

seq512-hd64-12h-8b | Q=512, K=512, head\_dim=64, heads=12, batch=8

<b>Kernel</b>	<b>Avg Forward Time (ms)</b>	<b>Relative to Triton Flash</b>
scaled_dot_prod_attention	9.28	5.08x slower
vectorized_torch	7.10	3.89x slower
vectorized_torch_compiled	0.80	2.27x faster
<b>flash_attention_triton</b>	<b>1.83</b>	<b>1.00x</b>