# Multi-Head Attention Forward Pass



# Multi-Head Attention Backward Pass



1

**Operations (Ops)**

| Abbrev | Name | Type / Shape | Notes |
|---|---|---|---|
| LN | Layer Normalization | op | Normalizes per token (per last dim $D$). |
| DO | Dropout | op | Training-time stochastic dropout. |
| + | Bias Add | op | Adds broadcast bias; see $\mathrm{BC}_{B,S}(\cdot)$. |
| T | Transpose | op | Context-dependent dims (e.g., $[B, N_H, S, D_h] \rightarrow [B, N_H, D_h, S]$). |
| R | Reshape / Split / Merge | op | Head split/merge: $[B, S, D] \leftrightarrow [B, N_H, S, D_h]$. |
| C | Concatenate | op | Join heads along last dim: $[B, S, N_H, D_h] \rightarrow [B, S, D]$. |
| SM | Scale (+ Mask) | op | Multiply by $1/\sqrt{D_h}$ and apply mask to scores. |
| S | Softmax | op | Softmax over key length ($S$) per head. |
| $\mathrm{BC}_{B,S}(\cdot)$ | Broadcast | op | Broadcast a length-$D$ (or $D_h$, $V$) bias to $[B, S, \cdot]$. |
| dS | Softmax Backward | op | Backprop through softmax over $S$. |
| dSM | Scale/Mask Backward | op | Backprop through scaling and masking. |
| dC | De-concatenate (Backward) | op | Split grads from concatenated heads. |
| dLN | LayerNorm Backward | op | Uses cached LN stats $(\mu, \sigma)$ and $\mathbf{X}$. |

**Data Tensors (Values)**

| Symbol | Name | Shape | Notes |
|---|---|---|---|
| $\mathbf{X}$ | Input hidden states | $[B, S, D]$ | Into MHA block (pre-LN). |
| $\mathbf{X}_{\mathrm{norm}}$ | LN output | $[B, S, D]$ | Result of $\mathrm{LN}(\mathbf{X})$. |
| $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ | Query/Key/Value | $[B, N_H, S, D_h]$ | From linear projections of $\mathbf{X}_{\mathrm{norm}}$. |
| $\widetilde{\mathbf{W}}_Q$ | Q weight | $[D, D]$ | Per-head realized by reshape; drawn as single matmul. |
| $\widetilde{\mathbf{W}}_K$ | K weight | $[D, D]$ | Same convention. |
| $\widetilde{\mathbf{W}}_V$ | V weight | $[D, D]$ | Same convention. |
| $\widetilde{\mathbf{W}}_O$ | Output-proj weight | $[D, D]$ | Maps concatenated heads back to model dim. |
| $\widetilde{\mathbf{b}}_O$ | Output bias | $[D]$ | Broadcast via $\mathrm{BC}_{B,S}$. |
| $\mathbf{A}$ | Attention scores | $[B, N_H, S, S]$ | $\mathbf{Q}\mathbf{K}^T/\sqrt{D_h}$ (plus mask). |
| $\mathbf{AS}$ | Attention weights | $[B, N_H, S, S]$ | $\mathrm{softmax}(\mathbf{A})$. |
| $\mathbf{AO}_{\mathrm{heads}}$ | Per-head outputs | $[B, N_H, S, D_h]$ | $\mathbf{AS} \cdot \mathbf{V}$. |
| $\mathbf{AO}_{\mathrm{cat}}$ | Concatenated heads | $[B, S, D]$ | After $C$. |
| $\mathbf{AO}_{\mathrm{lin}}$ | Linear output | $[B, S, D]$ | $\mathbf{AO}_{\mathrm{cat}}\widetilde{\mathbf{W}}_O$. |
| $\mathbf{AO}_{\mathrm{bias}}$ | Bias-added output | $[B, S, D]$ | $\mathbf{AO}_{\mathrm{lin}} + \widetilde{\mathbf{b}}_O$. |
| $\mathbf{A}_{\mathrm{out}}$ | MHA output | $[B, S, D]$ | After dropout; to next sublayer. |
| $\mathbf{dA}_{\mathrm{out}}$ | Grad wrt MHA output | $[B, S, D]$ | Backprop signal entering MHA. |
| $\mathbf{dQ}, \mathbf{dK}, \mathbf{dV}$ | Gradients for Q/K/V | $[B, N_H, S, D_h]$ | From attention-core backward. |
| $\mathbf{dK^T}$ | Grad of $K^T$ | $[B, N_H, D_h, S]$ | Before transpose/reshape to $\mathbf{dK}$. |
| $\mathbf{dAO}_{\mathrm{heads}}$ | Grad at heads | $[B, N_H, S, D_h]$ | Split from $\mathbf{dAO}_{\mathrm{cat}}$. |
| $\mathbf{dX}_{\mathrm{norm},Q}$ | Grad wrt $X_{\mathrm{norm}}$ (Q branch) | $[B, S, D]$ | Contribution via $W_Q^T$. |
| $\mathbf{dX}_{\mathrm{norm},K}$ | Grad wrt $X_{\mathrm{norm}}$ (K branch) | $[B, S, D]$ | Contribution via $W_K^T$. |
| $\mathbf{dX}_{\mathrm{norm},V}$ | Grad wrt $X_{\mathrm{norm}}$ (V branch) | $[B, S, D]$ | Contribution via $W_V^T$. |
| $\mathbf{dX}_{\mathrm{norm}}$ | Sum of above | $[B, S, D]$ | Input to dLN. |
| $\mathbf{dX}$ | Grad wrt input $X$ | $[B, S, D]$ | Output of dLN. |
| $\mathbf{d}\widetilde{\mathbf{W}}_Q$ | Q weight grad | $[D, D]$ | Standard matmul rule. |
| $\mathbf{d}\widetilde{\mathbf{W}}_K$ | K weight grad | $[D, D]$ | Standard matmul rule. |
| $\mathbf{d}\widetilde{\mathbf{W}}_V$ | V weight grad | $[D, D]$ | From reshaped $d\mathbf{V}$ and $X_{\mathrm{norm}}$. |
| $\mathbf{d}\widetilde{\mathbf{W}}_O$ | Output-proj grad | $[D, D]$ | From $\mathbf{AO}_{\mathrm{cat}}^T$ and $d\mathbf{AO}_{\mathrm{lin}}$. |
| $\mathbf{d}\widetilde{\mathbf{b}}_O$ | Output bias grad | $[D]$ | Sum over $B, S$ of $d\mathbf{AO}_{\mathrm{lin}}$. |

**Shape symbols:** $B$=batch size, $S$=sequence length, $D$=model dim, $N_H$=num heads, $D_h = D/N_H$.

**Implementation note:** Per-head $[D, D]$ drawings depict fused linears realized via reshape to $N_H \times D_h$.