



# ME 471/571

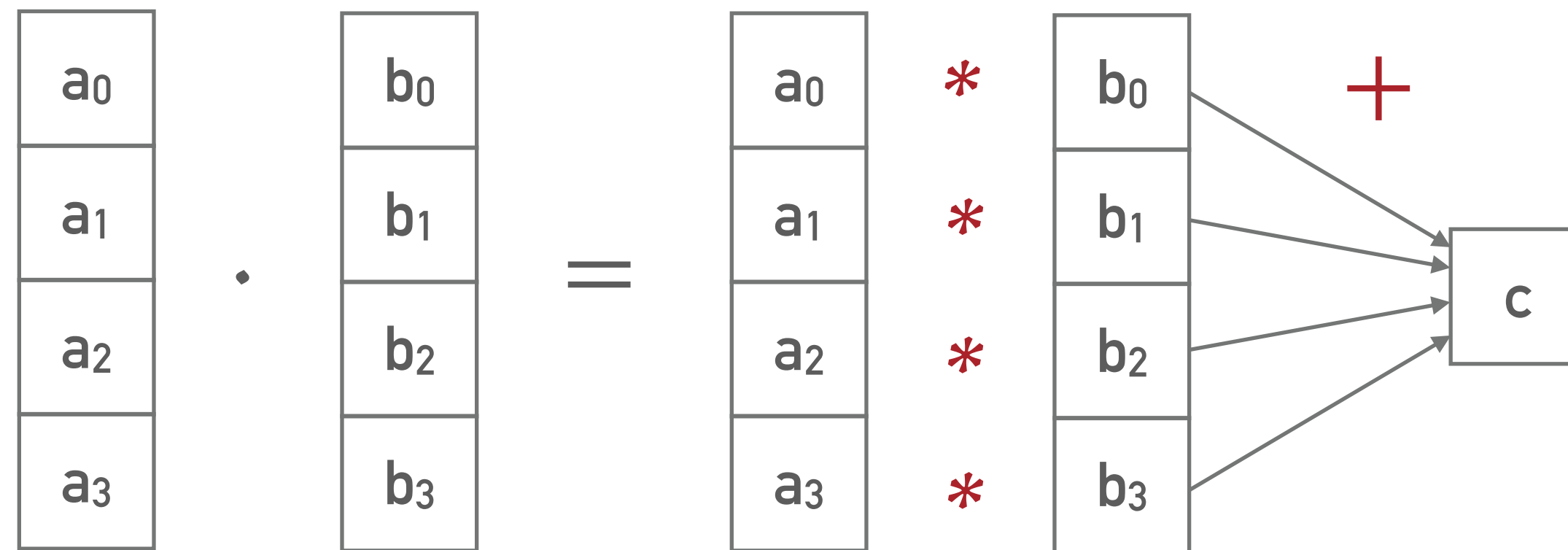
---

*Week 13 - Shared memory and thread synchronization*

# THREADS SYNCHRONIZATION

---

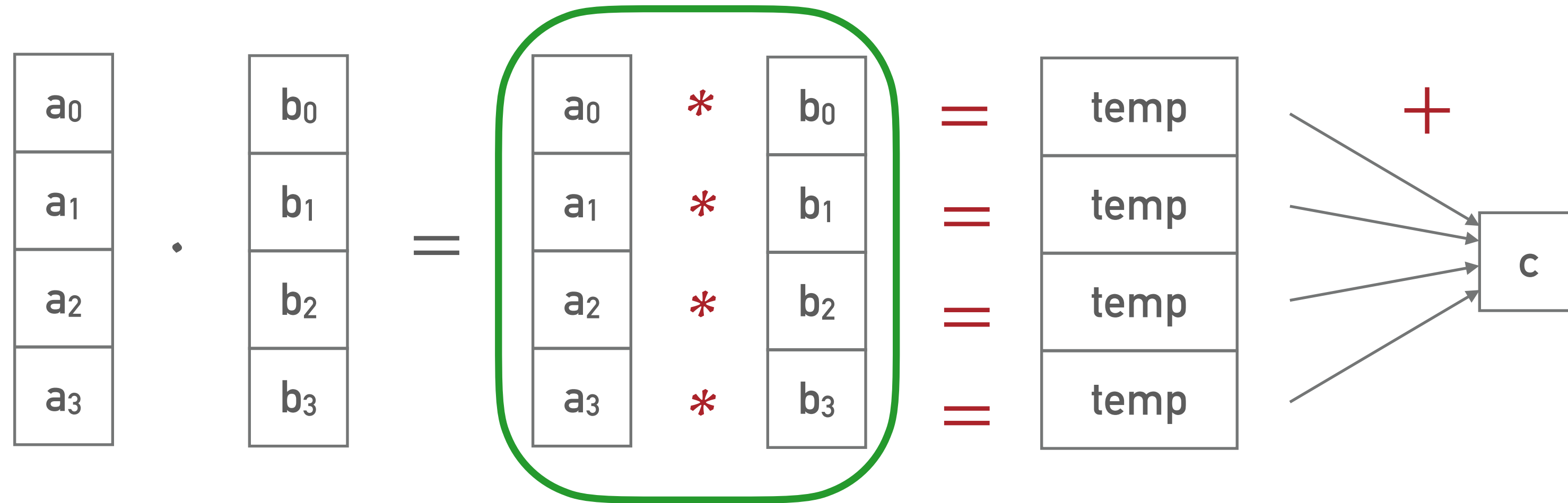
Consider a computation of dot product:



# THREADS SYNCHRONIZATION

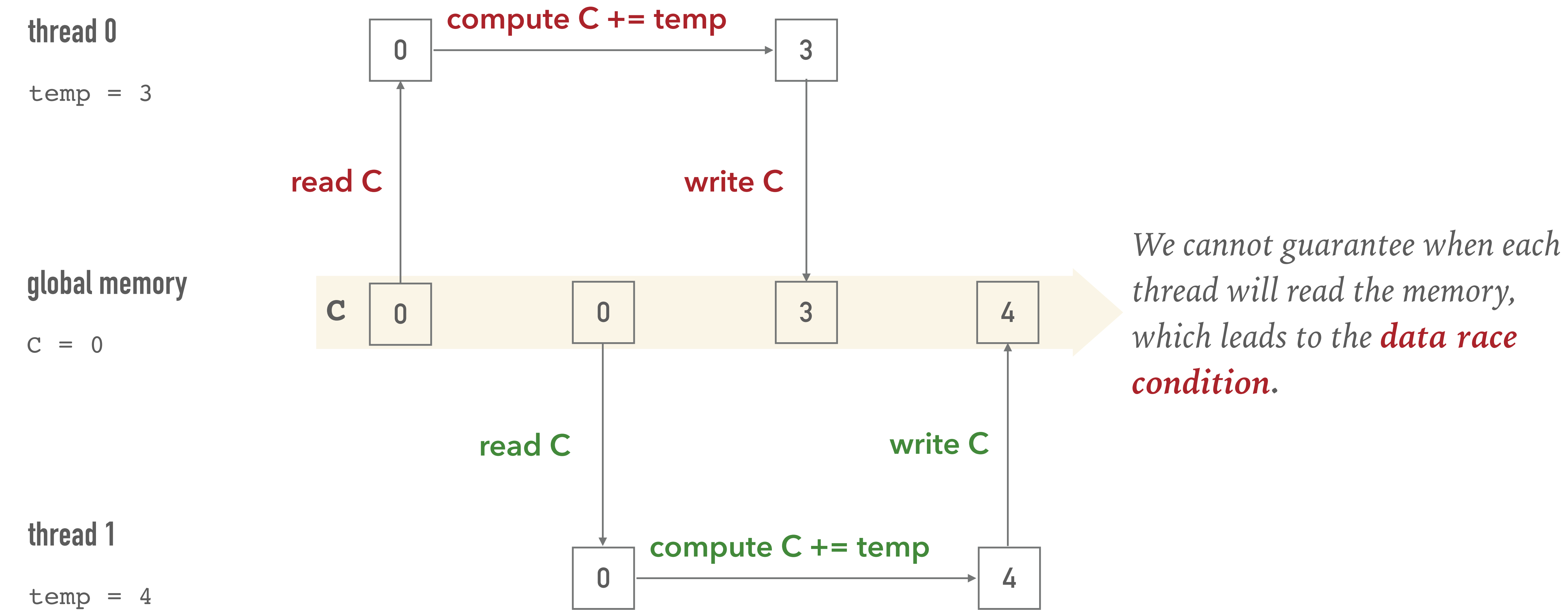
---

Parallel threads have no problem with the multiplication part:

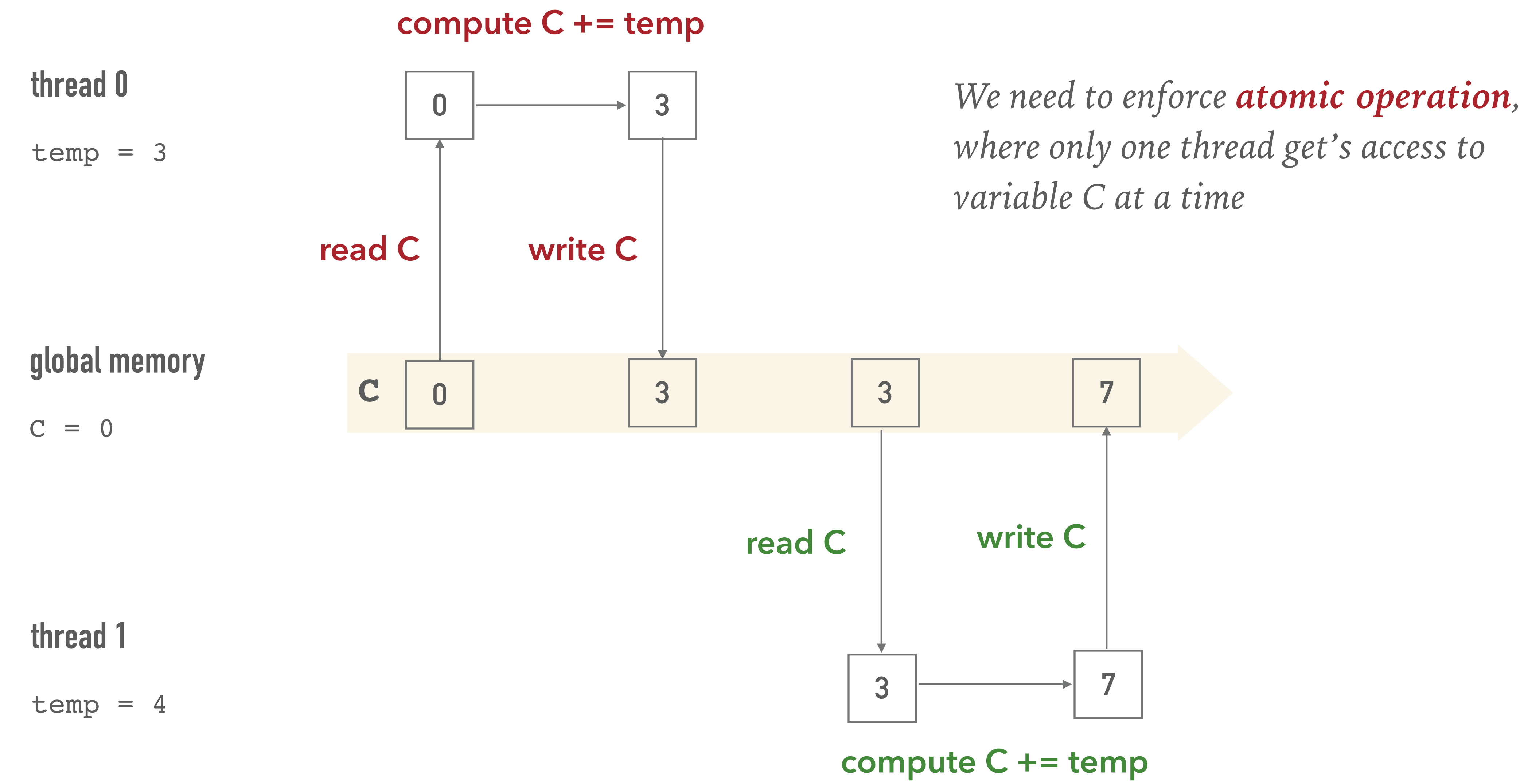


```
__global__ void dot( float *a, float *b, float *c, int N ) {  
    int idx = blockIdx.x * blockDim.x + threadIdx.x;  
    // Each thread computes a pairwise product  
    float temp += a[idx] * b[idx];  
    c += temp;  
}
```

# DATA RACE CONDITION



# DATA RACE CONDITION



# ATOMIC OPERATIONS

---

Atomic operations ensure that only one thread can access and modify a memory location. Other threads need to wait until atomic operation is completed.

- |             |                    |              |                         |
|-------------|--------------------|--------------|-------------------------|
| ➤ atomicAdd | <i>addition</i>    | ➤ atomicInc  | <i>increment</i>        |
| ➤ atomicSub | <i>subtraction</i> | ➤ atomicDec  | <i>decrement</i>        |
| ➤ atomicMin | <i>minimum</i>     | ➤ atomicExch | <i>exchange</i>         |
| ➤ atomicMax | <i>maximum</i>     | ➤ atomicCAS  | <i>compare and swap</i> |

atomicAdd\_block      *atomic with respect to block only*

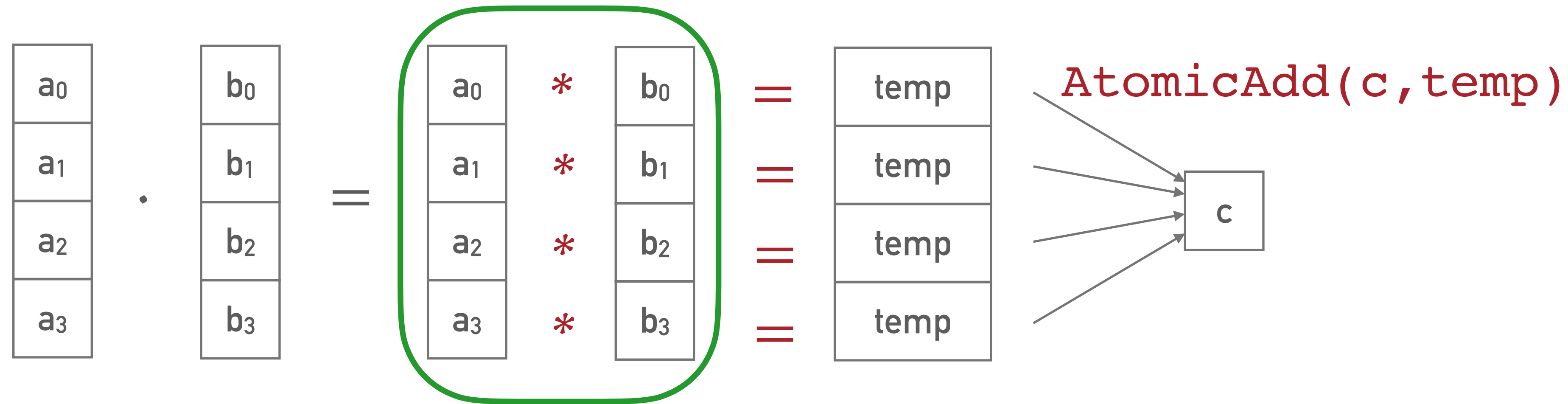
atomicAdd\_system      *atomic with respect to GPU and CPU (for unified memory)*



# THREADS SYNCHRONIZATION

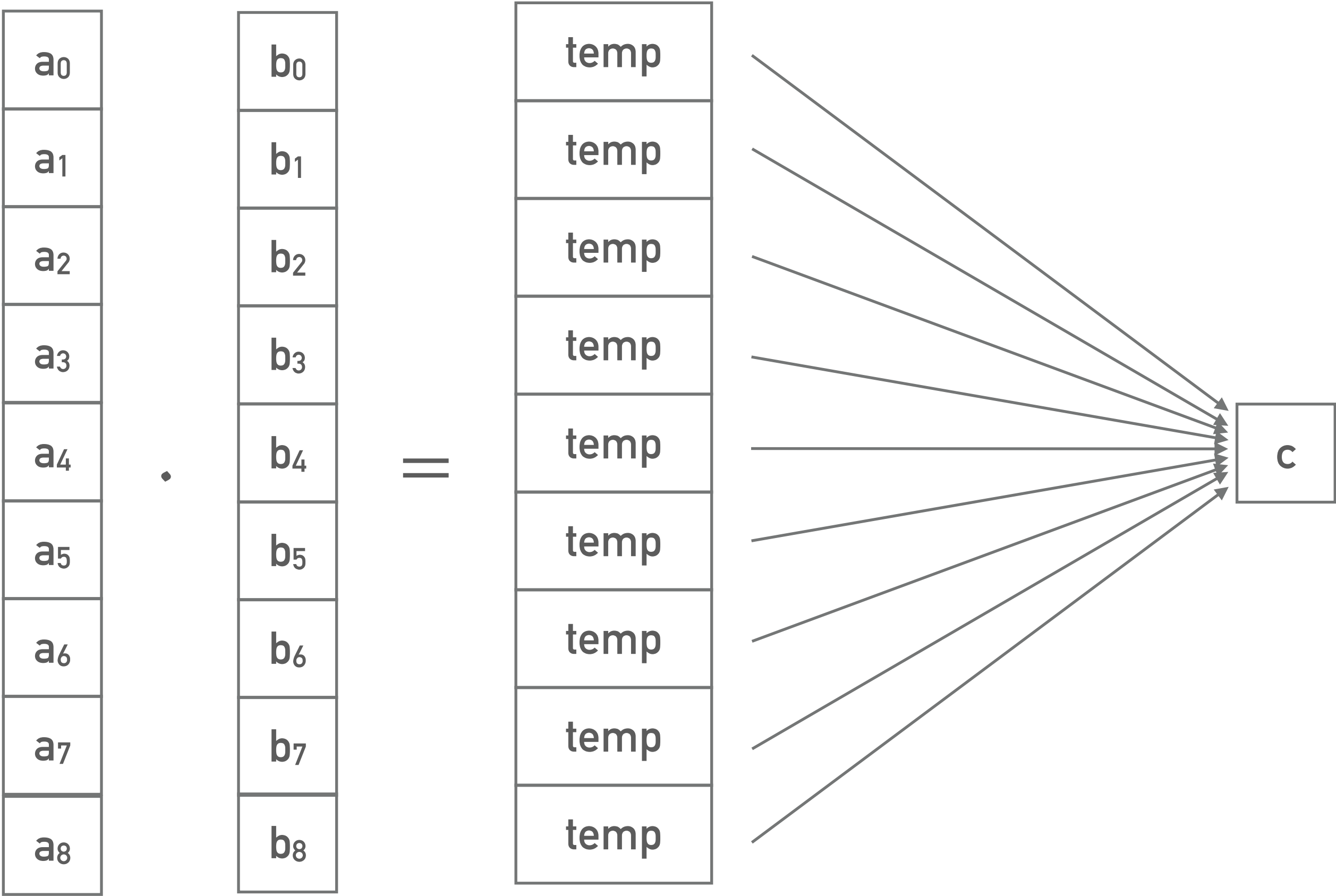
---

Parallel threads have no problem with the multiplication part:



```
__global__ void dot( float *a, float *b, float *c, int N ) {  
    int idx = blockIdx.x * blockDim.x + threadIdx.x;  
    // Each thread computes a pairwise product  
    float temp += a[idx] * b[idx];  
    AtomicAdd(c, temp);  
}
```

# BLOCKS AND THREADS





# BLOCKS AND THREADS

---

block 0

thread 0

a <sub>0</sub>
a <sub>1</sub>
a <sub>2</sub>

thread 1

b <sub>0</sub>
b <sub>1</sub>
b <sub>2</sub>

thread 2

temp
temp
temp

block 1

thread 0

a <sub>3</sub>
a <sub>4</sub>
a <sub>5</sub>

thread 1

b <sub>3</sub>
b <sub>4</sub>
b <sub>5</sub>

thread 2

temp
temp
temp

block 2

thread 0

a <sub>6</sub>
a <sub>7</sub>
a <sub>8</sub>

thread 1

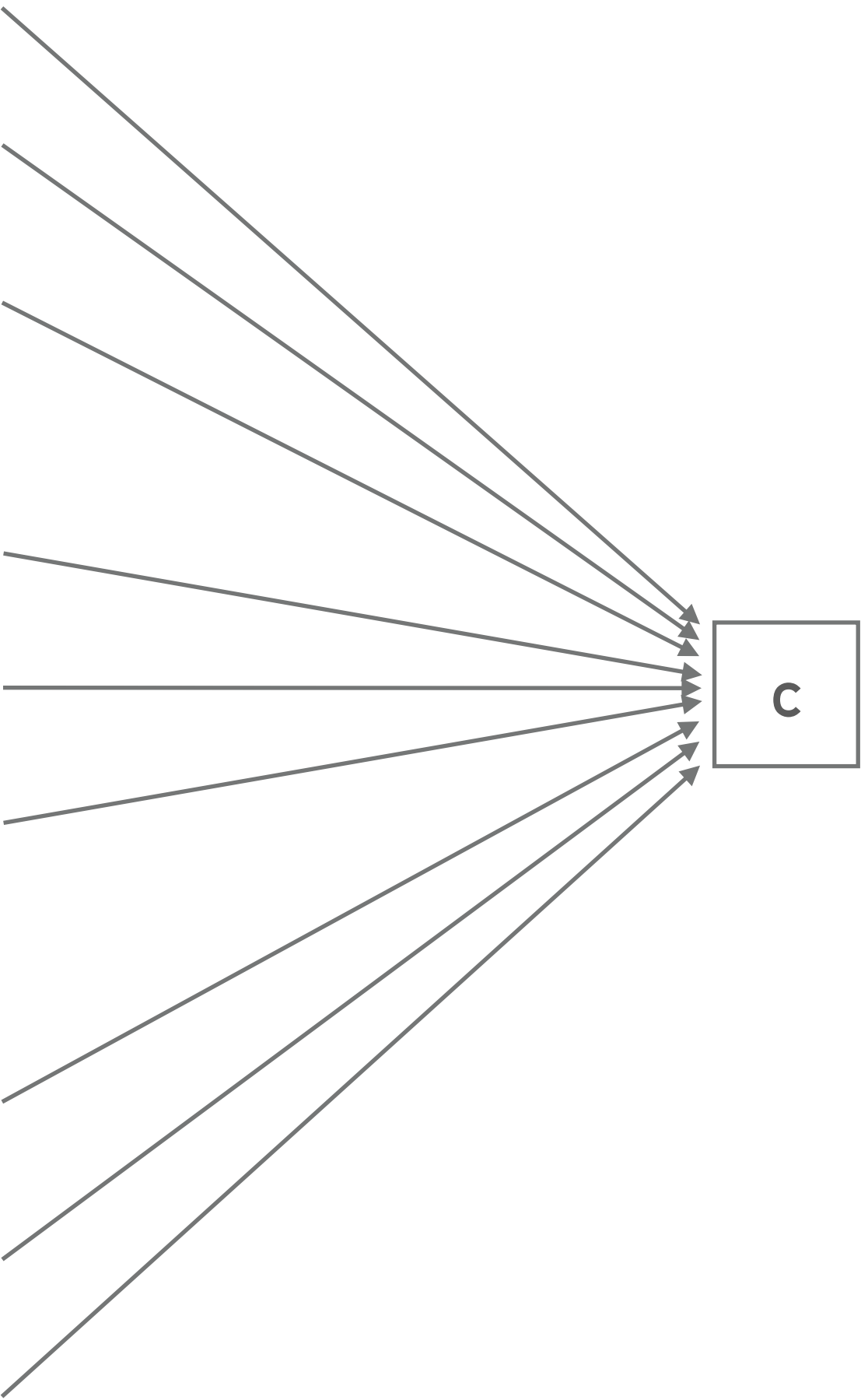
b <sub>6</sub>
b <sub>7</sub>
b <sub>8</sub>

thread 2

temp
temp
temp

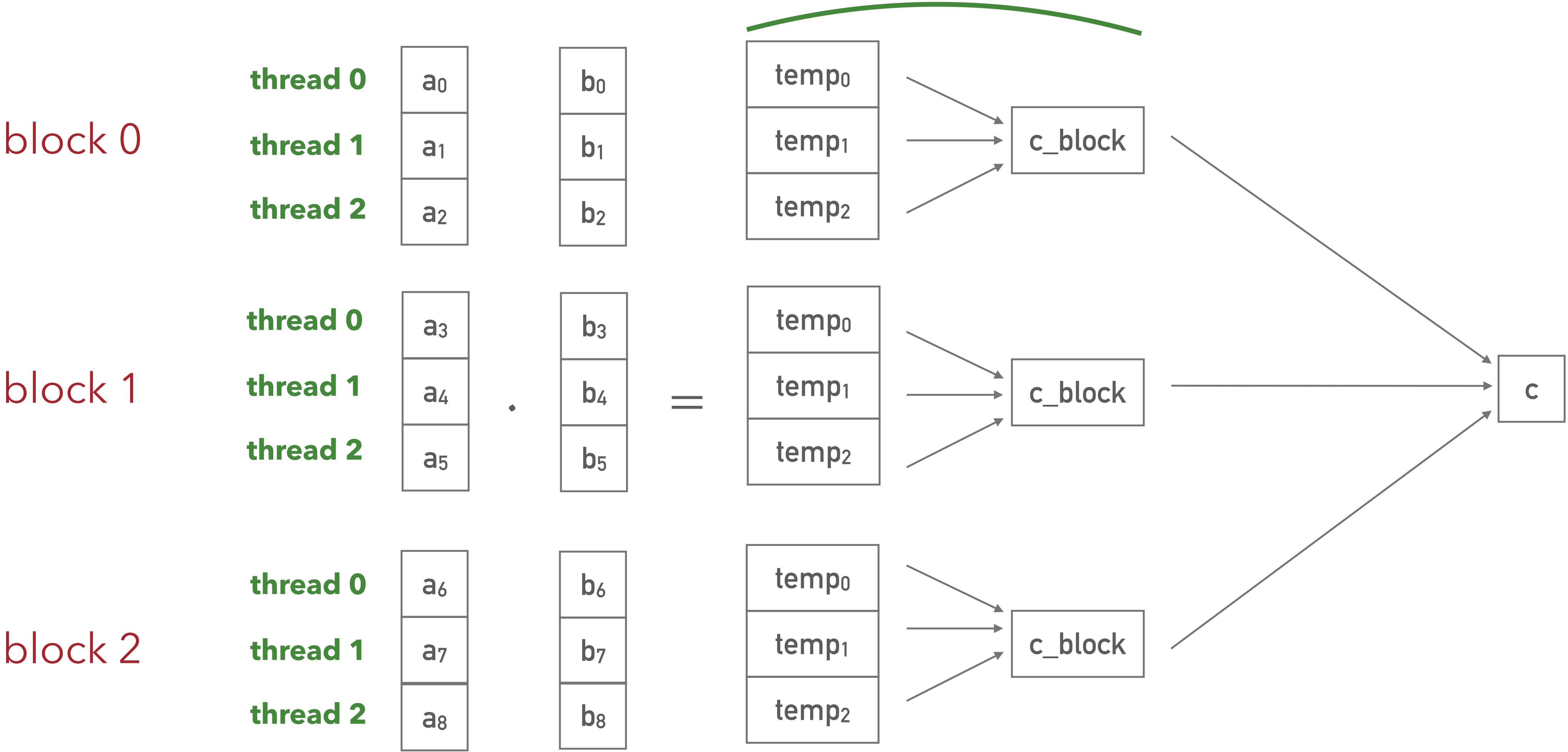
.

=



# BLOCKS AND THREADS

*shared memory within a block*



# THREADS SYNCHRONIZATION

---

Shared memory:

- is shared among the threads, but private to each block
- is extremely fast (think cache, but shared among threads)

```
__global__ void dot( int *a, int *b, int *c, int N ) {  
    int idx = blockIdx.x * blockDim.x + threadIdx.x;  
  
    // Each thread computes a pairwise product  
    __shared__ int temp[THREADS_PER_BLOCK];  
    temp[threadIdx.x] = a[idx] * b[idx];  
}
```

# THREADS SYNCHRONIZATION

---

```
__global__ void dot( int *a, int *b, int *c, int N ) {
    int idx = blockIdx.x * blockDim.x + threadIdx.x;

    // Each thread computes a pairwise product
    __shared__ int temp[TREADS_PER_BLOCK];
    temp[threadIdx.x] = a[idx] * b[idx];

    __syncthreads();

    // Thread 0 sums up the pairwise products
    if(threadIdx.x == 0) {
        float block_c = 0;
        for (int i = 0; i<N; i++)
            block_c += temp[i];
        atomicAdd(c,block_c);
    }
}
```