

# MATH 565: Lecture 29 (04/28/2026)

Today: \* Optimization in transformer networks

## Transformer Networks

**Motivation:** Feedforward NNs map fixed size inputs to outputs. They do not have built-in notions of

1. variable length inputs
2. sequential structure in inputs
3. dependencies between positions in a sequence

Transformer networks: work with a set of tokens as input, and allow each token to interact with every other token.

The goal of this lecture is to give a brief overview of transformer networks and how the results for NNs generalize to their setting.

## The Setup

\* Vocabulary ( $V$ ): set of all distinct tokens the model can process  
 $|V|$ : size of vocabulary  $|V| \approx 30,000 - 50,000$  → for LLMs

\* Embedding, positional encoding: each token is mapped to a point in  $\mathbb{R}^d$  via a learned embedding matrix  $E \in \mathbb{R}^{|V| \times d}$  where  $d$  is the model dimension, e.g.,  $d=768$  in GPT<sub>2</sub>,  $d=12,288$  in GPT-3, etc.

Embedding a token  $t$  is equivalent to looking up row  $t$  of  $E$ .

For a given input of  $T$  tokens, the order/position of the tokens is important. We use position encodings  $\bar{p}_1, \dots, \bar{p}_T \in \mathbb{R}^d$  and add them to the token embeddings to get the token matrix  $X \in \mathbb{R}^{T \times d}$ . Thus, the rows of  $X$  are  $d$ -dimensional representations of the  $T$  input tokens.

$\rightarrow \bar{x}_t = \bar{e}_t + \bar{p}_t, \quad \bar{e}_t = E_t.$   
 $\hookrightarrow t^{\text{th}} \text{ row of } E$

## Self-Attention Mechanism (with one attention head)

**Q:** For a given set of  $T$  tokens, for each token  $t \in T$ , self attention asks "given what  $t$  is about, which other tokens in the sequence are most relevant, and what info should I retrieve from them?"

Answer in three steps

1. Compare  $t$  to all other tokens to get a relevance score.
2. Normalize these scores to a probability using the softmax function.
3. Retrieve a weighted average of information from all tokens using these probabilities as weights.

We project  $\bar{x}_t$  ( $t^{\text{th}}$  row of  $X$  taken as a  $d$ -vector) in three ways:

\* query  $\bar{q}_t = W_Q^T \bar{x}_t$  (what token  $t$  is "looking for")

\* key  $\bar{k}_s = W_K^T \bar{x}_s$  (what token  $s$  "has to offer")

\* value  $\bar{v}_s = W_V^T \bar{x}_s$  (what token  $s$  contributes, if selected)

where  $W_Q, W_K, W_V \in \mathbb{R}^{d \times d_h}$  (matrices to be learned)

The relevance of token  $s$  to token  $t$  is measured by  $\bar{a}_{ts}^T \bar{k}_s / \sqrt{d_h}$  for  $d_h =$  head dimension ( $d_h = d$ , by default)

Softmax then normalizes these scores across all  $s$ .

Def  $\bar{a} = \text{softmax}(\bar{z})$  is defined by setting

$$a_i = \frac{e^{z_i}}{\sum_{k=1}^d e^{z_k}} ; a_i\text{'s are probabilities, as } \sum_{i=1}^d a_i = 1.$$

For 2-d case with  $z_1 = z, z_2 = 0$ , the softmax function reduces to the sigmoid function:  $\frac{e^z}{e^z + 1}$ .

attention weights  $\rightarrow$

$$a_{ts} = \frac{e^{(\bar{a}_t^T \bar{k}_s / \sqrt{d_h})}}{\sum_{s'} e^{(\bar{a}_t^T \bar{k}_s / \sqrt{d_h})}}$$

The output for token  $t$ , called the **attention output**, is

$$\text{Attn}_t = \sum_{s=1}^T a_{ts} \bar{v}_s$$

We write the attention for all tokens together as follows:

query  $Q = XW_Q$   
 key  $K = XW_K$   
 value  $V = XW_V$

Recall,  $W_Q, W_K, W_V \in \mathbb{R}^{d \times d_h}$

The attention output is

$Attn(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_h}} \right) V \in \mathbb{R}^{T \times d_h}$

↗ applied row-wise

$QK^T \in \mathbb{R}^{T \times T}$  has all pairwise compatibility scores between tokens.

$\text{softmax} \left( \frac{QK^T}{\sqrt{d_h}} \right)$  applied row-wise gives the stochastic attention matrix  $A \in \mathbb{R}^{T \times T}$ , and  $AV$  gives the weighted combination of value vectors.

Q: Why divide by  $\sqrt{d_h}$ ? ↗ normal distribution

If  $Q, K$  entries are  $\mathcal{N}(0, 1)$ , then  $QK^T$  entries have variance  $d_h$ . Hence, the  $\sqrt{d_h}$  scaling restores unit variance. Else, we could get  $\approx 0$  Jacobians (i.e., vanishing gradient problem).

### Multi-head Attention

Run  $h$  attention heads in parallel with independent projections and concatenate:

$MHA(X) = \text{cat}(Attn_1, \dots, Attn_h) \cdot W_O \quad | \quad W_O \in \mathbb{R}^{d \times d}$

$\sum_{i=1}^h d_{hi} = d$  typically.

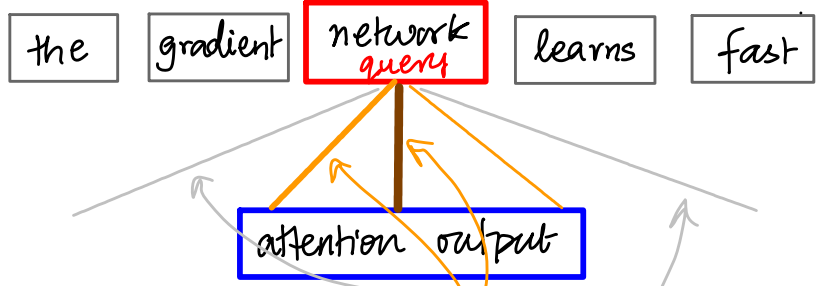
# Why MHA?

"The bank by the river offers high interest rates."

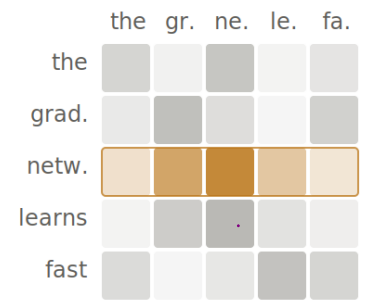
For  $t = \text{"bank"}$  needs attention from "river" and from "interest rates".

## Illustration of Attention Block

$T=5$   $L=3$  ("network")



attention weights: high, low (thickness + dark color)



softmax matrix

## The Transformer Block

One transformer block applies two operations in sequence, each with a residual (skip) connection and a layer normalization (LN).

$$\tilde{x} = LN(x + MHA(x))$$

$$x' = LN(\tilde{x} + FFN(\tilde{x}))$$

Layer Normalization (LN):

$$LN(\bar{z}) = \frac{\bar{z} - \mu}{\sigma} \odot \bar{r} + \bar{\beta} \quad \text{where}$$

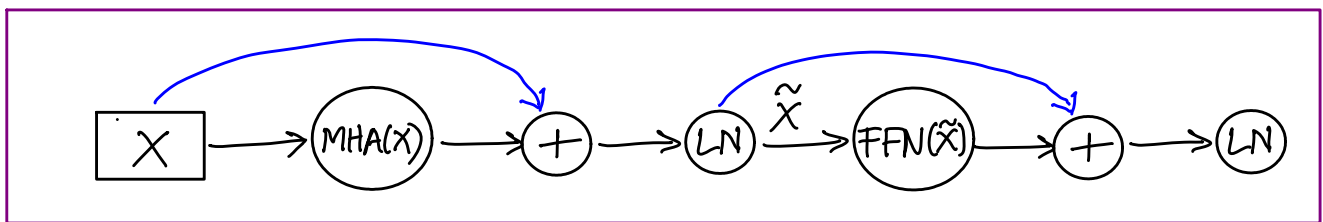
$\mu, \sigma$ : mean and standard deviation of  $\bar{z}$  (a single token)  
 $\bar{r}, \bar{\beta}$ : learned scale and shift parameter vectors.  
 $\bar{z} \in \mathbb{R}^d$

Pointwise Feedforward Network (FFN) function!

$$FFN(\bar{z}) = W_2 \Phi(W_1 \bar{z} + \bar{b}_1) + \bar{b}_2$$

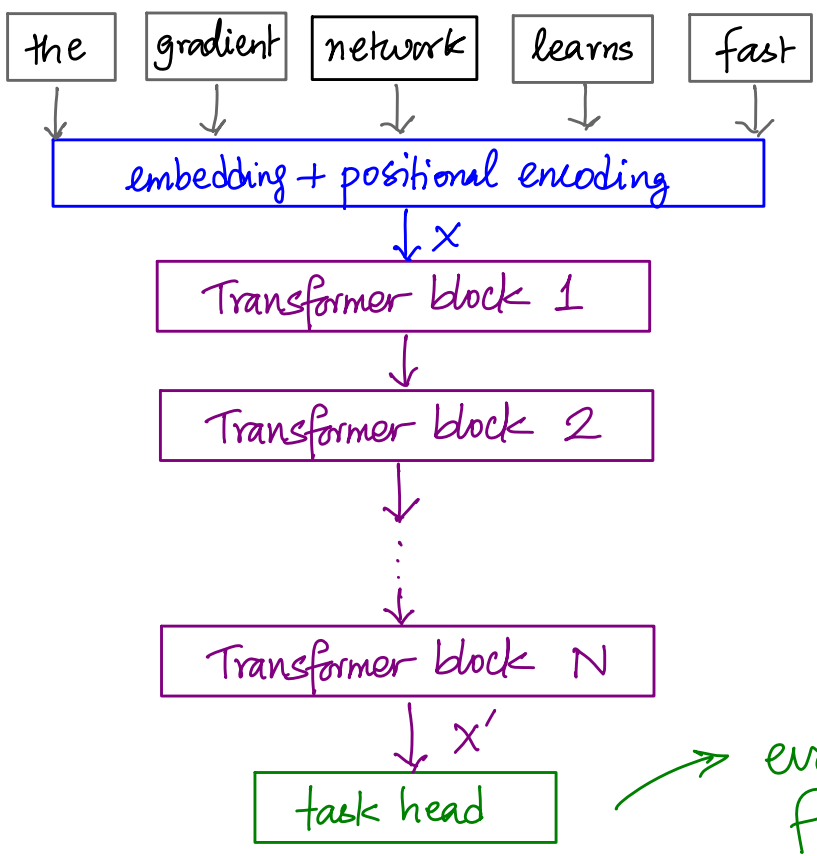
for  $W_1 \in \mathbb{R}^{d_{ff} \times d}$ ,  $W_2 \in \mathbb{R}^{d \times d_{ff}}$ ,  $d_{ff} \approx 4d$  is the feedforward dimension.

→ compare to  $\bar{z}^{(l)} = W^{(l)} \cdot \bar{y}^{(l-1)}$  in NNs, followed by  $\Phi(\bar{z}^{(l)})$ .



### The Transformer Stack

A full transformer network consists of N identical transformer blocks stacked between the embedding and a task specific output head. The same matrix X flows through each block, unchanged in shape.



→ evaluate loss function J

Backpropagation through the transformer network works just like in FeedForward neural networks.

## Theoretical Results

**Theorem 17** (Jacobian of softmax)

Let  $\bar{a} = \text{softmax}(\bar{z}) \in \mathbb{R}^T$  be defined by

$a_i = \frac{e^{\bar{z}_i}}{\sum_k e^{\bar{z}_k}}$ . Then, the softmax is differentiable and its Jacobian is given by

$$\mathcal{D}\bar{a} = \frac{\partial \bar{a}}{\partial \bar{z}} = \text{diag}(\bar{a}) - \bar{a}\bar{a}^T \in \mathbb{R}^{T \times T}$$

Proof

$$\text{For } i=j, \frac{\partial a_i}{\partial z_i} = \frac{(\sum_k e^{\bar{z}_k}) e^{\bar{z}_i} - e^{\bar{z}_i} \cdot e^{\bar{z}_i}}{(\sum_k e^{\bar{z}_k})^2} = a_i - a_i^2 = a_i(1 - a_i).$$

$$\text{For } i \neq j, \frac{\partial a_i}{\partial z_j} = \frac{(\sum_k e^{\bar{z}_k}) \cdot 0 - e^{\bar{z}_i} \cdot e^{\bar{z}_j}}{(\sum_k e^{\bar{z}_k})^2} = -a_i a_j.$$

$$\Rightarrow \mathcal{D}\bar{a} = \frac{\partial \bar{a}}{\partial \bar{z}} = \text{diag}(\bar{a}) - \bar{a}\bar{a}^T \quad \square$$

**Corollary 18** (Gradient through softmax) If  $J$  is a scalar loss function and  $\bar{g} = \frac{\partial J}{\partial \bar{a}}$  is the upstream gradient, then

$$\frac{\partial J}{\partial \bar{z}} = (\text{diag}(\bar{a}) - \bar{a}\bar{a}^T) \bar{g} = \bar{a} \circ \bar{g} - (\bar{a}^T \bar{g}) \bar{a}$$

Remark  $\text{rank}(\mathcal{D}\bar{a}) = T-1$ , since  $(\mathcal{D}\bar{a})\bar{1} = \bar{0}$  as  $(\text{diag}(\bar{a}) - \bar{a}\bar{a}^T)\bar{1} = \bar{a} - \bar{a}(\underbrace{\bar{a}^T\bar{1}}_{=1}) = \bar{0}$  as

$$\bar{a}^T\bar{1} = \sum a_k = 1$$

In contrast,  $\text{diag}(\Phi'(\bar{x}))$  has full rank in FFNNs.

Theorem 19 (gradient bound with residual connections)

For a transformer network with  $N$  blocks, let

$$\bar{s}^{(l)} = \frac{\partial J}{\partial \text{vec}(X^{(l)})}$$

be the vectorized gradient of  $J$

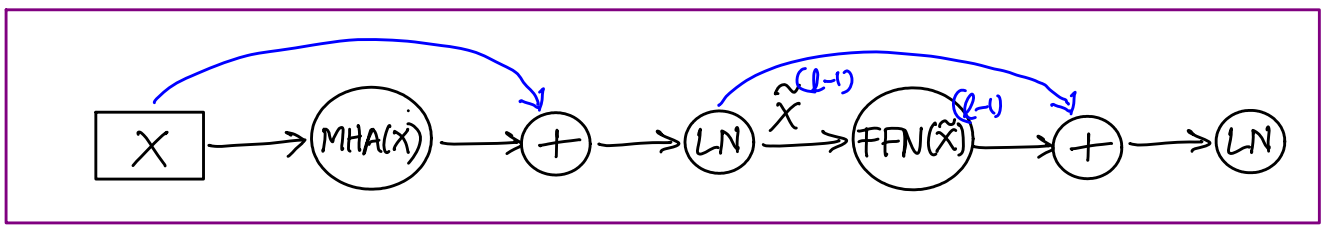
with respect to the token matrix after block  $l$  (we ignore layer normalization for simplicity). Suppose for each block  $l$  the sublayer function  $F^{(l)}$  (MHA or FFN) satisfies  $\|\mathcal{D}_{F^{(l)}}\|_2 \leq \rho$ . Then

$$\|\bar{s}^{(l)}\|_2 \leq (1+\rho)^{N-l} \|\bar{s}^{(N)}\|_2$$

Proof The residual connection acts as

$$X^{(l)} = \tilde{X}^{(l-1)} + F^{(l)}(\tilde{X}^{(l-1)})$$

Recall:



Backward pass through this relation gives (by chain rule at  $\oplus$  node) gives

$$\bar{s}^{(l-1)} = \frac{\partial J}{\partial X^{(l-1)}} = \underbrace{\frac{\partial J}{\partial X^{(l)}}}_{\bar{s}^{(l)}} \cdot \underbrace{\frac{\partial X^{(l)}}{\partial X^{(l-1)}}}_{=I} + \underbrace{\frac{\partial J}{\partial X^{(l)}}}_{\bar{s}^{(l)}} \cdot \underbrace{\frac{\partial X^{(l)}}{\partial F^{(l)}}}_{=I} \cdot \underbrace{\frac{\partial F^{(l)}}{\partial X^{(l-1)}}}_{=\mathcal{D}_{F^{(l)}}}$$

$$\begin{aligned} &= \bar{s}^{(l)} + \bar{s}^{(l)} \cdot I \cdot \mathcal{D}_{F^{(l)}} \\ \Rightarrow \bar{s}^{(l-1)} &= \bar{s}^{(l)} + \mathcal{D}_{F^{(l)}}^T \bar{s}^{(l)} \end{aligned}$$

Taking norms and applying submultiplicativity of spectral norm gives

$$\begin{aligned} \|\bar{s}^{(l-1)}\|_2 &\leq \|\bar{s}^{(l)}\|_2 + \|\mathcal{D}_{F^{(l)}}\|_2 \|\bar{s}^{(l)}\|_2 \\ &\leq (1 + \rho) \|\bar{s}^{(l)}\|_2 \end{aligned}$$

Applying the bound inductively gives

$$\|\bar{s}^{(l-1)}\|_2 \leq (1 + \rho)^{N-l} \|\bar{s}^{(N)}\|_2. \quad \square$$

Compare this result with that of Theorem 16, where  $\sigma \leq 1$  typically, resulting in exponential decay of gradients in early layers. But here, with  $\rho \geq 0$ ,  $(1 + \rho)^{N-l}$  represent a polynomial growth, i.e., there is no decay!

Thus, the residual connections solve the vanishing gradient problem!

## Optimization Challenges for Transformers

\* Learning rate warmup: transformers are highly sensitive to learning rates. In practice, we increase them linearly for several thousand steps and then decay.

\* Adam is mandatory: as plain SGD fails on transformers. The per-parameter adaptive scaling of Adam makes optimization tractable here.

\* gradient clipping:

$$\text{if } \|\nabla J\|_2 > \tau, \quad \nabla J \leftarrow \tau \frac{\nabla J}{\|\nabla J\|_2}$$

as gradients could become large for several reasons.