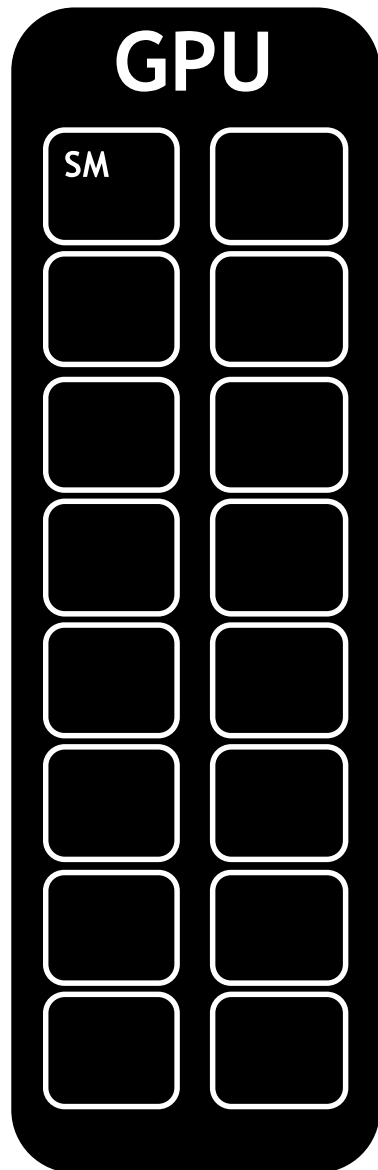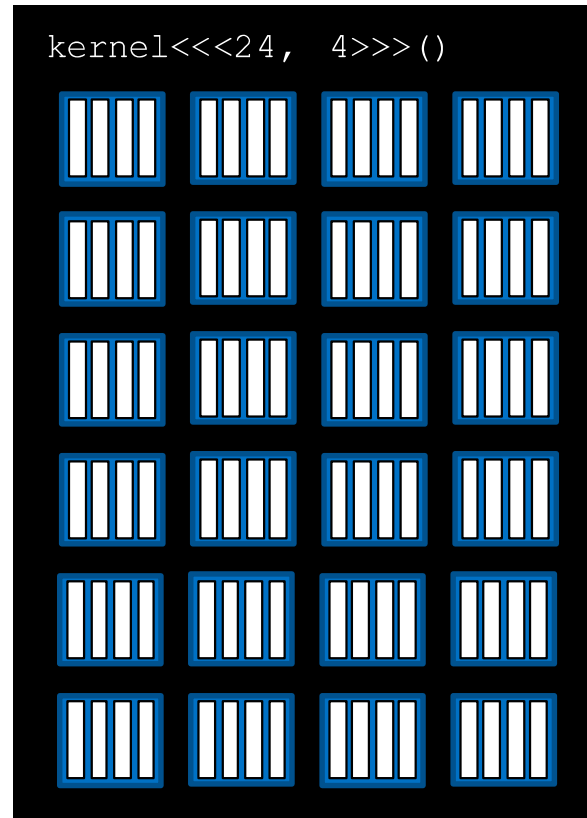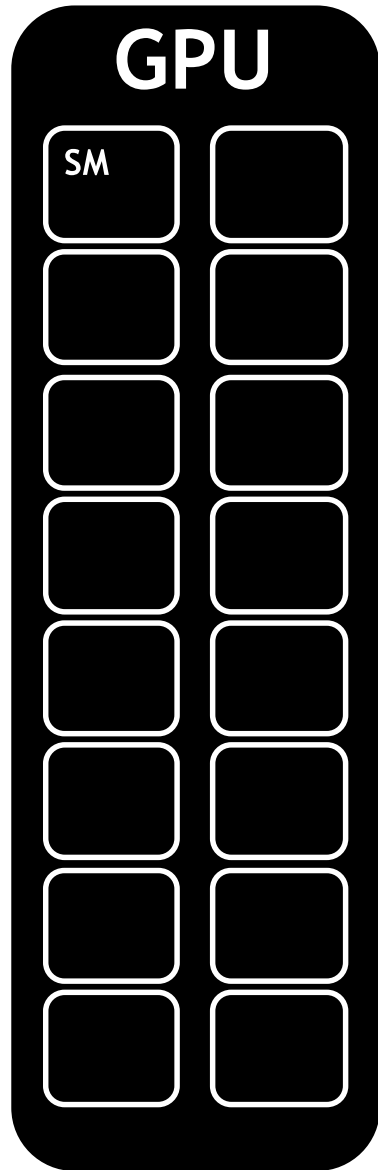# Streaming Multiprocessors

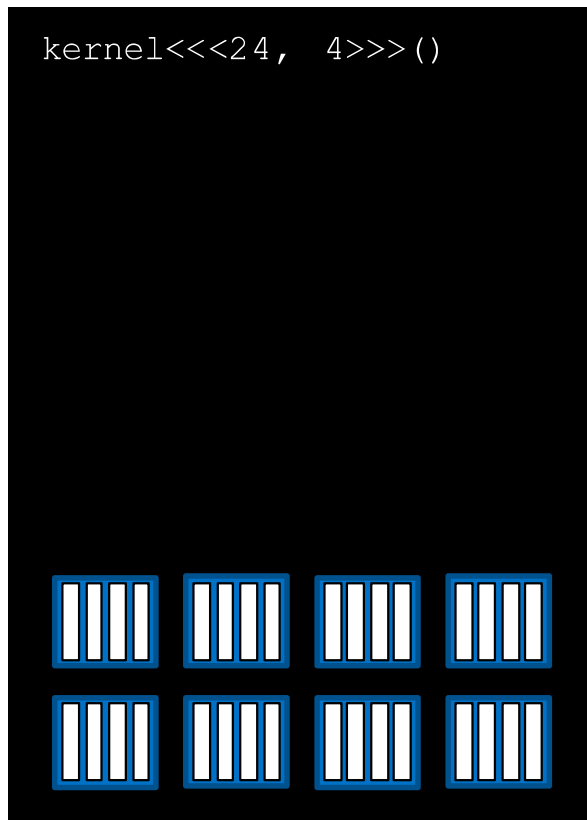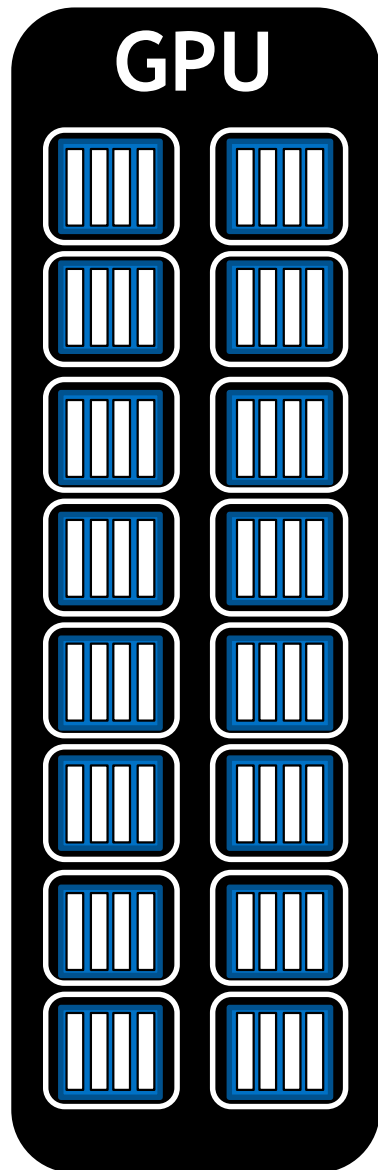**GPU**

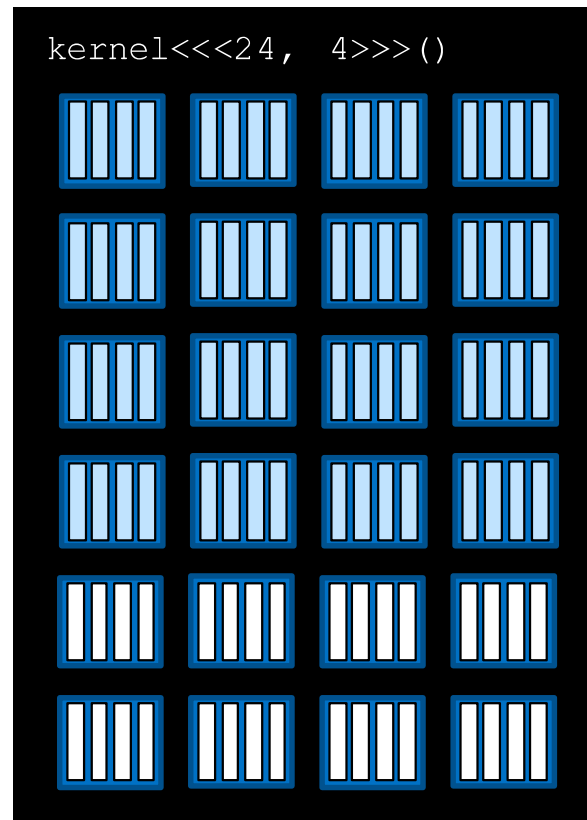NVIDIA GPUs contain functional units called **Streaming Multiprocessors**, or **SMs**

# GPU

SM

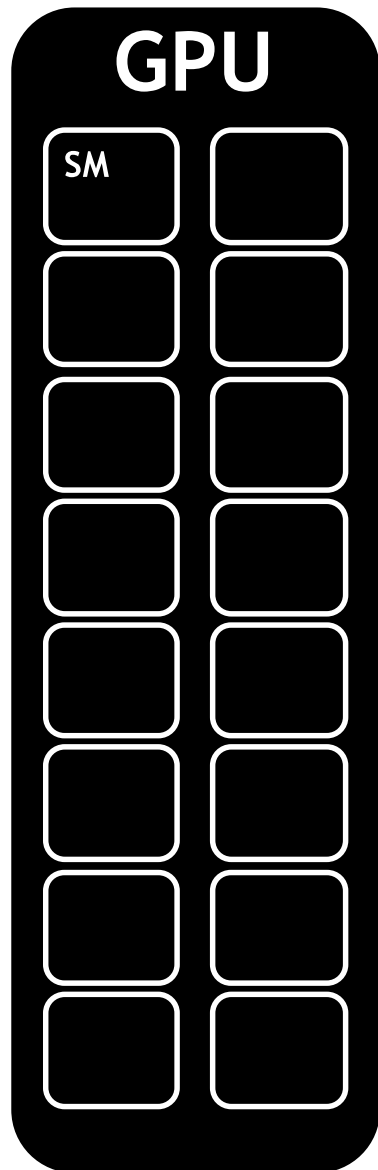NVIDIA GPUs contain functional units called **Streaming Multiprocessors**, or **SMs**

# GPU

SM

kernel<<<24, 4>>>()

Blocks of threads are scheduled to run on SMs

# GPU

```
kernel<<<24, 4>>>()
```

Depending on the number of SMs on a GPU, and the requirements of a block, more than one block can be scheduled on an SM

# GPU

SM

```
kernel<<<24, 4>>>()
```
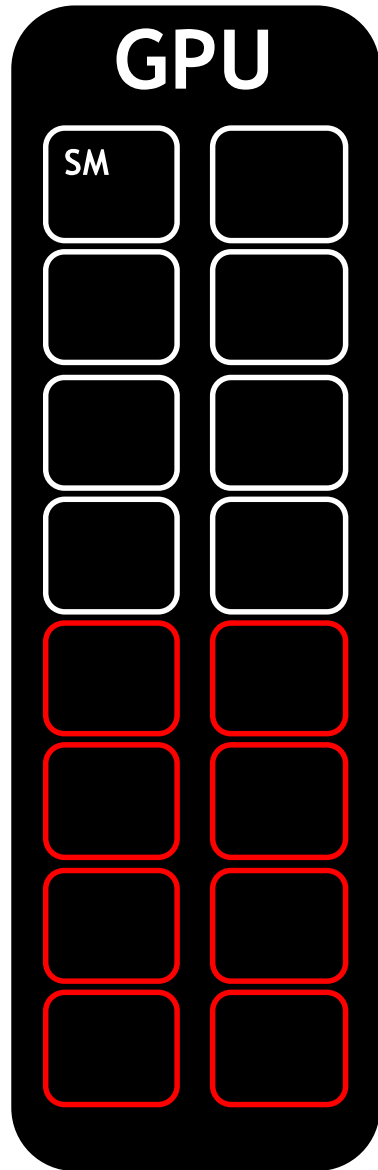
Depending on the number of SMs on a GPU, and the requirements of a block, more than one block can be scheduled on an SM

# GPU

kernel<<<24, 4>>>()

Grid dimensions divisible by the number of SMs on a GPU can promote full SM utilization

# GPU

SM

```
kernel<<<24, 4>>>()
```

# Unified Memory Behavior

When **UM** is allocated, it may not be resident initially on the CPU or the GPU

DATA
GPU
CPU

?

GPU

CPU

cudaMallocManaged()

Time

When some work asks for the memory for the first time, a **page fault** will occur

DATA
GPU
CPU

GPU

CPU

`cudaMallocManaged()`   `init()`

?

Time

The page fault will trigger the migration of the demanded memory

DATA

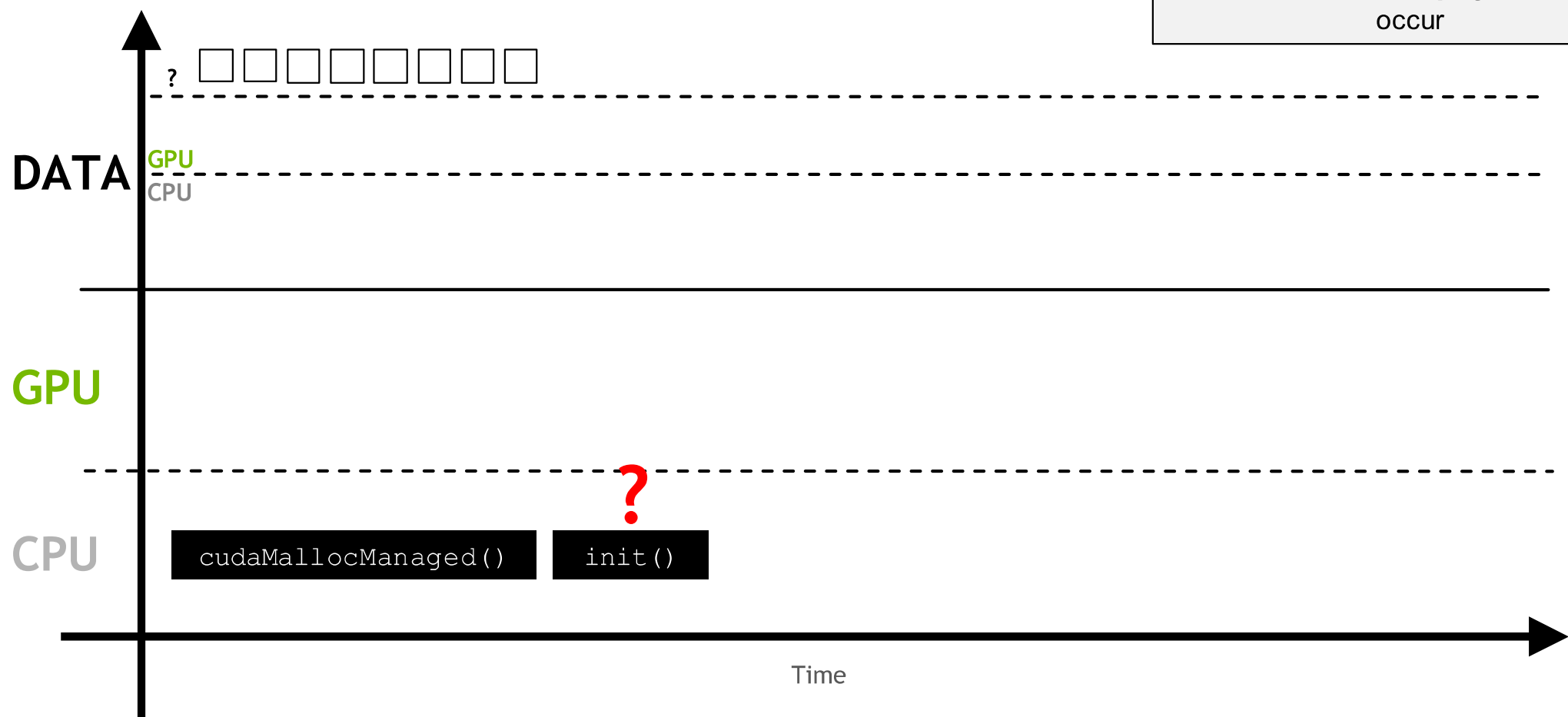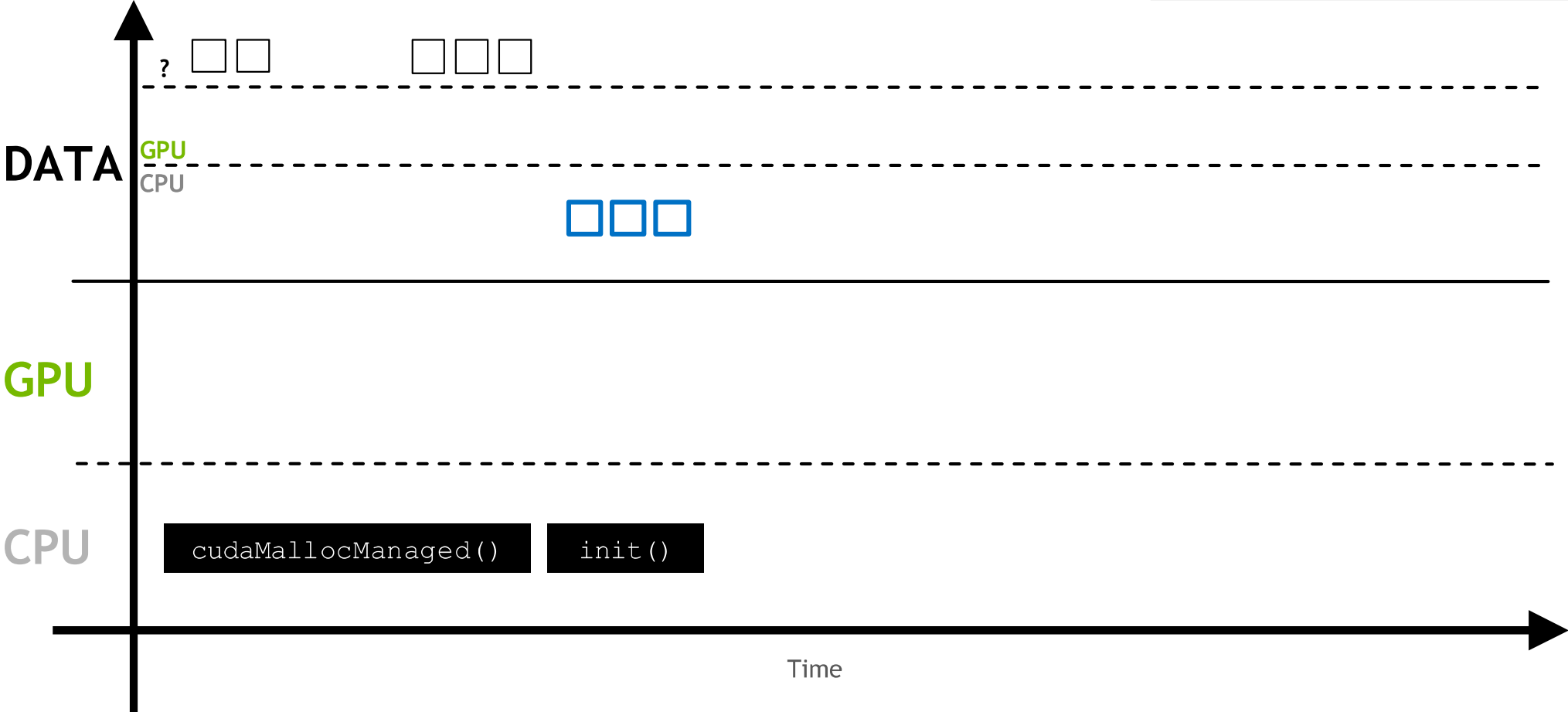GPU
CPU

GPU

CPU

cudaMallocManaged()    init()

Time

DATA
GPU
CPU

GPU

CPU

?

work<<<>>>()

cudaMallocManaged()    init()

Time

This process repeats anytime the memory is requested somewhere in the system where it is not resident

If it is known that the memory **will be** accessed somewhere it is not resident, asynchronous prefetching can be used

**DATA** GPU
CPU

**GPU**

```
work<<<>>>()
```

**CPU**

```
cudaMallocManaged()
```

```
init()
```

```
cudaMemPrefetchAsync(cpu)
```

Time

This moves the memory in larger batches, and prevents page faulting

**DATA** GPU
CPU

? □ □        □ □ □

□ □ □

**GPU**

`work<<<>>>()`

**CPU**

`cudaMallocManaged()`    `init()`    `cudaMemPrefetchAsync(cpu)`    `check()`

Time

www.nvidia.com/dli