

# Microkernel Construction

## I.11 – Small Address Spaces (Special Optimization for Untagged TLBs)

Lecture Summer Term 2017

Wednesday 15:45-17:15 R 131, 50.34 (INFO)

Jens Kehne | Marius Hillenbrand  
Operating Systems Group, Department of Computer Science



# Untagged TLB Context-Switch Costs

- Enter/exit kernel
- Switch thread
- Switch address space
  - Flush TLB
  - Refill TLB
- Refill L1 caches  
(P4 only, not P3/Core)

## 486 ... PIII

40 ... 200 cycles

≈ 10 cycles

≈ 50 ... 80 cycles

6 ... 96 TLB refills  
15 ... 40 cycles/refill  
≈ 100 ... 4000 cycles

## Pentium 4

150 ... 200 cycles

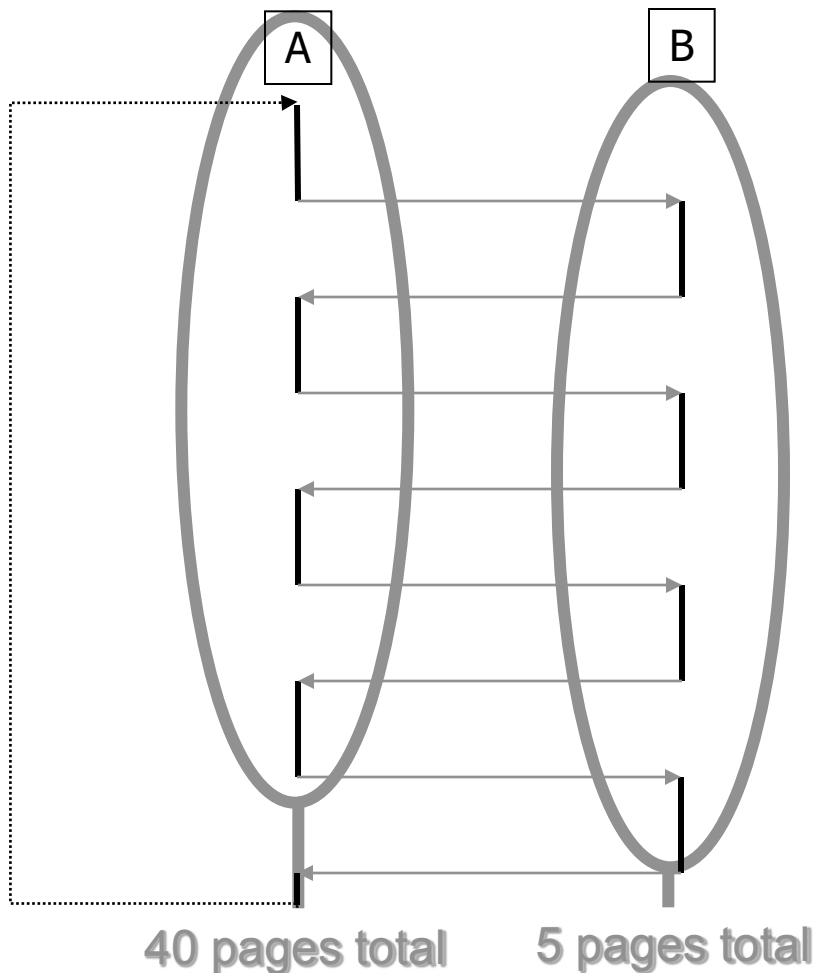
≈ 10 cycles

≈ 230 ... 250 cycles

6 ... 192 TLB refills  
15 ... 500 cycles/refill  
≈ 100 ... 96000 cycles

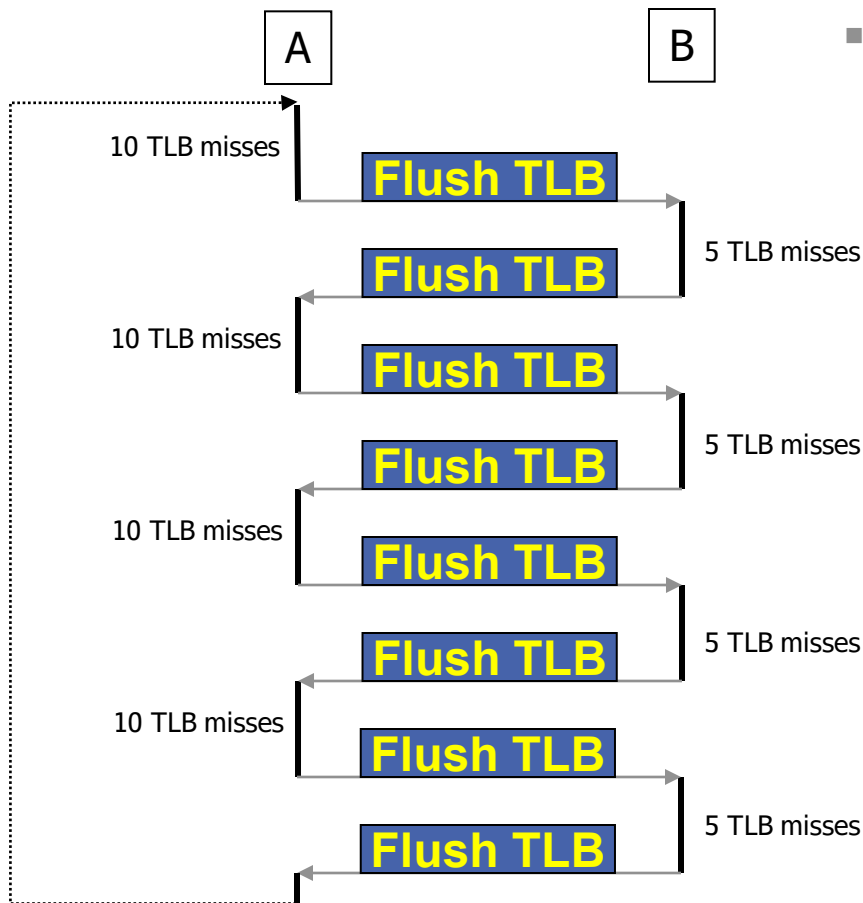
12 K trace cache  
8 K data cache  
15 ... 25 cycles/refill  
≈ 100 ... 16000 cycles

# Untagged TLB Context-Switch Costs



- Even when calling a thread with a very **small TLB working set**
  - Thread A frequently calls thread B
  - Working sets
    - Thread A: 4 different sets of 10 pages between B-calls
    - Thread B: always the same 5 pages

# Untagged TLB Context-Switch Costs



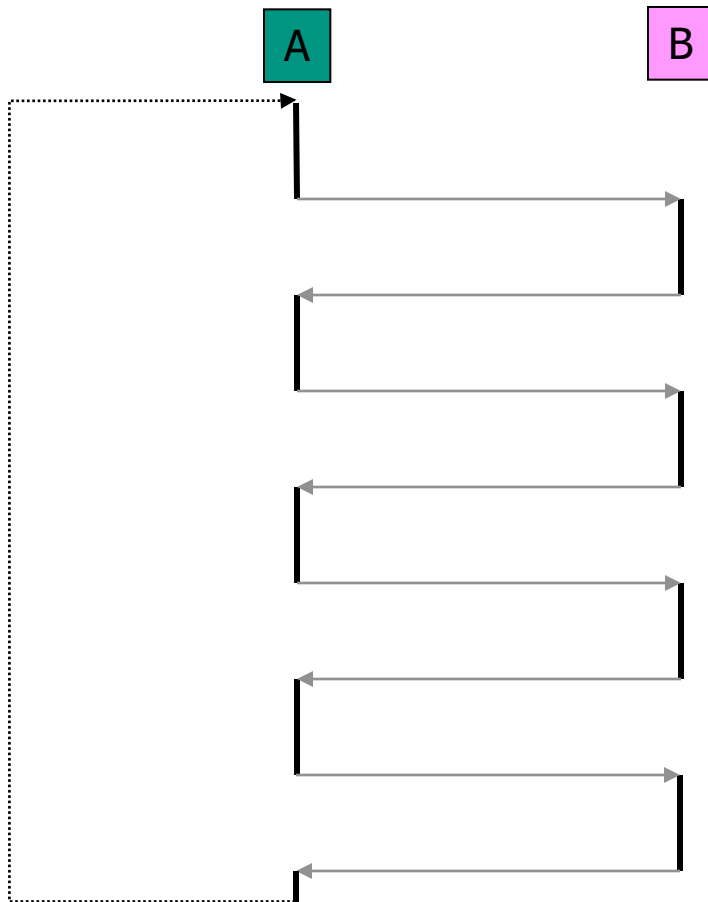
- Even when calling a thread with a very **small TLB working set**

- Thread A frequently calls thread B
- Working sets
  - Thread A: 4 different sets of 10 pages between B-calls
  - Thread B: always the same 5 pages

	[cycles]
■ 8 IPCs (w/o AS costs):	1440
■ 60 TLB misses:	900 ... 30000
■ 8 TLB flushes:	400 ... 2000

**Untagged TLB total: 2740 ... 33440**

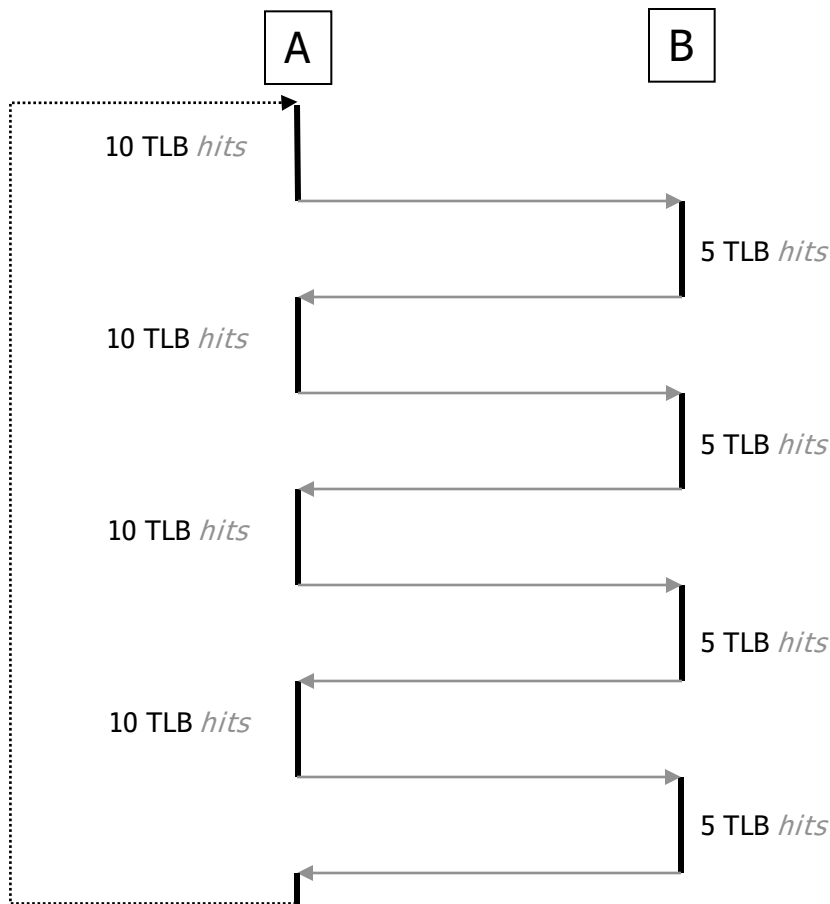
# Tagged TLB



- Tagged TLB:
  - Associate Address-Space ID with translations
  - No flushing during AS switch
  - Not available for x86 address spaces (only with virtualization extensions)

	ASID	VADDR	PADDR	sz	rwX	U/ K
	A					
	B					
	B					
	A					
	B					
	A					
	A					

# Untagged TLB Context-Switch Costs



- Even when calling a thread with a very **small TLB working set**

- Thread A frequently calls thread B
- Working sets
  - Thread A: 4 different sets of 10 pages between B-calls
  - Thread B: always the same 5 pages

- |                          | [cycles]      |
|--------------------------|---------------|
| ■ 8 IPCs (w/o AS costs): | 1440          |
| ■ 60 TLB misses:         | 900 ... 30000 |
| ■ 8 TLB flushes:         | 400 ... 2000  |

**Untagged TLB total: 2740 ... 33440**

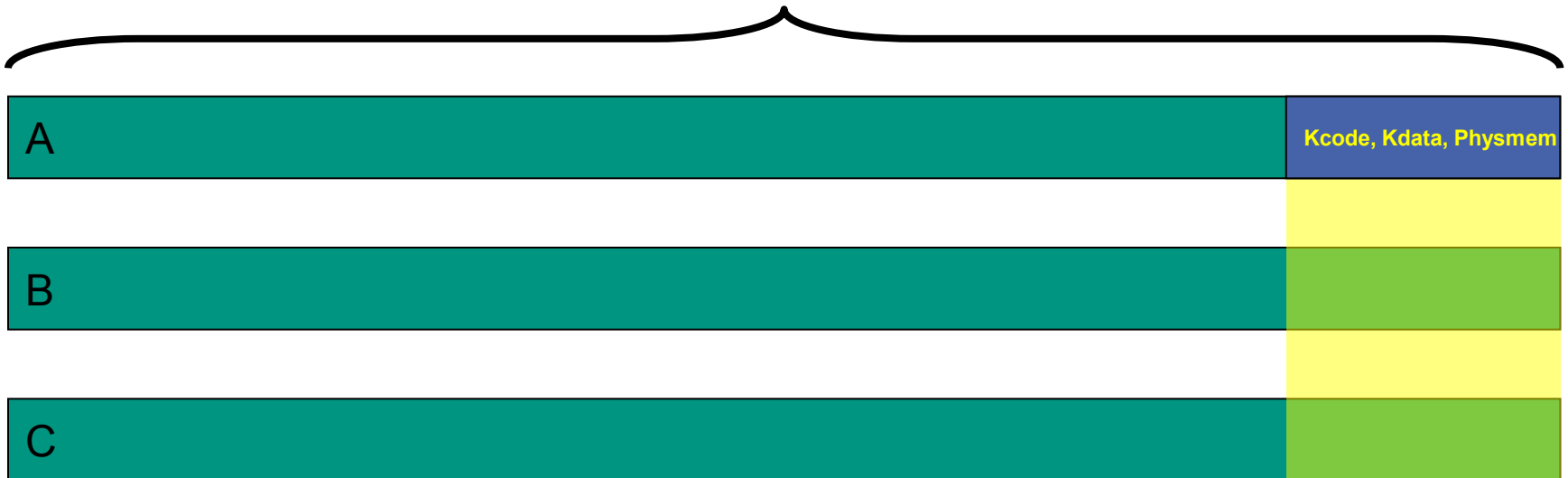
**Tagged TLB total:  $\approx$  1500**

# SMALL ADDRESS SPACES ON X86

How to emulate tagged TLBs?

# Address Spaces

CS, DS





# Address Spaces

kernel CS, DS

user CS, DS

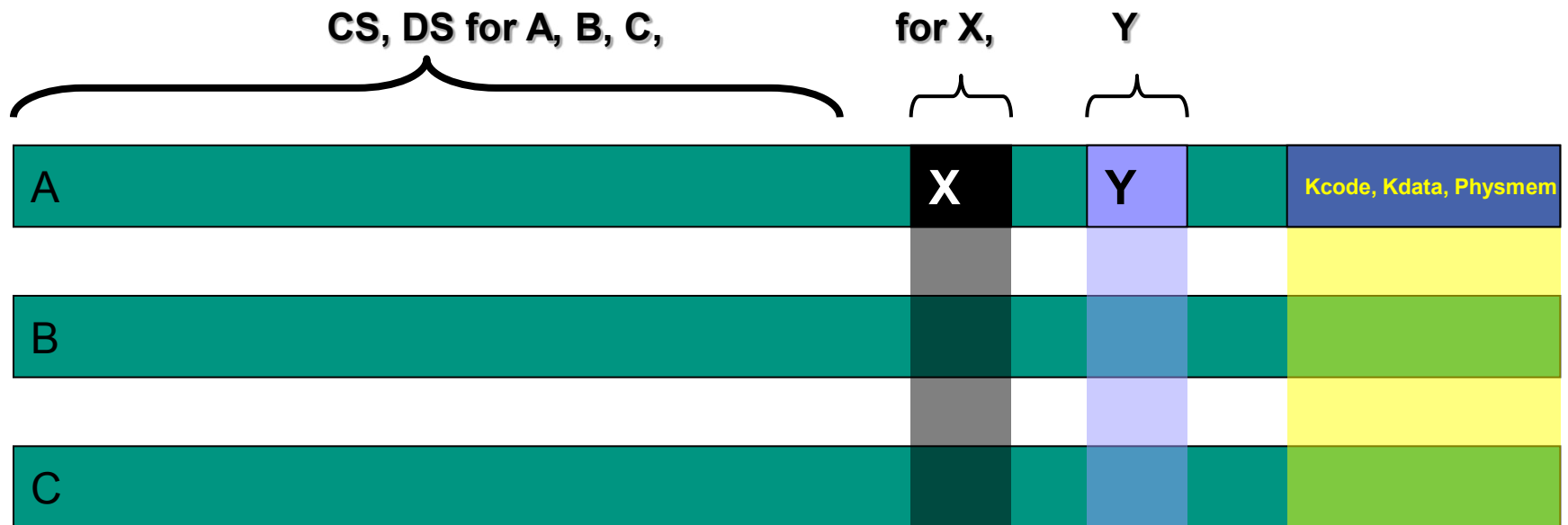


# Small Address Spaces

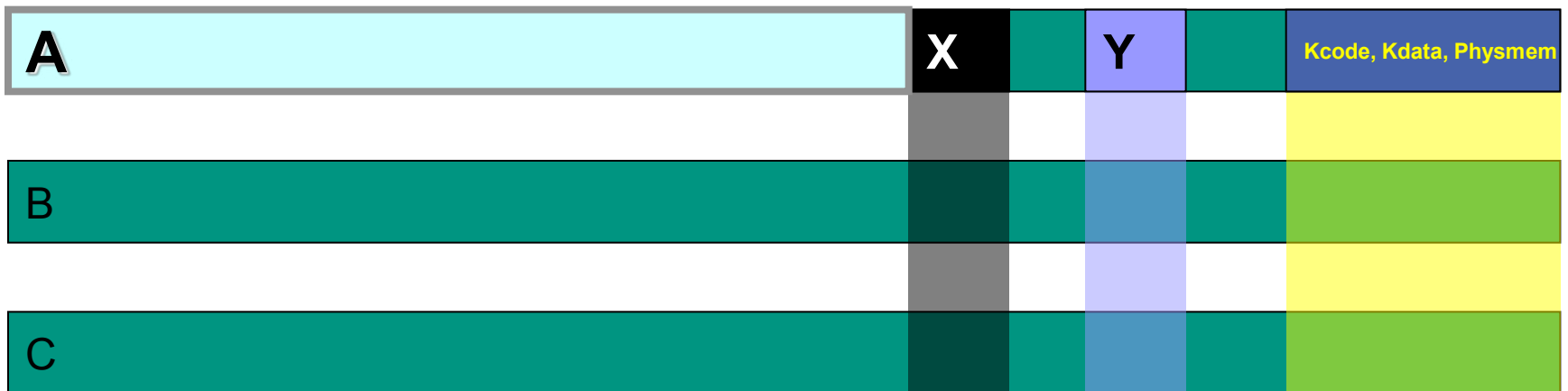
CS, DS for A, B, C



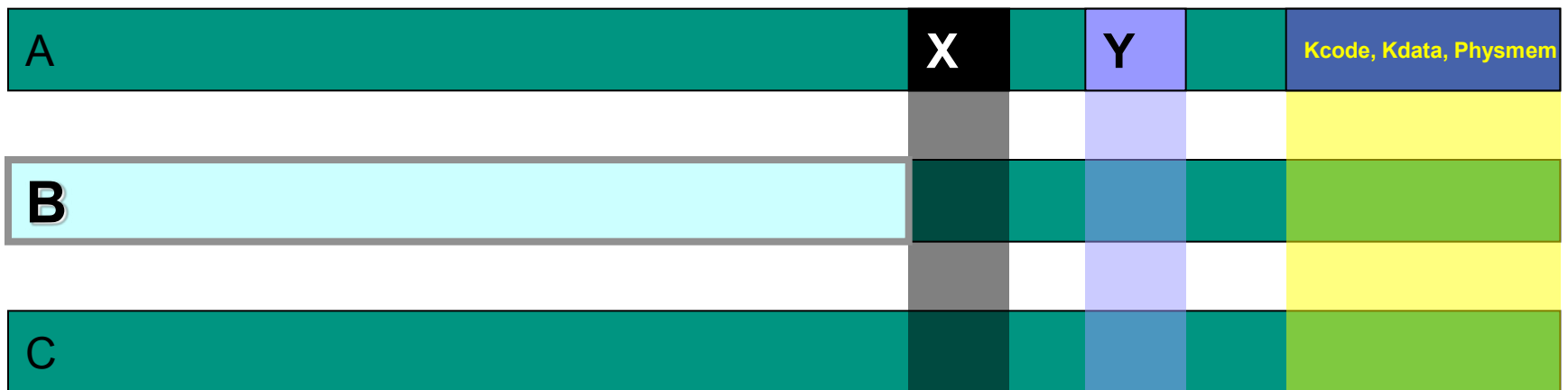
# Small Address Spaces



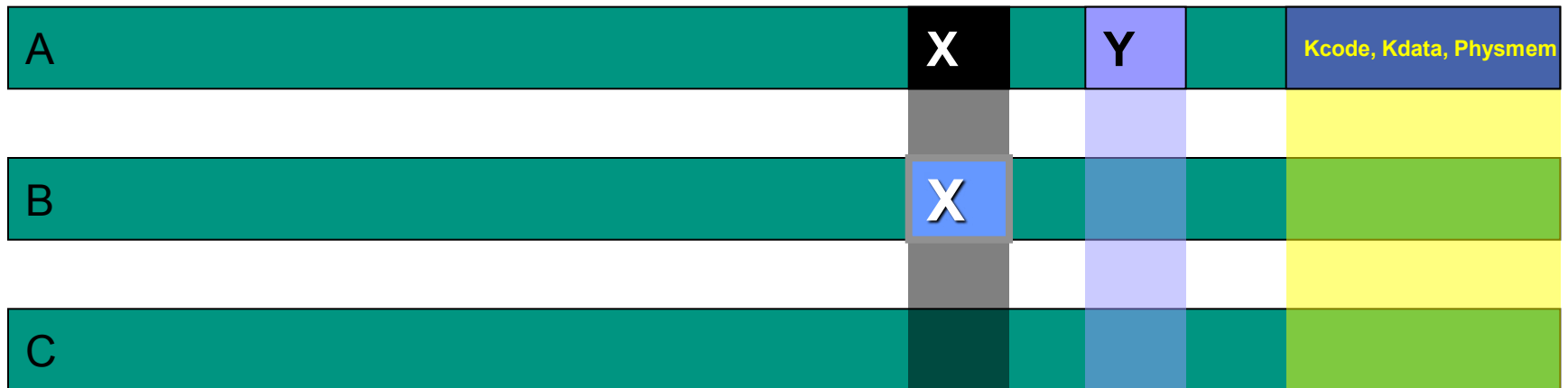
# Small Address Spaces



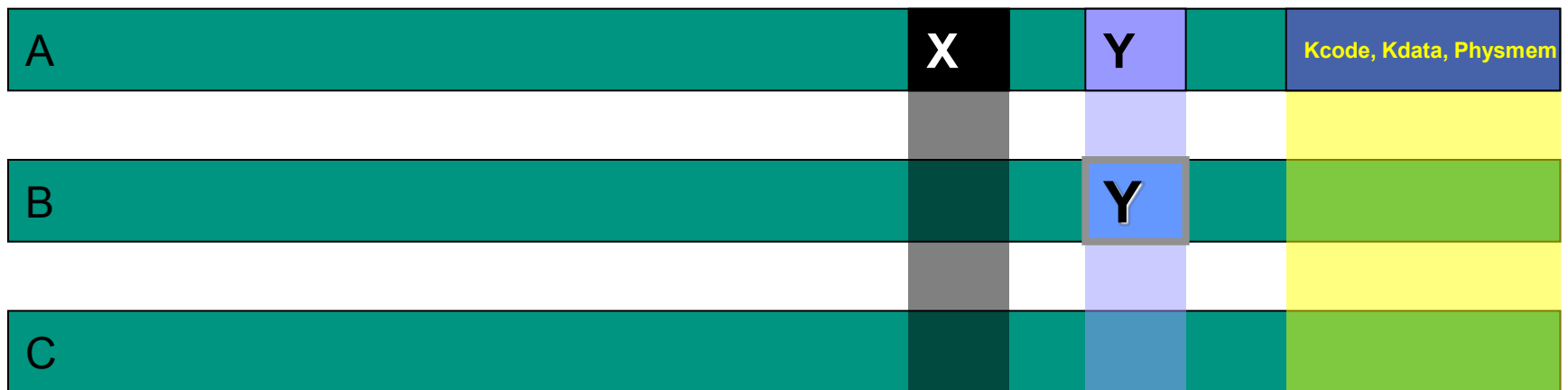
# Small Address Spaces



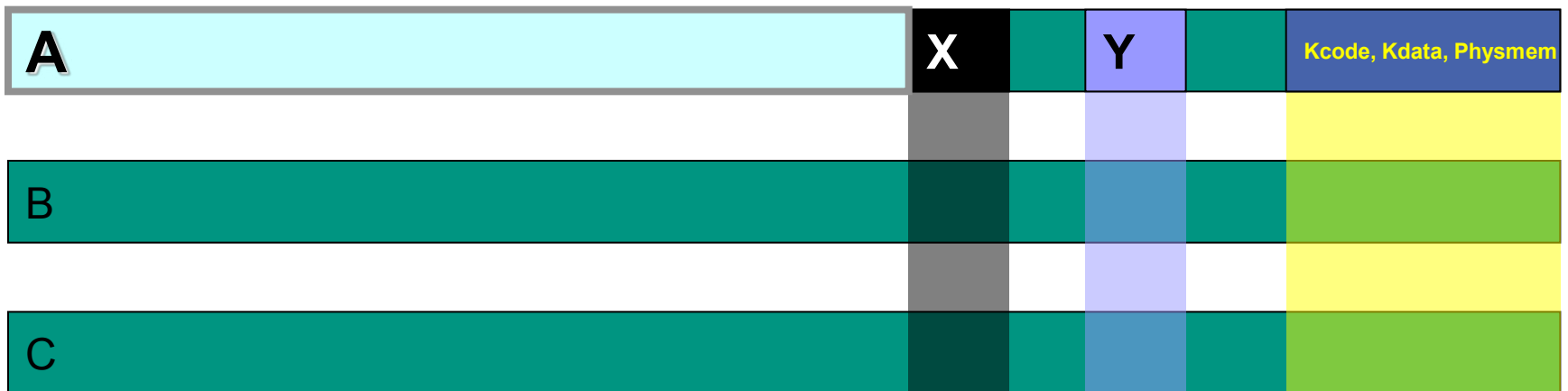
# Small Address Spaces



# Small Address Spaces

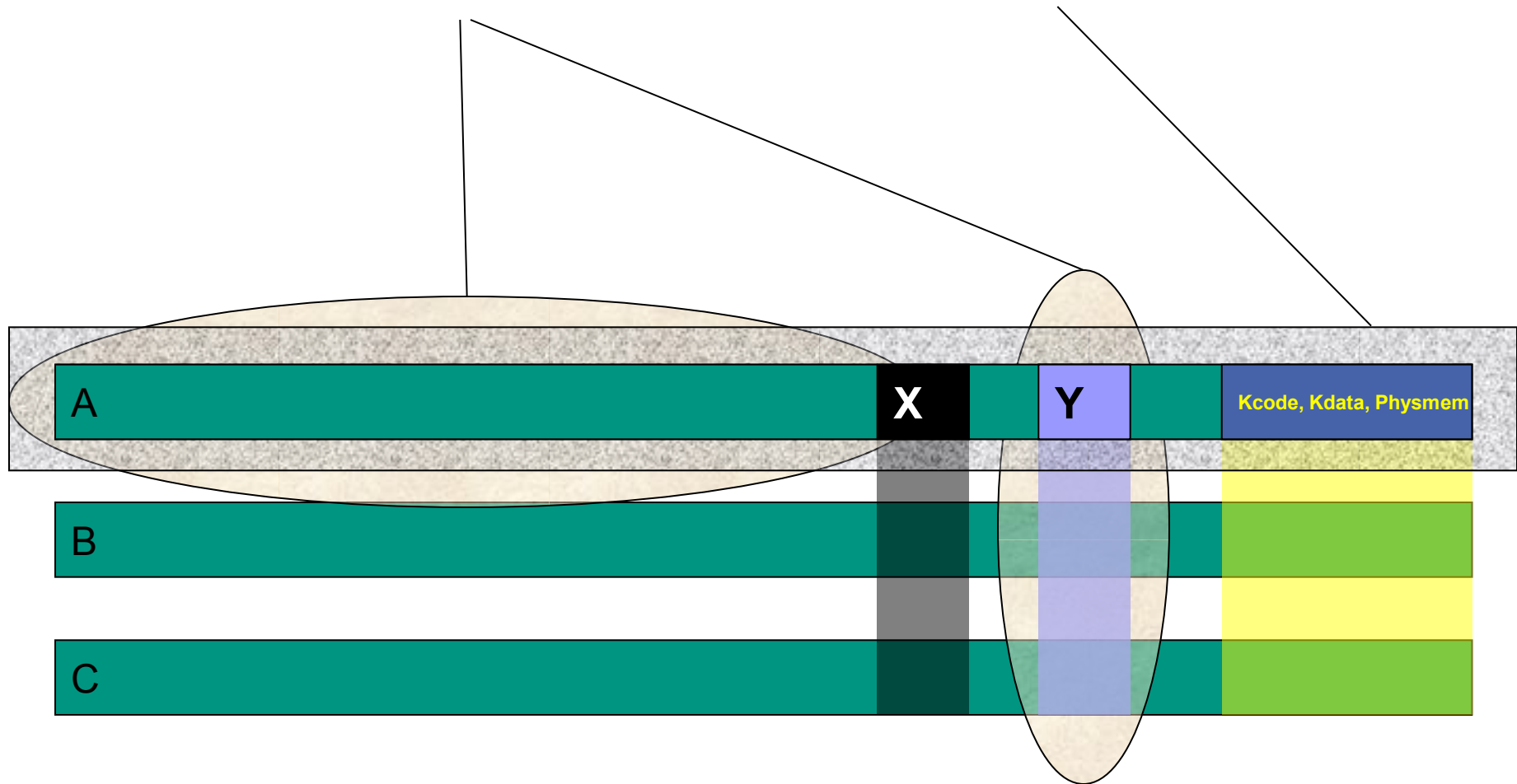


# Small Address Spaces

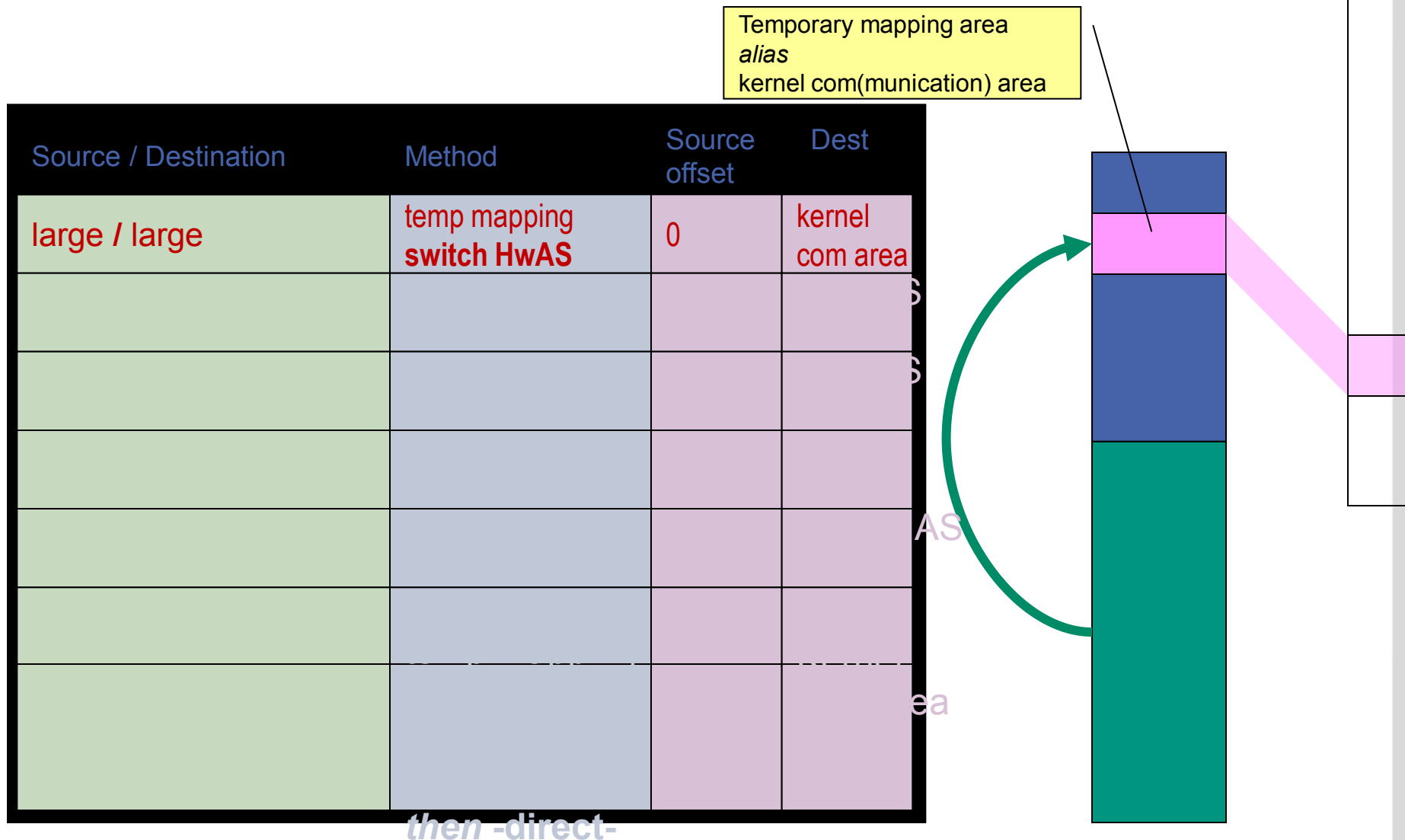




# Address Space (AS) vs. Hardware AS (HwAS)



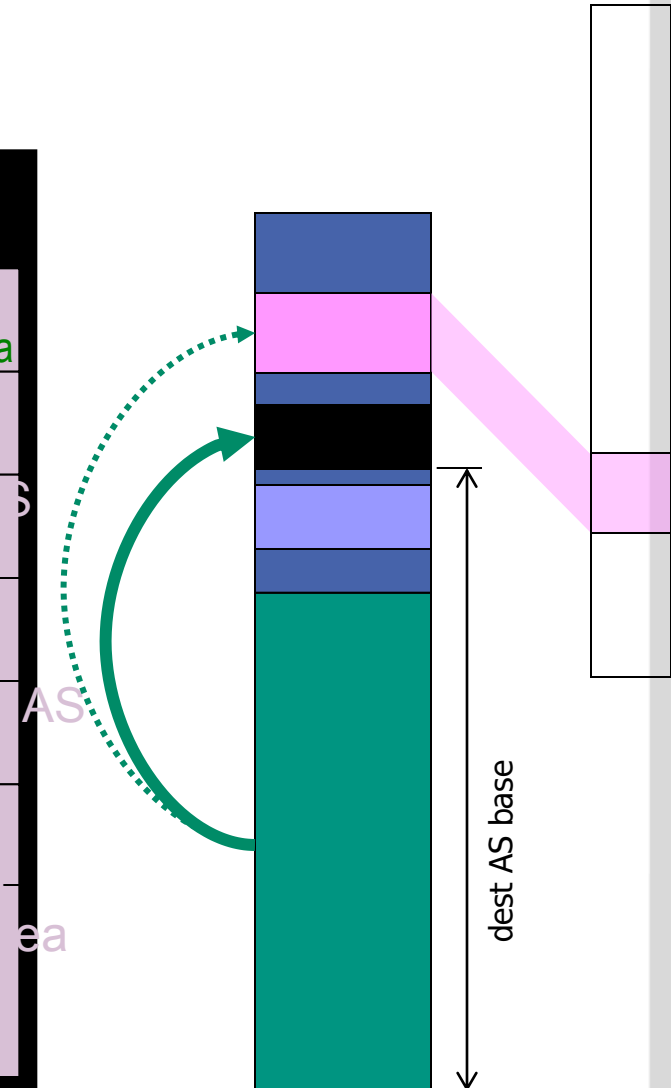
# Long-IPC Implementation Revisited



# Long-IPC Implementation Revisited

Source / Destination	Method	Source offset	Dest
large / large	temp mapping <b>switch HwAS</b>	0	kernel com area
large / small	direct <b>- no HwAS switch -</b>	0	dest AS base

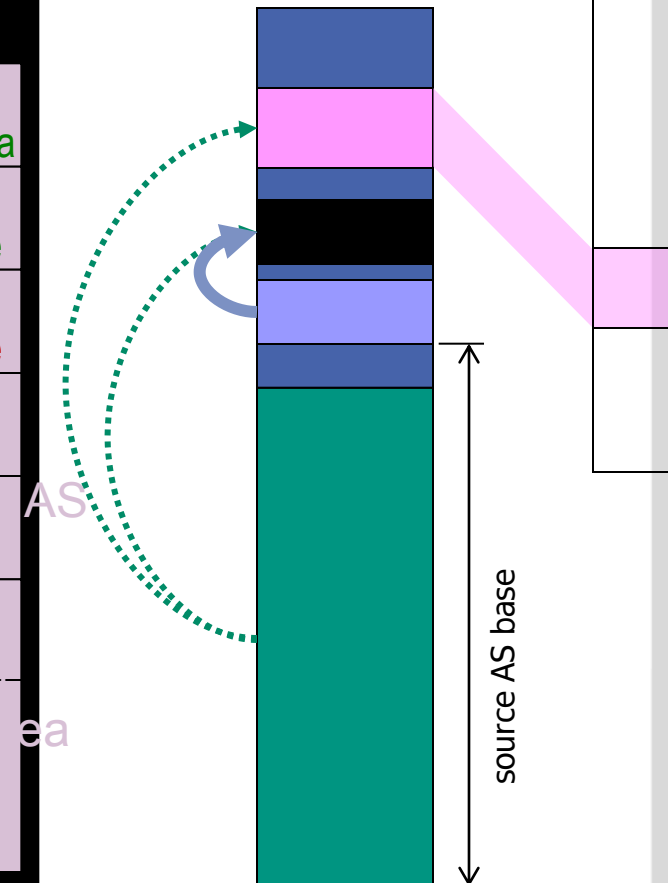
*then -direct-*



# Long-IPC Implementation Revisited

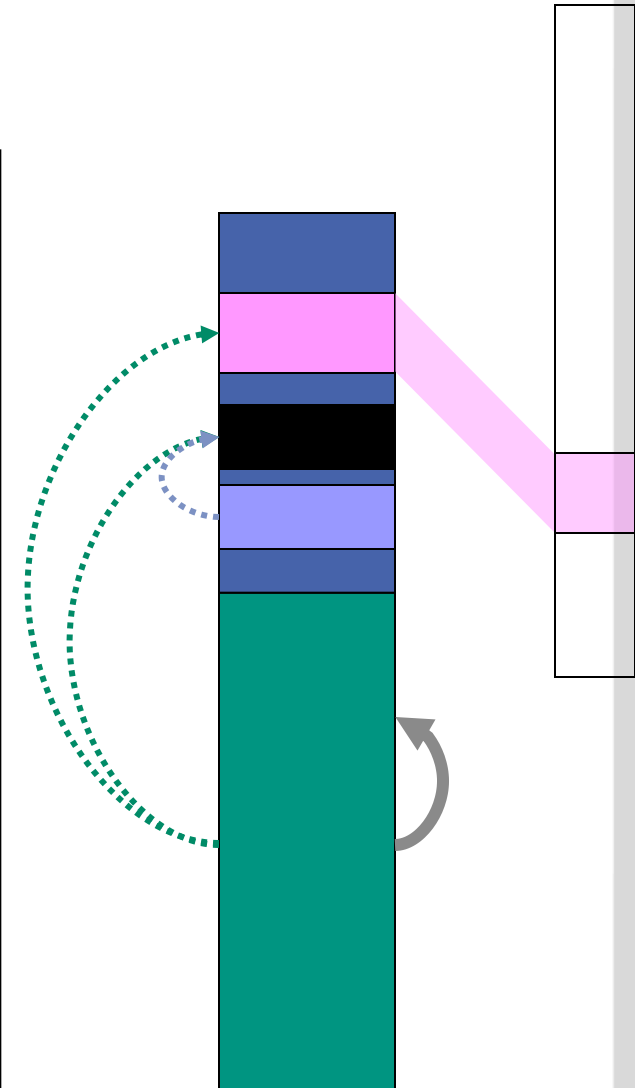
Source / Destination	Method	Source offset	Dest
large / large	temp mapping <b>switch HwAS</b>	0	kernel com area
large / small	direct <b>- no HwAS switch -</b>	0	dest AS base
small / small	direct <b>- no HwAS switch -</b>	source AS base	dest AS base

*then -direct-*



# Long-IPC Implementation Revisited

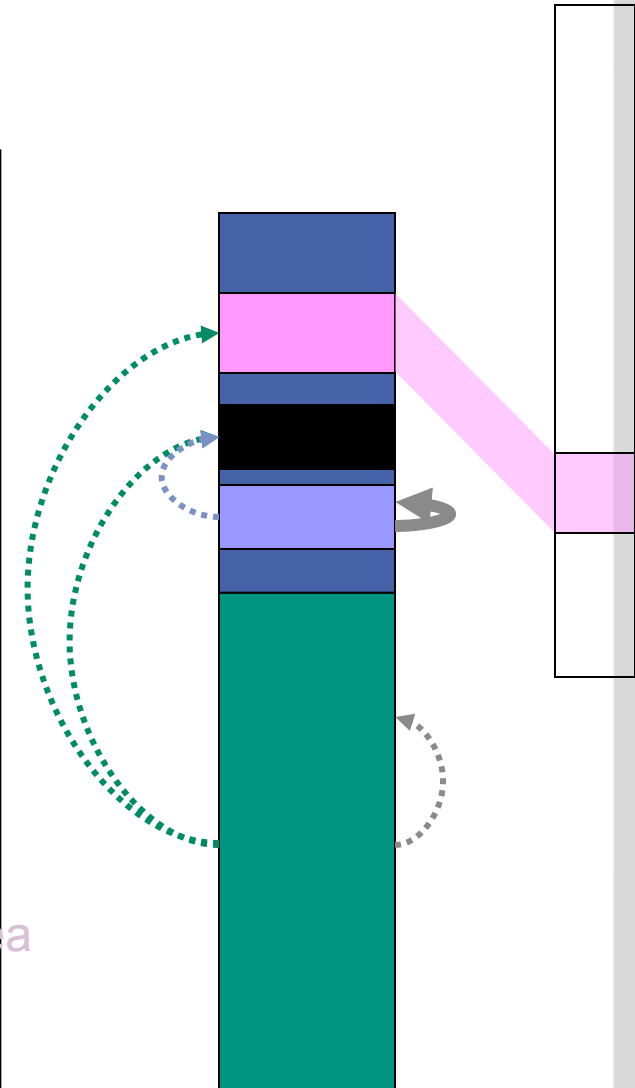
Source / Destination	Method	Source offset	Dest
large / large	temp mapping <b>switch HwAS</b>	0	kernel com area
large / small	direct <b>- no HwAS switch -</b>	0	dest AS base
small / small	direct <b>- no HwAS switch -</b>	source AS base	dest AS base
large / same	direct <b>- no HwAS switch -</b>	0	0



# Long-IPC Implementation Revisited

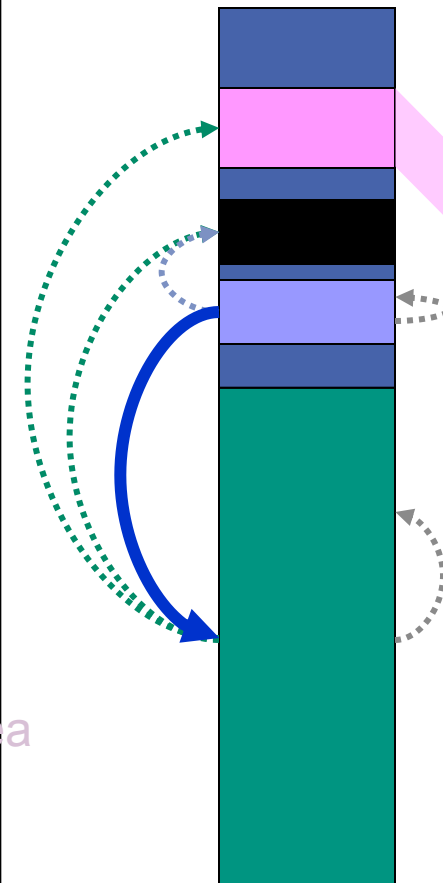
Source / Destination	Method	Source offset	Dest
large / large	temp mapping <b>switch HwAS</b>	0	kernel com area
large / small	direct <b>- no HwAS switch -</b>	0	dest AS base
small / small	direct <b>- no HwAS switch -</b>	source AS base	dest AS base
large / same	direct <b>- no HwAS switch -</b>	0	0
small / same	direct <b>- no HwAS switch -</b>	source AS base	source AS base

*then -direct-*



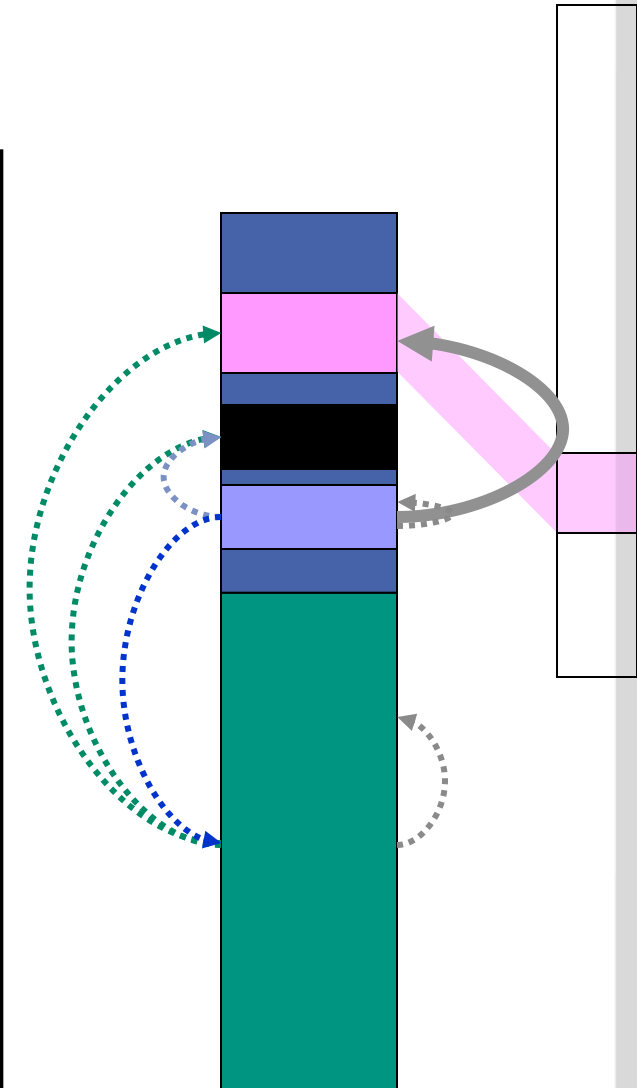
# Long-IPC Implementation Revisited

Source / Destination	Method	Source offset	Dest
large / large	temp mapping <b>switch HwAS</b>	0	kernel com area
large / small	direct <b>- no HwAS switch -</b>	0	dest AS base
small / small	direct <b>- no HwAS switch -</b>	source AS base	dest AS base
large / same	direct <b>- no HwAS switch -</b>	0	0
small / same	direct <b>- no HwAS switch -</b>	source AS base	source AS base
small / large = Current HwAS	direct <b>- no HwAS switch -</b>	source AS base	0



# Long-IPC Implementation Revisited

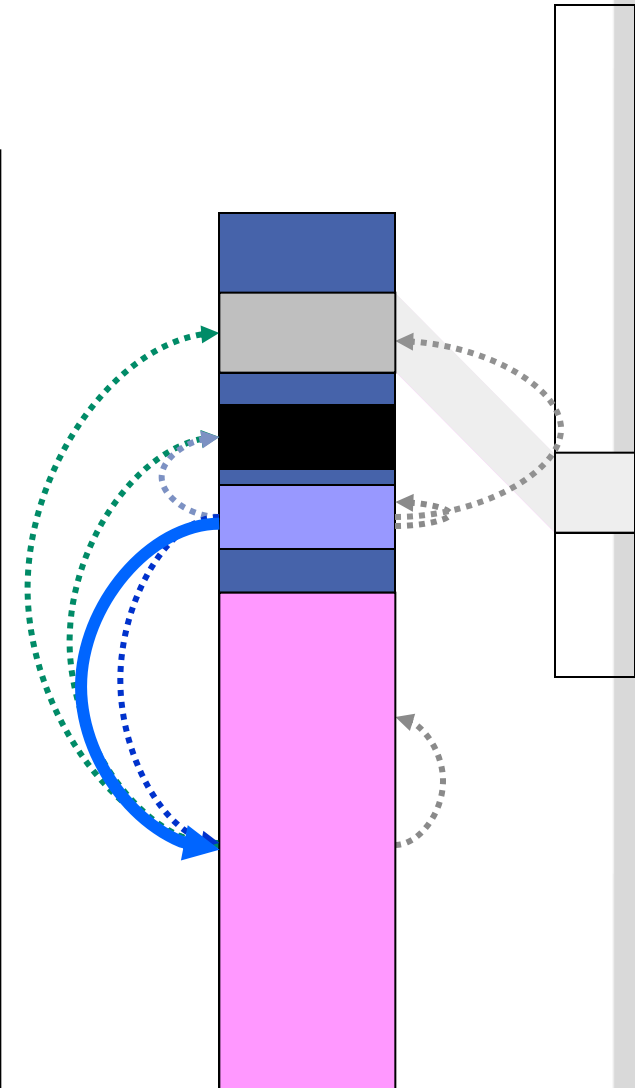
Source / Destination	Method	Source offset	Dest
large / large	temp mapping <b>switch HwAS</b>	0	kernel com area
large / small	direct <b>- no HwAS switch -</b>	0	dest AS base
small / small	direct <b>- no HwAS switch -</b>	source AS base	dest AS base
large / same	direct <b>- no HwAS switch -</b>	0	0
small / same	direct <b>- no HwAS switch -</b>	source AS base	Source AS base
small / large = Current HwAS	direct <b>- no HwAS switch -</b>	Source AS base	0
small / large $\neq$ Current HwAS	temp mapping <b>switch HwAS</b>	source AS base	kernel com area





# Long-IPC Implementation Revisited

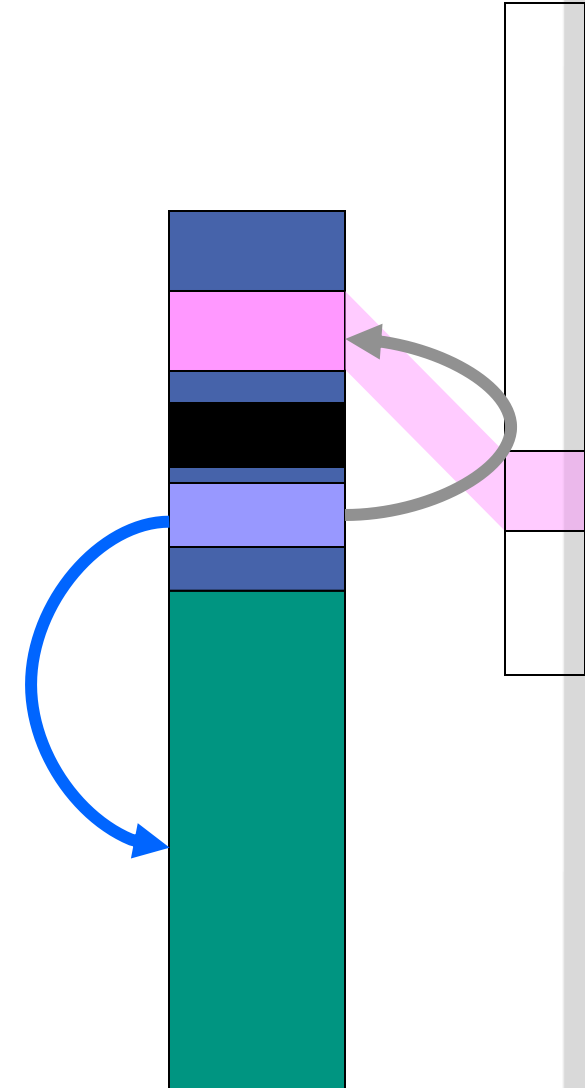
Source / Destination	Method	Source offset	Dest
large / large	temp mapping <b>switch HwAS</b>	0	kernel com area
large / small	direct <b>- no HwAS switch -</b>	0	dest AS base
small / small	direct <b>- no HwAS switch -</b>	source AS base	dest AS base
large / same	direct <b>- no HwAS switch -</b>	0	0
small / same	direct <b>- no HwAS switch -</b>	source AS base	source AS base
small / large = Current HwAS	direct <b>- no HwAS switch -</b>	source AS base	0
small / large $\neq$ Current HwAS	temp mapping <b>switch HwAS</b> – or – <b>first switch HwAS then direct</b>	source AS base	kernel com area – or – 0



# Long-IPC Implementation Revisited

- A **global bit** in page table entry indicates that TLB entry should not be flushed on TLB flushes
  - Used for not flushing kernel entries

Source / Destination	Method	Source	Dest
large / large			
large / small			
small / small			
large / same			
small / same	- no HwAS sw	base	base
small / large = Current HwAS	direct	source	0
	- no HwAS sw	base	
small / large $\neq$ Current HwAS	temp mapping switch HwAS – or – first switch HwAS then direct	source AS base	kernel com area – or – 0

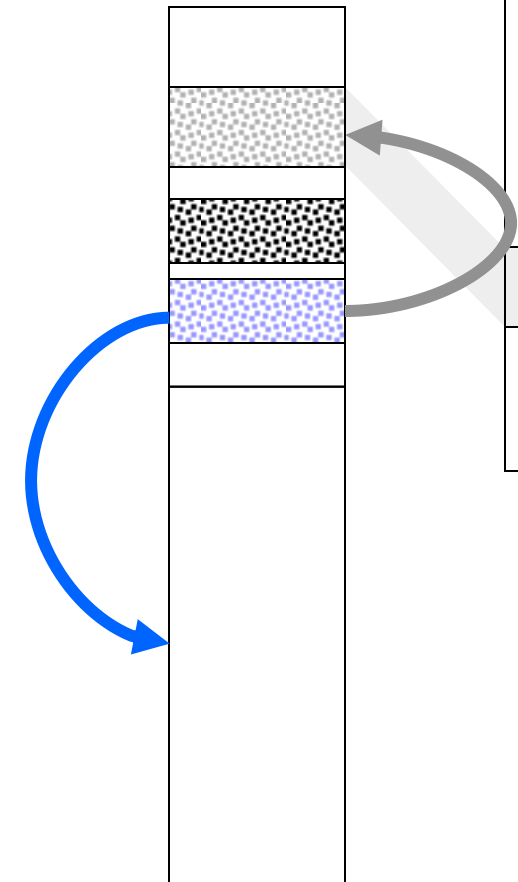


# Long-IPC Implementation Revisited

- A **global bit** in page table entry indicates that TLB entry should not be flushed on TLB flushes
  - Used for not flushing kernel entries

No “global bit”:  
All TLB entries flushed

Source / Destination	Method	Source	Dest
large / large			
large / small			
small / small			
large / same			
small / same			
small / large = Current HwAS	- no HwAS sw direct - no HwAS sw	base source base	base 0
small / large $\neq$ Current HwAS	temp mapping switch HwAS – or – first switch HwAS then direct	source AS base	kernel com area – or – 0

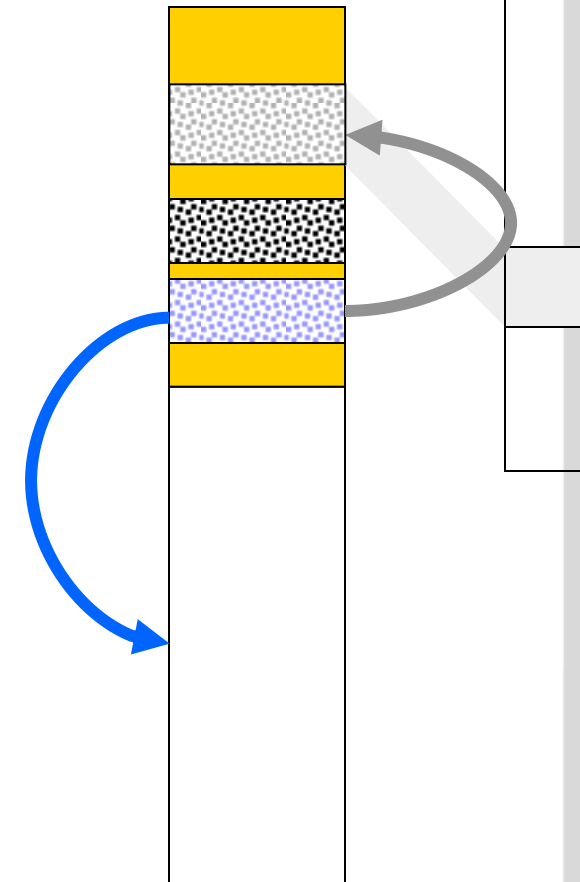


# Long-IPC Implementation Revisited

With “global bit”:  
Non-global TLB entries flushed

- A **global bit** in page table entry indicates that TLB entry should not be flushed on TLB flushes
  - Used for not flushing kernel entries

Source / Destination	Method	Source	Dest
large / large			
large / small			
small / small			
large / same			
small / same			
small / large = Current HwAS	- no HwAS sw direct - no HwAS sw	base source base	base 0
small / large $\neq$ Current HwAS	temp mapping switch HwAS – or – first switch HwAS then direct	source AS base	kernel com area – or – 0

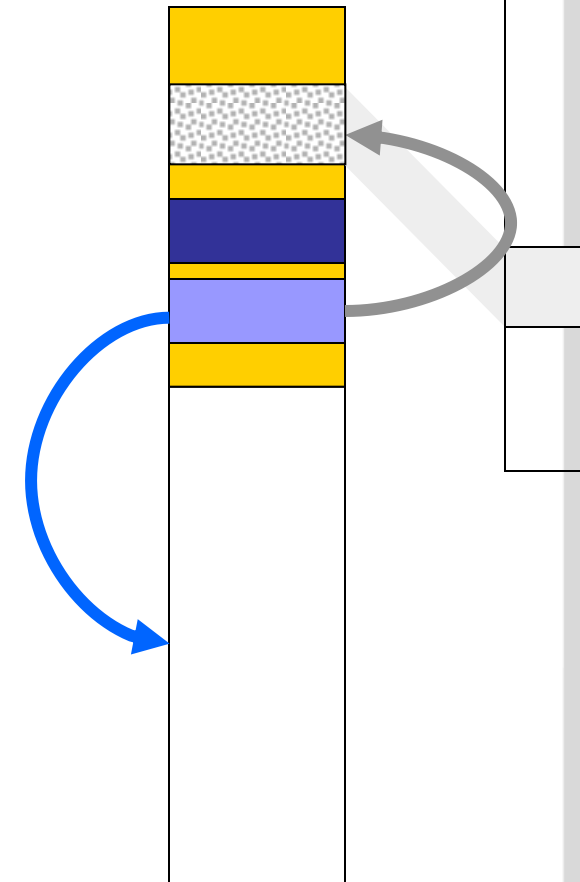


# Long-IPC Implementation Revisited

With “global bit”:  
Non-global TLB entries flushed

- A **global bit** in page table entry indicates that TLB entry should not be flushed on TLB flushes
  - Used for not flushing kernel entries
- small spaces are global over all hardware address spaces
  - Mark small spaces' PTEs as **global**

Source / Destination	Method	Source	Dest
large / large			
large / small			
small / small			
large / same			
small / same			
small / large = Current HwAS	- no HwAS sw direct - no HwAS sw	base source base	base 0
small / large $\neq$ Current HwAS	temp mapping switch HwAS - or - first switch HwAS then direct	source AS base	kernel com area - or - 0



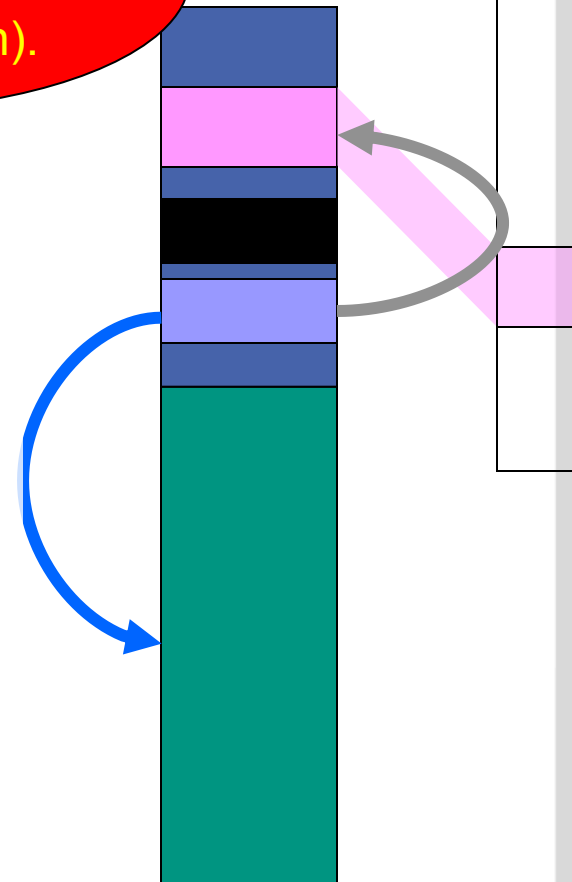
# Long-IPC Implementation Revisited

Source / Destination	Method		
large / large	temp mapping switch HwAS	0	com area
large / small	direct no HwAS switch -	0	dest AS base
small / small	temp mapping switch HwAS	source AS base	dest AS base
small / large = Current HwAS	temp mapping switch HwAS	0	0
small / large $\neq$ Current HwAS	temp mapping switch HwAS - or - first switch HwAS then direct	source AS base	kernel com area - or - 0

With “global PTE entries”,  
TLB misses only in dest space.

Without global pages,  
TLB misses in source and dest.

TLB misses on com area  
and in dest space  
(after the switch).



# Temporary Mapping Revisited

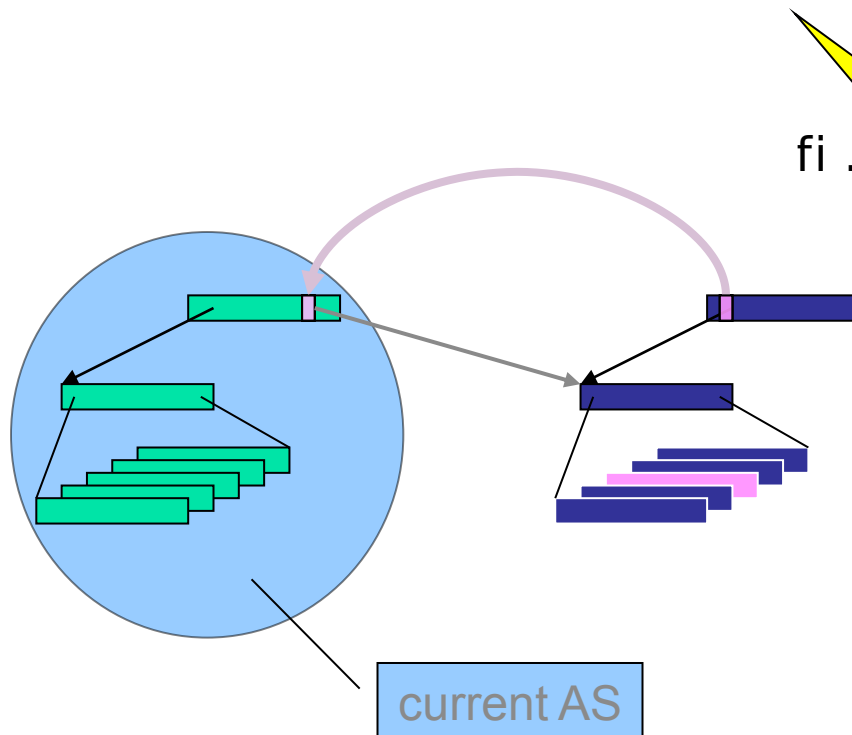
Leave thread:

if mytcb.partner  $\neq$  nilthread then

myPDE.TMarea := nil ;

if dest AS = my AS then

flush TLB



Optimization: avoids second TLB flush if subsequent thread and AS switch would flush TLB anyway

## Evicting the Temporary Mapping Area

- We must evict the temp. mapping from the TLB when switching from ... to ...
- Depending on whether small spaces use the temp. mapping / or not at all

from \ to	small AS	Active large AS	Inactive large AS
small AS	Yes / No	Yes / No <sup>1</sup>	No / No
large AS	Yes / Yes <sup>2</sup>	Yes	No

<sup>1</sup>: assuming the temp. mapping is invalidated when switching to the small space

<sup>2</sup>: to prevent T3 from using T1's mapping area after T1 (large AS A) → T2 (small AS B) → T3 (large AS A), possible due to <sup>1</sup>



# Temporary Mapping Revisited

Leave thread:

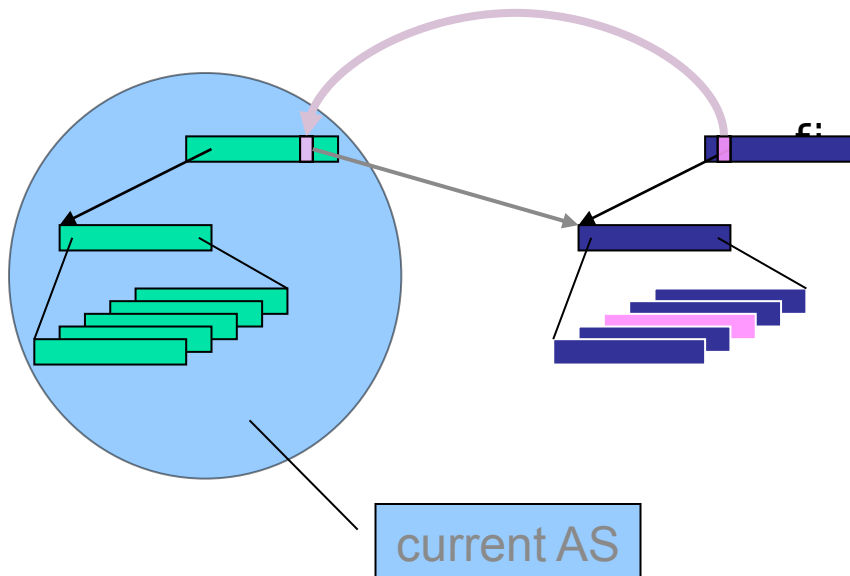
```
if mytcb.partner ≠ nilthread then
  myPDE.TMarea := tm,
```

Assuming small spaces use the temp. mapping area.

```
if (dest is small or
    dest HwAS = curr HwAS)
then
```

```
  flush TLB
```

```
fi
```



# Temporary Mapping Revisited

Leave thread:

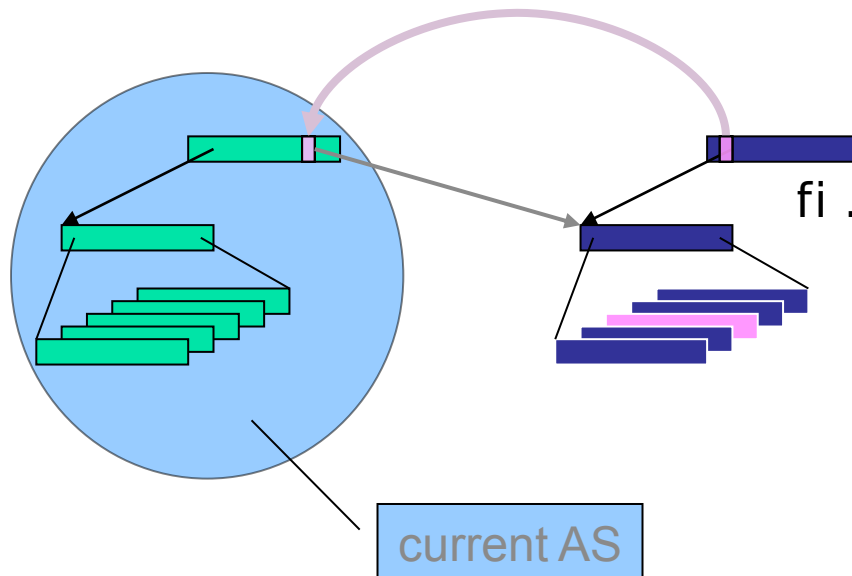
```
if mytcb.partner ≠ nilthread then
  myPDE.TMarea := tm ,
```

Assuming small spaces never use the temp. mapping area.

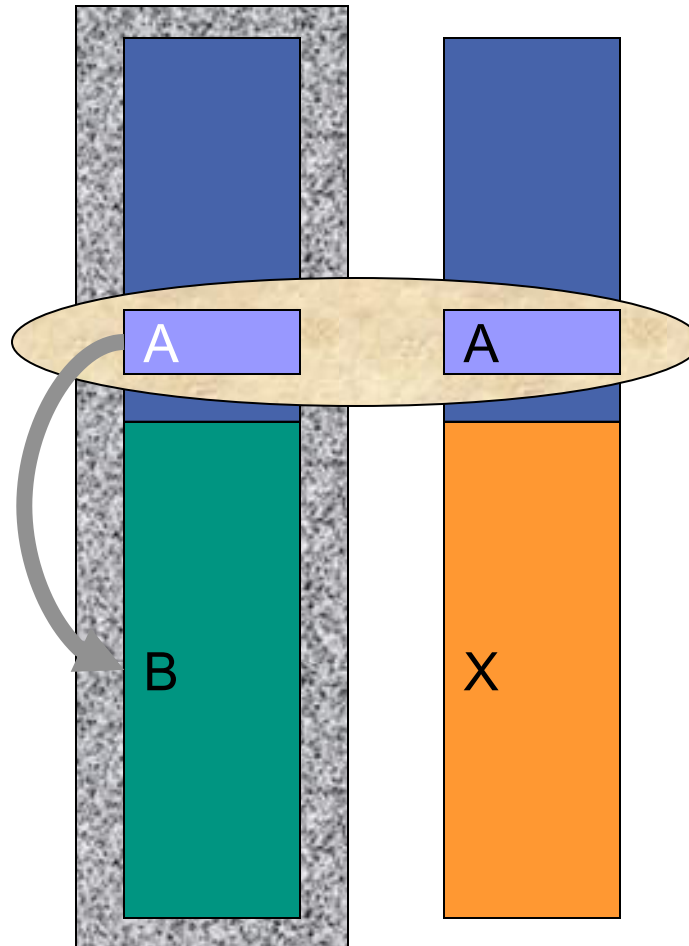
```
if (dest is small or
    dest HwAS = curr HwAS) and
    not curr is small
then
```

```
flush TLB
```

```
fi
```

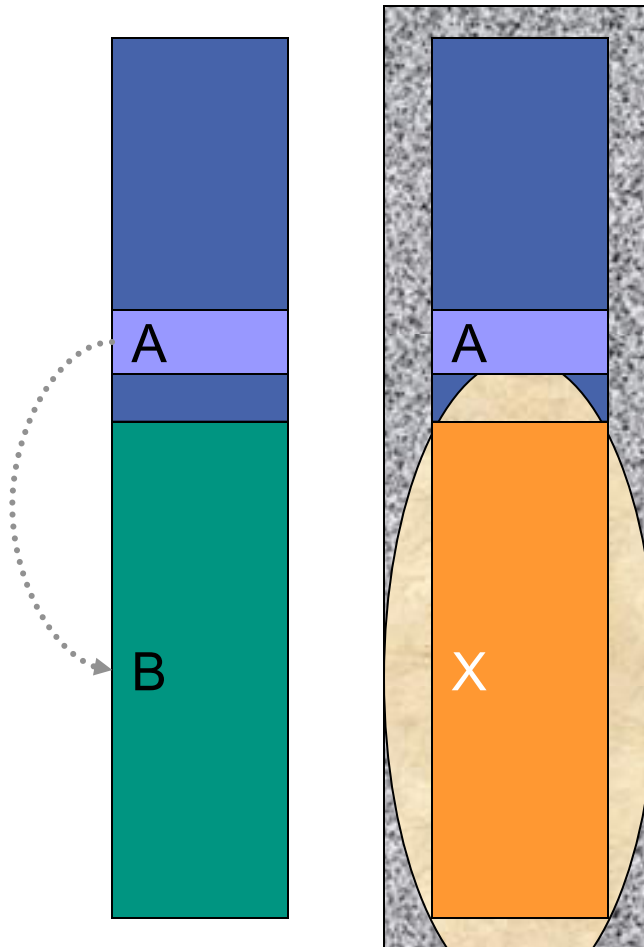


# Thread Switching Revisited



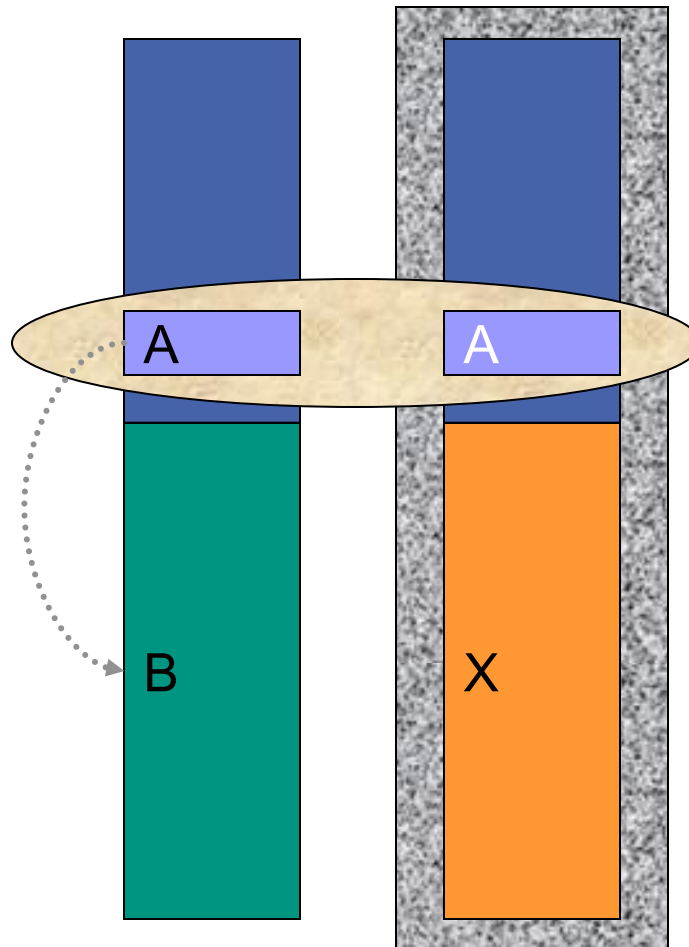
- **A** sends to **B**, executes in HwAS **B**
  - **A** is preempted or PF in **A**
  - Thread switch from **A** to **X**

# Thread Switching Revisited



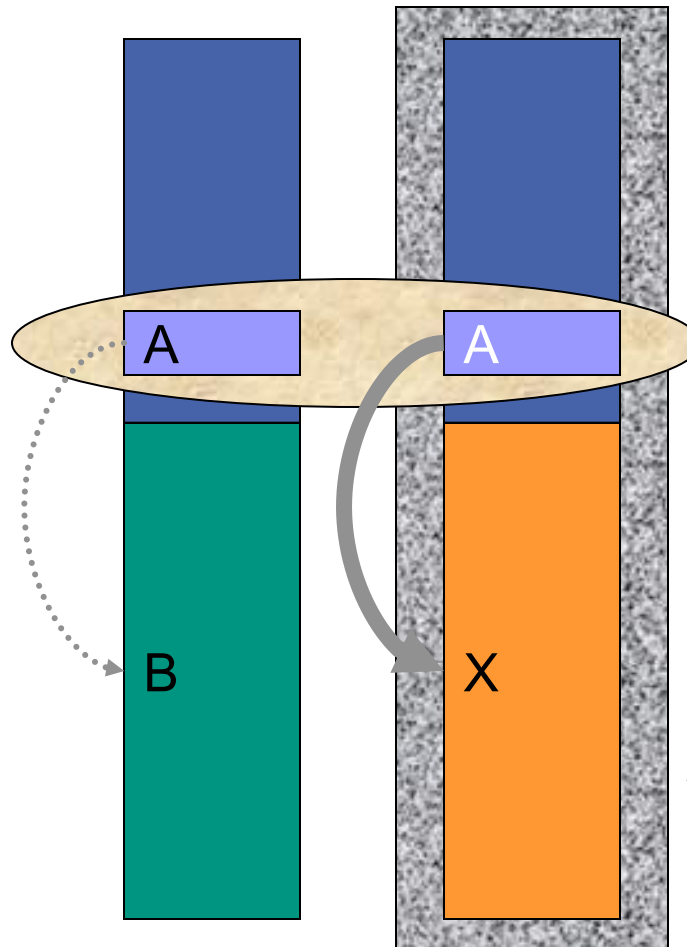
- **A** sends to **B**, executes in HwAS **B**
  - **A** is preempted or PF in **A**
  - Thread switch from **A** to **X**
  - Switch back from **X** to **A**

# Thread Switching Revisited



- **A** sends to **B**, executes in HwAS **B**
  - **A** is preempted or PF in **A**
  - Thread switch from **A** to **X**
  - Switch back from **X** to **A**

# Thread Switching Revisited

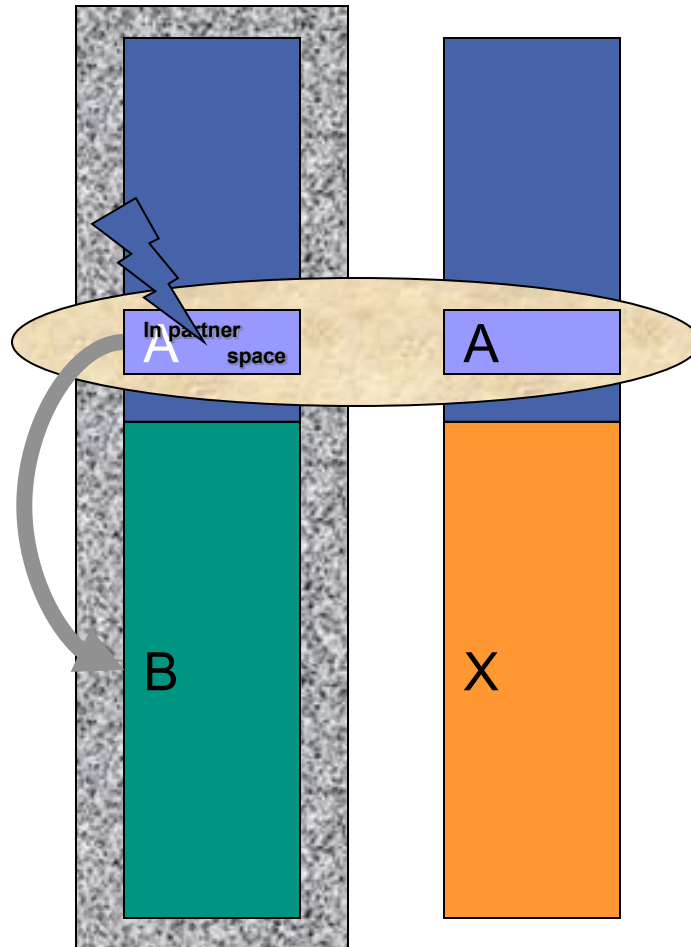


- **A** sends to **B**, executes in HwAS **B**
  - **A** is preempted or PF in **A**
  - Thread switch from **A** to **X**
  - Switch back from **X** to **A**
- Preemption or PF resumes in **A**, but **A** now executes in HwAS **X**



# Thread Switching Revisited

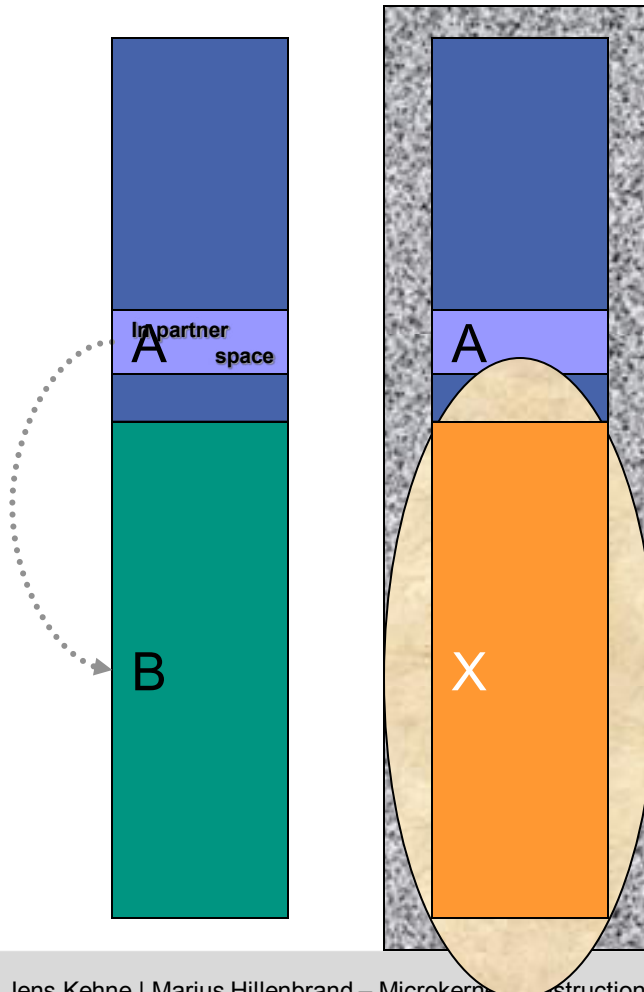
## A Solution



- **A** sends to **B**, executes in HwAS **B**
  - **A** is preempted or PF in **A**
    - Mark **A** “in partner space”

# Thread Switching Revisited

## A Solution

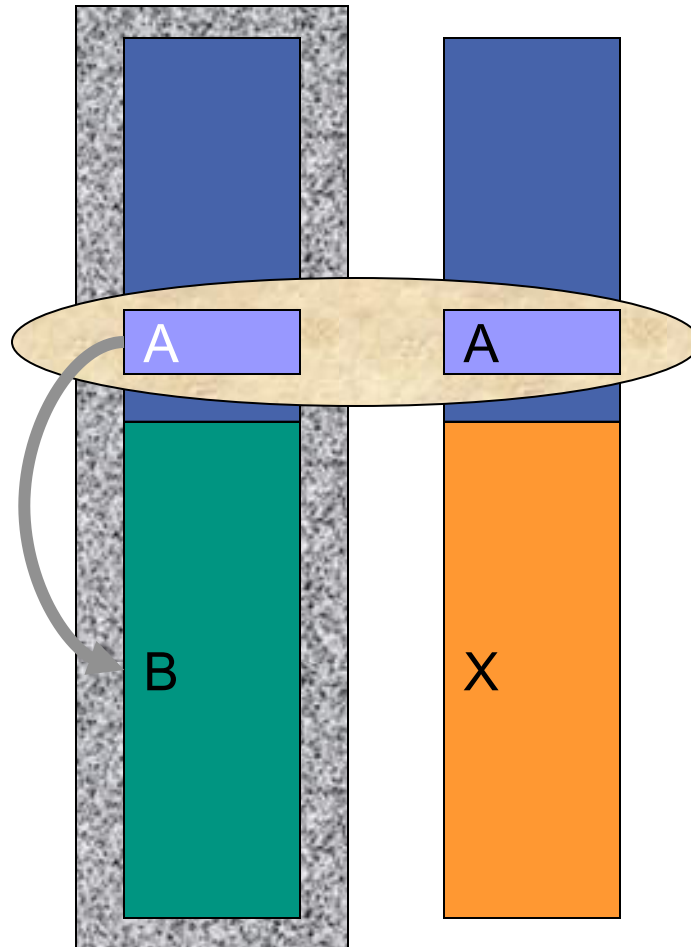


- **A** sends to **B**, executes in HwAS **B**
  - **A** is preempted or PF in **A**
    - Mark **A** “in partner space”
  - Thread switch from **A** to **X**



# Thread Switching Revisited

## A Solution



- **A** sends to **B**, executes in HwAS **B**
  - **A** is preempted or PF in **A**
    - Mark **A** “in partner space”
  - Thread switch from **A** to **X**
  - Switch back from **X** to **A**
    - Also switch HwAS to HwAS **B**

# SMALL ADDRESS SPACES AND FAST SYSTEM CALLS

Getting around automatic  
segment register reloading

# Fast System Calls

## ■ Optimized instructions

- sysenter/sysexit (Intel)
- syscall/sysret (AMD)

## ■ Faster than software interrupts

- Can avoid certain checks
- Unconditionally reload segment registers (page based protection)

	450 MHz PIII
Int / Iret	280
Sysenter / Sysexit	50

# Fast System Calls

## ■ Optimized instructions

- sysenter/sysexit (Intel)
- syscall/sysret (AMD)

## ■ Faster than software interrupts

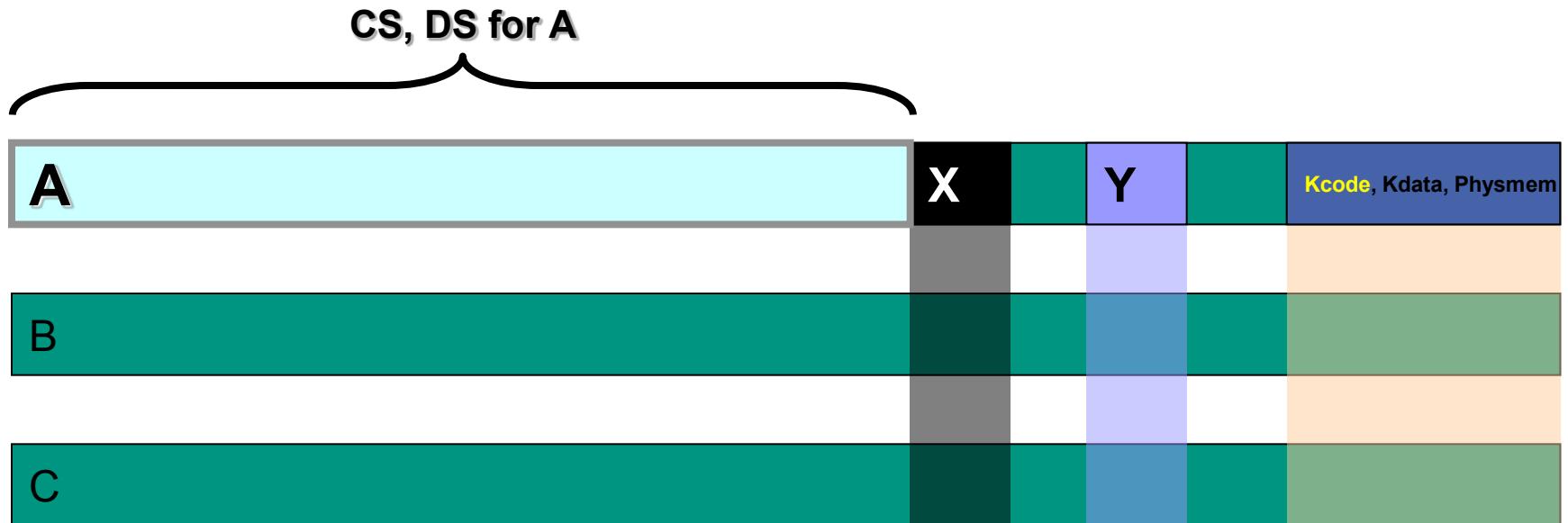
- Can avoid certain checks
- Unconditionally reload segment registers (page based protection)

	450 MHz PIII	1.5 GHz P4
Int / Iret	280	1600
Sysenter / Sysexit	50	140

Problematic for  
small spaces  
(segment based  
protection)

# Fast System Calls

## Automatic Segment Register Reloading



# Fast System Calls

## Automatic Segment Register Reloading

sysenter

kernel CS, DS

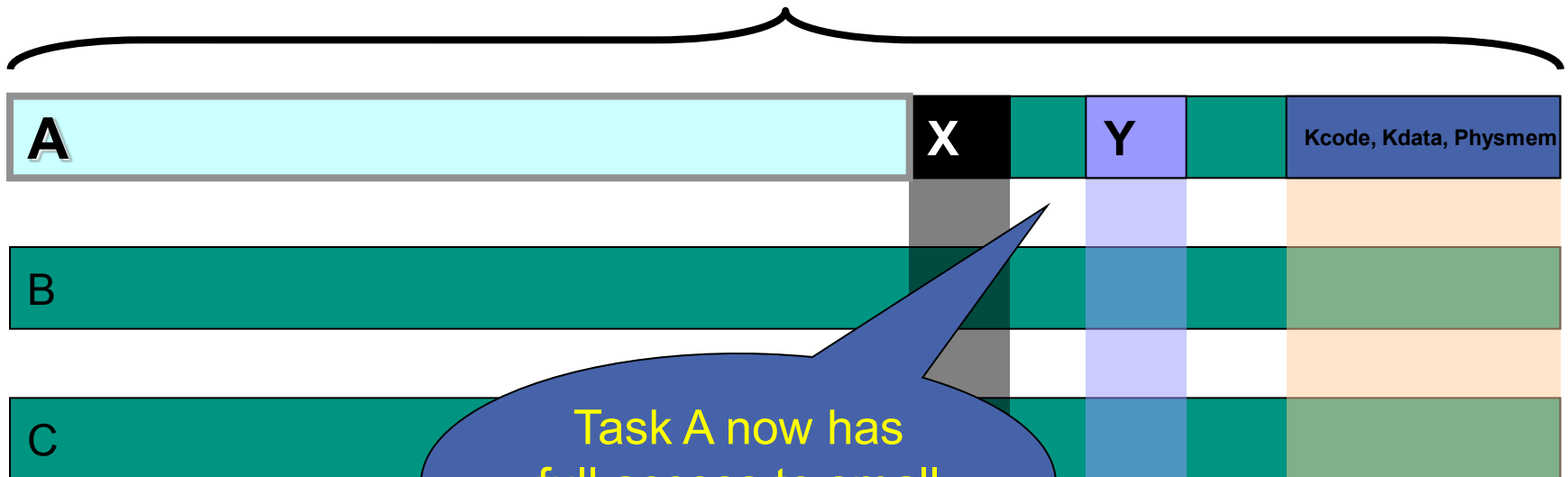


# Fast System Calls

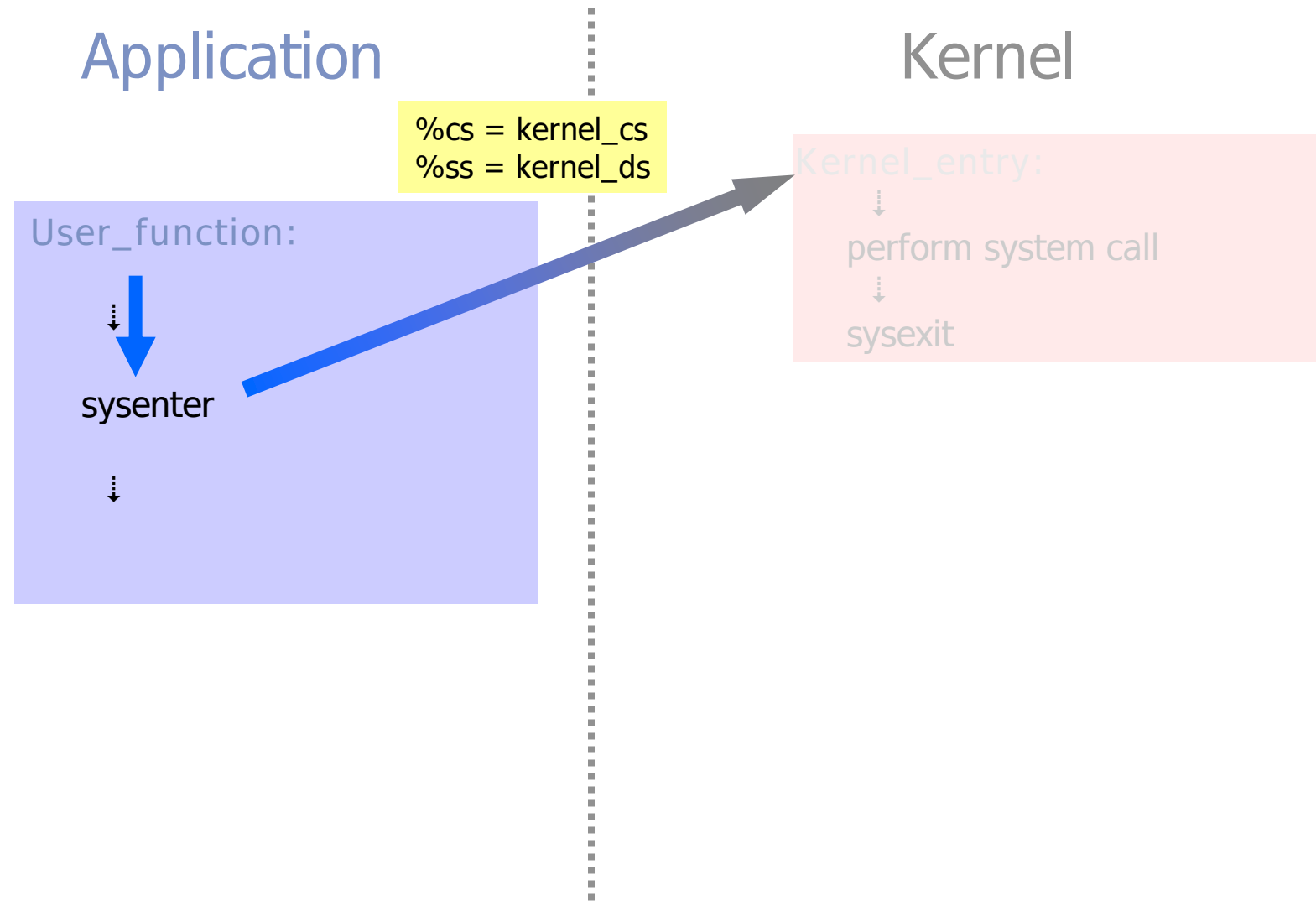
## Automatic Segment Register Reloading

sysexit

user CS, DS

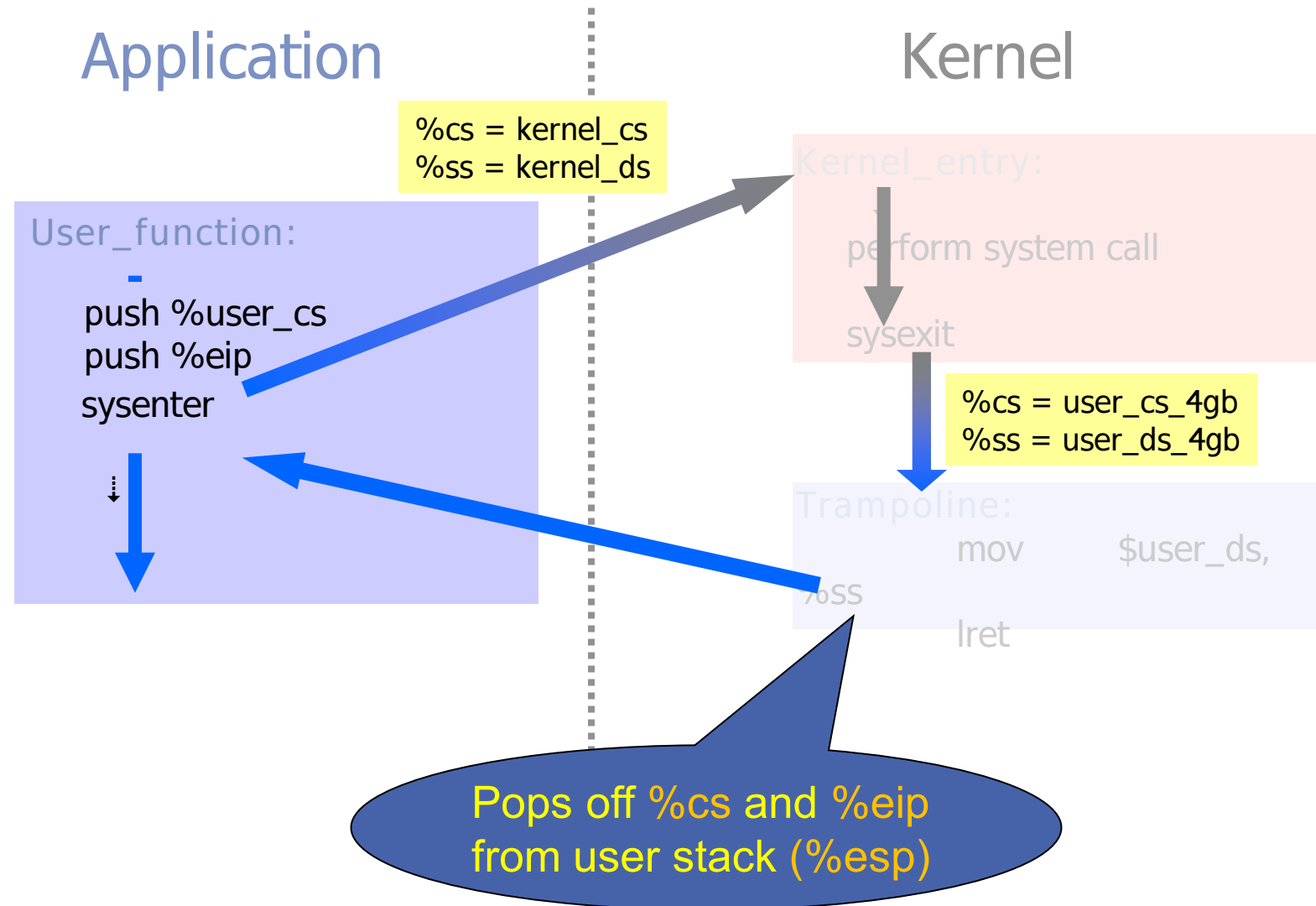


# Automatic Segment Register Reloading Solution – In-Kernel Trampoline

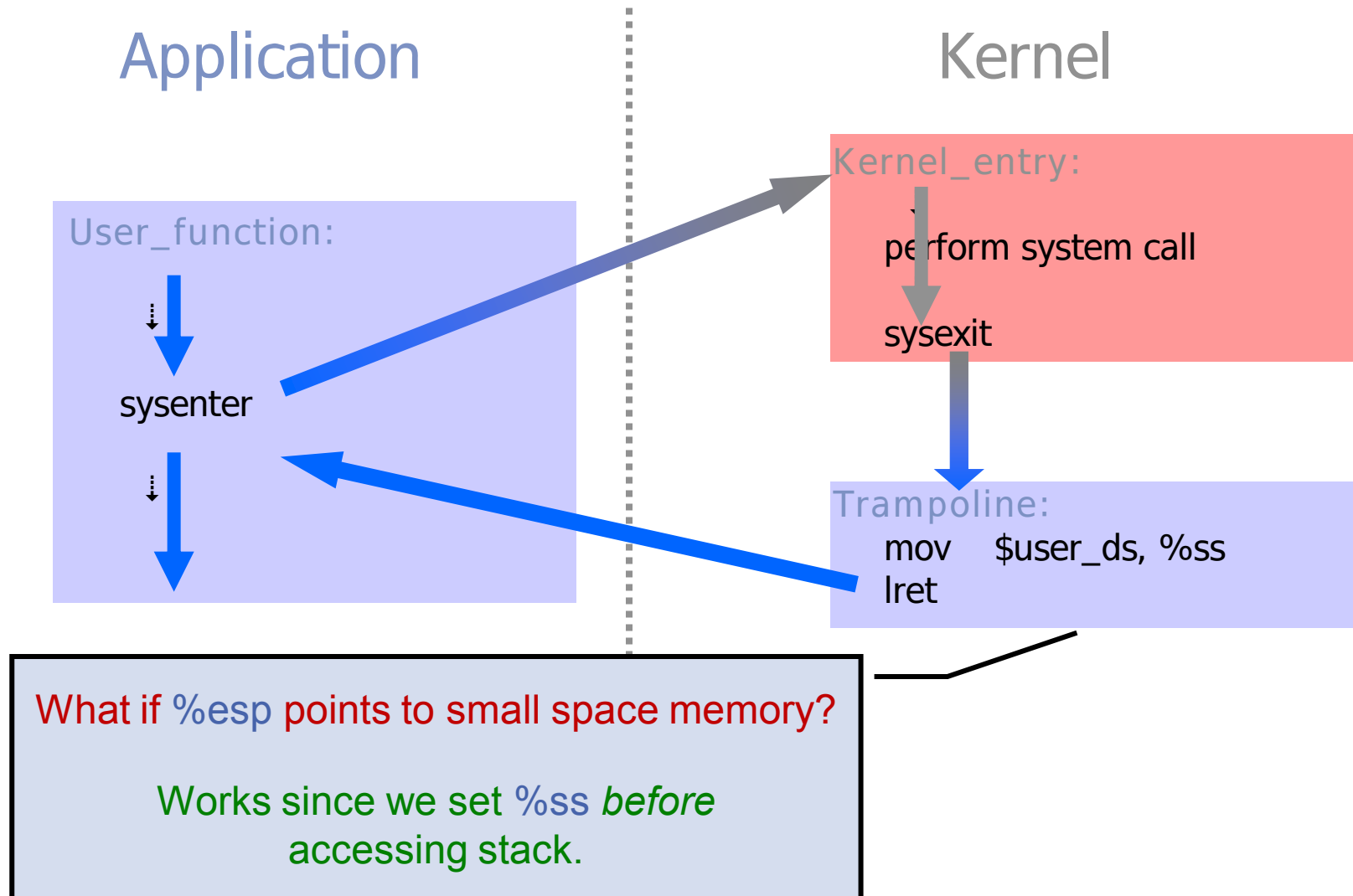




# Automatic Segment Register Reloading Solution – In-Kernel Trampoline



# Automatic Segment Register Reloading Solution – In-Kernel Trampoline



# Automatic Segment Register Reloading Solution – In-Kernel Trampoline

Application

Kernel

User\_function:

Kernel\_entry:

perform system call

sysexit

Trampoline:

mov \$user\_ds, %ss

iret

What if interrupts arrive or  
iret instruction raises exceptions  
(e.g., page-fault on stack access)?

- Exception gets handled
- Return to faulting instruction (iret)
  - Via iret (no nested fast syscall)
  - Restores %cs = user\_cs\_4gb and %ss = user\_ds.
  - Kernel must have set up %[defg]s before!
- Continue to completion