

kernel calls

main: add <<< N_1, N_2 >>> (---)

N_1 : gridDim.x

N_2 : blockDim.x

kernel: gridDim.x blockDim.x

general form of kernel call:

kernel <<< $X_1, X_2, [n, s]$ >>> (---)

X_1, X_2 : dim3 objects

$X_1, X_2 = (n_1, n_2, n_3)$
 $= (m_1, m_2, m_3)$

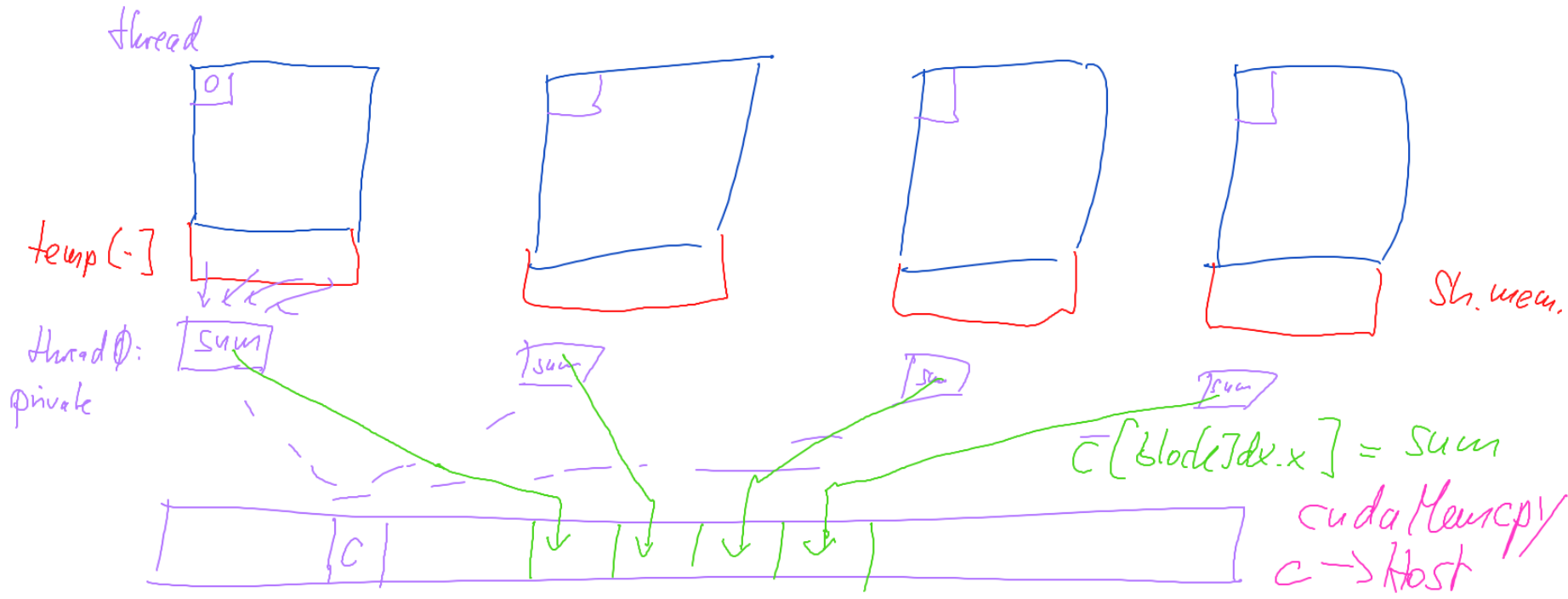
just integers N_1, N_2 : $X_1 = (N_1, 0, 0)$

$X_2 = (N_2, 0, 0)$

$n_{1-3} = \text{gridDim.x, y, z}$

$m_{1-3} = \text{blockDim.x, y, z}$

69000 > n: number of bytes to be allocated in shared memory for each block!
s: stream id \leftrightarrow grid identification (dynamic)



Global Memory

atomicAdd(c, sum)

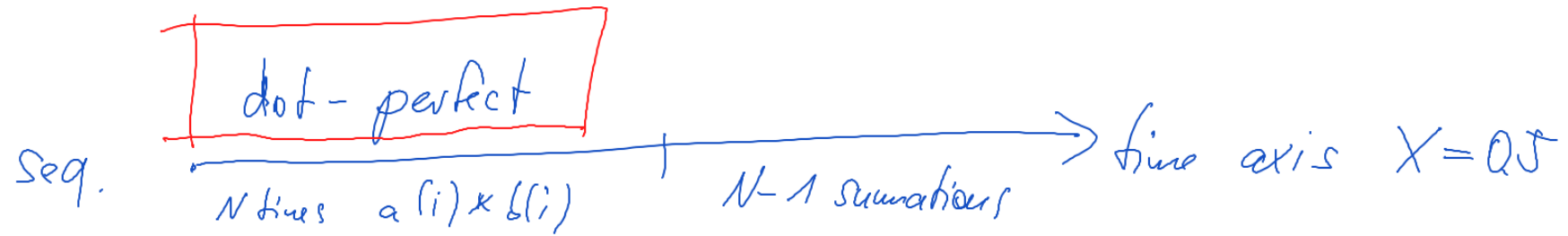
Declare $c[\leftarrow]$ — $\text{gridDim.x} = \frac{N}{\text{threads per block}}$

Version 1: `AtomicAdd(c, sum)`

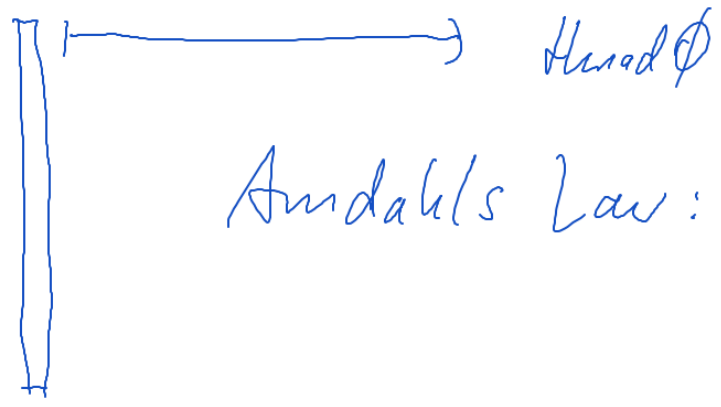
Sum Result of every block into one scalar c (global mem.)

Version 2: `c[blockIdx.x] = sum`

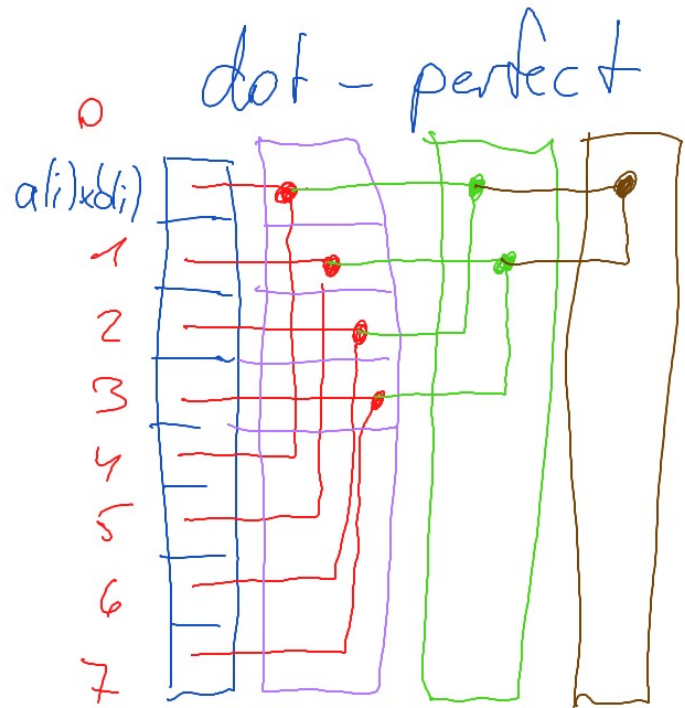
Store result of every block in a separate vector
element of c , sum on the host
Avoids race condition!



par.
~~⊗~~
 $a(i) \times b(i)$



Amdahl's Law: max. speedup 2



4 cycles to get all
 $a(i) \times b(i)$ into sum in thread \emptyset
 (previous one: 7 cycles needed!)

8 threads
 1 cycle