

Note the faster convergence and the  
crisper output. The convergence is  
noisier due to the large step sizes,  
but it's controlled thanks to the  
adaptation of  $\gamma$  and the averaging  
inside the minibatch.

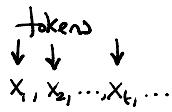
(set batchsize=20 to see it get  
smoother but a bit slower.)

# L25

Last time: RNNs

### ① Language modeling.

Language as a sequence of symbols



Objective: Maximize  $P(x_1, x_2, \dots, x_t, \dots)$  of data  
However, just like images, this is a high-dimensional space.

Any specific sequence has exponentially small probability!

Instead: Write  $P(x_1, x_2, \dots, x_t, \dots) \approx P(x_1)P(x_2|x_1)\dots P(x_t|x_{t-1}, \dots)$

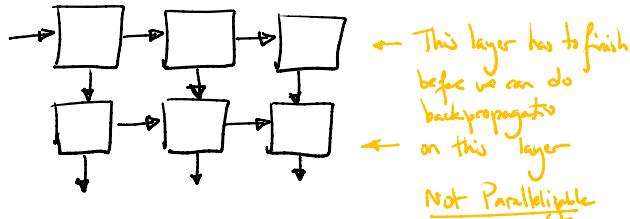
$$\text{Maximize } \sum_t \log P(x_t | x_1, x_2, \dots)$$

↑  
this probability is not as small!

- RNNs can be used to represent these partial probabilities  
How? Often, we let  $y_t$  (the output) be logits  
We go from logits to probabilities via a softmax:

$$P(x_{t+1} = i | x_1, x_2, \dots) = \frac{\exp(y_{t+1}[i])}{\sum_i \exp(y_{t+1}[i])}$$

- Issue with RNN: they're best with multiple layers



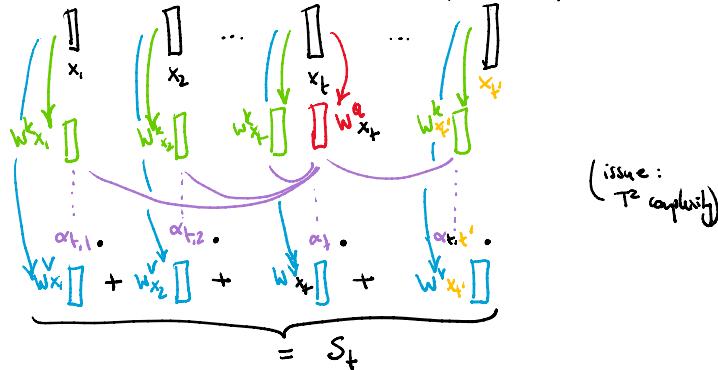
### ② (Self) Attention

- Instead of recurrence, why not directly assess what part information is relevant? (like in databases)
- Query: What the current symbol represents
- Key: The relevance of other sequence symbols to a query.
- Value: What each sequence query contributes

$$\begin{aligned}
 & \text{weights} \quad \text{embeding size} \\
 & Q_{t,h} = W^Q x_t \quad \left. \begin{array}{l} \text{embedding} \\ \text{attention weight} \end{array} \right\} \quad \alpha_{t,h} = \text{softmax} \left( \frac{1}{\sqrt{k}} \frac{Q_{t,h}^T K_{t,h}}{\text{inner product}} \right) \\
 & K_{t,h} = W^K x_t \quad \text{other} \\
 & V_{t,h} = W^V x_t \quad \text{(ignore all } h \text{ at first)}
 \end{aligned}$$

Then, to generate the new state variable, we combine values based on attention:  $S_t = W^O \text{concat}_h \sum_h \alpha_{t,h} V_{t,h}$

We could also have multiple heads, copies w/ different weights



Notes: Each head is given by the triple  $(W^Q, W^K, W^V)$

- The weights don't change by position.
- Position info can be placed inside  $x_t$  (positional embedding)
- Usually dimensions are maintained  $\# \text{heads} \times d = k$ ,  $W^O$  square.

### ③ Transformers

- Attention only tells us where the relevant information is
- We still need to act on it. We do it using:

(a) Add and normalize (Residual connection)

$$S' = \text{Layer Norm}(S + X)$$

$$S' = \gamma \frac{1}{\sqrt{d}} \sum_i S_i + K_i$$

$$\gamma^2 = \frac{1}{d} \sum_i (S_i + K_i)^2$$

$$S'_t = \gamma \frac{S_t + X_t - \mu + \beta}{\sqrt{d+\epsilon}}$$

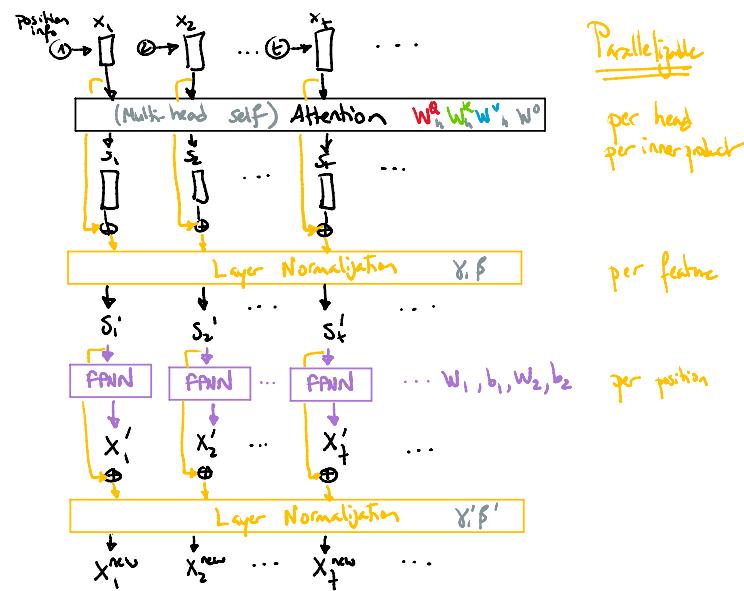
learnable learnable

(b) 2-Layer Feedforward Neural Network

$$X'_t = W^2 \text{ReLU}(W^1 S'_t + b') + b''$$

(c) Add and normalize (again!)

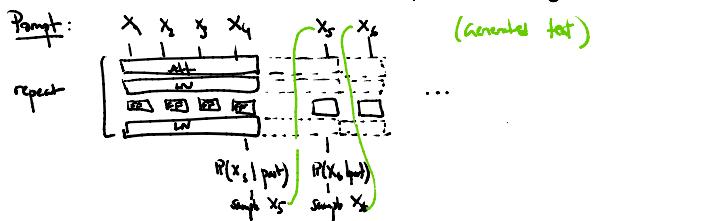
$$X'^{\text{new}} = \text{Layer Norm}(S' + X')$$



#### (4) Using a transformer:

##### (a) Decoder-only architecture (ChatGPT)

- After the prompt, we put outputs as inputs.
- Attention restricted to past outputs (since we generate forward)

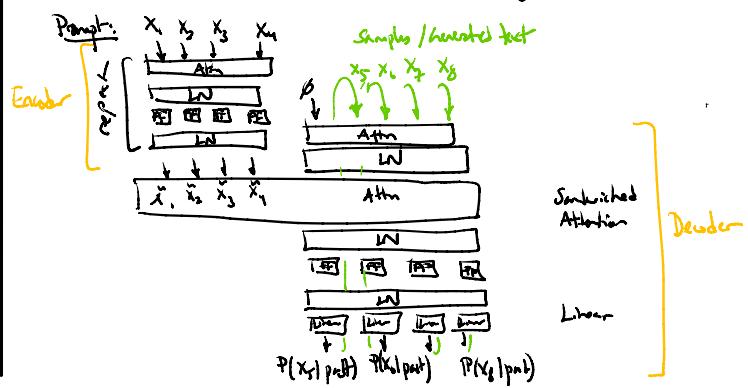


- Variations: Beam Search of most likely generation:

- Maintain k sequences, sample m (expand to km), prune back to top k

##### (b) Encoder-Decoder Architectures

- The outputs of an encoder transformer become the inputs of a decoder transformer, usually in sandwiched layer:



# L24

# Recurrent Neural Networks

Erdem Koyuncu

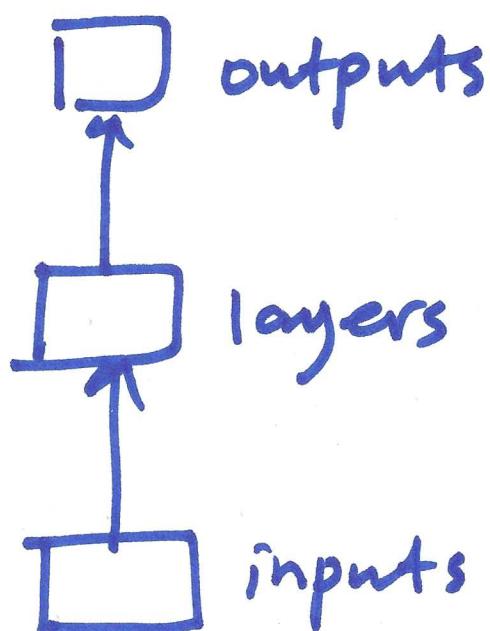
RNNs are used to process sequences (a sequence of images, words, letters, etc) in a way that takes into account the dependencies between the individual elements (or, the memory) of the sequence.

In regular (feedforward) neural networks, we provide an input to the network and observe the output after possibly many hidden layers. There is no feedback of the outputs back to the input.

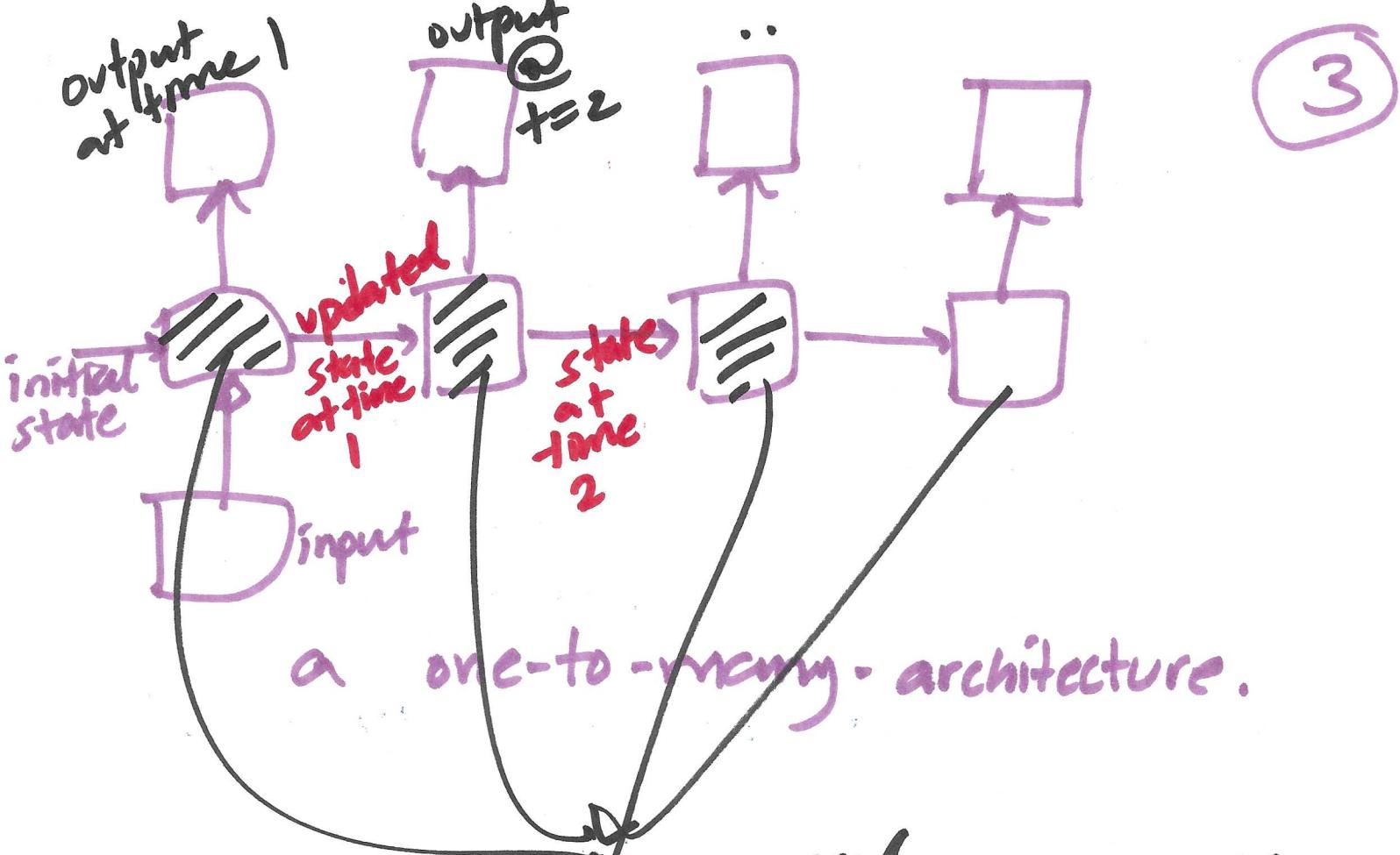
In RNNs, in general, we feed the input sequence one element at a time. Each input updates an internal state of the neural network,

which also serves as the RNN's memory. The output at any given time depends on the input at that time as well as the memory. One can imagine that the network feeds the internal states back to the network.

Ordinary (feedforward) NN:



a one-to-one architecture.



3

a one-to-many-architecture.

same neural network with same weights.

Ex: image description/captioning

Input can be an image:



|||

this is equivalent to

Output can be a sequence of words describing the image  
"person holding an apple"

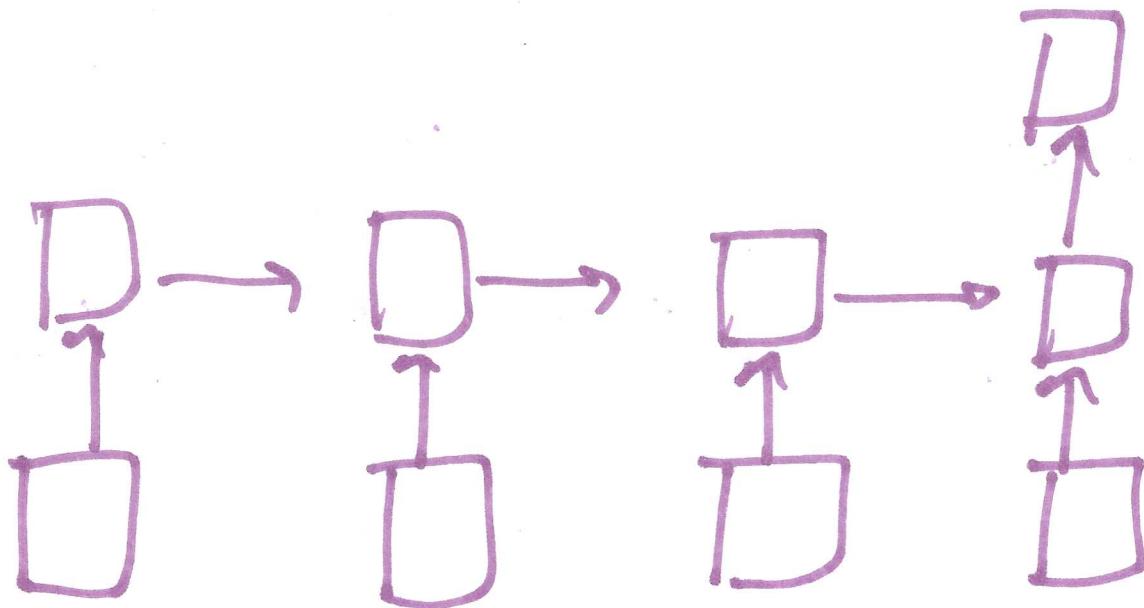
Outputs at different times will be different words.

unrolled 3 times.

state feedback



④



many-to-one architecture.

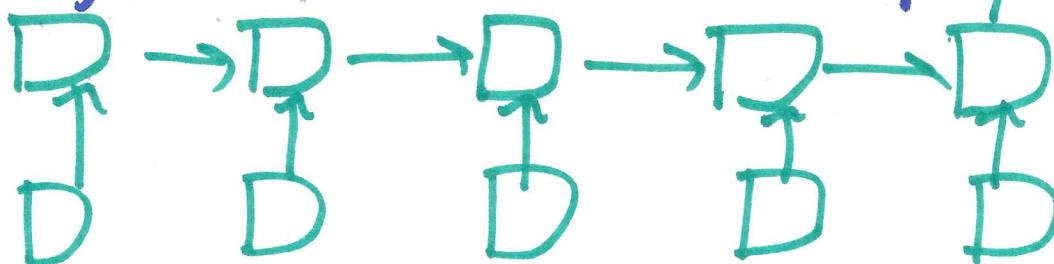
Ex: sentiment analysis

Input is a sequence of words

Output is whether that describes a good or bad sentiment (or any other class that we are interested in)

In most of these problems we will need an "END of sequence" word EOS.

Also, words are typically converted to vectors via an embedding algorithm, before they are input to the RNN.



"The" "movie" "was" "terrible" EOS

So, we are not inputting "words" of course but their embedded versions. 34

Another application of the many-to-one architecture could be a general time-series prediction problem:  $x_t$

(5)

