

# Recurrent Neural Network

Jun-ichi Sato

May 16, 2016

# Outline

- Review of last seminar
- Neural Networks
  - Activation function
  - Feed forward calculation
  - Backpropagation
- Recurrent Neural Networks
  - Structure
  - Calculation
- Example of RNN application

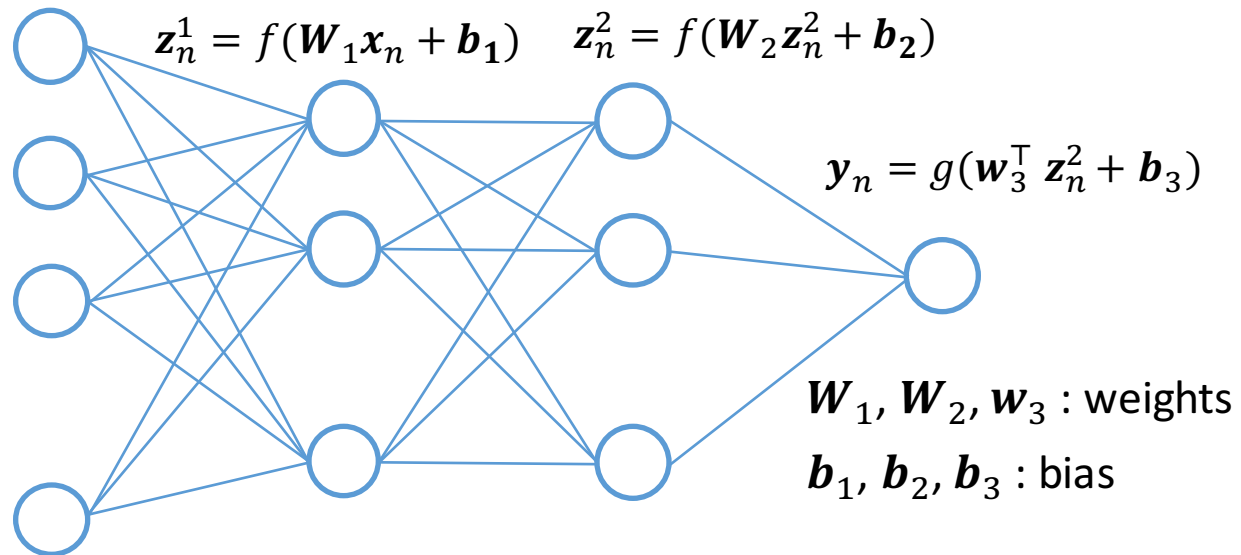
# Review of last seminar

- We proposed the method of **neural decoding** of c-VEP
- **Spatio-temporal inverse filter** for neural decoding
  - Linear model (LMSE, lasso)
  - Nonlinear model (Neural network)
- We adopt the feed **forward neural network** for the nonlinear model
- Feature works : more appropriate network for neural network
  - e.g. **Recurrent Neural Network**

# Neural Network

- **Feed-forward** structure
  - Input
  - Hidden
  - Output
- Model for **Classification** or **Regression**

$$\mathbf{x}_n = x_i[n + \tau + \tau_0]$$



# History of Neural Network

year	Events
1940s	Research started
1980s	<b>Backpropagation [1]</b>
1980s (late)	Convolutional neural network (CNN)
1990s	<b>Recurrent neural networks (RNN) [2]</b>
2000s	Deep belief networks (DBN) [3]

[1] D. E. Rumelhart and J. McClelland. "Parallel Distributed Processing: Explorations in the Microstructure of Cognition.", MIT Press, 1986

[2] C. Goller, "Learning task-dependent distributed representations by backpropagation through structure." Neural Networks, 1996., IEEE International Conference on. Vol. 1. IEEE, 1996.

[3] G. E. Hinton, "A fast learning algorithm for deep belief nets." Neural computation 18.7 (2006): 1527-1554.

# Neural Network in Computer Vision

## ImageNet Large-Scale Visual Recognition Challenge (ILSVRC)

- Classification problem
- 1000 class
- 1.2 million training

Year	Method	Accuracy
2010	Fast descriptor coding, large-scale SVM	72 %
2011	Fisher vector	75 %
2012	<b>Deep Learning (AlexNet)</b>	<b>85 %</b>
2013	<b>Deep Learning</b>	<b>89 %</b>
2014	<b>Deep Learning (GoogLeNet)</b>	<b>93 %</b>
2015	<b>Deep Learning</b>	<b>95.18 %</b>

- Human (hard studied) : 95 %

Beyond the human



# Feed forward Neural Network

- Each unit calculates the multi-inputs and the output
- Total input is calculated by

$$u = w_1x_1 + w_2x_2 + \cdots + w_nx_n + b$$

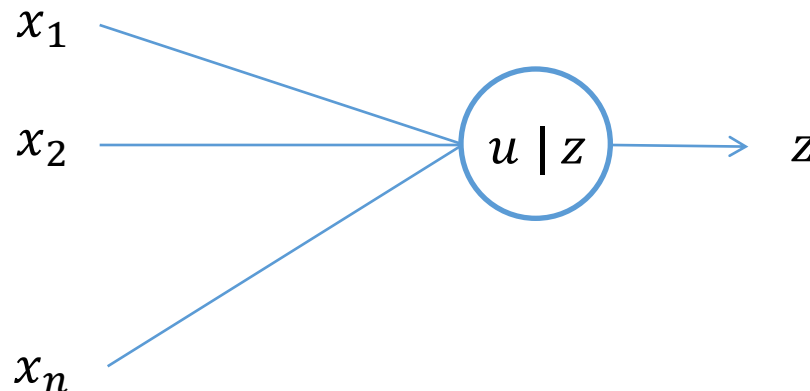
$w_i$  : weight of input

$b$  : bias

- Output is given by

$$z = f(u)$$

$f$  : Activation function



# Feed forward calculation of neural network

- Unit of first layer :  $i = 1, \dots, I$
- Unit of second layer :  $j = 1, \dots, J$
- The calculation is generalized by

$$u_j = \sum_i^I w_{ji} x_i + b_j$$
$$z_j = f(u_j)$$

- We can rewrite the equation by vector and matrix representation

$$\mathbf{u} = \mathbf{W}\mathbf{x} + \mathbf{b}$$

$$\mathbf{z} = \mathbf{f}(\mathbf{u})$$



# Multi-layer neural network

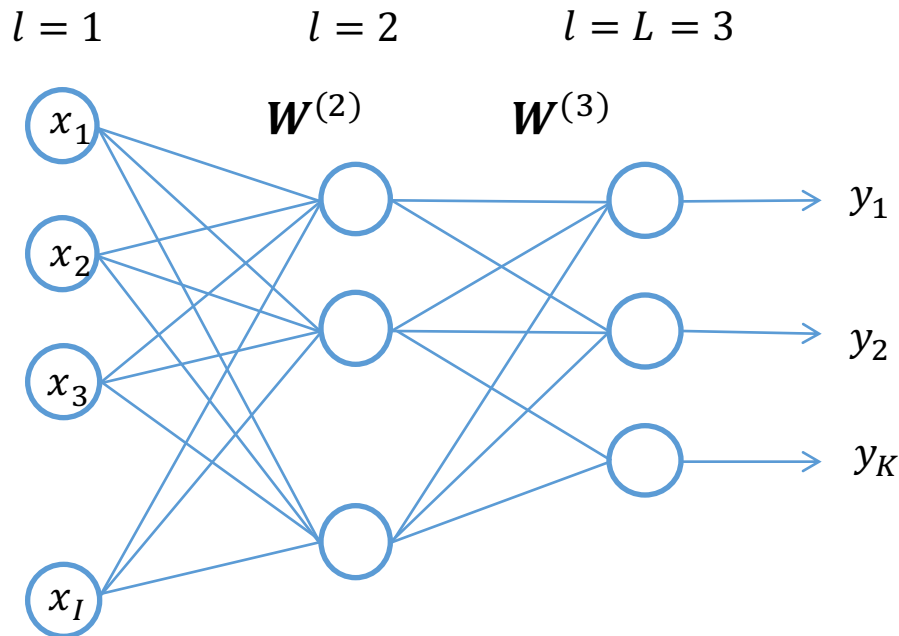
- $L = 3$  layer network

$$\mathbf{u}^{(l+1)} = \mathbf{W}^{(l+1)} \mathbf{z}^{(l)} + \mathbf{b}^{(l+1)}$$

$$\mathbf{z}^{(l+1)} = f(\mathbf{u}^{(l+1)})$$

- Output of the network is

$$\mathbf{y} \equiv \mathbf{z}^{(L)}$$



# Activation functions

- Logistic sigmoid function

- Most used
- domain :  $(-\infty, \infty)$ , range :  $(0, 1)$

$$f(u) = \frac{1}{1 + e^{-u}}$$

- Hyperbolic tangent

- Similar to logistic sigmoid function
- domain :  $(-\infty, \infty)$ , range :  $(-1, 1)$

$$f(u) = \frac{1}{1 + e^{-u}}$$

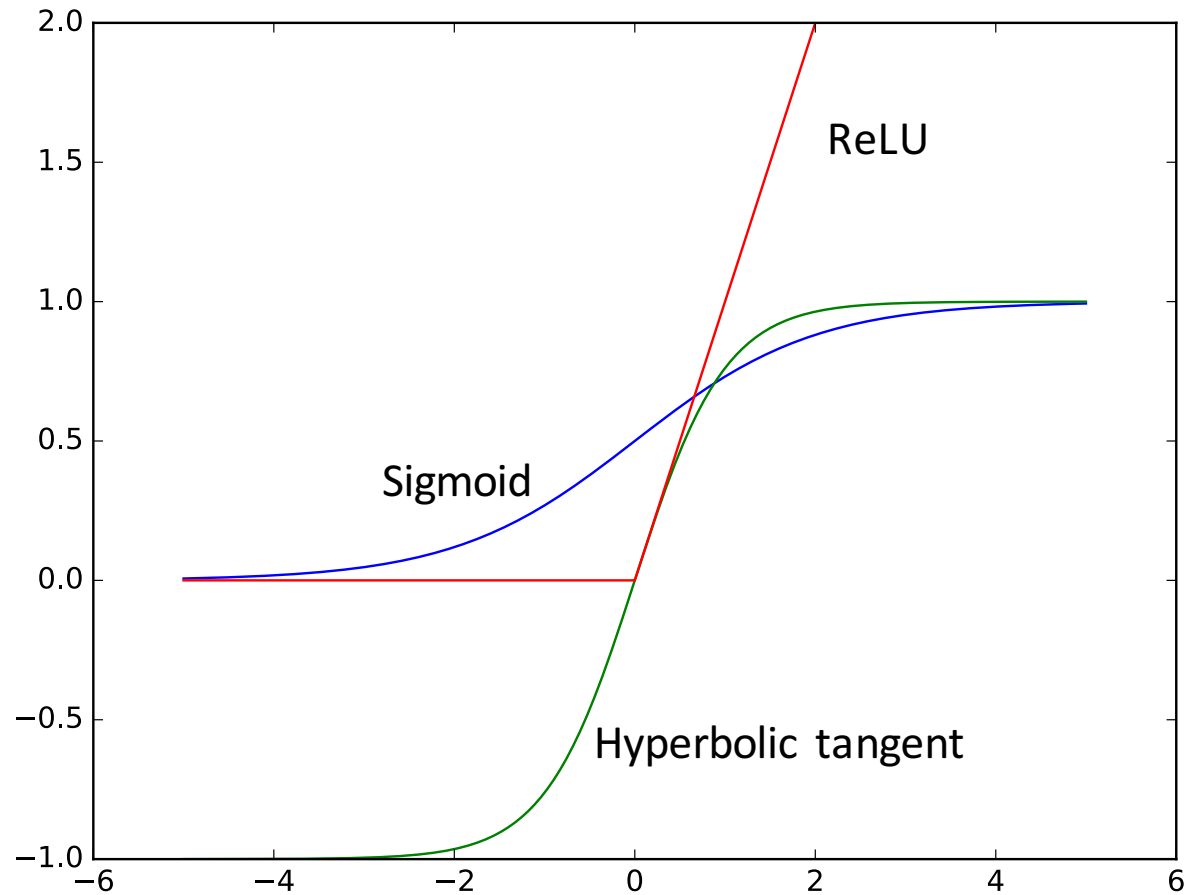
- Rectified linear function (ReLU)

- Used recently

$$f(u) = \max(u, 0)$$

- **Faster training** and **better results** than sigmoid and hyperbolic tangent

# Activation functions



# Design of output layer

- Linear activation function

$$f(u) = u$$

Problem	Activation function	Error function
Regression	Linear activation function	Squared error
Multi-class classification	Softmax function	Cross entropy

- Softmax function

$$y_k = \frac{\exp(u_k^L)}{\sum_{j=1}^K \exp(u_j^L)}$$

- Cross entropy

$$E(\mathbf{w}) = - \sum_{n=1}^N \sum_{k=1}^K d_{nk} \log y_k(\mathbf{x}_n; \mathbf{w})$$

- $\mathbf{d}_n$  : one-hot encoded vector (like  $[0, 0, 0, 1, 0]$ )

# Backpropagation

- Output layer error :

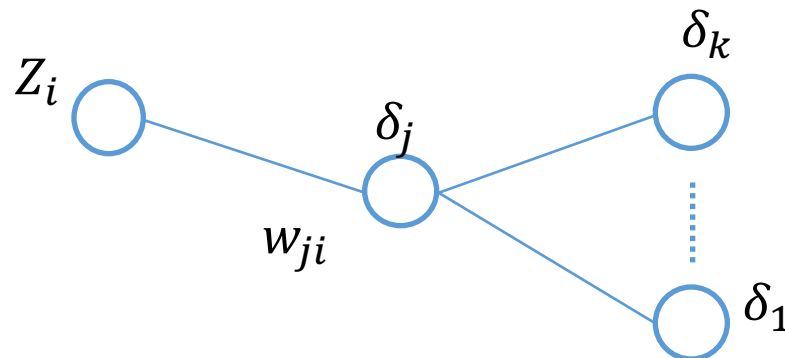
$$\delta_j \equiv \frac{\partial E}{\partial u_j}$$

- Backpropagation formula is given by

$$\delta_j = \sum_k w_{kj} \delta_k f'(u_j)$$

- Differential of  $E_n$  is given by

$$\frac{\partial E_n}{\partial w_{ji}} = \delta_j z_i$$



# Classification of sequential data

Sequential data :

- $x^1, x^2, \dots, x^T$
- $T$  is variable
- Ordered data
- E.g. voice, video, text

**Example : Predict the next word in the text**

We can get an idea of the quality of the learned **feature** vectors

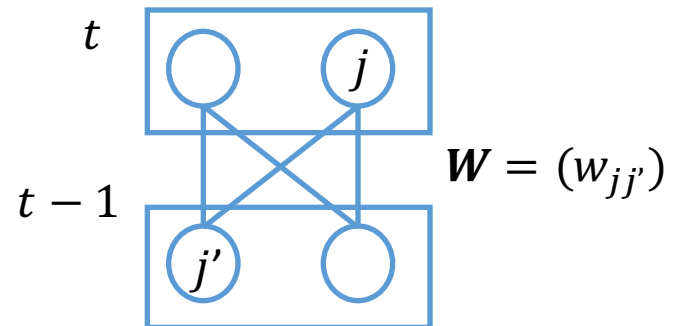
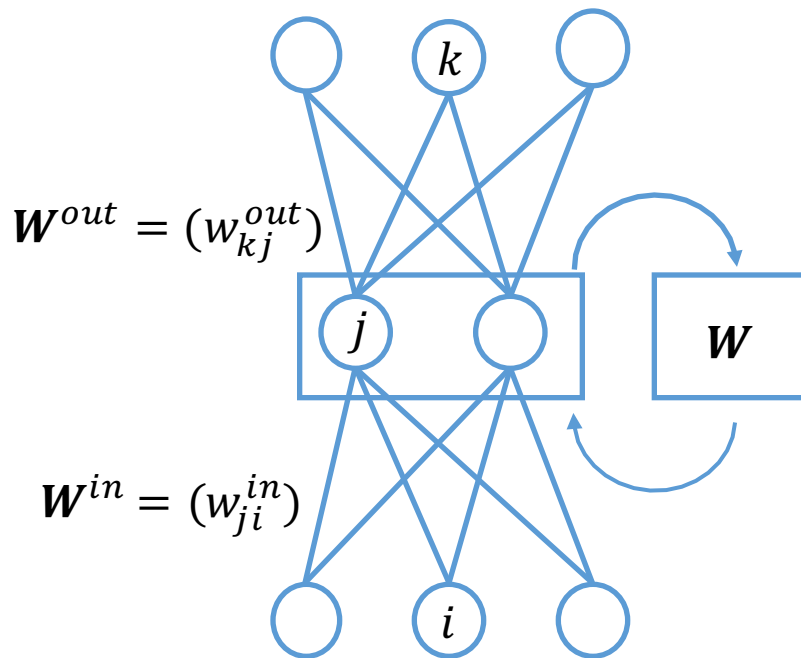


word	We	can	get	...	the	learned	?
input	$x^1$	$x^2$	$x^3$		$x^{t-1}$	$x^t$	$x^{t+1}$
output		$y^1$	$y^2$		$y^{t-2}$	$y^{t-1}$	$y^t$

Predict word  $y^t$  is influenced by previous words  $x^1, \dots, x^t$

# Recurrent Neural Network (RNN)

- RNN is a neural networks which has **closed circuits**
- RNN catches the **context** of a sequential data
- RNN receives  $x^t$  and outputs  $y^t$  for each time
- Theoretically, RNN outputs  $y^t$  from all inputs  $(x^1, \dots, x^t)$



# Design of output layer for RNN

- Same as feed forward network
- In classification problem, softmax function is used for activation function
- Error function is given by,

$$E(\mathbf{w}) = - \sum_{n=1}^N \sum_{t=1}^T \sum_{k=1}^K d_{nk}^t \log y_k^t(\mathbf{x}_n; \mathbf{w})$$

- $\mathbf{y}^1, \dots, \mathbf{y}^T$  : output sequence
- $\mathbf{d}^1, \dots, \mathbf{d}^T$  : target of the output
- $\mathbf{x}_n$  : input of  $n$  sample
- $d_{nk}^t$  : target of  $n$  sample at time  $t$



# RNN : Feed forward calculation

- Each unit is different state at each time  $t = 1, 2, \dots$

$\mathbf{x}^t = (x_i^t)$	input
$\mathbf{u}^t = (u_j^t)$	input of hidden layer
$\mathbf{z}^t = (z_j^t)$	output of hidden layer
$\mathbf{v}^t = (v_k^t)$	input of output layer
$\mathbf{y}^t = (y_k^t)$	output of output layer
$\mathbf{d}^t = (d_k^t)$	target of the output
$\mathbf{W}^{in} = (w_{ji}^{in})$	weight of the input to hidden layer
$\mathbf{W} = (w_{jj'})$	weight of feedback path
$\mathbf{W}^{out} = (w_{kj}^{out})$	weight of hidden to output

- $\mathbf{W}$  is constant value at any time

# RNN : Feed forward calculation

- Input of the hidden layer at  $t = t$

$$u_j = \sum_i w_{ji}^{(in)} x_i^t + \sum_{j'} w_{jj'} z_{j'}^{t-1}$$

- Sum of **inputs at  $t$**  and output of **hidden layer at  $t - 1$**
  - Bias is  $w_{j0}^{(in)}$
- Output of hidden layer

$$z_j^t = f(u_j^t)$$

$$\mathbf{z}^t = \mathbf{f}(\mathbf{W}^{(in)} \mathbf{x}^t + \mathbf{W} \mathbf{z}^{t-1})$$

- Output of output layer

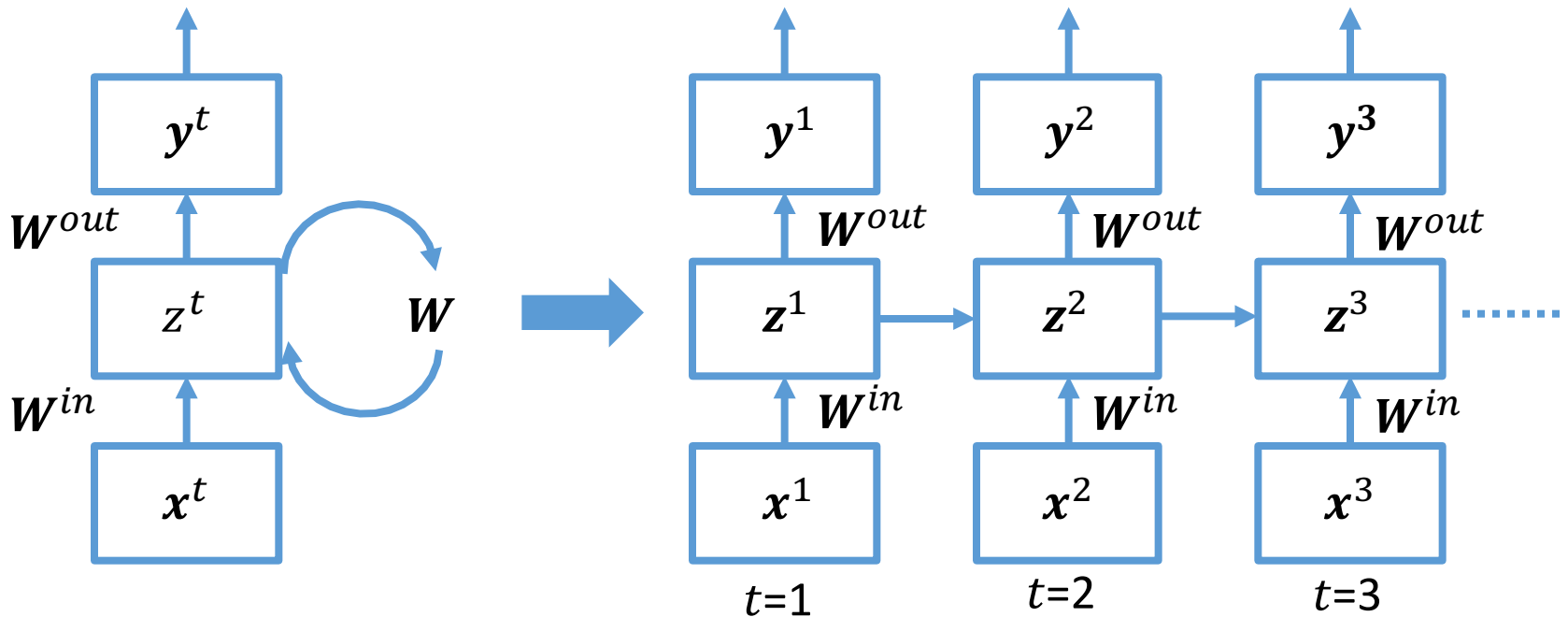
$$v_k^t = \sum_j w_{kj}^{(out)} z_j^t$$

$$\mathbf{y}^t = \mathbf{f}^{(out)}(\mathbf{v}^t) = \mathbf{f}^{(out)}(\mathbf{W}^{(out)} \mathbf{z}^t)$$

# RNN : Backpropagation (1)

## BPTT : backpropagation through time

- Replace RNN with feed forward NN
- Unfolding RNN through time



# RNN : Backpropagation (2)

- Backpropagation of feed forward networks

$$\delta_j^{(l)} = \sum_k w_{kj}^{(l+1)} \delta_k^{(l+1)} f'(u_j^{(l)})$$

- Output delta of unit  $k$  at  $t$

$$\delta_k^{out,t} \equiv \frac{\partial E}{\partial v_k^t}$$

- Hidden delta of unit  $j$  at  $t$

$$\delta_j^t \equiv \frac{\partial E}{\partial u_j^t}$$

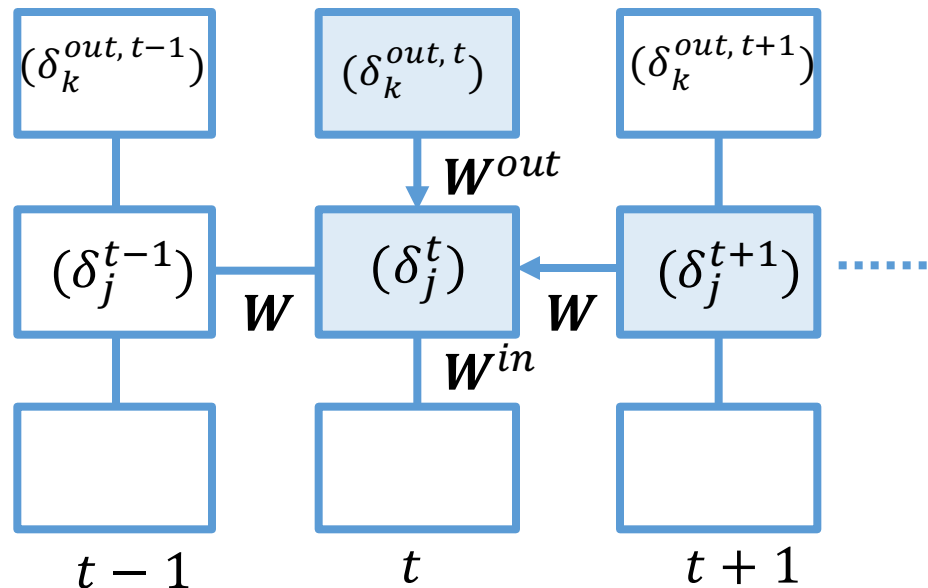
# RNN : Backpropagation (3)

- $\delta_j^t$  is given by

$$\delta_j^t = \left( \sum_k w_{kj}^{out} \delta_k^{out,t} + \sum_{j'} w_{jj'} \delta_{j'}^{t+1} \right) f'(u_j^t)$$

- Delta of  $T + 1$  is not able to calculate

$$\delta_j^{T+1} = 0$$



# RNN : Backpropagation (4)

- Differential of  $E$  are given by

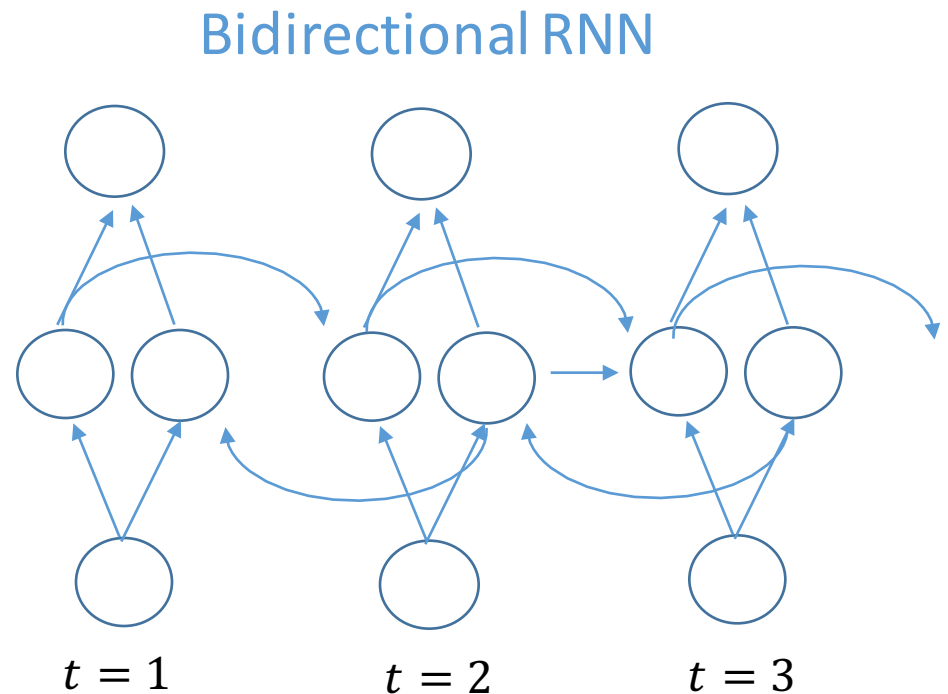
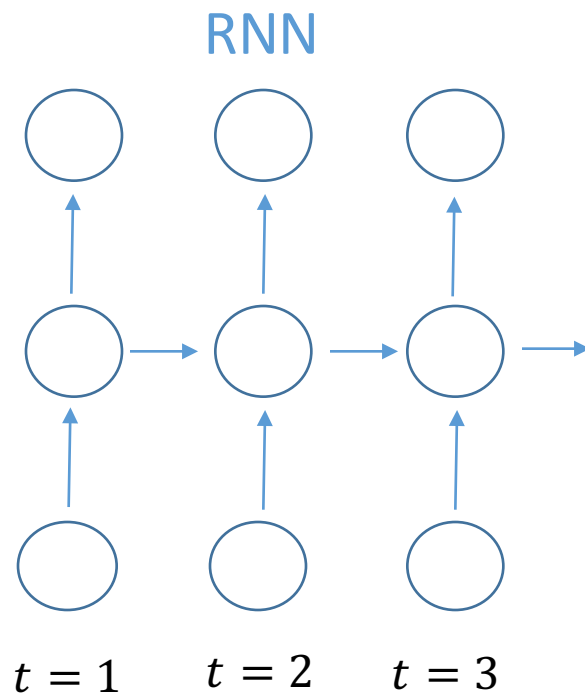
$$\frac{\partial E}{\partial w_{ji}^{in}} = \sum_{t=1}^T \frac{\partial E}{\partial u_j^t} \frac{\partial u_j^t}{\partial w_{ji}^{in}} = \sum_{t=1}^T \delta_j^t x_i^t$$

$$\frac{\partial E}{\partial w_{jj'}} = \sum_{t=1}^T \frac{\partial E}{\partial u_j^t} \frac{\partial u_j^t}{\partial w_{jj'}} = \sum_{t=1}^T \delta_j^t x_j^{t-1}$$

$$\frac{\partial E}{\partial w_{kj}^{out}} = \sum_{t=1}^T \frac{\partial E}{\partial v_k^t} \frac{\partial v_k^t}{\partial w_{kj}^{out}} = \sum_{t=1}^T \delta_j^t z_j^t$$

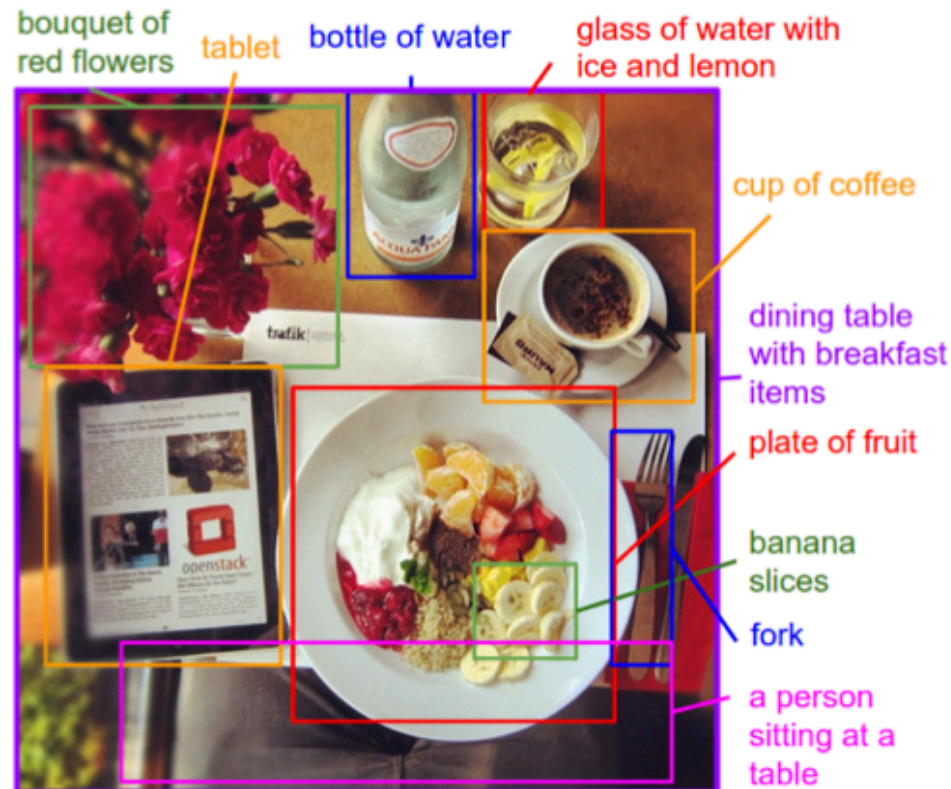
# Bidirectional RNN

- If all data is given at the same time,  
We can input the reversal data ( $t = T, T - 1 \dots, 1$ ) to the RNN
- B-RNN has better performance than normal RNN



# Example of RNN

- “Deep Visual-Semantic Alignments for Generating Image Descriptions”, K. Andrej and F. Li, CVPR, 2015
- Generate the description from the **image region**





# Proposed method

## Image representation : CNN

- CNN which is trained by **ImageNet**
- $I_b$ : pixel of image region
- Image representation is given by

$$v = W_m [CNN_{\theta_c}(I_b)] + b_m$$

## Sentence representation : BRNN

- $I_t$  : input vector, word of the number  $t$ . 1 of K representation
- $s_t$  : output vector

- **BRNN** is used to get output vector
- $W_w$  : weight of **word2vec**
- $f: x \mapsto \max(0, x)$  , ReLU

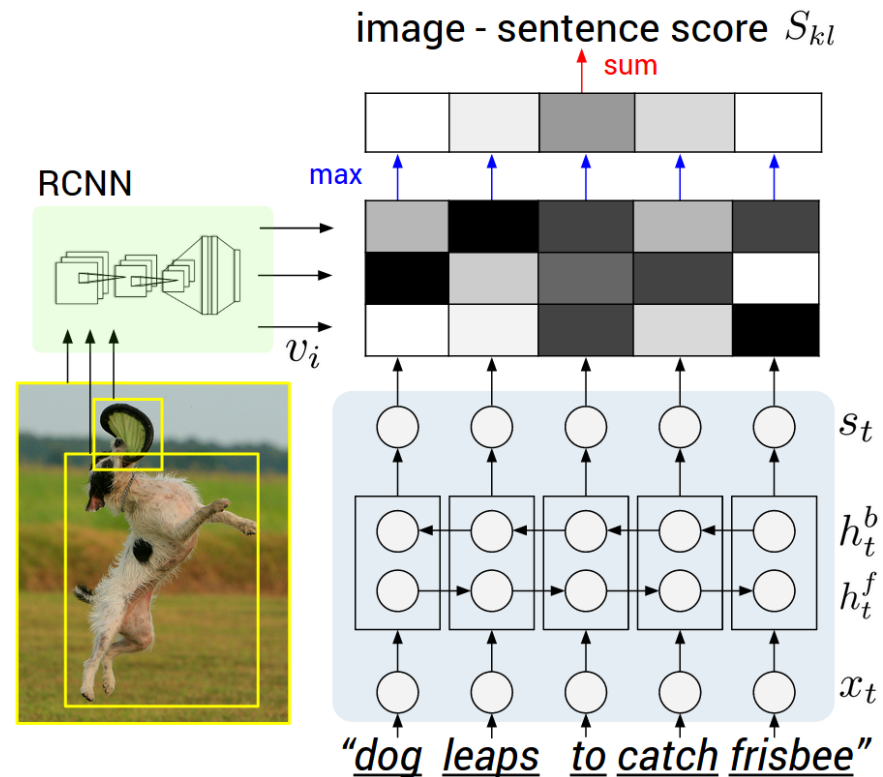
$$\begin{aligned}x_t &= W_w I_t \\e_t &= f(W_e x_t + b_e) \\h_t^f &= f(e_t + W_f h_{t-1}^f + b_f) \\h_t^b &= f(e_t + W_b h_{t+1}^b + b_b) \\s_t &= f(W_d (h_t^f + h_t^b) + b_d)\end{aligned}$$

# Score

- The score is given by

$$S_{kl} = \sum_{t \in g_l} \sum_{i \in g_k} \max(0, v_i^T s_t)$$

- $g_l$  : index set of sentence
- $g_k$  : index set of image region



# Result : Full frame

- “man in black shirt is playing a guitar”  
does not exist in training data sets
- “man in black shirt” and “is playing guitar” exists



man in black shirt is playing guitar.



construction worker in orange safety vest is working on road.

# Full frame : Failure

- Not young
- Not wakeboard



two young girls are playing with lego toy.

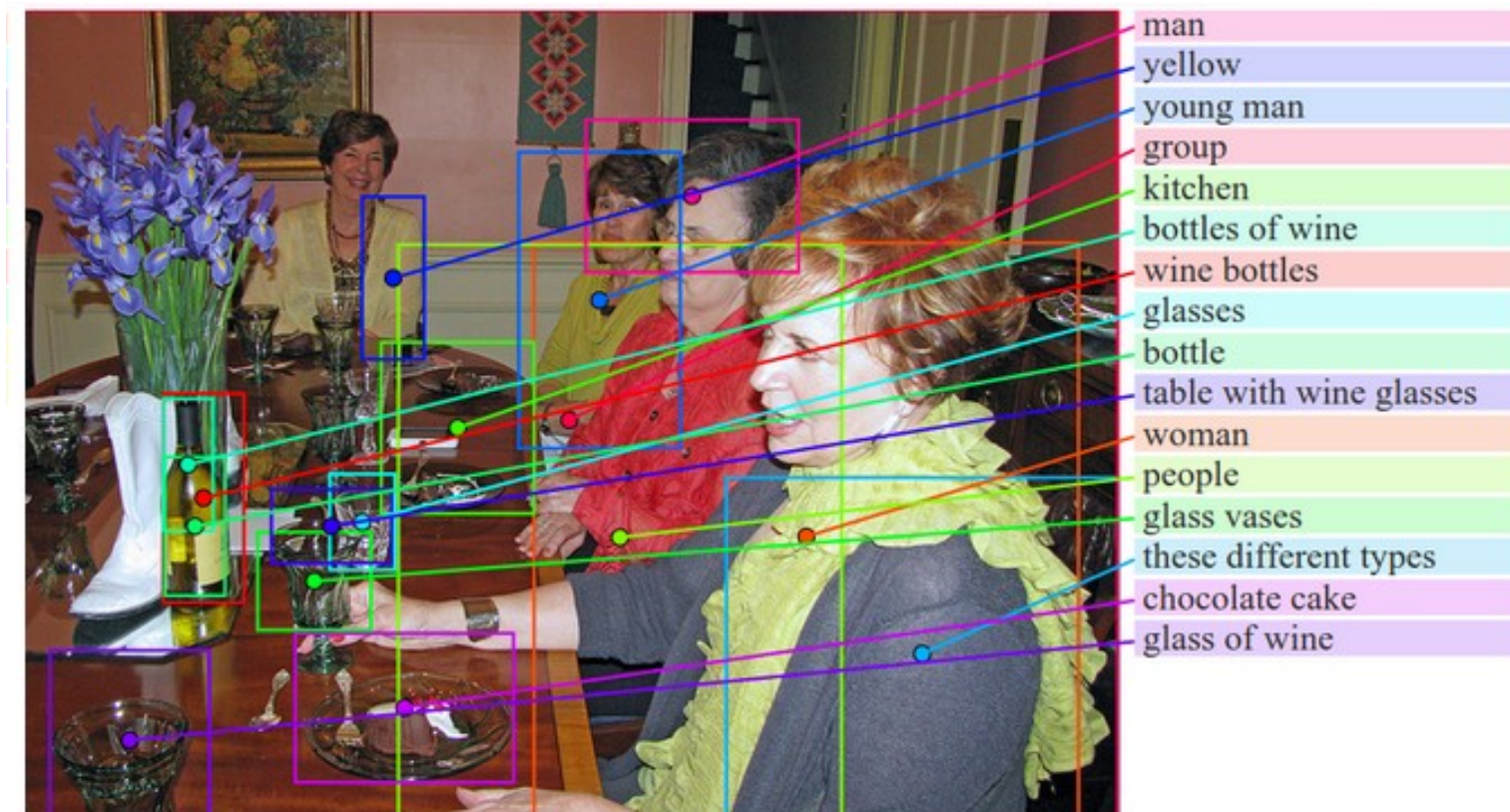


boy is doing backflip on wakeboard.



# Result : region-level descriptions

- “Table with wine glasses”
  - In training, it occurred in small region only 30 times.



# Future works

- To apply RNN to c-VEP BCI
- To study Long Short Term Memory (LSTM)
  - One of the RNN