# Sequence Modeling:
# Recurrent and ~~Recursive~~ Nets
# (10.1-10.5)

Jaehyoung Hong

# Recurrent neural network(RNN) is specialized for sequence-like input
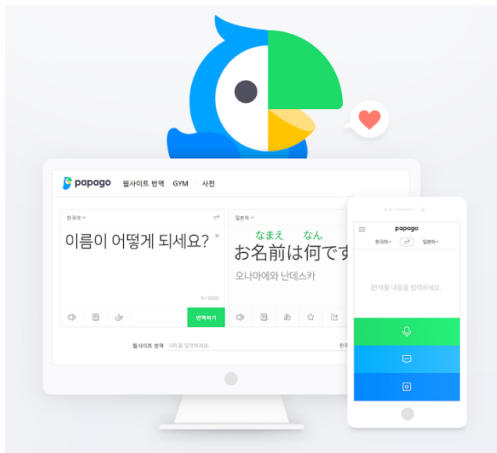
✓ CNN → Specialized for grid-like input

Image

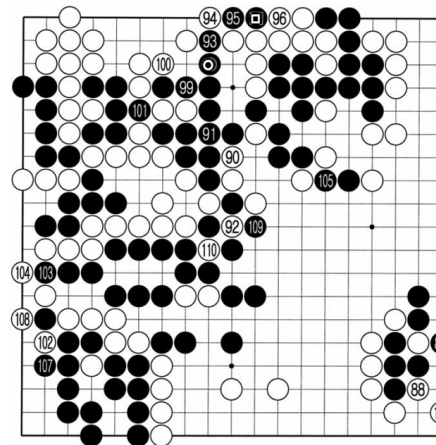| | M1 | M2 | M3 | M4 | M5 |
|---|---|---|---|---|---|
| | 1 | 3 | 2 | 5 | 4 |
| | 2 | 1 | 1 | 1 | 5 |
| | 3 | 2 | 3 | 1 | 5 |
| | 2 | 4 | 1 | 5 | 2 |

Ana    Movie 5    ★★★★☆

**Netflix problem
(Rating of movie)**

Fixed input (Given for one time-step)

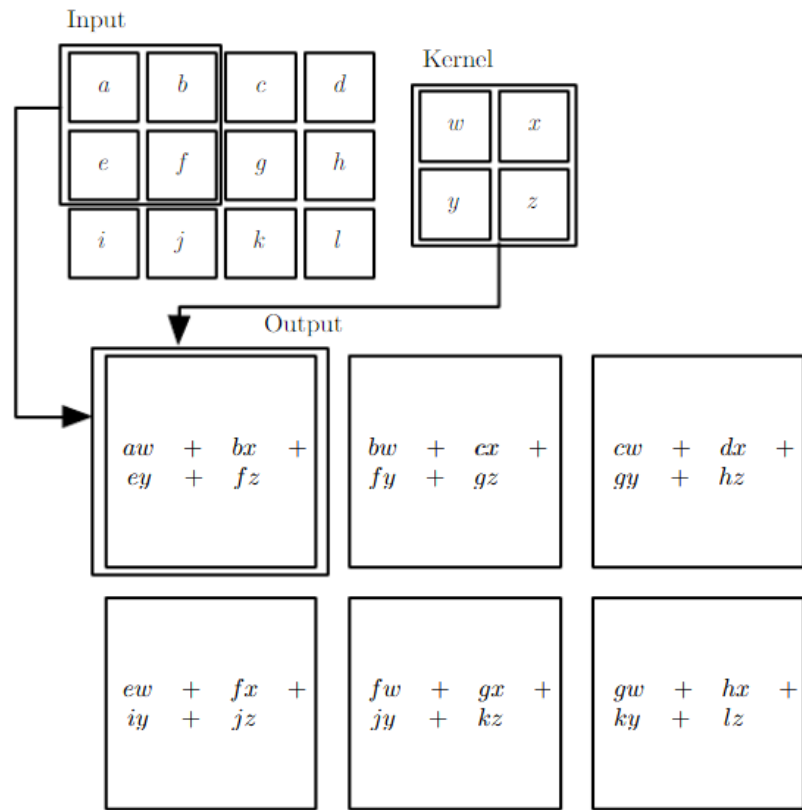✓ RNN → Specialized for sequence-like input

Language

Go

For many case, each time-point of sequence are correlated with each other
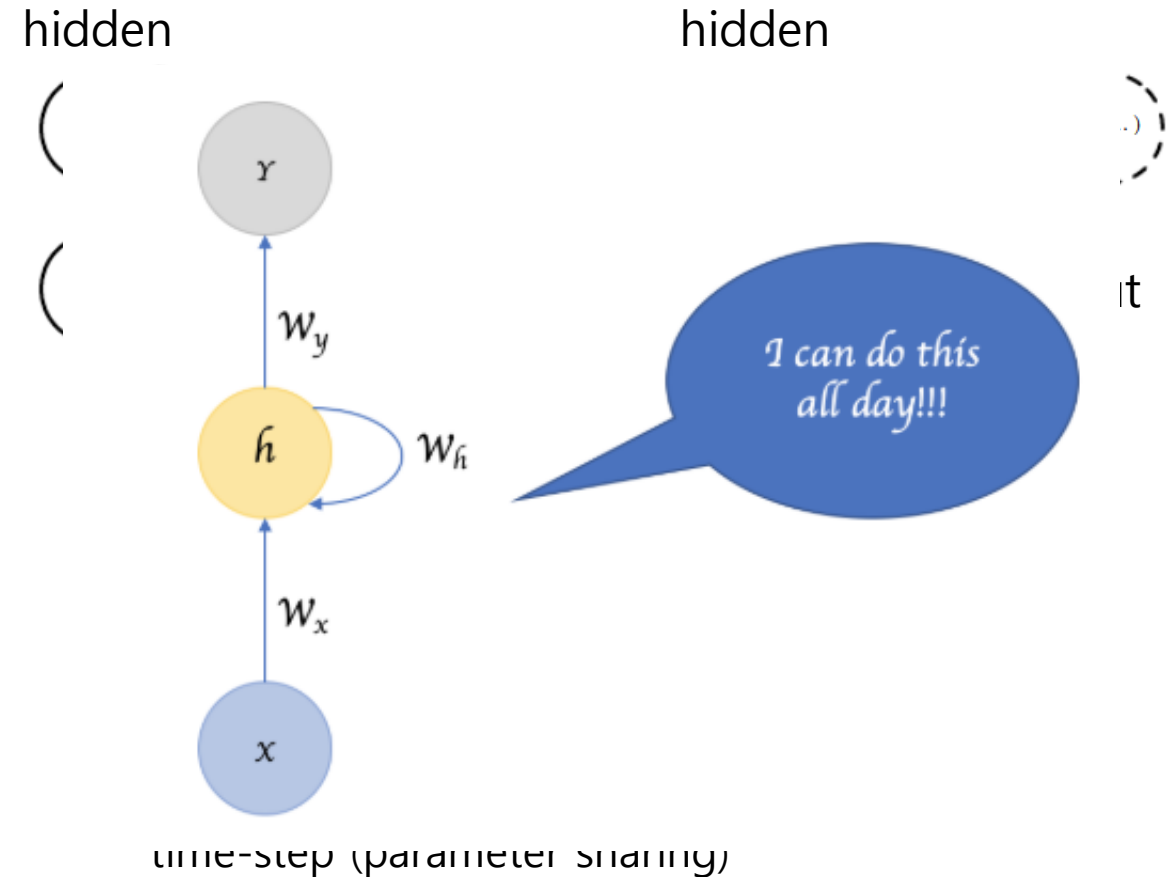
$x^{(t)}$ and $x^{(t+1)}$ are correlated

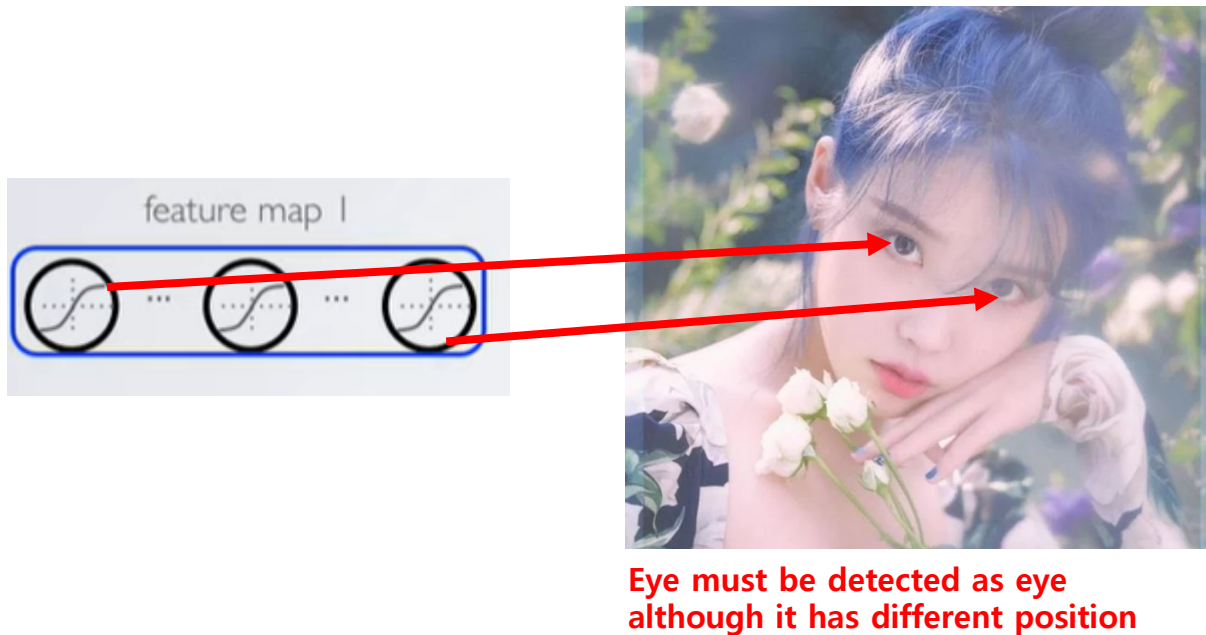# Like CNN, RNN can handle various length of sequence by using parameter sharing



CNN : Using same kernel (parameter sharing)

Handle various size of image

Handle various length of image

# By parameter sharing, RNN is invariant under position



feature map I

**Eye must be detected as eye although it has different position**

*I went to Nepal in 2009*

*In 2009, I went to Nepal*

CNN : Using same kernel (parameter sharing)

Same input gives same feature

RNN : Using same weight (parameter) for each time-step (parameter sharing)

Same word gives same meaning

# Description of RNN

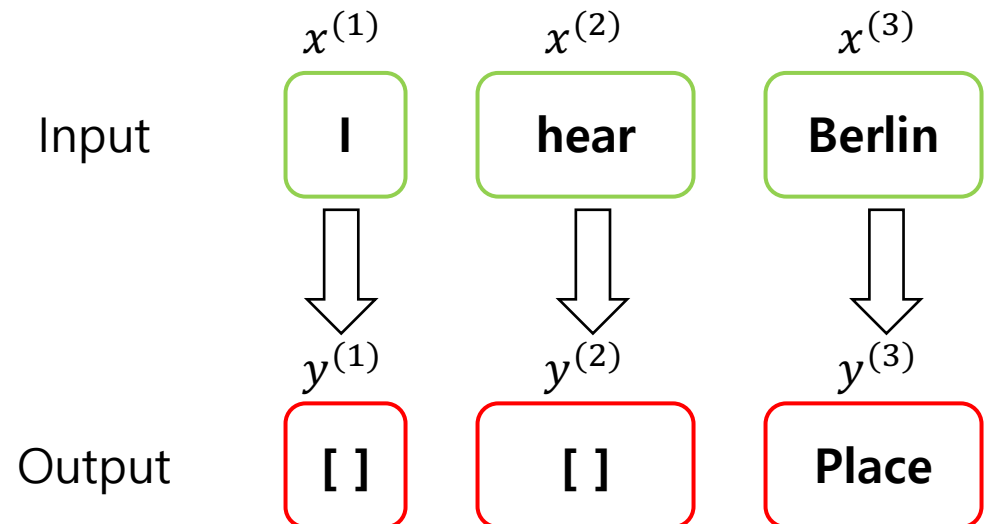# First, we consider sequence-to-sequence RNN

| Type of RNN | Illustration | Example |
|---|---|---|
| Many-to-many $T_x = T_y$ |  | Name entity recognition |

I hear **Berlin** is wonderful in the **winter**

Place — Time

<Name entity recognition>

| | $x^{(1)}$ | $x^{(2)}$ | $x^{(3)}$ |
|---|---|---|---|
| Input | **I** | **hear** | **Berlin** |
| | $y^{(1)}$ | $y^{(2)}$ | $y^{(3)}$ |
| Output | **[ ]** | **[ ]** | **Place** |

# Such RNN can be defined as recurrent form from $t = 1$ to $t = \tau$



Output

Input

Previous information influenced on current information only through hidden state ($h$)

RNN

$$t = 1 \text{ to } t = \tau$$

$$
\begin{aligned}
\boldsymbol{a}^{(t)} &= \boldsymbol{b} + \boldsymbol{W}\boldsymbol{h}^{(t-1)} + \boldsymbol{U}\boldsymbol{x}^{(t)}, \\
\boldsymbol{h}^{(t)} &= \tanh(\boldsymbol{a}^{(t)}), \\
\boldsymbol{o}^{(t)} &= \boldsymbol{c} + \boldsymbol{V}\boldsymbol{h}^{(t)}, \\
\hat{\boldsymbol{y}}^{(t)} &= \mathrm{softmax}(\boldsymbol{o}^{(t)}),
\end{aligned}
$$

Same parameter!!

Loss

$$
\begin{aligned}
& L\left(\{\boldsymbol{x}^{(1)},\ldots,\boldsymbol{x}^{(\tau)}\}, \{\boldsymbol{y}^{(1)},\ldots,\boldsymbol{y}^{(\tau)}\}\right) \\
&= \sum_t L^{(t)} \\
&= -\sum_t \log p_{\mathrm{model}}\left(y^{(t)} \mid \{\boldsymbol{x}^{(1)},\ldots,\boldsymbol{x}^{(t)}\}\right)
\end{aligned}
$$

Total loss is sum of losses over all the time steps

Gradient

$$
\frac{\partial \mathcal{L}^{(T)}}{\partial W} = \sum_{t=1}^{T} \frac{\partial \mathcal{L}^{(T)}}{\partial W}\bigg|_{(t)}
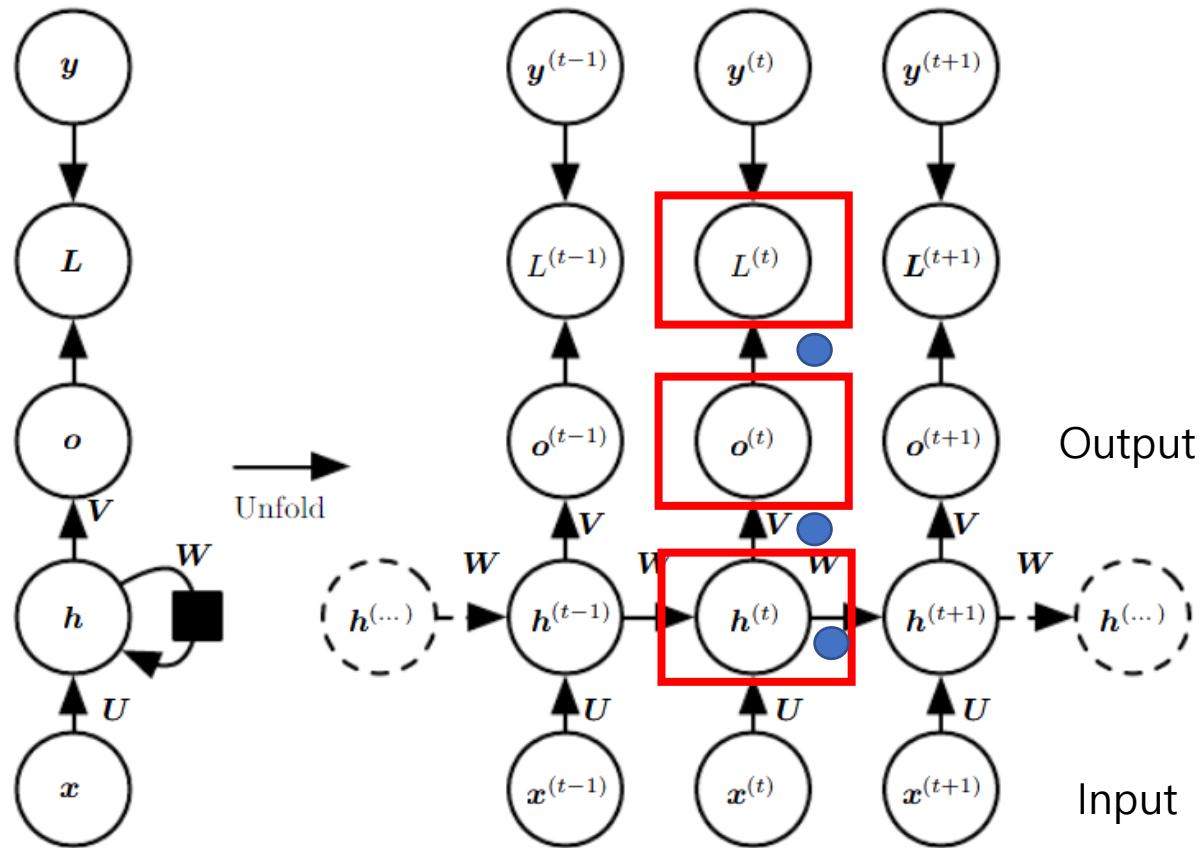$$

Runtime : $O(\tau)$

Memory : $O(\tau)$

(Backpropagation through time : BPTT)

# Example of calculation of BPTT



$$\begin{aligned} \boldsymbol{a}^{(t)} &= \boldsymbol{b} + \boldsymbol{W}\boldsymbol{h}^{(t-1)} + \boldsymbol{U}\boldsymbol{x}^{(t)}, \\ \boldsymbol{h}^{(t)} &= \tanh(\boldsymbol{a}^{(t)}), \\ \boldsymbol{o}^{(t)} &= \boldsymbol{c} + \boldsymbol{V}\boldsymbol{h}^{(t)}, \\ \hat{\boldsymbol{y}}^{(t)} &= \operatorname{softmax}(\boldsymbol{o}^{(t)}), \end{aligned}$$
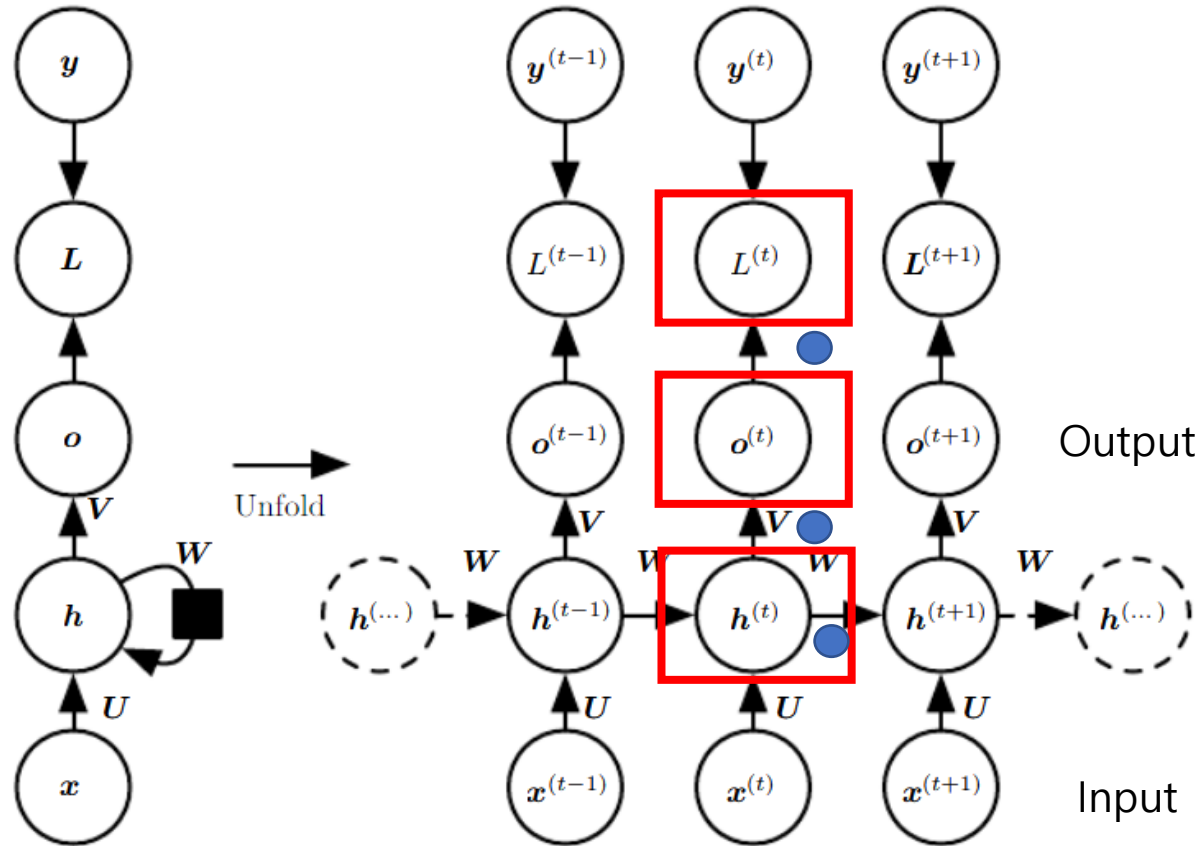
**Follow the arrow**

(1) $\dfrac{\partial L}{\partial L^{(t)}} = 1.$

(2) $(\nabla_{\boldsymbol{o}^{(t)}} L)_i = \dfrac{\partial L}{\partial o_i^{(t)}} = \dfrac{\partial L}{\partial L^{(t)}} \dfrac{\partial L^{(t)}}{\partial o_i^{(t)}} = \hat{y}_i^{(t)} - \mathbf{1}_{i=y^{(t)}}.$

(3) $\nabla_{\boldsymbol{h}^{(\tau)}} L = \boldsymbol{V}^{\top} \nabla_{\boldsymbol{o}^{(\tau)}} L.$

$$\nabla_{\boldsymbol{h}^{(t)}} L = \left(\dfrac{\partial \boldsymbol{h}^{(t+1)}}{\partial \boldsymbol{h}^{(t)}}\right)^{\top} (\nabla_{\boldsymbol{h}^{(t+1)}} L) + \left(\dfrac{\partial \boldsymbol{o}^{(t)}}{\partial \boldsymbol{h}^{(t)}}\right)^{\top} (\nabla_{\boldsymbol{o}^{(t)}} L)$$

$$= \boldsymbol{W}^{\top} \operatorname{diag}\left(1 - \left(\boldsymbol{h}^{(t+1)}\right)^2\right) (\nabla_{\boldsymbol{h}^{(t+1)}} L) + \boldsymbol{V}^{\top} (\nabla_{\boldsymbol{o}^{(t)}} L)$$

Output

Input

# Example of calculation of BPTT

$$
\begin{aligned}
\boldsymbol{a}^{(t)} &= \boldsymbol{b} + \boldsymbol{W}\boldsymbol{h}^{(t-1)} + \boldsymbol{U}\boldsymbol{x}^{(t)}, \\
\boldsymbol{h}^{(t)} &= \tanh(\boldsymbol{a}^{(t)}), \\
\boldsymbol{o}^{(t)} &= \boldsymbol{c} + \boldsymbol{V}\boldsymbol{h}^{(t)}, \\
\hat{\boldsymbol{y}}^{(t)} &= \mathrm{softmax}(\boldsymbol{o}^{(t)}),
\end{aligned}
$$

Follow the arrow

$$
\begin{aligned}
\nabla_{\boldsymbol{c}}L &= \sum_t \left(\frac{\partial \boldsymbol{o}^{(t)}}{\partial \boldsymbol{c}}\right)^{\top} \nabla_{\boldsymbol{o}^{(t)}}L = \sum_t \nabla_{\boldsymbol{o}^{(t)}}L, \\
\nabla_{\boldsymbol{b}}L &= \sum_t \left(\frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{b}^{(t)}}\right)^{\top} \nabla_{\boldsymbol{h}^{(t)}}L = \sum_t \mathrm{diag}\left(1 - \left(\boldsymbol{h}^{(t)}\right)^2\right) \nabla_{\boldsymbol{h}^{(t)}}L, \\
\nabla_{\boldsymbol{V}}L &= \sum_t \sum_i \left(\frac{\partial L}{\partial o_i^{(t)}}\right) \nabla_{\boldsymbol{V}^{(t)}} o_i^{(t)} = \sum_t (\nabla_{\boldsymbol{o}^{(t)}}L) \boldsymbol{h}^{(t)\top}, \\
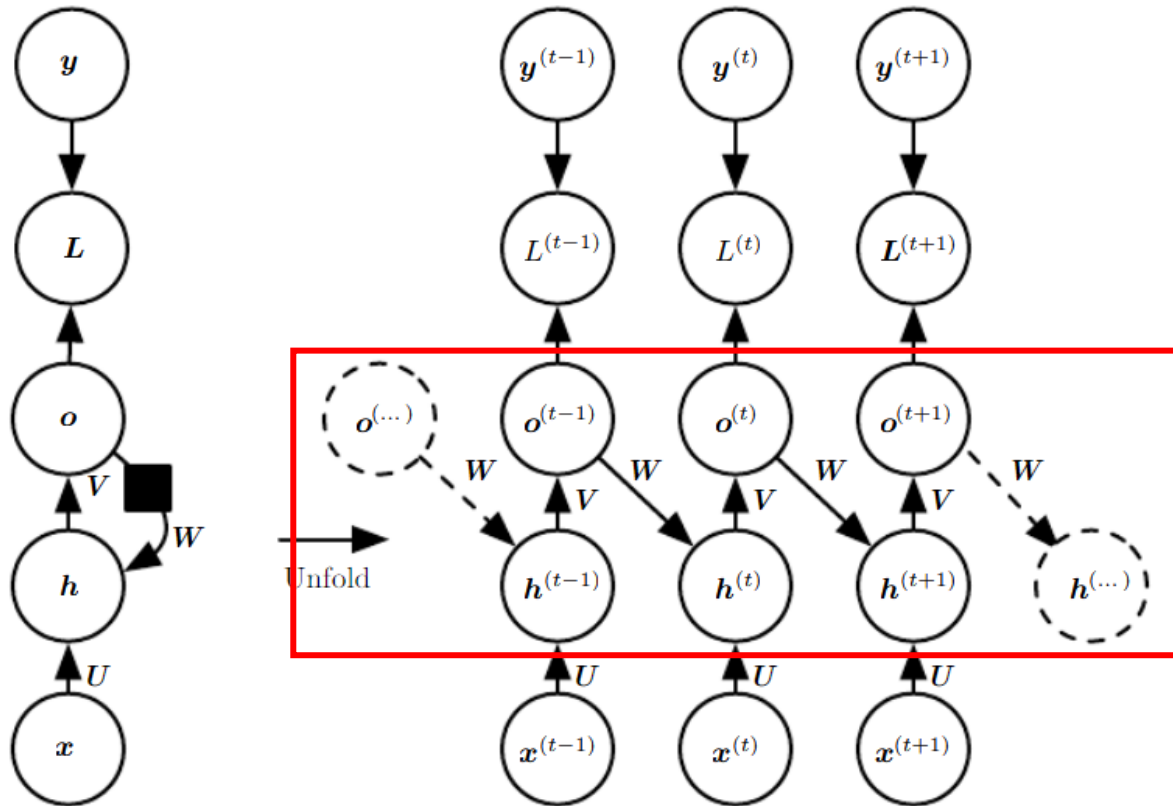\nabla_{\boldsymbol{W}}L &= \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}}\right) \nabla_{\boldsymbol{W}^{(t)}} h_i^{(t)} \\
&= \sum_t \mathrm{diag}\left(1 - \left(\boldsymbol{h}^{(t)}\right)^2\right)(\nabla_{\boldsymbol{h}^{(t)}}L) \boldsymbol{h}^{(t-1)\top}, \\
\nabla_{\boldsymbol{U}}L &= \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}}\right) \nabla_{\boldsymbol{U}^{(t)}} h_i^{(t)} \\
&= \sum_t \mathrm{diag}\left(1 - \left(\boldsymbol{h}^{(t)}\right)^2\right)(\nabla_{\boldsymbol{h}^{(t)}}L) \boldsymbol{x}^{(t)\top},
\end{aligned}
$$

**Bad Result**

$$
W = W^{(1)} = \cdots = W^{(\tau)},\ U = U^{(1)} = \cdots = U^{(\tau)},\ V = V^{(1)} = \cdots = V^{(\tau)}
$$

Output

Input

# We can use previous information from output to better trained the RNN (Teacher forcing)



$$\log p\left(\boldsymbol{y}^{(1)}, \boldsymbol{y}^{(2)} \mid \boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}\right)$$

$$= \log p\left(\boldsymbol{y}^{(2)} \mid \boldsymbol{y}^{(1)}, \boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}\right) + \log p\left(\boldsymbol{y}^{(1)} \mid \boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}\right).$$

Generally less powerful :
Output has less information than hidden

Teacher forcing : Use output as next input
Avoid BPTT !! (No hidden-to-hidden)

# We can use previous information from output to better trained the RNN (Teacher forcing)



Train time      Test time

Teacher forcing : Use output as next input

**Want to train?**

"Mary had a little lamb whose fleece was white as snow"

| $X$ | $\widehat{y}$ |
|---|---|
| [Start] | "a" |
| [Start], "a" | ?<br>(Any word can<br>get punished) |

<Without teacher forcing : Slow & incorrect>

| $X$ | $\widehat{y}$ |
|---|---|
| [Start] | "a" |
| [Start], "Mary" | ? |

With teacher forcing

# RNN which takes single vector as input



When input is $x$ rather than $x^{(t)}$

1. as an extra input at each time step, or

2. as the initial state $h^{(0)}$, or

3. both.

$x^T R$ is added as additional input to the hidden units

⇩

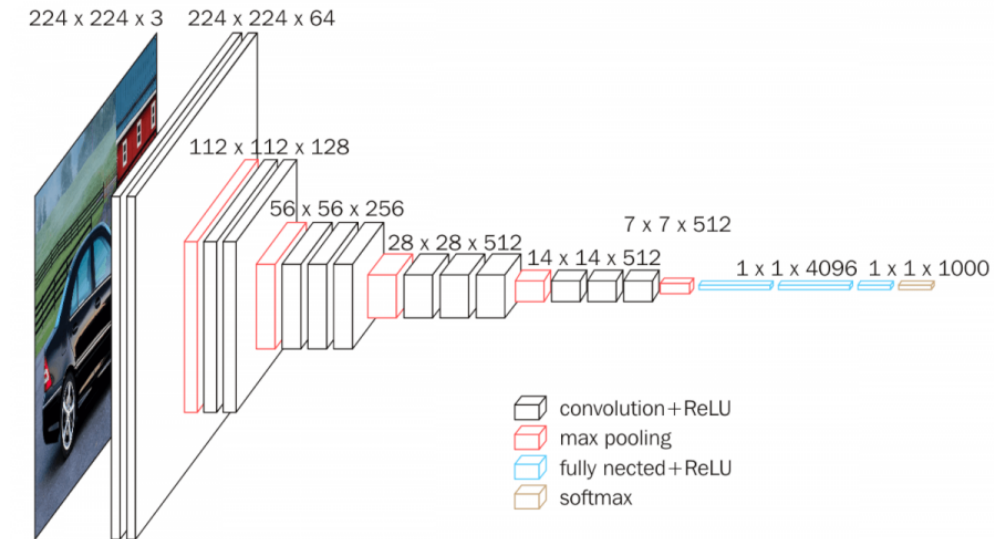$x^T R$ is new bias parameter

$$P(y|x) \rightarrow P(y|w), w = f(x)$$

# Variant type of RNN : Depending on Input and Output

| Type of RNN | Illustration | Example |
|---|---|---|
| One-to-one $T_x = T_y = 1$ |  | Traditional neural network |
| Many-to-many $T_x = T_y$ |  | Name entity recognition |





I hear Berlin is wonderful in the winter

# Variant type of RNN : Depending on Input and Output

# Various type of RNN : Bidirectional RNNs using future information



**Example**

✓ Speech recognition

✓ Handwriting recognition

✓ Bioinformatics

Vanilla RNN
$$(\text{Causal}: x^{(t)} \leftarrow x^{(1)}, \dots, x^{(t-1)})$$

vs

Bidirectional RNN
$$(\text{Around } t: x^{(t)} \leftarrow x^{(t+1)}, x^{(t-1)}, \dots)$$

❏ **Attention model** — This model allows an RNN to pay attention to specific parts of the input that is considered as being important, which improves the performance of the resulting model in practice. By noting $\alpha^{<t,t'>}$ the amount of attention that the output $y^{<t>}$ should pay to the activation $a^{<t'>}$ and $c^{<t>}$ the context at time $t$, we have:

$$c^{<t>} = \sum_{t'} \alpha^{<t,t'>} a^{<t'>} \qquad \text{with} \qquad \sum_{t'} \alpha^{<t,t'>} = 1$$

*Remark: the attention scores are commonly used in image captioning and machine translation.*

*Encoder*

*Deco*


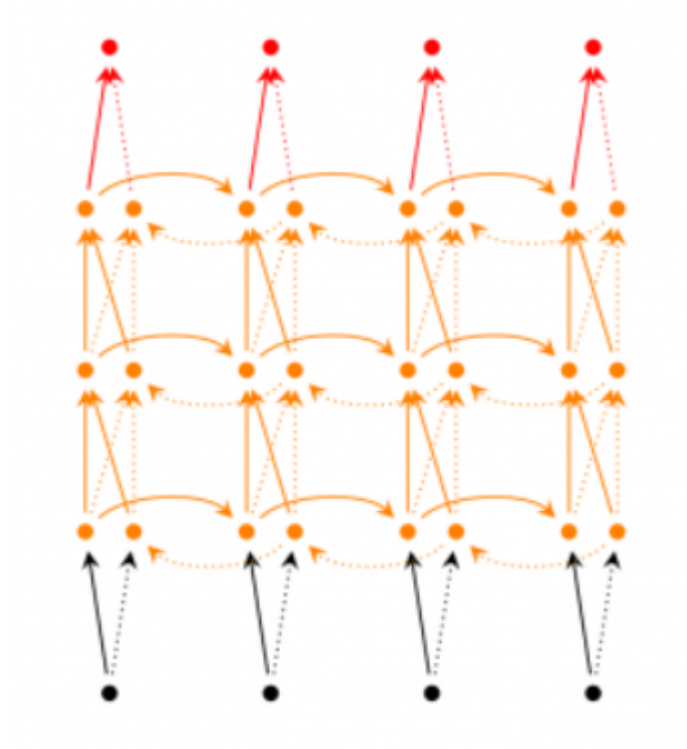
A cute teddy bear is reading Persian literature



A cute teddy bear is reading Persian literature

# RNN has opportunities to make deep RNN



Deep RNN