# *Recurrent Neural Network (RNN)*
# Software Implementation using *Python*

June 29th, 2017

Presenter: Ryunosuke Murakami

Demonstration: Maiko Arakawa

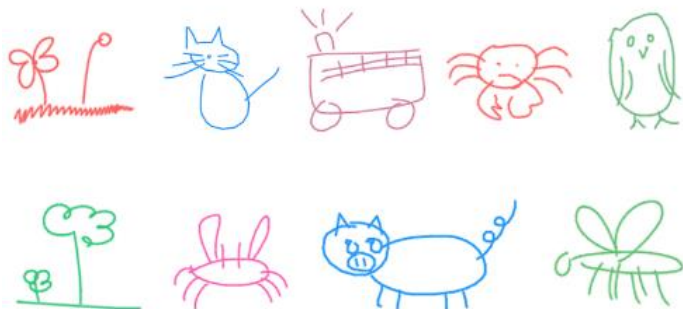Supporters: Hiroki Yomogita,

Tomohide Fukuchi, Yohei Shimmyo

# Outline

- Introduction

- Recurrent architecture

- Learning algorithm

    - Backpropagation in RNN.

    - Long-short term memory.

    - Connectionist Temporal Classification.

- Summary

- Demo

    - machine translation model

# What's RNN?

Recurrent Neural Network can handle
- – Voice data
- – Language data
- – Movie data



Sketch-rnn (Google)



DeepFix (Indian Institute of Science)

# Machine translation

*We can get an idea of the learned feature vectors by displaying.*

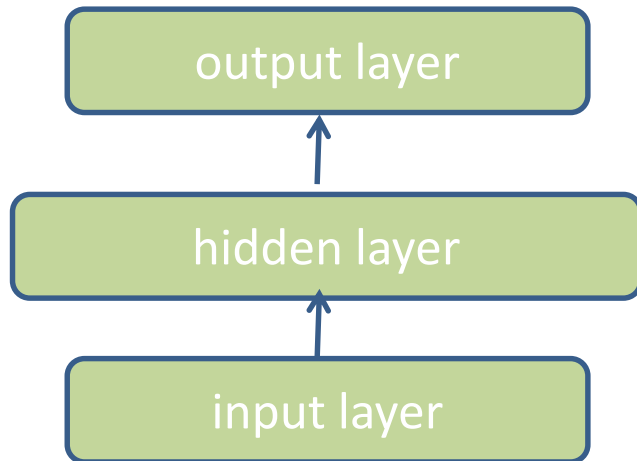| word | We | can | get | **...** | the | learned | ? |
|---|---|---|---|---|---|---|---|
| input | $x^1$ | $x^2$ | $x^3$ | | $x^{t-1}$ | $x^t$ | $x^t+1$ |
| output | | $y^1$ | $y^2$ | | $y^{t-2}$ | $y^{t-1}$ | $y^t$ |

Estimation task for next word from given inputs

$x^t$ : series of input

$y^t$ : *series of output*

# Difference of structure

**Feed-forward NN**
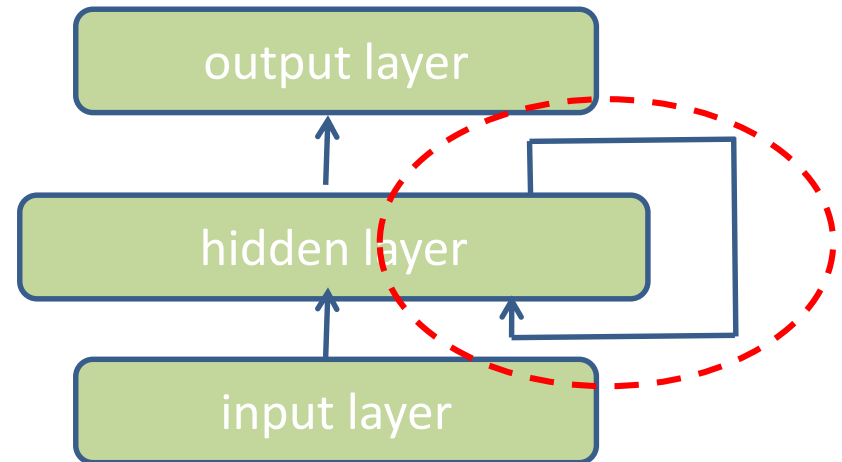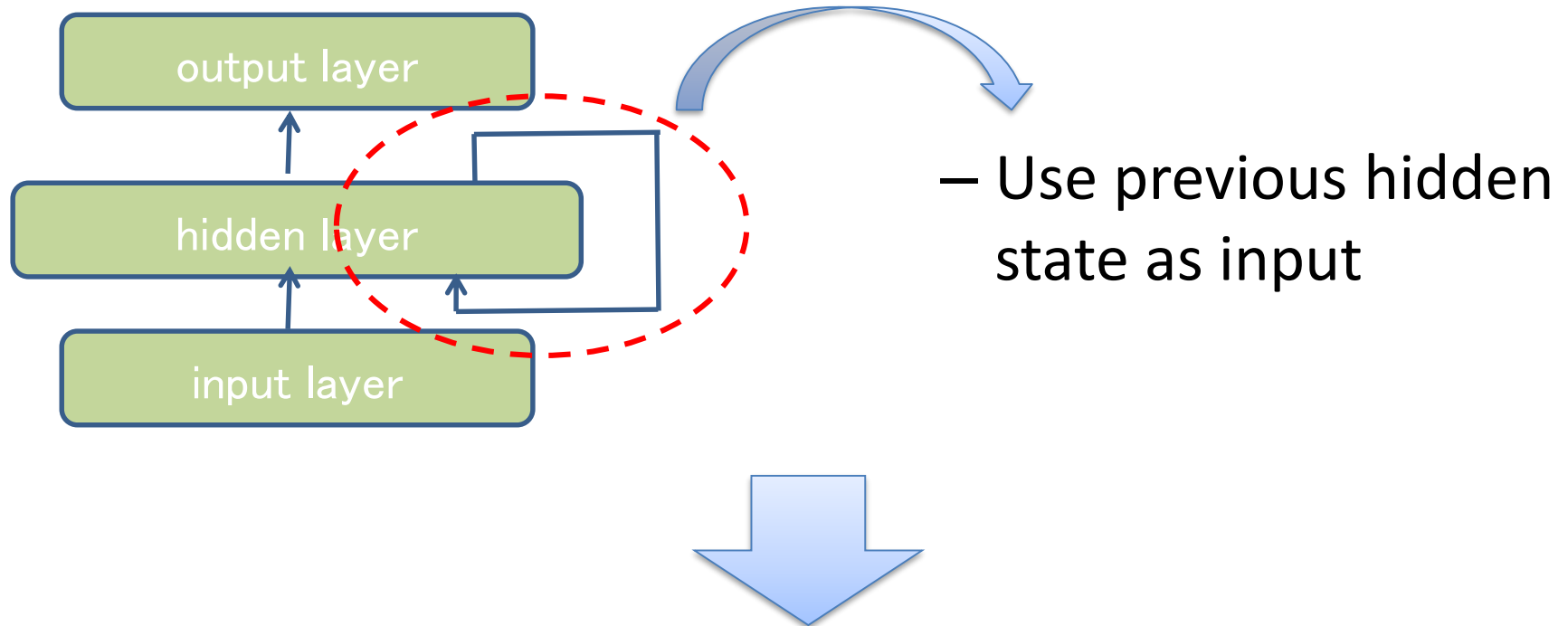
- Directed acyclic graph

**Recurrent NN**

- Directed cycle graph

# Why use RNN?

output layer

hidden layer

input layer

– Use previous hidden state as input
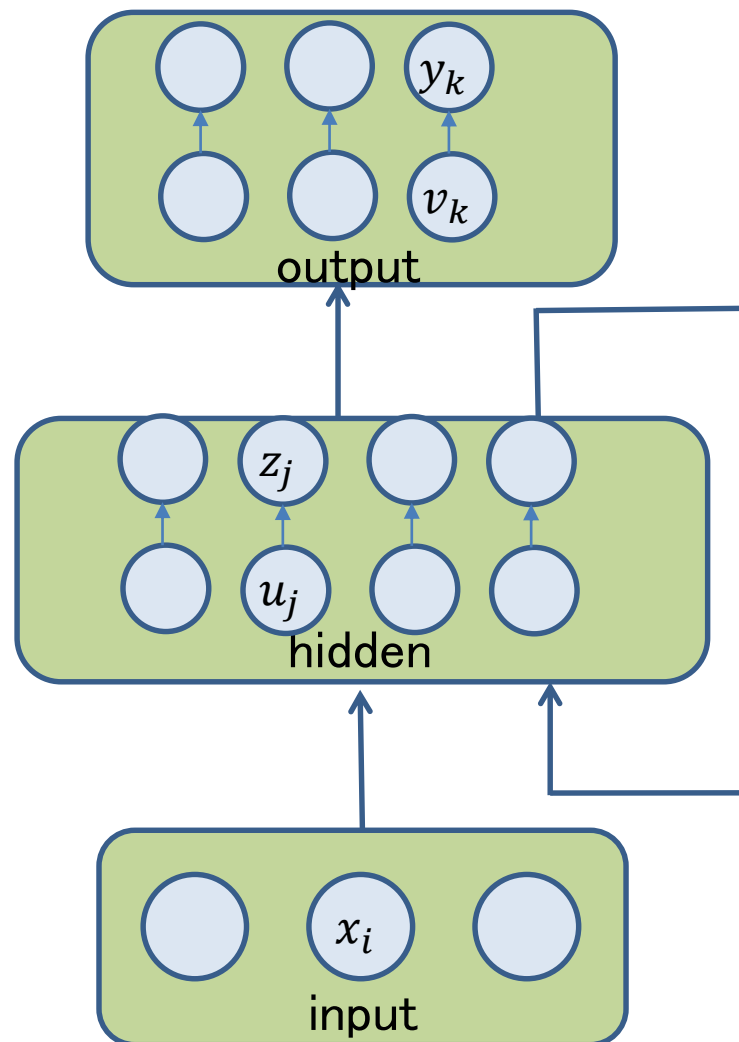
- Learn context and dependence between words
- High accuracy of word expectation

# Forward propagation of RNN

- Input of hidden layer
  - $u_j^t = \sum_i w_{ji}^{(in)} x_i^t + \sum_{j'} w_{jj'} z_{j'}^{t-1}$
- Output of hidden layer
  - $z_j^t = f(u_j^t)$
- Input of output layer
  - $v_k^t = \sum_j w_{kj}^{(out)} z_j^t$
- Output of output layer
  - $y_k^t = f(v_k^t)$

where t is time

# Forward propagation of RNN

- Input of hidden layer
  - $u_j^t = \sum_i w_{ji}^{(in)} x_i^t + \sum_{j'} w_{jj'} z_{j'}^{t-1}$
- Output of hidden layer
  - $z_j^t = f(u_j^t)$
- Input of output layer
  - $v_k^t = \sum_j w_{kj}^{(out)} z_j^t$
- Output of output layer
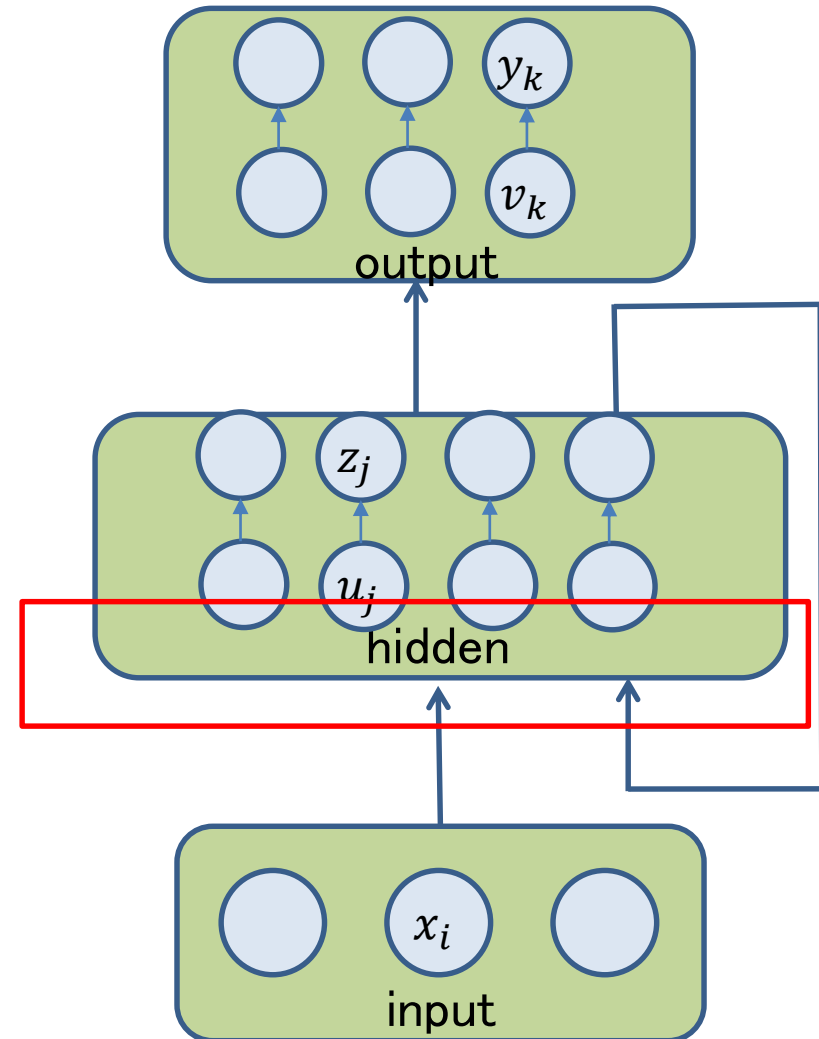  - $y_k^t = f(v_k^t)$

where t is time

# Detail of input of hidden layer
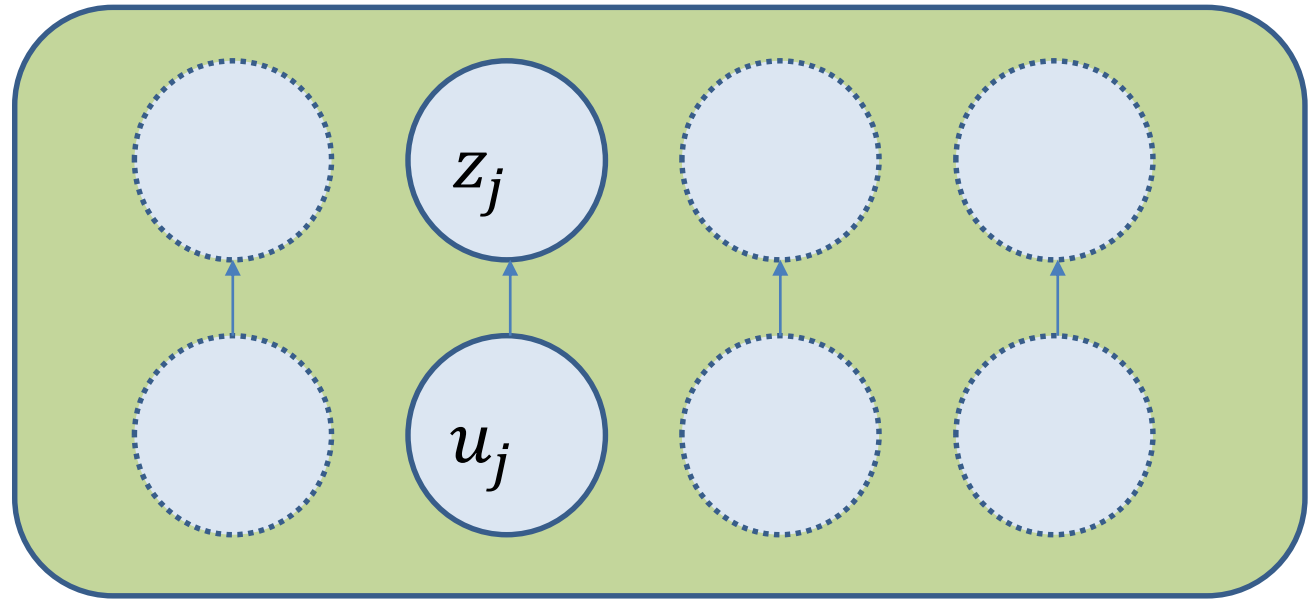
$$u_j^t = \sum_i w_{ji}^{(in)} x_i^t + \sum_{j'} w_{jj'} z_{j'}^{t-1}$$

# Detail of input of hidden layer



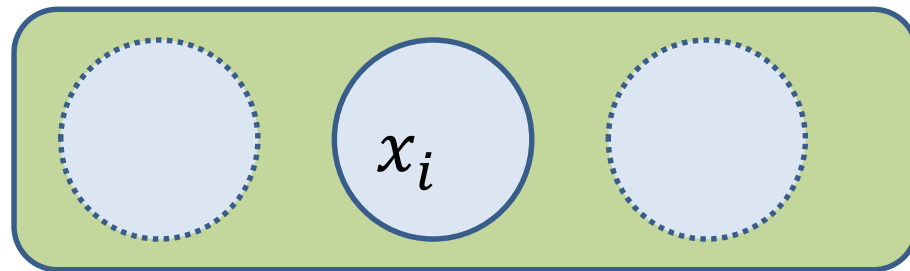$$u_j^t = \sum_i \textcolor{red}{w_{ji}^{(in)} x_i^t} + \sum_{j'} w_{jj'} z_{j'}^{t-1}$$

# Detail of input of hidden layer



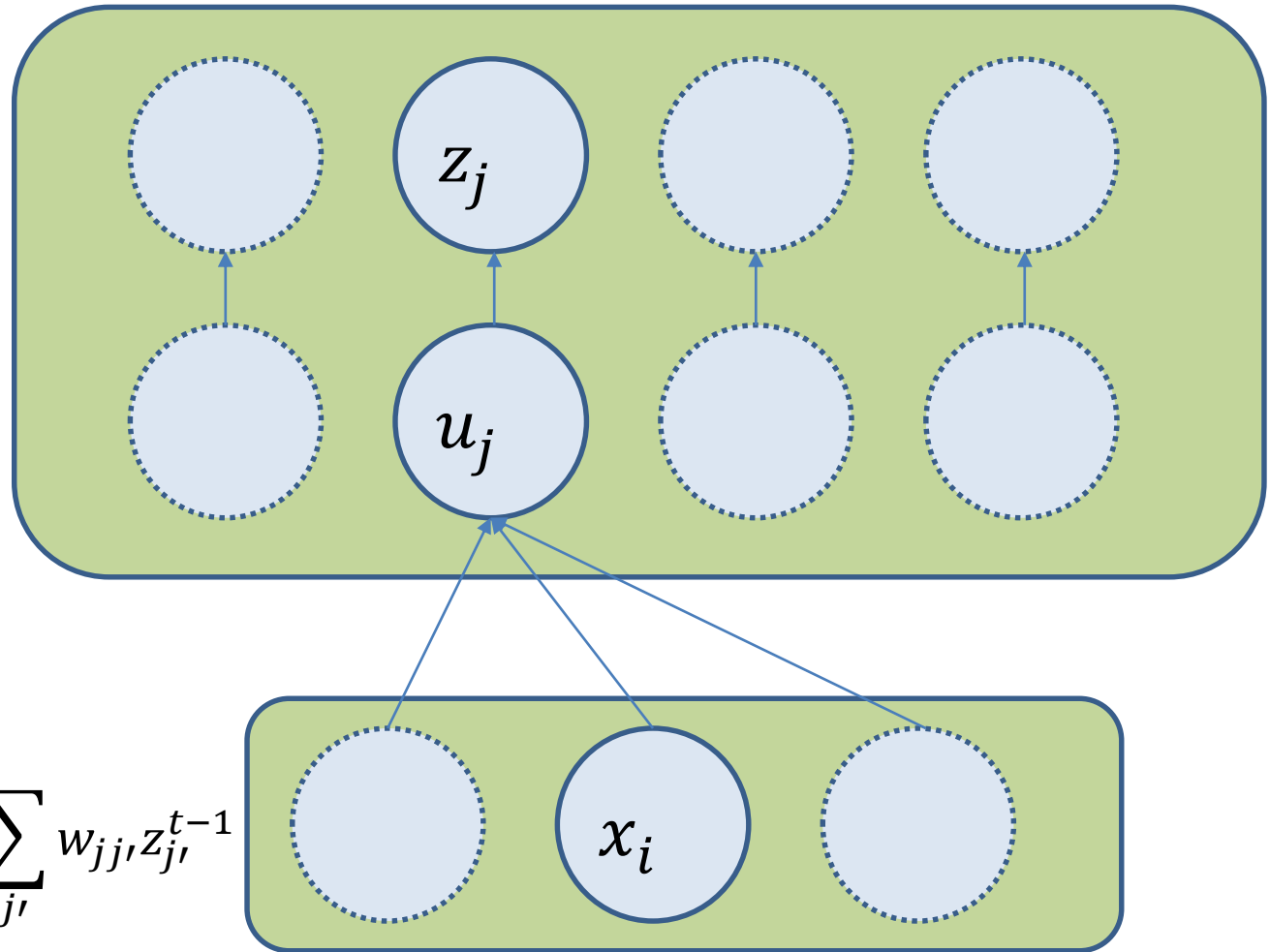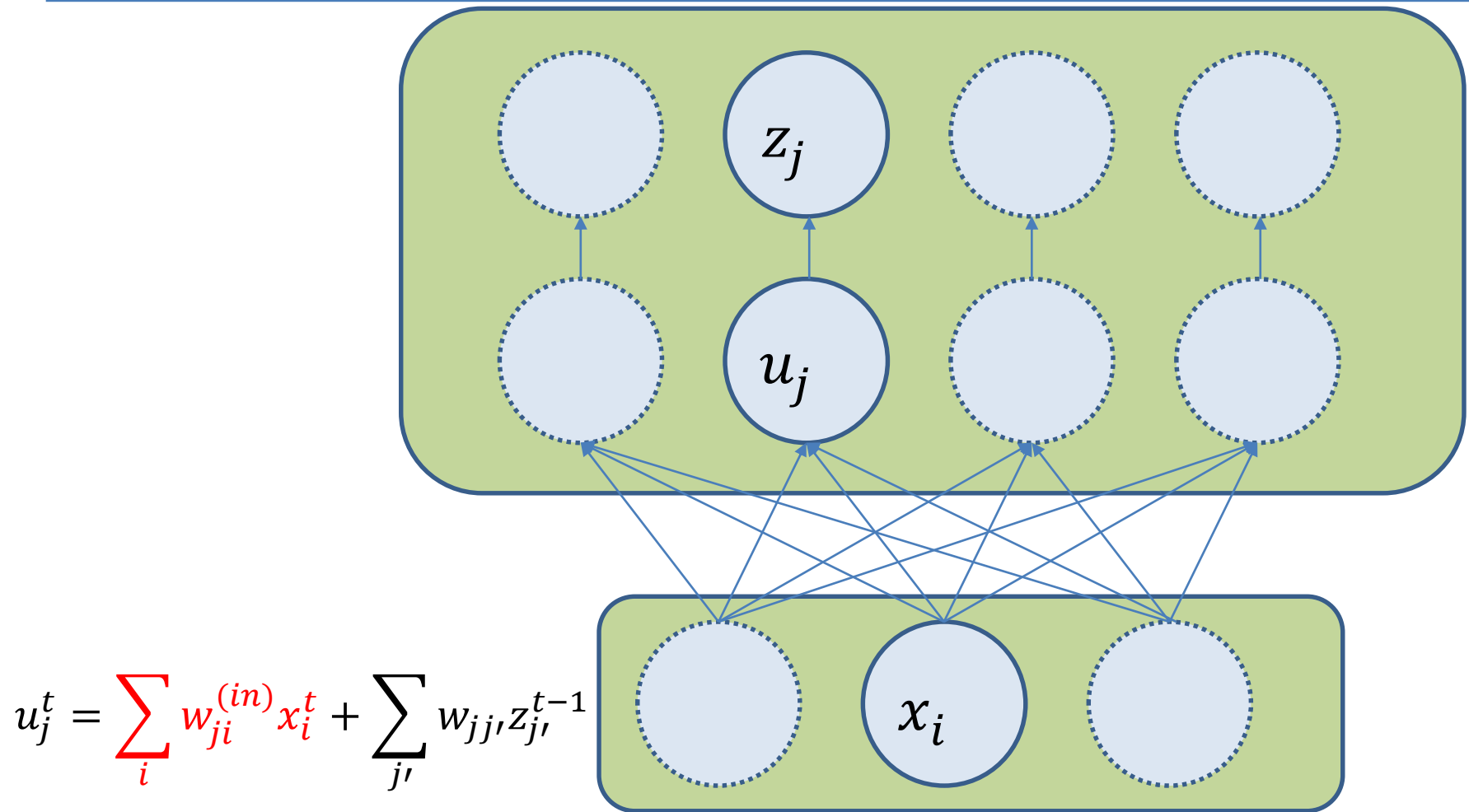$$u_j^t = \sum_i w_{ji}^{(in)} x_i^t + \sum_{j'} w_{jj'} z_{j'}^{t-1}$$

# Detail of input of hidden layer



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

# Detail of input of hidden layer



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

# Detail of input of hidden layer



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

# Forward propagation of RNN

- Input of hidden layer
  - $u_j^t = \sum_i w_{ji}^{(in)} x_i^t + \sum_{j'} w_{jj'} z_{j'}^{t-1}$
- Output of hidden layer
  - $z_j^t = f(u_j^t)$
- Input of output layer
  - $v_k^t = \sum_j w_{kj}^{(out)} z_j^t$
- Output of output layer
  - $y_k^t = f(v_k^t)$

where t is time

# Detail of output of hidden layer



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

# Forward propagation of RNN

- Input of hidden layer
  - $u_j^t = \sum_i w_{ji}^{(in)} x_i^t + \sum_{j'} w_{jj'} z_{j'}^{t-1}$
- Output of hidden layer
  - $z_j^t = f(u_j^t)$
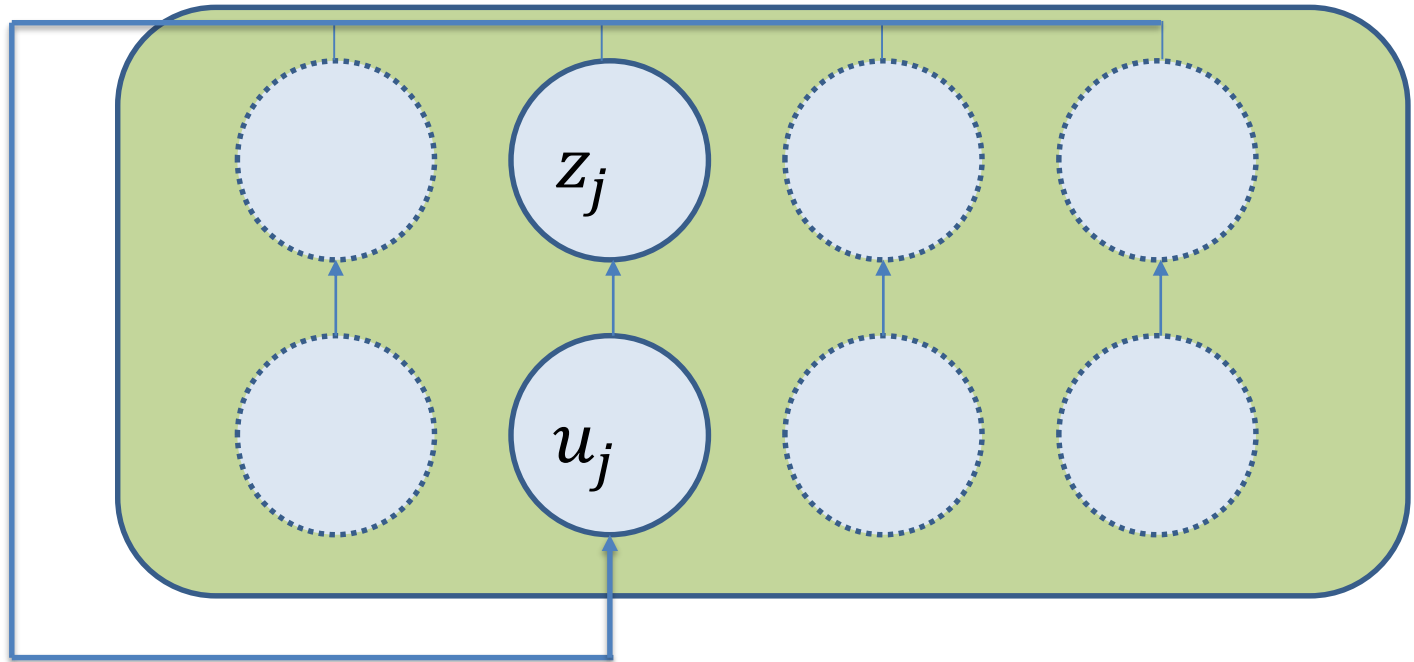- Input of output layer
  - $v_k^t = \sum_j w_{kj}^{(out)} z_j^t$
- Output of output layer
  - $y_k^t = f(v_k^t)$

where t is time

# Detail of input of output layer



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

output

$y_k$

$v_k$

$z_j$

$u_j$

# Forward propagation of RNN
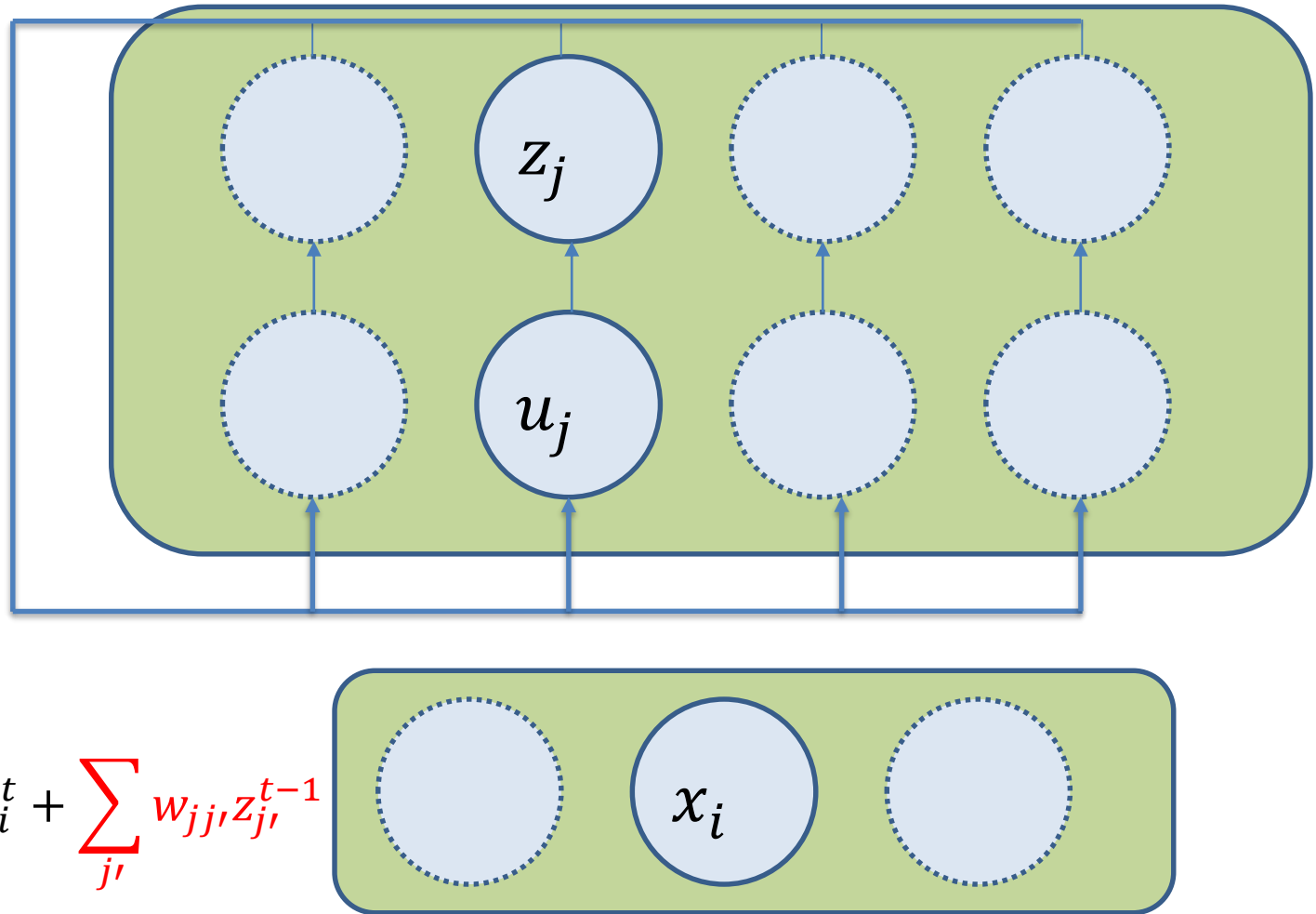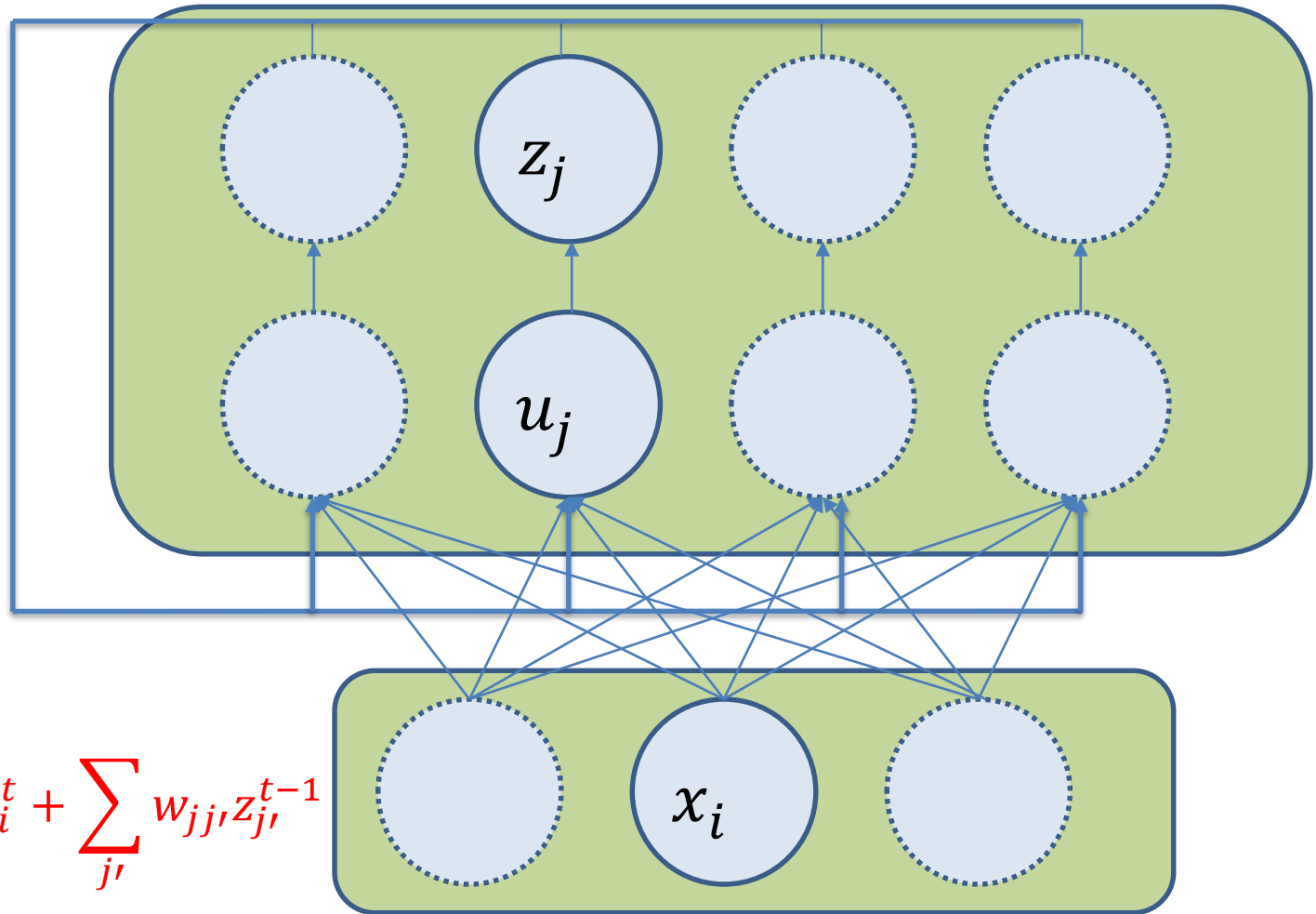
- Input of hidden layer
  - $u_j^t = \sum_i w_{ji}^{(in)} x_i^t + \sum_{j'} w_{jj'} z_{j'}^{t-1}$
- Output of hidden layer
  - $z_j^t = f(u_j^t)$
- Input of output layer
  - $v_k^t = \sum_j w_{kj}^{(out)} z_j^t$
- Output of output layer
  - $y_k^t = f(v_k^t)$

where t is time

# Detail of output on output layer



$$y_k^t = f(v_k^t)$$

# Forward propagation of RNN

- Input of hidden layer
  - $u_j^t = \sum_i w_{ji}^{(in)} x_i^t + \sum_{j'} w_{jj'} z_{j'}^{t-1}$
- Output of hidden layer
  - $z_j^t = f(u_j^t)$
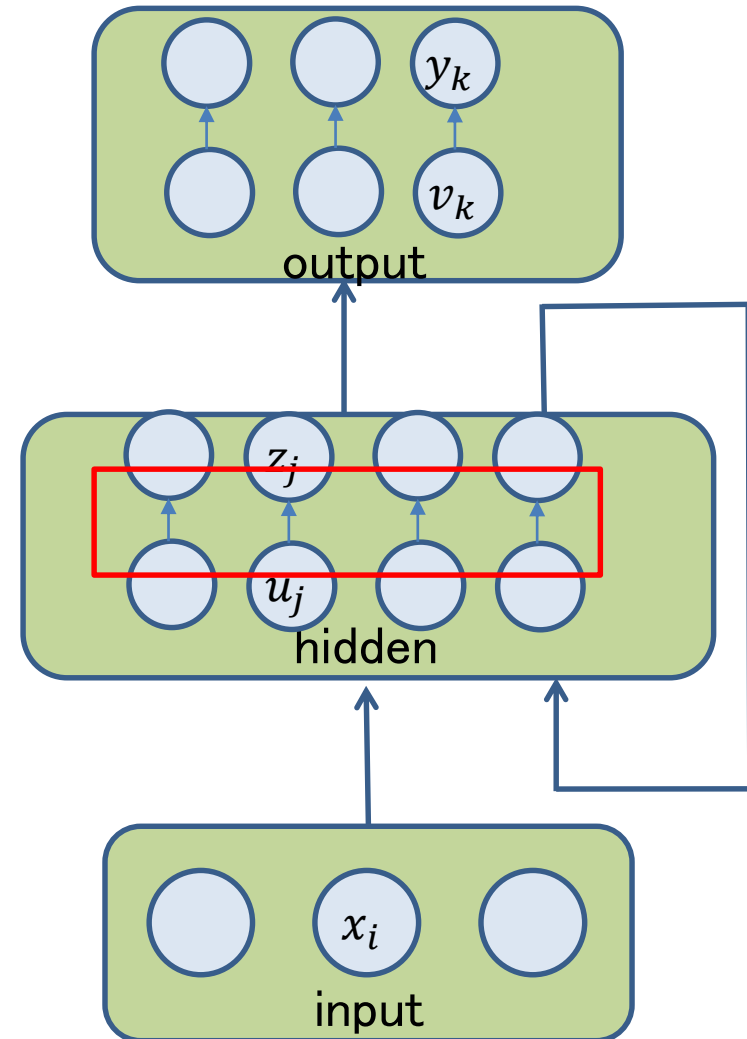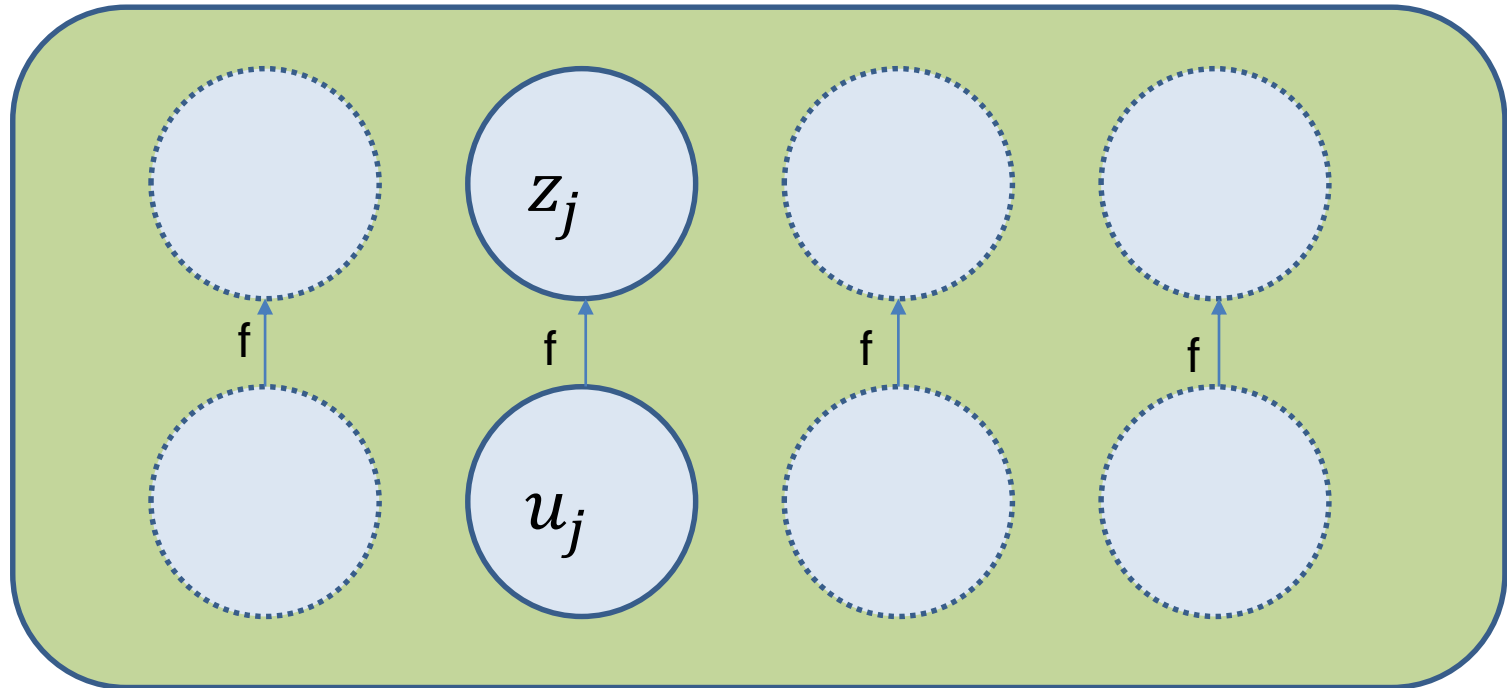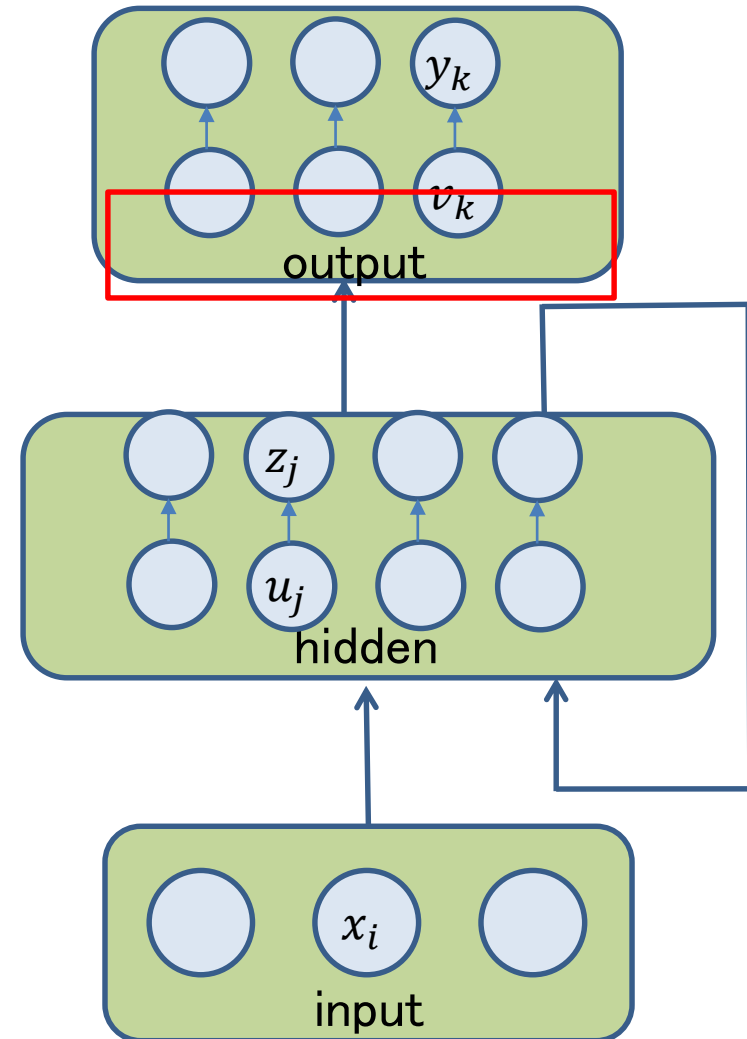- Input of output layer
  - $v_k^t = \sum_j w_{kj}^{(out)} z_j^t$
- Output of output layer
  - $y_k^t = f(v_k^t)$

where t is time

# Backpropagation of RNN

Backpropagation of RNN has 2 methods:

- Real time recurrent learning(RTRL)

- Backpropagation through time(BPTT)

**RTRL**

- memory efficient
- possible to learn real-time
- limited learning capability

**BPTT**

- fast to calculate
- high performance of learning capability
- Limited learning cycle

# Expand network through time direction



$y^t$

$W^{out}$   output

$z^t$   $W$

$W^{in}$   hidden

$x^t$

input

$t$

RNN network

$y^1$   $y^2$   $y^3$   $y^t$

$W^{out}$   $W^{out}$   $W^{out}$   $W^{out}$

$z^1$   $z^2$   $z^3$   $\cdots$   $z^t$

$W$   $W$   $W$

$W^{in}$   $W^{in}$   $W^{in}$   $W^{in}$

$x^1$   $x^2$   $x^3$   $x^t$

$t = 1$   $2$   $3$   $t$

Expanded network through time

# BPTT

- Find δ of output layer.

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

$l$: number of layers

- Find δ of hidden layer.

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

- Calculate gradient of error

$$\Delta w_{ij}(t) = -\eta \sum_{\tau=t_0} \delta_i(r+1) y_j(t)$$

$(\delta_k^{out,t-1})$  $(\delta_k^{out,t})$  $(\delta_k^{out,t+1})$

$W^{out}$

$(\delta_j^{t-1})$  $(\delta_j^{t})$  $(\delta_j^{t+1})$

$W$  $W$

$W^{in}$

$t-1$   $t$   $t+1$

# Problem on *RNN* learning

Length of the series data
determines the performance of RNN

The longer data length become,
the longer networks become

Numerous layers network causes
*vanishing gradient* problem

# What is *vanishing gradient*?

- vanishing gradient
  - Assume Sigmoid for activation functions
  - The more layers, the slower learning speed
  - The reason of slowing is "*gradient vanishment*"
- exploding gradient
  - *Gradient explode*

# Why gradient vanishes?

- The more layer, more number of  sigmoid is applied

# Why gradient vanishes?

- The more sigmoid is applied, the more its shape becomes flat -> gradient vanishes



y = sigmoid(x)



y = sigmoid(sigmoid(x))

output become too small

https://rechneronline.de/function-graphs/

# Long-short term memory(LSTM)

- *Deeper layer may cause vanishing gradient*
  - *RNN has deep network when it expand through time*
  - *It makes difficult to learn for long series data*
- *LSTM* is a solution for *vanishing gradient problem*
  - *memorize series data for long time*
- LSTM has only difference unit in hidden layer
  - Just replace the conventional unit with *memory-unit*

# Oblivion gate

Output layer

Input layer



- $unit\ f\ outputs\ g_j^{F,t}$
- $s_j^{t-1}$ is multiplied by $g_j^{F,t}$
- When the output is
  - Close to 0 > reset (oblivion)
  - Close to 1 > keep the state

$$g_j^{F,t} = f\left(u_j^{F,t}\right)$$
$$= f\left(\sum_i w_{ji}^{F,in}\, x_i^t + \sum_{j'} w_{jj'}\, z_{j'}^{t-1} + w_j^F s_j^{t-1}\right)$$

Input layer

hidden layer

# Input gate

Output layer | Input layer

$z_j^t$
$w_{ji}^{0,in}$ $w_{jj'}^0$

$g_j^{0,t}$

Output gate

d

e

$s_j^{t-1}$

$s_j^{t-1}$ $s_j^t$ $s_j^t$

$g_j^{F,t}$

f $s_j^{t-1}$

a

Oblivion gate

$g_j^{I,t}$

$w_{ji}^{F,in}$ $w_{jj'}^F$

c

Input gate

b

$w_{ji}^{in} w_{jj'}$

$w_{ji}^{I,in}$ $w_{jj'}^I$

- Unit $c$ outputs $g_j^{I,t}$
- unit $b$ receives input from outside and then multiplied by $g_j^{I,t}$
- The value is transmitted into memory cell

$$g_j^{I,t} = f\left(u_j^{I,t}\right)$$
$$= f\left(\sum_i w_{ji}^{I,in} x_i^t + \sum_{j'} w_{jj'}^I z_{j'}^{t-1} + w_j^I s_j^{t-1}\right)$$

Input layer | hidden layer

# Output gate

$z_j^t$

$g_j^{O,t}$

$w_{ji}^{O,in}$  $w_{jj'}^O$

Output gate

d

e

$s_j^{t-1}$

$s_j^{t-1}$  $s_j^t$

$s_j^t$

$g_j^{F,t}$

f

a

$s_j^{t-1}$

Oblivion gate

$g_j^{I,t}$

c

Input gate

$w_{ji}^{F,in}$ $w_{jj'}^F$

b

$w_{ji}^{in} w_{jj'}$

$w_{ji}^{I,in}$ $w_{jj'}^I$

- Unit $d$ outputs $g_j^{O,t}$
- When the $g_j^{O,t}$
  - Close to 1 > outputs transmitted into outsides
  - Close to 0 > block

$$z_j^t = g_j^{O,t} f(s_j^t)$$

$$g_j^{O,t} = f(u_j^{O,t})$$
$$= f\left(\sum_i w_{ji}^{O,in} x_i^t + \sum_{j'} w_{jj'}^O z_{j'}^{t-1} + w_j^O s_j^t\right)$$

# Memory cell $a$

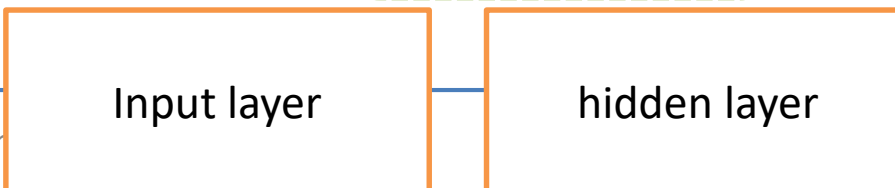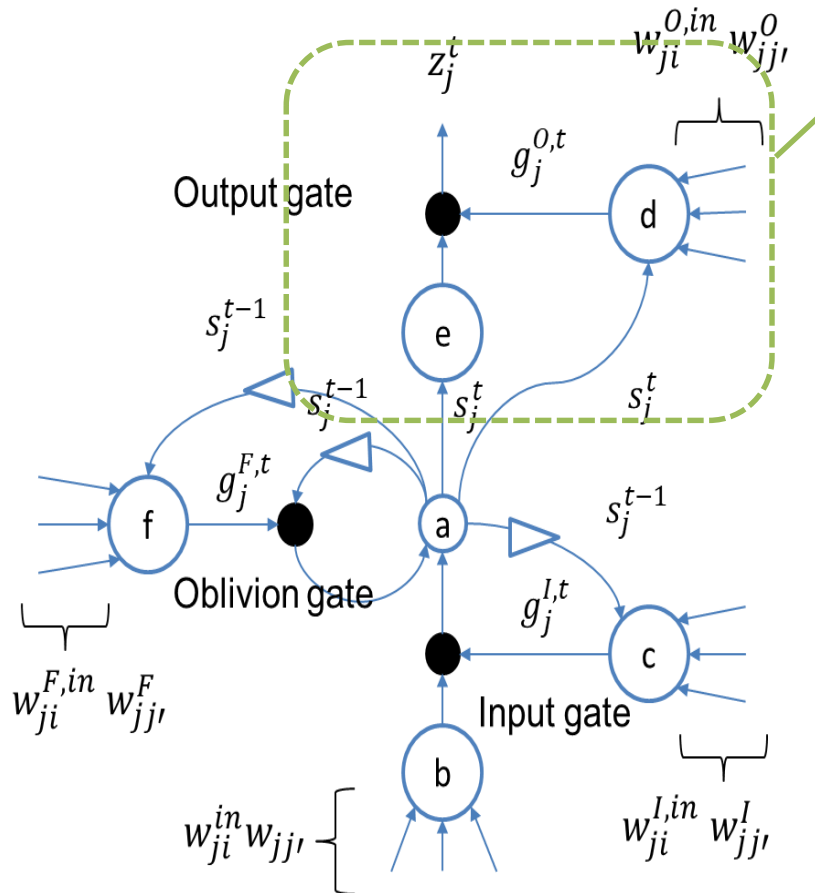- Memory cell (a) contains the state
- Memorization is realized by make its output return 1 time after
- In this time, the value is multiplied by output of unit (f)

$$s_j^t = g_j^{F,t} s_j^{t-1} + g_j^{I,t} f(u_j^t)$$

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

# Weakness of RNN

- RNN is unsuitable in the case: length of sequence data differs between input and output

- Input-length corresponds with output-length
  - RNN has to output $y_t$ from input $x_t$ in each time $t$

    *Ex) input: "abac", ouput: "xxyy"   but true answer = "xy"*

- Solution for this problem
  - Connectionist Temporal Classification

# Connectionist temporal classification

- CTC can solve classification tasks

  How to solve?

- Add output unit with '_' unit(vacant unit)
  - Expand the label by adding vacant value

# Ex. phoneme estimation

- $L = \{a, b\}$ ... *The set of labels to be recognized*
- $L' = \{a, b, \_\}$ ... *Label: L with vacant label '_'*

*There are countless numbers of redundant expression*

*ex. {a,_,_,b_,_} ... {a,a,a,_,_,_,b} ... {a,a,b,b,_,_,_,_} ...*

*l(Series data without redundancy) and $\pi$(with redundancy) have following "many-to-one" relation*

$$l = B(\pi) \qquad l = B(a\_\_b\_\_) = B(aaa\_\_\_b)$$

# Ex. phoneme estimation

The case: input length is 6

Adjust the length of the data

– Assume vacant label _

– Complement label with sequential vacant or same label

- $l = \text{``}ab\text{''}$

  *True label*

Length of input = 6

a,b,_,_,_,_
a,a,b,_,_,_
a,_,_,_,_,b
⋮
a,a,a,b,b,b

*Combination*

# Combination of estimates $\pi$



Each of the unit is applied *softmax* function for estimating probability

$t \rightarrow$

$l$     $l'$

$s$ $\downarrow$

$$\mathrm{B}^{-1} = \{\pi | B(\pi) = l\}$$

# Estimation



$$B^{-1} = \{\pi | B(\pi) = l\}$$

$$p(l|X) = \sum_{\pi \in B^{-1}(l)} p(\pi|X) \quad \text{sum of probability for total path}$$

$$p(\pi|X) = \prod_{t=1}^{T} y_{\pi t}^{t} \quad \text{Probability that the path } \pi \text{ is true (input is } X)$$

# $p(l|X)$ computation

- $p_l(ab|\,X) = $

$$p_l(a, b, \_, \_, \_, \_|X)$$
$$+$$
$$p_l(a,a,b,\_,\_,\_\,|X)$$
$$+$$
$$p_l(a,\_,\_,\_,\_,b\,|X)$$
$$\vdots$$
$$p_l(a,a,a,b,b,b\,|X)$$

$=$

$$y_a^1 \cdot y_b^2 \cdot y_-^3 \cdot y_-^4 \cdot y_-^5 \cdot y_-^6$$
$$+$$
$$y_a^1 \cdot y_a^2 \cdot y_b^3 \cdot y_-^4 \cdot y_-^5 \cdot y_-^6$$
$$+$$
$$y_a^1 \cdot y_-^2 \cdot y_-^3 \cdot y_-^4 \cdot y_-^5 \cdot y_b^6$$
$$\vdots$$
$$y_a^1 \cdot y_a^2 \cdot y_a^3 \cdot y_b^4 \cdot y_b^5 \cdot y_b^6$$

$p(\pi|X)$: each line of probability

## Highest value is estimated as correct $l$

# $p(l|X)$ computation

- $p(l|X) = \sum_{\pi \in B^{-1}(l)} p(\pi|X), \ ex. \, l = ab$

$p_l(a, b, \_, \_, \_, \_|X)$

$y_a^1 \cdot y_b^2 \cdot y_-^3 \cdot y_-^4 \cdot y_-^5 \cdot y_-^6$
+

$p_l(ab|X) \ needs \ huge \ amount \ of \ computation!$

- Highest value is the estimation of correct $l$

# Forward backward method

- Assume the set of path in time=$t$ and $s^{\text{th}}$ label
  - Former part of path to $t$ $\quad \pi_{1:t} = (\pi_1, \ldots, \pi_t)$
  - Latter part of path from $t$ $\quad \pi_{t:T} = (\pi_t, \ldots, \pi_T)$

- *Sum of probabilities*
  - *Former part of path* $p(\pi_{1:t})$ $\quad \alpha_{s,t}$
  - *Latter part of path* $\quad p(\pi_{t:T})$ $\quad \beta_{s,t}$

- *$\alpha$ and $\beta$ could be calculated recurrently*

# To calculate $p(l|X)$

$t \rightarrow$



$p(l|X) = $ sum of total path probability

# *Focus on the certain time*

# *Focus on the certain time*

$t \rightarrow$

1   2   3   4   5   6

$s$

$\downarrow$

1
2
3
4
5

path a

# *Focus on the certain time*

t →



path b

# *Focus on the certain time*



path c

path d

# *Focus on the certain time*



path e

# Forward backward method

$t \rightarrow$

$s$ ↓

$$p(l|X) = \text{path a + path b + } \dots \text{ path d}$$

# Forward backward method

*Sum of former part of path:* $\alpha_{s,t}$

$t$



*Latter part of path* $p(\pi_{t:T})$: $\beta_{s,t}$

$$p(l|X) = \sum_s \alpha_{s,t} + \beta_{s,t}$$

# Why forward backward is fast?

α$_{s,t}$ is calculated by recurrence relation

$\alpha(s-1, t-1)$

$\alpha(s, t-1)$

$y^t_{l'(s)}$

$$\alpha(s,t) = y^t_{l'(s)}\alpha(s,t) + y^t_{l'(s)}\alpha(s,t)$$

(pattern: s = '_')

# Why forward backward is fast?

$\alpha_{s,t}$ is calculated by recurrence relation



$$\alpha(s-2, t-1)$$

$$\alpha(s-1, t-1)$$

$$\alpha(s, t-1)$$

$$y^t_{l'(s)}$$

(b)

$$\alpha(s,t) = y^t_{l'(s)}\alpha(s-2, t-1) + y^t_{l'(s)}\alpha(s-1, t-1) + y^t_{l'(s)}\alpha(s, t-1)$$

## (pattern: s != '_')

t = 1

$$\alpha(1,1) = y^1_{l'(s)}$$

$$\alpha(2,1) = y^1_{l'(s)}$$

Accelerate the calculation of $\alpha$ with dynamic programming (as well as $\beta$)

# Summary

- RNN can memorize past state
- Length of series data determine performance
- Learning is conducted by expanding network through time-direction
- LSTM is solution for longer term memorization
- CTC can estimates likelihood label
  - Forward and backward method is efficient

# **Demo**

- Machine translation
  - Japanese to English

- Encoder decoder model

# Encoder-Decoder model



Tips:input sequence of sentence is reversion for better result

# Bilingual data

jp.txt

誰 が 一番 に 着く か 私 に は 分 かり ませ ん 。
十 中八九 彼 は 成功 する だろう 。
あなた の 銀行 口座 を 教え てい ただけ ます か 。

eng.txt

i can 't tell who will arrive first .

ten to one , he will succeed .

may we know your bank account ?

We use 10000 sentences for training
Word is separated by space

Encoder-Decoder translation model

# TRAINING

# Implementing Encoder-Decoder model

```
┌─────────────────┐      ┌─────────────────┐      ┌─────────────────┐
│  Read bilingual │  →   │  make word-id   │  →   │  Create model   │
│      data       │      │   dictionary    │      │                 │
└─────────────────┘      └─────────────────┘      └─────────────────┘
                                                           │
        ┌──────────────────────────────────────────────────┘
        ▼
┌─────────────────┐      ┌─────────────────┐
│ Optimize Setup  │  →   │ Learn for every │
│                 │      │    sentence     │
└─────────────────┘      └─────────────────┘
```

# Read data and make dictionary

**jp.txt**

```
jvocab = {}
jlines = open('jp.txt').read().split('¥n')
for i in range(len(jlines)):
    lt = jlines[i].split()
    for w in lt:
        if w not in jvocab:
            jvocab[w] = len(jvocab)

jvocab['<eos>'] = len(jvocab)
jv = len(jvocab)
```

**eng.txt**

```
evocab = {}
elines =
open('eng.txt').read().split('¥n')
for i in range(len(elines)):
    lt = elines[i].split()
    for w in lt:
        if w not in evocab:
            evocab[w] = len(evocab)

evocab['<eos>'] = len(evocab)
ev = len(evocab)
```

# Encoder-Decoder translation model graph

RNN
MyMT Class

$y_1$     $y_2$     …     <eos>
$y_{n+1}$

| | | | | softmax | softmax | | softmax |
|---|---|---|---|---|---|---|---|

W    W    W

LSTM → LSTM → … → LSTM → LSTM → LSTM → … → LSTM

$Embed_x$   $Embed_x$   …   $Embed_x$   $Embed_x$   $Embed_y$   …   $Embed_y$

$x_1$     $x_2$     …     $x_m$     $x_{m+1}$     $y_1$     …     $y_n$

```
class MyMT(chainer.Chain):
    def __init__(self, jv, ev, k):
        super(MyMT, self).__init__(
            embedx = L.EmbedID(jv, k),
            embedy = L.EmbedID(ev, k),
            H = L.LSTM(k, k),
            W = L.Linear(k, ev),
```

```python
def __call__(self, jline, eline):
    for i in range(len(jline)):
        wid = jvocab[jline[i]]
        x_k = self.embedx(Variable(np.array([wid], dtype=np.int32)))
        h = self.H(x_k)

    x_k = self.embedx(Variable(np.array([jvocab['<eos>']], dtype=np.int32)))
    tx = Variable(np.array([evocab[eline[0]]], dtype=np.int32))
    h = self.H(x_k)
    accum_loss = F.softmax_cross_entropy(self.W(h), tx)

    for i in range(len(eline)):
        wid = evocab[eline[i]]
        x_k = self.embedy(Variable(np.array([wid], dtype=np.int32)))
        next_wid = evocab['<eos>']  if (i == len(eline) - 1) else evocab[eline[i+1]]
        tx = Variable(np.array([next_wid], dtype=np.int32))
        h = self.H(x_k)
        loss = F.softmax_cross_entropy(self.W(h), tx)
        accum_loss += loss
    return accum_loss
```

# Create model and Setup optimizer

```
demb = 100
model = MyMT(jv, ev, demb)
optimizer = optimizers.Adam()
optimizer.setup(model)
```

# Learning

```
for epoch in range(100):
    for i in range(len(jlines)-1):
        jln = jlines[i].split()
        jlnr = jln[::-1]
        eln = elines[i].split()
        model.H.reset_state()
        model.zerograds()
        loss = model(jlnr, eln)
        loss.backward()
        loss.unchain_backward()  # truncate
        optimizer.update()
        print i, " finished"
    outfile = "mt-" + str(epoch) + ".model"
    serializers.save_npz(outfile, model)
```

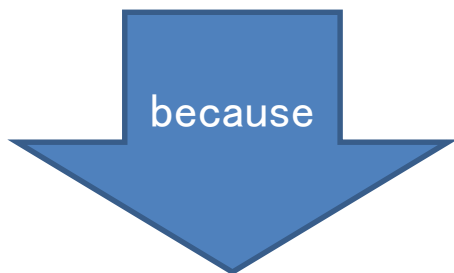Encoder-Decoder model

# TESTING

# Test Data

- Test data(100) + Training data(2000)
  = 2100 sentences

Uh...
Some sentences include unknown word.
Skip it.

# Result

- Accuracy is 0.077



because

- The number of training data is too small
- The low frequency word is forgotten.

# References

References

　　- 岡谷貴之. (2015)　「深層学習」

Demo

　　-「ChainerとRNNと機械翻訳」
　　<http://qiita.com/odashi_t/items/a1be7c4964fbea6a116e>

# Thank you for your listening!