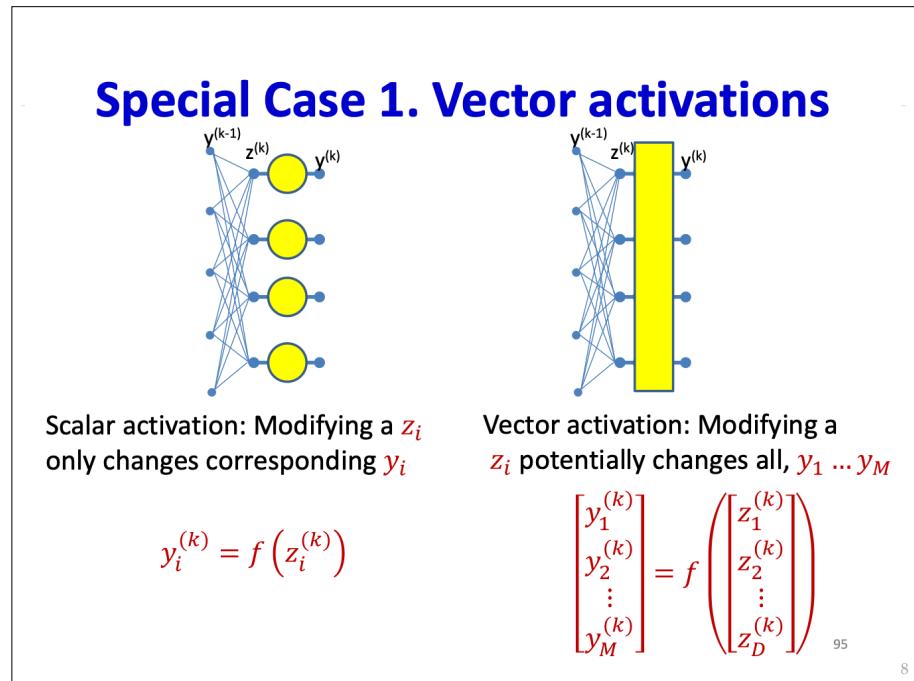
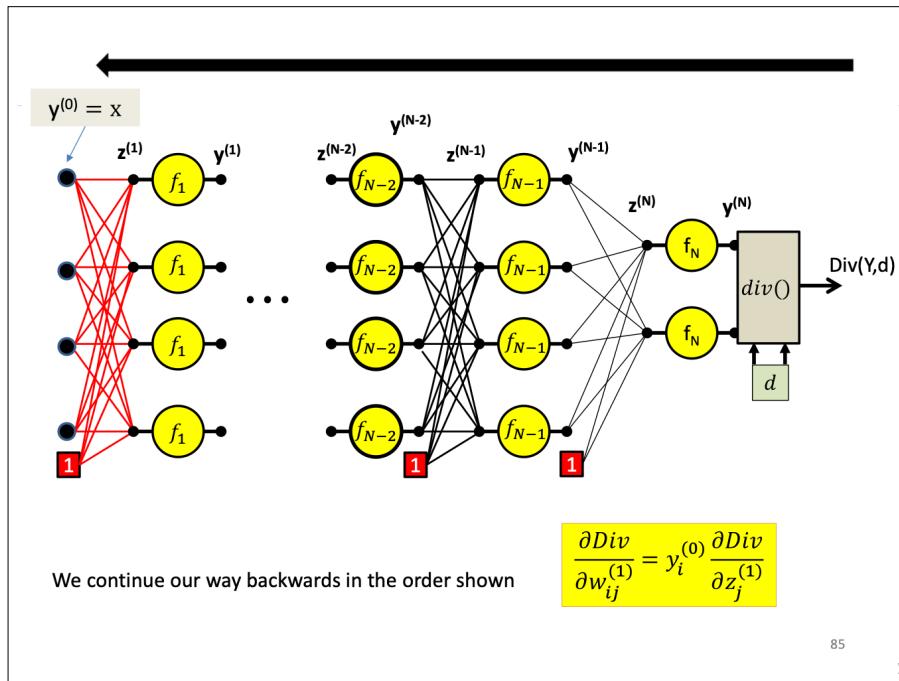
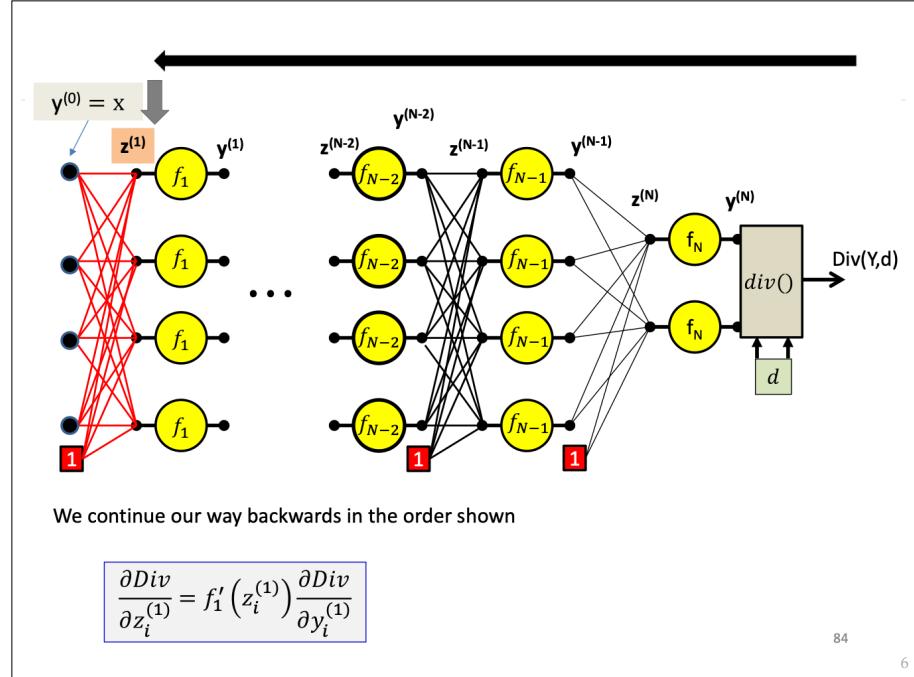
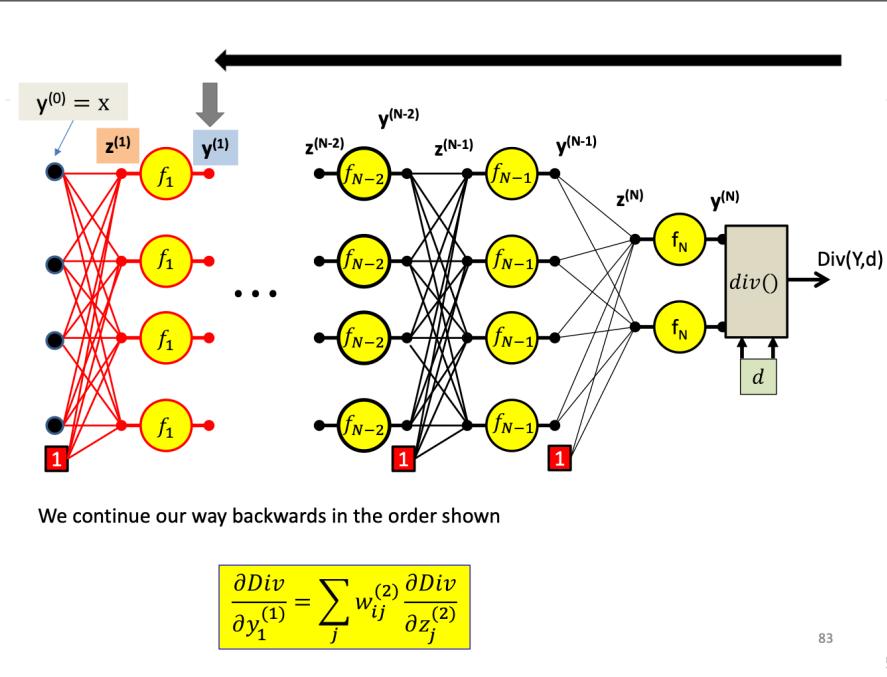


We continue our way backwards in the order shown

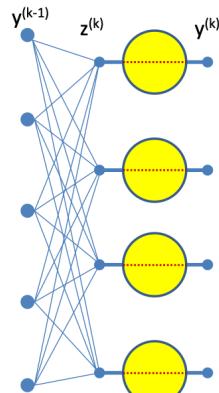
$$\frac{\partial \text{Div}}{\partial z_i^{(N-2)}} = f'_{N-2}(z_i^{(N-2)}) \frac{\partial \text{Div}}{\partial y_i^{(N-2)}}$$

82

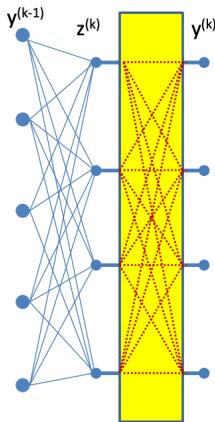
4



## “Influence” diagram



Scalar activation: Each  $z_i$  influences one  $y_i$

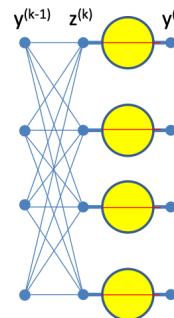


Vector activation: Each  $z_i$  influences all,  $y_1 \dots y_M$

96

9

## Scalar Activation: Derivative rule



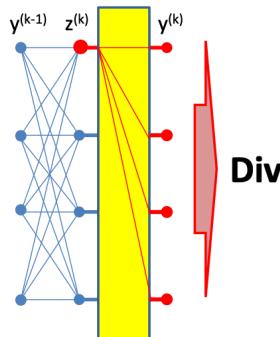
$$\frac{\partial \text{Div}}{\partial z_i^{(k)}} = \frac{\partial \text{Div}}{\partial y_i^{(k)}} \frac{\partial y_i^{(k)}}{\partial z_i^{(k)}}$$

- In the case of *scalar* activation functions, the derivative of the loss w.r.t to the input to the unit is a simple product of derivatives

97

0

## Derivatives of vector activation



$$\frac{\partial \text{Div}}{\partial z_i^{(k)}} = \sum_j \frac{\partial \text{Div}}{\partial y_j^{(k)}} \frac{\partial y_j^{(k)}}{\partial z_i^{(k)}}$$

Note: derivatives of scalar activations are just a special case of vector activations:

$$\frac{\partial y_j^{(k)}}{\partial z_i^{(k)}} = 0 \text{ for } i \neq j$$

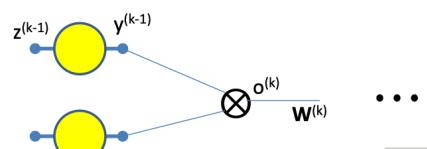
- For vector activations the derivative of the loss w.r.t. to any input is a sum of partial derivatives

– Regardless of the number of outputs  $y_j^{(k)}$

98

1

## Special Case 2: Multiplicative networks



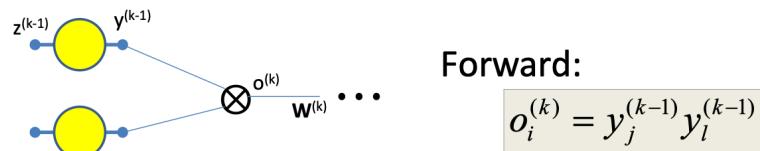
$$\text{Forward: } o_i^{(k)} = y_j^{(k-1)} y_l^{(k-1)}$$

- Some types of networks have *multiplicative* combination
  - In contrast to the *additive* combination we have seen so far
- Seen in networks such as LSTMs, GRUs, attention models, etc.

106

2

## Backpropagation: Multiplicative Networks



Backward:

$$\frac{\partial \text{Div}}{\partial o_i^{(k)}} = \sum_j w_{ij}^{(k+1)} \frac{\partial \text{Div}}{\partial z_j^{(k+1)}}$$

$$\frac{\partial \text{Div}}{\partial y_j^{(k-1)}} = \frac{\partial o_i^{(k)}}{\partial y_j^{(k-1)}} \frac{\partial \text{Div}}{\partial o_i^{(k)}} = y_l^{(k-1)} \frac{\partial \text{Div}}{\partial o_i^{(k)}}$$

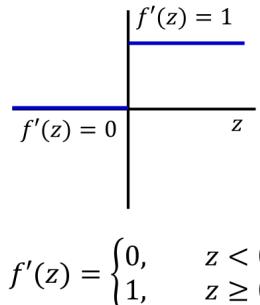
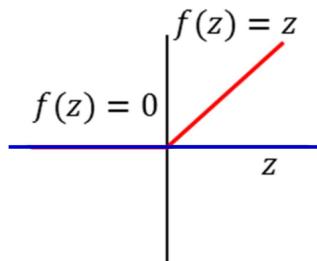
$$\frac{\partial \text{Div}}{\partial y_l^{(k-1)}} = y_j^{(k-1)} \frac{\partial \text{Div}}{\partial o_i^{(k)}}$$

- Some types of networks have *multiplicative* combination

107

3

## Non-differentiability: RELU



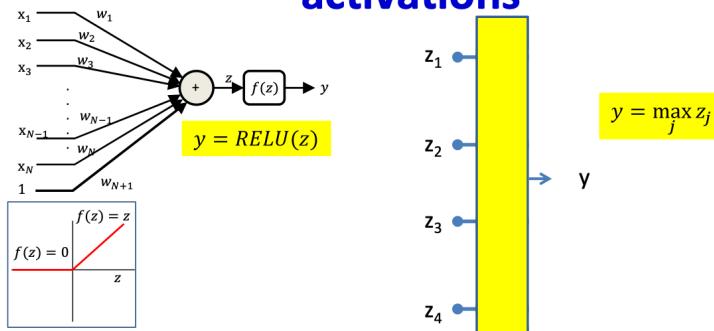
$$\Delta f(z) = \alpha \Delta z$$

- At 0 a *negative* perturbation  $\Delta z < 0$  results in no change of  $f(z)$ 
  - $\alpha = 0$
- A *positive* perturbation  $\Delta z > 0$  results in  $\Delta f(z) = \Delta z$ 
  - $\alpha = 1$
- Peering very closely, we can imagine that the curve is rotating continuously from slope = 0 to slope = 1 at  $z = 0$ 
  - So any slope between 0 and 1 is valid

113

5

## Special Case : Non-differentiable activations

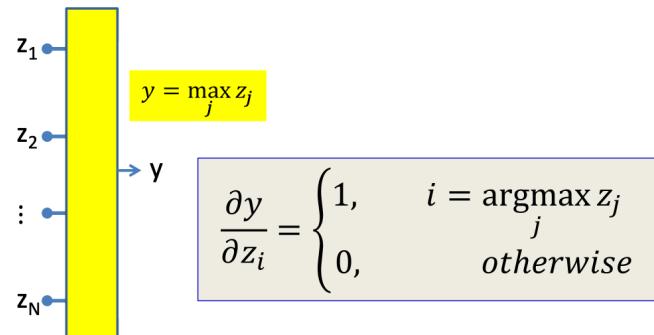


- Activation functions are sometimes not actually differentiable
  - E.g. The RELU (Rectified Linear Unit)
    - And its variants: leaky RELU, randomized leaky RELU
  - E.g. The "max" function
- Must use "subgradients" where available
  - Or "secants"

111

4

## Subgradients and the Max



- Vector equivalent of subgradient
  - 1 w.r.t. the largest incoming input
    - Incremental changes in this input will change the output
  - 0 for the rest
    - Incremental changes to these inputs will not change the output

115

6

## Overall Approach

- For each data instance
  - Forward pass:** Pass instance forward through the net. Store all intermediate outputs of all computation.
  - Backward pass:** Sweep backward through the net, iteratively compute all derivatives w.r.t weights
- Actual loss is the sum of the divergence over all training instances

$$\text{Loss} = \frac{1}{|\{X\}|} \sum_X \text{Div}(Y(X), d(X))$$

- Actual gradient is the sum or average of the derivatives computed for each training instance

$$\nabla_W \text{Loss} = \frac{1}{|\{X\}|} \sum_X \nabla_W \text{Div}(Y(X), d(X)) \quad W \leftarrow W - \eta \nabla_W \text{Loss}^T$$

120

7

## Training by BackProp

- Initialize weights  $\mathbf{W}^{(k)}$  for all layers  $k = 1 \dots K$
- Do: (*Gradient descent iterations*)
  - Initialize  $\text{Loss} = 0$ ; For all  $i, j, k$ , initialize  $\frac{d\text{Loss}}{dw_{i,j}^{(k)}} = 0$
  - For all  $t = 1:T$  (*Iterate over training instances*)
    - Forward pass:** Compute
      - Output  $\mathbf{Y}_t$
      - $\text{Loss} += \text{Div}(\mathbf{Y}_t, d_t)$
    - Backward pass:** For all  $i, j, k$ :
      - Compute  $\frac{d\text{Div}(\mathbf{Y}_t, d_t)}{dw_{i,j}^{(k)}}$
      - $\frac{d\text{Loss}}{dw_{i,j}^{(k)}} += \frac{d\text{Div}(\mathbf{Y}_t, d_t)}{dw_{i,j}^{(k)}}$
    - For all  $i, j, k$ , update:
- $w_{i,j}^{(k)} = w_{i,j}^{(k)} - \frac{\eta}{T} \frac{d\text{Loss}}{dw_{i,j}^{(k)}}$
- Until  $\text{Loss}$  has converged

121

8

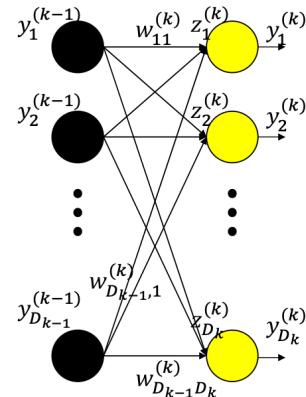
## Vector formulation

- For layered networks it is generally simpler to think of the process in terms of vector operations
  - Simpler arithmetic
  - Fast matrix libraries make operations *much* faster
- We can restate the entire process in vector terms
  - This is what is *actually* used in any real system

122

9

## Vector formulation



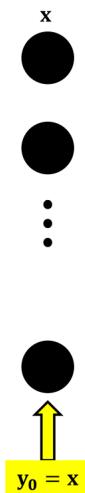
$$\mathbf{z}_k = \begin{bmatrix} z_1^{(k)} \\ z_2^{(k)} \\ \vdots \\ z_{D_k}^{(k)} \end{bmatrix} \quad \mathbf{y}_k = \begin{bmatrix} y_1^{(k)} \\ y_2^{(k)} \\ \vdots \\ y_{D_k}^{(k)} \end{bmatrix} \quad \mathbf{W}_k = \begin{bmatrix} w_{11}^{(k)} & w_{12}^{(k)} & \cdots & w_{D_{k-1}1}^{(k)} \\ w_{12}^{(k)} & w_{22}^{(k)} & \cdots & w_{D_{k-1}2}^{(k)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{1D_k}^{(k)} & w_{2D_k}^{(k)} & \cdots & w_{D_{k-1}D_k}^{(k)} \end{bmatrix} \quad \mathbf{b}_k = \begin{bmatrix} b_1^{(k)} \\ b_2^{(k)} \\ \vdots \\ b_{D_k}^{(k)} \end{bmatrix}$$

- Arrange the *inputs* to neurons of the  $k$ th layer as a vector  $\mathbf{z}_k$
- Arrange the outputs of neurons in the  $k$ th layer as a vector  $\mathbf{y}_k$
- Arrange the weights to any layer as a matrix  $\mathbf{W}_k$ 
  - Similarly with biases

123

0

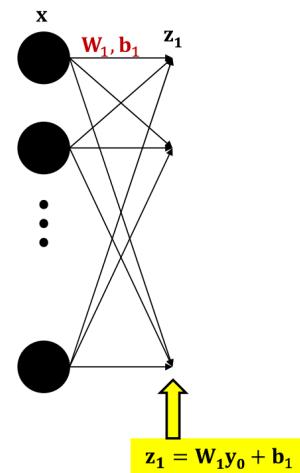
## The forward pass: Evaluating the network



125

1

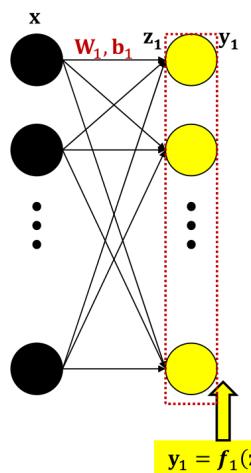
## The forward pass



126

2

## The forward pass



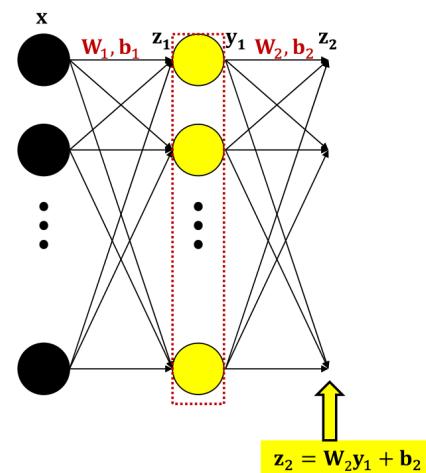
The Complete computation

$$y_1 = f_1(W_1 x + b_1)$$

127

3

## The forward pass



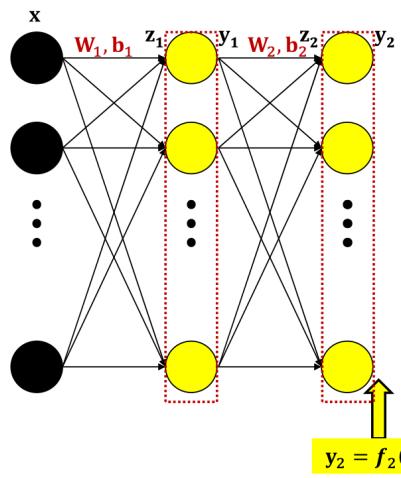
The Complete computation

$$y_1 = f_1(W_1 x + b_1)$$

128

4

## The forward pass



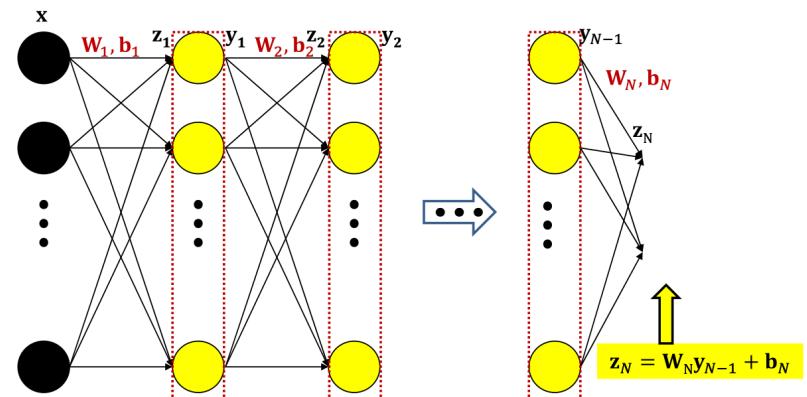
The Complete computation

$$y_2 = f_2(W_2 f_1(W_1 x + b_1) + b_2)$$

129

5

## The forward pass



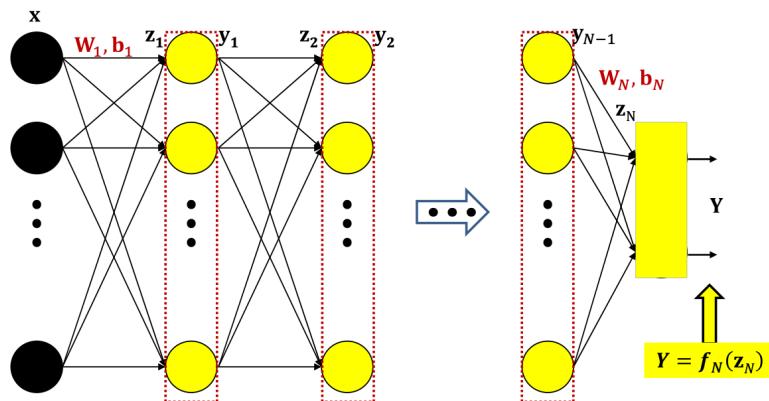
The Complete computation

$$z_N = W_N f_{N-1}(\dots f_2(W_2 f_1(W_1 x + b_1) + b_2) \dots) + b_N$$

130

6

## The forward pass



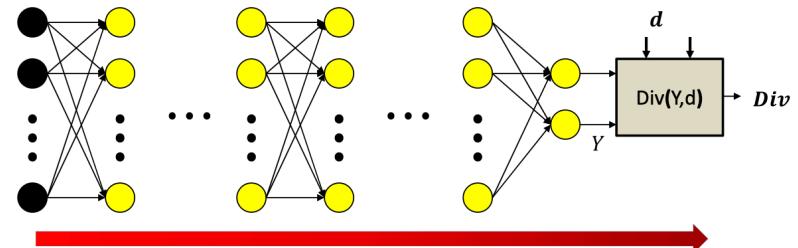
The Complete computation

$$Y = f_N(W_N f_{N-1}(\dots f_2(W_2 f_1(W_1 x + b_1) + b_2) \dots) + b_N)$$

131

7

## Forward pass



### Forward pass:

Initialize

$$\mathbf{y}_0 = \mathbf{x}$$

For  $k = 1$  to  $N$ :  $\mathbf{z}_k = \mathbf{W}_k \mathbf{y}_{k-1} + \mathbf{b}_k$      $\mathbf{y}_k = f_k(\mathbf{z}_k)$

Output

$$Y = \mathbf{y}_N$$

132

8

## The Backward Pass

- Have completed the forward pass
- Before presenting the backward pass, some more calculus...
  - Vector calculus this time

134

9

## Calculus Notes: The Jacobian

- The derivative of a vector function w.r.t. vector input is called a *Jacobian*
- It is the matrix of partial derivatives given below

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_M \end{bmatrix} = f \left( \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_D \end{bmatrix} \right)$$

Using vector notation

$$\mathbf{y} = f(\mathbf{z})$$

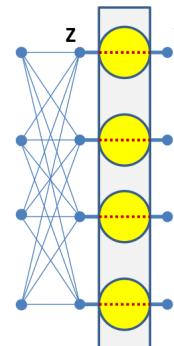
$$J_y(\mathbf{z}) = \begin{bmatrix} \frac{\partial y_1}{\partial z_1} & \frac{\partial y_1}{\partial z_2} & \dots & \frac{\partial y_1}{\partial z_D} \\ \frac{\partial y_2}{\partial z_1} & \frac{\partial y_2}{\partial z_2} & \dots & \frac{\partial y_2}{\partial z_D} \\ \dots & \dots & \ddots & \dots \\ \frac{\partial y_M}{\partial z_1} & \frac{\partial y_M}{\partial z_2} & \dots & \frac{\partial y_M}{\partial z_D} \end{bmatrix}$$

Check:  $\Delta \mathbf{y} = J_y(\mathbf{z}) \Delta \mathbf{z}$

137

0

## Jacobians can describe the derivatives of neural activations w.r.t their input



$$y_i = f(z_i)$$

$$J_y(\mathbf{z}) = \begin{bmatrix} f'(z_1) & 0 & \dots & 0 \\ 0 & f'(z_2) & \dots & 0 \\ \dots & \dots & \ddots & \dots \\ 0 & 0 & \dots & f'(z_M) \end{bmatrix}$$

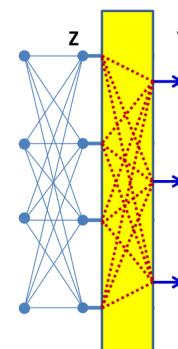
- For scalar activations (shorthand notation):

- Jacobian is a diagonal matrix
- Diagonal entries are individual derivatives of outputs w.r.t inputs

138

1

## For Vector activations



$$J_y(\mathbf{z}) = \begin{bmatrix} \frac{\partial y_1}{\partial z_1} & \frac{\partial y_1}{\partial z_2} & \dots & \frac{\partial y_1}{\partial z_D} \\ \frac{\partial y_2}{\partial z_1} & \frac{\partial y_2}{\partial z_2} & \dots & \frac{\partial y_2}{\partial z_D} \\ \dots & \dots & \ddots & \dots \\ \frac{\partial y_M}{\partial z_1} & \frac{\partial y_M}{\partial z_2} & \dots & \frac{\partial y_M}{\partial z_D} \end{bmatrix}$$

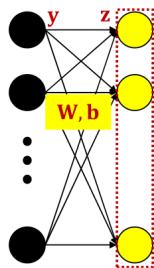
- Jacobian is a full matrix

- Entries are partial derivatives of individual outputs w.r.t individual inputs

139

2

## Special case: Affine functions



$$\mathbf{z} = \mathbf{W}\mathbf{y} + \mathbf{b}$$

↓

$$\nabla_{\mathbf{y}} \mathbf{z} = J_{\mathbf{z}}(\mathbf{y}) = \mathbf{W}$$

$$\nabla_{\mathbf{b}} \mathbf{z} = J_{\mathbf{z}}(\mathbf{b}) = \mathbf{I}$$

- Matrix **W** and bias **b** operating on vector **y** to produce vector **z**
- The Jacobian of **z** w.r.t **y** is simply the matrix **W**

140

3

## Vector Calculus Notes 2: Chain rule

- Chain rule for Jacobians:

- For vector functions of vector inputs:**

$$\mathbf{y} = \mathbf{y}(\mathbf{z}(\mathbf{x})) \rightarrow J_{\mathbf{y}}(\mathbf{x}) = J_{\mathbf{y}}(\mathbf{z})J_{\mathbf{z}}(\mathbf{x})$$

Check  $\Delta \mathbf{y} = J_{\mathbf{y}}(\mathbf{z})\Delta \mathbf{z}$

$\Delta \mathbf{z} = J_{\mathbf{z}}(\mathbf{x})\Delta \mathbf{x}$

$\Delta \mathbf{y} = J_{\mathbf{y}}(\mathbf{z})J_{\mathbf{z}}(\mathbf{x})\Delta \mathbf{x} = J_{\mathbf{y}}(\mathbf{x})\Delta \mathbf{x}$

Note the order: The derivative of the outer function comes first

143

4

## Extended Chain rule

$$x \xrightarrow{\nabla_x z_1} z_1 \xrightarrow{\nabla_{z_1} y_1} y_1 \xrightarrow{\nabla_{y_1} z_2} z_2 \xrightarrow{\nabla_{z_2} y_2} y_2 \xrightarrow{\nabla_{y_2} D} D$$

How do we compute the derivative of **D** w.r.t. **x**, **z<sub>1</sub>**, **y<sub>1</sub>**, **z<sub>2</sub>** and **y<sub>2</sub>**, from the local derivatives shown on the edges?

145

5

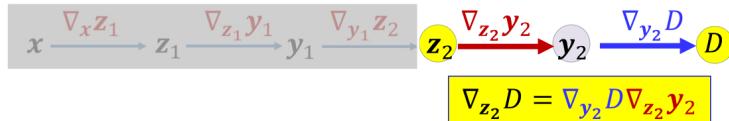
## Extended Chain rule

$$x \xrightarrow{\nabla_x z_1} z_1 \xrightarrow{\nabla_{z_1} y_1} y_1 \xrightarrow{\nabla_{y_1} z_2} z_2 \xrightarrow{\nabla_{z_2} y_2} y_2 \xrightarrow{\nabla_{y_2} D} D$$

146

6

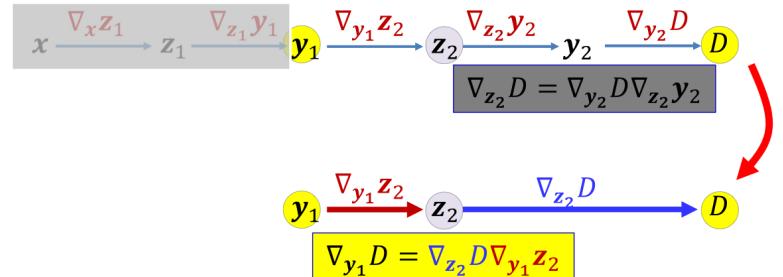
## Extended Chain rule



Note the order: The derivative of the outer function comes first

147

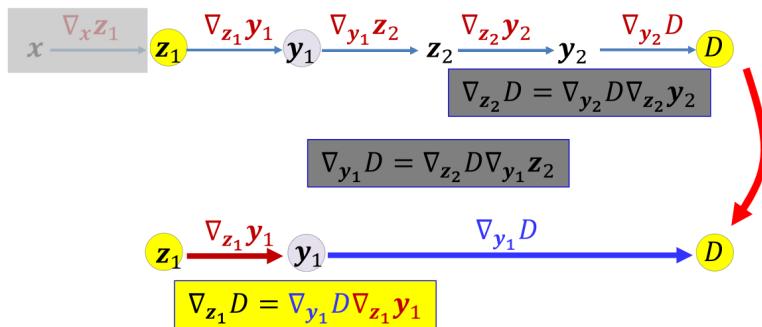
## Extended Chain rule



Note the order: The derivative of the outer function comes first

148

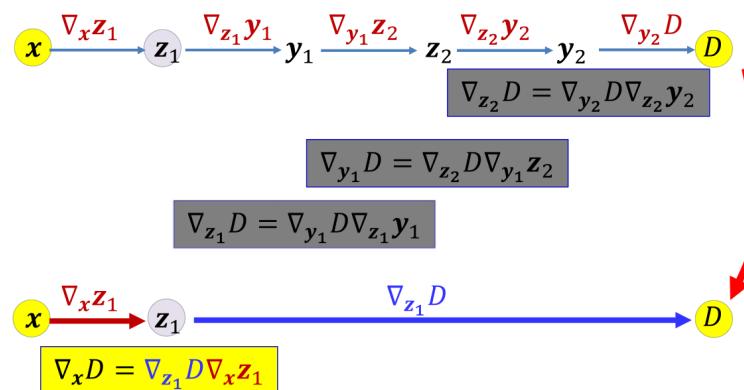
## Extended Chain rule



Note the order: The derivative of the outer function comes first

149

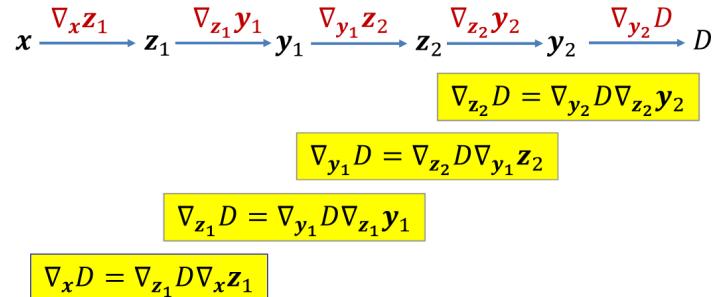
## Extended Chain rule



Note the order: The derivative of the outer function comes first

150

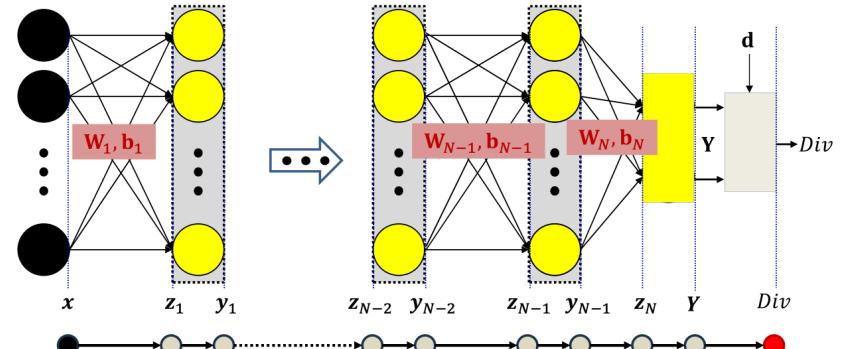
## Extended Chain rule



Note the order: The derivative of the outer function comes first

151

## The backward pass

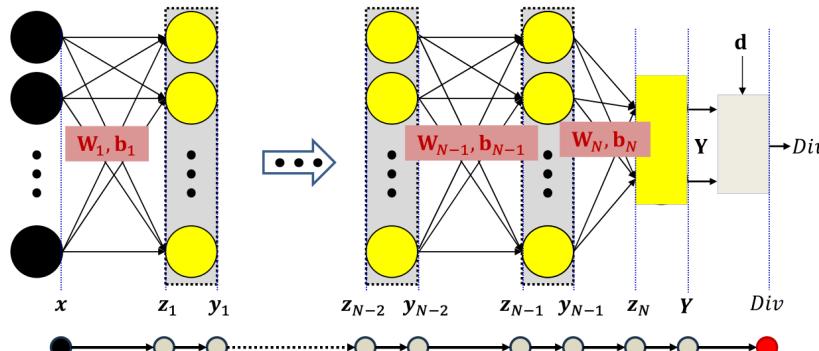


- The network again (with variables shown)...
- With the divergence we will minimize...
- And the entire influence diagram

167

12

## The backward pass

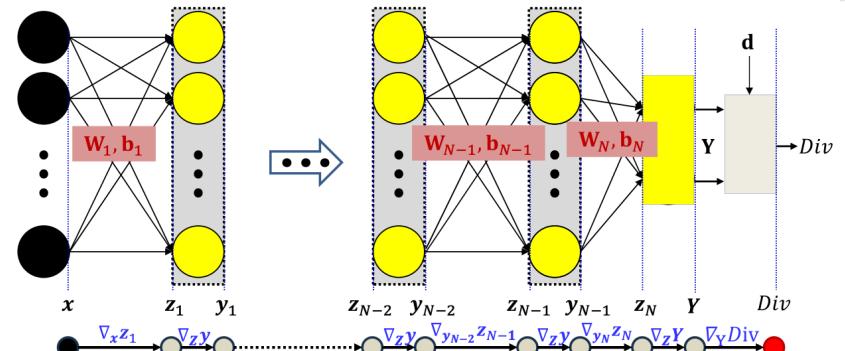


In the following slides we will also be using the notation  $\nabla_z y$  to represent the derivative of any  $y$  w.r.t any  $z$

Note that for activation functions, these are actually Jacobians

168

## The backward pass

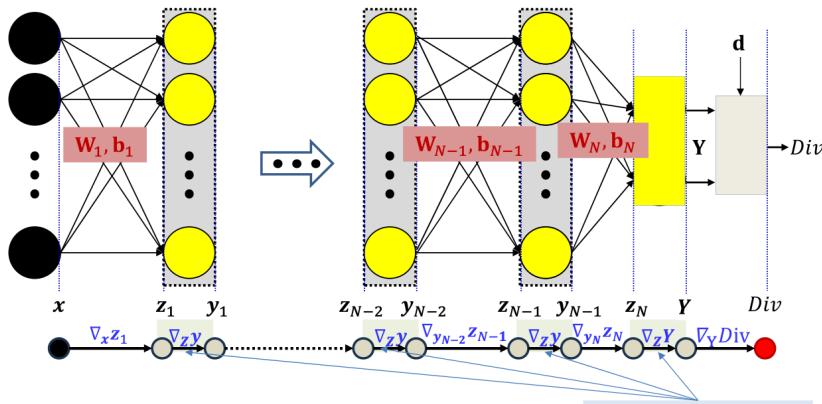


- The network again (with variables shown)...
- With the divergence we will minimize...
- And the entire influence diagram (with derivatives)
  - Variable subscripts not shown in  $\nabla_z y$  for brevity

169

14

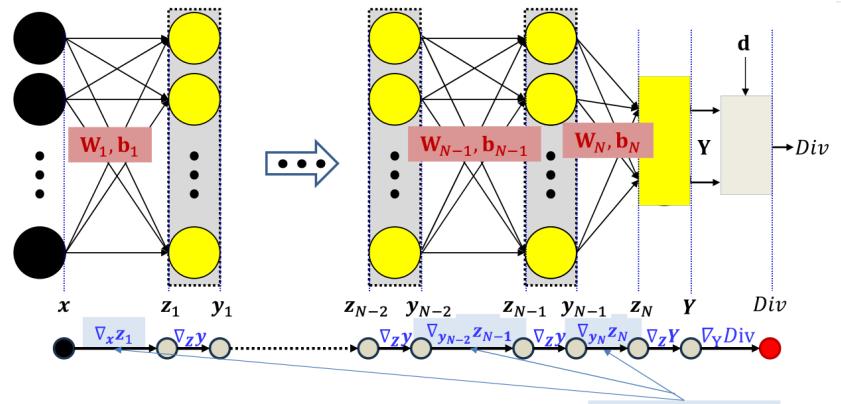
## The backward pass



- The network again (with variables shown)...
- With the divergence we will minimize...
- And the entire influence diagram (with derivatives)
  - Variable subscripts not shown in  $\nabla_z y$  for brevity

170  
15

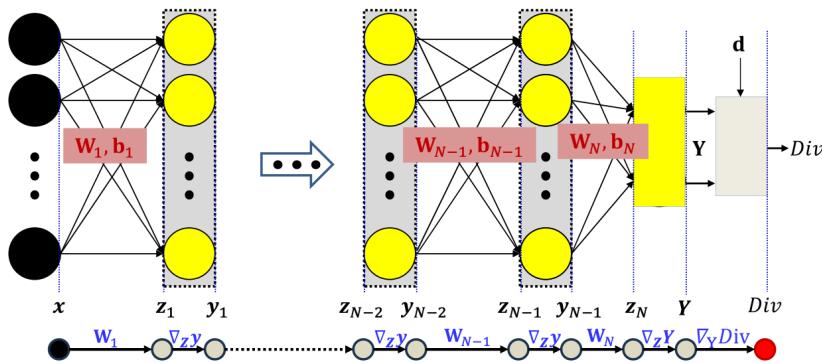
## The backward pass



- The network again (with variables shown)...
- With the divergence we will minimize...
- And the entire influence diagram (with derivatives)
  - Variable subscripts not shown in  $\nabla_z y$  for brevity

171  
16

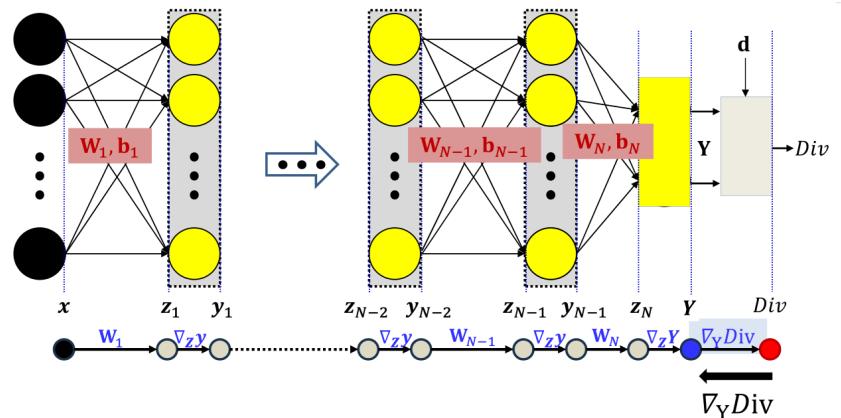
## The backward pass



- The network again (with variables shown)...
- With the divergence we will minimize...
- And the entire influence diagram (with derivatives)
  - Variable subscripts not shown in  $\nabla_z y$  for brevity

172  
17

## The backward pass

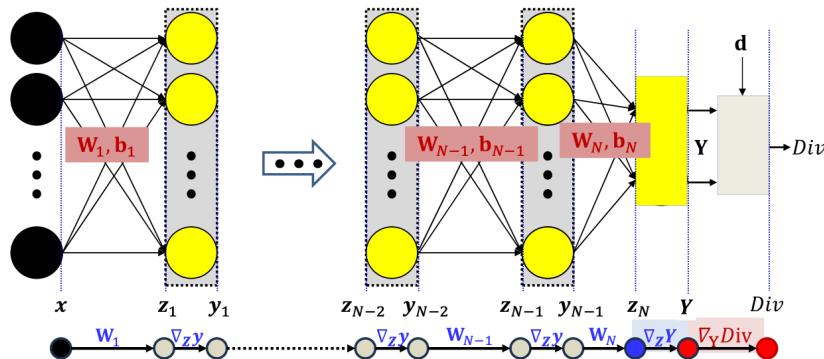


First compute the derivative of the divergence w.r.t.  $Y$ .  
The actual derivative depends on the divergence function.

N.B: The gradient is the transpose of the derivative

173  
18

## The backward pass

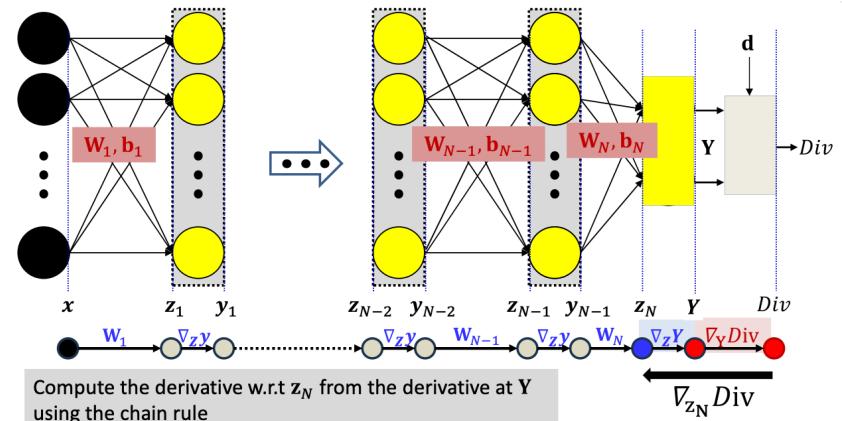


$$V_{z_N} \text{Div} = V_Y \text{Div} \cdot V_{z_N} Y$$

Already computed      New term

174

## The backward pass



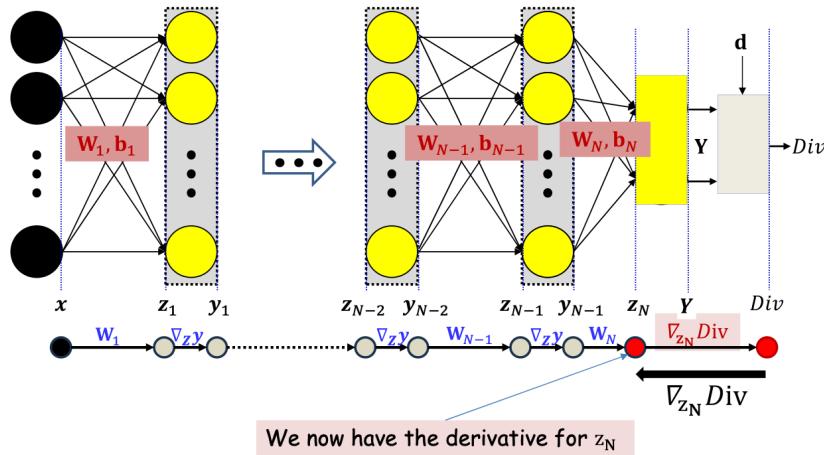
$$V_{z_N} \text{Div} = V_Y \text{Div} \cdot J_Y(z_N)$$

Already computed      Jacobian

$V_{z_N} Y$  is just the Jacobian of the activation function

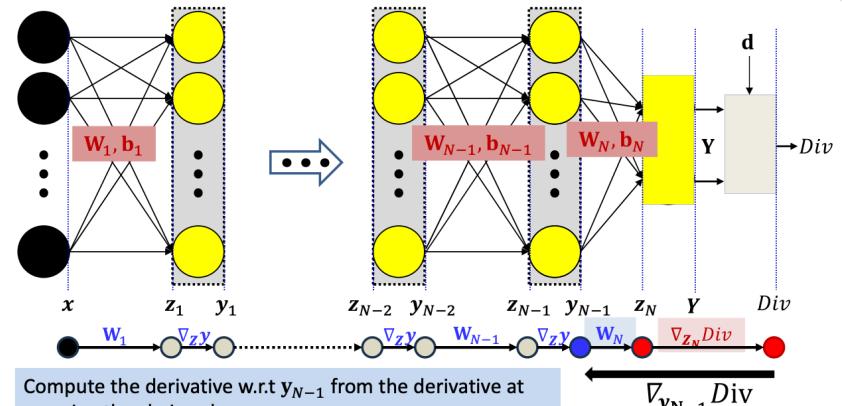
175

## The backward pass



176

## The backward pass

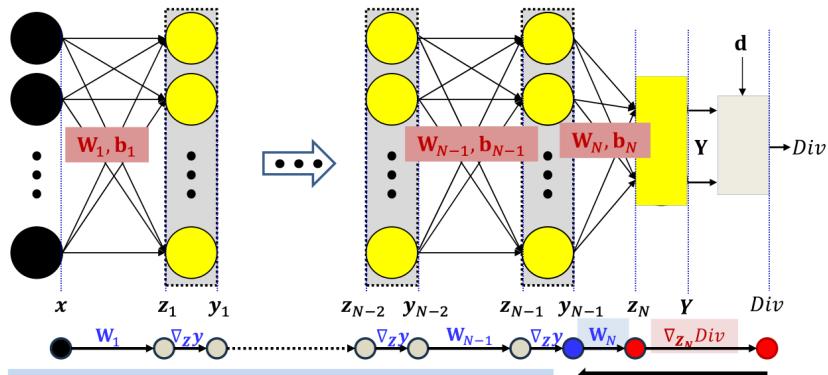


$$V_{y_{N-1}} \text{Div} = V_{z_N} \text{Div} \cdot W_N$$

$$z_N = W_N y_{N-1} + b_N \Rightarrow V_{y_{N-1}} z_N = W_N$$

177

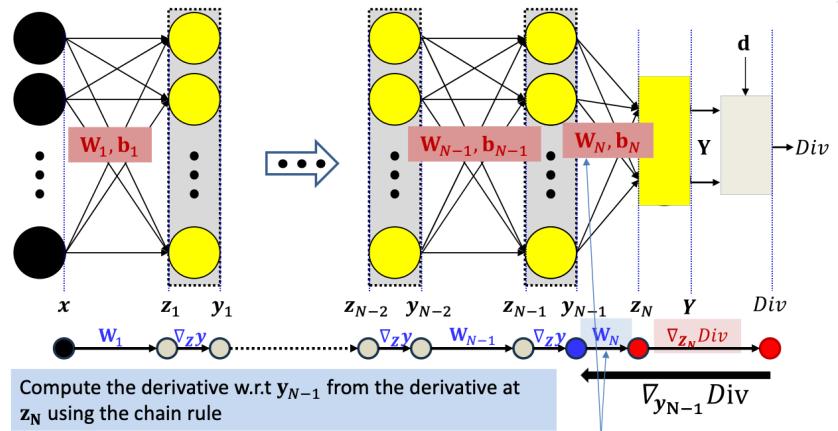
## The backward pass



$$z_N = W_N y_{N-1} + b_N \Rightarrow \nabla_{y_{N-1}} z_N = W_N$$

178

## The backward pass

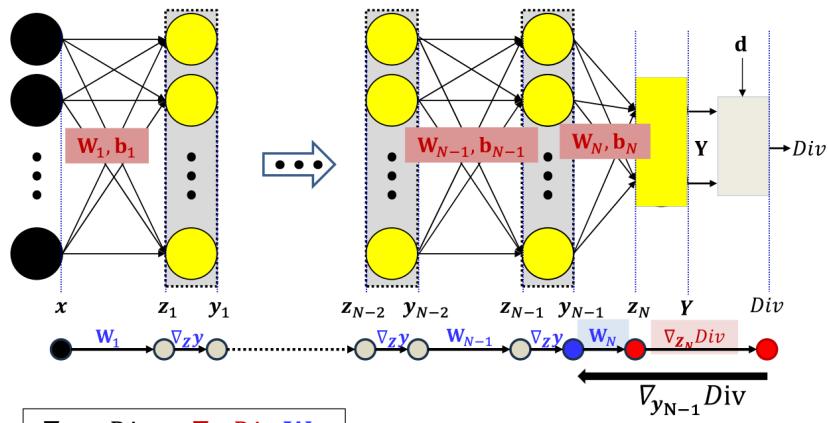


$$z_N = W_N y_{N-1} + b_N \Rightarrow \nabla_{y_{N-1}} z_N = W_N$$

179

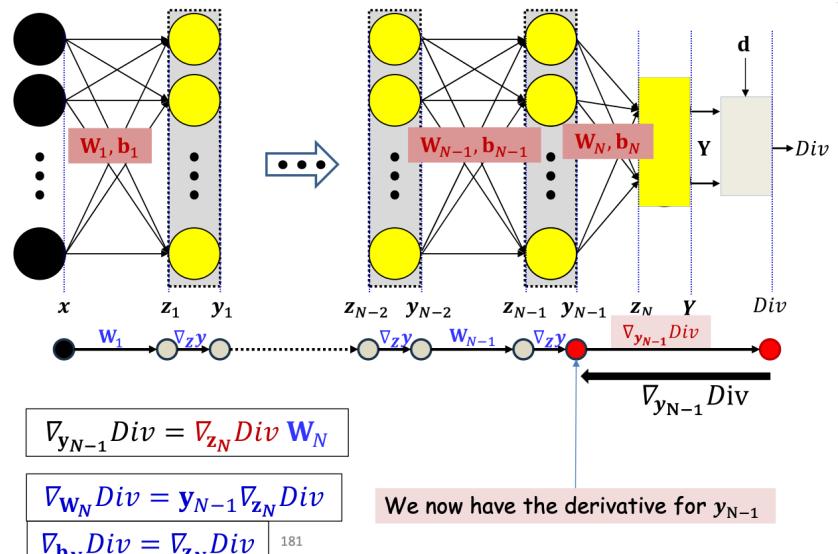
14

## The backward pass



180

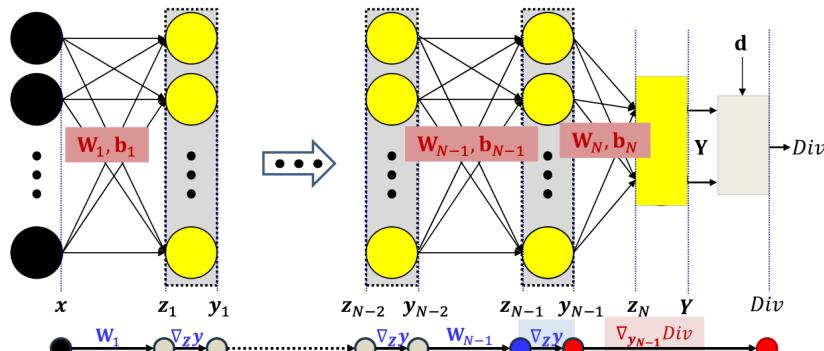
## The backward pass



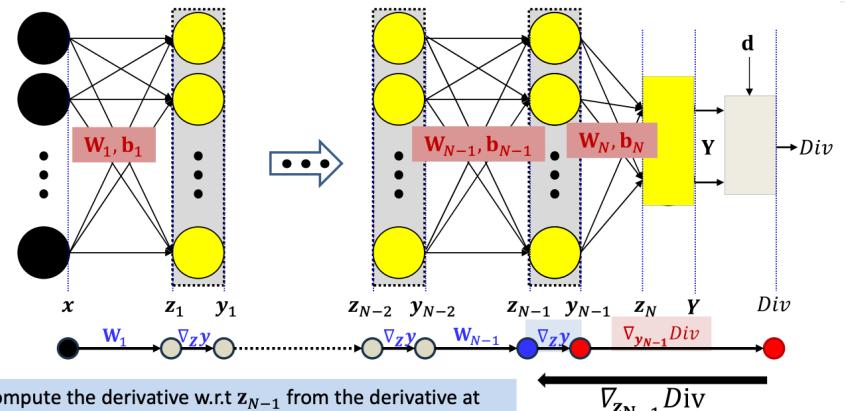
181

16

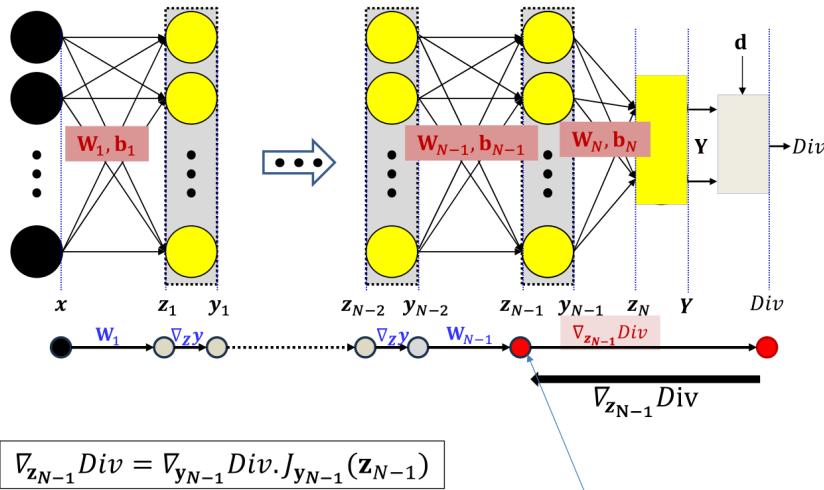
## The backward pass



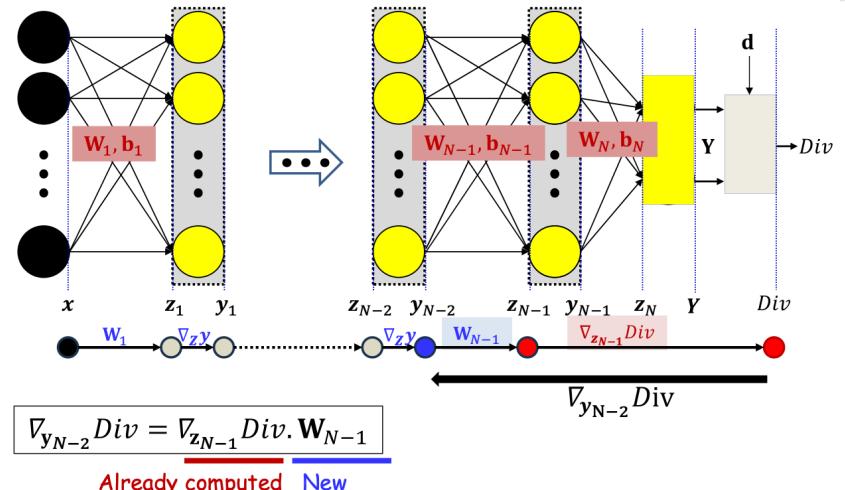
## The backward pass



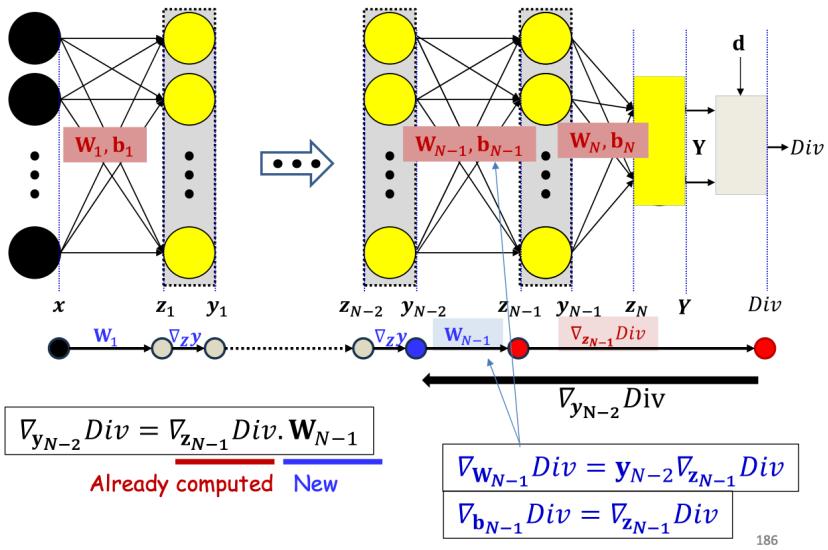
## The backward pass



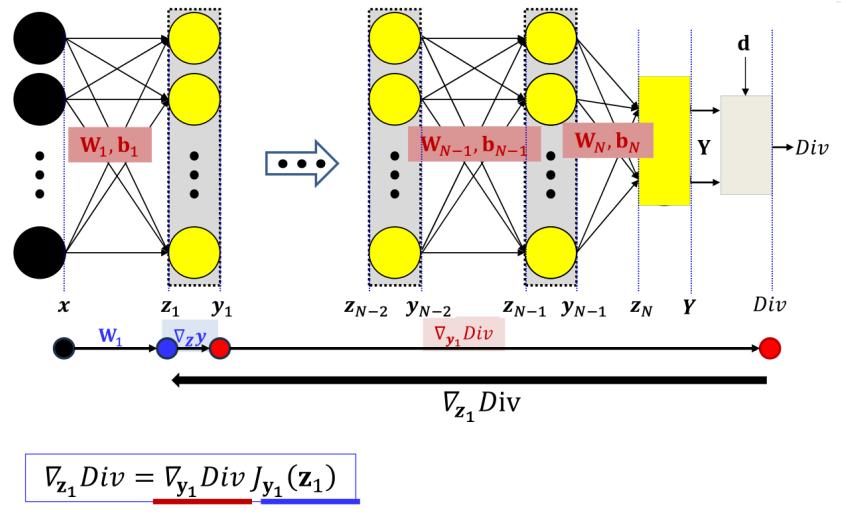
## The backward pass



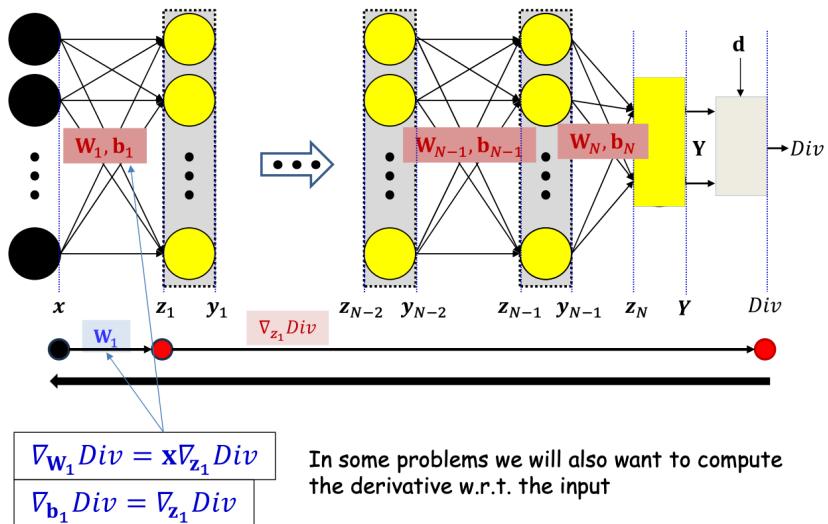
## The backward pass



## The backward pass



## The backward pass



## Neural network training algorithm

- Initialize all weights and biases ( $W_1, b_1, W_2, b_2, \dots, W_N, b_N$ )
- Do:
  - $LOSS = 0$
  - For all  $k$ , initialize  $\nabla_{W_k} Loss = 0, \nabla_{b_k} Loss = 0$
  - For all  $t = 1:T$  # Loop through training instances
    - Forward pass : Compute
      - Output  $Y(X_t)$
      - Divergence  $Div(Y_t, d_t)$
      - $Loss += Div(Y_t, d_t)$
    - Backward pass: For all  $k$  compute:
      - $\nabla_{y_k} \text{Div} = \nabla_{z_{k+1}} \text{Div} \cdot W_{k+1}$
      - $\nabla_{z_k} \text{Div} = \nabla_{y_k} \text{Div} J_{y_k}(z_k)$
      - $\nabla_{W_k} \text{Div}(Y_t, d_t) = y_{t-1} \nabla_{z_k} \text{Div}; \nabla_{b_k} \text{Div}(Y_t, d_t) = \nabla_{z_k} \text{Div}$
      - $\nabla_{W_k} Loss += \nabla_{W_k} \text{Div}(Y_t, d_t); \nabla_{b_k} Loss += \nabla_{b_k} \text{Div}(Y_t, d_t)$
    - For all  $k$ , update:
 
$$W_k = W_k - \frac{\eta}{T} (\nabla_{W_k} Loss)^T; \quad b_k = b_k - \frac{\eta}{T} (\nabla_{b_k} Loss)^T$$
  - Until  $Loss$  has converged

## Issues

- Convergence: How well does it learn
  - And how can we improve it
- How well will it generalize (outside training data)
- What does the output really mean?
- *Etc..*