# Multiprocessor Kernel Performance Profiling

**Alex Mirgorodskii**

`mirg@cs.wisc.edu`

Computer Sciences Department

University of Wisconsin

1210 W. Dayton Street

Madison, WI 53706-1685

USA

# Kperfmon: Overview

- Specify a *resource*
  - Almost any function or basic block in the kernel
- Apply a *metric* to the resource:
  - Number of entries to a function or basic block
  - Wall clock time, CPU time (virtual time)
  - All Sparc Hardware Counters: cache misses, branch mispredictions, instructions per cycle, ...
- Visualize the metric data in real time

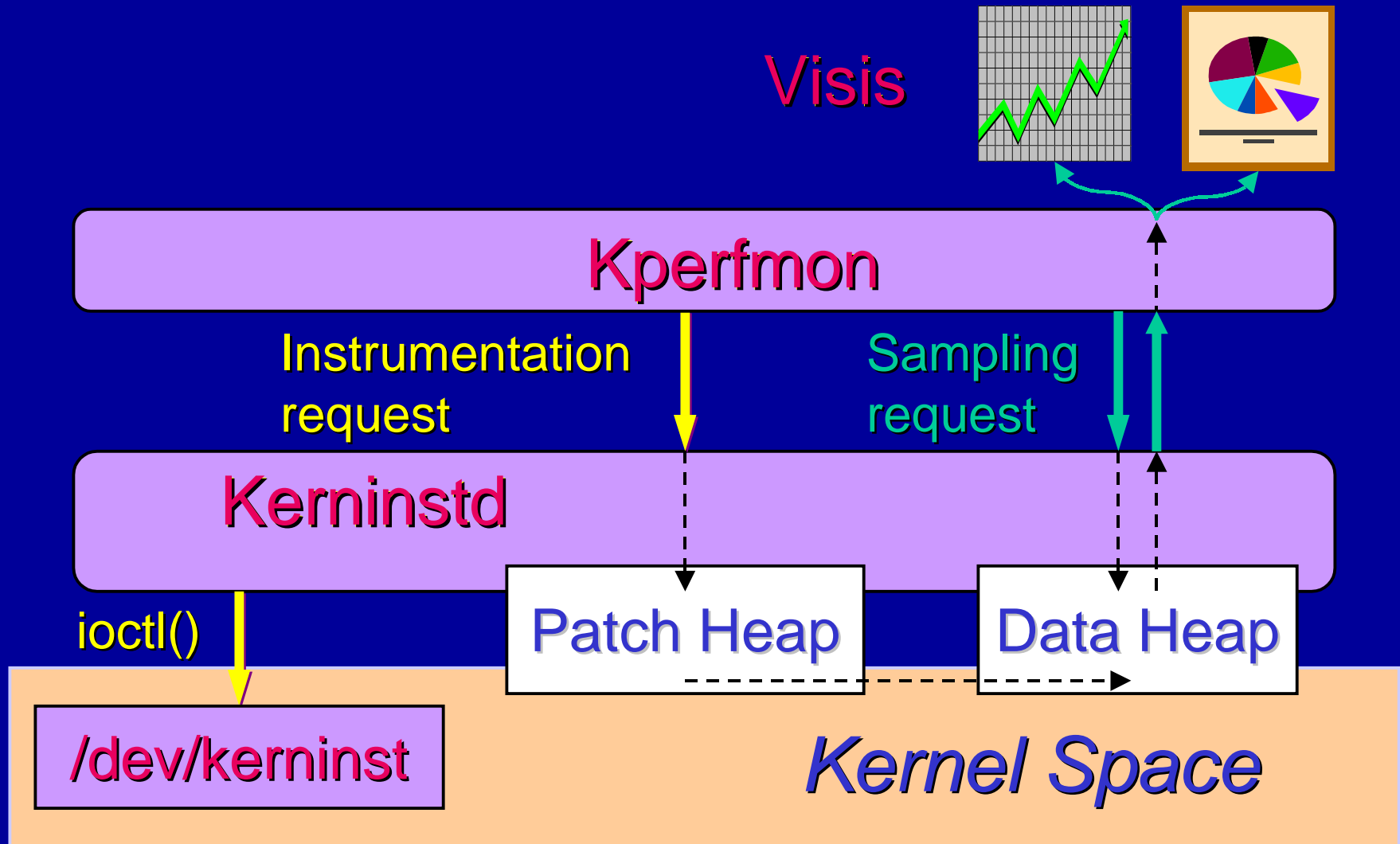# Kperfmon-MP: Goals

Modify uniprocessor Kperfmon to provide:

- Safe operation on SMP machines
    - Thread safety
    - Migration safety
- New feature: Per-CPU performance data
    - More detailed performance data
    - Reduce cache coherence traffic caused by the tool

# Kperfmon: Technology

- Use the *KernInst* framework to:
  - Insert measurement code in the kernel at run time
  - Sample accumulated metric values from the user space periodically

- No need for kernel recompilation
  - Works with stock SPARC Solaris 7 kernels
  - Supports both 32-bit and 64-bit kernels

- No need for rebooting
  - Important for 24 x 7 systems

# Kperfmon System

Visis

Kperfmon

Instrumentation
request

Sampling
request

Kerninstd

ioctl()

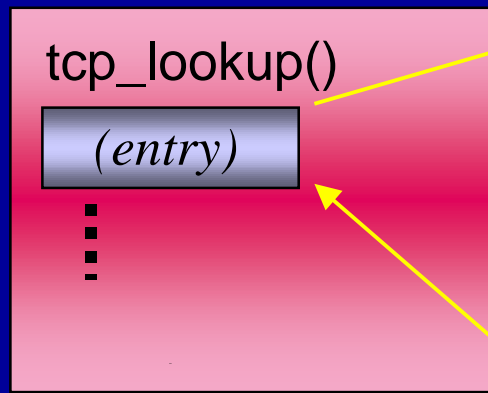Patch Heap

Data Heap

/dev/kerninst

*Kernel Space*

# Kperfmon instrumentation

- Counter primitive
  - Number of entries to a function or a basic block

- Wall clock timer primitive
  - Real time spent in a function

- CPU timer primitive
  - Excludes time while the thread was switched-out
  - Can count more than just timer ticks
    - All HW-counter metrics use this mechanism

# Non-MP Counter primitive

**Code Patch Area**

tcp_lookup()

*(entry)*

```
sethi hi(&cnt), r0
or r0, lo(&cnt), r0
ldx [r0], r1
retry:
    add r1, 1, r2
    casx [r0], r1, r2
cmp r1, r2
bne retry
mov r2, r1
```

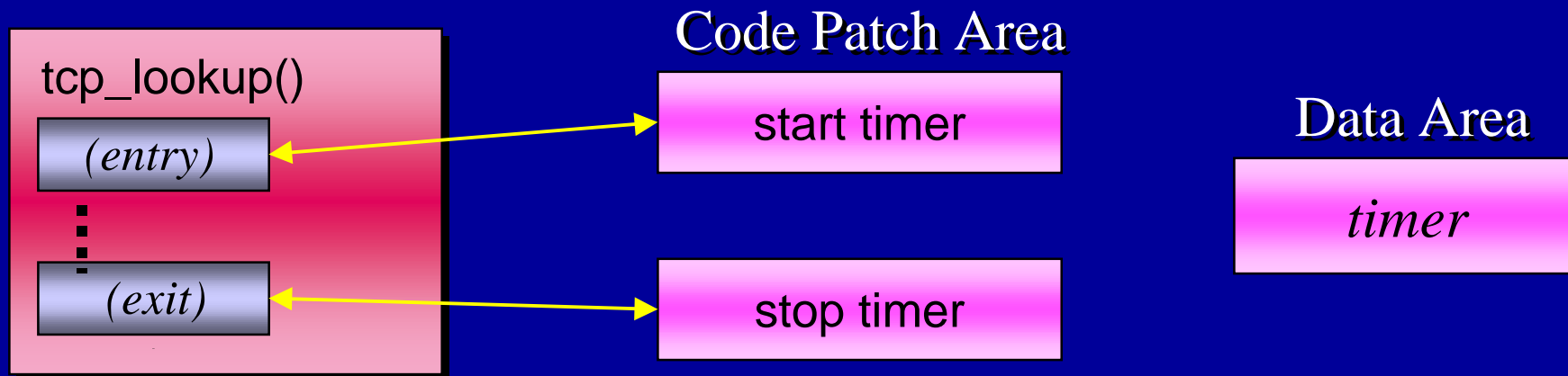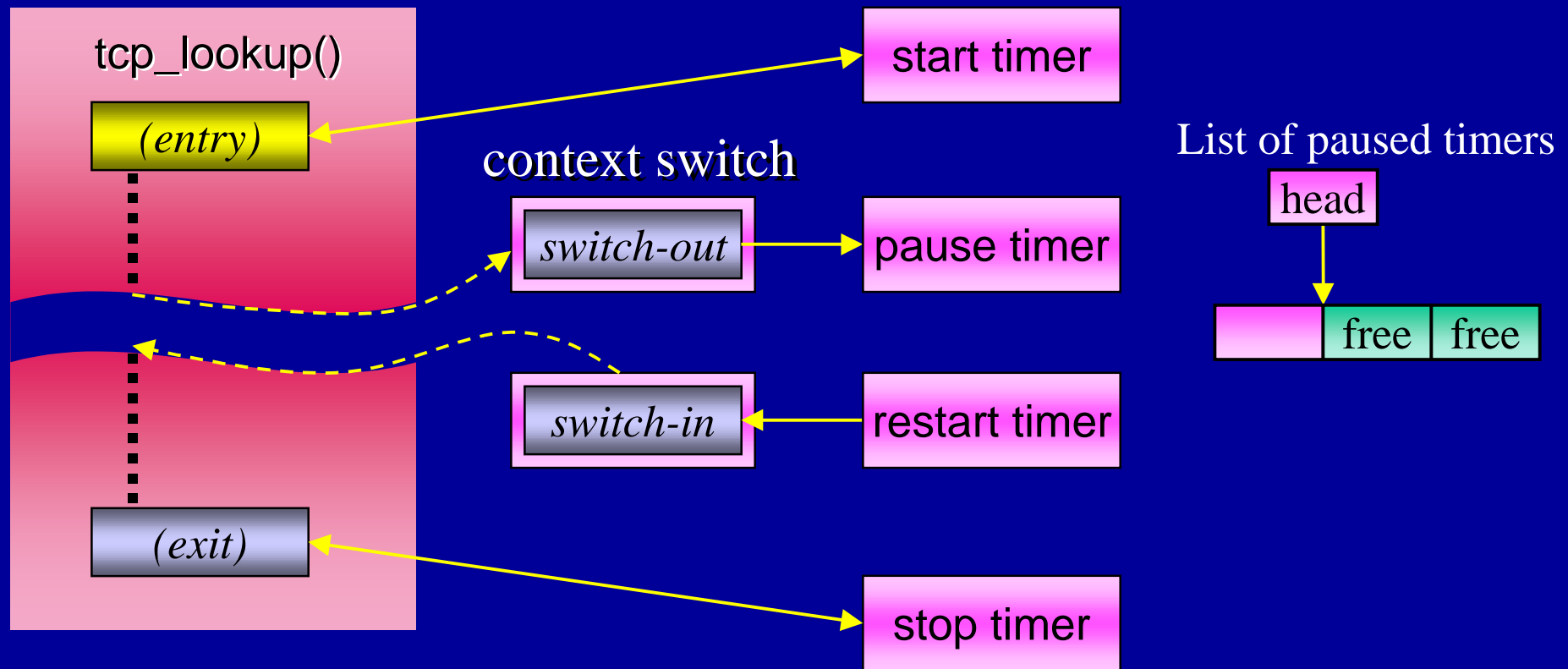*Relocated instruction*

ba,a tcp_lookup+4

**Data Area**

*cnt*

- Atomic, thread-safe update
- Lightweight
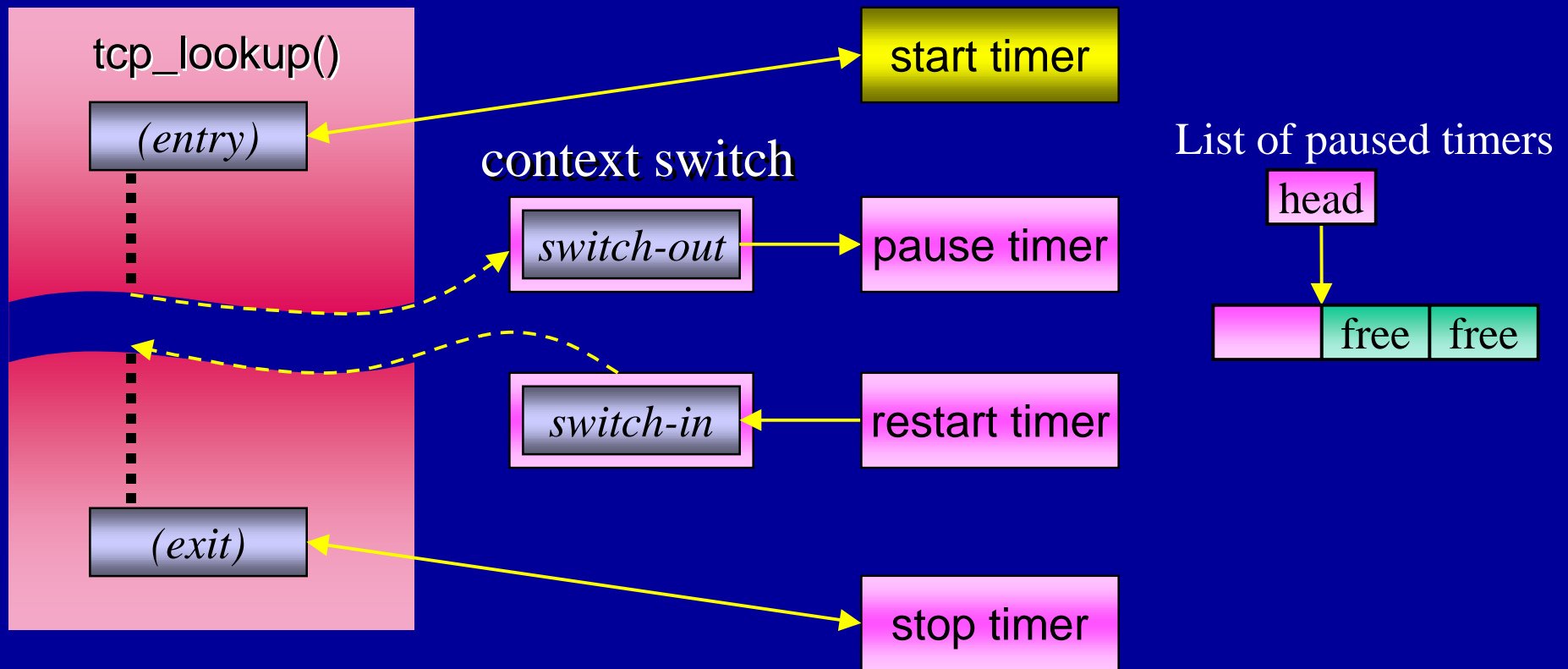- No register save/restore required

# Non-MP Wall clock timer primitive

tcp_lookup()

*(entry)*

*(exit)*

Code Patch Area

start timer

stop timer

Data Area

*timer*

- Inclusive (includes time in callees)
- Keeps accumulating if switched-out

# Non-MP CPU timer primitive

**tcp_lookup()**

*(entry)* → start timer

**context switch**

*switch-out* → pause timer

*switch-in* ← restart timer

*(exit)* → stop timer
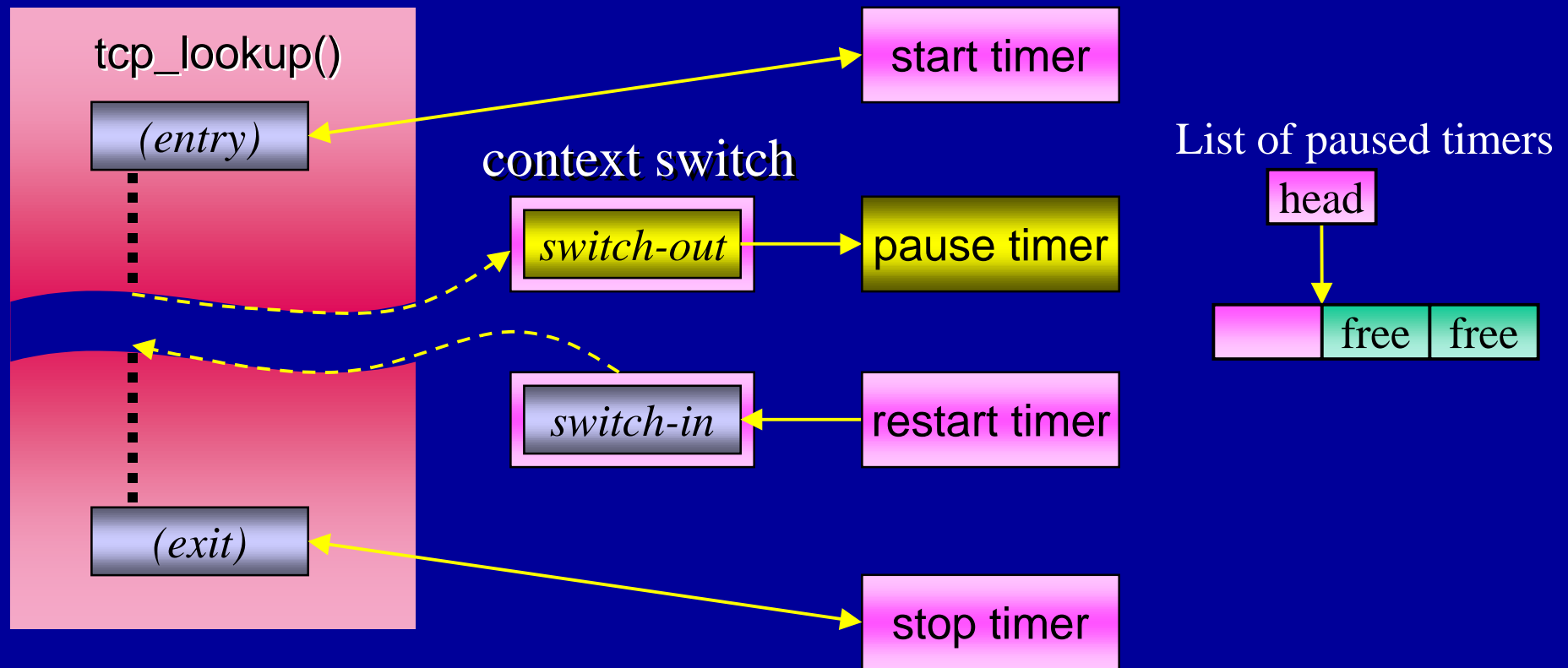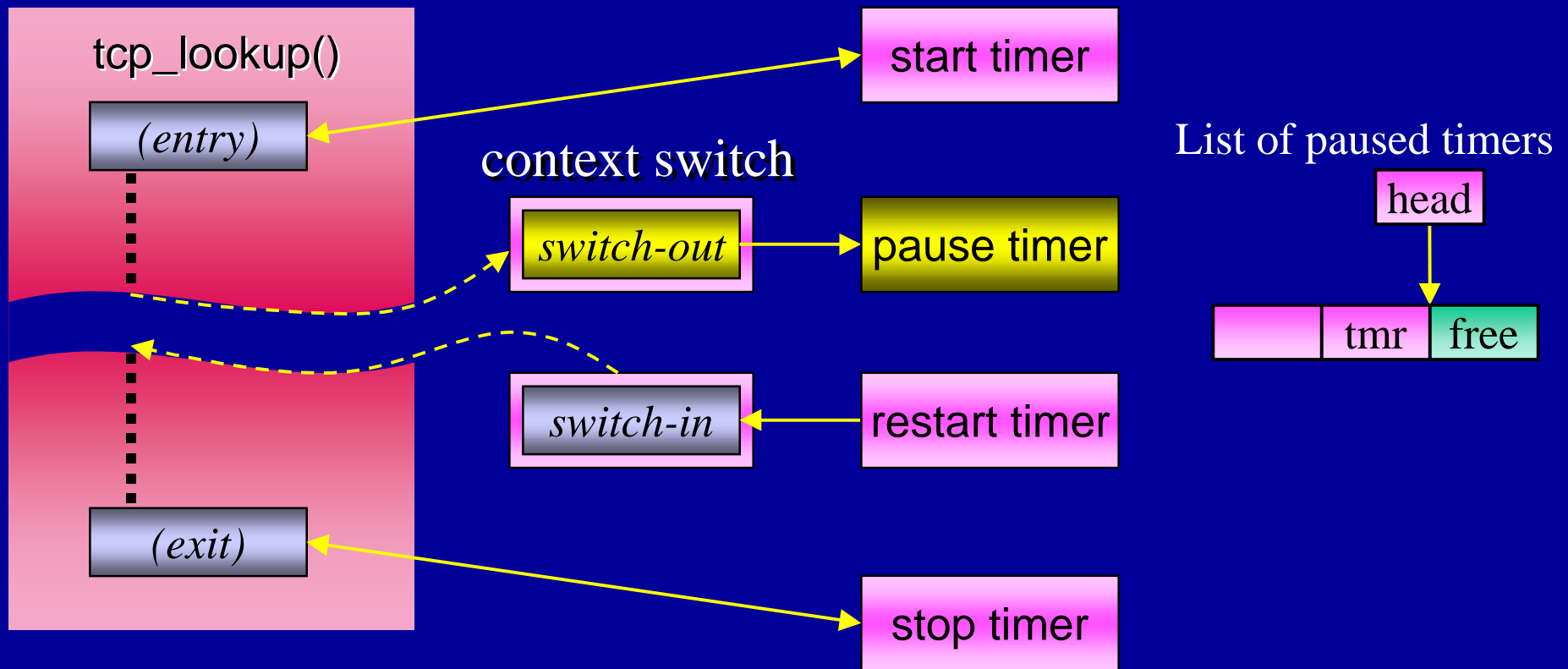
**List of paused timers**

head

free | free

- Exclude the time spent while switched out
  - Instrument context switch routines
- HW counter metrics are based on this mechanism

# Non-MP CPU timer primitive

**tcp_lookup()**

*(entry)*

start timer

**context switch**

*switch-out* → pause timer

*switch-in* ← restart timer

*(exit)* → stop timer
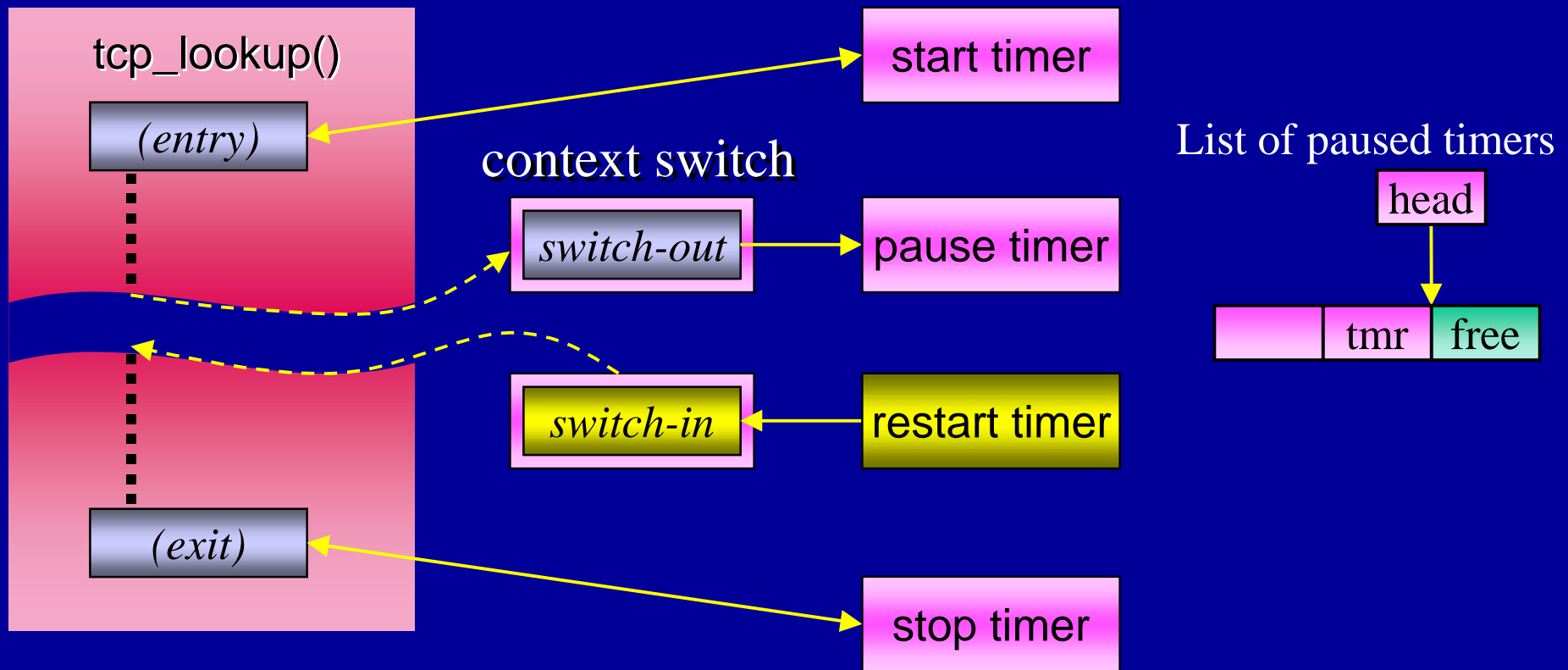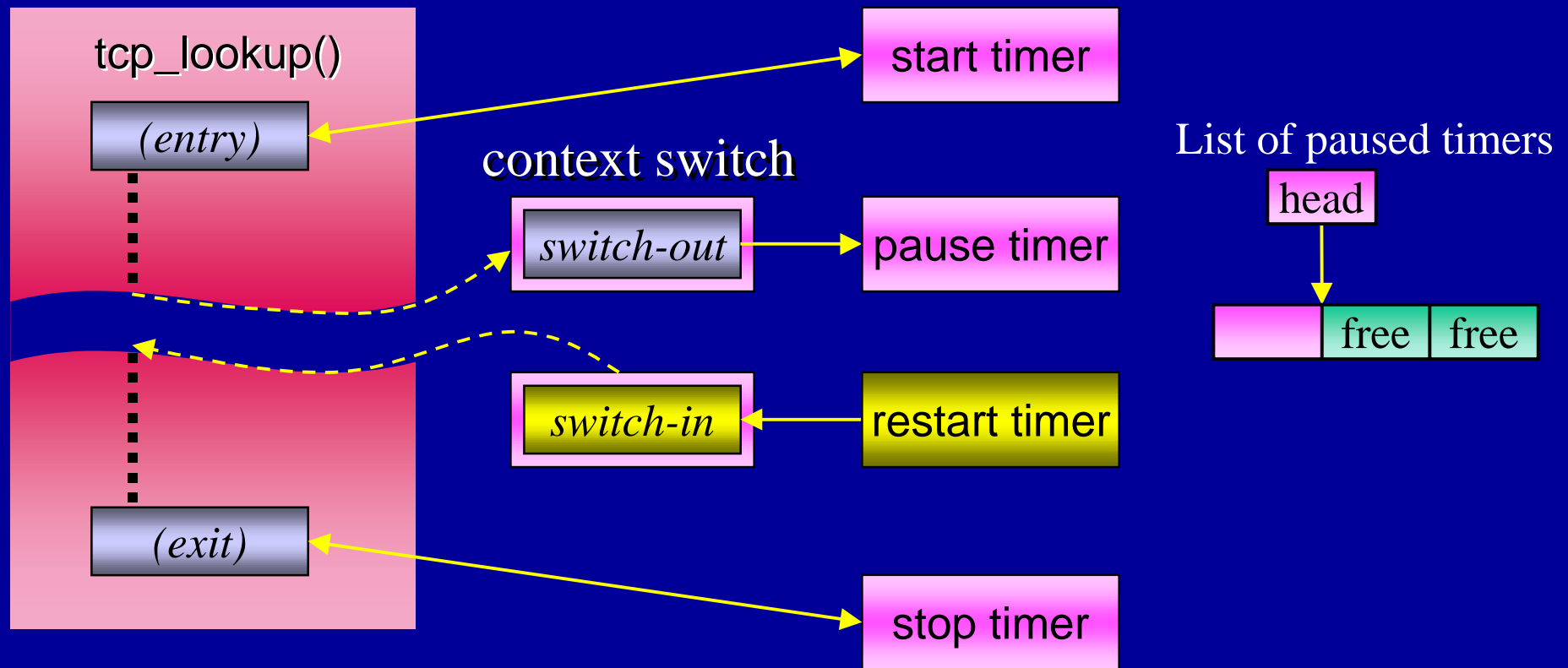
List of paused timers

head

free | free

- Exclude the time spent while switched out
  - Instrument context switch routines
- HW counter metrics are based on this mechanism

# Non-MP CPU timer primitive

**tcp_lookup()**

*(entry)* → start timer

**context switch**

*switch-out* → pause timer

*switch-in* ← restart timer

*(exit)* → stop timer
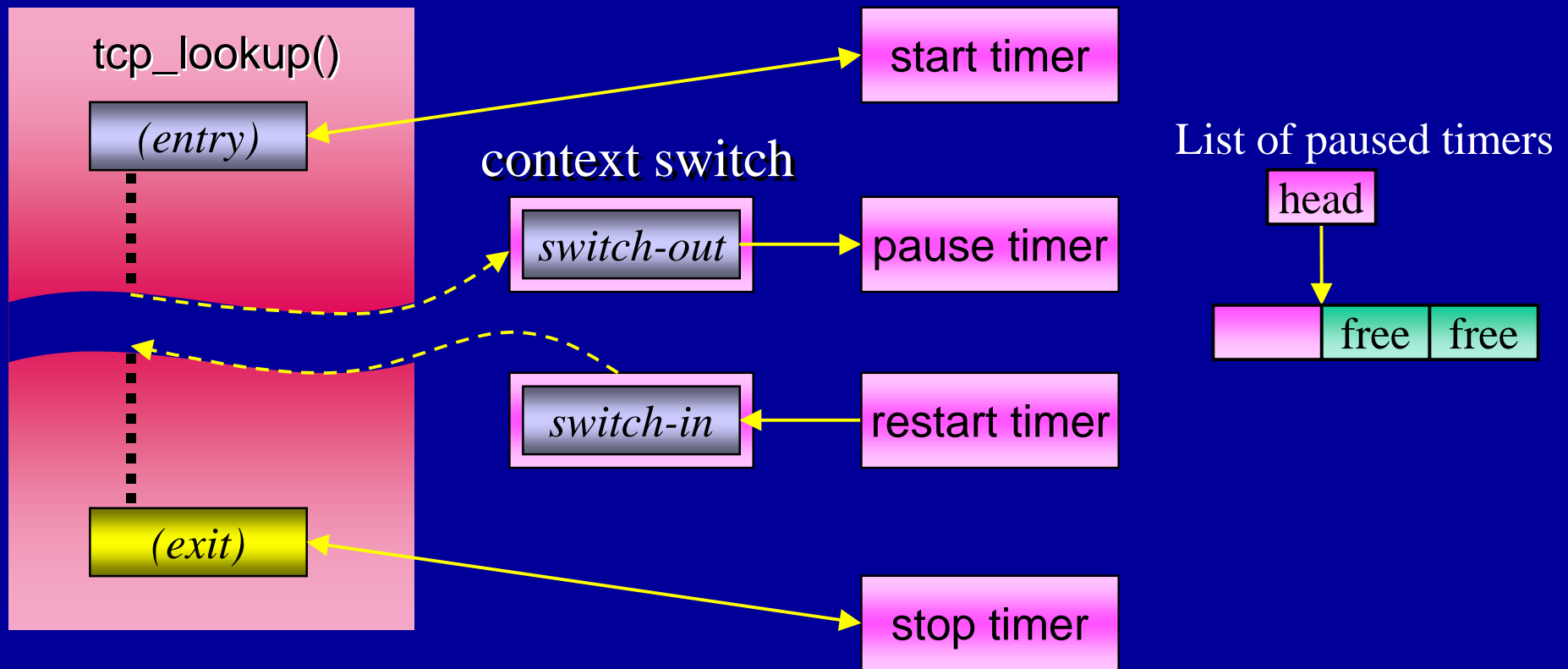
List of paused timers

head

free | free

- Exclude the time spent while switched out
  - Instrument context switch routines
- HW counter metrics are based on this mechanism

# Non-MP CPU timer primitive

**tcp_lookup()**

*(entry)*

start timer

**context switch**

*switch-out* → pause timer

*switch-in* ← restart timer

*(exit)*

stop timer

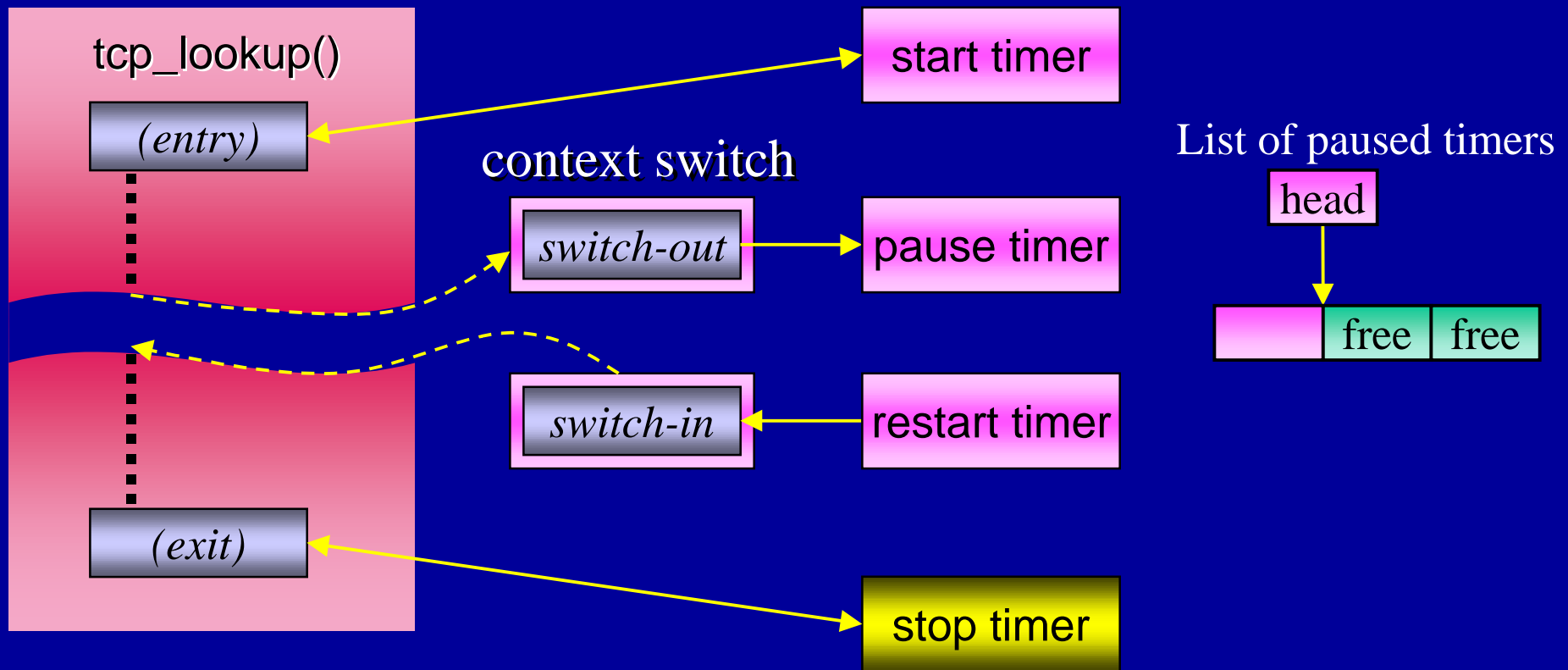List of paused timers

head

tmr | free

- Exclude the time spent while switched out
  - Instrument context switch routines
- HW counter metrics are based on this mechanism

# Non-MP CPU timer primitive

**tcp_lookup()**

*(entry)*

start timer

**context switch**

*switch-out* → pause timer

*switch-in* ← restart timer

*(exit)* → stop timer

List of paused timers

head

tmr | free

- Exclude the time spent while switched out
  – Instrument context switch routines
- HW counter metrics are based on this mechanism
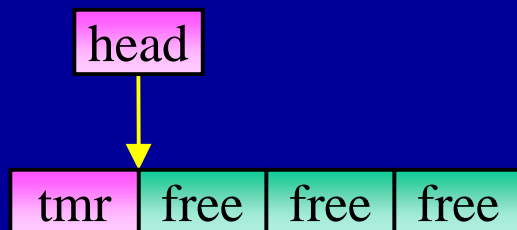
# Non-MP CPU timer primitive

tcp_lookup()

(entry)

start timer

context switch

switch-out → pause timer

switch-in ← restart timer

(exit)

stop timer

List of paused timers

head

free | free

- Exclude the time spent while switched out
  - Instrument context switch routines
- HW counter metrics are based on this mechanism

# Non-MP CPU timer primitive

**tcp_lookup()**

*(entry)*

**start timer**

**context switch**

*switch-out* → **pause timer**

*switch-in* ← **restart timer**

*(exit)* ← **stop timer**

List of paused timers

head

free | free

- Exclude the time spent while switched out
  - Instrument context switch routines
- HW counter metrics are based on this mechanism

# Non-MP CPU timer primitive

tcp_lookup()

*(entry)*

start timer

context switch

*switch-out* → pause timer

List of paused timers

head

free | free

*switch-in* ← restart timer

*(exit)*

stop timer

- Exclude the time spent while switched out
  - Instrument context switch routines
- HW counter metrics are based on this mechanism

# Kperfmon-MP: Goals

Modify uniprocessor Kperfmon to provide:

- Safe operation on SMP machines
  - Thread safety
  - Migration safety
- New feature: Per-CPU performance data
  - More detailed performance data
  - Reduce cache coherence traffic caused by the tool

# Thread Safety

Non-MP timer allocation routine

head

tmr | free | free | free

```
ld [head], R1
add R1, 4, R1
st R1, [head]
```

- Used on switch-out to save the paused timers
- Context switch is serial on uniprocessors
  - No thread safety problems there
- Context switches may be concurrent on SMPs!
  - Multiple threads are being scheduled simultaneously
  - The allocation code is no longer safe

# Thread Safety

head

tmr | free | free | free

MP timer allocation routine

```
alloc:
    ld [head], R1
    add R1, 4, R2
    cas [head], R1, R2
    cmp R1, R2
    bne alloc
```

- Context switches may be concurrent on SMPs
- Use the atomic cas instruction to ensure safety

# Per-CPU performance data

tcp_lookup()

*(entry)*

Code Patch Area

…
rd cpu#, r0
ldx cnt[r0], r1
add r1, 1, r2
casx r2, cnt[r0]
…

Data Area

*cnt-cpu0*

*cnt-cpu1*

*cnt-cpu31*

- Instrumentation code is shared by all CPUs
- Per-CPU copies of the primitive's data
  - Two copies are never placed in the same cache line

# Migration Between Primitives

CPU0

tcp_lookup()

*(entry)*

start timer

context switch

Data Area

switch-out

*timer-cpu0*

switch-in

*timer-cpu1*

*(exit)*

stop timer

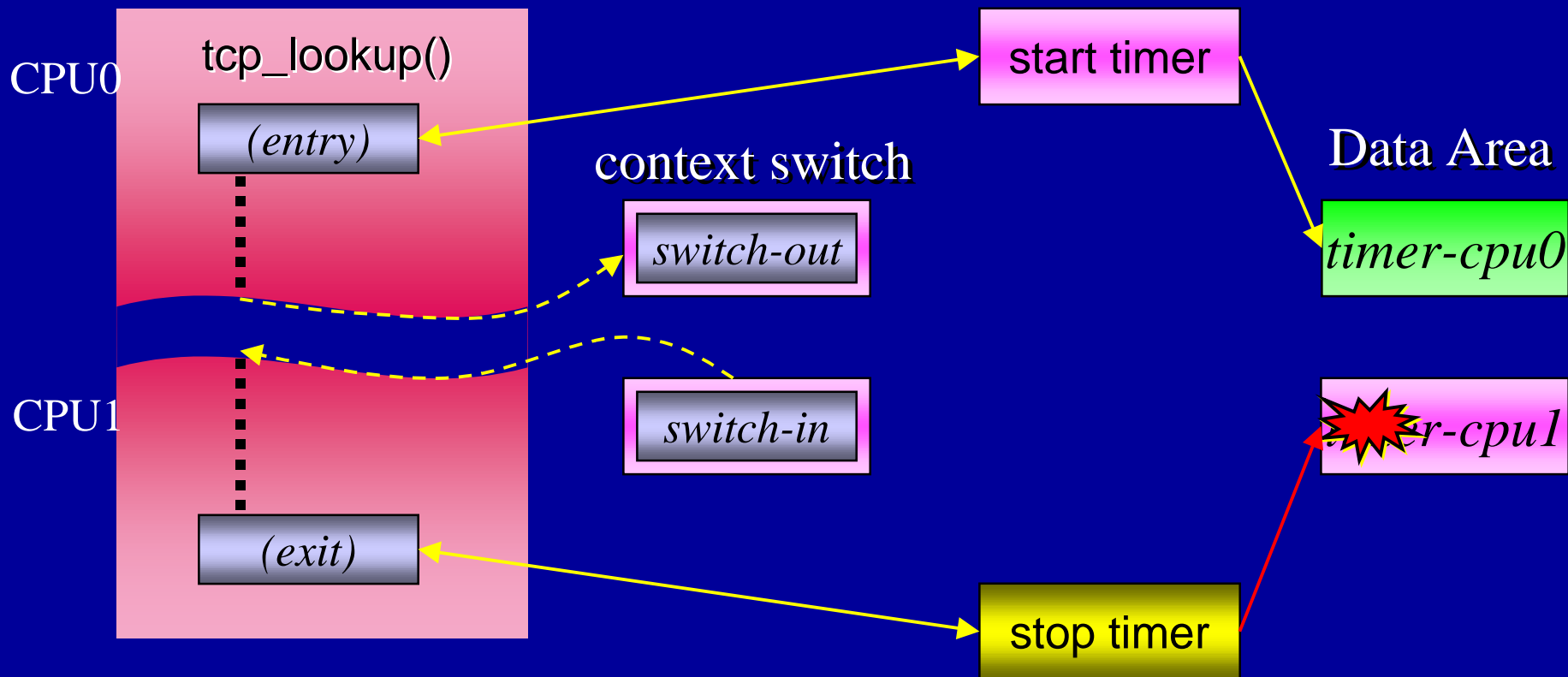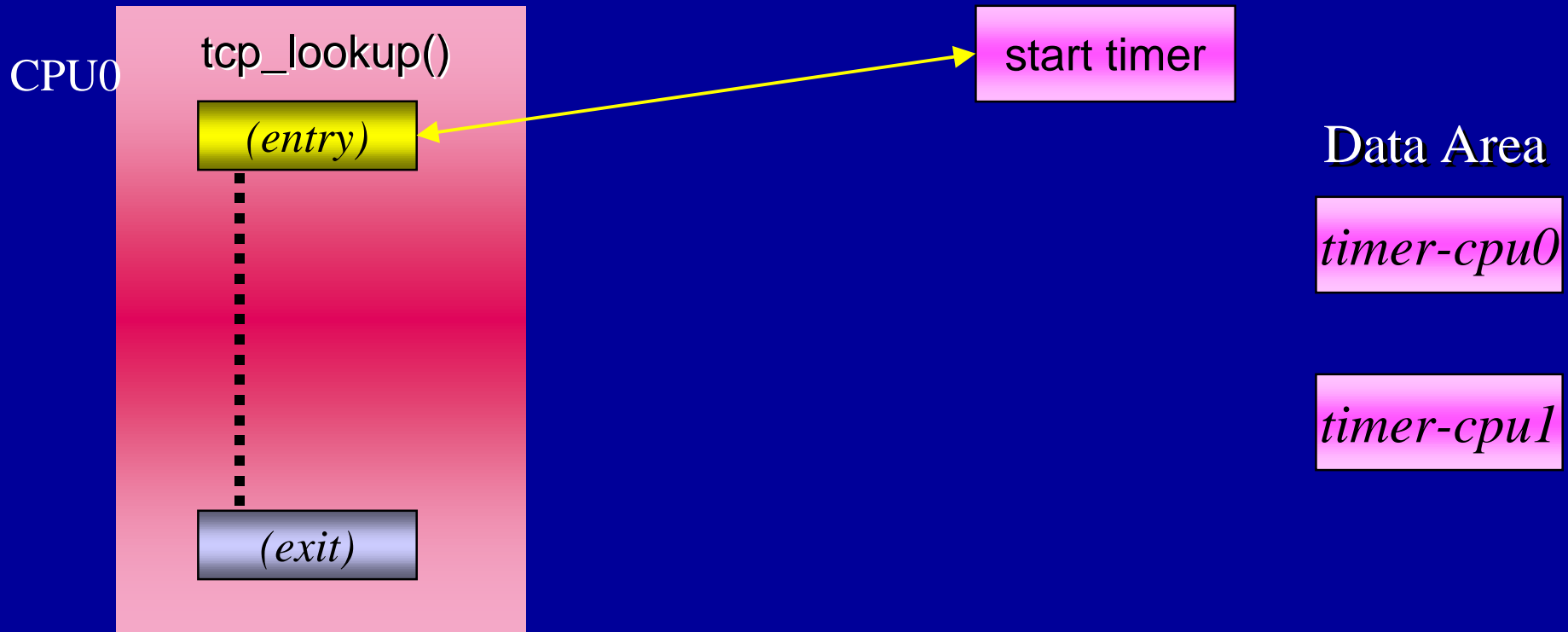- Wall timer started on CPU0, stopped on CPU1
- Counters and CPU timers are not affected

# Migration Between Primitives

CPU0

tcp_lookup()
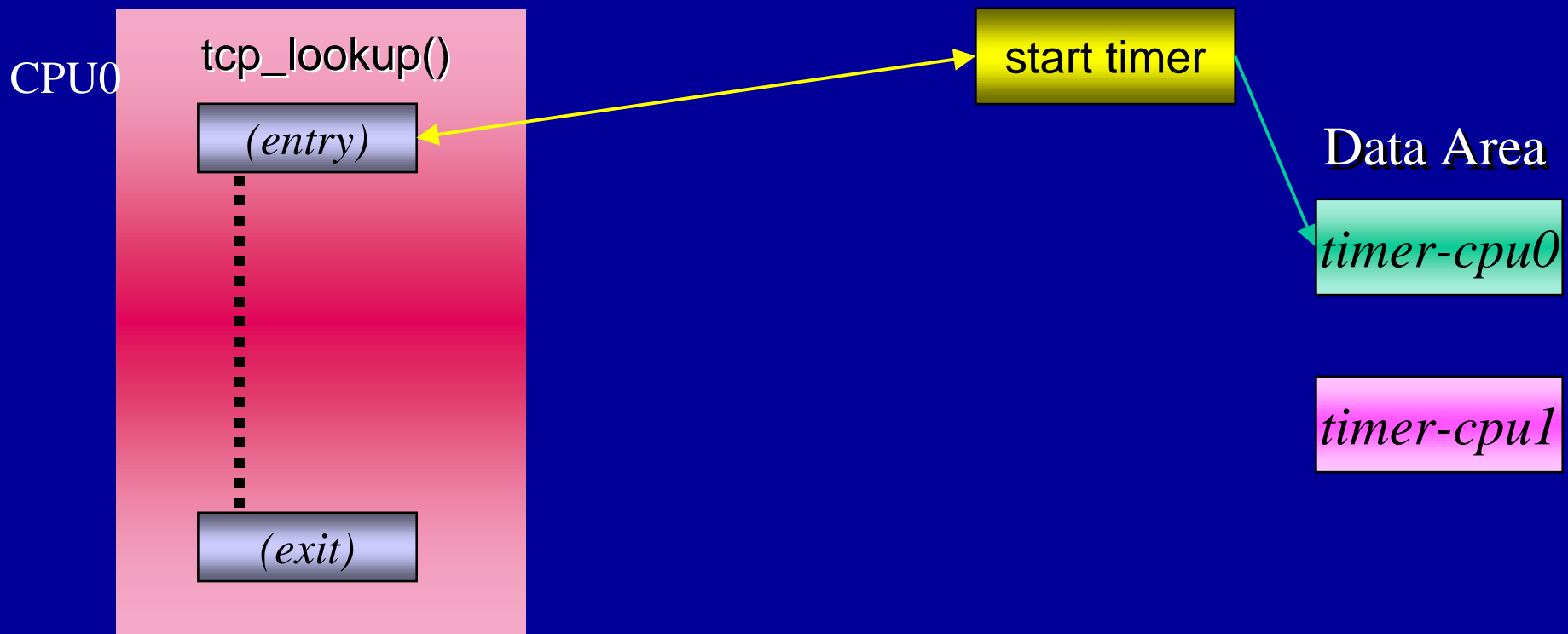
*(entry)*

*(exit)*

start timer

context switch

switch-out

switch-in

Data Area

*timer-cpu0*

*timer-cpu1*

stop timer

- Wall timer started on CPU0, stopped on CPU1
- Counters and CPU timers are not affected

# Migration Between Primitives

CPU0

tcp_lookup()

*(entry)*

start timer

context switch

Data Area

*switch-out*

*timer-cpu0*

*switch-in*

*timer-cpu1*

*(exit)*

stop timer

- Wall timer started on CPU0, stopped on CPU1
- Counters and CPU timers are not affected

# Migration Between Primitives

CPU0

tcp_lookup()

*(entry)*

start timer

Data Area

*timer-cpu0*

context switch

*switch-out*

CPU1

*switch-in*

*timer-cpu1*

*(exit)*

stop timer

- Wall timer started on CPU0, stopped on CPU1
- Counters and CPU timers are not affected

# Migration Between Primitives

CPU0

tcp_lookup()

*(entry)*

start timer

context switch

*switch-out*

Data Area

*timer-cpu0*

CPU1

*switch-in*

*timer-cpu1*

*(exit)*

stop timer

- Wall timer started on CPU0, stopped on CPU1
- Counters and CPU timers are not affected

# Migration Between Primitives

CPU0

**tcp_lookup()**

*(entry)*

**start timer**

context switch

*switch-out*

Data Area

*timer-cpu0*

CPU1

*switch-in*

*...er-cpu1*

*(exit)*

**stop timer**

- Wall timer started on CPU0, stopped on CPU1
- Counters and CPU timers are not affected

# Solution: virtualization

CPU0

tcp_lookup()

*(entry)*

start timer

Data Area

*timer-cpu0*

*timer-cpu1*

*(exit)*

- Implement wall timers on top of CPU timers!

# Solution: virtualization

CPU0

tcp_lookup()

(entry)

(exit)

start timer

Data Area

*timer-cpu0*

*timer-cpu1*

- Implement wall timers on top of CPU timers!

# Solution: virtualization



- Implement wall timers on top of CPU timers!

# Solution: virtualization



- Implement wall timers on top of CPU timers!

# Solution: virtualization



- Implement wall timers on top of CPU timers!

# Solution: virtualization



- Implement wall timers on top of CPU timers!

## kperfmon

### Kperfmon v0.4.2

#### Metrics

(Middle-click on a metric for details)

- ☐ entries to  ☐ exits from
- ☐ exits from (no unwinding tail calls)
- ☐ threadSeconds  ☐ latency/invoc
- ☐ vtime  ☐ vtime/invoc  ☐ walltime
- ☐ walltime/invoc
- ☐ D-$ VReads  ☐ D-$ VReads/invoc
- ☐ D-$ VReadHits
- ☐ D-$ VReadHits/invoc  ☐ D-$ VWrites
- ☐ D-$ VWrites/invoc  ☐ D-$ VWriteHits
- ☐ D-$ VWriteHits/invoc
- ☐ E-$ VRefs  ☐ E-$ VRefs/invoc

#### Kernel Code

root

Code

| genunix ▶ |
| hme (FEPS Ethernet Driver  v1.121 ) ▶ |
| inst_sync (instance binding syscall) ▶ |
| intpexec (exec mod for interp) ▶ |
| ip (IP Streams module) ▶ |
| ipc (common ipc code) ▶ |
| ipdcm (IP/Dialup v1.9) ▶ |
| iwscn (Workstation Redirection driver 'iwscn') ▶ |
| kb (streams module for keyboard) ▶ |
| kerninst (kerninst driver v0.4.1) ▶ |
| krtld |

#### Predicates

Pid(s): [            ]  Clear

#### Function Modifiers

◆ **Fn entry**     ◇ **Fn exit**

◇ **At Insn #**

☐ Just testing (no launcher)   Re-source .tcl

#### CPUs

☐ cpu 0  ☐ cpu 1  ☐ cpu 2  ☐ cpu 3  ☐ cpu 4  ☐ cpu 5  ☐ cpu 6  ☐ cpu 7
☐ cpu 8  ☐ cpu 9  ☐ cpu10  ☐ cpu11  ☐ cpu12  ☐ cpu13  ☐ cpu14  ☐ cpu15
■ sum

#### Disassemble a range of memory

From-addr: [            ]   To-addr: [            ]

Disassemble (kernel)     Disassemble (kerninstd)

☐ include ascii in disassembly

#### Disassemble Selected Fn/Block(s)

Disassemble (orig)

Disassemble (curr mem)

■ do code replacement   **Start a visi**

Kperfmon v0.4.2

| Metrics | Kernel Code |
|---|---|
| (Middle-click on a metric for details) | |

Table Visualization

Table Visualization

Para dyn

File   Actions   View

Phase: Current Phase

| | BranchMispred cpus | VStallTime E-$ frac | VMissRatio I-$ frac | VStallFraction entries to #/sec |
|---|---|---|---|---|
| afs/afs_close/sum all cpus | 9e-06 | 0.005 | 0.008 | 2 |
| afs/afs_getpage/sum all cpus | 9e-06 | 0.2 | 0.06 | 1e+01 |
| afs/afs_open/sum all cpus | 9e-06 | 0.2 | 0.07 | 3 |
| genunix/kmem_alloc/sum all cpus | 0.0001 | 0.3 | 0.1 | 4e+02 |
| tcp/tcp_rput_data/sum all cpus | 7e-06 | 0.1 | 0.2 | 2e+01 |
| ufs/ufs_getpage/sum all cpus | 0.002 | 0.3 | 0.06 | 2e+03 |
| ufs/ufs_lookup/sum all cpus | 0.001 | 0.1 | 0.05 | 2e+03 |
| ufs/ufs_read/sum all cpus | 0.01 | 0.3 | 0.04 | 2e+03 |
| ufs/ufs_readdir/sum all cpus | 0.0005 | 0.2 | 0.05 | 5e+01 |

From-addr:          To-addr:

| Disassemble (kernel) | Disassemble (kerninstd) | Disassemble (orig) |
|---|---|---|
| | | Disassemble (curr mem) |

☐ include ascii in disassembly

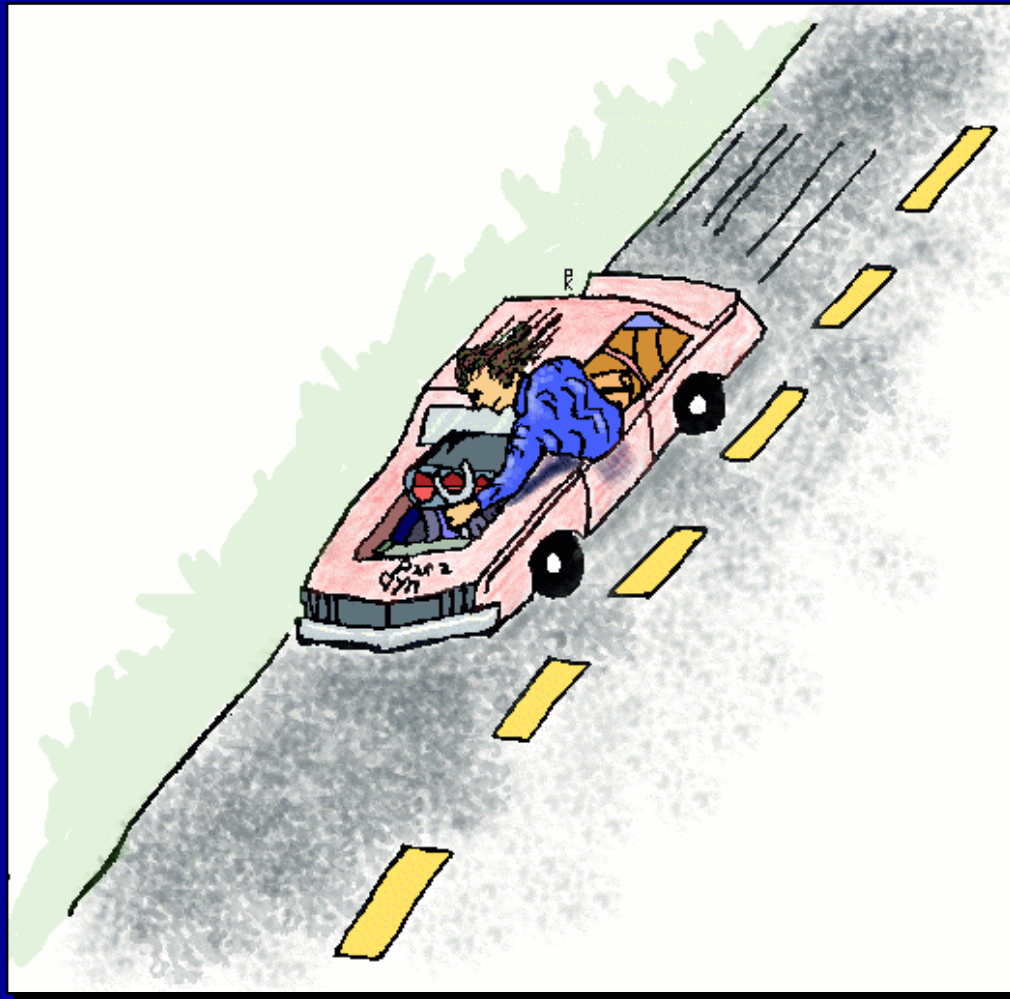☑ do code replacement    Start a visi

# Conclusion

- Techniques for correct MP profiling:
  - Atomic memory updates to ensure thread safety
  - Virtualized timers to handle thread migration
- Per-CPU data collection is important
  - Provides detailed performance information
  - Introduces fewer coherence cache misses

# Future Work

- New metrics
    - Locality of CPU assignments
    - Per-thread performance data
- Formal verification of instrumentation code for migration/preemption problems
- Ports to other architectures and OS'es

# The Big Picture



http://www.cs.wisc.edu/paradyn

# The Big Picture



- Demo: Wednesday, Room 6372
- Available for download on request
  - mailto: mirg@cs.wisc.edu
  - Public release in April