CS 6180

<u>Review from last time</u>

LSTMs ( specific type of RNNs)

<u>Machine Translation</u>

task of generating a sentence y from one language to another.

French  Il faut profiter du moment présent

$\downarrow$

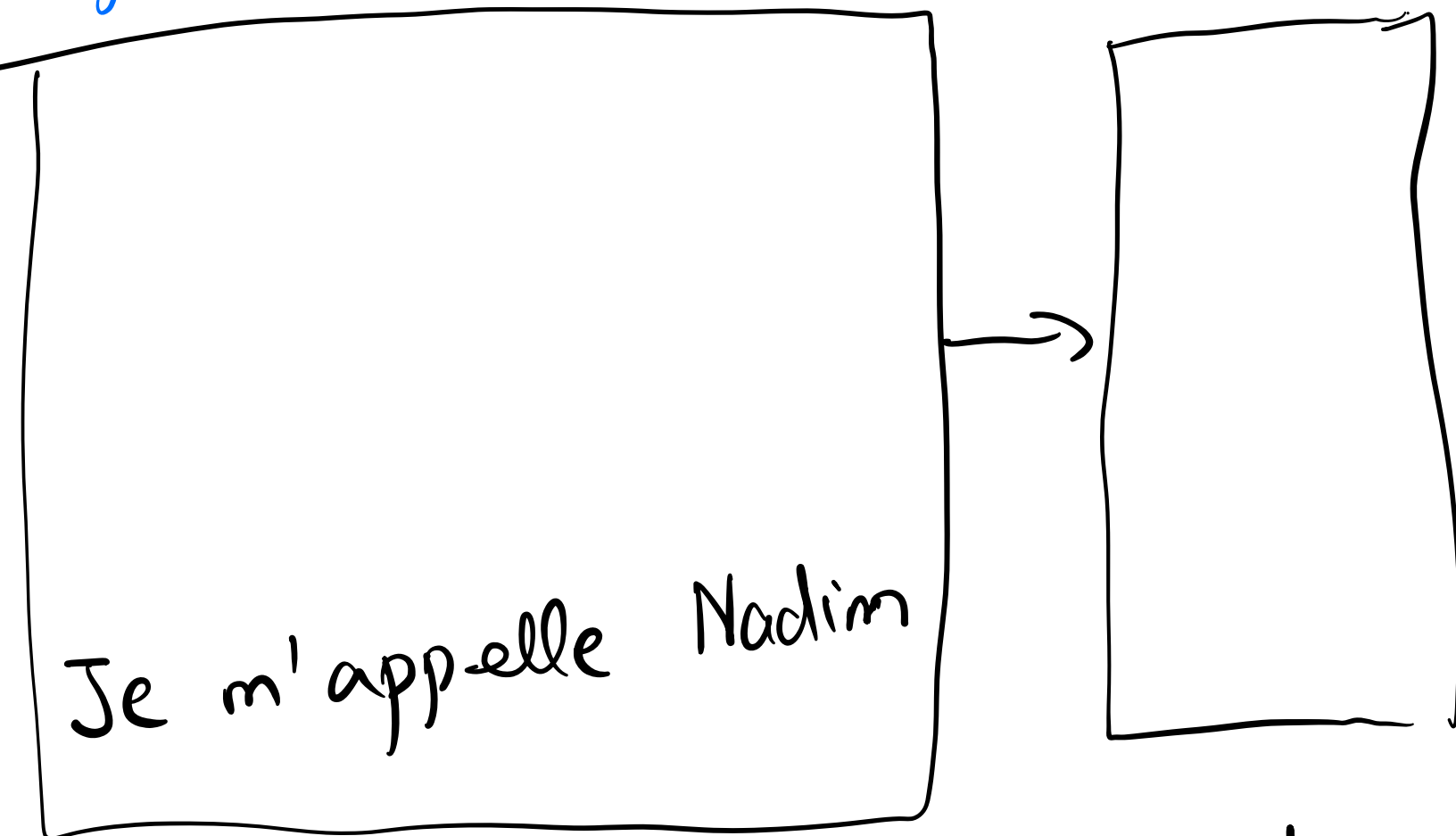English  You have to enjoy the present moment

(Carpe Diem) (Latin)

<u>Learn</u>  model

$$p(y \mid x)$$

(sentence in French)

# Goal

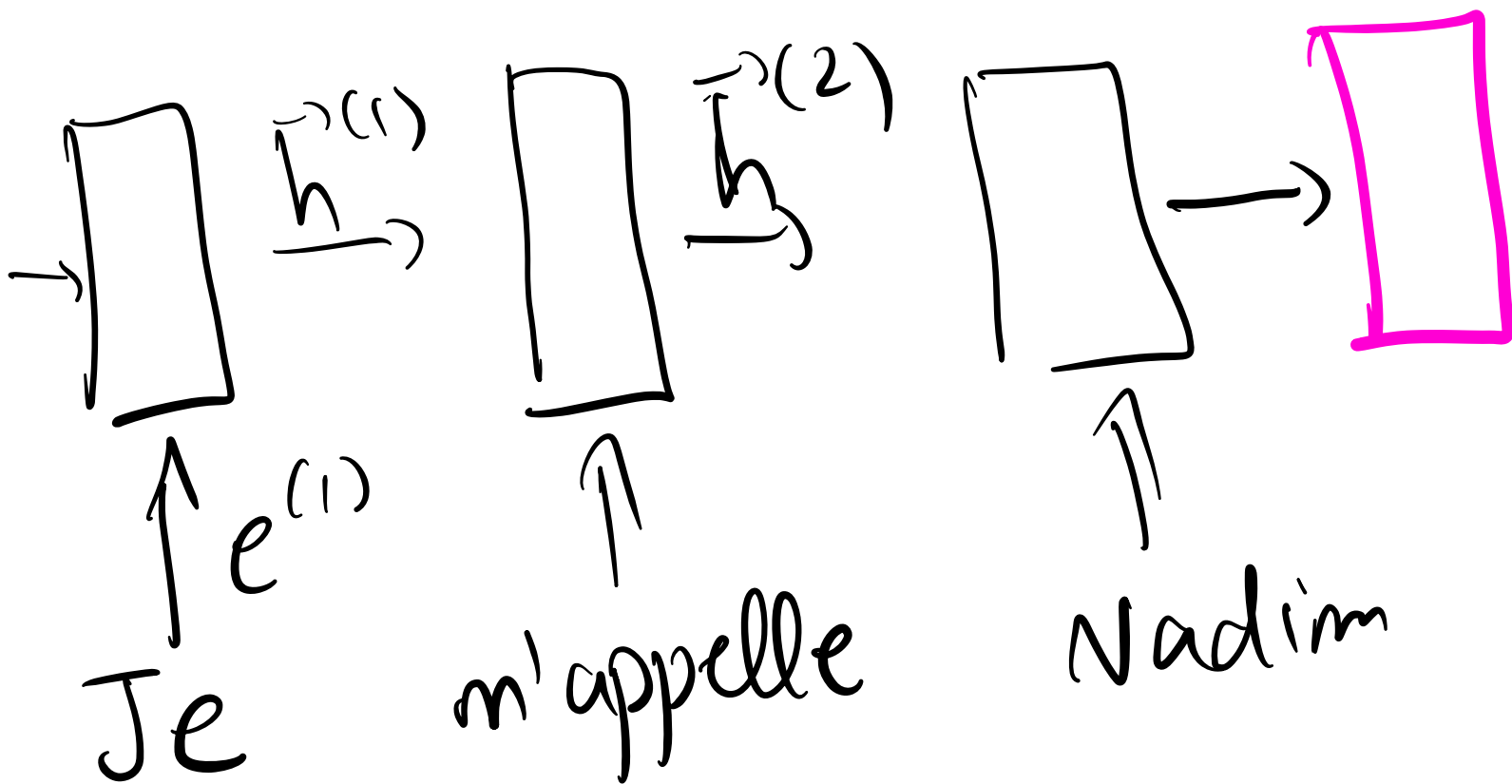maximize this probability

in English

## Neural Machine Translation

Je m'appelle Nadim

Encoder

Decoder

responsible of the translation

(Encoder RNN)



$\vec{h}^{(1)}$ $\vec{h}^{(2)}$

$e^{(1)}$

Je     n'appelle     Nadim

## Decoder RNN

my     name     is

start my name

Nadim
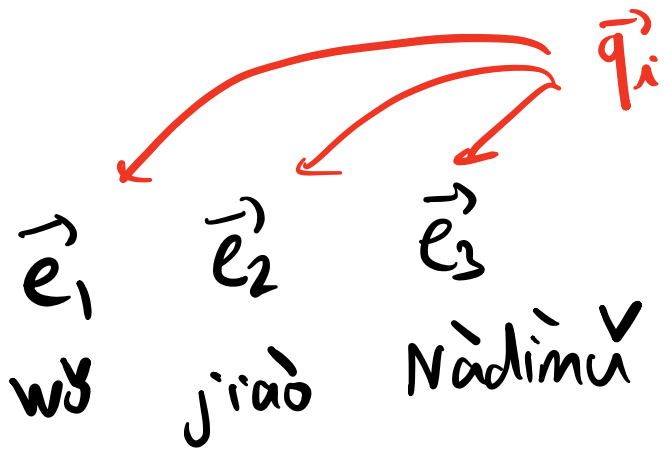
end

Nadim

IS

Nadim

* they were the best around 2017 for translation
* not parallelizable
computations can be expensive

* O(sequence length) info
  stored into one vector [ ]

⟹ Attention can help

(image on board)

* dot product between words in the source sentence and word to be generated

$\vec{q_i}$

$\vec{e_1}$    $\vec{e_2}$    $\vec{e_3}$

wǒ    jiào    Nàdìmǔ

$$\left.\begin{array}{c} \vec{q_i}^T \vec{e_1} \\ \vec{q_i}^T \vec{e_2} \\ \vec{q_i}^T \vec{e_3} \end{array}\right\}$$ attention scores
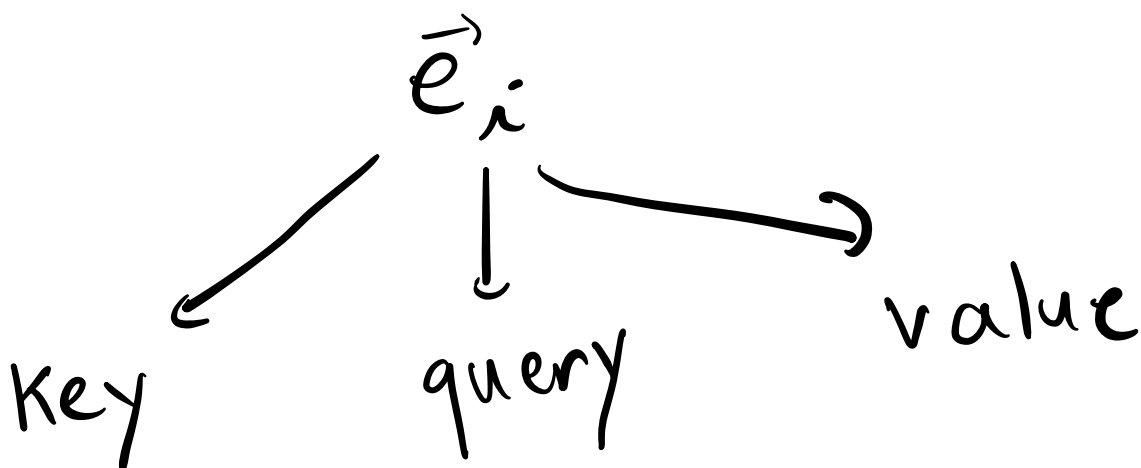
apply softmax to get weights →

$$\alpha_1 = \frac{\exp(\vec{q_i}^T \vec{e_1})}{\sum_j \exp(\vec{q_i}^T \vec{e_j})}$$

$$\alpha_2 = \frac{\exp(\vec{q_i}^T \vec{e_2})}{\sum_j \exp(\vec{q_i}^T \vec{e_j})}$$

$$\alpha_3 = \frac{\exp(\vec{q_i}^T \vec{e_3})}{\sum_j \exp(\vec{q_i}^T \vec{e_j})}$$

$$\text{Output} = \alpha_1 \vec{e}_1 + \alpha_2 \vec{e}_2$$
$$+ \alpha_3 \vec{e}_3$$

we are asking too much from the word embeddings

→ meaning of the word

→ similarity to other words

→ similarity in other languages

$\vec{e}_i$

Key    query    value

$$\vec{K}_i = \boxed{K}\, \vec{e}_i \qquad \vec{q}_i = \boxed{Q}\, \vec{e}_i \qquad \vec{v}_i = \boxed{V}\, \vec{e}_i$$

$d \times d \qquad\qquad d \times d \qquad\qquad d \times d$

matrices that we have to learn (parameters)

dot products :
$$\vec{q}_i{}^T \vec{K}_j$$
$$= (Q\vec{e}_i)^T (K\vec{e}_j)$$
$$= \vec{e}_i{}^T \boxed{Q^T K}\, \vec{e}_j$$

M

to train the model computationally, separating the Q and K leads to a much more efficient implementation.. (Let's see why)

$d_Q \times 1$

$$\vec{q}_i = \boxed{Q}\, \vec{e}_i \qquad\qquad \vec{K}_i = \boxed{K}\, \vec{e}_i$$

$d_Q \times 1 \quad \boxed{d_Q \times d} \qquad\qquad \boxed{d_Q \times d}$

what if $d_Q \ll d$?

$\boxed{2 d d_Q}$ total # of parameters for both Q and K

$vs$   $M = Q^T K$

$\underbrace{\underset{d \times d_Q}{} \quad \underset{d_Q \times d}{}}_{d \times d}$
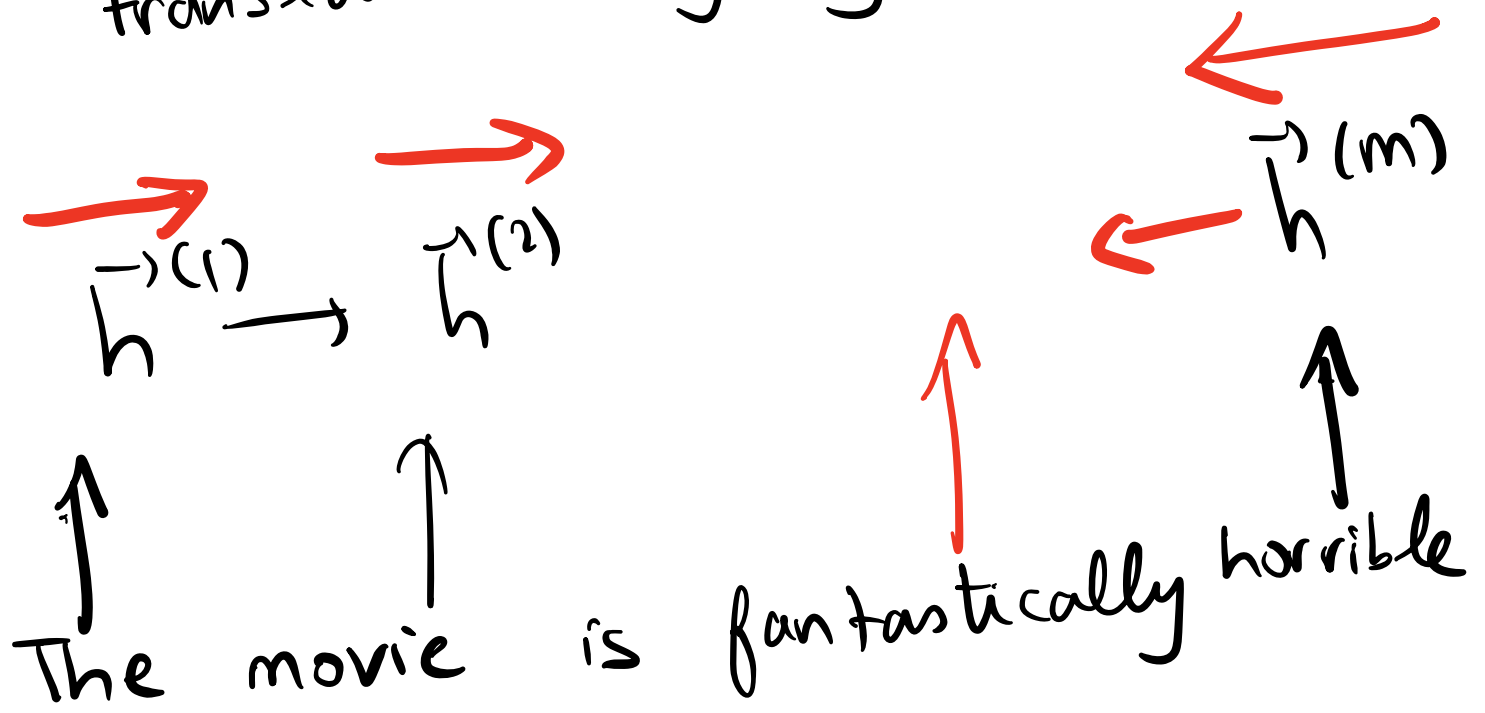
$d^2$ parameters

$2 d d_Q$   $vs$   $d^2$   $\sim$if $d_Q << d$

less parameters to keep track of
more efficient computationally.

(note to Nadim: add two other
screenshots)

Attention can be used for Language models
in general (not just translation)

translation (using only decoder)

$\overrightarrow{h}^{(1)} \longrightarrow \overrightarrow{h}^{(2)}$ ... $\overrightarrow{h}^{(m)}$

The movie is fantastically horrible

LSTM ($\longrightarrow$) left to right

LSTM ($\longleftarrow$) right to left

RNN

$\overrightarrow{h}^{(t)} = \sigma\left(W_h \overrightarrow{h}^{(t-1)} + W_e \overrightarrow{e}^{(t)} + \overrightarrow{b}_p\right)$

$\overleftarrow{RNN}$

$$\overleftarrow{\overrightarrow{h}}^{(t)} = \sigma\left(W_h \overrightarrow{\overleftarrow{h}}^{(t+1)} + W_e \overrightarrow{e}^{(t)} + \overrightarrow{b}\right)$$

$$\text{word } t = \begin{bmatrix} \overrightarrow{h}^{(t)} \\ \overleftarrow{h}^{(t)} \end{bmatrix}$$

Bi-directional LSTMs can apply for encoders Not for the decoders

NoA for Language model