# Three possible on-the-fly timing strategies
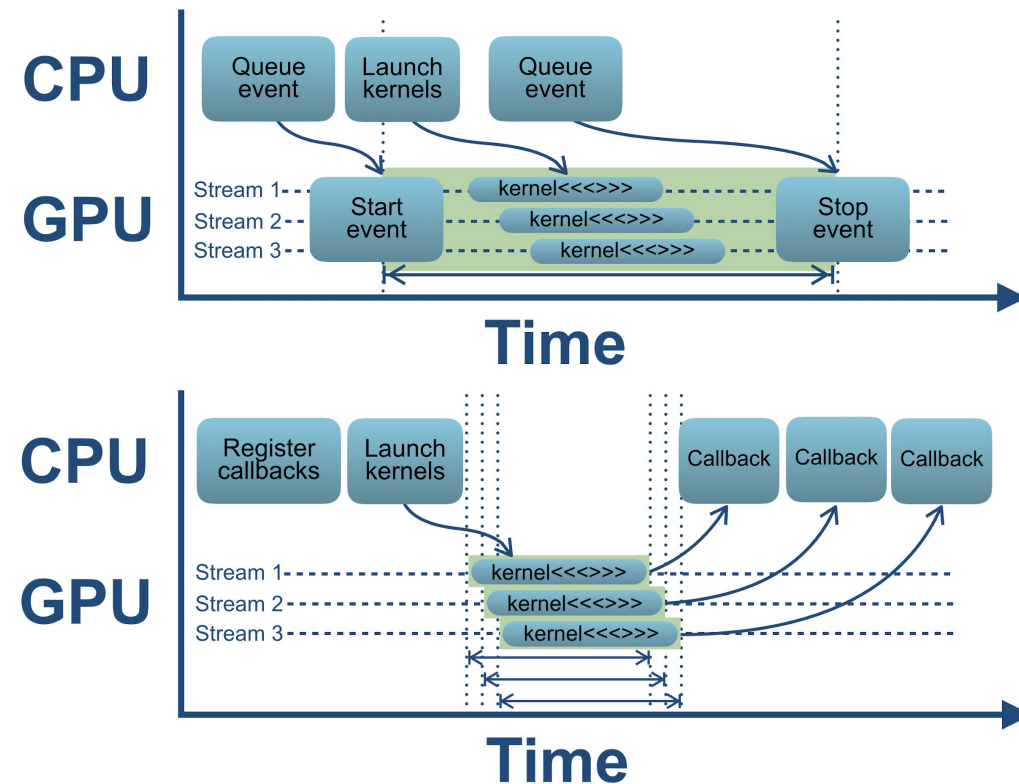
1. Count GPU clock cycles
   - Requires additional device to host transfers
   - Implementation may be invasive

2. CUDA Events
   - Can give ambiguous results

3. CUDA Profiling Tools Interface (CUPTI)
   - Buffer requests and delivery of timing information handled by CUPTI
   - Gives unambiguous kernel timings

```
1.   Start GPU timer
2.   Do GPU kernel work
3.   Stop GPU timer
4.   Send elapsed time to host
```

BERKELEY LAB

U.S. DEPARTMENT OF ENERGY | Office of Science

# Three possible on-the-fly timing strategies

1. ~~Count GPU clock cycles~~
   - ~~Requires additional device-to-host transfers~~
   - ~~Implementation may be invasive~~

2. ~~CUDA Events~~
   - ~~Can give ambiguous results~~

3. CUDA Profiling Tools Interface (CUPTI)
   - Buffer requests and delivery of timing information handled by CUPTI
   - Gives unambiguous kernel timings

```
1.    Start GPU timer
2.    Do GPU kernel work
3.    Stop GPU timer
4.    Send elapsed time to host
```

**CPU** | Queue event | Launch kernels | Queue event |

**GPU**
Stream 1 — Start event — kernel<<<>>>
Stream 2 — kernel<<<>>> — Stop event
Stream 3 — kernel<<<>>>

**Time**

**CPU** | Register callbacks | Launch kernels | | Callback | Callback | Callback |

**GPU**
Stream 1 —— kernel<<<>>>
Stream 2 —— kernel<<<>>>
Stream 3 —— kernel<<<>>>

**Time**

BERKELEY LAB

U.S. DEPARTMENT OF ENERGY | Office of Science