

(A) BurstAttention (B) Megatron-RingAttention

## Attention Kernel Performance

