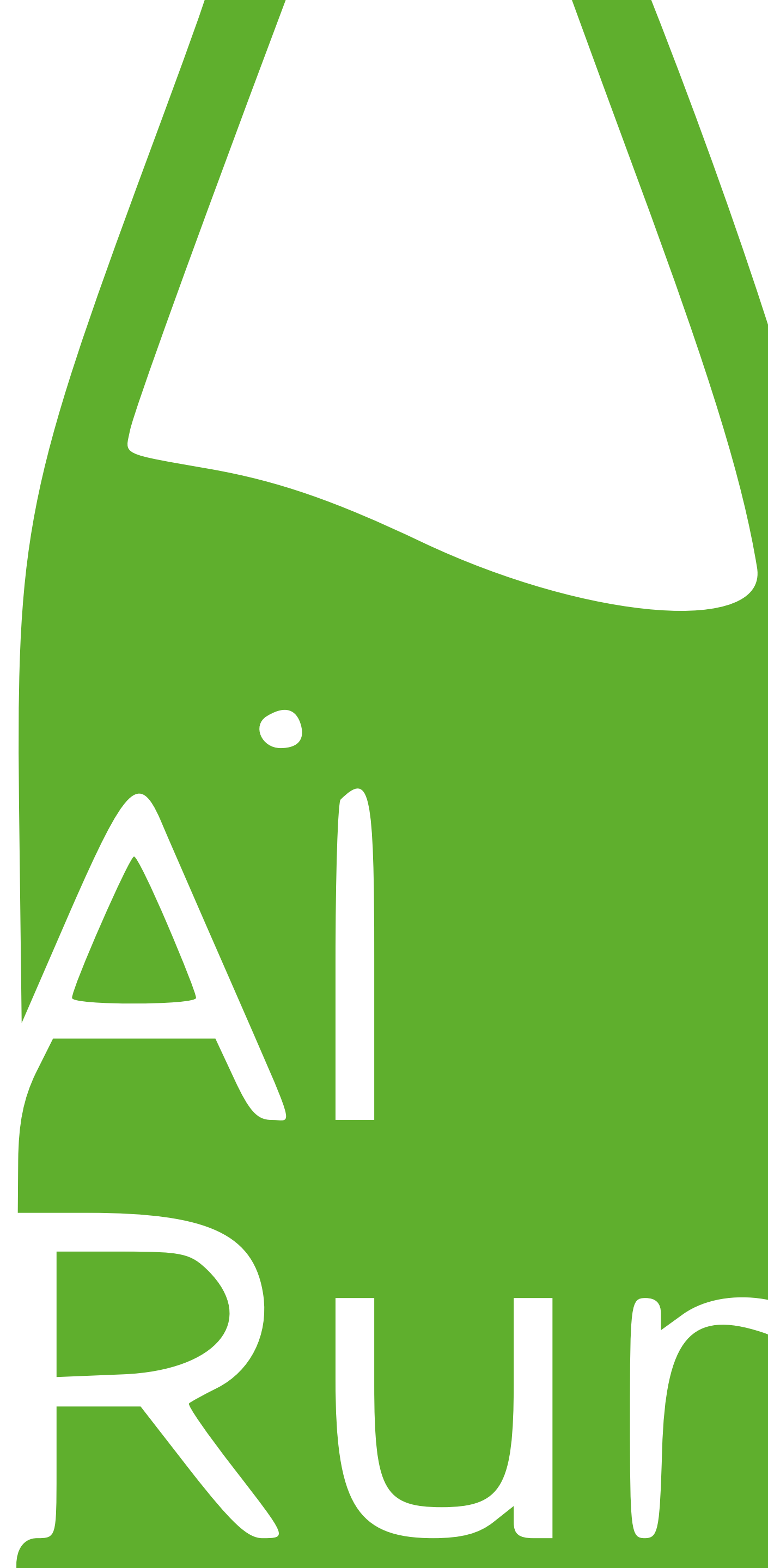


SEQ2SEQ И ATTENTION

Андреева Дарья
Data Scientist, X5 Tech



машинный перевод

1954г. – джорджтаунский эксперимент

до 2010г. – статистические модели

с 2016г. – машинный перевод

формально

source: $x = (x_1, \dots, x_n)$

target: $y = (y_1, \dots, y_m)$

C – корпус из пар (x, y)

$$P_{\theta}(Y|X) = \prod_{i=1}^n p(y_i | X, y_{<i}) \longrightarrow \max$$

$$\hat{y}_i = \operatorname{argmax}(p_{\theta}(y_i | X, \hat{y}_{<i}))$$

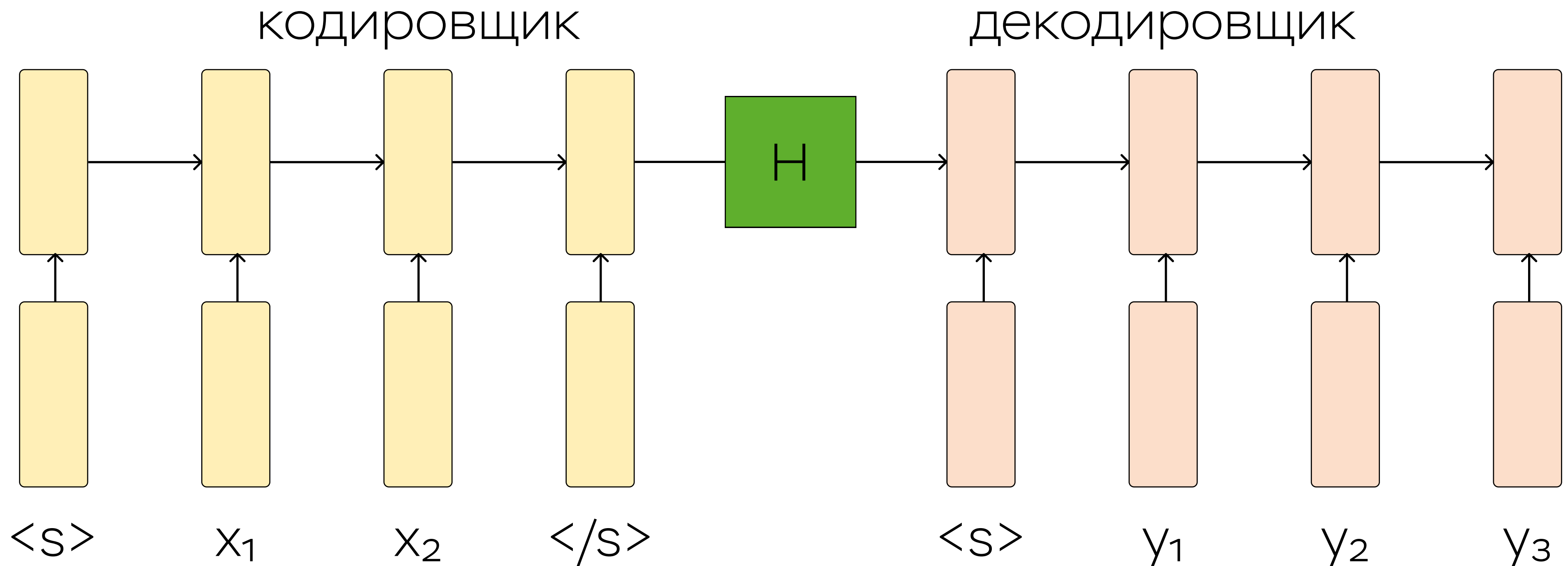
жадно семплируем?

КОДИРОВЩИК-ДЕКОДИРОВЩИК

source: $x = (x_1, \dots, x_n)$

target: $y = (y_1, \dots, y_m)$

C – корпус из пар (x, y)

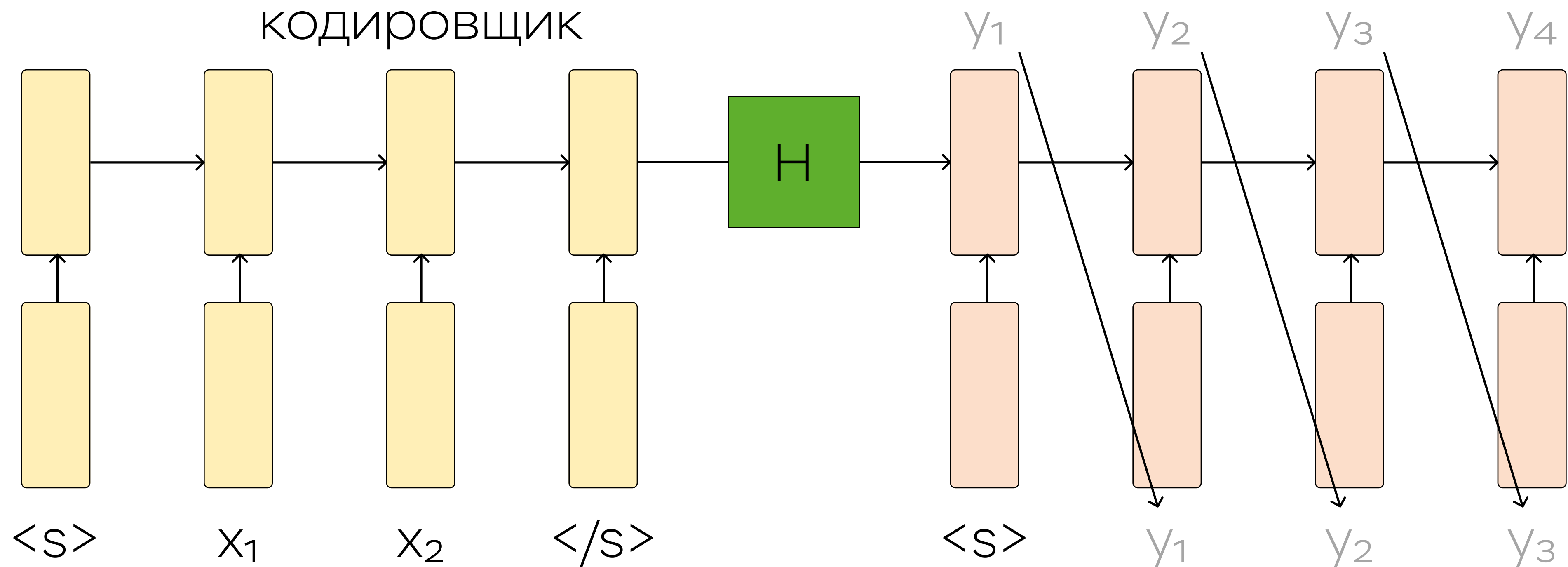


КОДИРОВЩИК-ДЕКОДИРОВЩИК

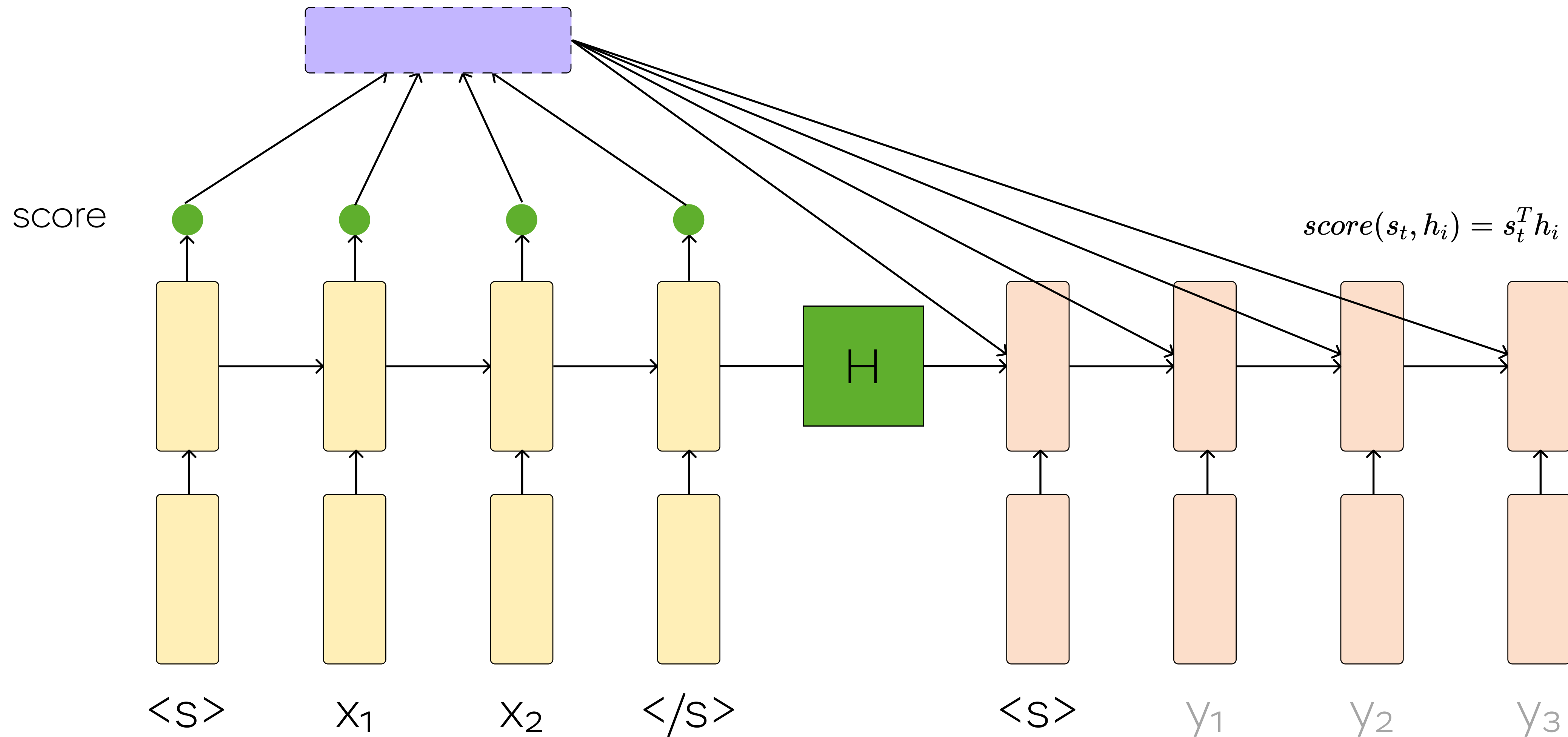
source: $x = (x_1, \dots, x_n)$

target: $y = (y_1, \dots, y_m)$

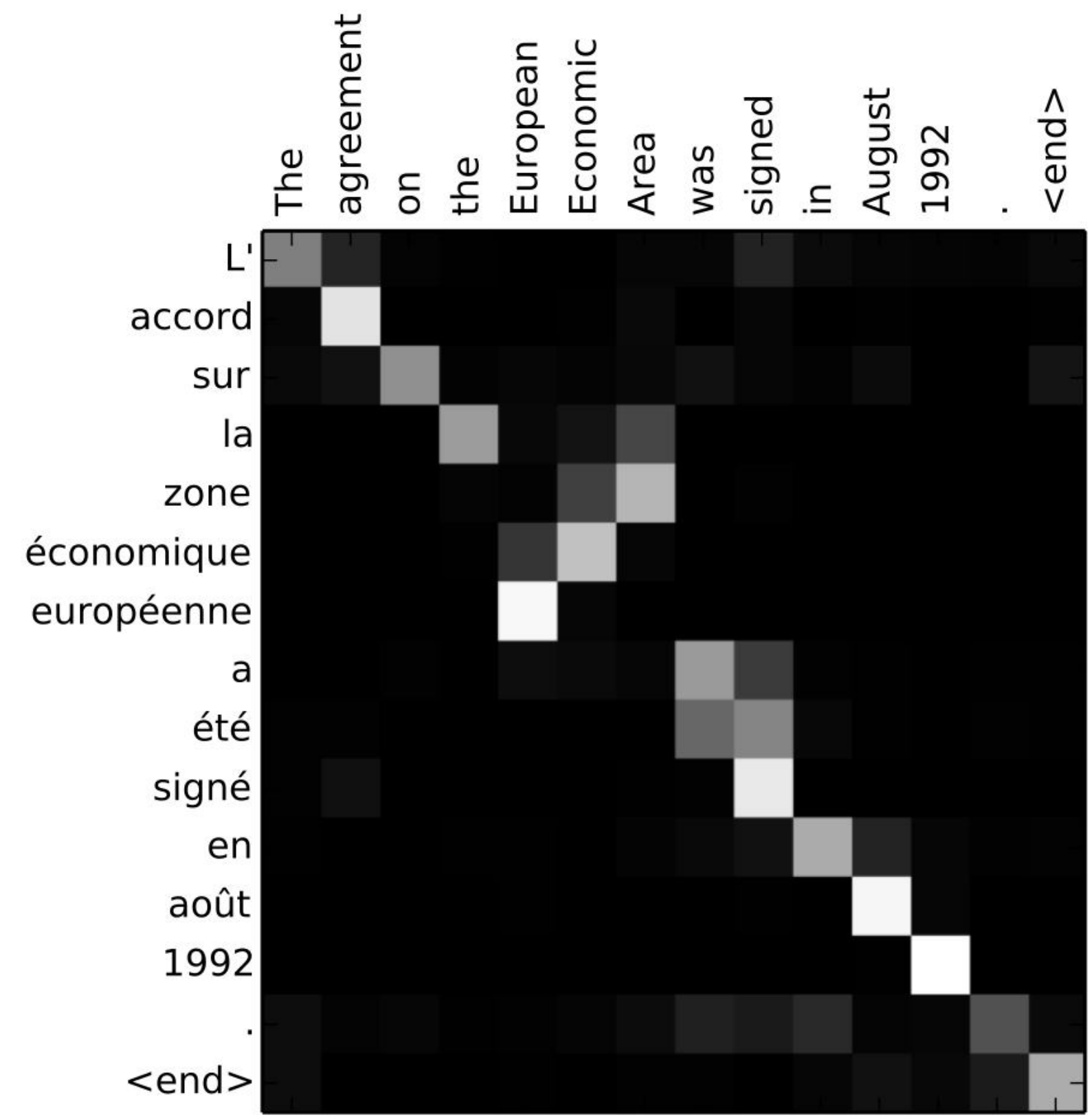
C – корпус из пар (x, y)



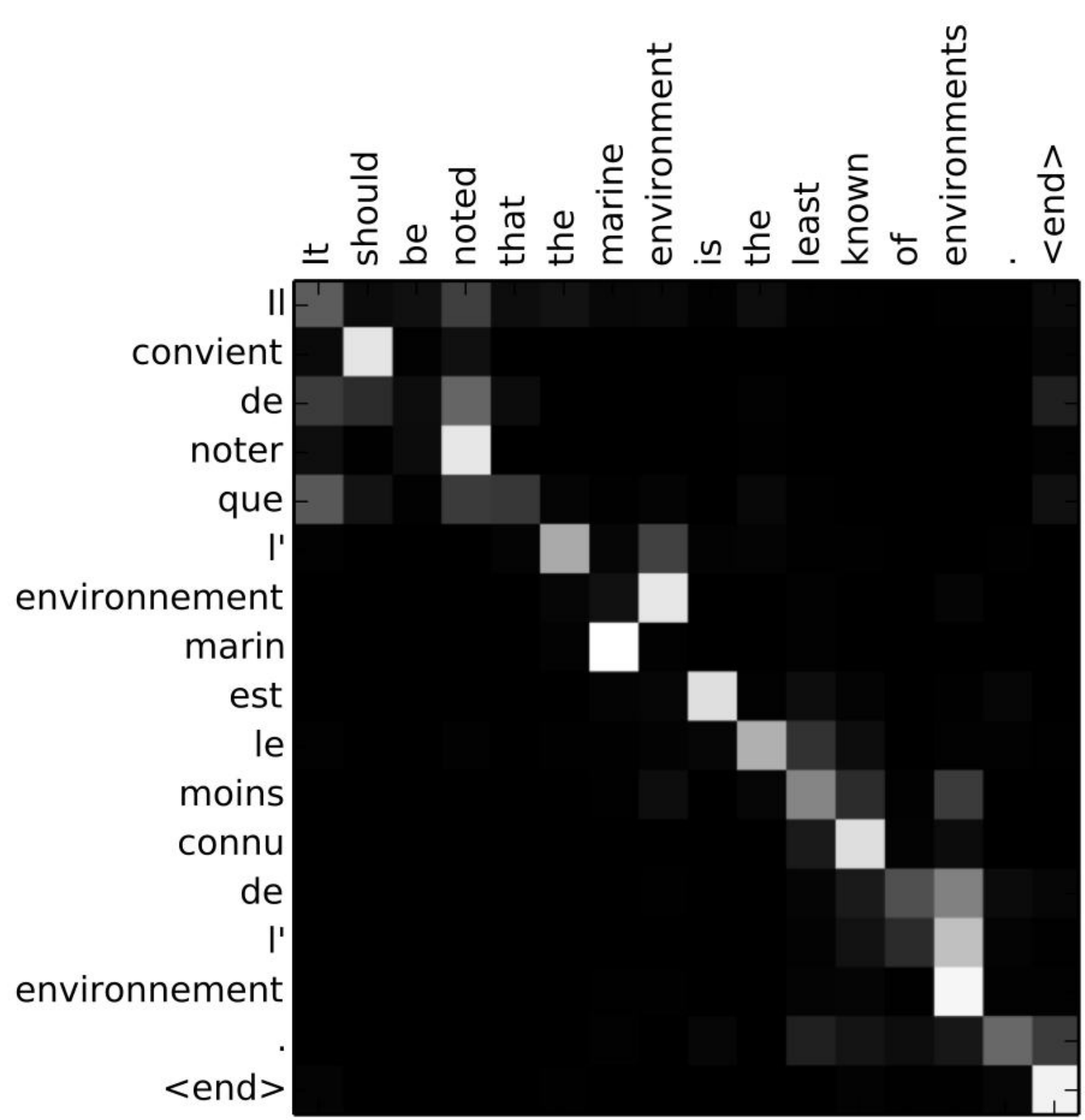
ВНИМАНИЕ



ВНИМАНИЕ



(a)



(b)

ATTENTION

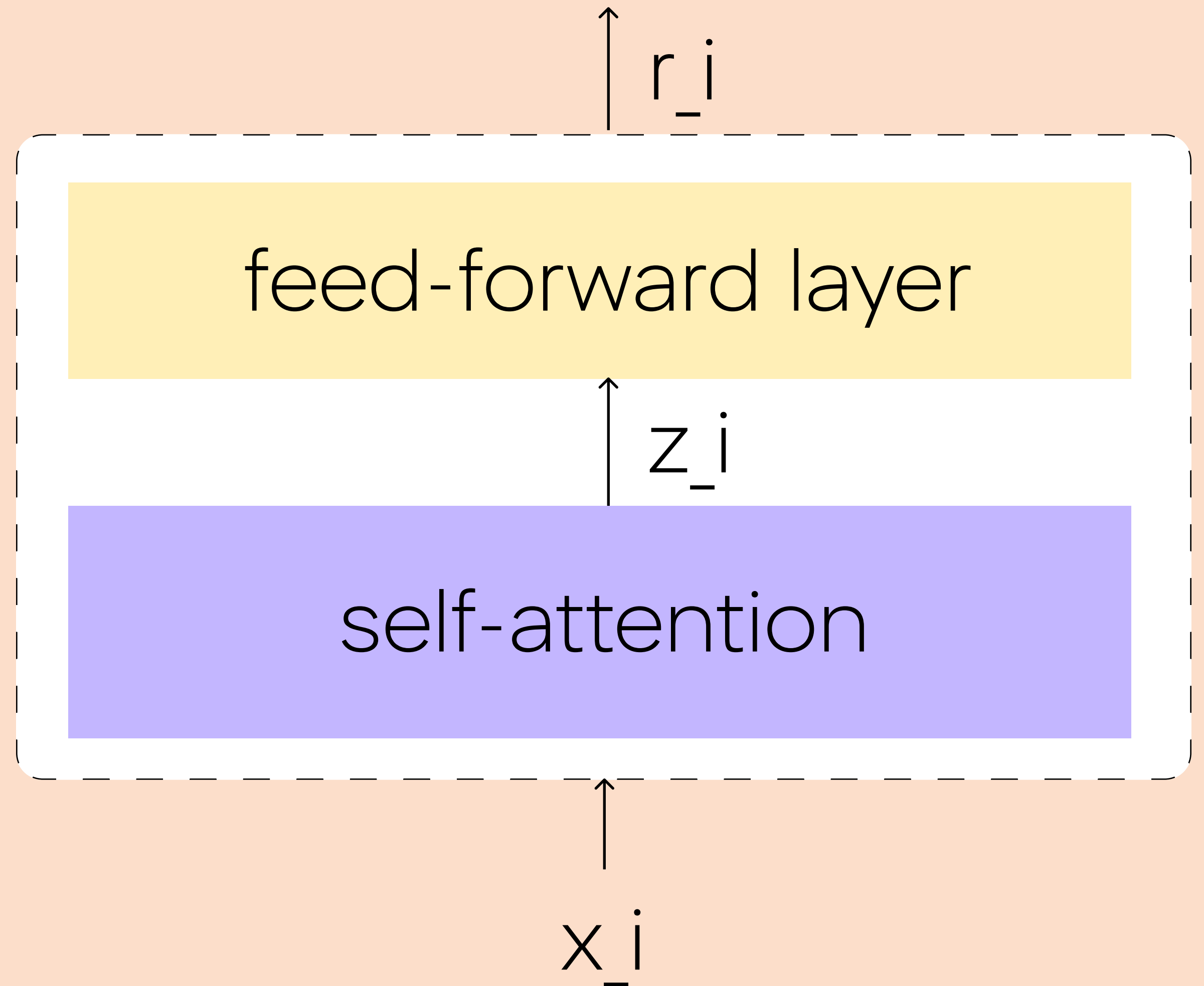
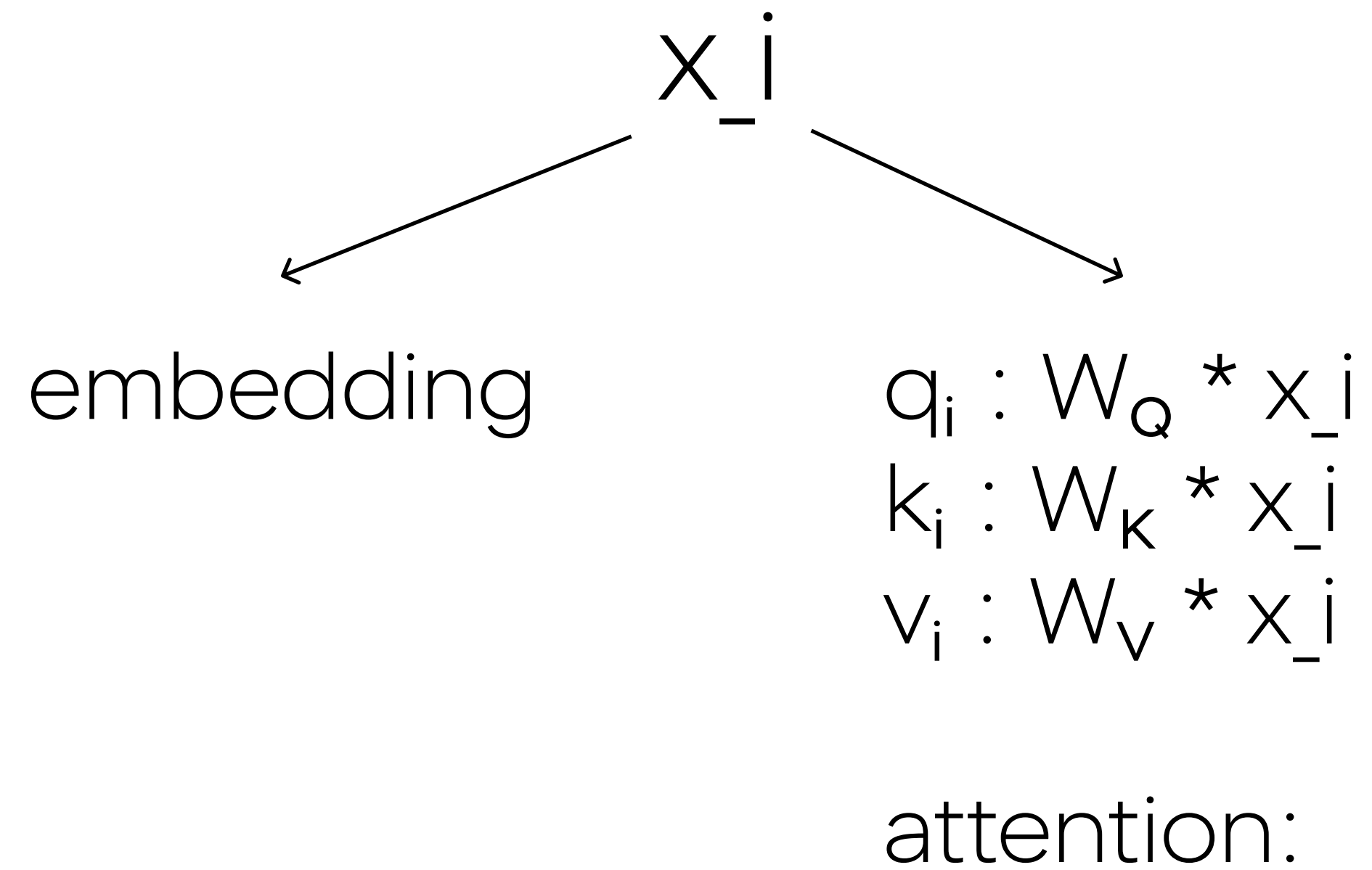
IS

ALL

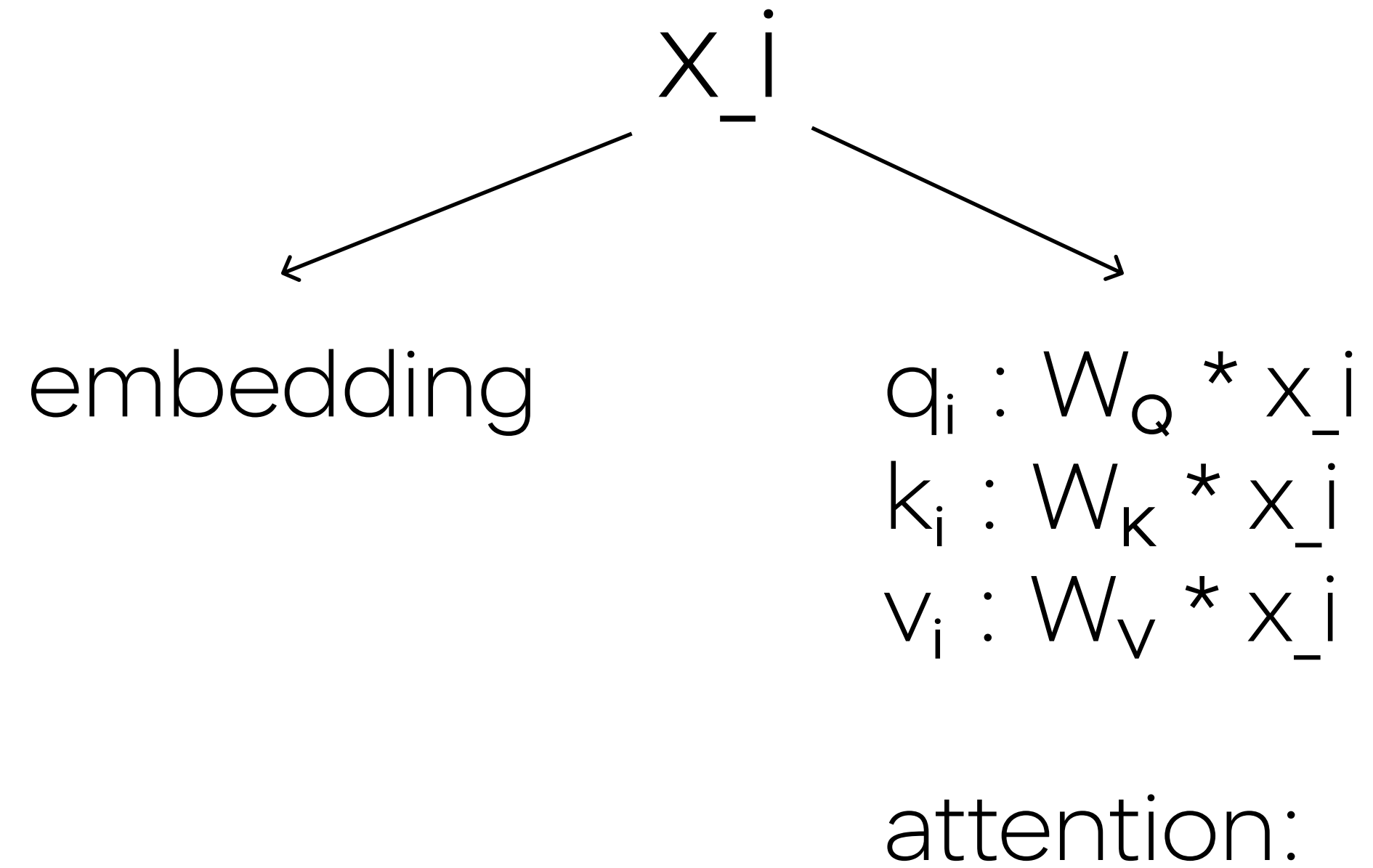
YOU

NEED

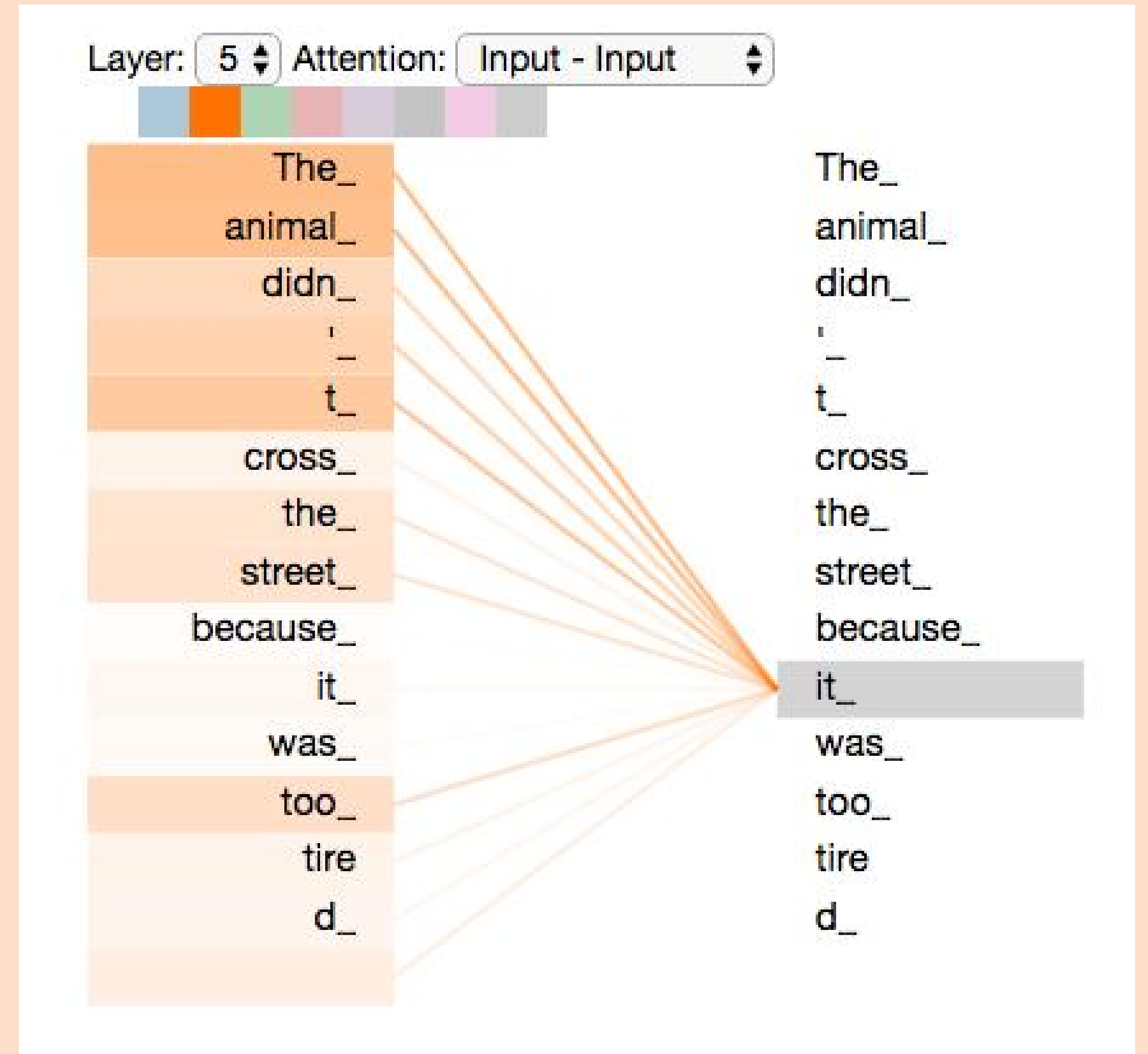
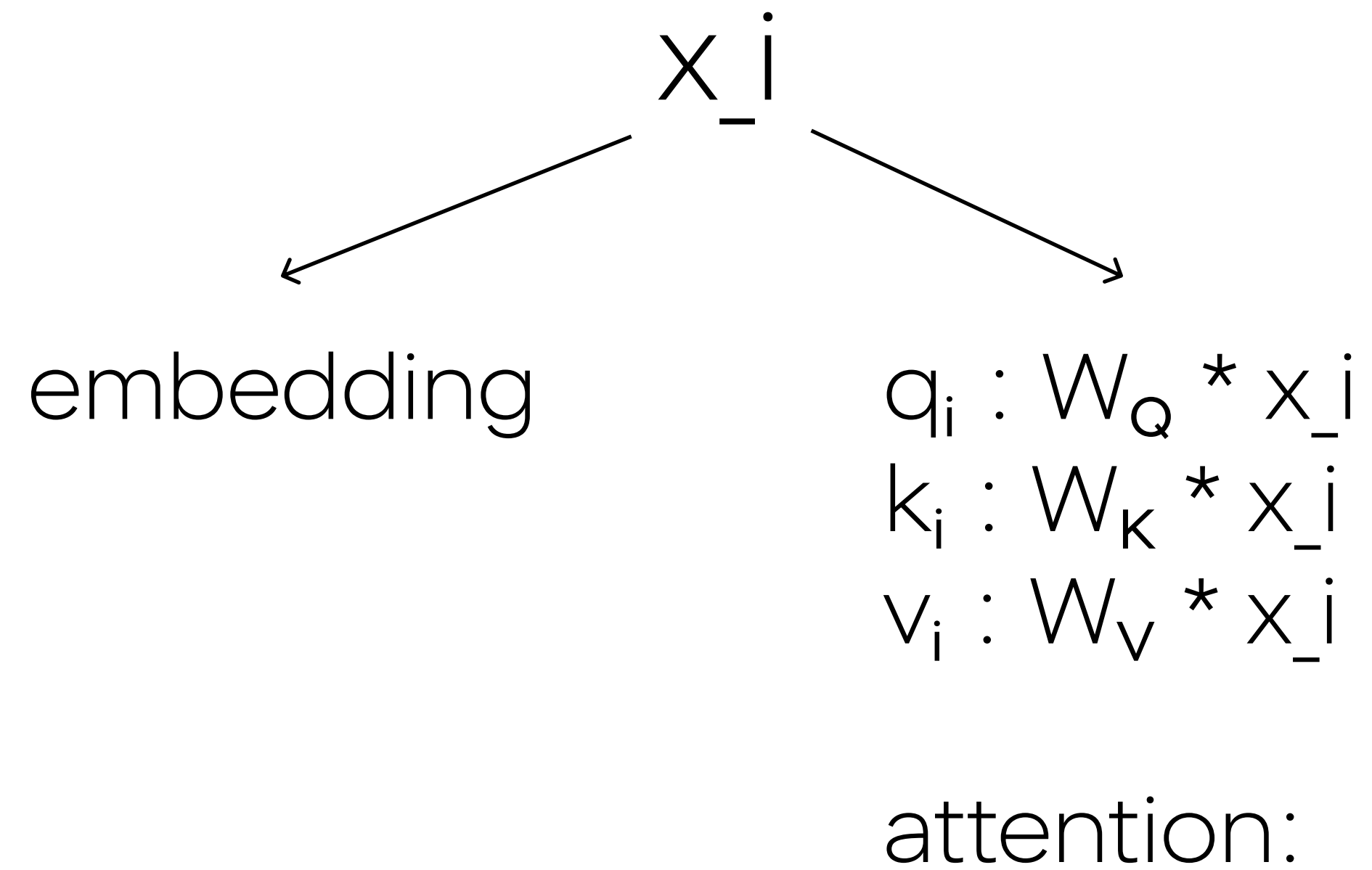
encoder



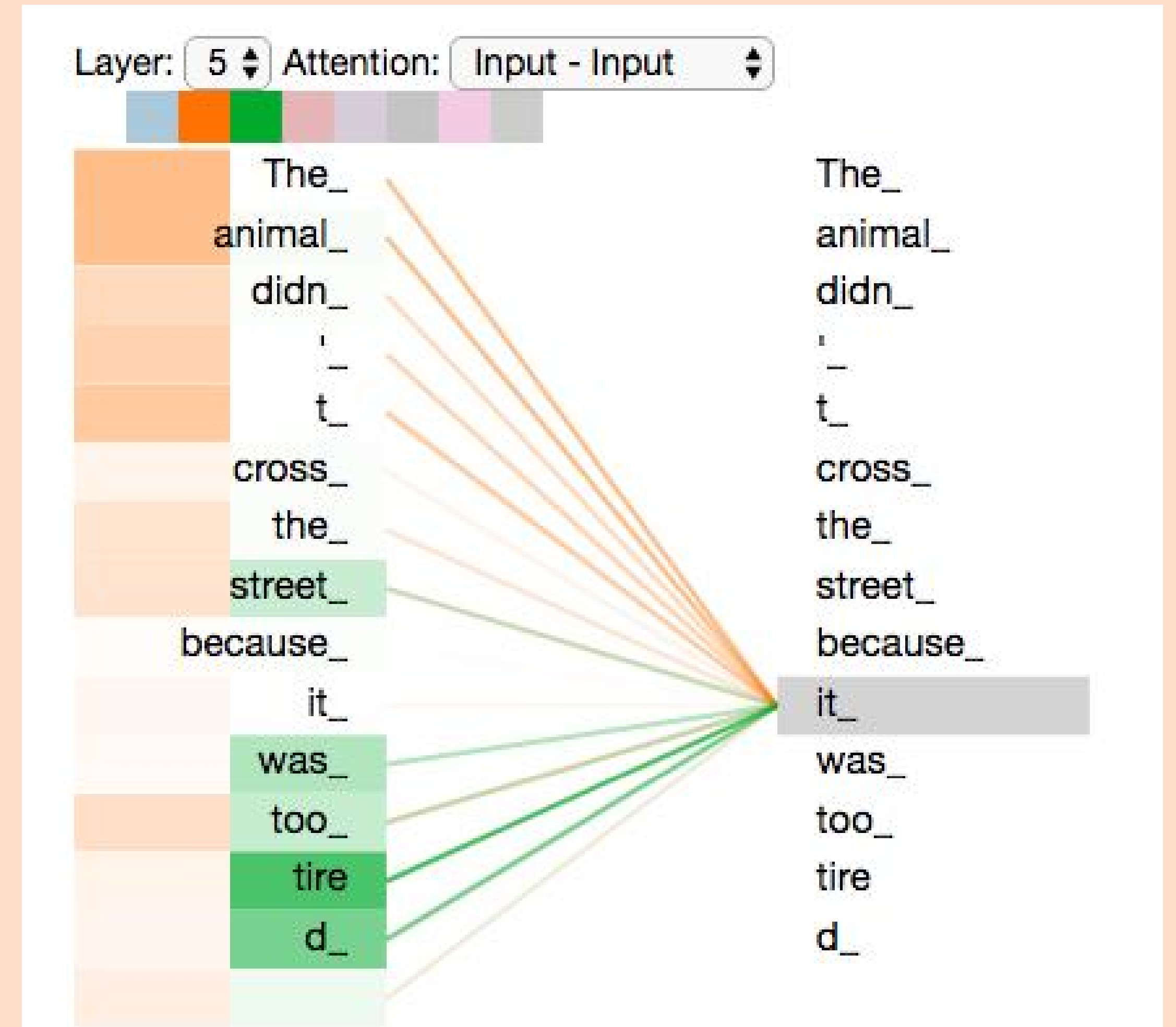
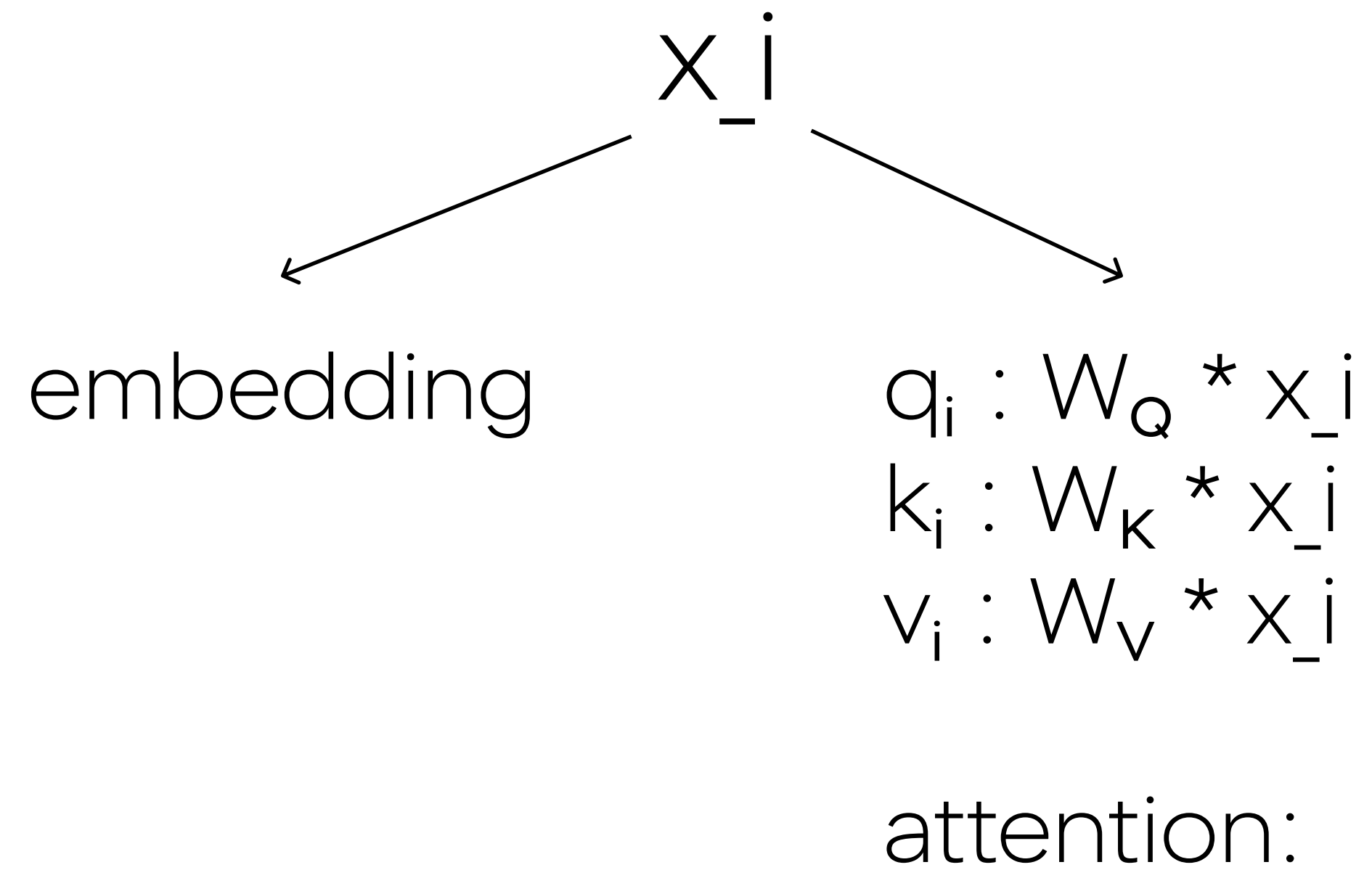
encoder



encoder

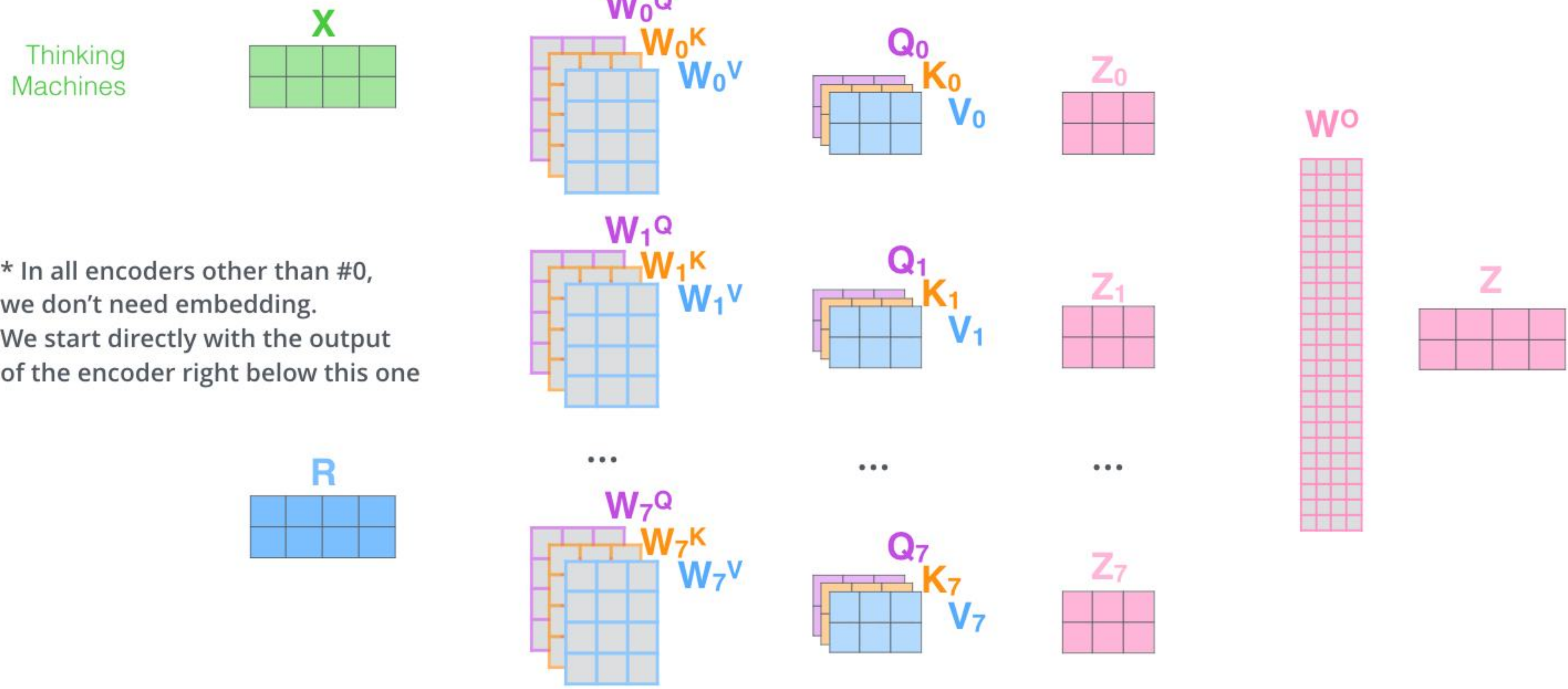


encoder

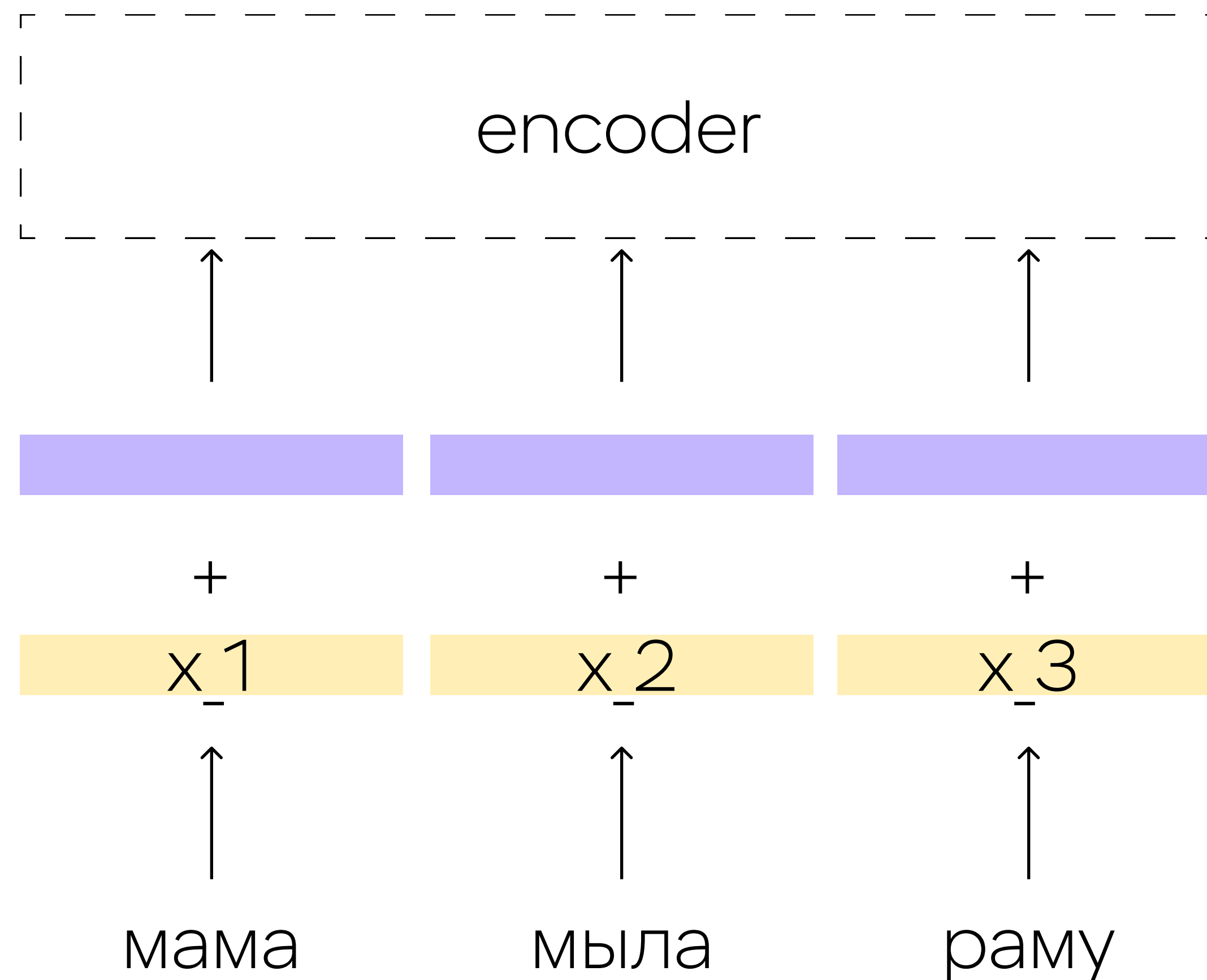


encoder

- 1) This is our input sentence*
- 2) We embed each word*
- 3) Split into 8 heads. We multiply X or R with weight matrices
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer



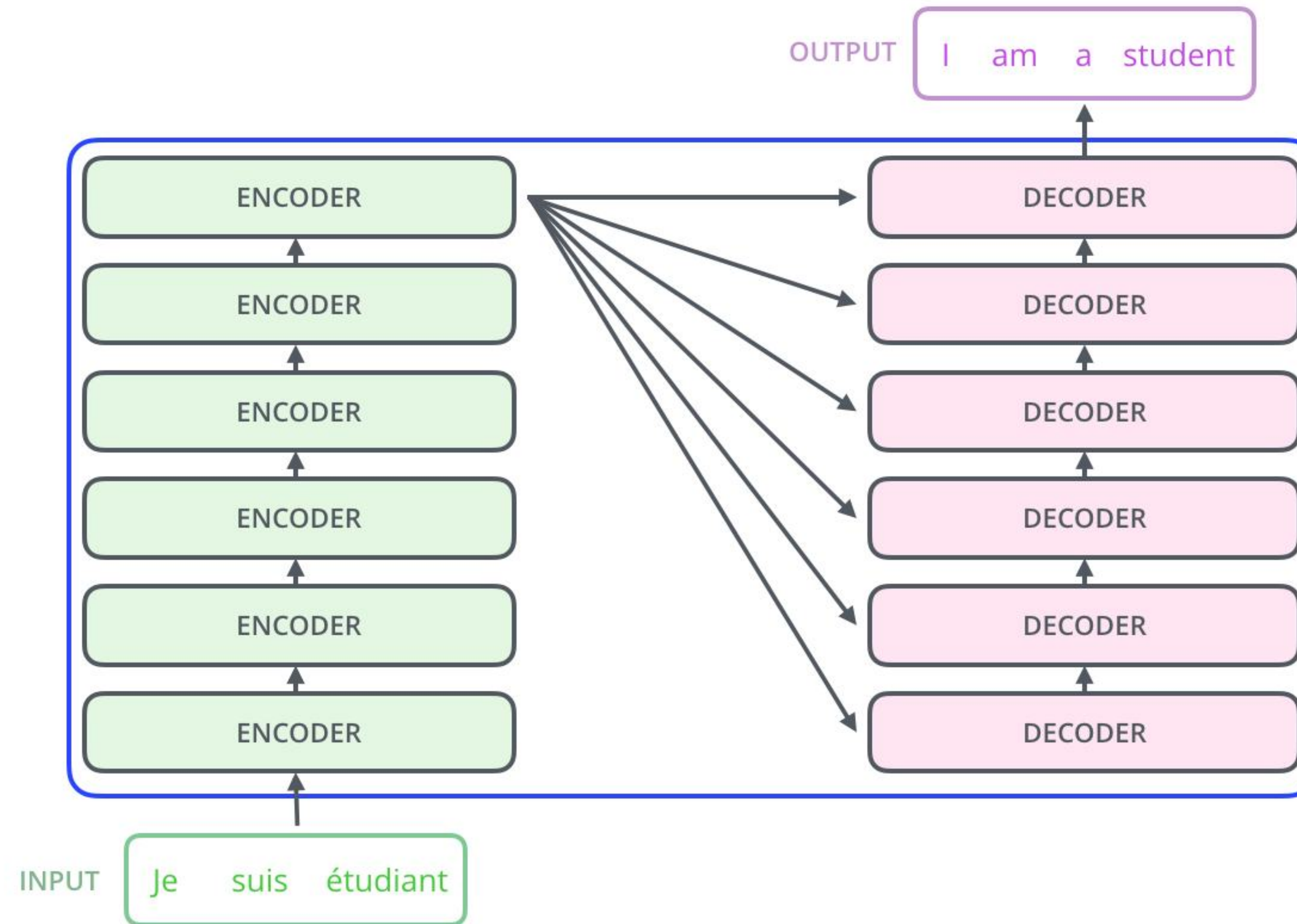
ПОЗИЦИОННЫЙ ЭМБЕДДИНГ



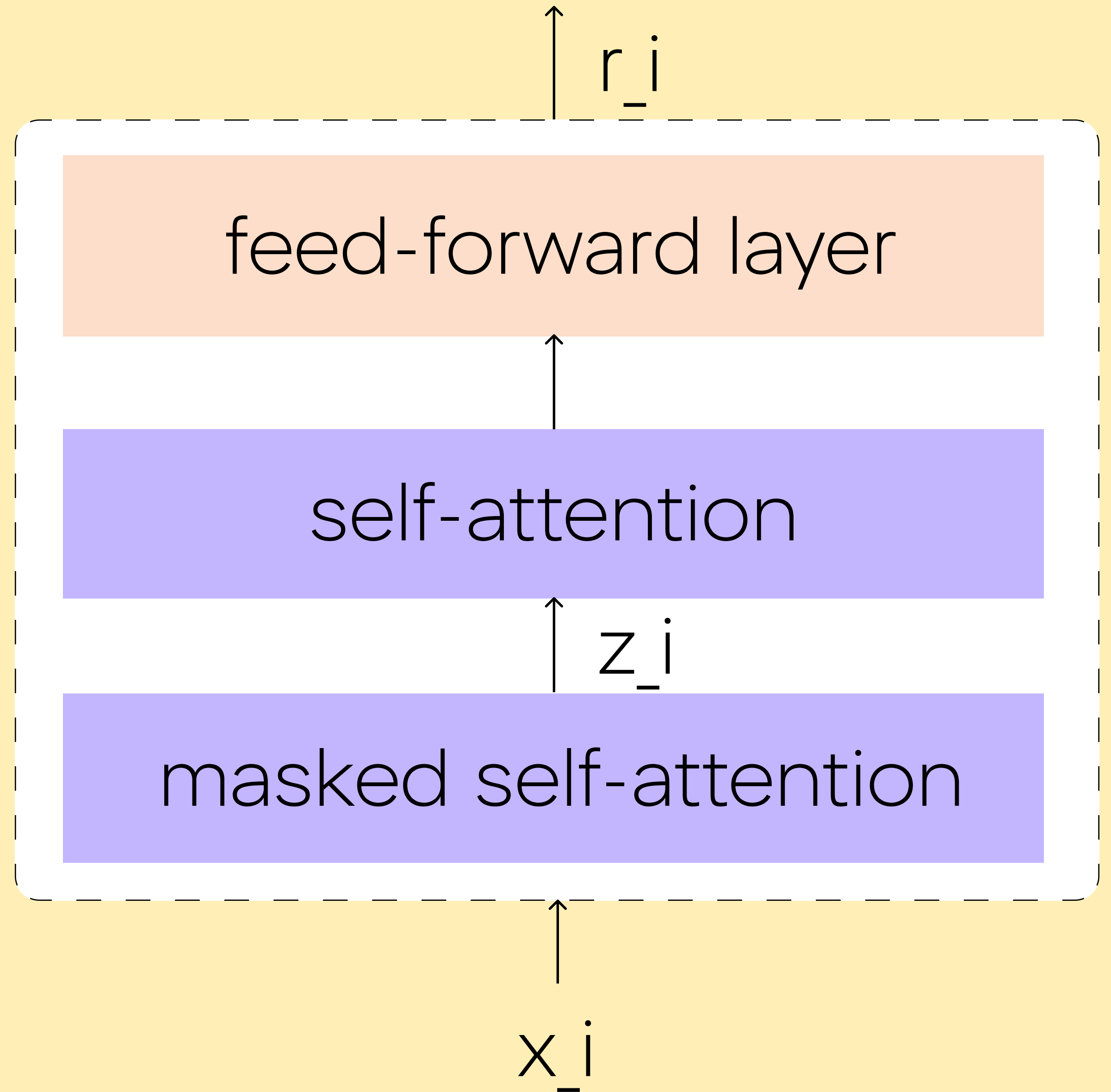
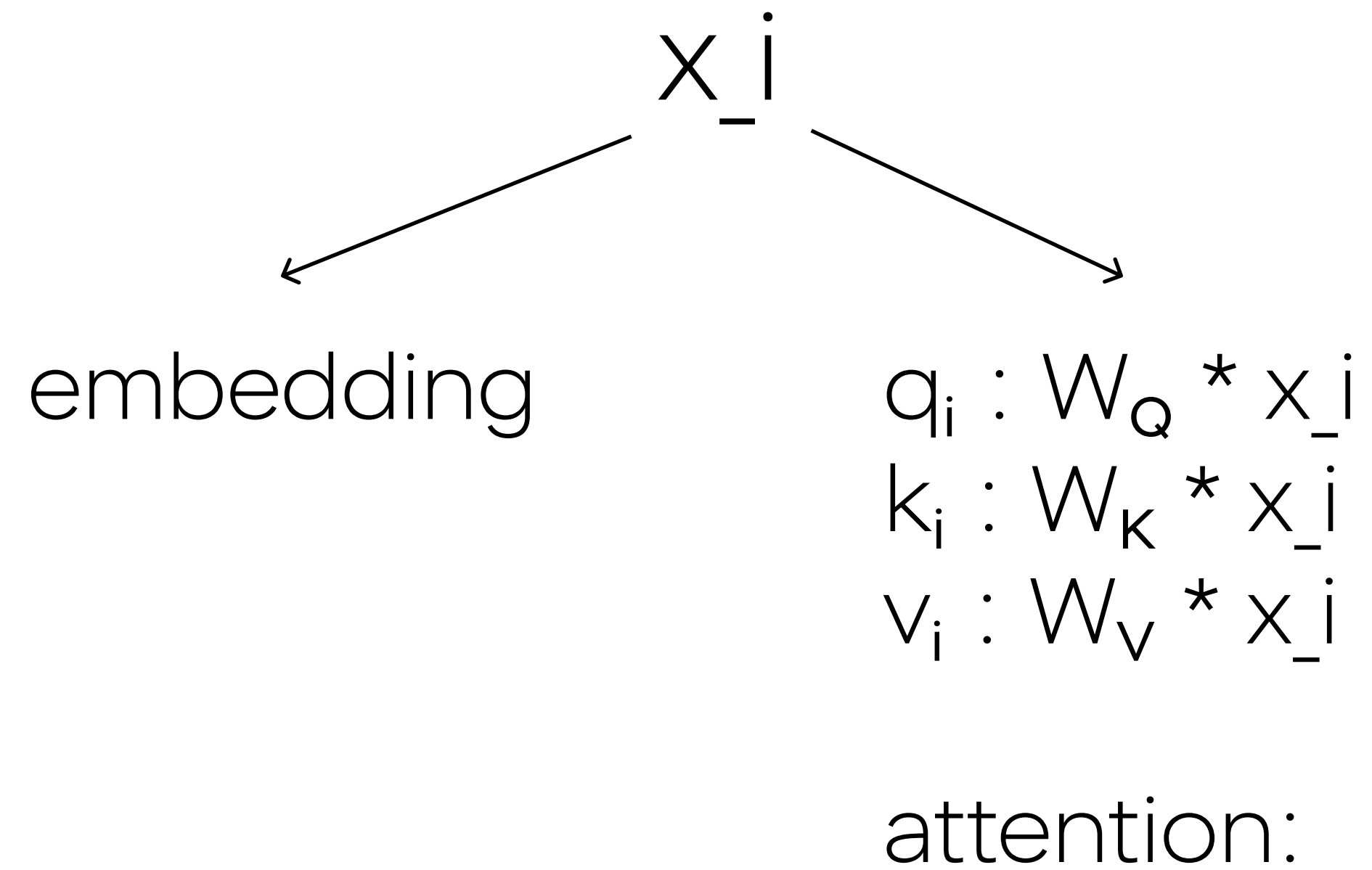
\swarrow
 $\sin(\text{angle_rates})$

\searrow
 $\cos(\text{angle_rates})$

overview



decoder



overview

