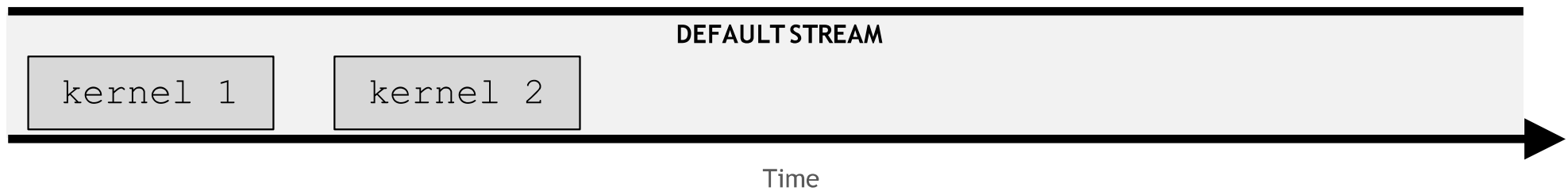# Concurrent CUDA Streams

# DEFAULT STREAM

Time

By default, CUDA kernels run in the **default stream**
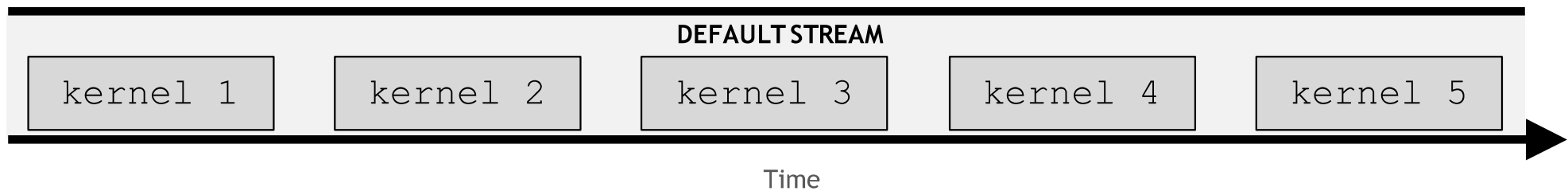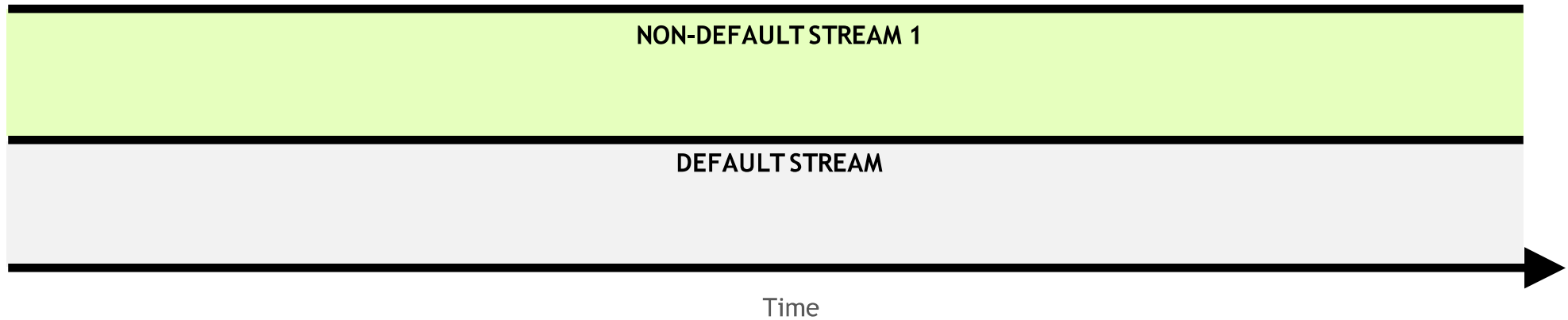
**DEFAULT STREAM**

kernel 1

Time

In any stream, including the default, an instruction in it (here a kernel launch) must complete before the next can begin
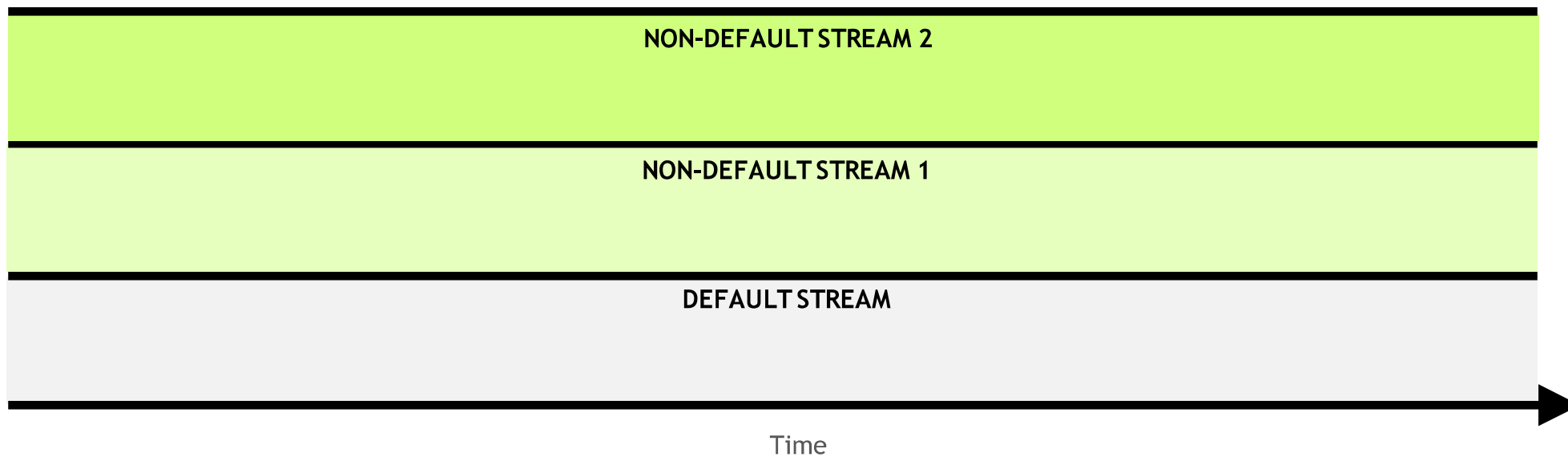
**DEFAULT STREAM**

| kernel 1 | kernel 2 |

Time

**Non-default streams** can also be created for kernel execution

**NON-DEFAULT STREAM 1**

**DEFAULT STREAM**

Time

**Non-default streams** can also be created for kernel execution

NON-DEFAULT STREAM 2

NON-DEFAULT STREAM 1

DEFAULT STREAM

Time

However, kernels in **different, non-default streams**, can interact concurrently

**NON-DEFAULT STREAM 2**

kernel 3

**NON-DEFAULT STREAM 1**

kernel 1

kernel 2

**DEFAULT STREAM**

Time

However, kernels in **different, non-default streams**, can interact concurrently

**NON-DEFAULT STREAM 2**

kernel 3    kernel 4

**NON-DEFAULT STREAM 1**

kernel 1    kernel 2

**DEFAULT STREAM**

Time

However, kernels in **different, non-default streams**, can interact concurrently

**NON-DEFAULT STREAM 2**

kernel 3  kernel 4  kernel 5

**NON-DEFAULT STREAM 1**

kernel 1  kernel 2

**DEFAULT STREAM**

Time

The default stream is special: **it blocks all kernels in all other streams**

NON-DEFAULT STREAM 2

NON-DEFAULT STREAM 1

kernel 1

DEFAULT STREAM

Time

The default stream is special: **it blocks all kernels in all other streams**

NON-DEFAULT STREAM 2

NON-DEFAULT STREAM 1

| kernel 1 | kernel 2 |

DEFAULT STREAM

Time

The default stream is special: **it blocks all kernels in all other streams**

**NON-DEFAULT STREAM 2**

kernel 3

**NON-DEFAULT STREAM 1**

kernel 1

kernel 2

**DEFAULT STREAM**

Time

The default stream is special: **it blocks all kernels in all other streams**

**NON-DEFAULT STREAM 2**

kernel 3    kernel 4

**NON-DEFAULT STREAM 1**

kernel 1    kernel 2

**DEFAULT STREAM**

Time

The default stream is special: **it blocks all kernels in all other streams**

NON-DEFAULT STREAM 2

kernel 3    kernel 4

NON-DEFAULT STREAM 1

kernel 1    kernel 2

DEFAULT STREAM

kernel 5

Time

The default stream is special: **it blocks all kernels in all other streams**

NON-DEFAULT STREAM 2

kernel 3    kernel 4    kernel 6

NON-DEFAULT STREAM 1

kernel 1    kernel 2

DEFAULT STREAM

kernel 5

Time

The default stream is special: **it blocks all kernels in all other streams**

NON-DEFAULT STREAM 2

kernel 3    kernel 4    kernel 6

NON-DEFAULT STREAM 1

kernel 1    kernel 2    kernel 7

DEFAULT STREAM

kernel 5

Time

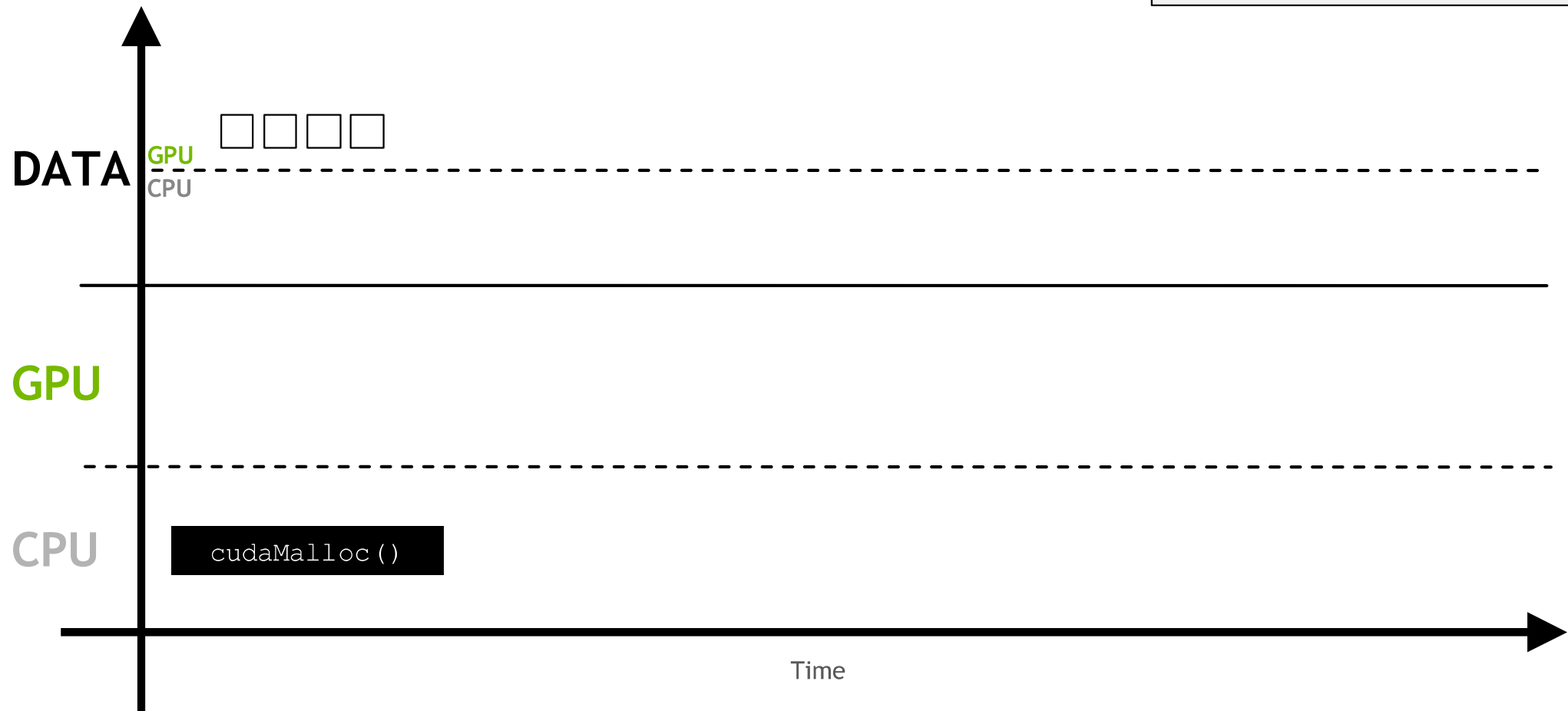The default stream is special: **it blocks all kernels in all other streams**

**NON-DEFAULT STREAM 2**

kernel 3    kernel 4    kernel 6

**NON-DEFAULT STREAM 1**

kernel 1    kernel 2    kernel 7

**DEFAULT STREAM**

kernel 5

Time

# Non-Unified Memory

Memory can be allocated directly to
the GPU with `**cudaMalloc**`

DATA
GPU
CPU

GPU

CPU
`cudaMalloc()`

Time

Memory can be allocated directly to the host with `cudaMallocHost`

DATA
GPU
CPU

GPU

CPU
cudaMallocHost()

Time

Memory allocated in either of these ways can be **copied** to other locations in the system with `cudaMemCpy`

**DATA**
GPU
CPU

**GPU**

**CPU**

cudaMallocHost()   cudaMemcpy(HtoD)

Time

DATA

GPU

CPU

Copying leaves 2 copies in of in the system

□
□
□
□

□ □ □ □

GPU

CPU

cudaMallocHost()    cudaMemcpy(HtoD)

Time

# cudaMemcpyAsync

This can allow the **overlapping** memory copies and computation

DATA GPU
CPU

GPU

CPU

`cudaMallocHost()` `cpy`

Time

This can allow the **overlapping** memory copies and computation

DATA
GPU
CPU

GPU

work

CPU

`cudaMallocHost()`  `cpy`

Time

This can allow the **overlapping** memory copies and computation

DATA GPU
CPU

GPU

work work

CPU

cudaMallocHost() cpy cpy

Time

This can allow the **overlapping** memory copies and computation

DATA
GPU
CPU

GPU

work   work   work

CPU

cudaMallocHost()   cpy   cpy   cpy

Time

This can allow the **overlapping** memory copies and computation

DATA
GPU
CPU

GPU

work  work  work  work

CPU

cudaMallocHost()  cpy  cpy  cpy  cpy

Time