

Outline :

- ① Applications of S2SVA (Many to many).
- ② Introduce Self Attention in the S2SVA.
- ③ Introduce : Attention is all you need.
- ④ Correcting part of Exam 2022.

1 Applications of S2S WA:

NMT

Output Sequence: Tom a été entarté $\langle \text{eos} \rangle$

s_i^{ty}

h_i^1 h_i^2 h_i^3 h_i^4 $h_i^{T_x-1}$ $h_i^{T_x}$

Input Sequence: Tom was put with a pie

Time Series Predictions

x_i^{T+1} \hat{y}_i^{T+2} \hat{y}_i^{T+H}

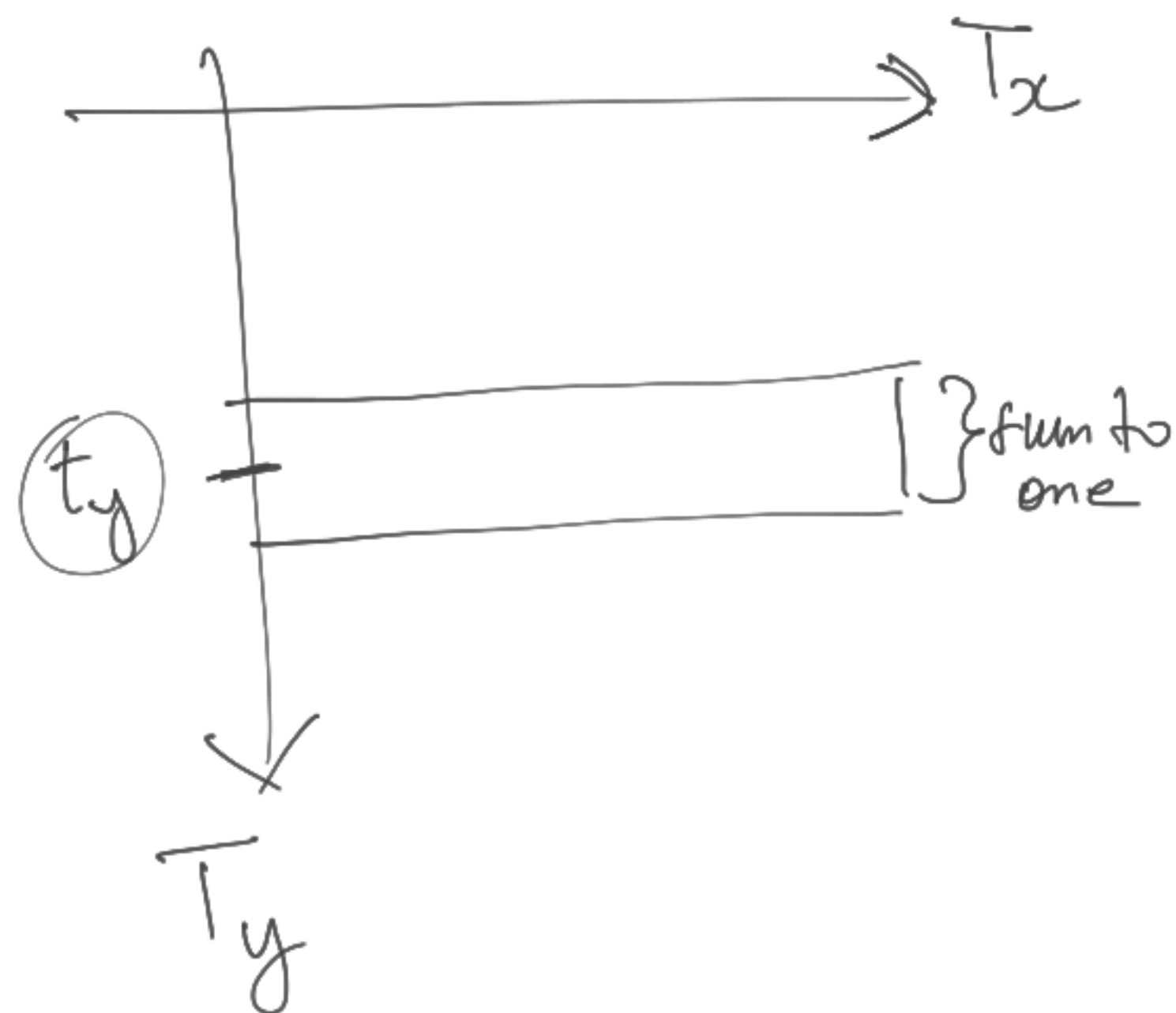
s_i^{ty}

x_i^1 x_i^2 x_i^T

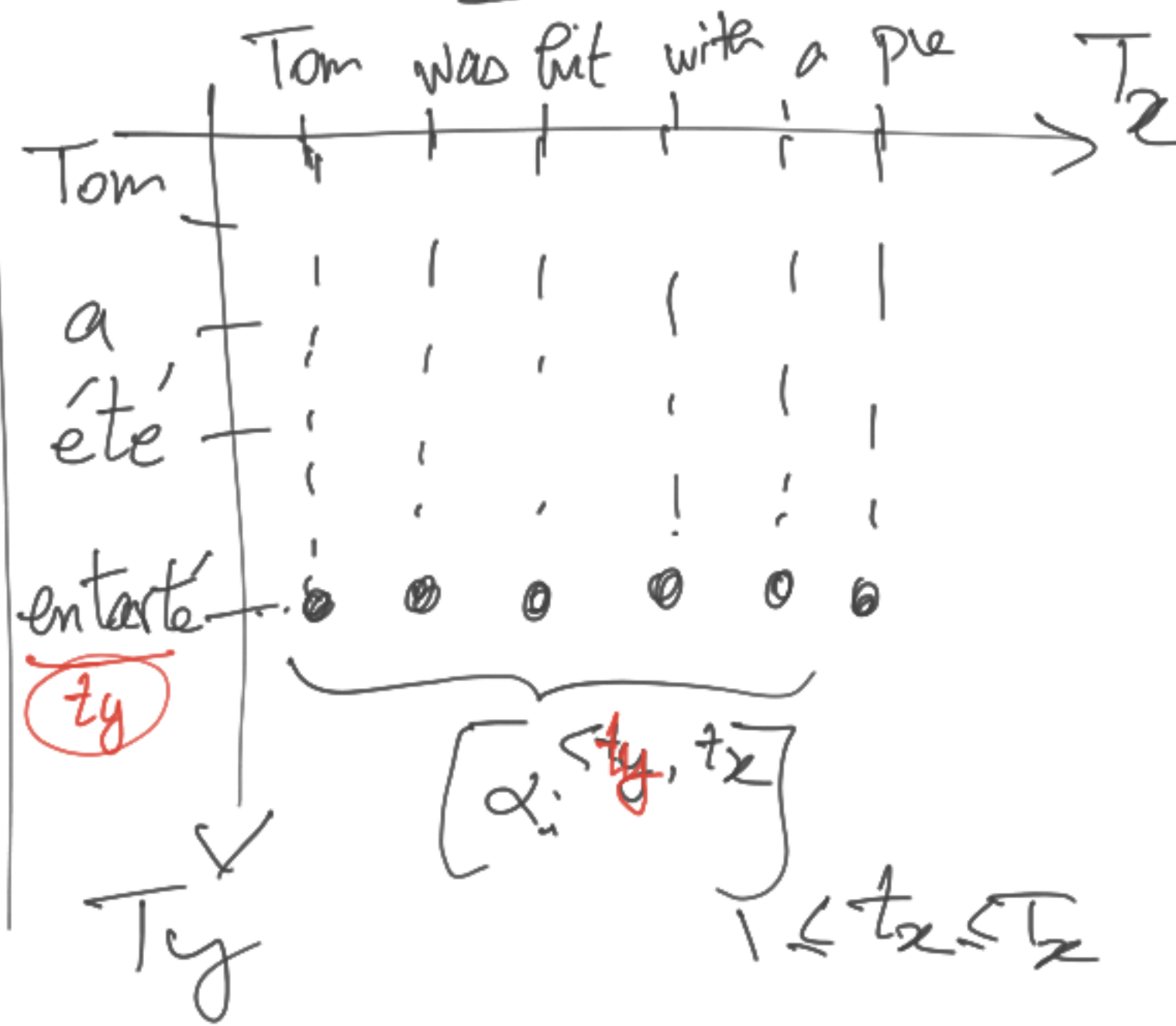
Alignment pb

Objective : $x_i^1 \dots x_i^{T_x}$

$\rightarrow \hat{x}_i^1 \dots \hat{x}_i^{T_y}$

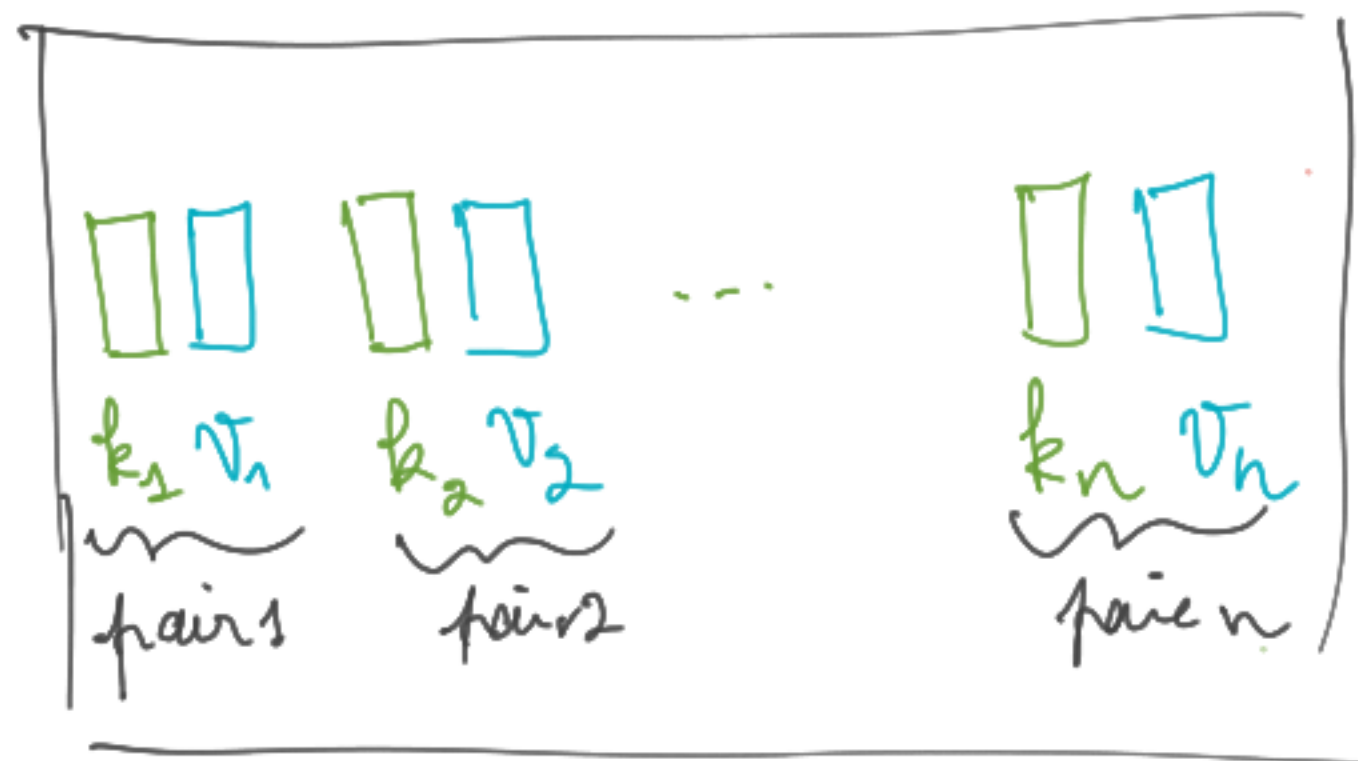


example



Attention:

$$\boxed{} \quad q \in \mathbb{R}^d$$



Database

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

||

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

$$\frac{\exp(q \cdot \underline{k_i})}{\underbrace{\sum_{j=1}^n \exp(q \cdot \underline{k_j})}_{\mathcal{Z}_i}} \quad \underline{v_i}$$



We can use other similarity measures.

Self Attention:

Objective: $x_1^1 \dots x_i^T$ (they represent embedding vectors in NLP)

↳ Create a sequence of contextual embedding vectors.

ex (28) : Tom a été entarté. (Tom was hit with a pie)

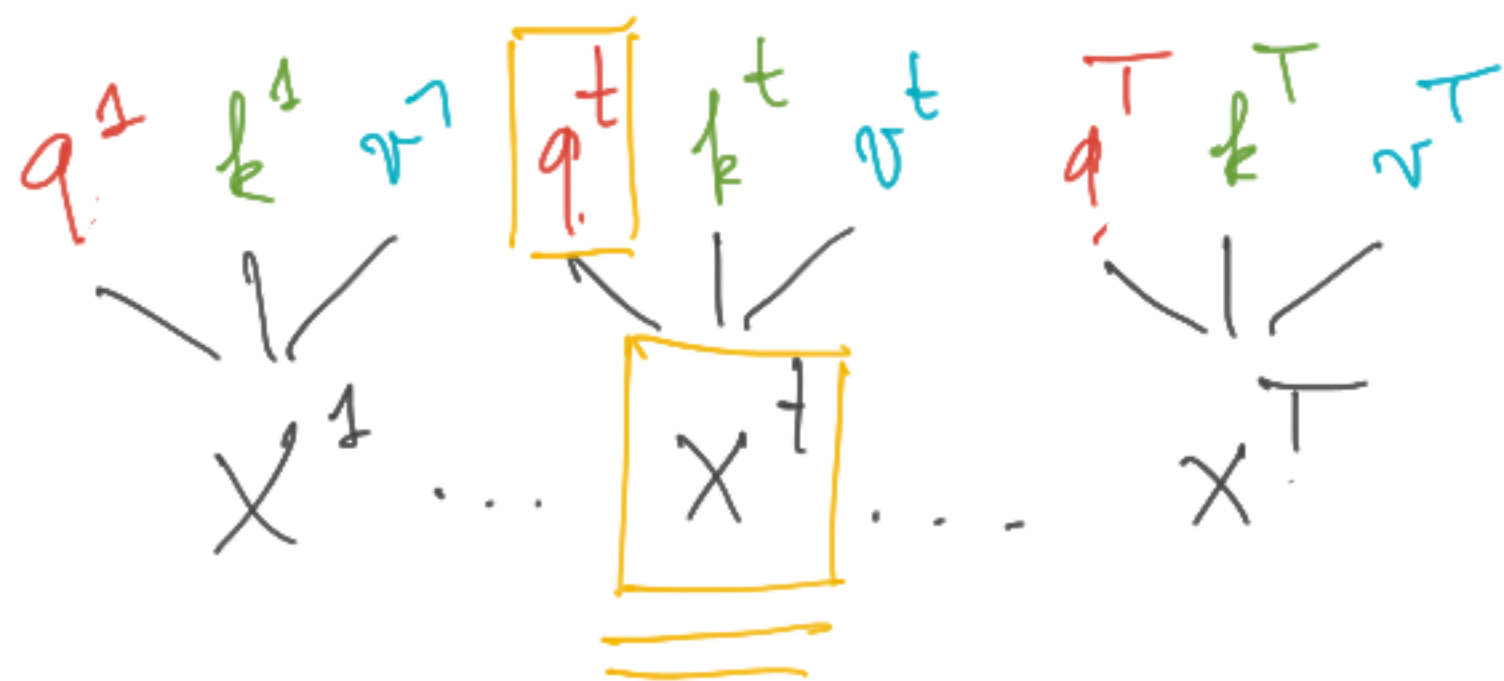
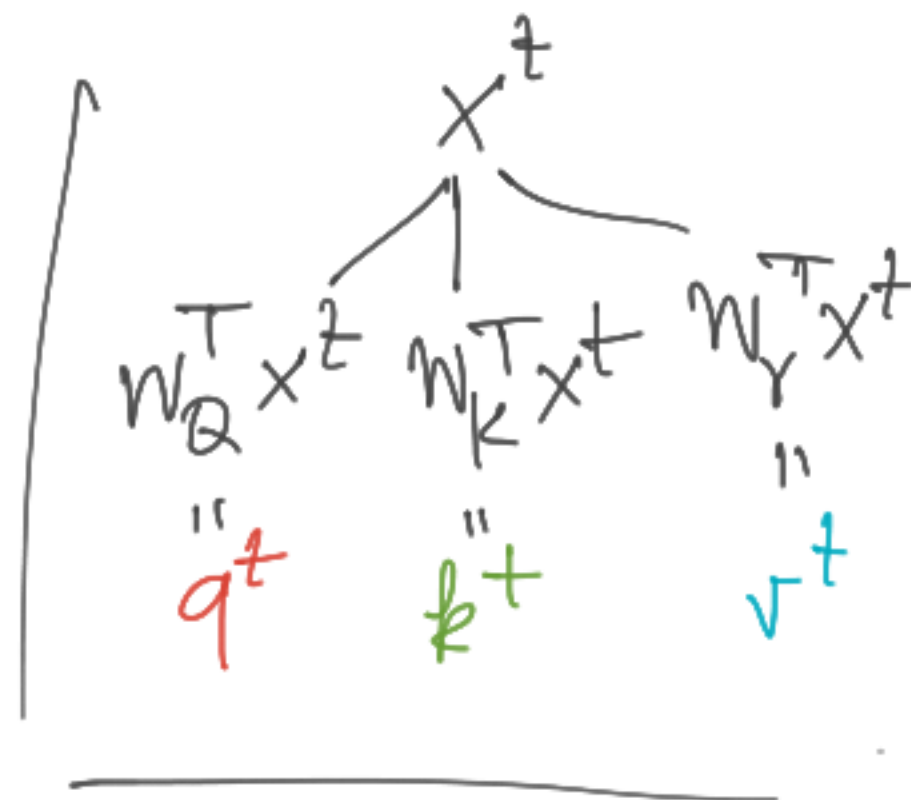
↳ Let été was il fera horriblement chaud.
Summer

Self Attention Layer

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

Contextual representation of x^t :

$$A(q^t, \{k^{t'}, v^{t'}\}_{t' \leq t}) = \sum_{t'=1}^T \frac{\exp\left(\frac{q^t \cdot k^{t'}}{\sqrt{d_k}}\right)}{\sum_{t''=1}^T \exp\left(\frac{q^t \cdot k^{t''}}{\sqrt{d_k}}\right)} v^{t'}$$

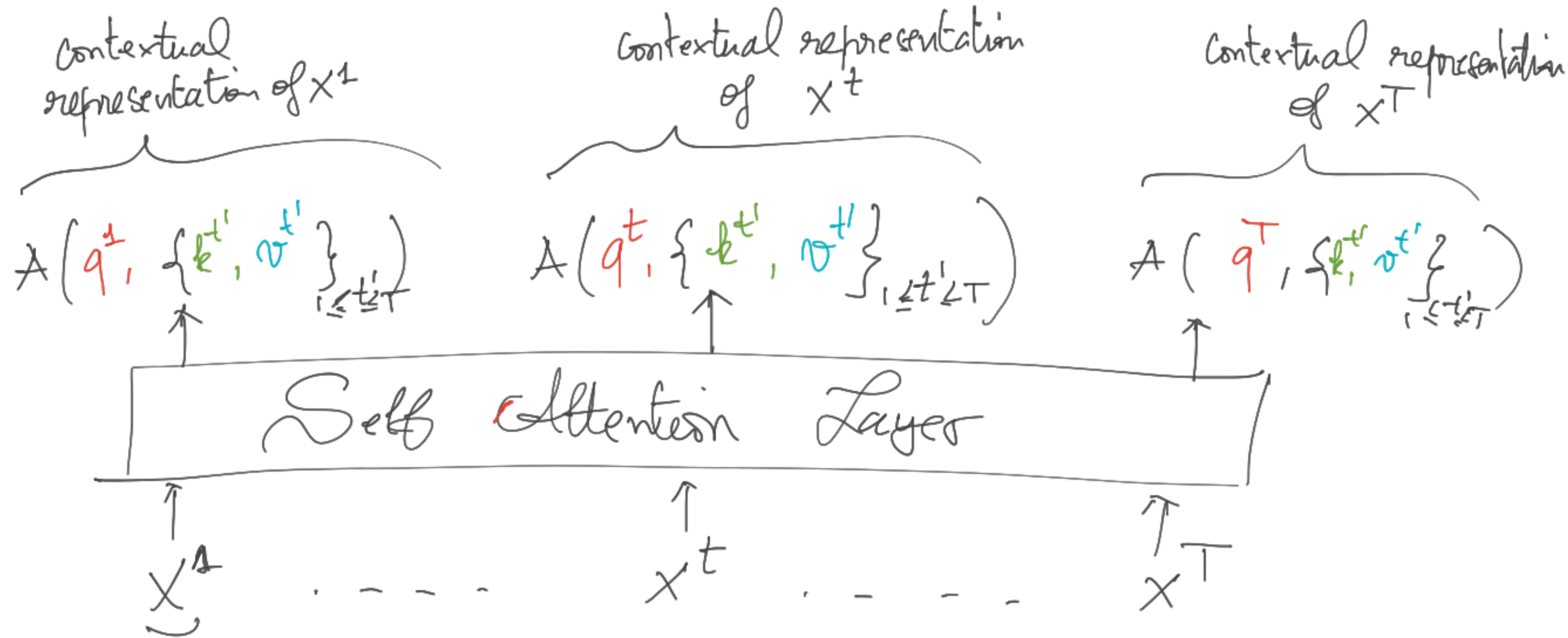


$$\alpha_{\langle t, t' \rangle}$$

: Impact of $x^{t'}$ in generation of the contextual representation of x^t .

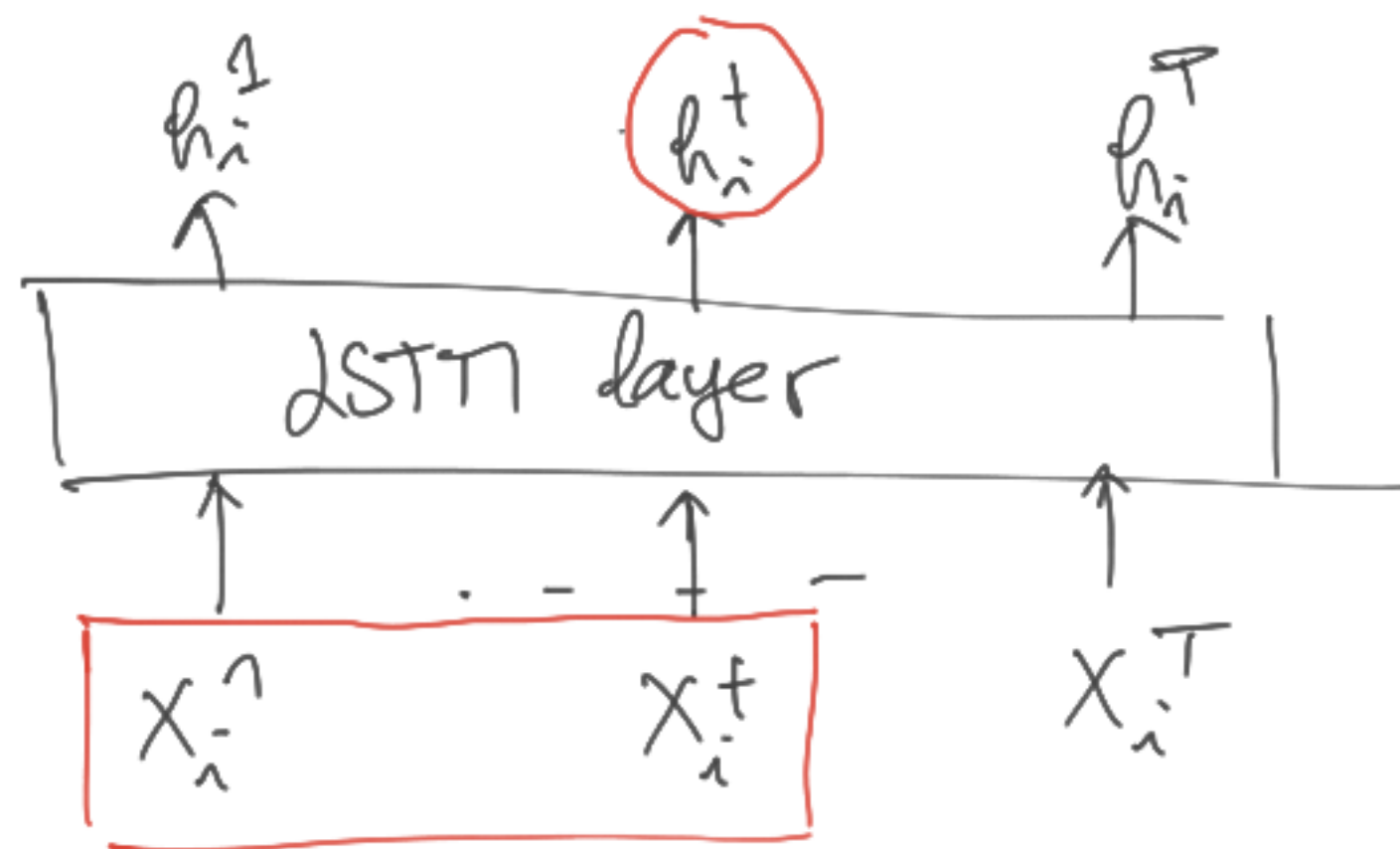
Self Attention Layer

Forward Propagation
Self Attention.
 $(N, T, D) \rightarrow (N, T, d_v)$



LSTM Layer

$$h_i^t = f(x_i^1, \dots, x_i^t)$$



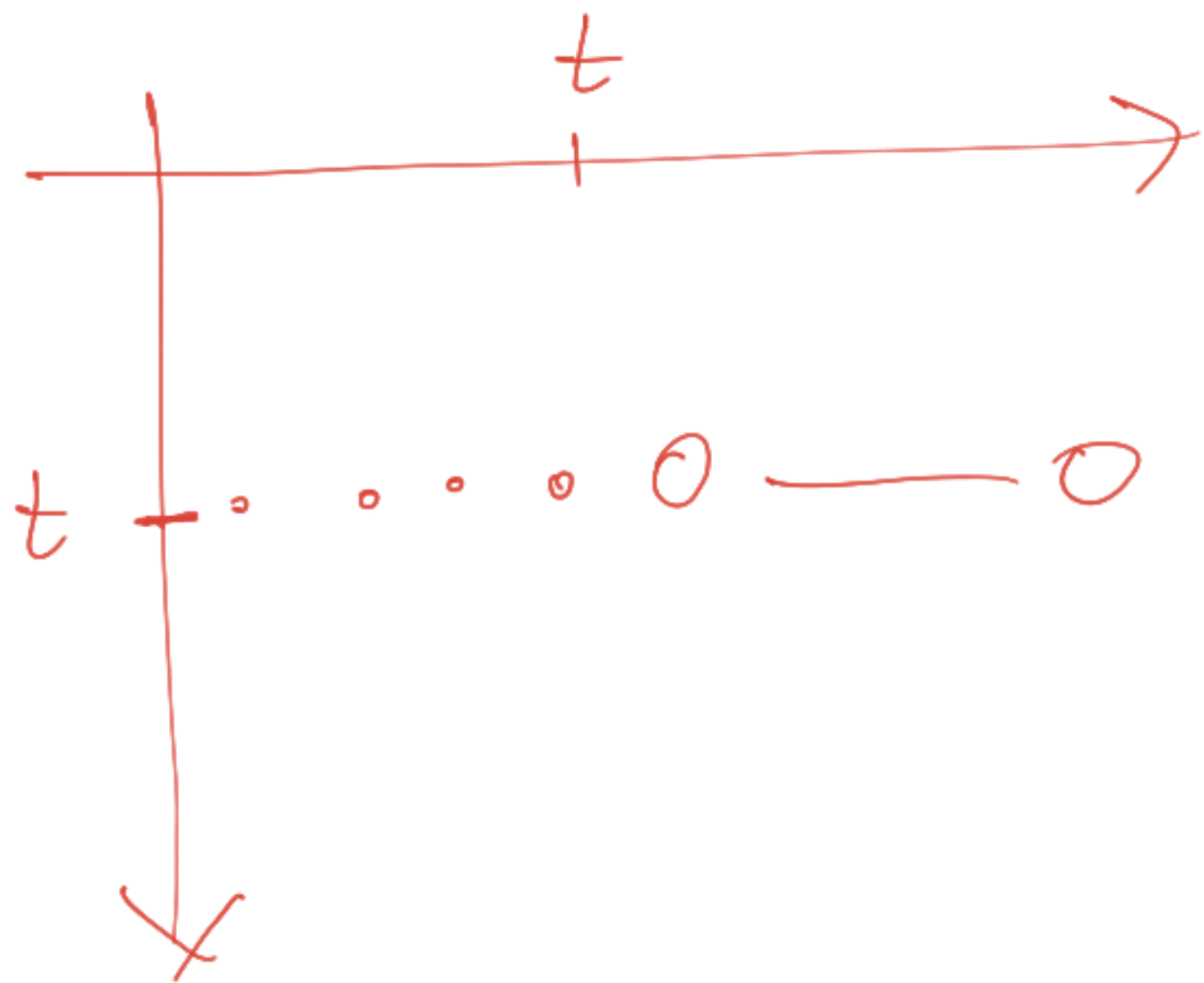
Self Attention Layer

$$A_i^t = f(x_i^1, \dots, x_i^T)$$

$$A_i^t = \sum_{t'=1}^T \underbrace{\alpha_i^{(t,t')}}_{\text{weight}} v_i^{t'}$$



→ Causal Self Attention Layer
(Masked Attention).



$$A(q^t, \{k^{t'}, v^{t'}\}_{t' \leq t})$$

$$\sum_{t'=1}^t \underbrace{\alpha^{<t, t'>}}_{\text{wavy line}} v^{t'}$$

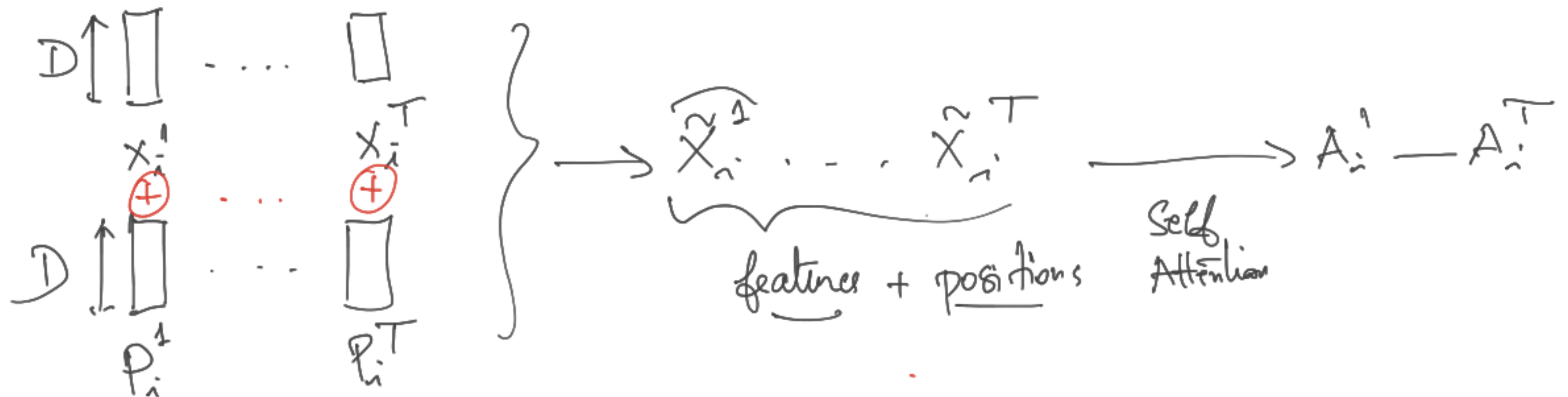
$$\underbrace{\forall t' \leq t \quad \alpha^{<t, t'>} = 0}_{\text{bracketed}} \quad \text{wavy line}$$

$$A(q^t, \{k^{t'}, v^{t'}\}_{t' \leq t}) = f(x^1 \dots x^t)$$

Positional Encoding:

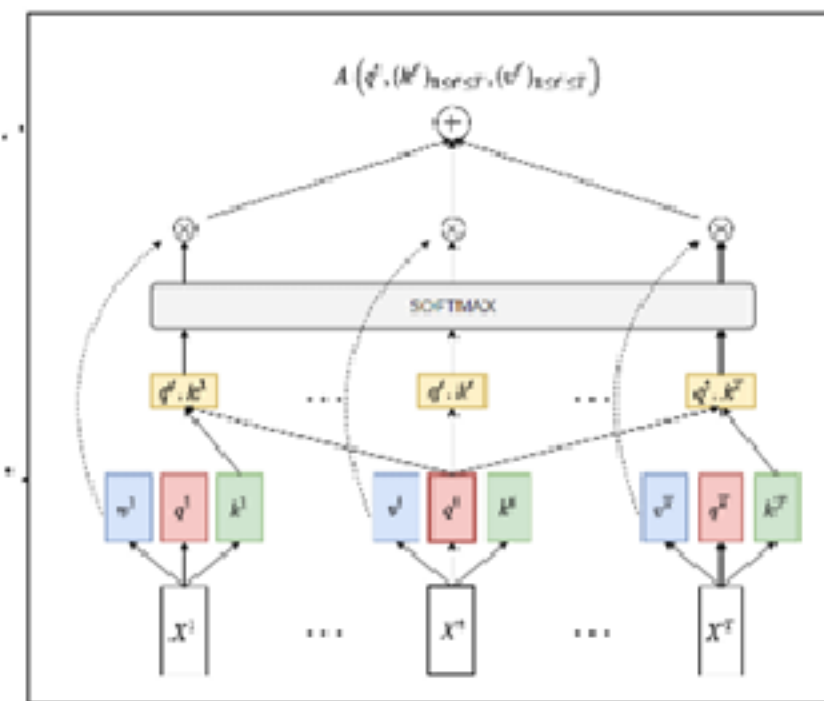
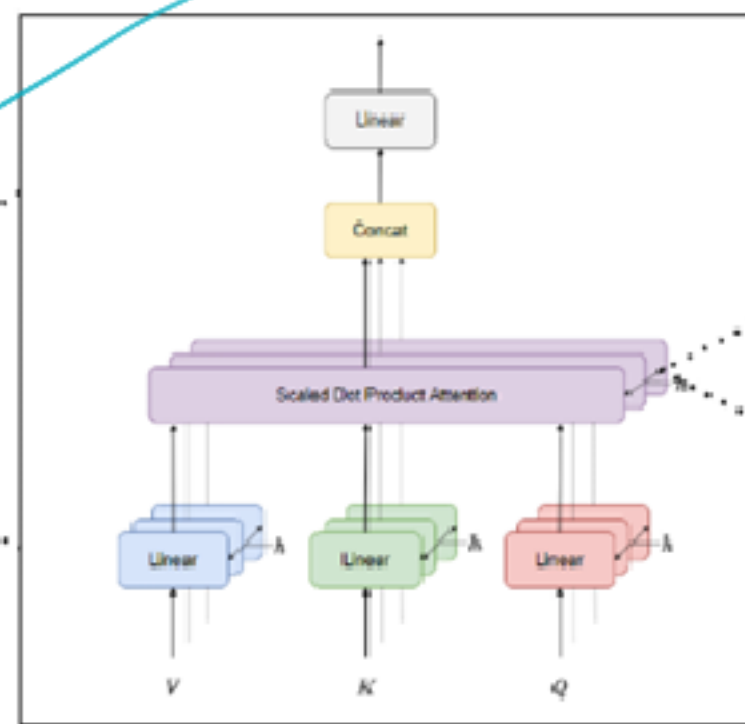
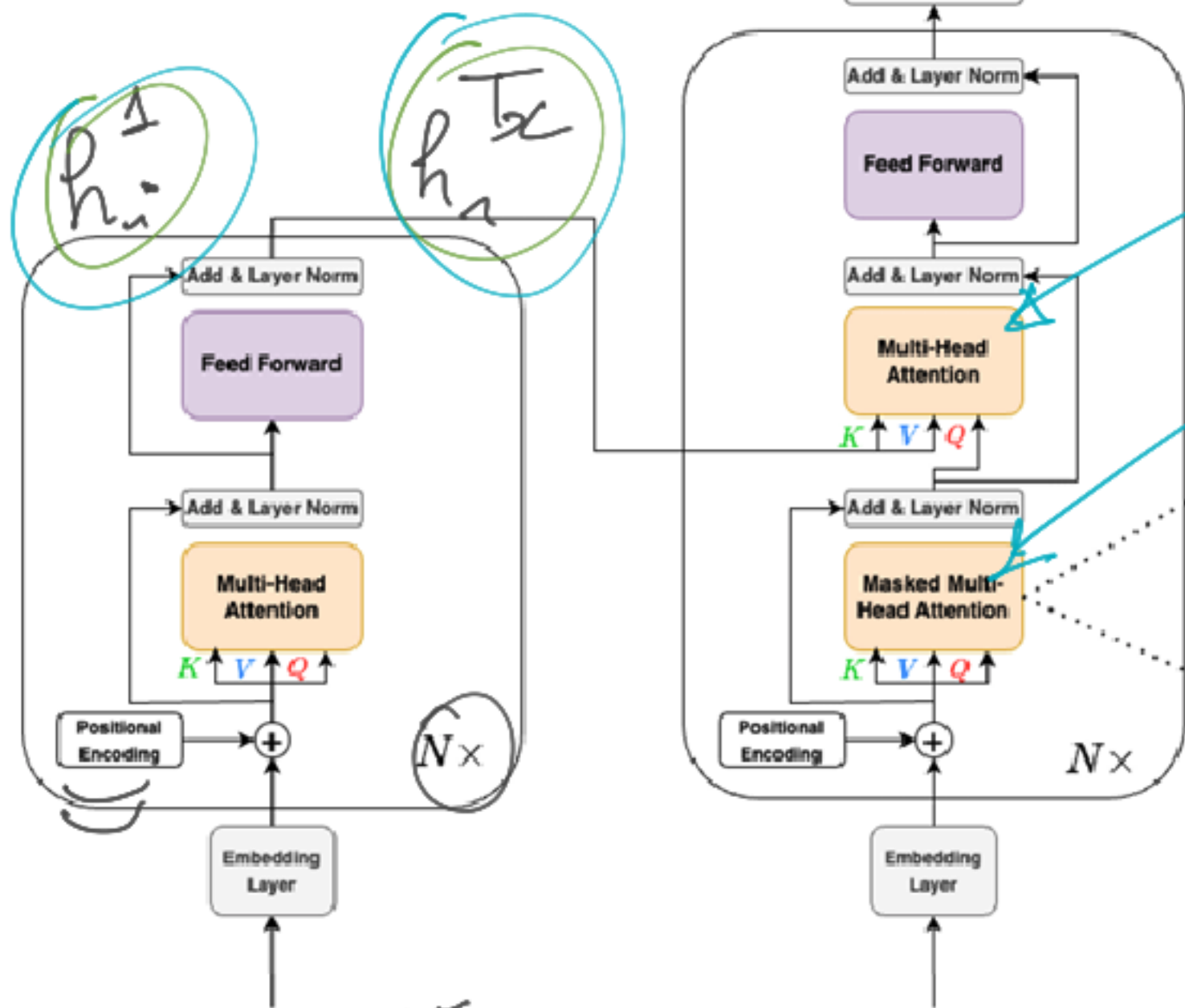
② Ph: $A(q^t, \{k^{t'}, v^{t'}\}_{1 \leq t' \leq T})$ doesn't take into consideration the order of $x_i^1 \dots x_i^T$.

③ Solⁿ:



objective Tom a été entanté <eos>) decoder
Output

Gross Attention Layer (same as before)
Causal Attention Layer
(Decoder \rightarrow query
Encoder \Rightarrow key, value)



$x_1^A \dots x_n^{Tx}$
Tom was hit with a pine

[<eos> Tom a été entanté] : Decoder Input

Exam 2022

Newsid	date	Stock Id	news	Sentiment
\downarrow	t	i		

Q10)
$$N = \sum_{t=1}^T \sum_{i=1}^{N_u} n_{t,i}$$

} $n_{t,i}$
news
related
to stock of
id i at
time t .

Q11

newsID	News
1	$w_1 = (w_1^1 \dots w_1^L)$
'	\vdots
'	
N	$w_N = (w_N^1 \dots w_N^L)$

→ describe the word_index to map words → integers
→ Padding to make all the sequences of the same size. (See Review 2 slide 2)

Q12

(N, L) tensor.

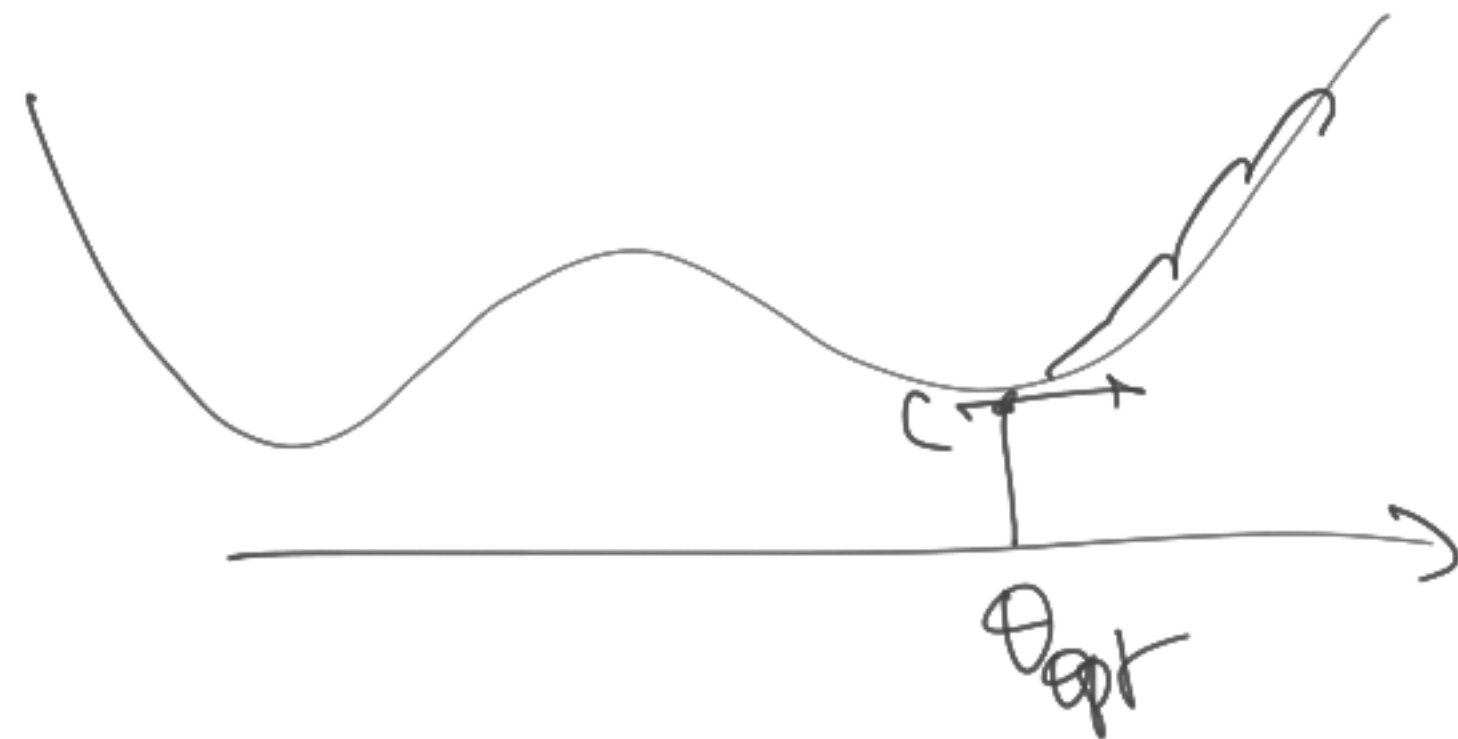
Q14



- How to gather the data to create X
- Explain the training of:

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

J ?



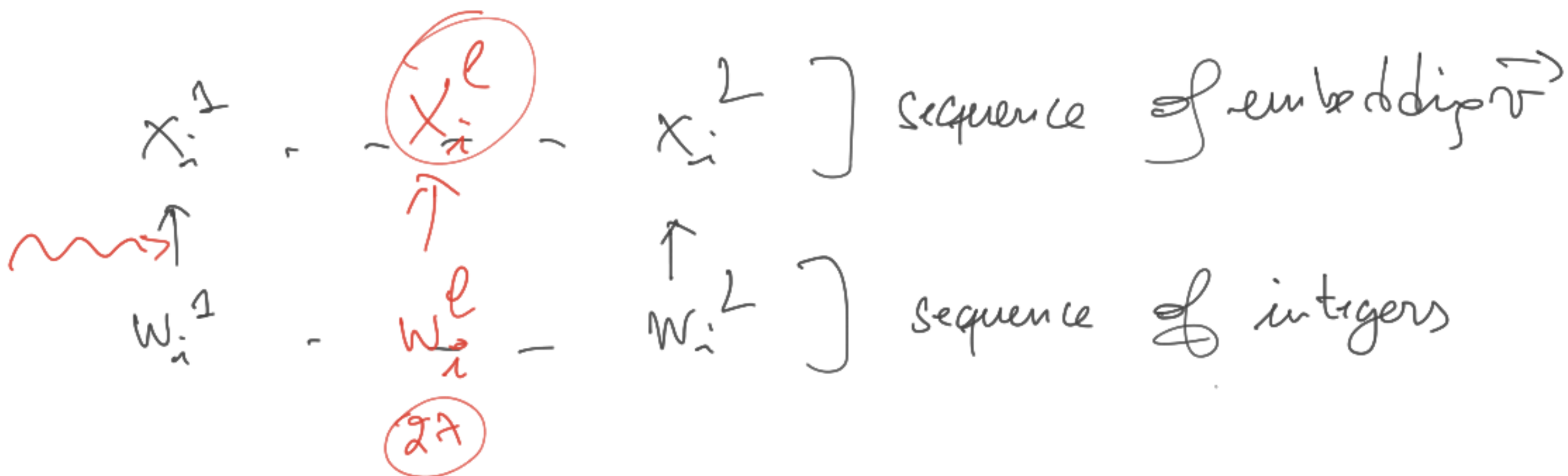
in θ_{opt} we have the embedding
matrix

$$E = \begin{bmatrix} \text{---} w_1 \text{---} \\ \vdots \\ \text{---} w_V \text{---} \end{bmatrix}$$

(Slide 9 Reviewed 2)

Q15 $W_e \in \mathbb{R}^{V \times D}$

Q16



$$x_i^l = W_e [w_i^l]$$

\uparrow
row # w_i^l in W_e .

Q17

$$\begin{bmatrix} \phi_i^{l-1} \\ c_i^{l-1} \\ x_i^l \end{bmatrix}$$

past memory \vec{v}

new obs

LSTM

$$(\phi_i^l, c_i^l)$$

$$\begin{bmatrix} W_F, b_F \\ W_I, b_I \\ W_O, b_O \\ W_C, b_C \end{bmatrix}$$

Gates

$$F_i^l = \sigma \left[W_F^T [\phi_i^{l-1}, x_i^l] + b_F \right]$$

$$I_i^l = \sigma \left[W_I^T [\phi_i^{l-1}, x_i^l] + b_I \right]$$

$$O_i^l = \sigma \left[W_O^T [\phi_i^{l-1}, x_i^l] + b_O \right]$$

Update equations

$$\tilde{c}_i^l = \tanh \left(W_C^T [\phi_i^{l-1}, x_i^l] + b_C \right)$$

$$c_i^l = F_i^l \odot c_i^{l-1} + I_i^l \odot \tilde{c}_i^l$$

$$\phi_i^l = O_i^l \odot \tanh(c_i^l)$$