

# Under the hood of autodiff

Neurons, fully connected networks, computational graphs,  
Jacobians and vector-Jacobian product

Karel Zimmermann

Czech Technical University in Prague

Faculty of Electrical Engineering, Department of Cybernetics

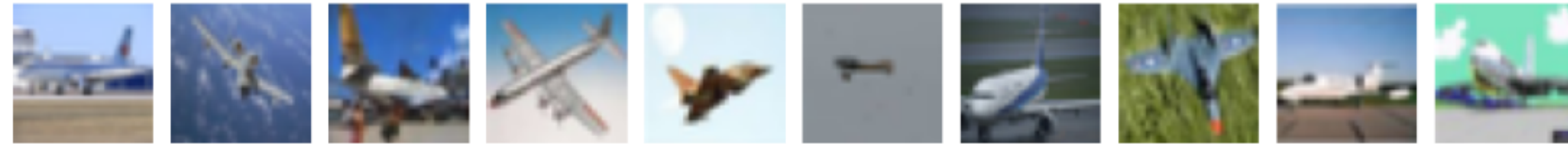


# Linear classifier and neuron

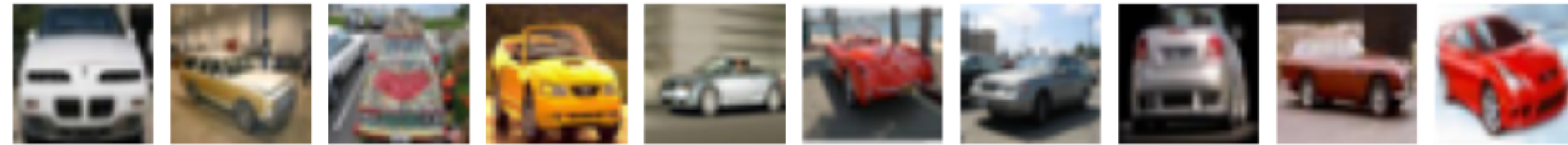
Labels

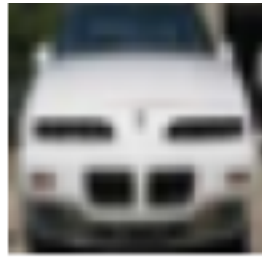
RGB images

+1

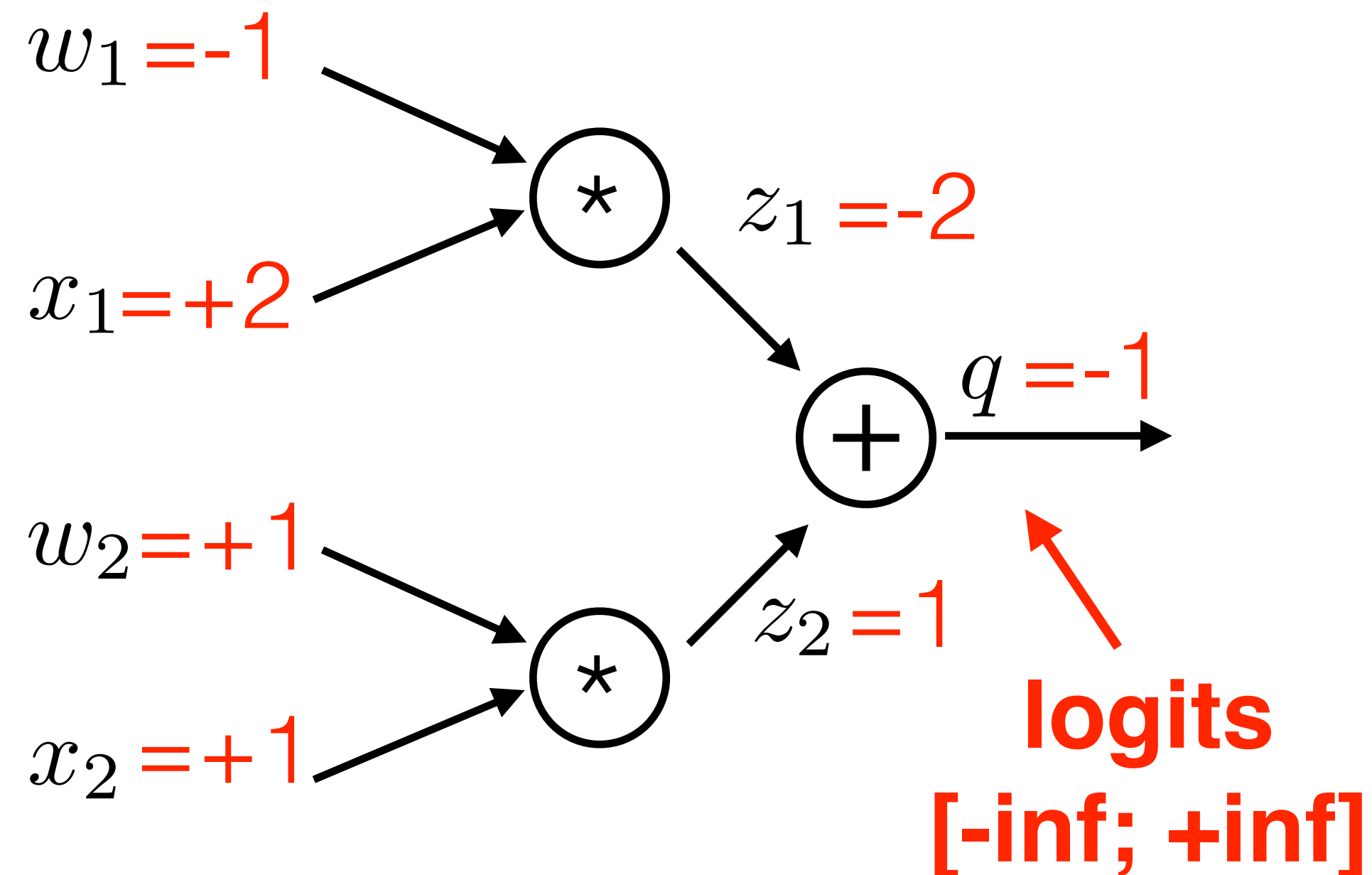


-1



```
def classify(  
    return  $\mathbf{w}^T \mathbf{x}$ 
```

Computational graph of linear classifier



# Linear classifier and neuron

Labels


RGB images

+1



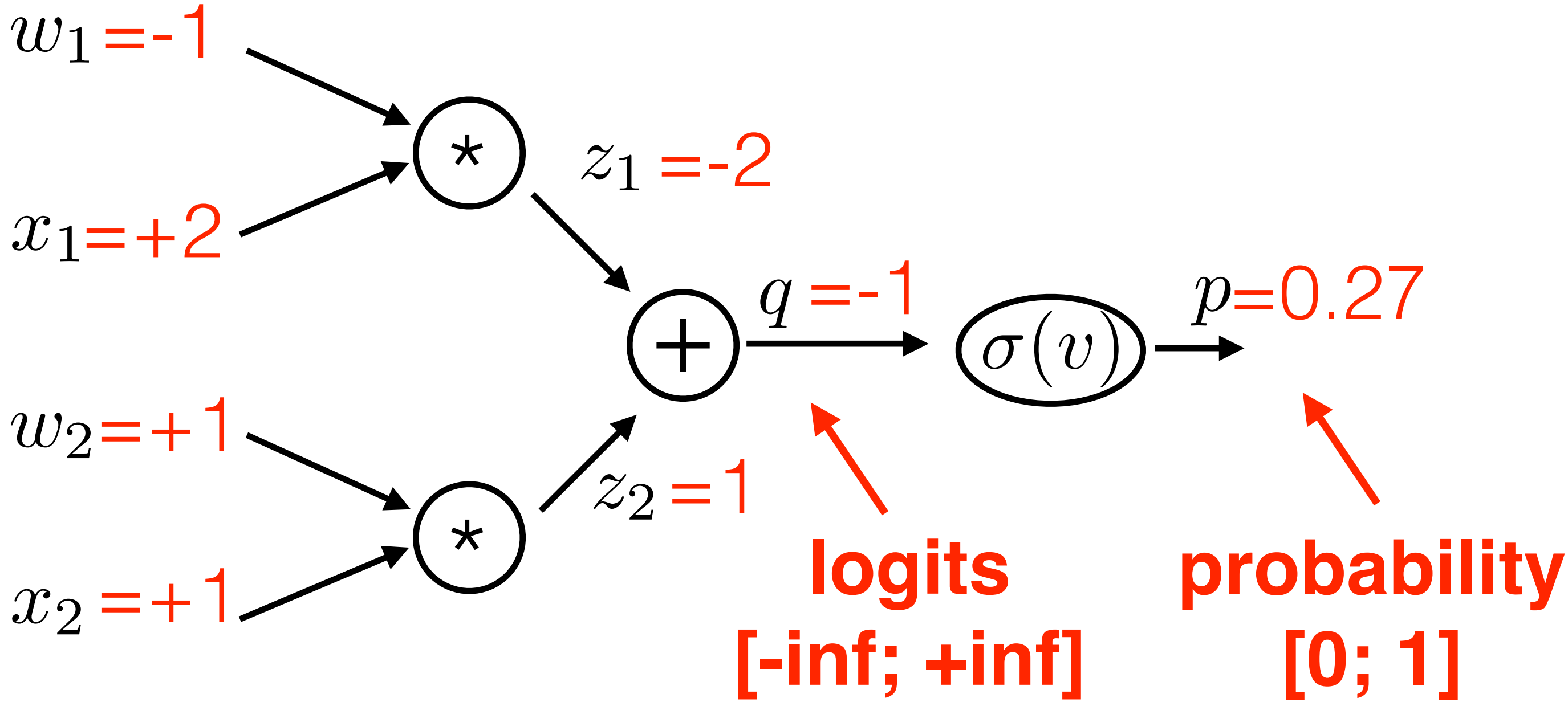
-1



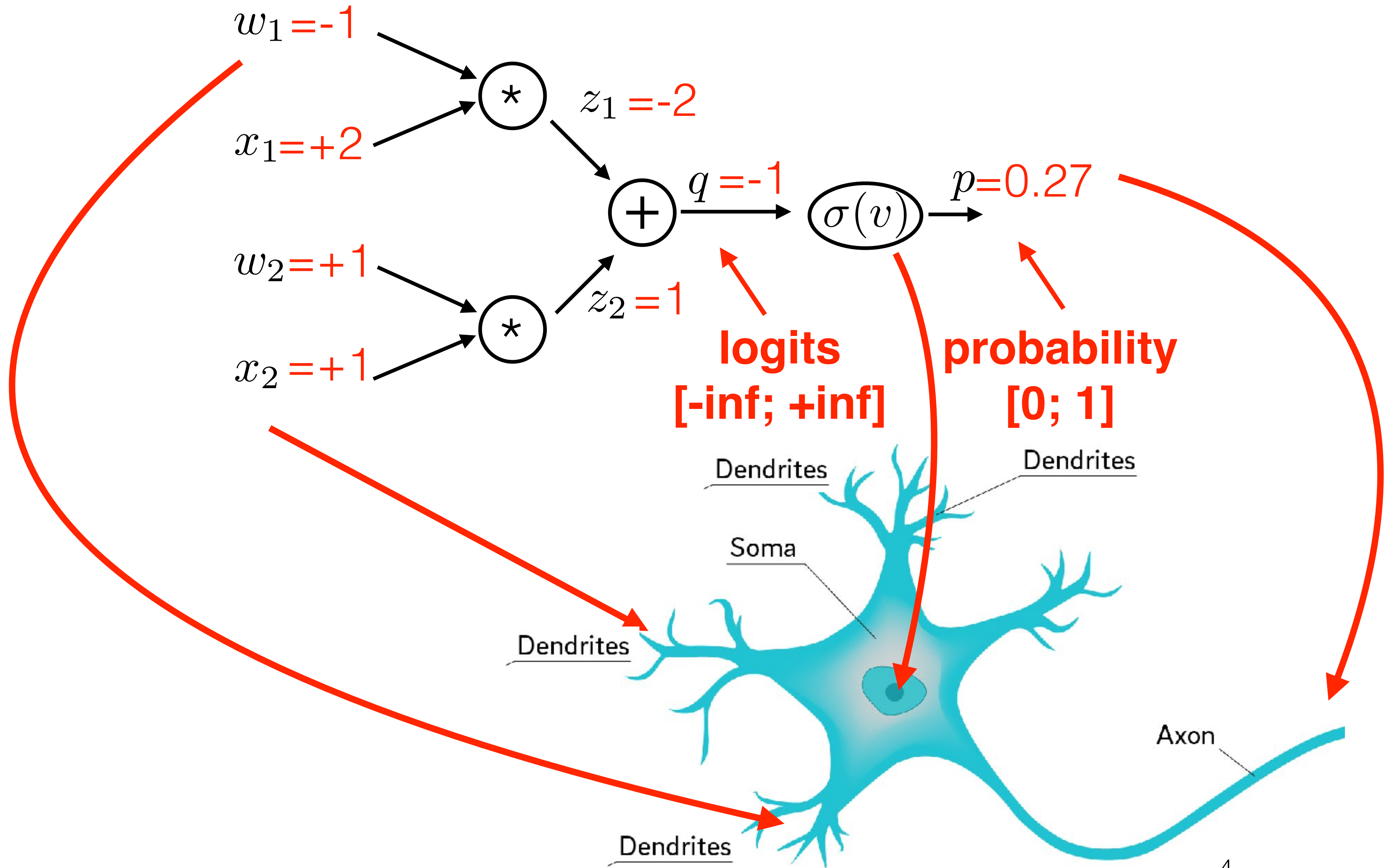
```
def classify(

Computational graph of neuron


```



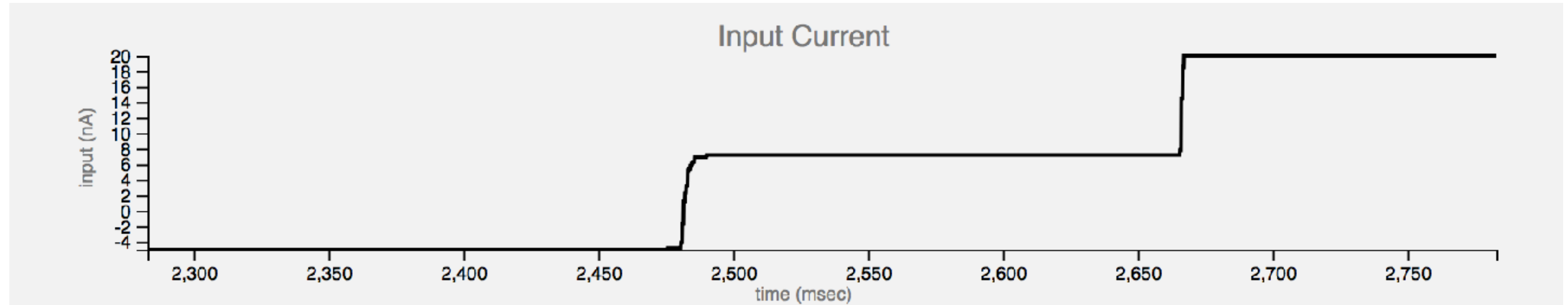
# Relation to biological neuron



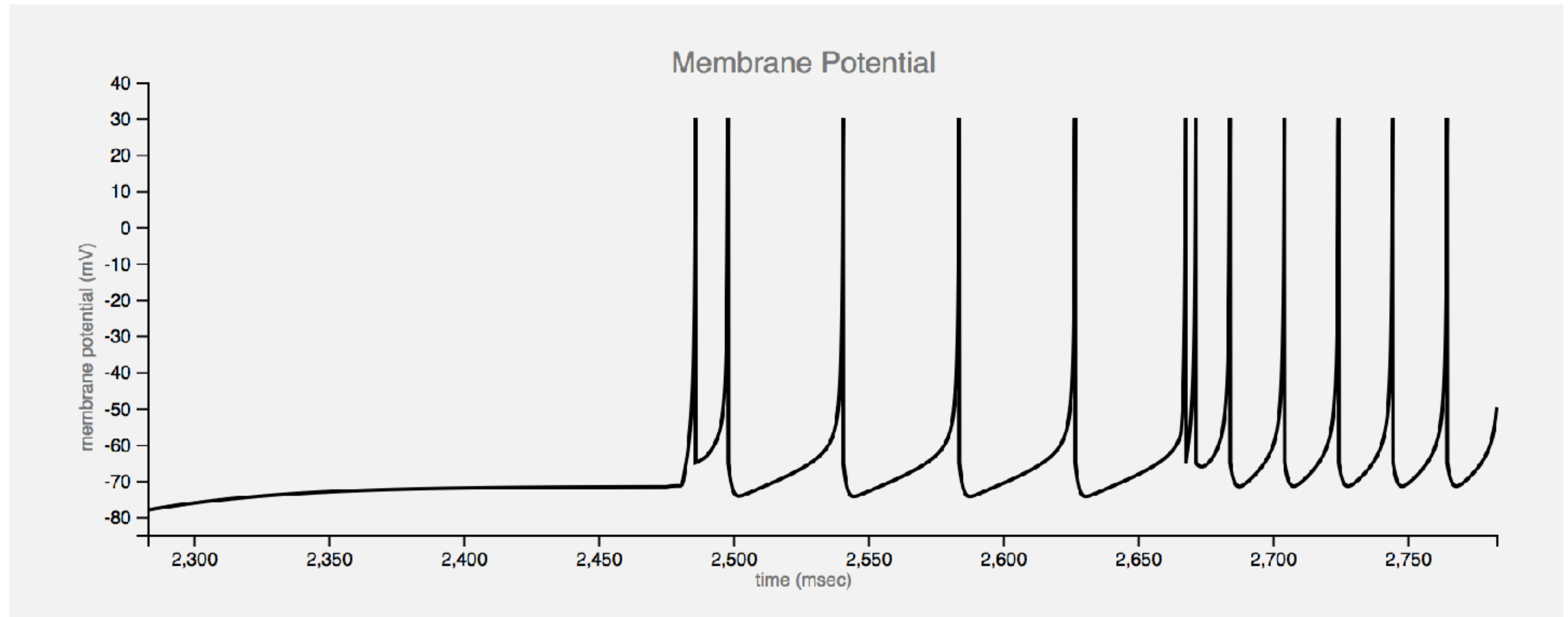


# Modeling dynamic neuron behaviour

Input:

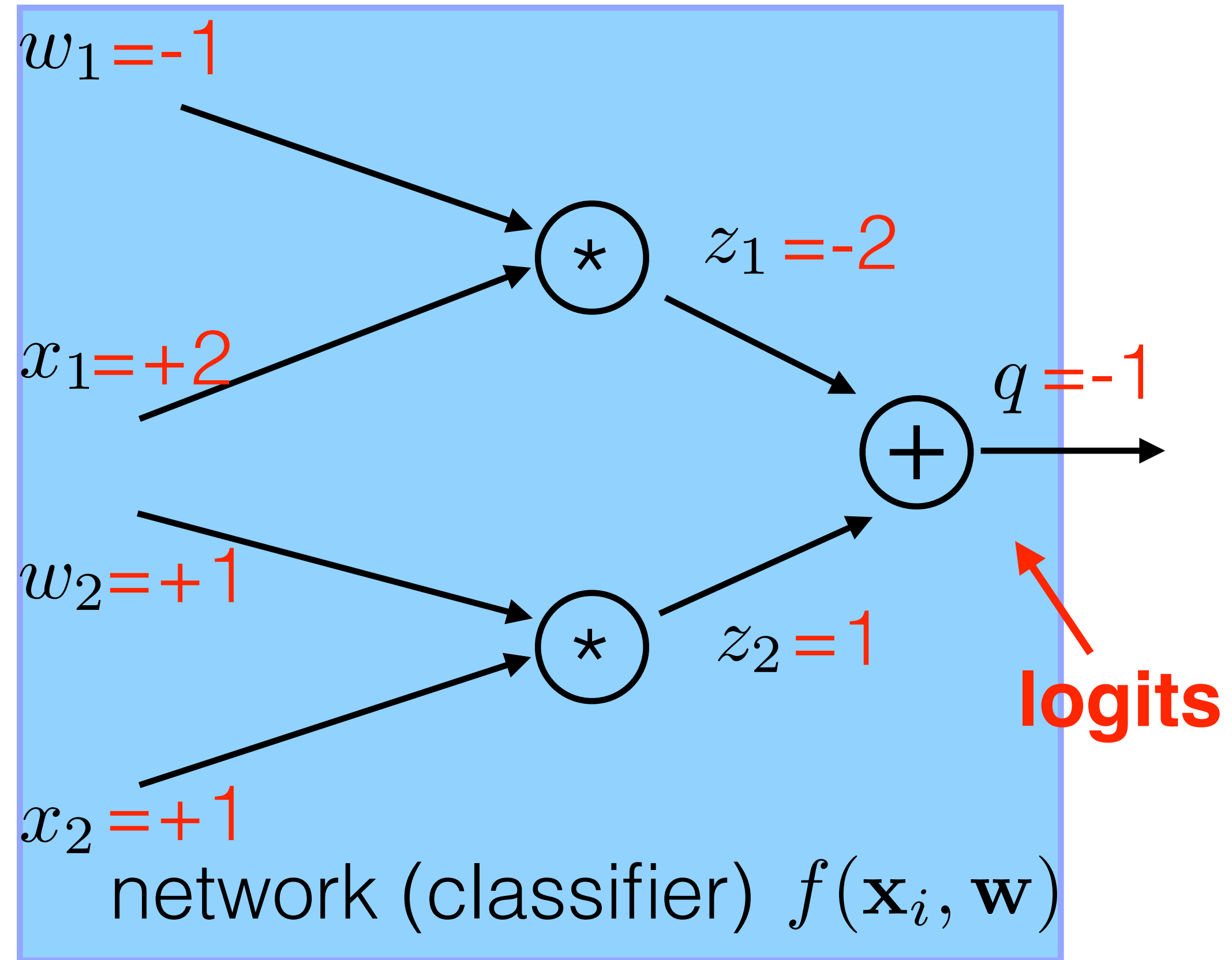


Output:

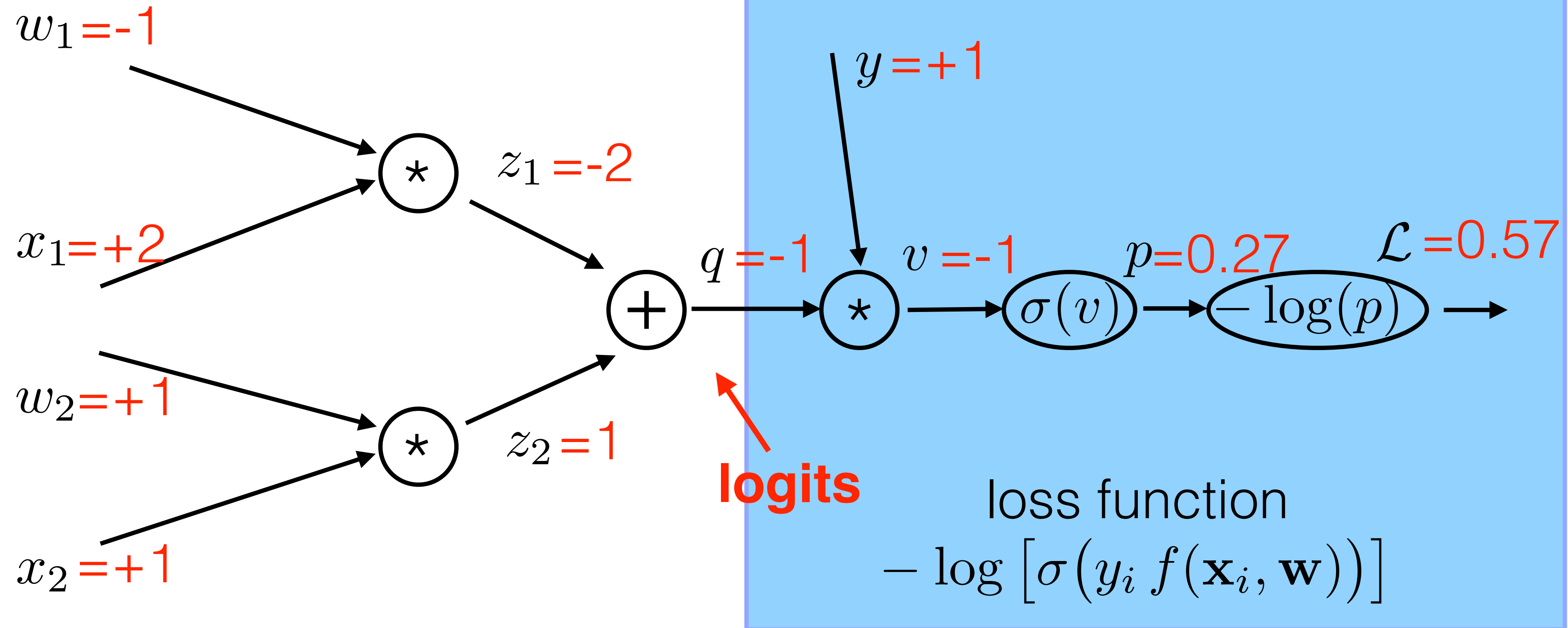


<http://jackterwilliger.com/biological-neural-networks-part-i-spiking-neurons/>

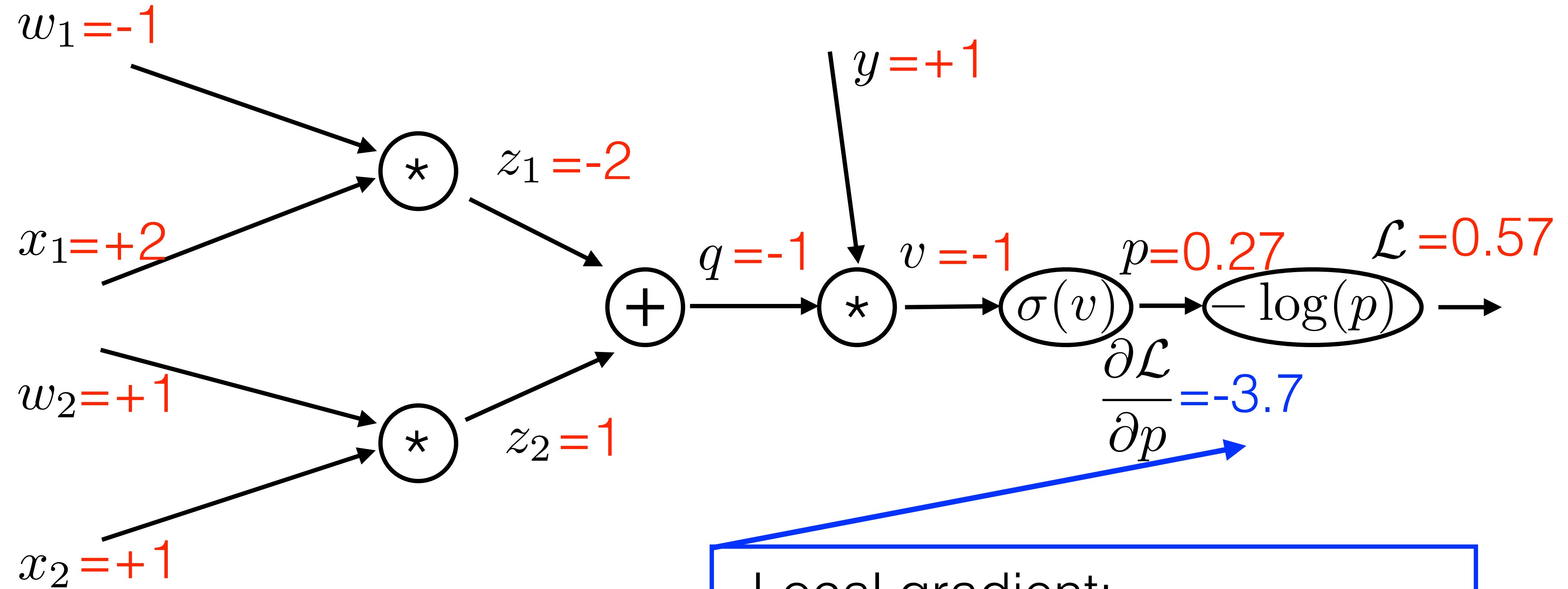
# Learning in computation graph



# Learning in computation graph



# Learning in computation graph

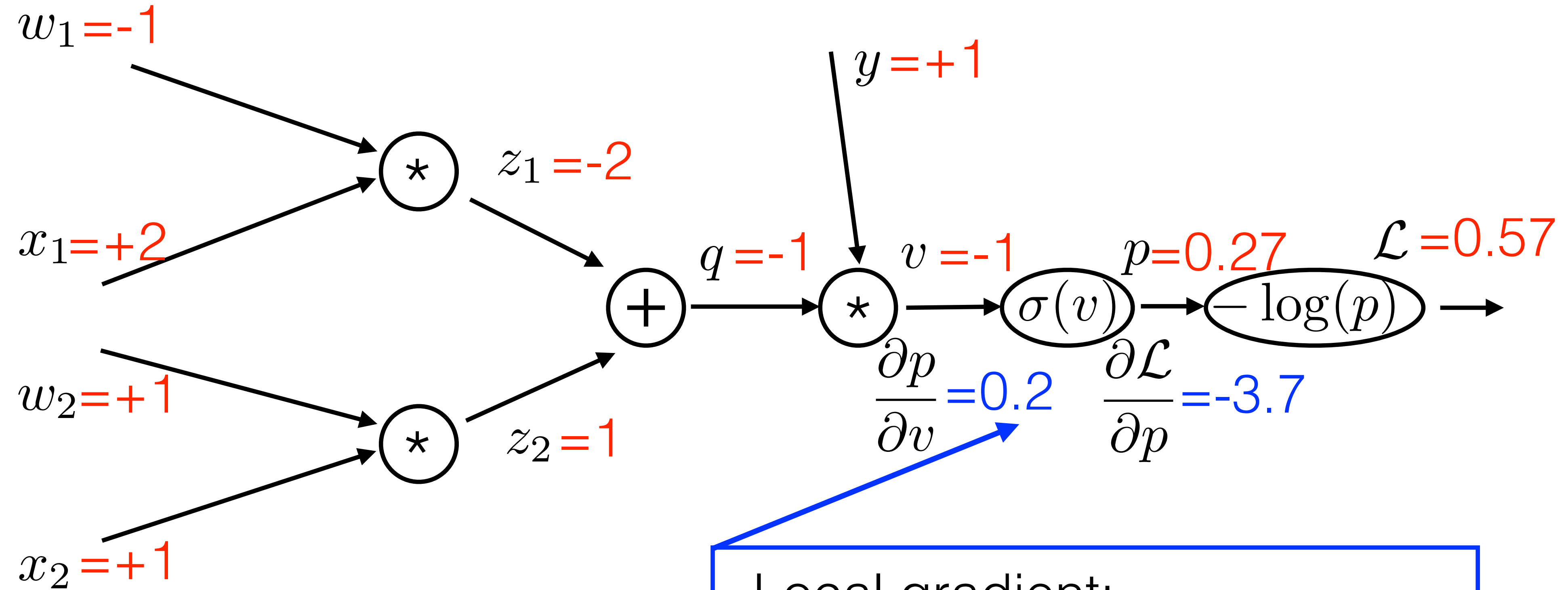


Local gradient:

$$\frac{\partial \mathcal{L}}{\partial p} = \frac{\partial(-\log(p))}{\partial p} = -\frac{1}{p}$$



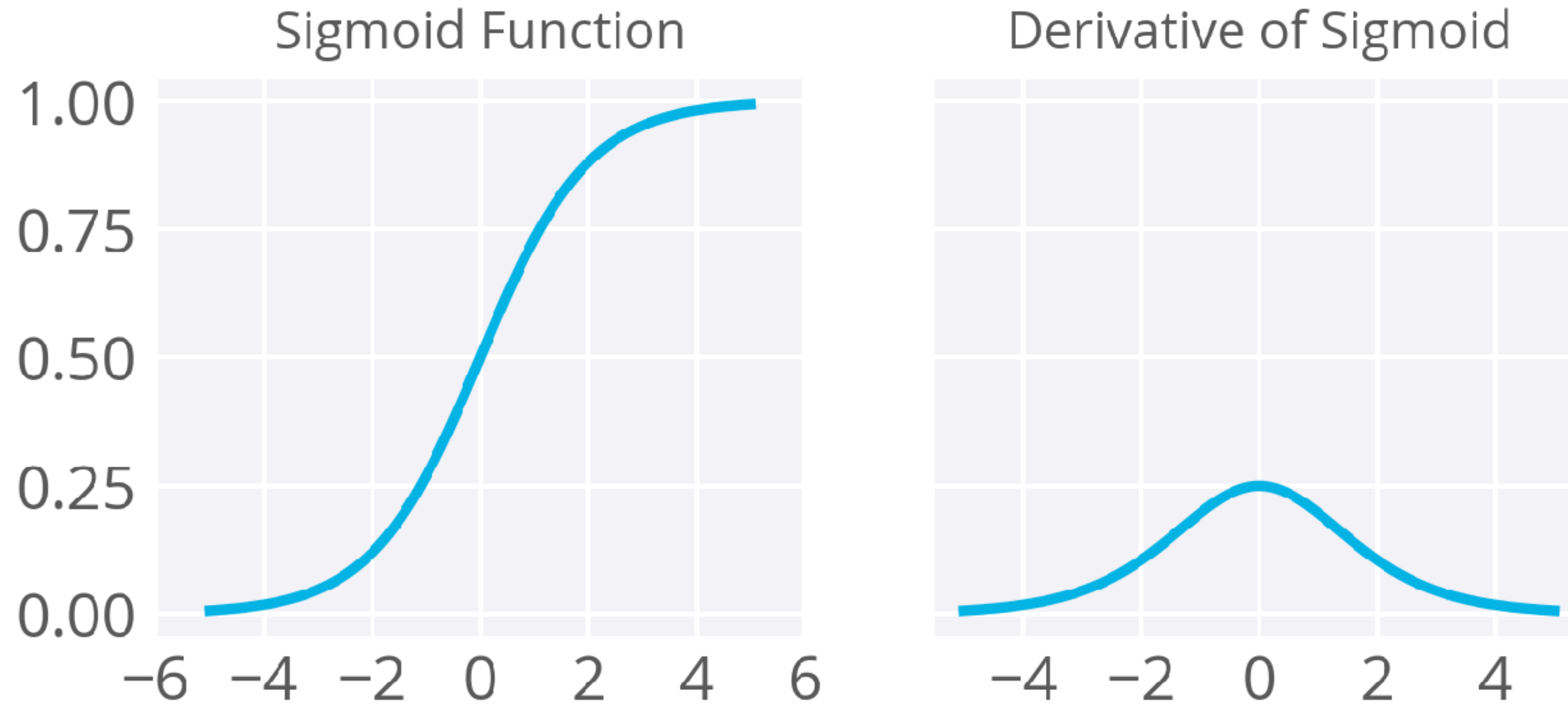
# Learning in computation graph



Local gradient:

$$\frac{\partial p}{\partial v} = \frac{\partial \sigma(v)}{\partial v} = \sigma(v)(1 - \sigma(v))$$

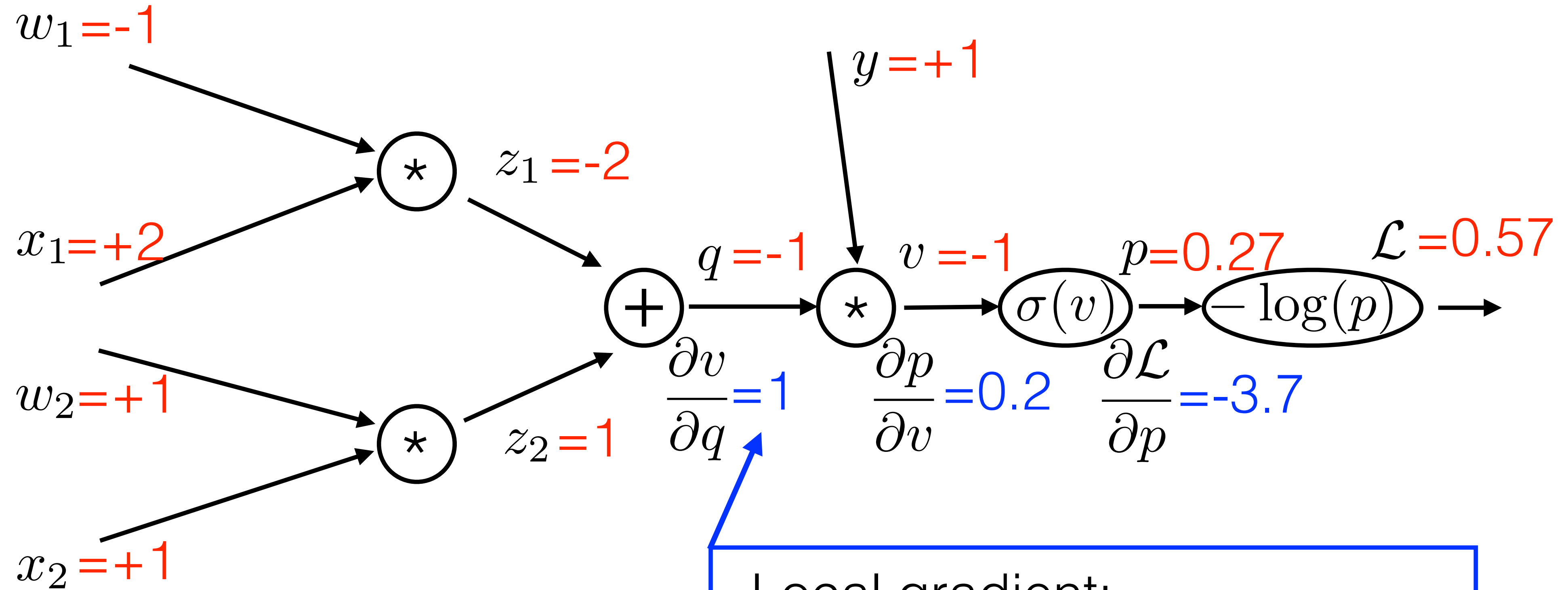
# Learning in computation graph



Local gradient:

$$\frac{\partial p}{\partial v} = \frac{\partial \sigma(v)}{\partial v} = \sigma(v)(1 - \sigma(v))$$

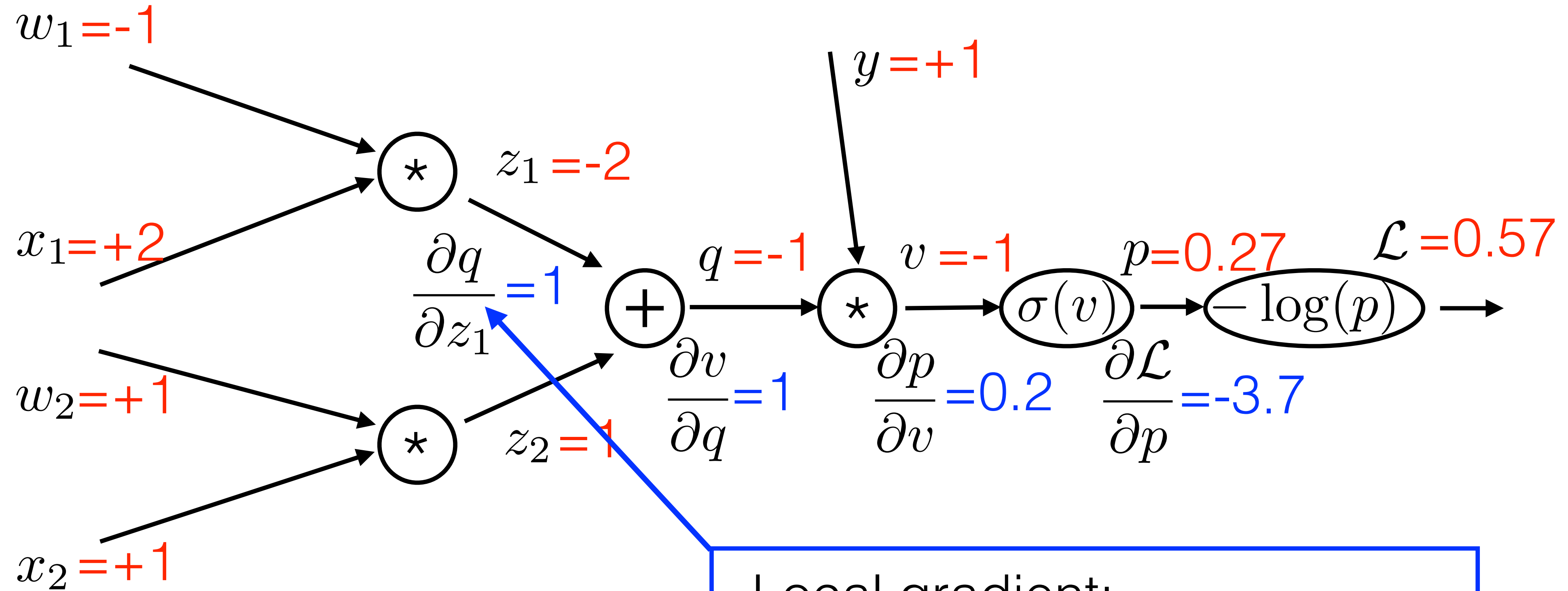
# Learning in computation graph



Local gradient:

$$\frac{\partial v}{\partial q} = \frac{\partial(yq)}{\partial q} = y$$

# Learning in computation graph

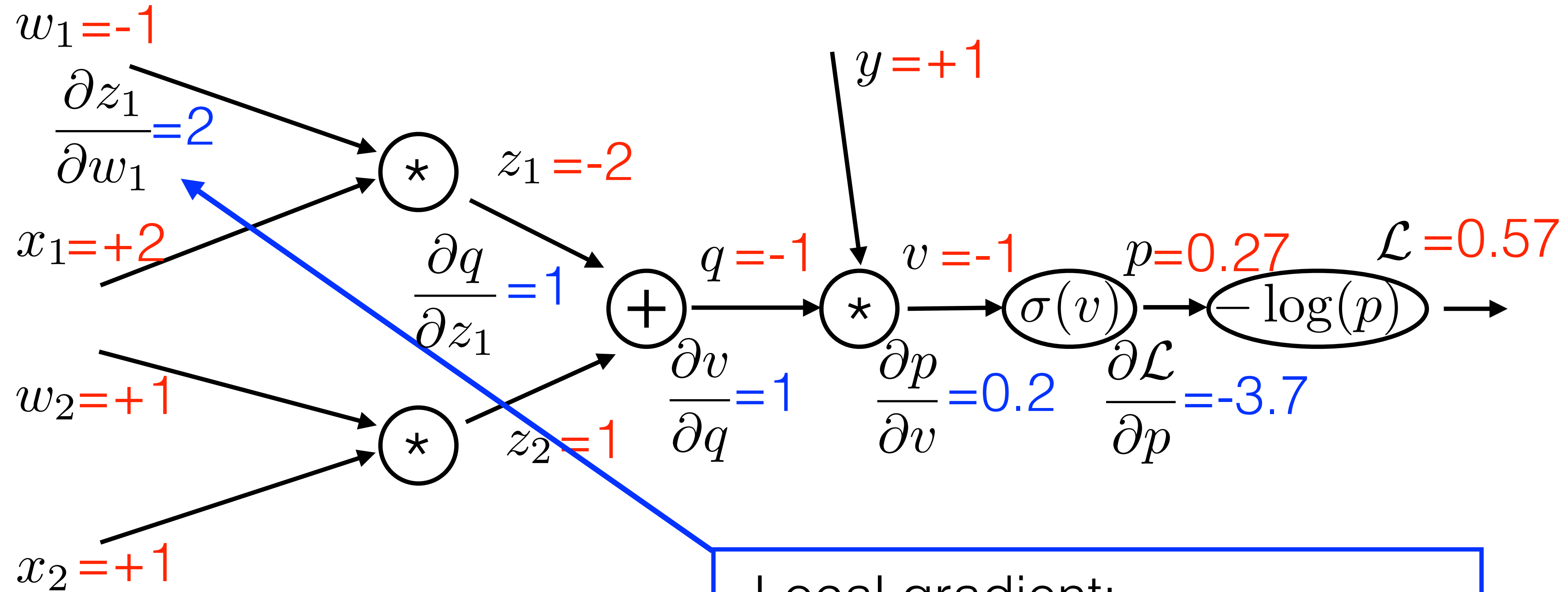


Local gradient:

$$\frac{\partial q}{\partial z_1} = \frac{\partial(z_1 + z_2)}{\partial z_1} = 1$$



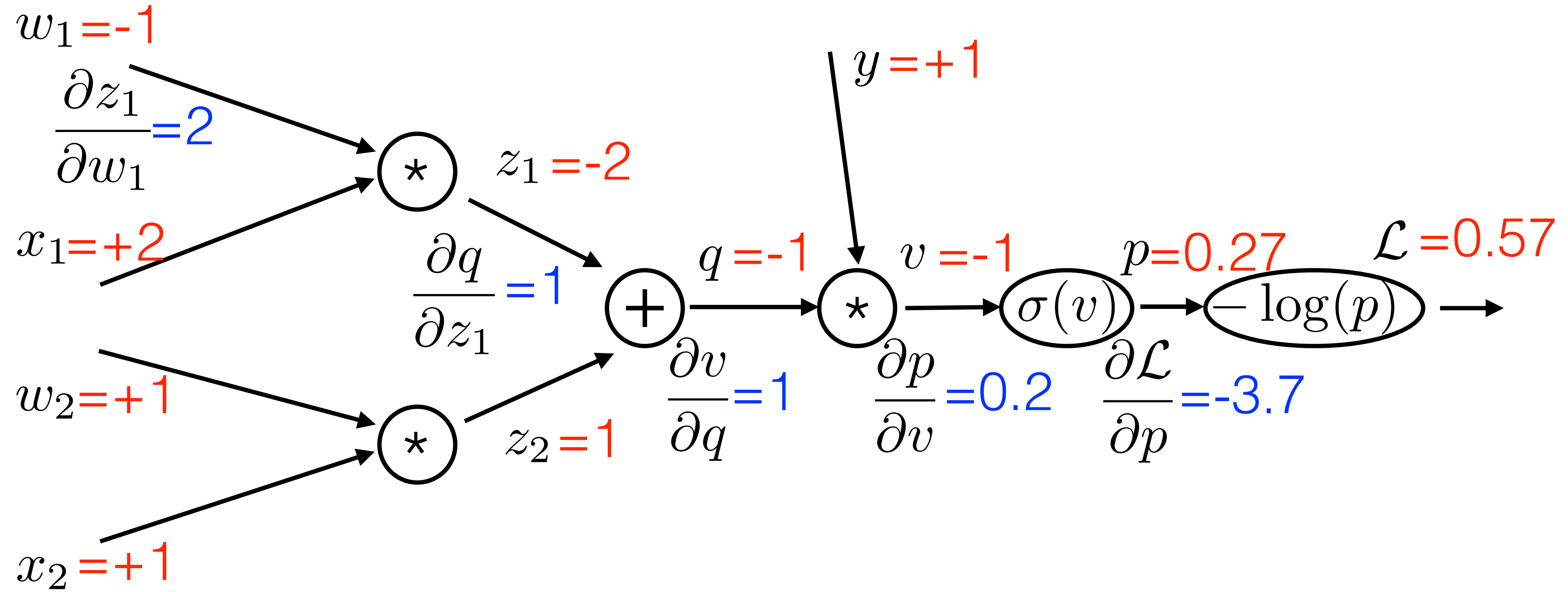
# Learning in computation graph



Local gradient:

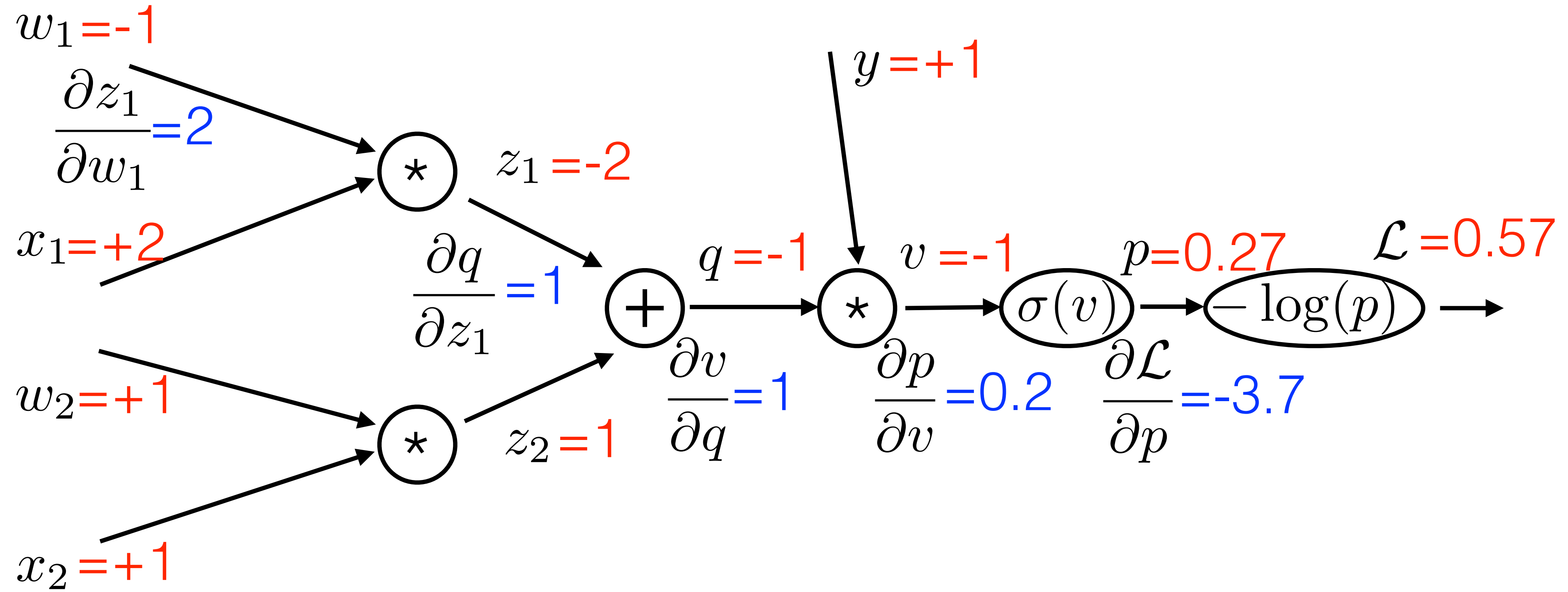
$$\frac{\partial z_1}{\partial w_1} = \frac{\partial (w_1 x_1)}{\partial w_1} = x_1$$

# Learning in computation graph



$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial p} \frac{\partial p}{\partial v} \frac{\partial v}{\partial q} \frac{\partial q}{\partial z_1} \frac{\partial z_1}{\partial w_1} = -3.7 * 0.2 * 1 * 1 * 2 = -1.48$$

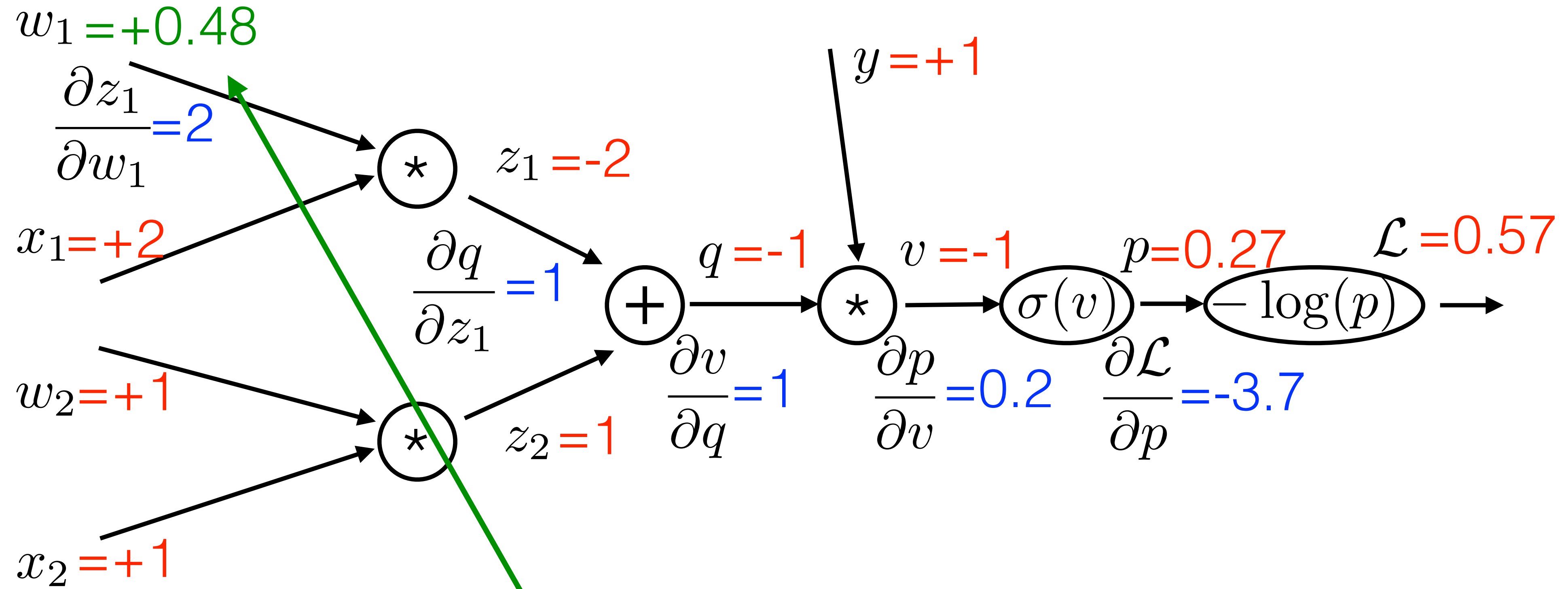
# Learning in computation graph



$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial p} \frac{\partial p}{\partial v} \frac{\partial v}{\partial q} \frac{\partial q}{\partial z_1} \frac{\partial z_1}{\partial w_1} = -3.7 * 0.2 * 1 * 1 * 2 = -1.48$$

$$w_1 = w_1 - \alpha \frac{\partial \mathcal{L}}{\partial w_1} = +0.48$$

# Learning in computation graph

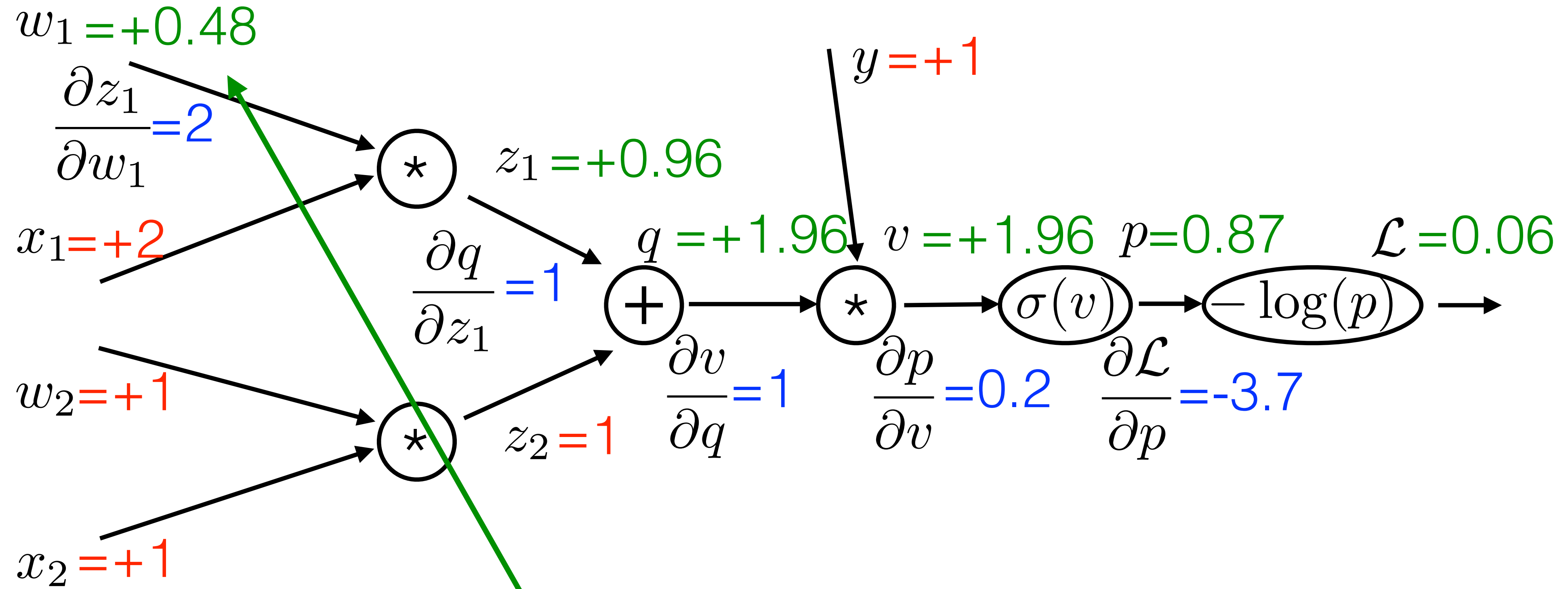


$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial p} \frac{\partial p}{\partial v} \frac{\partial v}{\partial q} \frac{\partial q}{\partial z_1} \frac{\partial z_1}{\partial w_1} = -3.7 * 0.2 * 1 * 1 * 2 = -1.48$$

$$w_1 = w_1 - \alpha \frac{\partial \mathcal{L}}{\partial w_1} = +0.48$$



# Learning in computation graph

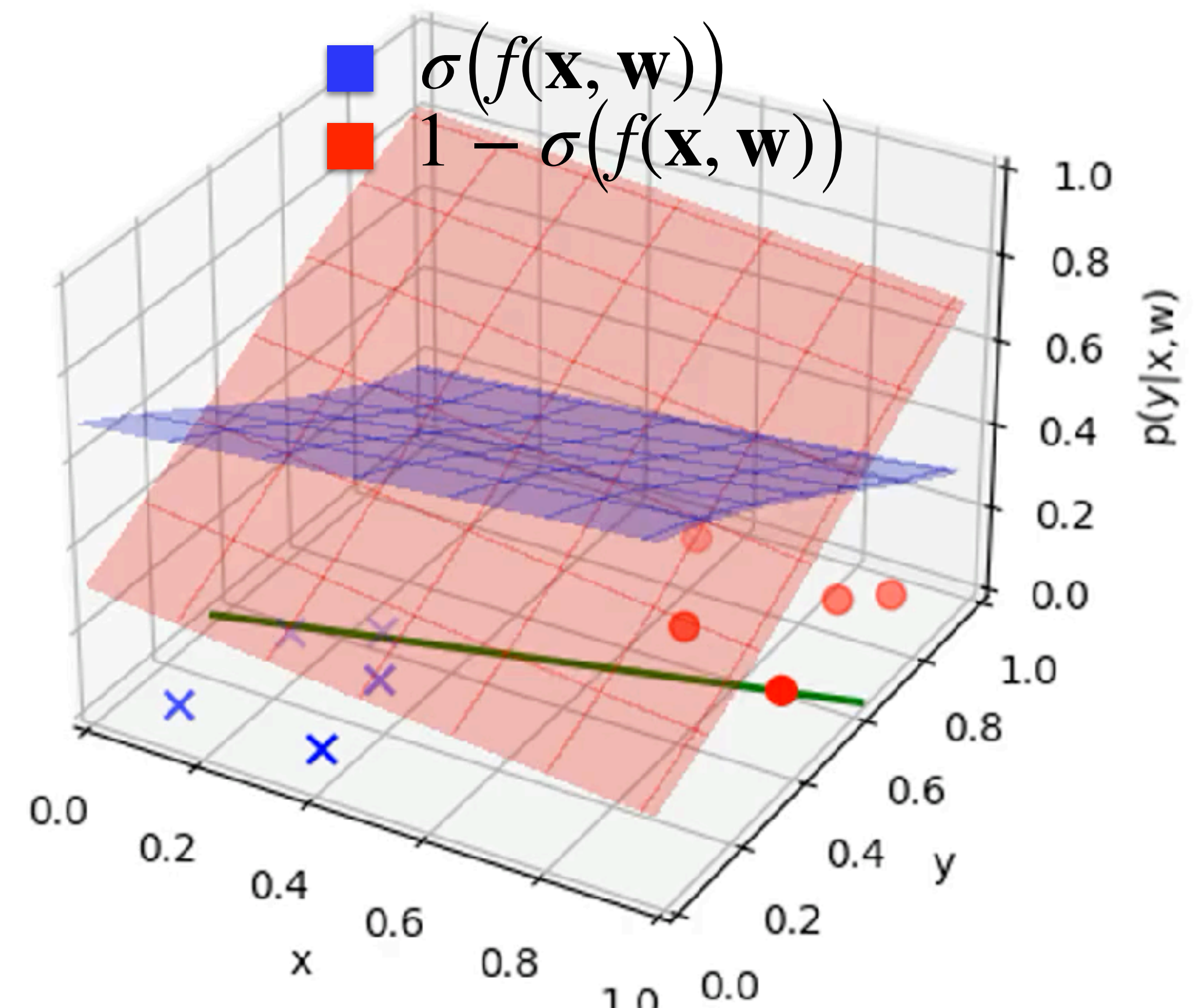
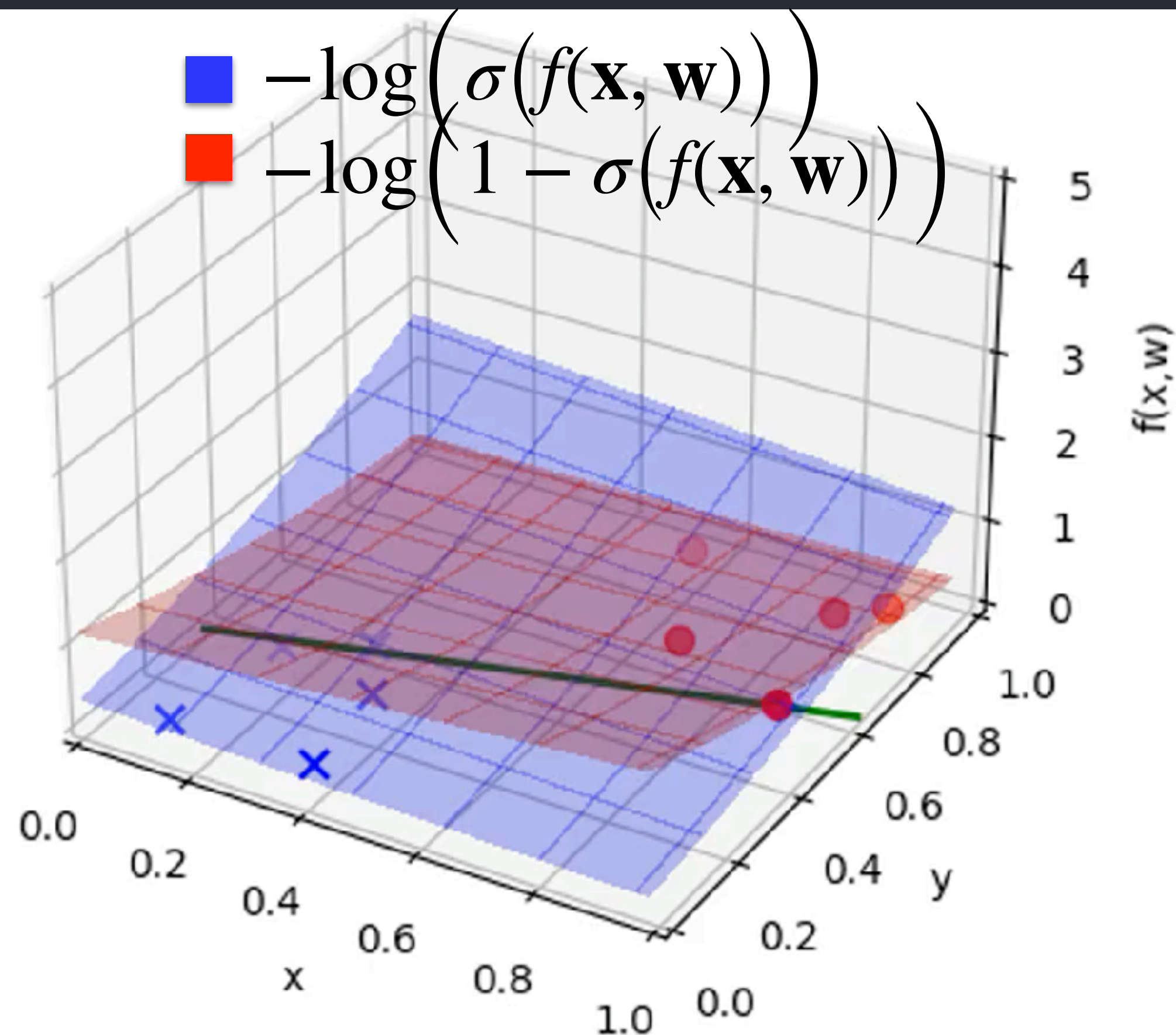


$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial p} \frac{\partial p}{\partial v} \frac{\partial v}{\partial q} \frac{\partial q}{\partial z_1} \frac{\partial z_1}{\partial w_1} = -3.7 * 0.2 * 1 * 1 * 2 = -1.48$$

$$w_1 = w_1 - \alpha \frac{\partial \mathcal{L}}{\partial w_1} = +0.48$$

# Numpy implementation

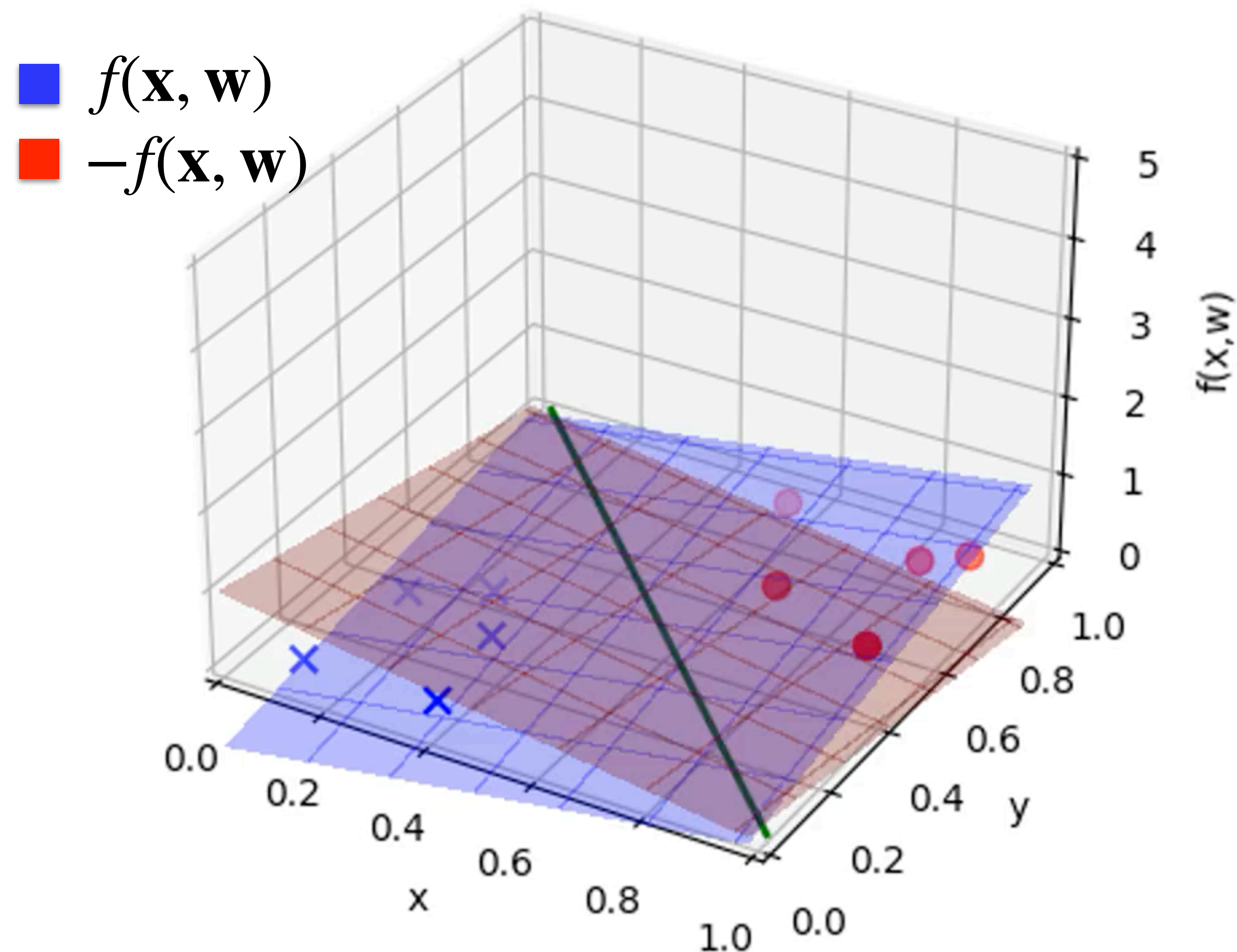
```
for i in range(30):  
    u = w[0] * x[:, 0] + w[1] * x[:, 1] + w[2]  
    p = sigmoid(u)  
    loss = (-np.log(p)*y + -np.log(1-p)*(1-y)).sum()  
    grad = -1/p * sigmoid(u) * (1-sigmoid(u)) * ...  
    w = w - 0.1 * grad
```



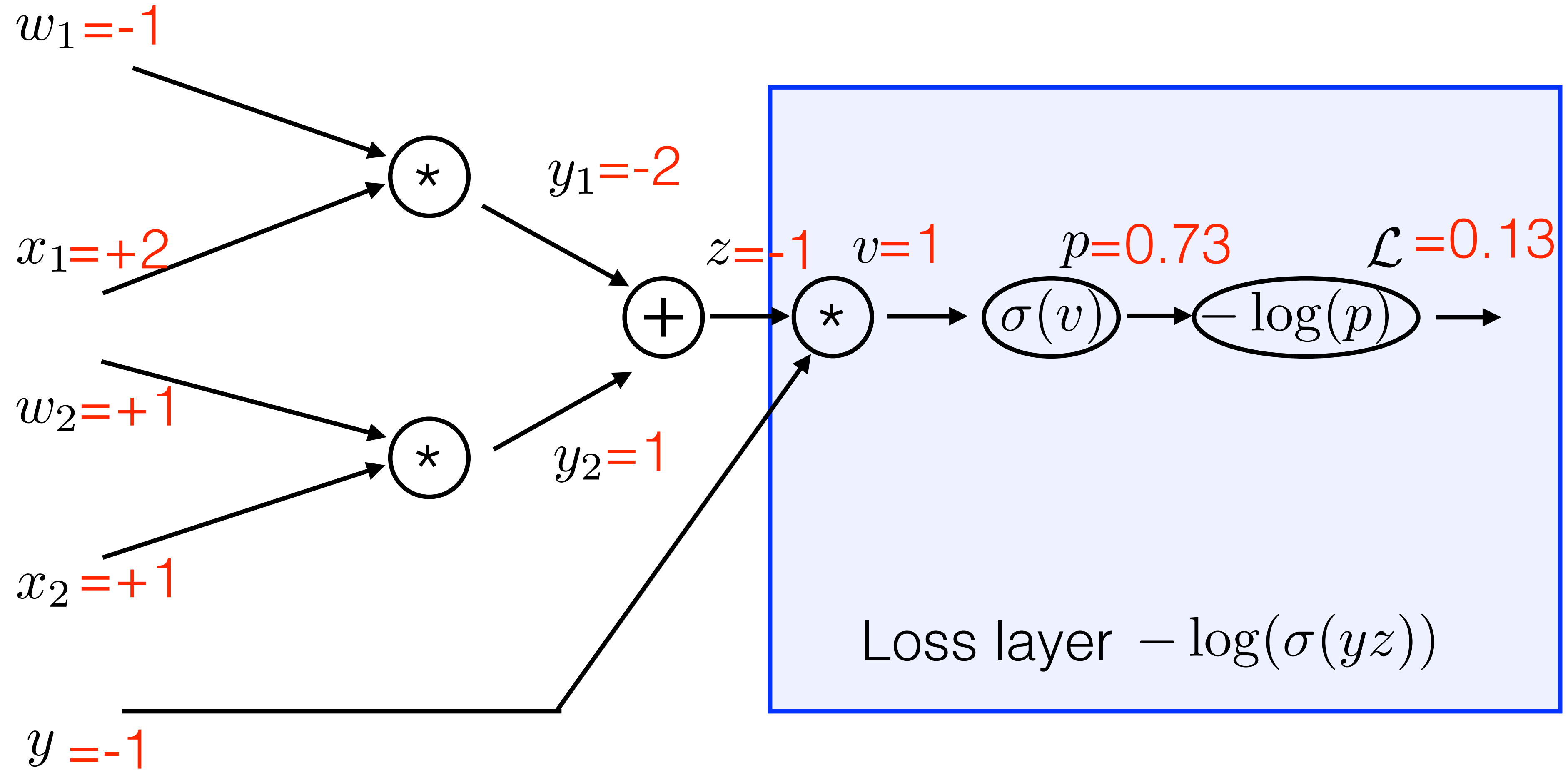


# Numpy implementation

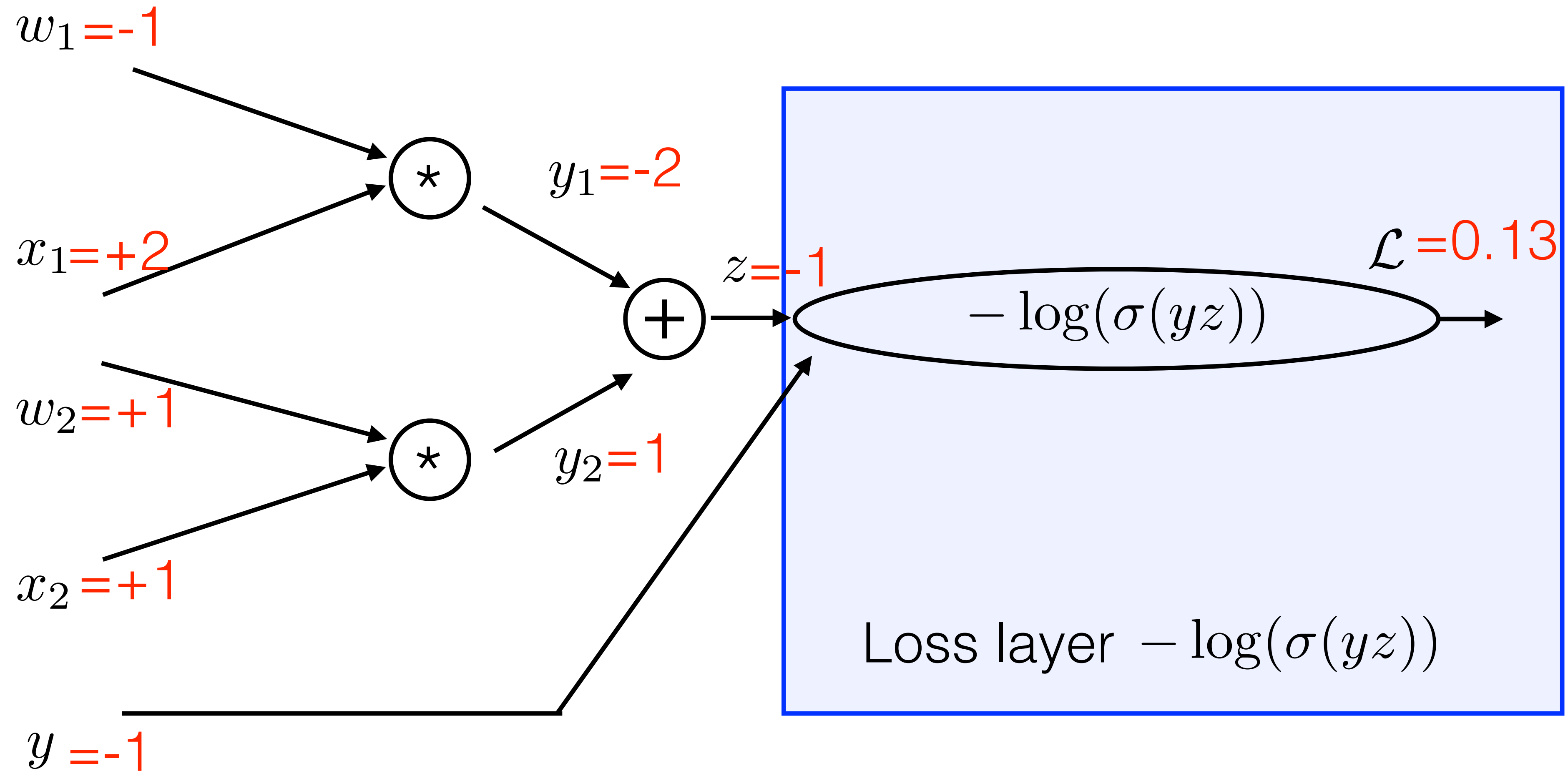
```
for i in range(30):  
    u = w[0] * x[:, 0] + w[1] * x[:, 1] + w[2]  
    loss = u*y + (1-u)*(1-y).sum()  
    grad = ...  
    w = w + 0.1 * grad
```



# Computational graph of the learning



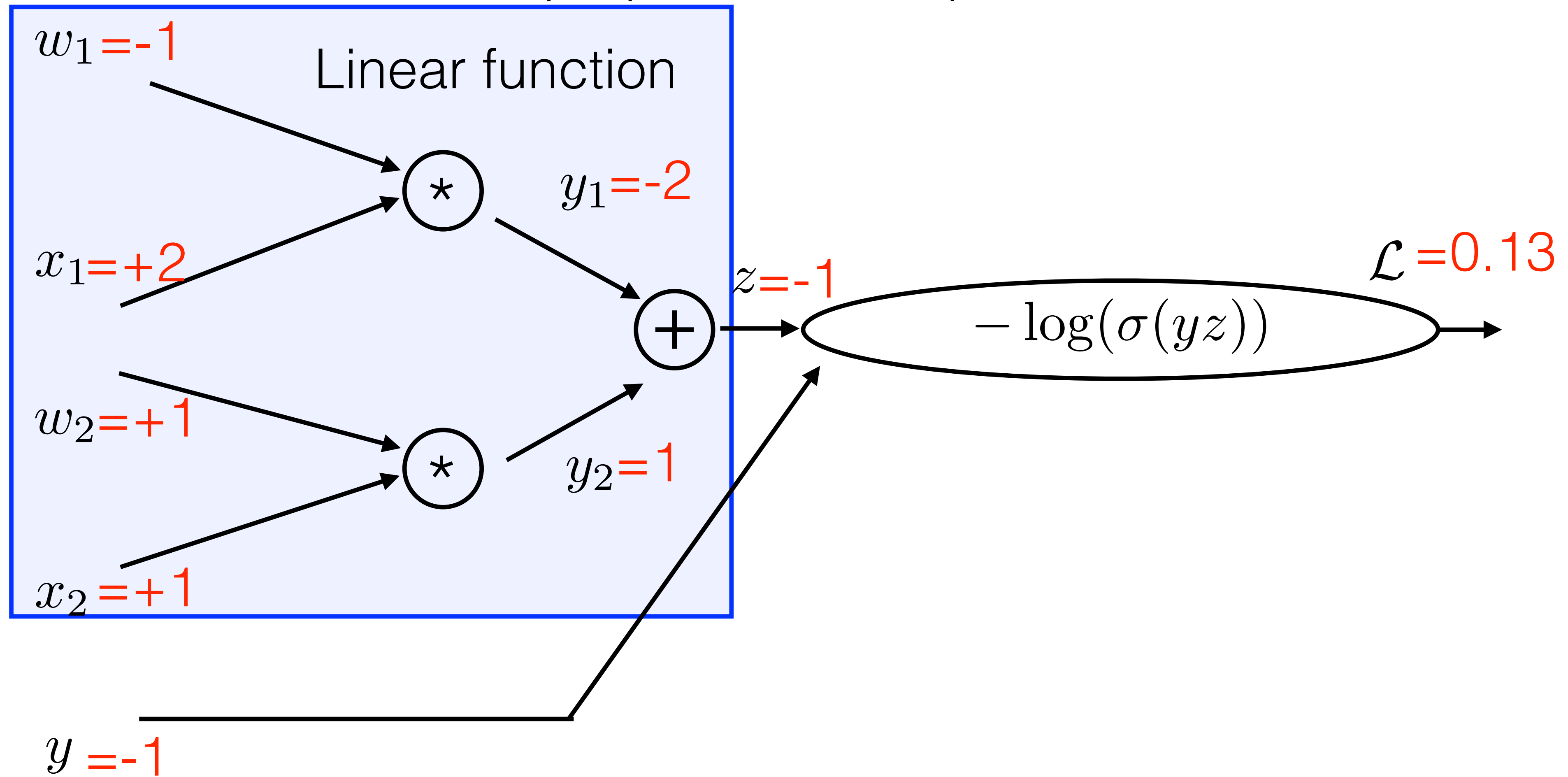
# Backprop in vector representation



This is the logistic loss!

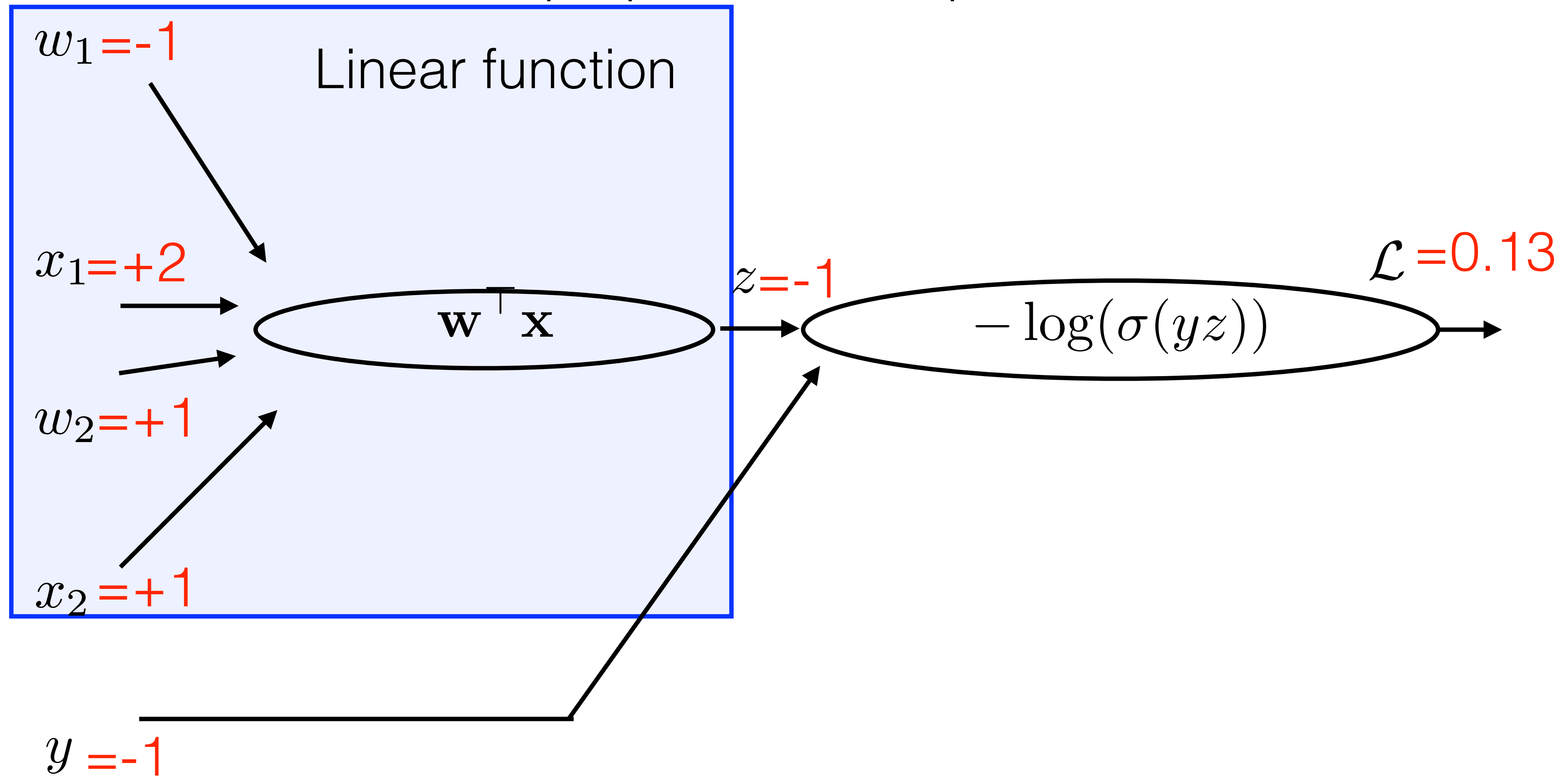
$$\mathcal{L}(y, z) = -\log(\sigma(yz)) = \log(1 + \exp(-yz))$$

# Backprop in vector representation

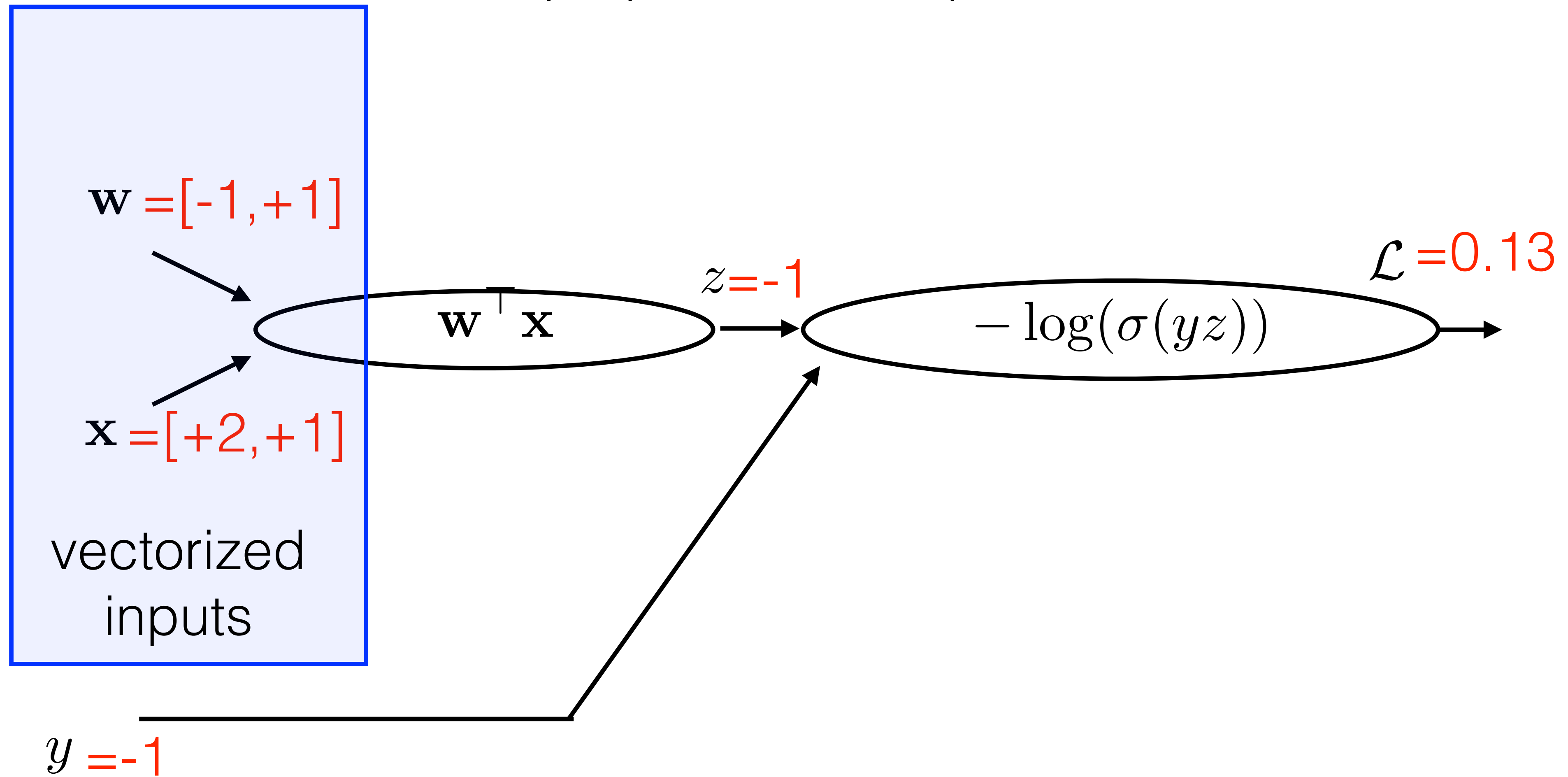




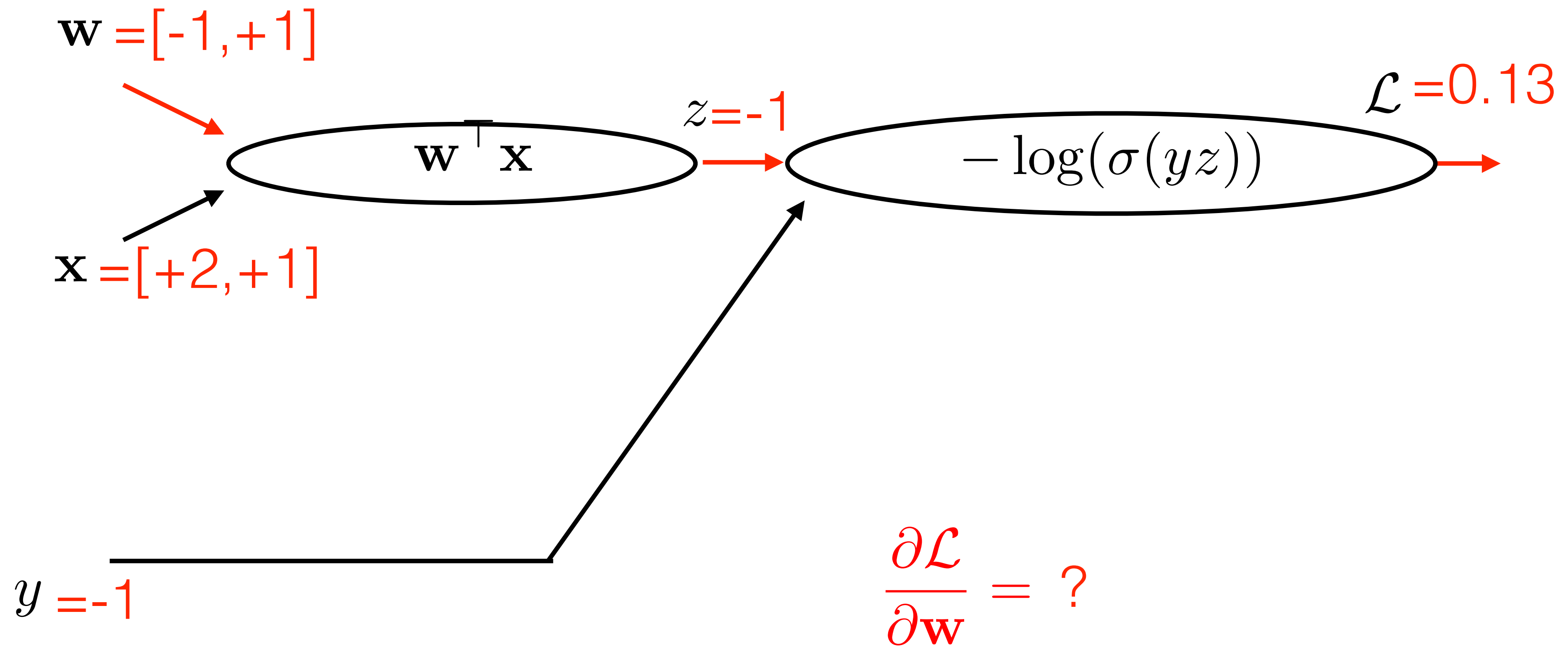
# Backprop in vector representation



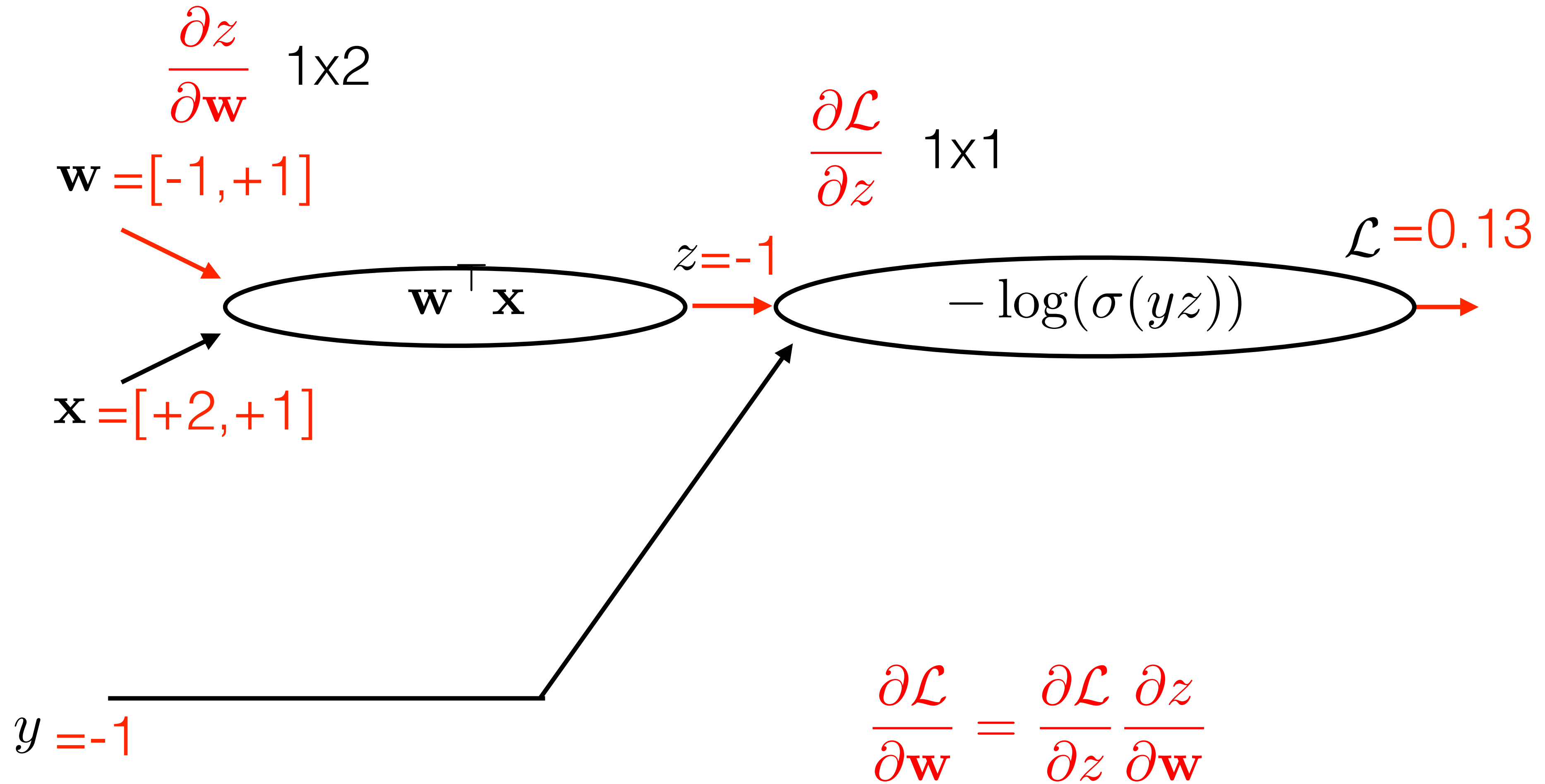
# Backprop in vector representation



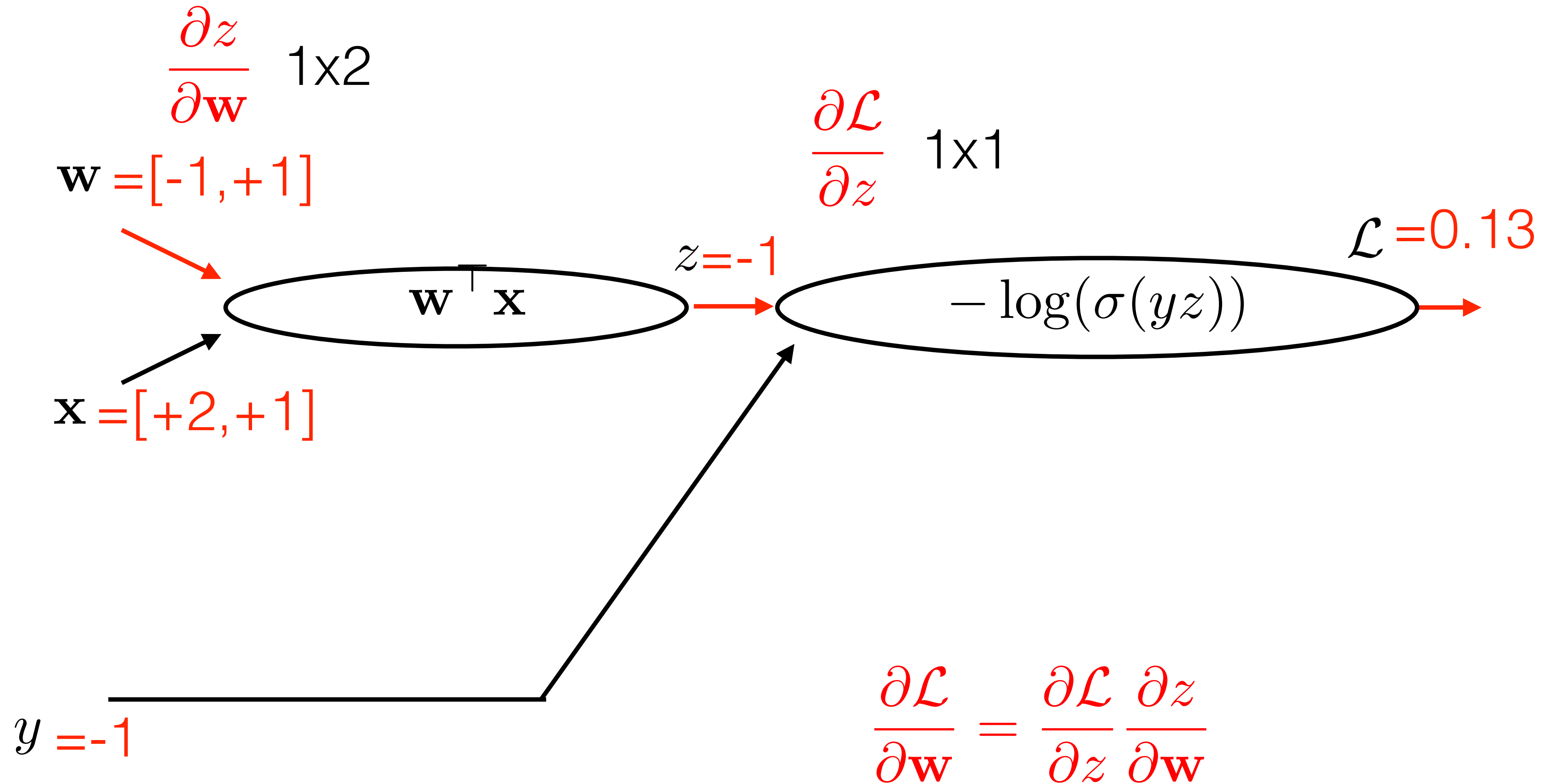
# Backprop in vector representation



# Backprop in vector representation

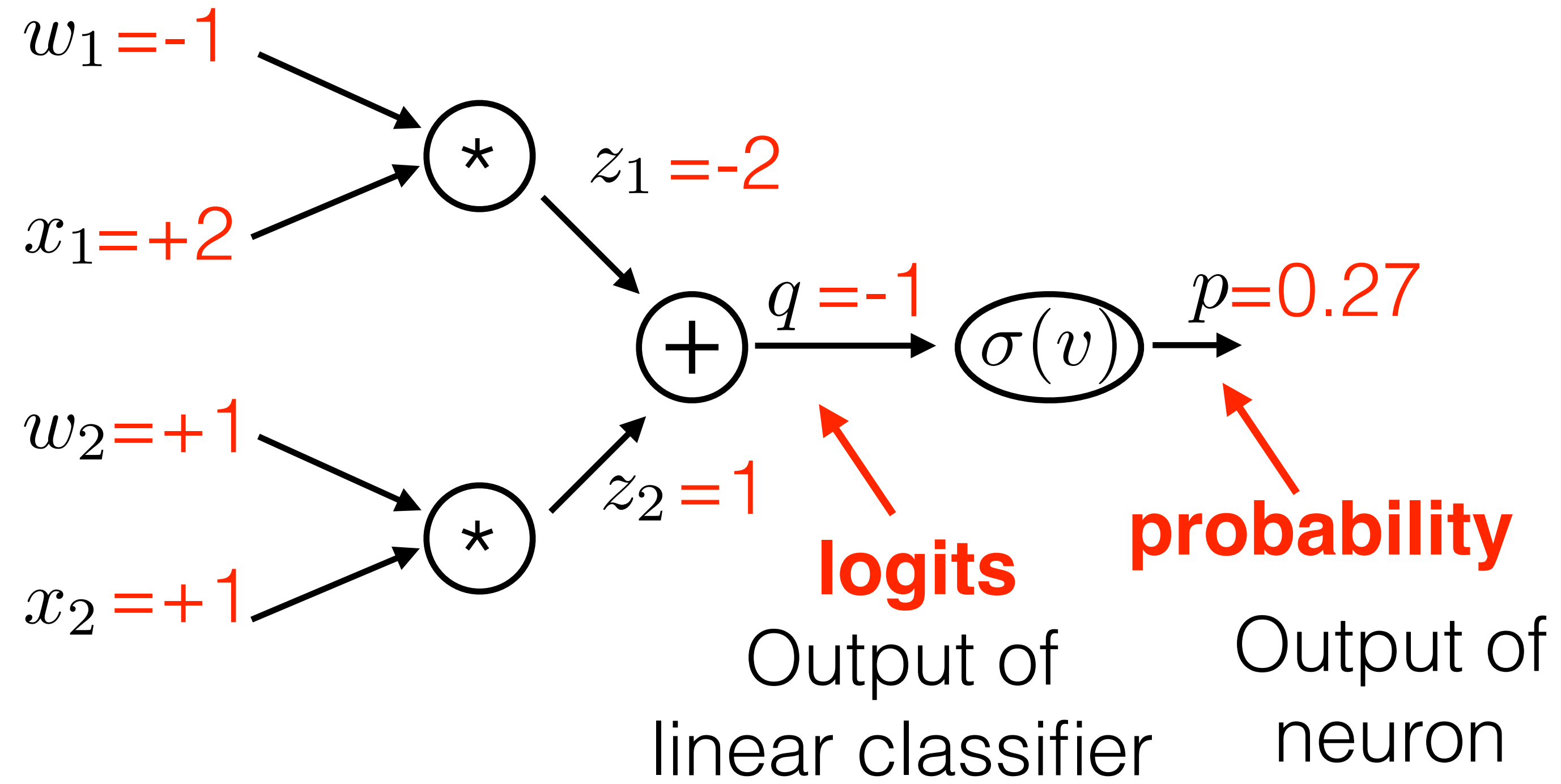


# Backprop in vector representation



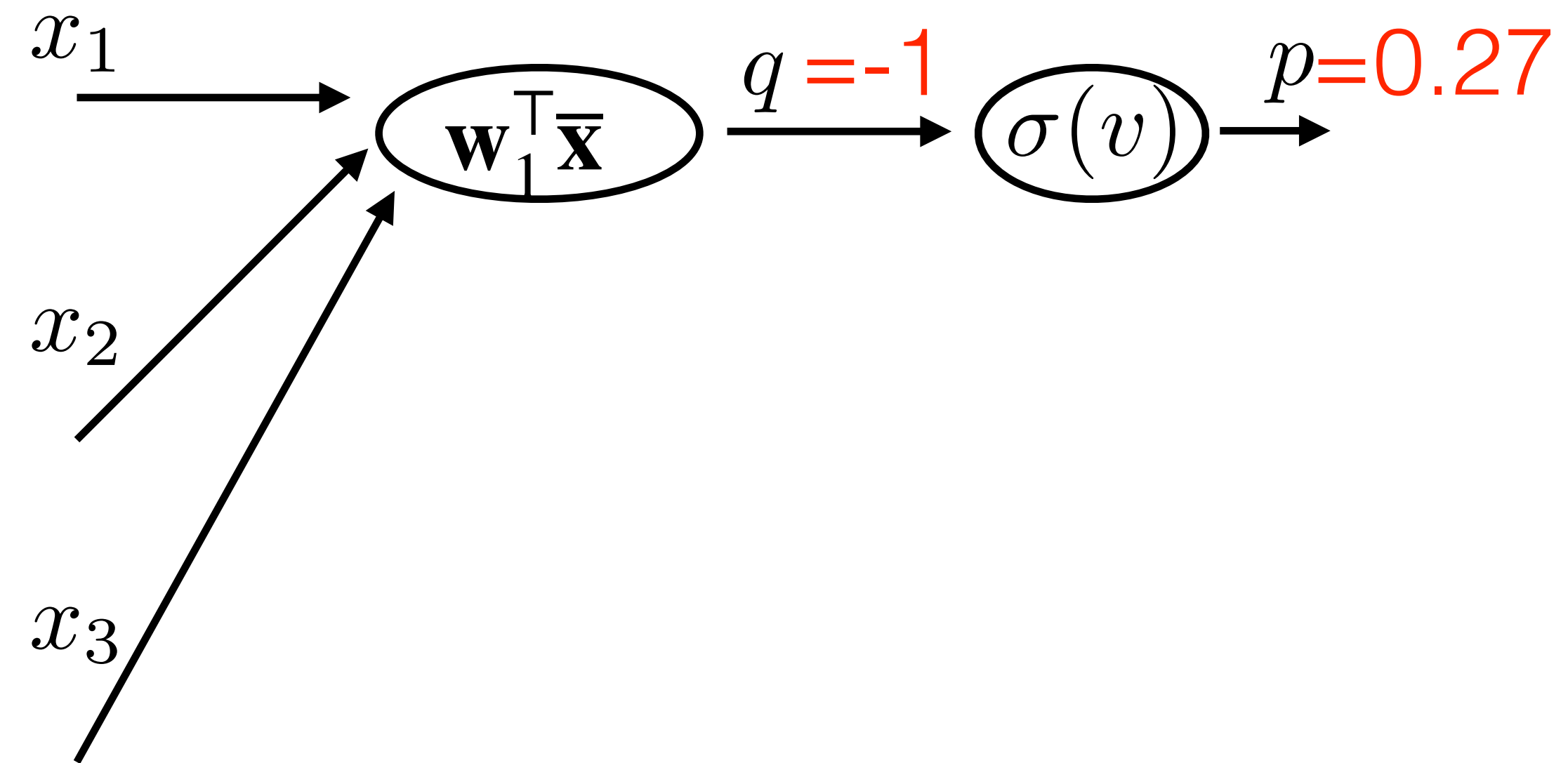
Learning from multiple training samples means summing up the partial derivative over all samples

# Neuron

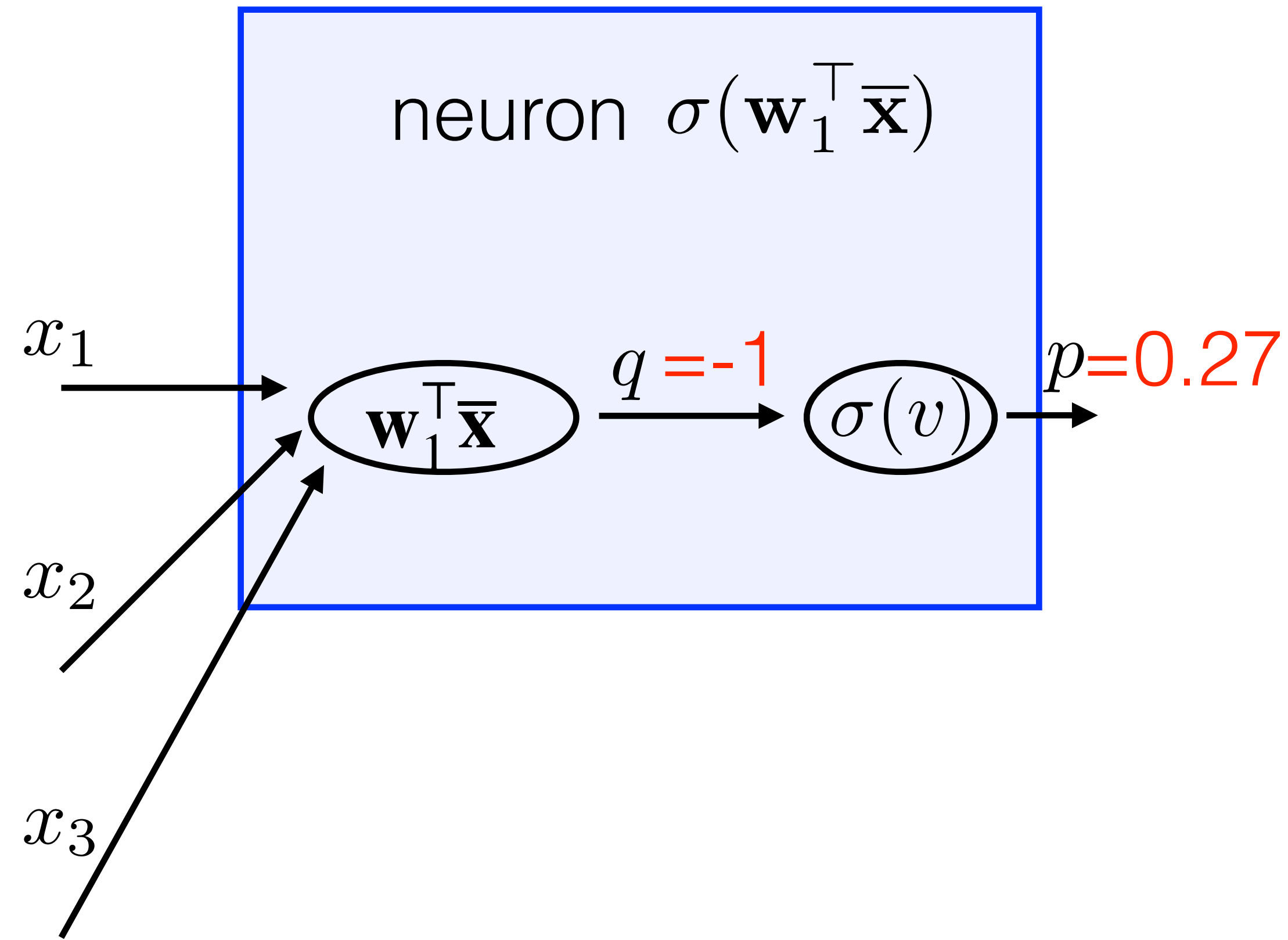




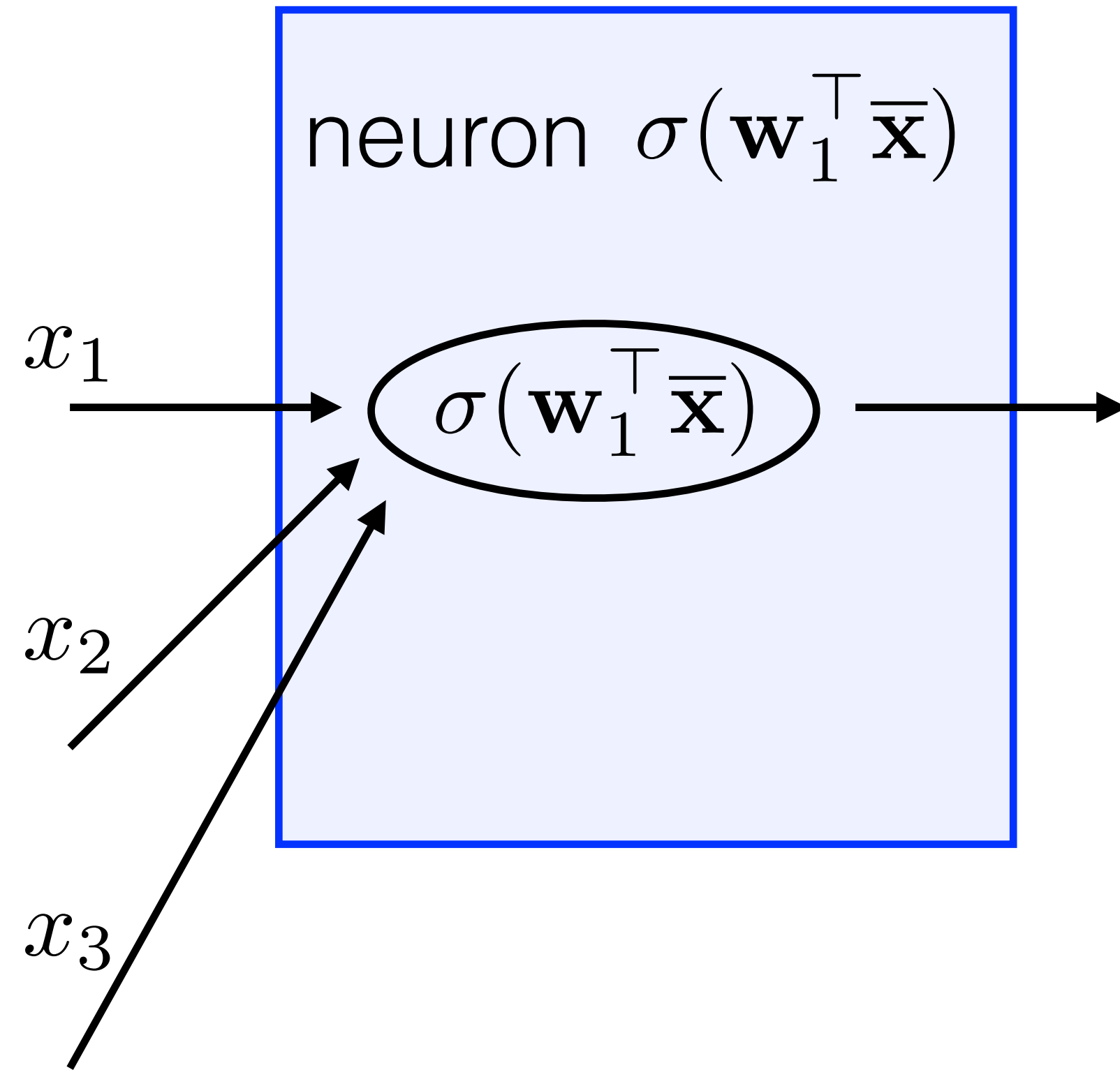
# Neuron



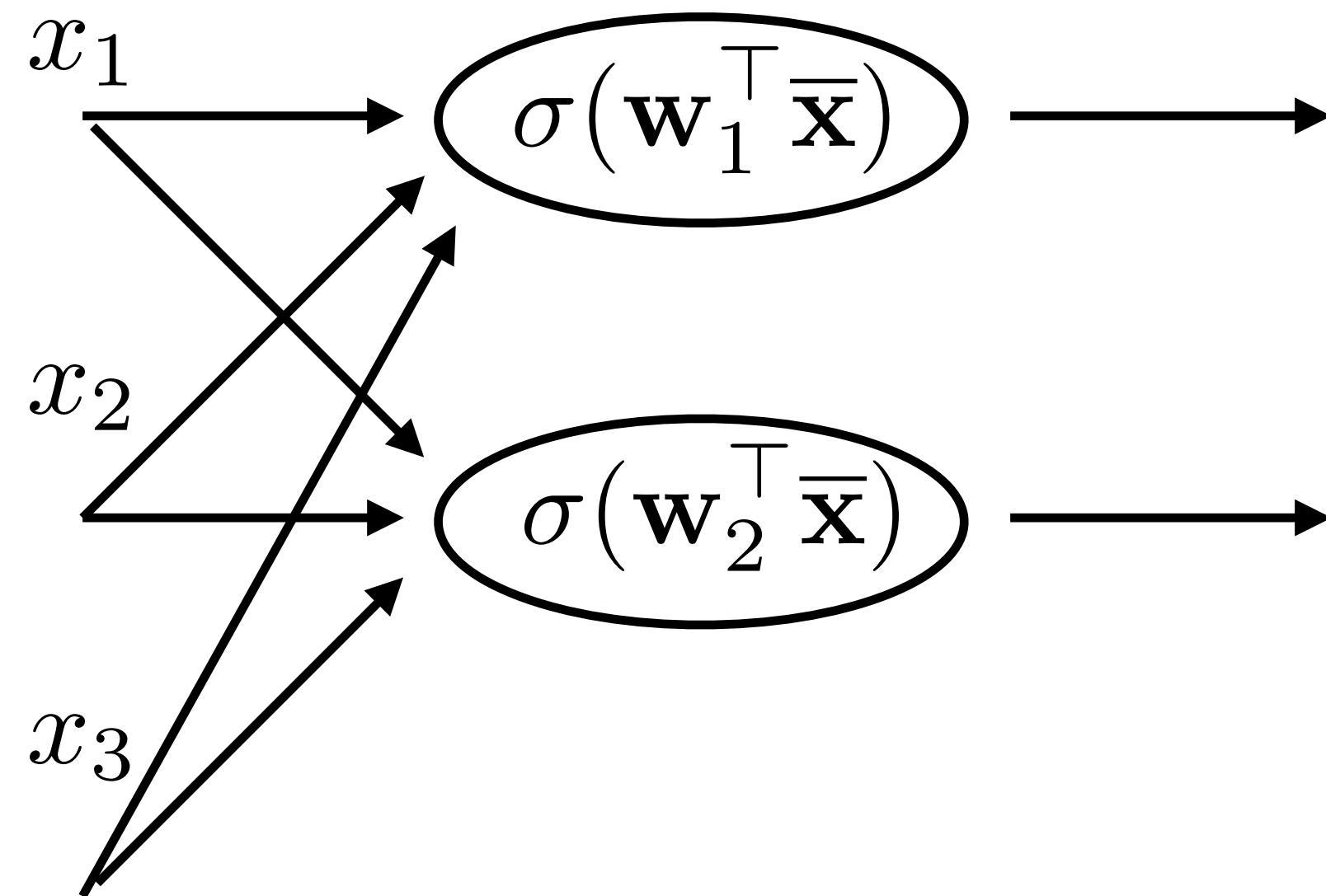
# Neuron



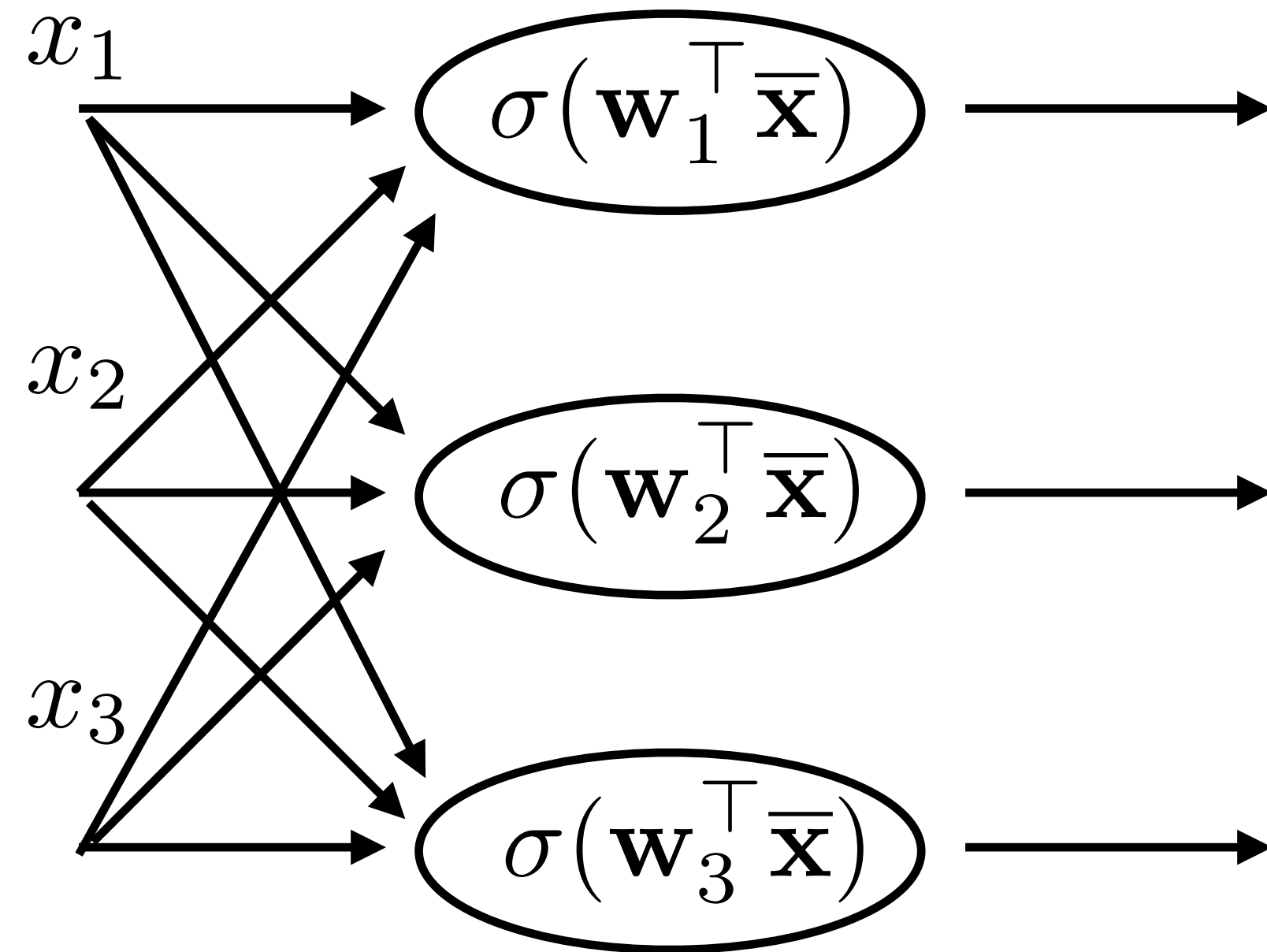
# Fully-connected neural network



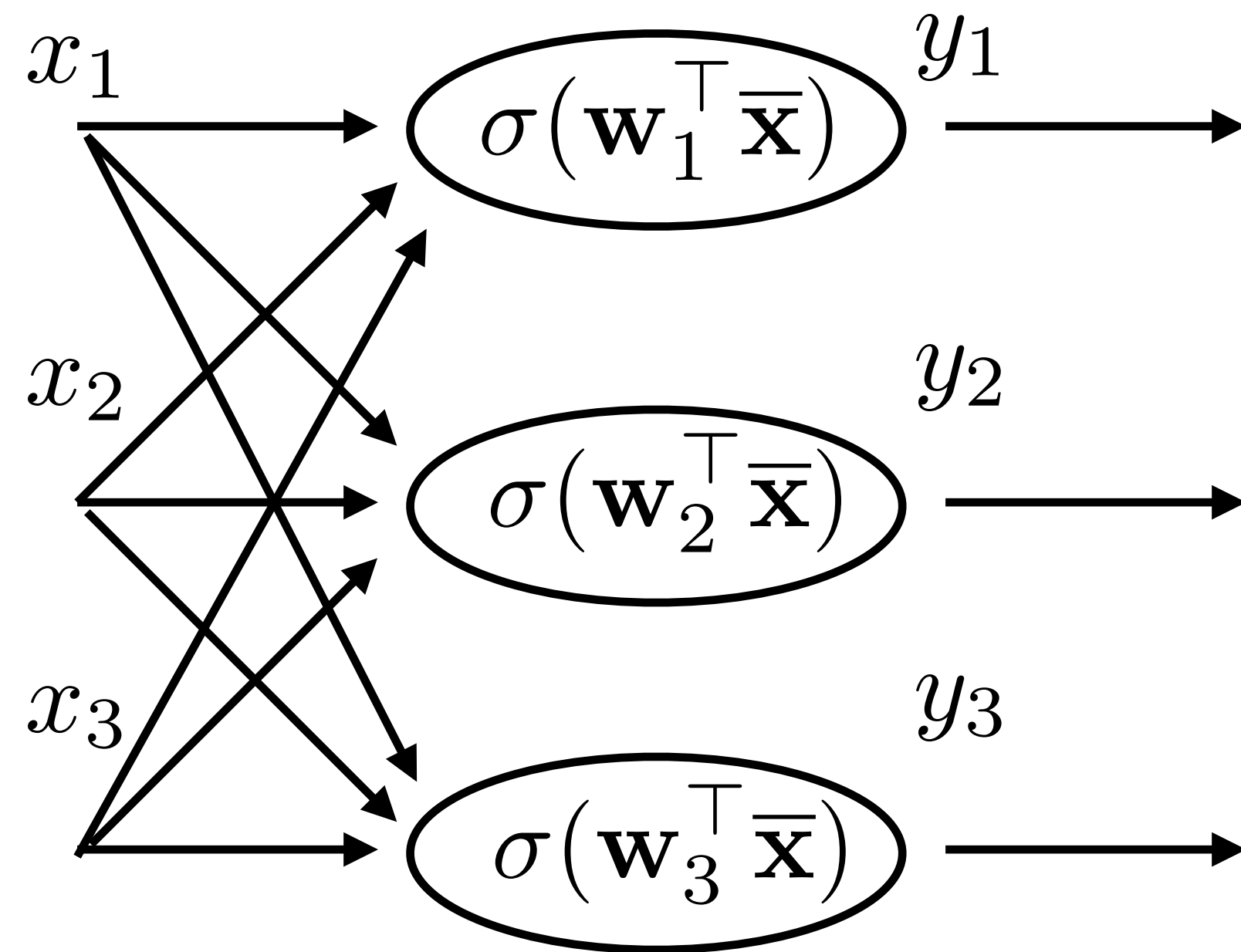
# Fully-connected neural network



# Fully-connected neural network

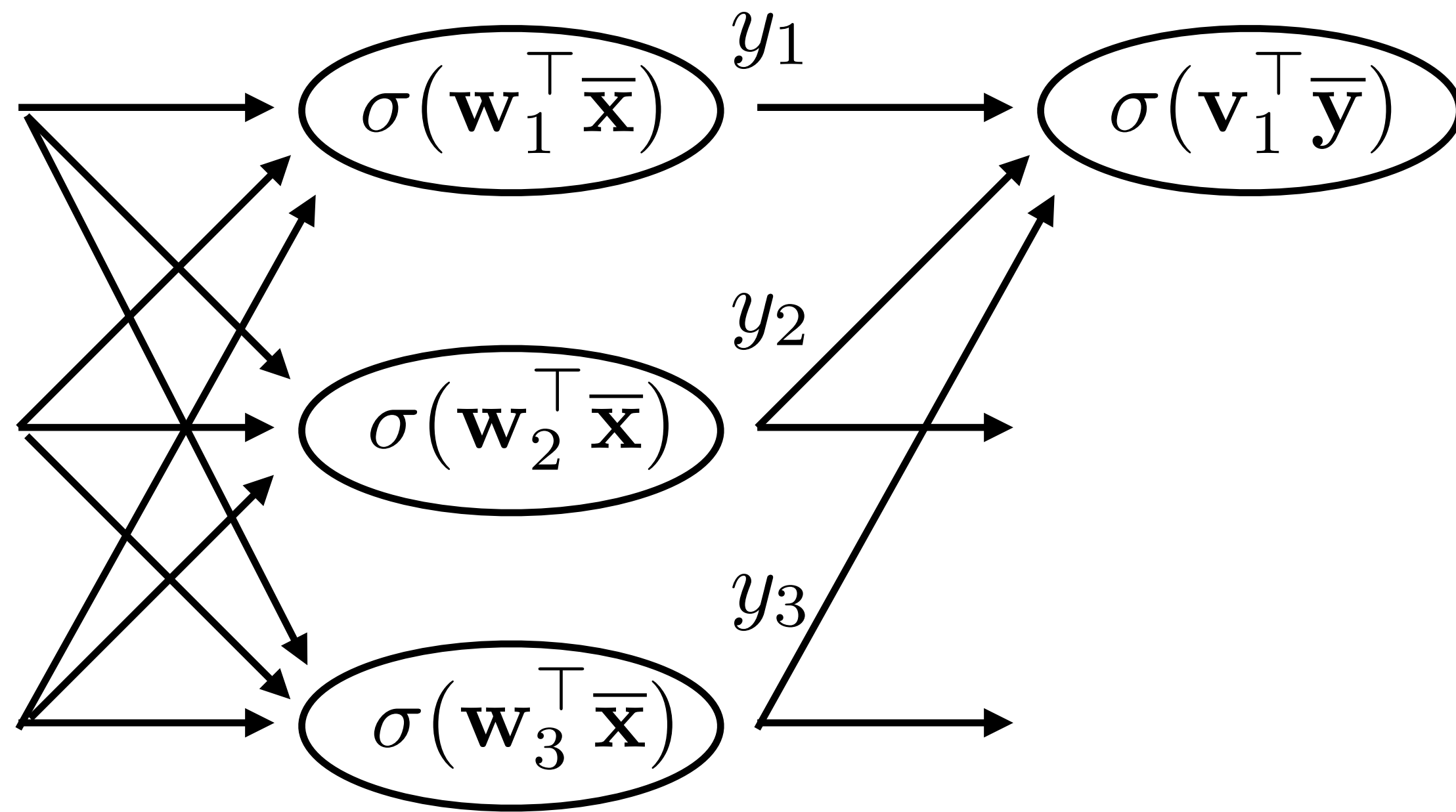


# Fully-connected neural network

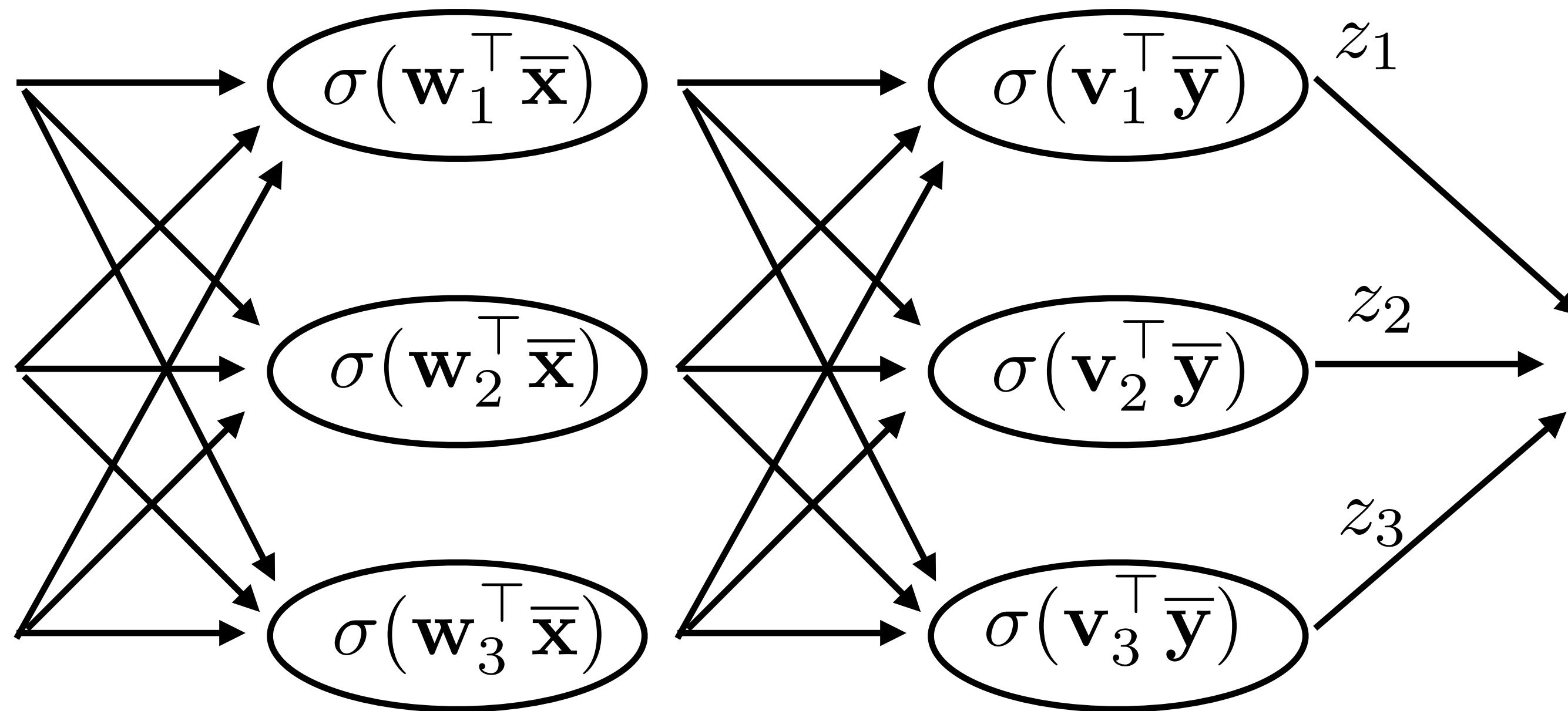




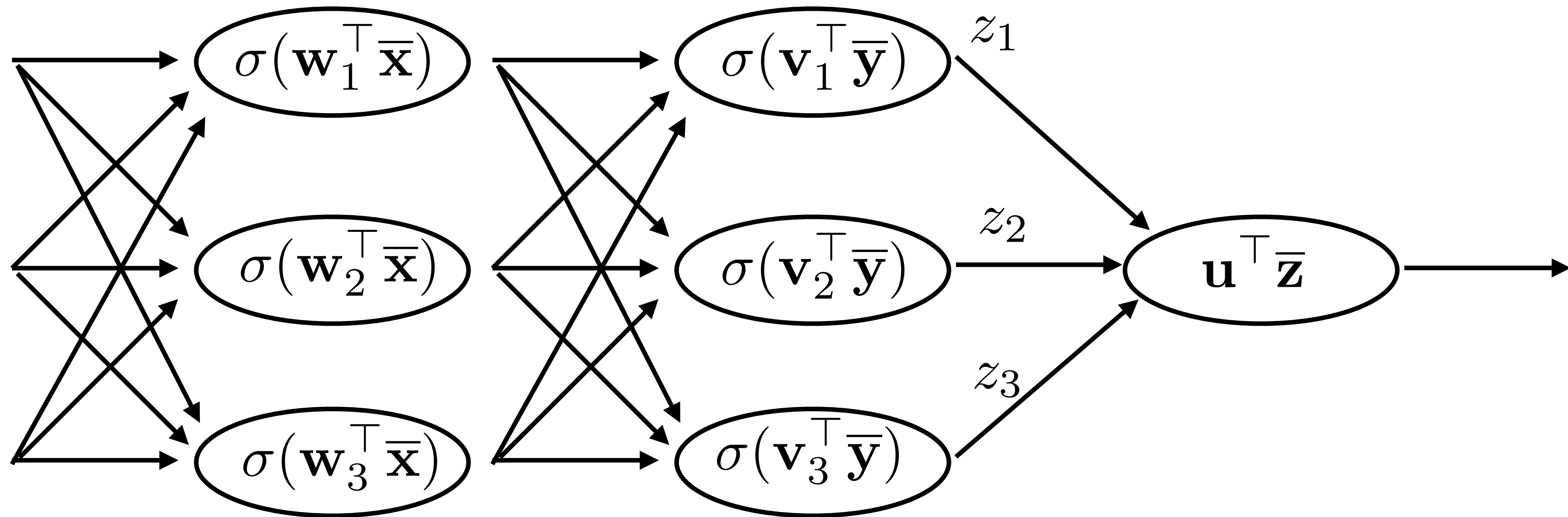
# Fully-connected neural network



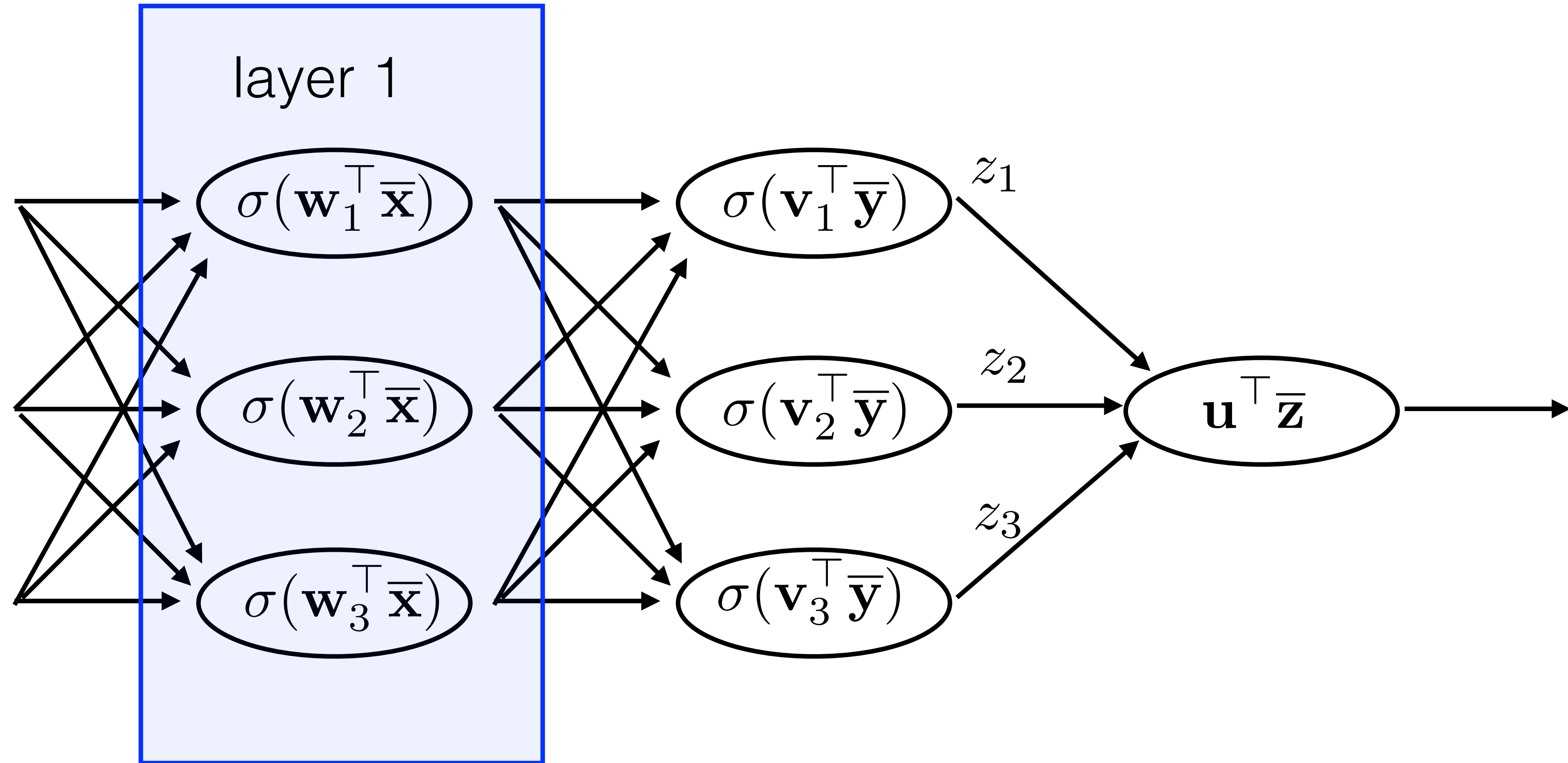
# Fully-connected neural network



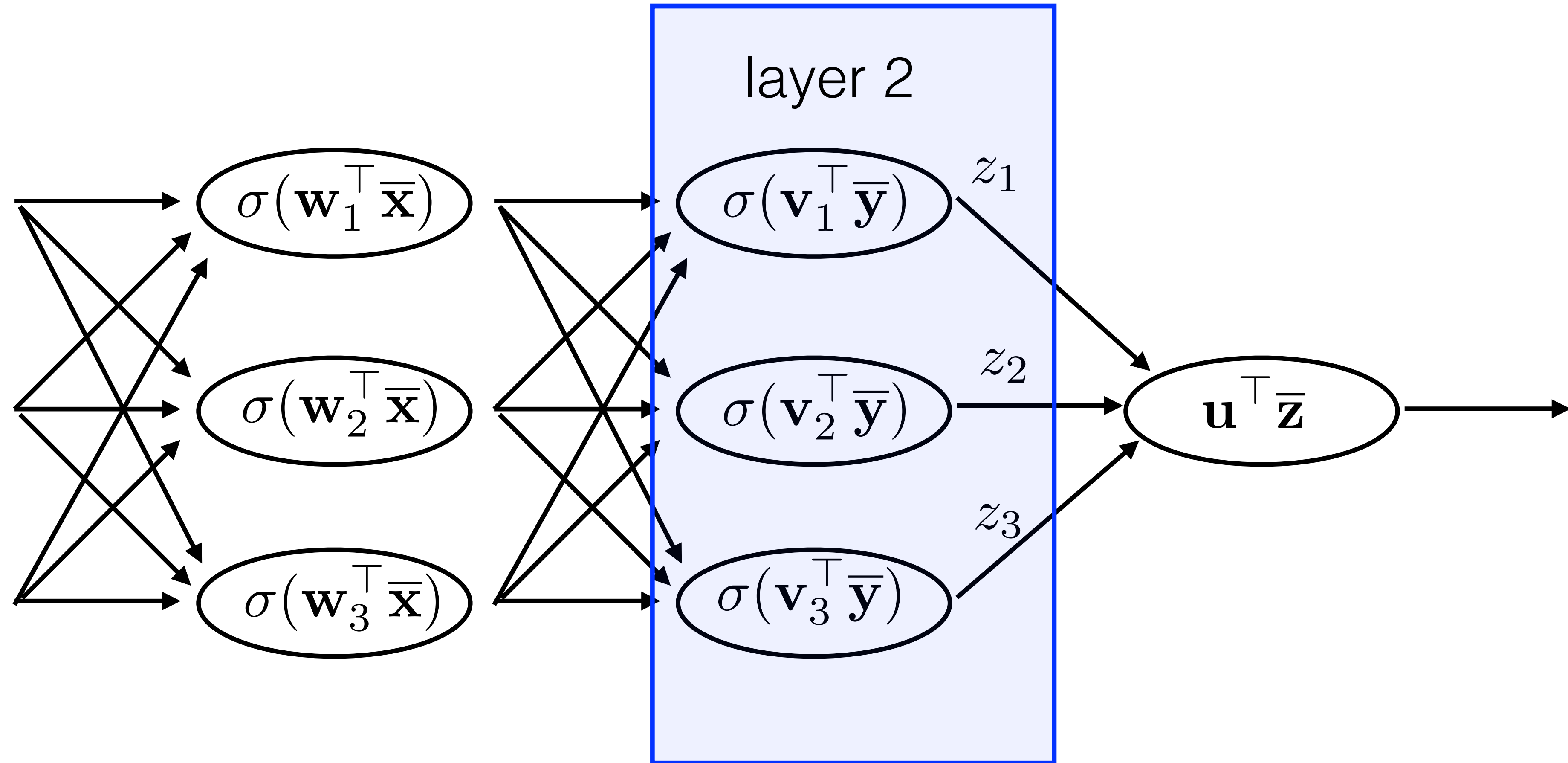
# Fully-connected neural network



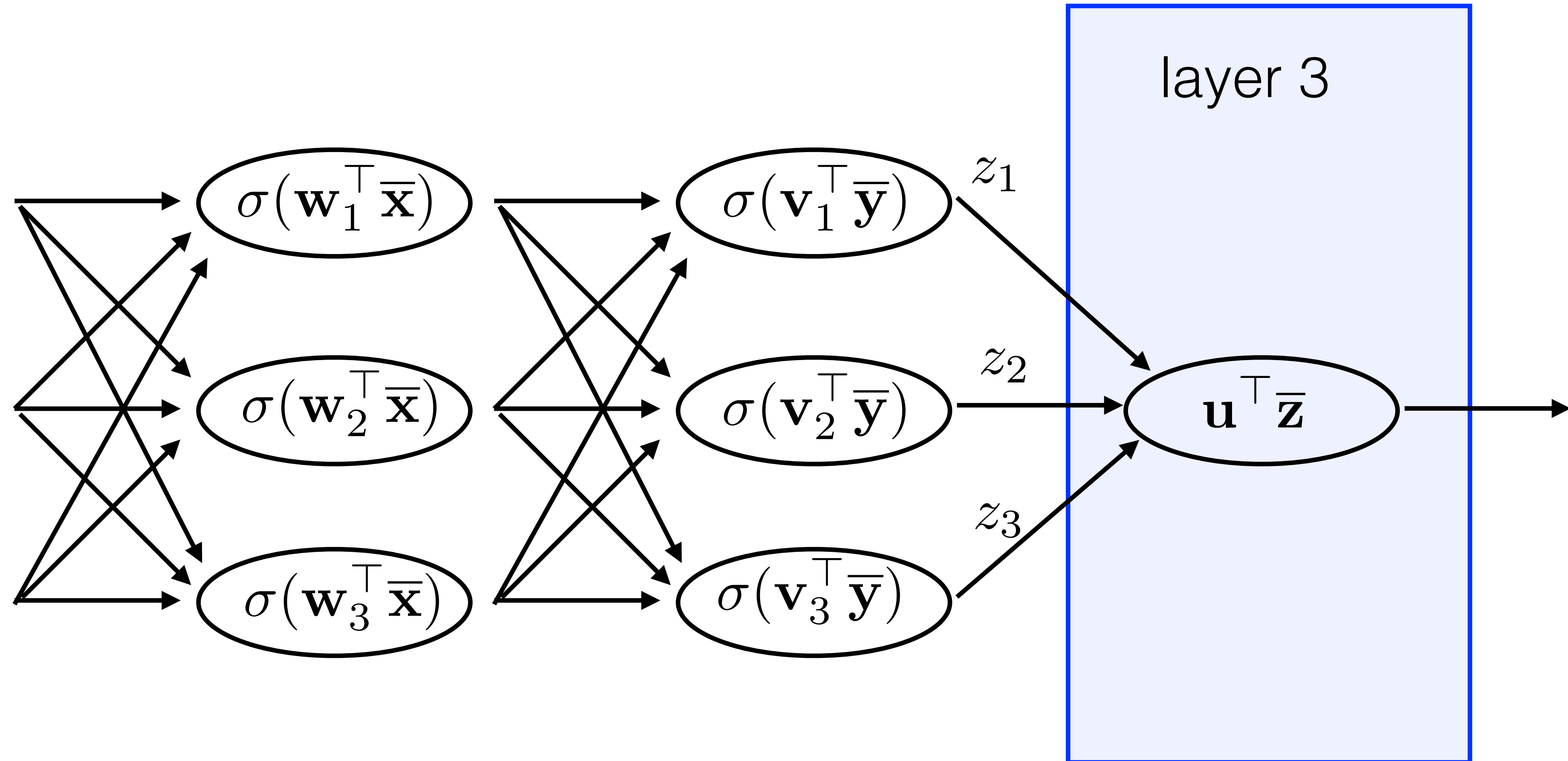
# Fully-connected neural network



# Fully-connected neural network

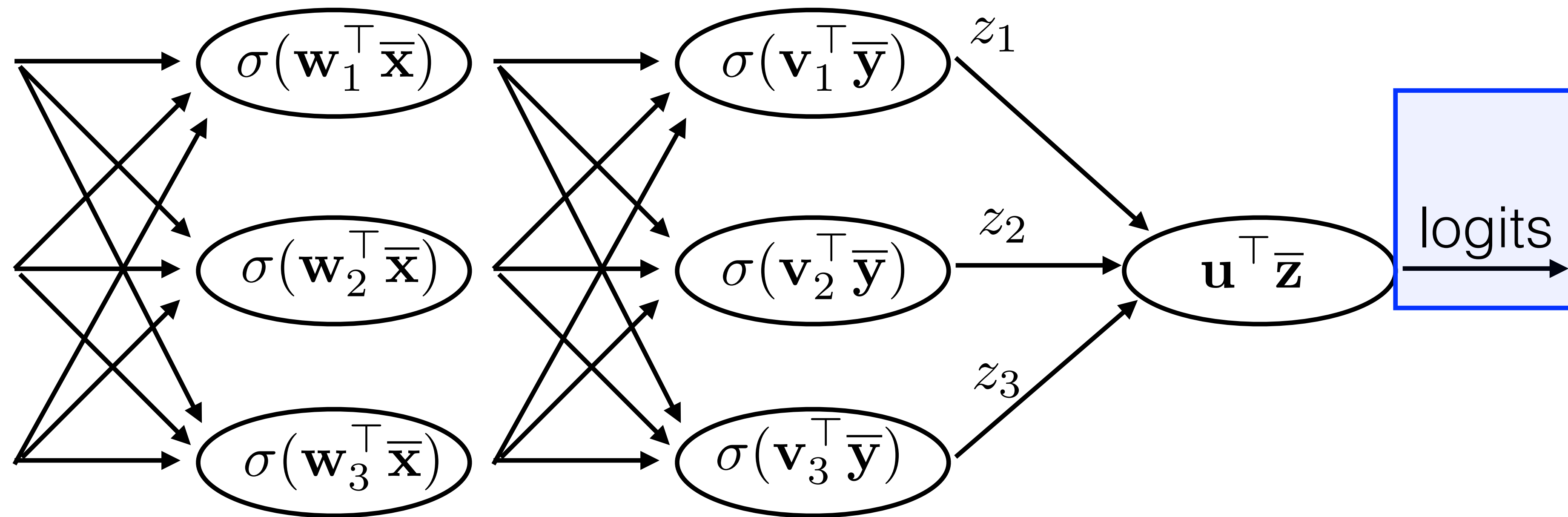


# Fully-connected neural network

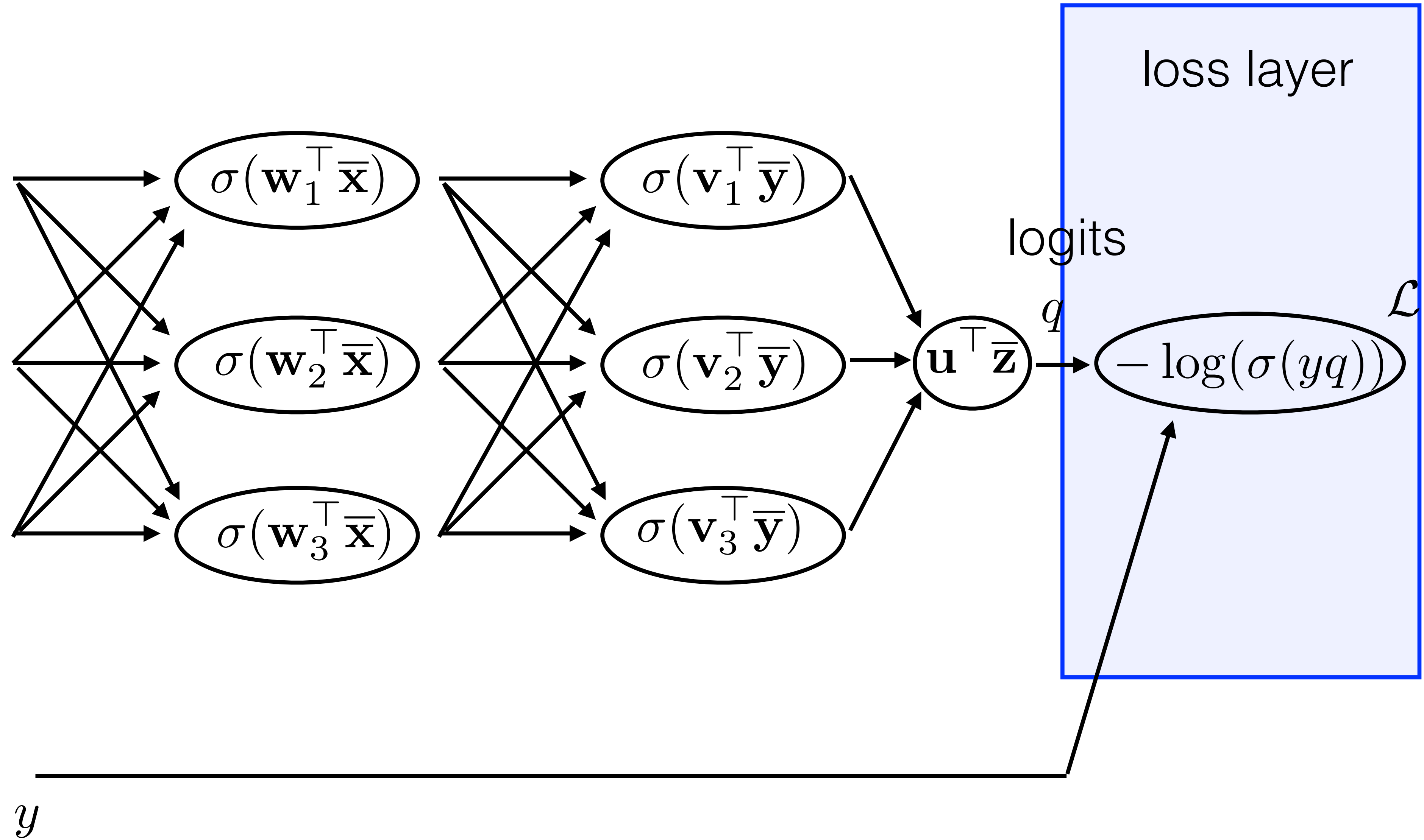




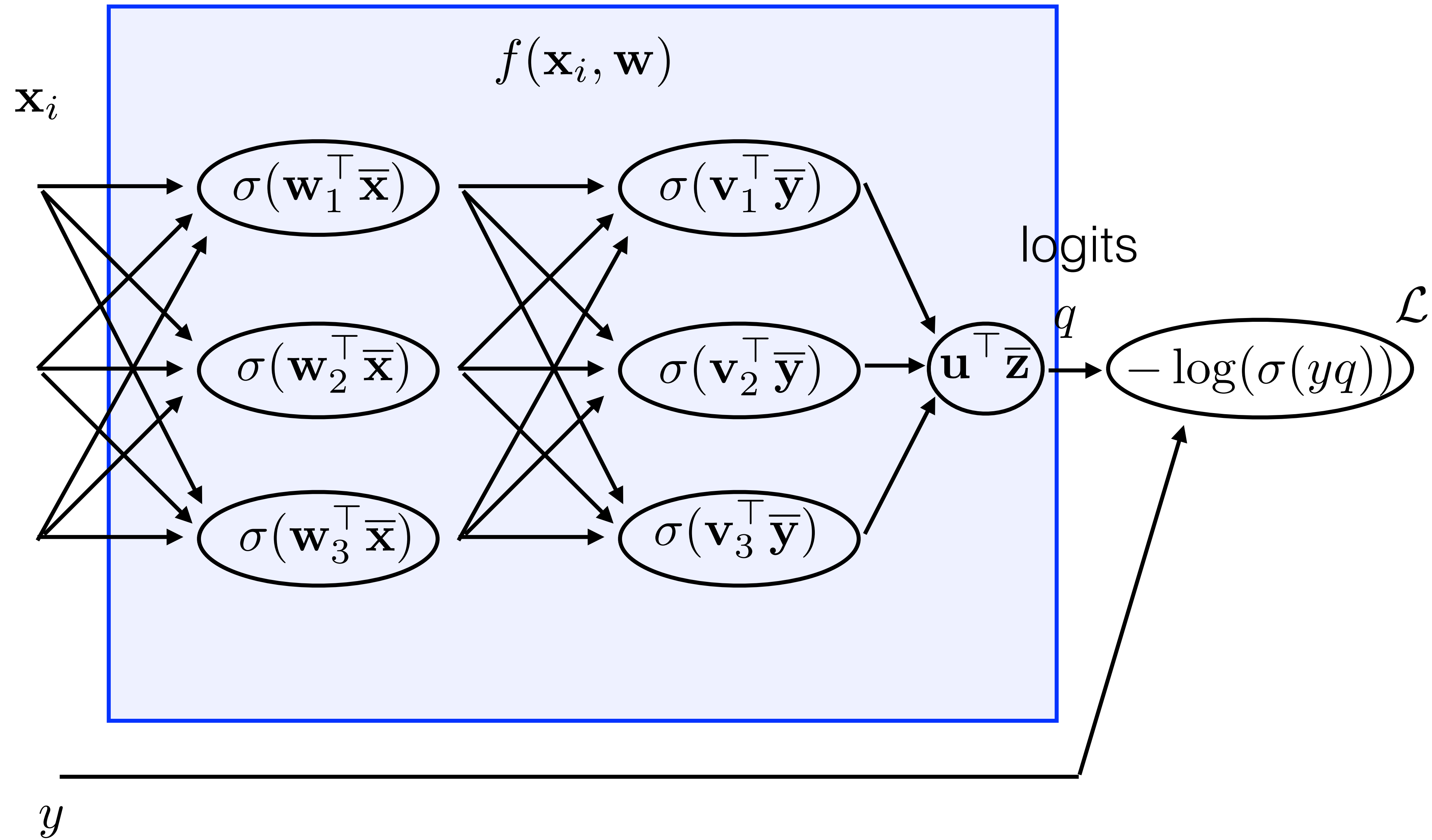
# Fully-connected neural network



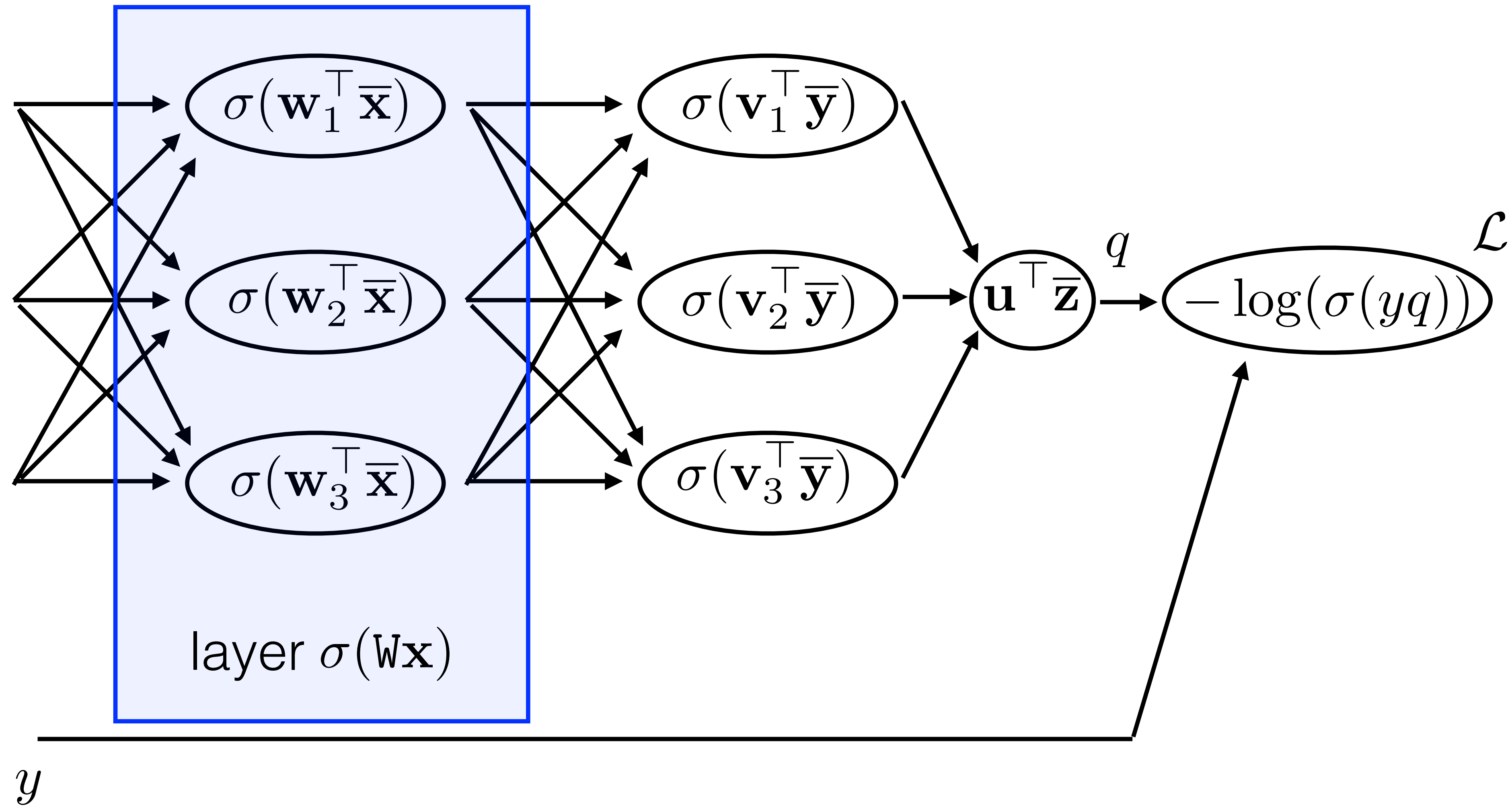
# Fully-connected neural network



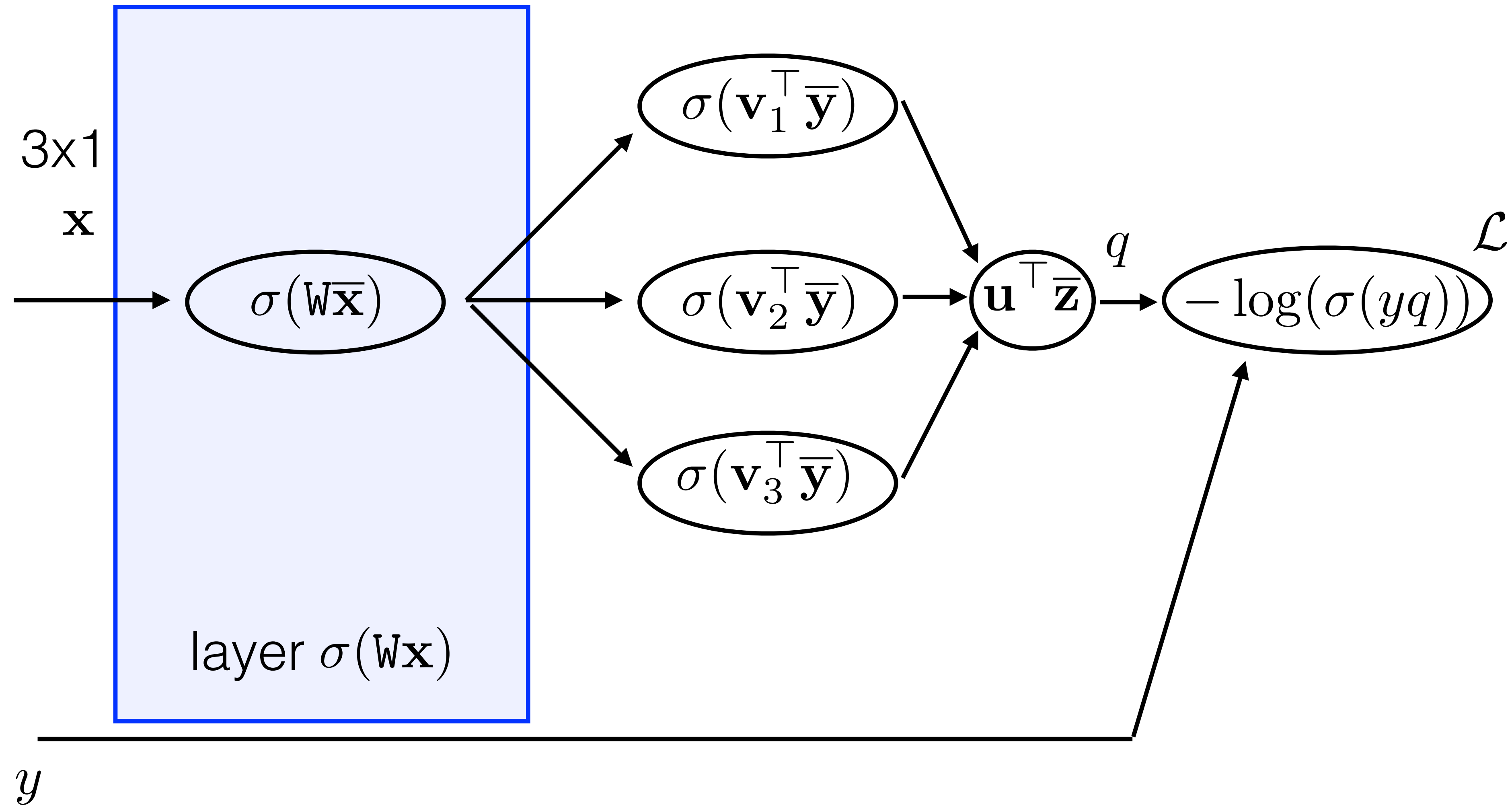
# Fully-connected neural network



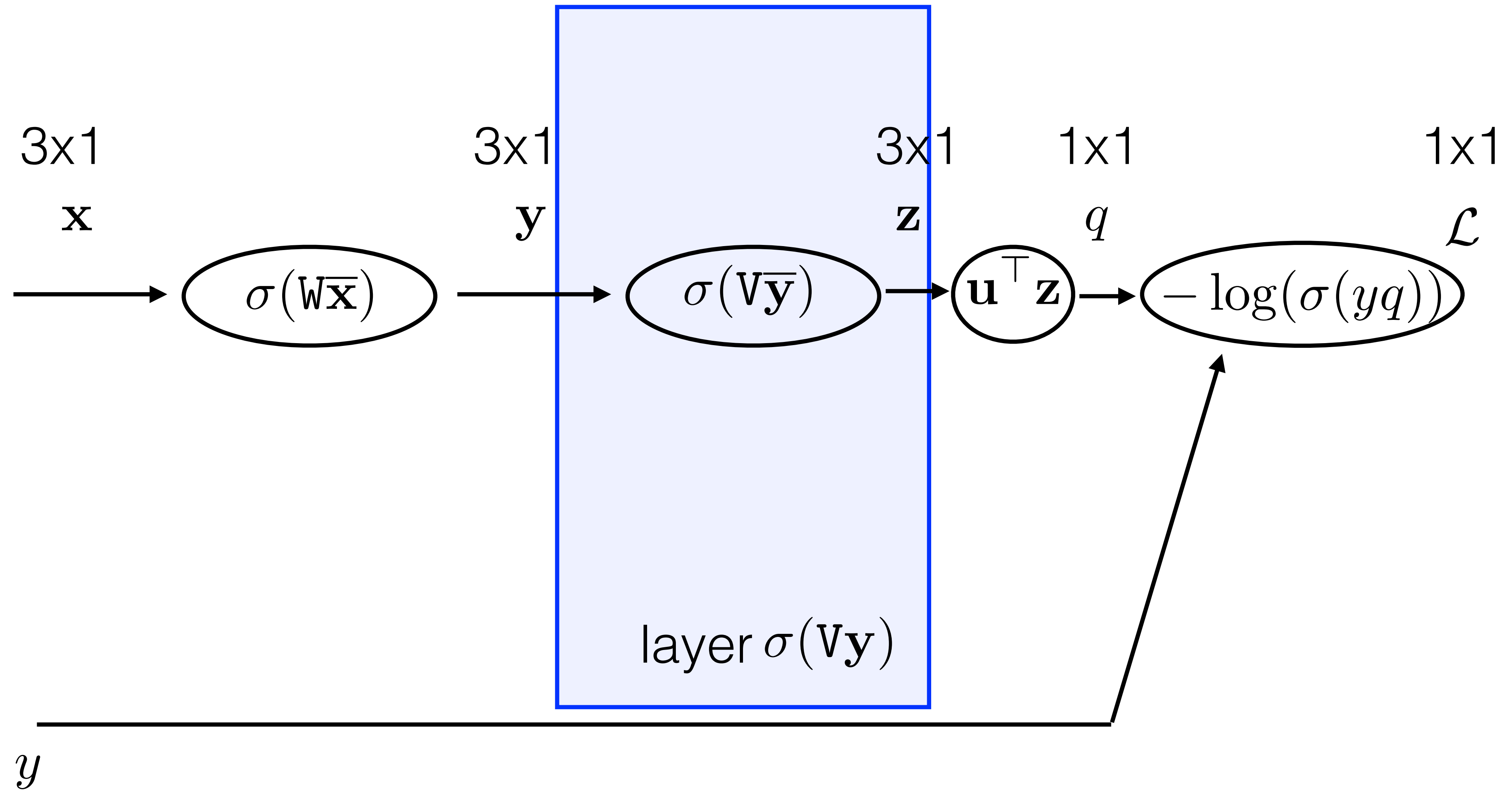
# Fully-connected neural network



# Fully-connected neural network



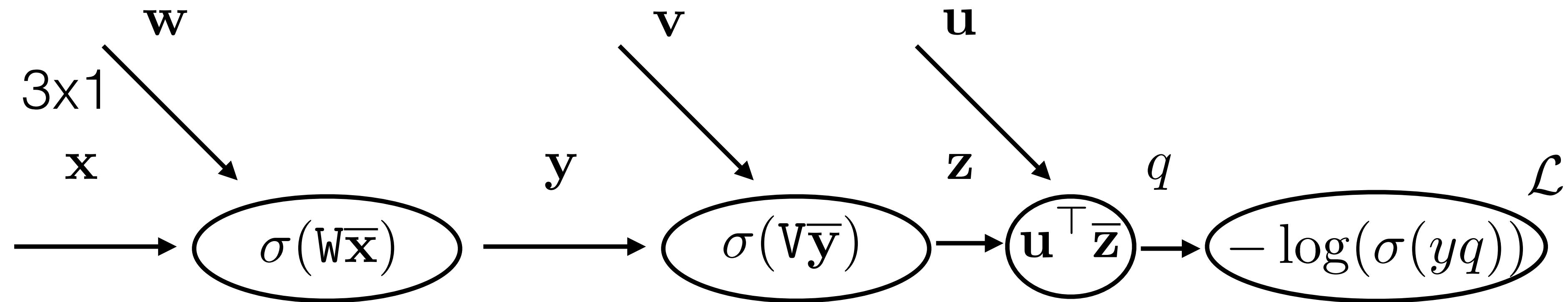
# Fully-connected neural network



# Fully-connected neural network

$$\mathbf{w} = \text{vec}(W)$$

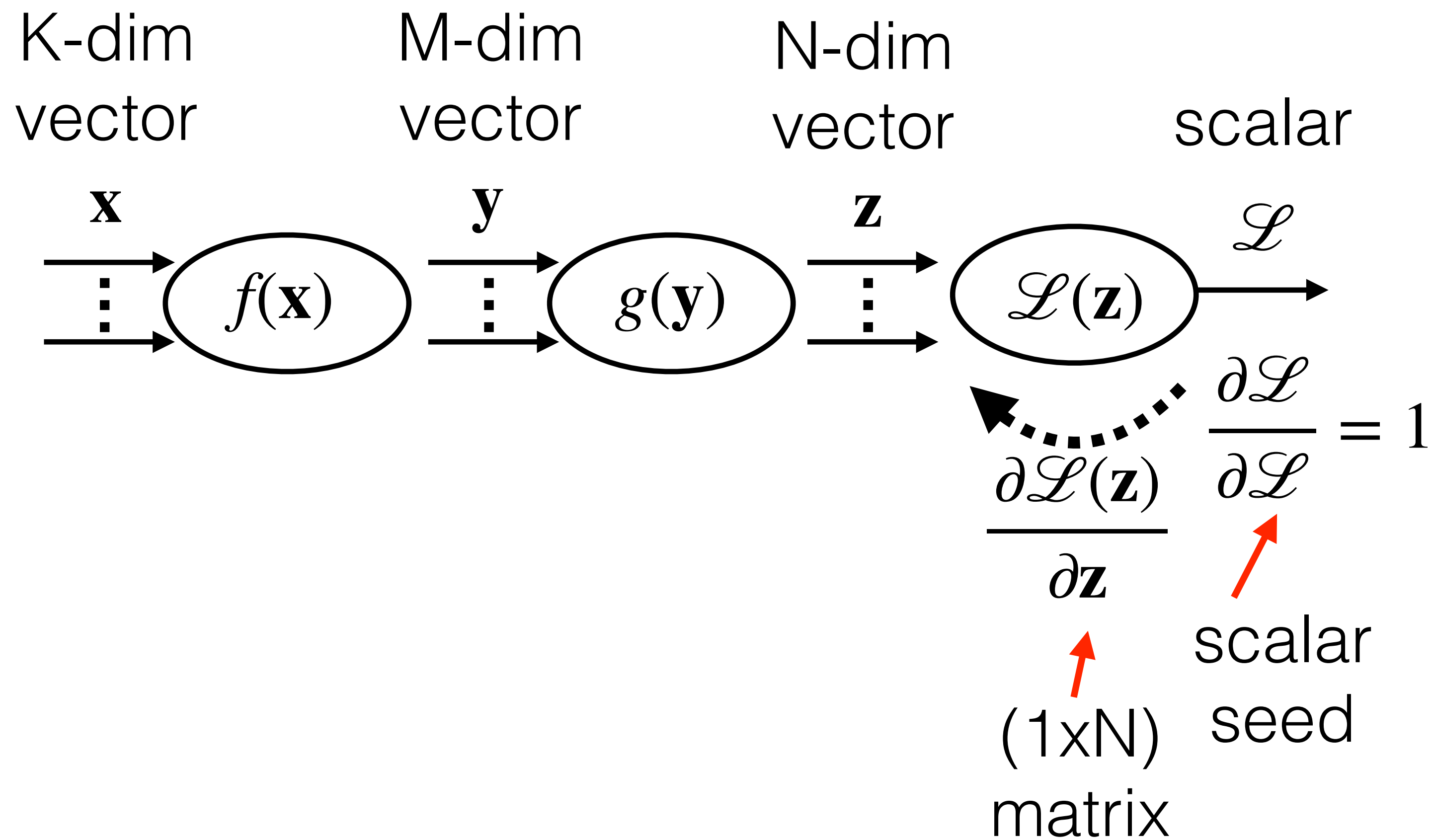
$$\mathbf{v} = \text{vec}(V)$$



Let's summarize the backpropagation  
and then apply on this example

$y$

# Chain-rule in computational graph and Jacobians



Layer:  $f(\mathbf{x}) : \mathbb{R}^K \rightarrow \mathbb{R}^M$

Jacobian:  $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} : \mathbb{R}^K \rightarrow \mathbb{R}^{M \times K}$

$$= \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial \mathbf{x}_1} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial \mathbf{x}_K} \\ \vdots & & \vdots \\ \frac{\partial f_M(\mathbf{x})}{\partial \mathbf{x}_1} & \dots & \frac{\partial f_M(\mathbf{x})}{\partial \mathbf{x}_K} \end{bmatrix} = \begin{bmatrix} \nabla f_1(\mathbf{x})^\top \\ \vdots \\ \nabla f_M(\mathbf{x})^\top \end{bmatrix}$$

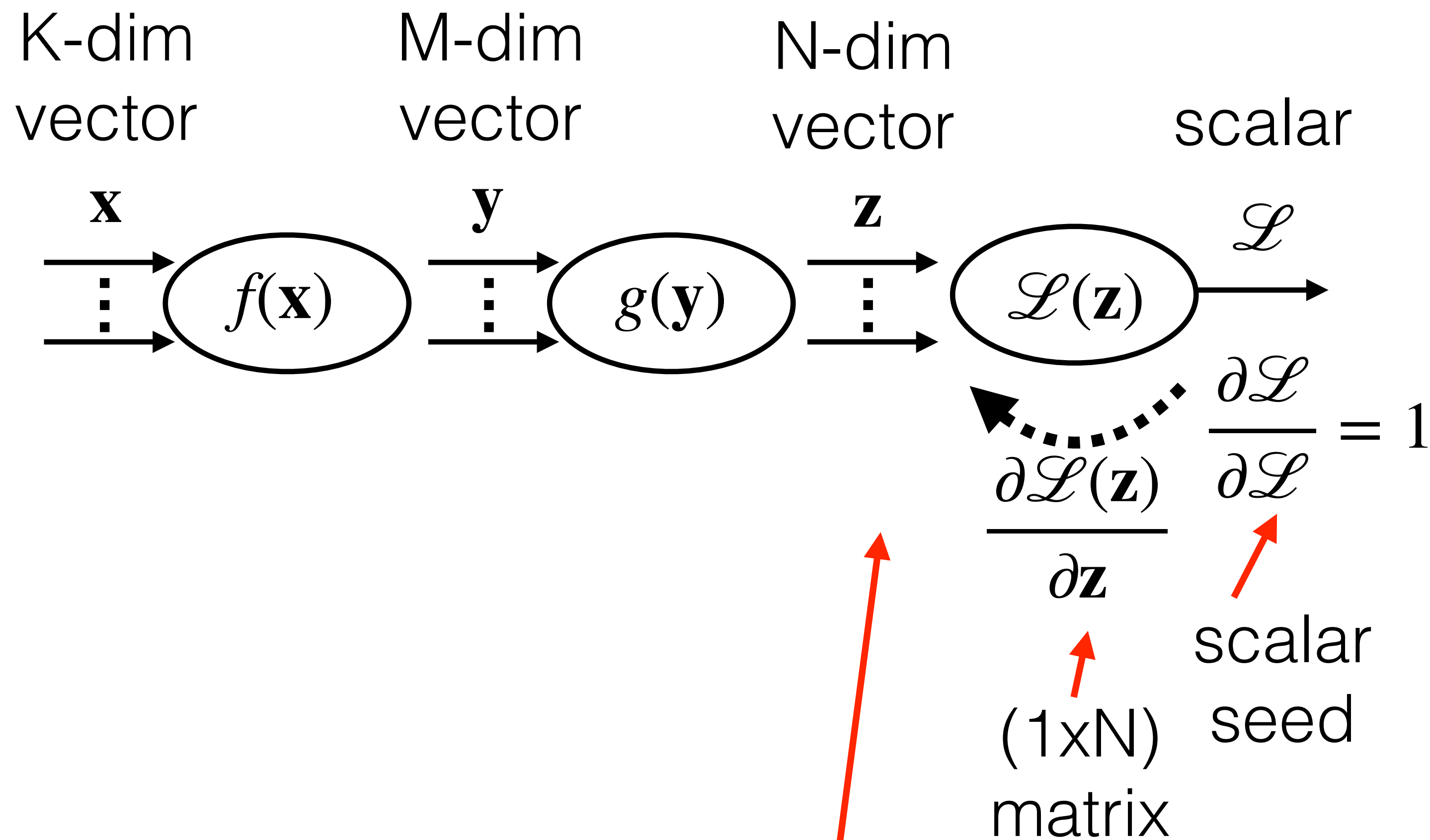
gradient's transpose

Loss jac wrt  $\mathbf{x}$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \begin{matrix} 1 \times K & & 1 \times 1 & & 1 \times N \\ \frac{\partial \mathcal{L}}{\partial \mathbf{x}} & = & 1 & \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \end{matrix}$$



# Chain-rule in computational graph and Jacobians



Layer:  $f(\mathbf{x}) : \mathbb{R}^K \rightarrow \mathbb{R}^M$   
 Jacobian:  $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} : \mathbb{R}^K \rightarrow \mathbb{R}^{M \times K}$

$$= \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial \mathbf{x}_1} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial \mathbf{x}_K} \\ \vdots & & \vdots \\ \frac{\partial f_M(\mathbf{x})}{\partial \mathbf{x}_1} & \dots & \frac{\partial f_M(\mathbf{x})}{\partial \mathbf{x}_K} \end{bmatrix} = \begin{bmatrix} \nabla f_1(\mathbf{x})^\top \\ \vdots \\ \nabla f_M(\mathbf{x})^\top \end{bmatrix}$$

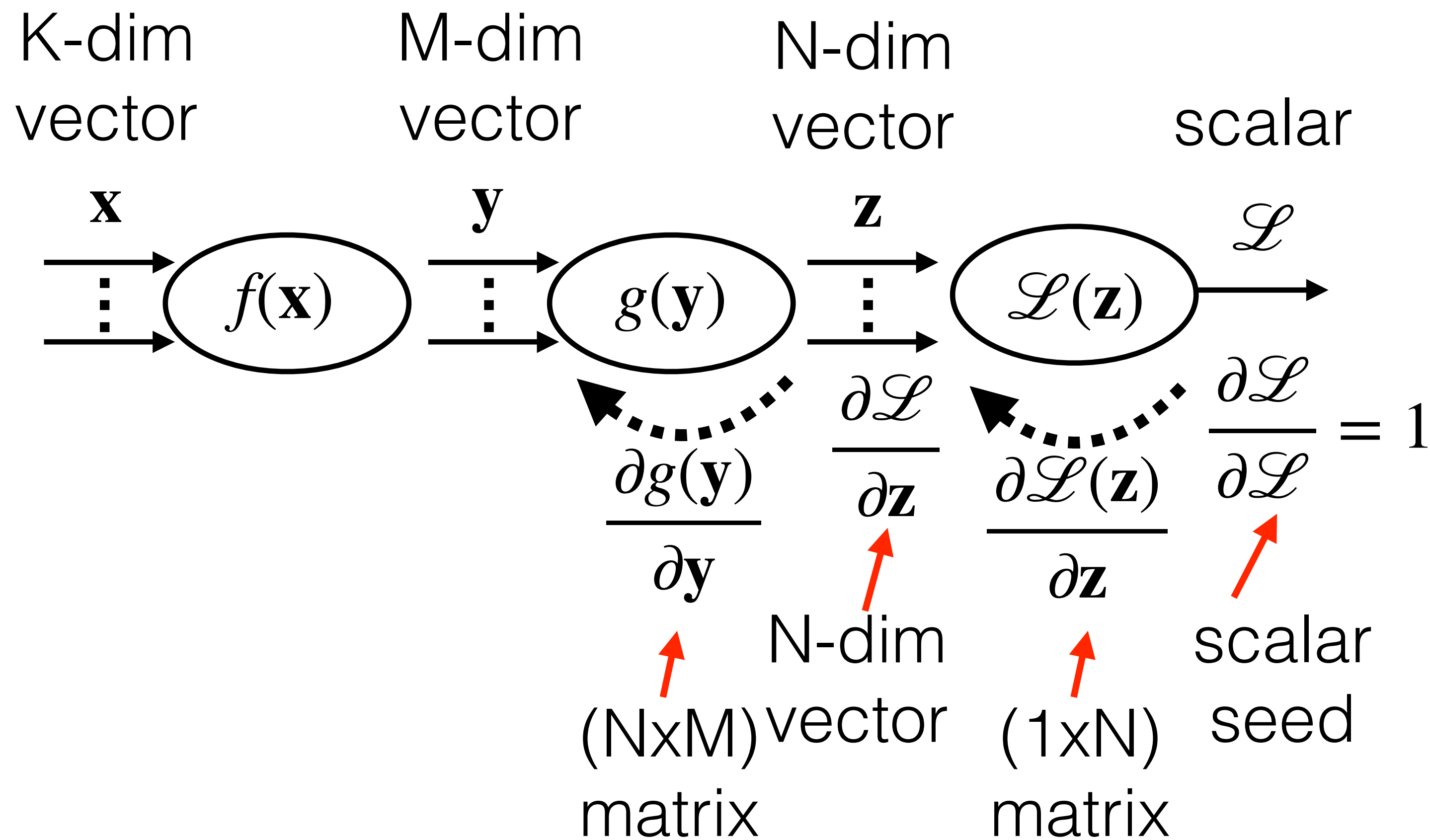
gradient's transpose

Loss jac wrt  $\mathbf{x}$   $\frac{\partial \mathcal{L}}{\partial \mathbf{x}}$

$$= \begin{bmatrix} 1 & \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \end{bmatrix} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \text{ N-dim vector}$$

Dimensions: 1xK, 1x1, 1xN

# Chain-rule in computational graph and Jacobians



Layer:  $f(\mathbf{x}) : \mathbb{R}^K \rightarrow \mathbb{R}^M$   
 Jacobian:  $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} : \mathbb{R}^K \rightarrow \mathbb{R}^{M \times K}$

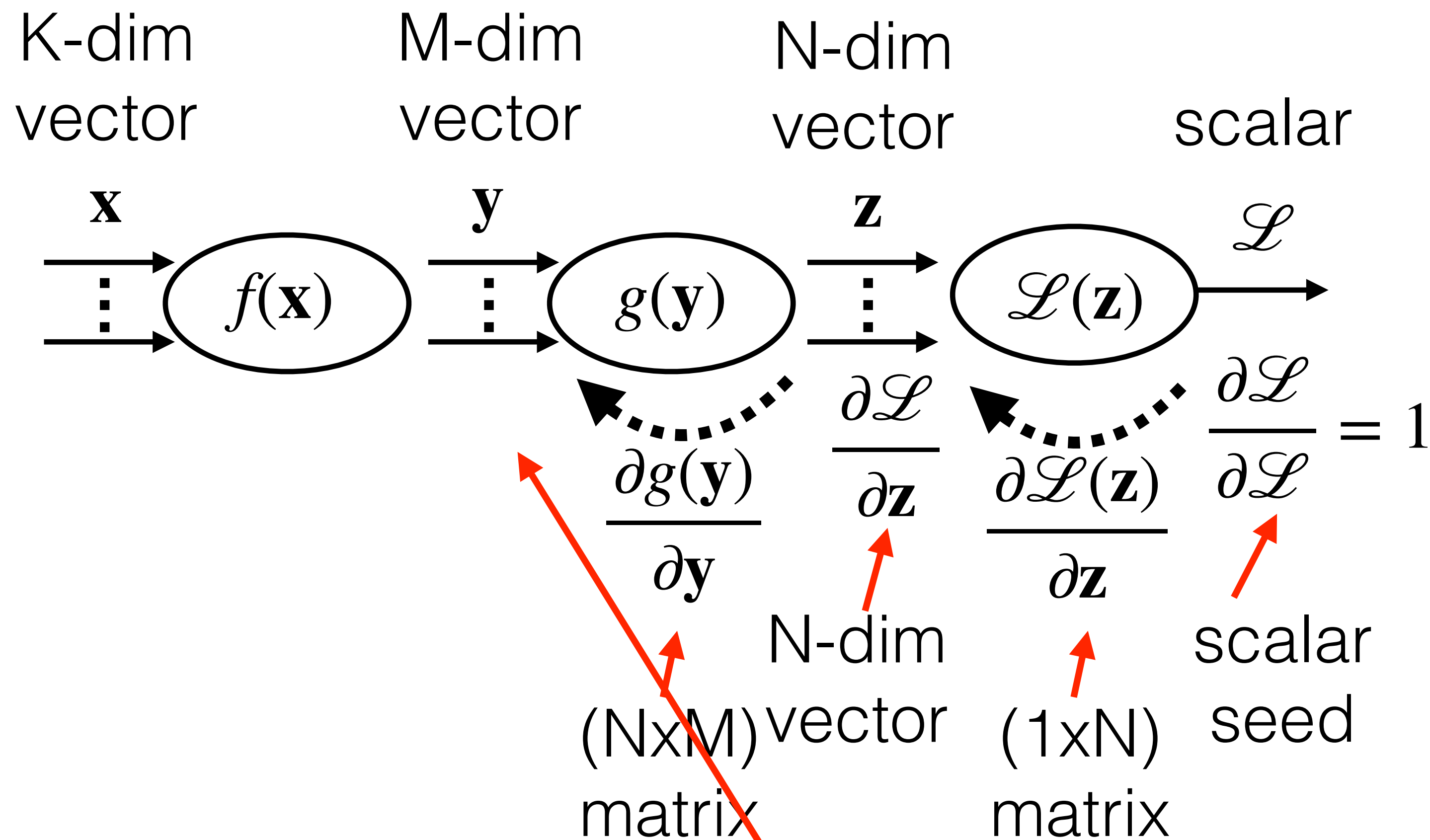
$$= \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial \mathbf{x}_1} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial \mathbf{x}_K} \\ \vdots & & \vdots \\ \frac{\partial f_M(\mathbf{x})}{\partial \mathbf{x}_1} & \dots & \frac{\partial f_M(\mathbf{x})}{\partial \mathbf{x}_K} \end{bmatrix} = \begin{bmatrix} \nabla f_1(\mathbf{x})^\top \\ \vdots \\ \nabla f_M(\mathbf{x})^\top \end{bmatrix}$$

gradient's transpose

Loss jac wrt  $\mathbf{x}$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \begin{matrix} 1 \times K \\ \frac{\partial \mathcal{L}}{\partial \mathbf{x}} \end{matrix} = \begin{matrix} 1 \times 1 \\ 1 \end{matrix} \begin{matrix} 1 \times N \\ \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \end{matrix} \begin{matrix} N \times M \\ \frac{\partial g(\mathbf{y})}{\partial \mathbf{y}} \end{matrix}$$

# Chain-rule in computational graph and Jacobians



Layer:  $f(\mathbf{x}) : \mathbb{R}^K \rightarrow \mathbb{R}^M$   
 Jacobian:  $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} : \mathbb{R}^K \rightarrow \mathbb{R}^{M \times K}$

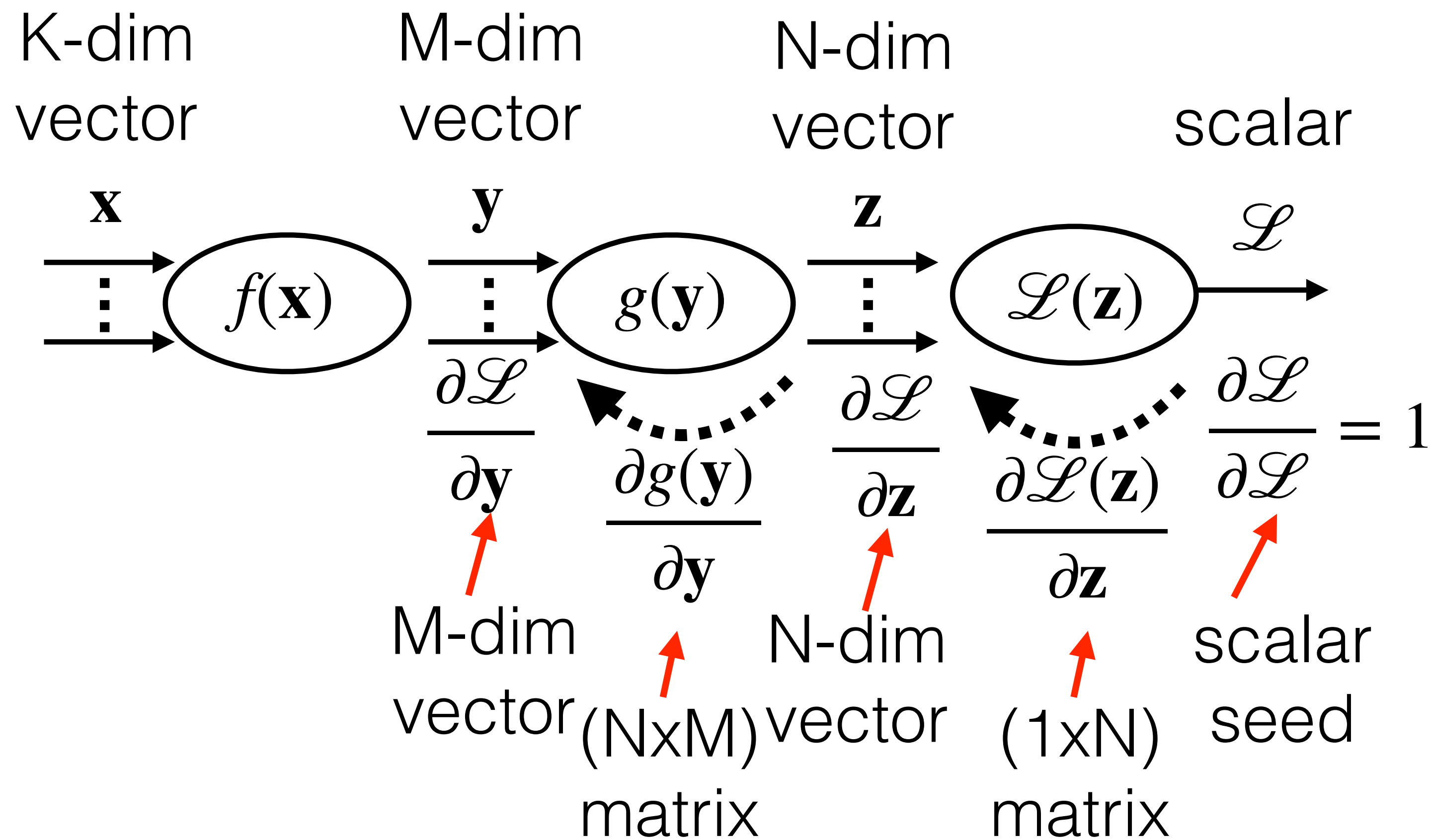
$$= \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial \mathbf{x}_1} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial \mathbf{x}_K} \\ \vdots & & \vdots \\ \frac{\partial f_M(\mathbf{x})}{\partial \mathbf{x}_1} & \dots & \frac{\partial f_M(\mathbf{x})}{\partial \mathbf{x}_K} \end{bmatrix} = \begin{bmatrix} \nabla f_1(\mathbf{x})^\top \\ \vdots \\ \nabla f_M(\mathbf{x})^\top \end{bmatrix}$$

gradient's transpose

Loss jac wrt  $\mathbf{x}$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \begin{bmatrix} 1 & \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \end{bmatrix} \begin{bmatrix} 1 \times 1 & 1 \times N \\ N \times M \end{bmatrix} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}} \quad \text{M-dim vector}$$

# Chain-rule in computational graph and Jacobians



Layer:  $f(\mathbf{x}) : \mathbb{R}^K \rightarrow \mathbb{R}^M$   
 Jacobian:  $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} : \mathbb{R}^K \rightarrow \mathbb{R}^{M \times K}$

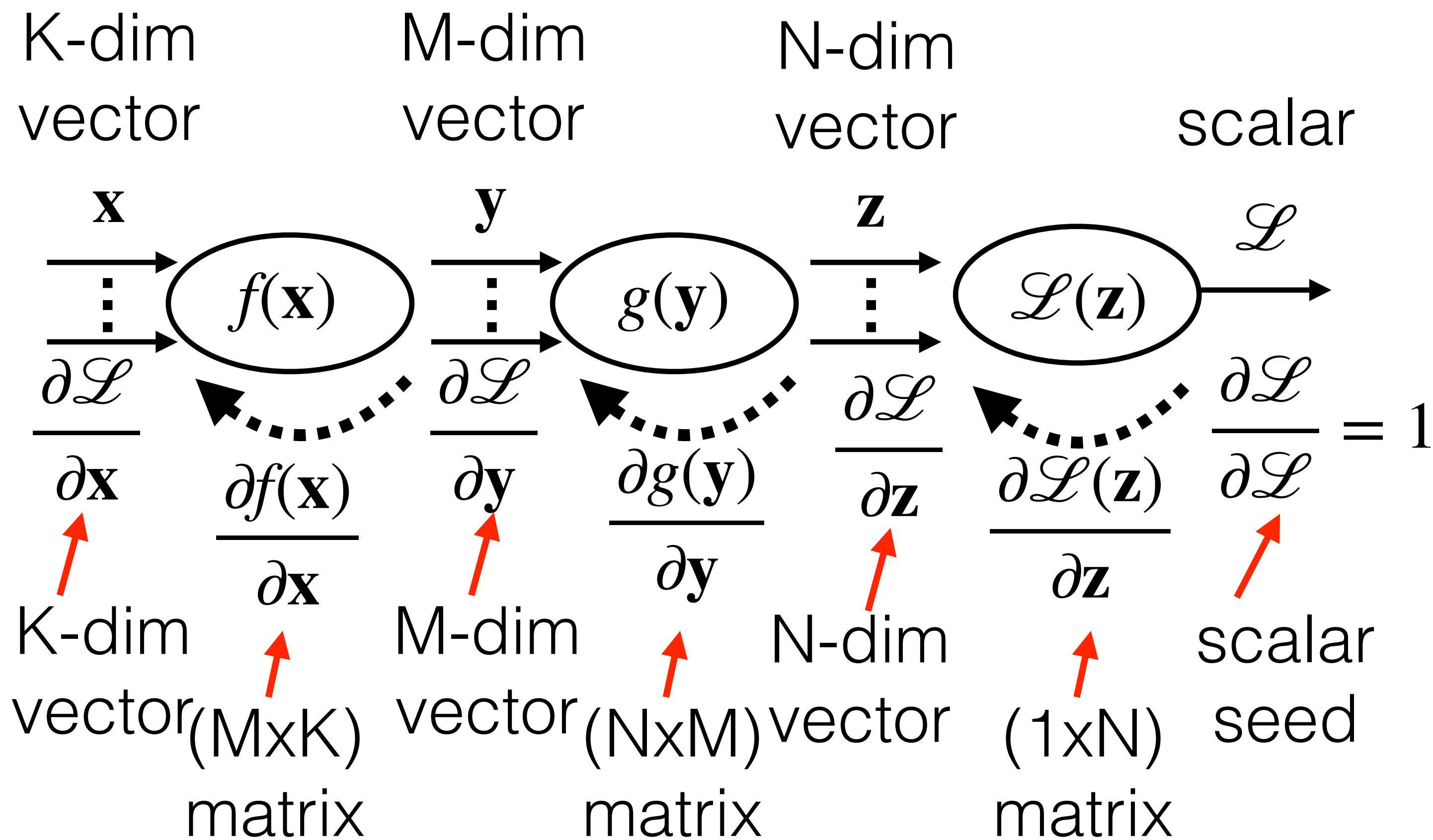
$$= \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial \mathbf{x}_1} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial \mathbf{x}_K} \\ \vdots & & \vdots \\ \frac{\partial f_M(\mathbf{x})}{\partial \mathbf{x}_1} & \dots & \frac{\partial f_M(\mathbf{x})}{\partial \mathbf{x}_K} \end{bmatrix} = \begin{bmatrix} \nabla f_1(\mathbf{x})^\top \\ \vdots \\ \nabla f_M(\mathbf{x})^\top \end{bmatrix}$$

gradient's transpose

Loss jac wrt  $\mathbf{x}$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \begin{matrix} 1 \times K \\ \frac{\partial \mathcal{L}}{\partial \mathbf{x}} \end{matrix} = \begin{matrix} 1 \times 1 \\ 1 \end{matrix} \begin{matrix} 1 \times N \\ \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \end{matrix} \begin{matrix} N \times M \\ \frac{\partial g(\mathbf{y})}{\partial \mathbf{y}} \end{matrix}$$

# Chain-rule in computational graph and Jacobians



Layer:  $f(\mathbf{x}) : \mathbb{R}^K \rightarrow \mathbb{R}^M$   
 Jacobian:  $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} : \mathbb{R}^K \rightarrow \mathbb{R}^{M \times K}$

$$= \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial \mathbf{x}_1} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial \mathbf{x}_K} \\ \vdots & & \vdots \\ \frac{\partial f_M(\mathbf{x})}{\partial \mathbf{x}_1} & \dots & \frac{\partial f_M(\mathbf{x})}{\partial \mathbf{x}_K} \end{bmatrix} = \begin{bmatrix} \nabla f_1(\mathbf{x})^\top \\ \vdots \\ \nabla f_M(\mathbf{x})^\top \end{bmatrix}$$

gradient's transpose

Loss jac wrt  $\mathbf{x}$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = 1 \cdot \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \cdot \frac{\partial g(\mathbf{y})}{\partial \mathbf{y}} \cdot \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$$

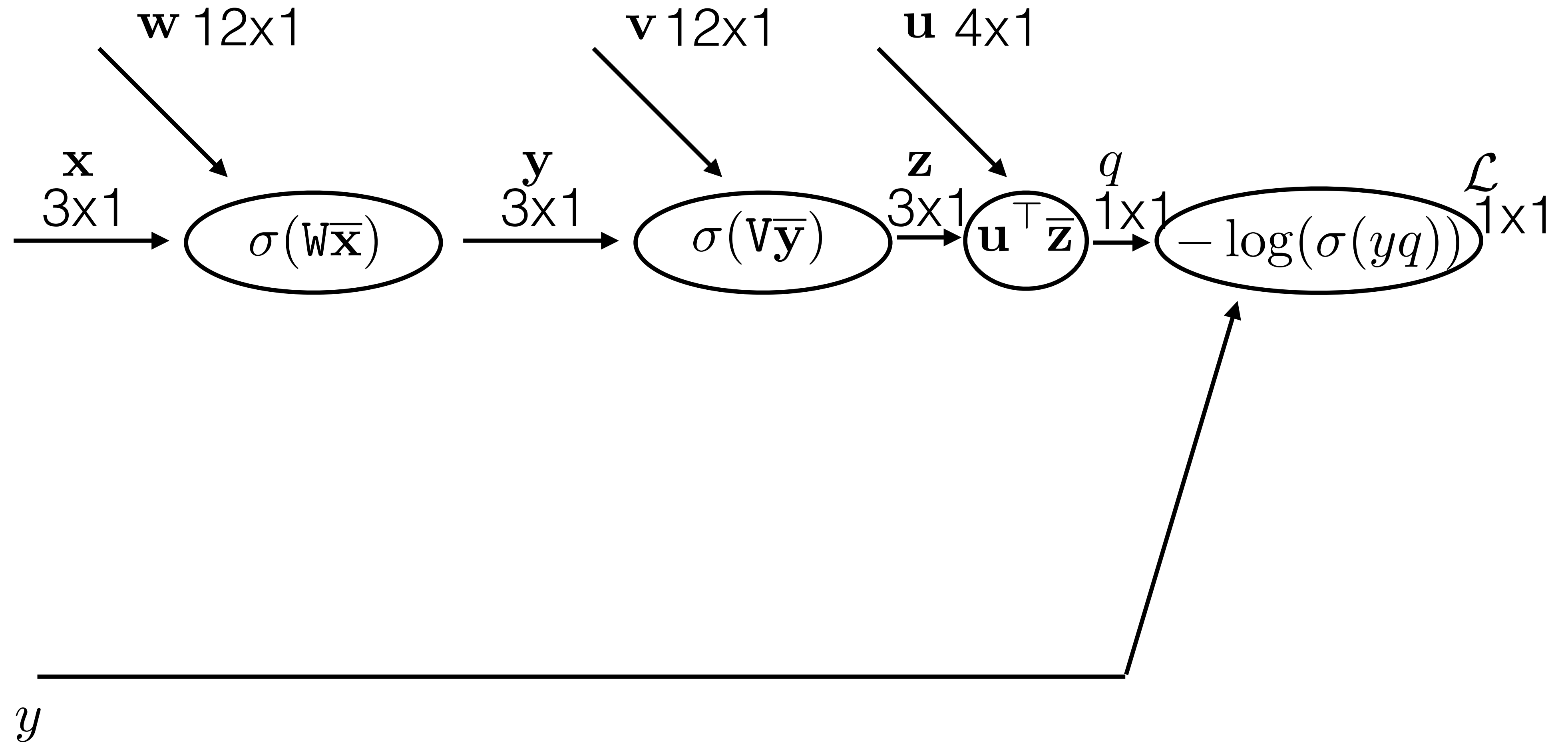
Dimensions:  $1 \times K = 1 \times 1 \cdot 1 \times N \cdot (N \times M) \cdot (M \times K)$

# Fully-connected neural network

$$\mathbf{w} = \text{vec}(W)$$

$$\mathbf{v} = \text{vec}(V)$$

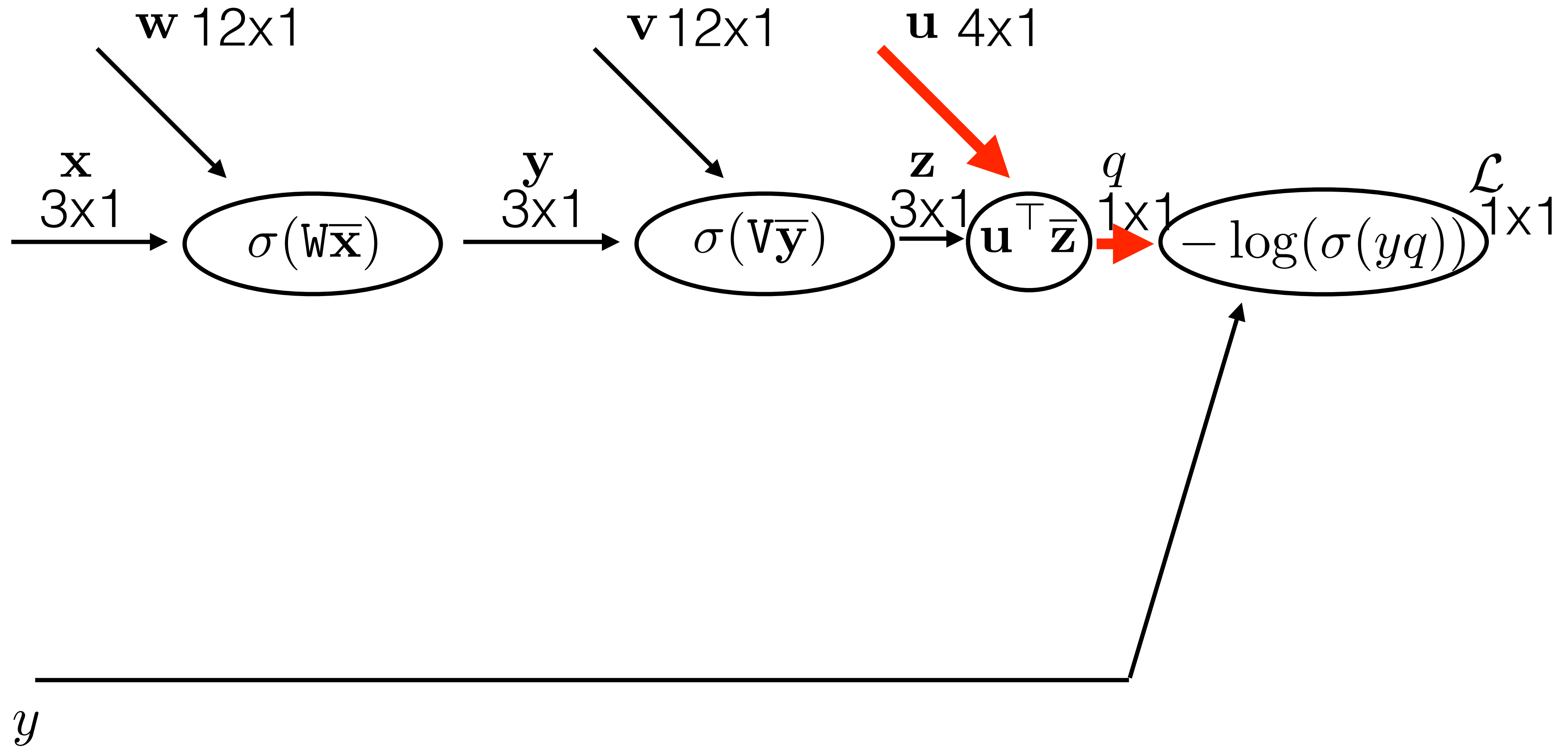
What is the  $\mathbf{w}$ -weight dimensionality?





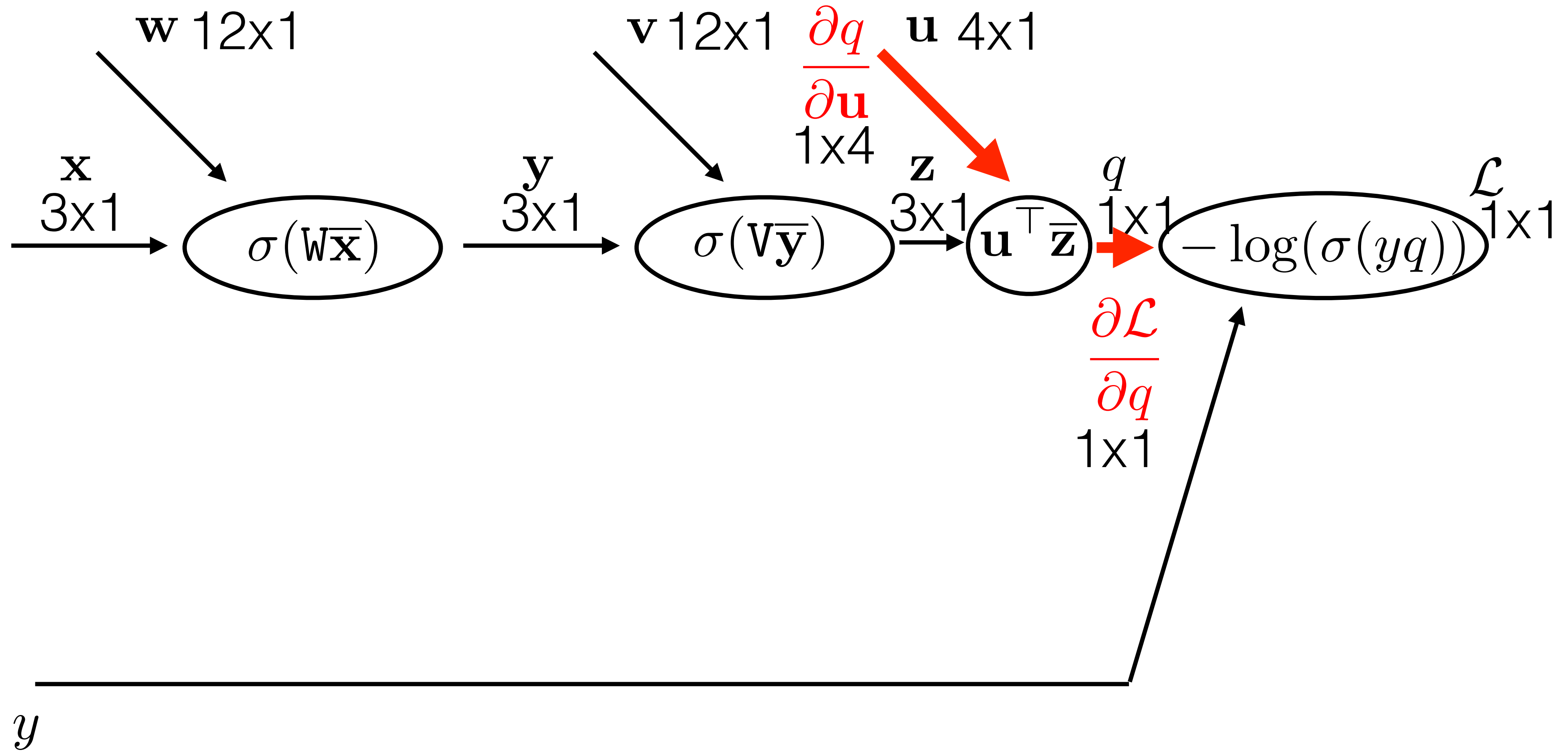
# Chainrule in fully-connected neural network

Jacobian wrt  $\mathbf{u}$  :  $\frac{\partial \mathcal{L}}{\partial \mathbf{u}} = ?$



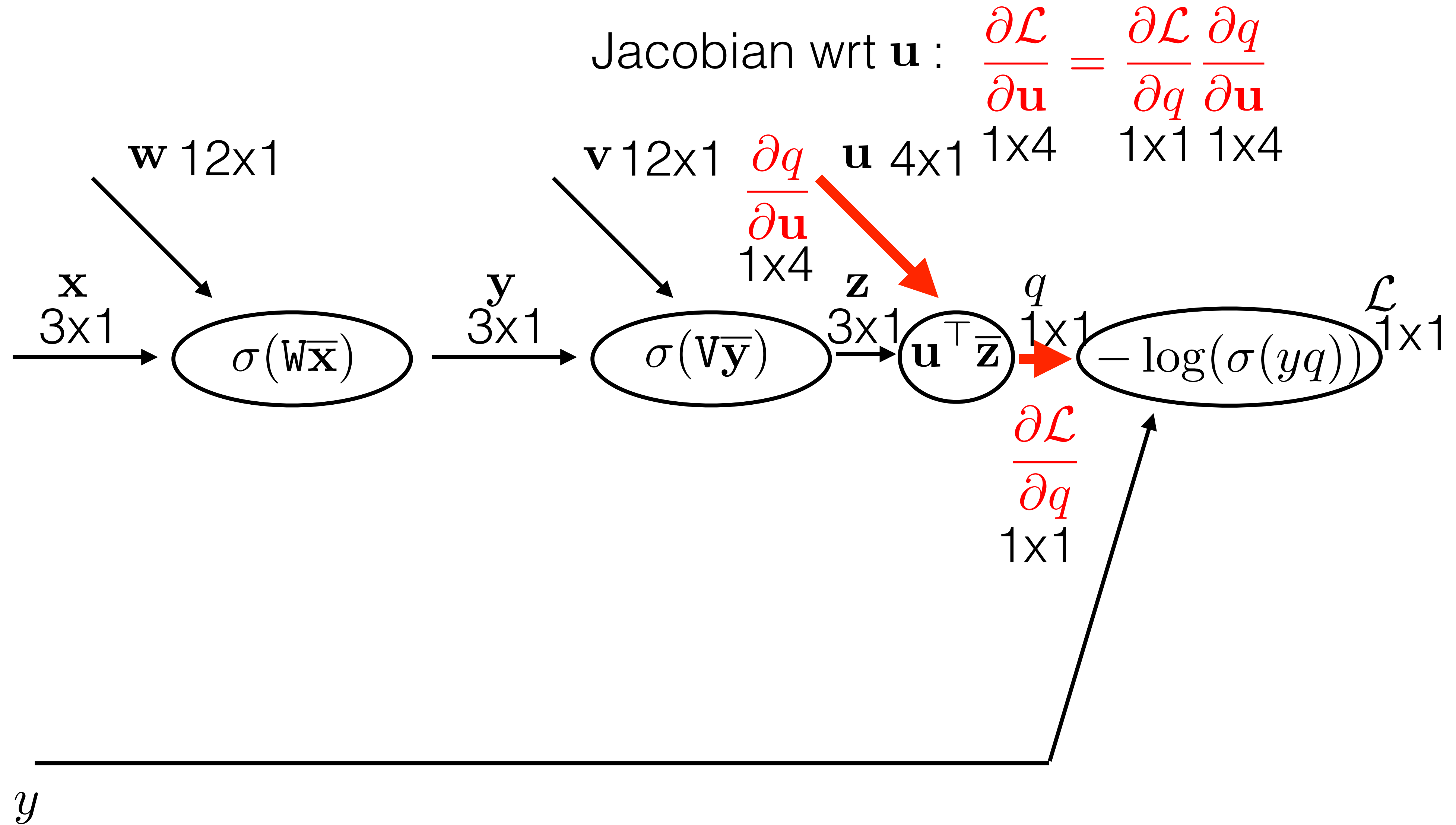
# Chainrule in fully-connected neural network

Jacobian wrt  $\mathbf{u}$  :  $\frac{\partial \mathcal{L}}{\partial \mathbf{u}} = \frac{\partial \mathcal{L}}{\partial q} \frac{\partial q}{\partial \mathbf{u}}$

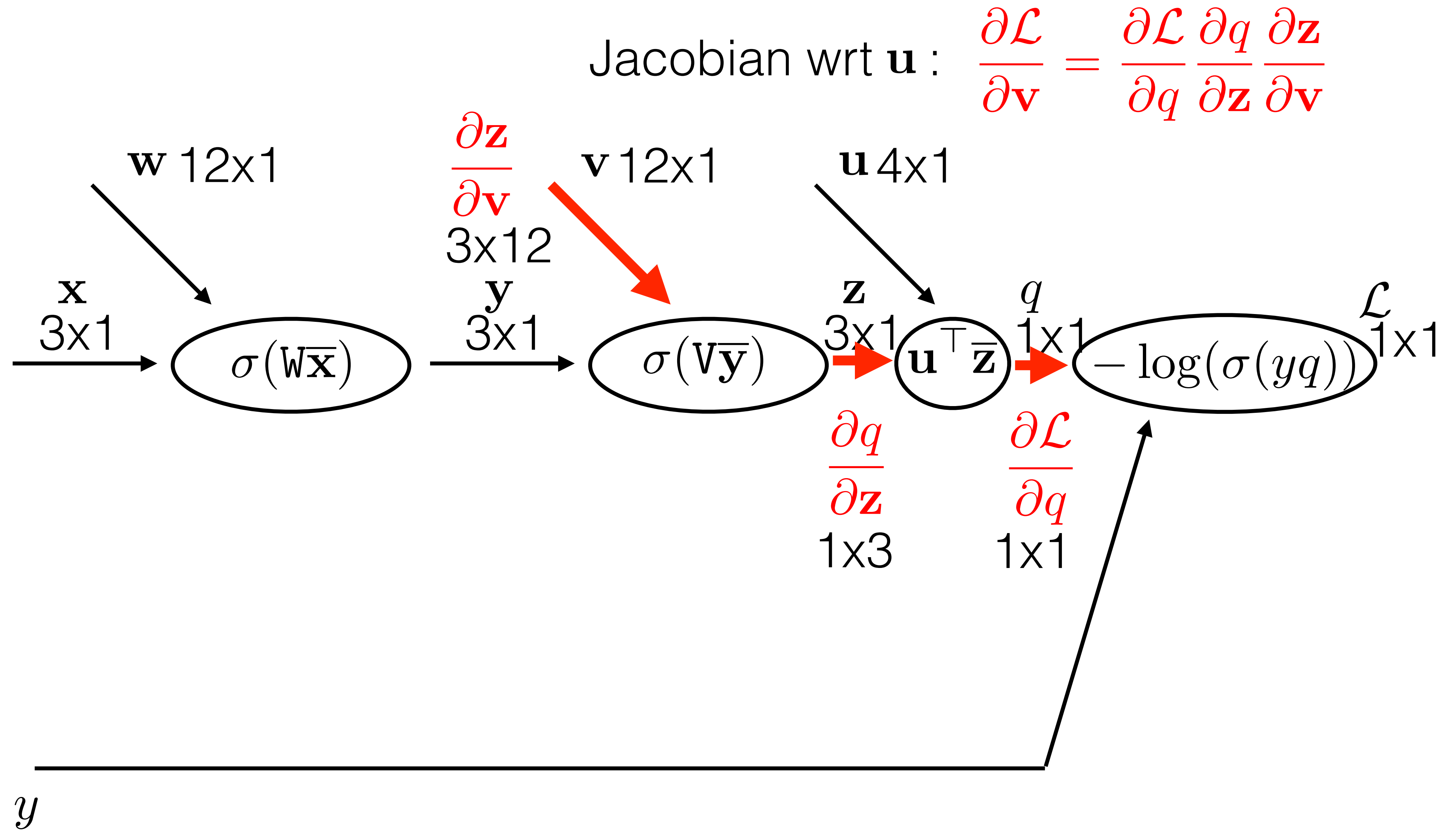




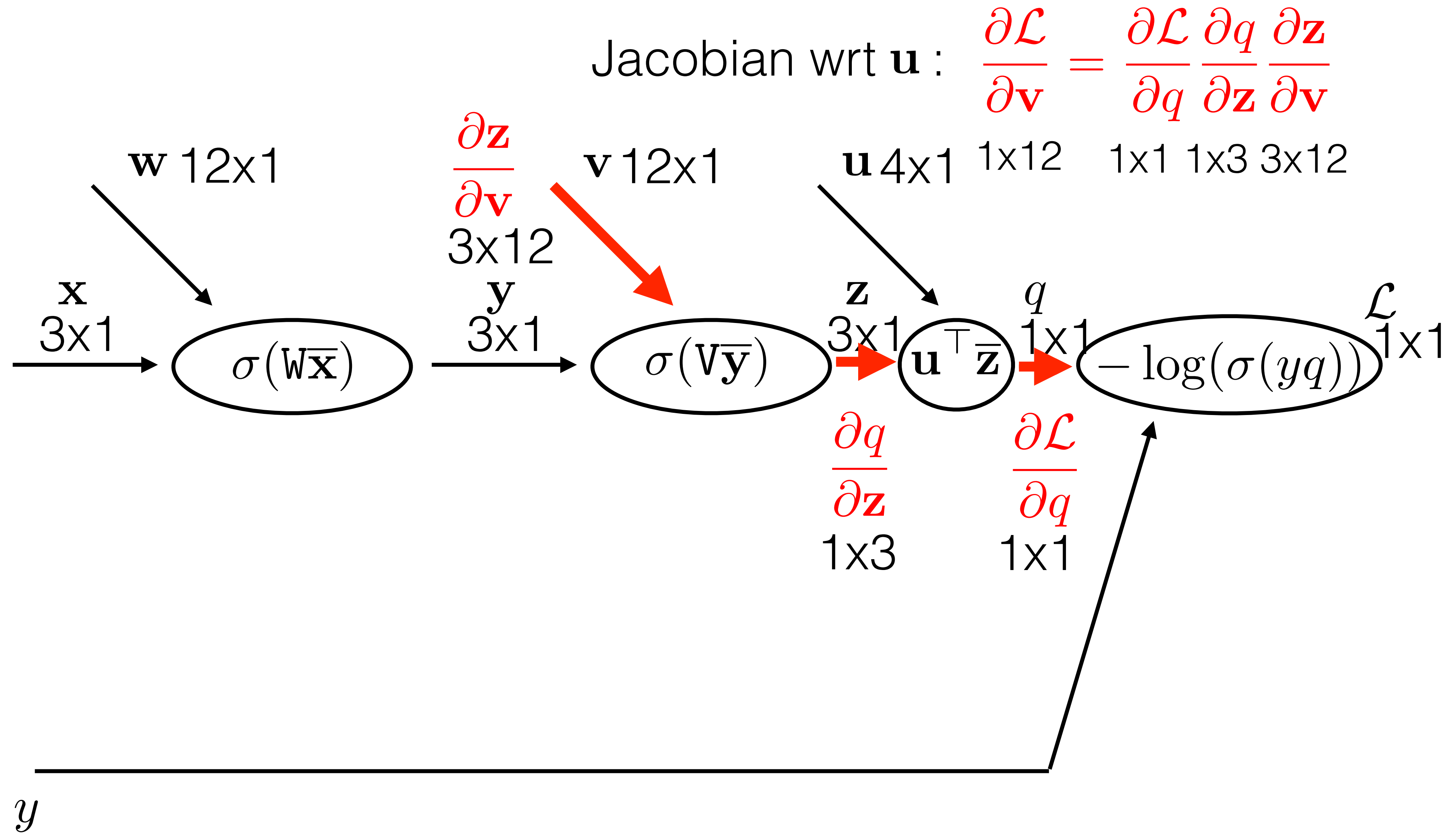
# Chainrule in fully-connected neural network



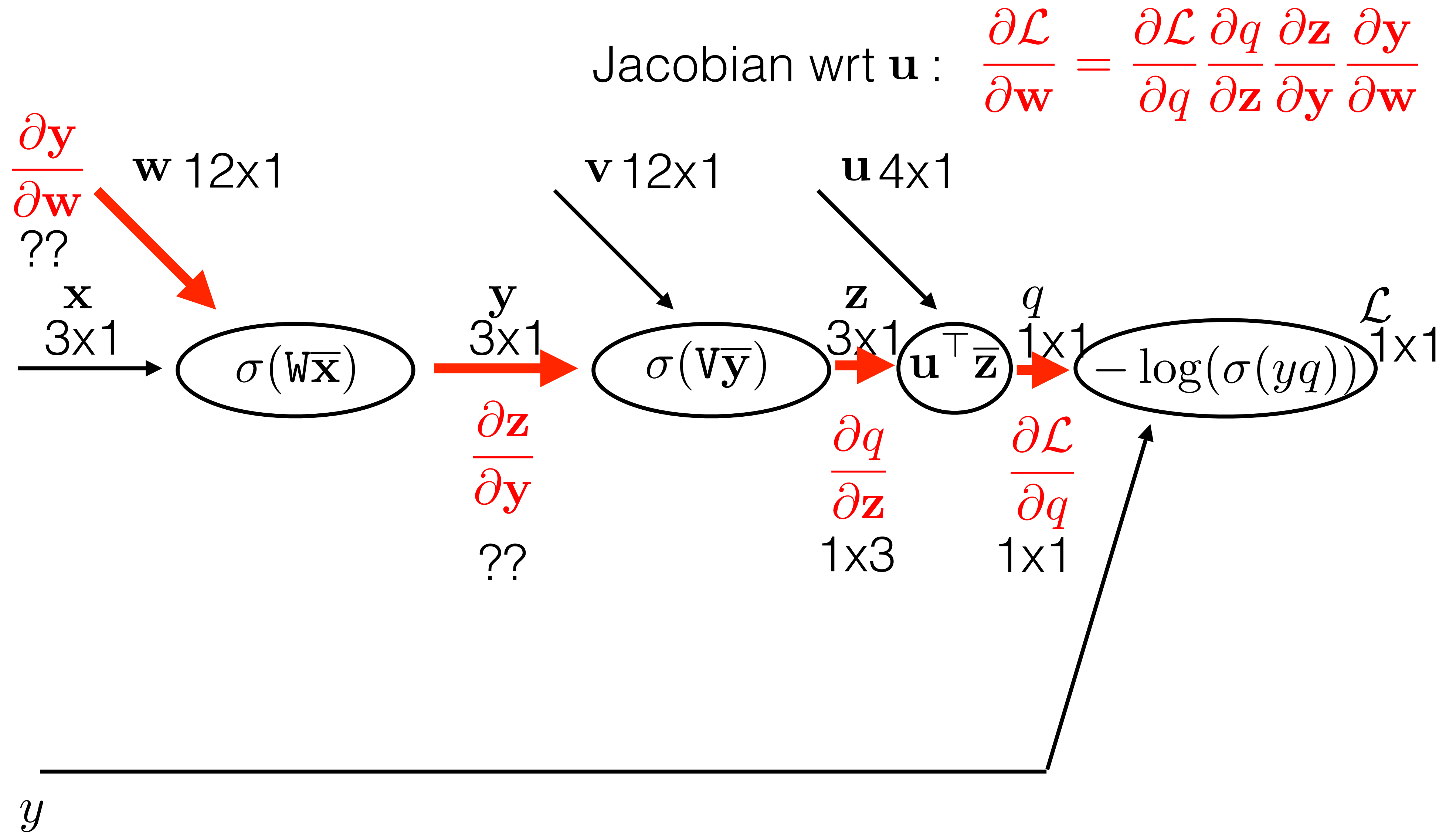
# Chainrule in fully-connected neural network



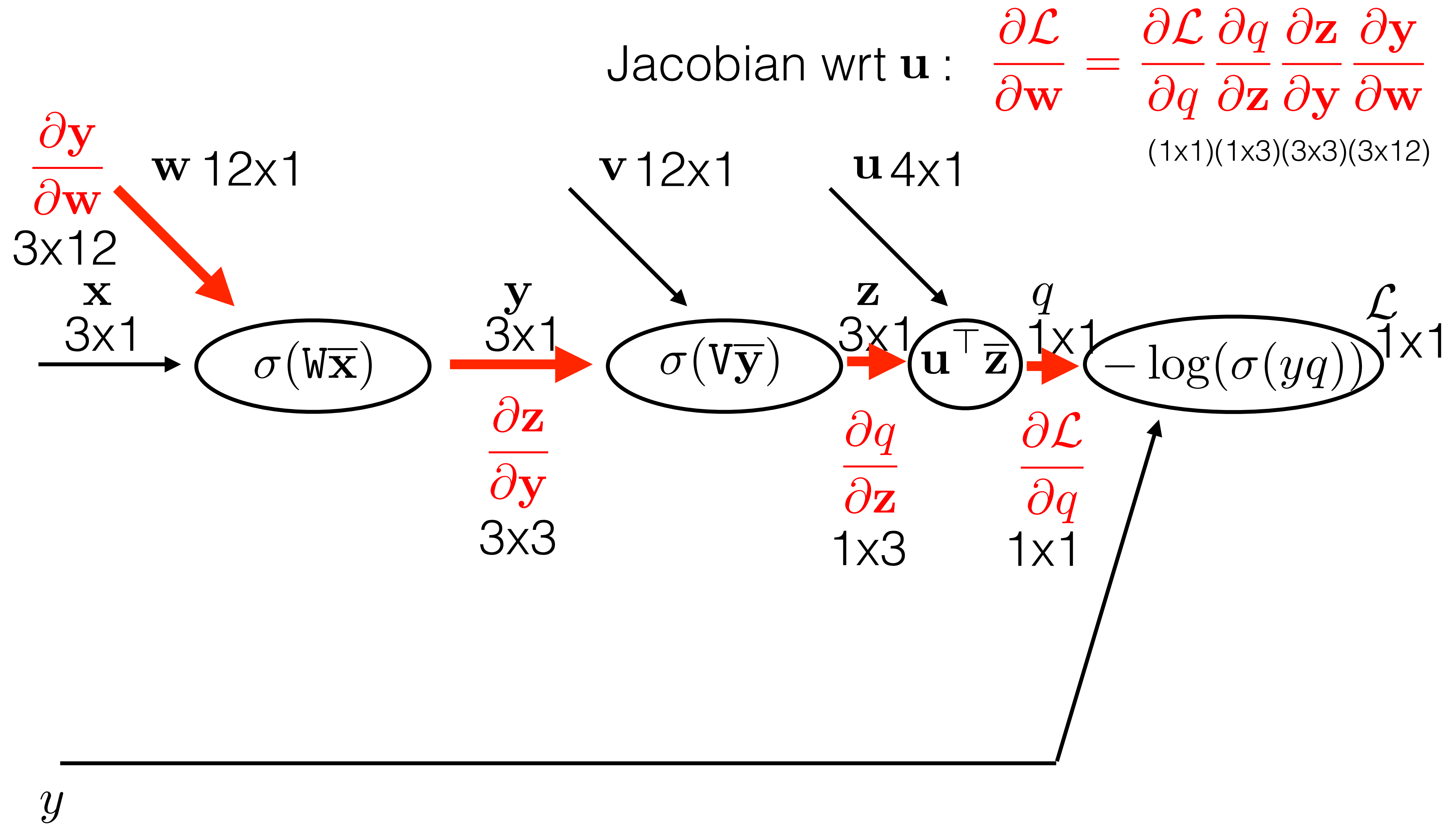
# Chainrule in fully-connected neural network



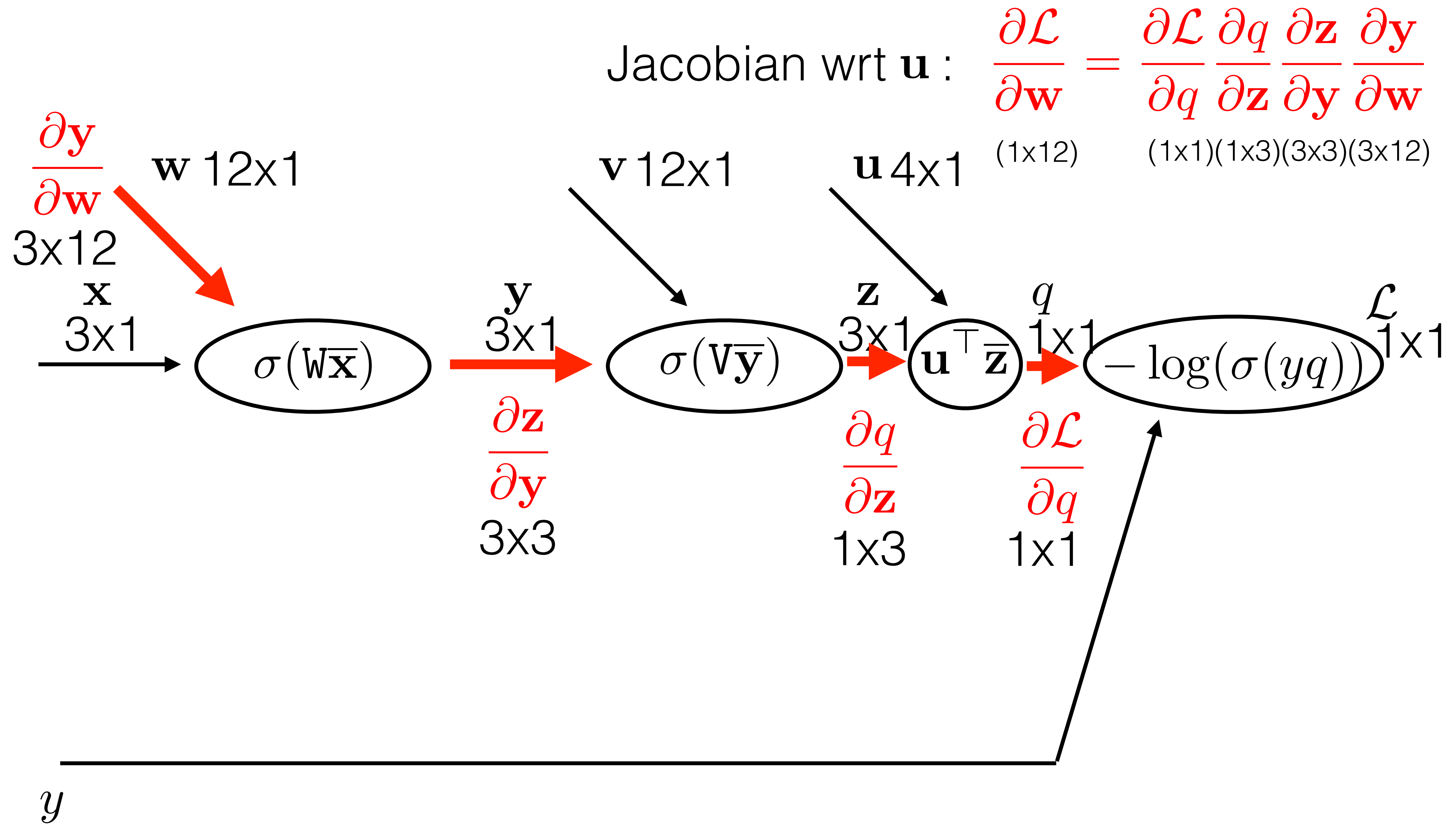
# Chainrule in fully-connected neural network



# Chainrule in fully-connected neural network

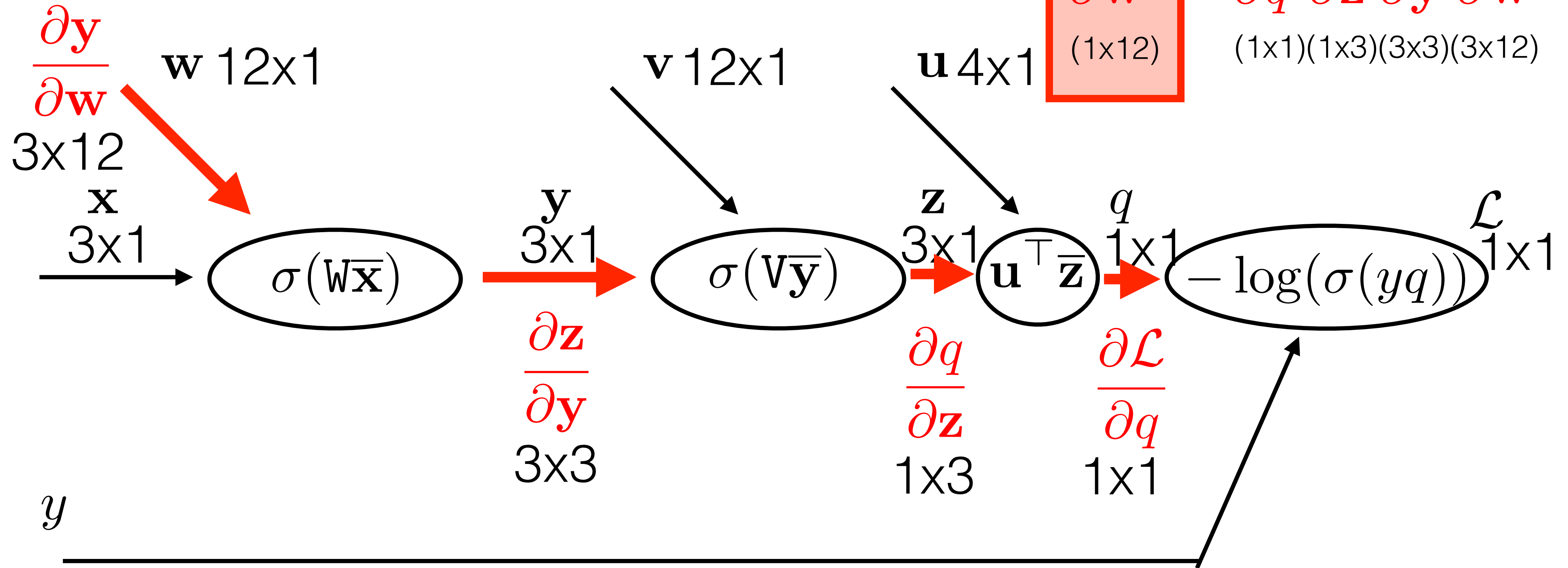


# Chainrule in fully-connected neural network



# Efficient implementation of autodiff?

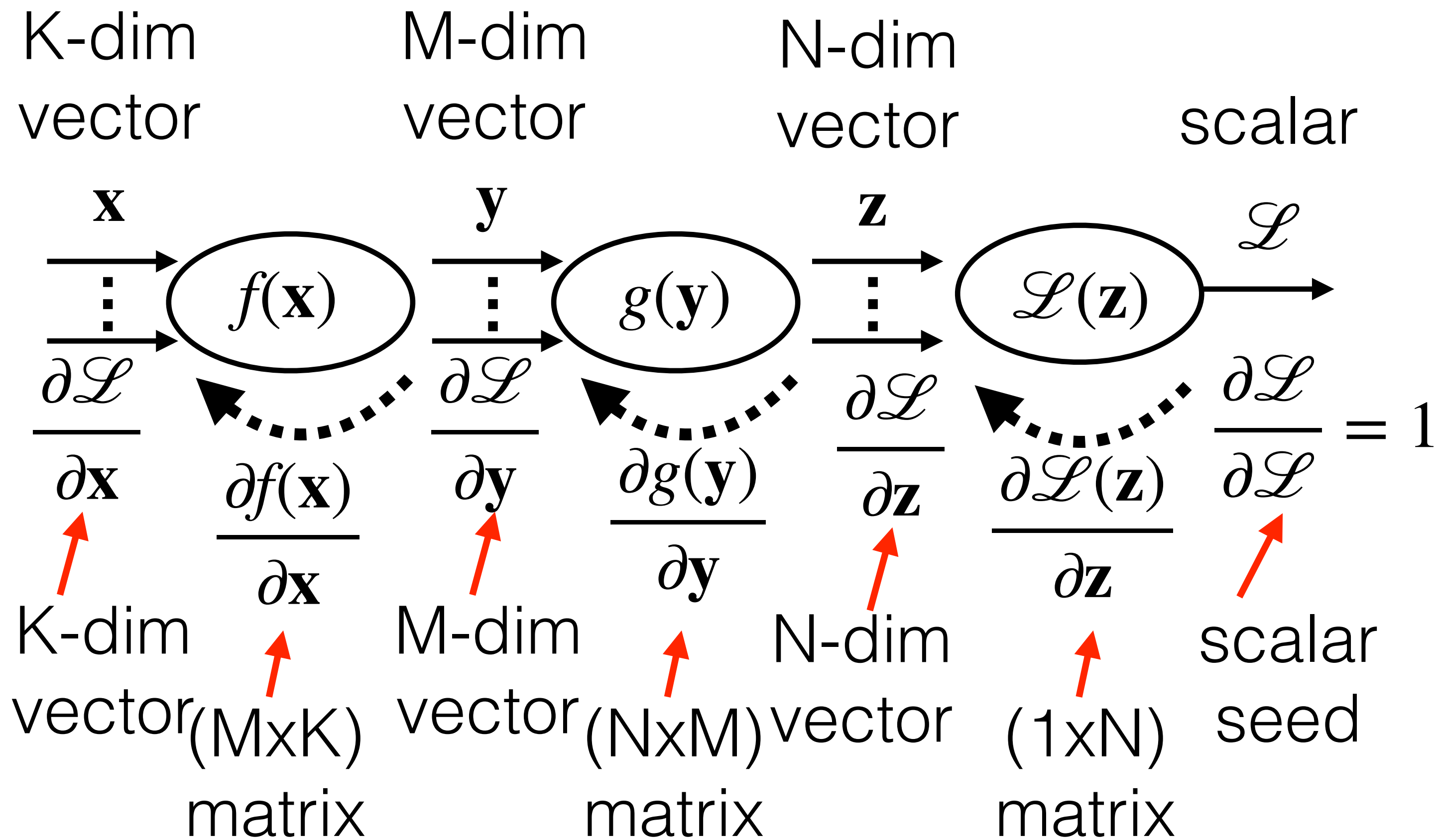
The main purpose is estimation of this



Forward differentiation (JVP): computes Jacobian of the loss one column at a time  
 Reverse differentiation (VJP): computes Jacobian of the loss one row at a time

**Which one would you choose (we want to compute 1x12 jacobian)?**

# Efficient implementation of autodiff?

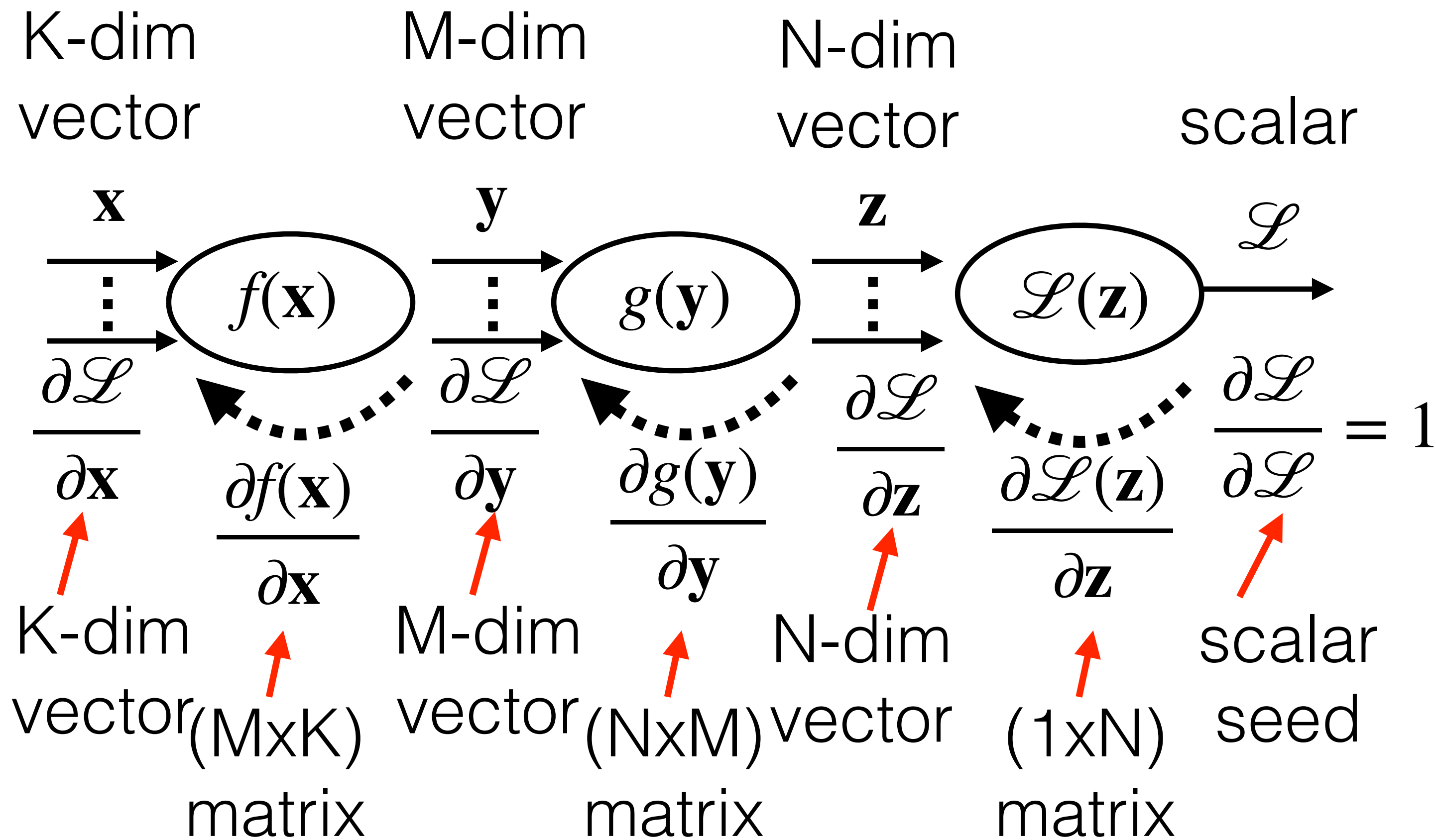


Loss jac wrt  $\mathbf{x}$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \begin{matrix} 1 \times 1 \\ 1 \end{matrix} \begin{matrix} 1 \times N \\ \frac{\partial L}{\partial \mathbf{z}} \end{matrix} \begin{matrix} N \times M \\ \frac{\partial g(\mathbf{y})}{\partial \mathbf{y}} \end{matrix} \begin{matrix} M \times K \\ \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \end{matrix}$$



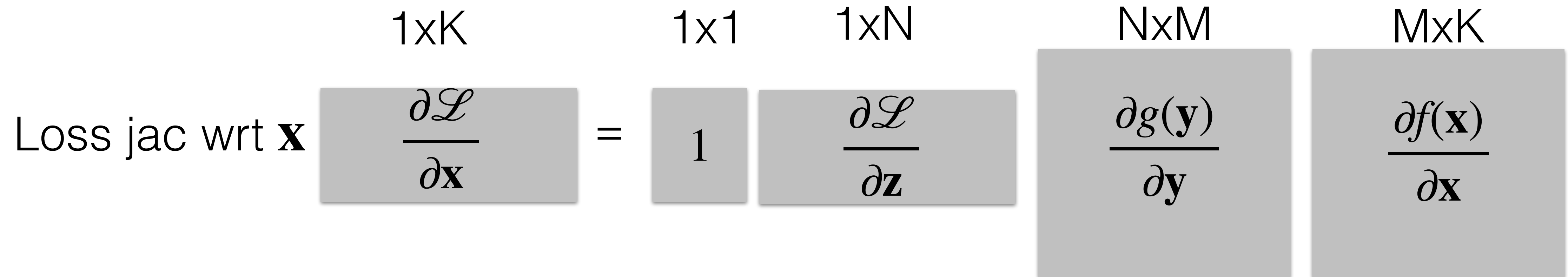
# Efficient implementation of autodiff?



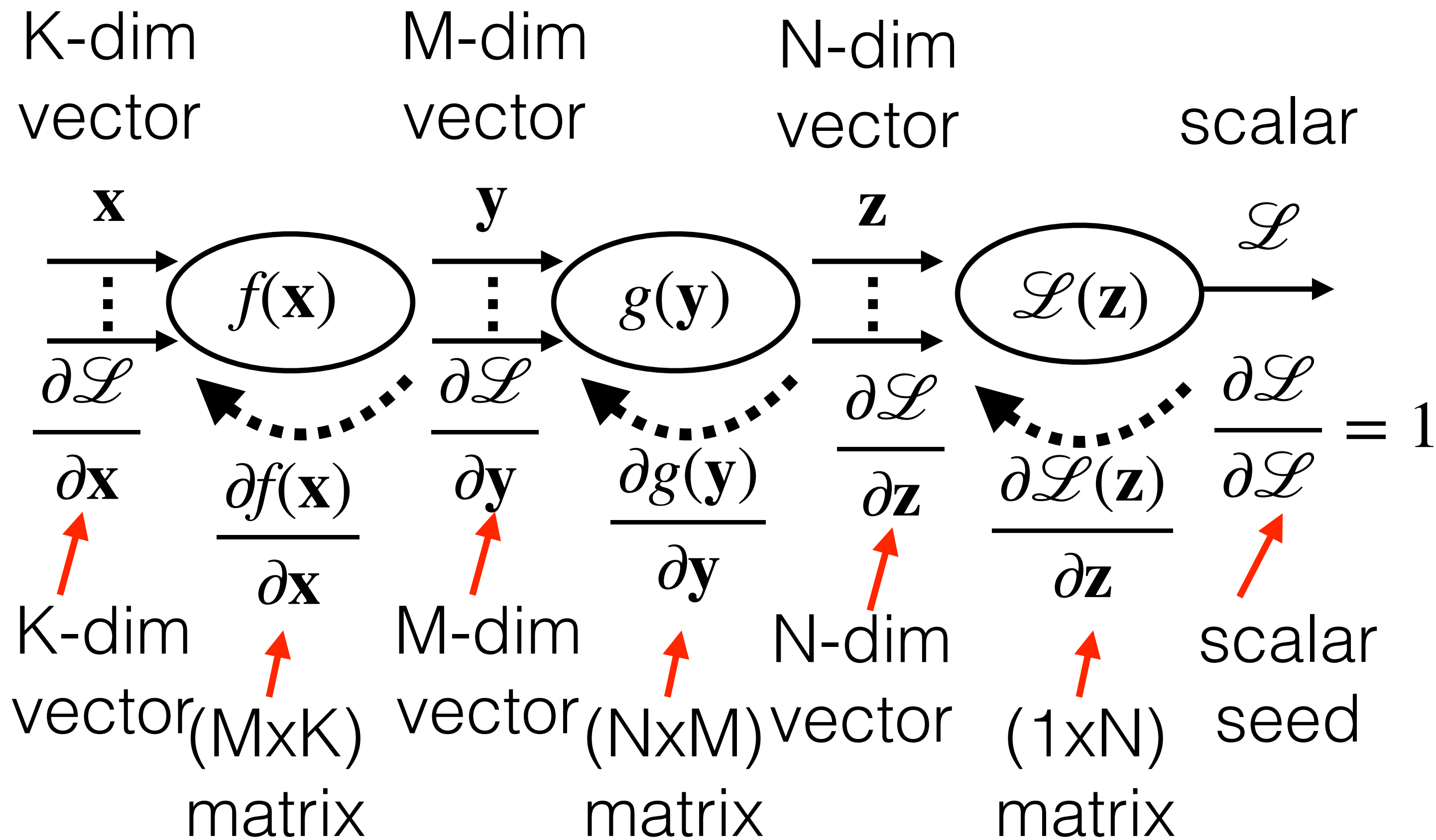
```
def vjp_L(v, z):
    return v.T * dL/dz
```

```
def vjp_g(v, y):
    return v.T * dg/dy
```

```
def vjp_f(v, x):
    return v.T * df/dx
```



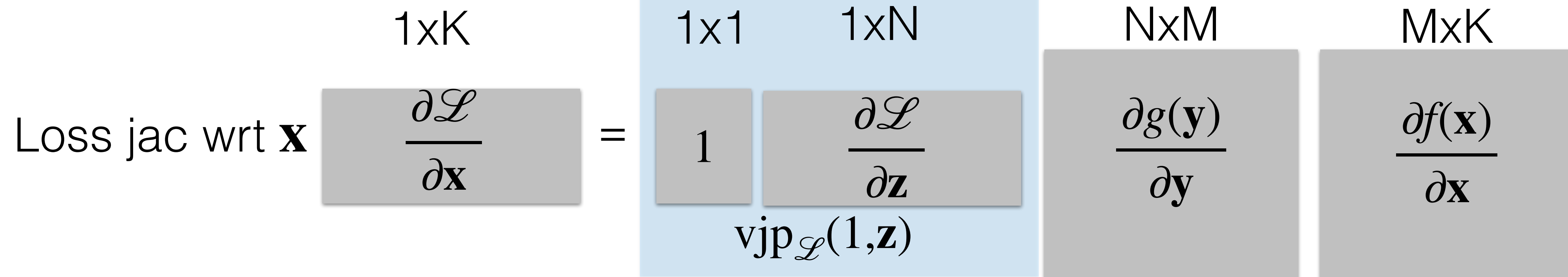
# Efficient implementation of autodiff?



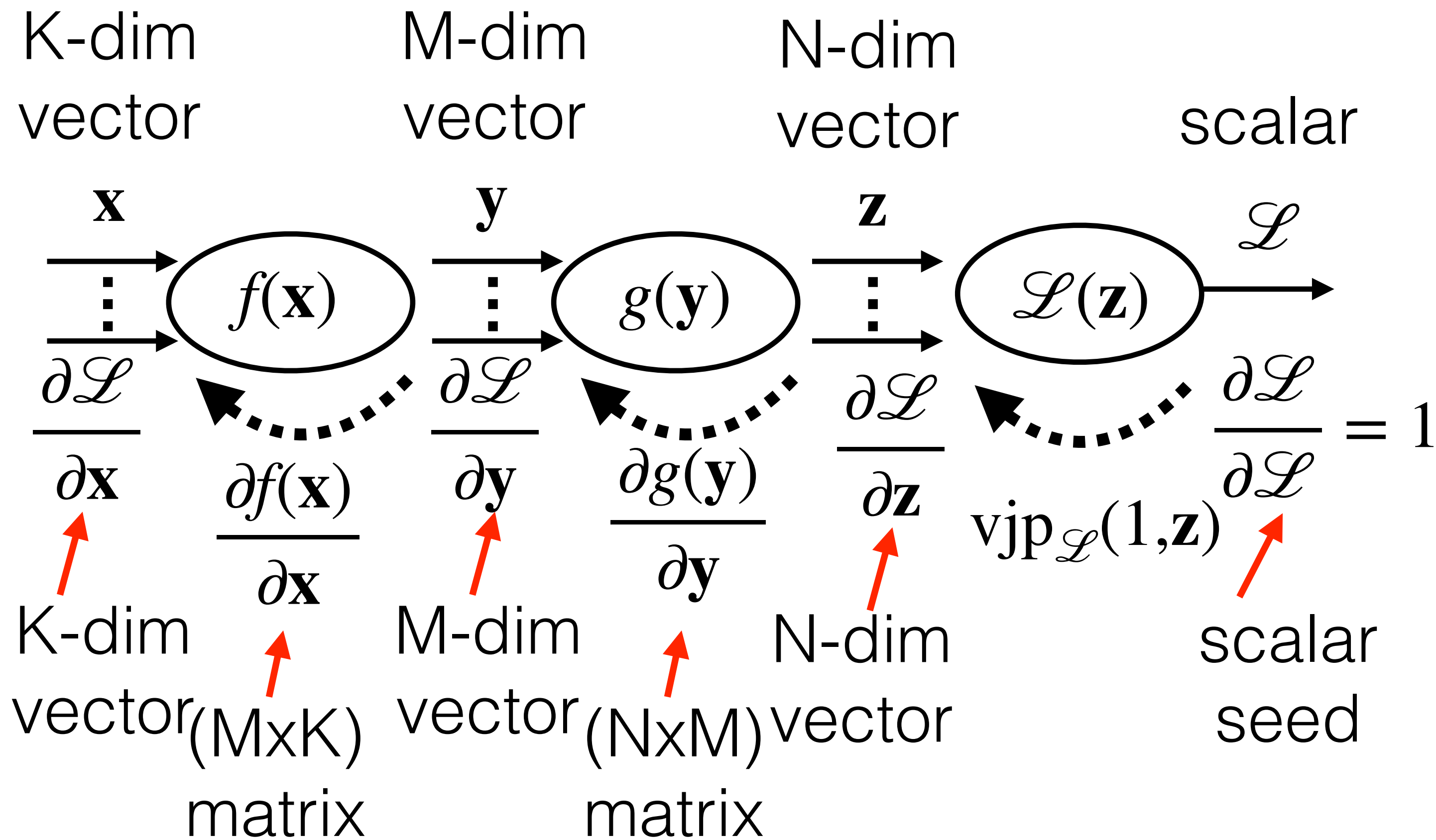
```
def vjp_L(v, z):
    return v.T * (dL/dz)
```

```
def vjp_g(v, y):
    return v.T * (dg/dy)
```

```
def vjp_f(v, x):
    return v.T * (df/dx)
```



# Efficient implementation of autodiff?



```
def vjp_L(v, z):
    return v.T * dL/dz
```

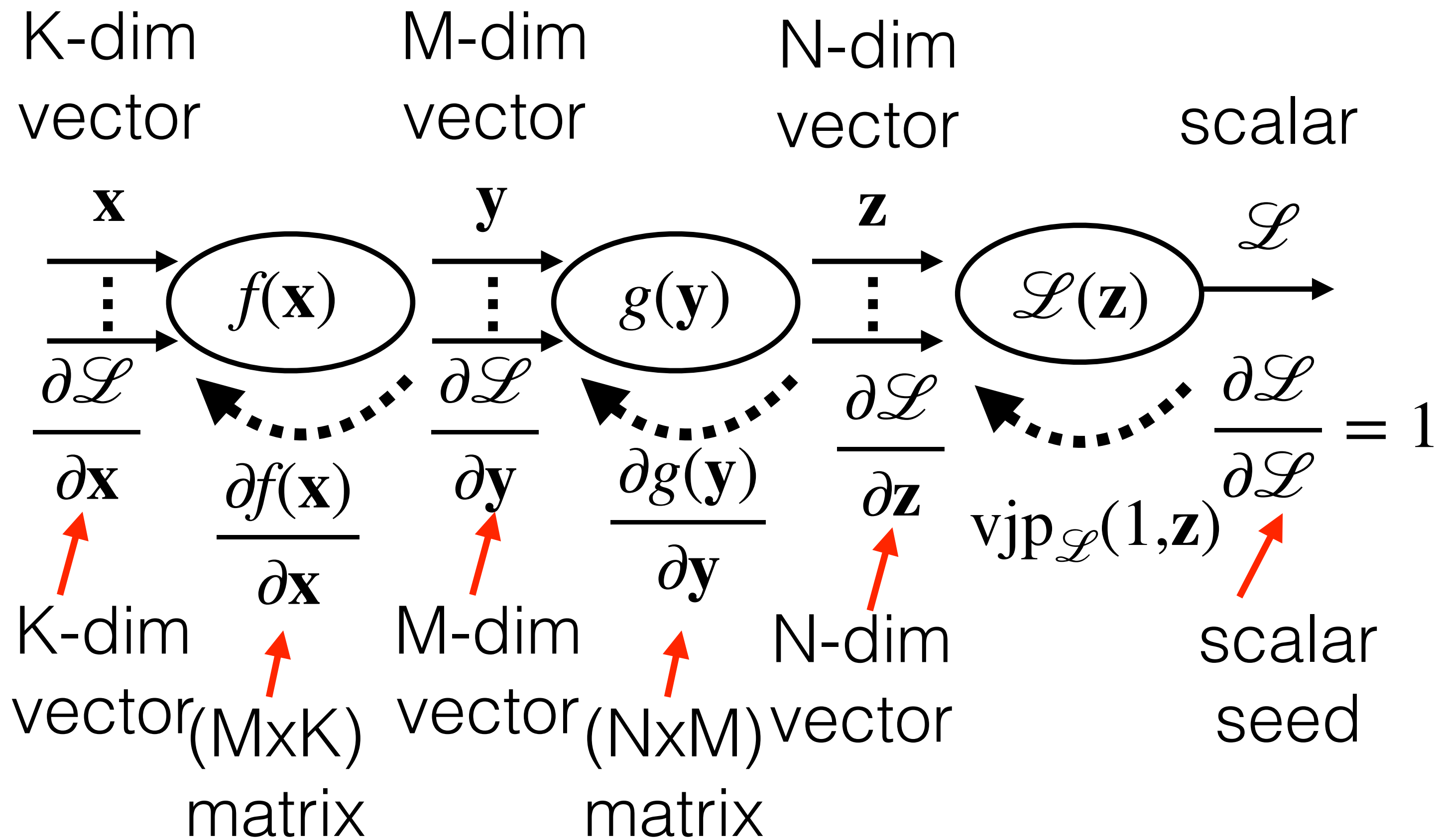
```
def vjp_g(v, y):
    return v.T * dg/dy
```

```
def vjp_f(v, x):
    return v.T * df/dx
```

Loss jac wrt  $\mathbf{x}$   $\frac{\partial \mathcal{L}}{\partial \mathbf{x}}$  =  $\text{vjp}_L(1, \mathbf{z})$   $\frac{\partial g(\mathbf{y})}{\partial \mathbf{y}}$   $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$

Dimensions:  $1 \times K$ ,  $N \times M$ ,  $M \times K$

# Efficient implementation of autodiff?



```
def vjp_L(v, z):
    return v.T * dL/dz
```

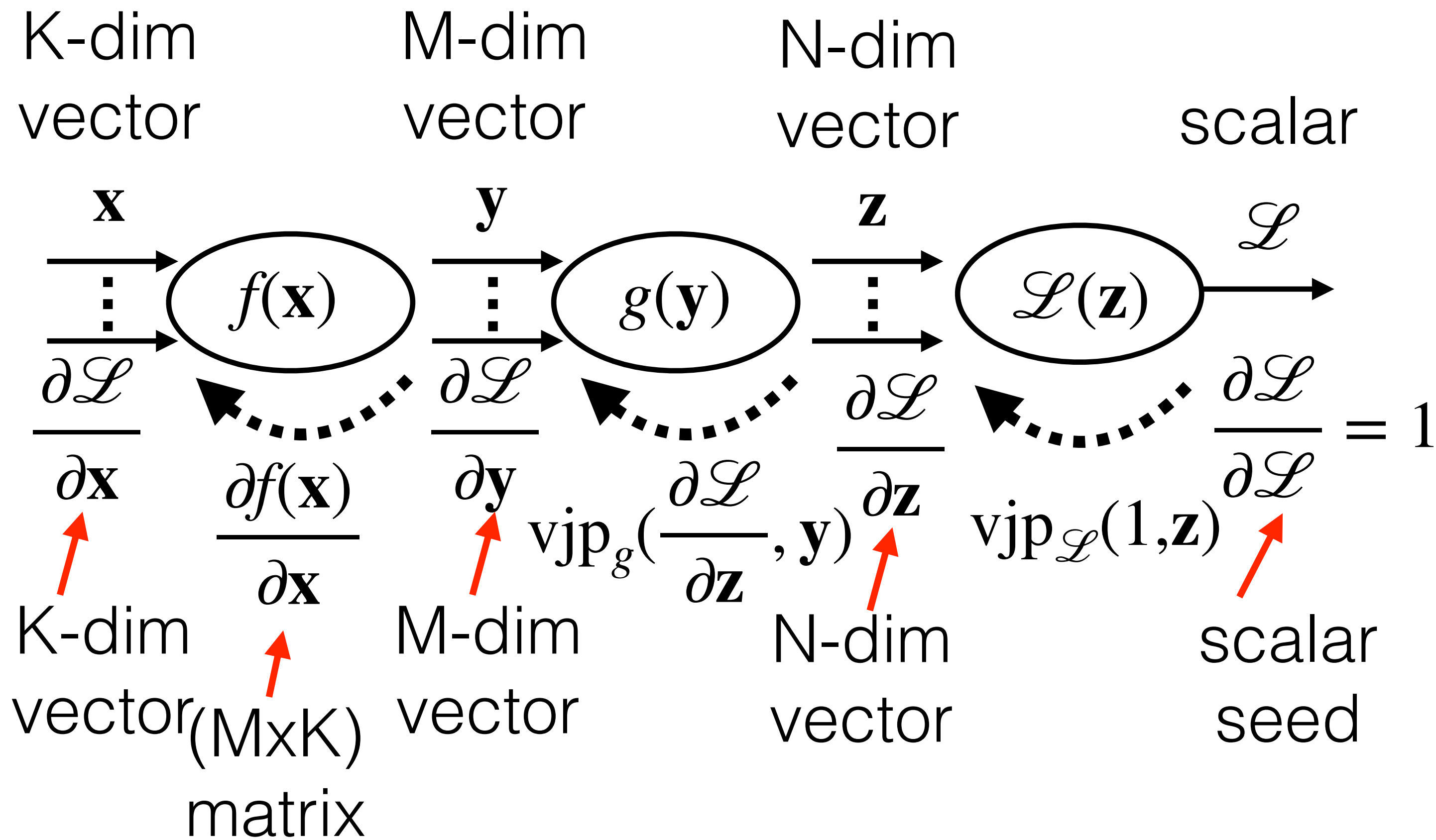
```
def vjp_g(v, y):
    return v.T * dg(y)/dy
```

```
def vjp_f(v, x):
    return v.T * df(x)/dx
```

Loss jac wrt  $\mathbf{x}$   $\frac{\partial \mathcal{L}}{\partial \mathbf{x}}$  =  $\text{vjp}_{\mathcal{L}}(1, \mathbf{z})$   $\frac{\partial g(\mathbf{y})}{\partial \mathbf{y}}$   $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$

Dimensions:  $1 \times K$ ,  $N \times M$ ,  $M \times K$

# Efficient implementation of autodiff?



```
def vjp_L(v, z):
    return v.T . dL/dz
```

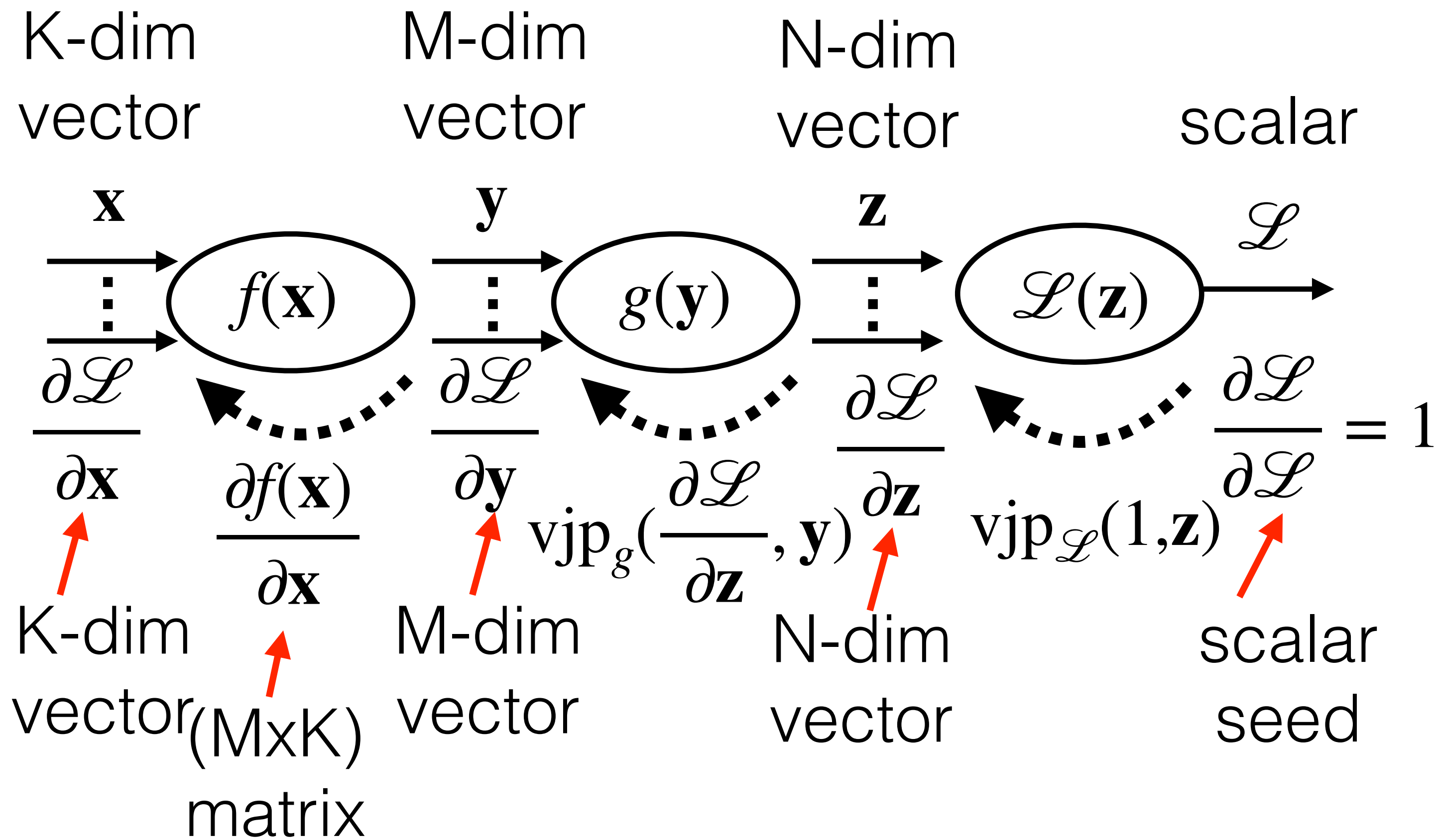
```
def vjp_g(v, y):
    return v.T . dg/dy
```

```
def vjp_f(v, x):
    return v.T . df/dx
```

Loss jac wrt  $\mathbf{x}$   $\frac{\partial \mathcal{L}}{\partial \mathbf{x}}$  (1xK) =  $\text{vjp}_g(\text{vjp}_L(1, \mathbf{z}), \mathbf{y})$  (MxK)  $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$  (MxK)



# Efficient implementation of autodiff?

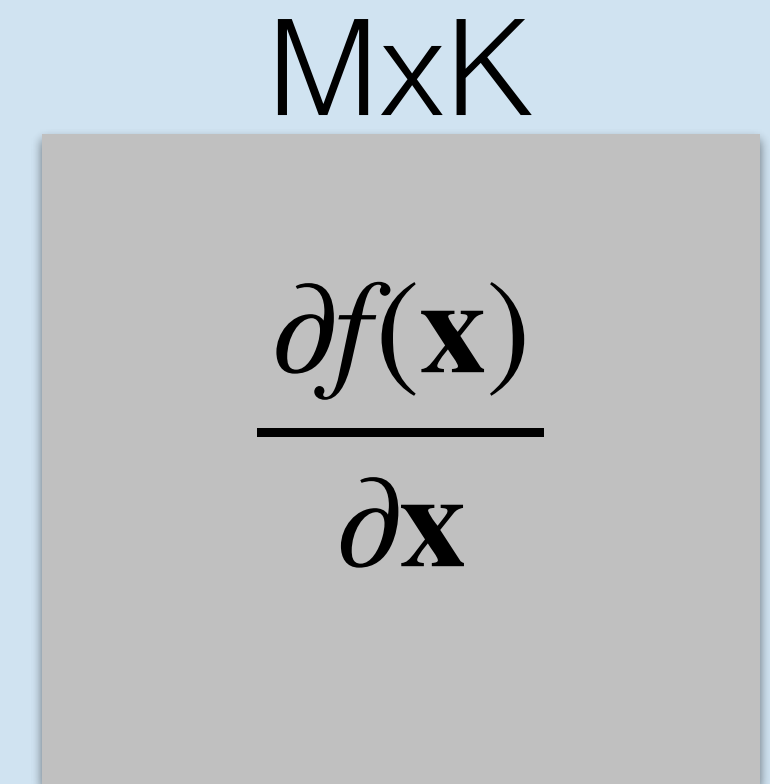


```
def vjp_L(v, z):
    return v.T * dL/dz
```

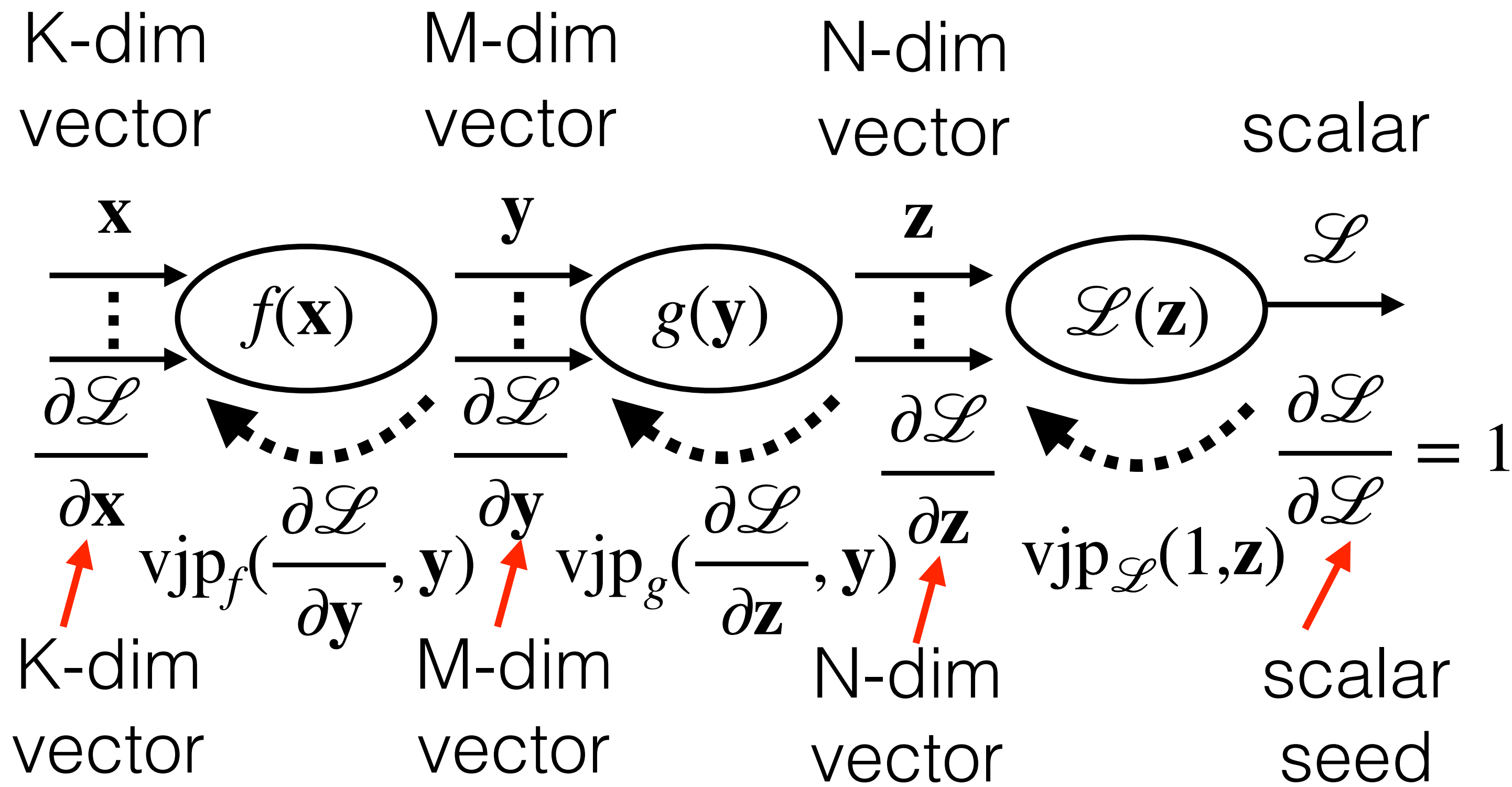
```
def vjp_g(v, y):
    return v.T * dg(y)/dy
```

```
def vjp_f(v, x):
    return v.T * df(x)/dx
```

Loss jac wrt  $\mathbf{x}$   $\frac{\partial \mathcal{L}}{\partial \mathbf{x}}$   $=$   $\text{vjp}_g(\text{vjp}_L(1, \mathbf{z}), \mathbf{y})$



# Efficient implementation of autodiff?



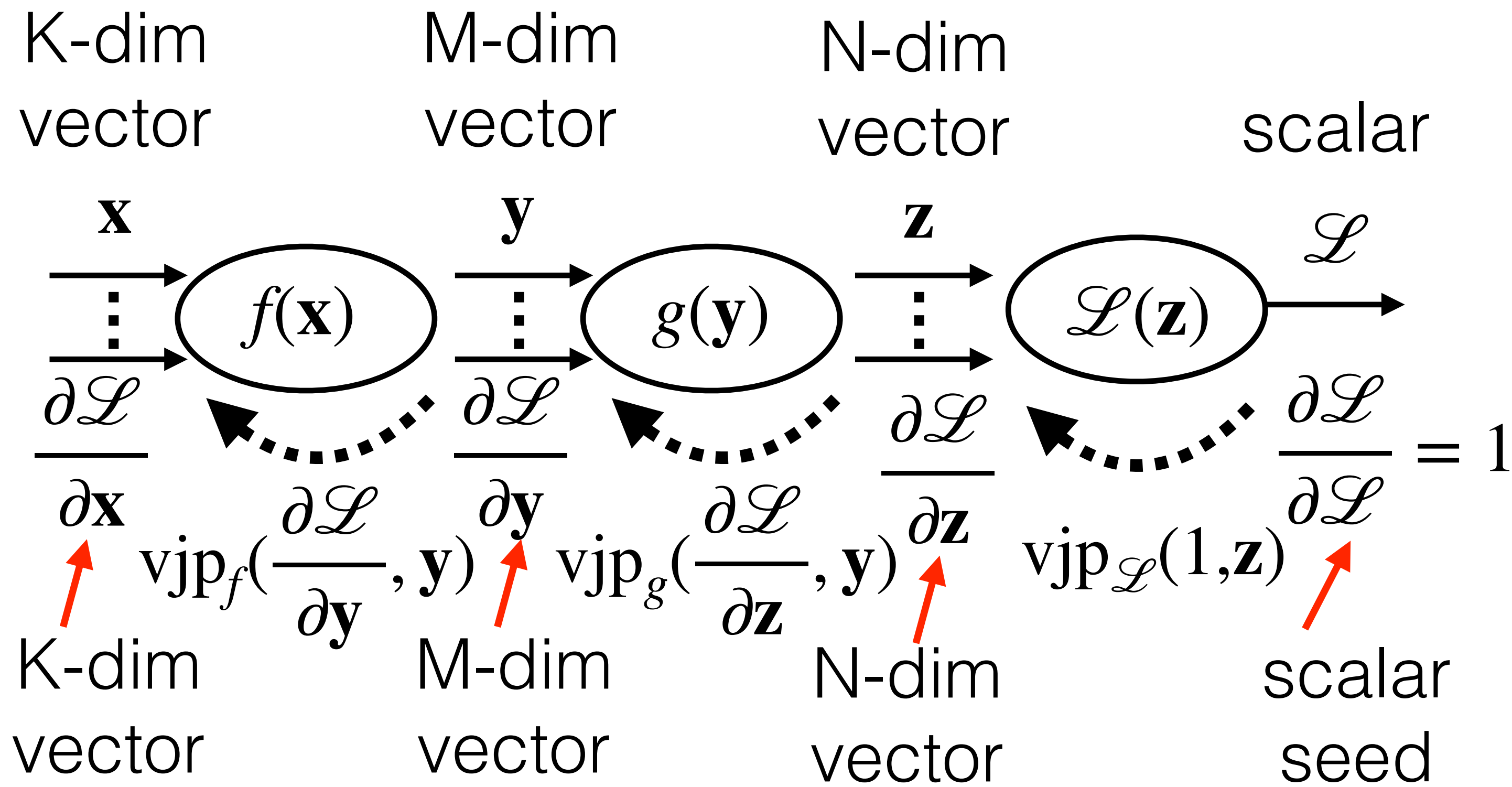
```
def vjp_L(v, z):
    return v.T . dL/dz
```

```
def vjp_g(v, y):
    return v.T . dg/dy
```

```
def vjp_f(v, x):
    return v.T . df/dx
```

Loss jac wrt  $\mathbf{x}$   $\frac{\partial \mathcal{L}}{\partial \mathbf{x}}$   $=$   $\text{vjp}_f(\text{vjp}_g(\text{vjp}_{\mathcal{L}}(1, \mathbf{z}), \mathbf{y}), \mathbf{x})$

# Efficient implementation of autodiff?



```
def vjp_L(v, z):
    return v.T * dL/dz
```

```
def vjp_g(v, y):
    return v.T * dg/dy
```

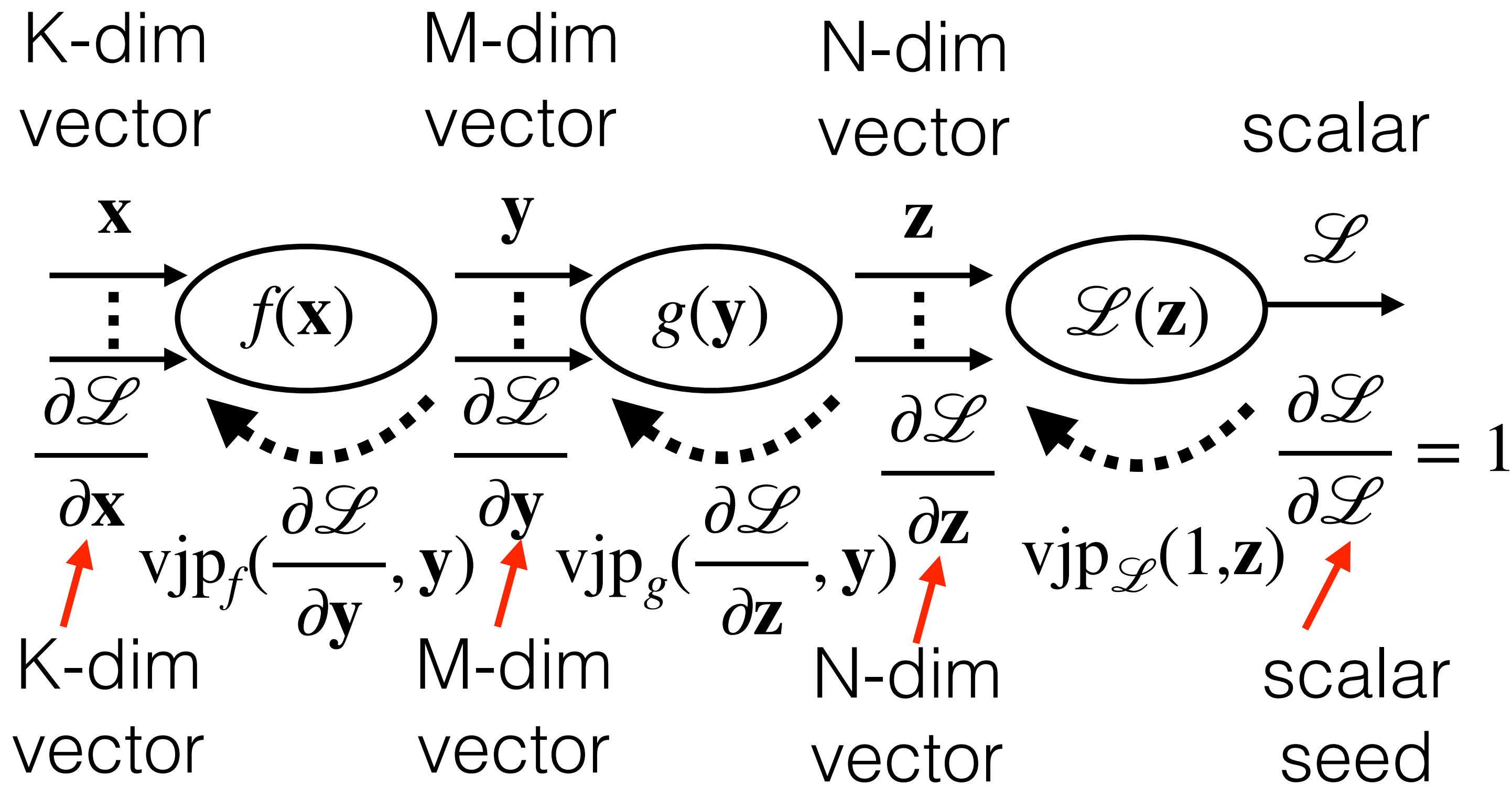
```
def vjp_f(v, x):
    return v.T * df/dx
```

## Reverse autodiff algorithm

input:	forward pass:	backward pass:	output:
$\mathbf{x}$	$\mathbf{y} = f(\mathbf{x})$	$\mathbf{v} = \text{vjp}_{\mathcal{L}}(1, \mathbf{z})$	$\mathbf{v} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}}$
	$\mathbf{z} = g(\mathbf{y})$	$\mathbf{v} = \text{vjp}_g(\mathbf{v}, \mathbf{y})$	
	$\mathcal{L} = \mathcal{L}(\mathbf{z})$	$\mathbf{v} = \text{vjp}_f(\mathbf{v}, \mathbf{x})$	



# Efficient implementation of autodiff?



```
def vjp_L(v, z):
    return v.T * dL/dz
```

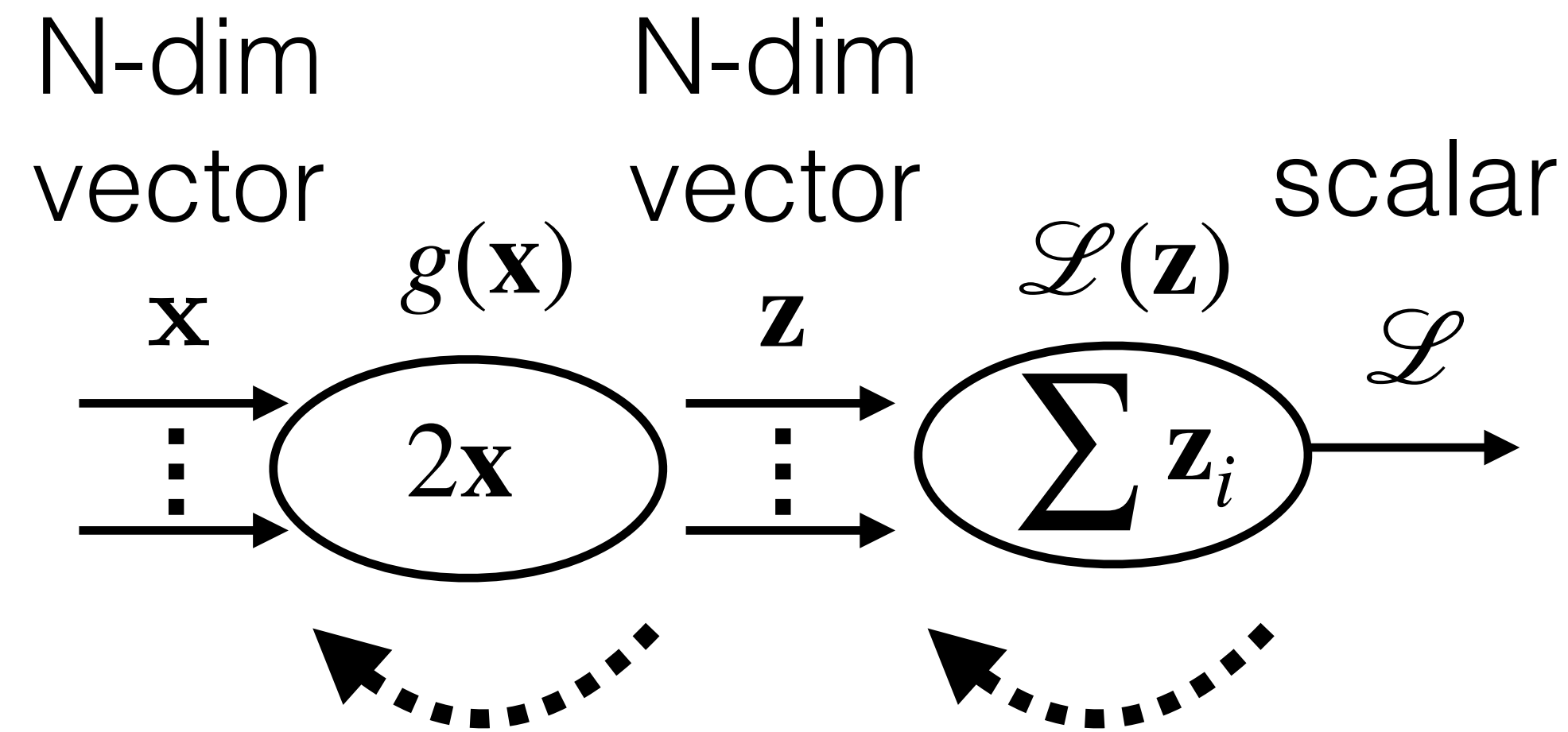
```
def vjp_g(v, y):
    return v.T * dg(y)/dy
```

```
def vjp_f(v, x):
    return v.T * df(x)/dx
```

Why the hell should I implement it in such a way?

- vjp usually implemented more efficiently (building the jacobian not required)
- preserves dimensionality of inputs in backward pass

- vjp usually implemented more efficiently (building the jacobian not required)



$$\frac{\partial g(\mathbf{y})}{\partial \mathbf{y}} = \begin{matrix} 2 & & & & \\ & 2 & & & \\ & & 2 & & \\ & & & \dots & \\ & & & & 2 \end{matrix}$$

NXN

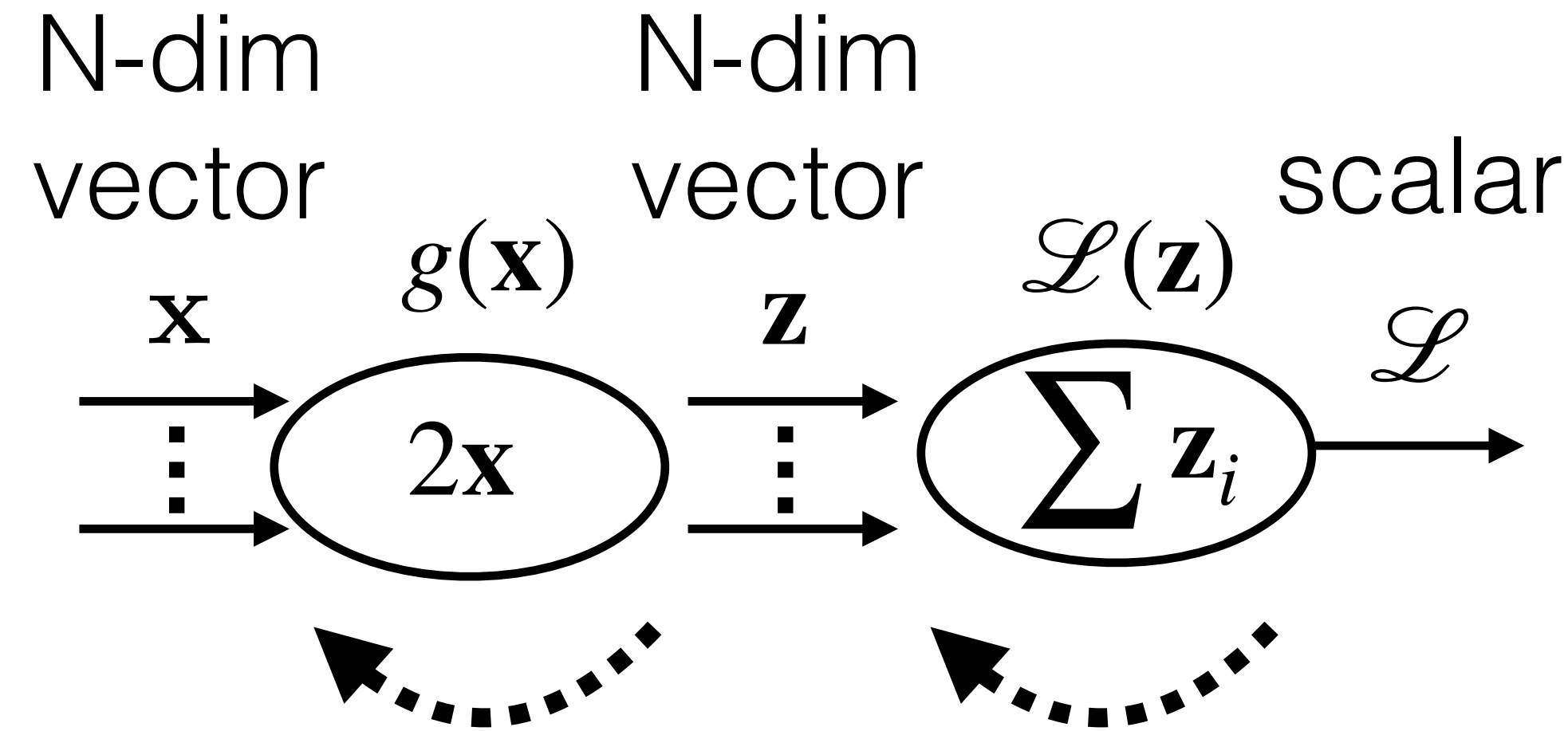
$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}} = \begin{matrix} 1, 1, 1, \dots, 1 \end{matrix}$$

1xN

```
def vjp_L(v, z):
    return v.T * (∂L(z) / ∂z)
```

```
def vjp_g(v, x):
    return v.T * (∂g(y) / ∂y)
```

- vjp usually implemented more efficiently (building the jacobian not required)



$$\frac{\partial g(\mathbf{y})}{\partial \mathbf{y}} = \begin{matrix} 2 \\ & 2 \\ & & 2 \\ & & & \dots \\ & & & & 2 \end{matrix}$$

NXN

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}} = \begin{matrix} 1, 1, 1, \dots, 1 \end{matrix}$$

1xN

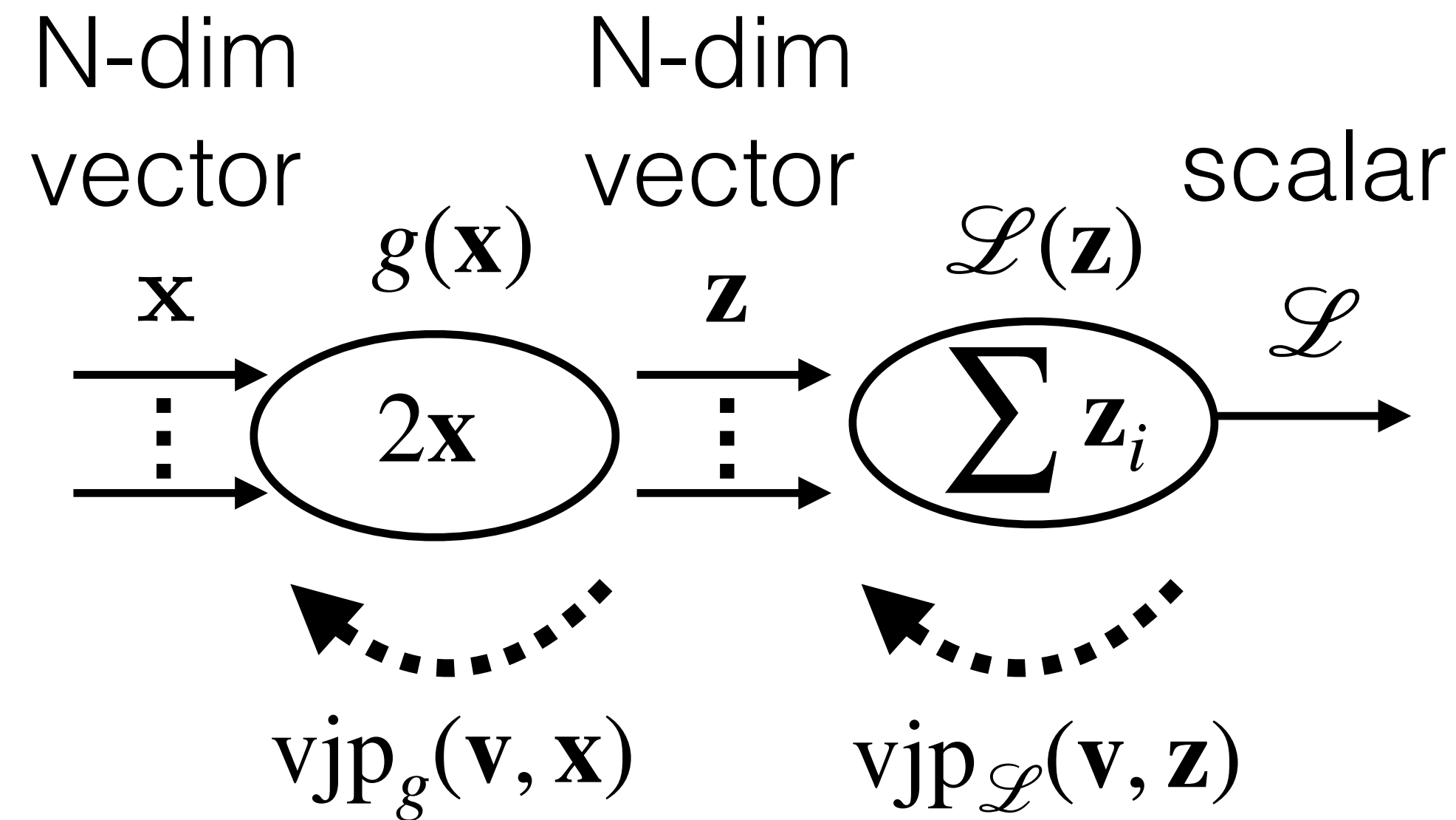
```
def vjp_L(v, z):
    return v.T . [1, 1, 1, ... 1]
```

```
def vjp_g(v, x):
    return v.T . \begin{matrix} 2 \\ & 2 \\ & & 2 \\ & & & \dots \\ & & & & 2 \end{matrix}
```

**Do I really need to construct the jacobian???**



- vjp usually implemented more efficiently (building the jacobian not required)



```
def vjp $\mathcal{L}$ (v, z):
    return tile(v, 1, N)
```

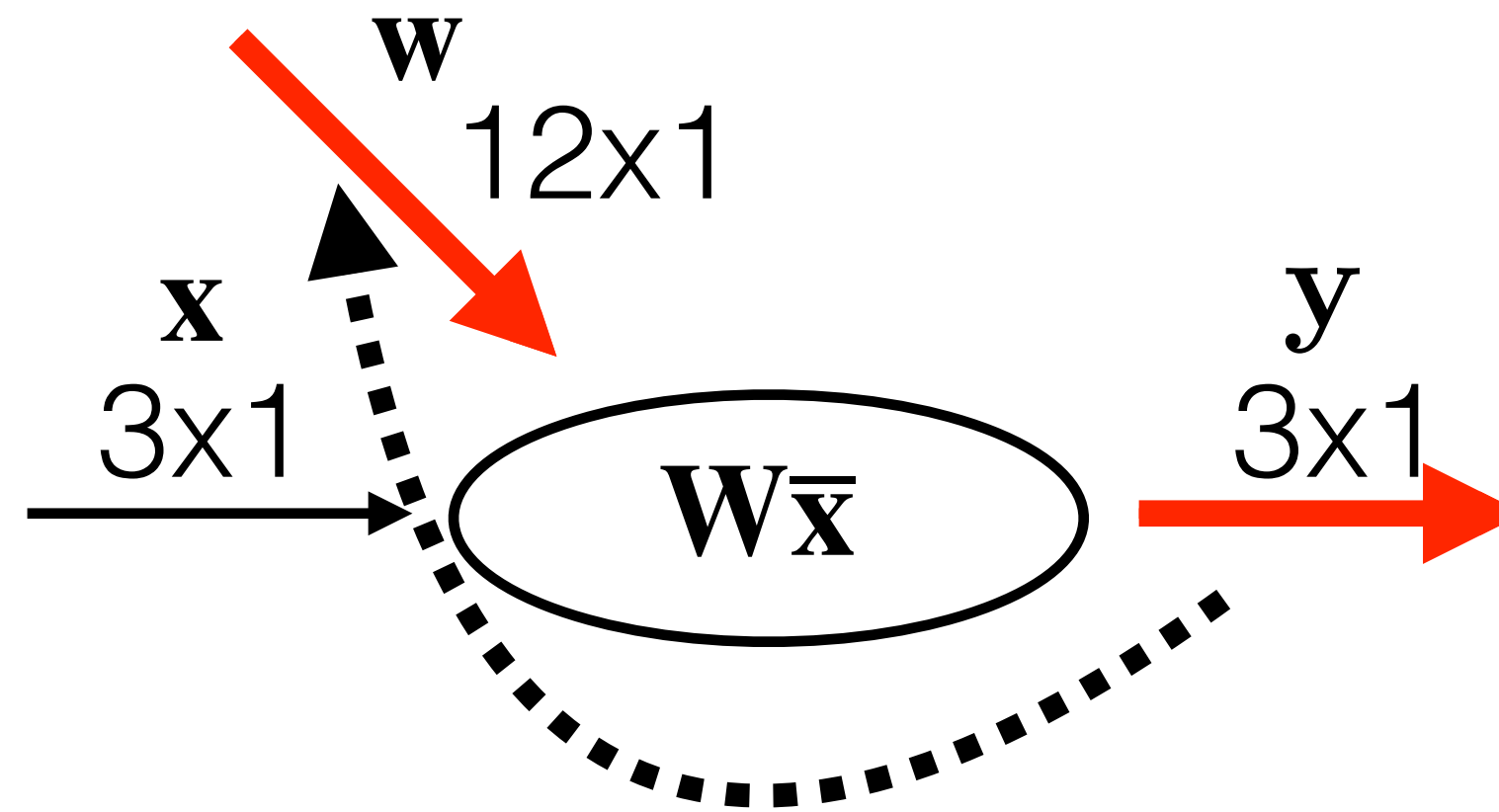
```
def vjp $g$ (v, x):
    return vT · 2
```

vjp is automatically added to tensors during the construction of comp.g.

```
grad_fn=<fBackward>
```

```
x = torch.tensor([[1, 2], [3, 4]], dtype=torch.float32, requires_grad=True)
z = x.sum()
print(z)
tensor(10, grad_fn=<SumBackward0>)
```

- vjp usually implemented more efficiently (building the jacobian not required)



```
def vjp_w(v, (x, w)) :
```

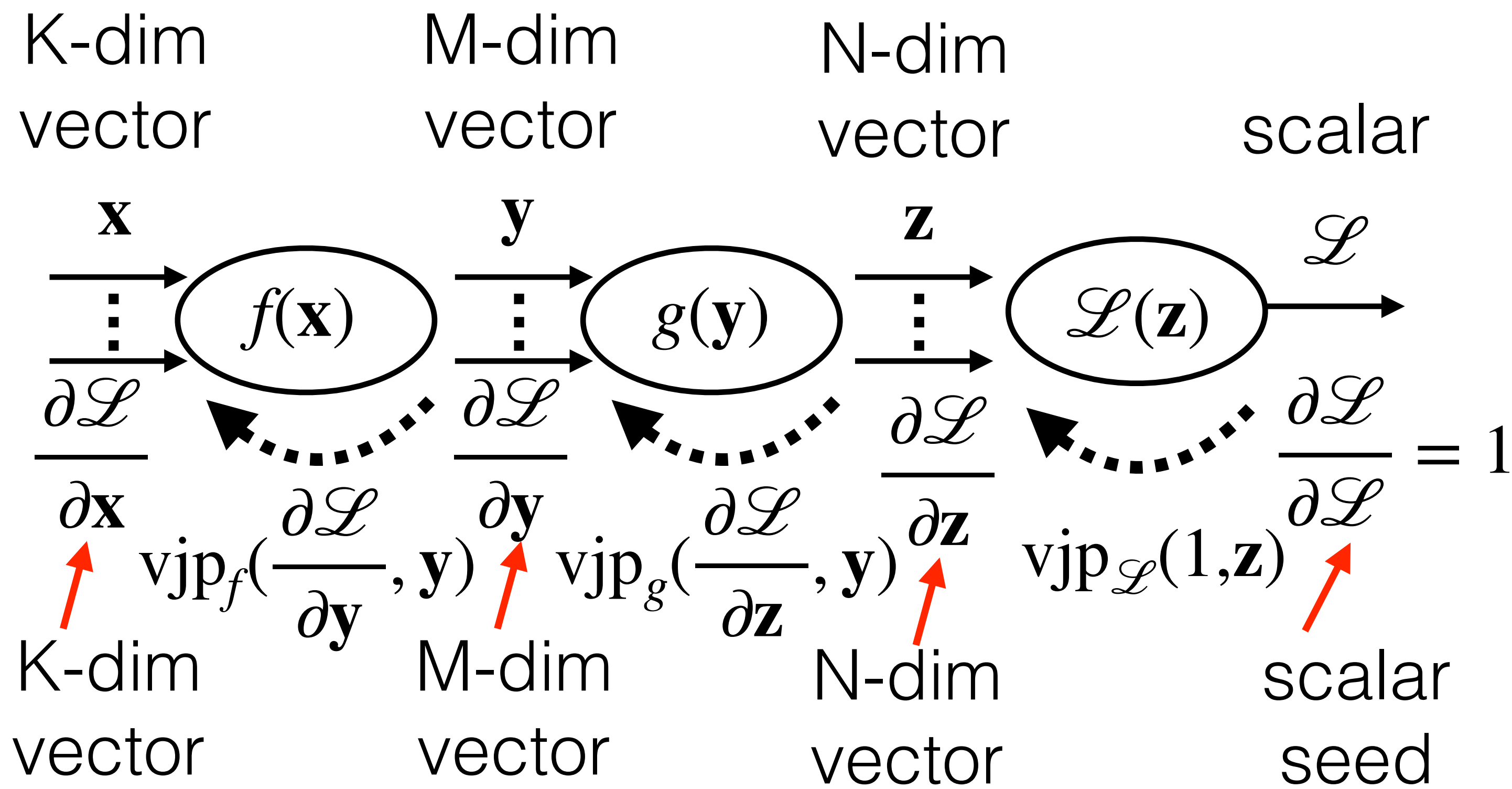
$$\text{return } \mathbf{v} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{w}} = \mathbf{v} \cdot \begin{bmatrix} -\bar{\mathbf{x}}^T - & & \\ & -\bar{\mathbf{x}}^T - & \\ & & -\bar{\mathbf{x}}^T - \end{bmatrix} = [v_1 \cdot \bar{\mathbf{x}}^T, v_2 \cdot \bar{\mathbf{x}}^T, v_3 \cdot \bar{\mathbf{x}}^T]$$

1x3

3x12

1x12

# Efficient implementation of autodiff?



```
def vjp_L(v, z):
    return v *  $\frac{\partial \mathcal{L}(\mathbf{z})}{\partial \mathbf{z}}$ 
```

```
def vjp_g(v, y):
    return v *  $\frac{\partial g(\mathbf{y})}{\partial \mathbf{y}}$ 
```

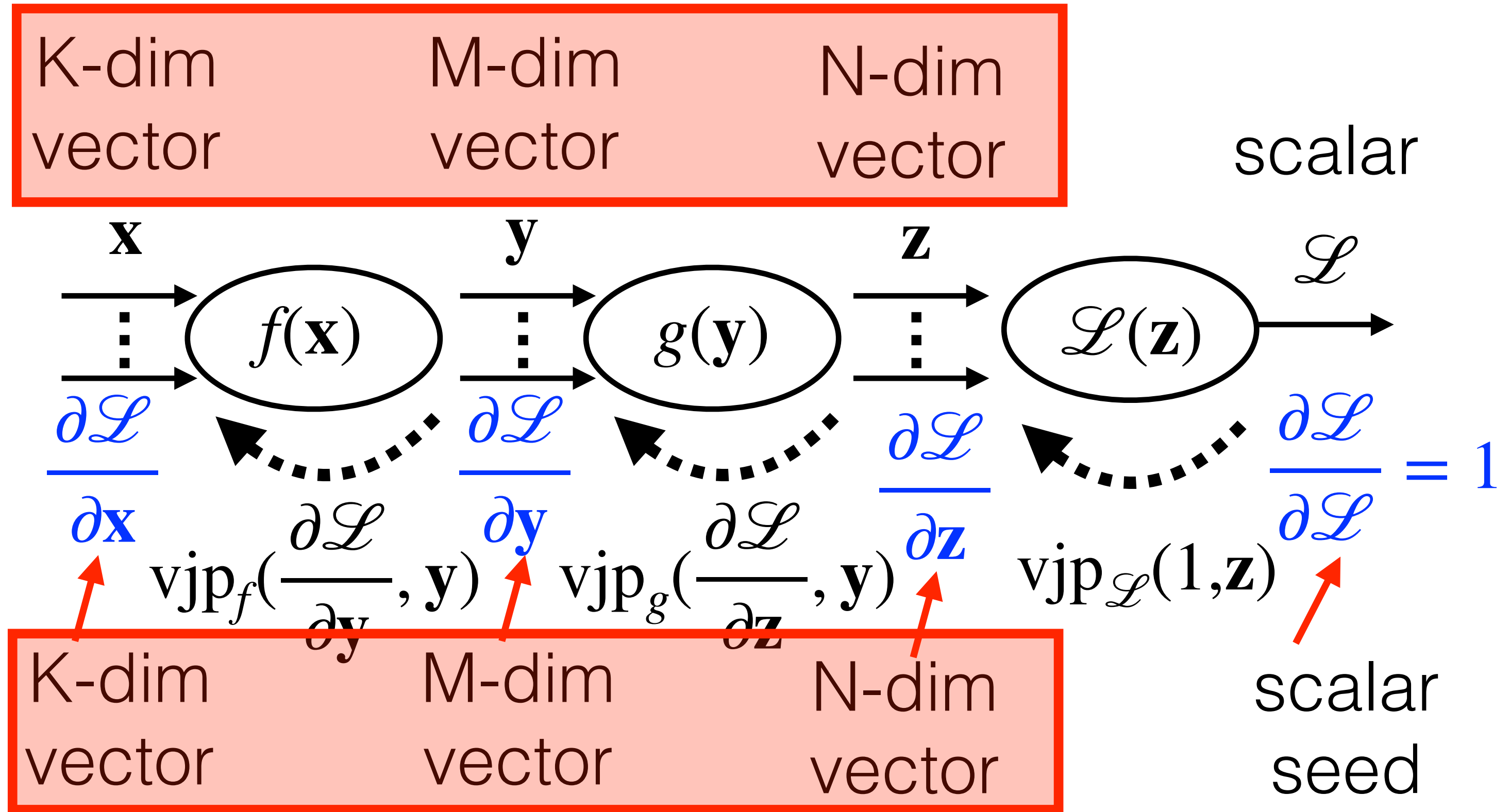
```
def vjp_f(v, x):
    return v *  $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$ 
```

Why the hell should I implement it in such a way?

- vjp usually implemented more efficiently (building the jacobian not required)
- preserves dimensionality of inputs in backward pass



# Efficient implementation of autodiff?



```
def vjp_L(v, z):
    return v *  $\frac{\partial \mathcal{L}(\mathbf{z})}{\partial \mathbf{z}}$ 
```

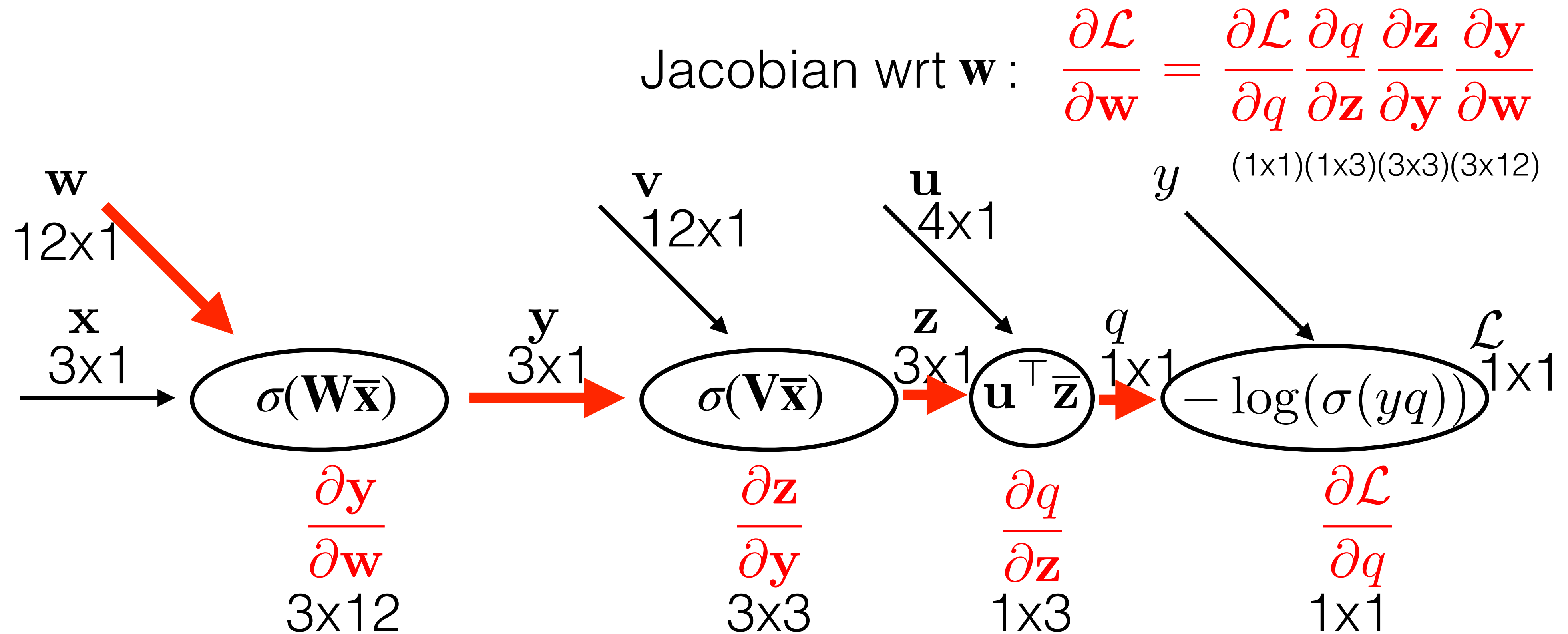
```
def vjp_g(v, y):
    return v *  $\frac{\partial g(\mathbf{y})}{\partial \mathbf{y}}$ 
```

```
def vjp_f(v, x):
    return v *  $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$ 
```

Why the hell should I implement it in such a way?

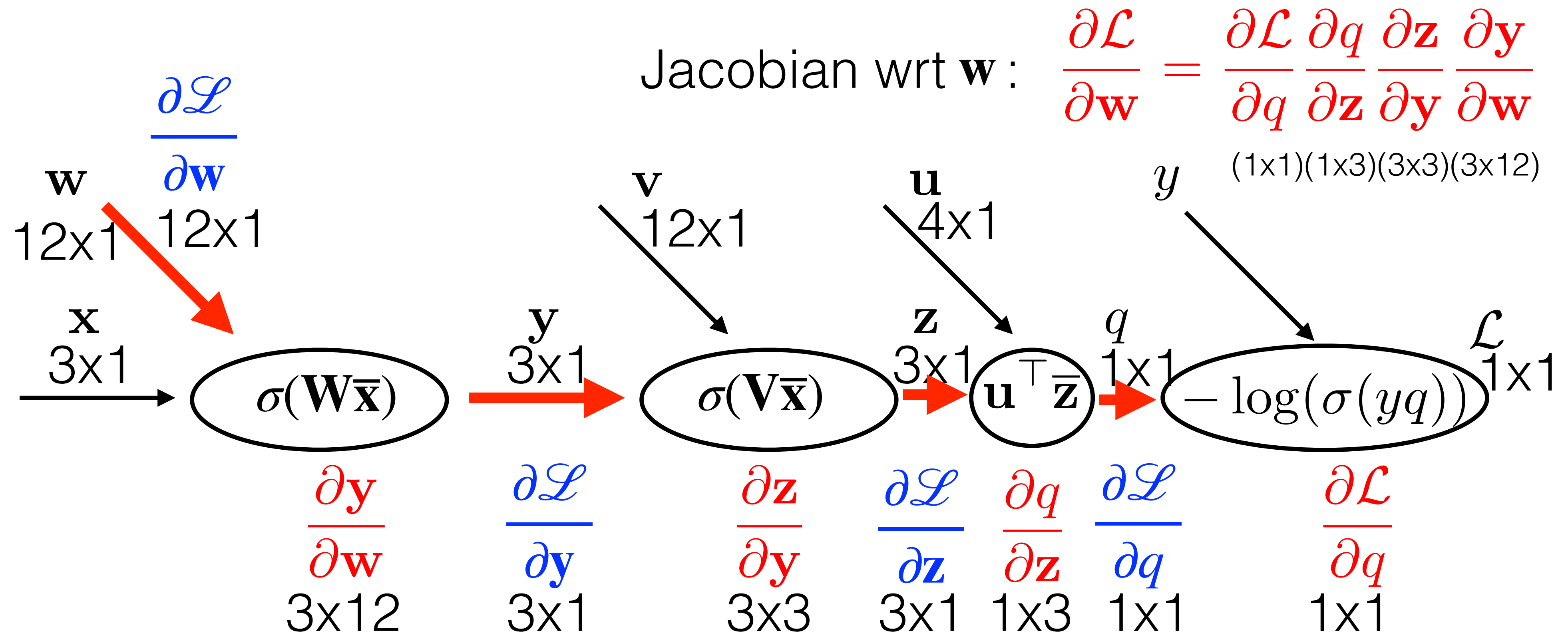
- vjp usually implemented more efficiently (building the jacobian not required)
- preserves dimensionality of inputs in backward pass

- preserves dimensionality of inputs in backward pass

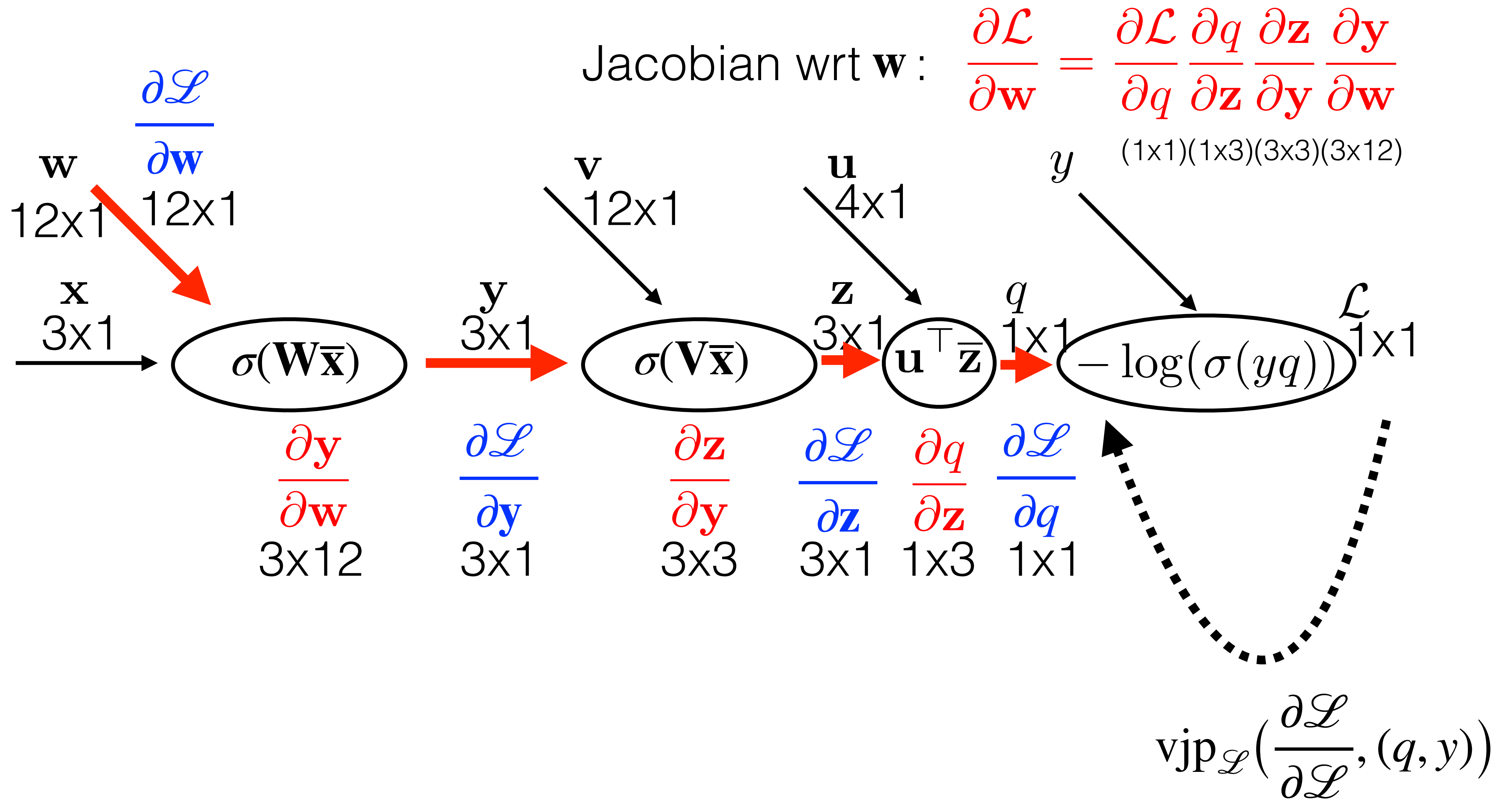


We avoided jacobian of multi-dimensional inputs by explicit vectorization

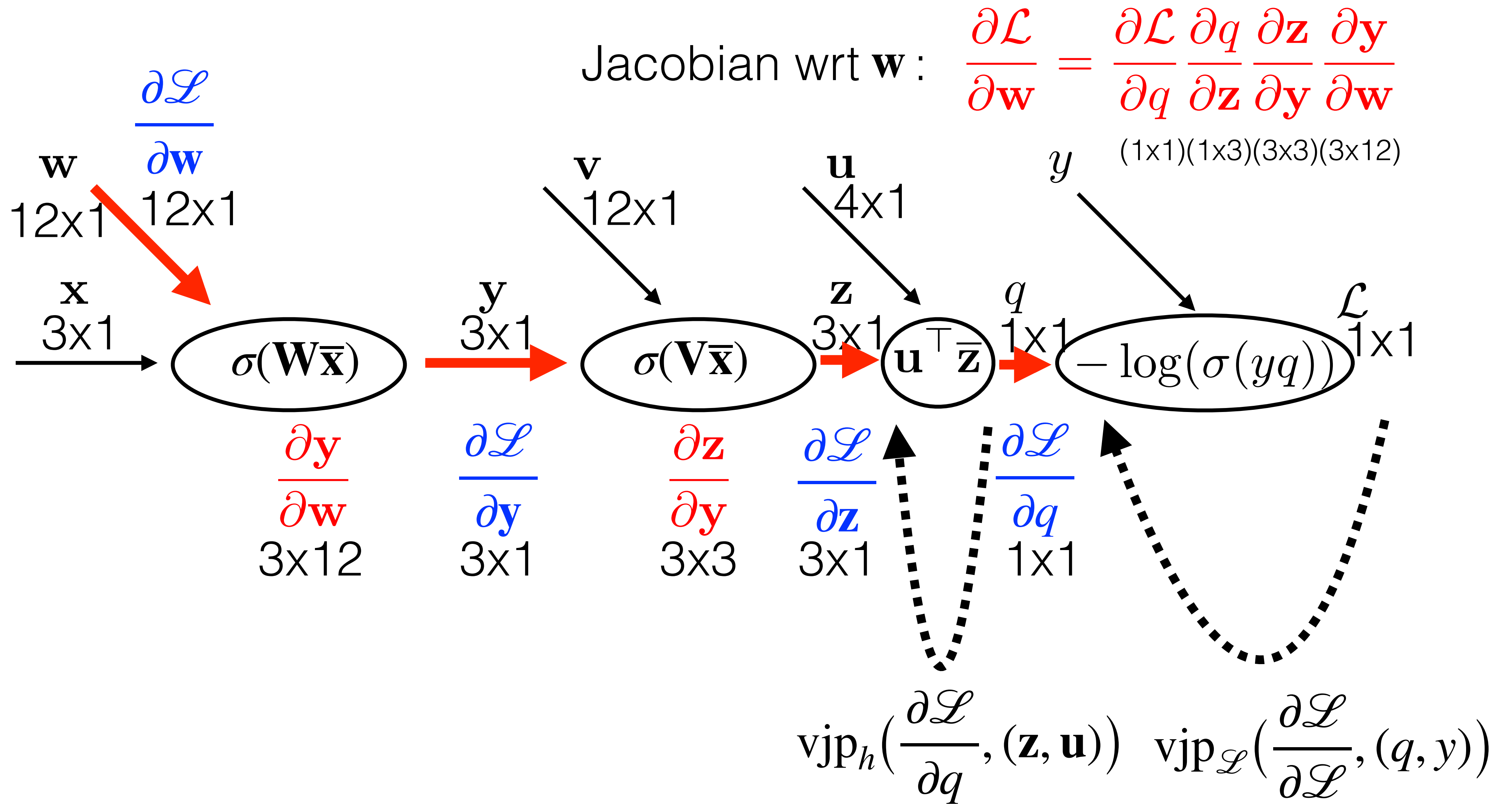
- preserves dimensionality of inputs in backward pass



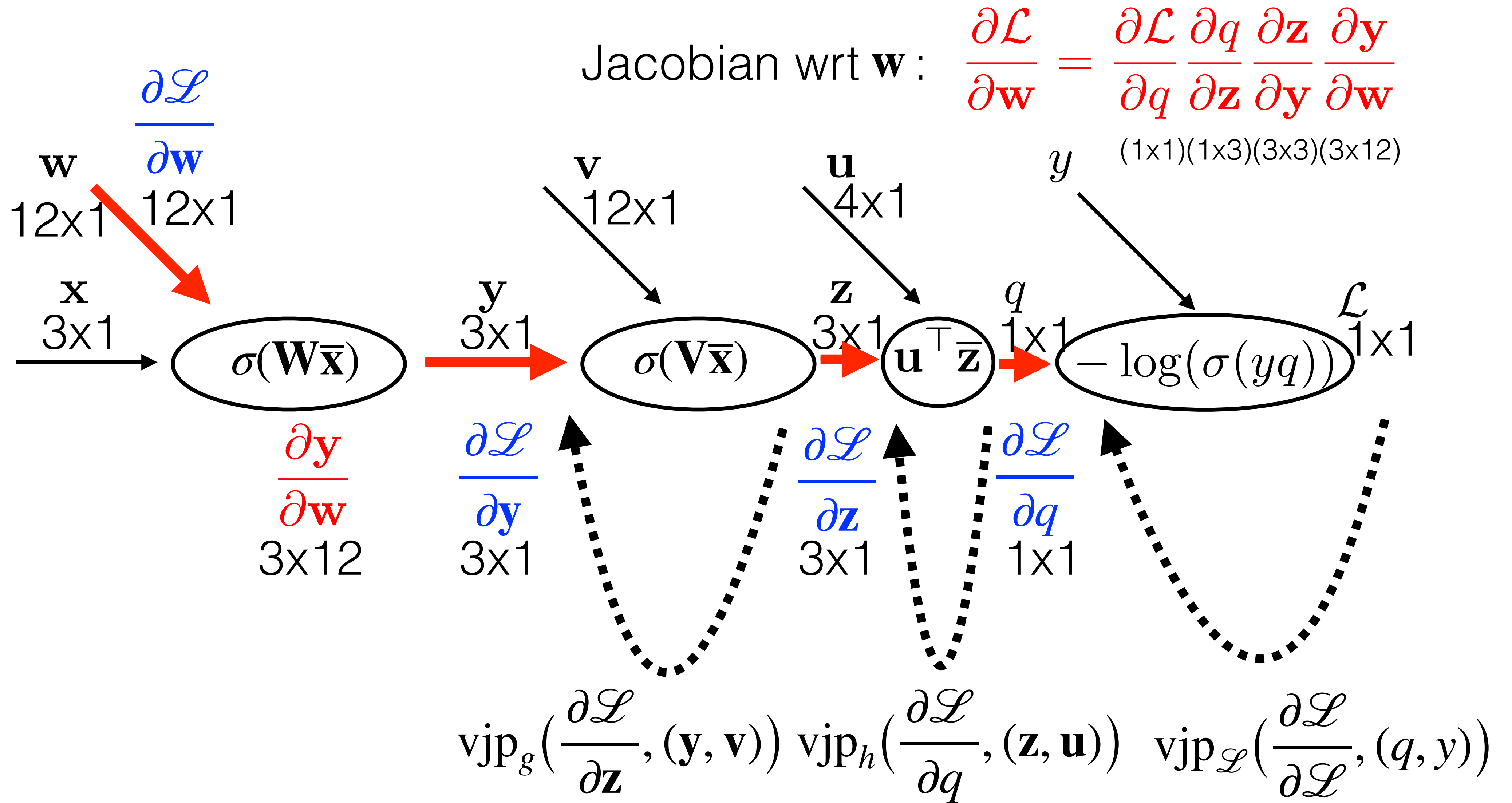
- preserves dimensionality of inputs in backward pass



- preserves dimensionality of inputs in backward pass

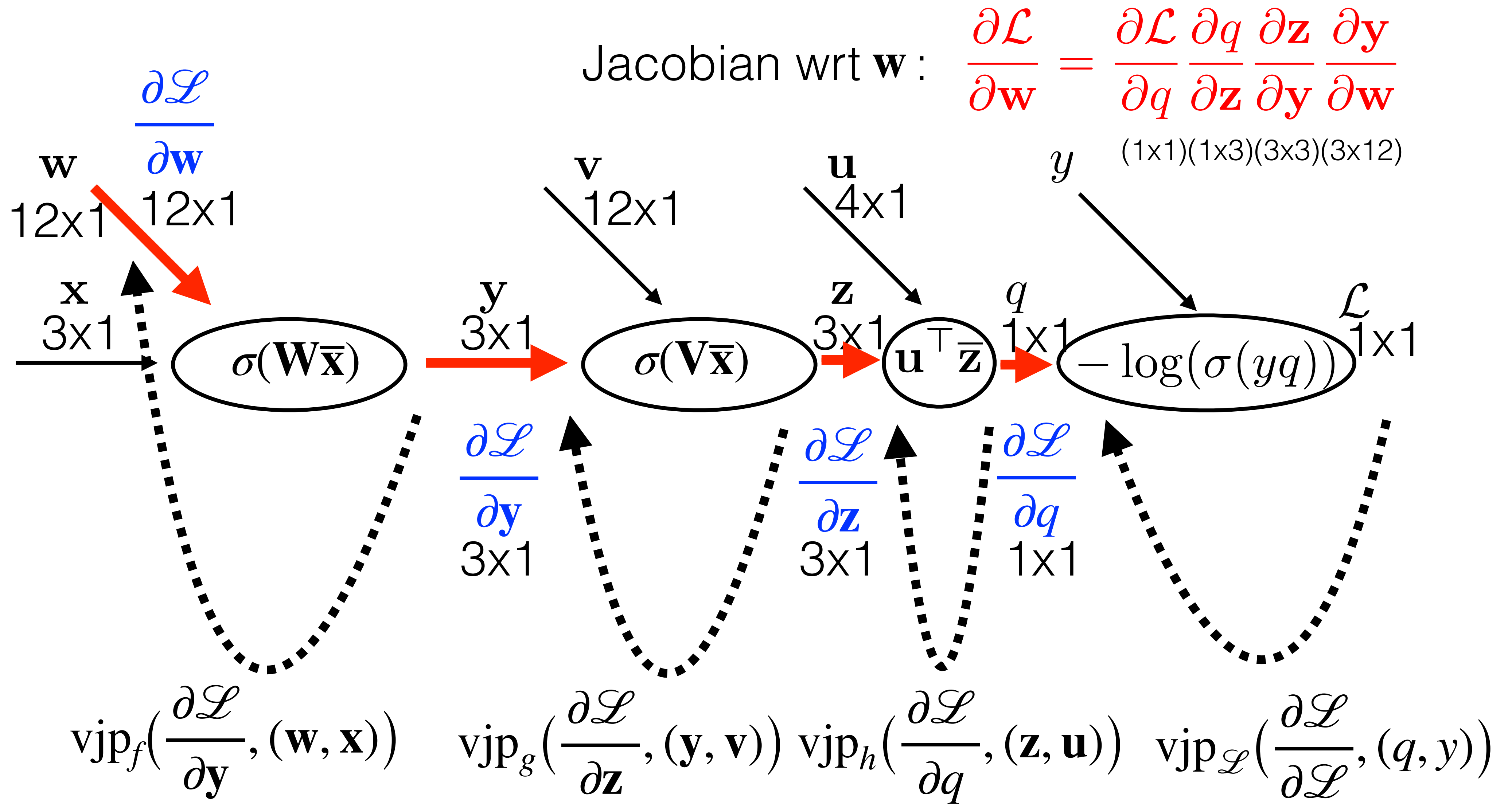


- preserves dimensionality of inputs in backward pass





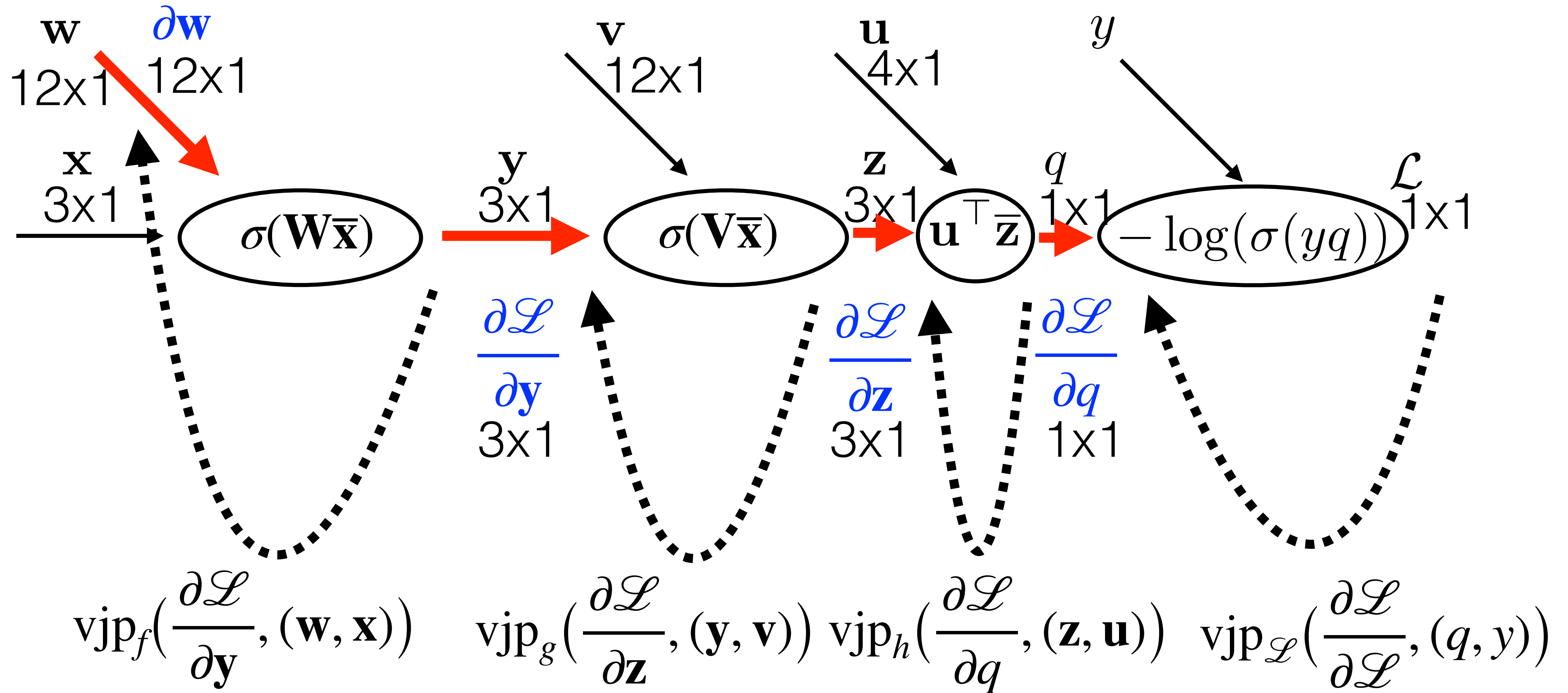
- preserves dimensionality of inputs in backward pass





- preserves dimensionality of inputs in backward pass

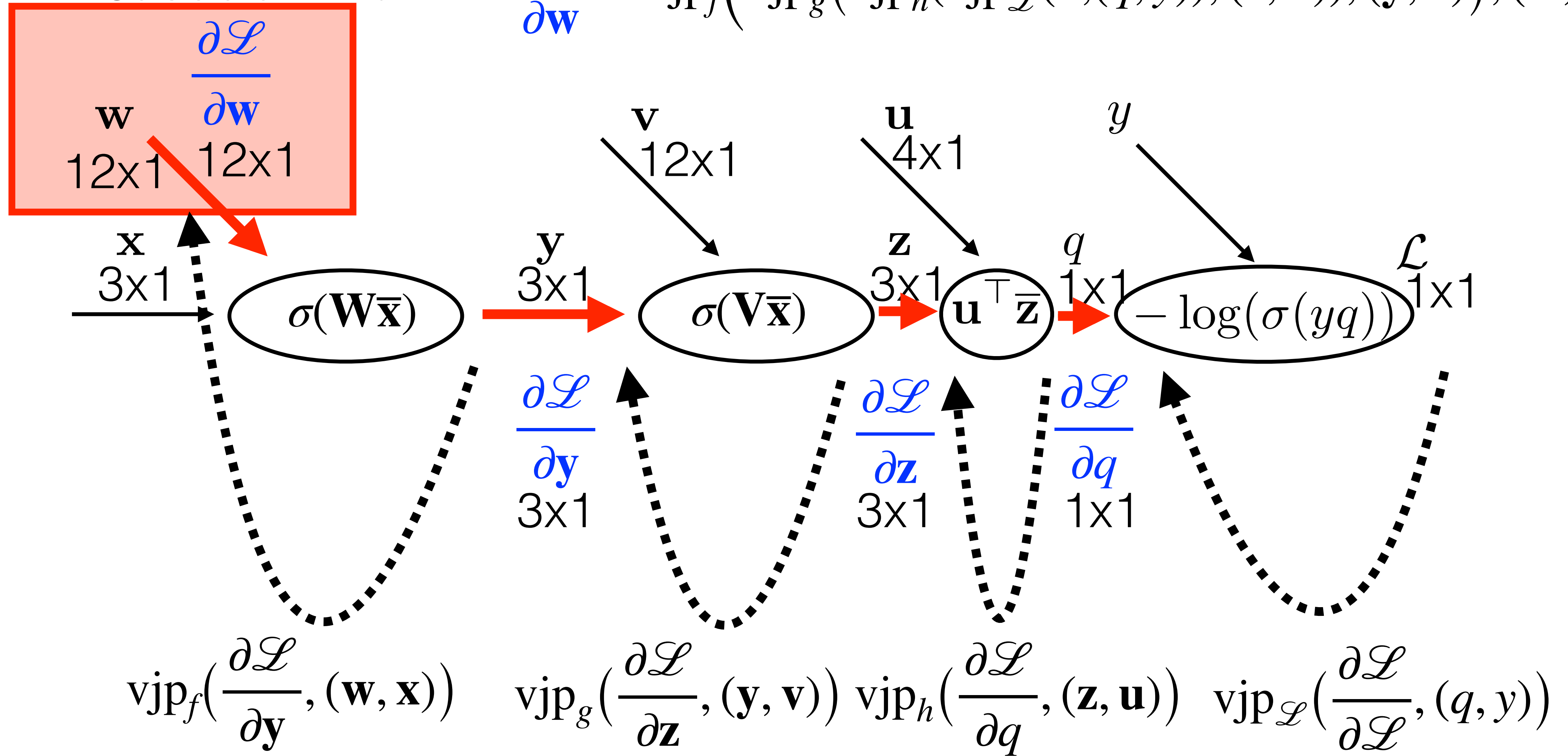
Jacobian wrt  $\mathbf{w}$  :  $\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \text{vjp}_f\left(\text{vjp}_g\left(\text{vjp}_h\left(\text{vjp}_{\mathcal{L}}(1, (q, y)), (\mathbf{z}, \mathbf{u})\right), (\mathbf{y}, \mathbf{v})\right), (\mathbf{w}, \mathbf{x})\right)$



There are no more jacobians

- preserves dimensionality of inputs in backward pass

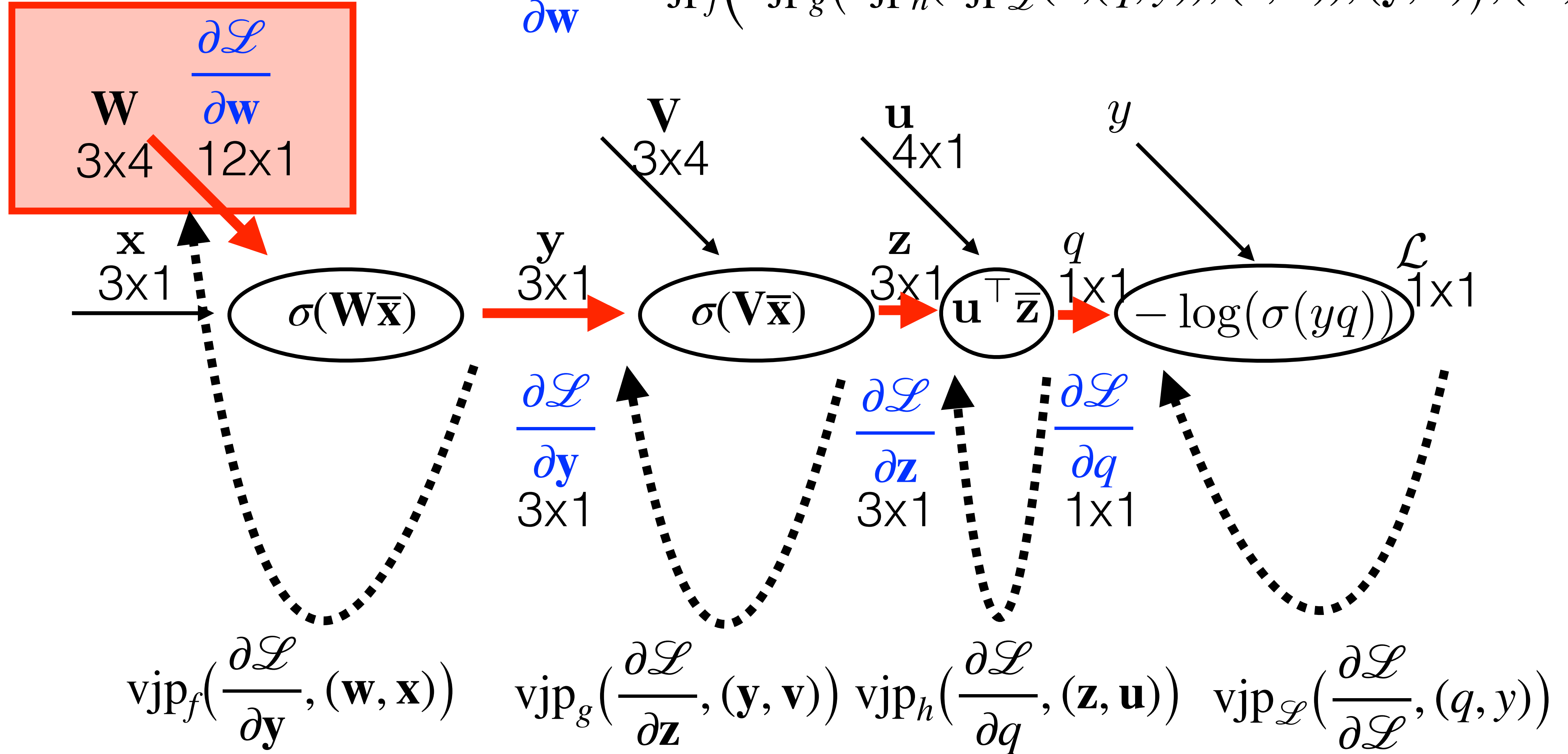
Jacobian wrt  $\mathbf{w}$  :  $\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \text{vjp}_f\left(\text{vjp}_g\left(\text{vjp}_h\left(\text{vjp}_{\mathcal{L}}(1, (q, y)), (\mathbf{z}, \mathbf{u})\right), (\mathbf{y}, \mathbf{v})\right), (\mathbf{w}, \mathbf{x})\right)$



We can (re-)implement vjp in order to preserve inputs resolution

- preserves dimensionality of inputs in backward pass

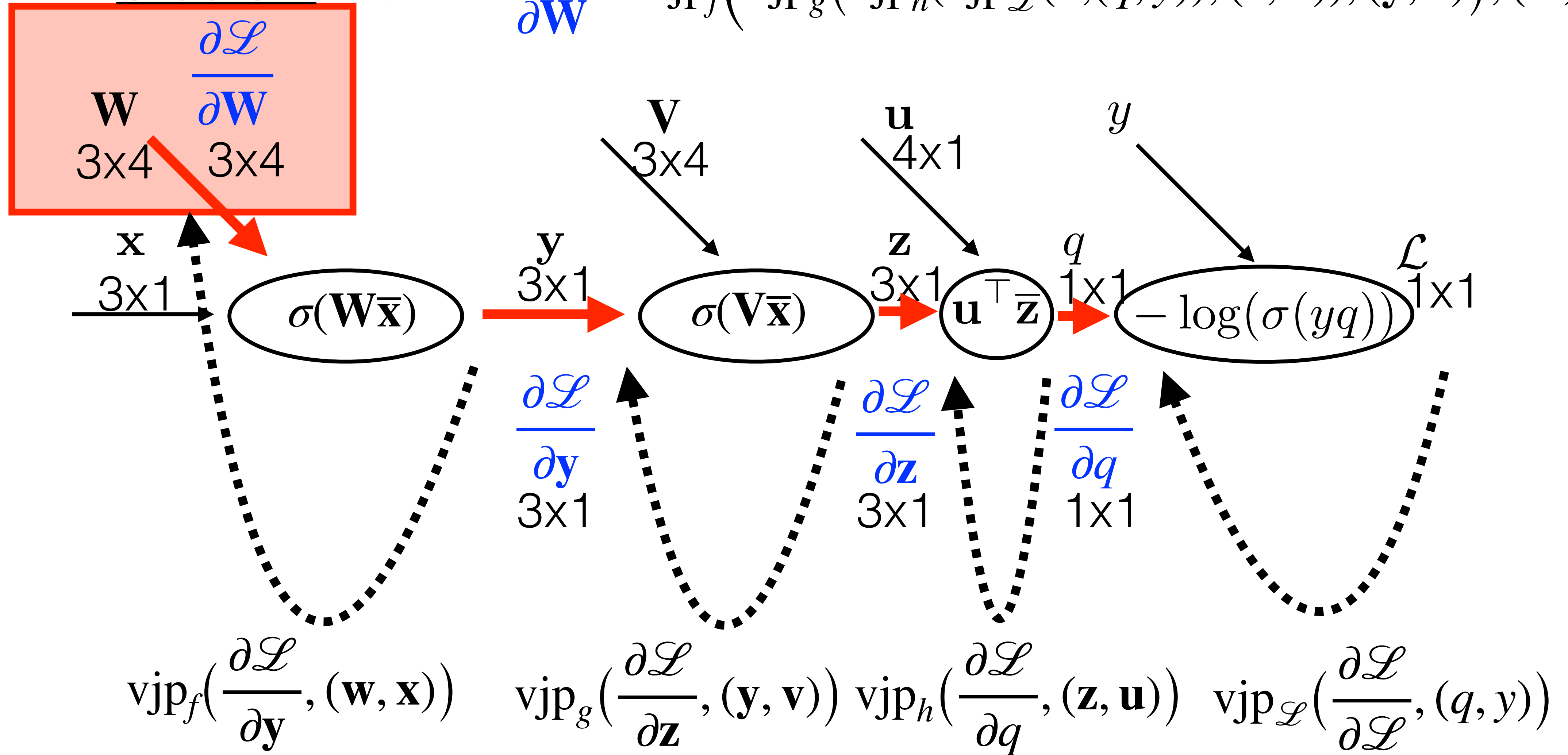
Jacobian wrt  $\mathbf{W}$ :  $\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \text{vjp}_f\left(\text{vjp}_g\left(\text{vjp}_h\left(\text{vjp}_{\mathcal{L}}(1, (q, y)), (\mathbf{z}, \mathbf{u})\right), (\mathbf{y}, \mathbf{v})\right), (\mathbf{w}, \mathbf{x})\right)$



Resulting matrix is reshaped jacobian/gradient of  $\mathcal{L}$  wrt  $\mathbf{W}$

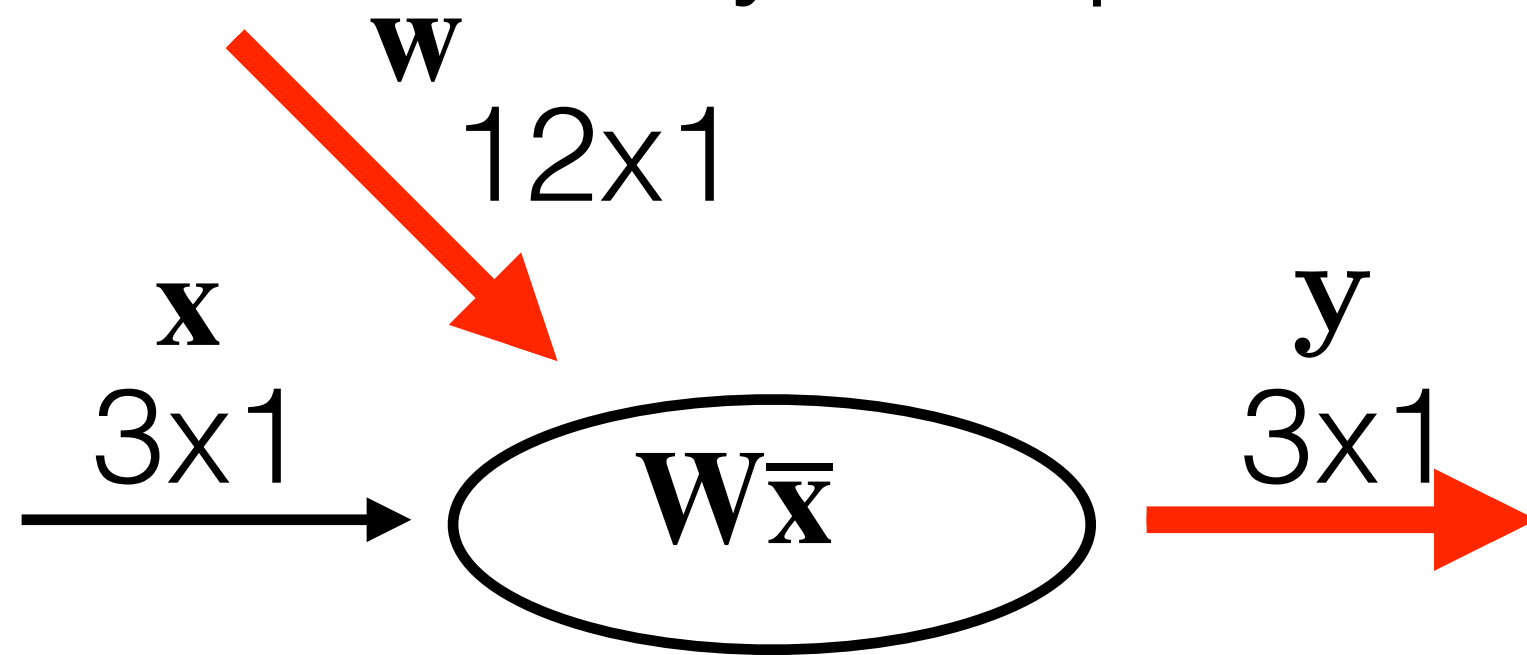
- preserves dimensionality of inputs in backward pass

**Gradient** wrt  $\mathbf{W}$  :  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \text{vjp}_f\left(\text{vjp}_g\left(\text{vjp}_h\left(\text{vjp}_{\mathcal{L}}(1, (q, y)), (\mathbf{z}, \mathbf{u})\right), (\mathbf{y}, \mathbf{v})\right), (\mathbf{w}, \mathbf{x})\right)$



Resulting matrix is reshaped jacobian/gradient of  $\mathcal{L}$  wrt  $\mathbf{W}$

- preserves dimensionality of inputs in backward pass



```
def vjp_w(v, (x, w)) :
```

$$\text{return } \mathbf{v} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{w}} = \mathbf{v} \cdot \begin{bmatrix} -\bar{\mathbf{x}}^T - & & \\ & -\bar{\mathbf{x}}^T - & \\ & & -\bar{\mathbf{x}}^T - \end{bmatrix} = \begin{bmatrix} v_1 \cdot \bar{\mathbf{x}}^T, & v_2 \cdot \bar{\mathbf{x}}^T, & v_3 \cdot \bar{\mathbf{x}}^T \end{bmatrix}$$

1x3

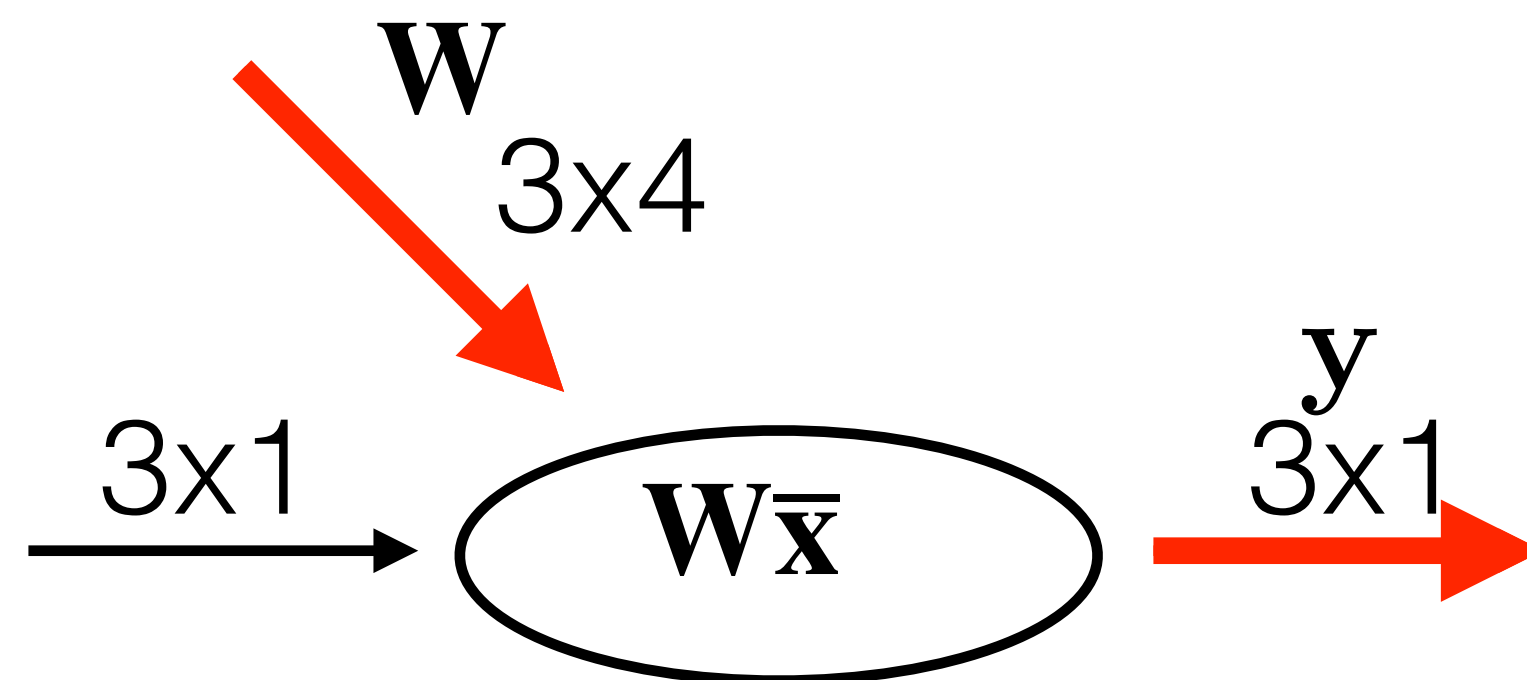
3x12

1x12

```
def vjp_w(v, (x, W)) :
```

$$\text{return } \begin{bmatrix} -v_1 \cdot \bar{\mathbf{x}}^T - \\ -v_2 \cdot \bar{\mathbf{x}}^T - \\ -v_3 \cdot \bar{\mathbf{x}}^T - \end{bmatrix}$$

3x4

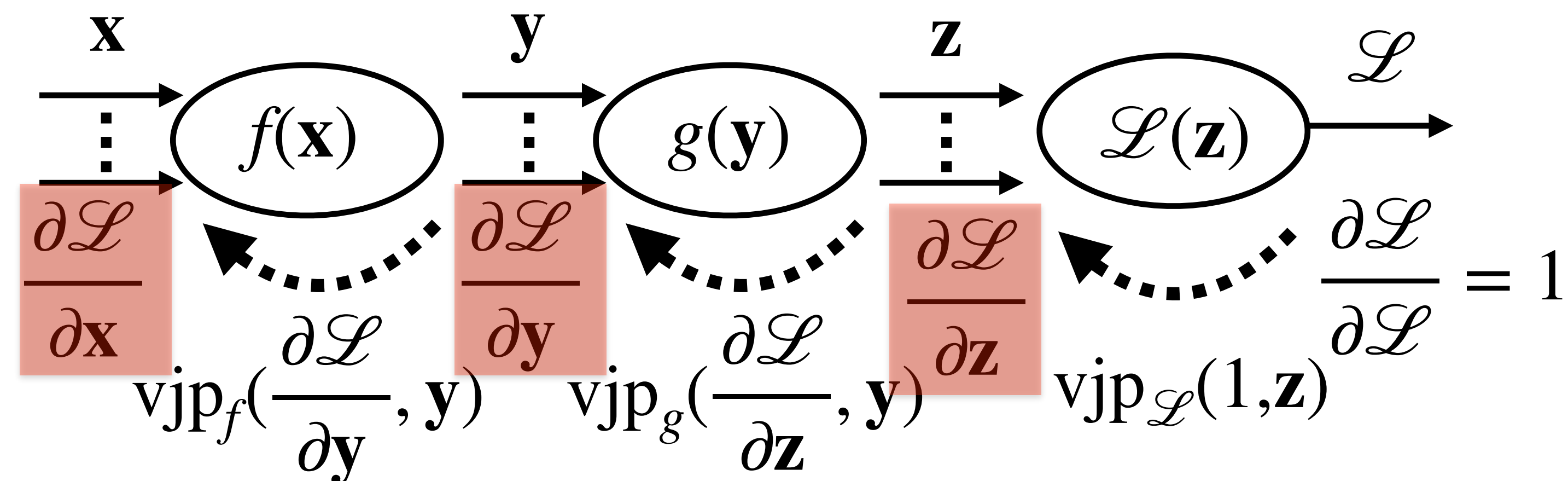




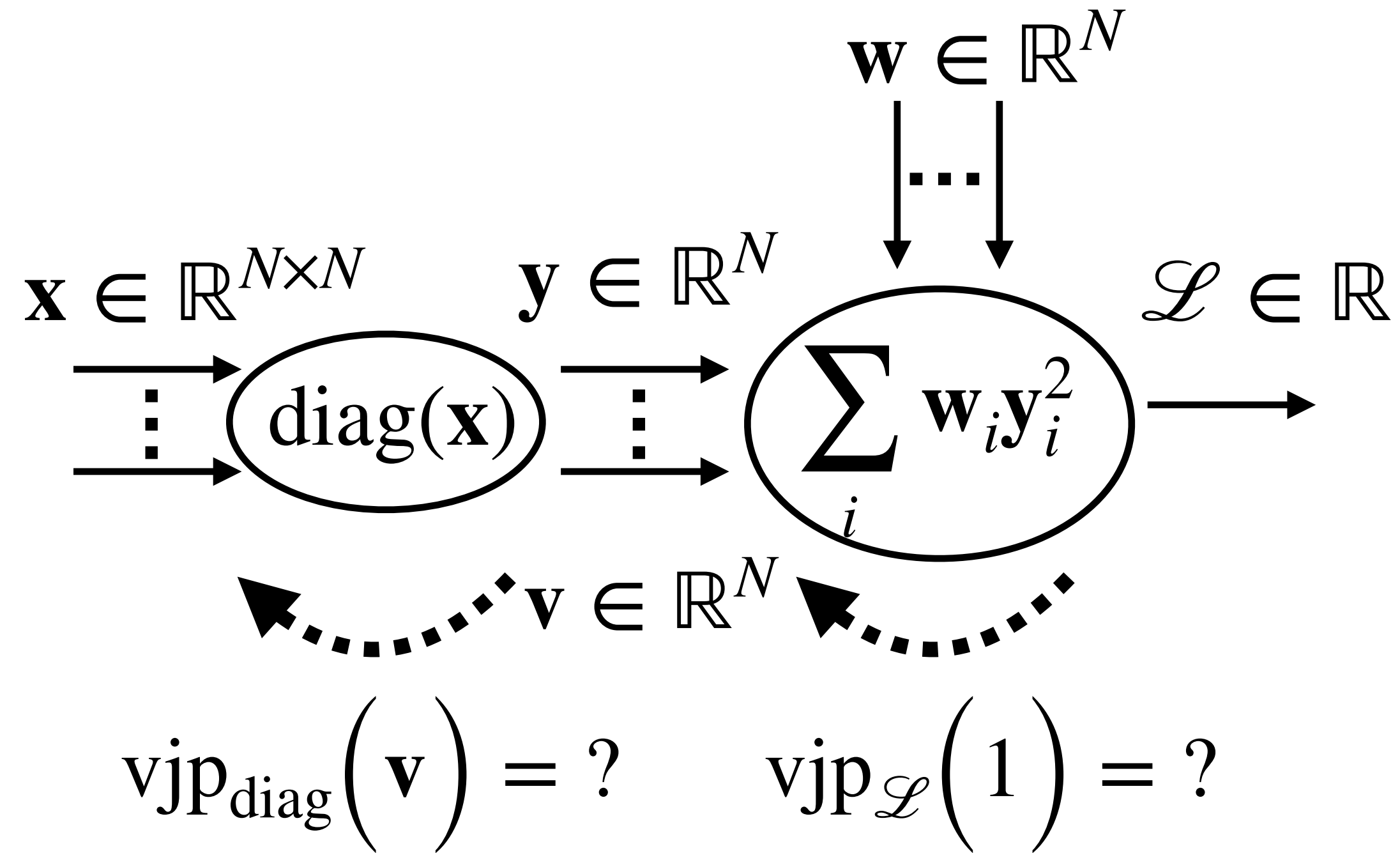
- preserves dimensionality of inputs in backward pass
- if  $g : \mathbb{R}^{m \times m} \rightarrow \mathbb{R}^{n \times n}$ , what is its jacobian?  $\frac{\partial g(\mathbf{y})}{\partial \mathbf{y}} : \mathbb{R}^{m \times m} \rightarrow \mathbb{R}^{n \times n \times m \times m}$
- vjp avoids multiplications of high-dimensional jacobians tensors
- vjp avoids explicit vectorization of higher-dimensional data structures (images)

```
u = torch.tensor([[1, 2], [3, 4]], dtype=torch.float32, requires_grad=True)
loss = torch.sigmoid(u).sum()
grad = torch.autograd.grad(loss, u)[0]
print(grad)
tensor([[0.1966, 0.1050],
        [0.0452, 0.0177]])
```

- edges of comp. graph are populated by gradients of loss wrt edge-variables



# Example: Jacobian vs vector-jacobian-product function

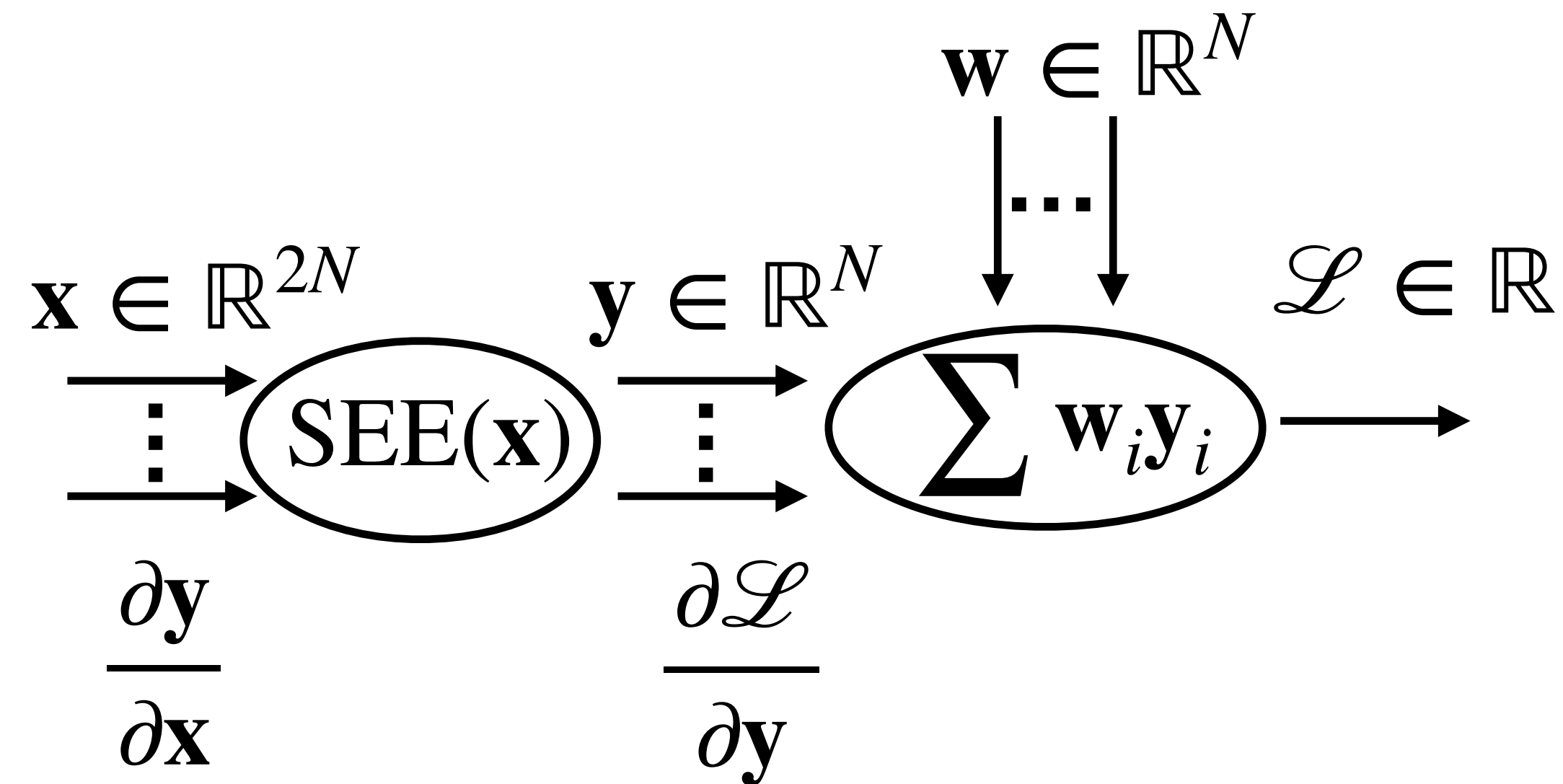


$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = ???$$



# Example: Jacobian vs vector-jacobian-product function

$\mathbf{y} = \text{SEE}(\mathbf{x}) = \text{Select Even Element from input vector } \mathbf{x}$



$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = ???$$

## Neural nets summary

- **Neural net** is a function created as concatenation of simpler functions (e.g. neurons or layers of neurons)
- **Fully connected neural net** is neural network created from neurons, where **all** outputs from previous layer are connected to **all** inputs of the following layer
- **Learning** is gradient optimization of a neural net concatenated with a loss-layer on a training set.
- **Gradient** evaluation is implemented as backward concatenation of vector-jacobian functions (neither symbolic differentiation nor numeric differentiation)
- **VJP** allows to evaluate gradient of scalar output in one backward pass, however the full jacobian of M-dim output requires M passes of vjp! => inefficient LM
- **Deep learning frameworks** (Pytorch, Tensorflow, Caffe) has many predefined layers and takes care about the efficient autodiff on GPU.
- **Deep learning = GPU + data + autodiff**
- **Spoiler alert:** Fully connected NN does not work on structural data (images, sound) well.

Dataset

Error

Linear

FCNN

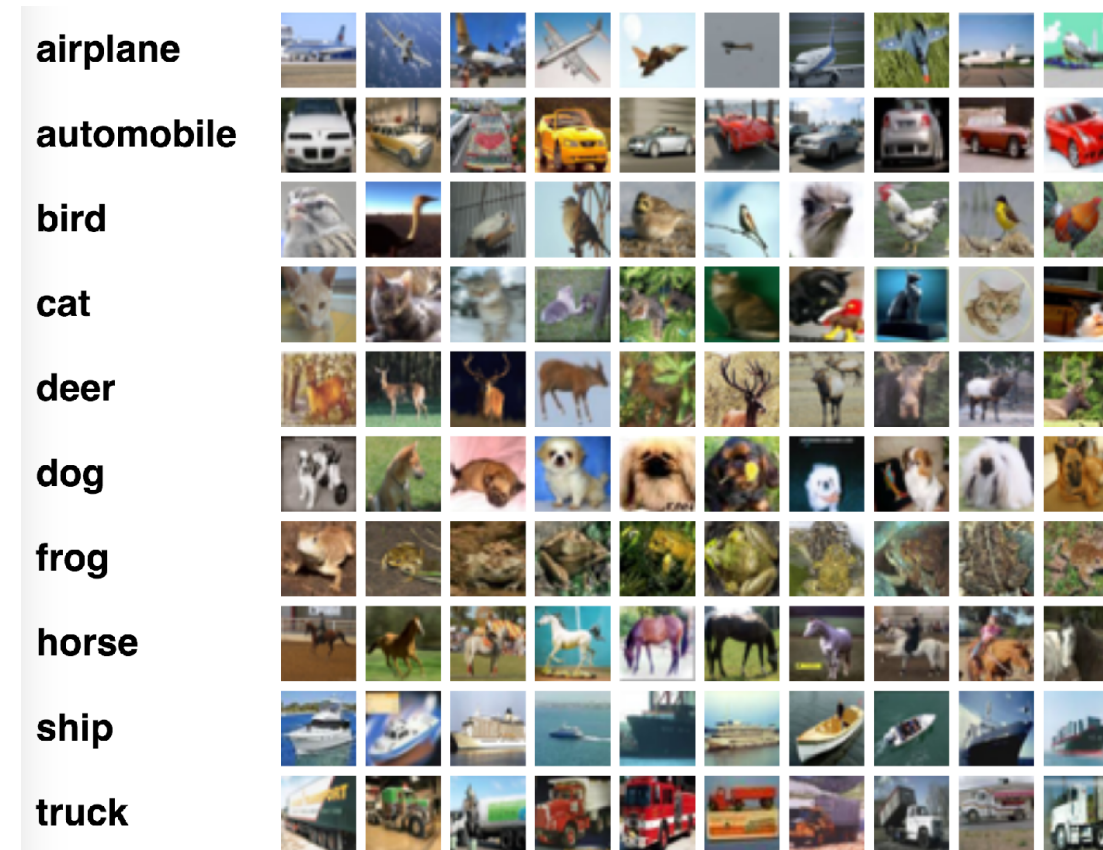
MNIST



8%

??

CIFAR-10



63%

??

<https://benchmarks.ai>

Dataset

Error

Linear

FCNN

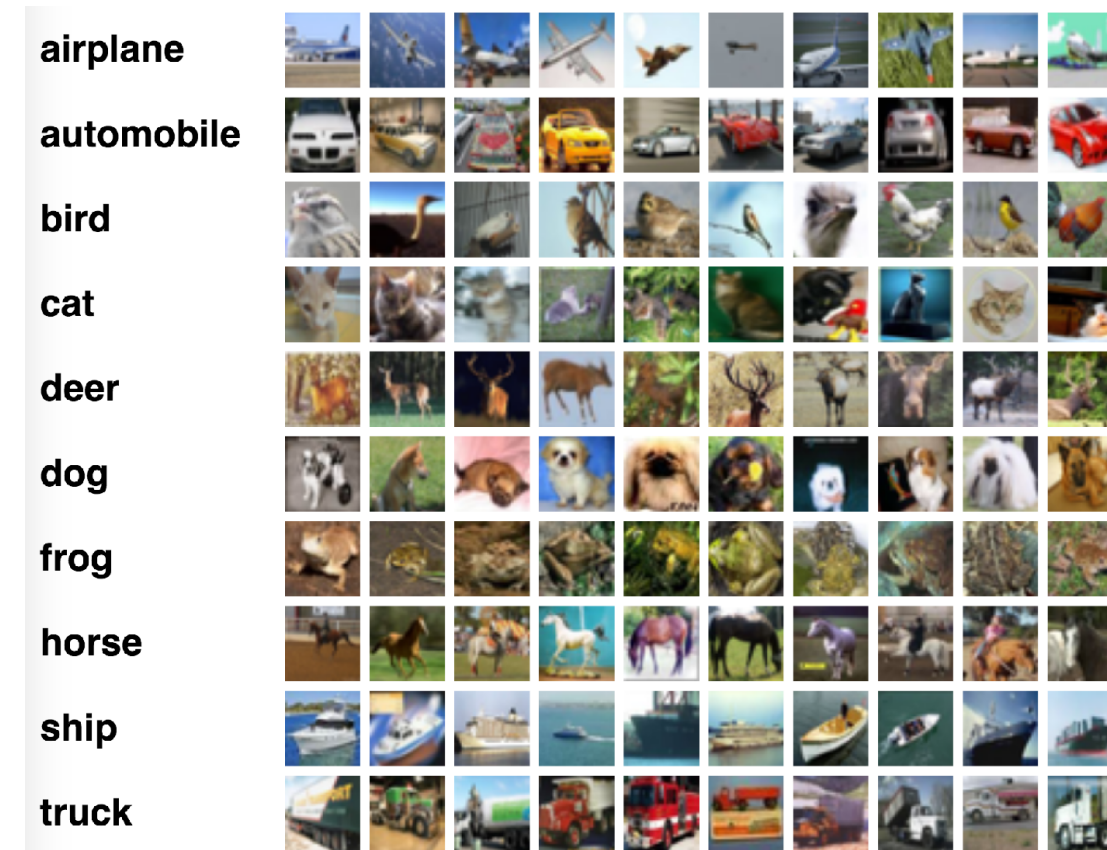
MNIST



8%

2%

CIFAR-10



63%

55%

<https://benchmarks.ai>



Dataset

Error

Linear

FCNN

ConvNet

MNIST

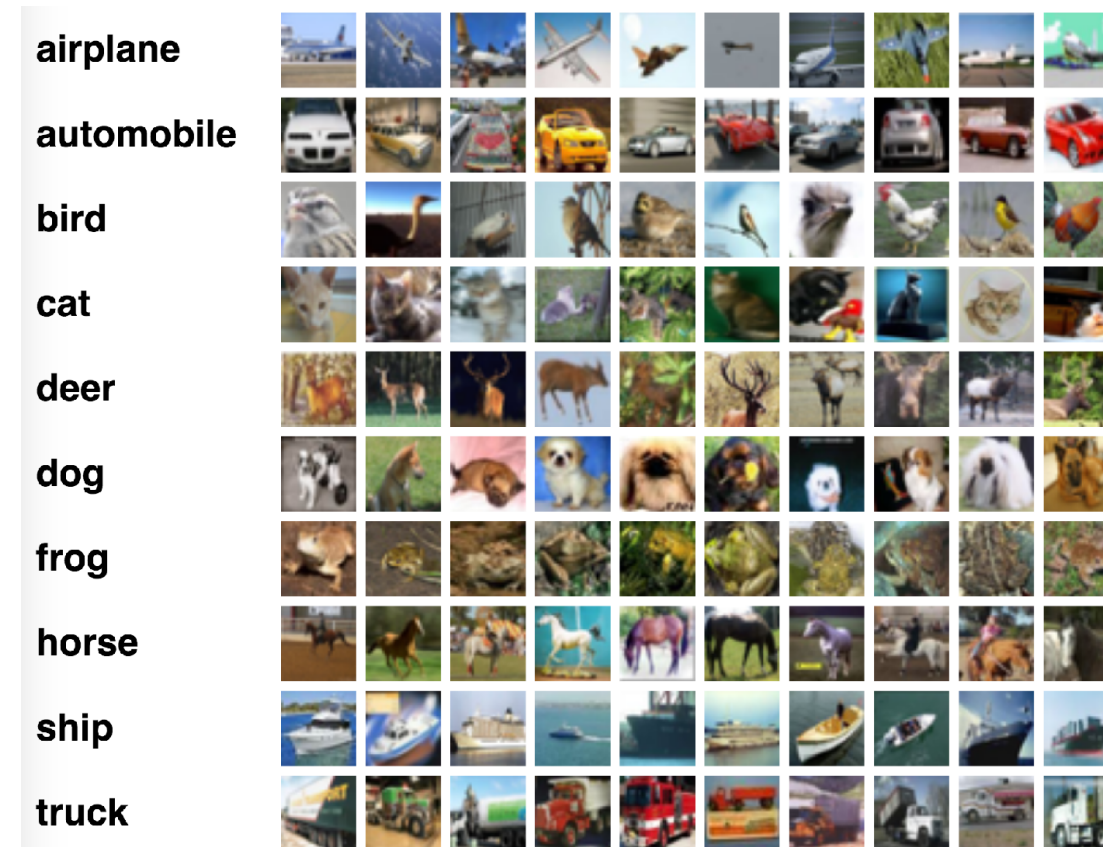


8%

2%

??

CIFAR-10



63%

55%

??

<https://benchmarks.ai>

# Dataset

## Error

Linear

FCNN

ConvNet

MNIST



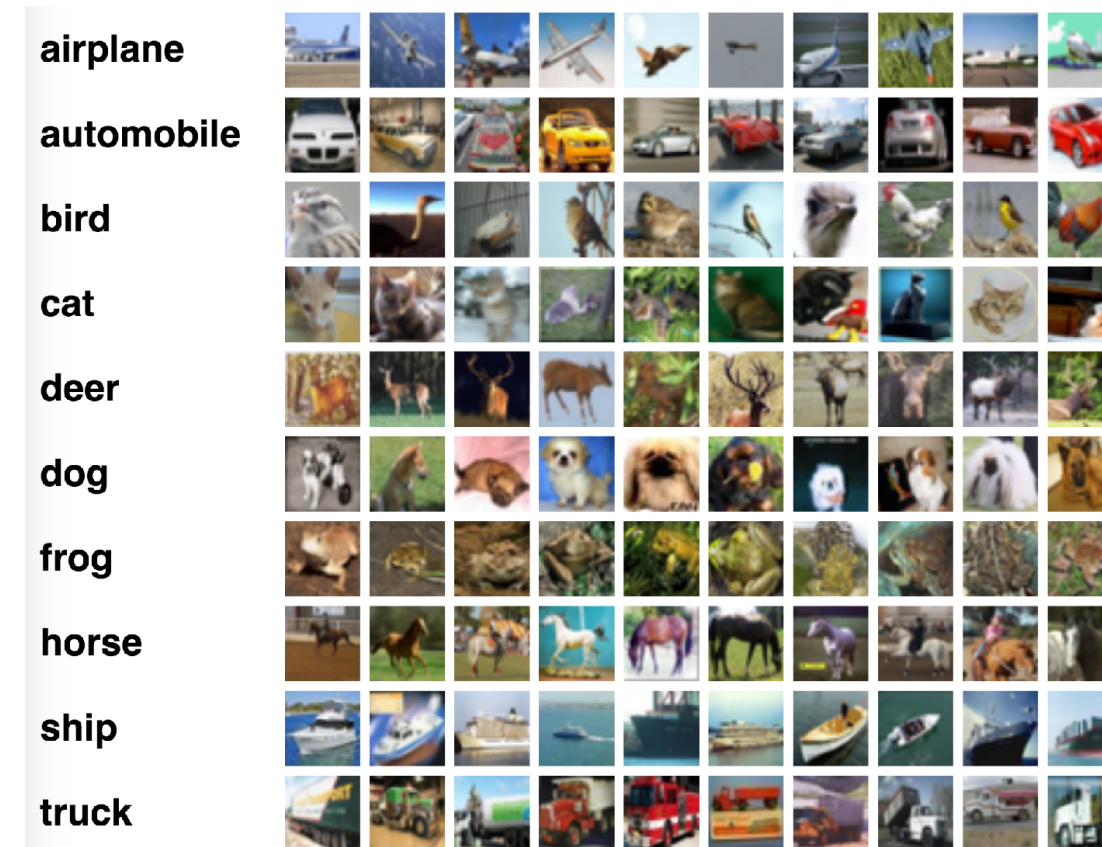
8%

2%

0.2%

[CVPR 2013]

CIFAR-10



63%

55%

1%

[EfficientNet,  
2018]

<https://benchmarks.ai>



# Dataset

## Error

Linear

FCNN

ConvNet

MNIST



8%

2%

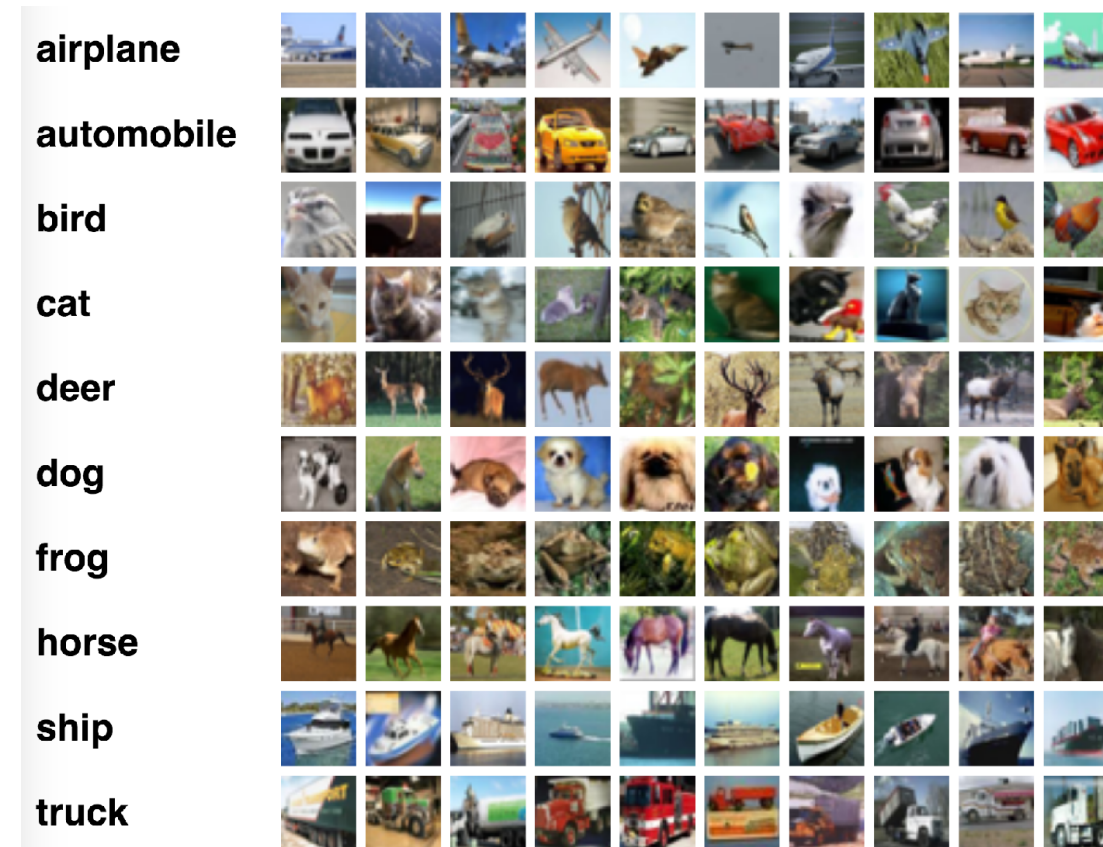
0.2%

underfit  
too  
simple  
arch.

overfit  
too  
complex  
arch.

[CVPR 2013]  
good fit ;-)  
cortex  
inspired  
architecture

CIFAR-10



63%

55%

1%

[EfficientNet,  
2018]

<https://benchmarks.ai>



## **Competencies required for the test T1**

- Ability to draw a computational graph.
- Compute vector-jacobian-product of a given mapping
- Compute backpropagation in computational graph