

Recurrent Neural Networks

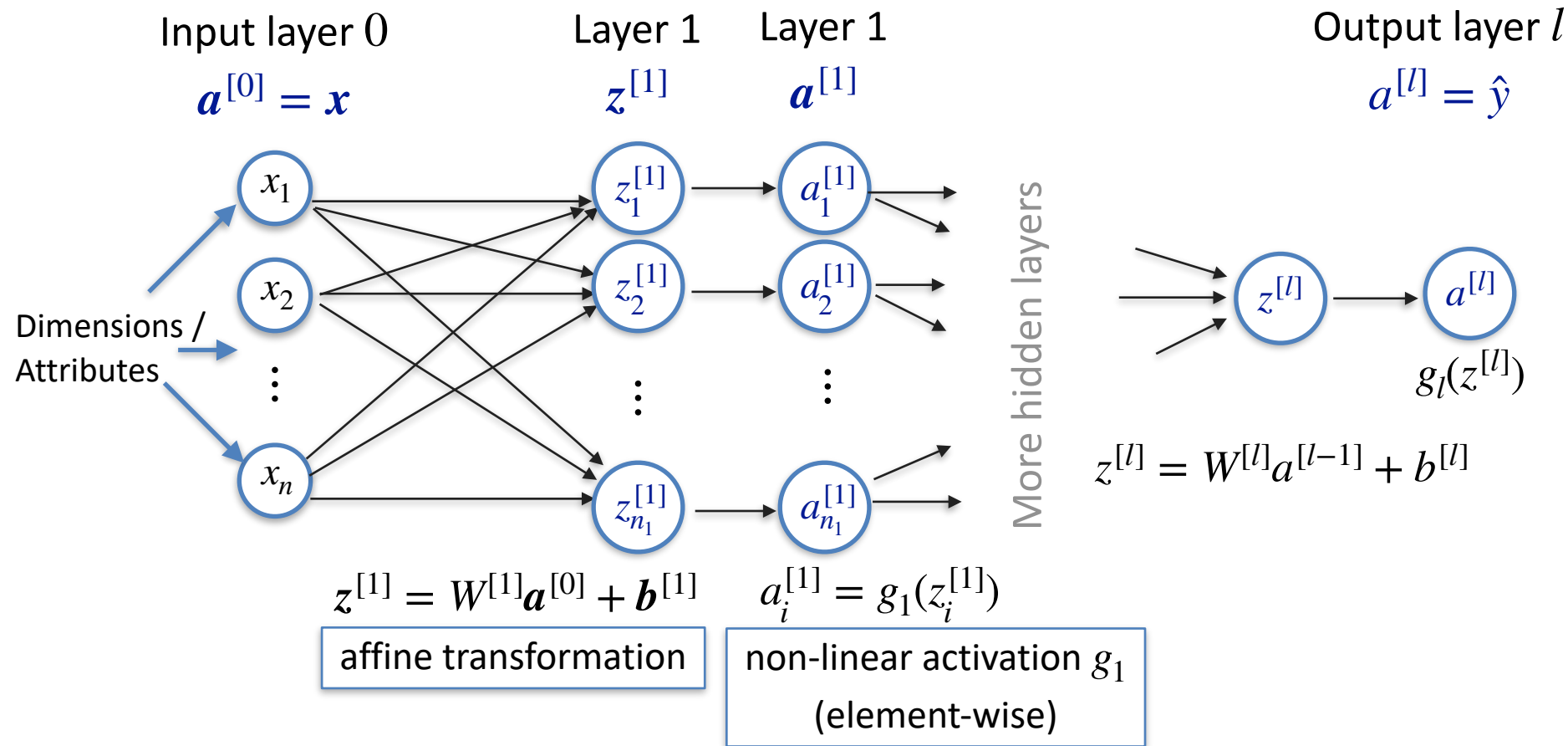
RNN, Types of RNN

Debapriyo Majumdar

debapriyo@isical.ac.in

Remember: deep feedforward networks

2



Each layer: an affine transformation, followed by a non-linear activation

- Data that are temporal or in the form of sequences
 - Text
 - Speech
 - Music
 - Video
 - DNA sequence, ...
- Each data point (and output, sometimes) is a sequence of input units
 - Inputs later in the sequence have dependency on previous inputs
- Early NLP techniques used the bag-of-words model
 - Ignored the sequence
 - Simulated somehow by set of sequences of words

Representing text in deep learning

4

21 46 34 32 42 7
sequential models are for text data

Assign an ID to each word

39 36 10 34 32 21 46
recurrent neural networks are for sequential models

Each sentence (input)
becomes a sequence of
integers

47 36 10 34 48 2 32 28
convolutional neural networks are generally used for images

28 31 30 23 3 49 16
images may be colored with three channels

16 5 1 47 36 12 4
ResNet is a convolutional neural network architecture

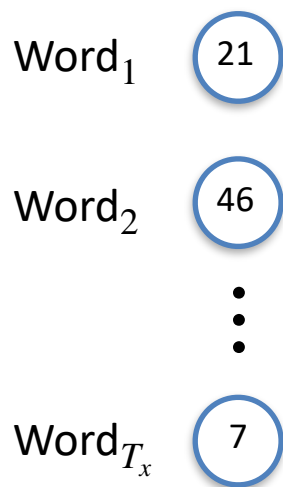
21 19 6 39 36 10
sequential inputs need recurrent neural networks

Representing text for deep learning

5

Consider a sentence of length T_x

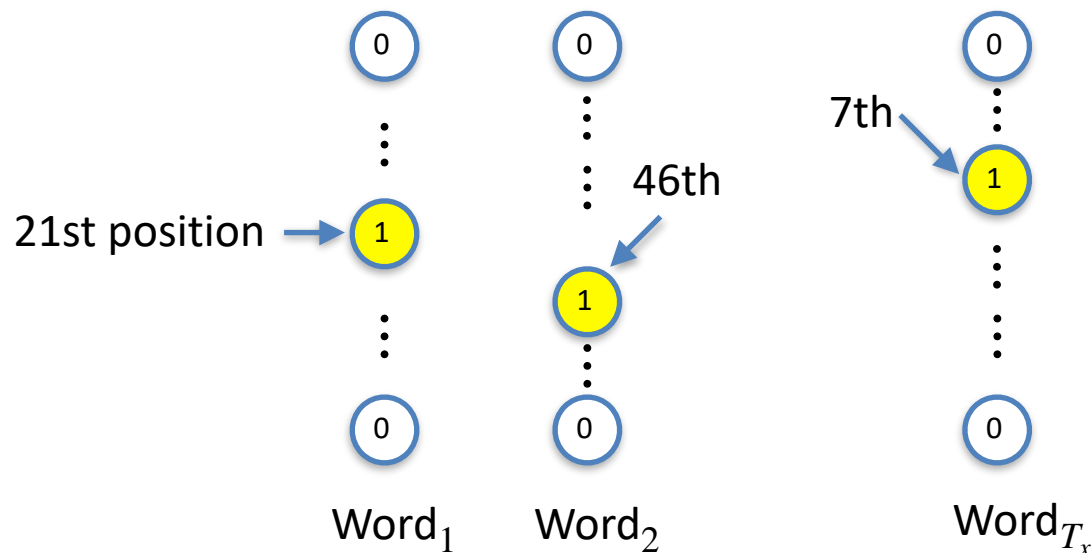
The integer sequence
as input vector



Problem: the numbers have an order that does not represent the words.

One-hot encoding

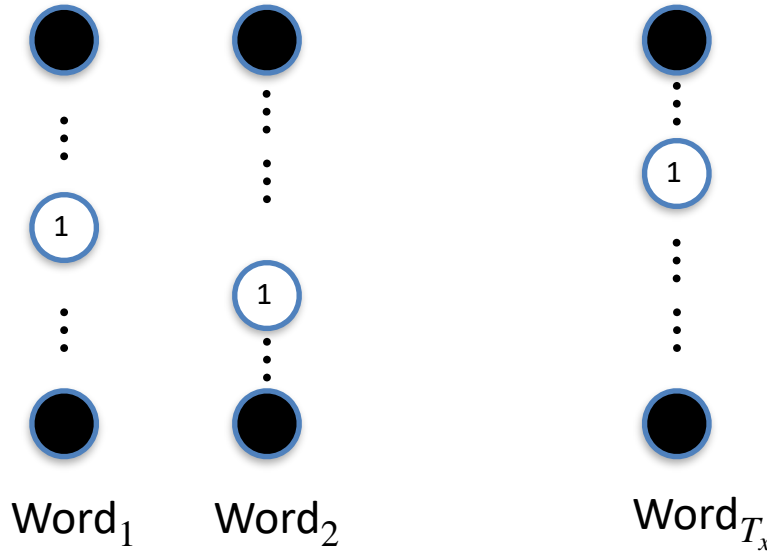
Each integer (word) as a vector



Each word is a vector with all zeros except one
Orthogonal to each other

Word embeddings

6

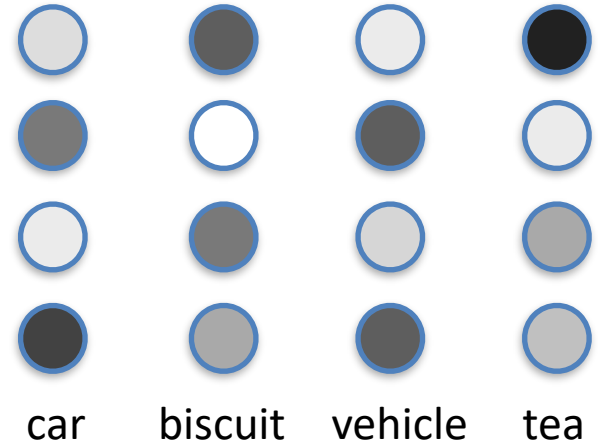


One-hot encoding

Each word is a vector with all zeros except one
Orthogonal to each other

But words have meanings
Some words are similar to each other

White = 1; Black = 0, Grey $\in (0,1)$.



Better idea

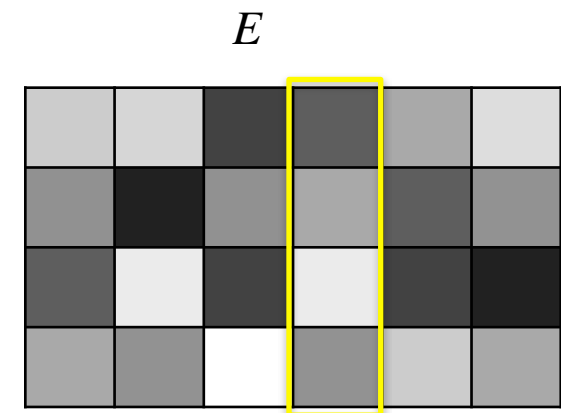
More **compact** vector representation
retaining **semantics** of words
Similar words \implies similar vectors

Learn word embeddings

The embedding matrix

7

White = 1; Black = 0, Grey $\in (0,1)$.



Embedding matrix for d
dimensional embedding

$$e_i = E o_i$$

o_i

0

0

0

1

0

0

$N \times 1$

One-hot vector for the i -th
word in the vocabulary

=

e_i

$d \times 1$

Embedding for the i -th
word in the vocabulary

Learning word embeddings

8

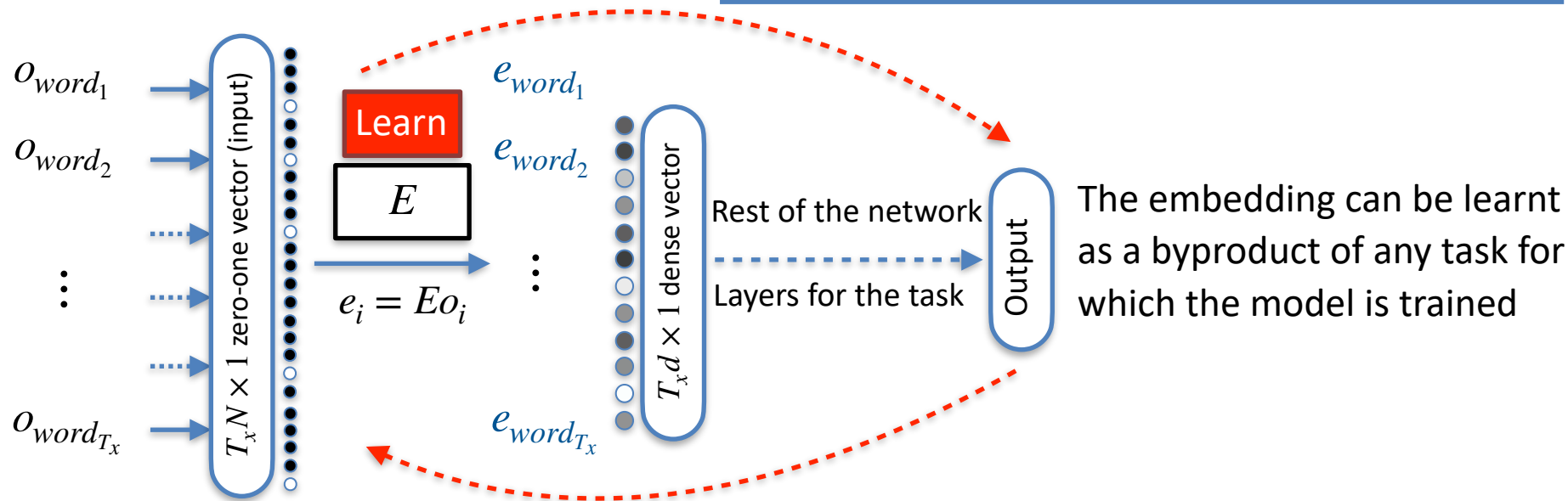
White = 1; Black = 0, Grey $\in (0,1)$.

Add an embedding layer in the beginning

```
inputs = keras.Input(shape=(None,), dtype="int32")
```

Embed each integer in a d -dimensional vector

```
x = layers.Embedding(N, d)(inputs)
```



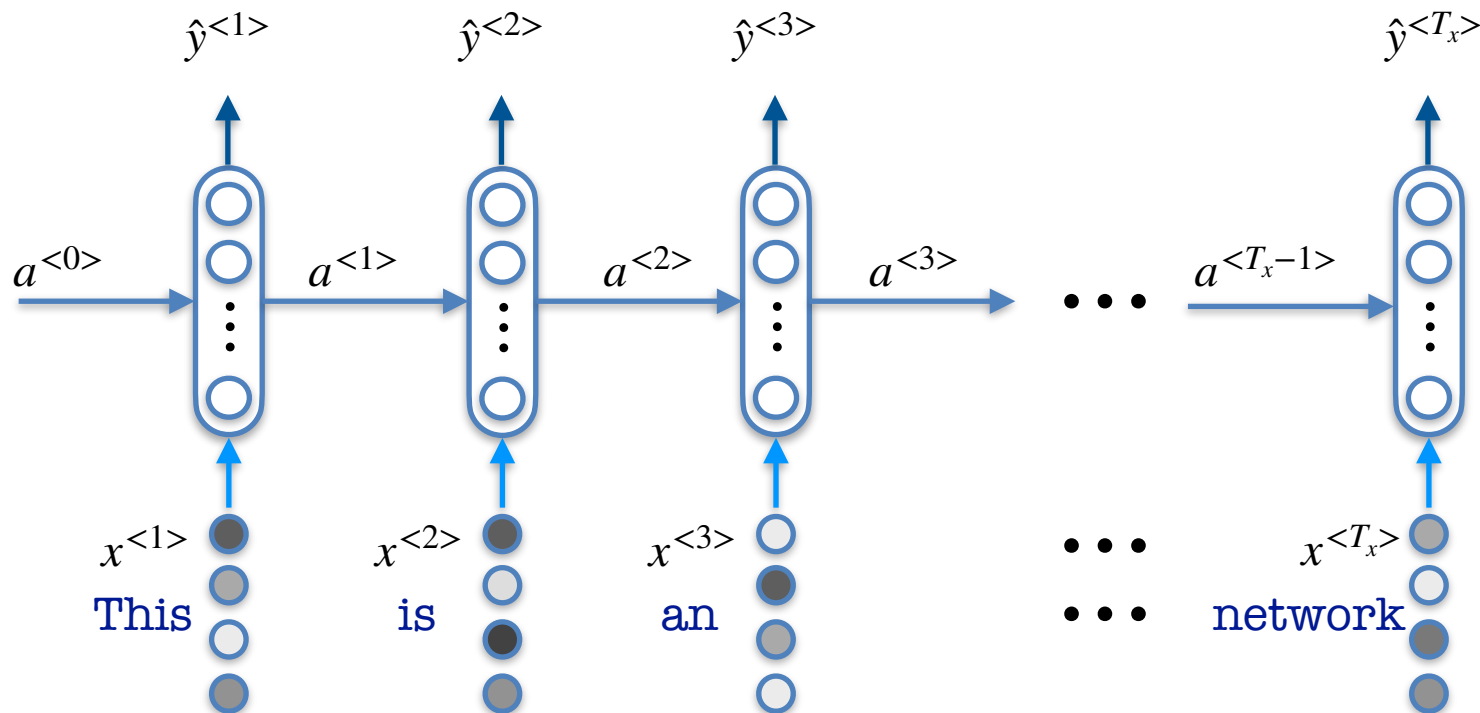
Input as one-hot vectors

We will learn more about word embeddings later.

The Recurrent Neural Network (RNN)

9

Output \hat{y} = the sequence $\hat{y}^{<1>}, \dots, \hat{y}^{<T_x>}$



Input x : This is an example of a recurrent neural network.

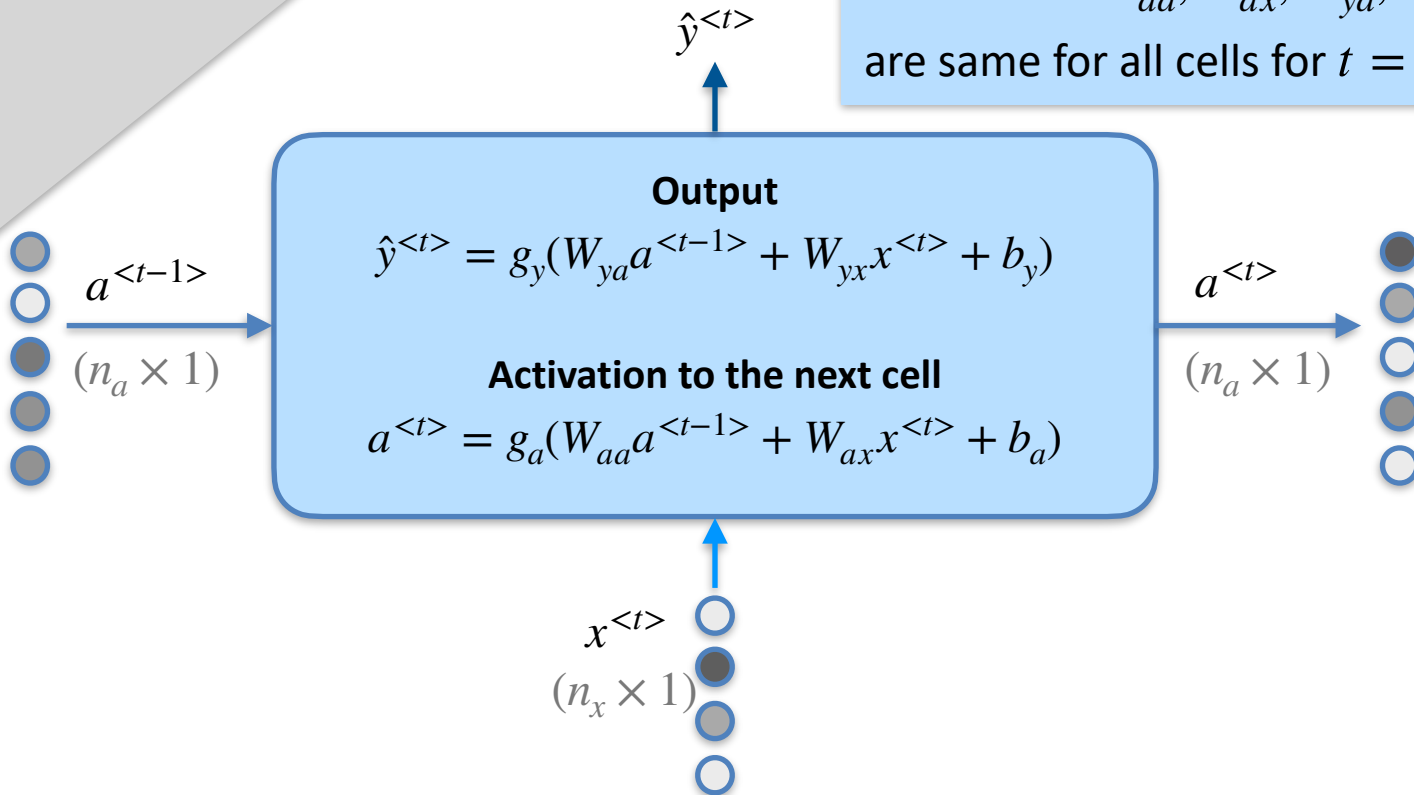
In this example, length of \hat{y} = length of $x = T_x$

But it can be different, as well as the architecture also may be somewhat different

An RNN Cell

10

What happens inside an RNN cell?



- Input $x^{<t>}$: some vector representation (say, one-hot) of the input element (e.g. a word)

An RNN Cell: Details of the Operations

11

$$a^{<t>} = g_a(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

Diagram illustrating the dimensions of the variables in the equation:

- $n_a \times 1$ (input vector dimension) points to $a^{<t>}$
- $n_a \times 1$ (input vector dimension) points to $a^{<t-1>}$
- $n_x \times 1$ (input vector dimension) points to $x^{<t>}$
- $n_a \times n_a$ (weight matrix dimension) points to W_{aa}
- $n_a \times n_x$ (weight matrix dimension) points to W_{ax}
- $n_a \times 1$ (bias vector dimension) points to b_a

 \equiv

$$a^{<t>} = g_a \left(\underbrace{\begin{bmatrix} \overbrace{W_{aa} \mid W_{ax}}^{W_a} \end{bmatrix}}_{n_a \times (n_a + n_x)} \underbrace{\begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix}}_{(n_a + n_x) \times 1} + b_a \right)$$

We can write

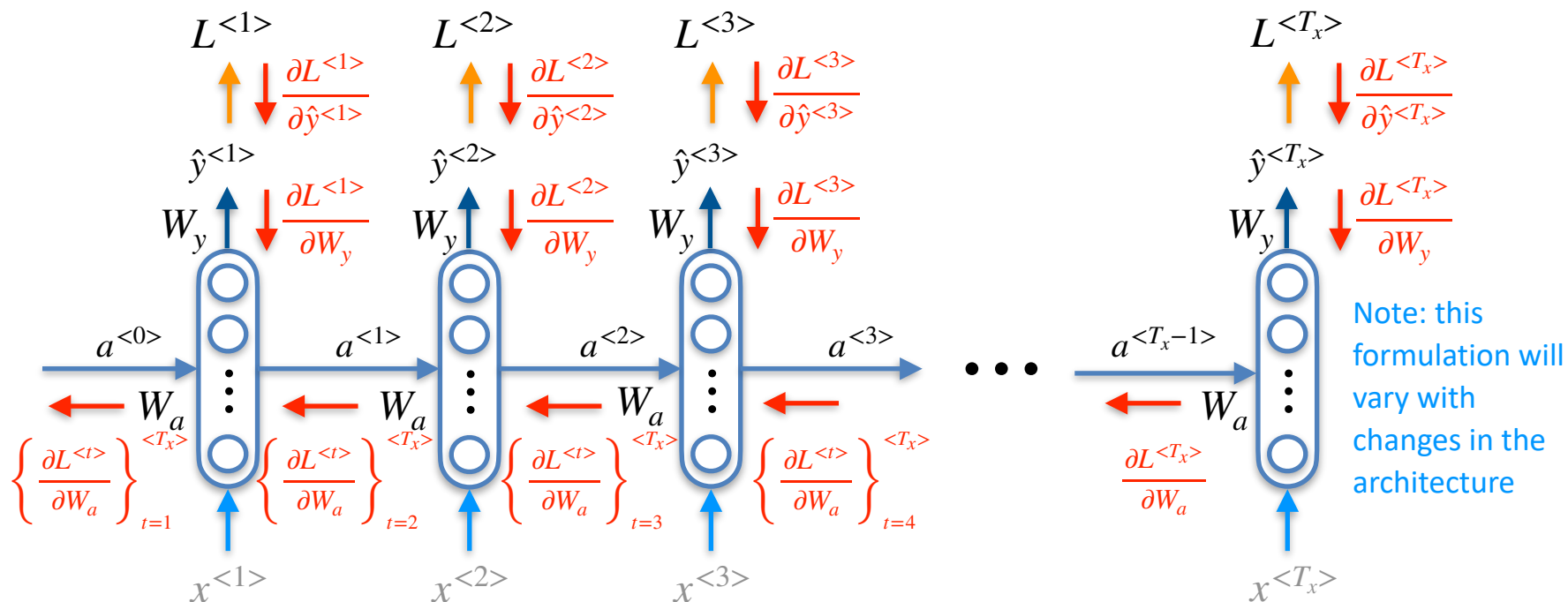
$$a^{<t>} = g_a(W_a[a^{<t-1>}, x^{<t>}] + b_a)$$

and, similarly

$$\hat{y}^{<t>} = g_y(W_y[a^{<t-1>}, x^{<t>}] + b_y)$$

Backpropagation through time

12



$$\text{Total loss: } L = \sum_{t=1}^{T_x} L^{<t>} = \sum_{t=1}^{T_x} L(\hat{y}^{<t>}, y^{<t>})$$

$$\frac{\partial L}{\partial W_y} = \sum_{t=1}^{T_x} \frac{\partial L^{<t>}}{\partial W_y}$$

$$\frac{\partial L}{\partial W_a} = \sum_{t=1}^{T_x} \frac{\partial L^{<t>}}{\partial W_a}$$

Similarly for b_y and b_a

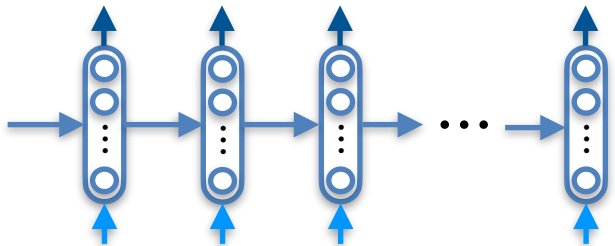
Input x is a single training example. Update W_y, b_y, W_a, b_a happens after the complete backpropagation (possibly for all examples of a mini-batch)

Different RNN Architectures

13

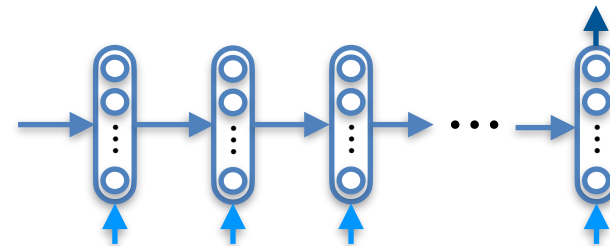
Many to Many

Named entity recognition, POS tagging, ...



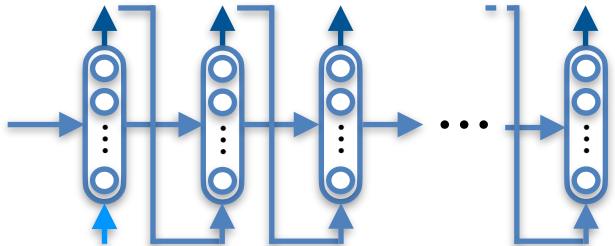
Many to One

Sentiment classification, ...



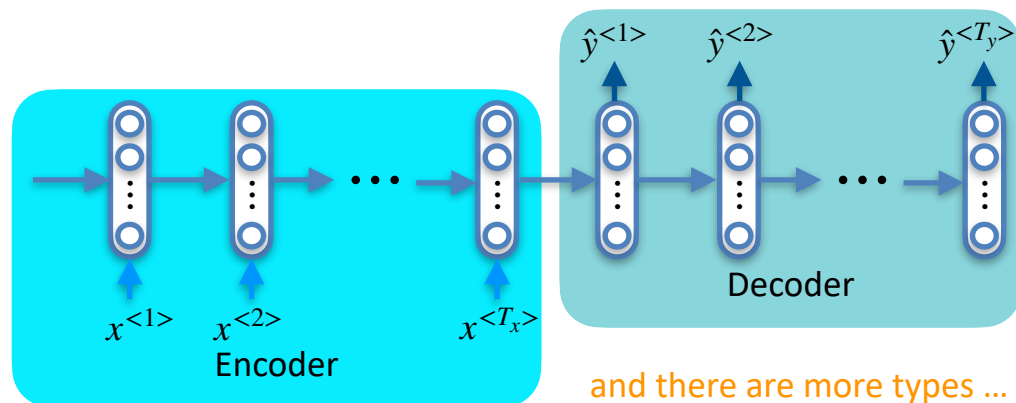
One to Many

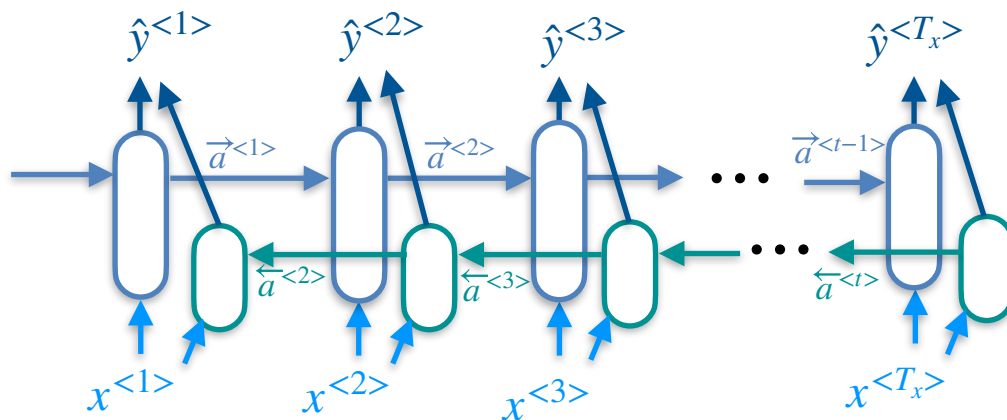
Music generation, text generation, ...



Many to Many (encoder - decoder)

Machine Translation, ...



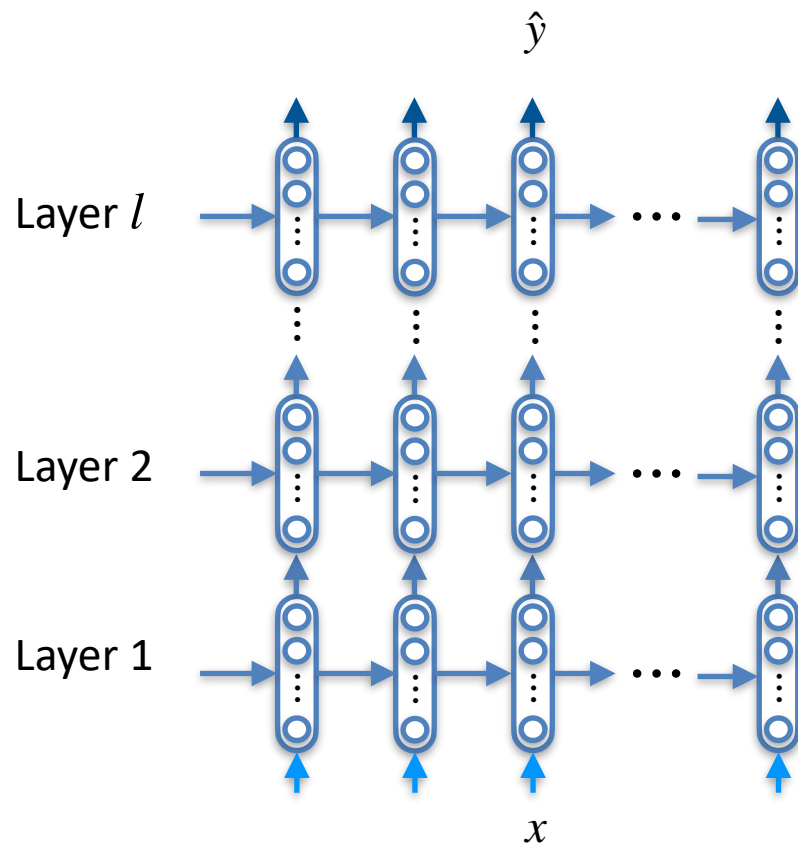


An acyclic graph

- Sometimes, earlier cells may require information from later cells to make sense
- Example: “Teddy bears are lovely gifts.” and “Teddy Roosevelt was a president of the United States of America.”
- Each output $\hat{y}^{<t>}$ can be obtained by both forward and backward cells
$$\hat{y}^{<t>} = g(W_y[\vec{a}^{<t>}, \overleftarrow{a}^{<t>}, x^{<t>}] + b_y)$$
- Usually the forward-RNN and backward-RNN are trained together by stacking the hidden state activations $\vec{a}^{<t>}$ and $\overleftarrow{a}^{<t>}$
- Bidirectional RNNs are widely used in practice

Deep RNN with multiple layers

15



- Multiple layers of RNN can be stacked to make a deep RNN
- Also called *stacked RNNs*
- Intuition: lower-level layers capture lower-level features, higher-level layers capture higher level features
- In practice, number of layers for multi-layer RNNs (2-6) are not as much as ConvNets (can be 150 - 200)
- Transformer based networks (e.g. BERT) can be up to 24 layers

- Andrew Ng's lectures on *Sequence Models*: www.coursera.org/learn/nlp-sequence-models
- Chris Manning, Abigail See and other TAs. *Natural Language Processing with Deep Learning*. Stanford University Course (CS224n), Winter 2019. web.stanford.edu/class/archive/cs/cs224n/cs224n.1194/
- Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016. www.deeplearningbook.org