

# Cs231n Lecture 10

## RNN, LSTM

2170051 송여진



## 목차

1. RNN(Recurrent Neural Networks)
2. LSTM(Long Short Term Memory)

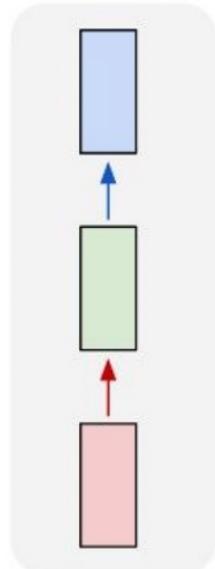


# 1. RNN(Recurrent Neural Networks)



# RNN vs Vanilla Neural Networks

one to one



Vanilla Neural Networks는 고정된 크기의 벡터를 입력으로 받고 고정된 크기의 벡터를 출력 즉 “고정된 양의 계산 단계” 를 사용하여 mapping을 수행

VS

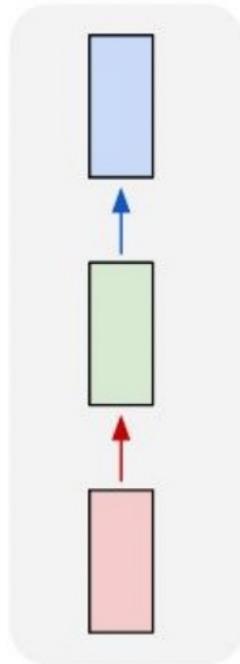
RNN은 vector sequence에 대해 연산을 수행 더 구체적으로 만들 수 있다

**Vanilla Neural Networks**

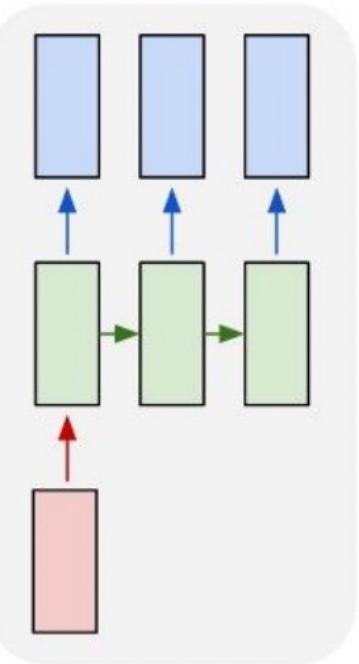


# RNN: Process Sequences

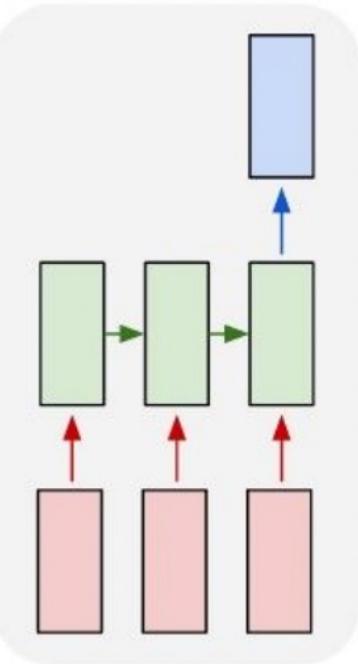
one to one



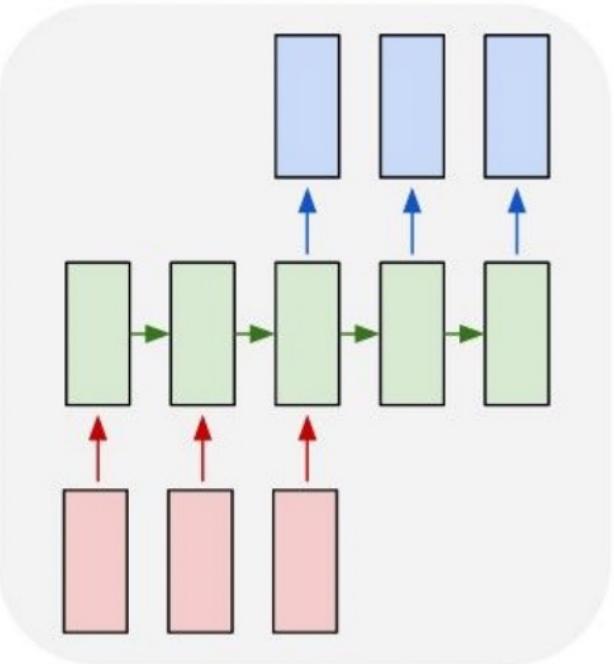
one to many



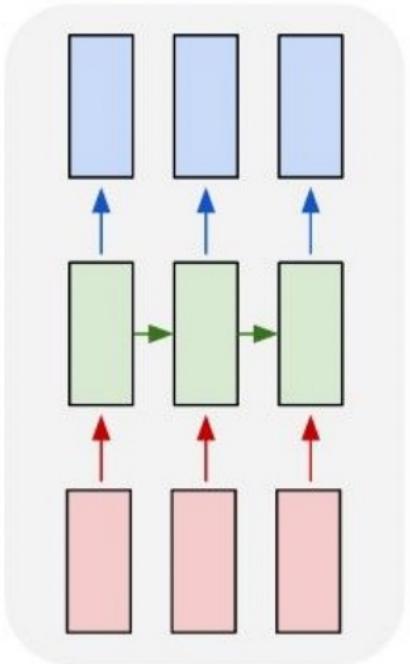
many to one



many to many



many to many



- █ Input vector
- █ Output vector
- █ RNN state

e.g. **Image classification**  
image -> sequence

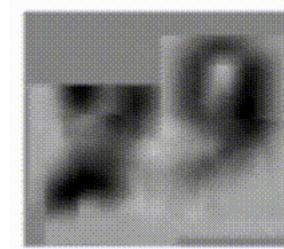
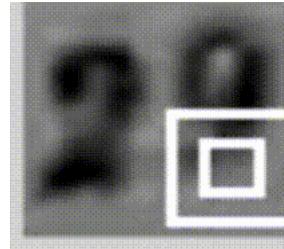
e.g. **Sequence prediction**  
sequence -> sequence

e.g. **Video classification on frame level**



# Sequential Processing of Non-Sequence Data

2	3	8	2	9	1	1	1	1	8
3	3	2	8	6	9	6	5	1	3
8	8	1	8	1	6	9	8	3	4
1	0	2	7	6	0	9	1	4	5
7	1	4	4	4	9	4	4	7	9
3	1	8	9	3	4	2	7	3	3
6	6	1	6	3	4	3	3	9	0
8	1	0	5	3	5	1	8	3	4
9	9	1	1	3	0	5	9	5	4
1	1	8	6	9	8	3	2	1	2

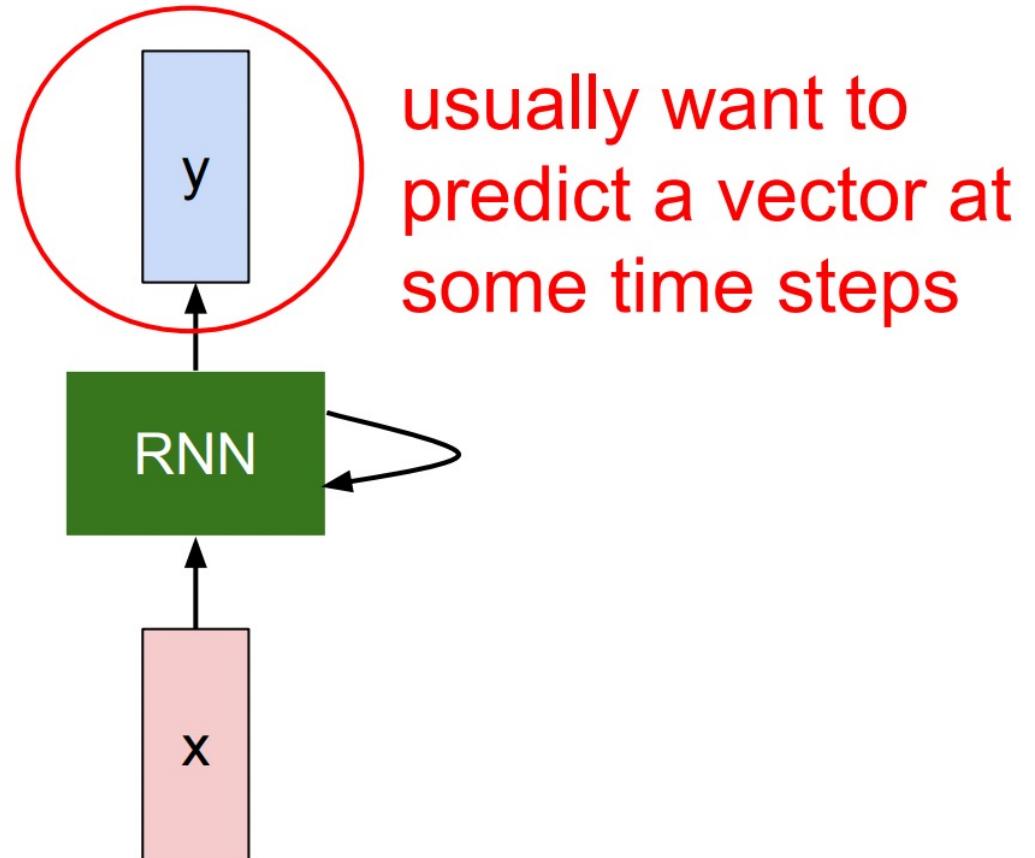


이미지의 숫자가 몇인지 분류하는 문제

-> 한번의 처리로 결정하지 않고  
여러 부분을 조금씩 살펴봄



# Recurrent Neural Network



우리가 원하는 것:  
특정 time steps의 벡터 예측



# Recurrent Neural Network

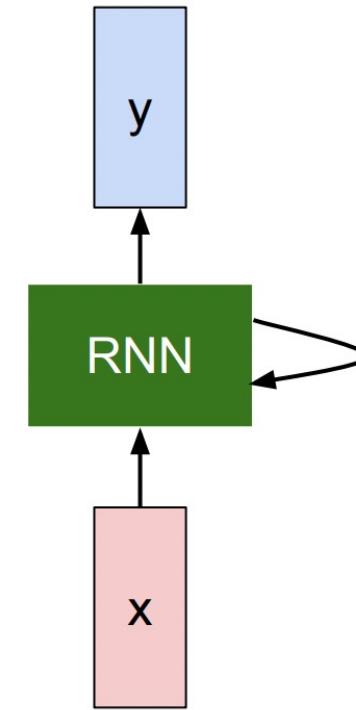
We can process a sequence of vectors  $\mathbf{x}$  by applying a **recurrence formula** at every time step:

Recurrence Function

$$h_t = f_W(h_{t-1}, x_t)$$

new state      old state      input vector at some time step

some function with parameters W

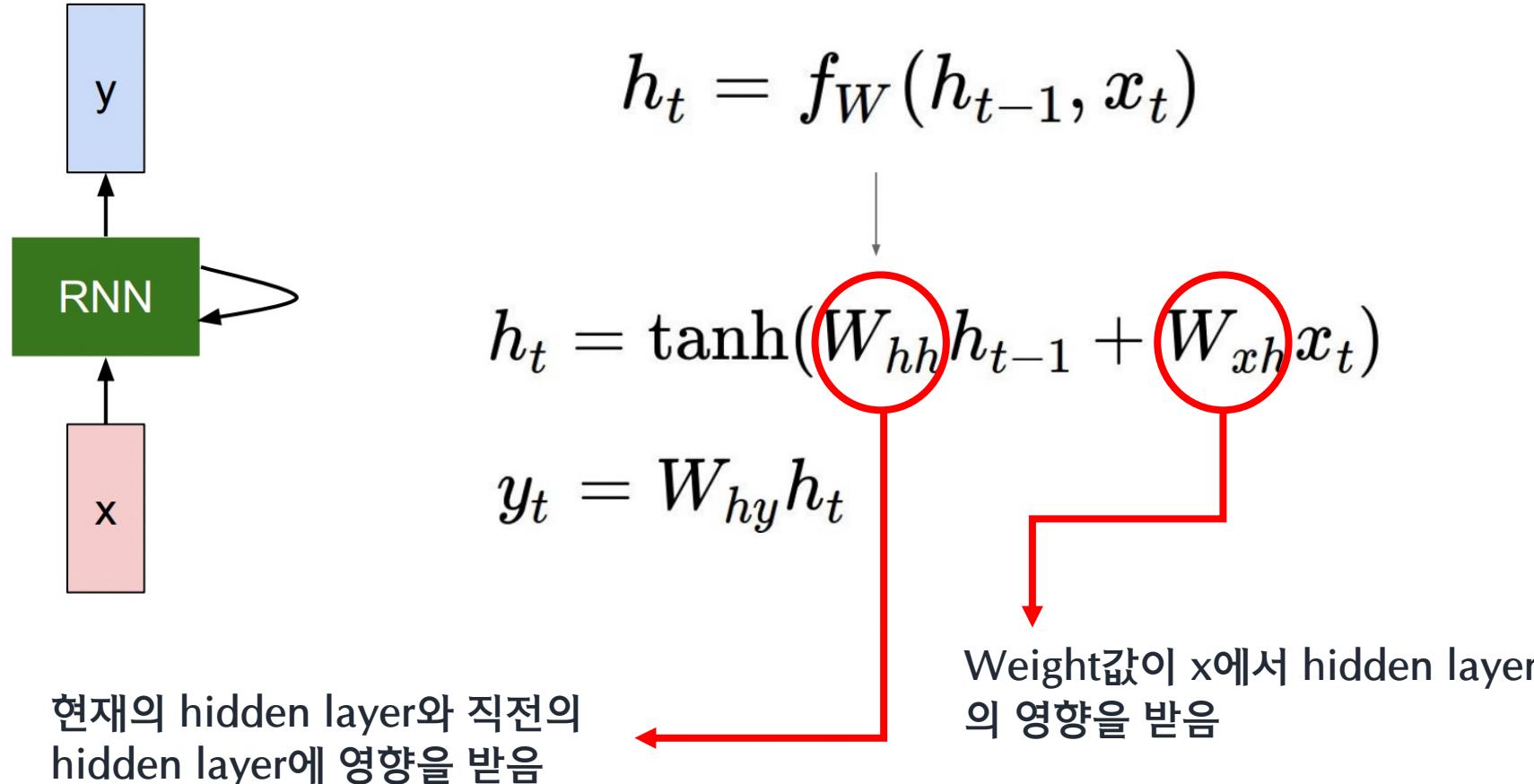


Notice: the same function and the same set of parameters are used at every time step.



# (Vanilla) Recurrent Neural Network

State가 하나의 "hidden" vector  $h$ 로 이루어짐.



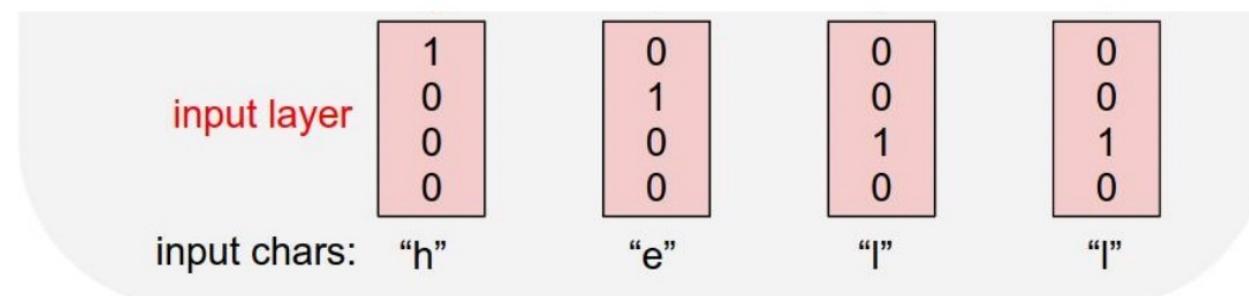
# Example: Character-level Language Model

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
“hello”

One-hot encoding

선택해야 하는 선택지의 개수만큼 차원을 가지면서, 각 선택지의 인덱스에 해당하는 원소에는 1, 나머지 원소에는 0의 값을 가지도록 하는 표현 방법



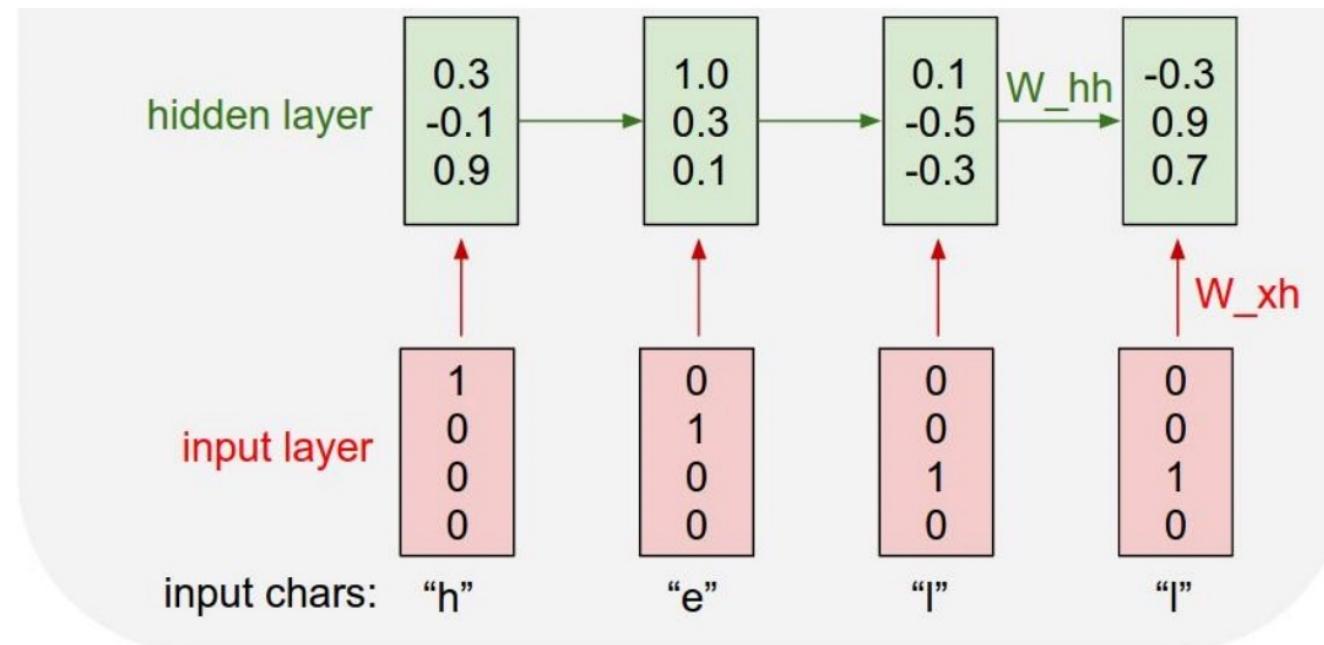
# Example: Character-level Language Model

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

Vocabulary:

[h,e,l,o]

Example training  
sequence:  
“hello”



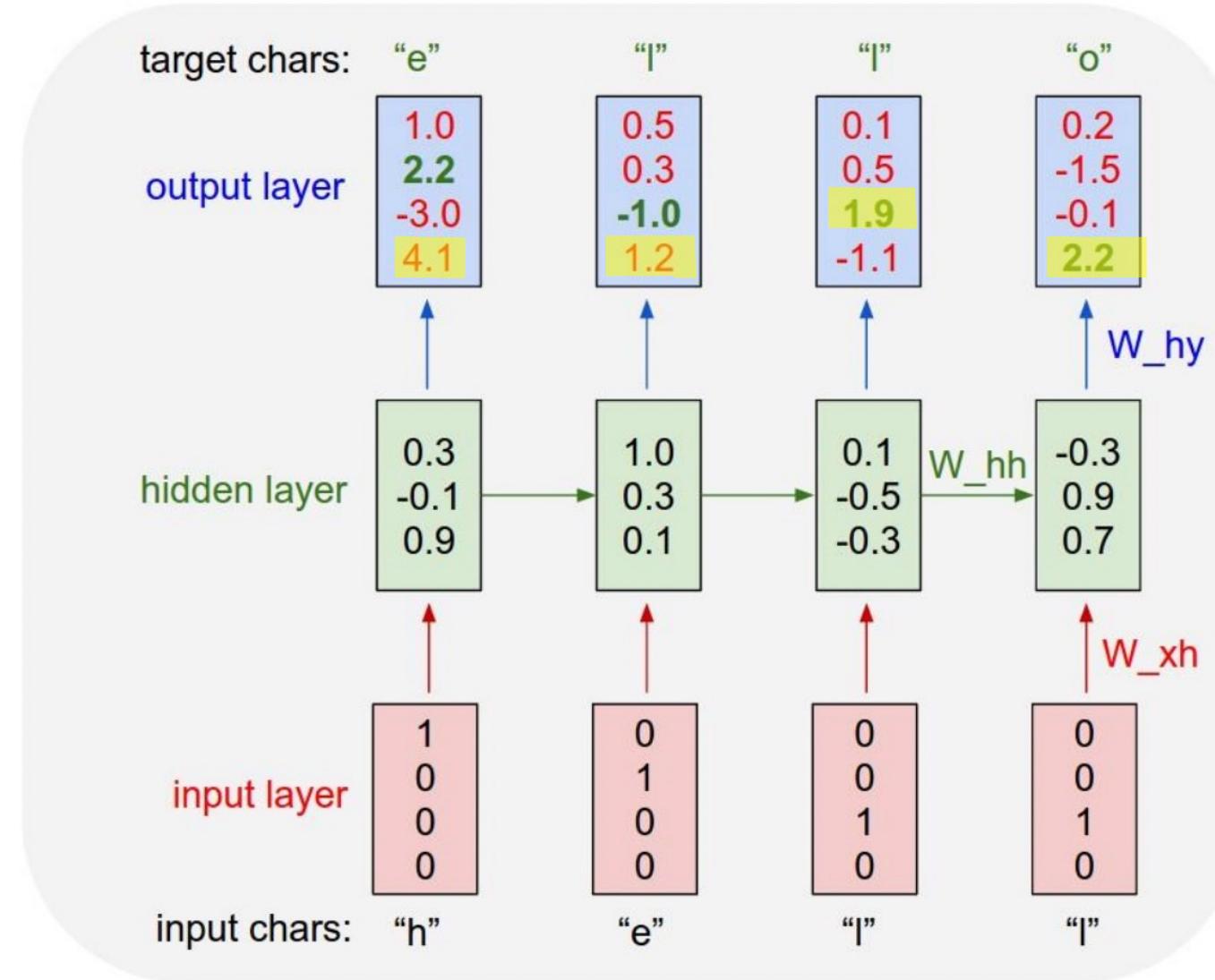
$W_{hh}$ : h가 h에 영향  
 $W_{xh}$ : x가 h에 영향  
 $W_{hy}$ : h가 y에 영향



# Example: Character-level Language Model

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
“hello”



1,2번째는 잘못된 예측  
을 하고 있음!!



정답과 비교해서  
Loss를 구해서  
Back propagation  
학습



# Example: Code

min-char-rnn.py gist: 112 lines of Python

```
"""
Minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)
BSD License
"""

import numpy as np

# data I/O
data = open('input.txt', 'r').read() # should be simple plain text file
chars = list(set(data))
data_size, vocab_size = len(data), len(chars)
print 'data has %d characters, %d unique.' % (data_size, vocab_size)
char_to_ix = { ch:i for i,ch in enumerate(chars) }
ix_to_char = { i:ch for i,ch in enumerate(chars) }

# hyperparameters
hidden_size = 100 # size of hidden layer of neurons
seq_length = 25 # number of steps to unroll the RNN for
learning_rate = 1e-1

# model parameters
Wxh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
bh = np.zeros((hidden_size, 1)) # hidden bias
by = np.zeros((vocab_size, 1)) # output bias

def lossFun(inputs, targets, hprev):
    """
    inputs,targets are both list of integers.
    hprev is Hx1 array of initial hidden state
    returns the loss, gradients on model parameters, and last hidden state
    """
    xs, hs, ys, ps = {}, {}, {}, {}
    hs[-1] = np.copy(hprev)
    loss = 0
    # forward pass
    for t in xrange(len(inputs)):
        xs[t] = np.zeros((vocab_size,1)) # encode in 1-of-k representation
        xs[t][inputs[t]] = 1
        hs[t] = np.tanh(np.dot(Wxh, xs[t]) + np.dot(Whh, hs[t-1]) + bh) # hidden state
        ys[t] = np.dot(Why, hs[t]) + by # unnormalized log probabilities for next chars
        ps[t] = np.exp(ys[t]) / np.sum(np.exp(ys[t])) # probabilities for next chars
        loss += -np.log(ps[t][targets[t],0]) # softmax (cross-entropy loss)
    # backward pass: compute gradients going backwards
    dWxh, dWhh, dWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
    dbh, dby = np.zeros_like(bh), np.zeros_like(by)
    dhnext = np.zeros_like(hs[0])
    for t in reversed(xrange(len(inputs))):
        dy = np.copy(ps[t])
        dy[targets[t]] -= 1 # backprop into y
        dWhy += np.dot(dy, hs[t].T)
        dby += dy
        dy = np.dot(Why.T, dy) + dhnext # backprop into h
        dWxh += np.dot(dy.T, xs[t])
        dWhh += np.dot(dy.T, hs[t-1])
        dbh += dy
    loss /= seq_length
    return loss, dWxh, dWhh, dWhy, dbh, dby, hs[-1], ps
```

```
def sample(h, seed_ix, n):
    """
    sample a sequence of integers from the model
    h is memory state, seed_ix is seed letter for first time step
    """
    x = np.zeros((vocab_size, 1))
    x[seed_ix] = 1
    ixes = []
    for t in xrange(n):
        h = np.tanh(np.dot(Wxh, x) + np.dot(Whh, h) + bh)
        y = np.dot(Why, h) + by
        p = np.exp(y) / np.sum(np.exp(y))
        ix = np.random.choice(range(vocab_size), p=p.ravel())
        x[ix] = 1
        ixes.append(ix)
    return ixes
```

```
n, p = 0, 0
mWxh, mWhh, mWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
mbh, mbv = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
while True:
    # prepare inputs (we're sweeping from left to right in steps seq_length long)
    if pseq_length+1 >= len(data) or n == 0:
        hprev = np.zeros((hidden_size,1)) # reset RNN memory
        p = 0 # go from start of data
        inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
        targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
    else:
        # sample from the model now and then
        if n % 100 == 0:
            sample_ix = sample(hprev, inputs[0], 200)
            txt = ''.join(ix_to_char[ix] for ix in sample_ix)
            print '----\n%s ----' % (txt,)

        # forward sequence characters through the net and fetch gradient
        loss, dWxh, dWhh, dWhy, hprev = lossFun(inputs, targets, hprev)
        smooth_loss = smooth_loss * 0.999 + loss * 0.001
        if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress

        # perform parameter update with Adagrad
        for param, dparam, mem in zip([Wxh, Whh, Why, bh, by],
                                      [dWxh, dWhh, dWhy, dbh, dby],
                                      [mWxh, mWhh, mWhy, mbh, mbv]):
            mem += dparam * dparam
            param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
            p += seq_length # move data pointer
            n += 1 # iteration counter
```

<https://gist.github.com/karpathy/d4dee566867f8291f0>



# Example: RNN Language model

tyntd-iafhatawiaoahrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e  
plia tklrgd t o idoe ns,smtt h ne etie h,hregtrs nigtike,aoaenns lng

↓ train more

"Tmont thithey" fomesscerliund  
Keushey. Thom here  
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwyl fil on aseterlome  
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

↓ train more

Aftair fall unsuch that the hall for Prince Velzonski's that me of  
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort  
how, and Gogition is so overelical and ofter.

↓ train more

"Why do what that day," replied Natasha, and wishing to himself the fact the  
princess, Princess Mary was easier, fed in had oftened him.  
Pierre aking his soul came to the packs and drove up his father-in-law women.



# Example: RNN Language model

The Stacks Project

home about tags explained tag lookup browse search bibliography recent comments blog add slogans

Browse chapters

Part	Chapter	online	TeX source	view pdf
Preliminaries	1. Introduction	<a href="#">online</a>	<a href="#">tex</a>	<a href="#">pdf &gt;</a>
	2. Conventions	<a href="#">online</a>	<a href="#">tex</a>	<a href="#">pdf &gt;</a>
	3. Set Theory	<a href="#">online</a>	<a href="#">tex</a>	<a href="#">pdf &gt;</a>
	4. Categories	<a href="#">online</a>	<a href="#">tex</a>	<a href="#">pdf &gt;</a>
	5. Topology	<a href="#">online</a>	<a href="#">tex</a>	<a href="#">pdf &gt;</a>
	6. Sheaves on Spaces	<a href="#">online</a>	<a href="#">tex</a>	<a href="#">pdf &gt;</a>
	7. Sites and Sheaves	<a href="#">online</a>	<a href="#">tex</a>	<a href="#">pdf &gt;</a>
	8. Stacks	<a href="#">online</a>	<a href="#">tex</a>	<a href="#">pdf &gt;</a>
	9. Fields	<a href="#">online</a>	<a href="#">tex</a>	<a href="#">pdf &gt;</a>
	10. Commutative Algebra	<a href="#">online</a>	<a href="#">tex</a>	<a href="#">pdf &gt;</a>

Parts

1. [Preliminaries](#)
2. [Schemes](#)
3. [Topics in Scheme Theory](#)
4. [Algebraic Spaces](#)
5. [Topics in Geometry](#)
6. [Deformation Theory](#)
7. [Algebraic Stacks](#)
8. [Miscellany](#)

Statistics

The Stacks project now consists of

- o 455910 lines of code
- o 14221 tags (56 inactive tags)
- o 2366 sections

Stacks Project에 있는 대수기하학 교과서를 RNN에 돌려봄



# Example: RNN Language model

For  $\bigoplus_{n=1,\dots,m} \mathcal{L}_{m,n} = 0$ , hence we can find a closed subset  $\mathcal{H}$  in  $\mathcal{H}$  and any sets  $\mathcal{F}$  on  $X$ ,  $U$  is a closed immersion of  $S$ , then  $U \rightarrow T$  is a separated algebraic space.

*Proof.* Proof of (1). It also start we get

$$S = \text{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by  $\coprod Z \times_U U \rightarrow V$ . Consider the maps  $M$  along the set of points  $\text{Sch}_{fppf}$  and  $U \rightarrow U$  is the fibre category of  $S$  in  $U$  in Section, ?? and the fact that any  $U$  affine, see Morphisms, Lemma ???. Hence we obtain a scheme  $S$  and any open subset  $W \subset U$  in  $\text{Sh}(G)$  such that  $\text{Spec}(R') \rightarrow S$  is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that  $f_i$  is of finite presentation over  $S$ . We claim that  $\mathcal{O}_{X,x}$  is a scheme where  $x, x', s'' \in S'$  such that  $\mathcal{O}_{X,x'} \rightarrow \mathcal{O}'_{X',x'}$  is separated. By Algebra, Lemma ?? we can define a map of complexes  $\text{GL}_{S'}(x'/S'')$  and we win.  $\square$

To prove study we see that  $\mathcal{F}|_U$  is a covering of  $\mathcal{X}'$ , and  $\mathcal{T}_i$  is an object of  $\mathcal{F}_{X/S}$  for  $i > 0$  and  $\mathcal{F}_p$  exists and let  $\mathcal{F}_i$  be a presheaf of  $\mathcal{O}_X$ -modules on  $\mathcal{C}$  as a  $\mathcal{F}$ -module. In particular  $\mathcal{F} = U/\mathcal{F}$  we have to show that

$$\widetilde{M}^\bullet = \mathcal{I}^\bullet \otimes_{\text{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1} \mathcal{F}$$

is a unique morphism of algebraic stacks. Note that

$$\text{Arrows} = (\text{Sch}/S)_{fppf}^{\text{opp}}, (\text{Sch}/S)_{fppf}$$

and

$$V = \Gamma(S, \mathcal{O}) \longrightarrow (U, \text{Spec}(A))$$

is an open subset of  $X$ . Thus  $U$  is affine. This is a continuous map of  $X$  is the inverse, the groupoid scheme  $S$ .

*Proof.* See discussion of sheaves of sets.  $\square$

The result for prove any open covering follows from the less of Example ???. It may replace  $S$  by  $X_{\text{spaces},\text{étale}}$  which gives an open subspace of  $X$  and  $T$  equal to  $S_{\text{Zar}}$ , see Descent, Lemma ???. Namely, by Lemma ?? we see that  $R$  is geometrically regular over  $S$ .

**Lemma 0.1.** Assume (3) and (3) by the construction in the description.

Suppose  $X = \lim |X|$  (by the formal open covering  $X$  and a single map  $\underline{\text{Proj}}_X(\mathcal{A}) = \text{Spec}(B)$  over  $U$  compatible with the complex

$$\text{Set}(\mathcal{A}) = \Gamma(X, \mathcal{O}_{X,\mathcal{O}_X}).$$

When in this case of to show that  $\mathcal{Q} \rightarrow \mathcal{C}_{Z/X}$  is stable under the following result in the second conditions of (1), and (3). This finishes the proof. By Definition ?? (without element is when the closed subschemes are catenary. If  $T$  is surjective we may assume that  $T$  is connected with residue fields of  $S$ . Moreover there exists a closed subspace  $Z \subset X$  of  $X$  where  $U$  in  $X'$  is proper (some defining as a closed subset of the uniqueness it suffices to check the fact that the following theorem

(1)  $f$  is locally of finite type. Since  $S = \text{Spec}(R)$  and  $Y = \text{Spec}(R)$ .

*Proof.* This is form all sheaves of sheaves on  $X$ . But given a scheme  $U$  and a surjective étale morphism  $U \rightarrow X$ . Let  $U \cap U = \coprod_{i=1,\dots,n} U_i$  be the scheme  $X$  over  $S$  at the schemes  $X_i \rightarrow X$  and  $U = \lim_i X_i$ .  $\square$

The following lemma surjective restrocomposes of this implies that  $\mathcal{F}_{x_0} = \mathcal{F}_{x_0} = \mathcal{F}_{x,\dots,0}$ .

**Lemma 0.2.** Let  $X$  be a locally Noetherian scheme over  $S$ ,  $E = \mathcal{F}_{X/S}$ . Set  $\mathcal{I} = \mathcal{J}_1 \subset \mathcal{I}_n$ . Since  $\mathcal{I}^n \subset \mathcal{I}^n$  are nonzero over  $i_0 \leq p$  is a subset of  $\mathcal{J}_{n,0} \circ \bar{A}_2$  works.

**Lemma 0.3.** In Situation ???. Hence we may assume  $q' = 0$ .

*Proof.* We will use the property we see that  $p$  is the next functor (??). On the other hand, by Lemma ?? we see that

$$D(\mathcal{O}_{X'}) = \mathcal{O}_X(D)$$

where  $K$  is an  $F$ -algebra where  $\delta_{n+1}$  is a scheme over  $S$ .  $\square$



# Example: Code

Proof. Omitted. □

**Lemma 0.1.** Let  $\mathcal{C}$  be a set of the construction.

Let  $\mathcal{C}$  be a gerber covering. Let  $\mathcal{F}$  be a quasi-coherent sheaves of  $\mathcal{O}$ -modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

*Proof.* This is an algebraic space with the composition of sheaves  $\mathcal{F}$  on  $X_{\text{étale}}$  we have

$$\mathcal{O}_X(\mathcal{F}) = \{\text{morph}_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where  $\mathcal{G}$  defines an isomorphism  $\mathcal{F} \rightarrow \mathcal{F}$  of  $\mathcal{O}$ -modules. □

**Lemma 0.2.** This is an integer  $\mathcal{Z}$  is injective.

*Proof.* See Spaces, Lemma ??.

**Lemma 0.3.** Let  $S$  be a scheme. Let  $X$  be a scheme and  $X$  is an affine open covering. Let  $\mathcal{U} \subset \mathcal{X}$  be a canonical and locally of finite type. Let  $X$  be a scheme. Let  $X$  be a scheme which is equal to the formal complex.

The following to the construction of the lemma follows.

Let  $X$  be a scheme. Let  $X$  be a scheme covering. Let

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

be a morphism of algebraic spaces over  $S$  and  $Y$ .

*Proof.* Let  $X$  be a nonzero scheme of  $X$ . Let  $X$  be an algebraic space. Let  $\mathcal{F}$  be a quasi-coherent sheaf of  $\mathcal{O}_X$ -modules. The following are equivalent

- (1)  $\mathcal{F}$  is an algebraic space over  $S$ .
- (2) If  $X$  is an affine open covering.

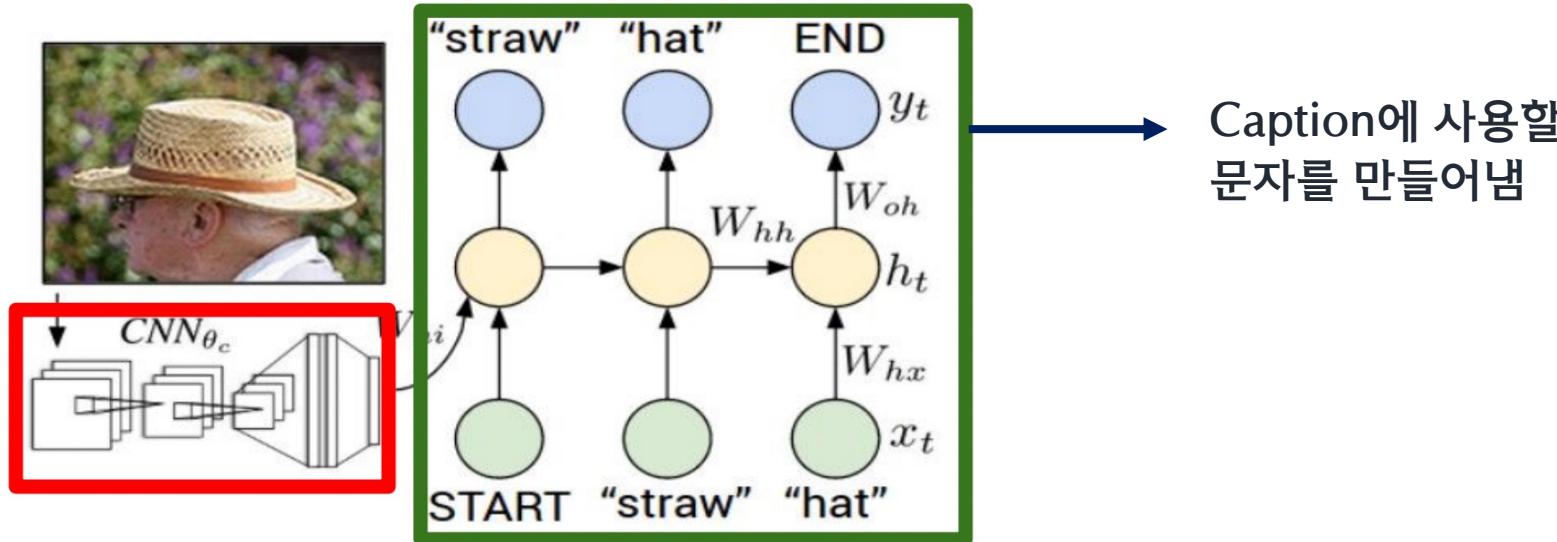
Consider a common structure on  $X$  and  $X$  the functor  $\mathcal{O}_X(U)$  which is locally of finite type. □

**Proof. Omitted**  
= 증명 생략



# Image Captioning

## Recurrent Neural Network



## Convolutional Neural Network

입력: 이미지

출력: caption(자연어, 가변 길이)



# Image Captioning

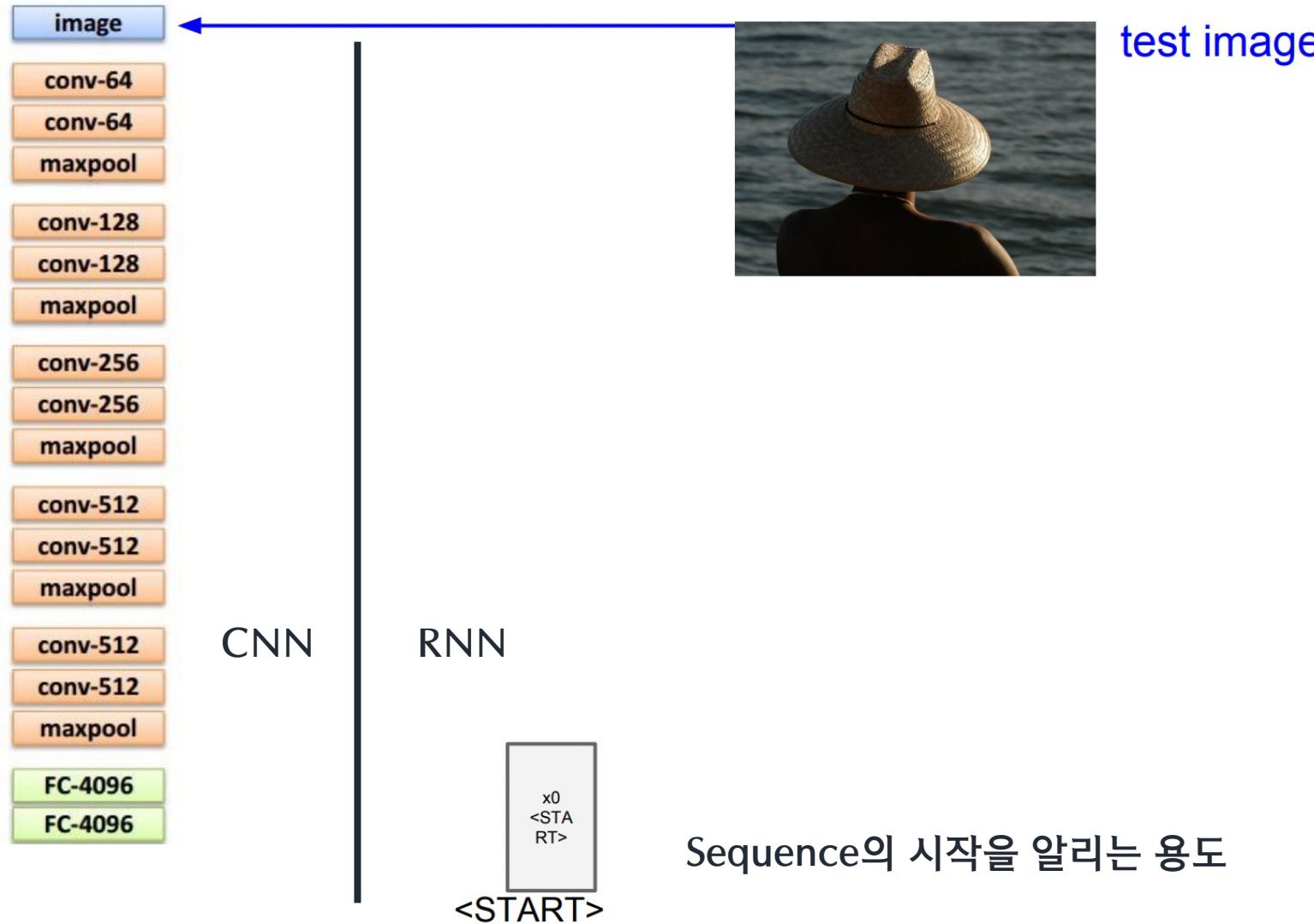


test image

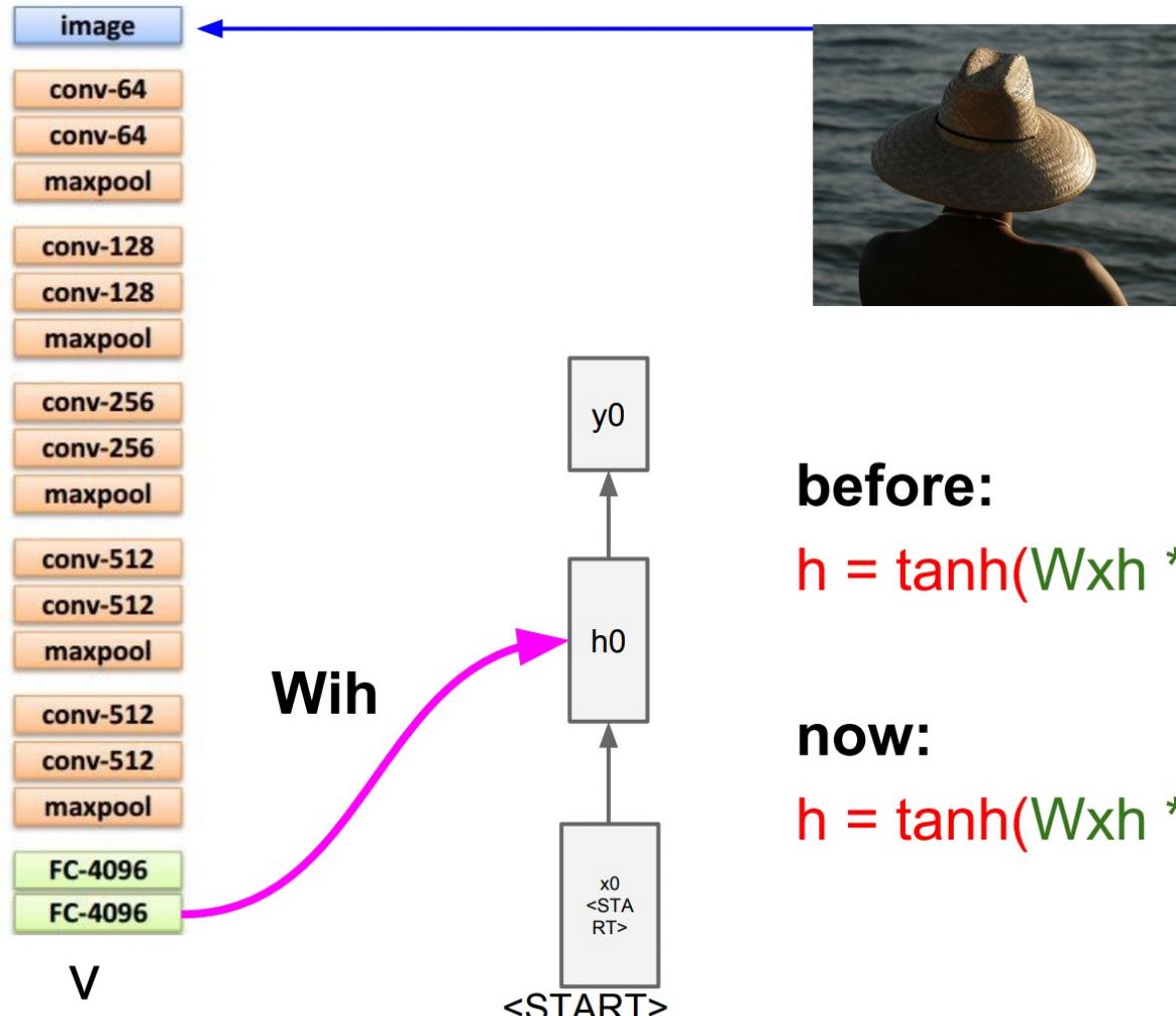
상위 2개의 Layer(FC-1000, softmax)를 삭제



# Image Captioning



# Image Captioning



test image

before:

$$h = \tanh(W_{xh} * x + W_{hh} * h)$$

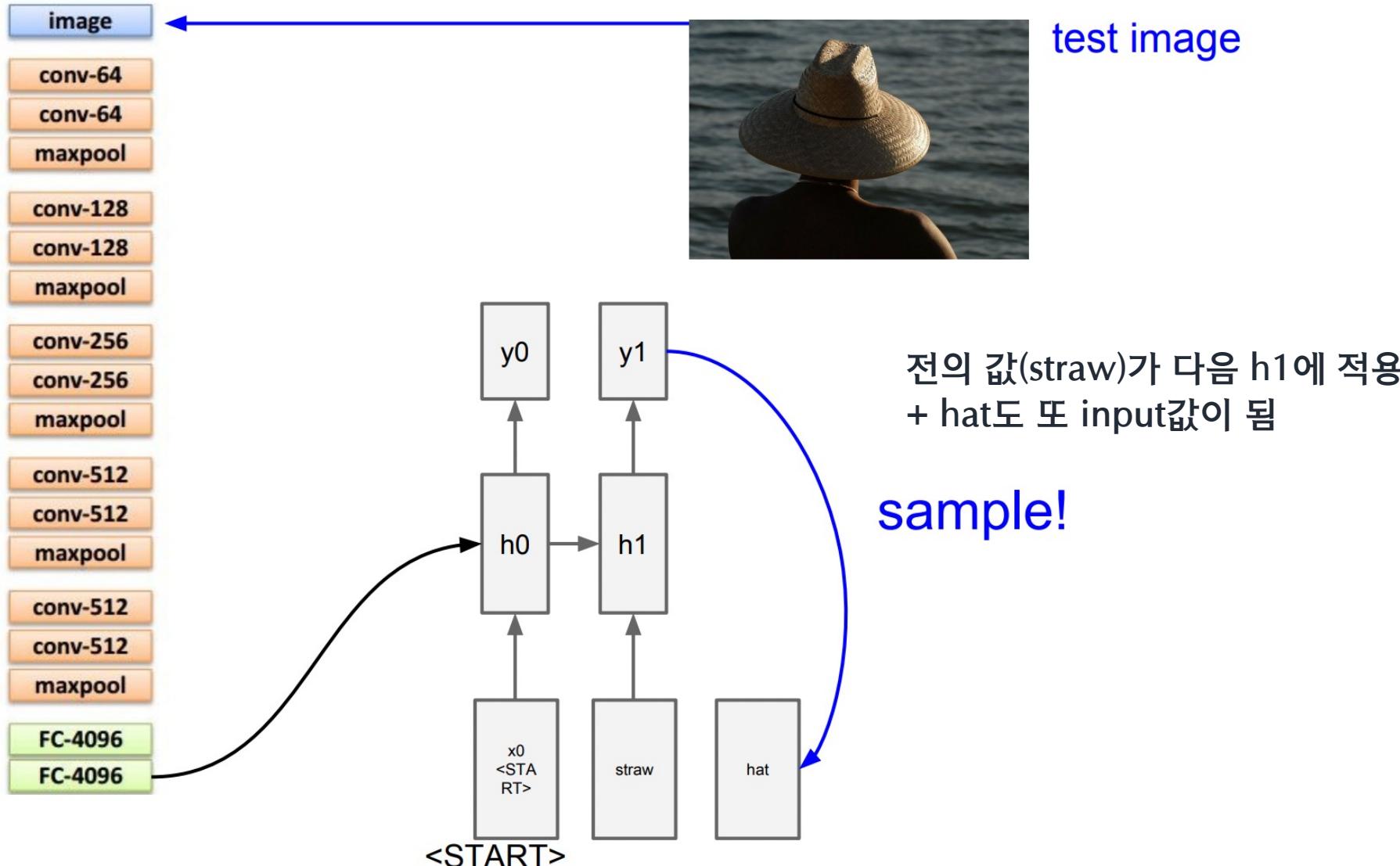
now:

$$h = \tanh(W_{xh} * x + W_{hh} * h + W_{vh} * v)$$

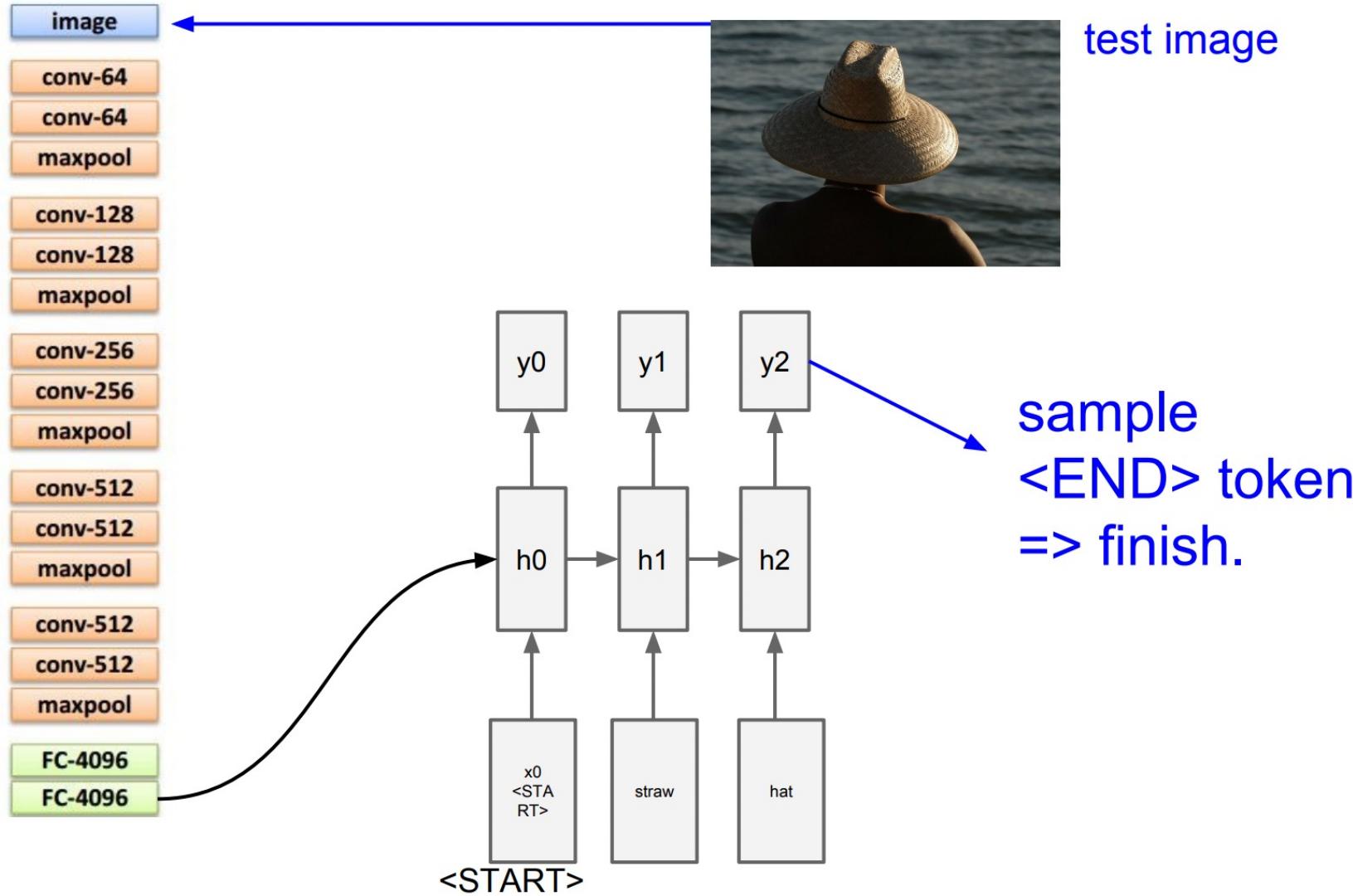
Image에서 온 가중치



# Image Captioning



# Image Captioning



# Image Captioning



*A cat sitting on a suitcase on the floor*



*A cat is sitting on a tree branch*



*A dog is running in the grass with a frisbee*



*A white teddy bear sitting in the grass*



*Two people walking on the beach with surfboards*



*A tennis player in action on the court*



*Two giraffes standing in a grassy field*



*A man riding a dirt bike on a dirt track*



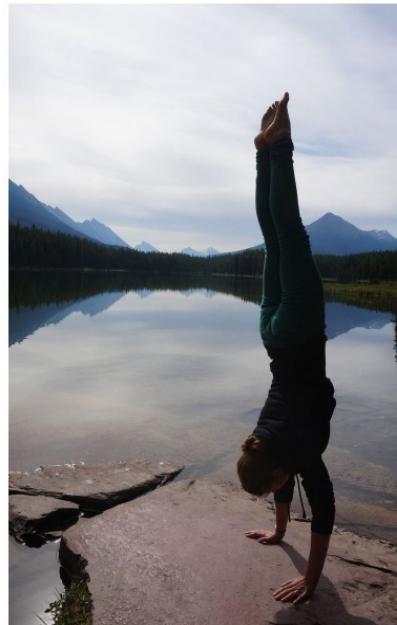
# Image Captioning: Failure Cases



*A woman is holding a cat in her hand*



*A person holding a computer mouse on a desk*



*A woman standing on a beach holding a surfboard*



*A bird is perched on a tree branch*



*A man in a baseball uniform throwing a ball*



# Multilayer RNNs

Layer가 여러 개인 RNN (hidden state가 여러 개)

RNN:

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$h \in \mathbb{R}^n \quad W^l [n \times 2n]$$

LSTM:

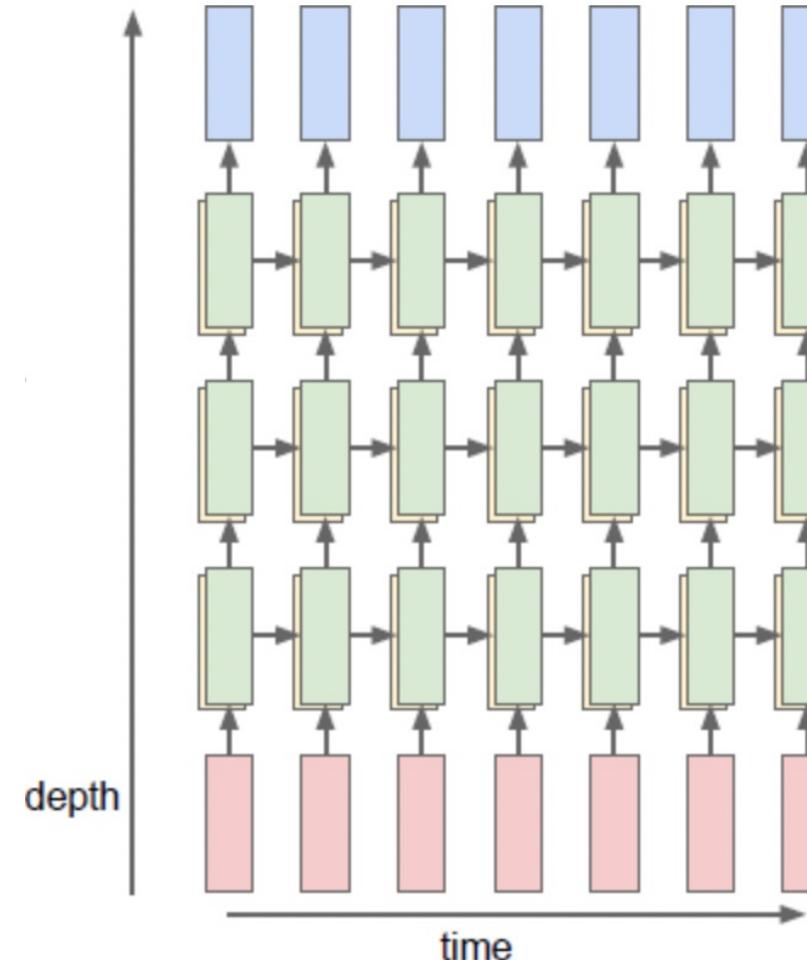
$$W^l [4n \times 2n]$$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

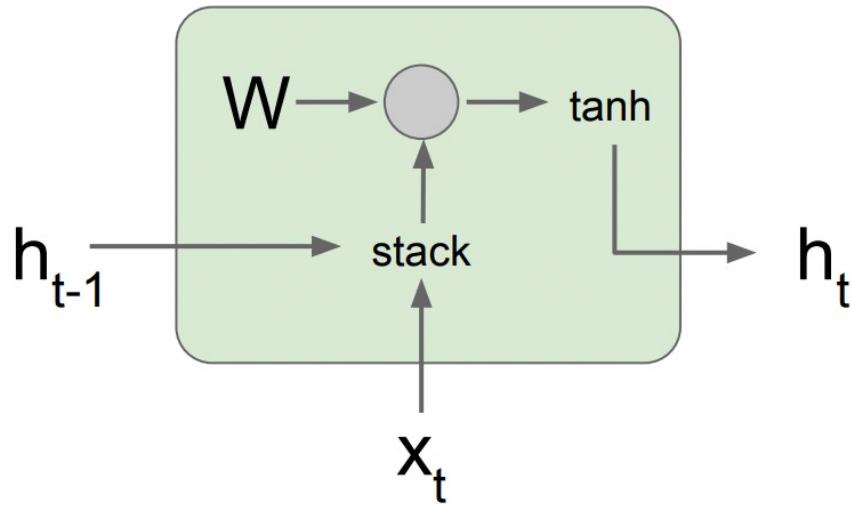
$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

3-Layer RNN



# Vanilla RNN Gradient Flow



$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\ &= \tanh \left( \begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \\ &= \tanh \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \end{aligned}$$

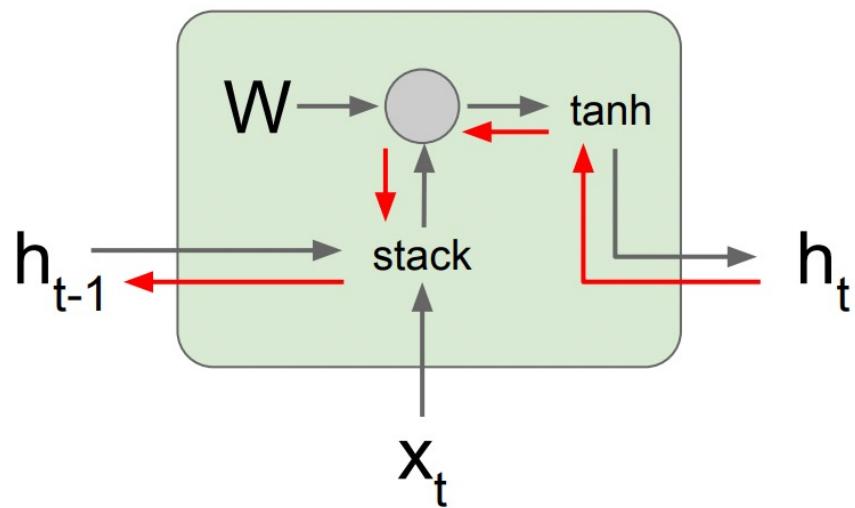
현재 값  $X(t)$ 와  $h(t-1)$ 을 쌓는다(stack)

가중치 행렬  $W$ 와 행렬곱 연산을 한다



# Vanilla RNN Gradient Flow

Backpropagation from  $h_t$  to  $h_{t-1}$  multiplies by  $W$   
(actually  $W_{hh}^T$ )



Background pass로 Gradient를 계산한다면?

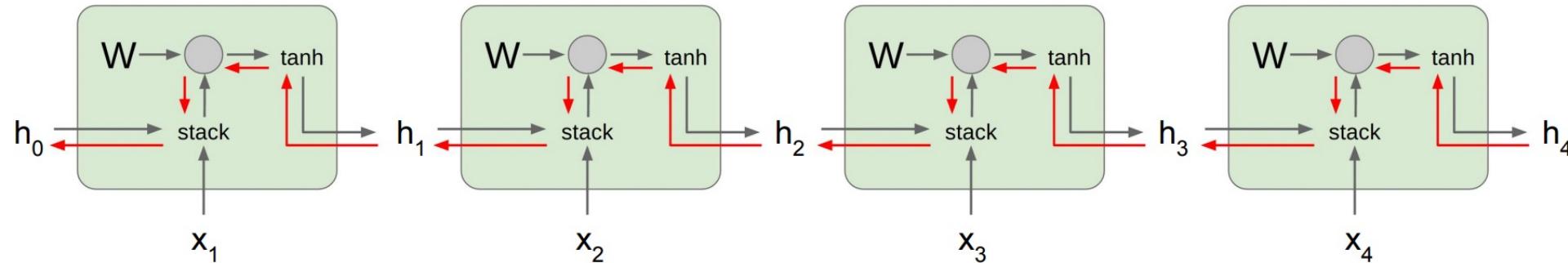


$W(hh)$ 의 Transpose matrix(전치 행렬)을 곱하게 됨

즉 1개의 cell을 통과할 때마다 가중치 행렬의 일부를  
곱하게 됨



# Vanilla RNN Gradient Flow



$h(0)$ 에 대한 gradient를 구하고자 한다면 결국 모든 RNN cells를 거쳐야 함

-> 엄청나게 많은 연산, 좋지 않음



# Vanilla RNN Gradient Flow

가중치가 벡터가 아닌 스칼라라고 생각해 본다면...

- 곱해지는 값이 1보다 크다  
-> 점점 값이 커짐(Exploding gradients) → 곱해지는 값이 “1”이어야 함!
- 곱해지는 값이 1보다 작다  
-> 점점 값이 작아져서 0이 됨(Vanishing gradients)



# Vanilla RNN Gradient Flow

Largest singular value > 1:  
**Exploding gradients**



**Gradient clipping:** Scale  
gradient if its norm is too big

```
grad_norm = np.sum(grad * grad)
if grad_norm > threshold:
    grad *= (threshold / grad_norm)
```

Largest singular value < 1:  
**Vanishing gradients**



Change RNN architecture

:LSTM



## 2. LSTM(Long Short Term Memory)



# Long Short Term Memory

## Vanilla RNN

$$h_t = \tanh \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

Hidden state 존재  
매 step마다 재귀적 방법으로  
hidden state 업데이트

## LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

한 cell마다 2개의 hidden state 존재



# Long Short Term Memory

## LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

C(t) - Cell State : LSTM 내부에만 존재하며 노출되지 않는 변수

H(t) - Hidden State



# Long Short Term Memory

f: Forget gate, Whether to erase cell

i: Input gate, whether to write to cell

g: Gate gate (?), How much to write to cell

o: Output gate, How much to reveal cell

F: 0이면 이전 cell state를 잊고 1이면 이전 cell state를 기억한다

I: 각 element에 cell state를 사용하고 싶다면 1, 아니라면 0

G: input cell을 얼마나 포함시킬지 결정

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

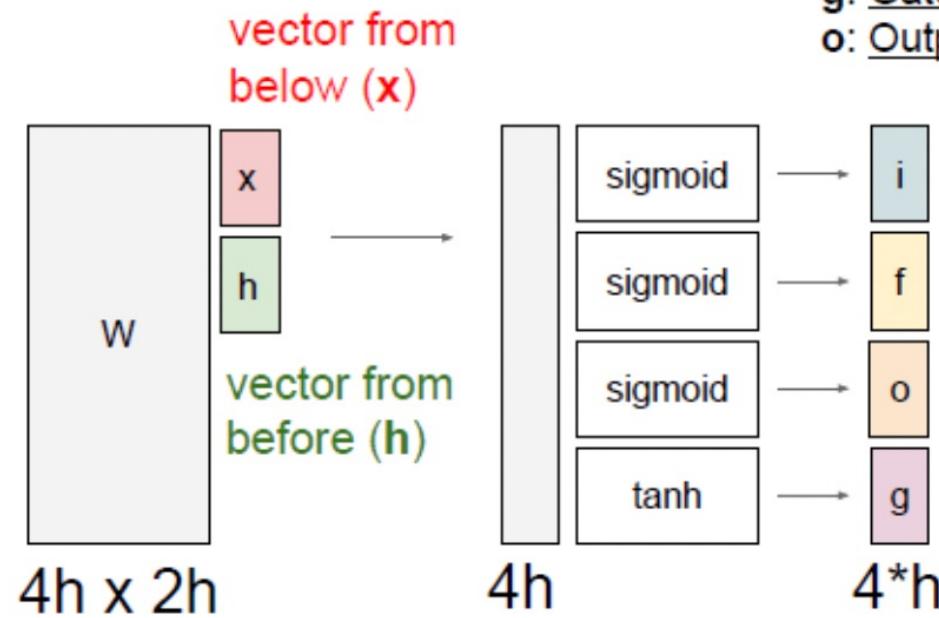
C(t)의 수식은 이전 cell state를  
계속 기억할지 말지 결정한다.



# Long Short Term Memory

## Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



- f: Forget gate, Whether to erase cell
- i: Input gate, whether to write to cell
- g: Gate gate (?), How much to write to cell
- o: Output gate, How much to reveal cell

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

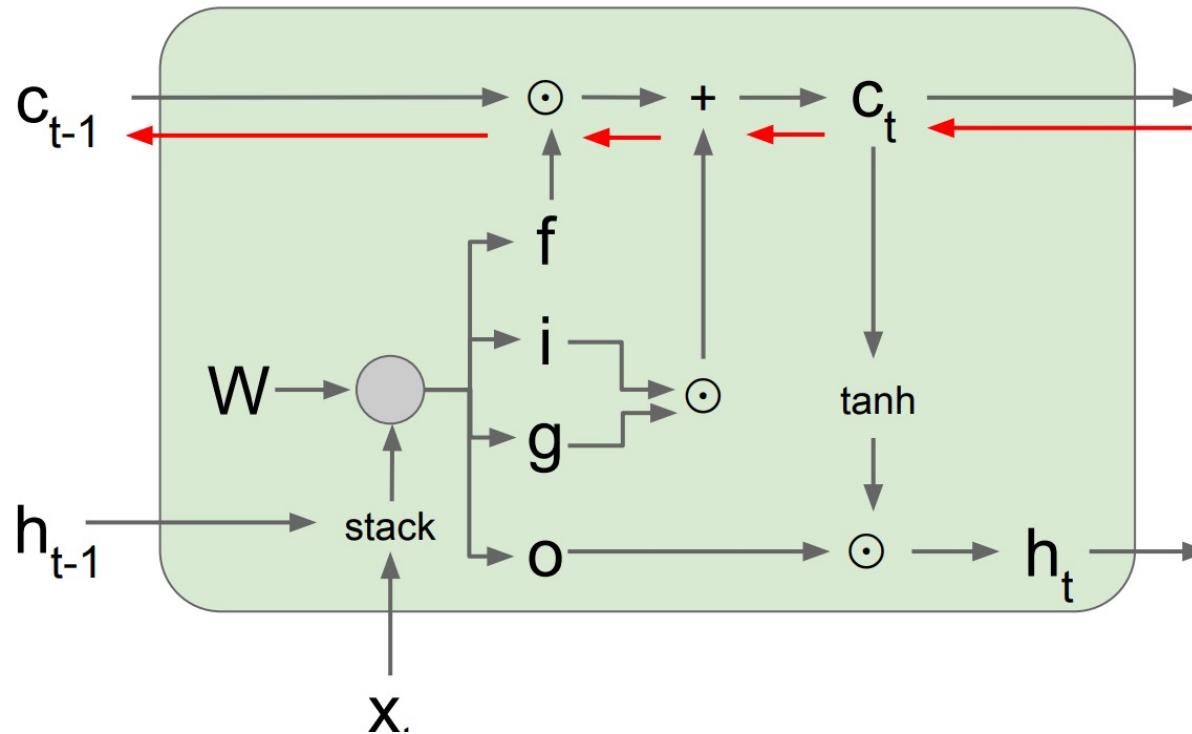
$$h_t = o \odot \tanh(c_t)$$

i,f,o는 sigmoid 함수를  
지나기 때문에  
[0,1]의 값을 가짐

G는 tanh를 지나기  
때문에  
[-1,1]의 값을 가짐



# LSTM Gradient Flow



Backpropagation from  $c_t$  to  $c_{t-1}$  only elementwise multiplication by  $f$ , no matrix multiply by  $W$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Gradient : Upstream Gradient \* Forget gate -> 문제 해결!



# LSTM Gradient Flow

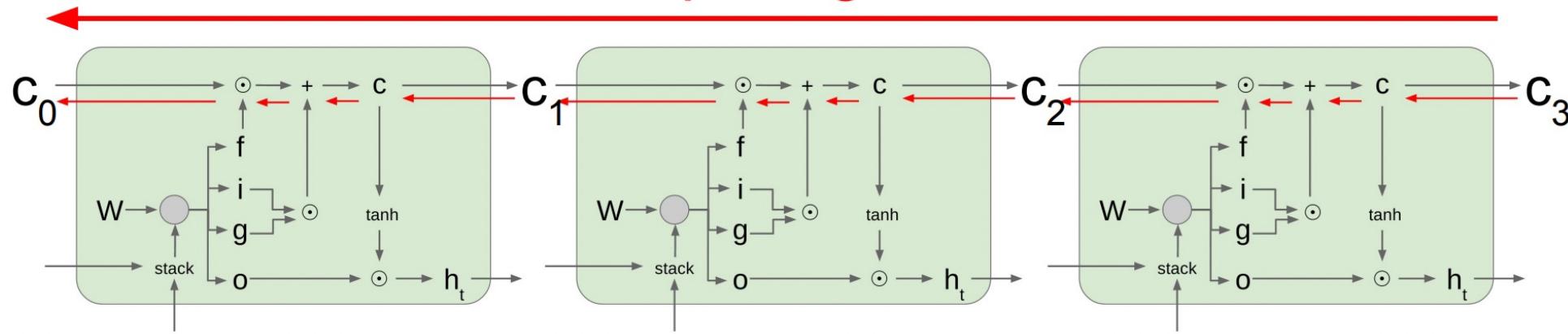
왜 더 좋을까? 🤔

1. Matrix multiplication이 아닌 element-wise(각 성분별로 계산)
2. element-wise를 통해 매 스텝 다른  $f$ 값과 곱해질 수 있다. 또한,  $f$ 는 sigmoid 함수에 의해 0과 1 사이의 값이므로 곱할수록 더 좋은 수치적 특성으로 변한다.
3. RNN처럼 매 단계 tanh 함수를 거칠 필요 없이 단 1번만 tanh를 거치면 됨!



# LSTM Gradient Flow

Uninterrupted gradient flow!



EWHA,  
THE FUTURE  
WE CREATE

감사합니다.

