

seq2048-hd128-16h-32b | Q=2048, K=2048, head_dim=128, heads=16, batch=32

Kernel	Avg Forward Time (ms)	Relative to Triton Flash
scaled_dot_prod_attention	ERROR	N/A
vectorized_torch	ERROR	N/A
vectorized_torch_compiled	182.95	1.48x slower
flash_attention_triton	123.57	1.00x