# Section 2: Learning

Backpropagation

Recurrent Neural Nets
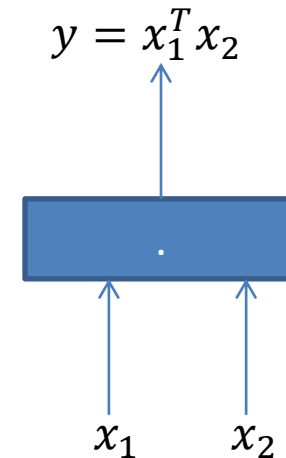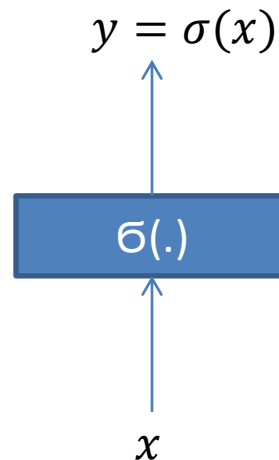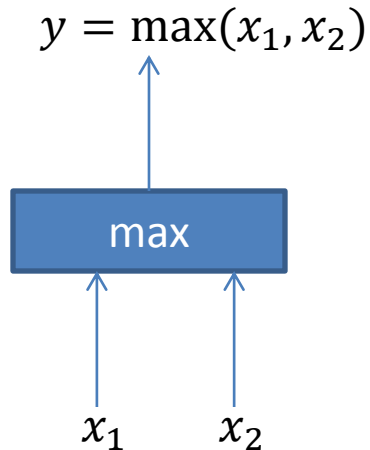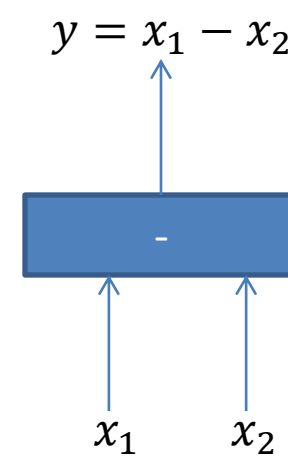
# Topics

- Review of backprop
  - Basic operations
  - Class example
- Additional complexity
  - Shared parameters
  - Dealing with vectors (optional)
- Recurrent Neural Net (RNN)
  - Motivation
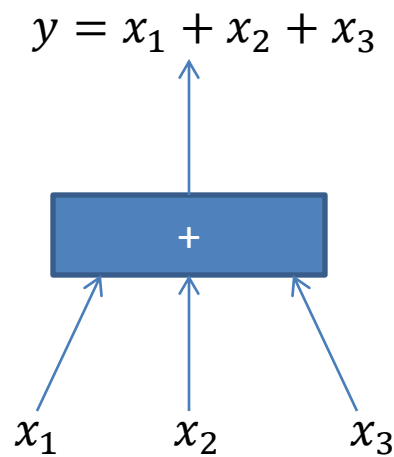  - Simple backprop (vector one optional)
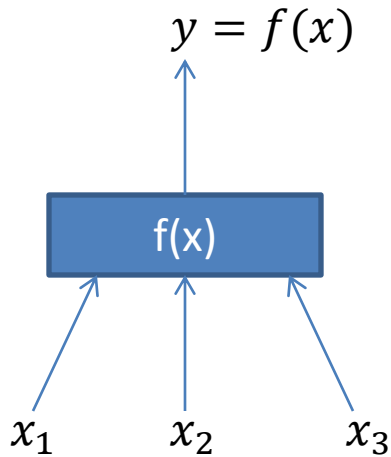  - Demo

# Topics

- **Review of backprop**
  - Basic operations
  - Class example
- Additional complexity
  - Shared parameters
  - Dealing with vectors (optional)
- Recurrent Neural Net (RNN)
  - Motivation
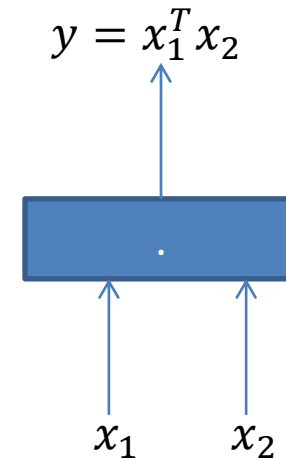  - Simple backprop (vector one optional)
  - Demo

# Review of backprop: Basic Operations

$y = f(x)$

$$f(x)$$

$x_1 \quad x_2 \quad x_3$

$y = x_1 + x_2 + x_3$

$$+$$

$x_1 \quad x_2 \quad x_3$

$y = x_1 - x_2$

$$-$$

$x_1 \quad x_2$

$y = \max(x_1, x_2)$

$$\max$$

$x_1 \quad x_2$

$y = \sigma(x)$

$$\sigma(.)$$

$x$

$y = x_1^T x_2$

$$.$$

$x_1 \quad x_2$

# Review of backprop: Basic Operations

$$y = f(x)$$

f(x)

$$\frac{\partial y}{\partial x_1} \quad \frac{\partial y}{\partial x_2} \quad \frac{\partial y}{\partial x_3}$$

$$x_1 \quad x_2 \quad x_3$$
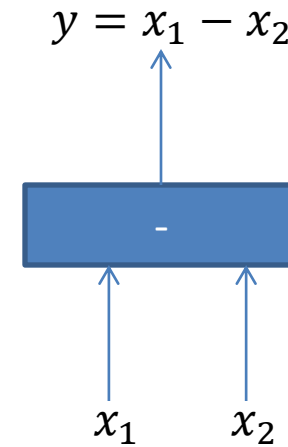
$$y = x_1 + x_2 + x_3$$

+

$$x_1 \quad x_2 \quad x_3$$

$$y = x_1 - x_2$$

-

$$x_1 \quad x_2$$

$$y = \max(x_1, x_2)$$

max

$$x_1 \quad x_2$$

$$y = \sigma(x)$$

б(.)

$$x$$

$$y = x_1^T x_2$$

.

$$x_1 \quad x_2$$

# Review of backprop: Basic Operations

$$y = f(x)$$



$$y = x_1 + x_2 + x_3$$

$$y = x_1 - x_2$$

$$\frac{\partial y}{\partial x_1} \qquad \frac{\partial y}{\partial x_2} \qquad \frac{\partial y}{\partial x_3}$$

f(x)

$x_1 \qquad x_2 \qquad x_3$

$$1 \qquad 1 \qquad 1$$

+

$x_1 \qquad x_2 \qquad x_3$

-

$x_1 \qquad x_2$

$$y = \max(x_1, x_2)$$

max

$x_1 \qquad x_2$

$$y = \sigma(x)$$

б(.)

$x$

$$y = x_1^T x_2$$

.

$x_1 \qquad x_2$

# Review of backprop:
# Basic Operations

$y = f(x)$

$$y = f(x)$$

$$y = x_1 + x_2 + x_3$$

$$y = x_1 - x_2$$

f(x)

$\frac{\partial y}{\partial x_1}$ $\frac{\partial y}{\partial x_2}$ $\frac{\partial y}{\partial x_3}$

$x_1$ $x_2$ $x_3$

+

1 1 1

$x_1$ $x_2$ $x_3$

-

1 $-1$

$x_1$ $x_2$

$$y = \max(x_1, x_2)$$

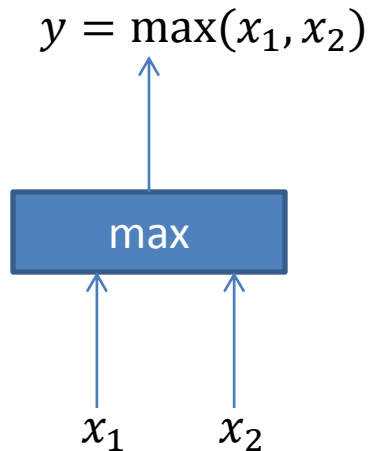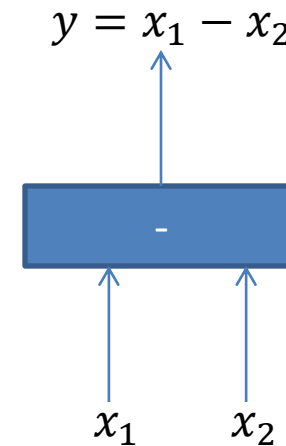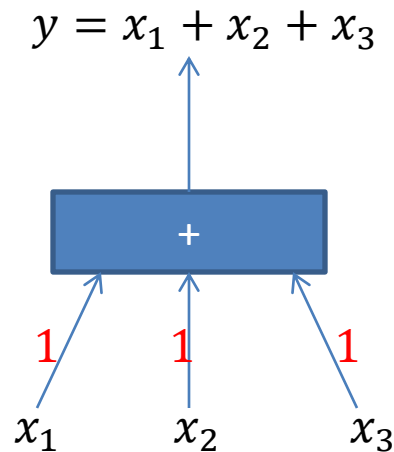$$y = \sigma(x)$$

$$y = x_1^T x_2$$

max

$x_1$ $x_2$

б(.)

$x$

.

$x_1$ $x_2$

# Review of backprop:
# Basic Operations

$y = f(x)$

$$\frac{\partial y}{\partial x_1} \quad \frac{\partial y}{\partial x_2} \quad \frac{\partial y}{\partial x_3}$$

f(x)

$x_1 \quad x_2 \quad x_3$

$y = x_1 + x_2 + x_3$

+

$1 \quad 1 \quad 1$

$x_1 \quad x_2 \quad x_3$

$y = x_1 - x_2$

-

$1 \quad -1$

$x_1 \quad x_2$

$y = \max(x_1, x_2)$

max

$1\{x_1 > x_2\} \quad 1\{x_2 > x_1\}$
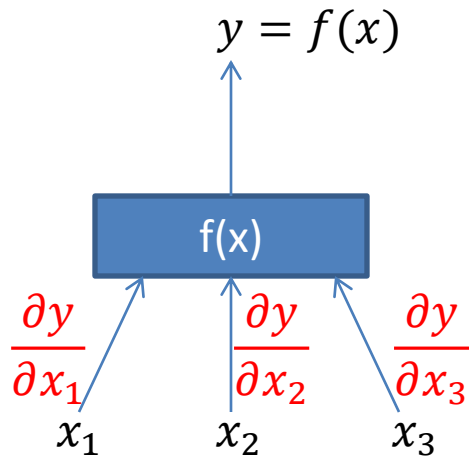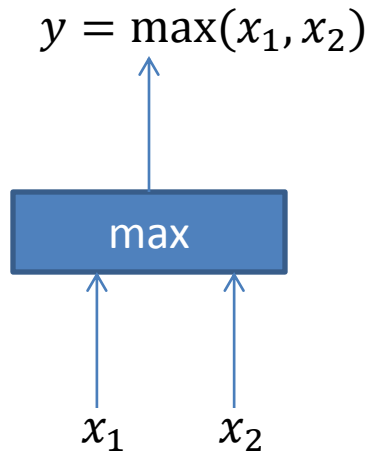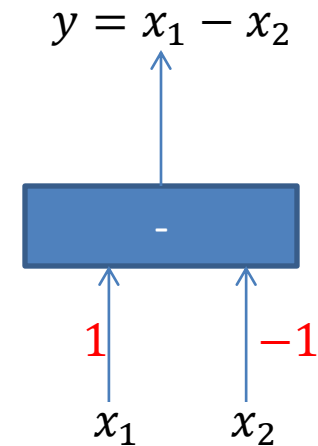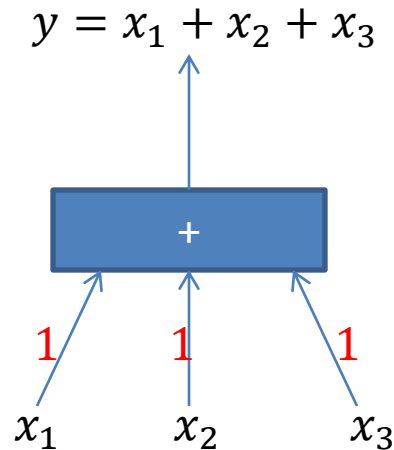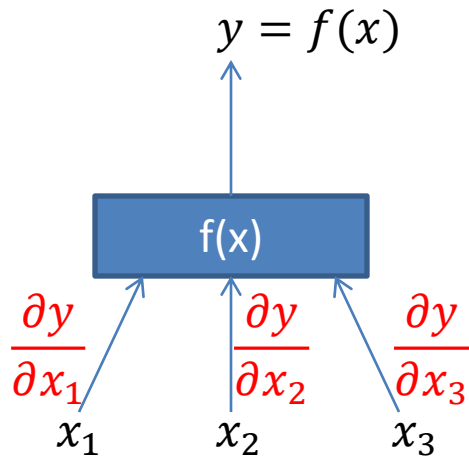
$x_1 \quad x_2$

$y = \sigma(x)$

б(.)

$x$

$y = x_1^T x_2$

.

$x_1 \quad x_2$

# Review of backprop: Basic Operations

$y = f(x)$

$$y = f(x)$$ is computed by the block f(x) with inputs $x_1$, $x_2$, $x_3$ and backprop gradients $\frac{\partial y}{\partial x_1}$, $\frac{\partial y}{\partial x_2}$, $\frac{\partial y}{\partial x_3}$.

$y = x_1 + x_2 + x_3$

The block $+$ has inputs $x_1$, $x_2$, $x_3$ with gradients $1$, $1$, $1$.

$y = x_1 - x_2$

The block $-$ has inputs $x_1$, $x_2$ with gradients $1$, $-1$.

$y = \max(x_1, x_2)$

The block max has inputs $x_1$, $x_2$ with gradients $1\{x_1 > x_2\}$, $1\{x_2 > x_1\}$.

$y = \sigma(x)$

The block $\sigma(.)$ has input $x$ with gradient $y(1 - y)$.

$y = x_1^T x_2$

The block $.$ has inputs $x_1$, $x_2$.

# Review of backprop:
# Basic Operations

$$y = f(x)$$

$$y = x_1 + x_2 + x_3$$

$$y = x_1 - x_2$$

f(x)

$\dfrac{\partial y}{\partial x_1}$  $\dfrac{\partial y}{\partial x_2}$  $\dfrac{\partial y}{\partial x_3}$

$x_1$  $x_2$  $x_3$

\+

1  1  1

$x_1$  $x_2$  $x_3$

\-

1  $-1$

$x_1$  $x_2$

$$y = \max(x_1, x_2)$$

$$y = \sigma(x)$$

$$y = x_1^T x_2$$

max

$1\{x_1 > x_2\}$  $1\{x_2 > x_1\}$

$x_1$  $x_2$

σ(.)

$y(1 - y)$

$x$

.

$x_2$  $x_1$

$x_1$  $x_2$

**Note**: Gradients may be vectors!

# Review of backprop: Example

x: input
p: predicted value
t: true value
L': (squared) loss
$w_1, w_2, v_1, v_2$: parameters to learn

$$y_1 = \sigma(w_1^T x)$$

$$y_2 = \sigma(w_2^T x)$$

$$p = \sigma(v_1 y_1 - v_2 y_2)$$

$$L' = (p - t)^2$$

# Review of backprop: Example



$$y_1 = \sigma(w_1^T x)$$
$$y_2 = \sigma(w_2^T x)$$
$$p = \sigma(v_1 y_1 - v_2 y_2)$$
$$L' = (p - t)^2$$

x: input
p: predicted value
t: true value
L': (squared) loss
$w_1, w_2, v_1, v_2$: parameters

# Review of backprop: Example



$$y_1 = \sigma(w_1^T x)$$

$$y_2 = \sigma(w_2^T x)$$

$$p = \sigma(v_1 y_1 - v_2 y_2)$$

$$L' = (p - t)^2$$

# Review of backprop: Example



$$y_1 = \sigma(w_1^T x)$$
$$y_2 = \sigma(w_2^T x)$$
$$p = \sigma(v_1 y_1 - v_2 y_2)$$
$$L' = (p - t)^2$$

14

# Review of backprop: Example



$$y_1 = \sigma(w_1^T x)$$
$$y_2 = \sigma(w_2^T x)$$
$$p = \sigma(v_1 y_1 - v_2 y_2)$$
$$L' = (p - t)^2$$

15

# Review of backprop: Example



$$y_1 = \sigma(w_1^T x)$$
$$y_2 = \sigma(w_2^T x)$$
$$p = \sigma(v_1 y_1 - v_2 y_2)$$
$$L' = (p - t)^2$$

# Review of backprop: Example



$$y_1 = \sigma(w_1^T x)$$
$$y_2 = \sigma(w_2^T x)$$
$$p = \sigma(v_1 y_1 - v_2 y_2)$$
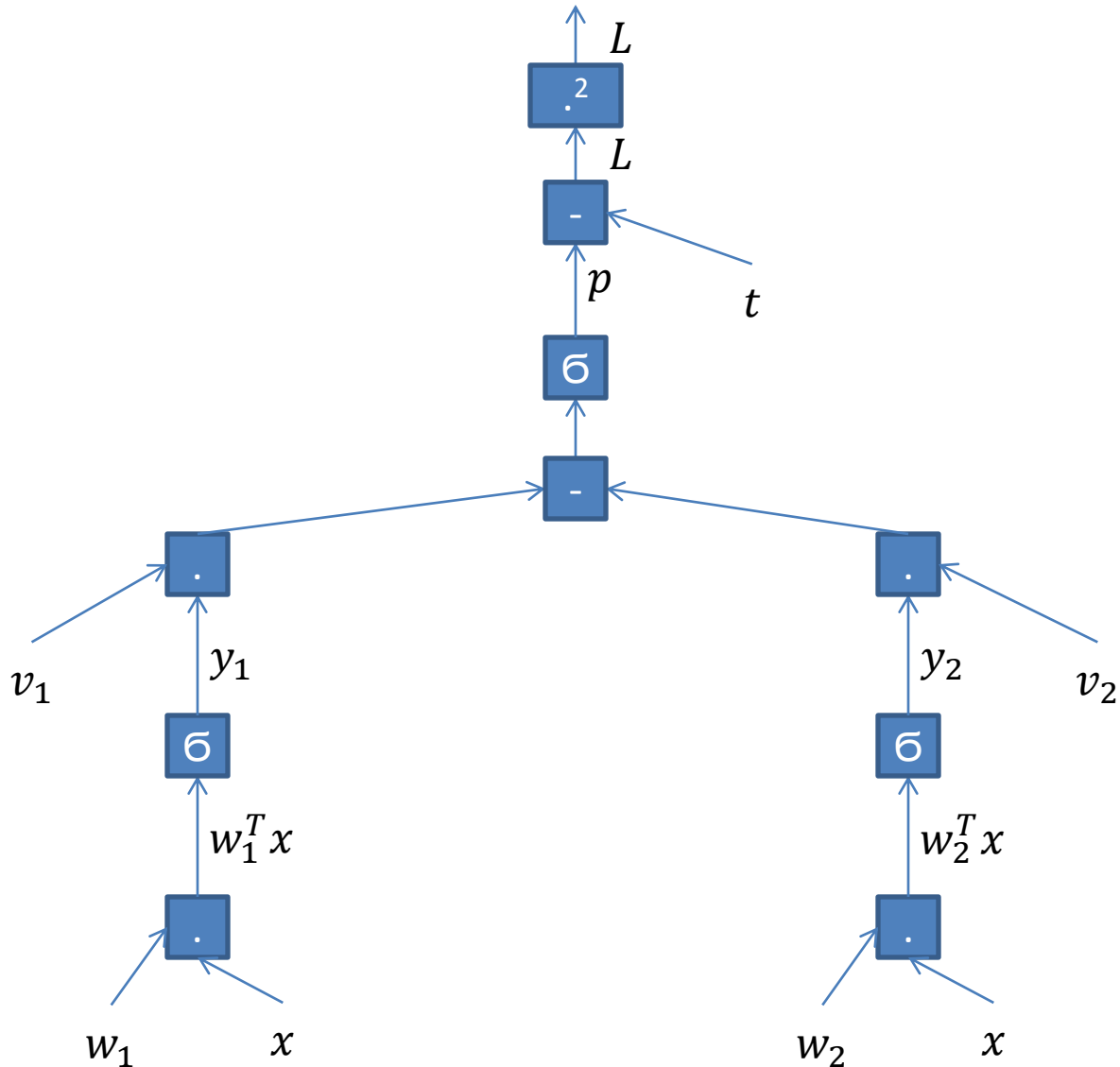$$L' = (p - t)^2$$

# Review of backprop: Example



$$y_1 = \sigma(w_1^T x)$$

$$y_2 = \sigma(w_2^T x)$$

$$p = \sigma(v_1 y_1 - v_2 y_2)$$

$$L' = (p - t)^2$$

18

# Topics

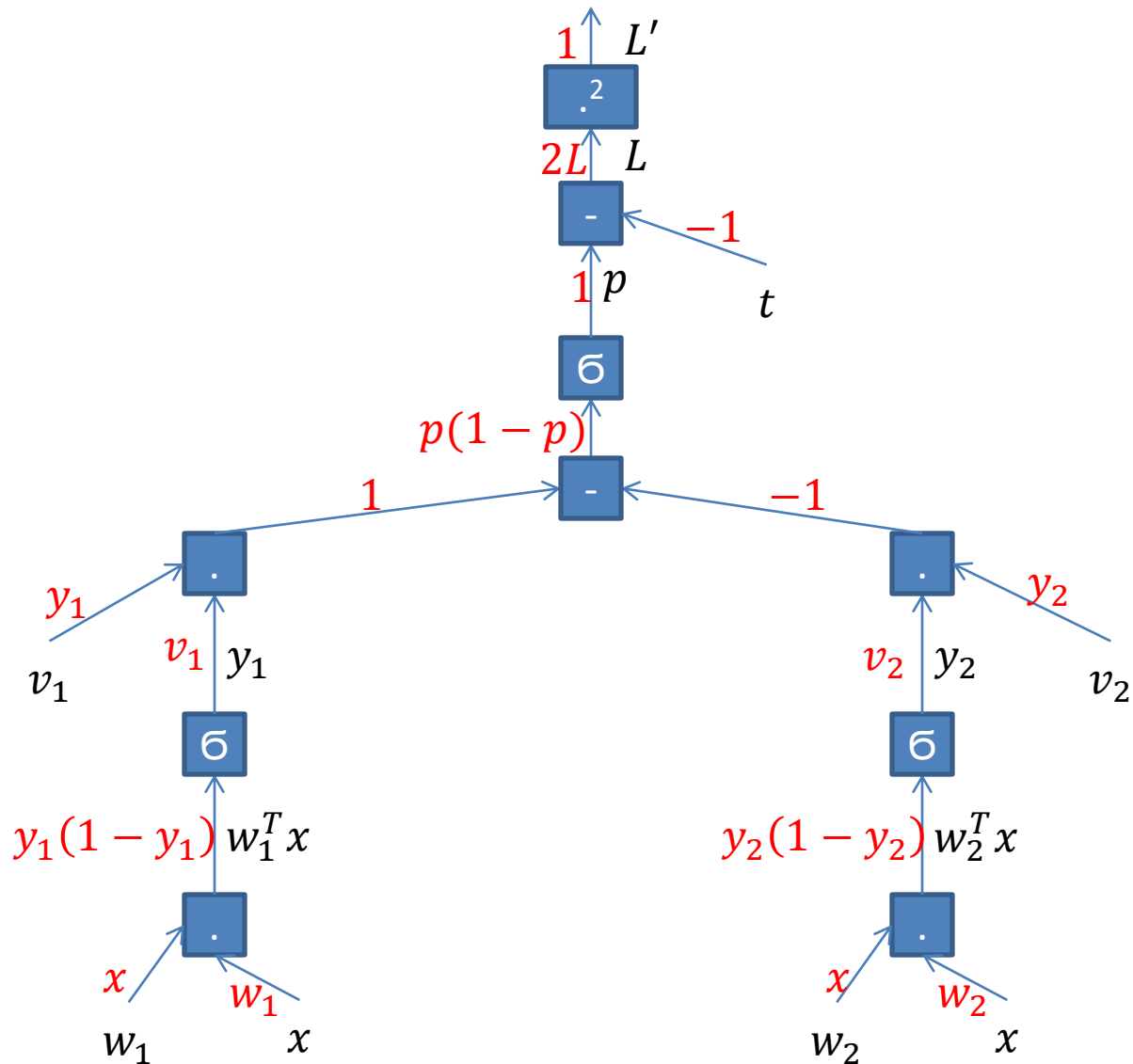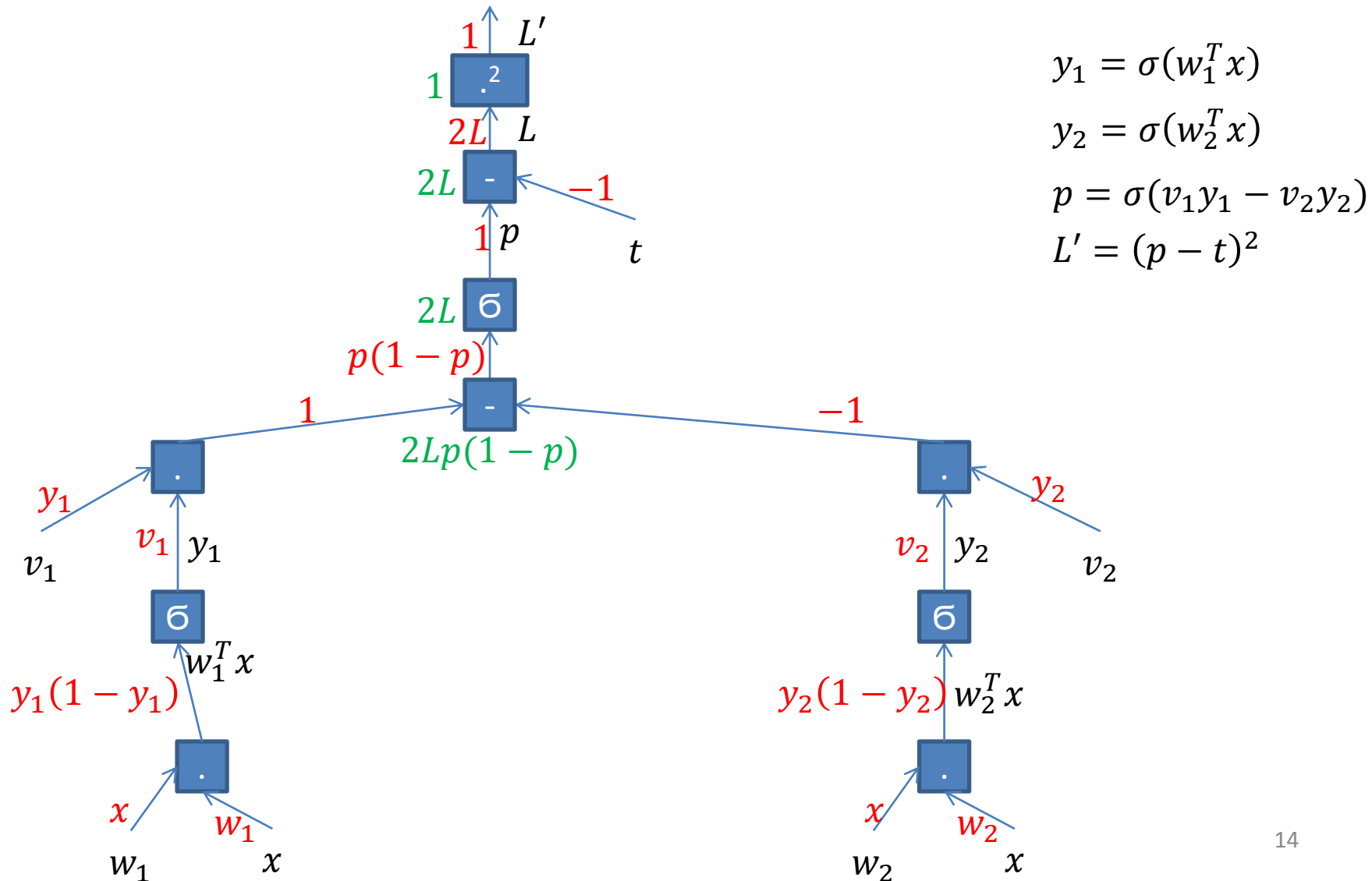- Review of backprop
  - Basic operations
  - Class example
- **Additional complexity**
  - **Shared parameters**
  - **Dealing with vectors (optional)**
- Recurrent Neural Net (RNN)
  - Motivation
  - Simple backprop (vector one optional)
  - Demo

# Additional complexity: Shared parameters

So far, each computed output goes to 1 unit only



$L'$

$1$

$.^2$

$2L$  $L$

$-$  $-1$

$1$  $p$  $t$

$\sigma$

$p(1-p)$

$1$  $+$  $1$

$.$  $.$

$y_1$  $y_2$

$v_1$  $v_2$

$v_1$  $y_1$  $v_2$  $y_2$

$\sigma$  $\sigma$

$y_1(1-y_1)$  $w_1^T x$  $y_2(1-y_2)$  $w_2^T x$

$.$  $.$

$x$  $x$

$w_1$  $w_2$

$w_1$  $x$  $w_2$  $x$

# Additional complexity: Shared parameters

So far, each computed output goes to 1 unit only

**What if?**

# Additional complexity: Shared parameters

So far, each computed output goes to 1 unit only

**What if?**

# Additional complexity: Shared parameters

So far, each computed output goes to 1 unit only

**ADD gradients**

# Additional complexity:
# Shared parameters

But what if same output goes to multiple units?

Example: $z = w^2(w - 3)$, compute $\frac{dz}{dw}$

# Additional complexity:
# Shared parameters

But what if same output goes to multiple units?

Example: $z = w^2(w - 3)$, compute $\frac{dz}{dw}$

By hand, use product rule:

$$\frac{dz}{dw} = \frac{dw^2}{dw}(w - 3) + w^2\frac{d(w - 3)}{dw}$$
$$= 2w(w - 3) \quad + w^2$$
$$= 3w^2 - 6w$$

# Additional complexity:
# Shared parameters

But what if same output goes to multiple units?

Example: $z = w^2(w - 3)$, compute $\frac{dz}{dw}$

By hand, use product rule:

$$\frac{dz}{dw} = \frac{dw^2}{dw}(w - 3) + w^2\frac{d(w - 3)}{dw}$$
$$= 2w(w - 3) \quad + w^2$$
$$= 3w^2 - 6w$$

Computational steps:
$x = w^2$
$y = w - 3$
$z = x.y$

# Additional complexity: Shared parameters

But what if same output goes to multiple units?

Example: $z = w^2(w - 3)$, compute $\frac{dz}{dw}$

By hand, use product rule:
$$\frac{dz}{dw} = \frac{dw^2}{dw}(w - 3) + w^2\frac{d(w - 3)}{dw}$$
$$= 2w(w - 3) \quad + w^2$$
$$= 3w^2 - 6w$$

Computational steps:
$x = w^2$
$y = w - 3$
$z = x.y$

# Additional complexity: Shared parameters

But what if same output goes to multiple units?

Example:     $\boldsymbol{z = w^2(w - 3)}$, compute $\frac{dz}{dw}$

By hand, use product rule:

$$\frac{dz}{dw} = \frac{dw^2}{dw}(w - 3) + w^2 \frac{d(w - 3)}{dw}$$

$$= 2w(w - 3) \quad + w^2$$

$$= 3w^2 - 6w$$

Computational steps:

$$x = w^2$$

$$y = w - 3$$

$$z = x.y$$

# Additional complexity: Shared parameters

But what if same output goes to multiple units?

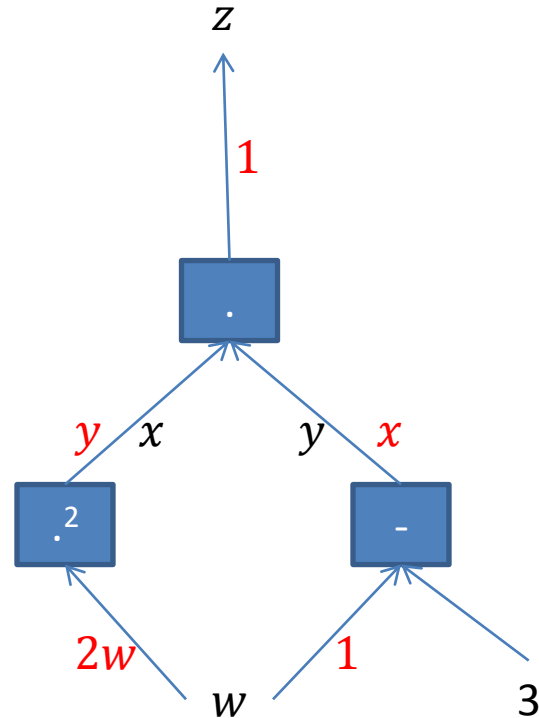Example: $z = w^2(w-3)$, compute $\frac{dz}{dw}$

By hand, use product rule:
$$\frac{dz}{dw} = \frac{dw^2}{dw}(w-3) + w^2 \frac{d(w-3)}{dw}$$
$$= 2w(w-3) \quad + w^2$$
$$= 3w^2 - 6w$$

Computational steps:
$x = w^2$
$y = w - 3$
$z = x.y$

# Additional complexity: Shared parameters

So far, each computed output goes to 1 unit only

**Similarly…**

# Additional complexity: Shared parameters

So far, each computed output goes to 1 unit only

**Similarly...**

# Additional complexity:
# Shared parameters

So far, each computed output goes to 1 unit only

**Similarly...**

# Additional complexity:
# Shared parameters

So far, each computed output goes to 1 unit only

**Similarly...**



$$\nabla_w L = \alpha + \beta$$

# Additional complexity:
# Shared parameters

- Why do we add when output/parameter replicated?

Proof

$$L(y', y'')$$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial y'} \cdot \frac{\partial y'}{\partial y} + \frac{\partial L}{\partial y''} \cdot \frac{\partial y''}{\partial y}$$

# Additional complexity:
# Shared parameters

- Why do we add when output/parameter replicated?

Proof

$$L(y', y'')$$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial y'} \cdot \frac{\partial y'}{\partial y} + \frac{\partial L}{\partial y''} \cdot \frac{\partial y''}{\partial y}$$

Since $y = y' = y''$, it follows:

# Additional complexity: Shared parameters

- Why do we add when output/parameter replicated?

Proof

$L(y', y'')$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial y'} \cdot \frac{\partial y'}{\partial y} + \frac{\partial L}{\partial y''} \cdot \frac{\partial y''}{\partial y}$$

Since $y = y' = y''$, it follows:

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial y'} + \frac{\partial L}{\partial y''}$$

# Additional complexity:
# Shared parameters

- Why do we add when output/parameter replicated?

Proof

$L(y', y'')$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial y'} \cdot \frac{\partial y'}{\partial y} + \frac{\partial L}{\partial y''} \cdot \frac{\partial y''}{\partial y}$$
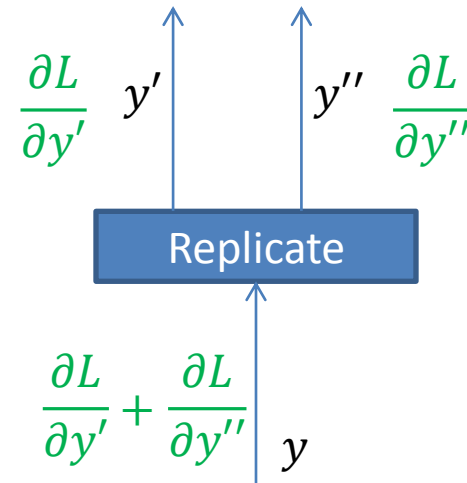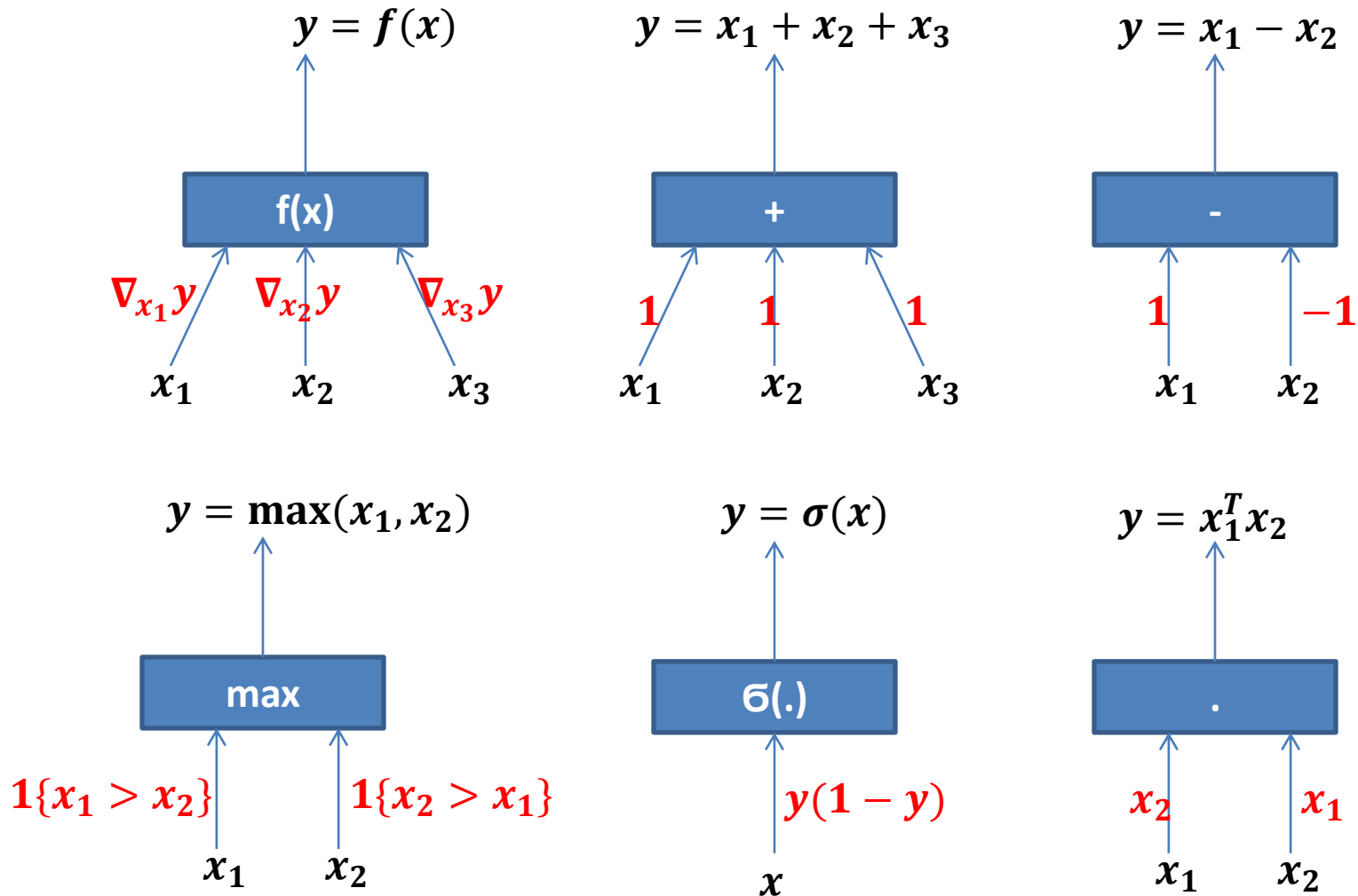
Since $y = y' = y''$, it follows:

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial y'} + \frac{\partial L}{\partial y''}$$

$\frac{\partial L}{\partial y'}$ $y'$ $y''$ $\frac{\partial L}{\partial y''}$

Replicate

$\frac{\partial L}{\partial y'} + \frac{\partial L}{\partial y''}$ $y$

# Additional complexity:
# Dealing with vectors (optional)

$$y = f(x)$$

$$y = x_1 + x_2 + x_3$$

$$y = x_1 - x_2$$

| f(x) |
| --- |

$\nabla_{x_1} y$   $\nabla_{x_2} y$   $\nabla_{x_3} y$

$x_1$   $x_2$   $x_3$

| + |
| --- |

$1$   $1$   $1$

$x_1$   $x_2$   $x_3$

| - |
| --- |

$1$   $-1$

$x_1$   $x_2$

$$y = \max(x_1, x_2)$$

$$y = \sigma(x)$$

$$y = x_1^T x_2$$

| max |
| --- |

$1\{x_1 > x_2\}$   $1\{x_2 > x_1\}$

$x_1$   $x_2$

| σ(.) |
| --- |

$y(1 - y)$

$x$

| . |
| --- |

$x_2$   $x_1$

$x_1$   $x_2$

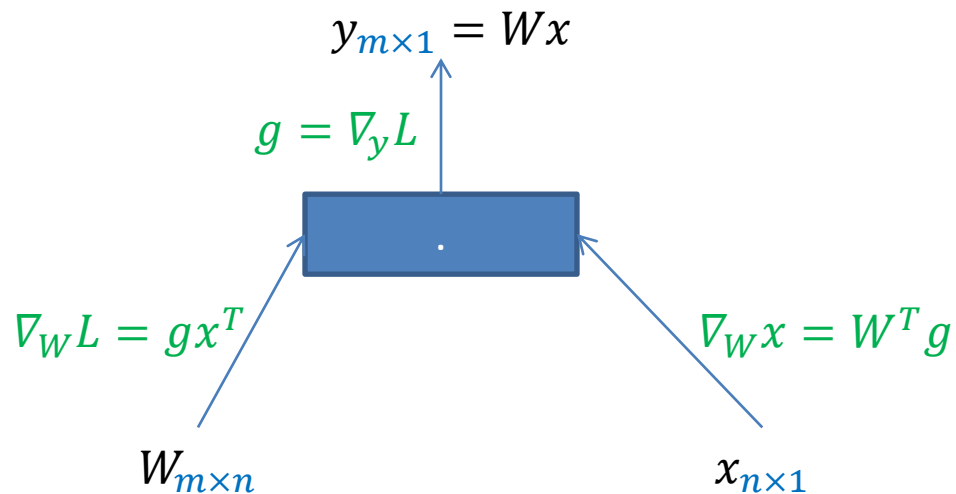**Note**: All operations, including backprop, are component-wise!

# Additional complexity:
# Dealing with vectors (optional)

$$y_{m \times 1} = Wx$$



$W_{m \times n}$        $x_{n \times 1}$

# Additional complexity:
# Dealing with vectors (optional)

$$y_{m \times 1} = Wx$$

$$g = \nabla_y L$$

$$\nabla_W L = g x^T$$

$$\nabla_W x = W^T g$$

$$W_{m \times n}$$

$$x_{n \times 1}$$

# Additional complexity:
# Dealing with vectors (optional)

$$y_{m \times 1} = Wx$$

$g$



$gx^T$

$W^T g$

$W_{m \times n}$

$x_{n \times 1}$

# Additional complexity:
# Dealing with vectors (optional)

$y_{m \times 1} = Wx$

$g$

$gx^T$ $W^T g$

$W_{m \times n}$ $x_{n \times 1}$

$L = \mathbf{1}^T \sigma(Wx)$

$L = z_1 + \ldots + z_m$

$\mathbf{1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{m \times 1}$

$z_{m \times 1} = \sigma(y)$

$\sigma$

$y_{m \times 1} = Wx$

$W_{m \times n}$ $x_{n \times 1}$

# Additional complexity:
# Dealing with vectors (optional)

$y_{m\times 1} = Wx$

$g$

$gx^T$

$W^T g$

$W_{m\times n}$

$x_{n\times 1}$

$L = \mathbf{1}^T \sigma(Wx)$

$L = z_1 + \ldots + z_m$

$\mathbf{1}$

$z_{m\times 1} = \sigma(y)$

$\sigma$

$y_{m\times 1} = Wx$

$W_{m\times n}$

$x_{n\times 1}$

# Additional complexity:
# Dealing with vectors (optional)

$$y_{m \times 1} = Wx$$

$g$

$gx^T$　　　　$W^T g$

$W_{m \times n}$　　　　$x_{n \times 1}$

$$L = \mathbf{1}^T \sigma(Wx)$$

$$1 \quad L = z_1 + \ldots + z_m$$

$\mathbf{z}$

$\mathbf{1}$

$\mathbf{1}$　$z_{m \times 1} = \sigma(y)$

$$y_{m \times 1} = Wx$$

$W_{m \times n}$　　　　$x_{n \times 1}$

44

# Additional complexity:
# Dealing with vectors (optional)

$y_{m \times 1} = Wx$

$g$

.

$gx^T$

$W^T g$

$W_{m \times n}$

$x_{n \times 1}$

$L = \mathbf{1}^T \sigma(Wx)$

$1 \quad L = z_1 + \dots + z_m$

$\mathbf{z}$

.

$\mathbf{1}$

$\mathbf{1}$

$z_{m \times 1} = \sigma(y)$

σ

$z^\circ(\mathbf{1} - z) = \begin{bmatrix} z_1(1 - z_1) \\ \vdots \\ z_2(1 - z_2) \end{bmatrix}_{m \times 1}$

$y_{m \times 1} = Wx$

.

$W_{m \times n}$

$x_{n \times 1}$

45

# Additional complexity:
# Dealing with vectors (optional)

$y_{m \times 1} = Wx$

$g$

$gx^T$

$W^T g$

$W_{m \times n}$

$x_{n \times 1}$

$L = \mathbf{1}^T \sigma(Wx)$

$1$

$L = z_1 + \ldots + z_m$

$\mathbf{z}$

$\mathbf{1}$

$\mathbf{1}$

$z_{m \times 1} = \sigma(y)$

$\sigma$

$\mathbf{z}°(\mathbf{1 - z})$

$y_{m \times 1} = Wx$

$W_{m \times n}$

$x_{n \times 1}$

# Additional complexity:
# Dealing with vectors (optional)

$$y_{m \times 1} = Wx$$

$$g$$



$$gx^T$$

$$W^T g$$

$$W_{m \times n}$$

$$x_{n \times 1}$$

$$L = \mathbf{1}^T \sigma(Wx)$$

$$1 \quad\quad L = z_1 + \ldots + z_m$$

$$\mathbf{z}$$

$$\mathbf{1}$$

$$\mathbf{1}$$

$$z_{m \times 1} = \sigma(y)$$

$$\sigma$$

$$\mathbf{z} \circ (\mathbf{1} - \mathbf{z}) \quad\quad y_{m \times 1} = Wx$$

$$.\,\mathbf{x}^T$$

$$\mathbf{W}^T.$$

$$W_{m \times n}$$

$$x_{n \times 1}$$

# Additional complexity:
# Dealing with vectors (optional)

$y_{m \times 1} = Wx$

$g$



$g x^T$

$W^T g$

$W_{m \times n}$

$x_{n \times 1}$

$L = \mathbf{1}^T \sigma(Wx)$

$1$

$L = z_1 + \ldots + z_m$

$1$

$\mathbf{z}$

$\mathbf{1}$

$1$

$z_{m \times 1} = \sigma(y)$

$1$

б

$\mathbf{z}°(\mathbf{1} - \mathbf{z})$

$y_{m \times 1} = Wx$

$\mathbf{1}°\mathbf{z}°(\mathbf{1} - \mathbf{z})$

$.\, \mathbf{x}^T$

$\mathbf{W}^T.$

$\mathbf{1}°\mathbf{z}°(\mathbf{1} - \mathbf{z})\mathbf{x}^T \, W_{m \times n}$

$x_{n \times 1}$

$\mathbf{W}^T \mathbf{1}°\mathbf{z}°(\mathbf{1} - \mathbf{z})$

# Additional complexity:
# Dealing with vectors (optional)

$$y_{m \times 1} = Wx$$

$g$

$gx^T$ $W^Tg$

$W_{m \times n}$ $x_{n \times 1}$

---

$1$ $1$ $L = z_1 + \ldots + z_m$

$+$

$1$ $\mathbf{1}$ $z_{m \times 1} = \sigma(y)$

$\sigma$

$\mathbf{1° z° (1 - z)}$ $\mathbf{z° (1 - z)}$ $y_{m \times 1} = Wx$

$.$

$\mathbf{1° z° (1 - z) x^T}$ $. \, \mathbf{x^T}$ $\mathbf{W^T} .$

$\mathbf{W^T 1° z° (1 - z)}$

$W_{m \times n}$ $x_{n \times 1}$

# Topics

- Review of backprop
  - Basic operations
  - Class example
- Additional complexity
  - Shared parameters
  - Dealing with vectors (optional)
- **Recurrent Neural Net (RNN)**
  - Motivation
  - Simple backprop (vector one optional)
  - Demo

# Recurrent Neural Net (RNN): Motivation

- Conventional Neural Networks

| Output y $\in R^n$ |
| :---: |

$\sigma(\mathrm{W}x), ReLU(Wx)$ or $Wx$

⋮

$\sigma(\mathrm{W}x), ReLU(Wx)$ or $Wx$

| Input $x \in R^d$ |
| :---: |

**Fixed input and output size**
- 1 input (*d* dimensional)
- 1 output (*n* dimensional)

**Fixed computing steps**
- Independent of input
- Static framework

**Typical applications**
- Image classification
- Regression

# Recurrent Neural Net (RNN): Motivation

- We desire variable input/output size, variable computational steps…

| one to one | one to many | many to one | many to many | many to many |



Applications:    Image captioning      Sentiment analysis      Machine translation      Text generation

Video classification

*Source*: http://karpathy.github.io/2015/05/21/rnn-effectiveness/

# Recurrent Neural Net (RNN): Motivation

- Image captioning (one-to-many)



"man in black shirt is playing guitar."

"man in blue wetsuit is surfing on wave."

"a young boy is holding a baseball bat."

*Source*: http://cs.stanford.edu/people/karpathy/deepimagesent/

# Recurrent Neural Net (RNN): Motivation

- ## Sentiment Analysis (many-to-one)

The action switches between past and present , but the material link is too tenuous to anchor the emotional connections that purport to span a 125-year divide .

Drops you into a dizzying , volatile , pressure-cooker of a situation that quickly snowballs out of control , while focusing on the what much more than the why .

The film is itself a sort of cinematic high crime , one that brings military courtroom dramas down very , very low .

**Classify as:**
0 - negative
1 - somewhat negative
2 - neutral
3 - somewhat positive
4 - positive

# Recurrent Neural Net (RNN): Motivation

- Machine translation (many-to-many)

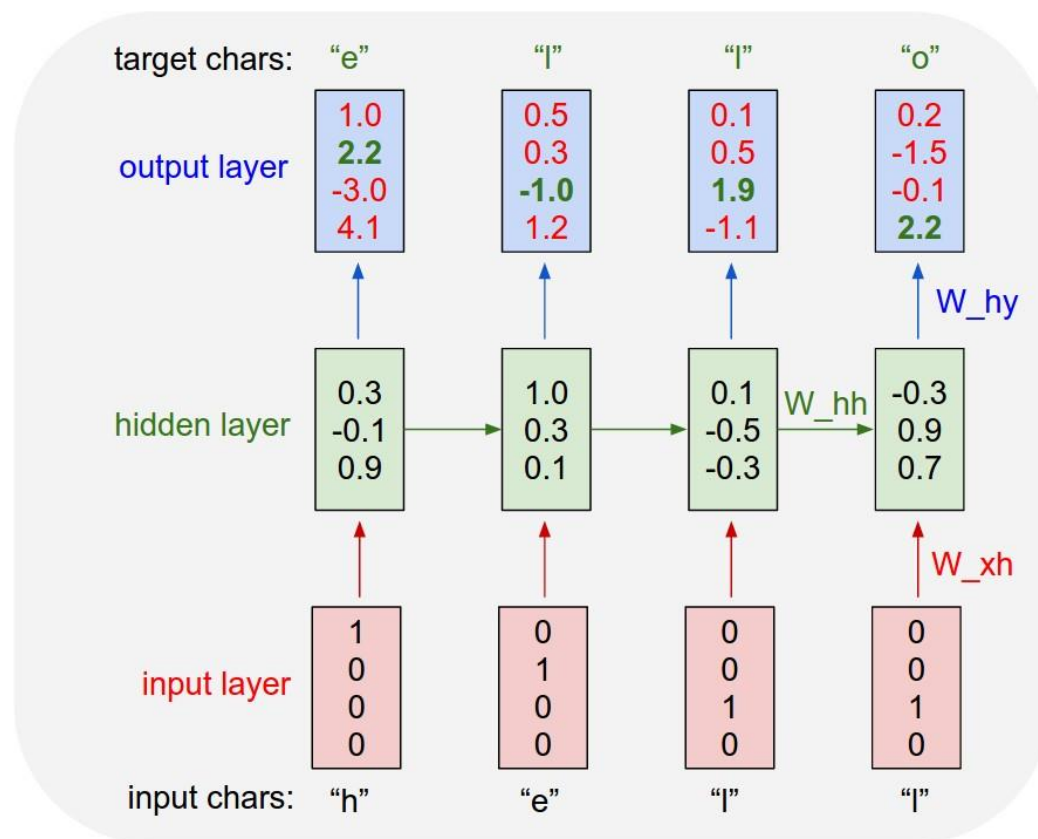*Source*: https://translate.google.com (above), http://cs224d.stanford.edu/lectures/CS224d-Lecture8.pdf (below)

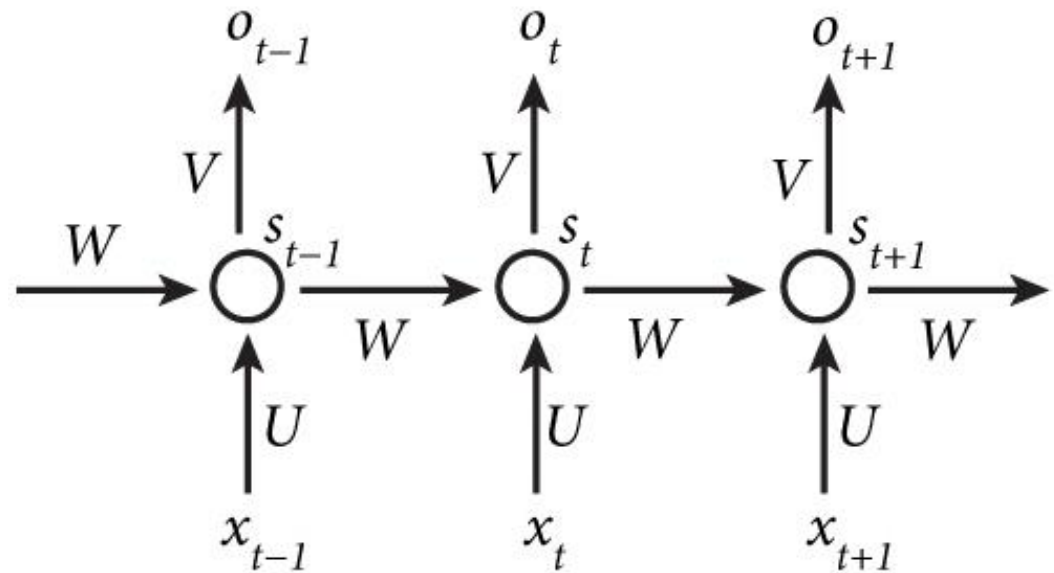# Recurrent Neural Net (RNN): Motivation
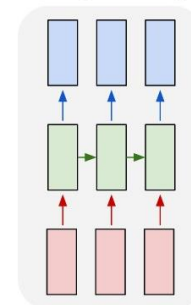
- Character-level language model (many-to-many)



Predict every next character. For example, consider "hello"; use "hell" to predict "ello"

# Recurrent Neural Net (RNN): Example

# Recurrent Neural Net (RNN): Example

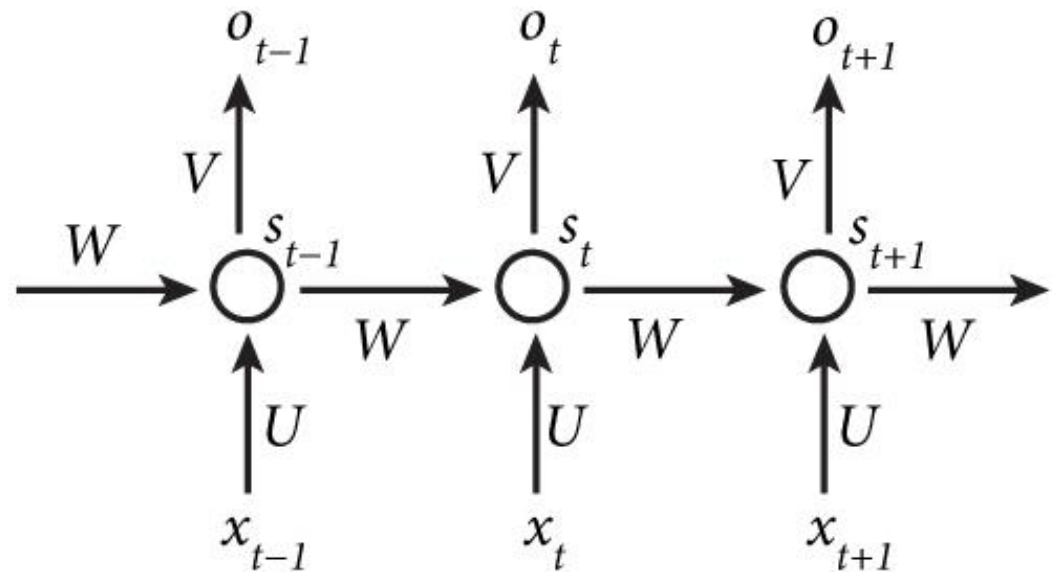$$s_t = \sigma(Ux_t + Ws_{t-1})$$

$$o_t = \sigma(Vs_t)$$

$x_t$: input
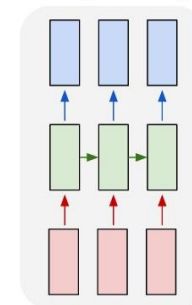$s_t$: hidden state   ($s_{-1} = 0$ )
$o_t$: output
$U, V, W$: parameters (matrices)

Hidden state stores past information
that may be relevant in future.
- Provides context
- Long range dependence



many to many

58

# Recurrent Neural Net (RNN): Example

## Character-level language model (many-to-many)

[DEMO]
**Script**: https://gist.github.com/karpathy/d4dee566867f8291f086
**Config**: hidden state has 100 dimensions, 93 different characters
**Input**: norvig.com/big.txt

**iter 0, loss: 113.314988**
I fechowta ecyoumepuave omas mmur a band chou os Carbinn yond here wa,k, oly soongy pas yin fou alinfo#gtid ed levenupksen Ia tbinl and.  Yury sleve lsok ufimeme conlanf youlsseg ve-;aud Mas finn ass w

**iter 185000, loss: 49.667141**
hing the hri, theme ummengi-hy linced. The candiccevinicas he visur. The her in war to Ereart in dnintorvaned wenced to as rewnighly restera he by appored bat riculing at hooke thiming a somews, and

**Observe**: Starts looking like English; new sentence starts with big letter; end with full stop; short words spelled correctly; long words still messed up

# Recurrent Neural Net (RNN): Example
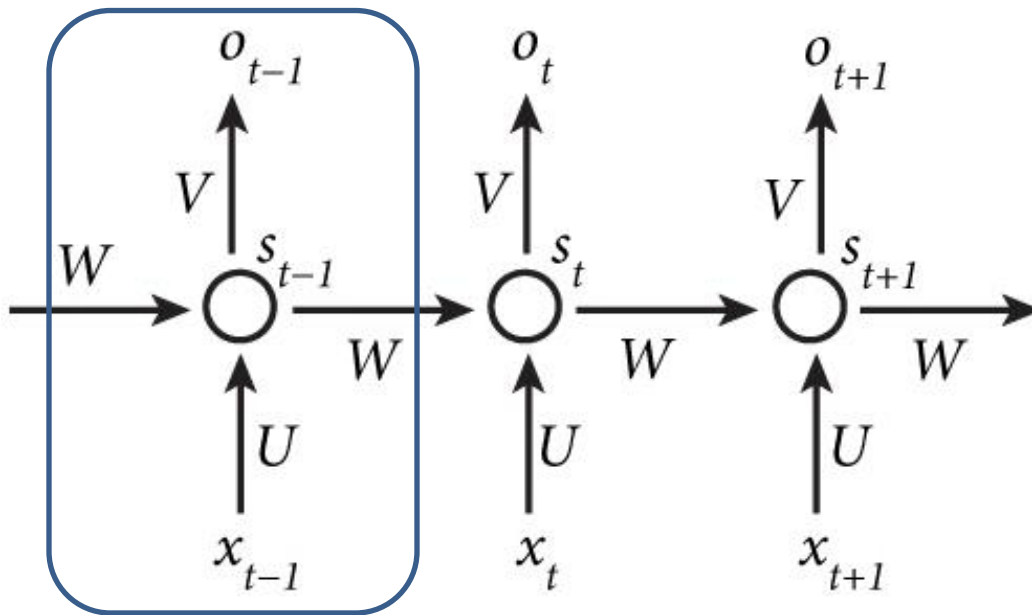
## Character-level language model (many-to-many)

**Paul Graham's essays**
Source: http://karpathy.github.io/2015/05/21/rnn-effectiveness/

*"The surprised in investors weren't going to raise money. I'm not the company with the time there are all interesting quickly, don't have to get off the same programmers. There's a super-angel round fundraising, why do you can do. If you have a different physical investment are become in people who reduced in a startup with the way to argument the acquirer could see them just that you're also the founders will part of users' affords that and an alternation to the idea. [2] Don't work at first member to see the way kids will seem in advance of a bad successful startup. And if you have to act the big company too."*

**Observe**: Learns spelling and grammar from scratch; learns to cite; says *"a company is a meeting to think to investors", starts understanding a bit.*

# Recurrent Neural Net (RNN): Example



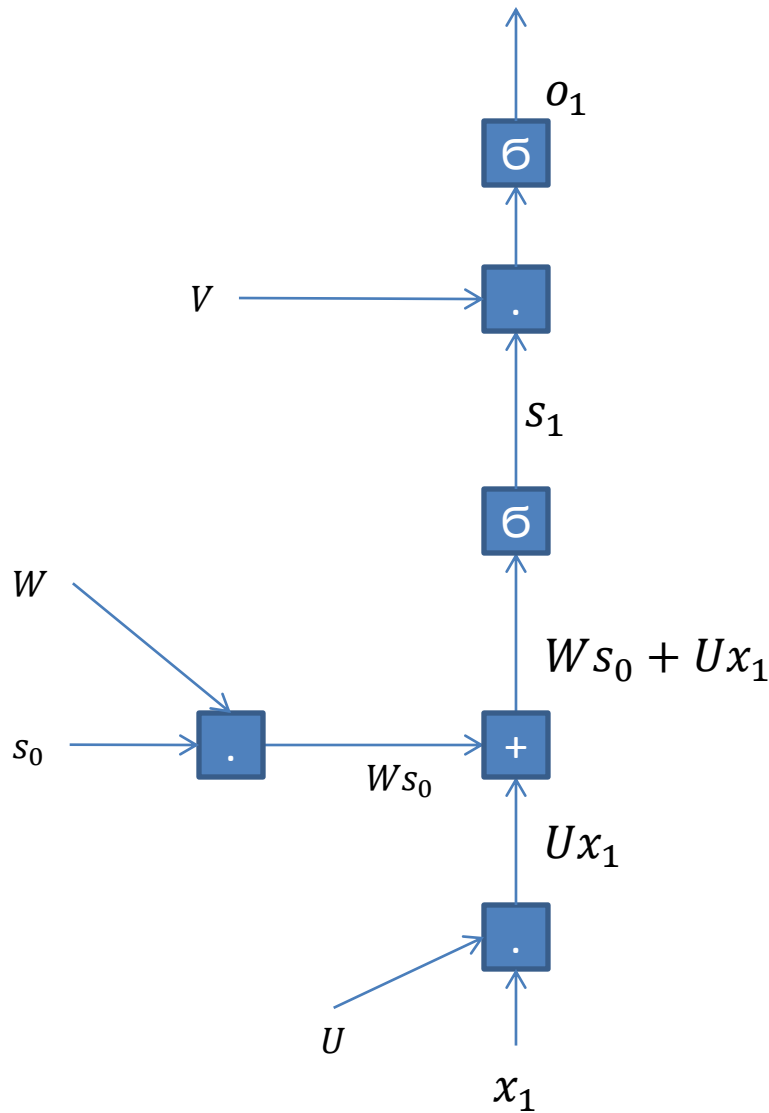$$s_t = \sigma(Ux_t + Ws_{t-1})$$
$$o_t = \sigma(Vs_t)$$

$x_t$: input
$s_t$: hidden state $(s_{-1} = 0)$
$o_t$: output
$U, V, W$: parameters (matrices)

# Recurrent Neural Net (RNN): Example



$$s_t = \sigma(Ux_t + Ws_{t-1})$$
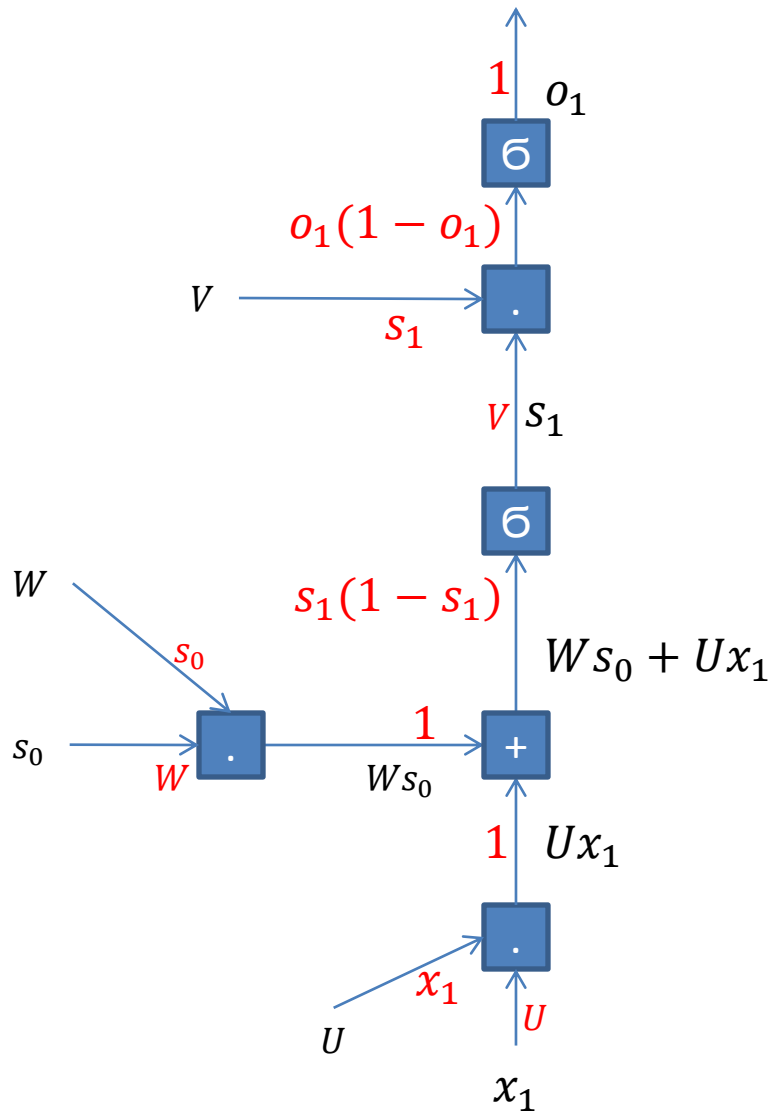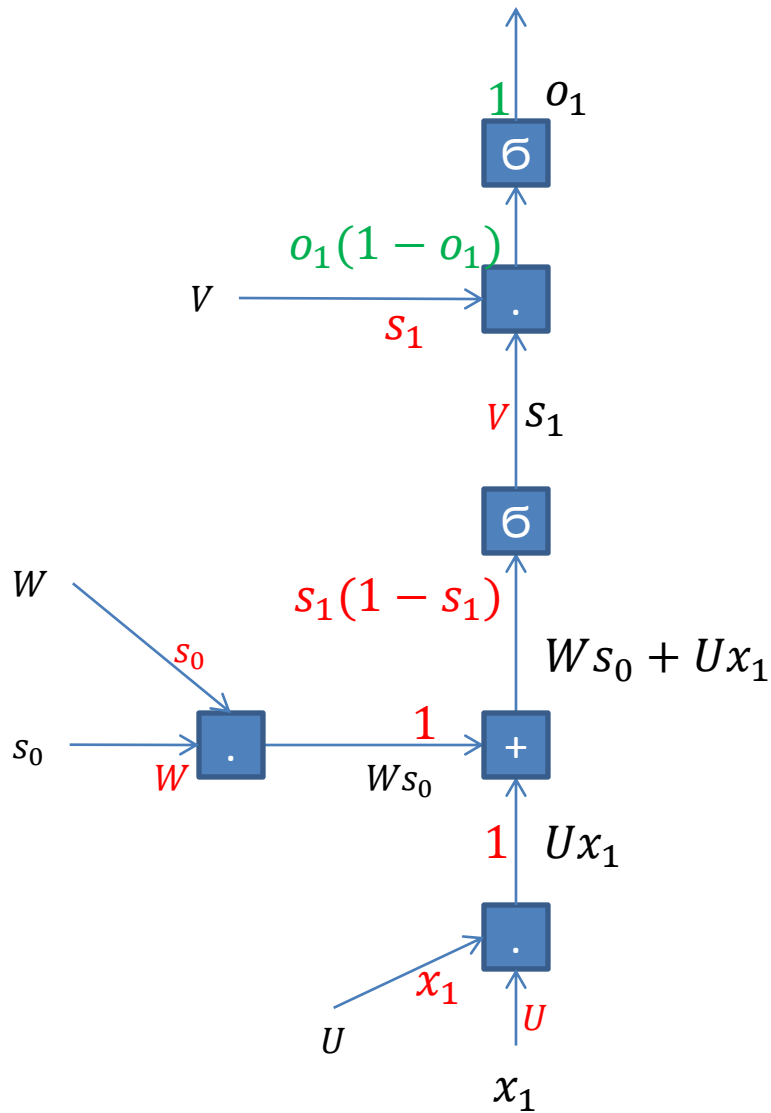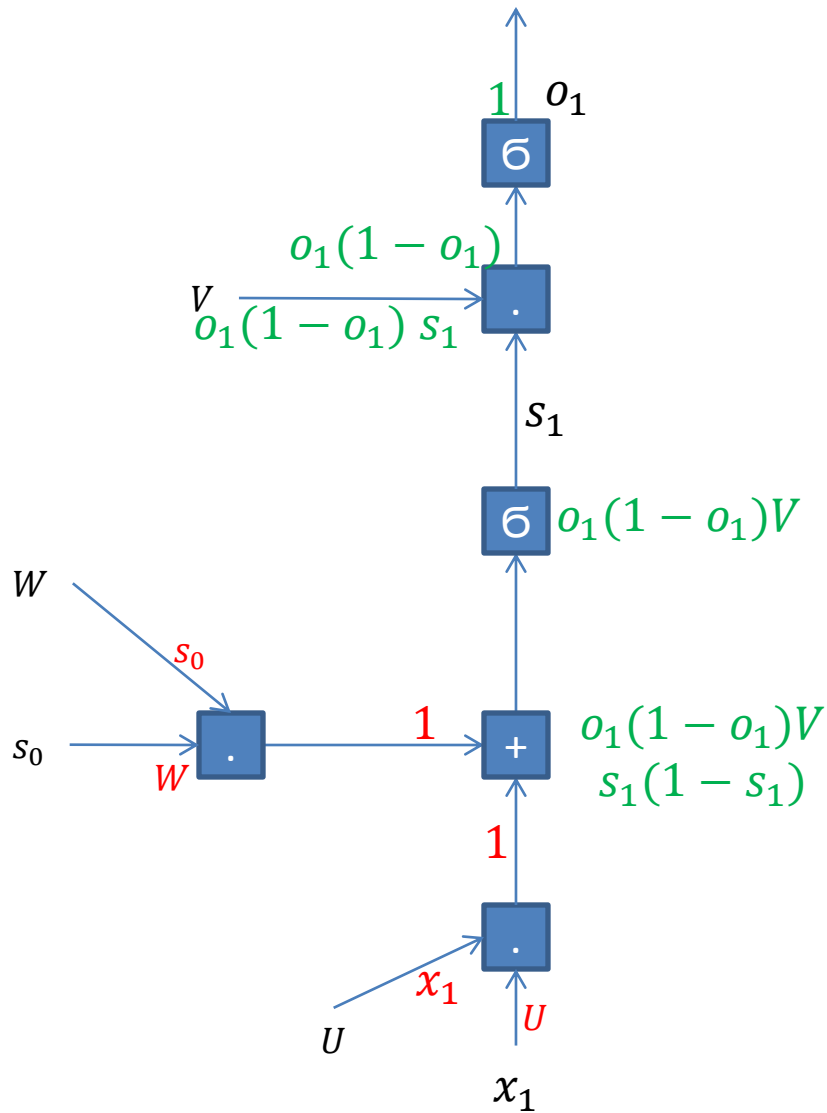
$$o_t = \sigma(Vs_t)$$

$x_t$: input
$s_t$: hidden state $(s_{-1} = 0)$
$o_t$: output
$U, V, W$: parameters

$L = \sum_t o_t$ : loss (assume)

# Recurrent Neural Net (RNN): Example



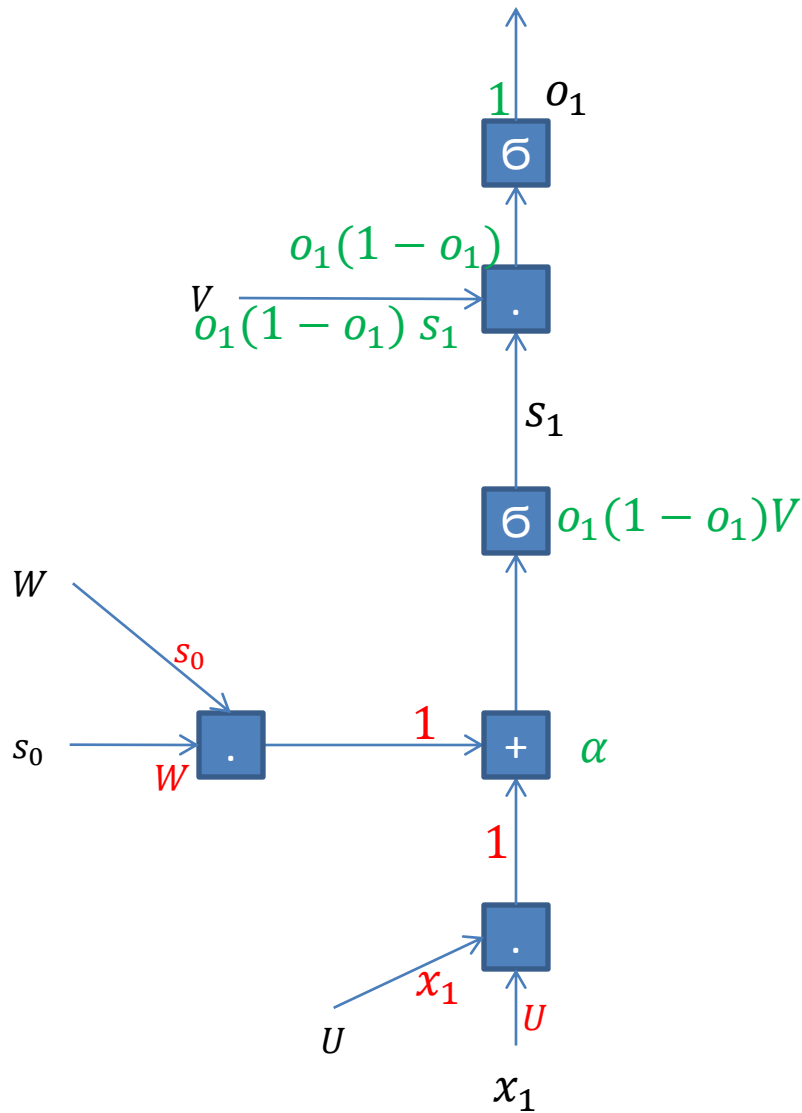$$s_t = \sigma(Ux_t + Ws_{t-1})$$

$$o_t = \sigma(Vs_t)$$

$x_t$: input
$s_t$: hidden state $(s_{-1} = 0)$
$o_t$: output
$U, V, W$: parameters

$$L = \sum_t o_t : \text{loss (assume)}$$

# Recurrent Neural Net (RNN): Example



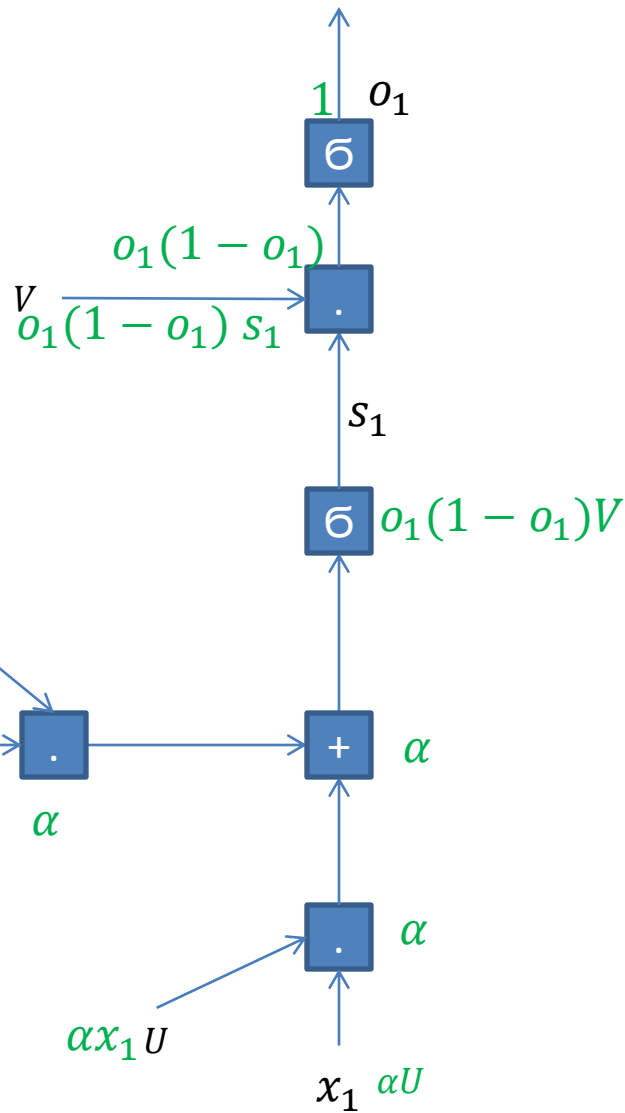$$s_t = \sigma(Ux_t + Ws_{t-1})$$
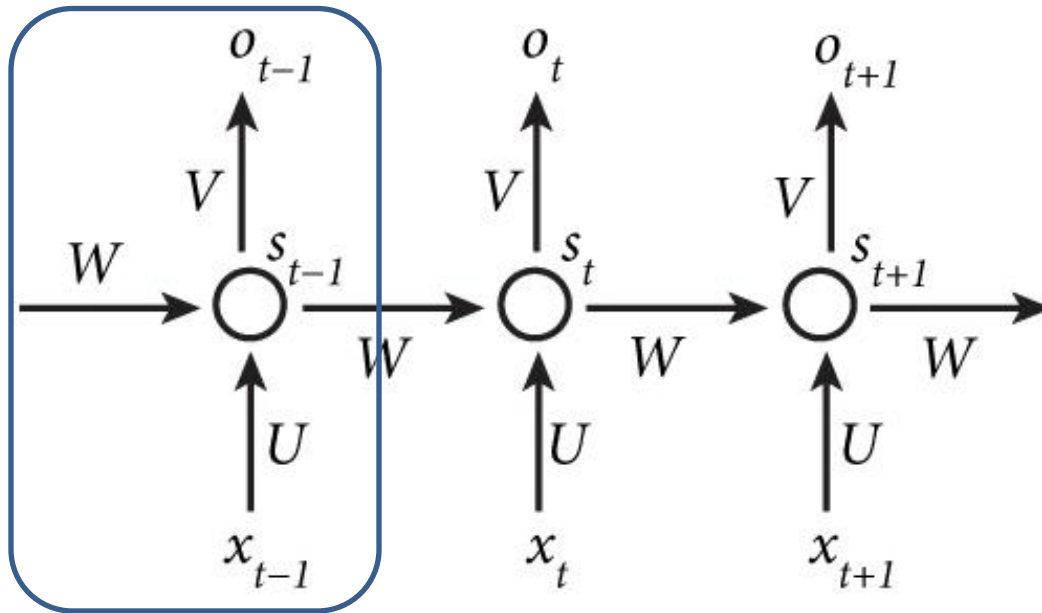
$$o_t = \sigma(Vs_t)$$

$x_t$: input
$s_t$: hidden state   $(s_{-1} = 0)$
$o_t$: output
$U, V, W$: parameters

$L = \sum_t o_t$ : loss (assume)

# Recurrent Neural Net (RNN): Example



$$s_t = \sigma(U x_t + W s_{t-1})$$

$$o_t = \sigma(V s_t)$$

$x_t$ : input

$s_t$ : hidden state   $(s_{-1} = 0)$

$o_t$ : output

$U, V, W$ : parameters

$$L = \sum_t o_t : \text{loss (assume)}$$

# Recurrent Neural Net (RNN): Example



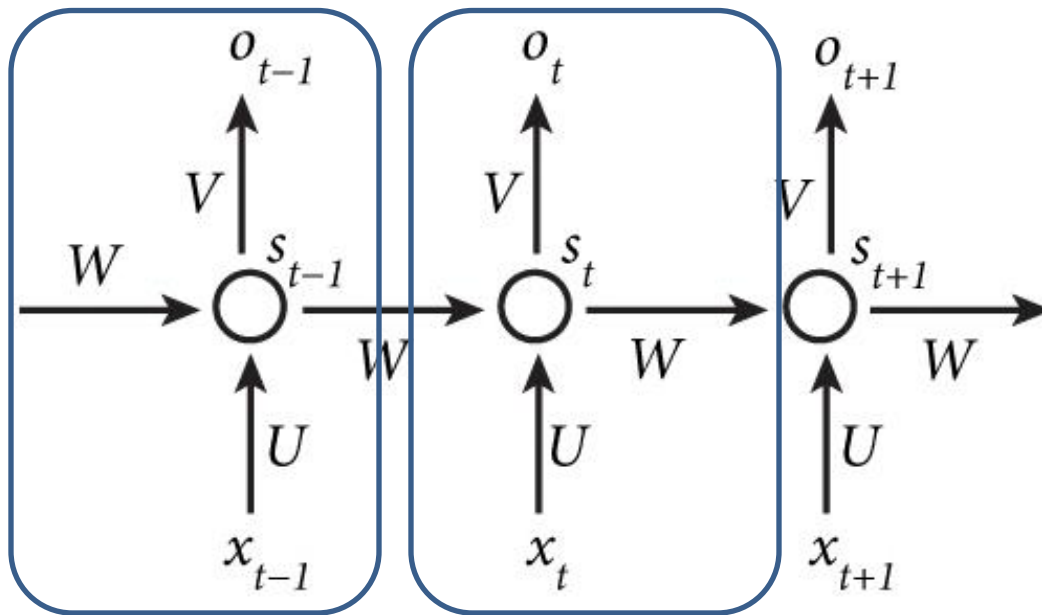$$s_t = \sigma(U x_t + W s_{t-1})$$

$$o_t = \sigma(V s_t)$$

$x_t$: input
$s_t$: hidden state   $(s_{-1} = 0)$
$o_t$: output
$U, V, W$: parameters

$L = \sum_t o_t$ : loss (assume)

# Recurrent Neural Net (RNN): Example



$$s_t = \sigma(Ux_t + Ws_{t-1})$$

$$o_t = \sigma(Vs_t)$$

$x_t$: input
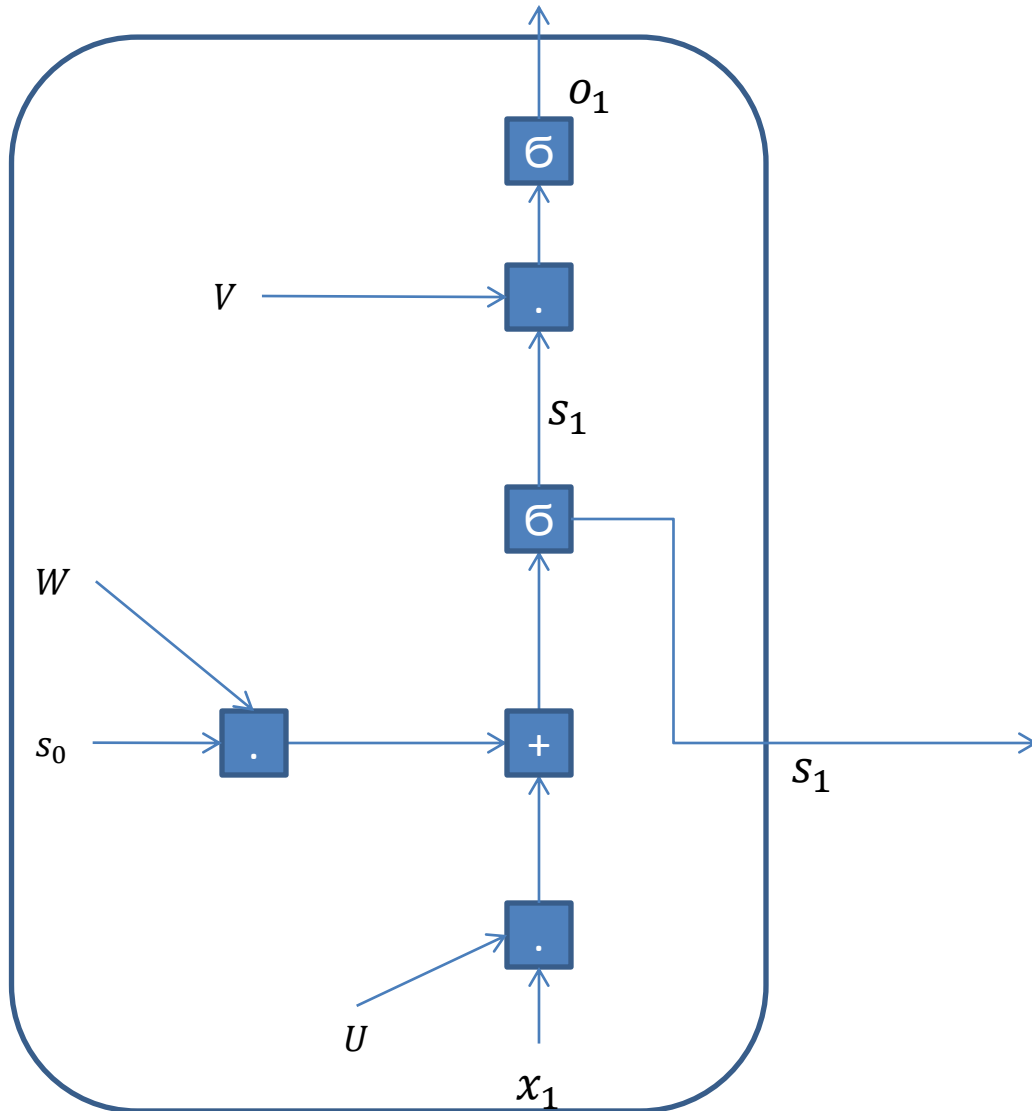$s_t$: hidden state $(s_{-1} = 0)$
$o_t$: output
$U, V, W$: parameters

$$L = \sum_t o_t : \text{loss (assume)}$$

# Recurrent Neural Net (RNN): Example



$$s_t = \sigma(Ux_t + Ws_{t-1})$$
$$o_t = \sigma(Vs_t)$$

$x_t$: input
$s_t$: hidden state  $(s_{-1} = 0)$
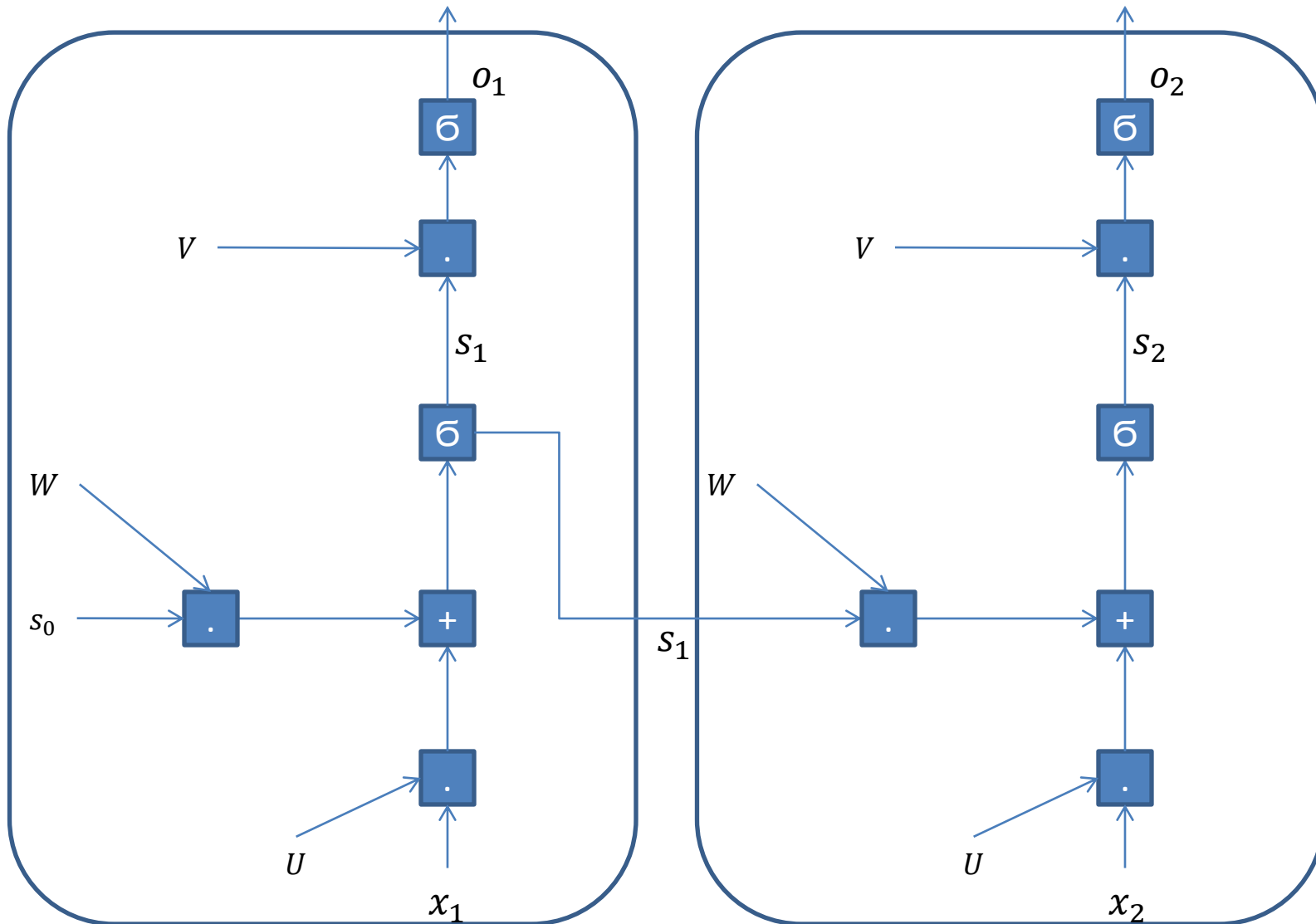$o_t$: output
$U, V, W$: parameters (matrices)

# Recurrent Neural Net (RNN): Example



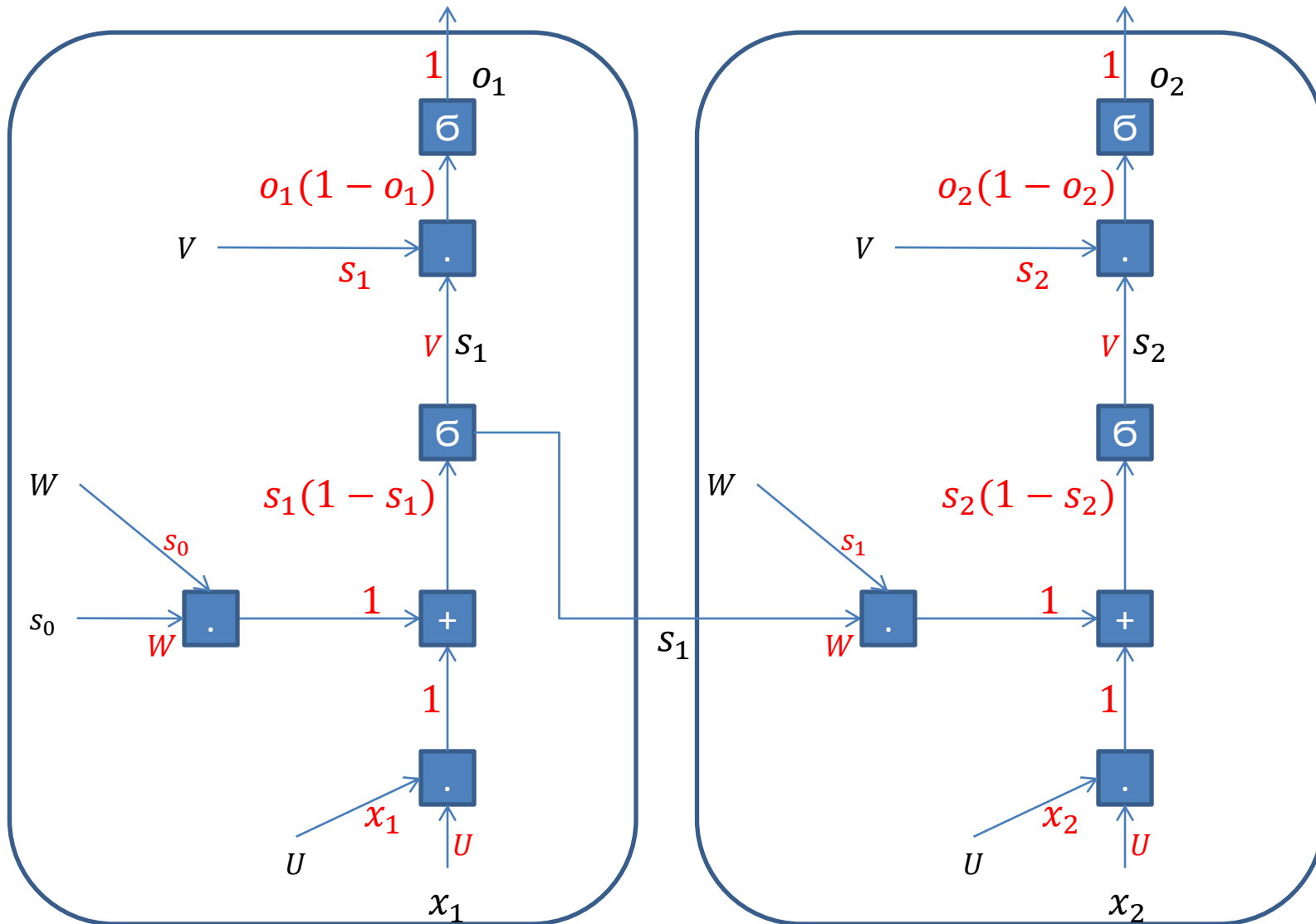$$s_t = \sigma(Ux_t + Ws_{t-1})$$

$$o_t = \sigma(Vs_t)$$

$x_t$: input
$s_t$: hidden state $\quad(s_{-1} = 0)$
$o_t$: output
$U, V, W$: parameters (matrices)

# Recurrent Neural Net (RNN): Example

$o_1$

$\sigma$

$V$ — .

$s_1$

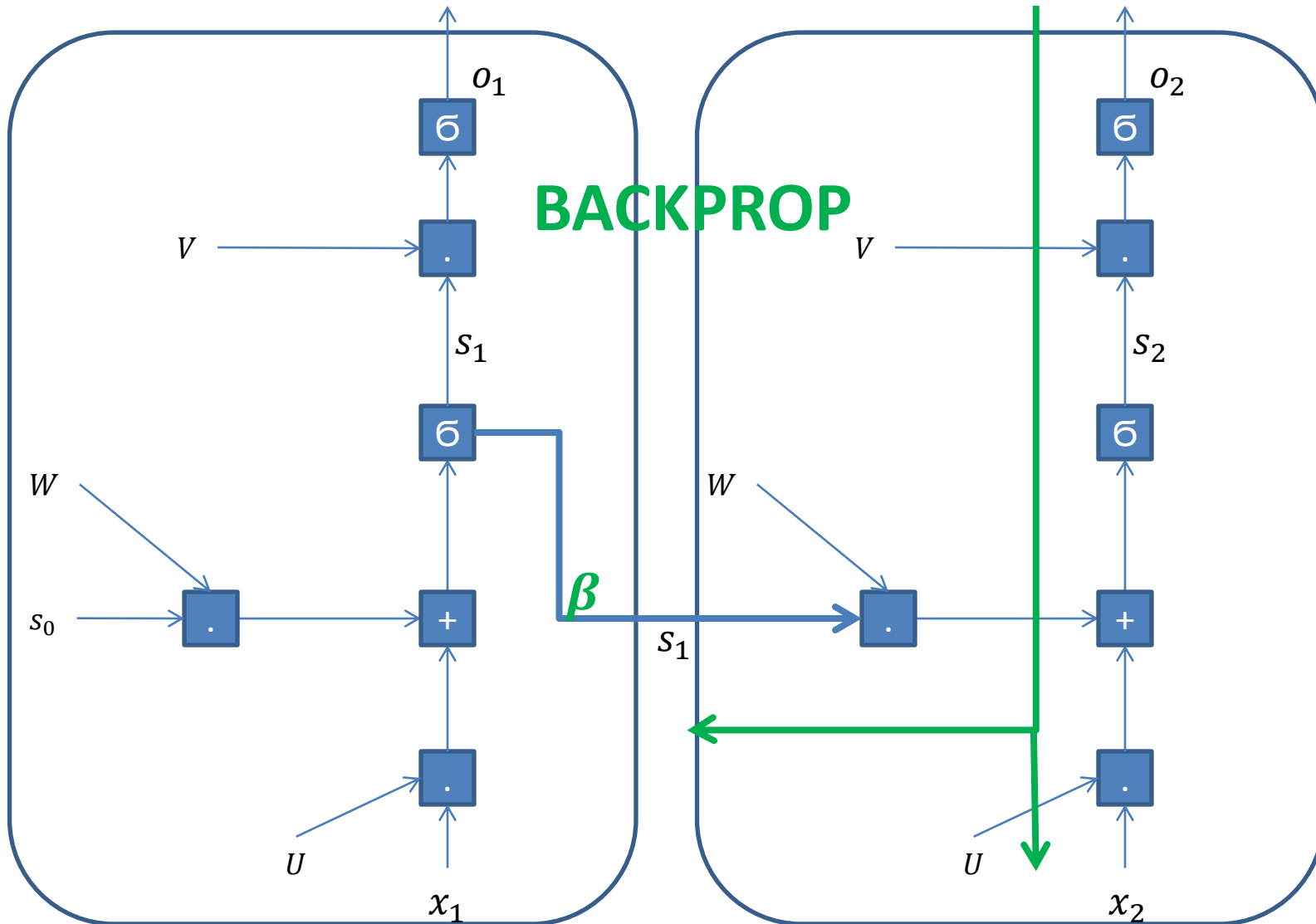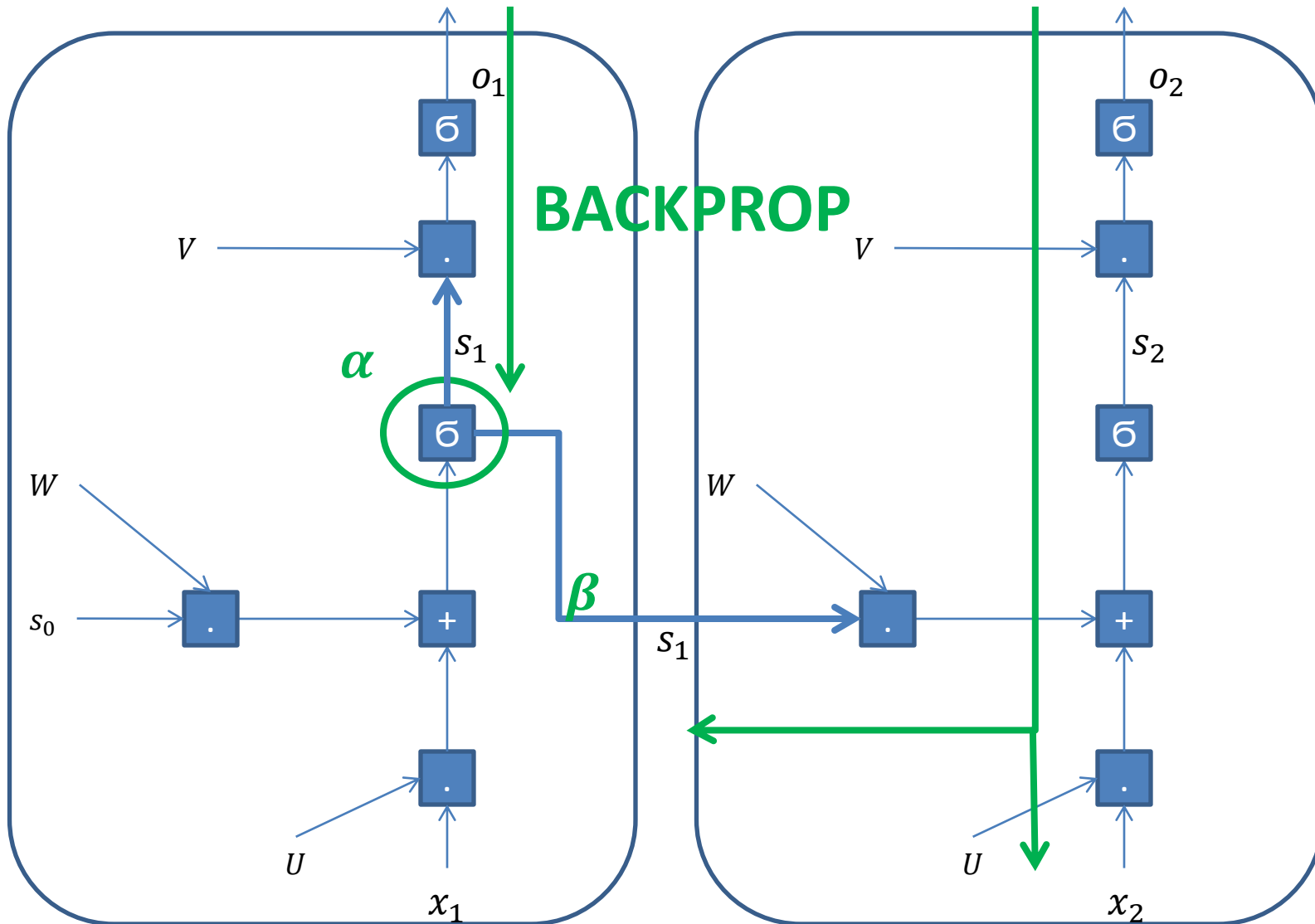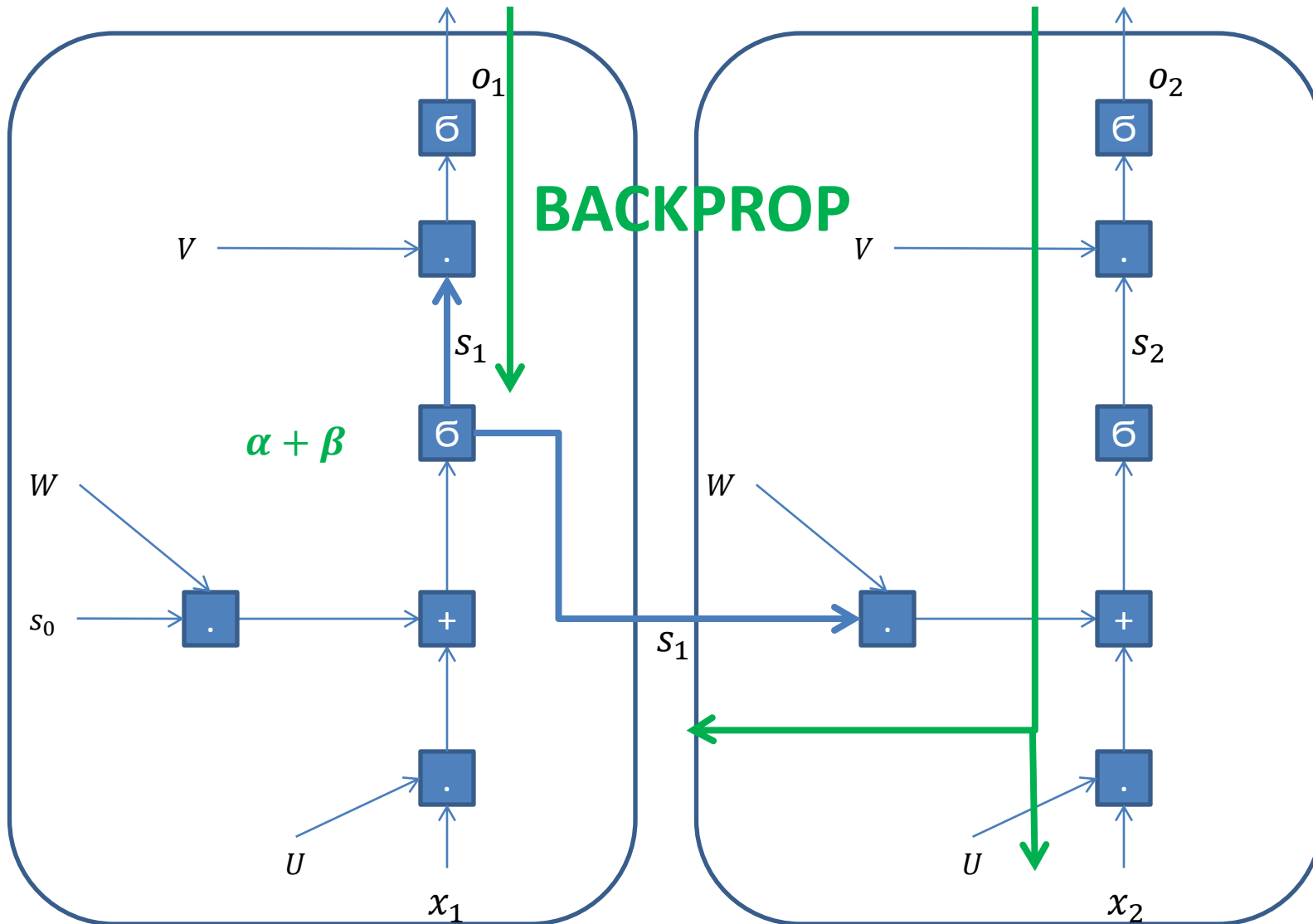$\sigma$
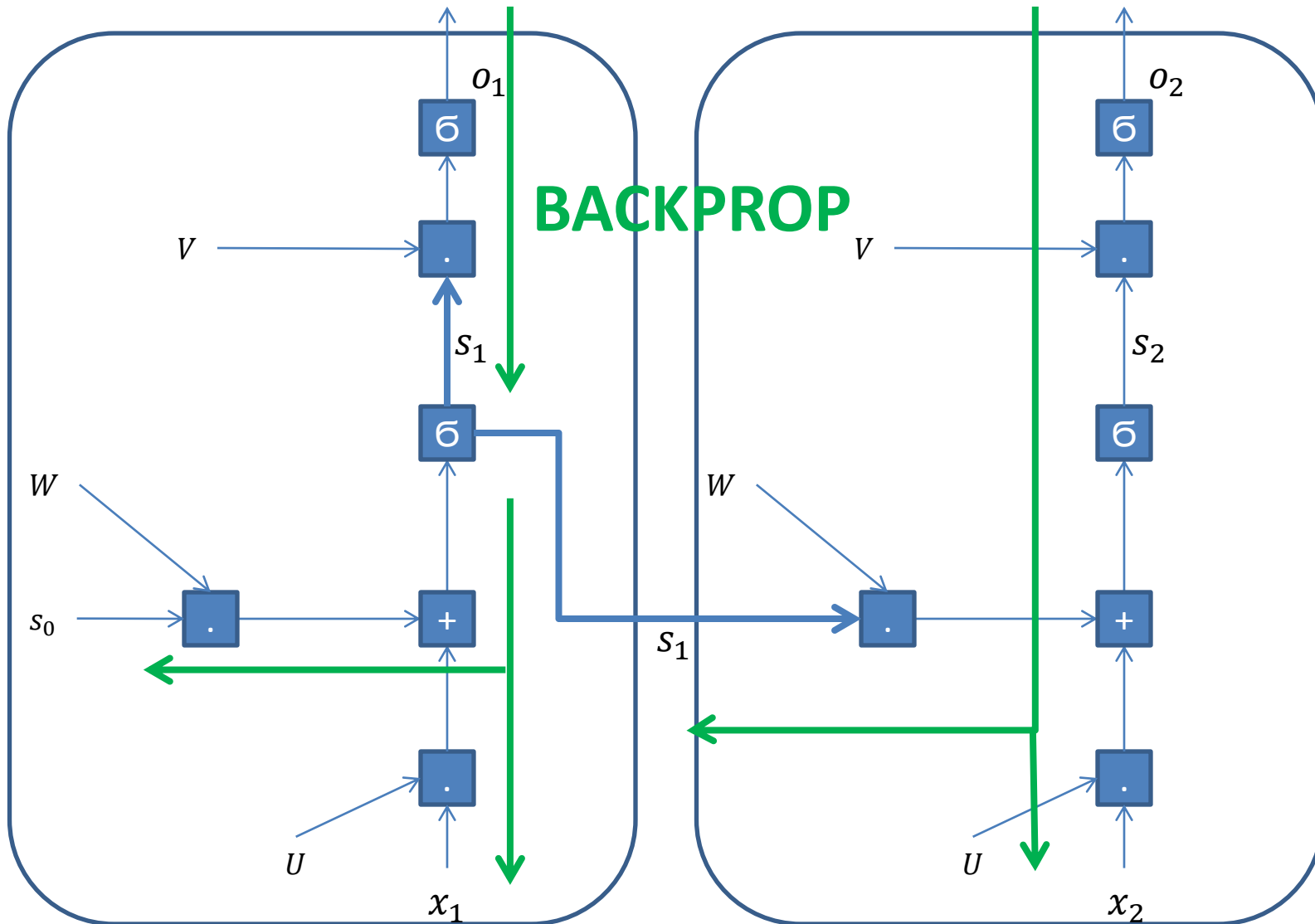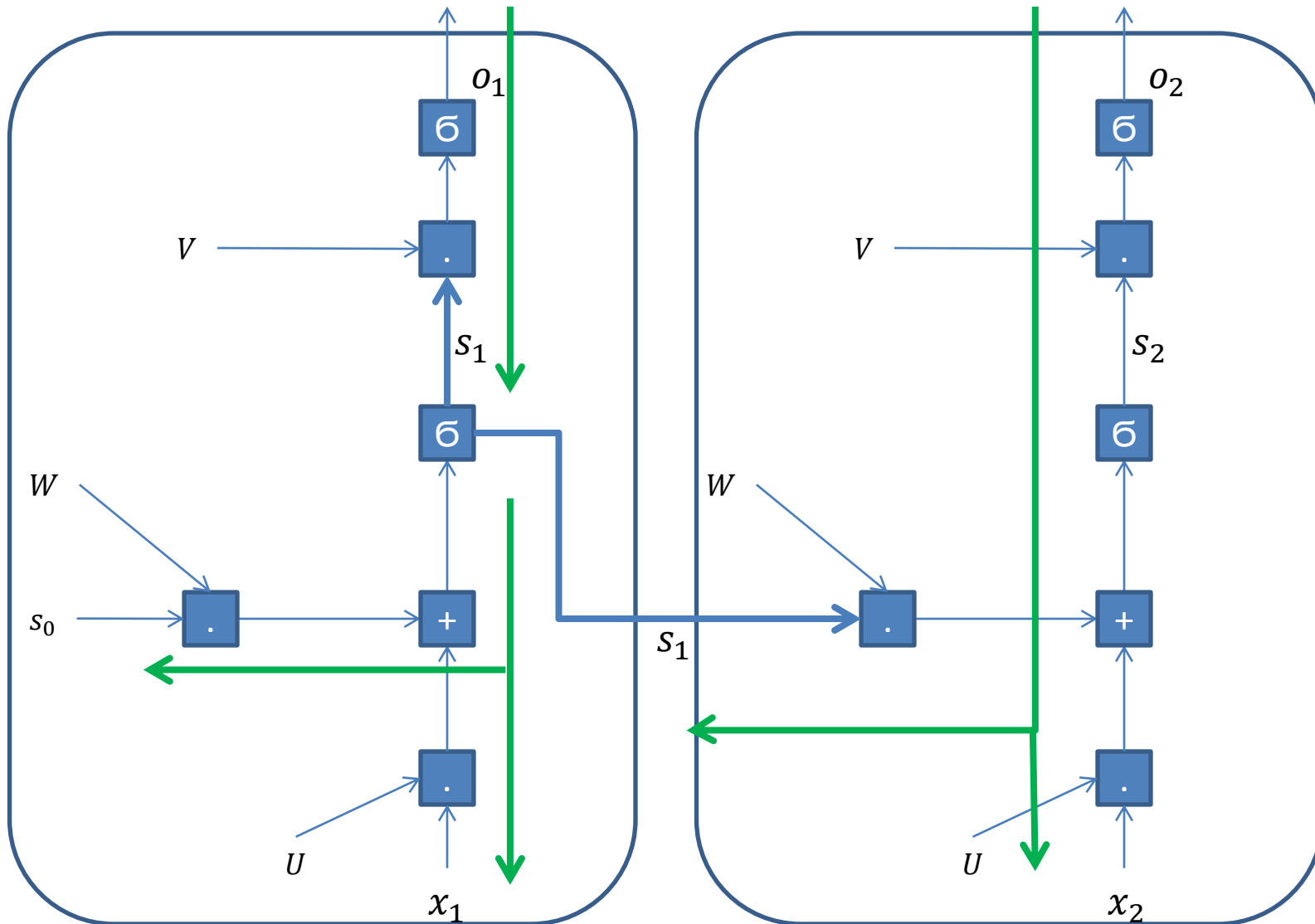
$W$

$s_0$ — .  +

$U$

$x_1$

# Recurrent Neural Net (RNN): Example

# Recurrent Neural Net (RNN): Example

# Recurrent Neural Net (RNN): Example

# Recurrent Neural Net (RNN): Example



**BACKPROP**

$o_1$

$V$

$s_1$

$W$

$s_0$

$U$

$x_1$

$\beta$

$s_1$

$o_2$

$V$

$s_2$

$W$

$U$

$x_2$

# Recurrent Neural Net (RNN): Example



**BACKPROP**

$o_1$

$V$

$\alpha$

$s_1$

$W$

$s_0$

$U$

$x_1$

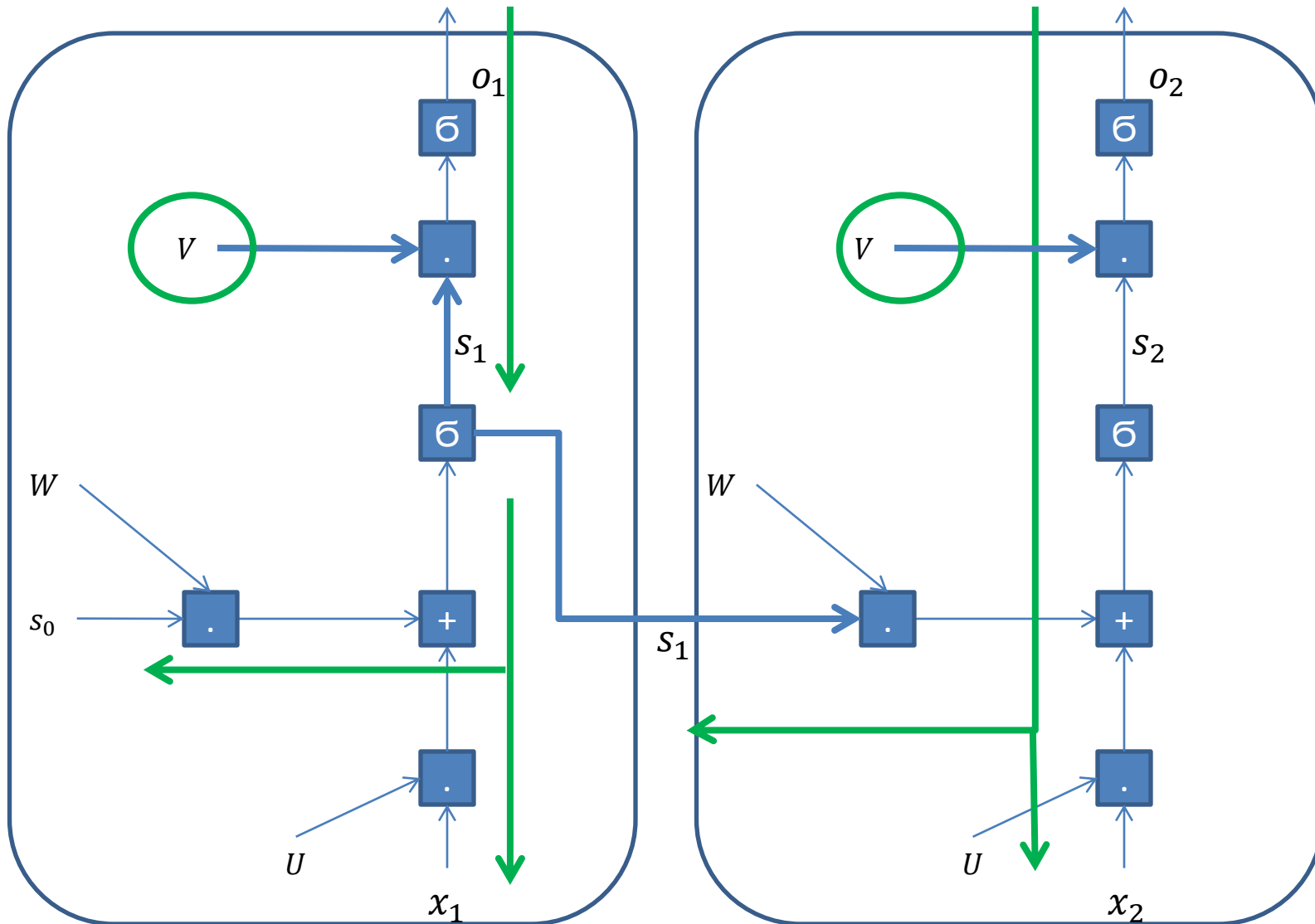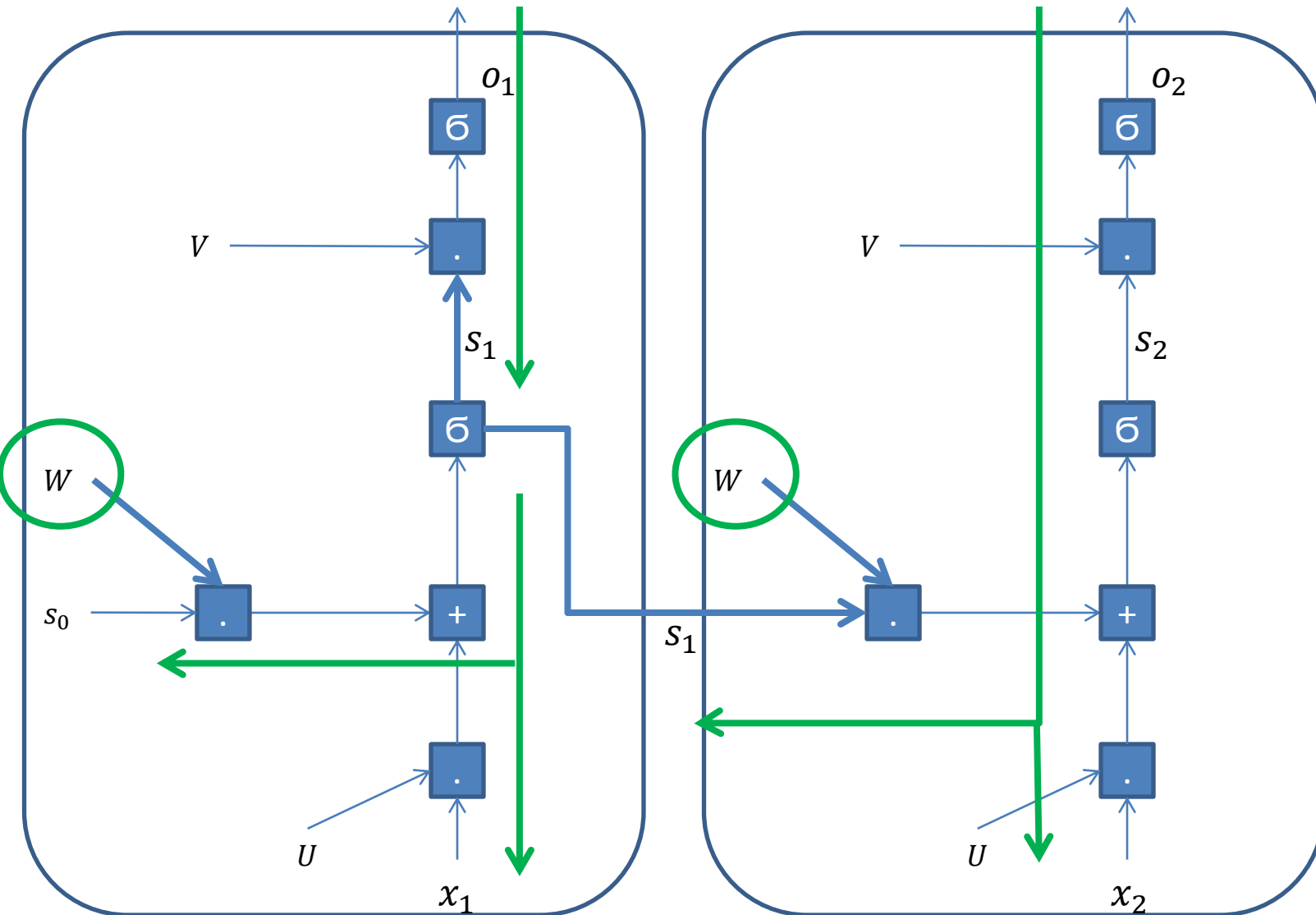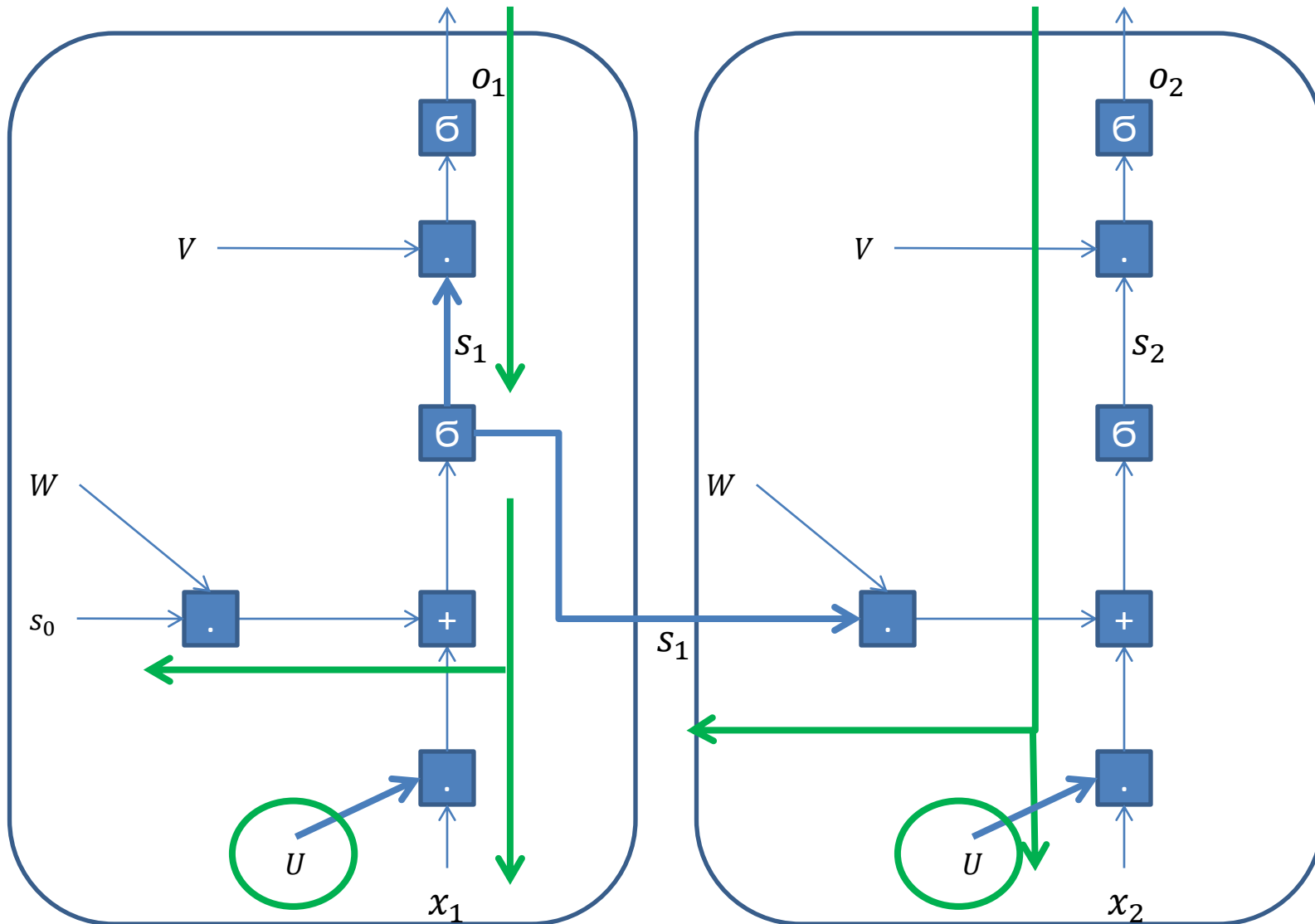$\beta$

$s_1$

$o_2$

$V$

$s_2$

$W$

$U$

$x_2$

# Recurrent Neural Net (RNN): Example

# Recurrent Neural Net (RNN): Example

# Recurrent Neural Net (RNN): Example

# Recurrent Neural Net (RNN): Example

# Recurrent Neural Net (RNN): Example

# Recurrent Neural Net (RNN): Example

# Recurrent Neural Net (RNN): Example (optional)

We assumed U,V and W were scalars

Work through the backprop when U, V and W are matrices.