Lecture 8: Attention and Transformers

Stanford CS231n 10th Anniversary

Lecture 8 - 1

Administrative

- Assignment 2 released yesterday (4/23)
- Project proposals are due tomorrow (4/25)

Stanford CS231n 10th Anniversary

Lecture 8 - 2

Last Time: Recurrent Neural Networks



one to many



many to many

many to many



Stanford CS231n 10th Anniversary

Lecture 8 - 3

Today: Attention + Transformers

Attention: A new primitive that operates on sets of vectors



Transformer: A neural network architecture that uses attention everywhere



Stanford CS231n 10th Anniversary

Lecture 8 - 4

Today: Attention + Transformers

Attention: A new primitive that operates on sets of vectors



Transformers are used everywhere today!

But they developed as an offshoot of RNNs so let's start there

Transformer: A neural network architecture that uses attention everywhere



Stanford CS231n 10th Anniversary

Lecture 8 - 5

Sequence to Sequence with RNNs: Encoder - Decoder

Input: Sequence x_1, \dots, x_T **Output**: Sequence $y_1, \dots, y_{T'}$ A motivating example for today's discussion – machine translation! English \rightarrow Italian

Lecture 8 - 6

April 24, 2025

Encoder:
$$h_t = f_W(x_t, h_{t-1})$$



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

Input: Sequence $x_1, ..., x_T$ **Output**: Sequence $y_1, ..., y_{T'}$

Encoder: $h_t = f_W(x_t, h_{t-1})$ From final hidden state predict: **Initial decoder state** s_0 **Context vector** c (often c=h_T)



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

Stanford CS231n 10th Anniversary



Decoder: $s_t = g_U(y_{t-1}, s_{t-1}, c)$

Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

Stanford CS231n 10th Anniversary

Lecture 8 - 8 April 24, 2025



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

Stanford CS231n 10th Anniversary

Lecture 8 - 9 April 24, 2025



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

Stanford CS231n 10th Anniversary

Lecture 8 - 10 April 24, 2025



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

Stanford CS231n 10th Anniversary

Lecture 8 - **11** April 24, 2025



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

Stanford CS231n 10th Anniversary

Lecture 8 - **12** April 24, 2025

Input: Sequence $x_1, ..., x_T$ **Output**: Sequence $y_1, ..., y_{T'}$

Encoder: $h_t = f_W(x_t, h_{t-1})$ From final hidden state: Initial decoder state s_0





Lecture 8 - 13

April 24, 2025

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Compute (scalar) **alignment scores** $e_{t,i} = f_{att}(s_{t-1}, h_i)$ (f_{att} is a Linear Layer)

April 24, 2025

Lecture 8 - 14



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015



Compute (scalar) **alignment scores** $e_{t,i} = f_{att}(s_{t-1}, h_i)$ (f_{att} is a Linear Layer)

Lecture 8 - 15

Normalize alignment scores to get **attention weights** $0 < a_{t,i} < 1$ $\sum_i a_{t,i} = 1$

April 24, 2025

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015



Compute (scalar) **alignment scores** $e_{t,i} = f_{att}(s_{t-1}, h_i)$ (f_{att} is a Linear Layer) vediamo Normalize alignment scores to get **attention weights** $0 < a_{t,i} < 1$ $\sum_i a_{t,i} = 1$ Compute context vector as **weighted sum of hidden states** $c_t = \sum_i a_{t,i} h_i$

April 24, 2025

Lecture 8 - 16

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Stanford CS231n 10th Anniversary

Lecture 8 - 17 April 24, 2025



Lecture 8 - 18

April 24, 2025

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015





Repeat: Use s_1 to compute new context vector c_2

Compute new alignment scores $e_{2,i}$ and attention weights $a_{2,i}$

April 24, 2025

Lecture 8 - 20

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Stanford CS231n 10th Anniversary

Lecture 8 - 21 April 24, 2025



Lecture 8 - 22

April 24, 2025

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

 S_0

Use a different context vector in each timestep of decoder

- Input sequence not bottlenecked through single vector
- At each timestep of decoder, context vector "looks at" different parts of the input sequence





April 24, 2025

Lecture 8 - 23

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015



Example: English to Visualize attention weights a_{ti} French translation agreement European Economic igned August 1992 <end> Area vas The he S accord sur la zone économique européenne а été signé en août 1992 <end>

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Stanford CS231n 10th Anniversary

Lecture 8 - 24

Example: English to French translation

Input: "The agreement on the European Economic Area was signed in August 1992."

Output: "L'accord sur la zone économique européenne a été signé en août 1992."



April 24, 2025

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Stanford CS231n 10th Anniversary

Lecture 8 - 25





Stanford CS231n 10th Anniversary

Lecture 8 - 27



Stanford CS231n 10th Anniversary

Lecture 8 - 28

Query vectors (decoder RNN states) and data vectors (encoder RNN states) get transformed to

output vectors (Context states).

Each **query** attends to all **data** vectors and gives one **output** vector



There's a general operator hiding here:



Stanford CS231n 10th Anniversary

Lecture 8 - 29

Inputs: Query vector: q [D_Q]



Stanford CS231n 10th Anniversary

Lecture 8 - 30

 $\label{eq:linear} \begin{array}{l} \underline{Inputs} \\ \textbf{Query vector: } \textbf{q} \ [D_Q] \\ \textbf{Data vectors: } \textbf{X} \ [N_X \ x \ D_X] \end{array}$



Stanford CS231n 10th Anniversary

Lecture 8 - 31

Inputs: **Query vector**: $\mathbf{q} [D_Q]$ **Data vectors**: $\mathbf{X} [N_X \times D_X]$



Lecture 8 - 32

April 24, 2025

<u>Computation</u>: Similarities: $e[N_X] e_i = f_{att}(q, X_i)$

Inputs: **Query vector**: $\mathbf{q} [D_Q]$ **Data vectors**: $\mathbf{X} [N_X \times D_X]$



Lecture 8 - 33

April 24, 2025

<u>Computation</u>: Similarities: $e[N_X] e_i = f_{att}(q, X_i)$ Attention weights: $a = softmax(e) [N_X]$

Inputs: **Query vector**: $\mathbf{q} [D_Q]$ **Data vectors**: $\mathbf{X} [N_X \times D_X]$



Lecture 8 - 34

April 24, 2025

<u>Computation</u>: **Similarities**: $e[N_X] e_i = f_{att}(\mathbf{q}, \mathbf{X}_i)$ **Attention weights**: $a = \text{softmax}(e) [N_X]$ **Output vector**: $\mathbf{y} = \sum_i a_i \mathbf{X}_i$ [D_X]

Inputs: **Query vector**: $\mathbf{q} [D_Q]$ **Data vectors**: $\mathbf{X} [N_X \times D_X]$



<u>Computation</u>: **Similarities**: $e[N_X] e_i = f_{att}(\mathbf{q}, \mathbf{X}_i)$ **Attention weights**: $a = softmax(e) [N_X]$ **Output vector**: $\mathbf{y} = \sum_i a_i \mathbf{X}_i$ [D_X]

Let's generalize this!

Stanford CS231n 10th Anniversary

Lecture 8 - 35 April 24, 2025

Inputs: **Query vector**: $\mathbf{q} [D_Q]$ **Data vectors**: $\mathbf{X} [N_X \times D_X]$



<u>Computation</u>: Similarities: $e[N_X] = \mathbf{q} \cdot \mathbf{X}_i$ **Attention weights**: $a = \text{softmax}(e) [N_X]$ **Output vector**: $\mathbf{y} = \sum_i a_i \mathbf{X}_i$ [D_X]

Changes

Lecture 8 - 36

Use dot product for similarity

April 24, 2025
Inputs: **Query vector**: $\mathbf{q} [D_Q]$ **Data vectors**: $\mathbf{X} [N_X \times D_X]$



<u>Computation</u>: **Similarities**: $e[N_X] = \mathbf{q} \cdot \mathbf{X}_i / \sqrt{D_Q}$ **Attention weights**: $a = \operatorname{softmax}(e) [N_X]$ **Output vector**: $\mathbf{y} = \sum_i a_i \mathbf{X}_i [D_X]$

Changes

Use scaled dot product for similarity

Stanford CS231n 10th Anniversary

Lecture 8 - 37 April 24, 2025

<u>Inputs</u>: Query vector: $q [D_Q]$ Data vectors: $X [N_X \times D_X]$

Large similarities will cause softmax to saturate and give vanishing gradients Recall $a \cdot b = |a||b| \cos(angle)$ Suppose that a and b are constant vectors of dimension D Then $|a| = (\sum_i a^2)^{1/2} = a \sqrt{D}$

<u>Computation</u>: **Similarities**: $e[N_X] = \mathbf{q} \cdot \mathbf{X}_i / \sqrt{D_Q}$ **Attention weights**: $a = \operatorname{softmax}(e) [N_X]$ **Output vector**: $\mathbf{y} = \sum_i a_i \mathbf{X}_i [D_X]$



Changes

Use **scaled** dot product for similarity

Stanford CS231n 10th Anniversary

Lecture 8 - 38

Inputs: **Query vector**: $\mathbf{Q} [N_Q \times D_Q]$ **Data vectors**: $\mathbf{X} [N_X \times D_X]$



Computation:

Similarities: $E = QX^T / \sqrt{D_Q} [N_Q \times N_X]$ $E_{ij} = Q_i \cdot X_j / \sqrt{D_Q}$ Attention weights: $A = \text{softmax}(E, \text{dim}=1) [N_Q \times N_X]$ Output vector: $Y = AX [N_Q \times D_X]$ $Y_i = \sum_i A_{ij} X_i$

Changes

- Use scaled dot product for similarity
- Multiple query vectors

Stanford CS231n 10th Anniversary

Lecture 8 - 39

 $\label{eq:linear_line$



Keys:
$$\mathbf{K} = \mathbf{XW}_{\mathbf{K}} [N_X \times D_Q]$$

Values: $\mathbf{V} = \mathbf{XW}_{\mathbf{V}} [D_X \times D_{\mathbf{V}}]$
Similarities: $\mathbf{E} = \mathbf{QK}^T / \sqrt{D_Q} [N_Q \times N_X]$
 $\mathbf{E}_{ij} = \mathbf{Q}_i \cdot \mathbf{K}_j / \sqrt{D_Q}$
Attention weights: $\mathbf{A} = \operatorname{softmax}(\mathbf{E}, \operatorname{dim}=1) [N_Q \times \mathbf{D}_Q]$
Output vector: $\mathbf{Y} = \mathbf{AV} [N_Q \times \mathbf{D}_X]$
 $\mathbf{Y}_i = \sum_j A_{ij} \mathbf{V}_j$



Changes

 N_X]

- Use scaled dot product for similarity
- Multiple query vectors
- Separate key and query

Stanford CS231n 10th Anniversary

Lecture 8 - 40

 $\label{eq:linear_line$

Computation:

Keys: $\mathbf{K} = \mathbf{XW}_{\mathbf{K}} [N_X \times D_Q]$ Values: $\mathbf{V} = \mathbf{XW}_{\mathbf{V}} [D_X \times D_V]$ Similarities: $\mathbf{E} = \mathbf{QK}^T / \sqrt{D_Q} [N_Q \times N_X]$ $\mathbf{E}_{ij} = \mathbf{Q}_i \cdot \mathbf{K}_j / \sqrt{D_Q}$ Attention weights: $\mathbf{A} = \operatorname{softmax}(\mathbf{E}, \operatorname{dim}=1) [N_Q \times N_Y]$ Output vector: $\mathbf{Y} = \mathbf{AV} [N_Q \times D_X]$ $\mathbf{Y}_i = \sum_i A_{ij} \mathbf{V}_i$

X₁



 X_3



Stanford CS231n 10th Anniversary

Lecture 8 - 41

Inputs:

Query vector: Q $[N_Q \times D_Q]$ Data vectors: X $[N_X \times D_X]$ Key matrix: W_K $[D_X \times D_Q]$ Value matrix: W_V $[D_X \times D_V]$

Computation:

Keys: $\mathbf{K} = \mathbf{XW}_{\mathbf{K}} [\mathbf{N}_{\mathbf{X}} \times \mathbf{D}_{\mathbf{Q}}]$ Values: $\mathbf{V} = \mathbf{XW}_{\mathbf{V}} [\mathbf{D}_{\mathbf{X}} \times \mathbf{D}_{\mathbf{V}}]$ Similarities: $\mathbf{E} = \mathbf{Q}\mathbf{K}^{\top} / \sqrt{D_Q} [\mathbf{N}_{\mathbf{Q}} \times \mathbf{N}_{\mathbf{X}}]$ $\mathbf{E}_{ij} = \mathbf{Q}_i \cdot \mathbf{K}_j / \sqrt{D_Q}$ Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1) [\mathbf{N}_{\mathbf{Q}} \times \mathbf{N}_{\mathbf{Y}}]$ Output vector: $\mathbf{Y} = \mathbf{AV} [\mathbf{N}_{\mathbf{Q}} \times \mathbf{D}_{\mathbf{X}}]$ $\mathbf{Y}_i = \sum_i \mathbf{A}_{ij} \mathbf{V}_i$





Stanford CS231n 10th Anniversary

Lecture 8 - 42

Inputs:

Query vector: Q $[N_Q \times D_Q]$ Data vectors: X $[N_X \times D_X]$ Key matrix: W_K $[D_X \times D_Q]$ Value matrix: W_V $[D_X \times D_V]$

Computation:

Keys: $K = XW_{K} [N_{X} \times D_{Q}]$ Values: $V = XW_{V} [D_{X} \times D_{V}]$ Similarities: $E = QK^{T} / \sqrt{D_{Q}} [N_{Q} \times N_{X}]$ $E_{ij} = Q_{i} \cdot K_{j} / \sqrt{D_{Q}}$ Attention weights: $A = \text{softmax}(E, \text{dim}=1) [N_{Q} + M_{Q} \times D_{X}]$ Output vector: $Y = AV [N_{Q} \times D_{X}]$



Stanford CS231n 10th Anniversary

Inputs:

Query vector: Q $[N_Q \times D_Q]$ Data vectors: X $[N_X \times D_X]$ Key matrix: W_K $[D_X \times D_Q]$ Value matrix: W_V $[D_X \times D_V]$

Computation:

Keys: $K = XW_{K} [N_{X} \times D_{Q}]$ Values: $V = XW_{V} [D_{X} \times D_{V}]$ Similarities: $E = QK^{T} / \sqrt{D_{Q}} [N_{Q} \times N_{X}]$ $E_{ij} = Q_{i} \cdot K_{j} / \sqrt{D_{Q}}$ Attention weights: $A = \text{softmax}(E, \text{dim}=1) [N_{Q} \times N_{X}]$ Output vector: $Y = AV [N_{Q} \times D_{X}]$ Softmax normalizes each column: each **query** predicts a distribution over the **keys**



Stanford CS231n 10th Anniversary

Lecture 8 - 44



Stanford CS231n 10th Anniversary

Lecture 8 - 45

Cross-Attention Layer

Inputs:

Query vector: Q $[N_Q \times D_Q]$ Data vectors: X $[N_X \times D_X]$ Key matrix: W_K $[D_X \times D_Q]$ Value matrix: W_V $[D_X \times D_V]$ Each **query** produces one **output**, which is a mix of information in the **data** vectors

Computation:

 $\begin{array}{lll} \mbox{Keys:} & \mbox{K} = \mbox{XW}_{K} & [\mbox{N}_{X} \times \mbox{D}_{Q}] \\ \mbox{Values:} & \mbox{V} = \mbox{XW}_{V} & [\mbox{D}_{X} \times \mbox{D}_{V}] \\ \mbox{Similarities:} & \mbox{E} = \mbox{QK}^{\top} / \sqrt{D_{Q}} & [\mbox{N}_{Q} \times \mbox{N}_{X}] \\ & \mbox{E}_{ij} = \mbox{Q}_{i} \cdot \mbox{K}_{j} / \sqrt{D_{Q}} \\ \mbox{Attention weights:} & \mbox{A} = \mbox{softmax}(\mbox{E}, \mbox{dim} = 1) & [\mbox{N}_{Q} \times \mbox{N}_{X}] \\ \mbox{Output vector:} & \mbox{Y} = \mbox{AV} & [\mbox{N}_{Q} \times \mbox{D}_{X}] \\ & \mbox{Y}_{i} = \sum_{j} A_{ij} \mbox{V}_{j} \end{array}$



Stanford CS231n 10th Anniversary

Lecture 8 - 46

Inputs:

Input vectors: X [N x D_{in}] Key matrix: W_{K} [D_{in} x D_{out}] Value matrix: W_{V} [D_{in} x D_{out}] Query matrix: W_0 [D_{in} x D_{out}]

Computation:

Queries: $\mathbf{Q} = \mathbf{XW}_{\mathbf{Q}}$ [N x D_{out}] Keys: $K = XW_{K}$ [N x D_{out}] Values: $V = XW_V$ [N x D_{out}] Similarities: E = QK^T $/\sqrt{D_0}$ [N x N] $E_{ii} = Q_i \cdot K_i / \sqrt{D_o}$ Attention weights: A = softmax(E, dim=1) [N x N] **Output vector**: $Y = AV [N \times D_{out}]$ $\mathbf{Y}_{i} = \sum_{i} A_{ii} \mathbf{V}_{i}$

Each input produces one output, which is a mix of information from all inputs

Shapes get a little simpler:

- N input vectors, each D_{in}
- Almost always $D_Q = D_V = D_{out}$



Stanford CS231n 10th Anniversary

Lecture 8 - 47

Inputs:

Input vectors: X [N x D_{in}] Key matrix: W_K [$D_{in} x D_{out}$] Value matrix: W_V [$D_{in} x D_{out}$] Query matrix: W_Q [$D_{in} x D_{out}$] Each **input** produces one **output**, which is a mix of information from all **inputs**

Computation:



Stanford CS231n 10th Anniversary

Inputs:

Input vectors: X [N x D_{in}] Key matrix: W_K [$D_{in} x D_{out}$] Value matrix: W_V [$D_{in} x D_{out}$] Query matrix: W_Q [$D_{in} x D_{out}$] Each **input** produces one **output**, which is a mix of information from all **inputs**

Computation:

Queries: $Q = XW_Q$ [N x D_{out}] Keys: $K = XW_K$ [N x D_{out}] Values: $V = XW_V$ [N x D_{out}] Similarities: $E = QK^T / \sqrt{D_Q}$ [N x N] $E_{ij} = Q_i \cdot K_j / \sqrt{D_Q}$ Attention weights: A = softmax(E, dim=1) [N x Output vector: Y = AV [N x D_{out}] $Y_i = \sum_i A_i V_i$



Stanford CS231n 10th Anniversary

Lecture 8 - 49

Inputs:

Input vectors: X [N x D_{in}] Key matrix: W_K [$D_{in} \times D_{out}$] Value matrix: W_V [$D_{in} \times D_{out}$] Query matrix: W_Q [$D_{in} \times D_{out}$] Each **input** produces one **output**, which is a mix of information from all **inputs**

Computation:

Queries: $\mathbf{Q} = \mathbf{XW}_{\mathbf{Q}} [N \times D_{out}]$ Keys: $\mathbf{K} = \mathbf{XW}_{\mathbf{K}} [N \times D_{out}]$ Values: $\mathbf{V} = \mathbf{XW}_{\mathbf{V}} [N \times D_{out}]$ Similarities: $\mathbf{E} = \mathbf{QK}^{\mathsf{T}} / \sqrt{D_Q} [N \times N]$ $E_{ij} = \mathbf{Q}_i \cdot \mathbf{K}_j / \sqrt{D_Q}$ Attention weights: $\mathbf{A} = \operatorname{softmax}(\mathbf{E}, \operatorname{dim}=1) [N \times N]$ Output vector: $\mathbf{Y} = \mathbf{AV} [N \times D_{out}]$ $\mathbf{Y}_i = \sum \mathbf{A}_i \mathbf{V}_i$ Normalize over each column: each **query** computes a distribution over **keys**



Stanford CS231n 10th Anniversary

Lecture 8 - 50

Inputs:

Input vectors: X [N x D_{in}] Key matrix: W_K [$D_{in} x D_{out}$] Value matrix: W_V [$D_{in} x D_{out}$] Query matrix: W_Q [$D_{in} x D_{out}$] Each **input** produces one **output**, which is a mix of information from all **inputs**

Computation:



Stanford CS231n 10th Anniversary

Lecture 8 - 51



Inputs:

Input vectors: X [N x D_{in}] Key matrix: W_{K} [D_{in} x D_{out}] Value matrix: W_{V} [D_{in} x D_{out}] Query matrix: W_{Q} [D_{in} x D_{out}]

Computation:

Stanford CS231n 10th Anniversary

Lecture 8 - 52

Consider permuting inputs:

Queries, keys, and values will be the same but permuted



Input vectors: X [N x D_{in}] Key matrix: W_K [$D_{in} x D_{out}$] Value matrix: W_V [$D_{in} x D_{out}$] Query matrix: W_Q [$D_{in} x D_{out}$]

Computation:

$Product(\rightarrow), Sum(\uparrow)$ V_3 V_2 Softmax(1) K_3 K₁ K_2 Q_2 Q_3 Q_1 X_2 **X**₃ X_1

Stanford CS231n 10th Anniversary

Lecture 8 - 53

Inputs:

Input vectors: X [N x D_{in}] Key matrix: W_{K} [D_{in} x D_{out}] Value matrix: W_{V} [D_{in} x D_{out}] Query matrix: W_{Q} [D_{in} x D_{out}]

Computation:

Consider permuting inputs:

Queries, keys, and values will be the same but permuted

Similarities are the same but permuted



Stanford CS231n 10th Anniversary

Lecture 8 - 54

Inputs:

Input vectors: X [N x D_{in}] Key matrix: W_{K} [D_{in} x D_{out}] Value matrix: W_{V} [D_{in} x D_{out}] Query matrix: W_{Q} [D_{in} x D_{out}]

Computation:

Consider permuting inputs:

Queries, keys, and values will be the same but permuted

Similarities are the same but permuted

Attention weights are the same but permuted



Stanford CS231n 10th Anniversary

Lecture 8 - 55

Inputs:

Computation:

Queries: $Q = XW_{O}$ [N x D_{out}]

Keys: $K = XW_{K}$ [N x D_{out}]

Values: $V = XW_V$ [N x D_{out}]

Similarities: $E = QK^T / \sqrt{D_O} [N \times N]$

Consider permuting inputs:

Queries, keys, and values will be the same but permuted

Similarities are the same but permuted

Attention weights are the same but permuted

Outputs are the same but permuted

 $E_{ii} = Q_i \cdot K_i / \sqrt{D_0}$ Attention weights: A = softmax(E, dim=1) [N x N] **Output vector**: $Y = AV [N \times D_{out}]$ $\mathbf{Y}_{i} = \sum_{i} A_{ii} \mathbf{V}_{i}$

Input vectors: X [N x D_{in}] Key matrix: W_{K} [D_{in} x D_{out}] Value matrix: W_{V} [D_{in} x D_{out}] Query matrix: W_o [D_{in} x D_{out}]

 V_2 A_{3,2} $A_{1,2}$ $A_{2,2}$ Softmax(1) K_3 E_{1.3} E_{2.3} $E_{3.3}$ K₁ E_{3.1} E_{1,1} $E_{2,1}$ K_2 E_{1.2} E_{3,2} E_{2.2} Q_3 Q_2 Q_1 X_2 X_3

 V_3

٧٦

Stanford CS231n 10th Anniversary

Lecture 8 - 56

April 24, 2025

 X_1

 Y_3

A_{2.3}

 $A_{2,1}$

 Y_2

 $Product(\rightarrow), Sum(\uparrow)$

A_{1,3}

A_{1,1}

 $A_{3,3}$

A_{3.1}

Inputs:

Input vectors: X [N x D_{in}] Key matrix: W_{K} [D_{in} x D_{out}] Value matrix: W_{V} [D_{in} x D_{out}] Query matrix: W_{Q} [D_{in} x D_{out}]

Self-Attention is permutation equivariant: $F(\sigma(X)) = \sigma(F(X))$

This means that Self-Attention

Computation:

Y_3 Y_2 $Product(\rightarrow), Sum(\uparrow)$ V_3 A_{1,3} A_{2.3} $A_{3,3}$ V₁ A_{3.1} A_{1.1} $A_{2,1}$ V_2 A_{3,2} $A_{1,2}$ $A_{2,2}$ Softmax(1) K_3 E_{1.3} E_{2.3} $E_{3.3}$ K₁ E_{3.1} E_{1,1} $E_{2,1}$ K_2 E_{1.2} E_{3,2} E_{2.2} Q_2 Q_3 Q_1 X_2 X_3 X_1

Stanford CS231n 10th Anniversary

Lecture 8 - 57

Inputs:

Input vectors: X [N x D_{in}] Key matrix: W_K [$D_{in} \times D_{out}$] Value matrix: W_V [$D_{in} \times D_{out}$] Query matrix: W_Q [$D_{in} \times D_{out}$]

Computation:

Problem: Self-Attention does not know the order of the sequence



Softmax(1)



Stanford CS231n 10th Anniversary

Lecture 8 - 58

Inputs:

Input vectors: X [N x D_{in}] Key matrix: W_K [$D_{in} \times D_{out}$] Value matrix: W_V [$D_{in} \times D_{out}$] Query matrix: W_Q [$D_{in} \times D_{out}$]

Computation:

Problem: Self-Attention does not know the order of the sequence

Solution: Add positional encoding to each input; this is a vector that is a fixed function of the index



Softmax(1)



April 24, 2025

Stanford CS231n 10th Anniversary

Lecture 8 - 59

Masked Self-Attention Layer

Don't let vectors "look ahead" in the sequence

Inputs:

Input vectors: X $[N \times D_{in}]$ Key matrix: W_K $[D_{in} \times D_{out}]$ Value matrix: W_V $[D_{in} \times D_{out}]$ Query matrix: W_Q $[D_{in} \times D_{out}]$

Override similarities with -inf; this controls which inputs each vector is allowed to look at.



Softmax(1)



Computation:

Stanford CS231n 10th Anniversary

Lecture 8 - 60

Masked Self-Attention Layer

Don't let vectors "look ahead" in the sequence

Inputs:

Input vectors: X $[N \times D_{in}]$ Key matrix: W_K $[D_{in} \times D_{out}]$ Value matrix: W_V $[D_{in} \times D_{out}]$ Query matrix: W_Q $[D_{in} \times D_{out}]$

Computation:

Queries:
$$Q = XW_Q$$
 [N x Dout]where you want to next wordKeys: $K = XW_K$ [N x Dout]next wordValues: $V = XW_V$ [N x Dout]similarities: $E = QK^T / \sqrt{D_Q}$ [N x N]Similarities: $E = QK^T / \sqrt{D_Q}$ [N x N] $E_{ij} = Q_i \cdot K_j / \sqrt{D_Q}$ Attention weights: $A = \text{softmax}(E, \text{dim}=1)$ [N x N]Output vector: $Y = AV$ [N x Dout] $Y_i = \sum_j A_{ij}V_j$

Override similarities with -inf; this controls which inputs each vector is allowed to look at.

Used for language modeling where you want to predict the next word



Softmax(1)

April 24, 2025



Stanford CS231n 10th Anniversary

Lecture 8 - 61

Run H copies of Self-Attention in parallel

Inputs:

```
Input vectors: X [N x D<sub>in</sub>]
Key matrix: W_{K} [D<sub>in</sub> x D<sub>out</sub>]
Value matrix: W_{V} [D<sub>in</sub> x D<sub>out</sub>]
Query matrix: W_{Q} [D<sub>in</sub> x D<sub>out</sub>]
```

Computation:



Stanford CS231n 10th Anniversary

Lecture 8 - 62

Run H copies of Self-Attention in parallel

Inputs:

Input vectors: X [N x D_{in}] Key matrix: W_K [D_{in} x D_{out}] Value matrix: W_V [D_{in} x D_{out}] Query matrix: W_Q [D_{in} x D_{out}]

Computation:

```
Queries: \mathbf{Q} = \mathbf{XW}_{\mathbf{Q}} [\mathbf{N} \times \mathbf{D}_{out}]

Keys: \mathbf{K} = \mathbf{XW}_{\mathbf{K}} [\mathbf{N} \times \mathbf{D}_{out}]

Values: \mathbf{V} = \mathbf{XW}_{\mathbf{V}} [\mathbf{N} \times \mathbf{D}_{out}]

Similarities: \mathbf{E} = \mathbf{Q}\mathbf{K}^{T} / \sqrt{D_{Q}} [\mathbf{N} \times \mathbf{N}]

\mathbf{E}_{ij} = \mathbf{Q}_{i} \cdot \mathbf{K}_{j} / \sqrt{D_{Q}}

Attention weights: \mathbf{A} = \operatorname{softmax}(\mathbf{E}, \operatorname{dim}=1) [\mathbf{N}

Output vector: \mathbf{Y} = \mathbf{AX} [\mathbf{N} \times \mathbf{D}_{out}]

\mathbf{Y}_{i} = \sum_{i} A_{ij} \mathbf{V}_{i}
```

H = 3 independent self-attention layers (called heads), each with their own weights



Stanford CS231n 10th Anniversary

Lecture 8 - 63

Run H copies of Self-Attention in parallel

Inputs:

Input vectors: X [N x D_{in}] Key matrix: W_K [$D_{in} \times D_{out}$] Value matrix: W_V [$D_{in} \times D_{out}$] Query matrix: W_Q [$D_{in} \times D_{out}$]

Computation:

Queries: $\mathbf{Q} = \mathbf{XW}_{\mathbf{Q}}$ [N x D_{out}] Keys: $\mathbf{K} = \mathbf{XW}_{\mathbf{K}}$ [N x D_{out}] Values: $\mathbf{V} = \mathbf{XW}_{\mathbf{V}}$ [N x D_{out}] Similarities: $\mathbf{E} = \mathbf{QK}^{\mathsf{T}} / \sqrt{D_Q}$ [N x N] $\mathbf{E}_{ij} = \mathbf{Q}_i \cdot \mathbf{K}_j / \sqrt{D_Q}$ Attention weights: A = softmax(E, di

Output vector: Y = AX [N x D_{out}]

Stack up the H independent outputs for each input X

H = 3 independent self-attention layers (called heads), each with their own weights











Stanford CS231n 10th Anniversary

Lecture 8 - 64

Run H copies of Self-Attention in parallel

Inputs:

Input vectors: X [N x D_{in}] Key matrix: W_{K} [$D_{in} \times D_{out}$] Value matrix: W_{V} [$D_{in} \times D_{out}$] Query matrix: W_{Q} [$D_{in} \times D_{out}$]

Computation:

Queries: $\mathbf{Q} = \mathbf{XW}_{\mathbf{Q}}$ [N x D_{out}] Keys: $\mathbf{K} = \mathbf{XW}_{\mathbf{K}}$ [N x D_{out}] Values: $\mathbf{V} = \mathbf{XW}_{\mathbf{V}}$ [N x D_{out}] Similarities: $\mathbf{E} = \mathbf{Q}\mathbf{K}^{T} / \sqrt{D_{Q}}$ [N x N] $\mathbf{E}_{ij} = \mathbf{Q}_{i} \cdot \mathbf{K}_{j} / \sqrt{D_{Q}}$ Attention weights: A = softmax(E, di

Output vector: Y = AX [N x D_{out}]

 $\mathbf{Y}_{i} = \sum_{j} A_{ij} \mathbf{V}_{j}$

Output projection fuses data from each head

Stack up the H independent outputs for each input X

H = 3 independent self-attention layers (called heads), each with their own weights









E_{2,2}

Stanford CS231n 10th Anniversary

Lecture 8 - 65

Run H copies of Self-Attention in parallel

Inputs:

Input vectors: X [N x D] Key matrix: W_{K} [D x HD_H] Value matrix: W_{V} [D x HD_H] Query matrix: W_{Q} [D x HD_H] Output matrix: W_{O} [HD_H x D]

Computation:

Queries:
$$Q = XW_Q$$
 [H x N x D_H]

Keys:
$$K = XW_{K}$$
 [H x N x D_H]

Values:
$$V = XW_V [H \times N \times D_H]$$

Similarities: $E = \mathbf{QK^T} / \sqrt{D_Q} [H \times N \times N]$

```
Attention weights: A = \text{softmax}(E, \text{dim}=1) [H x N x N]
Head outputs: Y = AV [H x N x D_H] => [N x HD_H]
Outputs: O = YW_O [N x D]
```

Each of the H parallel layers use a qkv dim of D_{H} = "head dim"

Usually $D_H = D / H$, so inputs and outputs have the same dimension





Stanford CS231n 10th Anniversary

Lecture 8 - 66

Run H copies of Self-Attention in parallel

Inputs:

```
Input vectors: X [N x D]
Key matrix: W_{K} [D x HD<sub>H</sub>]
Value matrix: W_{V} [D x HD<sub>H</sub>]
Query matrix: W_{Q} [D x HD<sub>H</sub>]
Output matrix: W_{O} [HD<sub>H</sub> x D]
```

Computation:

```
Queries: Q = XW_Q [H x N x D<sub>H</sub>]
```

```
Keys: K = XW_{K} [H x N x D<sub>H</sub>]
Values: V = XW_{V} [H x N x D<sub>H</sub>]
```

Similarities: E = $\mathbf{Q}\mathbf{K}^{\mathsf{T}} / \sqrt{D_0} [\mathbf{H} \times \mathbf{N} \times \mathbf{N}]$

```
Attention weights: A = softmax(E, dim=1) [H x N x N]
Head outputs: Y = AV [H x N x D_H] => [N x HD_H]
Outputs: O = YW_O [N x D]
```

In practice, compute all H heads in parallel using batched matrix multiply operations.

Used everywhere in practice.





Stanford CS231n 10th Anniversary

Lecture 8 - 67

Inputs:

```
Input vectors: X [N x D]
Key matrix: W_{K} [D x HD<sub>H</sub>]
Value matrix: W_{V} [D x HD<sub>H</sub>]
Query matrix: W_{Q} [D x HD<sub>H</sub>]
Output matrix: W_{O} [HD<sub>H</sub> x D]
```

Computation:

```
Queries: Q = XW_Q [H x N x D<sub>H</sub>]

Keys: K = XW_K [H x N x D<sub>H</sub>]

Values: V = XW_V [H x N x D<sub>H</sub>]

Similarities: E = QK^T / \sqrt{D_Q} [H x N x N]

Attention weights: A = \text{softmax}(E, \text{dim}=1) [H x N x N]

Head outputs: Y = AV [H x N x D<sub>H</sub>] => [N x HD<sub>H</sub>]

Outputs: O = YW_O [N x D]
```

Stanford CS231n 10th Anniversary

Lecture 8 - 68

Inputs:

Input vectors: X [N x D] Key matrix: W_{K} [D x HD_H] Value matrix: W_{V} [D x HD_H] Query matrix: W_{Q} [D x HD_H] Output matrix: W_{O} [HD_H x D]

Computation:

Queries: $Q = XW_Q$ [H x N x D_H] Keys: $K = XW_K$ [H x N x D_H] Values: $V = XW_V$ [H x N x D_H]

Similarities: $\mathbf{E} = \mathbf{Q}\mathbf{K}^{\mathsf{T}} / \sqrt{D_Q} [\mathbf{H} \times \mathbf{N} \times \mathbf{N}]$

```
Attention weights: A = \text{softmax}(E, \text{dim}=1) [H \times N \times N]
Head outputs: Y = AV [H \times N \times D_H] => [N \times HD_H]
Outputs: O = YW_O [N \times D]
```

<u>QKV Projection</u> [N x D] [D x 3HD_H] => [N x 3HD_H] Split and reshape to get Q, K, V each of shape [H x N x D_H]

Stanford CS231n 10th Anniversary

Lecture 8 - 69

Inputs:

Input vectors: X [N x D] Key matrix: W_{K} [D x HD_H] Value matrix: W_{V} [D x HD_H] Query matrix: W_{Q} [D x HD_H] Output matrix: W_{O} [HD_H x D]

Computation:

Queries: $\mathbf{Q} = \mathbf{XW}_{\mathbf{Q}}$ [H x N x D_H]

Keys: $K = XW_{K}$ [H x N x D_H]

Values: $V = XW_V [H \times N \times D_H]$

Similarities: $\mathbf{E} = \mathbf{Q}\mathbf{K}^{\top} / \sqrt{D_Q} [\mathbf{H} \times \mathbf{N} \times \mathbf{N}]$

Attention weights: $A = \text{softmax}(E, \text{dim}=1) [H \times N \times N]$ Head outputs: $Y = AV [H \times N \times D_H] => [N \times HD_H]$ Outputs: $O = YW_O [N \times D]$

<u>QKV Projection</u> [N x D] [D x 3HD_H] => [N x 3HD_H] Split and reshape to get Q, K, V each of shape [H x N x D_H]

2. <u>QK Similarity</u> [H x N x D_H] [H x N x D_H] => [H x N x N]

Stanford CS231n 10th Anniversary

Lecture 8 - 70

Inputs:

Input vectors: X [N x D] Key matrix: W_{K} [D x HD_H] Value matrix: W_{V} [D x HD_H] Query matrix: W_{Q} [D x HD_H] Output matrix: W_{O} [HD_H x D]

Computation:

Queries: $Q = XW_Q$ [H x N x D_H]

```
Keys: K = XW_{K} [H \times N \times D_{H}]
Values: V = XW_{V} [H \times N \times D_{H}]
```

Similarities: $\mathbf{E} = \mathbf{Q}\mathbf{K}^{\mathsf{T}} / \sqrt{D_O} [\mathbf{H} \times \mathbf{N} \times \mathbf{N}]$

```
Attention weights: A = \text{softmax}(F, \text{dim}=1) [H x N x N]
Head outputs: Y = AV [H x N x D<sub>H</sub>] => [N x HD<sub>H</sub>]
```

- <u>QKV Projection</u>

 [N x D] [D x 3HD_H] => [N x 3HD_H]
 Split and reshape to get Q, K, V each of shape [H x N x D_H]
- 2. <u>QK Similarity</u> [H x N x D_H] [H x N x D_H] => [H x N x N]
- 3. <u>V-Weighting</u> [H x N x N] [H x D x D_H] => [H x N x D_H] Reshape to [N x HD_H]

Stanford CS231n 10th Anniversary

Lecture 8 - 71

Inputs:

Input vectors: X [N x D] Key matrix: W_{K} [D x HD_H] Value matrix: W_{V} [D x HD_H] Query matrix: W_{Q} [D x HD_H] Output matrix: W_{O} [HD_H x D]

Computation:

```
Queries: Q = XW_Q [H x N x D<sub>H</sub>]
```

Keys: $K = XW_{K} [H \times N \times D_{H}]$ Values: $V = XW_{V} [H \times N \times D_{H}]$

Similarities: $\mathbf{E} = \mathbf{Q}\mathbf{K}^{\mathsf{T}} / \sqrt{D_0} [\mathbf{H} \times \mathbf{N} \times \mathbf{N}]$

```
Attention weights: A = \text{softmax}(E, \text{dim}=1) [H x N x N]
```

```
Head outputs: Y = AV [H \times N \times D_H] => [N \times HD_H]
Outputs: O = YW_O [N \times D]
```

<u>QKV Projection</u>

 [N x D] [D x 3HD_H] => [N x 3HD_H]
 Split and reshape to get Q, K, V each of shape [H x N x D_H]

- 2. <u>QK Similarity</u> [H x N x D_H] [H x N x D_H] => [H x N x N]
- 3. <u>V-Weighting</u> [H x N x N] [H x D x D_H] => [H x N x D_H] Reshape to [N x HD_H]
- 4. <u>Output Projection</u> [N x HD_H] [HD_H x D] => [N x D]

Stanford CS231n 10th Anniversary

Lecture 8 - 72
Inputs:

Input vectors: X [N x D] Key matrix: W_{K} [D x HD_H] Value matrix: W_{V} [D x HD_H] Query matrix: W_{Q} [D x HD_H] Output matrix: W_{O} [HD_H x D]

Computation:

Queries: $Q = XW_Q$ [H x N x D_H]

Keys: $K = XW_{K} [H \times N \times D_{H}]$ Values: $V = XW_{V} [H \times N \times D_{H}]$

Similarities: $\mathbf{E} = \mathbf{Q}\mathbf{K}^{\mathsf{T}} / \sqrt{D_0} [\mathbf{H} \times \mathbf{N} \times \mathbf{N}]$

```
Attention weights: A = \text{softmax}(E, \text{dim}=1) [H \times N \times N]
Head outputs: Y = AV [H \times N \times D_H] => [N \times HD_H]
Outputs: O = YW_O [N \times D]
```

- <u>QKV Projection</u> [N x D] [D x 3HD_H] => [N x 3HD_H] Split and reshape to get Q, K, V each of shape [H x N x D_H]
- 2. <u>QK Similarity</u> [H x N x D_H] [H x N x D_H] => [H x N x N]
- 3. <u>V-Weighting</u> [H x N x N] [H x D x D_H] => [H x N x D_H] Reshape to [N x HD_H]
- 4. <u>Output Projection</u> [N x HD_H] [HD_H x D] => [N x D]

Q: How much <u>compute</u> does this take as the number of vectors N increases?

Stanford CS231n 10th Anniversary

Lecture 8 - 73

Inputs:

Input vectors: X [N x D] Key matrix: W_{K} [D x HD_H] Value matrix: W_{V} [D x HD_H] Query matrix: W_{Q} [D x HD_H] Output matrix: W_{O} [HD_H x D]

Computation:

Queries: $Q = XW_Q$ [H x N x D_H]

Keys: $K = XW_{K}$ [H x N x D_H] Values: $V = XW_{V}$ [H x N x D_H]

Similarities: $\mathbf{E} = \mathbf{Q}\mathbf{K}^{\mathsf{T}} / \sqrt{D_0} [\mathbf{H} \times \mathbf{N} \times \mathbf{N}]$

```
Attention weights: A = \text{softmax}(E, \text{dim}=1) [H \times N \times N]
Head outputs: Y = AV [H \times N \times D_H] => [N \times HD_H]
Outputs: O = YW_O [N \times D]
```

<u>QKV Projection</u>

 [N x D] [D x 3HD_H] => [N x 3HD_H]
 Split and reshape to get Q, K, V each of shape [H x N x D_H]

 <u>QK Similarity</u>

 [H x N x D_H] [H x N x D_H] => [H x N x N]
 Weighting

Lecture 8 - 74

Q: How much <u>compute</u> does this take as the number of vectors N increases? **A**: $O(N^2)$

April 24, 2025

Inputs:

Input vectors: X [N x D] Key matrix: W_{K} [D x HD_H] Value matrix: W_{V} [D x HD_H] Query matrix: W_{Q} [D x HD_H] Output matrix: W_{O} [HD_H x D]

Computation:

Queries: $Q = XW_Q$ [H x N x D_H]

Keys: $K = XW_{K} [H \times N \times D_{H}]$ Values: $V = XW_{V} [H \times N \times D_{H}]$

Similarities: $\mathbf{E} = \mathbf{Q}\mathbf{K}^{\mathsf{T}} / \sqrt{D_0} [\mathbf{H} \times \mathbf{N} \times \mathbf{N}]$

```
Attention weights: A = \text{softmax}(E, \text{dim}=1) [H \times N \times N]
Head outputs: Y = AV [H \times N \times D_H] => [N \times HD_H]
Outputs: O = YW_O [N \times D]
```

- <u>QKV Projection</u> [N x D] [D x 3HD_H] => [N x 3HD_H] Split and reshape to get Q, K, V each of shape [H x N x D_H]
- 2. <u>QK Similarity</u> [H x N x D_H] [H x N x D_H] => [H x N x N]
- 3. <u>V-Weighting</u> [H x N x N] [H x D x D_H] => [H x N x D_H] Reshape to [N x HD_H]
- 4. <u>Output Projection</u> [N x HD_H] [HD_H x D] => [N x D]

Q: How much <u>memory</u> does this take as the number of vectors N increases?

Stanford CS231n 10th Anniversary

Lecture 8 - 75

Inputs:

Input vectors: X [N x D] Key matrix: W_{K} [D x HD_H] Value matrix: W_{V} [D x HD_H] Query matrix: W_{Q} [D x HD_H] Output matrix: W_{O} [HD_H x D]

Computation:

Queries: $Q = XW_Q$ [H x N x D_H]

Keys: $K = XW_{K}$ [H x N x D_H] Values: $V = XW_{V}$ [H x N x D_H]

Similarities: $\mathbf{E} = \mathbf{Q}\mathbf{K}^{\mathsf{T}} / \sqrt{D_0} [\mathbf{H} \times \mathbf{N} \times \mathbf{N}]$

```
Attention weights: A = \text{softmax}(E, \text{dim}=1) [H \times N \times N]
Head outputs: Y = AV [H \times N \times D_H] => [N \times HD_H]
Outputs: O = YW_O [N \times D]
```

 <u>QKV Projection</u> [N x D] [D x 3HD_H] => [N x 3HD_H] Split and reshape to get Q, K, V each of shape [H x N x D_H]
 <u>QK Similarity</u> [H x N x D_H] [H x N x D_H] => [H x N x N]
 <u>V-Weighting</u>

 $[H \times N \times N] [H \times D \times D_{H}] => [H \times N \times D_{H}]$ Reshape to [N × HD_H]

Q: How much <u>memory</u> does this take as the number of vectors N increases? **A**: $O(N^2)$

Stanford CS231n 10th Anniversary

Lecture 8 - 76

Inputs:

Input vectors: X [N x D] Key matrix: W_{K} [D x HD_H] Value matrix: W_{V} [D x HD_H] Query matrix: W_{Q} [D x HD_H] Output matrix: W_{O} [HD_H x D]

Computation:

Queries: $Q = XW_Q$ [H x N x D_H]

Keys: $K = XW_{K} [H \times N \times D_{H}]$ Values: $V = XW_{V} [H \times N \times D_{H}]$

Similarities: $\mathbf{E} = \mathbf{Q}\mathbf{K}^{\mathsf{T}} / \sqrt{D_0} [\mathbf{H} \times \mathbf{N} \times \mathbf{N}]$

```
Attention weights: A = \text{softmax}(E, \text{dim}=1) [H \times N \times N]
Head outputs: Y = AV [H \times N \times D_H] => [N \times HD_H]
Outputs: O = YW_O [N \times D]
```

If N=100K, H=64 then HxNxN attention weights take 1.192 TB! GPUs don't have that much memory... 1. QKV Projection $[N \times D] [D \times 3HD_{H}] => [N \times 3HD_{H}]$ Split and reshape to get Q, K, V each of shape $[H \times N \times D_{H}]$ QK Similarity $[H \times N \times D_{H}]$ $[H \times N \times D_{H}] => [H \times N \times N]$ 3. V-Weighting $[H \times N \times N] [H \times D \times D_{H}] => [H \times N \times D_{H}]$ Reshape to $[N \times HD_{H}]$ 4. Output Projection

4. <u>Output Projection</u> [N x HD_H] [HD_H x D] => [N x D]

> **Q:** How much <u>memory</u> does this take as the number of vectors N increases? **A**: $O(N^2)$

Stanford CS231n 10th Anniversary

Lecture 8 - 77

Inputs:

Input vectors: X [N x D] Key matrix: W_{κ} [D x HD_H] Value matrix: W_V [D x HD_H] Query matrix: W_0 [D x HD_H] Output matrix: W₀ [HD_H x D] possible

Flash Attention algorithm computes 2+3 at the same time without storing the full attention matrix!

1.

Makes large N

Computation:

Queries: $Q = XW_0$ [H x N x D_H]

Keys:
$$K = XW_{K} [H \times N \times D_{H}]$$

values: $V = X W_V$ [H X N X D_H]

```
Similarities: E = QK^T / \sqrt{D_O} [H \times N \times N]
```

```
Attention weights: A = softmax(E, dim=1) [H x N x N]
Head outputs: Y = AV [H \times N \times D_{H}] => [N \times HD_{H}]
Outputs: O = YW_0 [N x D]
```

If N=100K, H=64 then HxNxN attention weights take 1.192 TB! GPUs don't have that much memory... **QKV** Projection $[N \times D] [D \times 3HD_{H}] => [N \times 3HD_{H}]$ Split and reshape to get Q, K, V each of shape [H x N x D_{H}]

- 2. QK Similarity $[H \times N \times D_{H}] [H \times N \times D_{H}] => [H \times N \times N]$
- 3. V-Weighting $[H \times N \times N] [H \times D \times D_{H}] => [H \times N \times D_{H}]$ Reshape to $[N \times HD_{H}]$
- 4. **Output Projection** $[N \times HD_{H}] [HD_{H} \times D] => [N \times D]$

Lecture 8 - 78

Q: How much <u>memory</u> does this take as the number of vectors N increases? A: O(N) with Flash Attention

Dao et al, "FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness", 2022

April 24, 2025

Stanford CS231n 10th Anniversary

Lecture 8 - 79

Recurrent Neural Network



Works on 1D ordered sequences

(+) Theoretically good at long sequences: O(N) compute and memory for a sequence of length N
(-) Not parallelizable. Need to compute hidden states sequentially

Stanford CS231n 10th Anniversary

Lecture 8 - 80



Recurrent Neural Network

Works on 1D ordered sequences

(+) Theoretically good at long sequences: O(N) compute and memory for a sequence of length N
(-) Not parallelizable. Need to compute hidden states sequentially Convolution



Works on N-dimensional grids

(-) Bad for long sequences: need to stack many layers to build up large receptive fields
(+) Parallelizable, outputs can be computed in parallel

Stanford CS231n 10th Anniversary

Lecture 8 - 81



Recurrent Neural Network

Works on 1D ordered sequences

(+) Theoretically good at long sequences: O(N) compute and memory for a sequence of length N
(-) Not parallelizable. Need to compute hidden states sequentially Convolution



Works on N-dimensional grids

(-) Bad for long sequences: need to stack many layers to build up large receptive fields

(+) Parallelizable, outputs can be computed in parallel

Self-Attention



Works on sets of vectors

(+) Great for long sequences; each output depends directly on all inputs
(+) Highly parallel, it's just 4 matmuls
(-) Expensive: O(N²) compute, O(N) memory for sequence of length N

Stanford CS231n 10th Anniversary

Lecture 8 - 82



Transformer Block

Input: Set of vectors x



April 24, 2025

Vaswani et al, "Attention is all you need," NeurIPS 2017

Stanford CS231n 10th Anniversary

Lecture 8 - 84

Transformer Block

Input: Set of vectors x





Vaswani et al, "Attention is all you need," NeurIPS 2017

Stanford CS231n 10th Anniversary

Lecture 8 - 85

Transformer Block

Input: Set of vectors x



Vaswani et al, "Attention is all you need," NeurIPS 2017

Stanford CS231n 10th Anniversary

Lecture 8 - 86

Transformer Block

Input: Set of vectors x

Recall Layer Normalization:



Ba et al, 2016



Vaswani et al, "Attention is all you need," NeurIPS 2017

Stanford CS231n 10th Anniversary

Lecture 8 - 87

Transformer Block

Input: Set of vectors x

Usually a two-layer MLP; classic setup is D => 4D => D

Also sometimes called FFN (Feed-Forward Network)

April 24, 2025



Lecture 8 - 88

Vaswani et al, "Attention is all you need," NeurIPS 2017

Transformer Block

Input: Set of vectors x



Vaswani et al, "Attention is all you need," NeurIPS 2017

Stanford CS231n 10th Anniversary

Lecture 8 - <u>89</u>

Transformer Block

Input: Set of vectors x



Vaswani et al, "Attention is all you need," NeurIPS 2017

Stanford CS231n 10th Anniversary

Lecture 8 - <u>90</u>

Transformer Block

Input: Set of vectors x Output: Set of vectors y

Self-Attention is the only interaction between vectors

LayerNorm and MLP work on each vector independently

Highly scalable and parallelizable, most of the compute is just 6 matmuls:

4 from Self-Attention 2 from MLP

Vaswani et al, "Attention is all you need," NeurIPS 2017

Stanford CS231n 10th Anniversary



Lecture 8 - 91

Transformer Block

Input: Set of vectors x Output: Set of vectors y

Self-Attention is the only interaction between vectors

LayerNorm and MLP work on each vector independently

Highly scalable and parallelizable, most of the compute is just 6 matmuls:

4 from Self-Attention 2 from MLP

Vaswani et al, "Attention is all you need," NeurIPS 2017

Stanford CS231n 10th Anniversary

A **Transformer** is just a stack of identical Transformer blocks!

They have not changed much since 2017... but have gotten a lot bigger



April 24, 2025

Lecture 8 - 92

Transformer Block

Input: Set of vectors x Output: Set of vectors y

Self-Attention is the only interaction between vectors

LayerNorm and MLP work on each vector independently

Highly scalable and parallelizable, most of the compute is just 6 matmuls:

4 from Self-Attention 2 from MLP

Vaswani et al, "Attention is all you need," NeurIPS 2017

Stanford CS231n 10th Anniversary

A **Transformer** is just a stack of identical Transformer blocks!

They have not changed much since 2017... but have gotten a lot bigger

<u>Original</u>: [Vaswani et al, 2017] 12 blocks, D=1024, H=16, N=512 213M params



Lecture 8 - 93

Transformer Block

Input: Set of vectors x Output: Set of vectors y

Self-Attention is the only interaction between vectors

LayerNorm and MLP work on each vector independently

Highly scalable and parallelizable, most of the compute is just 6 matmuls:

4 from Self-Attention 2 from MLP

A **Transformer** is just a stack of identical Transformer blocks!

They have not changed much since 2017... but have gotten a lot bigger

<u>Original</u>: [Vaswani et al, 2017] 12 blocks, D=1024, H=16, N=512 213M params

<u>GPT-2</u>: [Radford et al, 2019] 48 blocks, D=1600, H=25, N=1024 1.5B params



Vaswani et al, "Attention is all you need," NeurIPS 2017

Stanford CS231n 10th Anniversary

Lecture 8 - 94

Transformer Block

Input: Set of vectors x **Output**: Set of vectors y

Self-Attention is the only interaction between vectors

LayerNorm and MLP work on each vector independently

Highly scalable and parallelizable, most of the compute is just 6 matmuls:

4 from Self-Attention 2 from MLP

A **Transformer** is just a stack of identical Transformer blocks!

They have not changed much since 2017... but have gotten a lot bigger

<u>Original</u>: [Vaswani et al, 2017] 12 blocks, D=1024, H=16, N=512 213M params

<u>GPT-2</u>: [Radford et al, 2019] 48 blocks, D=1600, H=25, N=1024 1.5B params

<u>GPT-3</u>: [Brown et al, 2020] 96 blocks, D=12288, H=96, N=2048 175B params



Vaswani et al, "Attention is all you need," NeurIPS 2017

Stanford CS231n 10th Anniversary

Lecture 8 - 95

Learn an <u>embedding matrix</u> at the start of the model to convert words into vectors.

Given vocab size V and model dimension D, it's a lookup table of shape $[V \times D]$



Stanford CS231n 10th Anniversary

Lecture 8 - 96

Learn an <u>embedding matrix</u> at the start of the model to convert words into vectors.

Given vocab size V and model dimension D, it's a lookup table of shape $[V \times D]$

Use masked attention inside each transformer block so each token can only see the ones before it



Stanford CS231n 10th Anniversary

Lecture 8 - 97

Learn an <u>embedding matrix</u> at the start of the model to convert words into vectors.

Given vocab size V and model dimension D, it's a lookup table of shape $[V \times D]$

Use masked attention inside each transformer block so each token can only see the ones before it

At the end, learn a <u>projection matrix</u> of shape $[D \times V]$ to project each D-dim vector to a V-dim vector of scores for each element of the vocabulary.



Stanford CS231n 10th Anniversary

Lecture 8 - 98

Learn an <u>embedding matrix</u> at the start of the model to convert words into vectors.

Given vocab size V and model dimension D, it's a lookup table of shape [V x D]

Use masked attention inside each transformer block so each token can only see the ones before it

At the end, learn a <u>projection matrix</u> of shape $[D \times V]$ to project each D-dim vector to a V-dim vector of scores for each element of the vocabulary.

Train to predict next token using softmax + cross-entropy loss



Stanford CS231n 10th Anniversary

Lecture 8 - 99



Input image: e.g. 224x224x3

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

Stanford CS231n 10th Anniversary

Lecture 8 - 100



Input image: e.g. 224x224x3





Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021 Break into patches e.g. 16x16x3

Stanford CS231n 10th Anniversary

Lecture 8 - 101



Input image: e.g. 224x224x3



Break into patches
 Flatten and apply a linear
 e.g. 16x16x3
 transform 768 => D

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

Stanford CS231n 10th Anniversary

Lecture 8 - 102



Input image: e.g. 224x224x3



Break into patchesFlatten and apply a lineare.g. 16x16x3transform 768 => D

Q: Any other way to describe this operation?

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

Stanford CS231n 10th Anniversary

Lecture 8 - <u>103</u>



Input image: e.g. 224x224x3



Q: Any other way to describe this operation?

A: 16x16 conv with stride 16, 3 input channels, D output channels

Lecture 8 - 104

April 24, 2025

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021 Break into patches e.g. 16x16x3

chesFlatten and apply a linear<3</td>transform 768 => D



Input image: e.g. 224x224x3



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

Break into patches e.g. 16x16x3

transform 768 => D

D-dim vector per patch are the input vectors to the Transformer

April 24, 2025

Lecture 8 - 105



Input image: e.g. 224x224x3



Use positional encoding to tell the transformer the 2D position of each patch

April 24, 2025

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021 Break into patches e.g. 16x16x3

s Flatten and apply a linear transform 768 => D D-dim vector per patch are the input vectors to the Transformer

Lecture 8 - 106



Input image: e.g. 224x224x3



Don't use any masking; each image patch can look at all other image patches

Use positional encoding to tell the transformer the 2D position of each patch

April 24, 2025

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021 Break into patches e.g. 16x16x3

Flatten and apply a linear transform 768 => D D-dim vector per patch are the input vectors to the Transformer

Lecture 8 - 107



Input image: e.g. 224x224x3



Transformer gives an output vector per patch

Don't use any masking; each image patch can look at all other image patches

Use positional encoding to tell the transformer the 2D position of each patch

April 24, 2025

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

Break into patches e.g. 16x16x3

 Flatten and apply a linear transform 768 => D D-dim vector per patch are the input vectors to the Transformer

Lecture 8 - 108
Vision Transformers (ViT)



Input image: e.g. 224x224x3

Dosovitskiy et al, "An Image is Worth

16x16 Words: Transformers for Image

Recognition at Scale", ICLR 2021



transform 768 => D

Average pool NxD vectors to 1xD, apply a linear layer D=>C to predict class scores



Transformer gives an output vector per patch

Don't use any masking; each image patch can look at all other image patches

Use positional encoding to tell the transformer the 2D position of each patch

April 24, 2025

D-dim vector per patch are the input vectors to the Transformer

Lecture 8 - 109

Stanford CS231n 10th Anniversary

e.g. 16x16x3

Tweaking Transformers

The Transformer architecture has not changed much since 2017.

But a few changes have become common:



Stanford CS231n 10th Anniversary

Lecture 8 - 110

Pre-Norm Transformer

Layer normalization is outside the residual connections

Kind of weird, the model can't actually learn the identify function



Baevski & Auli, "Adaptive Input Representations for Neural Language Modeling", arXiv 2018

Stanford CS231n 10th Anniversary

Lecture 8 - 111

Pre-Norm Transformer

Layer normalization is outside the residual connections

Kind of weird, the model can't actually learn the identify function

Solution: Move layer normalization before the Self-Attention and MLP, inside the residual connections. Training is more stable.



Baevski & Auli, "Adaptive Input Representations for Neural Language Modeling", arXiv 2018

Stanford CS231n 10th Anniversary

Lecture 8 - 112

RMSNorm

Replace Layer Normalization with Root-Mean-Square Normalization (RMSNorm)

Input: x [shape D] **Output**: y [shape D] **Weight**: γ [shape D]

$$y_{i} = \frac{x_{i}}{RMS(x)} * \gamma_{i}$$
$$RMS(x) = \sqrt{\varepsilon + \frac{1}{N} \sum_{i=1}^{N} x_{i}^{2}}$$

Training is a bit more stable



Zhang and Sennrich, "Root Mean Square Layer Normalization", NeurIPS 2019

Stanford CS231n 10th Anniversary

Lecture 8 - 113

SwiGLU MLP

Classic MLP:

Input: X [N x D] Weights: W₁ [D x 4D] W₂ [4D x D] Output: Y = σ (XW₁)W₂ [N x D]



Lecture 8 - 114

April 24, 2025

Shazeer, "GLU Variants Improve Transformers", 2020

SwiGLU MLP

Classic MLP:

Input: X [N x D] Weights: W₁ [D x 4D] W₂ [4D x D] Output: Y = σ (XW₁)W₂ [N x D]

SwiGLU MLP:

Input: X [N x D] Weights: W₁, W₂ [D x H] W₃ [H x D] Output: $Y = (\sigma(XW_1) \odot XW_2)W_3$

Setting H = 8D/3 keeps same total params

Shazeer, "GLU Variants Improve Transformers", 2020

Stanford CS231n 10th Anniversary



Lecture 8 - 115

SwiGLU MLP

Classic MLP:

Input: X [N x D] Weights: W₁ [D x 4D] W₂ [4D x D] Output: Y = σ (XW₁)W₂ [N x D]

SwiGLU MLP:

Input: X [N x D] Weights: W₁ , W₂ [D x H] W₃ [H x D] Output: $Y = (\sigma(XW_1) \odot XW_2)W_3$ We offer no explanation as to why these architectures seem to work; we attribute their success, as all else, to divine benevolence.



Shazeer, "GLU Variants Improve Transformers", 2020

same total params

Setting H = 8D/3 keeps

Stanford CS231n 10th Anniversary

Lecture 8 - <u>116</u>

Mixture of Experts (MoE)

Learn E separate sets of MLP weights in each block; each MLP is an *expert*

 $W_1: [D x 4D] \implies [E x D x 4D]$ $W_2: [4D x D] \implies [E x 4D x D]$



April 24, 2025

Shazeer et al, "Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer", 2017

Stanford CS231n 10th Anniversary

Lecture 8 - 117

Mixture of Experts (MoE)

Learn E separate sets of MLP weights in each block; each MLP is an *expert*

 $W_1: [D x 4D] => [E x D x 4D]$ $W_2: [4D x D] => [E x 4D x D]$

Each token gets *routed* to A < E of the experts. These are the *active experts*.

Increases params by E, But only increases compute by A



April 24, 2025

Lecture 8 - 118

Shazeer et al, "Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer", 2017

Mixture of Experts (MoE)

Learn E separate sets of MLP weights in each block; each MLP is an *expert*

 $W_1: [D x 4D] => [E x D x 4D]$ $W_2: [4D x D] => [E x 4D x D]$

Each token gets *routed* to A < E of the experts. These are the *active experts*.

Increases params by E, But only increases compute by A

All of the biggest LLMs today (e.g. GPT4o, GPT4.5, Claude 3.7, Gemini 2.5 Pro, etc) almost certainly use MoE and have > 1T params; but they don't publish details anymore



April 24, 2025

Lecture 8 - 119

Tweaking Transformers

The Transformer architecture has not changed much since 2017.

But a few changes have become common:

- **Pre-Norm:** Move normalization inside residual
- **RMSNorm**: Different normalization layer
- SwiGLU: Different MLP architecture
- Mixture of Experts (MoE): Learn E different MLPs, use A < E of them per token. Massively increase params, modest increase to compute cost.



Stanford CS231n 10th Anniversary

Lecture 8 - 120

Summary: Attention + Transformers

Attention: A new primitive that operates on sets of vectors



Transformers are the backbone of all large Al models today!

Used for language, vision, speech, ...

Transformer: A neural network architecture that uses attention everywhere



Stanford CS231n 10th Anniversary

Lecture 8 - 121 April 24, 2025

Next Time: Detection, Segmentation, Visualization

Lecture 8 - 122

April 24, 2025