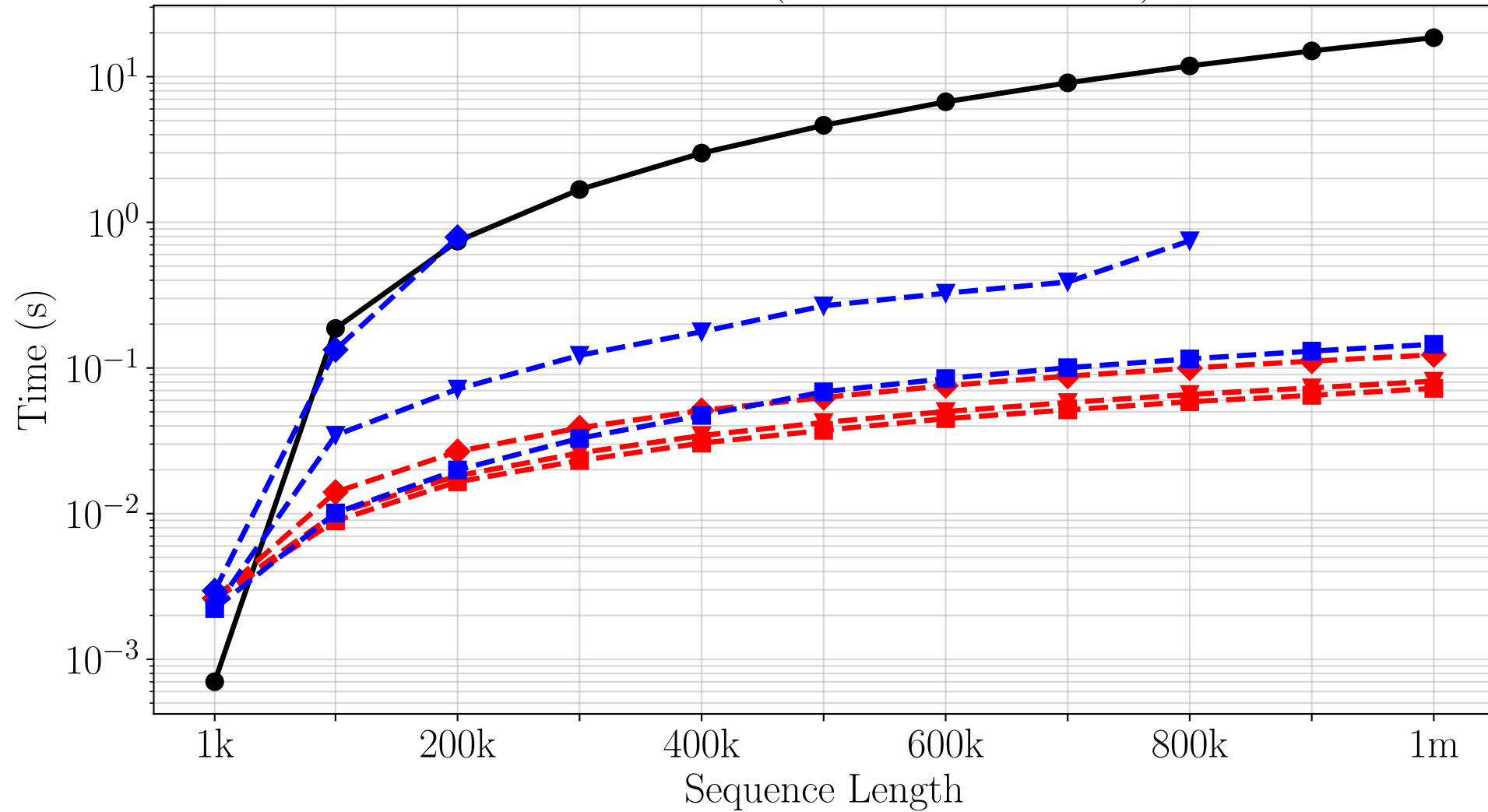
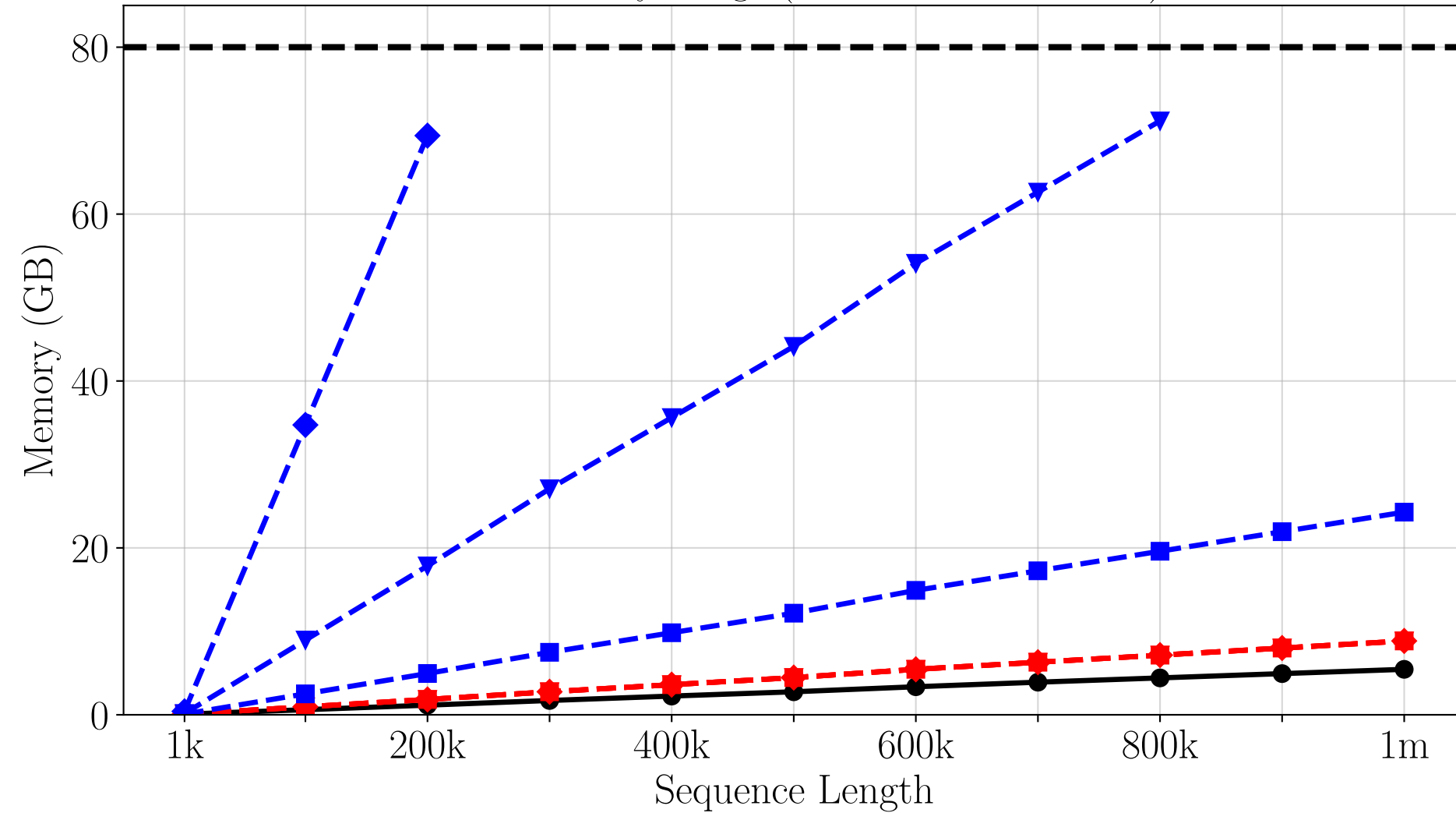


Execution Time (Forward + Backward)



Peak Memory Usage (Forward + Backward)



-■- FLARE (128 latents) (ours)    -◆- FLARE (2048 latents) (ours)    -▼- PhysAttention (512 slices)    -●- Softmax Attention  
 -▼- FLARE (512 latents) (ours)    -■- PhysAttention (128 slices)    -◆- PhysAttention (2048 slices)