

Lecture 1 → Lecture 3



$x \in \mathbb{R}^D$



- Binary Classification
- Multiclass classification
- Regression.

① PS 3

→ Preprocessing \Rightarrow

→ Feature engineering.

↙ One hot encoding ⊗
↘ Scaling

→ Model

L2

Tree based Models

⇓
HP

- depth
- # trees
- Impurity
- etc.

L4

NN

⇓
HP

- # neurons / layer
- # layers
- Act fcts ...



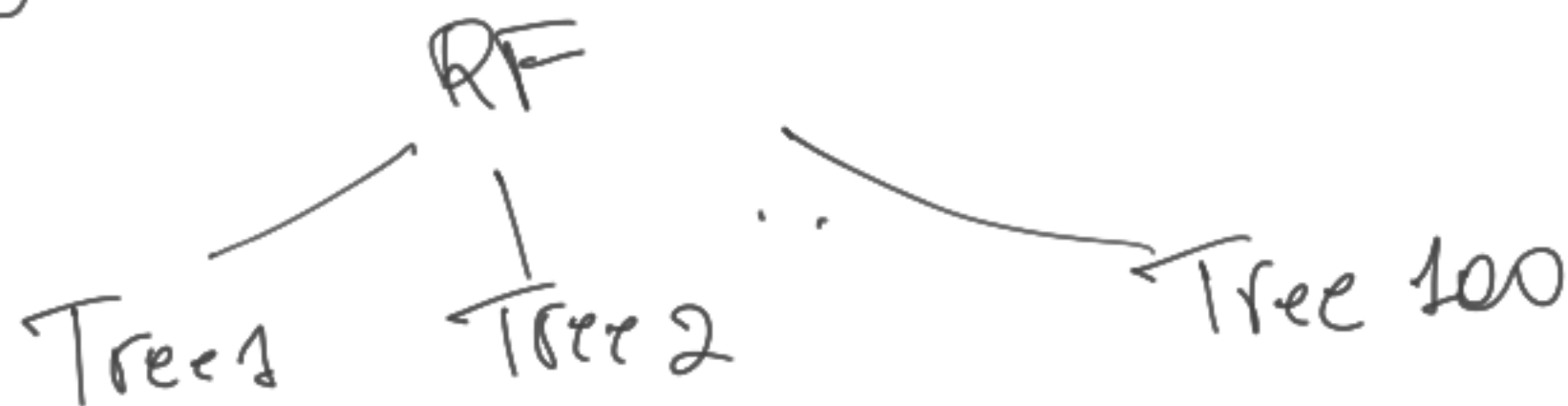
$x \in \mathbb{R}^D$

→ Pick HP :

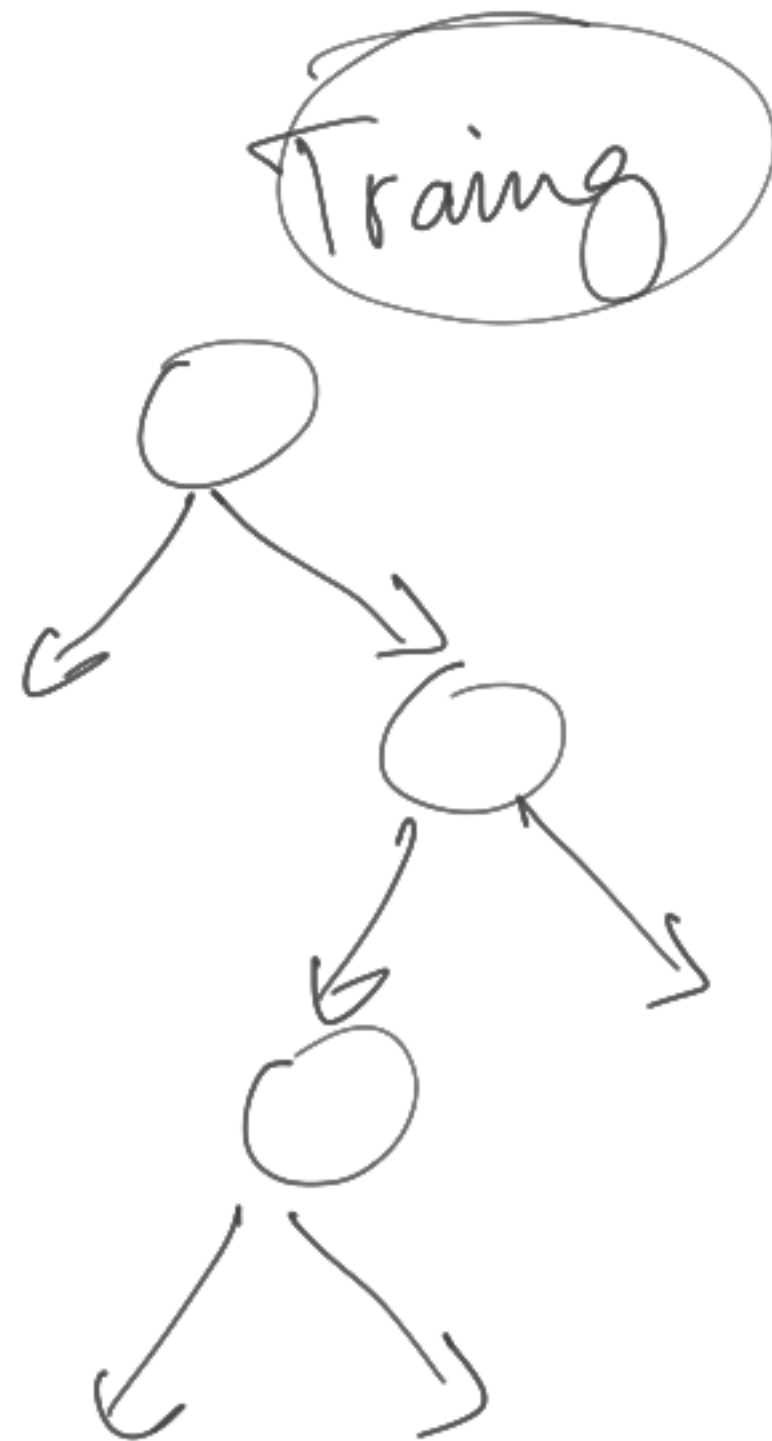


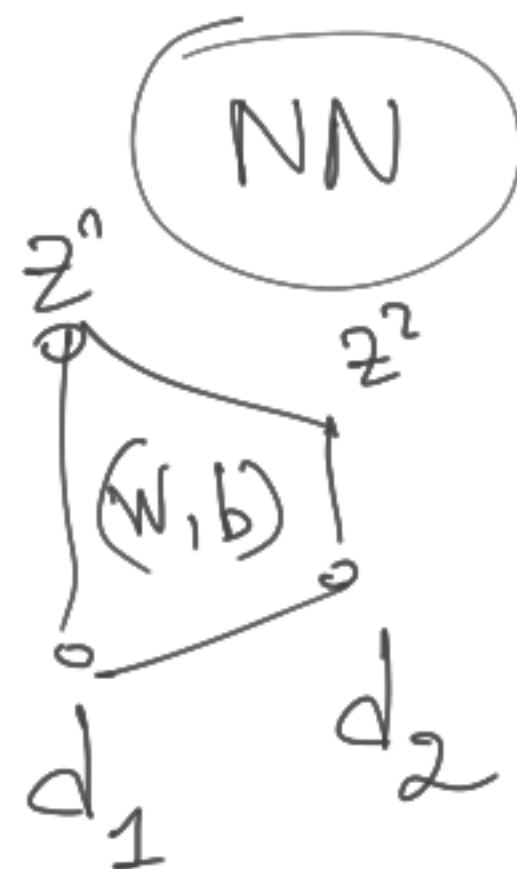
"Test"

$h = (100 \text{ trees, Entropy, } 200)$



④ evaluation metrics (21)





$$z^2 = \sigma(w^T z^1 + b)$$



Part 2 Introducing Sequentially

PSU, Imdb dataset

doc 1



y_1

2



V-dim

.

.

.

doc n

- -

- -

y_n

Limitation in PSIR

→ Words represented as indices (L1)

→ Non sequential model. (L2)

Lecture 5 : Limitation 1

words \longrightarrow indices \longrightarrow D-dimensional representation
of words

$$\mathcal{L} = \left[\begin{array}{c} \longleftarrow w_1 \longrightarrow \\ \vdots \\ \longleftarrow w_N \longrightarrow \end{array} \right]_{V \times D}$$

Glove recap

Dataloader

docs

.

,

,

docs_n



X

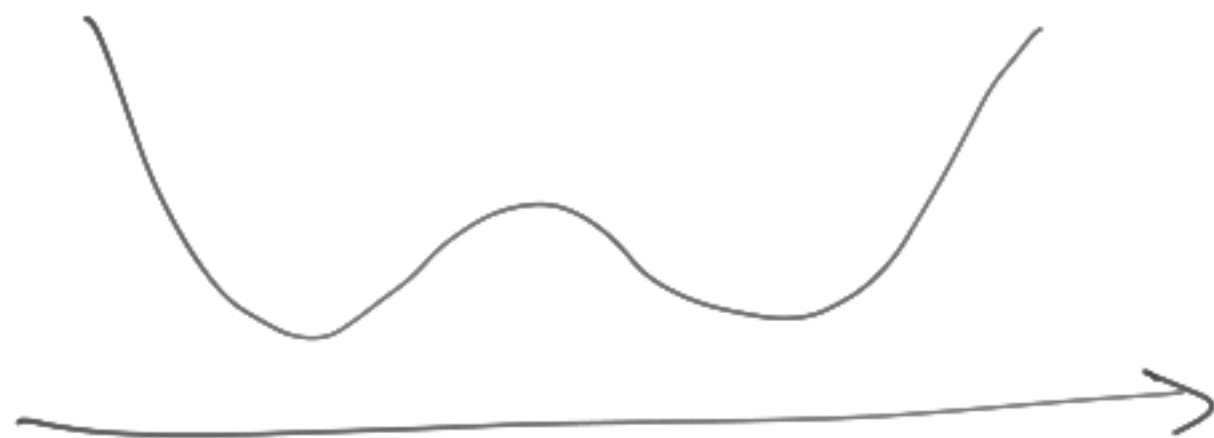
co-occurrence

matrix

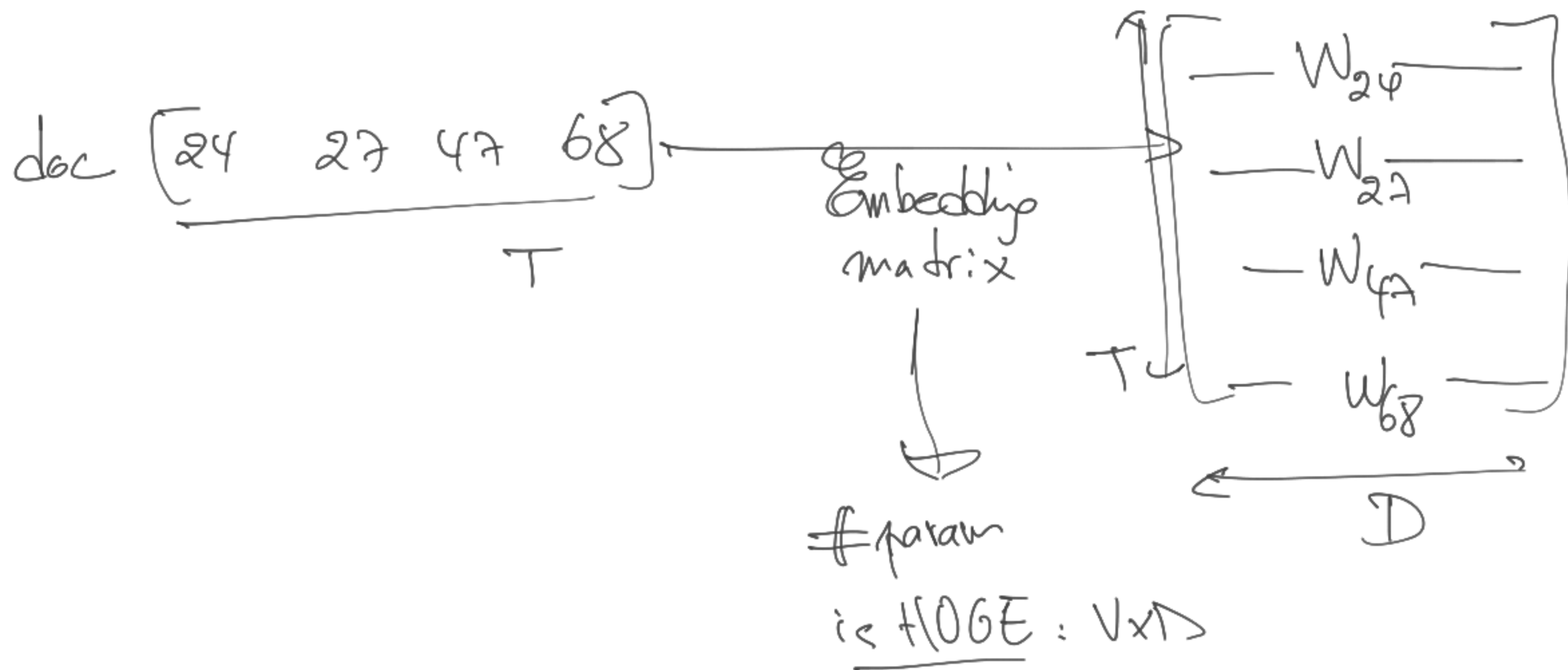


$$\log x_{ij} \sim w_i^T \tilde{w}_j$$

$$J = \sum_{i,j} (\log x_{ij} - w_i^T \tilde{w}_j)^2$$

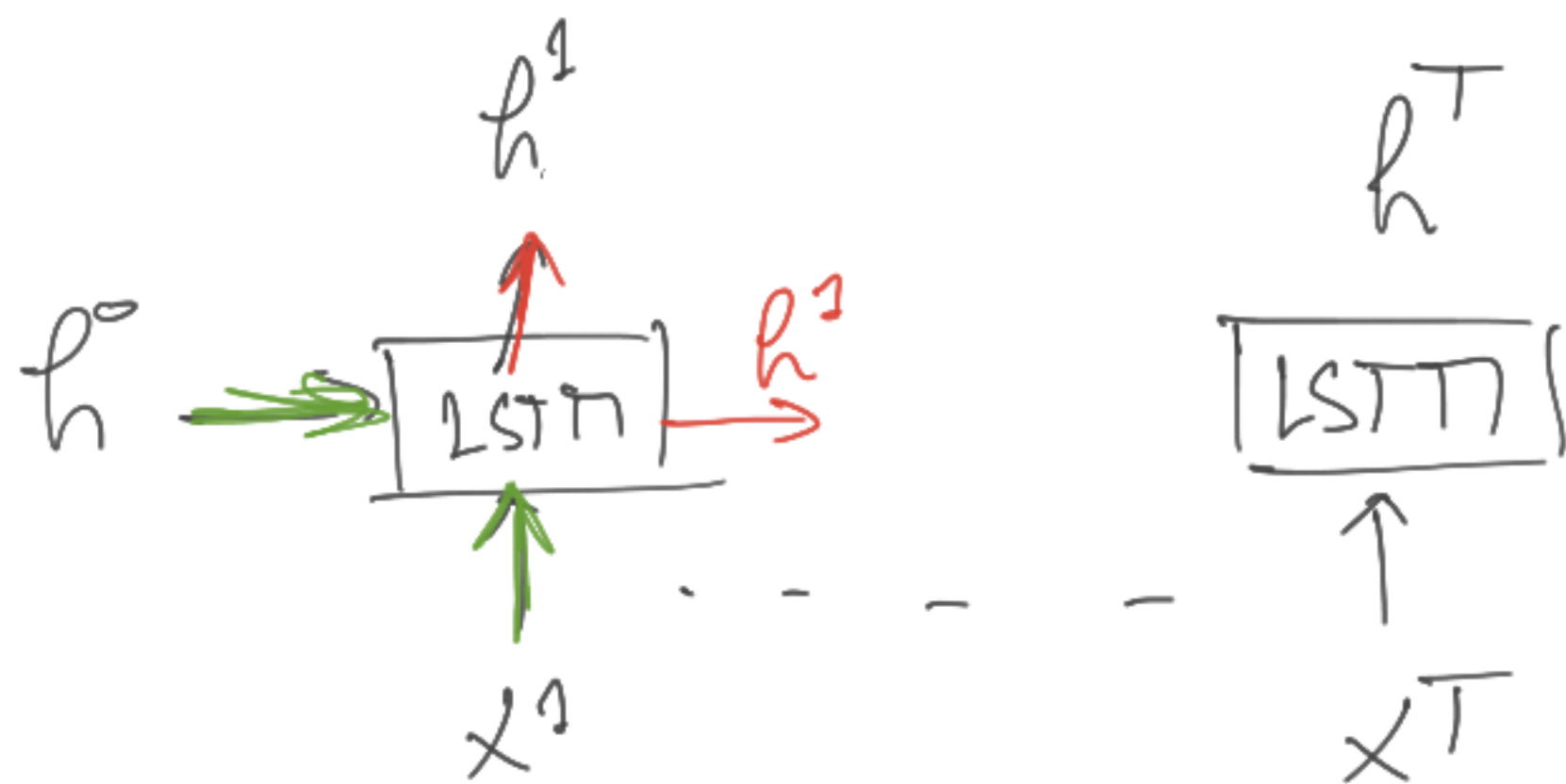


Embedding Layer (26)

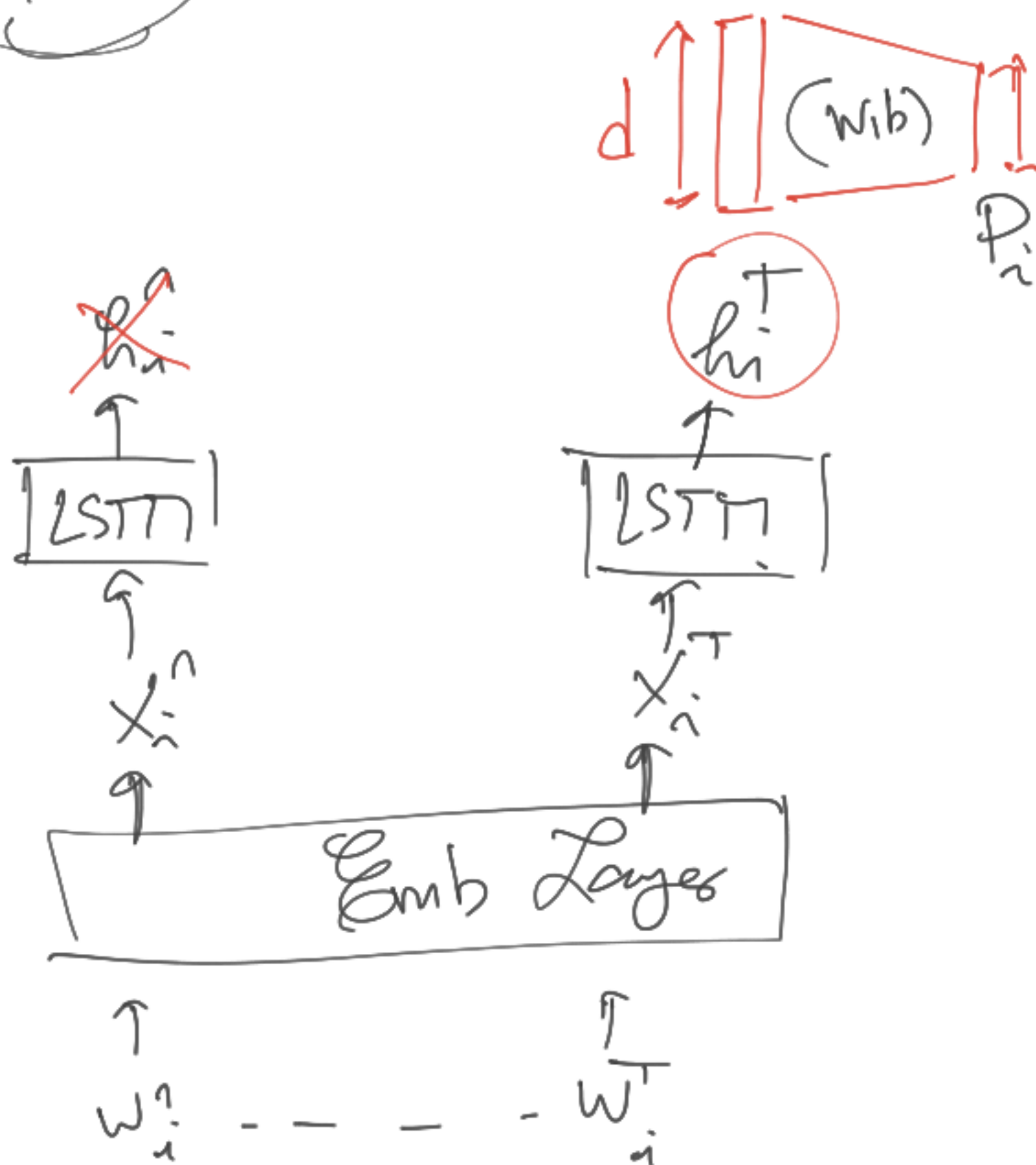


LSTM model

$x^1 \dots x^T \Rightarrow h^1 \dots h^T$



PSA



$K=5$ categories

$$J = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K (-y_i^k \log P_i^k)$$

Lecture 8

Many to One

→ Many to many

NLP

Tom a été entarté



Inter Say. Tom was hit with a pie

Finance

y^{T+1} y^{T+2} ... y^{T+H}



x^1 . . . x^T

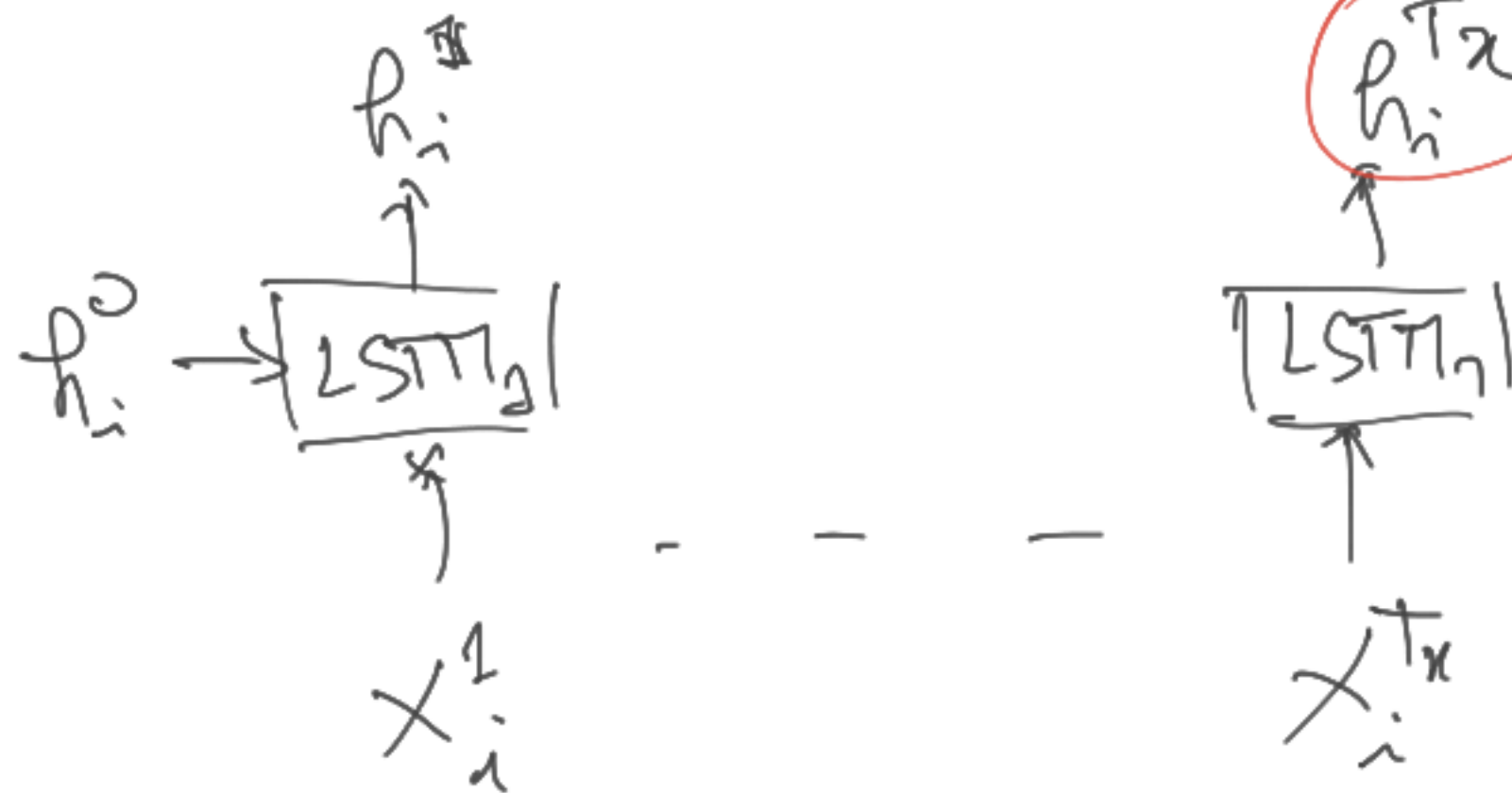
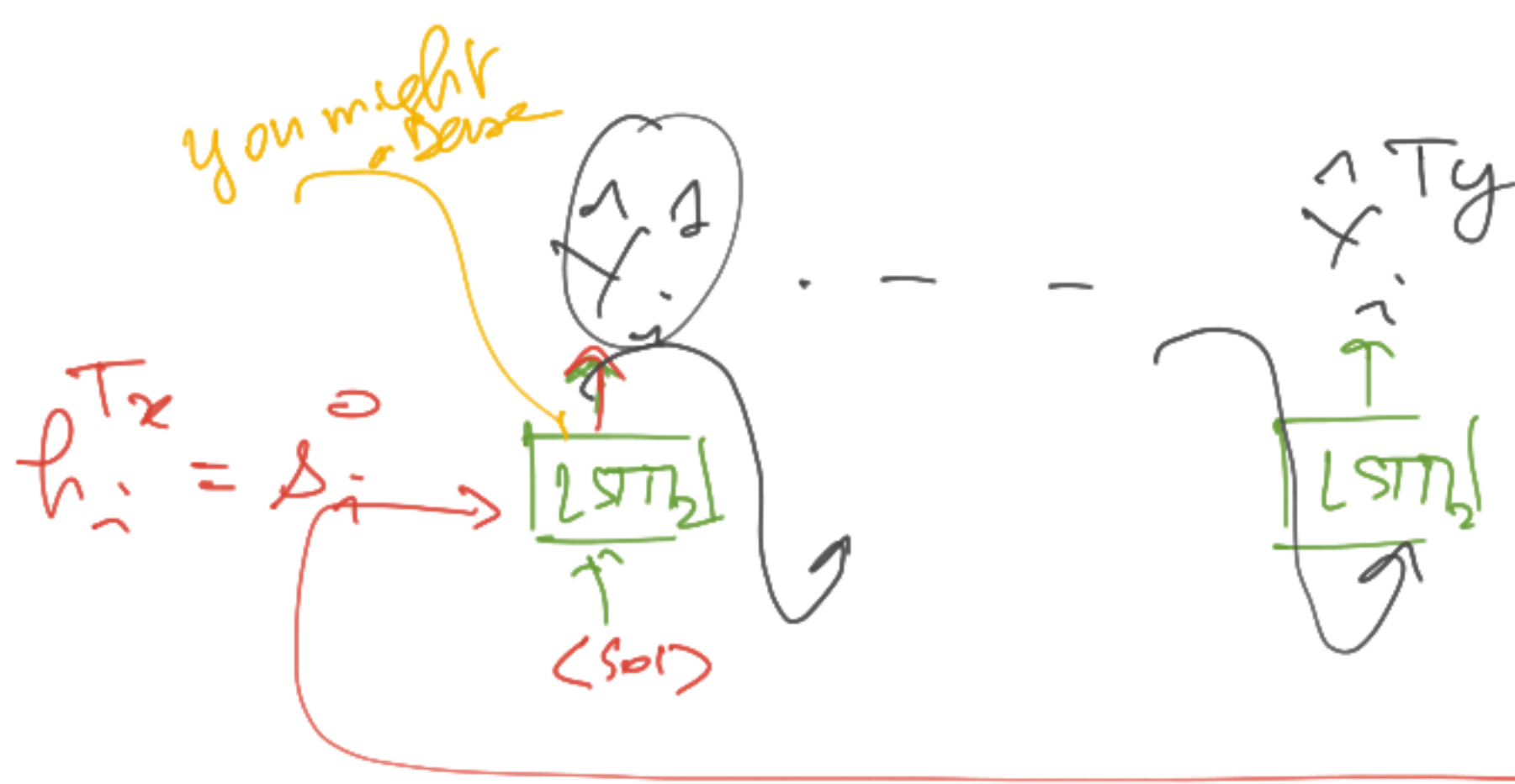


Diagram illustrating two horizontal lines. The top line has an arrow pointing right labeled T_y . The bottom line has an arrow pointing right labeled T_x .

Equation for T_y :

$$\forall t_y \quad \hat{y}_i^{t_y} \in \mathbb{R}^k$$

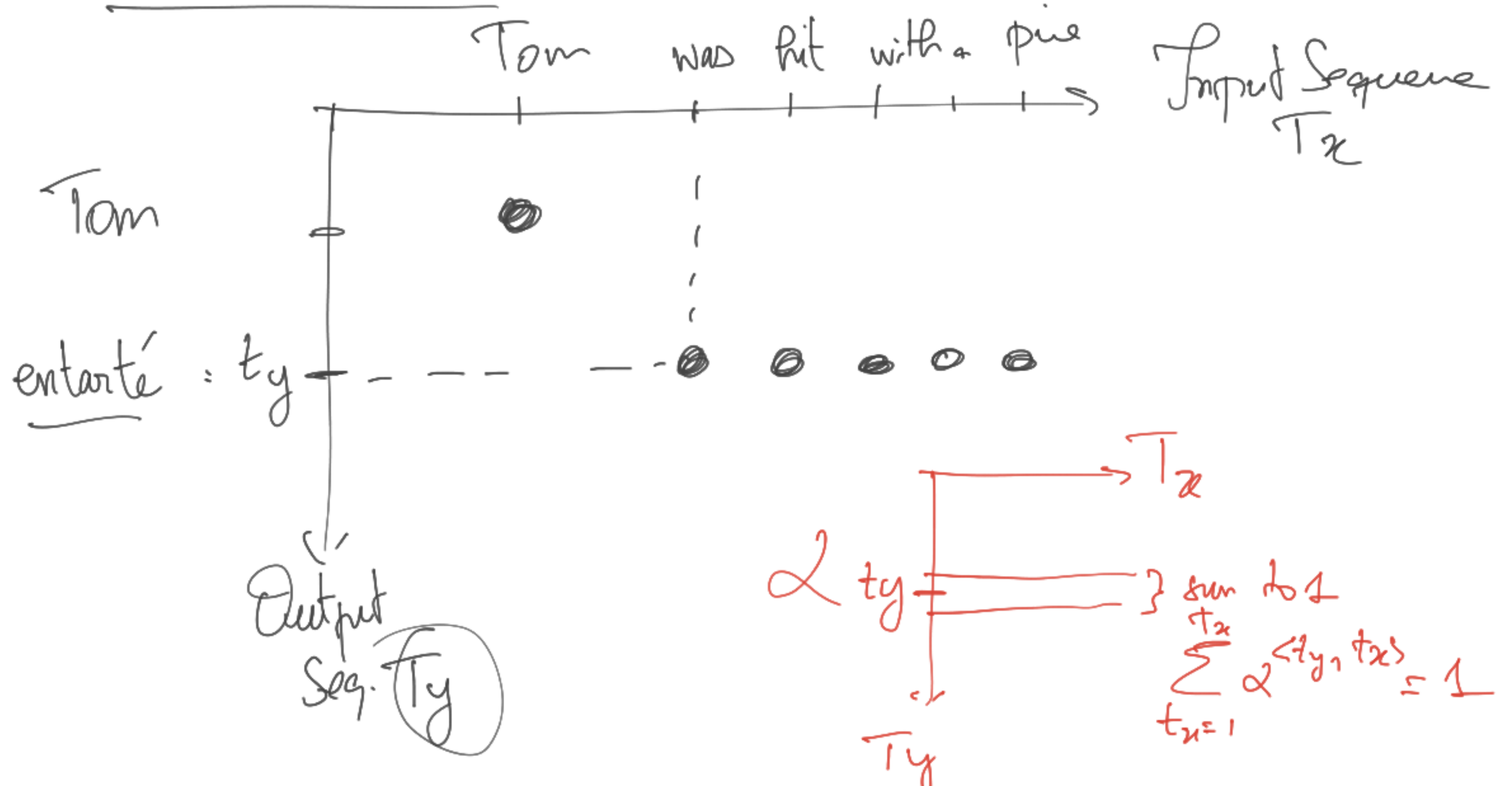
Equation for T_x :

$$\begin{bmatrix} 0 & 1 & 0 \end{bmatrix} = \gamma_i^{t_y} \in \mathbb{R}^k$$

Equation for the sum:

$$\sum_{t_y=1}^{T_y} \sum_{k=1}^K \left(-\gamma_i^{t_y}(k) \log \hat{\gamma}_i^{t_y}(k) \right)$$

Alignment ph



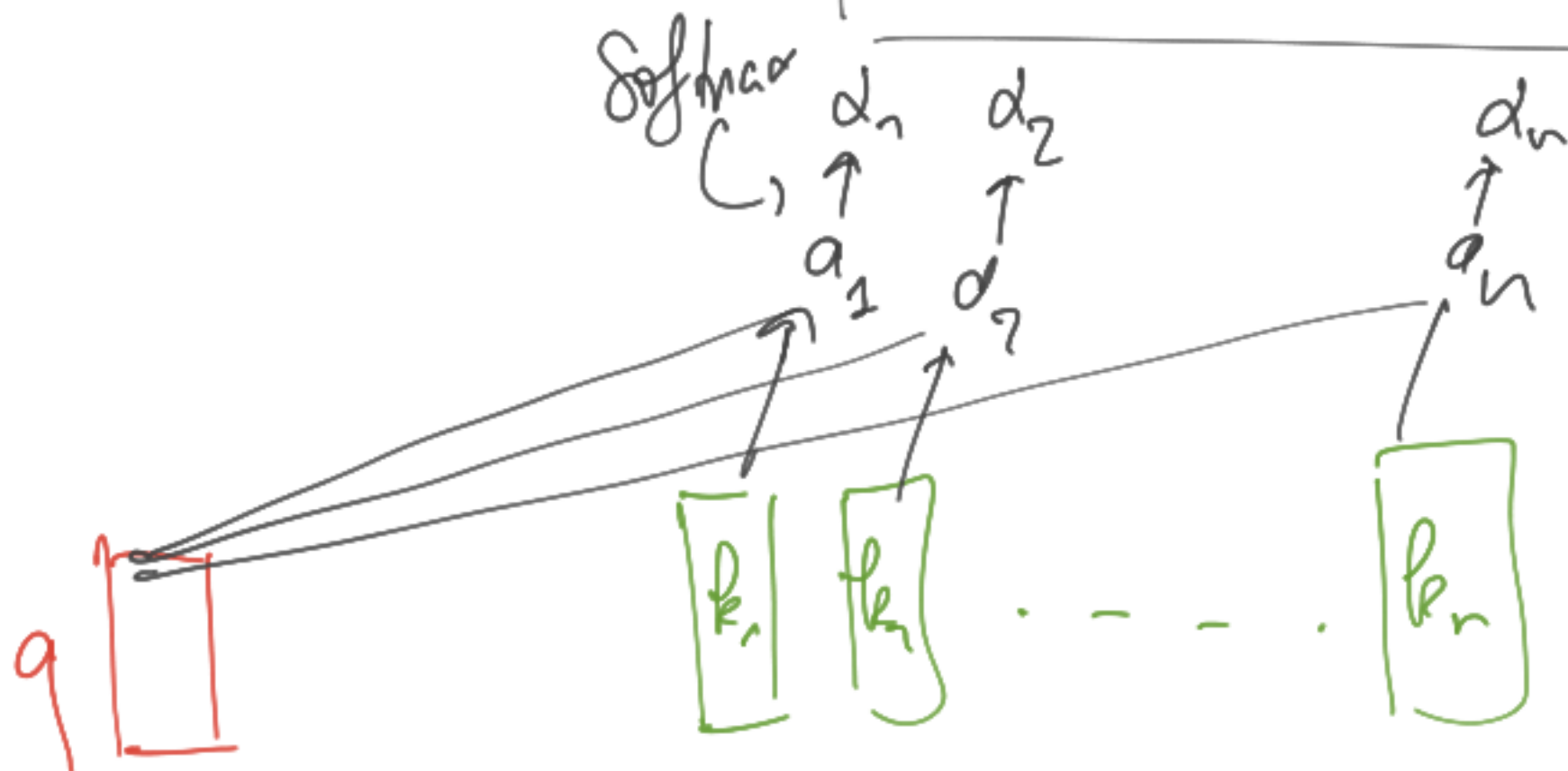
Attention ph

Intuition

YouTube database

title k_i

video v_i

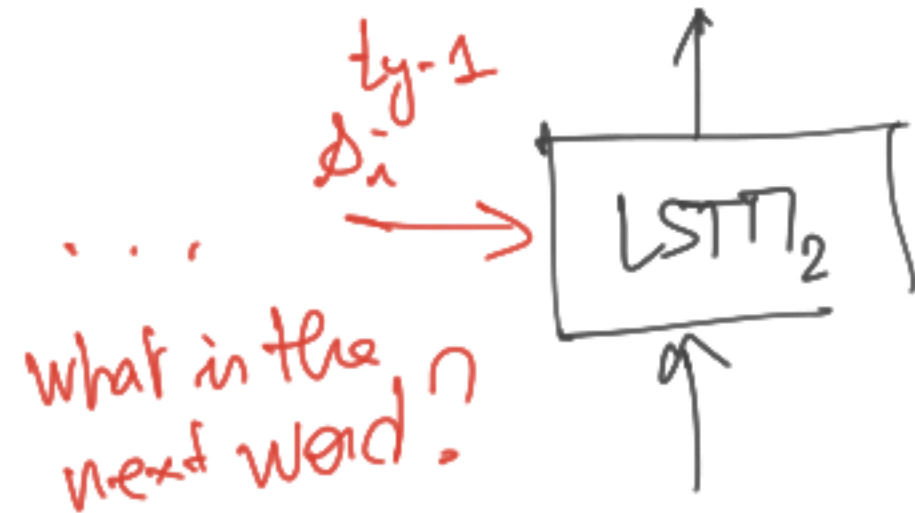


$$A(q, \{k_i, v_i\}_{1 \leq i \leq n})$$

$$\sum_{i=1}^n \alpha_i v_i$$

1st application of Attention (Not in exam)

objective: put attention on "hit with a fire".



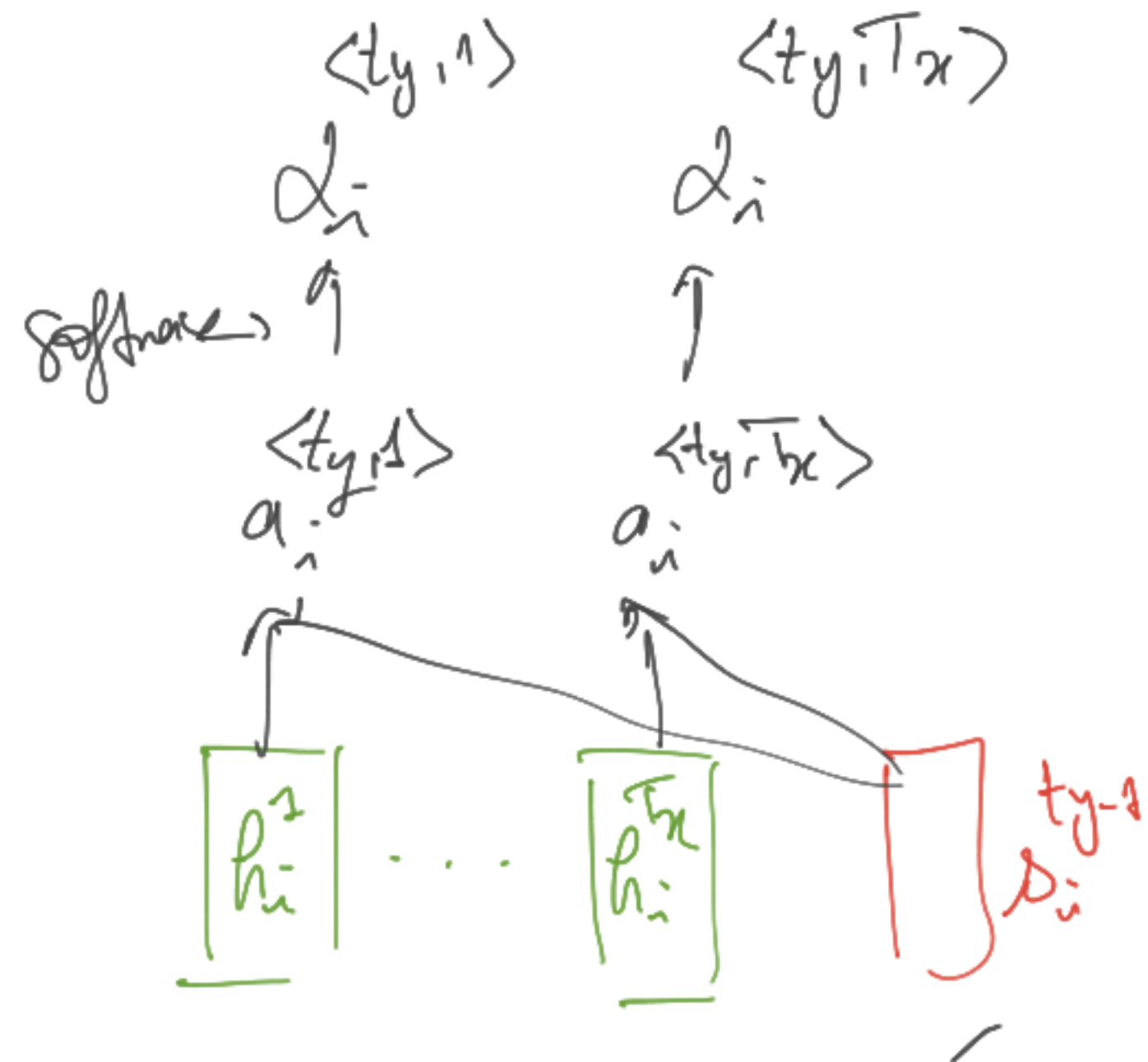
$$\Delta(\delta_i^{t_{y-1}}, \{h_i^{t_x}, h_i^{t_x}\}_{t_x \leq T_x})$$

$h_i^{t_1}$

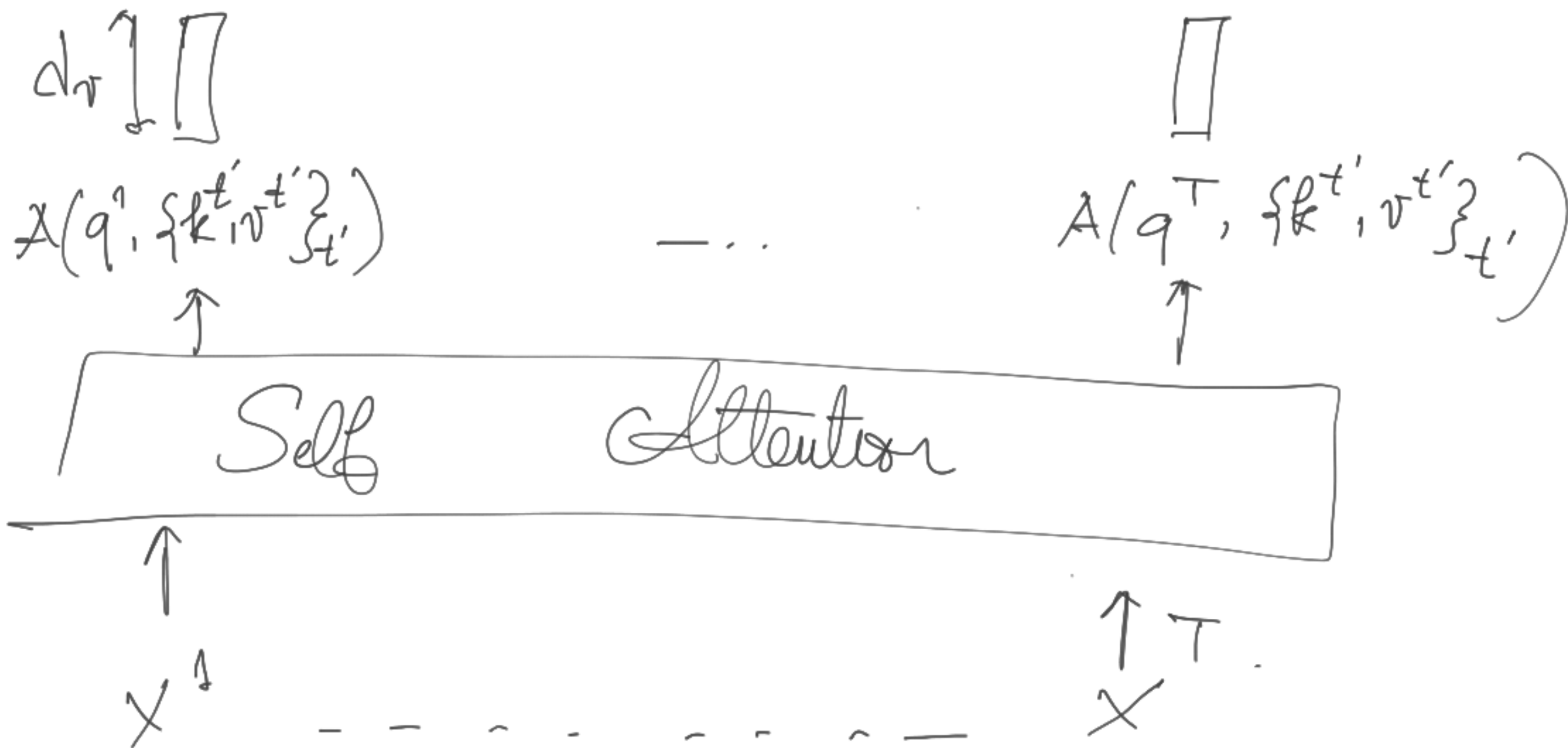
$x_{i, \text{Tom}}$

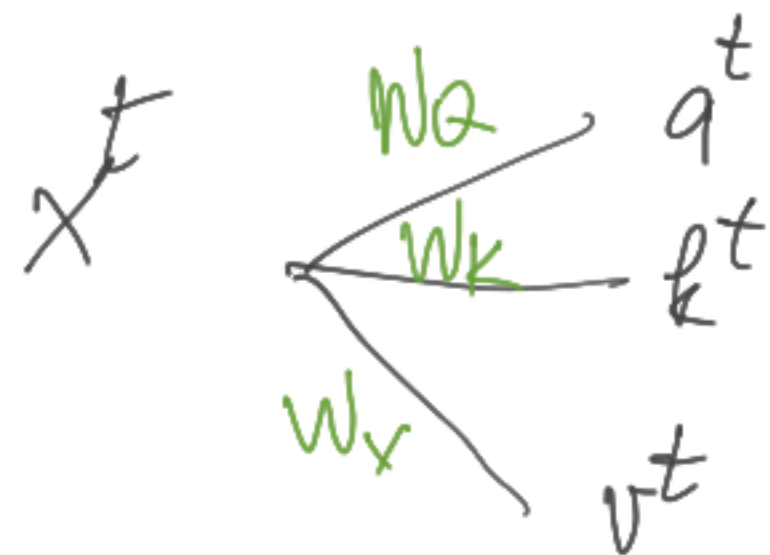
$h_i^{T_x}$

$x_{i, \text{pic}}^{T_x}$

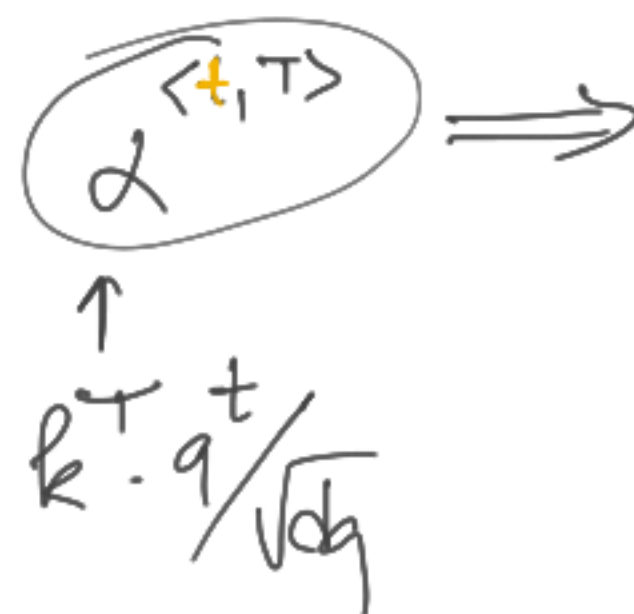
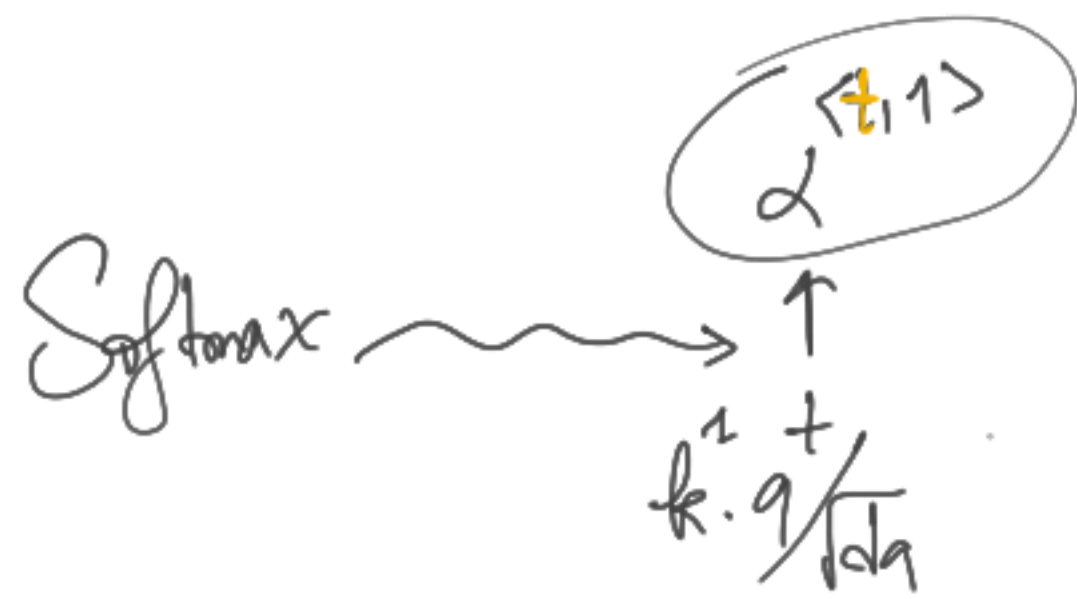


2nd Application of Attention: Contextual embedding vectors:





for all $t \in [1, T]$

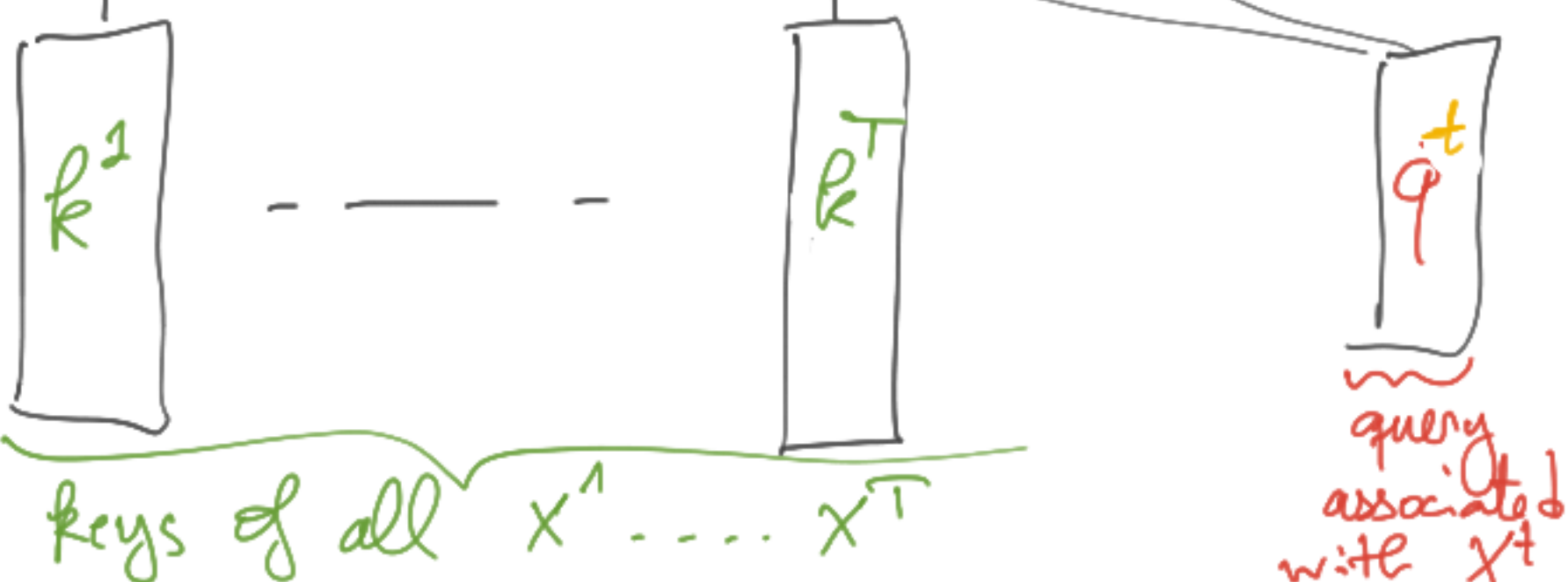


$$A(q^+, \{k^{t'}, v^{t'}\}_{t' \in [1, T]})$$

$$\sum_{t'=1}^T \alpha^{<t, t'>} v^{t'}$$

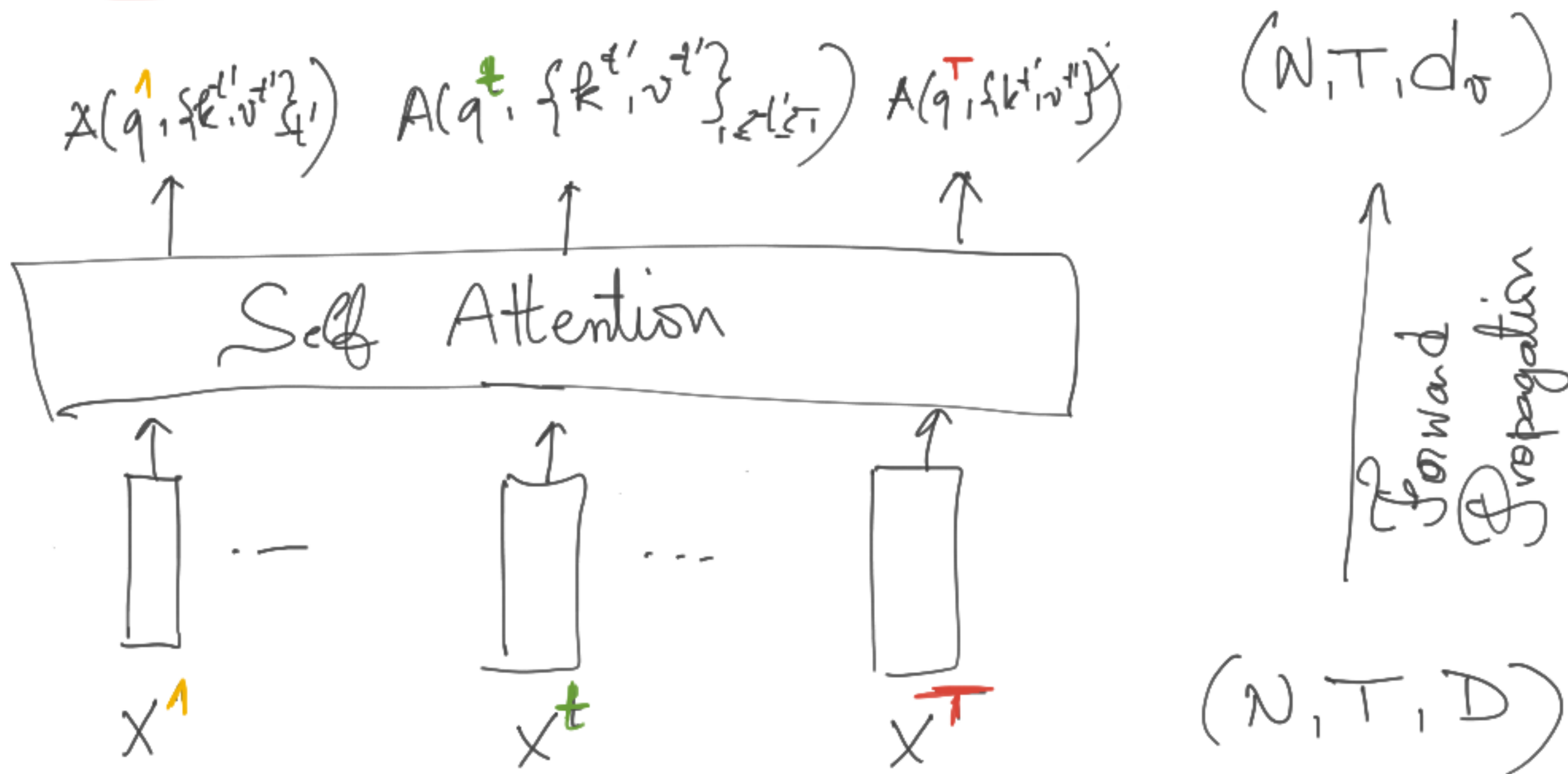
values associated with $x^{t'}$

Contextual representation of x^t



The Self Attention Layer

Parameters: $W_Q \in \mathbb{R}^{D \times d_q}$; $W_K \in \mathbb{R}^{D \times d_k}$, $W_V \in \mathbb{R}^{D \times d_v}$



Optional Session on Monday, May 15th

- Coding the Self Attention Layer
- Introduce the Transformer architecture and large language models
- Correct part of exam 2022.