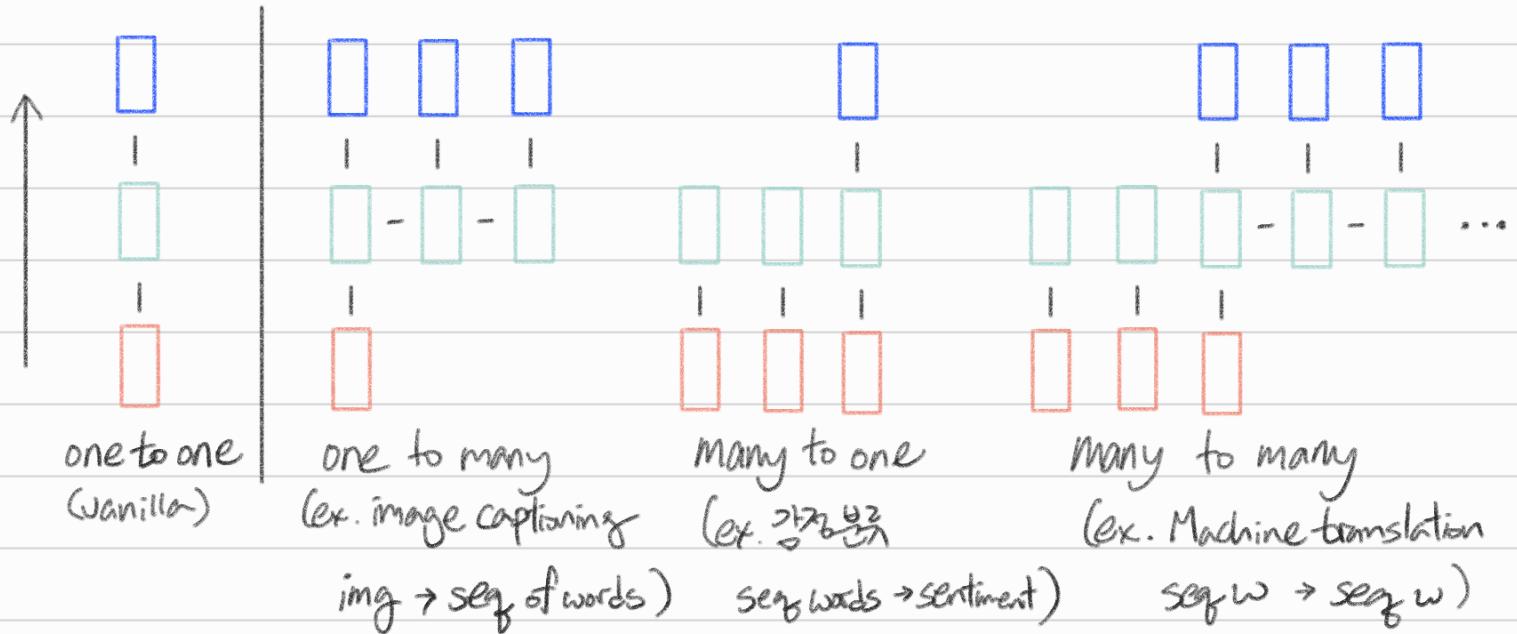
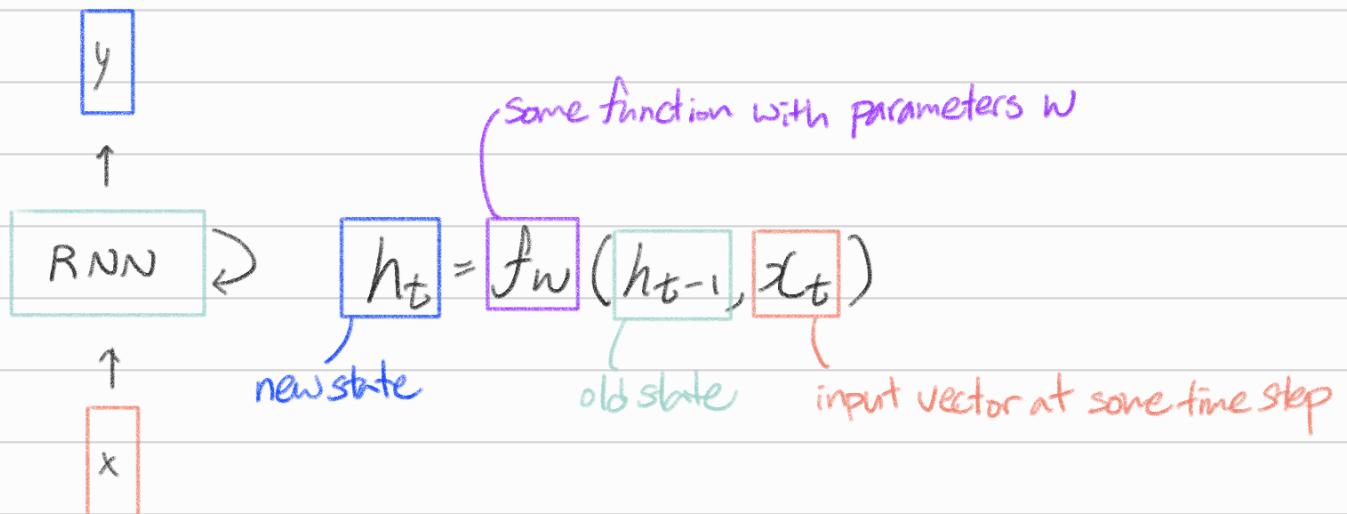


Recurrent Neural Networks

We want to have more flexibility in the types of data that our models can process ↴



RNN: recurrence formula @ every time step



* Same function and the same set of parameters are used at every time step

Vanilla RNN :

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

현재 입력값 (x_t)에 W_{xh} 곱한 것

+

과거 값 (x_{t-1})에 W_{hh} 곱한 것 + 편향값 (b)

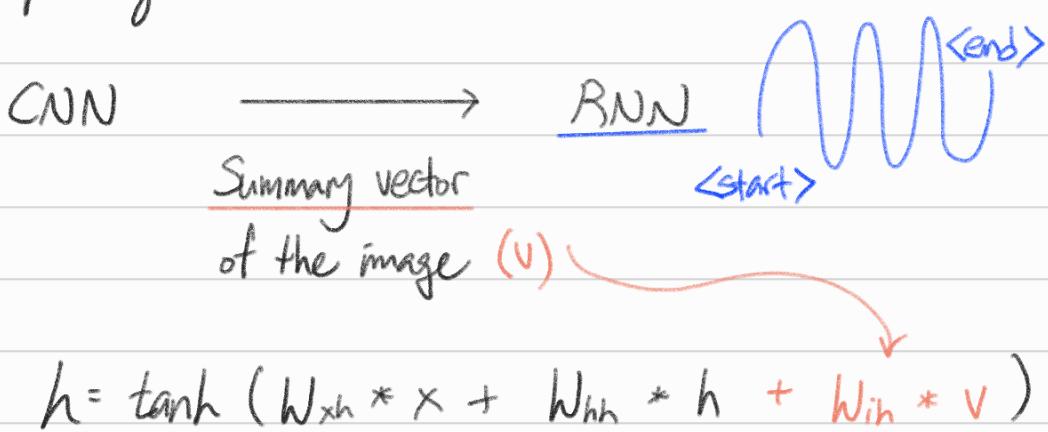
BPTT (Backpropagation through time)

: forward through entire sequence to compute loss,
then backward through entire seq to compute gradient

Truncated " : 오래 걸리니 짧게 (근사치 사용)

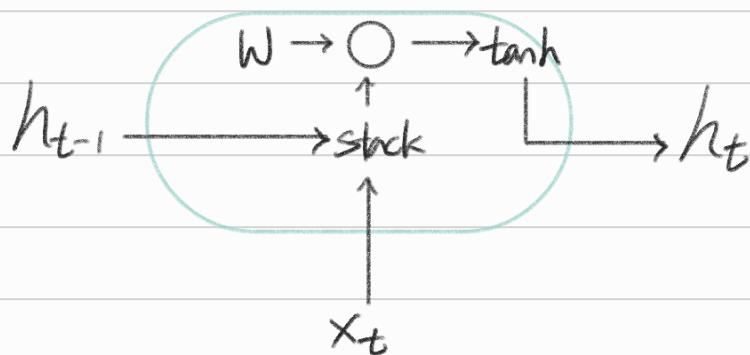
"현재와 과거를 함께 고려해 미래를 예측"

Image Captioning



Visual Question Answering : RNNs with Attention

Vanilla RNN Gradient Flow



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

$$= \tanh (W (h_{t-1} x_t))$$

$h_0 \ h_1 \ h_2 \ h_3 \ h_4 \ h_5$

←

Backprop

SE layers only
정지값

Computing gradient of h_0 involves many factors of W , Vanishing Gr
같은 값을 반복해서 곱하면 1이 아닌 경우 gradient가 0으로 수렴하거나
explode

explode 하면 'Gradient clipping' (if its norm is too big)

Vanish 하면 change RNN architecture

exploding Gr

≡ GRU

$$\begin{pmatrix} i : \text{input gate} \\ f : \text{forget gate} \\ o : \text{output gate} \\ g : \text{gate gate} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ \tanh \end{pmatrix} N(h_{t-1}, x_t) \quad \begin{matrix} \text{Sigmoid } 0 \leq \leq 1 \\ -1 \leq \leq 1 \end{matrix}$$

$$\text{Cell state: } C_t = f \odot C_{t-1} + i \odot g$$

$$\text{hidden state: } h_t = o \odot \tanh(C_t)$$

\odot : multiply element-wise

BACKPROP:

