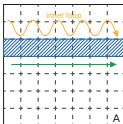
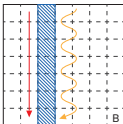


## Matrix memory layout:



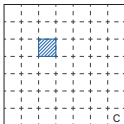
consecutive values  
in memory

x

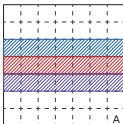


non consecutive values  
in memory

=

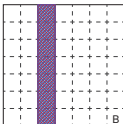


## Naïve kernel:



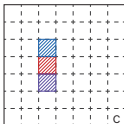
threads access non-consecutive  
values => cannot coalesce

x



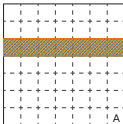
all threads access same  
values => within-warp broadcast

=



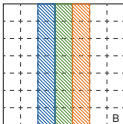
no benefit in putting these  
threads in same warp

## Coalescing kernel:



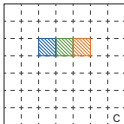
all threads access same  
values => within-warp broadcast

x



threads access consecutive  
values => can coalesce

=



make sure these threads end up  
in same warp to exploit coalescing