



Operators	Buffered Activation	Memory Usage
$X_{n1} = RMSNorm(X_{in})$	$X_{in}$ $\sigma_{in}^2$	$(b, s, d)[w]$ $(b, s)[w]$
$Q = X_{n1}(W_Q + A_Q B_Q)$ $K = X_{n1}(W_K + A_K B_K)$ $V = X_{n1}(W_V + A_V B_V)$	$X_{n1}$ $(X_{n1} A_Q), (X_{n1} A_K)$ $(X_{n1} A_V)$	$(b, s, d)[w]$ $3 \times (b, s, r)[w]$
$Q = RoPE(Q, cos, sin)$ $K = RoPE(K, cos, sin)$	$cos$ $sin$	$2 \times (s, d/h)[w]$
$S = QK^T, A = Softmax(S)$ $O = AV$ <span style="color: red;">w/o FlashAttn</span>	$Q, K, V$ $A$	$3 \times (b, s, d)[w]$ $(b, h, s, s)[w]$
$O = FlashAttn(Q, K, V)$ <span style="color: green;">w FlashAttn</span>	$Q, K, V$	$3 \times (b, s, d)[w]$
$X_{mid} = O(W_O + A_O B_O)$	$O$ $(O A_O)$	$(b, s, d)[w]$ $(b, s, r)[w]$
$X_{n2} = RMSNorm(X_{mid})$	$X_{mid}$ $\sigma_{mid}^2$	$(b, s, d)[w]$ $(b, s)[w]$
$X_G = X_{n2}(W_G + A_G W_G)$ $X_U = X_{n2}(W_U + A_U W_U)$	$X_{n2}$ $(X_{n2} A_G), (X_{n2} A_U)$	$(b, s, d)[w]$ $2 \times (b, s, r)[w]$
$X_{SiLU} = SiLU(X_G)$	$X_G$	$(b, s, d_f)[w]$
$X_D = X_{SiLU} \odot X_U$	$X_{SiLU}$ $X_U$	$(b, s, d_f)[w]$ $(b, s, d_f)[w]$
$X_{out} = X_D(W_D + A_D B_D)$	$X_D$ $(X_D A_D)$	$(b, s, d_f)[w]$ $(b, s, r)[w]$
<b>Estimated Total Size(bit)</b>		$(8d + 4d_f)bsw$
+HyCLoRA@raw quant		$(8d + 4d_f)bsw_q$
+HyCLoRA@inter + intra		$(8d + 2d_f)bsw_q$