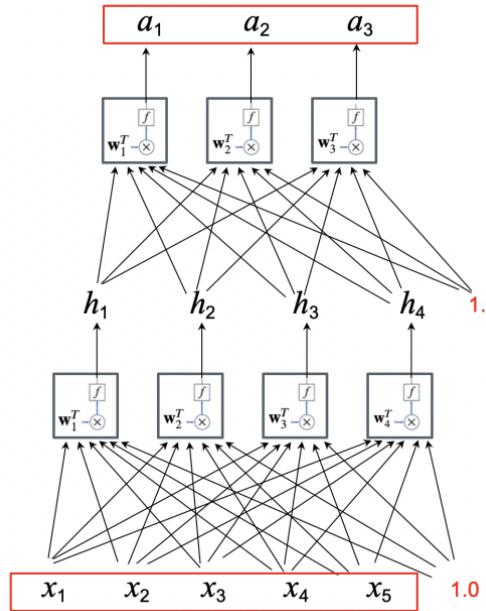


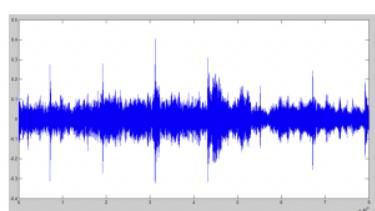
## Neural Networks for Sequence Data: RNNs, GRUs, LSTMs; Deep RNNs

## 1. Sequence Data: A problem for FFNNs

- a. A Feed-Forward Neural Network learns a function from vectors to vectors:



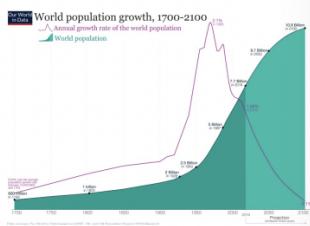
- b. You do not want it to learn the training set (overfitting)
- c. Note: The Order of the inputs does not matter at all!
- d. Not everything cannot be turned into vector to vector function (order of the input and output does not matter → completely symmetric, which is not an advantage when order matters)
- e. However, many kinds of data are inherently sequential, often as a time series
- f. Words in a sentence:
  - i. The matters order
  - ii. Order the matters
  - iii. Matters order the
  - iv. The order matters → clearly this is the best sentence with appropriate order
- g. Audio:



h. Stock prices:



i. Population:



j. We have already tried one way to deal with sequences, using the Markov Assumption: N-grams = short subsequences of words..... Can we do better?

## 2. Sequence Data: A problem for FFNNs in NLP

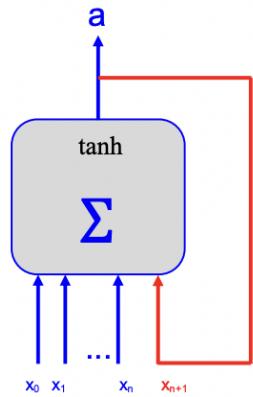
- a. So far, our representations for words and documents have either
  - i. Completely ignored sequencing (BOW, TF, TFIDF, mean of embeddings, cosine similarity)
  - ii. Accounted for very short sequences (N-grams)
- b. You care about what words are around it but we do not care about the order
- c. But (modern) natural languages are strongly sequential, and exhibit long-range dependencies:
  - i. The students in CS 505 love NLP. vs. The student in CS 505 loves NLP. (verb tense)
  - ii. The students in Professor Snyder's CS 505 class love NLP.
  - iii. The students in Professor Snyder's CS 505 (Natural Language Processing) class — at least when an assignment is not due that night — love NLP. (hyphen – open and close)
- d. Or just look at first sentence above! – parenthesis, verb, commas, either \_\_\_\_\_ or \_\_\_\_\_
  - i. All verbs are consistently in past tense. ...have ... ignored .... Accounted ....
  - ii. Consistent use of commas in list ... , ... , ... , ....
  - iii. Matching parentheses: ... ( ..... ) ...
  - iv. Miscellaneous: ... either .... or for very short sequences
- e. Clearly we need better methods to process long sequences
- f. And this is not just a problem with individual sentences
- g. With an extended discourse, many different things have to be remembered:
  - i. I grew up in France. I learn to cycle when I was very young but only learned to swim as an adult. I also love to cook and bake. I can make a

mean cake. I speak several languages but because of where I grew up, I am most fluent in \_\_\_\_\_.

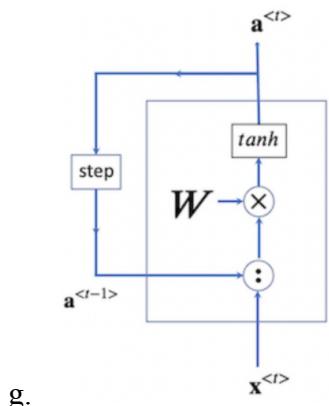
- ii. The answer should be French according to the text

### 3. Recurrent Neural Networks

- a. The solution is recursion (Though we call it “recurrence” in NNs)
- b. A basic recurrent neuron is the same with one important change: the outputs are fed back into the input layer:



- c. Another, less significant change, is that tanh is generally used instead of sigmoid or relu
- d. The most significant difference is that it goes tanh, activation comes out, feeds the activation in the next cycle (we cycle the information using a loop)
- e. Here, we have a sequence of vectors
- f. In a recurrent layer, all the output activations are fed back and concatenated with the inputs



- g. Now the input is a sequence of data vectors

$$\mathbf{x}^{<1>} , \mathbf{x}^{<2>} , \dots , \mathbf{x}^{<T>} ,$$

processed in a loop for time steps  $t = 1, 2, \dots, T$

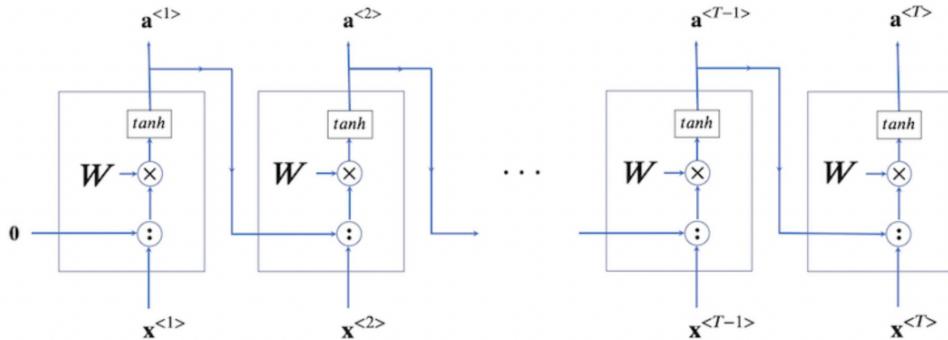
$$\begin{aligned}
 \mathbf{a}^{<0>} &= \mathbf{0} \\
 \textbf{for } t &= 1, 2, \dots, T: \\
 \mathbf{a}^{<t>} &= \tanh(W(\mathbf{a}^{<t-1>} : \mathbf{x}^{<t>}))
 \end{aligned}$$

concatenation

Weights for recurrent activations      Weights for data inputs

$$W = \begin{bmatrix} \mathbf{w}'_1 & \mathbf{w}_1 \\ \mathbf{w}'_2 & \mathbf{w}_2 \\ \vdots & \vdots \\ \mathbf{w}'_n & \mathbf{w}_n \end{bmatrix} = \begin{bmatrix} w'_{1,1} & w'_{1,2} & \cdots & w'_{1,n} & w_{1,1} & w_{1,2} & \cdots & w_{1,k} \\ w'_{2,1} & w'_{2,2} & \cdots & w'_{2,n} & w_{2,1} & w_{2,2} & \cdots & w_{2,k} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ w'_{n,1} & w'_{n,2} & \cdots & w'_{n,n} & w_{n,1} & w_{n,2} & \cdots & w_{n,k} \end{bmatrix}$$

- i. You have the weights coming from recurrent activations and data inputs
- j. You will commonly see diagrams unrolled through time

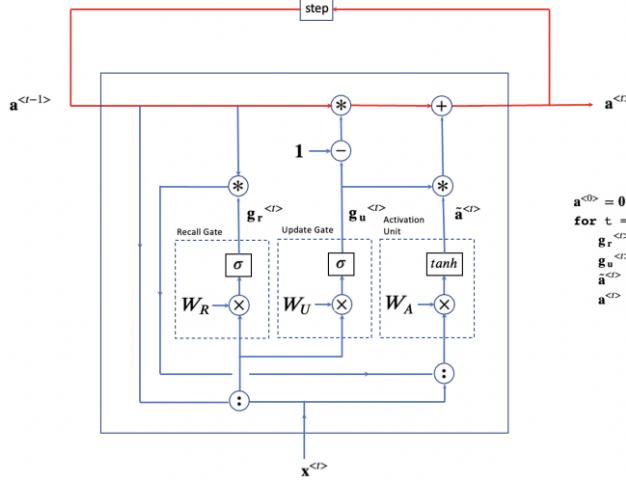


- k. It is the same  $W$  (one set of weights that is trained through the whole process)
- l. It is trying to remember what has happened in the previous parts of the stage
- m. But make sure you understand that this is one layer, with one set of weights  $W$

$$\begin{aligned}
 \mathbf{a}^{<0>} &= \mathbf{0} \\
 \textbf{for } t &= 1, 2, \dots, T: \\
 \mathbf{a}^{<t>} &= \tanh(W(\mathbf{a}^{<t-1>} : \mathbf{x}^{<t>}))
 \end{aligned}$$

#### 4. Recurrent Neural Networks: GRUs

- a. A Gated Recurrence Unit (GRU) adds a recurrent activation path which acts as a memory. Two “sub-neurons” have been added
- b. Recall Gate: learns how to read from memory
- c. Update Gate: learns how to write to memory



```

 $a^{<t>} = 0$ 
for t = 1, 2, ..., T:
     $g_r^{<t>} = \sigma(W_R(a^{<t-1>}; x^{<t>}))$ 
     $g_u^{<t>} = \sigma(W_U(a^{<t-1>}; x^{<t>}))$ 
     $\tilde{a}^{<t>} = \tanh(W_A((g_r^{<t>} * a^{<t-1>}); x^{<t>}))$ 
     $a^{<t>} = g_u^{<t>} * \tilde{a}^{<t>} + (1 - g_u^{<t>}) * a^{<t-1>}$ 

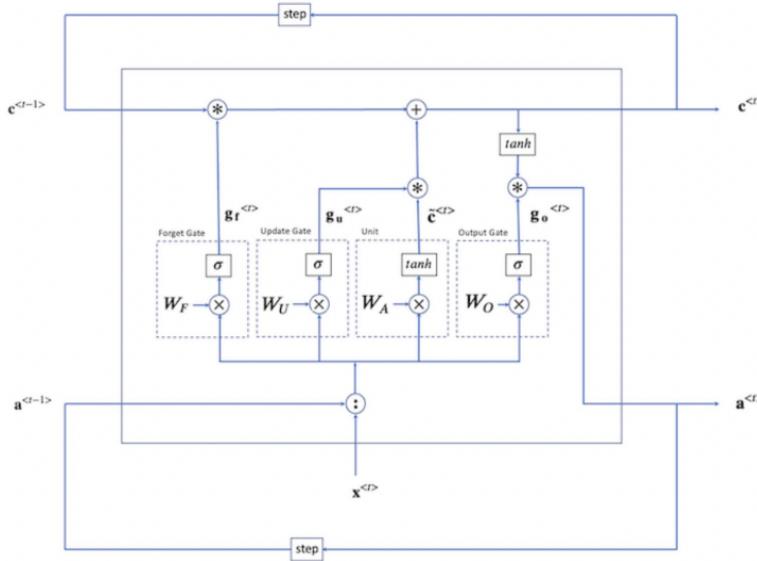
```

d.

i.  $W_U$  – Weights from the memory

## 5. Recurrent Neural Network: LSTM

- a. A Long Short-Term Memory (LSTM) is another, earlier design, still very much used: The memory cycle is separate from the activation, and an output gate determines how memory is used to form the activation



```

 $0> = 0$ 
 $0> = 0$ 
for t = 1, 2, ..., T:
     $g_f^{<t>} = \sigma(W_F(a^{<t-1>}; x^{<t>}))$ 
     $g_o^{<t>} = \sigma(W_O(a^{<t-1>}; x^{<t>}))$ 
     $g_u^{<t>} = \sigma(W_U(a^{<t-1>}; x^{<t>}))$ 
     $\tilde{c}^{<t>} = \tanh(W_A(a^{<t-1>}; x^{<t>}))$ 
     $c^{<t>} = g_u^{<t>} * \tilde{c}^{<t>} + g_f^{<t>} * c^{<t-1>}$ 
     $a^{<t>} = g_o^{<t>} * \tanh(c^{<t>})$ 

```

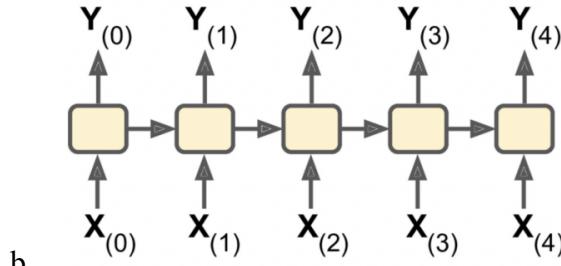
b.

## 6. RNNs: Which one is best for NLP?

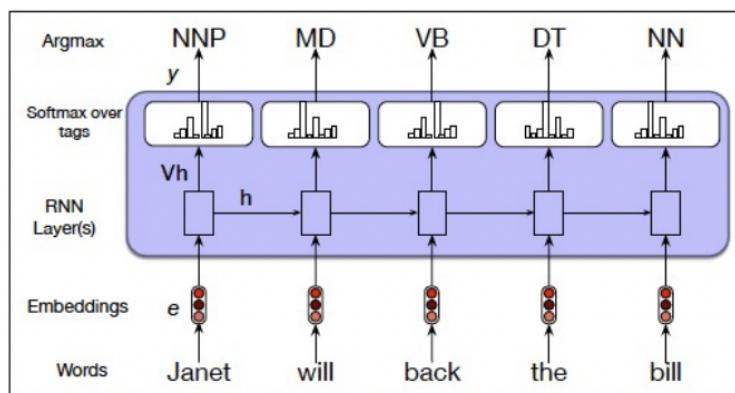
- a. Historically, the LSTM design was first, and the GRU was designed as a simplified version of the LSTM. The very basic RNN is also used.
- b. However, the LSTM has 4 weight matrices, compared with 3 for the GRU, 2 for the Simple GRU, and one for the FFNN.
- c. Adding more gates in general improves performance, but has an obvious impact on the complexity of training.
- d. For very large networks with large data sets, the GRU is generally used, and the LSTM seems to work better for smaller projects such as you would do in an academic course.

## 7. RNN Network Architectures

- a. There are many ways to configure an RNN. The simplest is a sequence-to-sequence

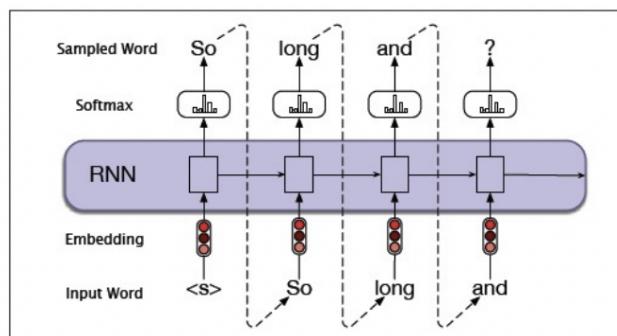


- b.
- c. NOTE: From now on, we'll show every RNN unrolled through time, though you should always remember that there is a for loop controlling the whole process.)
- d. Example 1: Part of speech tagging



**Figure 9.7** Part-of-speech tagging as sequence labeling with a simple RNN. Pre-trained word embeddings serve as inputs and a softmax layer provides a probability distribution over the part-of-speech tags as output at each time step.

- i.
- ii. Whatever goes out to the top is passed to the next part (remember what has happened before)
- e. Example 2 (vector to sequence): A sentence generator using a trained RNN can generate sentences by picking the most likely next word in each step, until it generates the end-of-sentence token </s>

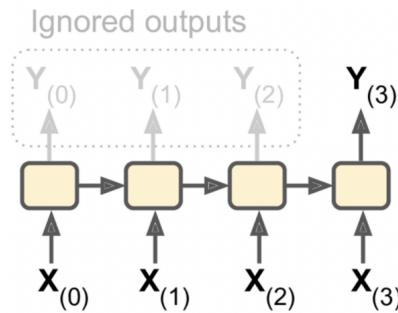


**Figure 9.9** Autoregressive generation with an RNN-based neural language model.

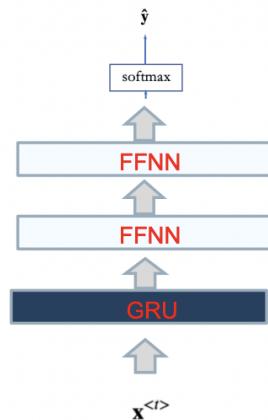
- i.
- ii. Something else is passed too with the word (the RNN)

## 8. RNN Network Architectures

- a. Another configuration is sequence-to-vector:

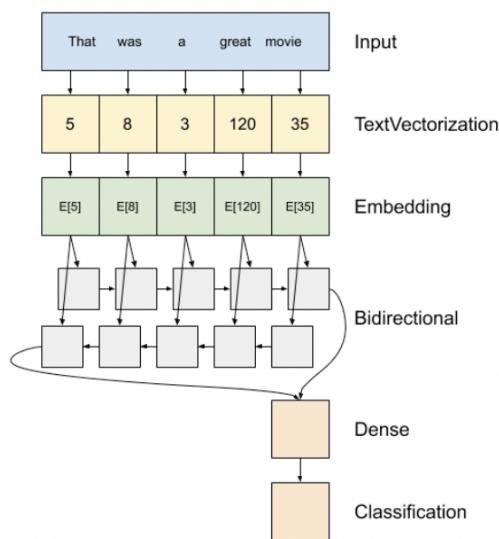


- b. With sequence-to-vector layers, it is very common to add FF layers downstream, especially for classification:

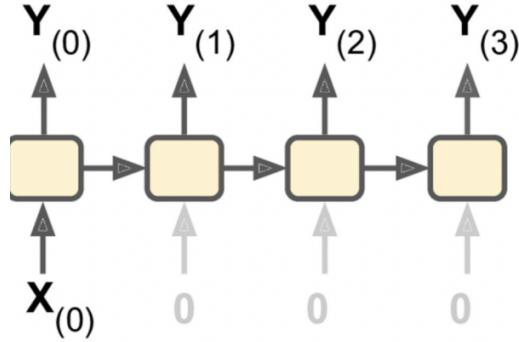


- c. Example: Text Classification

### Example: Text Classification



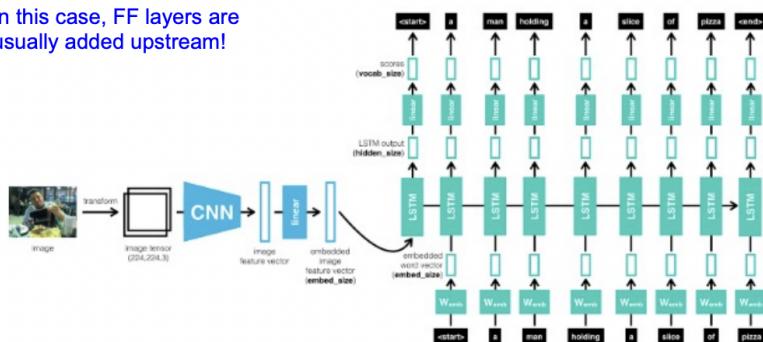
d. A third possibility is vector-to-sequence:



e. Example: Image Captioning

- i. In this case, FF layers are usually added upstream!

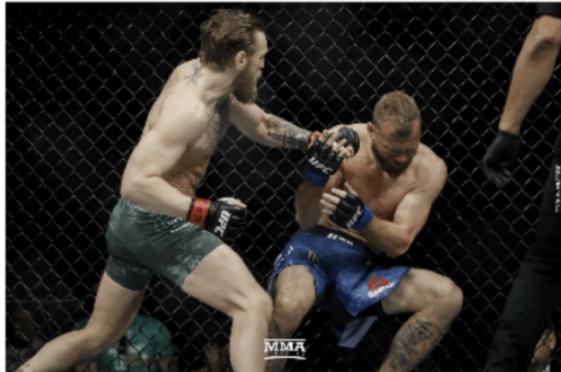
In this case, FF layers are usually added upstream!



ii.

f. Example: Image Captioning (hopefully accurate)

a man in a white shirt playing tennis



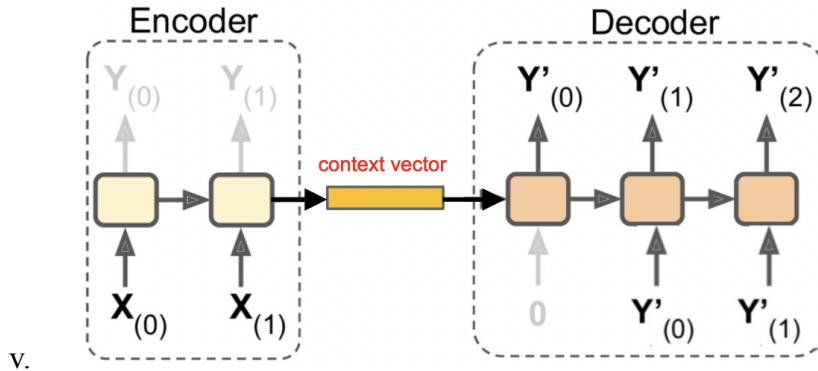
i.

- ii. There are two men
- iii. There is no white shirt
- iv. Not playing tennis

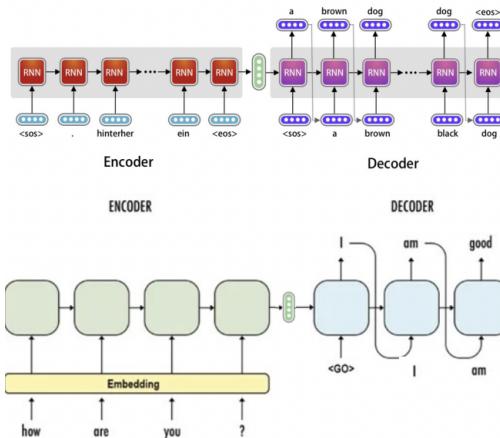
g. Finally, the SOTA

- i. An encoder (sequence-to-vector) and
- ii. A decoder (vector-to-sequence)
- iii. And passes a context vector between them

- iv. The bar gets passed to the decoder (sequence to vector and vector to sequence)

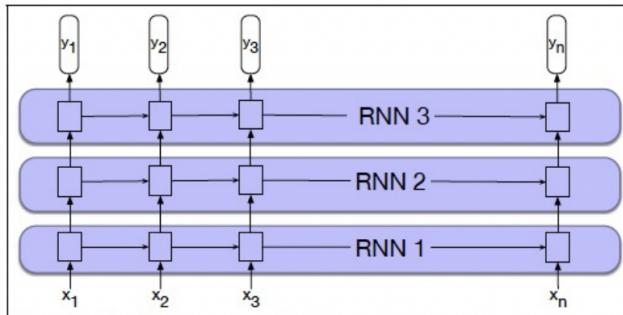


- v. h. Examples: Machine Translation, Conversation Agents:



- i. ii. Always keep in mind that these are unrolled in time  
 iii. Note that the input to the decoder at each time step is the previous output plus additional activations

- i. One more idea: Stacked Recurrent Layers (multiple recurrent layers)



**Figure 9.10** Stacked recurrent networks. The output of a lower level serves as the input to higher levels with the output of the last network serving as the final output.

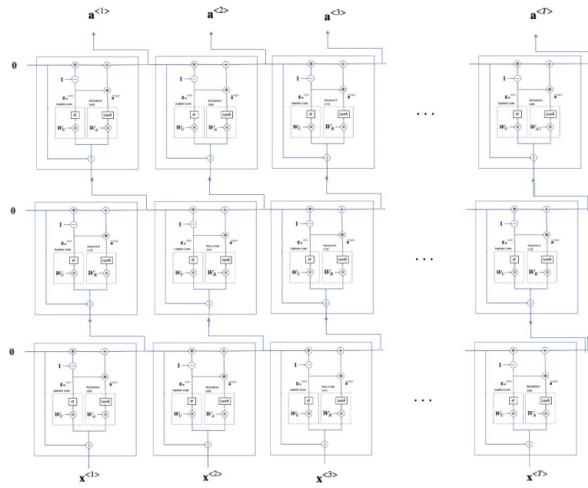
- j. k. Deep Networks

- i. Generally, networks for sequence data such as text have recurrent layers processing the sequence, and feed forward layers interpreting and producing output such as a classification



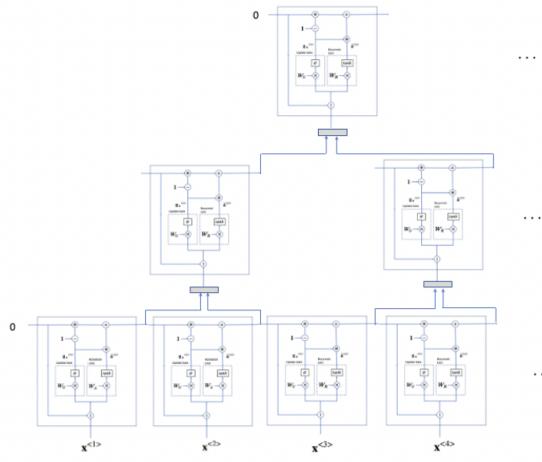
## 9. Recurrent Neural Networks: RNNs, GRUs, LSTMs

- a. Unrolling a deep RNN network reveals a very complicated design

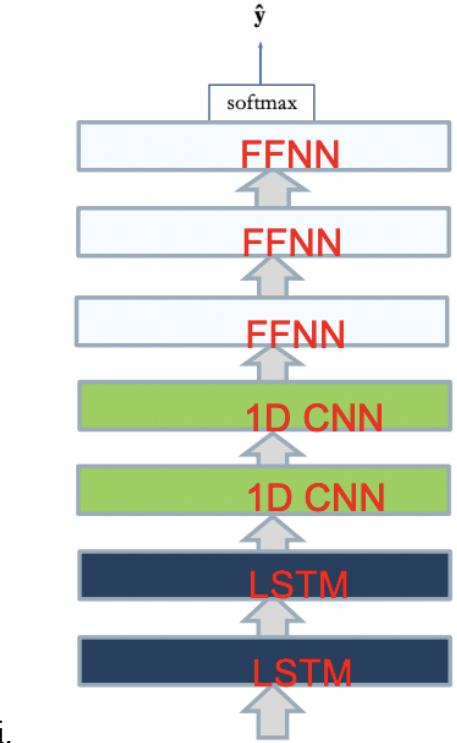


b.

- c. Many different designs have been proposed, with advantages and disadvantages
- d. Idea 1: Tree-structured network which combines lower levels using some aggregating function (weighted) sum, perhaps controlled by a gate

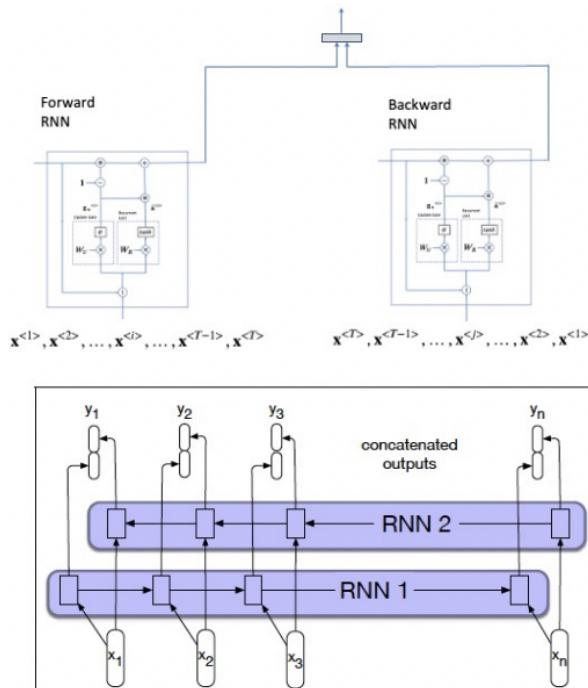


- e. Idea 2: Apply 1D Convolutions to the RNN layers



i.

- f. Idea 3: Bidirectional RNN (BRNN): Combine result of running two RNNs on forward and reverse sequence simultaneously. Results are fed to the next layer, usually by concatenation.

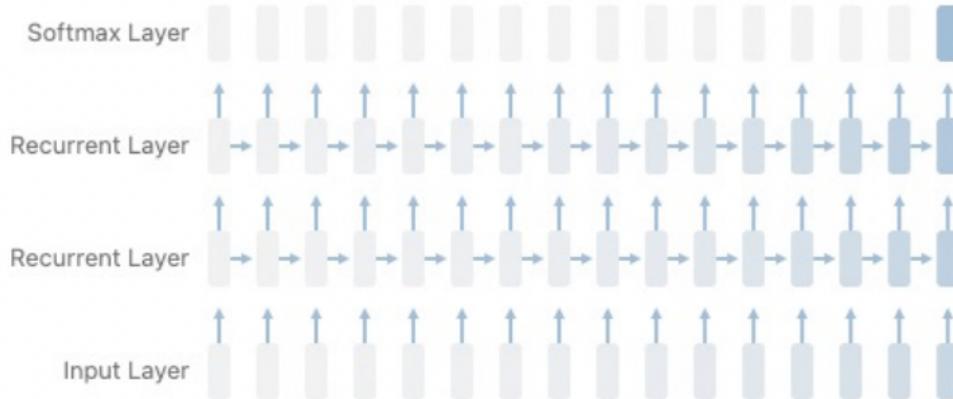


**Figure 9.11** A bidirectional RNN. Separate models are trained in the forward and backward directions, with the output of each model at each time point concatenated to represent the bidirectional state at that time point.

i.

## 10. RNN Network Architectures

- a. Unfortunately, all this complexity comes with a cost!
- b. Deep recurrent networks have a large number of parameters and take a long time to train.
- c. Plus, we have the Vanishing Gradients Problem: unrolling through time makes the network very large and preserving information (through weights) over long distances is a problem:



**Vanishing Gradient:** where the contribution from the earlier steps becomes insignificant in the gradient for the vanilla RNN unit.

- d.
- e. There is a symmetric problem, called the Exploding Gradients Problem!
- f. One solution is encoder decoder
- g. Another solution is normalize between layers

## 11. Fun with Character-Level RNN Text Generation

- a. From “The Unreasonable Effectiveness of Recurrent Neural Networks

### The evolution of samples while training

First, it's fun to look at how the sampled text evolves while the model trains. For example, I trained an LSTM of Leo Tolstoy's War and Peace and then generated samples every 100 iterations of training. At iteration 100 the model samples random jumbles:

```
tyntd-iafhatawiaoahrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e  
plia tkldrgd t o idoe ns,smtt h ne etie h,hregtrs nigtike,aoaenns lng
```

However, notice that at least it is starting to get an idea about words separated by spaces. Except sometimes it inserts two spaces. It also doesn't know that comma is almost always followed by a space. At 300 iterations we see that the model starts to get an idea about quotes and periods:

```
"Tmont thithey" fomesscerliund  
Keushey. Thom here  
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwy fil on aseterlome  
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."
```

- b.

The words are now also separated with spaces and the model starts to get the idea about periods at the end of a sentence. At iteration 500:

```
we counter. He stutn co des. His stanted out one ofler that concossions and was  
to gearang reay Jotrets and with fre colt otf paitt thin wall. Which das stimn
```

the model has now learned to spell the shortest and most common words such as "we", "He", "His", "Which", "and", etc. At iteration 700 we're starting to see more and more English-like text emerge:

```
Aftair fall unsuch that the hall for Prince Velzonski's that me of  
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort  
how, and Gogition is so overelical and ofter.
```

C.

At iteration 1200 we're now seeing use of quotations and question/exclamation marks. Longer words have now been learned as well:

```
"Kite vouch!" he repeated by her  
door. "But I would be done and quarts, feeling, then, son is people...."
```

Until at last we start to get properly spelled words, quotations, names, and so on by about iteration 2000:

```
"Why do what that day," replied Natasha, and wishing to himself the fact the  
princess, Princess Mary was easier, fed in had oftened him.  
Pierre aking his soul came to the packs and drove up his father-in-law women.
```

The picture that emerges is that the model first discovers the general word-space structure and then rapidly starts to learn the words; First starting with the short words and then eventually the longer ones. Topics and themes that span multiple words (and in general longer-term dependencies) start to emerge only much later.

d.