





# COURSE 5: SEQUENCE MODELS

## Week 1

Examples of sequence data:

1. Speech recognition:   $\rightarrow$  "Hello hi."
2. Music generation:   $\rightarrow$  
3. Sentiment classification: "This movie sucks"  $\rightarrow$  ★☆☆☆☆
4. DNA sequence analysis: AGCCCTGATGAGGAAGTAG  $\rightarrow$  AGCCCTGATGAGGAAGTAG
5. Machine Translation: Voulez-vous chanter avec moi?  $\rightarrow$  Do you want to sing with me?
6. Video Activity recognition:   $\rightarrow$  Running
7. Name Entity recognition: Yesterday, Harry  $\rightarrow$  Harry, Hermione met Hermione

Notation:

$x$ : Harry Potter and Hermione invented a new spell.  
 $x^{(1)} \quad x^{(2)} \quad x^{(3)} \quad \dots \quad x^{(t)} \quad \dots \quad x^{(9)}$   
 $T_x = 8$

$y$ :  
 $y^{(1)} \quad y^{(2)} \quad y^{(3)} \quad \dots \quad y^{(t)} \quad \dots \quad y^{(9)}$   
 $T_y = 8$

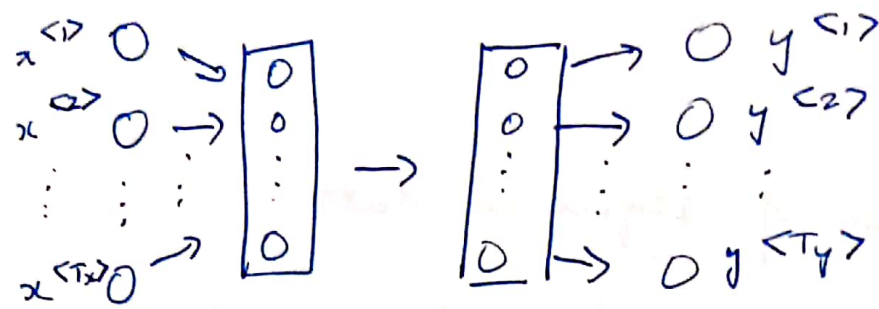
For training example  $i$ ,  $x^{(i)}$ ,  $y^{(i)}$ ,  $T_x^{(i)}$ ,  $T_y^{(i)}$   
The words are represented as vectors

Vocabulary

a	1
action	2
and	367
happy	4015
price	6630
Zulu	10,000

Then  $x^{(1)}$   $x^{(2)}$   $\leftarrow$  one-hot representation  
 $\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$  -4015  $\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$  -6630  
 $\uparrow$  10,000  
and so on for all 8 words

Why not a standard network?

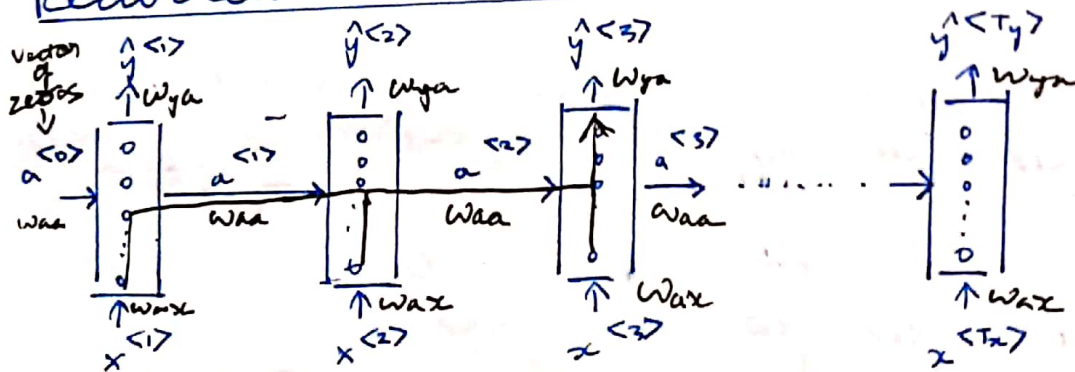


Problems:

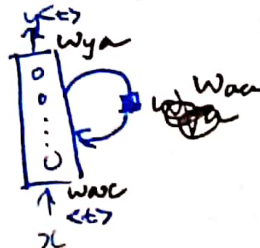
- Inputs, outputs can be different lengths in different examples (since the length of sentences vary).
- Doesn't share features learned across different positions of text (It can't share features learnt if Harry comes at different places in different sentences)

So solution is RNN.

## Recurrent Neural Network



- The parameters are shared
- Activations of previous words affects the prediction of the current word (as shown by black path)
- This is sometimes also represented as:



→ Problem: It only uses information from previous words (example: it doesn't use 4, 5, 6... Tx in the diagram for  $\hat{y}^{(1)}$ )

name	He said, "Teddy"	Roosevelt was a great President."
not a name	He said, "Teddy"	bears are on sale!"

So we can't determine by only knowing previous words, so we use Bidirectional RNN (BRNN)

Forward Propagation:

$$a^{(0)} = \vec{0}$$

$$a^{(1)} = g(w_{aa} a^{(0)} + w_{ax} x^{(1)} + b_a) \leftarrow \text{tanh is a common choice since it removes}$$

$$\hat{y}^{(1)} = g(w_{ya} a^{(1)} + b_y)$$

vanishing gradient problem.  
ReLU is used later often.

Sigmoid/softmax based on problem

$$a^{(t)} = g(w_{aa} a^{(t-1)} + w_{ax} x^{(t)} + b_a)$$

$$\hat{y}^{(t)} = g(w_{ya} a^{(t)} + b_y)$$

This can be written as,

$$a^{(t)} = g(w_a [a^{(t-1)}, x^{(t)}] + b_a)$$

$$\hat{y}^{(t)} = g(w_y a^{(t)} + b_y)$$

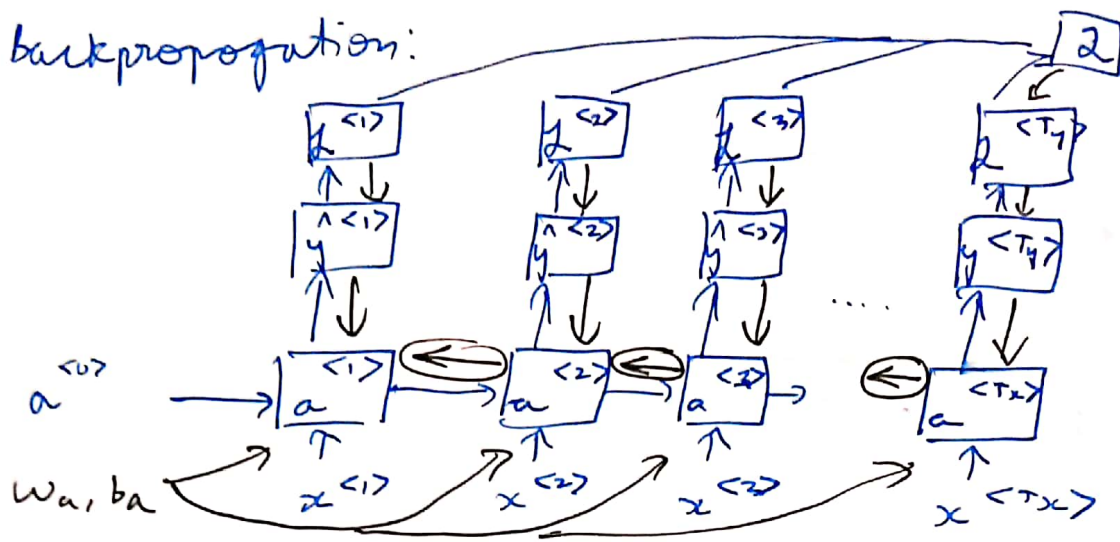
If  $w_{aa}$  is (100, 100) and  $w_{ax}$  is (100, 10,000)

$$w_a = \begin{matrix} \xrightarrow{100} \\ \begin{bmatrix} w_{aa} & w_{ax} \end{bmatrix} \\ \xleftarrow{100} \quad \xleftarrow{10,000} \end{matrix} \text{ and } [a^{(t-1)}, x^{(t)}] = \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} \begin{matrix} \xrightarrow{100} \\ \xleftarrow{10,000} \end{matrix}$$

$$\therefore [w_{aa} : w_{ax}] \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} = w_{aa} a^{(t-1)} + w_{ax} x^{(t)}$$



backpropagation:



← is backpropagation

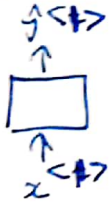
← is called backpropagation through time

$$L^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -y^{(t)} \log \hat{y}^{(t)} - (1 - y^{(t)}) \log (1 - \hat{y}^{(t)})$$

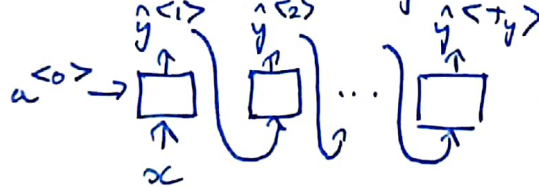
$$L(\hat{y}, y) = \sum_{t=1}^{T_x} L^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

Different Types of RNN:

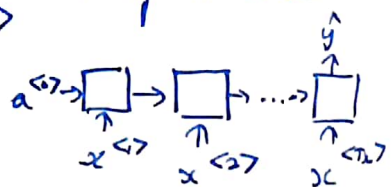
One-to-one:



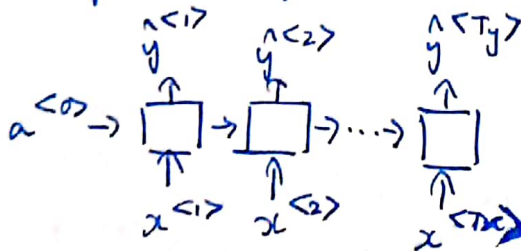
One to many:



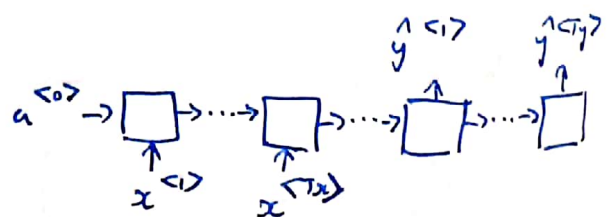
Many to one



Many to many:



Many to many:



Encoding  
(French)

Decoding  
(English)

Used in translation

What is language modelling?

Given a sentence, it'll tell you the probability of that sentence occurring.

1. The apple and pear salad  $P = 3.2 \times 10^{-13}$

✓ 2. The apple and pear salad  $P = 5.7 \times 10^{-10}$

$$P(\text{sentence}) = P(y^{<1>}, y^{<2>}, \dots, y^{<T_y>}) = ?$$

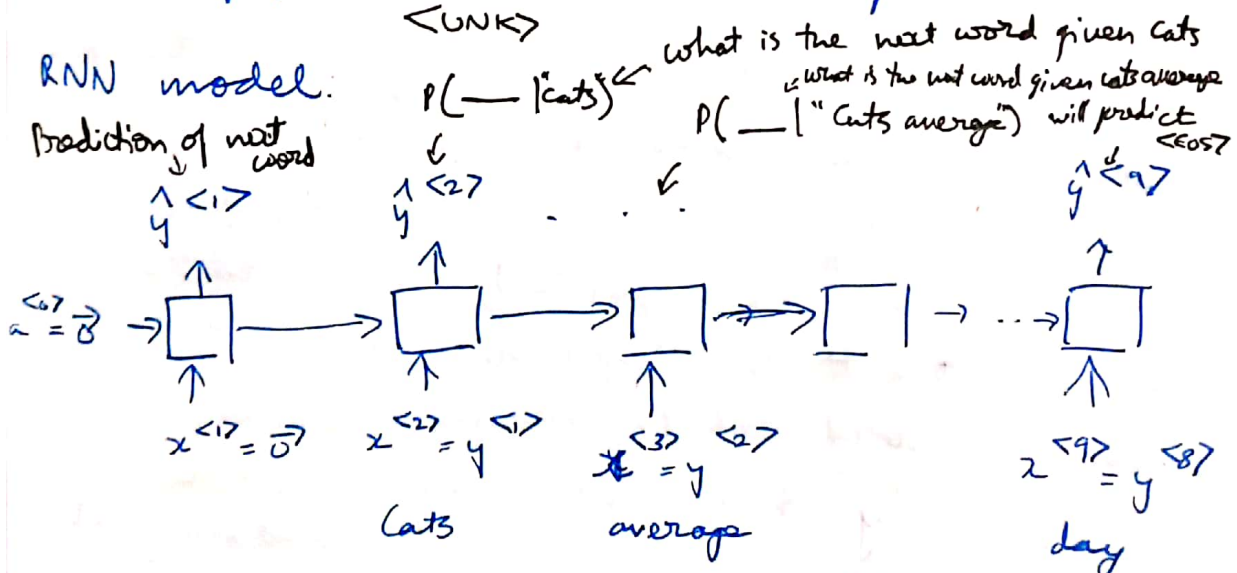
Language modelling with an RNN:

- Training set: large corpus of english text
- We tokenise the sentence first (convert to the number in the vocabulary list)
- Cats average 15 hours of sleep a day.  $\langle \text{EOS} \rangle$   
 $y^{<1>} \quad y^{<2>} \quad \dots \quad y^{<8>} \quad y^{<9>}$ 

end of sentence  
↓
- The egyptian Maam is a breed of cat.  $\langle \text{EOS} \rangle$   
 $\langle \text{UNK} \rangle$

RNN model:

Prediction of next word



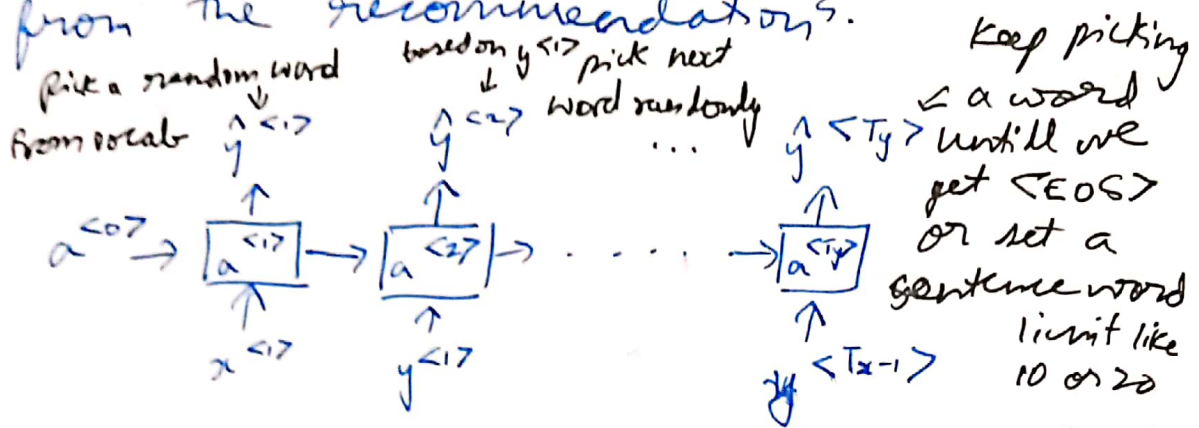
$$\mathcal{L}(\hat{y}^{<t>}, y^{<t>}) = - \sum_i y_i^{<t>} \log \hat{y}_i^{<t>} \quad \therefore \text{Given a sentence}$$

$$\mathcal{L} = \sum \mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

$$\begin{aligned} P(y^{<1>}, y^{<2>}, y^{<3>}) &= \\ P(y^{<1>}) P(y^{<2>} | y^{<1>}) &= \\ P(y^{<3>} | y^{<1>}, y^{<2>}) &= \end{aligned}$$

Sampling a sequence from a trained RNN:

Here we make it generate a sentence by making it randomly pick a word from the recommendations.



- Use np.random.choice to pick
- We can use a character level language model (with a vocabulary of characters) but it won't be good.

Vanishing gradients with RNNs:

If we have sentences like,

The cat, which already ate ..., was full  
The cats, which ..., were full

Since we saw in deeper NN, it's harder for a word at the beginning to affect a deeper word due to vanishing gradients. Usually only 3-4 words before have a strong influence on the current word using our current model. We will solve this.

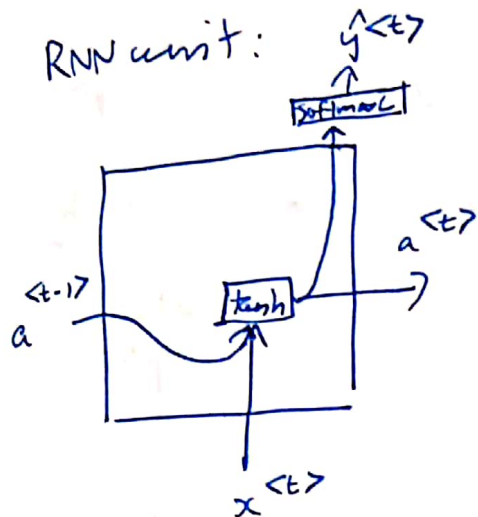
- If exploding gradient occurs, we will get NaN errors. For this we need to do gradient clipping (setting a threshold for gradients so they don't get very large values)



# Gated Recurrent Unit (GRU):

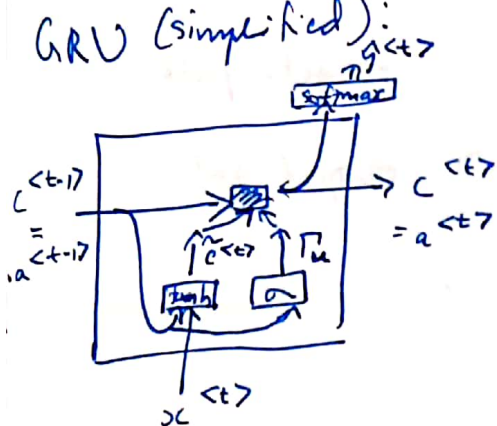
This helps with the vanishing gradient problem and makes it much better capturing long range connections.

RNN unit:



$$a^{<t>} = g(h[a^{<t-1>}, x^{<t>}] + b_a)$$

GRU (simplified):



$c$  = memory cell

$$\rightarrow c^{<t>} = a^{<t>}$$

$$\rightarrow \tilde{c}^{<t>} = \tanh(w_c[c^{<t-1>}, x^{<t>}] + b_c)$$

$$\rightarrow \Gamma_u = \sigma(w_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$\rightarrow c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$\Gamma_u = 1 \quad \Gamma_u = 0 \quad \Gamma_u = 0 \dots \dots \Gamma_u = 1$$

The cat which already ate was full  
 → The memory cell remembers that cat is singular till gate is activated again  
 → The gate decides when the  $c^{<t>}$  value changes, otherwise it won't affect it (since  $\Gamma_u = 0, \Gamma_u * \tilde{c}^{<t>} = 0$ )

Full GRU:

$$\rightarrow \tilde{c}^{<t>} = \tanh(w_c [\Gamma_r^{<t-1>} * c^{<t-1>}, x^{<t>}] + b_c) \quad \tilde{h}$$

$$\rightarrow \Gamma_u = \sigma(w_u [c^{<t-1>}, x^{<t>}] + b_u) \quad u$$

$$\rightarrow \Gamma_r = \sigma(w_r [c^{<t-1>}, x^{<t>}] + b_r) \quad r$$

$$\rightarrow c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>} \quad h$$

Long Short Term Memory (LSTM): < better option

$$\rightarrow \tilde{c}^{<t>} = \tanh(w_c [a^{<t-1>}, x^{<t>}] + b_c)$$

$$\rightarrow \Gamma_u = \sigma(w_u [a^{<t-1>}, x^{<t>}] + b_u) \quad \text{update gate}$$

$$\rightarrow \Gamma_f = \sigma(w_f [a^{<t-1>}, x^{<t>}] + b_f) \quad \text{forget gate}$$

$$\rightarrow \Gamma_o = \sigma(w_o [a^{<t-1>}, x^{<t>}] + b_o) \quad \text{output gate}$$

$$\rightarrow c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$\rightarrow a^{<t>} = \Gamma_o * \tanh(c^{<t>})$$

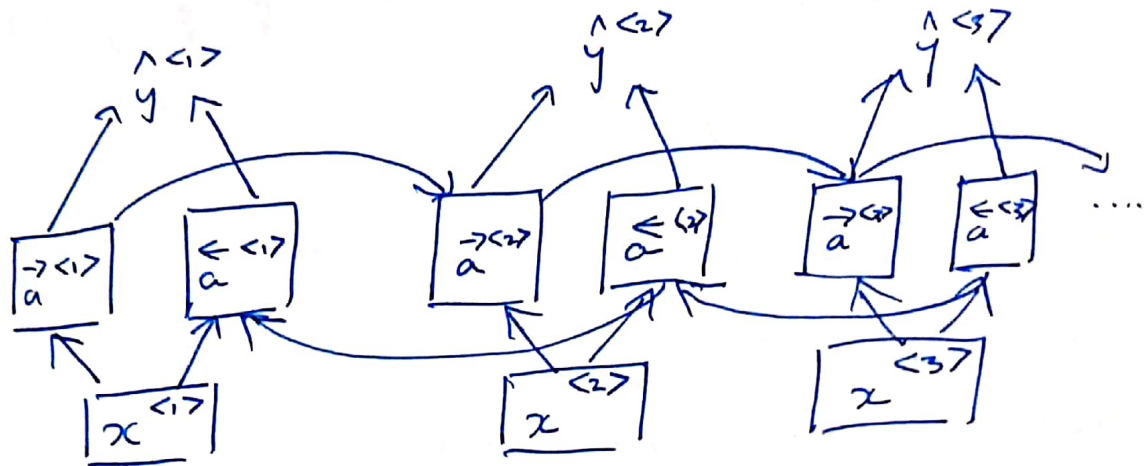
Check video for diagram of LSTM

$$y^{<t>}_{\text{pred}} = \text{softmax}(w_y a^{<t>} + b_y)$$



# Bidirectional RNN (BRNN)

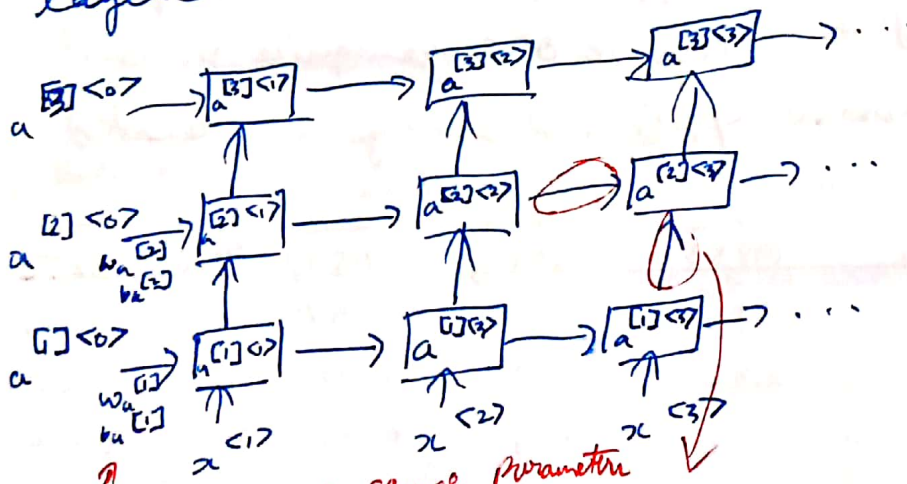
To address the issue of RNN ~~and~~ using information only from the previous words, we use bidirectional RNN.



$$\hat{y}^{<t>} = g(w_y [a^{<t> \rightarrow}, a^{<t> \leftarrow}] + b_y)$$

## Deep RNN example

You can make a deep RNN by stacking more layers horizontally.



for this entire row same parameter

Ex:  $a^{[2]<3>} = g(w_a^{[2]} [a^{[2]<2>}, a^{[2]<3>}] + b_a^{[2]})$