

seq1024-hd128-16h-32b | Q=1024, K=1024, head\_dim=128, heads=16, batch=32

<b>Kernel</b>	<b>Avg Forward Time (ms)</b>	<b>Relative to Triton Flash</b>
scaled_dot_prod_attention	97.68	3.10x slower
vectorized_torch	108.24	3.44x slower
vectorized_torch_compiled	31.83	1.01x slower
<b>flash_attention_triton</b>	<b>31.47</b>	<b>1.00x</b>