# Lecture 19
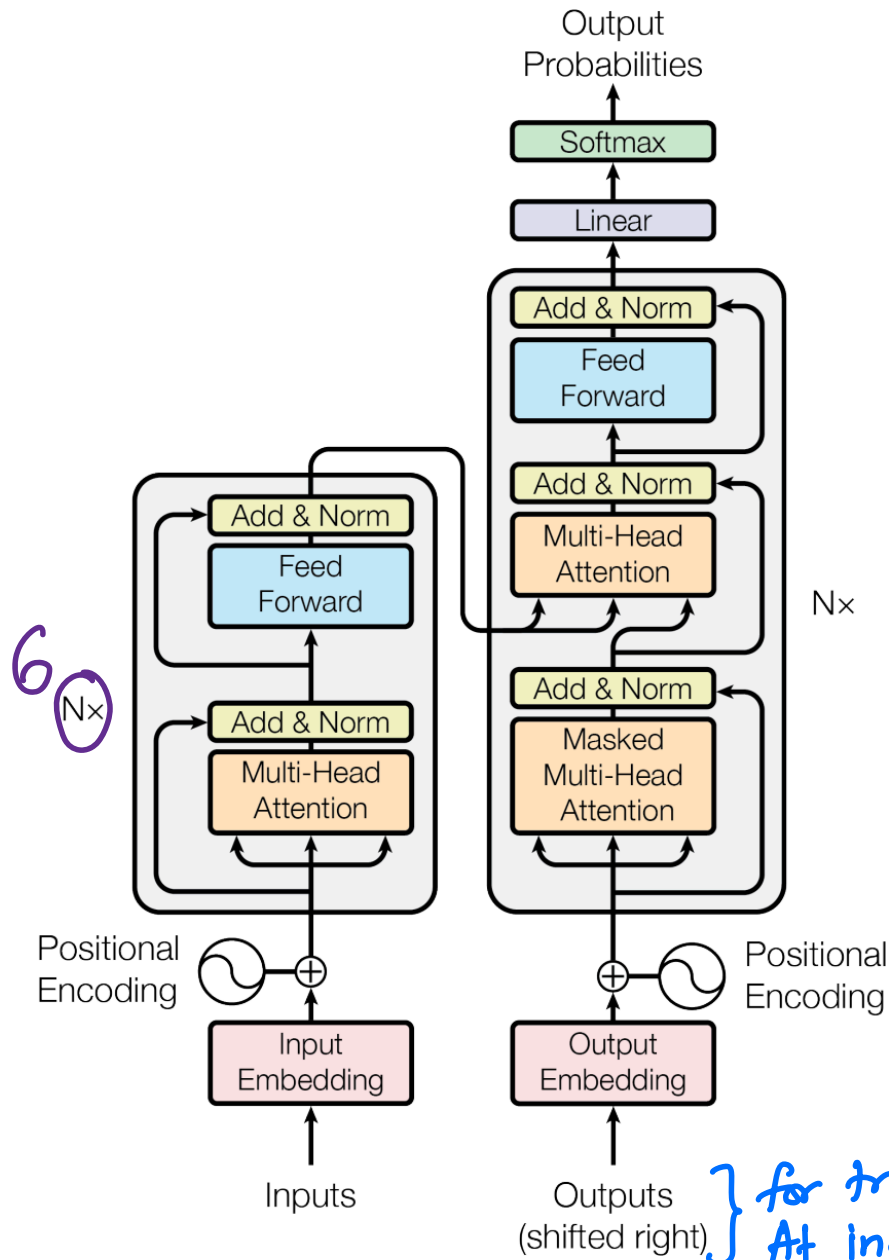
- Attention recap.
- Transformer architecture
- Embedding layer
- Positional encoding.

# Attention is all you need.



Figure 1: The Transformer - model architecture.

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Nx

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

6 Nx

Positional Encoding

Input Embedding

Output Embedding

Positional Encoding

Inputs

Outputs (shifted right)

- Normalization layers are now typically before the attention block.

decoder

6

⋮

1 2

1

} for training.
At inference, just put in <start>.

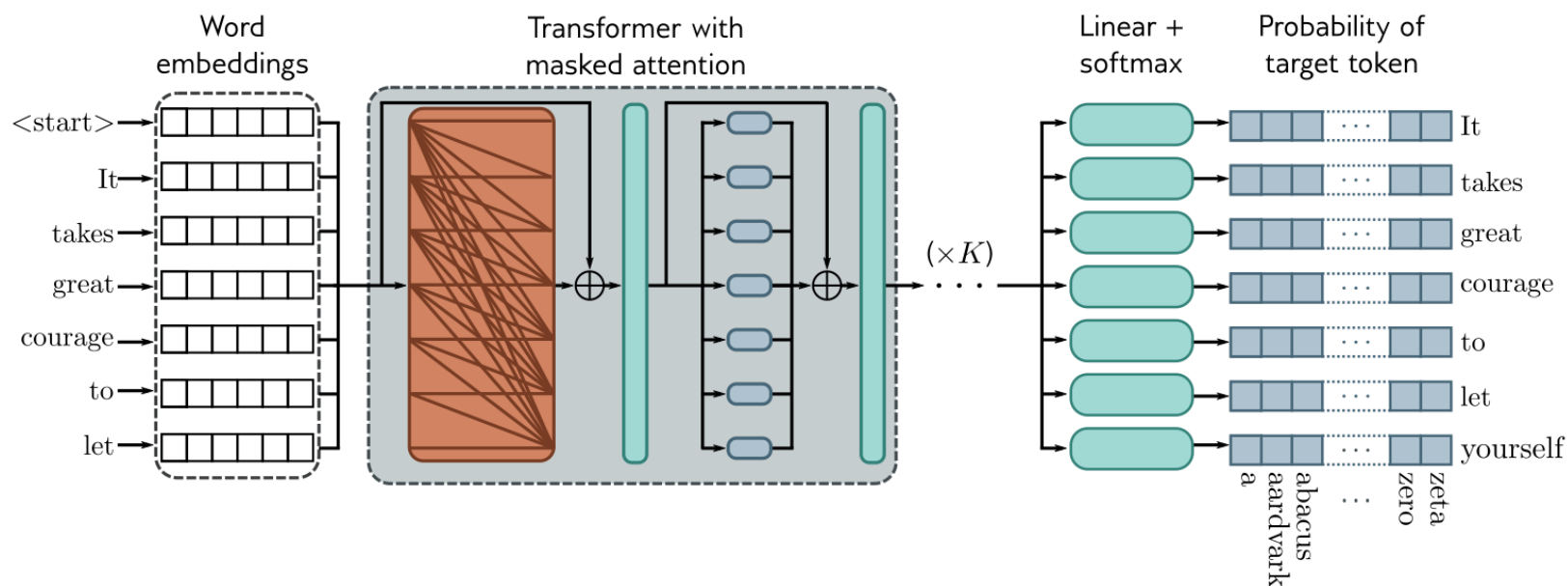# GPT-3 architecture (Decoder only).



**Figure 12.12** Training GPT3-type decoder network. The tokens are mapped to word embeddings with a special <start> token at the beginning of the sequence. The embeddings are passed through a series of transformer layers that use masked self-attention. Here, each position in the sentence can only attend to its own embedding and those of tokens earlier in the sequence (orange connections). The goal at each position is to maximize the probability of the following ground truth token in the sequence. In other words, at position one, we want to maximize the probability of the token It; at position two, we want to maximize the probability of the token takes; and so on. Masked self-attention ensures the system cannot cheat by looking at subsequent inputs. The autoregressive task has the advantage of making efficient use of the data since every word contributes a term to the loss function. However, it only exploits the left context of each word.

# Embedding Layer.

- Use Learned embeddings.

- How?

Assume fixed vocabulary size, e.g. 32,000.
Embedding dimension = 1000 say.

Input is a sequence of tokens from your
vocabulary.

$\longrightarrow$ One hot encoded vectors in $\mathbb{R}^{32,000}$.

$\longrightarrow$ Need 1000 dimensional vectors.

$\longrightarrow$ Use a linear layer and learn this.

- At output: Need to get scores that are compatible with a softmax.
Need to take 1000-dim vector to 32,000 score $\longrightarrow$ Softmax
$\longrightarrow$ Again Linear layer $\longrightarrow$ Probability.
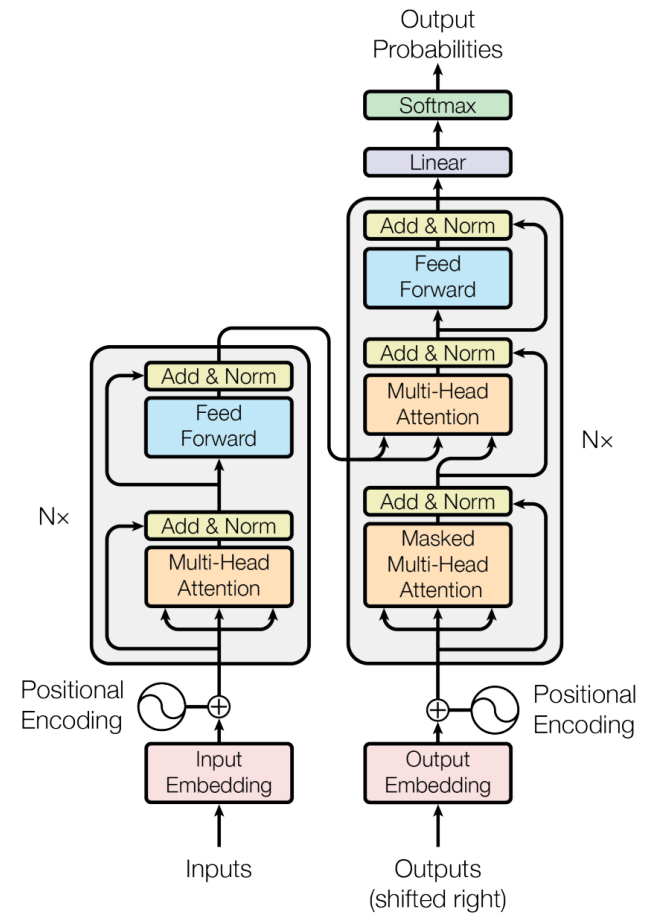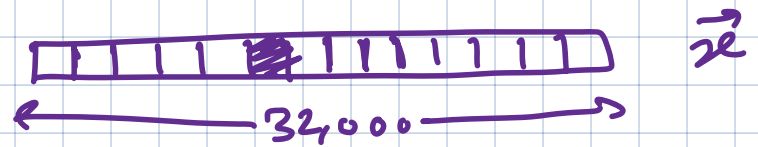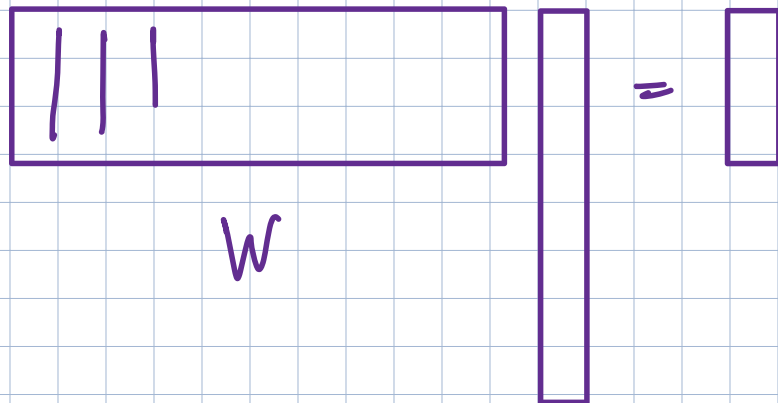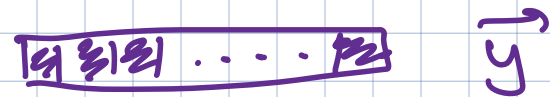


Figure 1: The Transformer - model architecture.

Input: RMS-norm = $\dfrac{1}{\sqrt{32,000}}$



$\vec{x}$
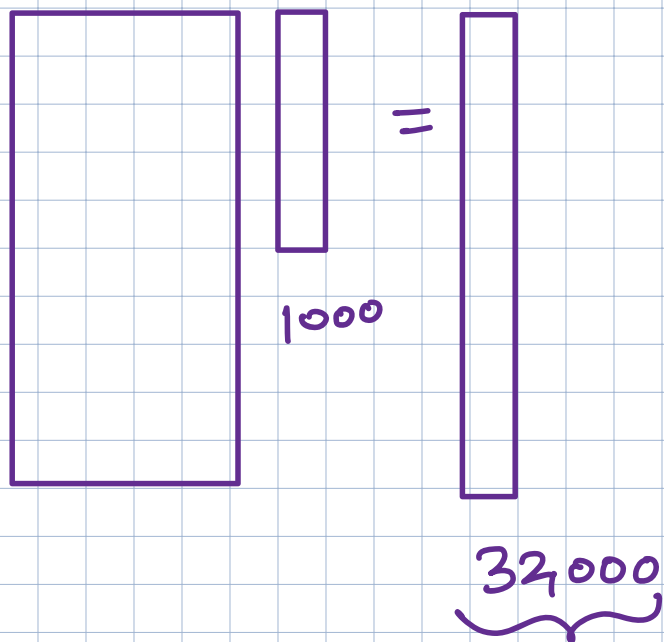
We want RMS-norm for input to attention heads to be 1.

$W\vec{x} = \vec{y}$



$\vec{y}$



W

↑
one hot
— so picks out one column of $W$.

So want each column of $W$ to be RMS norm 1.

For "de"-embedding:

$$\square = \square$$

1000

32000

$$\overset{\curvearrowright}{Q}\,\vec{\alpha} = \vec{\beta}.$$

$$\|\vec{\alpha}\|_{RMS} = 1.$$

want to make sure that entries are not too big.

Every entry of $Q \approx \dfrac{1}{\sqrt{1000}}$

$$q_1^2 \cdot 1 + q_2^2 \cdot 1 + \cdots + q_n^2 \cdot 1 = 1$$

$$\sum q_i^2 = 1 \implies 1000\, q_i^2 = 1 \implies q_i^2 = \frac{1}{1000} \implies q_i \approx \frac{1}{\sqrt{1000}}.$$

# Positional Encoding.

① Attention is all you need.
   Add positional encodings.
      Sinusoidal pattern.

   $\longrightarrow$ No longer used.

② NoPE.

   $\rightarrow$ No position encoding.

③ Learned Relative positional encoding.

Basic attention:     Score $s_i = \dfrac{\langle \vec{q_t}, \vec{k_i} \rangle}{\sqrt{d}}$

$q_t$ at time $t$
$k_i$ is key at pos $i$

$P_{t,i} = \dfrac{e^{s_i}}{\sum\limits_{j} e^{s_j}}$

$$P_{t,i} = \frac{e^{S_i + \boxed{\phantom{xx}}}}{\sum\limits_{j} e^{S_j + \boxed{\phantom{xx}}}}$$

bshift $\rightarrow$ t−i

$$P_{t,i} = \frac{e^{S_i + \boxed{b_{t-i}}}}{\sum\limits_{j} e^{S_j + \boxed{b_{t-j}}}}$$

learn these!

This changes the attention matrix!

# ④ RoPE

Modify attention calculation to capture relative position.

$$\langle \vec{q}_t, \vec{k}_i \rangle$$

$$d=2.$$

$$\begin{bmatrix} q_{t1} \\ q_{t2} \end{bmatrix}$$

$$R_t = \begin{bmatrix} \cos \omega_t & -\sin \omega_t \\ \sin \omega_t & \cos \omega_t \end{bmatrix}$$

$$\langle R_t \vec{q}_t, R_i \vec{k}_i \rangle$$

$$= \vec{q}_t^{\,T} R_t^T R_i \vec{k}_i$$

2d

2d

2d