

The Transformer

1. Self-attention

input sequence: $x_1, x_2 \dots x_T \quad x_t \in \mathbb{R}^d$

Queries: $q_t = Qx_t \quad Q \in \mathbb{R}^{d_q \times d}$

Keys: $k_t = Kx_t \quad K \in \mathbb{R}^{d_k \times d}$

Values: $v_t = Vx_t \quad V \in \mathbb{R}^{d_v \times d}$

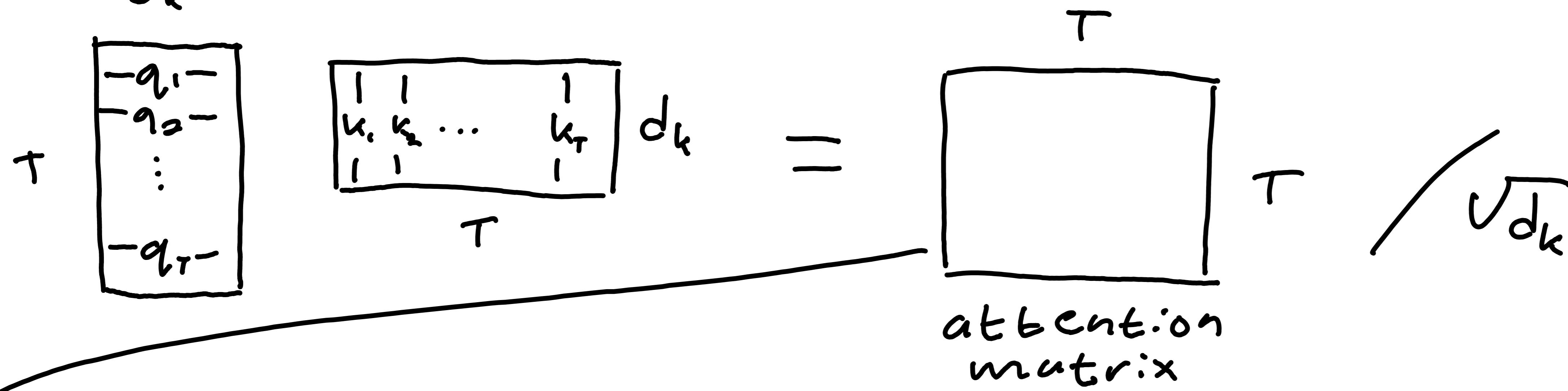
Comparison function: \rightarrow trainable params

$$a_{ij} = \alpha(q_i, k_j) = \frac{q_i^T k_j}{\sqrt{d_k}}$$

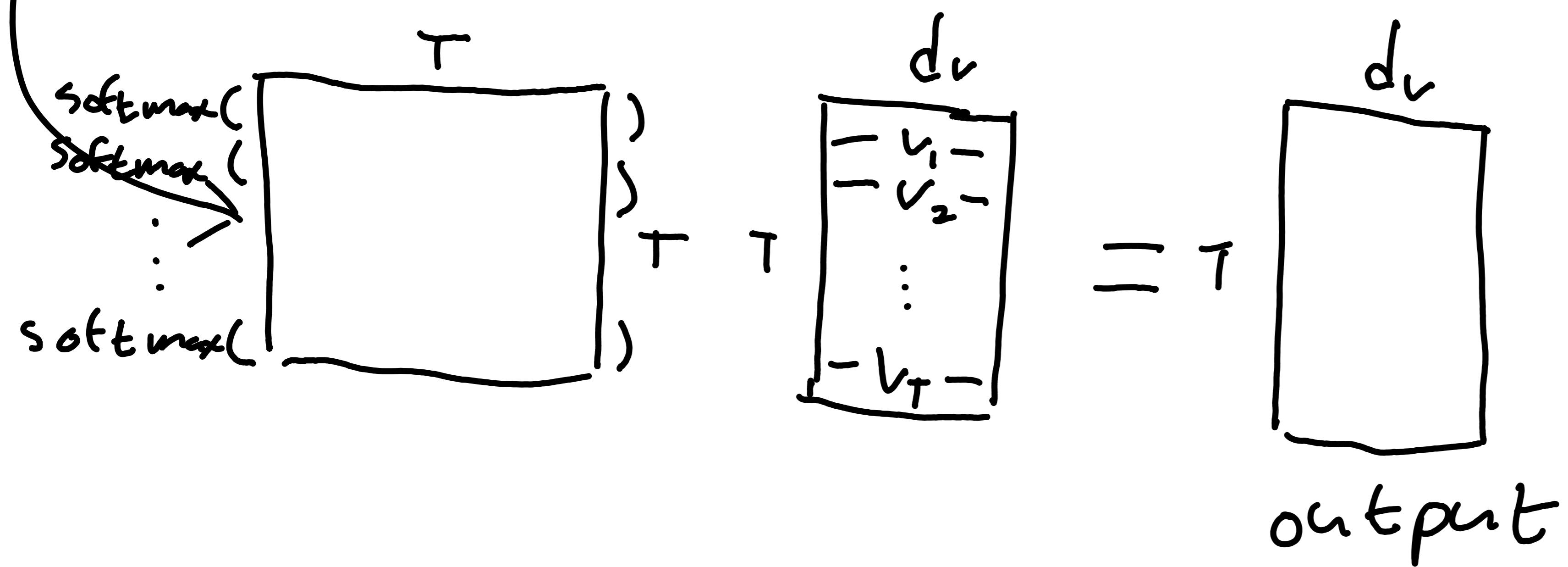
Reduction function:

softmax-weighted average of v_1, \dots, v_T

comparison function:



Reduction function



2. Position Embeddings

Self-attention is orderless (set operation)

- a) Learned vector for each sequence position
embedding
(can't extrapolate)
- b) Relative embeddings for each sequence offset up to some maximum offset

how?

i) $\alpha(q_i, k_j) = \frac{q_i^T (k_j + p_{ij})}{\sqrt{d_k}}$ learned vector $\in \mathbb{R}^{d_k}$

ii) $\alpha(q_i, k_j) = \frac{q_i^T k_j}{\sqrt{d_k}} + p_{ij}$ learned scalar

iii) $\alpha(q_i, k_j) = \frac{q_i^T k_j}{\sqrt{d_k}} + \lambda(i - j)$
↳ learned scalar

3. Multi-headed self-attention

Rather than having a single self-attention, we have H in parallel, then we aggregate the result.

$$h^{(i)} = \text{self attention}(Q^{(i)}, K^{(i)}, V^{(i)})$$

↳ sequence

↳ different parameters for each head

$$\text{output} = O \begin{bmatrix} h^{(1)} \\ h^{(2)} \\ \vdots \\ h^{(H)} \end{bmatrix}$$

↳ output matrix $R^{d_v \times H d_v}$

4. Layer norm

Given a length- T sequence $h_1, h_2, \dots, h_T \in \mathbb{R}^d$

$$\mu_t = \frac{1}{d} \sum_{i=1}^d h_{t,i}$$

$$\sigma_t = \sqrt{\frac{1}{d} \sum_{i=1}^d (h_{t,i} - \mu_t)^2}$$

$$\text{Layer Norm}(h)_t = \frac{h_t - \mu_t}{\sigma_t} + \beta \quad \xrightarrow{\text{learned parameters}}$$

5. Residual connections

$$f(x) \rightarrow f(x) + x$$

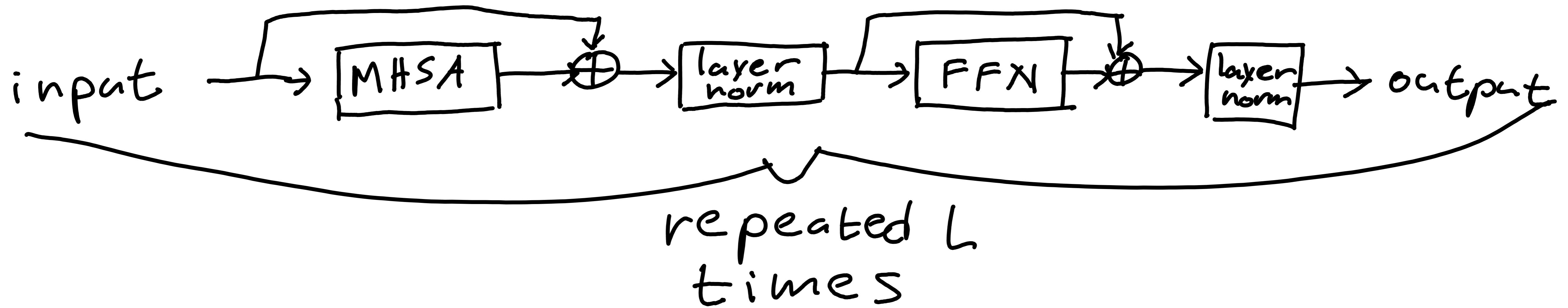
6. "position-wise" feed-forward networks

$$W_2 \phi(W_1 x_t + b_1) + b_2$$

↳ nonlinearity, typically ReLU

7. Embeddings

Transformer layer block



- 1) Transform using self-attention (mixes over time)
- 2) Residual connection + layer norm
for stability
- 3) Transform each position independently using FFN

Transformer block modifications

1. Layer norm position (pre- or post-)

pre-LN :

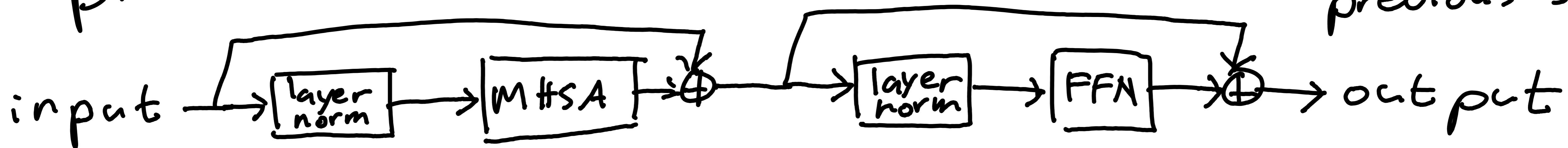


diagram from
previous slide

2. Removal of biases

$$W_2 \phi(W, x_t)$$

FFN

$$\frac{\partial}{\partial t} (h_t - \mu_t)$$

RMSNorm

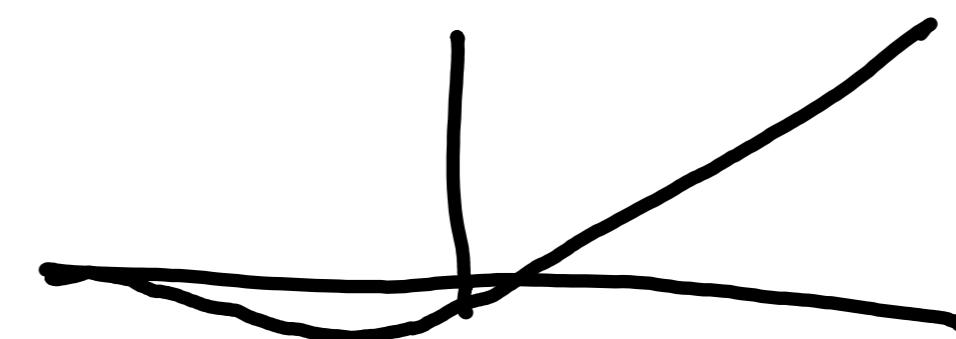
3. FFN non linearity \rightarrow gated linear unit (GLU)

$$W_3 (\phi(W, x_t) \odot W_2 x_t)$$

↳ often "swish"

$$\phi(x) = x \sigma(x)$$

↳ sigmoid

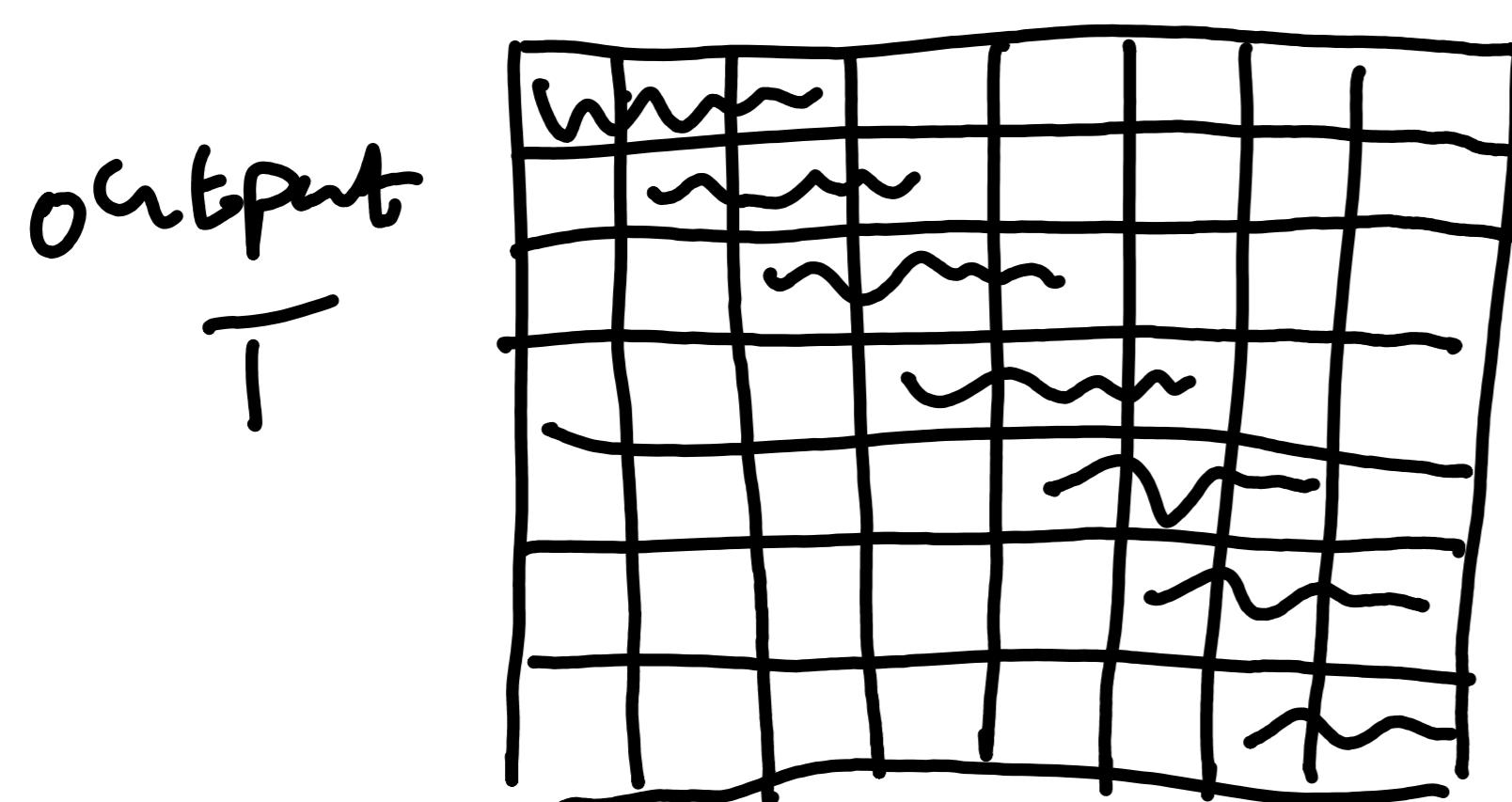


4. Sub-quadratic self-attention

Notes:

- i) For typical Transformer dimensions and modest sequence lengths ($T < 1000$) compute for FFN \gg compute for MHA
- ii) Can get around quadratic memory cost by being clever (Flash Attention)

e.g. "local" self-attention \rightarrow only compute attention over a neighbourhood around the current position



5. "parallel" block

standard: $x + \text{FFN}(\text{LN}(x + \text{MHSA}(x)))$
(pre-LN)

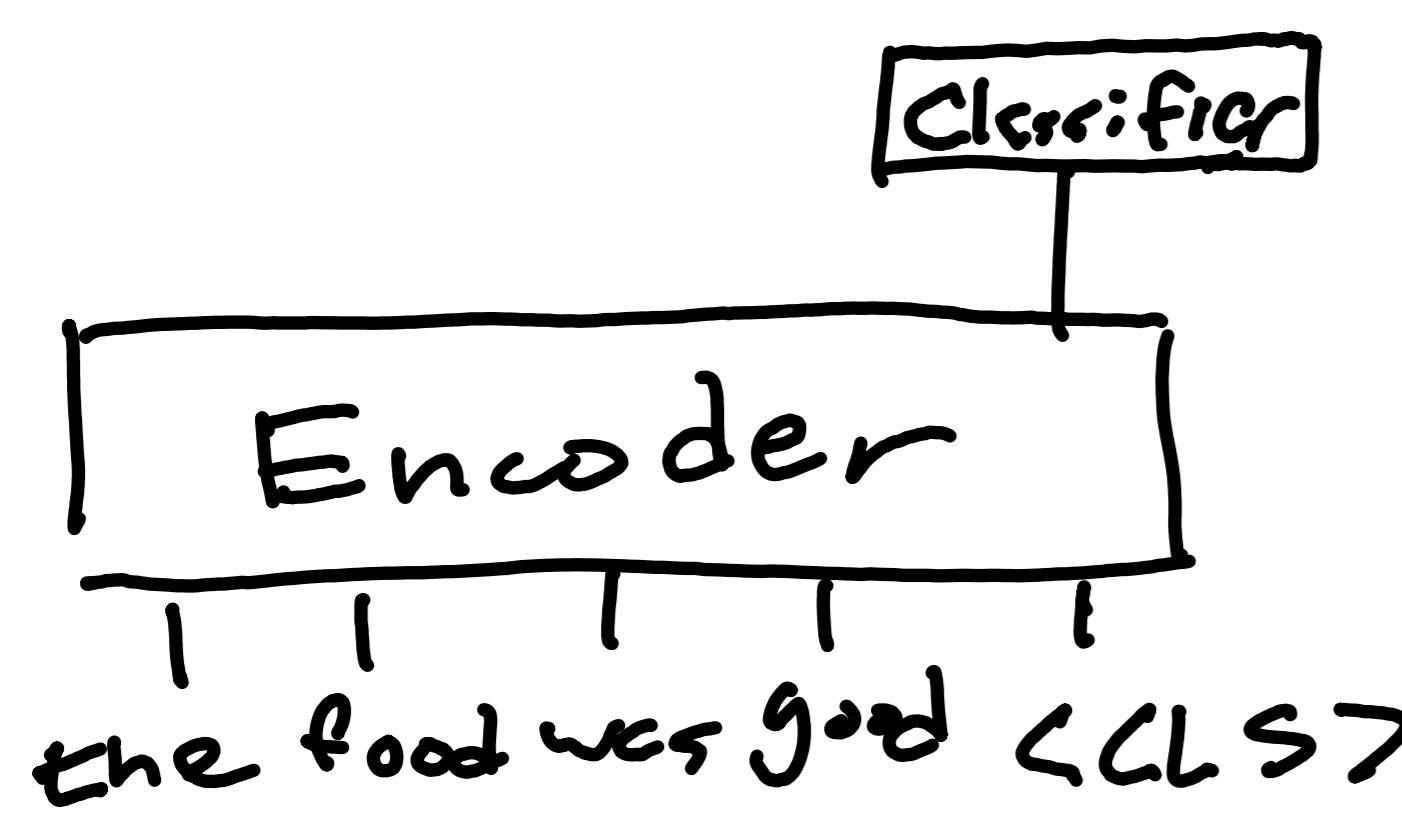
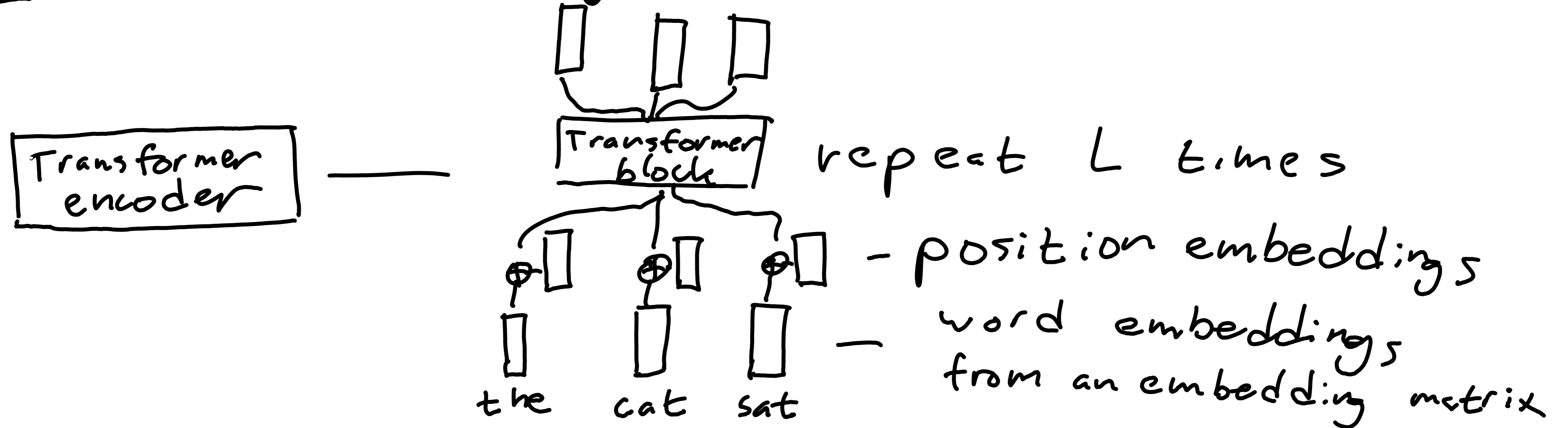
parallel: $x + \text{FFN}(\text{LN}(x)) + \text{MHSA}(\text{LN}(x))$

Doesn't harm performance *

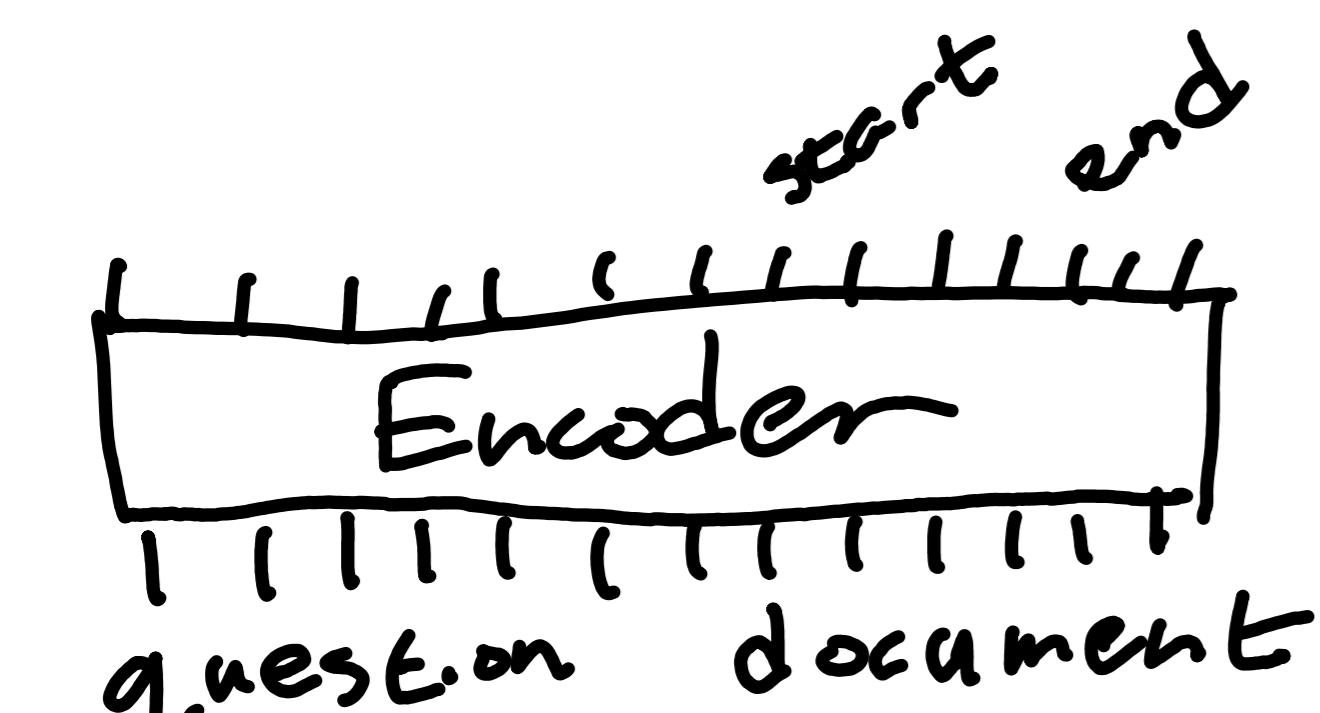
Improves GPU utilization

* maybe

Encoder-only Transformer

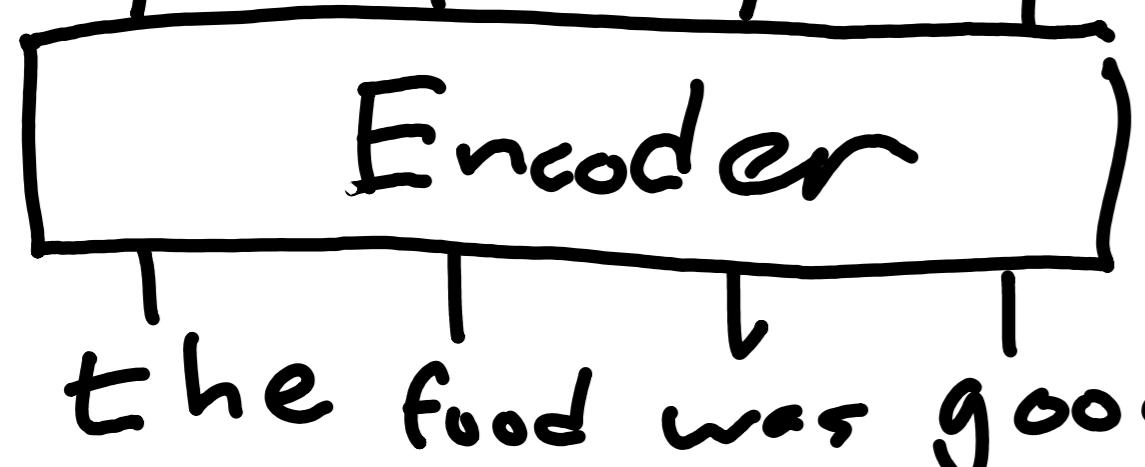


sequence classification



span extraction

article noun verb adj



sequence labeling

Why pre-train? (first step of transfer learning)

1. Transformer has a limited inductive bias.
→ might help to learn structure before applying the model.
2. Often we can learn structure on a large, data-rich (possibly unsupervised) task.
3. Empirically, pre-training can help
 - improve performance
 - using less labeled data
 - with faster convergence

How to pre-train? (on unlabeled data)

For text, use language modeling!

$$P(x_t | x_{t-1}, \dots, x_1)$$

This only conditions on the past.

So, use "masked" language modeling.

$$P(x_t | x_T, x_{T-1}, \dots, \underbrace{x_{t+1}, x_{t-1}, \dots, x_1}_{\text{entire input sequence other than } x_t})$$

MLM

entire input
sequence other
than x_t

How does BERT implement MLM?

Bidirectional encoder
representations from
Transformers

the cat sat on the mat

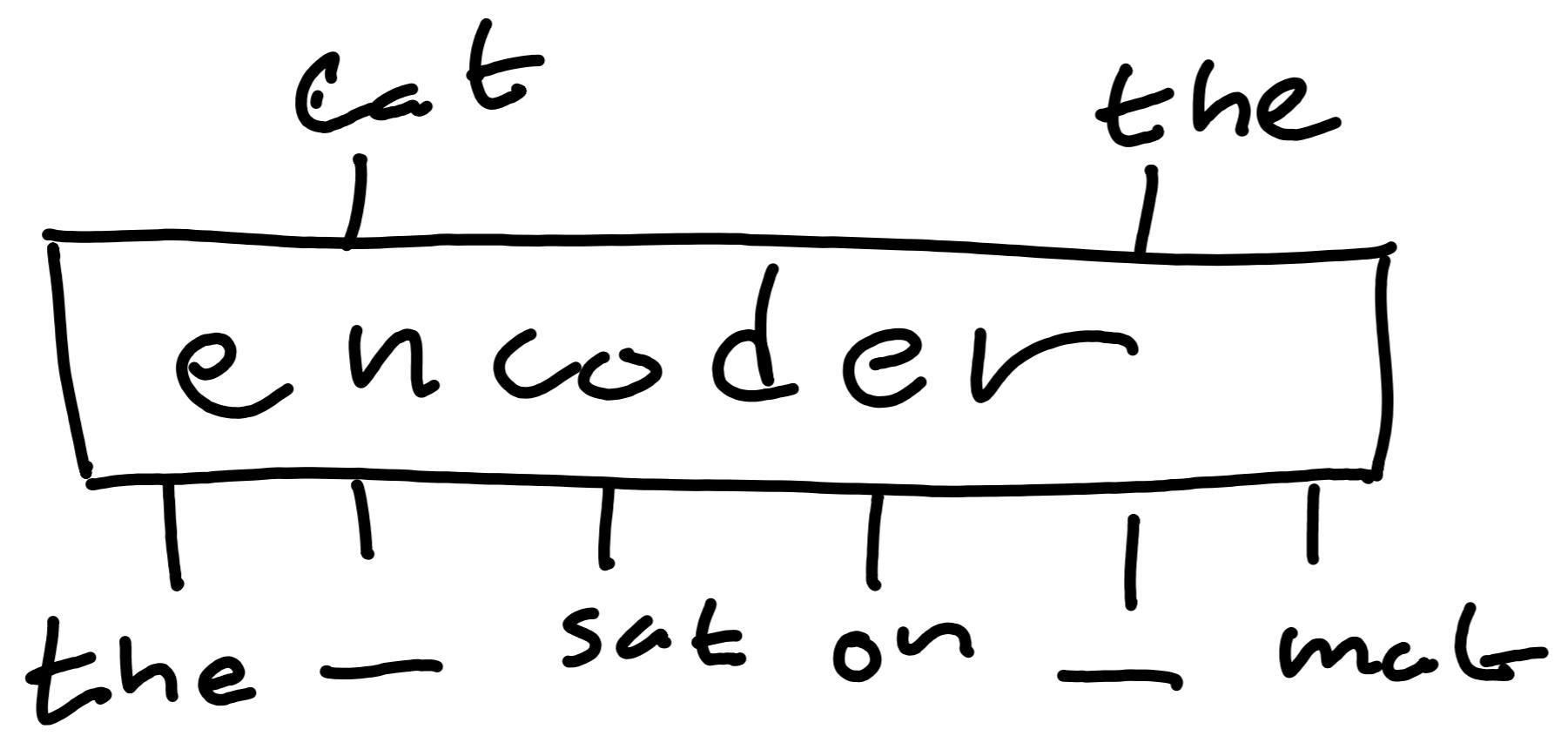


the — sat on — mat

$$P(x_2 = \text{cat} | \text{the} - \text{sat} \text{ on} - \text{mat})$$

$$P(x_5 = \text{the} | \text{the} - \text{sat} \text{ on} - \text{mat})$$

} pre-train
to predict
these



basically sequence
labelling on the
masked-out positions!