

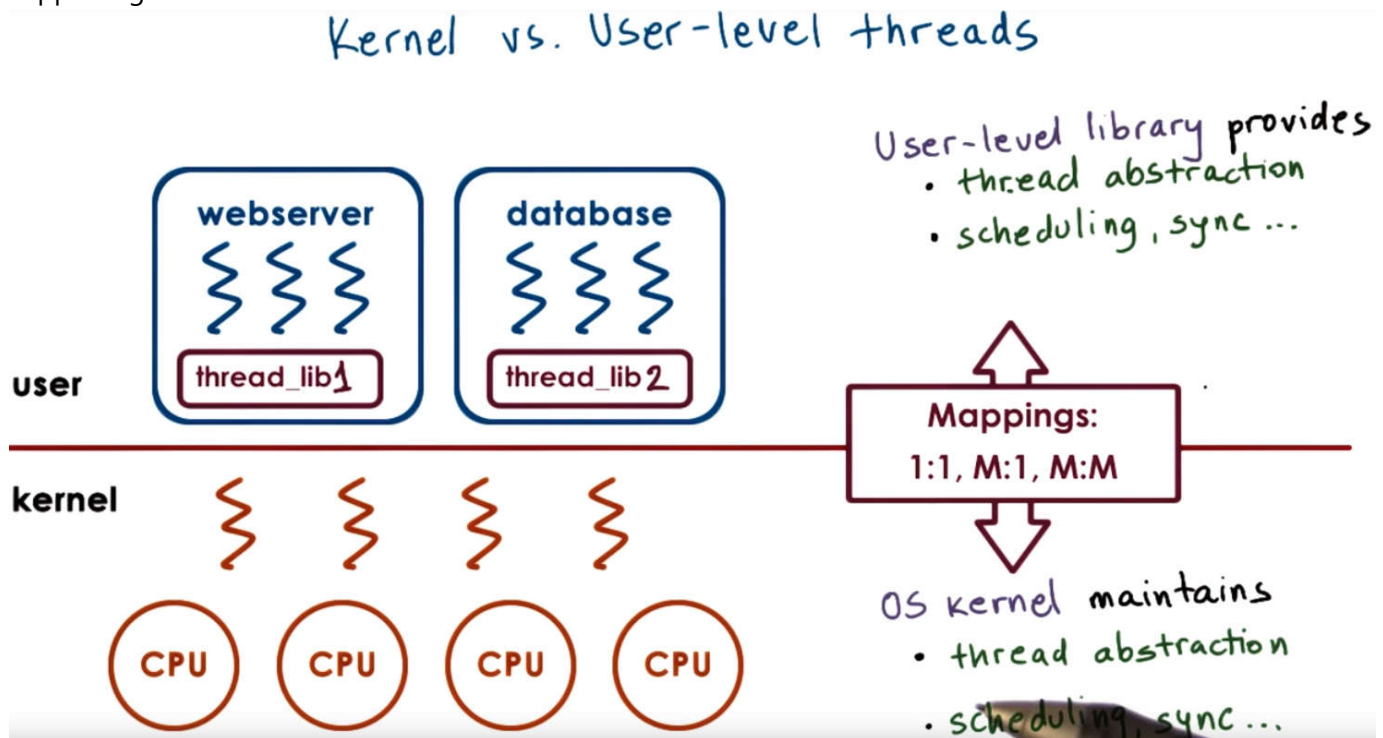
# P2L4 Thread Design Considerations

## Goal

- Data structures and mechanisms of kernel vs. user-level threads
- Two notification mechanisms supported by OSs:
  - Threads and interrupts
  - Threads and signal handling
- How threading systems evolve over time

## 1. Kernel vs. User-level threads

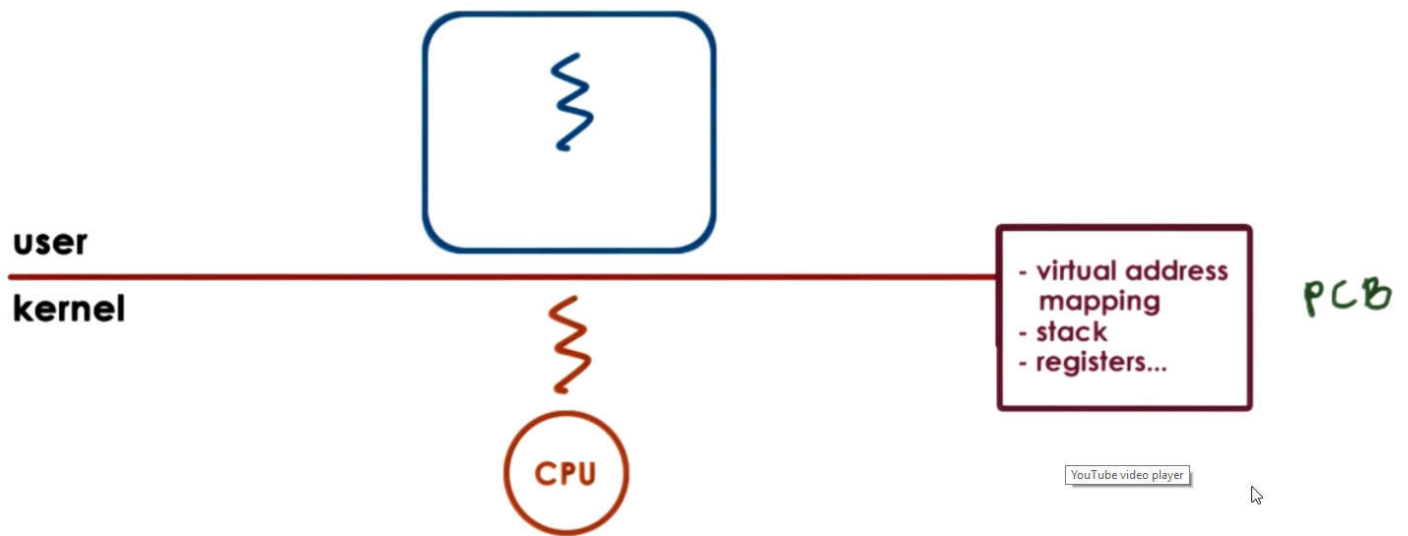
- supporting thread at the OS level means that the OS kernel itself is multithreaded



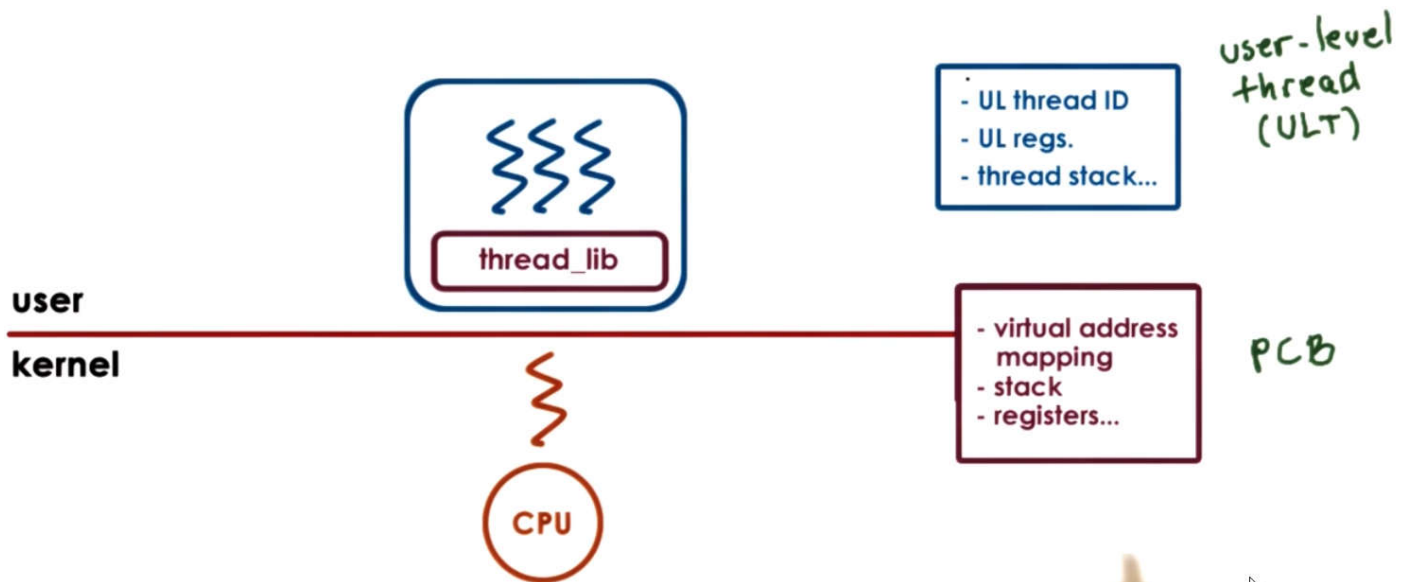
- user level threads
  - thread libs support data structures that's needed to implement the thread abstraction
  - provide all scheduling synchronization and other mechanisms
  - different threads can use entirely different thread libs.
- kernel level threads
  - OS kernel maintains thread abstraction, scheduling sync.
  - Support mapping between user and kernel threads.

### 1.1 Thread-related Data Structures

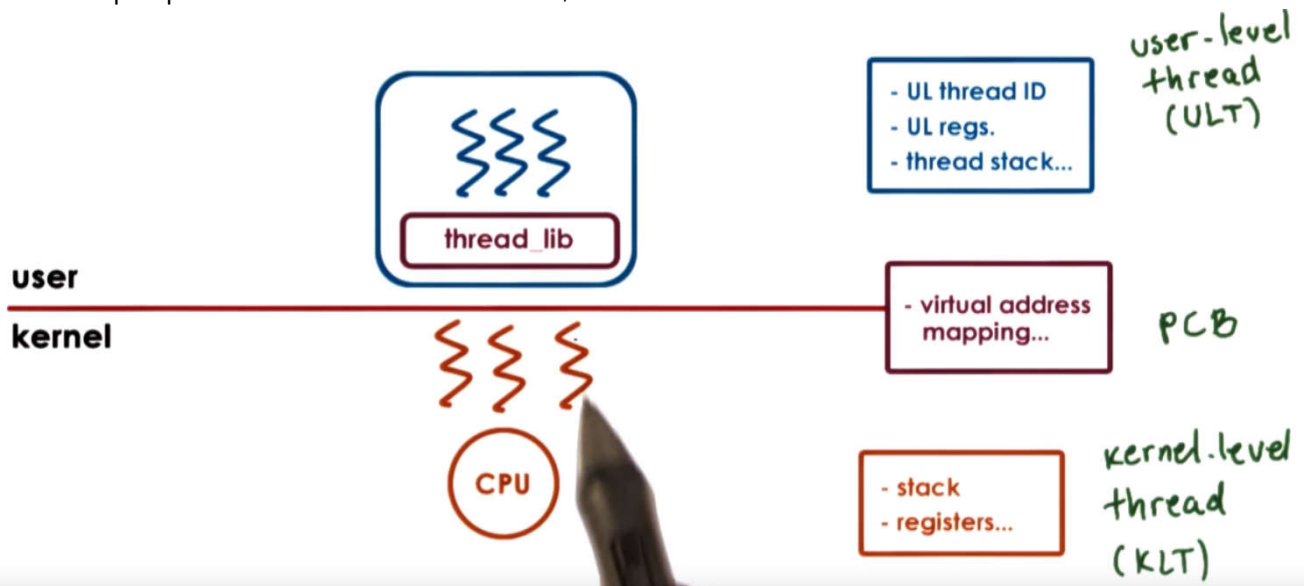
- Single thread on single process



- Multi user-threads on single kernel-thread process

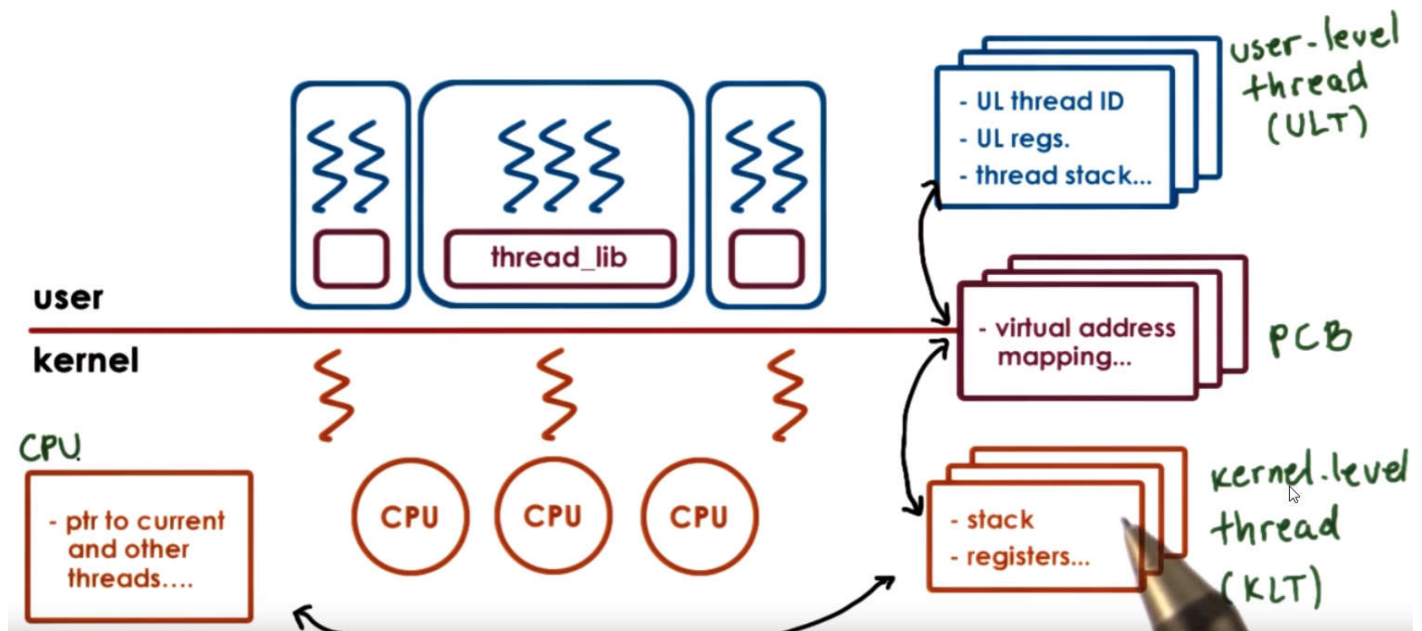


- Multiple kernel-level threads associated with the process with multiple user-thread
  - From the perspective of the user-level threads, the kernel level threads look like virtual CPUs.



## Relationships among ULT, PCB and KLT

- Both the ULT and KLT has to know what is the address space within which that thread executes
- If there are multiple CPUs we have to maintain a data structure to represent the CPU
- CPU data structure has a relationship with KLTs

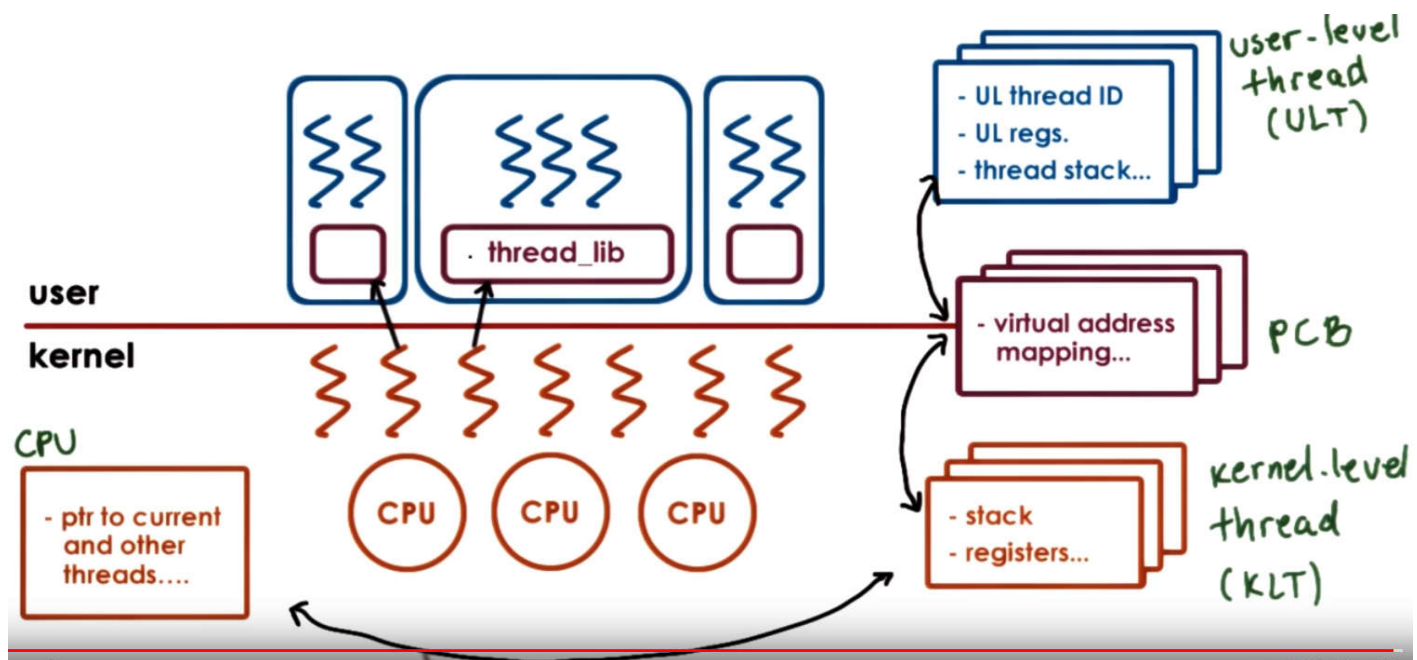


When the kernel is multi-threaded we can have multiple kernel-level threads supporting a single user-level process.

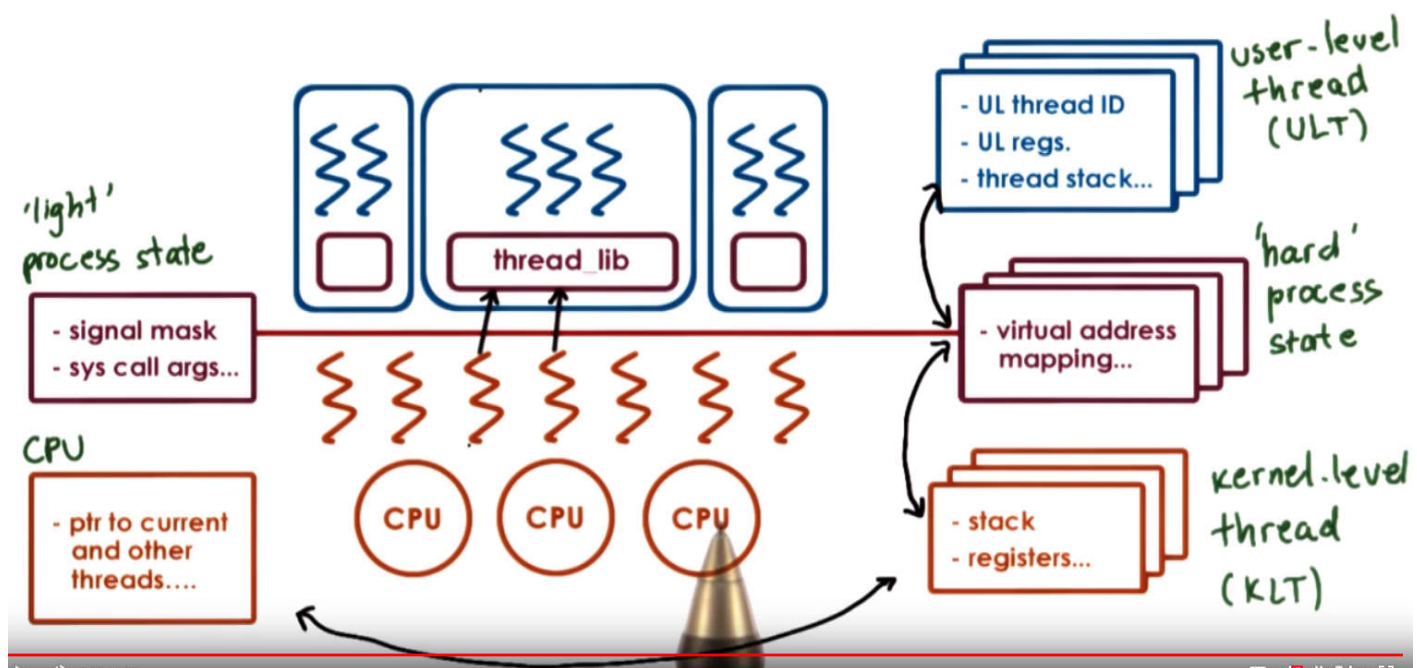
When the kernel needs to schedule/ context switch among kernel-level threads that belong to different processes, it can quickly determine that they point to a different process control block, hence different virtual address mappings.

So the kernel needs to completely invalidate the existing address mappings and restore new ones.

The kernel will save the entire PCB of the thread to be switched off and restore the PCB of the thread to be loaded to run.



## 1.2 Hard and Light Process State



When there are multiple kernel-level threads supporting one process, hence they belong to the same virtual address mapping scheme and the kernel is doing a context switch between them, there are portions of the PCB that they share the same data and part of the PCB that are specific to each kernel-level threads (e.g. signal masks, sys call args).

- **hard process state:** Information in PCB that the kernel-level threads in the process share— virtual address mapping etc.
- **light process state:** Information in PCB that are only relevant to a particular kernel-level thread in the process, and the user-level threads that are mapped to the kernel-level thread.

**Rationale for splitting single PCB into multiple smaller data structures**

## Rationale for Multiple Datastructures

### Single PCB

- large continuous data structure
- private for each entity
- saved and restored on each context switch
- update for any changes



### multiple data structures

- smaller data structures
- easier to share
- on context switch only save and restore what needs to change
- user-level library need only update portion of the state

## Thread Structure Quiz



1. What is the name of the kernel thread structure (name of C struct)?

kthread\_worker

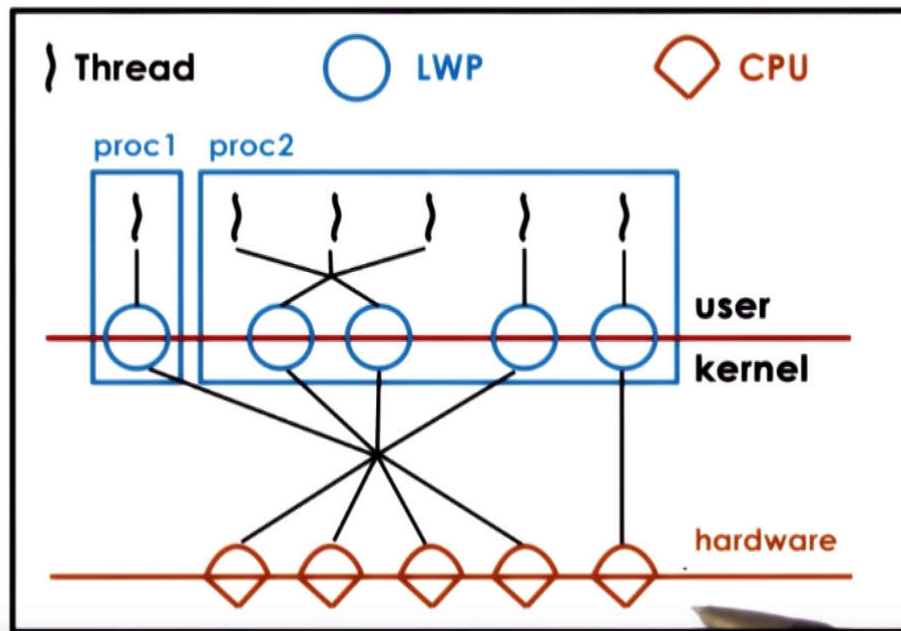
2. What is the name of the data structure - contained in the above structure - that describes the process the kernel thread is running (name of C struct)?

task\_struct

## 1.3 Sun OS 5.0 Threading Model

- Multi-kernel threads
- Both single and multi user-level threads
- 1 to 1 and many to many user-kernel level thread mapping
- Each kernel-level threads have a light weight process data structure representing the virtual CPUs onto which it's going to be scheduling the user-level threads
- kernel level scheduler managing the kernel level threads



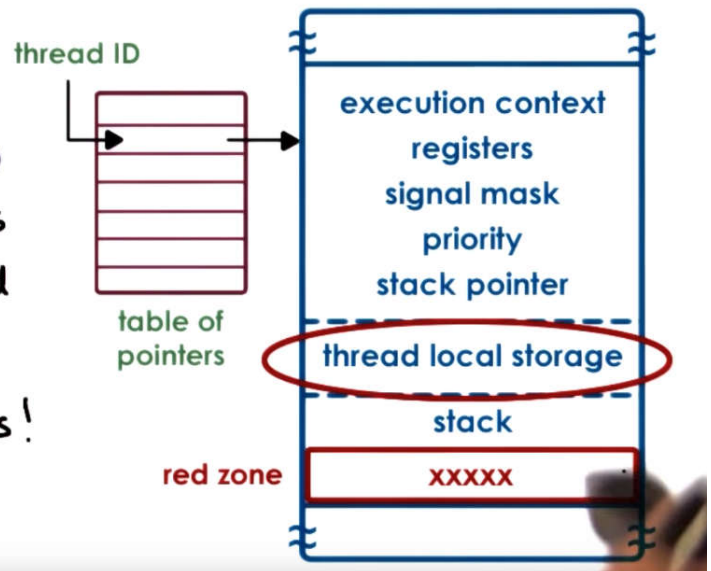


## 1.4 User-level thread data structures

### User-level Thread Data Structures

"Implementing Lightweight Threads", by Stein & Shah

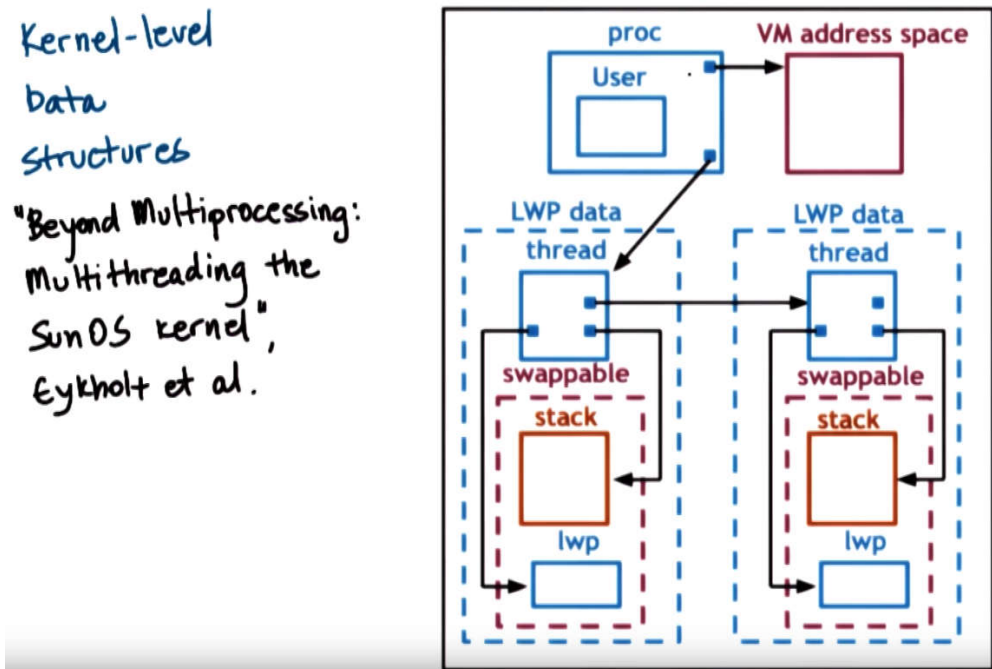
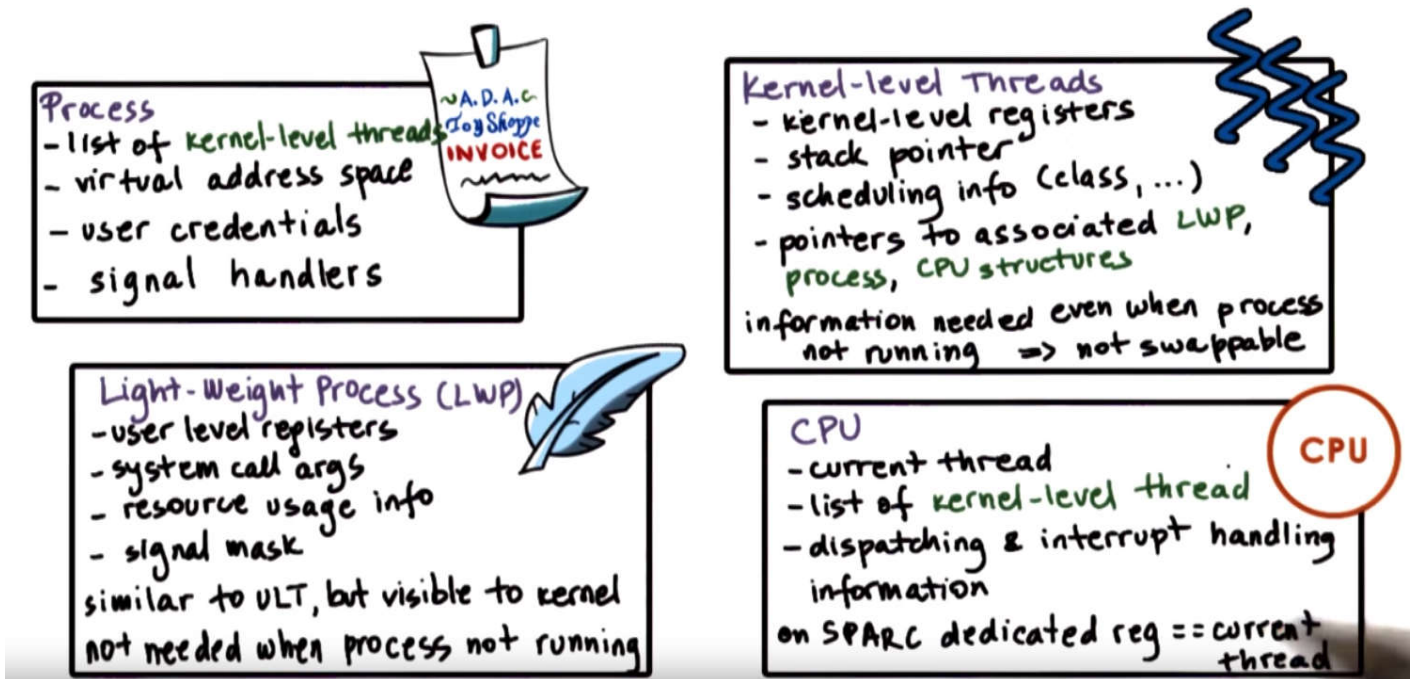
- not POSIX threads, but similar
- thread creation  $\Rightarrow$  thread ID (tid)
  - tid  $\Rightarrow$  index into table of pointers
  - table pointers point to per thread data structure
- stack growth can be dangerous!
  - solution  $\Rightarrow$  red zone



- Having thread ID point to a table entry, we can store some info about the thread in the table entry. This could help us avoid dereferencing a thread id pointer just to find it points to corrupted memory.
- Thread local storage: include the variables defined in thread functions that are known in the compile time so the compiler can allocate private storage on a per-thread basis for each of them.
- Avoid stack overflow, wracking other data structures by separating them with the non-dereferencable red zone.

## 1.5 Kernel-level data structures

- At OS level, the kernel tracks resource uses on a per kernel thread basis, maintained in the lightweight process that corresponds to that kernel level thread.

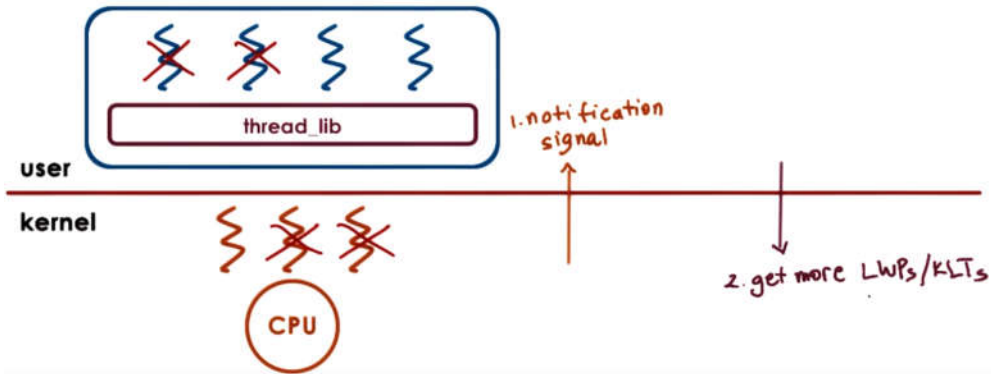


## 2. Basic Thread Management Interactions

- Example Case: The process originally has two KLTs, but both of them are blocked due to an I/O operation. The two other ULTs are unable, so the kernel can send a signal to ULT and then give it an extra KLT to run the runnable threads.

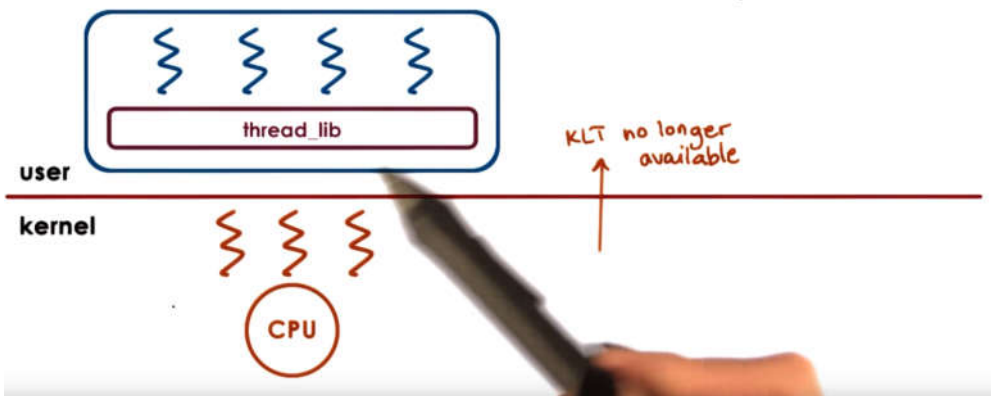
## Basic Thread Management Interactions

User-level library does not know what is happening in the kernel



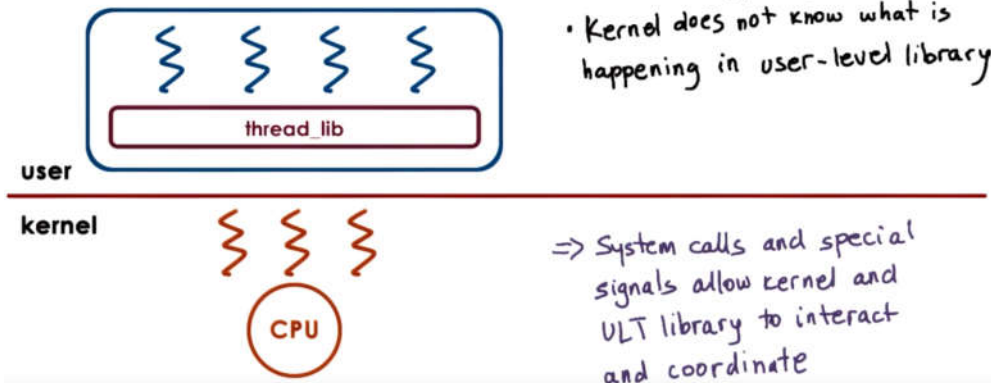
- Then when the blocking is done, the kernel will tell the ULTs that the extra KLT is no longer available.

User-level library does not know what is happening in the kernel



-Summary:

- User-level library does not know what is happening in the kernel
- Kernel does not know what is happening in user-level library



## 2.1 Thread Management Visibility

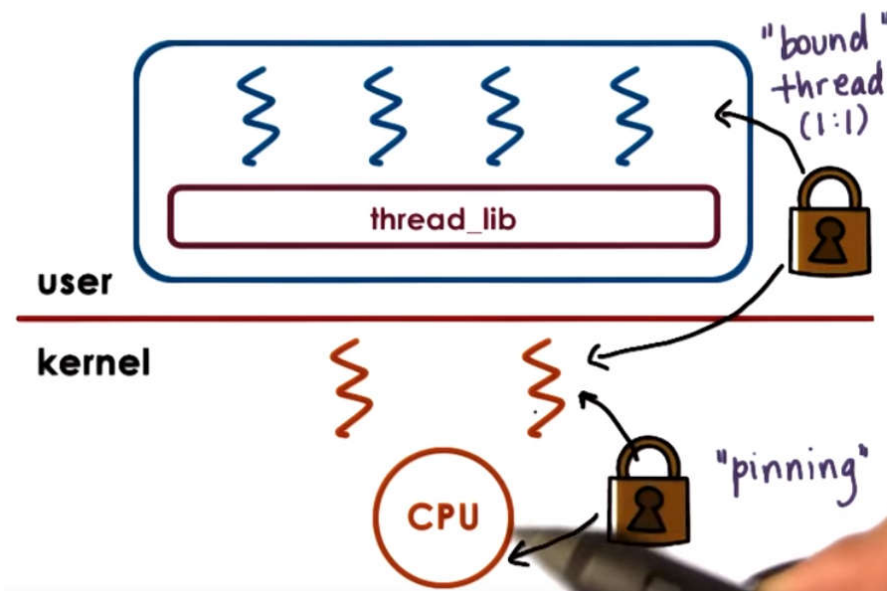
- Bound ULT to KLT



UL library sees  
- ULTs  
- available KLTs

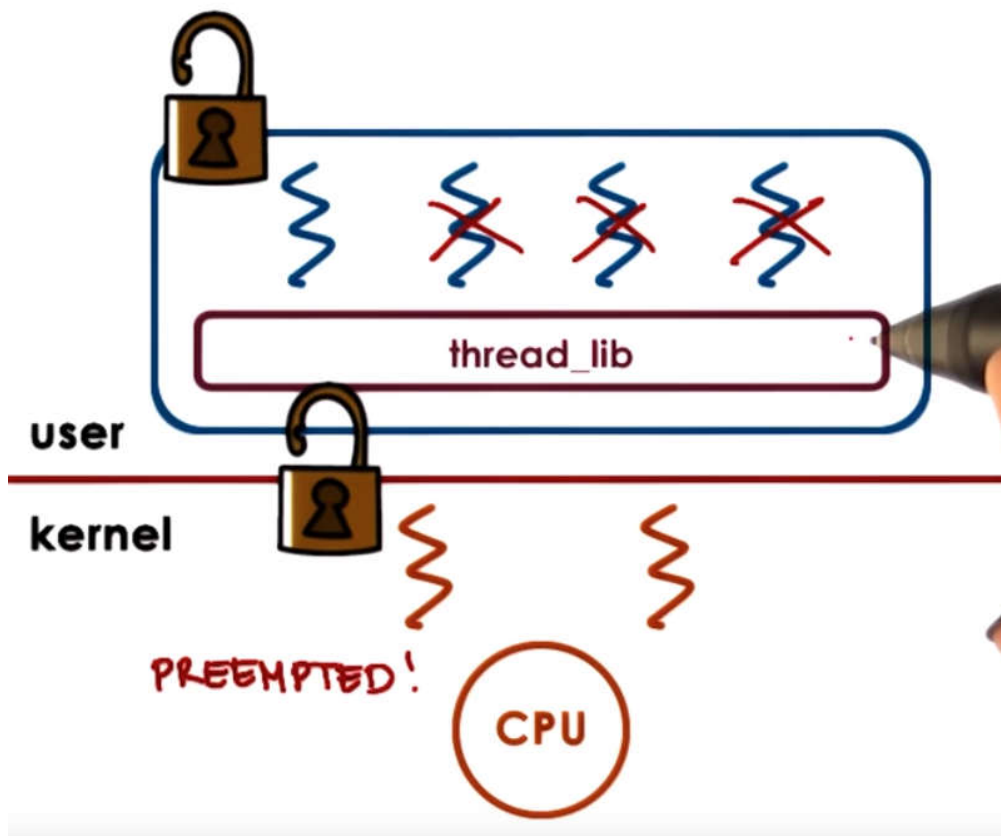
Kernel sees...

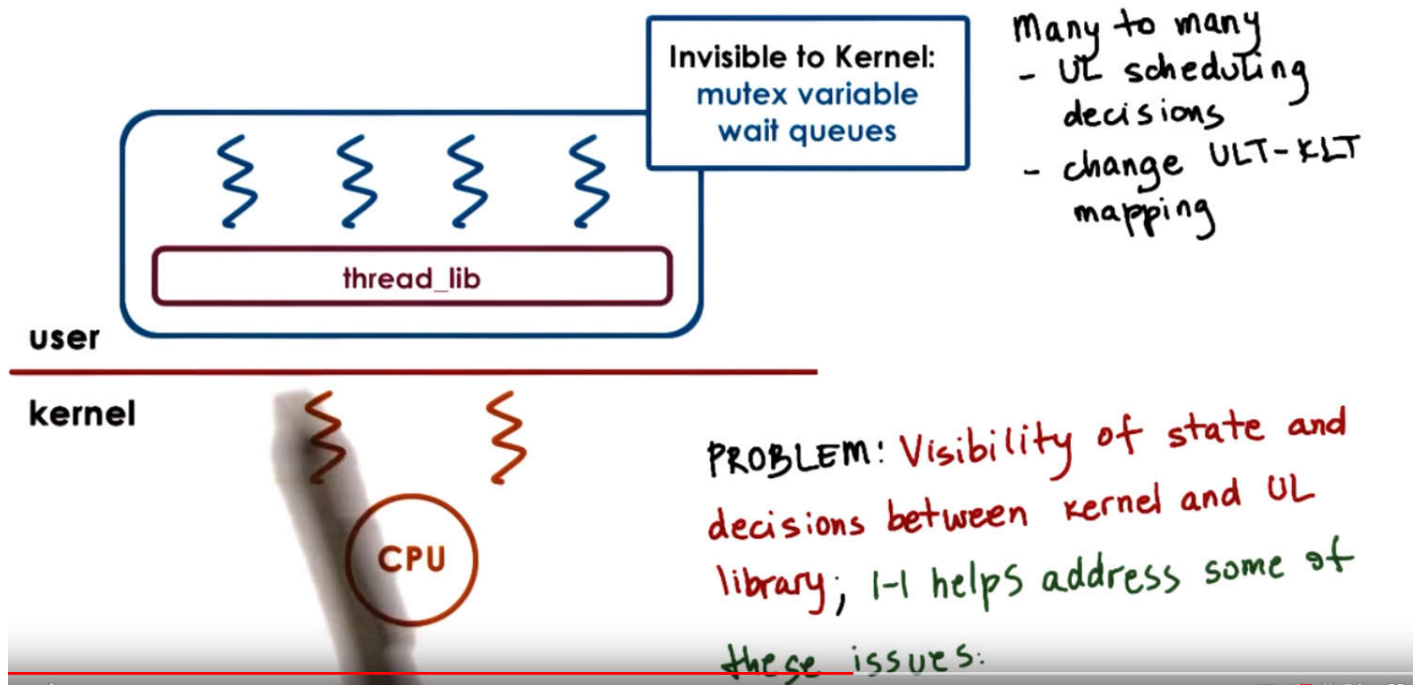
- KLTs  
- CPUs  
- KL scheduler



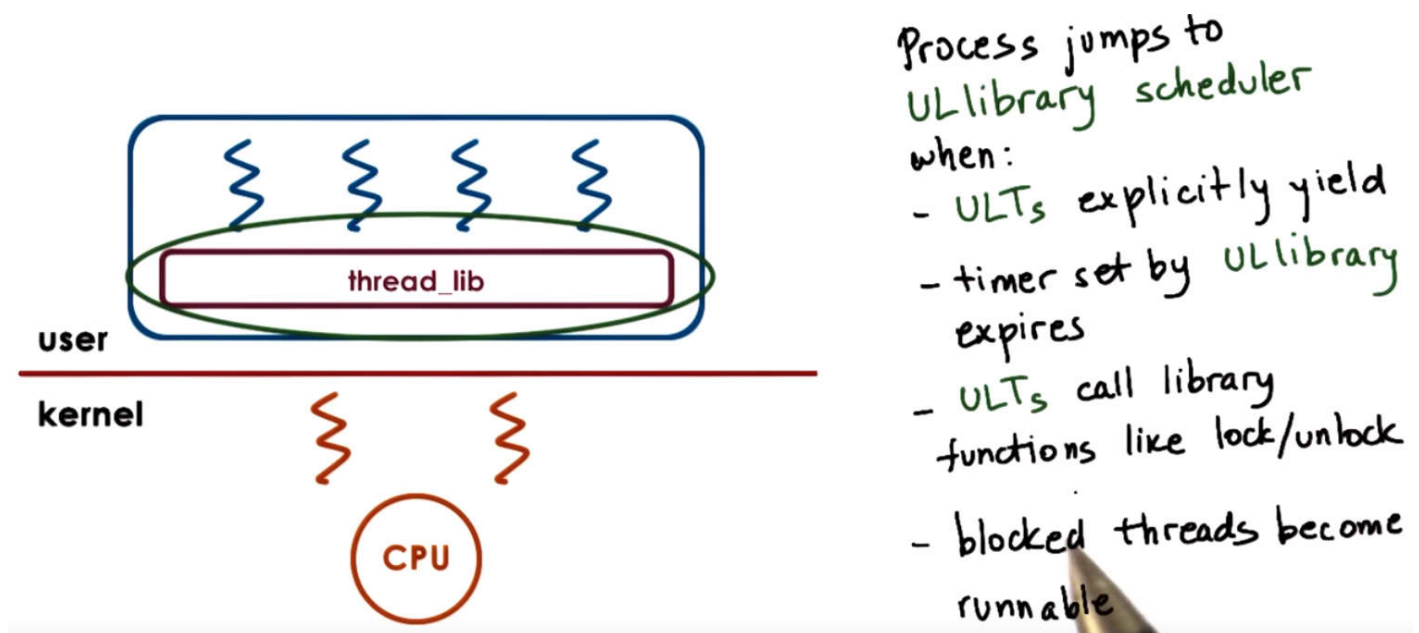
- Problem Case:

When the kernel cannot see ULT scheduling decisions, it might switch a thread that has the mutex lock off the CPU and then iterate through other threads just to find that they are blocked at the mutex lock acquisition.

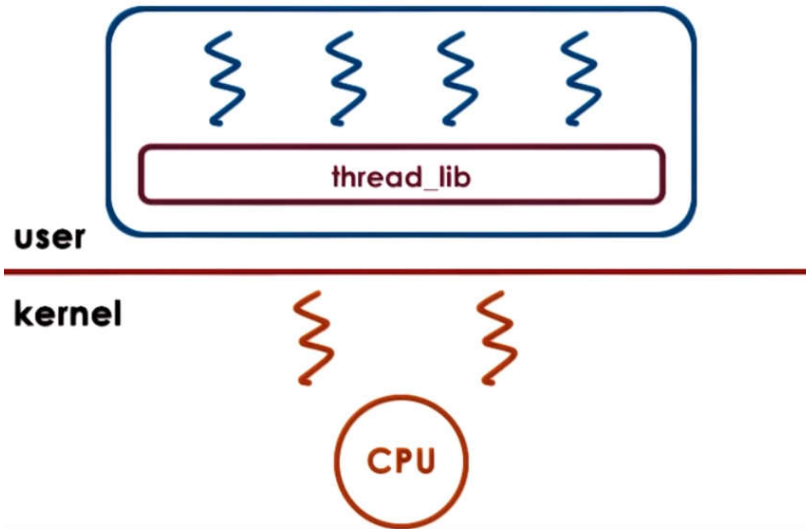




## How does the UL Lib Run?



UL library scheduler...  
- runs on ULT operations  
- runs on signals from timer or kernel



## 2.2 Issues on Multiple CPUs

Question: Why we cannot directly modify register of one CPU when executing on another CPU? --P2L4  
Lesson 13

Question: What does preempt mean?

When there are multiple CPUs, we need interaction between KLTs running on different CPUs.

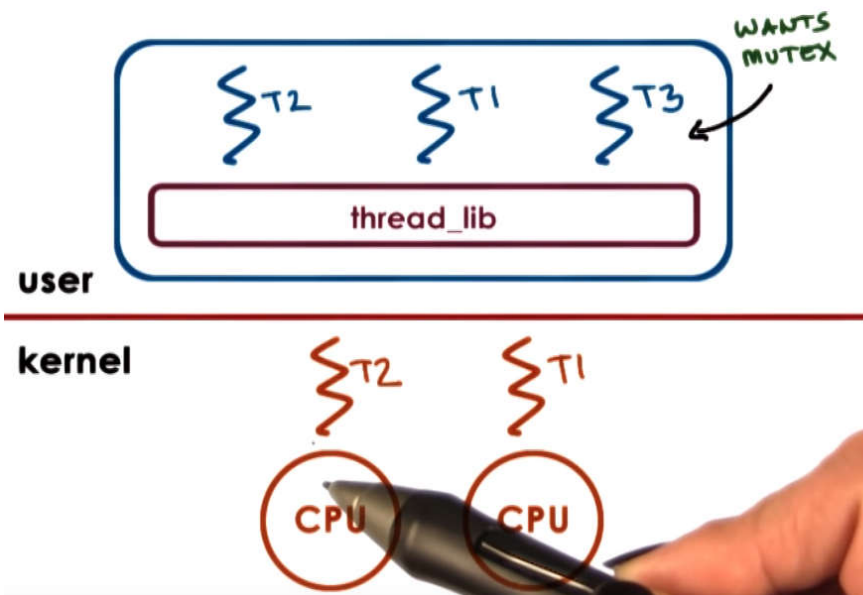
Thread Priority

$T3 > T2 > T1$

T2 unlocks

⇒ make T3 runnable

⇒ preempt T1, but  
on another CPU!



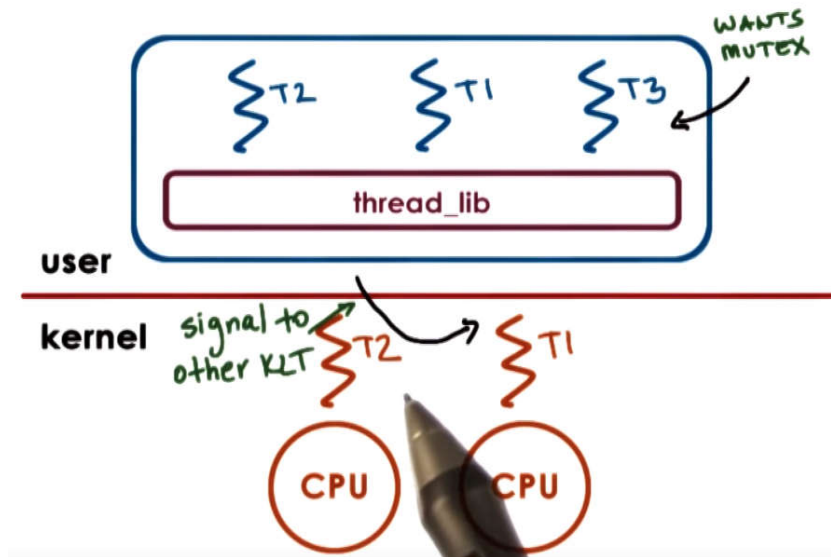
Thread Priority

$T_3 > T_2 > T_1$

T2 unlocks

⇒ make T3 runnable

⇒ send signal to other thread, on other CPU to run library code locally.



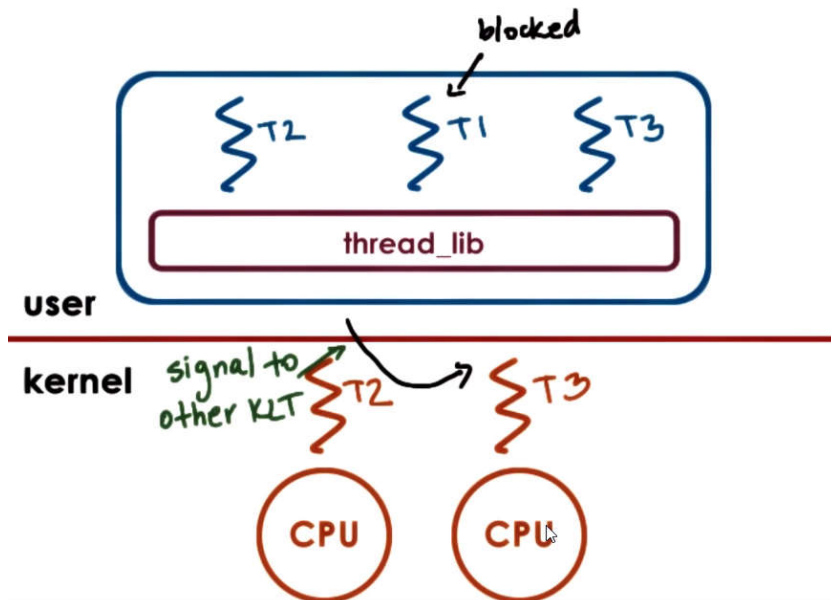
Thread Priority

$T_3 > T_2 > T_1$

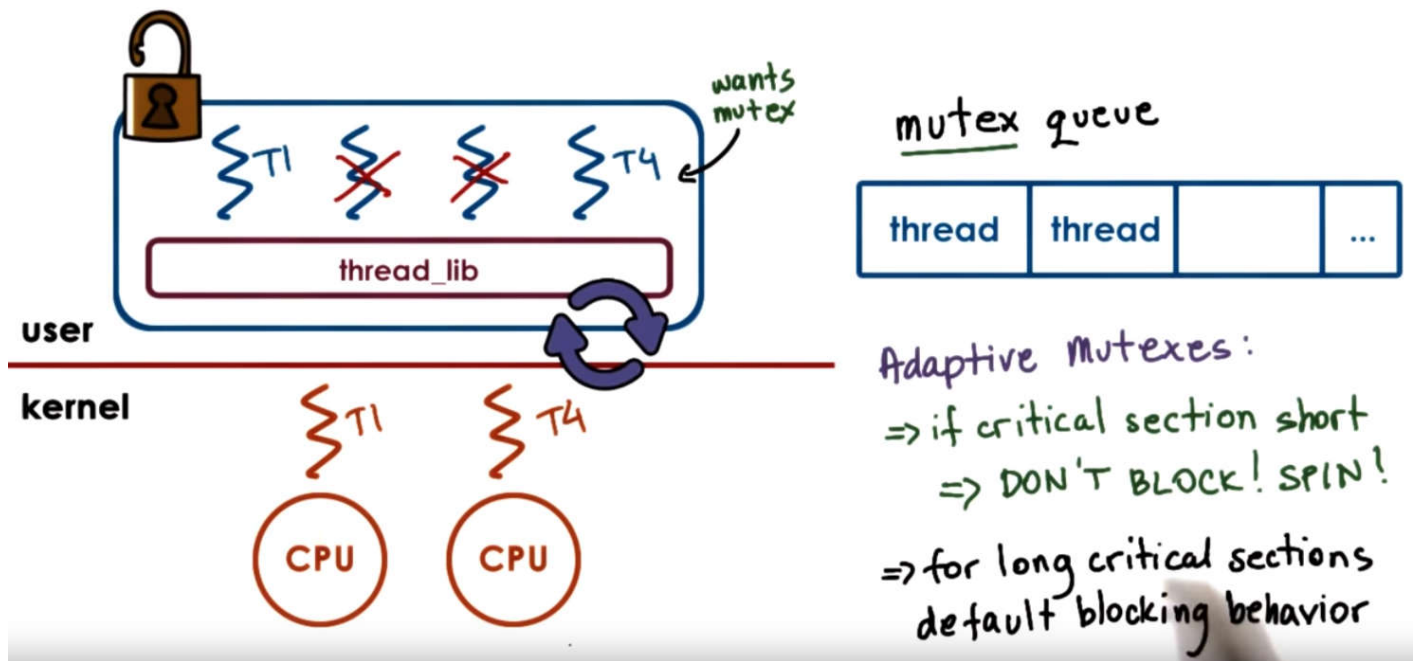
T2 unlocks

⇒ make T3 runnable

⇒ send signal to other thread, on other CPU to run library code locally.



## 2.3 Synchronization Related Issues.



When critical section is short, T1 might release the mutex lock before the context switching is completed on T4. So it's better off leave T4 spinning on CPU 2 and burn a few cycles.

This is only possible on multi-CPU context.

## 2.4 Destroying Threads

### Destroying threads

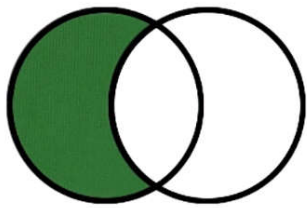
- Instead of destroying... reuse threads
- when a thread exits...
  - put on a "death row"
  - periodically destroyed by reaper thread
  - otherwise thread structures / stacks are reused => performance gains!



## 3. Communication Mechanisms

### 3.1 Interrupts vs. Signals





## Interrupts vs. Signals

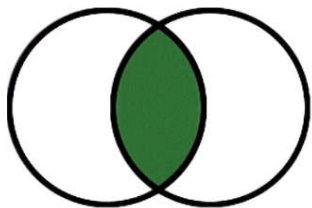
### Interrupts

- events generated externally by components other than the CPU (I/O devices, timers, other CPUs)
- determined based on the physical platform
- appear asynchronously

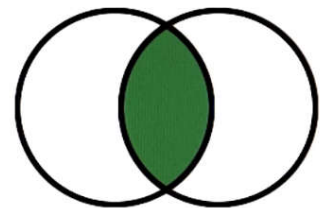
### Signals

- events triggered by the CPU & software running on it
- determined based on the operating system
- appear synchronously or asynchronously

- **synchronously:** means it's triggered by some specific actions that took place on the CPU, a synchronous signal is generated in response to that action. (such as an attempt to access unallocated memory.)



## Interrupts vs. Signals



### Interrupts and Signals

- have a unique ID depending on the hardware or OS
- can be masked and disabled/suspended via corresponding mask
  - per-CPU interrupt mask, per-process signal mask
- if enabled, trigger corresponding handler
  - interrupt handler set for entire system by OS
  - signal handlers set on per process basis, by process

## Visual Metaphor

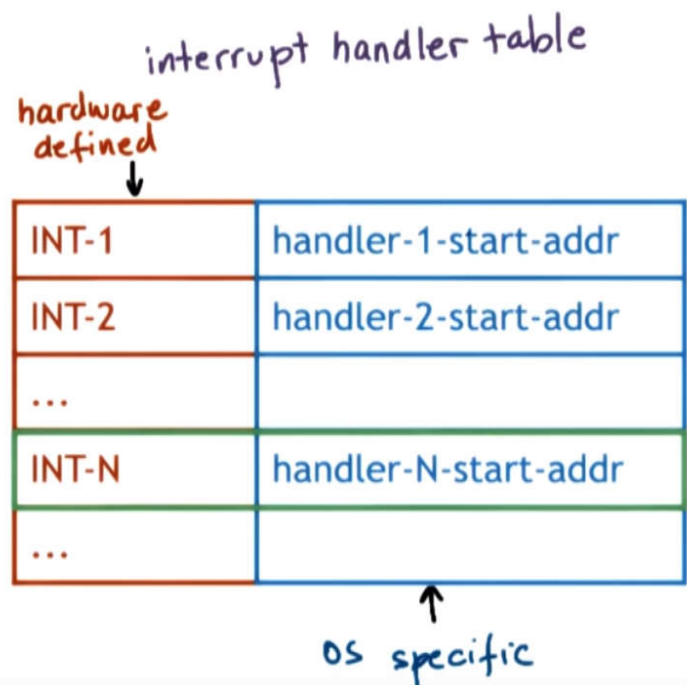
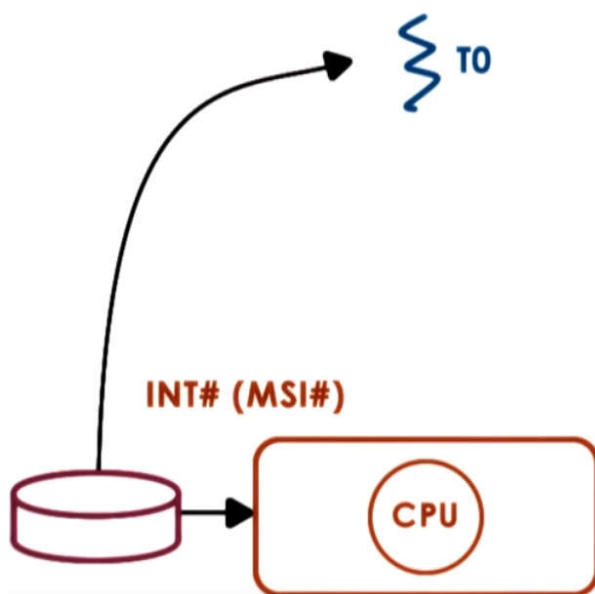
"An interrupt is like ... a snowstorm warning"  
"A signal is like ... a 'battery is low' warning"

- handled in specific ways
  - interrupt and signal handlers
- can be ignored
  - interrupt / signal mask
- expected or unexpected
  - appear sync. or async.

- handled in specific ways
  - safety protocols, hazard plans...
- can be ignored
  - continue working
- expected or unexpected
  - happen regularly or irregularly

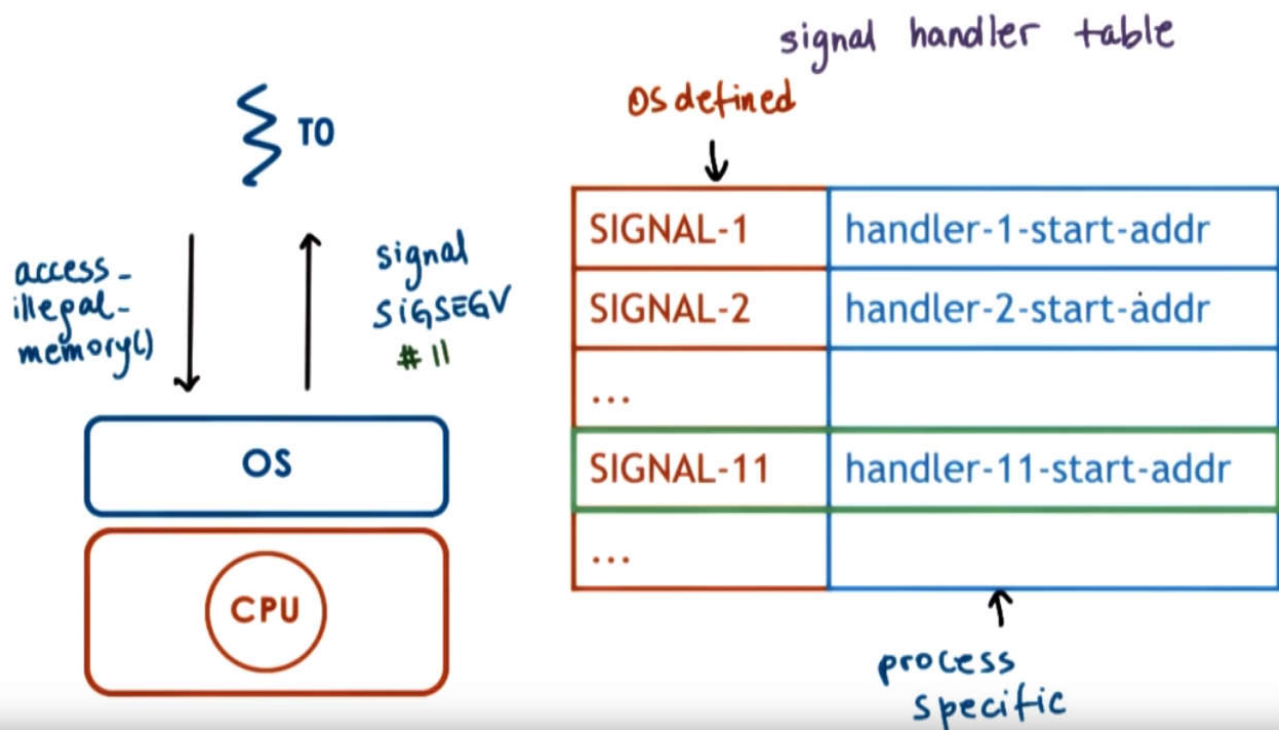
## Interrupt Handling

### Interrupts



## Signal Handling

# Signals



# Signals

## Handlers/Actions

### Default Actions

- Terminate, Ignore, Terminate and Core Dump, Stop or Continue

### Process Installs Handler

- signal(), sigaction()
- for most signals, some cannot be "caught"

## Synchronous

- SIGSEGV (access to protected mem)
- SIGFPE (divide by zero)
- SIGKILL (kill, id)  
can be directed to a specific thread

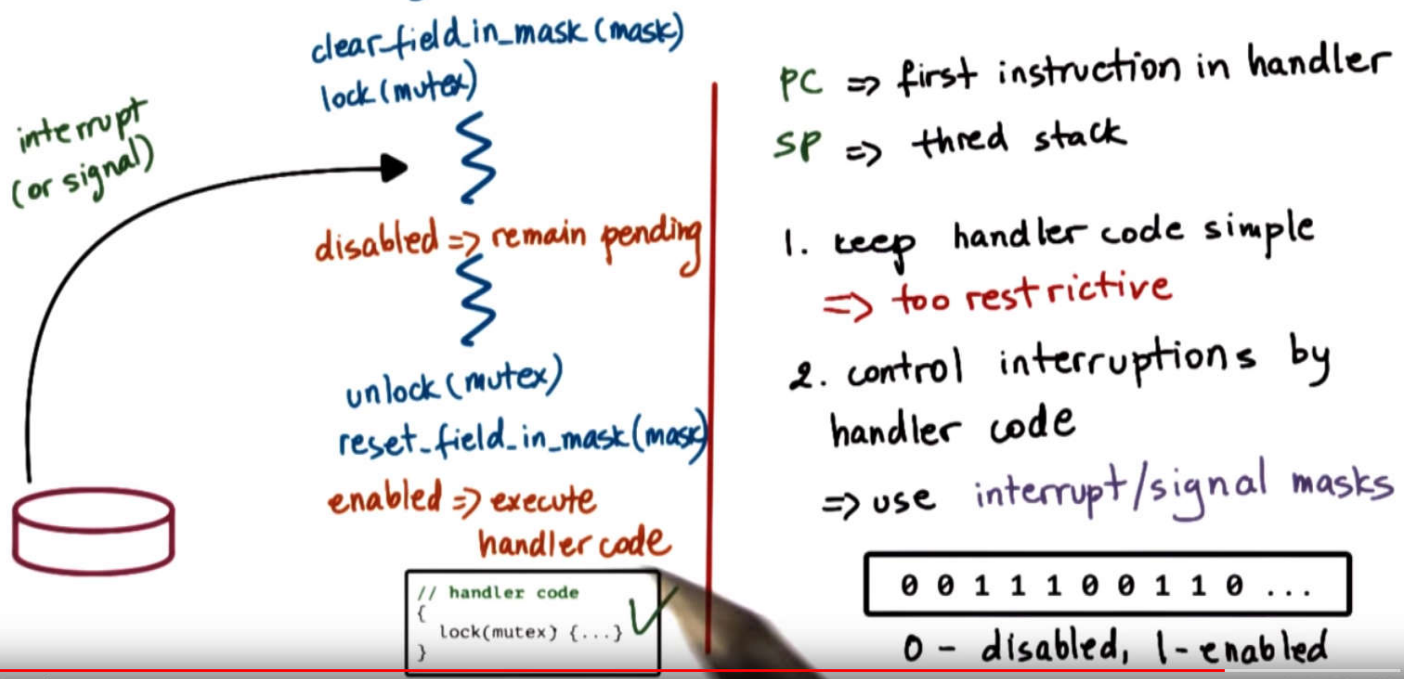
## Asynchronous

- SIGKILL (kill)
- SIGALRM

## 3.2 Disabling Interrupts or Signals

What happens when an interrupt/signal appears:

## Why Disable Interrupts or Signals ?



- When the signal/interrupt handler code needs to acquire a mutex to proceed, but the mutex is currently owned by the thread's code, there is a dead lock situation.
- Solution to the above problem: signal/interrupt masks, enabling us to enable/disable signal/interrupts.
- The mask is a sequence of bits that corresponds to specific signal or interrupts, 0 – disabled, 1 – enabled.
- In the thread code, when the thread needs to acquire the mutex, disable the possible signal/interrupts, after it's done with the critical section, unmask the signals.
- While a signal/interrupt is pending, other signals/interrupts that occurred will be pending too, once the mask is reset, the handler routine will typically be executed only once.

## More on Masks

Interrupt masks are per CPU

if mask **disables** interrupt  $\Rightarrow$  hardware interrupt routing mechanism will not deliver interrupt to CPU

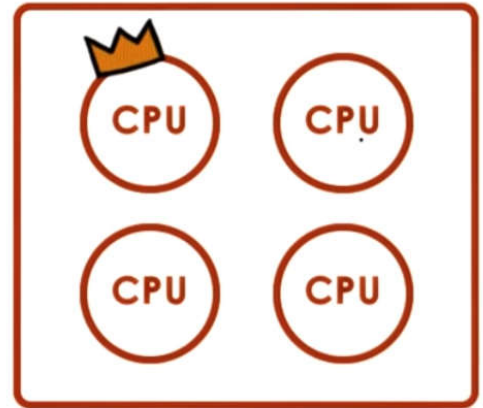
Signal masks are per execution context (ULT on top of KLT)

if mask **disables** signal  $\Rightarrow$  kernel sees mask and will not interrupt corresponding thread

### 3.3 Interrupts on Multicore Systems

#### Interrupts on Multicore Systems

- Interrupts can be directed to any CPU that has them enabled
- may set interrupt on just a single core
  - $\Rightarrow$  avoids overheads & perturbations on all other cores



### 3.4 Types of Signals

#### Types of Signals

##### One-Shot Signals

- "n signals pending == 1 signal pending" : at least once
- must be explicitly re-enabled

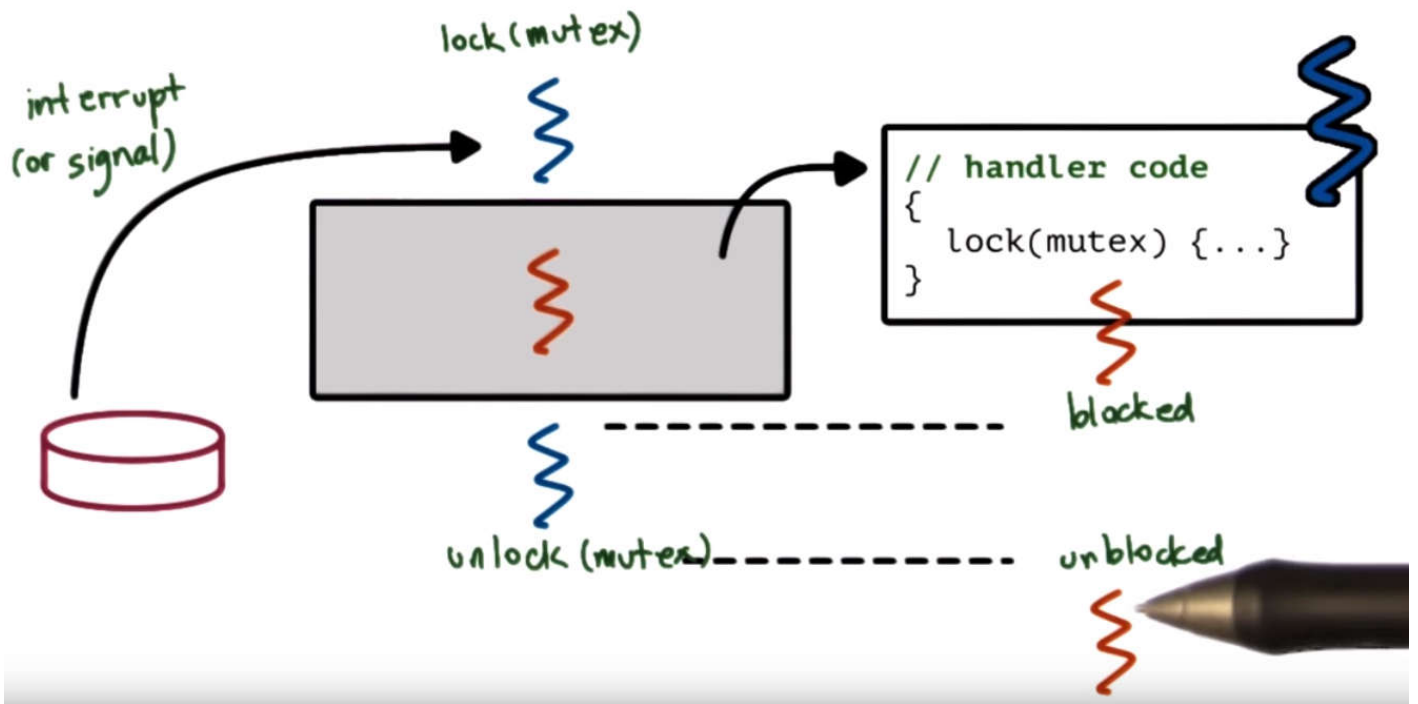
##### Real Time Signals

- "if n signals raised, then handler is called n times"

### 3.5 Handling Interrupts as Threads

- Problem: When an interrupt/signal is sent to a thread that has a mutex lock, also the interrupt handler needs to acquire the lock, there is a dead lock situation.
- One possible solution: Handle the interrupt as a thread, switch the original thread back till the mutex is released then unblock the interrupt handler "thread".





- However, dynamic thread creation is expensive:

... but, dynamic thread creation is expensive!

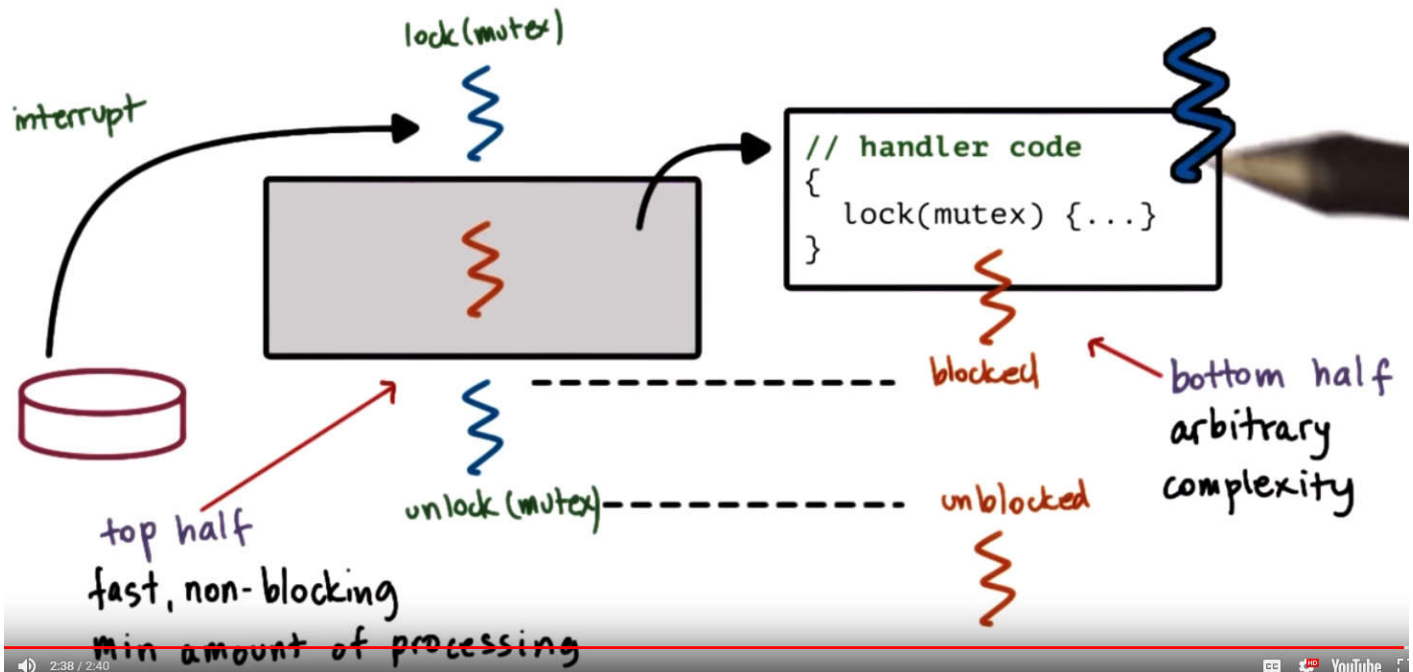
### Dynamic Decision

- if handler doesn't lock => execute on interrupted thread's stack
- if handler can block => turn into real thread

### Optimization

- precreate & preinitialize thread structures for interrupt routines

## Top vs. Bottom Half



Rationale behind treating interrupt as threads:

## Performance: Bottom Line

### Overall Cost

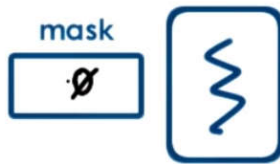
- overhead of 40 SPARC instructions per interrupt
- saving of 12 instructions per mutex
  - no changes in interrupt mask, level...
- fewer interrupts than mutex lock/unlock operations

⇒ It's a WIN!

Optimize for the common case!

## 3.6 Threads and Signal Handling

# Threads and Signal Handling



disable  $\Rightarrow$  clear signal mask  
signal occurs - what to do  
with the signal?

**user**

---

**kernel**



Problem: ULT and KLT has separate masks, when ULT disables a signal how to let the KLT that it maps to know how to deal with the signal.

**Case 1: ULT mask = KLT mask = 1**

Case 1:



**user**

---

**kernel**

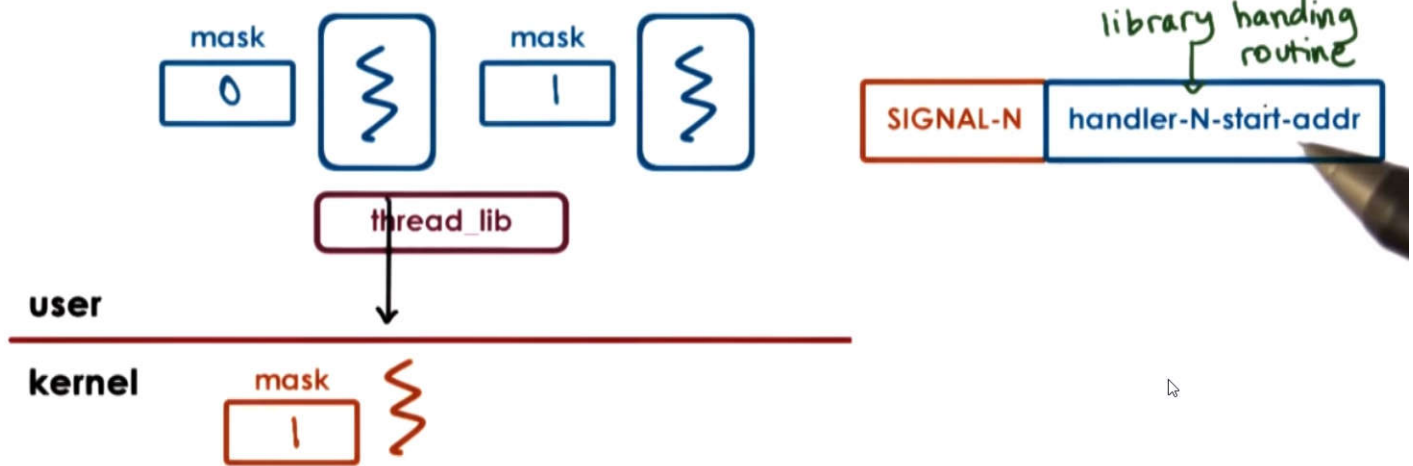


No problem here.

**Case 2: ULT mask = 0, KLT mask = 1, another KLT mask = 1**

Case 2:

ULT mask = 0  
KLT mask = 1  
another ULT mask = 1

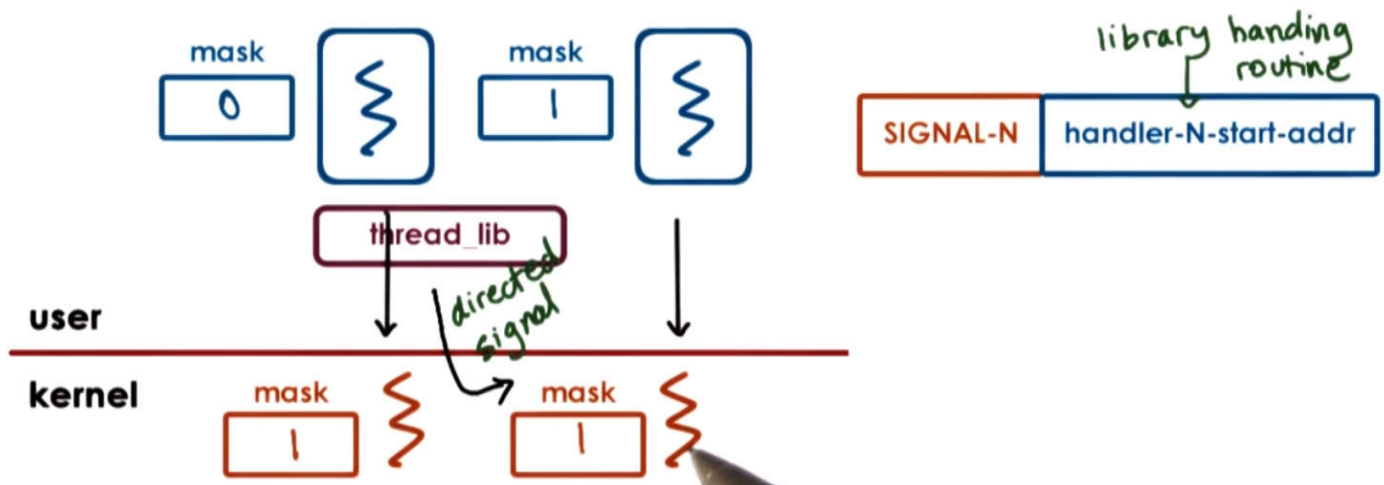


The user-level thread lib can see which ULT has the signal enabled.

Case 3: 2 KLTs mask = 1, 1 ULT mask = 1, 1 ULT mask = 0

Case 3:

ULT mask = 0  
KLT mask = 1  
another ULT mask = 1  
KLT mask = 1

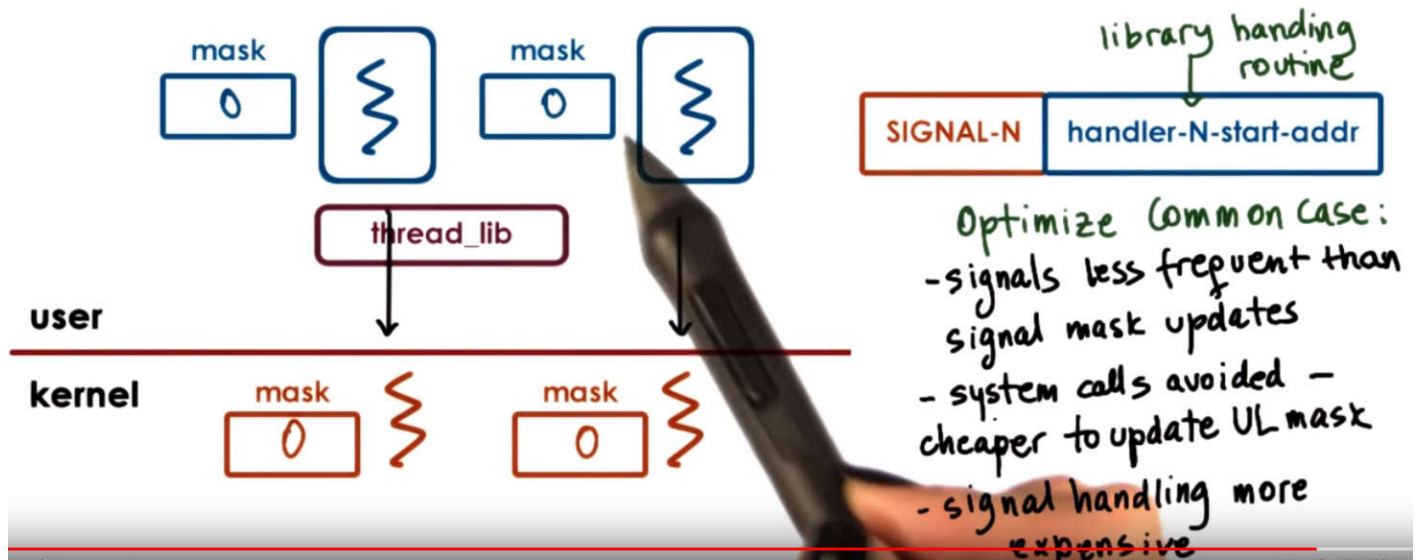


The thread lib will direct the signal to the KLT mapped to the ULT that has mask enabled.

Case 4: KLT mask = 1, all ULT mask = 0

Case 4:

ULT mask = 0  
KLT mask = 1  
all ULT mask = 0



Thread lib will perform a system call to change underlying KLT mask to 0. Then the OS will reissue the signal to other threads in the process, if the thread also disables the signal, then it'll change the other thread's underlying KLT to 0.

If one of the ULT mask is updated to enable a signal, the thread lib will have to perform a system call to update the mask in KLT.

## 4. Threading Supports in Linux—Task

- Task: Main execution abstraction
  - execution context of a kernel level thread.

Task Struct in Linux

- Main execution abstraction => task
  - kernel level thread
- Single-threaded process => 1 task
- Multi-threaded process => many tasks



# Task Struct in Linux

```
struct task_struct {  
    // ...  
    pid_t pid;  
    pid_t tgid;  
    int prio;  
    volatile long state;  
    struct mm_struct *mm;  
    struct files_struct *files;  
    struct list_head tasks;  
    int on_cpu;  
    cpumask_t cpus_allowed;  
    // ...  
}
```

Task Creation: clone

*clone (function, stack\_ptr, sharing\_flags, args)*

sharing_flags	Meaning when set	Meaning when cleared
CLONE_VM	Create a new thread	Create a new process
CLONE_FS	Share unmask, root, and working dirs	Do not share unmask, root, and working dirs
CLONE_FILES	Share the file descriptors	Copy the file descriptors
CLONE_SIGHAND	Share the signal handler table	Copy the sig. handler table
CLONE_PID	New thread gets old PID	New thread gets own PID
CLONE_PARENT	New thread has same parent as caller	New thread's parent is caller

- fork is implemented via clone

## Linux Threads model

Native POSIX Threads Library (NPTL) "1:1 model"

- kernel sees each ULT info
- kernel traps are cheaper
- more resources: memory, large range of IDs

Older LinuxTreads "M-M model"

- similar issues to those described in Solaris papers

- kernel traps: user to kernel level crossing