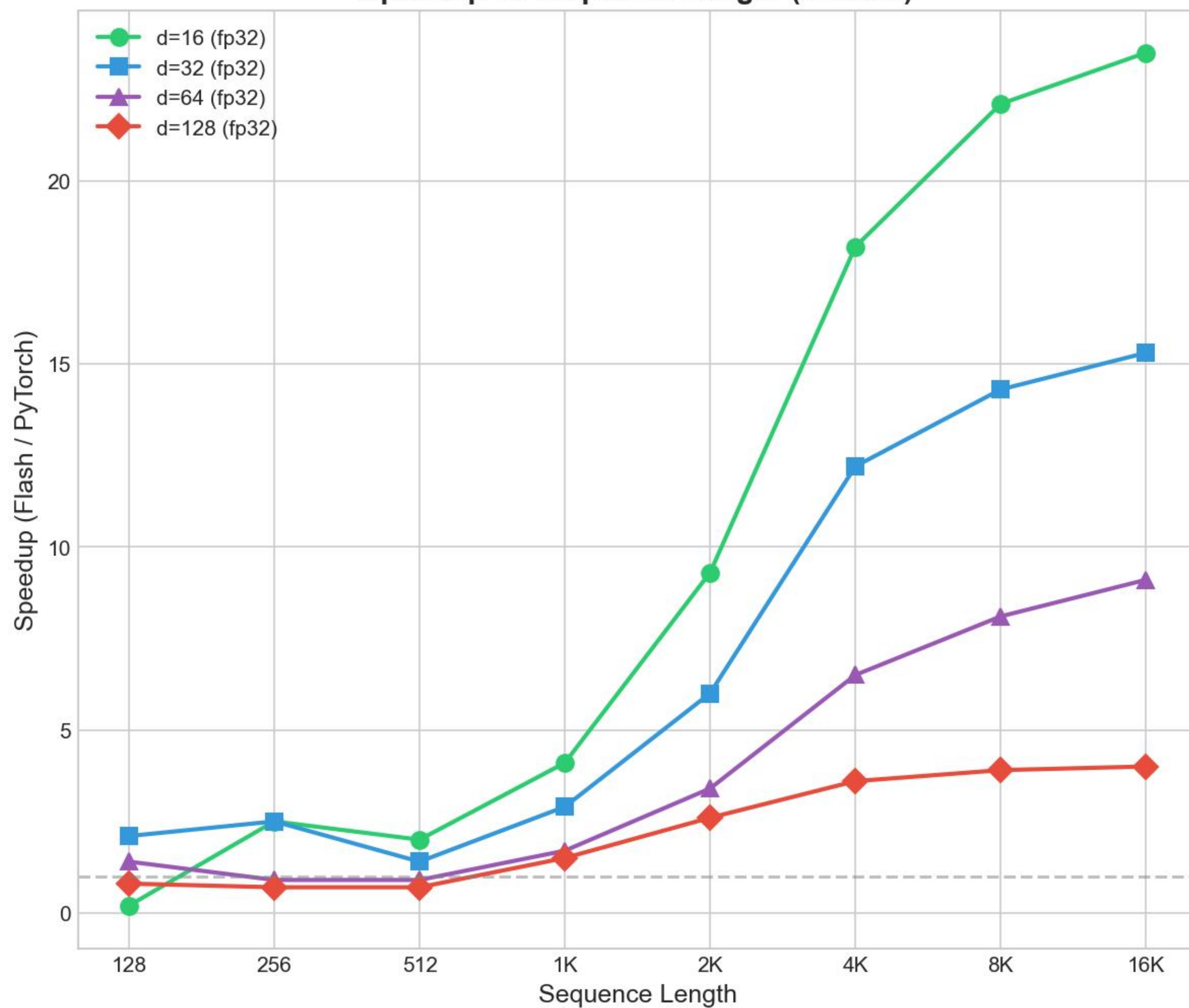
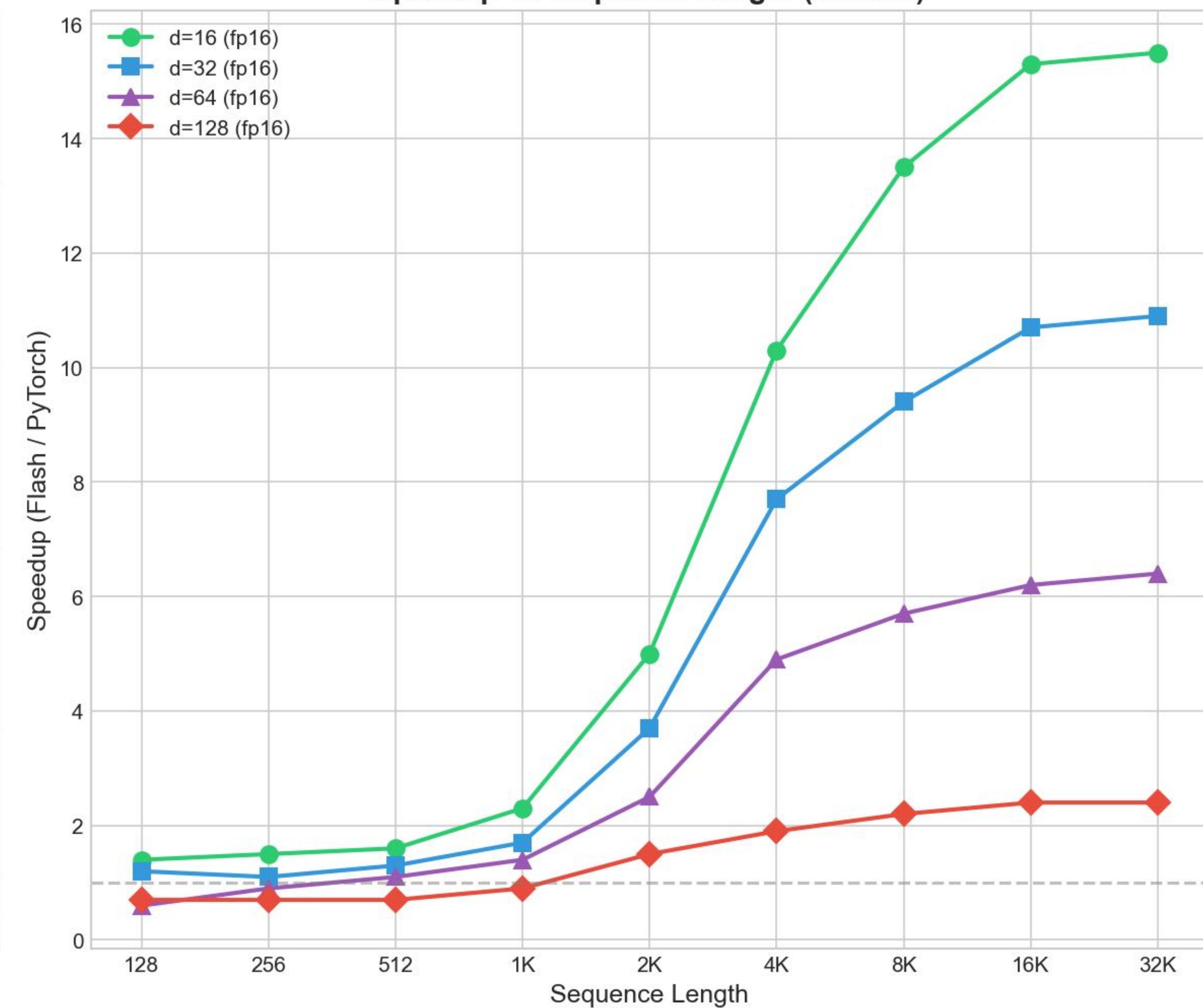


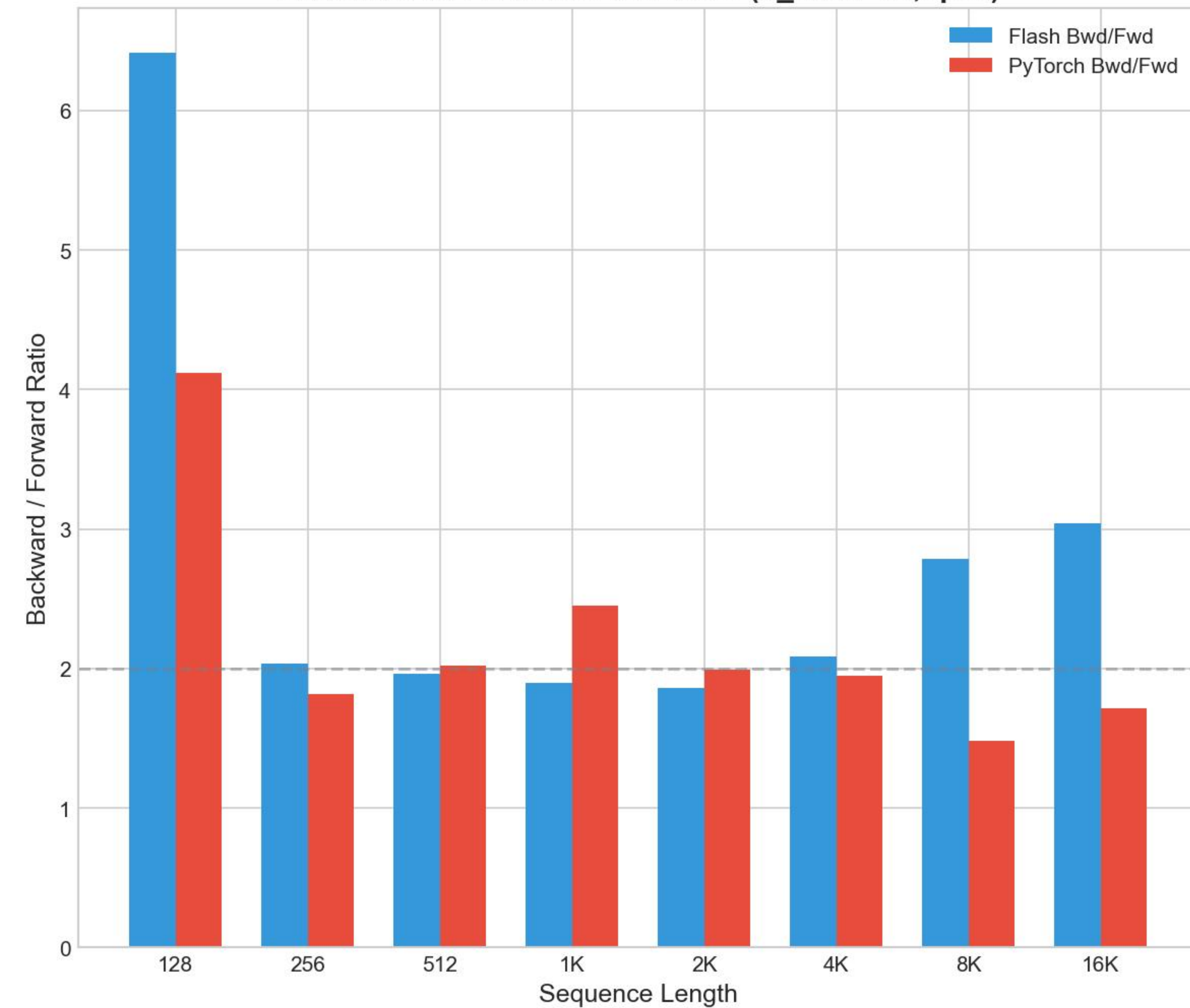
Speedup vs Sequence Length (Float32)



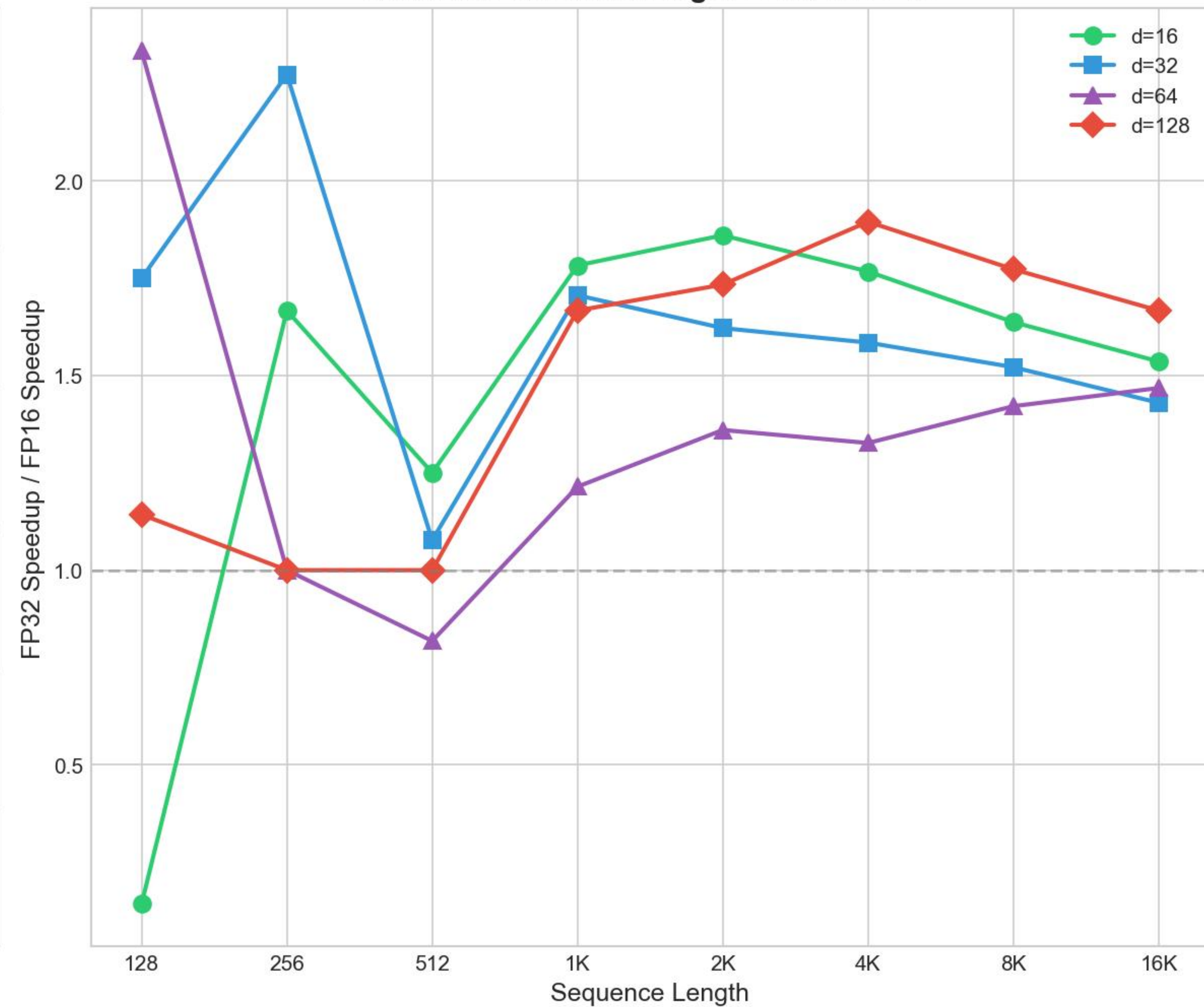
Speedup vs Sequence Length (Float16)



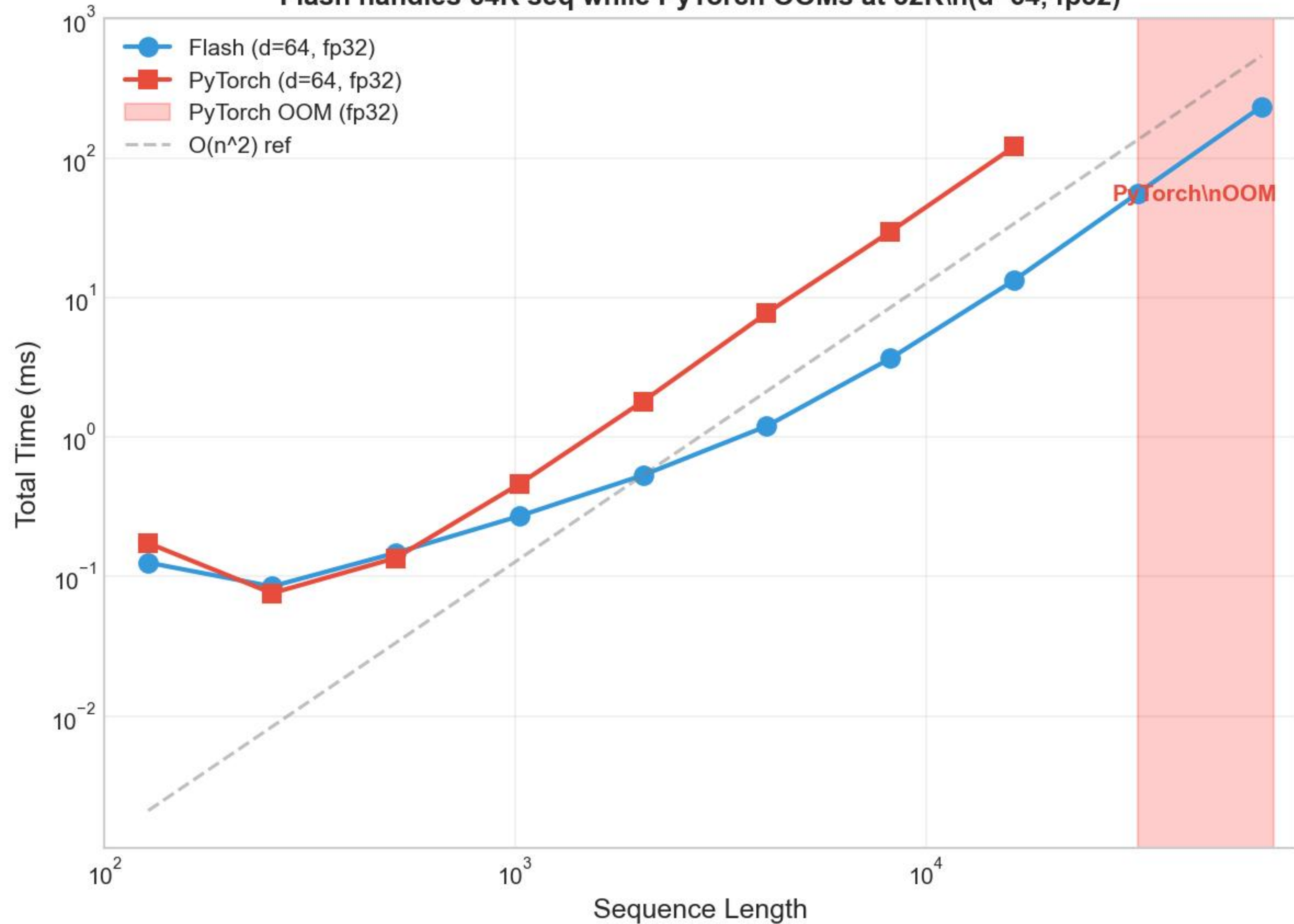
Backward vs Forward Pass Ratio (d_head=64, fp32)



Relative Flash Advantage: FP32 vs FP16



Flash handles 64K seq while PyTorch OOMs at 32K\n(d=64, fp32)



seq_len=65536: Flash works, PyTorch OOMs

