# Multicore

QNX

# Topics:

→ **Overview**
**Scheduling**
**Synchronization**
**Conclusion**

# SMP:

- is short for Symmetrical MultiProcessing
  - now generally called multicore
  - still used in naming the kernel
- means a system that has more than one processor/CPU tightly coupled
  - independent processors
  - shared hardware, RAM, bus, etc
- multicore processors, where the processors share one chip, are a common modern case of this

# To use multicore, nothing needs to be done:

– the kernel is already multicore aware as of SDP7.0

  • you don't have to write any special code

– **procnto-smp-instr** is designed to support single and multi-core processors

  • if run on a single-core processor, will run in single-core mode

  • on a multi-core processor, will run in multi-core mode

  • the BSP must support multi-core as well

Topics:

**Overview**

→ **Scheduling**
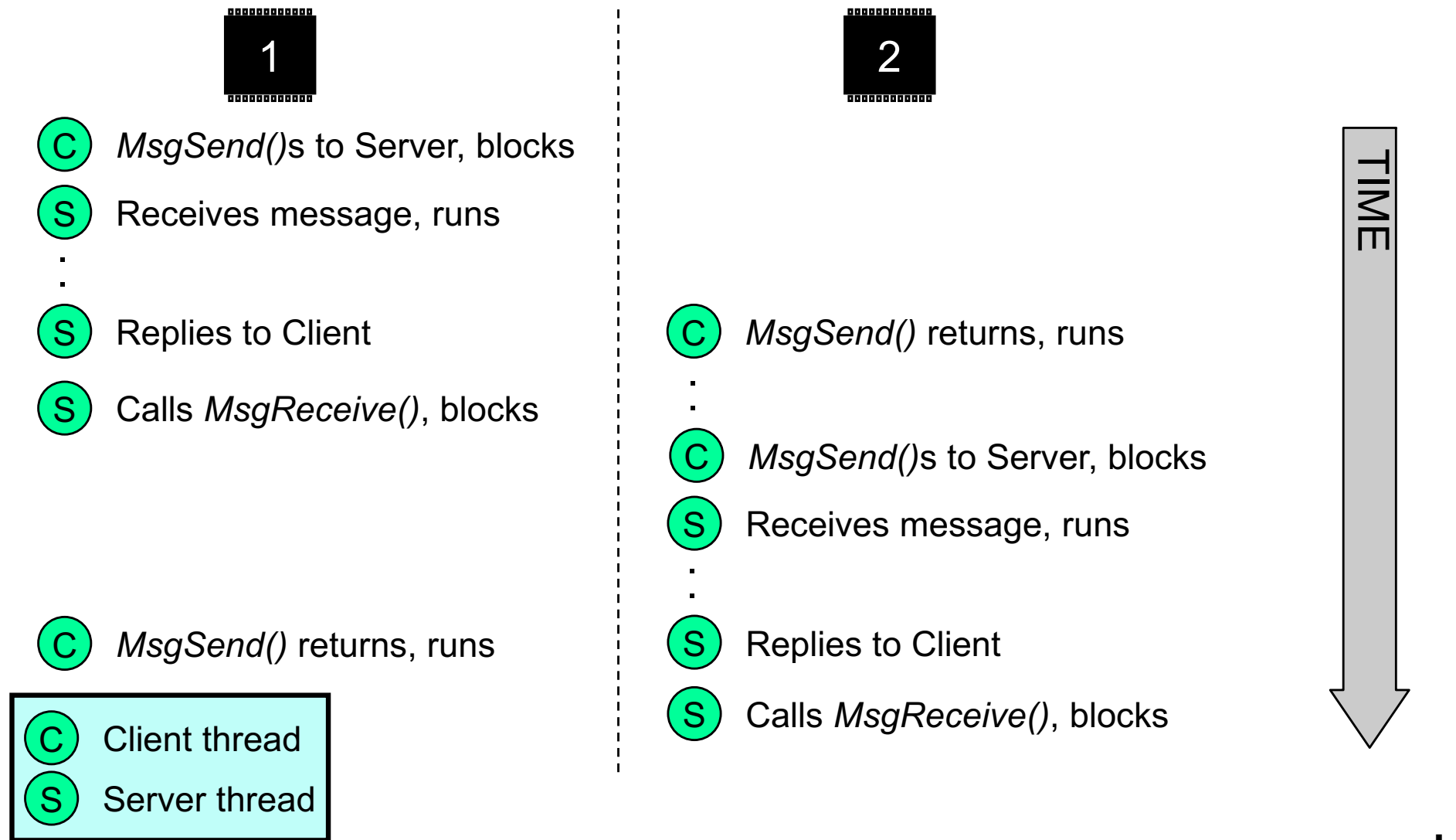
**Synchronization**

**Conclusion**

# With the multicore scheduler:

- when a thread becomes runnable, usually:
  - we find the core with the lowest priority running thread
  - if the newly runnable thread is higher priority than that thread, run the new thread on that core
  - this behavior can be configured with:

    ```
    SchedCtl( SCHED_CONFIGURE, … );
    ```

- we guarantee that the highest priority thread that is ready to use the CPU is running
  - this is the same thread that would be scheduled uniprocessor
- but, sometimes lower priority thread(s) may be running while higher priority thread(s) are ready
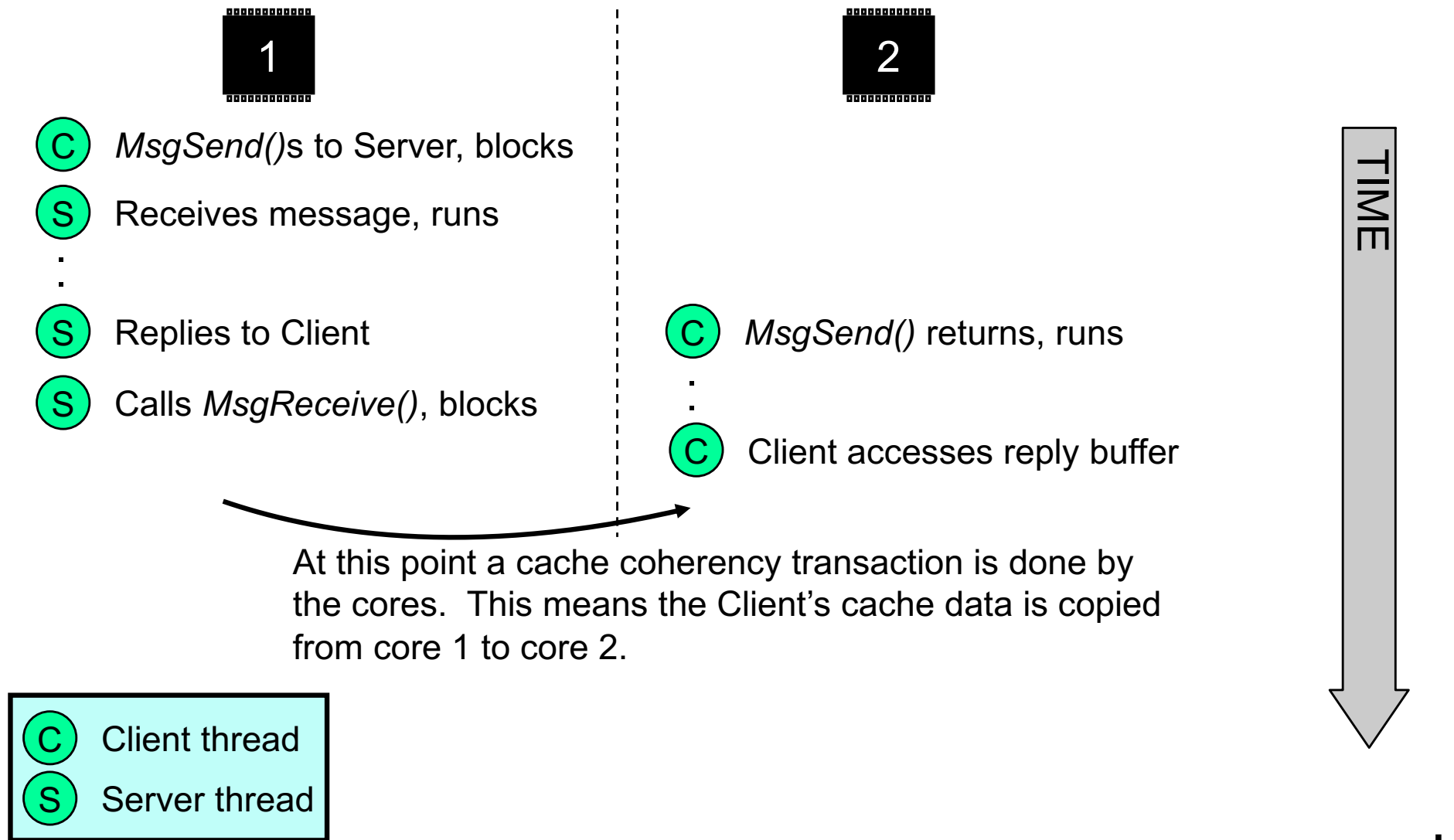  - a trade-off between latency and throughput…

# Bad throughput (not actual behavior):

**1**

**2**

(C) *MsgSend()*s to Server, blocks

(S) Receives message, runs

(S) Replies to Client

(S) Calls *MsgReceive()*, blocks

(C) *MsgSend()* returns, runs

(C) *MsgSend()* returns, runs

(C) *MsgSend()*s to Server, blocks

(S) Receives message, runs

(S) Replies to Client

(S) Calls *MsgReceive()*, blocks

TIME

(C) Client thread
(S) Server thread

**Multicore**
2020/09/18 R09

QNX

# Cache implications of our bad scenario:

**1**

**2**

**C** *MsgSend()*s to Server, blocks

**S** Receives message, runs

**S** Replies to Client

**C** *MsgSend()* returns, runs

**S** Calls *MsgReceive()*, blocks

**C** Client accesses reply buffer

TIME

At this point a cache coherency transaction is done by the cores.  This means the Client's cache data is copied from core 1 to core 2.

**C** Client thread

**S** Server thread

**Multicore**
2020/09/18 R09
**8**

QNX

# Actual behavior:

**1**

**C** *MsgSend()*s to Server, blocks

**S** Receives message, runs

**S** Replies to Client

**S** Calls *MsgReceive()*, blocks

**C** *MsgSend()* returns, runs

Once a server has replied to a client it typically goes back to a *MsgReceive()* to wait for another message.

So even though the Client becomes ready here...

... we don't schedule the client until…

The server blocks. Then we run the Client.

The end result is that:
- they don't thrash from one core to another and
- cache coherency transactions aren't needed.

TIME

**C** Client thread
**S** Server thread

**⊞ QNX**™

# How multicore scheduling is done in QNX:

– the scheduling state is stored in memory

- the data is shared by all cores

– when a thread changes state from RUNNING on a core

- the kernel runs another thread on that core

– when a thread changes state to a runnable state on a particular core:

- the kernel, running on that core, calculates a new scheduling state
- if it needs to run a different thread than is currently running on the current core, it does so
- if a different thread than is currently running on another core needs to be scheduled on that other core, it issues an InterProcessor Interrupt (IPI) to that core and exits the kernel

– the other core enters the kernel due to the IPI

- looks at the scheduling data, and schedules a (new) thread
- updates the scheduling state data
- perhaps issues another IPI

# How we keep the client on the same core:

**1**

(C) *MsgSend()*s to Server, blocks

(S) Receives message, runs

The kernel updates the scheduling data to mark the client READY

(S) Replies to Client

... but, we don't issue the IPI to schedule the client on another core…

(S) Calls *MsgReceive()*, blocks

So, usually, when the server blocks…

(C) *MsgSend()* returns, runs

…we run the Client on the same core.

TIME

(C) Client thread

(S) Server thread

**QNX**

**`SchedCtl(SCHED_CONFIGURE, …)`** takes a structure with two parameters:

- **`low_latency_priority`**
  - if a thread becomes runnable and is
    - higher priority than the current thread on the current core
    - higher priority than this setting
  - it preempts on the current core immediately
- **`migrate_priority`**
  - if a thread
    - gets preempted on this core
    - is higher than this priority
    - there is a lower-priority thread running on another core
  - then migrate it immediately to another core
- both default to **`INT_MAX`**

# A thread scheduling optimization:

– normally the kernel puts a newly scheduled thread on the core with the lowest priority thread

– if there are several such cores available, then the kernel tries to run that thread on the same core that it last ran on
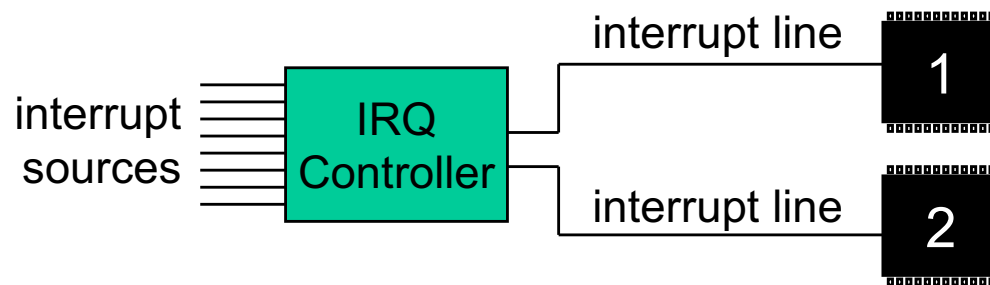


Thread G
RUNNING
priority 2

Thread B
RUNNING
priority 20

Thread E
RUNNING
priority 10

Thread H
RUNNING
priority 2

Thread D
un-blocks
priority 17

In this example, let's pretend that when Thread D last ran it was on core 4. If Thread D becomes ready to run now with the cores in use as in this diagram, the kernel will run it on core 4.

# Where do interrupt handlers run?

– controlled by a mapping of interrupt sources to cores
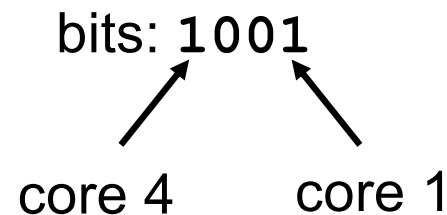  - for some controllers the `startup` code can set it up, for others it cannot be changed



– interrupt handlers run on the core that gets the interrupt
– if interrupt handling causes an event to be delivered then the recipient thread will typically want to use the same core
  - can be controlled using processor affinity...

:: QNX

# Processor affinity:

- you can arrange it so that the kernel will run a thread only on a certain core or set of cores
- to do this, use the *ThreadCtl()* function:

```
ThreadCtl(_NTO_TCTL_RUNMASK, (void *) runmask);
```

where each bit in **runmask** represents a core

bits: **1001**

core 4        core 1

Example:

```
int runmask;

runmask = 0x9; /* 0x9 = 1001, so run core 1 or 4*/
ThreadCtl(_NTO_TCTL_RUNMASK, (void *)runmask);
```

QNX

# Why use processor affinity?

- you can keep non-realtime threads on a specific core
  - in general this is not necessary since the highest priority ready thread will always pre-empt a lower priority thread

- to make better use of the CPU cache?
  - most useful after interrupt handling
  - chances are that some of the thread's code/data is in a core's cache
  - soft affinity is usually enough

# Where does the kernel run?

– to answer this you need to know when the kernel runs, it runs when:

- a hardware interrupt is generated
  - we already saw how interrupts are assigned to cores
  - timeslicing of threads is done by the kernel as a part of handling the timer interrupt
- a fault or exception occurs (e.g. a bus error, segment violation)
  - kernel handles it on the core that the fault occurred on
- a kernel call is made
  - kernel runs on the same core as the caller
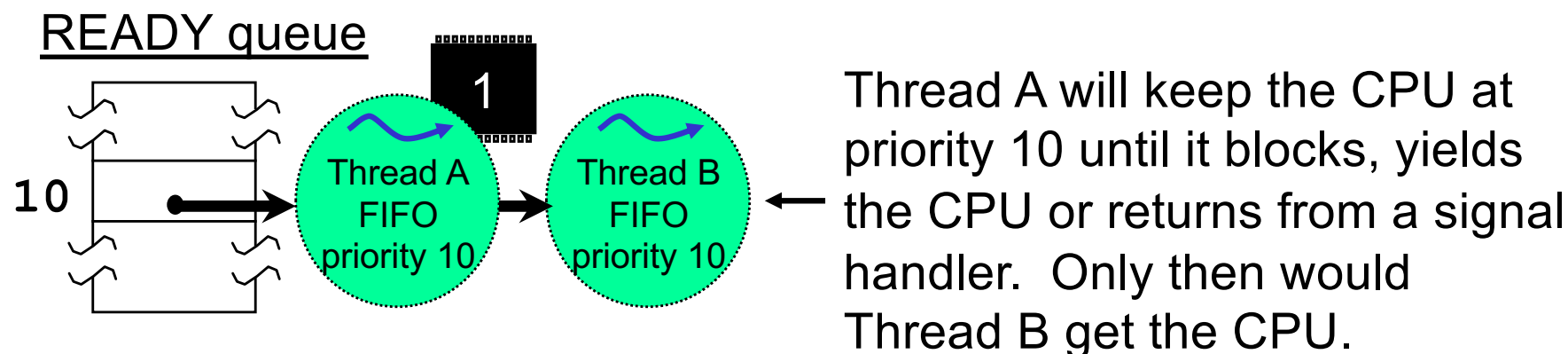
**QNX**

# Topics:
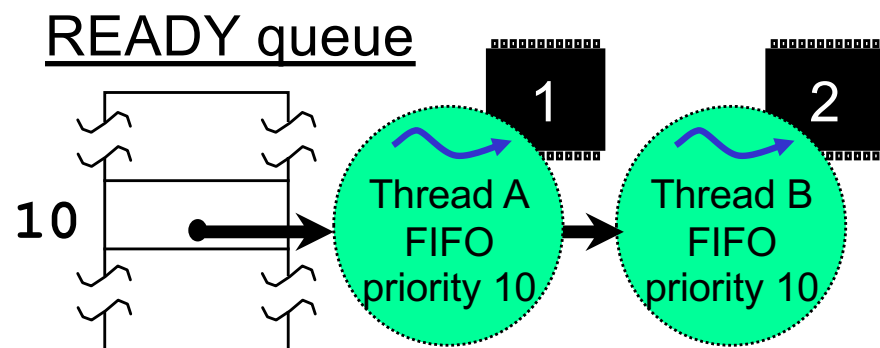
**Overview**

**Scheduling**

→ **Synchronization**

**Conclusion**

# Synchronizing between threads:

– on a uniprocessor system synchronization can sometimes be done using FIFO scheduling

READY queue

**1**

Thread A FIFO priority 10

Thread B FIFO priority 10

Thread A will keep the CPU at priority 10 until it blocks, yields the CPU or returns from a signal handler.  Only then would Thread B get the CPU.

– on an multicore system they may run at the same time

READY queue

**1**

**2**

Thread A FIFO priority 10

Thread B FIFO priority 10

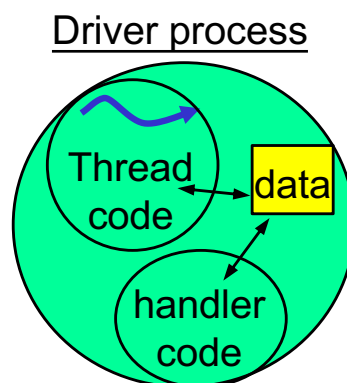So this type of synchronization cannot be used on multicore systems.

# Standard thread synchronization functions:

– barriers, mutexes, condvars, semaphores and all of their derivatives can still be used for synchronization on multicore

- we guarantee that all the thread sync functions will work as expected, providing proper behaviour in multicore systems

# Synchronizing between a thread and an interrupt handler:

– sometimes a thread and an interrupt handler within a driver process will have to share data
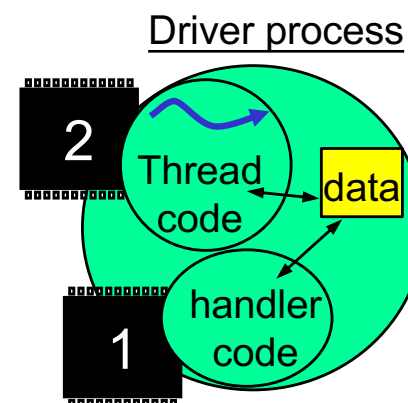
Driver process



– using an event to schedule a thread completely avoids this problem

• *InterruptAttachEvent()* rather than *InterruptAttach()*

QNX

# The problem:

- on a uniprocessor system:

  - an interrupt handler can pre-empt a thread at anytime

    - this is because interrupt handlers have higher priority than any thread

  - a thread cannot pre-empt an interrupt handler

- on an multicore system:

  - the interrupt handler and thread can both be running at the same time!

Driver process

2

Thread code

data

handler code

1

:: **QNX**™

# Solution 1 of 2 - use the atomic*() functions:

| | |
|---|---|
| `atomic_add` | does += some value |
| `atomic_add_value` | *atomic_add()*, and returns original value |
| `atomic_clr` | does &= ~some value |
| `atomic_clr_value` | *atomic_clr()*, and returns original value |
| `atomic_set` | does \|= some value |
| `atomic_set_value` | *atomic_set()*, and returns original value |
| `atomic_sub` | does -= some value |
| `atomic_sub_value` | *atomic_sub()*, and returns original value |
| `atomic_toggle` | does ^= some value |
| `atomic_toggle_value` | *atomic_toggle()*, and returns original value |

These functions:

- are guaranteed to complete correctly despite pre-emption or interruption, and between cores
- can be used between threads
- can be used between threads and an interrupt handler

# Solution 2 of 2 - use an exclusion lock:

- declare a spinlock variable such that both your thread and interrupt handler can access it:

  ```
  intrspin_t spinlock;
  ```

  - it must initially contain zeros
  - if it is not initialized that way, use *memset()* to zero it

- then, in both your thread and your interrupt handler do the following to access the data:

  ```
  InterruptLock(&spinlock);
  /* access your data */
  InterruptUnlock(&spinlock);
  ```

☞ you typically keep the time between the above lines as short as possible

☞ you must have called *ThreadCtl(_NTO_TCTL_IO, 0)* for these functions to work

# InterruptDisable()/InterruptEnable():

– always available

– only disables on the calling core

– will not provide the expected protection with multicore

Always use *InterruptLock()/InterruptUnlock()*

– even on a uniprocessor system

- driver, when moved to multicore, will work
- overhead is small enough to not matter

# Topics:

**Overview**

**Scheduling**

**Synchronization**

→ **Conclusion**

## Conclusion

You learned:

– how the multicore scheduler works

  • and how to tune its behavior

– how and why to use processor affinity

– how to synchronize between threads and interrupt handlers on multicore