

Producer warp groups

*load indices
to shared
① memory*

*load indices
to register
② memory*

*async load
keys/values
③ at indices*

*load indices
to shared
memory*

*load indices
to register
memory*

**④ $qk^t + \text{softmax} +$
 value comp.**

**④ $qk^t + \text{softmax} +$
 value comp.**

Consumer warp groups

.....
time