

CS 6180 10/08

HW2 released Yaaaay

Review from last time

Attention

Translation

source language \rightarrow target language

sequence to sequence model

wǒ jiào
Nàdìmx

Encoder
LSTM



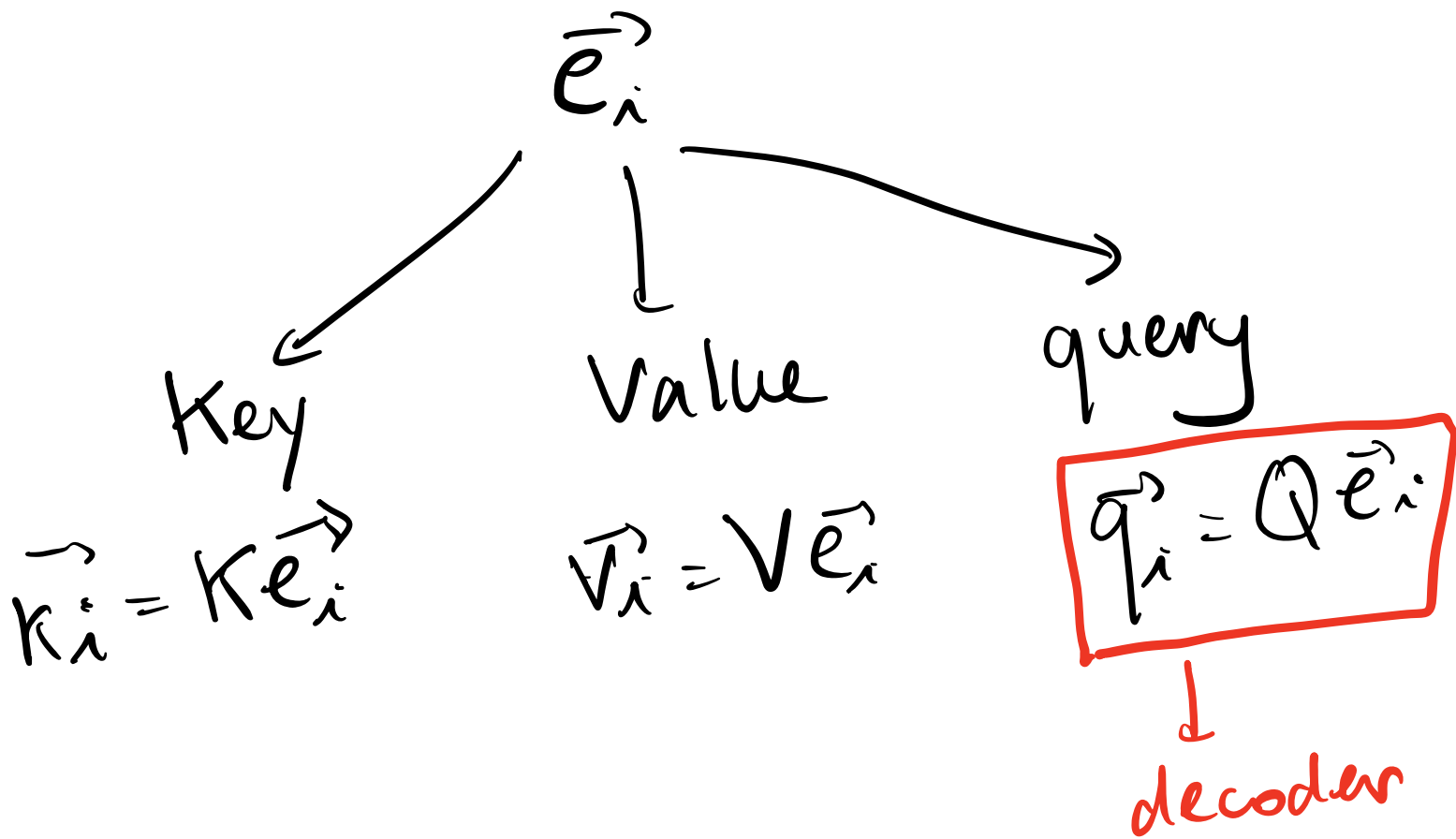
My name is
Nadim

Decoder
LSTM

Attention (Cross-Attention)
(between different models)

put attention on the source sentence

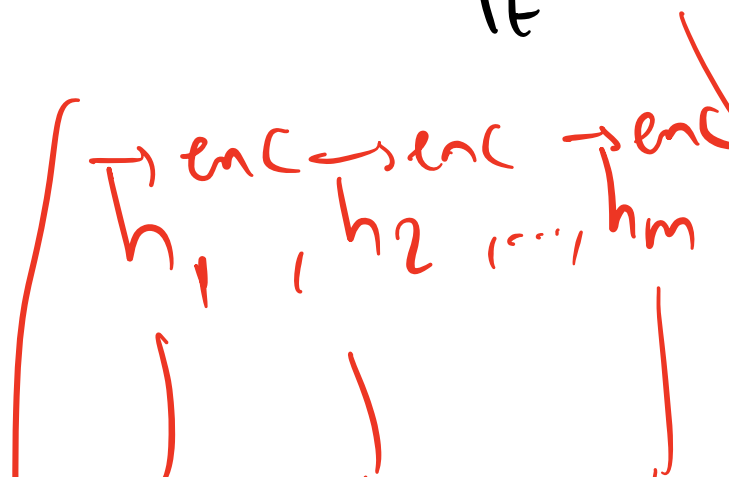
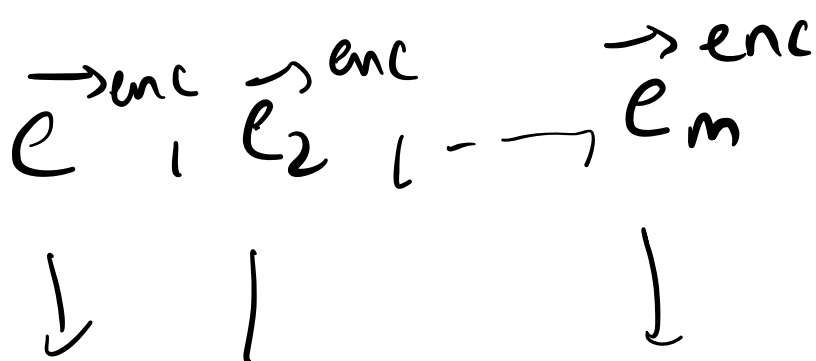
to generate a word in the target sentence.



Key representation

My name is
 \uparrow
 query
 \vec{q}_t

source sentence



$\vec{K}_1 \quad \vec{K}_2 \quad \dots \quad \vec{K}_m$

$\vec{K}_1 \quad \vec{K}_2 \quad \vec{K}_m$

similarity between Keys and queries

$$\begin{array}{ccccccc} \vec{K}_1^T \vec{q}_t & \vec{K}_2^T \vec{q}_t & \dots & \vec{K}_m^T \vec{q}_t & & & \\ \text{"} & \text{"} & & \text{"} & & & \text{(attention scores)} \\ s_{1t} & s_{2t} & & s_{mt} & & & \end{array}$$

softmax to scale the scores

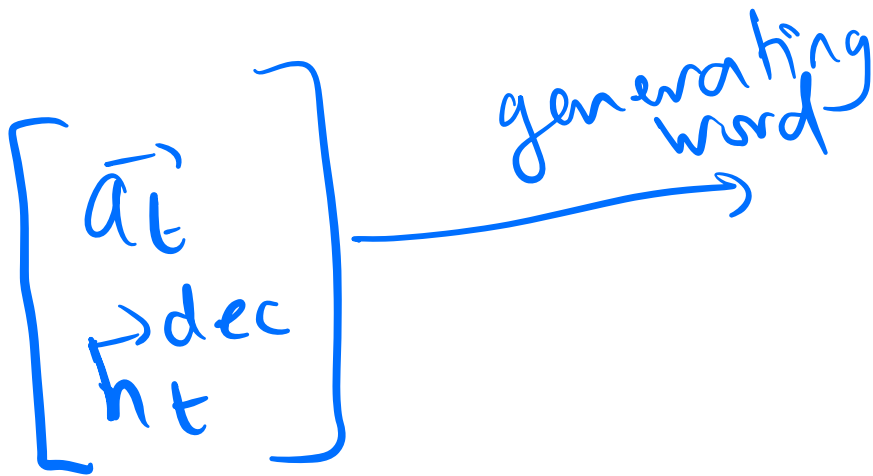
$$\frac{\exp(\vec{K}_i^T \vec{q}_t)}{\sum_{j=1}^m \exp(\vec{K}_j^T \vec{q}_t)} = \alpha_{it}$$

scaled
attention
scores

(attention
distribution)

attention output

$$\vec{a}_t = \sum_{i=1}^m \alpha_{it} \vec{V}_i$$



taking into account words in the original sentence (with attention)

$$\vec{h}_t^{\text{dec}}, \vec{c}_t^{\text{dec}} = \text{LSTM}(\vec{h}_{t-1}^{\text{dec}}, \vec{c}_{t-1}^{\text{dec}}, \vec{e}_t)$$

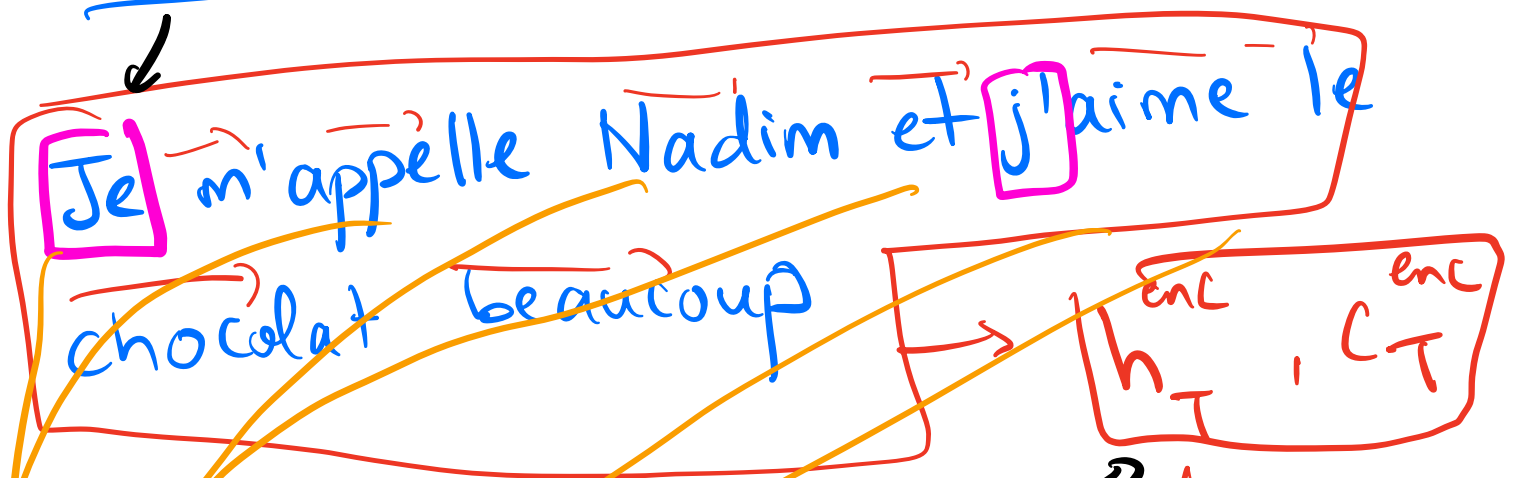
replace

$$\begin{bmatrix} \bar{e}_t \\ d_t \end{bmatrix}$$

how much info was propagating
from the source sentence to
the target one

$O(\text{length of source sentence})$
between relevant words in the
source sentence and target
(previously)

with attention $O(1)$



interaction
distance

input
decoder

between the word I want to
generate now and a relevant
source word is in the order
of the length of the source
sentence

My

With attention

interaction distance = 1

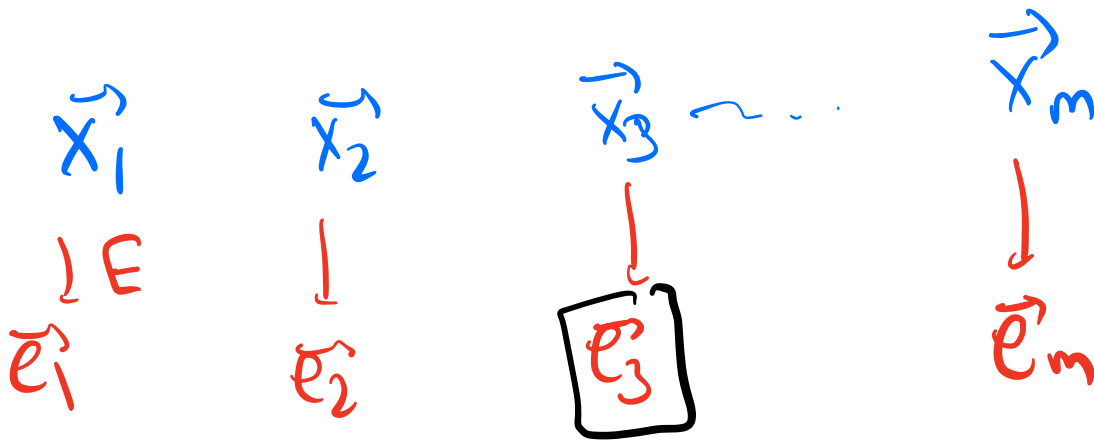
(getting some direct info from every source word)

Attention also helps with propagating info from previous words. That's why we looked at RNNs.

Why not give up on RNNs (LSTMs)

and build models that are entirely based on attention?

Self attention: attention framework
WITHIN the same sentence.



$$\begin{aligned} \vec{e}_i &\rightarrow \vec{q}_i = Q \vec{e}_i \\ &\rightarrow \vec{k}_i = K \vec{e}_i \\ &\rightarrow \vec{v}_i = V \vec{e}_i \end{aligned}$$

$$s_{ij} = \vec{q}_i^T \vec{k}_j \quad \alpha_{ij}, \quad \vec{a}_i = \sum_j \alpha_{ij} \vec{v}_j$$

* Need to make ^{sure} we cannot peek at future words

Come up with a general framework for every word in the sentence

* Order (position) of words is not taken into account.

Zu Ko made his uncle
some tea.

vs

His uncle made Zu Ko some
tea.

Same words in both sentences but

different meanings.

* Need non-linearities
to be able to handle
complex data.

* Let's fix first the
issue of word
position

Let's add position
vectors


$$\begin{array}{lcl} \vec{x}_1 & \rightarrow & \vec{e}_1 + \vec{p}_1 \\ \vec{x}_2 & \rightarrow & \vec{e}_2 + \vec{p}_2 \end{array}$$

$$\vec{X}_m \rightarrow \vec{C}_m \perp \vec{P}_m$$

Zuko made his uncle tea at Zuko's place.

vs

His uncle made Zuko tea at Iroh's place.

if we're only using self attention
the output at  should be the same.

We don't want the model to learn that. Adding position parameters to the model will give it flexibility to learn the difference between the sentences.

W

*Need to make sure that we cannot peek at future words.

Masking

t=4

Zuko made his uncle fear. at Zuko's

s_{1t} s_{2t} s_{3t} s_{4t} s_{5t} s_{6t} s_{7t}

need to
not
contribute
to my attention
output

$$\vec{a}_t = \sum_{i=1}^7 \alpha_{it} \vec{v}_i$$

$$\alpha_{6t} = \frac{\exp(S_{6t})}{\sum_{i=1}^7 \exp(S_{it})}$$

↓

0

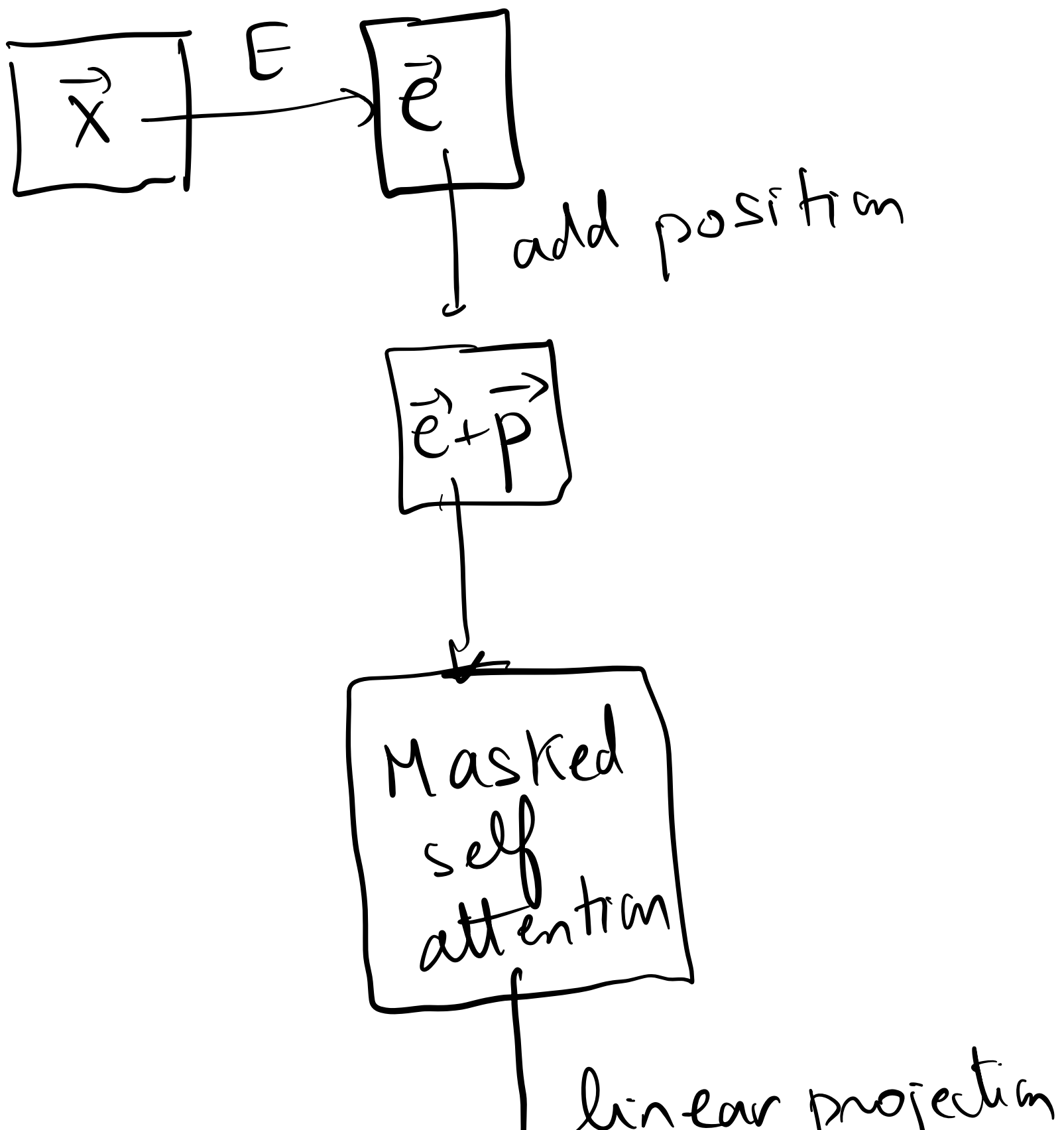
set $S_{6t} = -\infty$

$S_{7t} = -\infty$

$$S_{it} = \begin{cases} \vec{q}_t^T \vec{K}_i & \text{if } i \leq t \\ -\infty & \text{if } i > t \end{cases}$$

ensures that we are not peeking

Way for masking



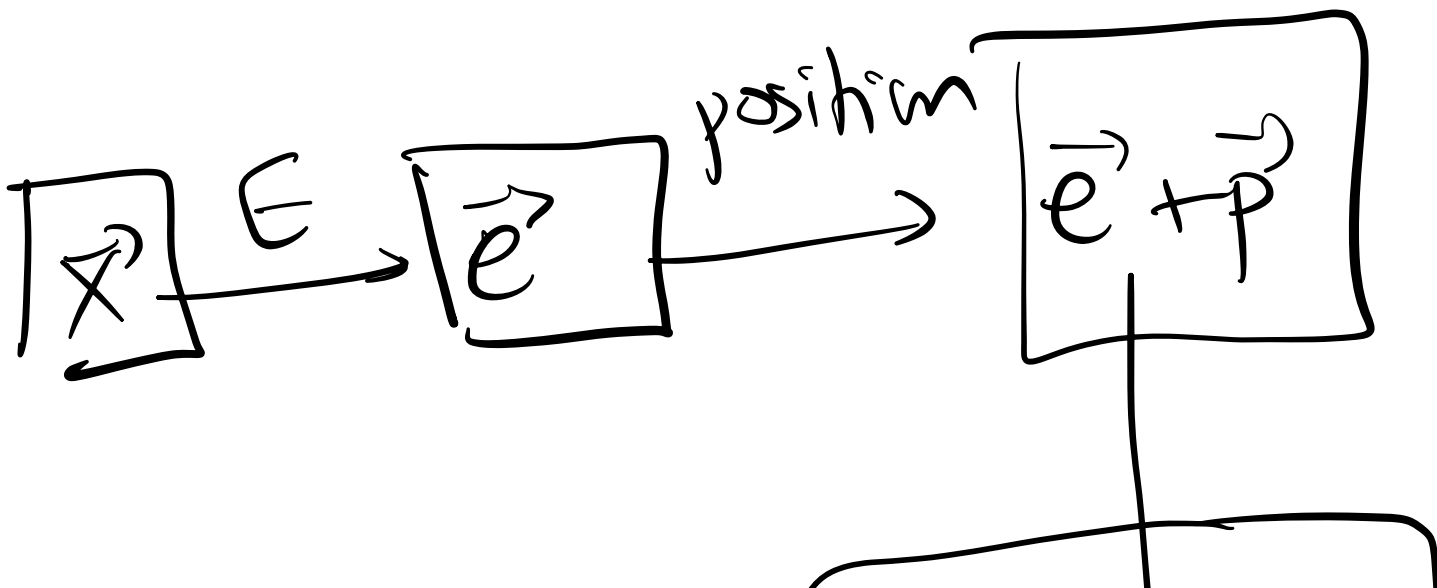
softmax

Multi layer perceptron

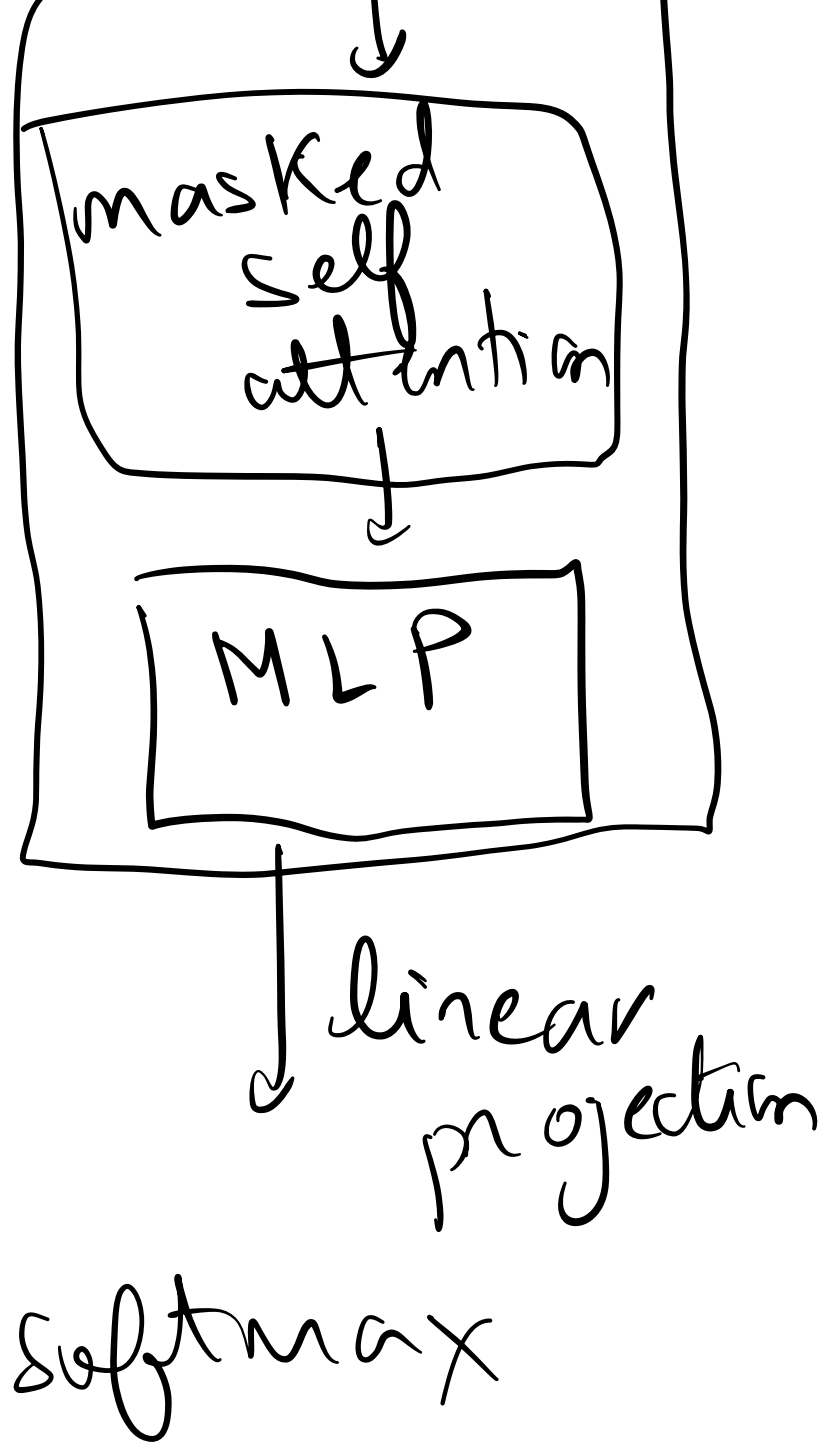
$$\vec{h}^{(1)} = \sigma(W^{(1)} \vec{x} + \vec{b}^{(1)})$$

$$\vec{h}^{(2)} = \sigma(W^{(2)} \vec{h}^{(1)} + \vec{b}^{(2)})$$

⋮



repeat →
such
unit
several
times



(Almost)

Transformer (Decoder)
architecture

3 main components to
add to get the full
transformer architecture

* have more than
one attention
head (Q, V, K)

→ meaning
→ grammar
→ syntax

1
2
3
4