# CUDA SHARED MEMORY

**NVIDIA Corporation**

- Difference between *host* and *device*

  *Host*         CPU

  *Device*       GPU

- Using **\_\_global\_\_** to declare a function as device code

  - Executes on the device

  - Called from the host (or possibly from other device code)

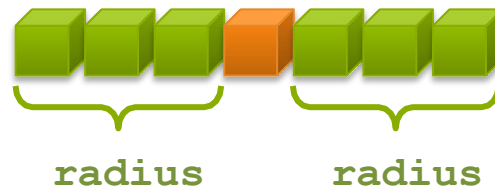- Passing parameters from host code to a device function

- Basic device memory management

  - **`cudaMalloc()`**

  - **`cudaMemcpy()`**

  - **`cudaFree()`**

- Launching parallel kernels

  - Launch **`N`** copies of **`add()`** with **`add<<<N,1>>>(…);`**

  - Use **`blockIdx.x`** to access block index

# Shared Memory

- Shared memory is equivalent to a user-managed cache: The application explicitly allocates and accesses it.
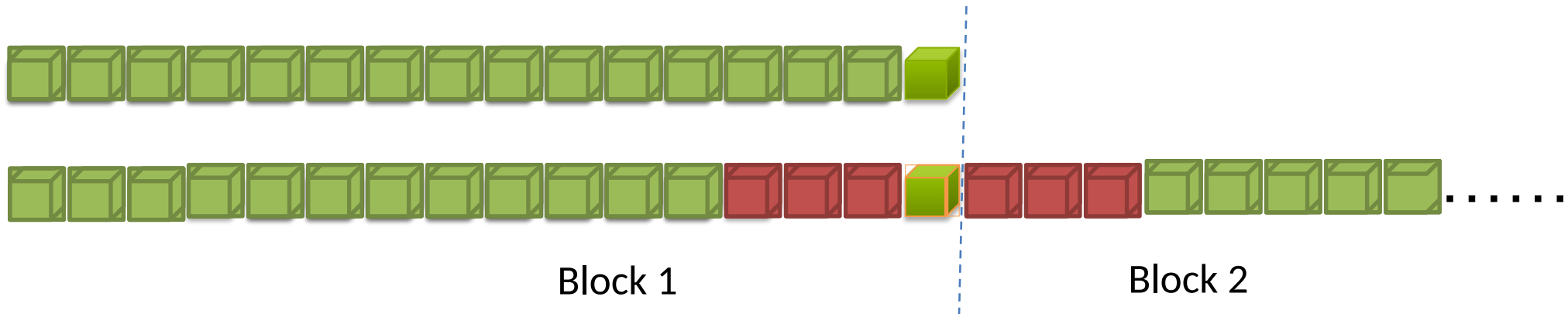
> Recall: how does cache work in OS and computer organization courses

- Launching parallel kernels

  - Launch **N** copies of **add()** with **add<<<N,1>>>(...);**

  - Use **blockIdx.x** to access block index

# 1D STENCIL

- Consider applying a 1D stencil to a 1D array of elements

    - Each output element is the sum of input elements within a radius

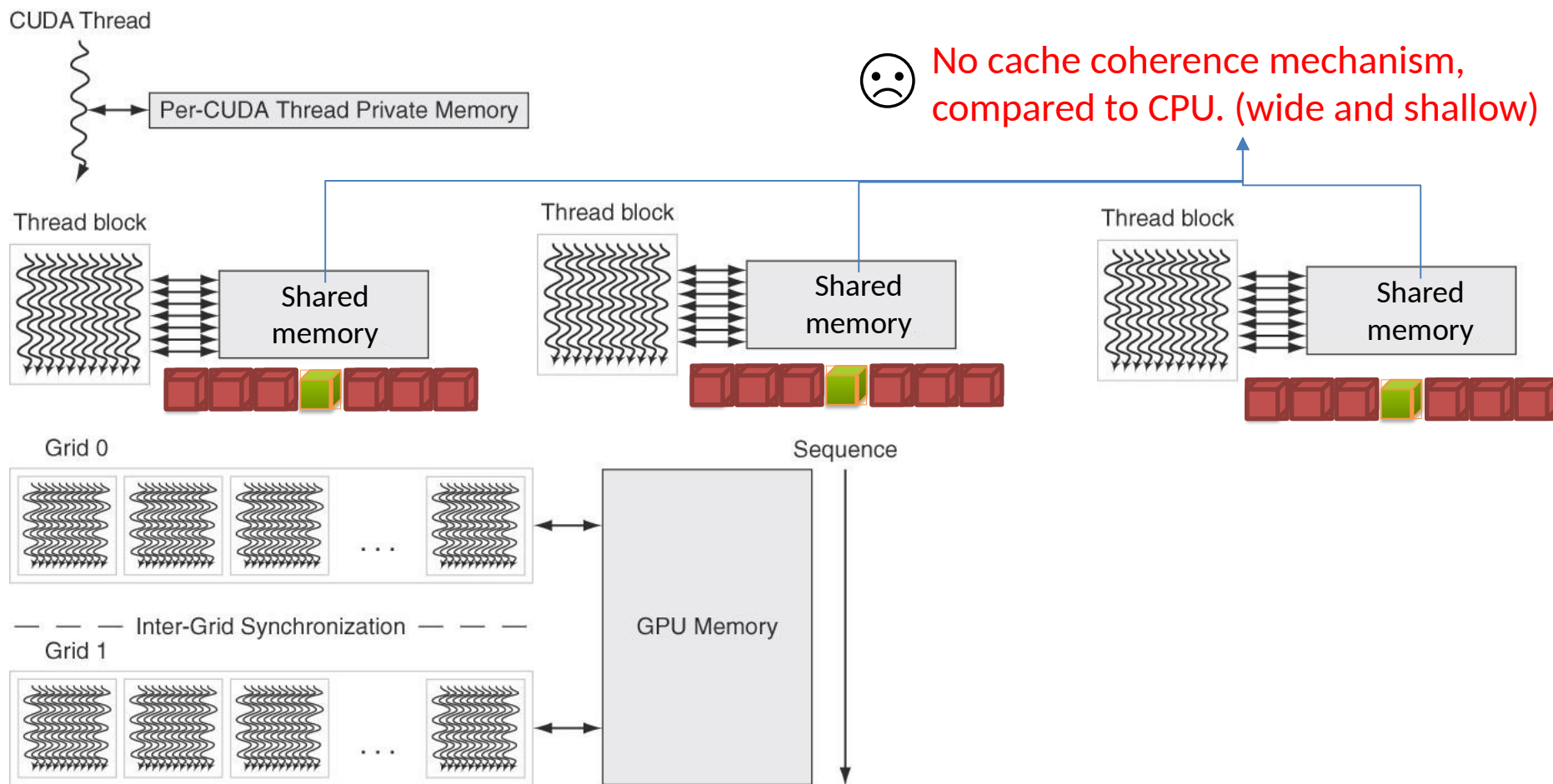- If radius is 3, then each output element is the sum of 7 input elements:



radius        radius

- Each thread processes one output element

  - **`blockDim.x`** elements per block

- Input elements are read several times

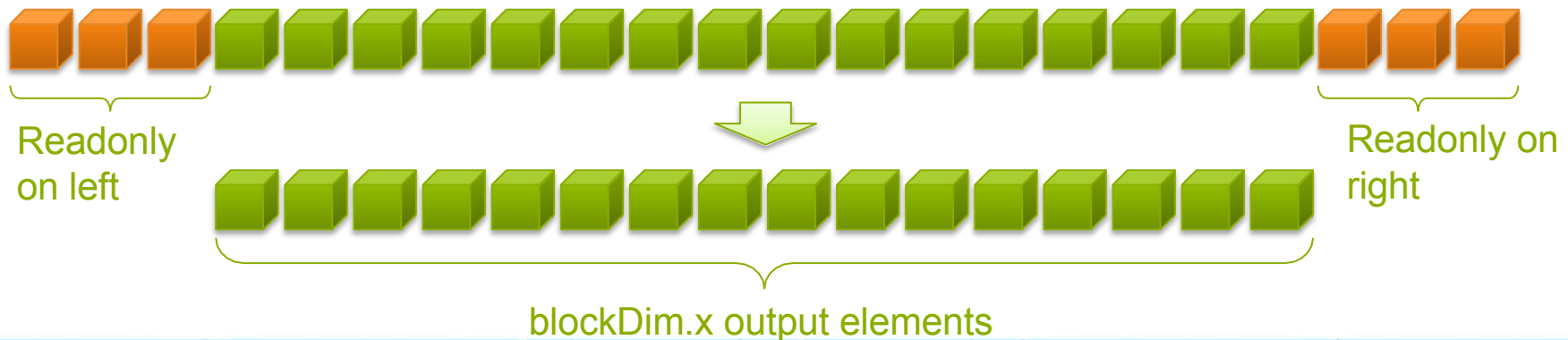  - With radius 3, each input element is read seven times



Block 1                              Block 2

- Terminology: within a block, threads share data via <span style="color:orange">shared memory</span>

- Shared memory is equivalent to a user-managed cache:

  - The application explicitly allocates and accesses it.

- Extremely fast on-chip memory, user-managed

- Declare using `__shared__`, allocated a variable per block that:
  - Resides in the shared memory space of a thread block,
  - Has the lifetime of the block,
  - Has a distinct object per block,
  - Is only accessible from all the CUDA threads within the block,

- Typical programming pattern:

  - Load data from device memory to shared memory

  - Synchronize with all the other threads of the block so that each thread can safely read shared memory locations that were populated by different threads,

  - Process the data in shared memory

  - Synchronize again if necessary to make sure that shared memory has been updated with the results,

  - Write the results back to device memory.

CUDA Thread

Per-CUDA Thread Private Memory

☹ No cache coherence mechanism, compared to CPU. (wide and shallow)

Thread block

Shared memory

Thread block

Shared memory

Thread block

Shared memory

Grid 0

Sequence

GPU Memory

Inter-Grid Synchronization

Grid 1

- Cache data in shared memory
  - Read (`blockDim.x` + 2 * radius) input elements from global memory to shared memory
  - Compute `blockDim.x` output elements
  - Write `blockDim.x` output elements to global memory

- Each block needs a halo of `radius` elements at each boundary

Readonly on left

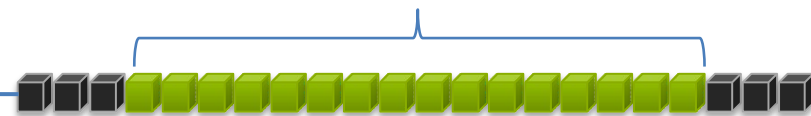Readonly on right

blockDim.x output elements

```
__global__ void stencil_1d(int *in, int *out) {
  __shared__ int temp[BLOCK_SIZE + 2 * RADIUS];
  int gindex = threadIdx.x + blockIdx.x * blockDim.x;
  int lindex = threadIdx.x + RADIUS;

  // Read input elements into shared memory
  temp[lindex] = in[gindex];
  if (threadIdx.x < RADIUS) {
    temp[lindex - RADIUS] = in[gindex - RADIUS];
    temp[lindex + BLOCK_SIZE] =in[gindex+BLOCK_SIZE];
  }
}
```

A thread block[0-15] calculates stencil

The 3 CUDA threads move left 3 orange elements from in[] to temp[]
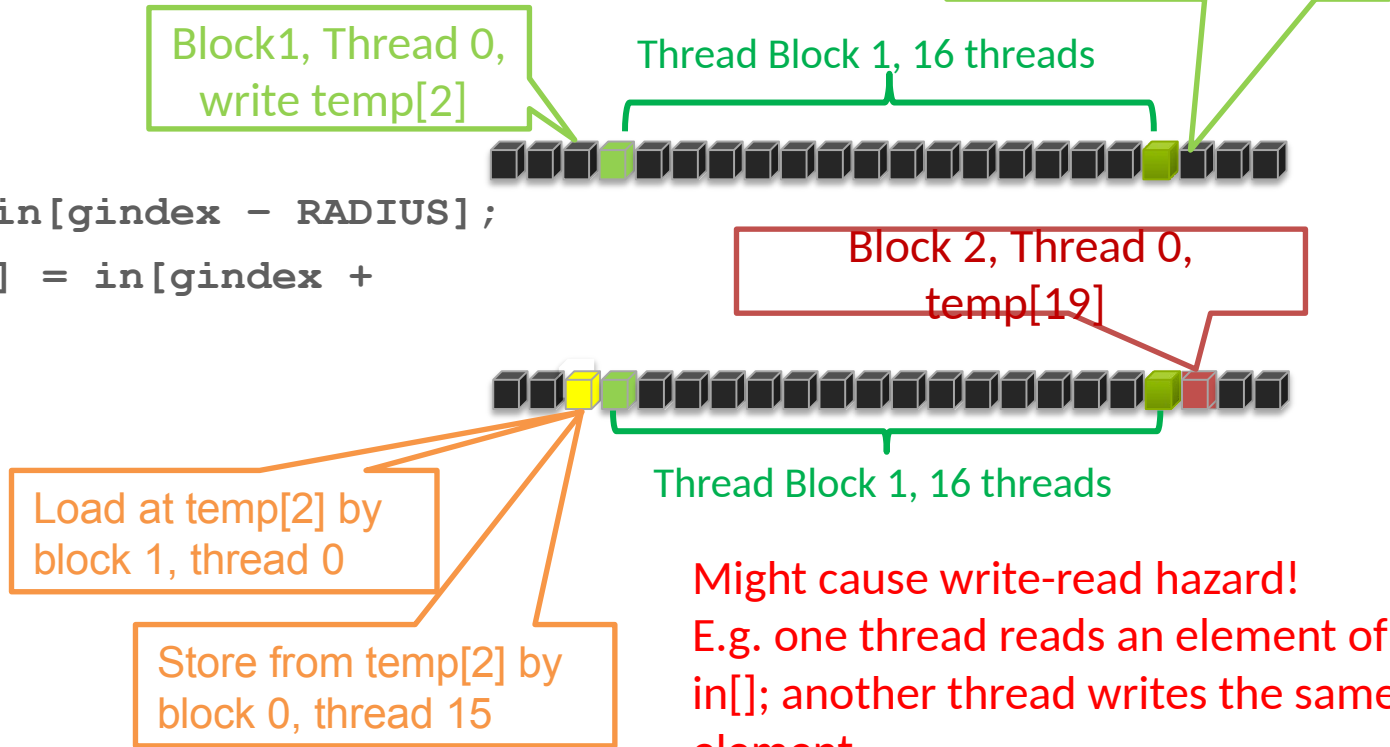
The 3 CUDA threads move right 3 orange elements from in[] to temp[]

```
// Apply the stencil in shared memory
int result = 0;
 for (int offset = -RADIUS ; offset <= RADIUS ; offset++)
   result += temp[lindex + offset];


 // Store the result
 out[gindex] = result;
}
```

- The stencil example will not work...
- Suppose thread 15 in Block 1 reads the temp[19] for calculating stencil after thread 0 in block 2 has fetched temp[19] from global memory in[].

Block 1, Thread 15, write temp[19]

```
temp[lindex] = in[gindex];
if (threadIdx.x < RADIUS) {
    temp[lindex - RADIUS] = in[gindex - RADIUS];
    temp[lindex + BLOCK_SIZE] = in[gindex +
    BLOCK_SIZE];
}
int result = 0;
result += temp[lindex + 1];
```

Block1, Thread 0, write temp[2]

Thread Block 1, 16 threads

Block 2, Thread 0, temp[19]

Load at temp[2] by block 1, thread 0

Thread Block 1, 16 threads

Store from temp[2] by block 0, thread 15

Might cause write-read hazard! E.g. one thread reads an element of in[]; another thread writes the same element.

13

- **`void _syncthreads();`**

- Synchronizes all threads within a CUDA block, how to explain in GPU architecture

  - Used to prevent RAW / WAR / WAW hazards

- All threads must reach the barrier

  - In conditional code, the condition must be uniform across the block

```
__global__ void stencil_1d(int *in, int *out) {
    __shared__ int temp[BLOCK_SIZE + 2 * RADIUS];
    int gindex = threadIdx.x + blockIdx.x * blockDim.x;
     int lindex = threadIdx.x + radius;


    // Read input elements into shared memory
    temp[lindex] = in[gindex];
    if (threadIdx.x < RADIUS) {
        temp[lindex - RADIUS] =
        in[gindex - RADIUS];

        temp[lindex +
        BLOCK_SIZE] = in[gindex
        + BLOCK_SIZE];
    }
    // Synchronize (ensure all
```

```
// Apply the stencil
int result = 0;
for (int offset = -RADIUS ; offset <= RADIUS ; offset++)
    result += temp[lindex + offset];


// Store the result
out[gindex] = result;
}
```

- Use**___ shared___** to declare a variable/array in shared memory

  - Data is shared between threads in a block

  - Not visible to threads in other blocks

- Use**___ syncthreads()** as a barrier

  - Use to prevent data hazards

# FUTURE SESSIONS

- CUDA GPU architecture and basic optimizations

- Atomics, Reductions, Warp Shuffle

- Using Managed Memory

- Concurrency (streams, copy/compute overlap, multi-GPU)

- Analysis Driven Optimization

- Cooperative Groups

# FURTHER STUDY

- Shared memory:

  - https://devblogs.nvidia.com/using-shared-memory-cuda-cc/

- CUDA Programming Guide:

  - https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#shared-memory

- CUDA Documentation:

  - https://docs.nvidia.com/cuda/index.html

    https://docs.nvidia.com/cuda/cuda-runtime-api/index.html (runtime API)