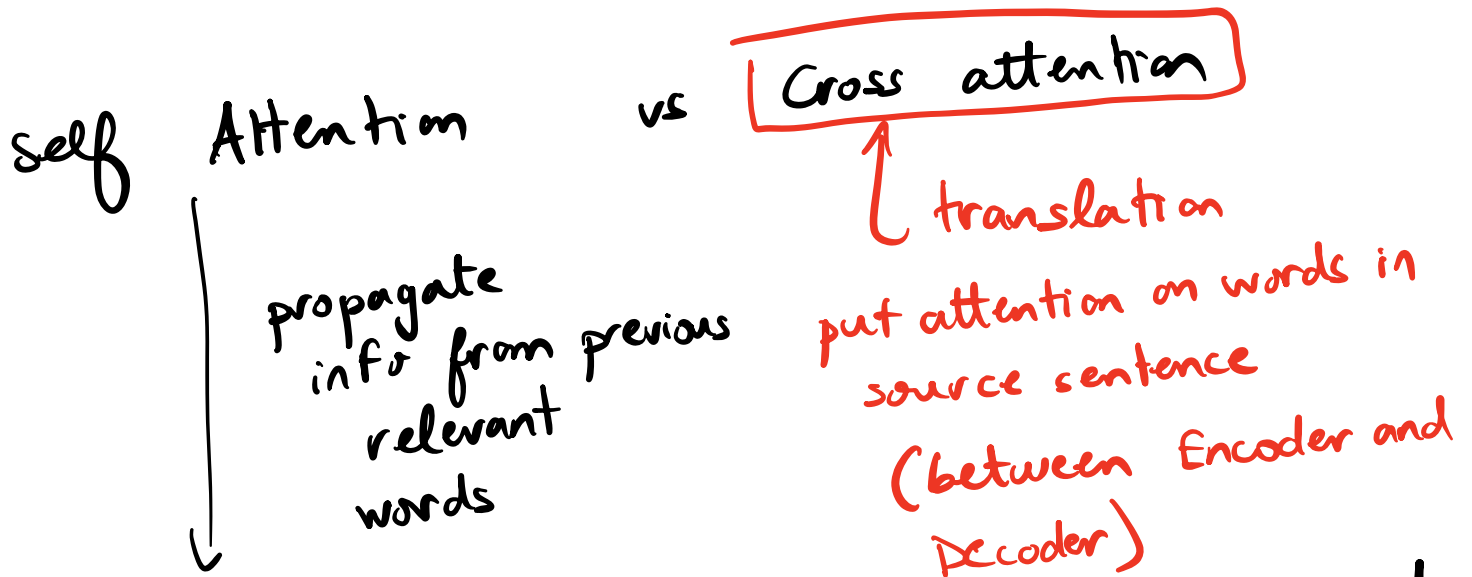


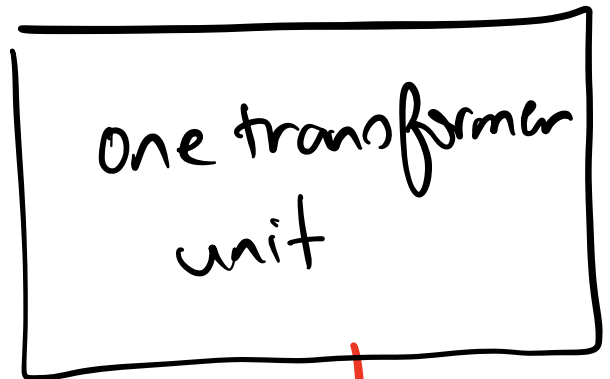
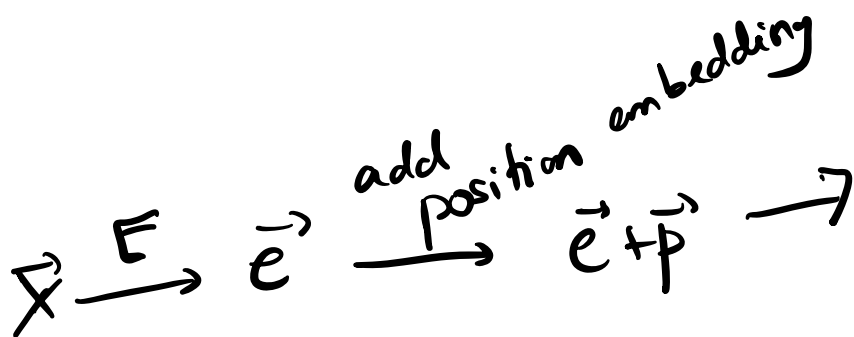
CS 6180 10/15

Review from last time

Transformer architecture (GPT, a lot of generative models)



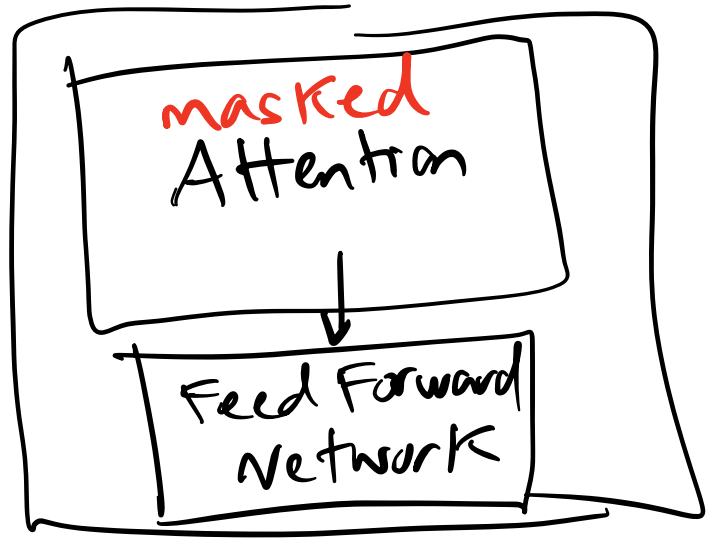
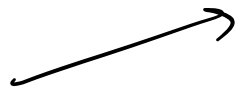
Let's build some models that are purely based on Attention.



Zuko made his uncle tea at Zuko's place.

His uncle made Zuko tea at Iroh's place.

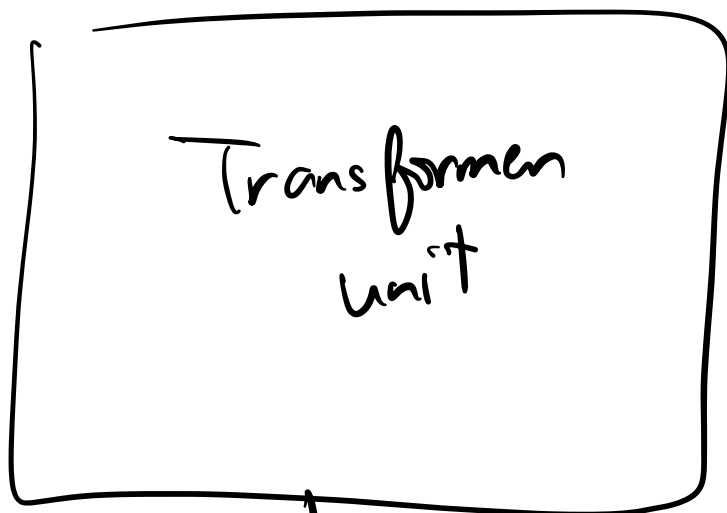
Transformer  
unit



Decoder  
model

masking : Set attention scores of future words to  $-\infty$  so that their contribution will be 0 in the softmax.

$$\vec{x} \rightarrow \vec{e} \rightarrow \vec{e} + \vec{p} \rightarrow$$



project back  
to the  
space of  
word vectors  
and apply softmax

Add & Norm

(component will be added to  
the transformer unit)

\* Layer Normalization

example:

$$z = Wx + b$$

$$x \sim N(0, 1)$$

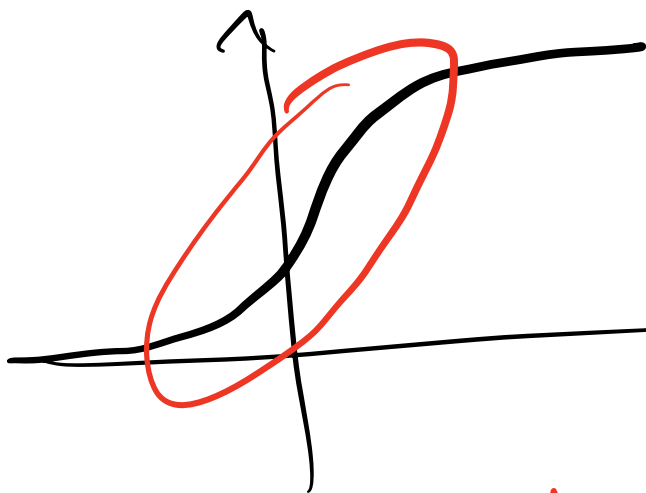
$$\begin{cases} h = \sigma(z) \end{cases}$$

ex1

$$W = 1$$

$$b = 0$$

$$z \sim N(0, 1)$$



regime where the  
gradients are not  
zero

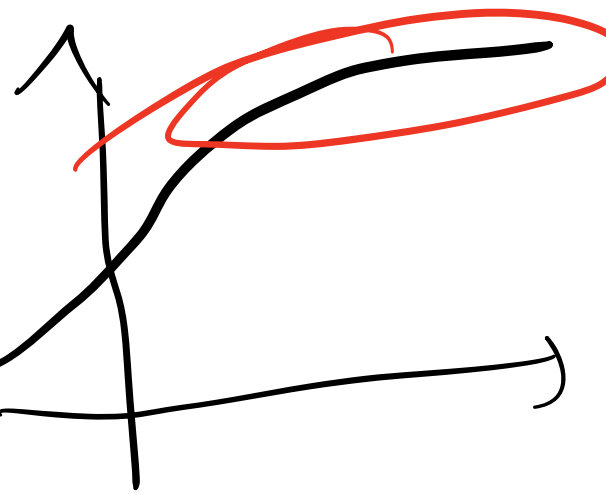
$\Rightarrow$  learning will be  
happening  $\Rightarrow$  yay

ex2

$$W = 3$$

$$b = 2$$

$$z \sim N(2, 9)$$



gradients  
are closer to  
0

learning  
will be  
slow

might even  
get  
interrupted.

⇒ vanishing  
gradients

⇒ Let's normalize

$$\hat{z} = \frac{z - \mu}{\sqrt{\sigma^2}}$$

$$E(\hat{z}) = 0$$

$$\text{Var}(\hat{z}) = 1$$

$\mu$  = average over the features

$\sigma^2$  = variance over the features

$\Rightarrow$  will go back to the "nice" area of the sigmoid so that the gradients won't vanish.

## Limitations

$\times$  if  $\sigma^2$  is small  
 $\Rightarrow$  dividing by a very small number  
 $\Rightarrow$  exploding values occurring with  $\hat{z}$

fix: 
$$\hat{z} = \frac{z - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\epsilon \sim 10^{-5} - 10^{-6}$$



helps with  
stability of the  
normalization so that  
we don't have  
exploding  $\hat{z}$ .

\* We might need to  
end up at the  
regimes where the  
gradients are close to  
0 (especially when we  
have almost converged)

Fix:  $y = \gamma \hat{z} + \beta$

learnable parameters

$$\gamma = \sigma$$

$$\hat{z} = \frac{z - \mu}{\sigma}$$

$$\beta = \mu$$

$$\hat{z} \sigma + \mu = z$$

\* at the beginning of the model

$\Rightarrow$  the model will try to keep the features standardized  
( $\gamma = 1, \beta = 0$ )

\* As the model is learning and has learned quite a bit (almost converged)  
model now has flexibility (through  $\gamma$  and  $\beta$ ) to go back to



the unstandardized features.  
(yay)

mainly helps with training

we have flowing gradients  
and could also end up with the  
0 gradients when converged.

## Residual addition

Add & Norm

(Attention is all  
you need, original  
transformer)

another  
version of  
Add & Norm

Used in GPT  
models  
for example)

## Approach 1

AH : attention heads

LN : layer normalization

FFN : feed forward  
Network

$x \rightarrow$  transformer unit

$$y_1 = \text{LN}(x + \text{ATT}(x))$$

$$y_2 = \text{LN}(y_1 + \text{FFN}(y_1))$$

residual ( adding  $x$  to  $\text{ATT}(x)$  )  
"  $y_1$  to  $\text{FFN}(y_1)$  )

making sure we have not  
lost any previous info and  
would keep on adding  
new info.

## Approach 2

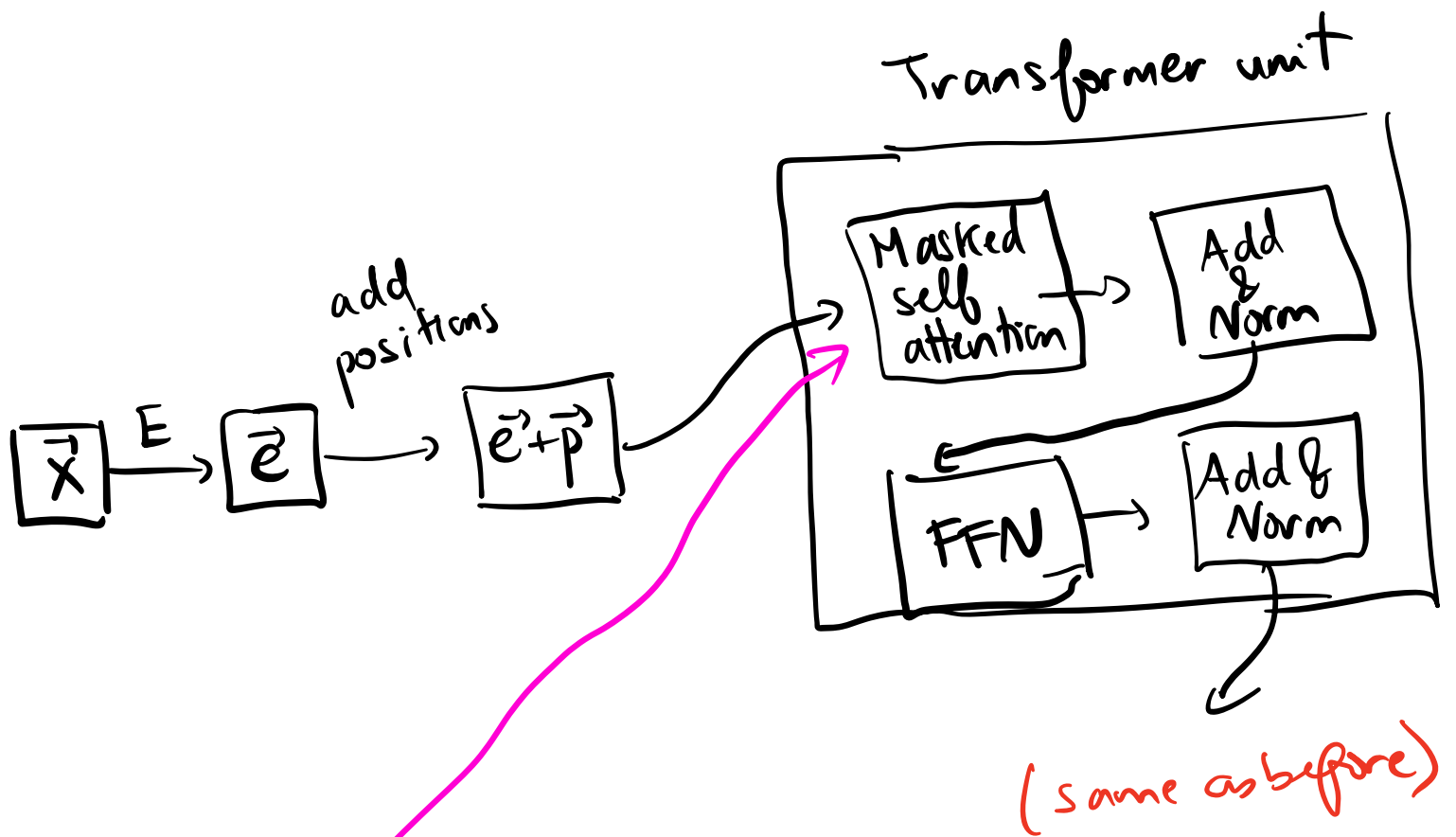
$$y_1 = x + \text{Act}(\text{LN}(x))$$

$$y_2 = y_1 + \text{FFN}(\text{LN}(y_1))$$

gradients flow  
more directly

this is more stable  
for deep networks

(more used now)



Instead of only one attention head  
 why not use many so that each  
 attention head could attend to  
 different features of the sentence?  
 (meaning, grammar, pronoun, ...)

Attention head

$$\vec{K}_i = K \vec{e}_i ; \quad \vec{q}_i = Q \vec{e}_i ; \quad \vec{V}_i = V \vec{e}_i$$

$$S_{ij} = \vec{q}_i^T \vec{K}_j$$

$$j=1, 2, \dots, m$$

$$\text{softmax}(\vec{S}_i)$$

"  $d_{ij}$

$$= \begin{bmatrix} \frac{\exp(S_{i1})}{\sum_k \exp(S_{ik})} \\ \vdots \\ \vdots \end{bmatrix}$$

$$\vec{O}_i = \sum_j d_{ij} \vec{V}_j$$

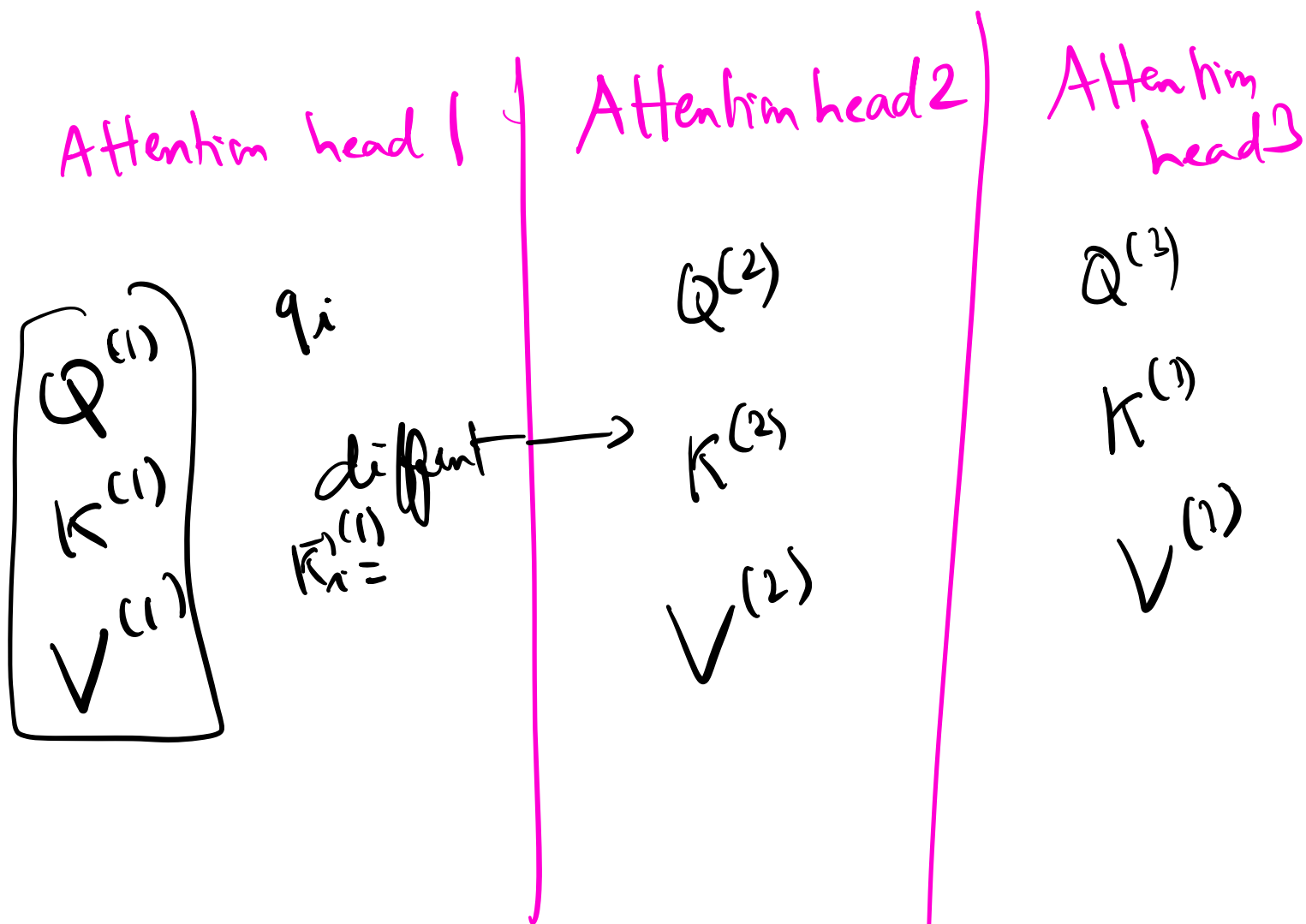
$$\vec{O}_i^{(1)}$$

$$\vec{O}_i^{(2)}$$

$$\vec{O}_i^{(3)}$$

parameters

$$\underbrace{\begin{bmatrix} \sigma_i^{(1)} & \sigma_i^{(2)} & \sigma_i^{(3)} \end{bmatrix}}_{d \times 3} \underbrace{V}_{3 \times 1}$$



each attention head will focus on a different task

↳ (intuition hope)

how do we know what they're  
focusing on? Open research  
question

Mechanistic  
interpretability

$$O^{(l)} = \text{softmax} \left( \frac{(XQ)(XK)^T}{\sqrt{d}} \right) XV$$

if dot products  
become  
very large  
→ let's scale by  
dimension

Transformers, LSTMs, RNNs

# Autoregressive models

what are autoregressive models?

$$\begin{aligned} &P(x_1, x_2, x_3, \dots, x_m) \\ &= P(x_1) \cdot P(x_2 | x_1) \\ &\quad P(x_3 | x_1, x_2) \\ &\quad \vdots \\ &\quad P(x_m | x_1, x_2, \dots, x_{m-1}) \end{aligned}$$

autoregressive

property

where

$x_m$

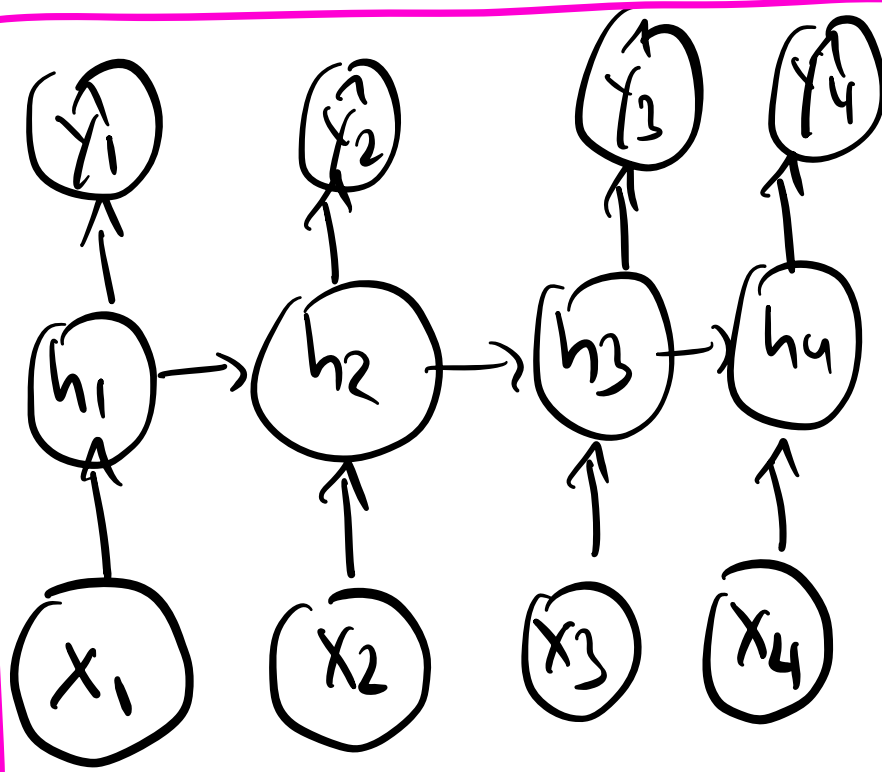


depends on the past words  
that occurred.

→ literature on time  
series (observations from  
previous times to predict  
an output at the next  
time step.)

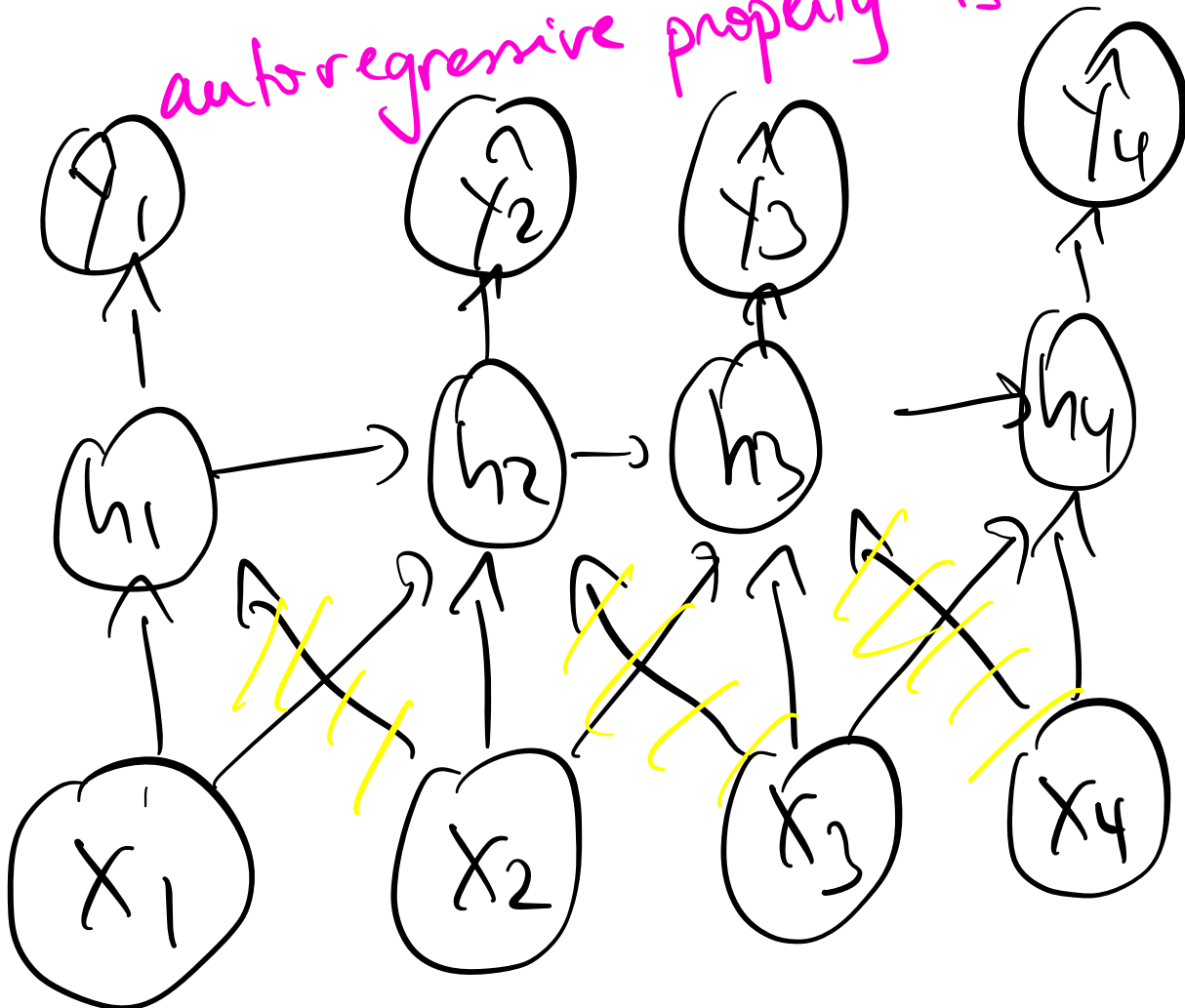
$$p(X_i | X_1, \dots, X_{i-1})$$

can use a neural network  
to model this function  
with a certain # of  
parameters



RNN

autoregressive property is satisfied



CNN

not autoregressive  
model.

removing these 3 edges  
makes it autoregressive.

$$\begin{array}{r} 1-4 \\ 1-2 \quad 2 \\ 2-3 \quad 1 \\ 3-4 \quad 6 \\ \hline 21 \end{array}$$

priority  
for you  
all

attention  
head 1

: Grammar

$\vec{K}_i$

$\vec{q}_i$   
↓

$\vec{V}_i$

attention  
head 2:

meaning

↓  
teach  
taught  
↑