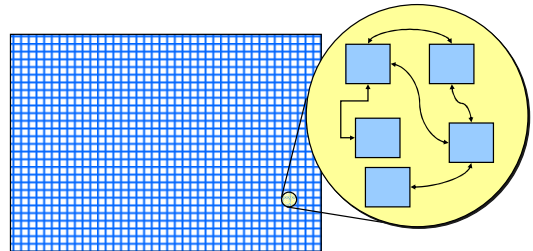


Microkernel Construction

IPC Implementation

IPC Importance



General IPC Algorithm

- Validate parameters
- Locate target thread
 - if unavailable, deal with it
- Transfer message
 - untyped - short IPC
 - typed message - long IPC
- Schedule target thread
 - switch address space as necessary
- Wait for IPC

IPC - Implementation

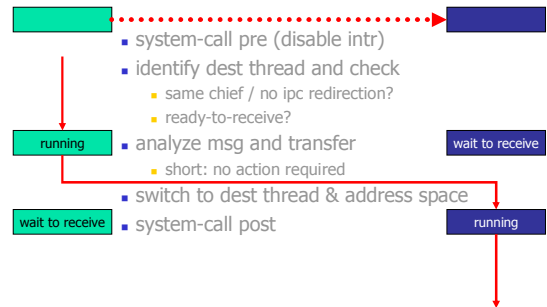
Short IPC

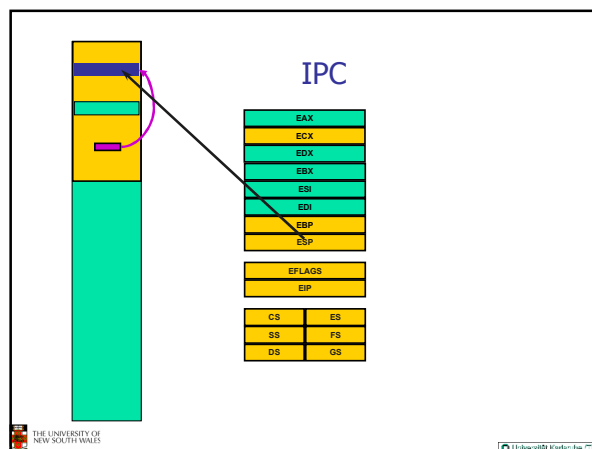
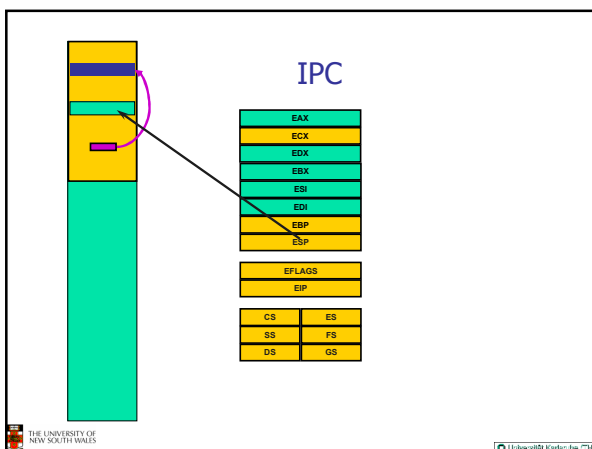
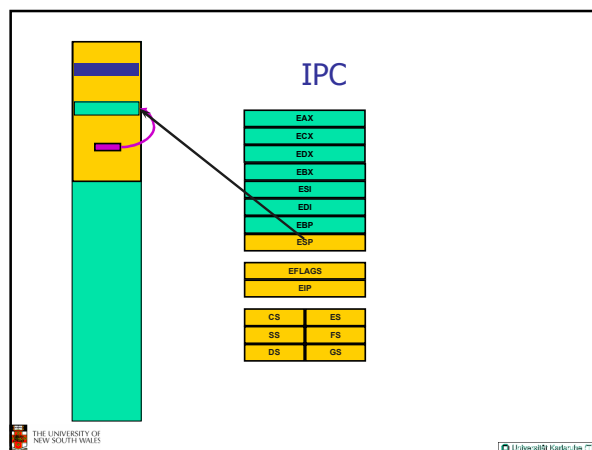
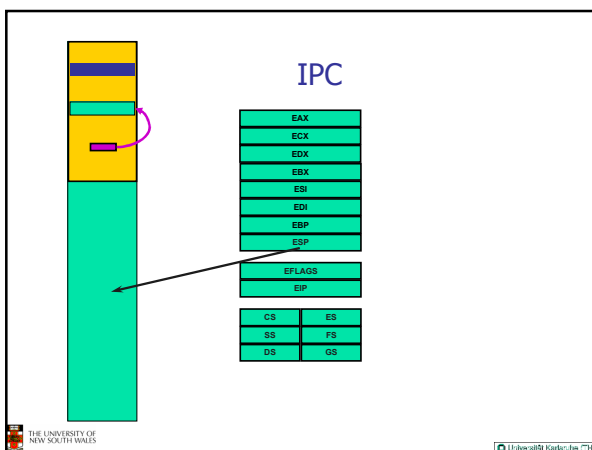
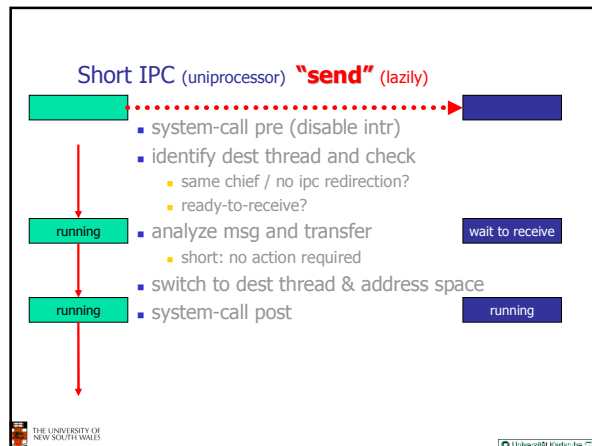
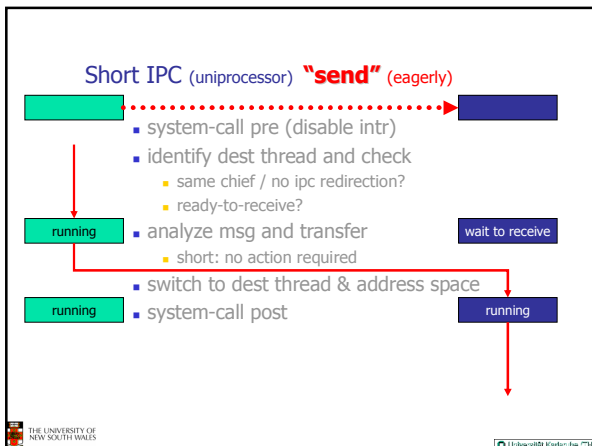
Short IPC (uniprocessor)

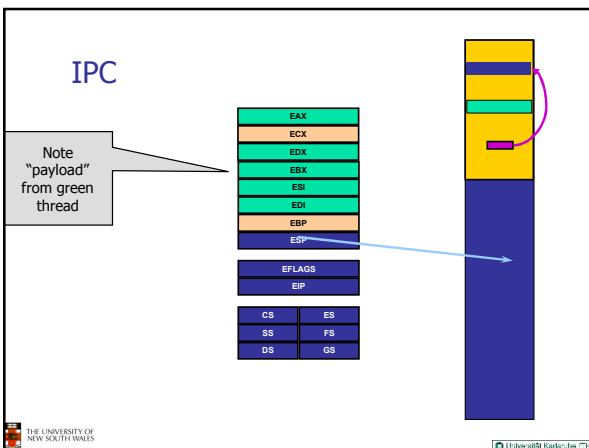
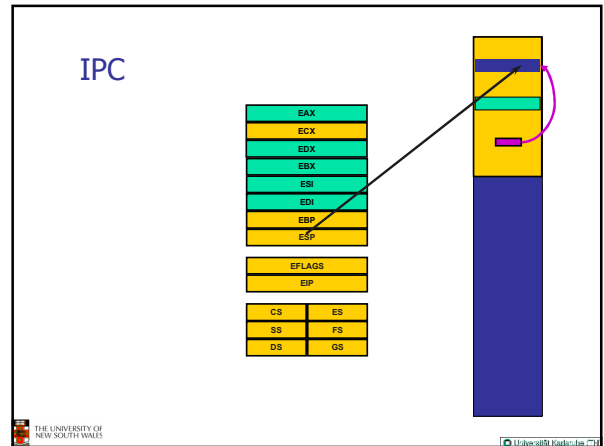
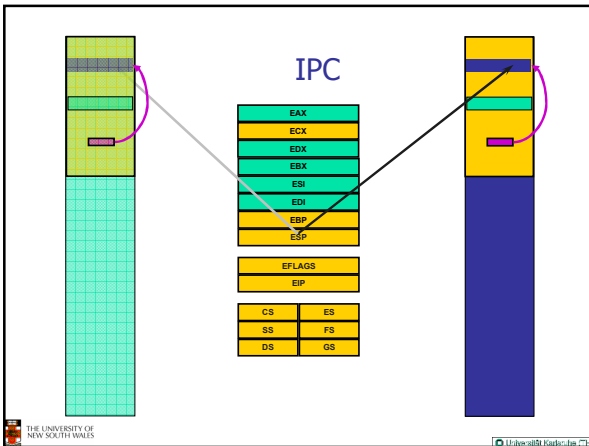
- system-call preamble (disable intr)
- identify dest thread and check
 - same chief / no ipc redirection?
 - ready-to-receive?
- analyze msg and transfer
 - short: no action required
- switch to dest thread & address space
- system-call postamble

The critical path

Short IPC (uniprocessor) "call"







- ### Implementation Goal
- Most frequent kernel op: short IPC
 - thousands of invocations per second
 - Performance is critical:
 - structure IPC for speed
 - **structure entire kernel to support fast IPC**
 - What affects performance?
 - cache line misses
 - TLB misses
 - memory references
 - pipe stalls and flushes
 - instruction scheduling

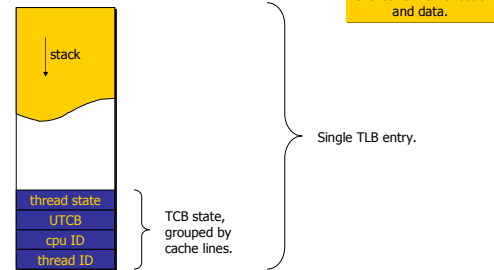
- ### Fast Path
- Optimize for common cases
 - write in assembler
 - non-critical paths written in C++
 - but still fast as possible
 - Avoid high-level language overhead:
 - function call state preservation
 - poor code "optimizations"
 - We want every cycle possible!

- ### IPC Attributes for Fast Path
- untyped message
 - single runnable thread after IPC
 - must be valid IPC call
 - switch threads, originator blocks
 - send phase:
 - the target is waiting
 - receive phase:
 - the sender is not ready to couple, causing us to block
 - no receive timeout

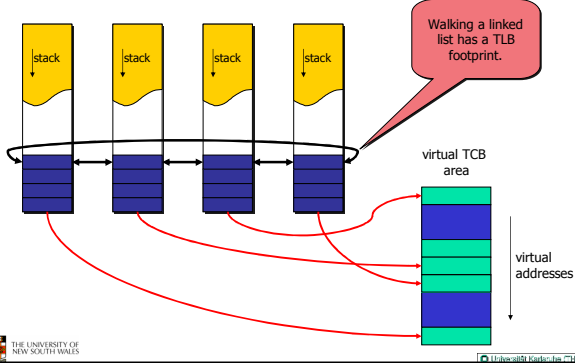
Avoid Memory References!!!

- Memory references are slow
 - avoid in IPC:
 - ex: use lazy scheduling
 - avoid in common case:
 - ex: timeouts
- Microkernel should minimize indirect costs
 - cache pollution
 - TLB pollution
 - memory bus

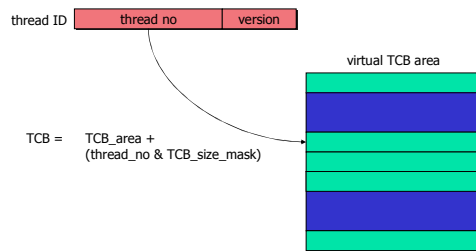
Optimized Memory



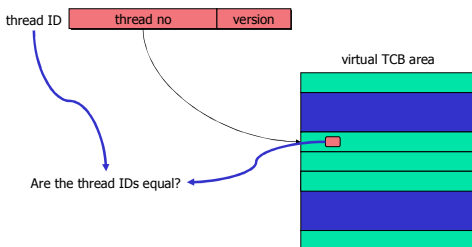
TLB Problem



Avoid Table Lookups



Validate Thread ID



Branch Elimination

```

slow = ~receiver->thread_state +
      (timeouts & 0xffff) +
      sender->resources +
      receiver->resources;
if( slow )
    enter_slow_path();
    
```

- Reduces branch prediction footprint.
- Avoids mispredicts & stalls & flushes.
 - Increases latency for slow path

Common case: -1

Common case: 0

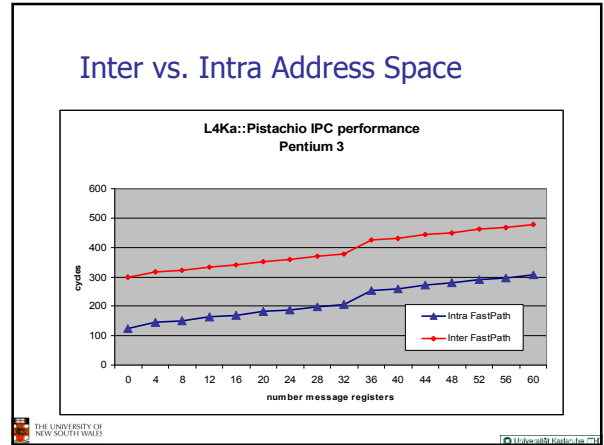
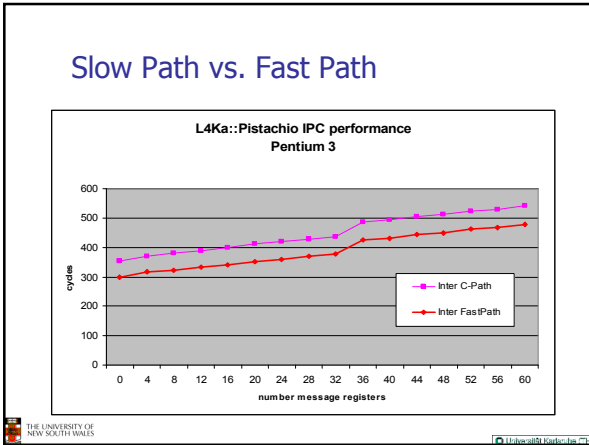
TCB Resources

- One bit per resource
- Fast path checks entire word
 - if not 0, jump to resource handlers

The diagram shows a 'Resources bitfield' with two bits. The left bit is connected to a green box labeled 'Copy area'. The right bit is connected to a yellow box labeled 'Debug registers'.

Message Transfer

The graph plots 'cycles per IPC' (y-axis, 180-280) against 'message registers copied' (x-axis, 0-60). Two lines are shown: 'instructions per IPC' (green) and 'cycles per IPC' (red). A callout indicates 'up to 10 physical registers' at x=10. Another callout shows a 'virtual register copy loop' between x=10 and x=40. A text box notes 'Many cycles wasted on pipe flushes for privileged instructions.' Text on the right specifies 'IBM PowerPC 750, 500 MHz, 32 registers'.



IPC - Implementation

Long IPC

Long IPC (uniprocessor)

- system-call preamble (disable intr)
 - Preemptions possible! (end of timeslice, device interrupt...)
- identify dest thread and check
 - same chief
 - ready-to-recv?
- analyze msg and transfer
 - long/map:
 - Pagefaults possible! (in source and dest address space)

— transfer message —

- switch to dest thread & address space
- system-call postamble

Long IPC (uniprocessor)

- system-call pre (disable intr)
- identify dest thread and check
 - same chief
 - ready-to-receive?
- analyze msg and transfer
 - long/map:
 - lock both partners
- transfer message -
- unlock both partners
- switch to dest thread & address space
- system-call post

Preemptions possible!
(end of timeslice, device interrupt...)

Pagefaults possible!
(in source and dest address space)

Long IPC (uniprocessor)

- system-call pre (disable intr)
- identify dest thread and check
 - same chief
 - ready-to-receive?
- analyze msg and transfer
 - long/map:
 - lock both partners
 - enable intr
 - transfer message -
 - disable intr
 - unlock both partners
- switch to dest thread & address space
- system-call post

Preemptions possible!
(end of timeslice, device interrupt...)

Pagefaults possible!
(in source and dest address space)

Long IPC (uniprocessor)

- system-call pre (disable intr)
- identify dest thread and check
 - same chief
 - ready-to-receive?
- analyze msg and transfer
 - long/map:
 - lock both partners
 - enable intr
 - transfer message -
 - disable intr
 - unlock both partners
- switch to dest thread & address space
- system-call post

IPC - mem copy

- Why is it needed? Why not share?
 - Security
 - Need own copy
 - Granularity
 - Object small than a page or not aligned

copy in - copy out

- copy into kernel buffer

copy in - copy out

- copy into kernel buffer
- switch spaces

copy in - copy out

- copy into kernel buffer
- switch spaces
- copy out of kernel buffer

costs for n words

- $2 \times 2n$ r/w operations
- $3 \times n/8$ cache lines
 - $1 \times n/8$ overhead cache misses (small n)
 - $4 \times n/8$ cache misses (large n)

The diagram shows two vertical bars representing memory spaces. The left bar has a green bottom section and a yellow top section. The right bar has a blue bottom section and a yellow top section. A red dashed arrow points from a green box in the left bar to a blue box in the right bar. A yellow box is shown in the top of both bars, representing a kernel buffer. A vertical pink line connects the green box to the yellow box in the left bar, and another vertical pink line connects the yellow box to the blue box in the right bar.

temporary mapping

The diagram shows two vertical bars. The left bar has a green bottom section and a yellow top section. The right bar has a blue bottom section and a yellow top section. A red dashed arrow points from a green box in the left bar to a blue box in the right bar. A yellow box is shown in the top of both bars, representing a kernel buffer. A vertical pink line connects the green box to the yellow box in the left bar, and another vertical pink line connects the yellow box to the blue box in the right bar.

temporary mapping

- select dest area (4+4 M)

The diagram shows two vertical bars. The left bar has a green bottom section and a yellow top section. The right bar has a blue bottom section and a yellow top section. A red dashed arrow points from a green box in the left bar to a blue box in the right bar. A yellow box is shown in the top of both bars, representing a kernel buffer. A vertical pink line connects the green box to the yellow box in the left bar, and another vertical pink line connects the yellow box to the blue box in the right bar.

temporary mapping

- select dest area (4+4 M)
- map into source AS (kernel)

The diagram shows two vertical bars. The left bar has a green bottom section and a yellow top section. The right bar has a blue bottom section and a yellow top section. A red dashed arrow points from a green box in the left bar to a blue box in the right bar. A yellow box is shown in the top of both bars, representing a kernel buffer. A vertical pink line connects the green box to the yellow box in the left bar, and another vertical pink line connects the yellow box to the blue box in the right bar.

temporary mapping

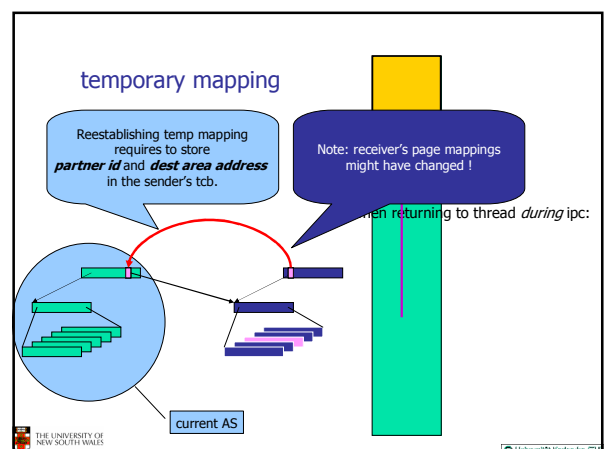
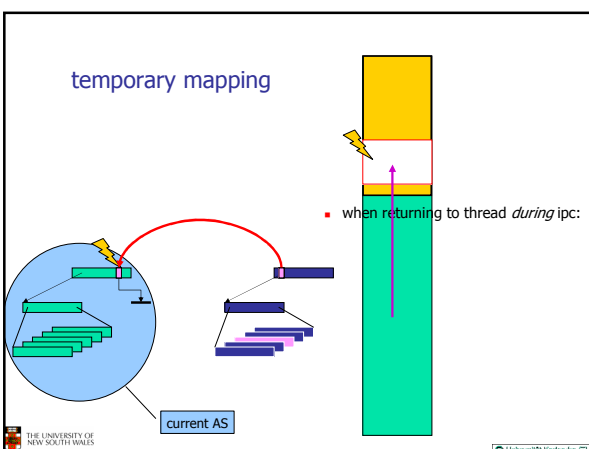
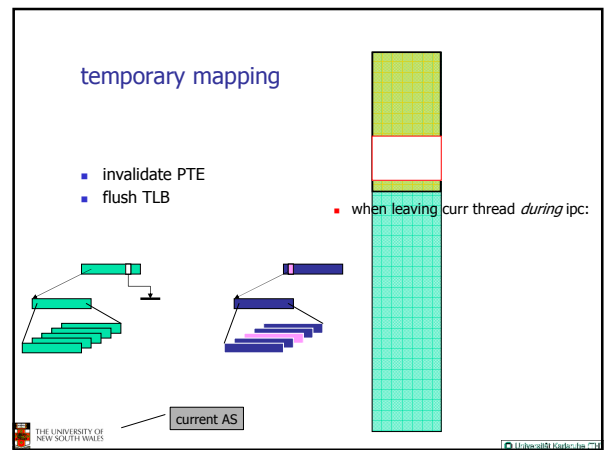
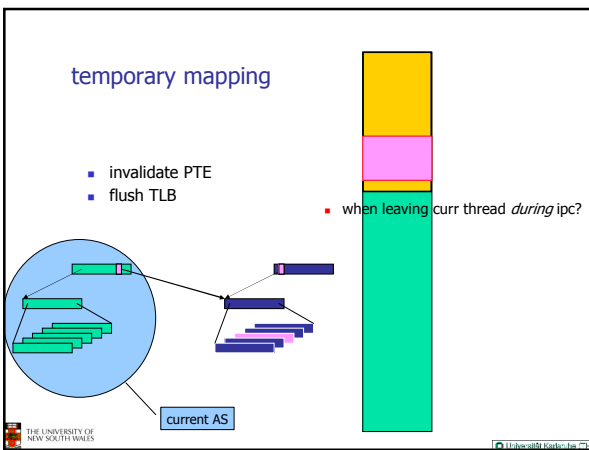
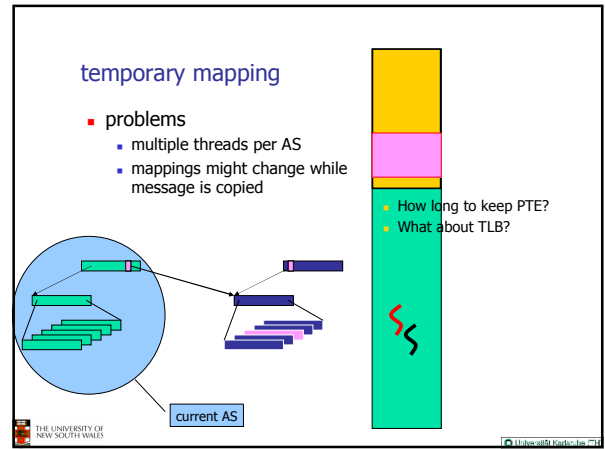
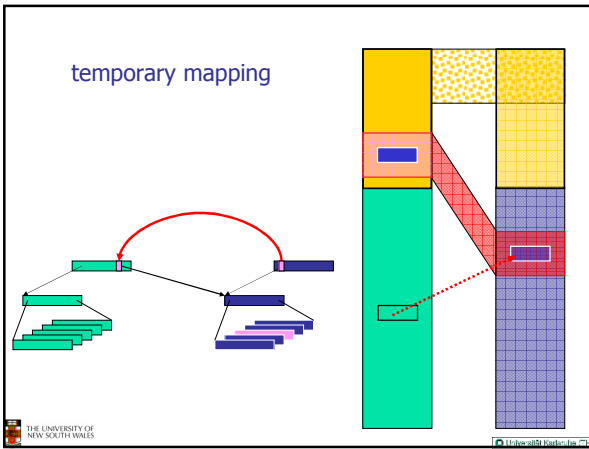
- select dest area (4+4 M)
- map into source AS (kernel)
- copy data

The diagram shows two vertical bars. The left bar has a green bottom section and a yellow top section. The right bar has a blue bottom section and a yellow top section. A red dashed arrow points from a green box in the left bar to a blue box in the right bar. A yellow box is shown in the top of both bars, representing a kernel buffer. A vertical pink line connects the green box to the yellow box in the left bar, and another vertical pink line connects the yellow box to the blue box in the right bar.

temporary mapping

- select dest area (4+4 M)
- map into source AS (kernel)
- copy data
- switch to dest space

The diagram shows two vertical bars. The left bar has a green bottom section and a yellow top section. The right bar has a blue bottom section and a yellow top section. A red dashed arrow points from a green box in the left bar to a blue box in the right bar. A yellow box is shown in the top of both bars, representing a kernel buffer. A vertical pink line connects the green box to the yellow box in the left bar, and another vertical pink line connects the yellow box to the blue box in the right bar.



temporary mapping

```

Start temp mapping:
mytcb.partner := partner ;
mytcb.waddr := dest 8M area base ;
myPDE.TMarea := destPDE.destarea .

Leave thread:
if mytcb.waddr ≠ nil then
  myPDE.TMarea := nil ;
  if dest AS = my AS then
    flush TLB
  fi
fi

Close temp mapping:
mytcb.waddr := nil .
myPDE.TMarea := nil ??
    
```

optimization only:
avoids second TLB flush if subsequent thread switch would flush TLB anyhow

THE UNIVERSITY OF NEW SOUTH WALES

temporary mapping

- Alternative method:

```

Leave thread:
if mytcb.waddr ≠ nil then
  myPDE.TMarea := nil ;
  flush TLB ;
  TLB flushed := true
fi .

Thread switch :
...
if TLB just flushed
  then TLB flushed := false
  else flush TLB
fi ;
PT root := ...
    
```

Requires separation of TLB flush and load PT root !

Does therefore not work reasonably on x86.

Load PT root implicitly includes TLB flush on x86.

current AS

THE UNIVERSITY OF NEW SOUTH WALES

temporary mapping

- Page Fault Resolution:

current AS

THE UNIVERSITY OF NEW SOUTH WALES

temporary mapping

- Page Fault Resolution:

current AS

THE UNIVERSITY OF NEW SOUTH WALES

temporary mapping

- Page Fault Resolution:

current AS

THE UNIVERSITY OF NEW SOUTH WALES

temporary mapping

- Page Fault Resolution:

```

TM area PF:
if myPDE.TMarea = destPDE.destarea then
  tunnel to (partner) ;
  access dest area ;
  tunnel to (my)
fi ;
myPDE.TMarea := destPDE.destarea .
    
```

current AS

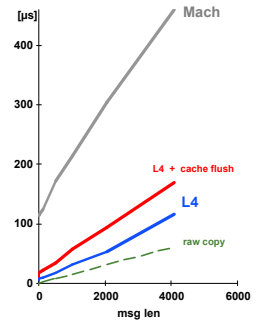
THE UNIVERSITY OF NEW SOUTH WALES

Cost estimates

	Copy in - copy out	Temporary mapping
<i>R/W operations</i>	$2 \times 2n$	$2n$
<i>Cache lines</i>	$3 \times n/8$	$2 \times n/8$
<i>Small n overhead cache misses</i>	$n/8$	0
<i>Large n cache misses</i>	$5 \times n/8$	$3 \times n/8$
<i>Overhead TLB misses</i>	0	$n / \text{words per page}$
<i>Startup instructions</i>	0	50

486 IPC costs

- Mach: copy in/out
- L4: temp mapping



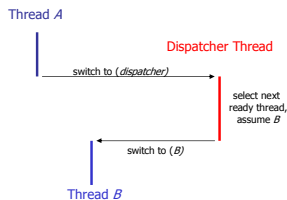
Dispatching

Dispatching topics:

- thread switch
 - (to a specific thread)
 - to next thread to be scheduled
 - (to nil)
 - implicitly, when ipc blocks
- priorities
- preemption
 - time slices
 - wakeups, interruptions
- timeouts and wake-ups
 - time

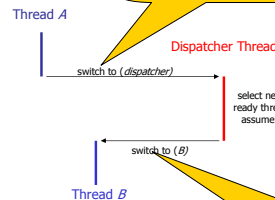
Switch to ():

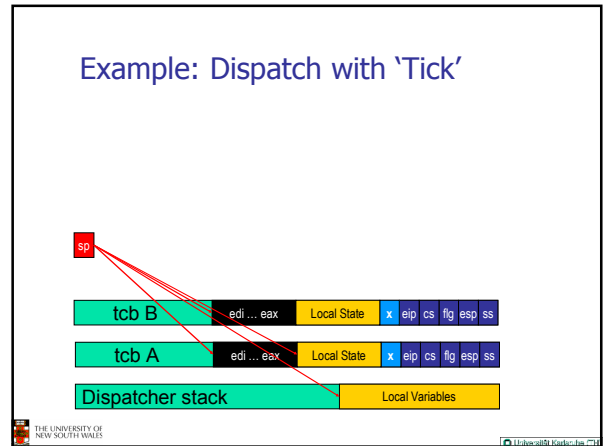
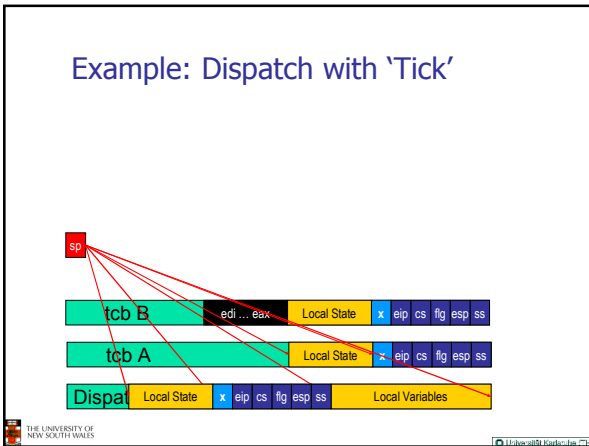
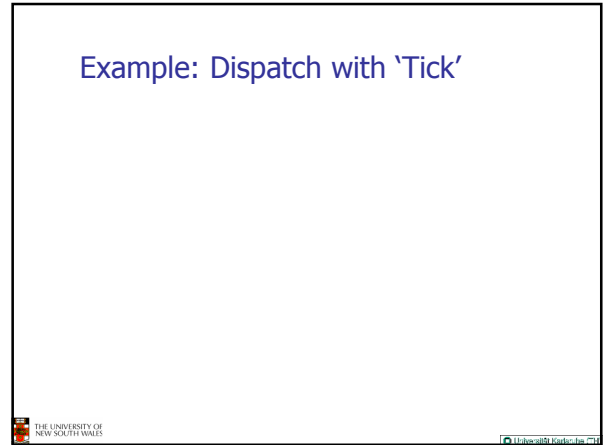
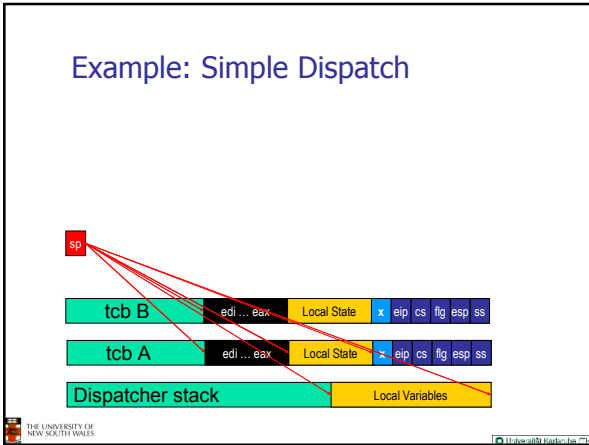
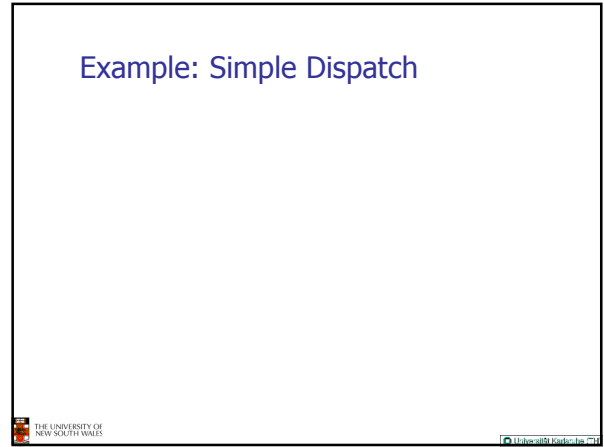
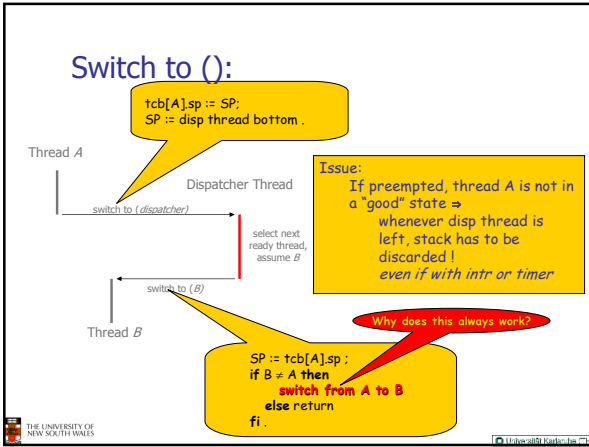
- Smaller stack per thread
- Dispatcher is preemptable
 - Improved interrupt latency if dispatching is time consuming



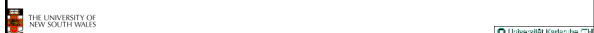
Switch to ():

- Optimizations :
 - disp thread is special
 - no user mode,
 - no own AS required
 - Can avoid AS switch
 - no id required
 - Freedom from tcb layout conventions
 - almost stateless (see priorities)
 - No need to preserve internal state between invocations
 - External state must be consistent
- costs (A → B)
 - costs (A → disp → B)
 - costs (select next)
 - costs (A → disp → A) are low

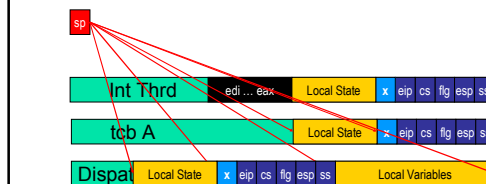




Example: Dispatch with Interrupt



Example: Dispatch with Interrupt

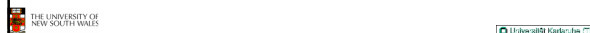
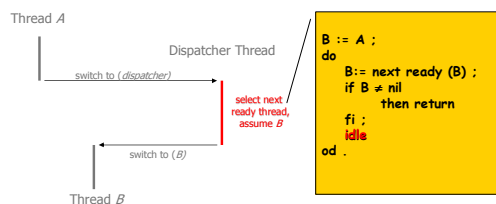


Example: Dispatch with Interrupt



Switch to ():

- dispatcher thread is also **idle thread**

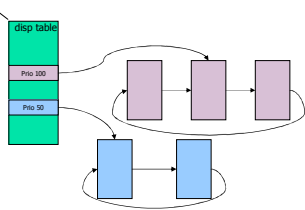


Priorities

- 0 (lowest) ... 255
- hard priorities
- round robin per prio
- dynamically changeable
- ready tcb list per prio
- 'current tcb' per list

```

do
  p := 255;
  do
    if currentp = nil
    then B := currentp;
      return
    fi;
    p -= 1
  until p < 0 od;
  idle
od
    
```

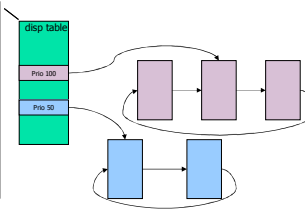


Priorities

- Optimization
 - keep highest active prio

```

do
  if current[highest active p] = nil
  then B := current[highest active p];
    return
  elif highest active p > 0
  then highest active p -= 1
  else
    idle
  fi
od
    
```



Priorities, Preemption

highest active p := max (new p, highest active p)

```

do
  if current[highest active p] ≠ nil
  then B := current[highest active p];
  return
  elif highest active p > 0
  then highest active p -= 1
  else
  idle
  fi
od .
  
```

disp table
 Prio 110
 Prio 100
 Prio 50

p=110
intr/wakeup

THE UNIVERSITY OF NEW SOUTH WALES

Priorities, Preemption

- What happens when a prio falls empty?

```

do
  if current[highest active p] ≠ nil
  then round robin if necessary;
  B := current[highest active p];
  return
  elif highest active p > 0
  then highest active p -= 1
  else
  idle
  fi
od .
  
```

disp table
 Prio 110
 Prio 100
 Prio 50

Remaining time slice > 0?

THE UNIVERSITY OF NEW SOUTH WALES

Priorities, Preemption

- What happens when a prio falls empty?

```

do
  if current[highest active p] ≠ nil
  then round robin if necessary;
  B := current[highest active p];
  return
  elif highest active p > 0
  then highest active p -= 1
  else
  idle
  fi
od .
  
```

disp table
 Prio 110
 Prio 100
 Prio 50

Remaining time slice > 0?

THE UNIVERSITY OF NEW SOUTH WALES

Preemption

- Preemption, time slice exhausted

```

do
  if current[highest active p] ≠ nil
  then round robin if necessary;
  B := current[highest active p];
  return
  elif highest active p > 0
  then highest active p -= 1
  else
  idle
  fi
od .
  
```

disp table
 Prio 110
 Prio 100
 Prio 50

Remaining time slice = 0

THE UNIVERSITY OF NEW SOUTH WALES

Preemption

- Preemption, time slice exhausted

```

do
  if current[highest active p] ≠ nil
  then round robin if necessary;
  B := current[highest active p];
  return
  elif highest active p > 0
  then highest active p -= 1
  else
  idle
  fi
od .
  
```

disp table
 Prio 110
 Prio 100
 Prio 50

Remaining time slice = 0

THE UNIVERSITY OF NEW SOUTH WALES

Lazy Dispatching

Thread state toggles frequently (per ipc)

- ready ↔ waiting
- delete/insert ready list is expensive
- therefore: delete *lazily* from ready list

disp table
 Prio 110
 Prio 100
 Prio 50

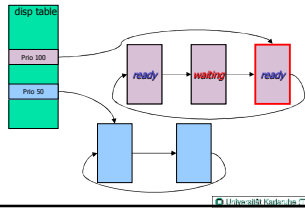
ready ready ready

THE UNIVERSITY OF NEW SOUTH WALES

Lazy Dispatching

Thread state toggles frequently (per ipc)

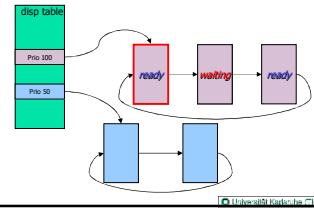
- ready ↔ waiting
 - delete/insert ready list is expensive
 - therefore: delete *lazily* from ready list



Lazy Dispatching

Thread state toggles frequently (per ipc)

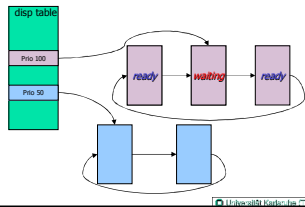
- ready ↔ waiting
 - delete/insert ready list is expensive
 - therefore: delete *lazily* from ready list



Lazy Dispatching

Thread state toggles frequently (per ipc)

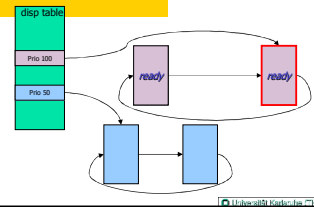
- ready ↔ waiting
 - delete/insert ready list is expensive
 - therefore: delete *lazily* from ready list



Lazy Dispatching

Thread state toggles frequently (per ipc)

- ready ↔ waiting
 - delete/insert ready list is expensive
 - therefore: delete *lazily* from ready list
 - Whenever reaching a non-ready thread,
 - delete it from list
 - proceed with next



Lazy Dispatching

```

do
  round robin if necessary:
  if current(highest active p) ≠ nil
    then B := current(highest active p); return
  elif highest active p > 0
    then highest active p -= 1
  else
    idle
  fi
od .

round robin if necessary:
while curr(hi act p) ≠ nil do
  if curr(hi act p).state ≠ ready
    then deletes from list (curr(hi act p))
  elif curr(hi act p).rem ts = 0
    then curr(hi act p).rem ts := new ts
  else leave round robin if necessary
  fi :
  curr(hi act p) := next ;
od .
    
```

