

Full Transformer Forward + Backprop

Sunday, 26 October 2025 08:36

FORWARD

Transformer.forward():

$$x = \text{Tokenise}(x)$$

$$X_{\text{one-hot}} = \text{One-hot}(x)$$

$$X_{\text{emb}} = x W_{\text{emb}}$$

$$X_{\text{pos}} = X_{\text{emb}} + \text{pos-enc}(x)$$

$$X = X_{\text{pos}}$$

FOR EACH LAYER:

$$x = \text{Layer.forward}(x)$$

$$X_{\text{norm-out}} = \text{RMSNorm.forward}(x)$$

$$out = X_{\text{norm-out}} W_{\text{out-emb}}$$

$$\text{Pred} = \text{softmax}(out)$$

SHAPES

$$(B, L)$$

$$(B, L, d_{\text{vocab}})$$

$$(B, L, d_{\text{vocab}}) (d_{\text{vocab}}, d_{\text{model}}) \rightarrow (B, L, d_{\text{model}})$$

$$(B, L, d_{\text{model}})$$

↑ Take softmax($\text{pred}[E]$) to get

Probability distribution over next token

$$(B, L, d_{\text{model}})$$

Layer.forward(x)

SHAPES

$$x_{\text{norm}} = \text{RMSNorm.forward}(x)$$

$$(B, L, d_{\text{model}})$$

$$X_{\text{multi-att}} = x + \text{MultiHead.forward}(x_{\text{norm}})$$

$$(B, L, d_{\text{model}}) + (B, L, d_{\text{model}})$$

$$X_{\text{multi-att-norm}} = \text{RMSNorm.forward}(X_{\text{multi-att}})$$

$$(B, L, d_{\text{model}})$$

$$\text{layer.out} = X_{\text{multi-att}} + \text{FeedForward.forward}(X_{\text{multi-att-norm}})$$

$$(B, L, d_{\text{model}}) + (B, L, d_{\text{model}})$$

return layer.out

$$(B, L, d_{\text{model}})$$

MultiHead.forward(x)

FOR EACH HEAD:

SHAPES

$$x_{\text{head}} = \text{Head.forward}(x)$$

$$(B, L, d_v)$$

$$X_{\text{head-out}} = \text{Concat}(x_{\text{head}})$$

$$(B, L, h \cdot d_v)$$

$$out = X_{\text{head-out}} \cdot W_0$$

$$(B, L, h \cdot d_v) (h \cdot d_v, d_{\text{model}}) \rightarrow (B, L, d_{\text{model}})$$

return out

$$(B, L, d_{\text{model}})$$

FeedForward.forward(X_H)

SHAPES

$$FF = X_H W_H + b_H,$$

$$(B, L, d_{\text{model}}) (d_{\text{model}}, d_{\text{ff}}) + (d_{\text{ff}}) \rightarrow (B, L, d_{\text{ff}})$$

Head forward forward (X_H)

SHAPES

$$FF_1 = X_H W_{H_1} + b_{H_1}$$

$$(B, L, d_{model}) (d_{model}, d_H) \rightarrow (B, L, d_H)$$

$$FF'_1 = \text{ReLU}(FF_1)$$

$$(B, L, d_H)$$

$$FF_2 = FF'_1 W_{H_2} + b_{H_2}$$

$$(B, L, d_H) (d_H, d_{model}) \rightarrow (B, L, d_{model})$$

return FF_2

(B, L, d_{model})

Head forward (x)

SHAPES

$$K = x W_K$$

$$(B, L, d_{model}) (d_{model}, d_K) \rightarrow (B, L, d_K)$$

$$Q = x W_Q$$

$$(B, L, d_{model}) (d_{model}, d_K) \rightarrow (B, L, d_K)$$

$$V = x W_V$$

$$(B, L, d_{model}) (d_{model}, d_V) \rightarrow (B, L, d_V)$$

$$X_{\text{att_lin}} = Q K^\top / \sqrt{d_K}$$

$$(B, L, d_K) (B, d_K, L) \rightarrow (B, L, L)$$

$$X_{\text{att_mask}} = X_{\text{att_lin}} + \text{Mask}$$

$$X_{\text{att}} = \text{Softmax.forward}(X_{\text{att_mask}})$$

$$(B, L, L)$$

$$X_{\text{att}_V} = X_{\text{att}} V$$

$$(B, L, L) (B, L, d_V) \rightarrow (B, L, d_V)$$

return X_{att_V}

Softmax.forward (x):

$$y = e^x / \sum_{i=1}^L e^{x_i}$$

$$(B, L, L)$$

return y

RMSNorm.forward (x):

Whatever RMSNorm
was init. to
Usually d_{model}

$$\text{rms_norm} = x / \sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2}$$

$$(B, L, d_{model})$$

$$\text{rms} = \text{rms_norm} * g$$

$$(B, L, d_{model}) * (d_{model})$$

return rms

BACKWARD

$$\mathcal{L} = \text{CrossEntropyLoss}(\text{pred})$$

Transformer.backward (Pred, labels):	SHAPES
$\frac{dt}{dout} = \text{outputs} - \text{labels}$	$(B, L, d_{vocab}) - (B, L, d_{vocab}) \rightarrow (B, L, d_{vocab})$
$\frac{dt}{dX_{norm_out}} = \frac{dt}{dout} \cdot W_{vocab}^T$	$(B, L, d_{vocab}) (d_{vocab}, d_{model}) \rightarrow (B, L, d_{model})$
$\frac{dt}{dW_{vocab}} = X_{norm_out}^T \cdot \frac{dt}{dout}$	$(B, d_{model}, L) (B, L, d_{vocab}) \rightarrow (B, d_{model}, d_{vocab})$
$\frac{dt}{layer_out} = \text{RMSNorm. backward}(\frac{dt}{dX_{norm_out}})$	(B, L, d_{model})
Prev-grad = $\frac{dL}{layer_out}$	
FOR EACH reversed (LAYER):	
Prev-grad = Layer.backward (Prev-grad)	(B, L, d_{model})
$\frac{dL}{dW_{vocab}} = X_{one-hot}^T \cdot \text{Prev-grad}$	$(B, d_{vocab}, L) (B, L, d_{model}) \rightarrow (B, d_{vocab}, d_{model})$
$\xrightarrow{(B, L, d_{model})}$	
Layer.backward (Prev-grad):	SHAPES
$\frac{dt}{dX_{multi_att_norm}} = \text{FeedForward.backward(Prev-grad)}$	(B, L, d_{model})
$\frac{dt}{dX_{multi_att}} = \text{RMSNorm.backward}(\frac{dt}{dX_{multi_att_norm}}) + \text{Prev-grad}$	$(B, L, d_{model}) + (B, L, d_{model})$
$\frac{dt}{dX_{norm}} = \text{Multi Head. backward}(\frac{dt}{dX_{multi_att}})$	(B, L, d_{model})
$\frac{dt}{dX} = \text{RMSNorm.backward}(\frac{dt}{dX_{norm}}) + \frac{dt}{dX_{multi_att}}$	$(B, L, d_{model}) + (B, L, d_{model})$
return $\frac{dt}{dX}$	
$\xrightarrow{(B, L, d_{model})}$	
FeedForward.backward (Prev-grad):	SHAPES
$\frac{dt}{dW_{H_2}} = FF_i^T \cdot \text{Prev-grad}$	$(B, dH, L) (B, L, d_{model}) \rightarrow (B, dH, d_{model})$
$\frac{dt}{db_{H_2}} = \text{Prev-grad}$	(B, L, d_{model})
$\frac{dt}{dFF_i} = \text{Prev-grad} \cdot W_{H_2}^T$	$(B, L, d_{model}) (d_{model}, dH) \rightarrow (B, L, dH)$
$\frac{dt}{dFF_i} = \text{ReLU. backward}(\frac{dt}{dFF_i})$	(B, L, dH)
$\frac{dt}{dW_{H_1}} = X_H^T \cdot \frac{dt}{dFF_i}$	$(B, d_{model}, L) (B, L, dH) \rightarrow (B, d_{model}, dH)$
$\frac{dt}{db_{H_1}} = \frac{dt}{dFF_i}$	(B, L, dH)
$\frac{dt}{dX_H} = \frac{dt}{dFF_i} \cdot W_{H_1}^T$	$(B, L, dH) (dH, d_{model}) \rightarrow (B, L, d_{model})$

$$\frac{dL}{dX_H} = \frac{dL}{dFF_i} \cdot W_H^T$$

return $\frac{dL}{dX_H}$

$$(B, L, d_H) (d_H, d_{model}) \rightarrow (B, L, d_{model})$$

Multi-Head. backward (prev-grad):

$$\frac{dL}{dW_0} = X_{\text{head-out}}^T \cdot \text{prev-grad}$$

$$\frac{dL}{dX_{\text{head-out}}} = \text{prev-grad} \cdot W_0^T$$

SHAPES

$$(B, h \cdot dv, L) (B, L, d_{model}) \rightarrow (B, h \cdot dv, d_{model})$$

$$(B, L, d_{model}) (d_{model}, h \cdot dv) \rightarrow (B, L, h \cdot dv)$$

FOR EACH HEAD: (B, L, dv)

$$\frac{dL}{dX} += \text{Head. backward} \left(\frac{dL}{dX_{\text{head-out}}} \right)^\text{[head]}$$

$$(B, L, d_{model})$$

$$\frac{dL}{dX} /= \text{num_heads}$$

return $\frac{dL}{dX}$

(B, L, dv)

Head. backward (prev-grad):

$$\frac{dL}{dV} = X_{\text{att}}^T \cdot \text{prev-grad}$$

$$\frac{dL}{dX_{\text{att}}} = \text{prev-grad} \cdot V^T$$

SHAPES

$$(B, L, L) (B, L, dv) \rightarrow (B, L, dv)$$

$$(B, L, dv) (B, dv, L) \rightarrow (B, L, L)$$

$$\frac{dL}{dX_{\text{att-mask}}} = \text{softmax. backward} \left(\frac{dL}{dX_{\text{att}}} \right)$$

$$(B, L, L)$$

$$\frac{dL}{dX_{\text{att-in}}} = \text{mask} \left(\frac{dL}{dX_{\text{att-mask}}} \right)$$

$$(B, L, L)$$

$$\frac{dL}{dQ} = \frac{1}{\sqrt{d_k}} \frac{dL}{dX_{\text{att-in}}} \cdot K$$

$$(B, L, L) (B, L, d_k) \rightarrow (B, L, d_k)$$

$$\frac{dL}{dS} = \frac{1}{\sqrt{d_k}} \frac{dL}{dX_{\text{att-in}}}^T \cdot Q$$

$$(B, L, L) (B, L, d_k) \rightarrow (B, L, d_k)$$

$$\frac{dL}{dW_Q} = X^T \cdot \frac{dL}{dQ}$$

$$(B, d_{model}, L) (B, L, d_k) \rightarrow (B, d_{model}, d_k)$$

$$\frac{dL}{dW_K} = X^T \cdot \frac{dL}{dK}$$

$$(B, d_{model}, L) (B, L, d_k) \rightarrow (B, d_{model}, d_k)$$

$$\frac{dL}{dW_V} = X^T \cdot \frac{dL}{dV}$$

$$(B, d_{model}, L) (B, L, dv) \rightarrow (B, d_{model}, dv)$$

$$\frac{dL}{dX} = \underbrace{\left(\frac{dL}{dK} \cdot W_K^T \right)}_{(B, L, d_k) (d_k, d_{model})} + \underbrace{\left(\frac{dL}{dQ} \cdot W_Q^T \right)}_{(B, L, dv) (dv, d_{model})} + \underbrace{\left(\frac{dL}{dV} \cdot W_V^T \right)}_{(B, L, dv) (d_{model}, dv)} \cdot \frac{1}{3}$$

$$(B, L, d_{model})$$

$$(B, L, d_k) (d_k, d_{model}) \rightarrow (B, L, d_{model})$$

$$(B, L, dv) (dv, d_{model}) \rightarrow (B, L, d_{model})$$

$$\text{return } \frac{df}{dx}$$

Softmax backward (prev-grad):

$$\text{Let } C = \sum_{i=1}^n e^{x_i}, \text{ so } y_i = \frac{e^{x_i}}{C}$$

$$\text{Use quotient rule: } \nabla \frac{du}{dx} = u \frac{dv}{dx} - v \frac{du}{dx}$$

$$v^2$$

Each y_i depends on all x (due to $C = \sum_{i=1}^n e^{x_i}$)

\therefore 2 cases:

Case 1: $i = k$

$$u = e^{x_i}, \frac{du}{dx_i} = e^{x_i}, v = C, \frac{dv}{dx_i} = e^{x_i}$$

$$\begin{aligned}\frac{dy_i}{dx_i} &= \frac{Ce^{x_i} - e^{2x_i}}{C^2} \\ &= \frac{e^{x_i}}{C} - \frac{e^{2x_i}}{C^2} \\ &= y_i - y_i^2 \\ &= y_i(1 - y_i)\end{aligned}$$

Case 2: $i \neq k$

$$u = e^{x_k}, \frac{du}{dx_k} = 0, v = C, \frac{dv}{dx_k} = e^{x_k}$$

$$\begin{aligned}\frac{dy_i}{dx_k} &= \frac{(0 \cdot C) - e^{x_i} e^{x_k}}{C^2} \\ &= -\frac{e^{x_i + x_k}}{C^2} \\ &= -\frac{e^{x_i}}{C} \frac{e^{x_k}}{C} \\ &= -y_i y_k\end{aligned}$$

$$\frac{dy_i}{dx_k} = \begin{cases} y_i(1 - y_i) & \text{if } i = k \\ -y_i y_k & \text{if } i \neq k \end{cases}$$

RMS Norm. backward (prev-grad):

$$\frac{df}{dg} = \text{rms_norm} * \text{prev_grad}$$

$$(\mathcal{B}, L, d_{\text{rms}}) * (\mathcal{B}, L, d_{\text{rms}}) \rightarrow (\mathcal{B}, L, d_{\text{rms}})$$

$$\tilde{dg} = \text{rms_norm} * \text{Prev_grad}$$

$$(\mathcal{B}, L, d_{\text{rms}}) * (\mathcal{B}, L, d_{\text{rms}}) \rightarrow (\mathcal{B}, L, d_{\text{rms}})$$

$$\frac{d\mathcal{L}}{dx_i}$$

$$\text{Let } r = \sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2}, \text{ so } \text{rms}_i = \frac{x_i}{r} * g_i$$

$$\frac{d_{\text{rms}_i}}{dx_i} = \frac{d}{dx_i} \left(\frac{x_i}{r} \right) * g_i, \text{ but rms}_i \text{ depends on all } x \text{ due to } r$$

$\therefore 2$ cases:

Case 1: $i = k$:

$$\begin{aligned} U &= x_i, \frac{du}{dx_i} = 1, V = r, \frac{dv}{dx_i} = \frac{d}{dx_i}(r) \\ \rightarrow \frac{dr}{dx_i} &: \frac{dr}{du} \frac{du}{dx_i}, V = \frac{1}{n} \sum_{i=1}^n x_i^2, r = V^{\frac{1}{2}} \\ &= \frac{1}{2} V^{-\frac{1}{2}} \cdot \frac{2x_i}{n} \\ &= \frac{1}{2r} \cdot \frac{2x_i}{n} \\ &= \frac{x_i}{rn} \end{aligned}$$

$$\begin{aligned} \frac{d_{\text{rms}_i}}{dx_i} &\approx \frac{r - \frac{x_i^2}{rn}}{r^2} \\ &= \left(\frac{1}{r} - \frac{x_i^2}{r^3 n} \right) * g_i \end{aligned}$$

Case 2: $i \neq k$

$$U = x_i, \frac{du}{dx_k} = 0, V = r, \frac{dv}{dx_k} = \frac{x_k}{rn}$$

$$\begin{aligned} \frac{d_{\text{rms}_i}}{dx_k} &= \frac{-x_i \frac{x_k}{rn}}{r^2} \\ &= -\frac{x_i x_k}{r^3 n} * g_i \end{aligned}$$

$$\frac{d_{\text{rms}_i}}{dx_k} = \begin{cases} \left(\frac{1}{r} - \frac{x_i^2}{r^3 n} \right) * g_i & \text{if } i=k \\ -\frac{x_i x_k}{r^3 n} * g_i & \text{if } i \neq k \end{cases}$$

$$\frac{d\mathcal{L}}{dx_k} = \text{Prev_grad} \cdot \frac{d_{\text{rms}_i}}{dx_k}$$

$$(\mathcal{B}, L, d_{\text{rms}}) (\mathcal{B}, L, d_{\text{rms}}, d_{\text{rms}}) \rightarrow (\mathcal{B}, L, d_{\text{rms}})$$