Introduction to Computer Vision



Lecture 14 Self-Attention & Transformer

Prof. He Wang

Embodied Perception and InteraCtion Lab

Spring 2025



Logistics

- Assignment 4 (Point Cloud Learning, Detection & RNN)
 - Released on 5/24
 - Due on 6/8 11:59PM

Attention + Transformer

The slides are borrowed and modified from Stanford CS 231N.

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



Recap: Sequence to Sequence with RNNs

Decoder: $s_t = g_U(y_{t-1}, s_{t-1}, c)$ Input: Sequence $x_1, \dots x_T$ **Output**: Sequence $y_1, \ldots, y_{T'}$ vediamo From final hidden state predict: **Y**1 **Initial decoder state** s₀ **Encoder:** $h_t = f_W(x_t, h_{t-1})$ **Context vector** c (often c=h_T) h₂ h₁ hз h4 **S**0 **S**1 С **X**1 **X**2 **X**3 **X**4 **y**0 [START] the skv we see

Recap: Sequence to Sequence with RNNs



Recap: Sequence to Sequence with RNNs



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





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





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** $0 < a_{t,i} < 1$ $\sum_i a_{t,i} = 1$

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR2015





Repeat: Use s_1 to compute new context vector c_2

Compute new alignment scores $e_{2,1}$ and attention weights $a_{2,1}$





Use a different context vector in each timestep of decoder

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







<end>

<pue>

992

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."





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."



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."



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:



Inputs: Query vector: q[D_Q]



Inputs:

Query vector: q[D_Q] Data vectors: X[N_X x D_X]



Inputs:

Query vector: q[D_Q] Data vectors: X[N_X x D_X]



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

Inputs:

Query vector: q[D_Q] Data vectors: X[N_X x D_X]



 $\label{eq:computation} \begin{array}{l} \hline \textbf{Computation} \\ \textbf{Similarities:} e [N_X] e_i = f_{att}(\textbf{q}, \textbf{X}_i) \\ \textbf{Attention weights:} a = softmax(e) [N_X] \end{array}$

Inputs:

Query vector: q[D_Q] Data vectors: X[N_X x D_X]



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

Inputs:

Query vector: q[D_Q] Data vectors: X[N_X x D_X]



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

Let's generalize this!

Inputs:

Query vector: q[D_X] Data vectors: X[N_X x D_X]



Computation: Similarities: $e[N_X] e_i = q \cdot X_i$ Attention weights: $a = softmax(e) [N_X]$ Output vector: $y = \sum_i a_i X_i [D_x]$

- Use dot product for similarity

Inputs: Query vector: q[D_X] Data vectors: X[N_X x D_X]



$$\begin{array}{l} \hline \textbf{Computation}:\\ \textbf{Similarities}: e [N_X] & e_i = \textbf{q} \cdot \textbf{X}_i \; \sqrt{D_X}\\ \textbf{Attention weights}: a = softmax(e) \; [N_X]\\ \textbf{Output vector}: \textbf{y} = \sum_i a_i \; \textbf{X}_i \; [D_X] \end{array}$$

Changes

- Use **scaled** dot product for similarity

Inputs: Query vector: q[D_X] Data vectors: X[N_X x 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 $|\mathbf{a}| = (\sum a^2)^{1/2} = \mathbf{a}\sqrt{D}$

Computation: Similarities: $e[N_X] e_i = q \cdot X_i \sqrt{D_X}$ Attention weights: $a = \text{softmax}(e)[N_X]$ Output vector: $y = \sum_i a_i X_i [D_x]$



Changes

- Use scaled dot product for similarity



Computation:



 $E_{ij} = \mathbf{Q}_i \cdot \mathbf{X}_j / \sqrt{D_X}$ we as Attention weights: A = softmax(E, dim=1) [N_Q x N_X] Output vector: Y = AX[N_Q x D_X] Y_i = \sum_i A_{ij} X_i

Similarities: $E = QX^T / \sqrt{D_X} [N_Q \times N_X]$

Changes

- Use scaled dot product for similarity
- Multiple query vectors

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: $\mathbf{V} = \mathbf{XW}_{\mathbf{V}} [\mathbf{N}_{\mathbf{X}} \times \mathbf{D}_{\mathbf{V}}]$ Similarities: $\mathbf{E} = \mathbf{Q}\mathbf{K}^{\mathsf{T}} / \sqrt{D_{Q}} [\mathbf{N}_{\mathbf{Q}} \times \mathbf{N}_{\mathbf{X}}]$ $\mathbf{E}_{\mathbf{i}\mathbf{i}} = \mathbf{Q}_{\mathbf{i}} \cdot \mathbf{K}_{\mathbf{i}} [/D_{Q}]$

Attention weights: A = softmax(E, dim=1) $[N_Q \times N_X]$ Output vector: Y = AV $[N_Q \times D_V]$

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



Changes

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

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





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: $V = XW_V [N_X \times D_V]$ 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: $Y = AV[N_Q \times D_V]$ $Y_i = \sum_i A_{ii}V_i$





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



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


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: $W_V[D_X \times D_V]$

Computation:

Keys: $K = XW_{K} [N_{X} \times D_{Q}]$ Values: $V = XW_{V} [N_{X} \times D_{V}]$ 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: $Y = AV[N_{Q} \times D_{V}]$ $Y_{i} = \sum_{i} A_{ii} V_{i}$



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:

Keys: $K = XW_{K} [N_{X} \times D_{Q}]$ Values: $V = XW_{V} [N_{X} \times D_{V}]$ 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: $Y = AV [N_{Q} \times D_{V}]$ $Y_{i} = \sum_{i} A_{ii} V_{i}$



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: $W_{Q}[D_{n} x D_{out}]$

Computation:

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

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 $D_Q = D_V = D_{out}$



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

Each input produces

one output, which is

a mix of information

from all inputs



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

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



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

Each input produces

one output, which is

a mix of information

from all inputs

Normalize over each column: each query computes a distribution over keys



Compute output vectors as linear combinations of value vectors

Each input produces

one output, which is

a mix of information

from all inputs

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:



Y₃

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:



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:

Consider permuting inputs:

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



Consider permuting inputs:

Queries, keys, and values

permuted

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:
$$Q = XW_Q$$
 [N x D_{out}]
Keys: $K = XW_K$ [N x D_{out}]
Values: $V = XW_V$ [N x D_{out}]
Similarities: $E = QK^T / \sqrt{O_Q}$ [N x N]
 $E_{ij} = Q_i \cdot K_j / \sqrt{O_Q}$

Attention weights: A = softmax(E, dim=1) [N x N] **Output vector**: $Y = AV [N \times D_{out}]$ $Y_i = \sum_i A_{ii} V_i$

 $Product(\rightarrow), Sum(\uparrow)$ will be the same but permuted V_3 Similarities are the same but V_2 Softmax(1) K₃ E3, 1 E1.1 E2.1 K_1 E1.2 E2,2 E3.2 K_2 E3,3 E1,3 E2,3 Q_2 Q₁ Q_3 X_2 X_1 **X**3

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:
$$\mathbf{Q} = \mathbf{XW}_{\mathbf{Q}}$$
 [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: $\mathbf{E} = \mathbf{QK}^{\mathsf{T}} / \sqrt{O_{\mathbf{Q}}}$ [N x N]
 $\mathbf{E}_{ij} = \mathbf{Q}_{i} \cdot \mathbf{K}_{j} / \sqrt{O_{\mathbf{Q}}}$

Attention weights: A = softmax(E, dim=1) [N x N] Output vector: Y = AV [N x D_{out}] Y_i = $\sum_{j} A_{ij} V_{j}$

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



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:

Attention weights: A = softmax(E, dim=1) [N x N] Output vector: Y = AV [N x D_{out}] Y_i = $\sum_{j} A_{ij} V_{j}$

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



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:

Attention weights: A = softmax(E, dim=1) [N x N] Output vector: Y = AV [N x D_{out}] Y_i = $\sum_{j} A_{ij} V_{j}$

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

This means that Self-Attention works on **sets of vectors**



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:

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



Softmax(1)



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:

Attention weights: A = softmax(E, dim=1) [N x N] Output vector: Y = AV [N x D_{out}] Y_i = $\sum_{j} A_{ij} V_{j}$

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)



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:

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







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 D_{out}] Keys: $K = XW_K$ [N x D_{out}] Values: $V = XW_V$ [N x D_{out}] Similarities: $E = QK^T / \sqrt{O_Q}$ [N x N] $E_{ij} = Q_i \cdot K_j / \sqrt{D_Q}$

Attention weights: A = softmax(E, dim=1) [N x N] Output vector: $Y = AV [N \times D_{out}]$ $Y_i = \sum_j A_{ij}V_j$

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)



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



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: $\mathbf{Q} = \mathbf{XW}_{\mathbf{Q}}$ [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: $\mathbf{E} = \mathbf{QK}^{\mathsf{T}} / \sqrt{O_{\mathbf{Q}}}$ [N x N] $\mathbf{E}_{ij} = \mathbf{Q}_{i} \cdot \mathbf{K}_{j} / \sqrt{D_{Q}}$

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

Attention weights: A = softmax(E, dim=1) [N x N] Output vector: Y = AX [N x D_{out}]





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: $\mathbf{Q} = \mathbf{XW}_{\mathbf{Q}}$ [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: $\mathbf{E} = \mathbf{QK}^{T} / \sqrt{O_{\mathbf{Q}}}$ [N x N] $\mathbf{E}_{ij} = \mathbf{Q}_{i} \cdot \mathbf{K}_{j} / \sqrt{D_{\mathbf{Q}}}$

Attention weights: A = softmax(E, dim=1) [N x N] Output vector: Y = AX [N x D_{out}]

 $\mathbf{Y}_i = \sum_j A_{ij} \mathbf{V}_j$

Stack up the H independent outputs for each input X

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





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: $\mathbf{Q} = \mathbf{XW}_{\mathbf{Q}}$ [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: $\mathbf{E} = \mathbf{QK}^{\top} / \sqrt{O_{\mathbf{Q}}}$ [N x N] $\mathbf{E}_{ij} = \mathbf{Q}_{i} \cdot \mathbf{K}_{j} / \sqrt{D_{\mathbf{Q}}}$

Attention weights: A = softmax(E, dim=1) [N x N] Output vector: Y = AX [N x D_{out}]

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





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:



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





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_H] Keys: $K = XW_K$ [H x N x D_H] Values: $V = XW_V$ [H x N x D_H] Similarities: $E = QK^T / \sqrt{O_Q}$ [H x N x 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 = YW₀ [N x D]

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

Used everywhere in practice.





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{O_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]
```

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: $E = QK^T / \sqrt{O_Q} [H \times N \times N]$ Attention weights: $A = softmax(E, dim=1) [H \times N \times N]$ 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]

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

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: $O = YW_0$ [N x 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>

 $[H \times N \times D_H] [H \times D_H \times N] \Longrightarrow [H \times N \times N]$

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: $E = QK^T / \sqrt{O_Q}$ [H x N x 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: $O = YW_Q$ [N x 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 D_H x N] => [H x N x N]

 $[H \times N \times N] [H \times N \times D_{H}] \Rightarrow [H \times N \times D_{H}]$ Reshape to [N x HD_H]

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] 4. Values: $V = XW_V$ [H x N x D_H] Similarities: $E = QK^T / \sqrt{O_Q}$ [H x N x 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]

<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]

 QK Similarity

 $[H \times N \times D_H] [H \times D_H \times N] \Longrightarrow [H \times N \times N]$

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]HKeys: K = XW_K [H x N x D_H]4. <u>9</u>Values: V = XW_V [H x N x D_H][Similarities: E = $QK^T / \sqrt{O_Q}$ [H x N x 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]

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

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

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] 4. Values: $V = XW_V$ [H x N x D_H] Similarities: $E = QK^T / \sqrt{O_Q}$ [H x N x 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]

 <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 D_H x N] => [H x N x N]
 <u>V-Weighting</u> [H x N x N] [H x N x D_H] => [H x N x D_H] Reshape to [N x HD_H]
 <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? **A**: $O(N^2)$

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<sub>Q</sub> [D x HD<sub>H</sub>]
Output matrix: W<sub>O</sub> [HD<sub>H</sub> x D]
```

Computation:

Queries: $Q = XW_Q$ [H x N x D_H]FKeys: K = XW_K [H x N x D_H]4. <u>Q</u>Values: V = XW_V [H x N x D_H][Similarities: E = $QK^T / \sqrt{O_Q}$ [H x N x 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]

- <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. $\frac{QK \text{ Similarity}}{[H \times N \times D_H] [H \times D_H \times N]} => [H \times N \times N]$
- 3. <u>V-Weighting</u>
 [H x N x N] [H x N 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?

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<sub>Q</sub> [D x HD<sub>H</sub>]
Output matrix: W<sub>O</sub> [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>] 4.

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

Similarities: E = QK^T / \sqrt{O_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]
```

- <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> $[H \times N \times D_H] [H \times D_H \times N] => [H \times N \times N]$
- 3. <u>V-Weighting</u>
 [H x N x N] [H x N 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? **A**: $O(N^2)$

If N=100K, H=64 then HxNxN attention weights take 1.192 TB! GPUs don't have that much memory...

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: WQ [D x HD_H] Output matrix: WO [HD_H 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>] 4.

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

Similarities: E = QK^T / \sqrt{O_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]
```

- <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> $[H \times N \times D_H] [H \times D_H \times N] => [H \times N \times N]$
- 3. <u>V-Weighting</u>
 [H x N x N] [H x N 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? **A**: $O(N^2)$

If N=100K, H=64 then HxNxN attention weights take 1.192 TB! GPUs don't have that much memory...

<u>Inputs</u> : Input vectors: X [N x D] Key matrix: W _K [D x HD _H] Value matrix: W _V [D x HD _H]	Flash Attention algorithm computes 2+3 at the same time without storing the full attention matrix!	 <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]
Query matrix: W _Q [D x HD _H] Output matrix: W _O [HD _H x D] <u>Computation</u> :	2 Makes large N possible 3	 <u>QK Similarity</u> [H x N x D_H] [H x D_H x N] => [H x N x N] <u>V-Weighting</u> [H x N x N] [H x N x D_H] => [H x N x D_H Pershape to [N x HD_H]
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: $E = QK^T / \sqrt{O_Q}$ [H	H] A A X N X N] (E dim=1) [H X N X	 A Contract of the contract of the
Head outputs : $Y = AV [H \times N \times D_H] => [N \times HD_H]$ Outputs : $O = YW_0 [N \times D]$ as the number of vectors N increases A: O(N) with Flash Attention		

Three Ways of Processing Sequences

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

Three Ways of Processing Sequences

Recurrent Neural Network

Convolution



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



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
Three Ways of Processing Sequences





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



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

	Y 1	¥2	Yз
[Product(→),目Sum(个)		
	t		
V3 -	A _{1,3}	A _{2,3}	A _{3,3}
V2 -	$A_{\scriptscriptstyle 1,2}$	A _{2,2}	A _{3,2}
V1 -	A _{1,1}	A _{2,1}	A _{3,1}
	Softmax(↑)		
	•		
K3 -	E _{1,3}	E _{2,3}	E _{3,3}
K2 -	E _{1,2}	E _{2,2}	E _{3,2}
К1 —	E _{1,1}	E _{2,1}	E _{3,1}
	+	+	1
	Q1	Q2	Q3
	+	-	
	Х1	X2	Хз

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

Three Ways of Processing Sequences



Transformer Block

Input: Set of vectors x



Transformer Block

Input: Set of vectors x

All vectors interact through (multiheaded) Self-Attention



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

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) MLP independently MLP MLP MLP MLP on each vector Layer normalization Layer Normalization normalizes all vectors **Residual connection** Self-Attention All vectors interact through (multiheaded) Self-Attention X_1 X_2 **X**₃ **X**₄

Transformer Block

Input: Set of vectors x



Transformer Block

Input: Set of vectors x



Transformer Block

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

Self-Attention is the only interaction between vectors

LayerNormand 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



Transformer Block

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

Self-Attention is the only interaction between vectors

LayerNormand 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



Transformer Block

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

Self-Attention is the only interaction between vectors

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



Transformer Block

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

Self-Attention is the only interaction between vectors

LayerNormand 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>Origina</u>l: [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



Transformer Block

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

Self-Attention is the only interaction between vectors

LayerNormand 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>Origina</u>l: [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 etal, 2020] 96 blocks, D=12288, H=96, N=2048 175B params



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 $[V \times D]$



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



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.



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.

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





Input image: e.g. 224x224x3



Input image: e.g. 224x224x3



Break into patches e.g. 16x16x3



Input image: e.g. 224x224x3



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



Input image: e.g. 224x224x3



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

Q: Any other way to describe this operation?



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

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



Input image: e.g. 224x224x3





Input image: e.g. 224x224x3



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

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

D-dim vector perpatch are the input vectors to the Transformer



Input image: e.g. 224x224x3



e.g. 16x16x3

transform 768 => D

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

Don't use any

masking; each

look at all other

image patches

Use positional

encoding to tell

the transformer

the 2D position

of each patch



Input image: e.g. 224x224x3



the transformer the 2D position of each patch

Transformer gives an output

Don't use any

masking; each

look at all other

image patches

Use positional

encoding to tell

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

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



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

e.g. 16x16x3

transform 768 => D

D-dim vector perpatch are the input vectors to the Transformer

Tweaking Transformers

The Transformer architecture has not changed much since 2017.

But a few changes have become common:



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

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

RMSNorm

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

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



SwiGLU MLP

Classic MLP:

Input: $X [N \times D]$ Weights: $W_1 [D \times 4D]$ $W_2 [4D \times D]$ Output: $Y = \delta(XW_1) W_2 [N \times D]$



SwiGLU MLP

Classic MLP:

Input: $X [N \times D]$ Weights: $W_1 [D \times 4D]$ $W_2 [4D \times D]$ Output: $Y = \delta(XW_1) W_2 [N \times D]$

SwiGLU MLP:

Input: $X [N \times D]$ Weights: $W_1, W_2[D \times H]$ $W_3[D \times H]$ Output: $Y = (\delta(XW_1) \odot XW_2)W_3$

Setting H = 8D/3 keeps same total params



SwiGLU MLP

Classic MLP:

Input: $X [N \times D]$ Weights: $W_1 [D \times 4D]$ $W_2 [4D \times D]$ Output: $Y = \delta(XW_1) W_2 [N \times D]$

SwiGLU MLP:

```
Input: X [N \times D]
Weights: W_1, W_2[D \times H]
W_3[D \times H]
Output:
Y = (\delta(XW_1) \odot XW_2)W_3
```

Setting H = 8D/3 keeps same total params We offer no explanation as to why these architectures seem to work; we attribute their success, as all else, to divine benevolence.



Mixture of Experts (MoE)

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

W₁: [D x 4D] => [E x D x 4D] W₂: [4D x D] => [E x 4D x D]


Mixture of Experts (MoE)

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

W₁: [D x 4D] => [E x D x 4D] W₂: [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



Mixture of Experts (MoE)

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

W₁: [D x 4D] => [E x D x 4D] W₂: [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



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 themper token. Massively increase params, modest increase to compute cost.



Summary: Attention + Transformers

Attention: A new primitive that operates on sets of vectors



Transformers are the backbone of all large AI models today!

Used for language, vision, speech, ...

Transformer: A neural network architecture that uses attention everywhere



Large Multi-modal Models

Large Multi-modal Models (BLIP)





ITC: Image-Text Contrastive Learning; ITM: Image-Text Match Learning; LM: Language Modelling

Li, Junnan, et al. "Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation." International Conference on Machine Learning. PMLR, 2022.

Large Multi-modal Models (BLIP)

Learning Framework of BLIP



Large Multi-modal Models (BLIP)

• Evaluation of the effect of the captioner (C) and filter (F) for dataset bootstrapping.

Pre-train	Boot	tstrap	Vision	Retrieval-F	T (COCO)	Retrieval-Z	ZS (Flickr)	Caption-FT	C (COCO)	Caption-ZS	S (NoCaps)
dataset	C	F	backbone	TR@1	IR@1	TR@1	IR@1	B@4	CIDEr	CIDEr	SPICE
COCO+VG +CC+SBU (14M imgs)	$\begin{array}{c} \mathbf{X} \\ \mathbf{X} \\ \mathbf{V}_B \\ \mathbf{V}_B \end{array}$	× ✓ B ×	ViT-B/16	78.4 79.1 79.7 80.6	60.7 61.5 62.0 63.1	93.9 94.1 94.4 94.8	82.1 82.8 83.6 84.9	38.0 38.1 38.4 38.6	127.8 128.2 128.9 129.7	102.2 102.7 103.4 105.1	13.9 14.0 14.2 14.4
COCO+VG +CC+SBU +LAION	× ✓ _B ✓ _L	× ✓ _B ✓ _L	ViT-B/16	79.6 81.9 81.2	62.0 64.3 64.1	94.3 96.0 96.0	83.6 85.0 85.5	38.8 39.4 39.7	130.1 131.4 133.3	105.4 106.3 109.6	14.2 14.3 14.7
(129M imgs)	× ✓ L	× ✓ _L	ViT-L/16	80.6 82.4	64.1 65.1	95.1 96.7	85.5 86.7	40.3 40.4	135.5 136.7	112.5 113.2	14.7 14.8

Downstream tasks include image-text retrieval and image captioning with finetuning (FT) and zero-shot (ZS) settings.

• Effect of sharing parameters between the captioner and filter.

Captioner &	Noise	Retrieval-	FT (COCO)	Retrieval	-ZS (Flickr)	Caption-	FT (COCO)	Caption-ZS (NoCaps)		
Filter	ratio	TR@1	IR@1	TR@1	IR@1	B@4	CIDEr	CIDEr	SPICE	
Share parameters	8%	79.8	62.2	94.3	83.7	38.4	129.0	103.5	14.2	
Decoupled	25%	80.6	63.1	94.8	84.9	38.6	129.7	105.1	14.4	

Large Multi-modal Models (BLIP-2)

• Pretraining Pipeline



Li, Junnan, et al. "Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models." *arXiv preprint arXiv:2301.12597* (2023).

Large Multi-modal Models (BLIP-2)

• Pretraining Pipeline

Model architecture of Q-Former and BLIP-2's first-stage vision-language representation learning objectives.



Q: query token positions; T: text token positions.

The self-attention masking strategy for each objective to control query-text interaction.



Large Multi-modal Models (BLIP-2)

• Results on Zero-shot Vision-Language Models

Models	#Trainable Open- Params sourced		Visual Question Answering VQAv2 (test-dev) VQA acc.	Image C NoCaj CIDEr	aptioning ps (val) SPICE	Image-Text Retrieval Flickr (test) TR@1 IR@1		
BLIP (Li et al., 2022)	583M	✓	-	113.2	14.8	96.7	86.7	
SimVLM (Wang et al., 2021b)	1.4B	X	-	112.2	-	-	-	
BEIT-3 (Wang et al., 2022b)	1.9B	×	-	-	-	94.9	81.5	
Flamingo (Alayrac et al., 2022)	10.2B	×	56.3	-	-	-	-	
BLIP-2	188M	✓	65.0	121.6	15.8	97.6	89.7	

Overview of BLIP-2 results on various zero-shot vision-language tasks.

Comparison with state-of-the-art methods on zero-shot visual question answering.

Models	#Trainable Params	#Total Params	V val	QAv2 test-dev	OK-VQA test	GQA test-dev
117 776	22.024	2001	10.5		5.0	60
VL-15 _{no-vqa}	224M	269M	13.5	-	5.8	6.3
FewVLM (Jin et al., 2022)	740M	785M	47.7	-	16.5	29.3
Frozen (Tsimpoukelli et al., 2021)	40M	7.1B	29.6	-	5.9	-
VLKD (Dai et al., 2022)	406M	832M	42.6	44.5	13.3	-
Flamingo3B (Alayrac et al., 2022)	1.4B	3.2B	-	49.2	41.2	-
Flamingo9B (Alayrac et al., 2022)	1.8B	9.3B	-	51.8	44.7	-
Flamingo80B (Alayrac et al., 2022)	10.2B	80B	-	56.3	50.6	-
BLIP-2 ViT-L OPT _{2.7B}	104M	3.1B	50.1	49.7	30.2	33.9
BLIP-2 ViT-G OPT _{2.7B}	107M	3.8B	53.5	52.3	31.7	34.6
BLIP-2 ViT-G OPT _{6.7B}	108M	7.8B	54.3	52.6	36.4	36.4
BLIP-2 ViT-L FlanT5 _{XL}	103M	3.4B	62.6	62.3	39.4	44.4
BLIP-2 ViT-G FlanT5 _{XL}	107M	4.1B	<u>63.1</u>	<u>63.0</u>	40.7	44.2
BLIP-2 ViT-G FlanT5 _{XXL}	108M	12.1B	65.2	65.0	<u>45.9</u>	44.7

Large Multi-modal Models (InstructBLIP)

Model Architecture



Dai, Wenliang, et al. "InstructBLIP: Towards General-purpose Vision-Language Models with Instruction Tuning." *arXiv* preprint arXiv:2305.06500 (2023).

Large Multi-modal Models (InstructBLIP)

• Comparison between BLIP-2 and InstructBLIP

	SaianaaOA			A-OKVQA					
	IMG	OCR-VQA	OKVQA	Direct	Answer	nswer Multi-			
				Val	Test	Val	Test		
Duraniana COTA	LLaVA [25]	GIT [42]	PaLM-E(562B) [9]	[15]	[36]	[15]	[<mark>36</mark>]		
Previous SOIA	89.0	70.3	66.1	56.3	61.6	73.2	73.6		
BLIP-2 (FlanT5 _{XXL})	89.5	72.7	54.7	57.6	53.7	80.2	76.2		
InstructBLIP (FlanT5 _{XXL})	90.7	73.3	55.5	57.1	54.8	81.0	76.7		
BLIP-2 (Vicuna-7B)	77.3	69.1	59.3	60.0	58.7	72.1	69.0		
InstructBLIP (Vicuna-7B)	79.5	72.8	62.1	64.0	62.1	75.7	73.4		

Large Multi-modal Models (Frozen)

Model Architecture



Tsimpoukelli, Maria, et al. "Multimodal few-shot learning with frozen language models." Advances in Neural Information Processing Systems 34 (2021): 200-212.

Large Multi-modal Models (Frozen)

Inference-time Interface •



(b) 1-shot outside-knowledge VQA

(c) Few-shot image classification

Large Multi-modal Models (Frozen)

Experiment Results

n-shot Acc.	n=0	n=1	n=4	au
Frozen	29.5	35.7	38.2	×
Frozen scratch	0.0	0.0	0.0	X
Frozen finetuned	24.0	28.2	29.2	X
Frozen train-blind	26.2	33.5	33.3	×
Frozen _{VQA}	48.4	_	–	✓
Frozen VQA-blind	39.1	-	-	1
Oscar [23]	73.8	_	_	1

n-shot Acc.	n=0	n=1	n=4	au
Frozen	5.9	9.7	12.6	X
Frozen 400mLM	4.0	5.9	6.6	X
Frozen finetuned	4.2	4.1	4.6	X
Frozen train-blind	3.3	7.2	0.0	×
Frozen _{VQA}	19.6	_	_	X
Frozen VQA-blind	12.5	-	_	×
MAVEx [42]	39.4	_	_	1

Table 1: Transfer from Conceptual Captions to VQAv2. The τ column indicates whether a model uses training data from the VQAv2 training set. The row denoted *Frozen* train-blind is the blind baseline described in subsection 4.1. *Frozen* VQA is a

Table 2: Transfer from Conceptual Captions to OKVQA. The τ column indicates if a model uses training data from the OKVQA training set. *Frozen* does not train on VQAv2 except in the baseline row, and it never trains on OKVQA.

Large Multi-modal Models (Flamingo)

Architecture Overview



Alayrac, Jean-Baptiste, et al. "Flamingo: a visual language model for few-shot learning." Advances in Neural Information Processing Systems 35 (2022): 23716-23736.

Large Multi-modal Models (Flamingo)

• GATED XATTN-DENSE layers



Large Multi-modal Models (Flamingo)

• Experiment Results

<u></u>																		
Method	FT	Shot	OKVQA (I)	VQAv2 (I)	COCO (I)	MSVDQA (V)	VATEX (V)	VizWiz (I)	Flick30K (I)	MSRVTTQA (V)	iVQA (V)	YouCook2 (V)	STAR (V)	VisDial (I)	TextVQA (I)	NextQA (I)	HatefulMemes (I)	Rare Act (V)
Zero/Few shot SOTA	×	(N)	[34] 43.3	[114] 38.2	[124] 32.2	[58] 35.2	-	-	-	[58] 19.2	[135] 12.2		[143] 39.4	[79] 11.6	-	7 .]	[85] 66.1	[85] 40.7
	~	(X)	(16)	(4)	(0)	(0)	40.1	28.0	60.6	(0)	(0)	550	(0)	(0)	20.1	21.2	(0)	(0)
Elaminas 2D	\$	0	41.2	49.2	75.0	27.5	40.1	28.9	72.0	14.0	32.1	55.8	39.0	40.1	30.1	21.5	53.1	58.4
Fumingo-5B	x	32	45.9	57.1	99.0	42.6	59.2	45.5	71.2	25.6	37.7	76.7	41.6	47.3	30.6	26.1	56.3	-
	×	0	44.7	51.8	79.4	30.2	39.5	28.8	61.5	13.7	35.2	55.0	41.8	48.0	31.8	23.0	57.0	57.9
Flamingo-9B	×	4	49.3	56.3	93.1	36.2	51.7	34.9	72.6	18.2	37.7	70.8	42.8	50.4	33.6	24.7	62.7	-
	×	32	51.0	60.4	106.3	47.2	57.4	44.0	72.8	29.4	40.7	77.3	41.2	50.4	32.6	28.4	63.5	-
	×	0	50.6	56.3	84.3	35.6	46.7	31.6	67.2	17.4	40.7	60.1	39.7	52.0	35.0	26.7	46.4	60.8
Flamingo	×	4	57.4	63.1	103.2	41.7	56.0	39.6	75.1	23.9	44.1	74.5	42.4	55.6	36.5	30.8	68.6	-
Tunningo	×	32	57.8	67.6	113.8	<u>52.3</u>	65.1	49.8	75.4	31.0	45.3	86.8	42.2	55.6	37.9	33.5	70.0	
Pretrained			54.4	80.2	143.3	47.9	76.3	57.2	67.4	46.8	35.4	138.7	36.7	75.2	54.7	25.2	79.1	
FT SOTA	~	(X)	[34] (10K)	[140] (444K)	[124] (500K)	[28] (27K)	[153] (500K)	[65] (20K)	[150] (30K)	[51] (130K)	[135] (6K)	[132] (10K)	[128] (46K)	[79] (123K)	[137] (20K)	[129] (38K)	[62] (9K)	-



Introduction to Computer Vision

Next week: Lecture 15, Generative Model

Embodied Perception and InteraCtion Lab

Spring 2025

