

The State of COMPOSITE: a Customizable Component-Based OS for Predictable, Reliable, and Scalable Computation

Gabriel Parmer

gparmer@gwu.edu

The George Washington University (GWU)

OSPRT 2013

Researchers include

Qi Wang, Jiguo Song, Jakob Kaivo,
Andrew Sweeney, John Wittrock, ...

Embedded Systems

Past:

- single, simple task
- uni-processor
- fault-tolerance ignored (reboot), or custom

Present/Future:

- consolidation
- certification
- multi-/many-core
- increased faults due to shrinking manufacturing processes

Embedded OSes

Past:

- single memory protection domain
- threads, FP scheduling, semaphores, mailboxes, timing
- FreeRTOS, OSEK, ...

Challenges of the Present/Future:

- spatial + temporal isolation
- system composition from independently certifiable pieces
- intra- and inter-task parallelism
- reliability built-in

Embedded OSes

Past:

- single memory protection domain
- threads, FP scheduling, semaphores, mailboxes, timing
- FreeRTOS, OSEK, ...

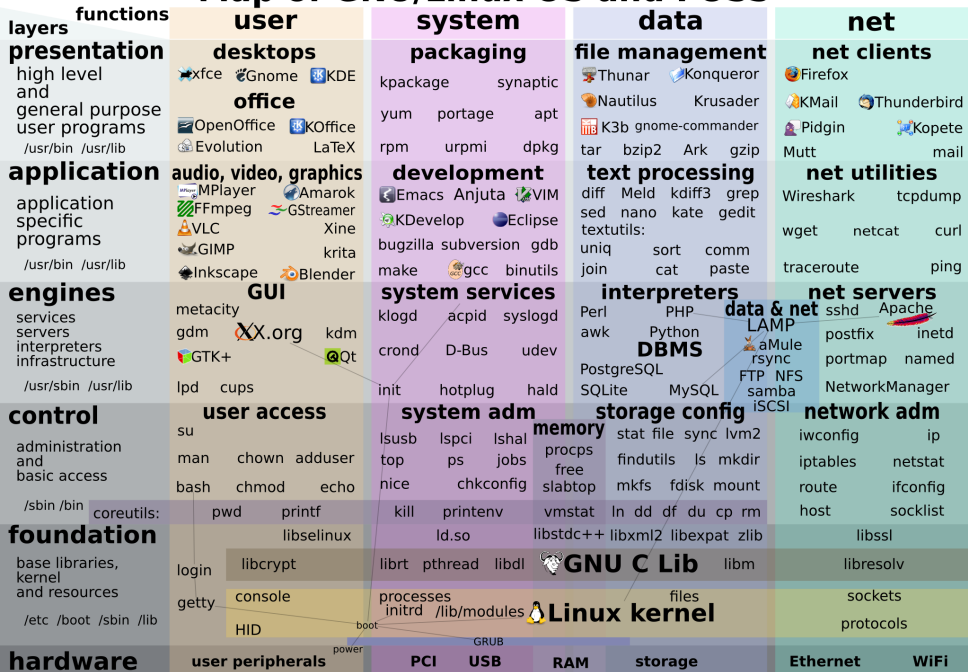
Challenges of the Present/Future:

- spatial + temporal isolation
- system composition from independently certifiable pieces
- intra- and inter-task parallelism
- reliability built-in

Challenge: predictability

Challenge: maintaining system simplicity

Map of GNU/Linux OS and FOSS



The COMPOSITE Component-Based OS

System policies/abstractions are *components*

- user-level
- minimal unit of spatial isolation

Low-level functions are components

- scheduling
- memory mapping
- I/O processing

Threads orthogonal to components

- thread migration
- concurrent/parallel components

Components interact via *invocation* of exported function

- contractually specified interfaces
- function call semantics

System = Components + Composition

Composition

- complex behavior from simple(ish) pieces
- gluing components together \rightarrow raise level of abstraction

Complex functionality from simple pieces...sound familiar?

Hint: Thompson & Ritchie

System = Components + Composition

Composition

- complex behavior from simple(ish) pieces
- gluing components together → raise level of abstraction

Complex functionality from simple pieces...sound familiar?

Hint: Thompson & Ritchie

```
wget -O - www.ecrts.org | grep 'ospert' | wc -l
```

System = Components + Composition

Composition

- complex behavior from simple(ish) pieces
- gluing components together → raise level of abstraction

Complex functionality from simple pieces...sound familiar?

Hint: Thompson & Ritchie

```
wget -O - www.ecrts.org | grep 'ospert' | wc -l
```

```
wget = c "bin/wget" "-O - www.ecrts.org"  
grep = c "bin/grep" "ospert"  
wc    = c "bin/wc" "-l"  
sys   = deps [ (cat, [grep, POSIX]),  
               (grep, [wc, POSIX]) ]
```

System = Components + Composition

Composition

■ complete

■ gluing

Complex function

abstraction

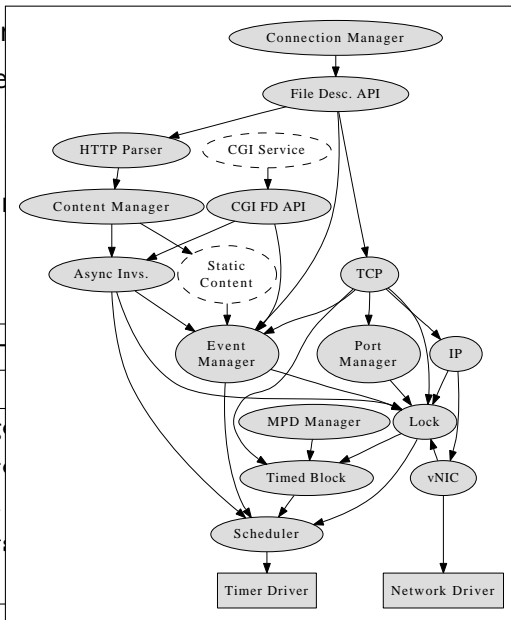
familiar?

wget -O -

' | wc -l

wget
gr
wc
sys

"



System = Components + Composition

Composition Challenges:

- complex
- end-to-end predictability
- gluing
- dependent-task structure to mirror components?
- trade between component concurrency, and memory

abstraction

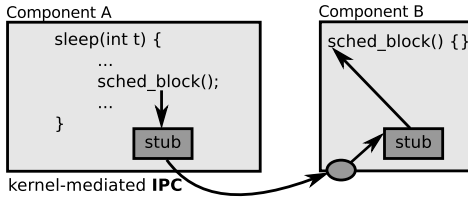
familiar?

```
wget -O - www.ecrts.org | grep 'ospert' | wc -l
```

```
wget = c "bin/wget" "-O - www.ecrts.org"  
grep = c "bin/grep" "ospert"  
wc    = c "bin/wc" "-l"  
sys   = deps [ (cat, [grep, POSIX]),  
               (grep, [wc, POSIX]) ]
```

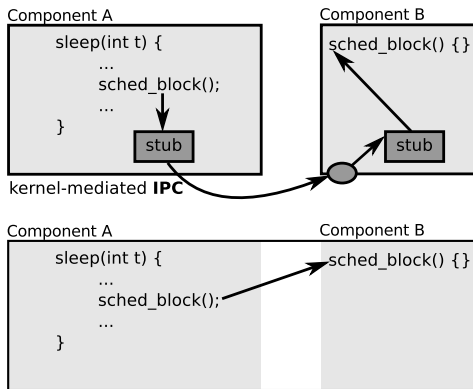
But people understand components...what else?

All problems can be solved by another level of indirection. – Dijkstra



But people understand components...what else?

All problems can be solved by another level of indirection. – Dijkstra



Mutable Protection Domains

- generalizes other system structures (μ kern, exokern, ..)

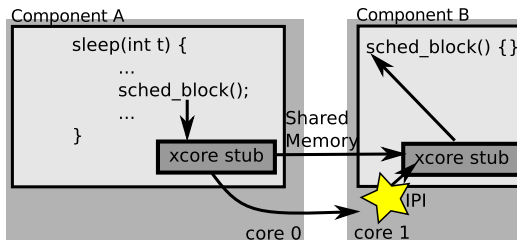
Predictable Parallel Computation

Parallel systems are here, what do we do with them?

- Inter-task parallelism: simple until
 - shared resources
 - schedulability: partitioned + bin-packing
- Intra-task parallelism:
 - fork/join (OpenMP) schedulability
 - general abstractions + mechanisms for parallelism
 - harness hidden parallelism in concurrent systems

think: `wget www.ecrts.org& wget www.rtss.org&`

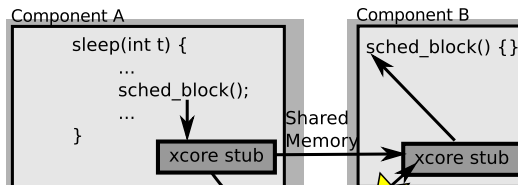
Many-core Composite: MC²



Inter-component parallelism:

- bin-packing overheads for partitioned systems
- *cut* a task across cores
 - synchronous communication across cores
- specialized mechanisms for cross-core thread activation
 - intra-component: 4x faster than Linux (WC)
 - inter-component: harness non-blocking, async APIs

Many-core Composite: MC²

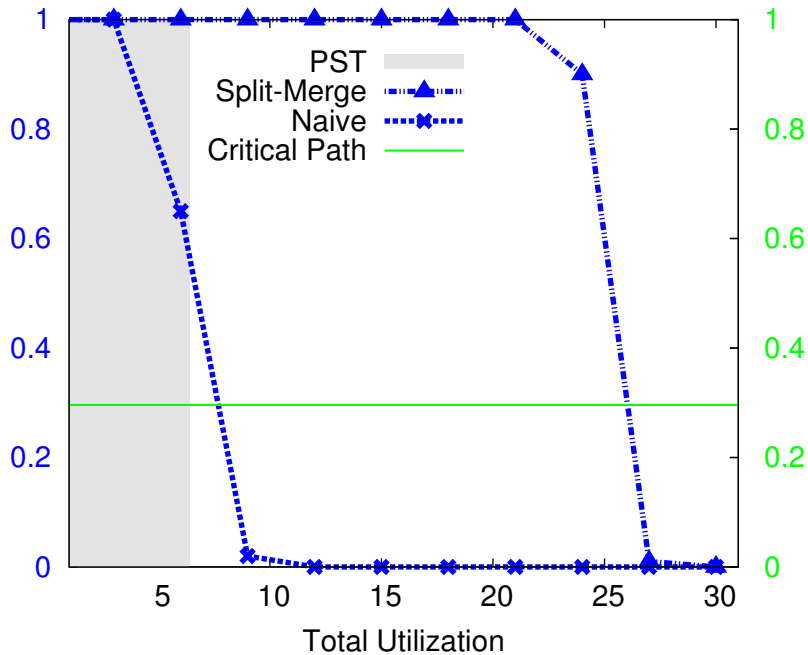


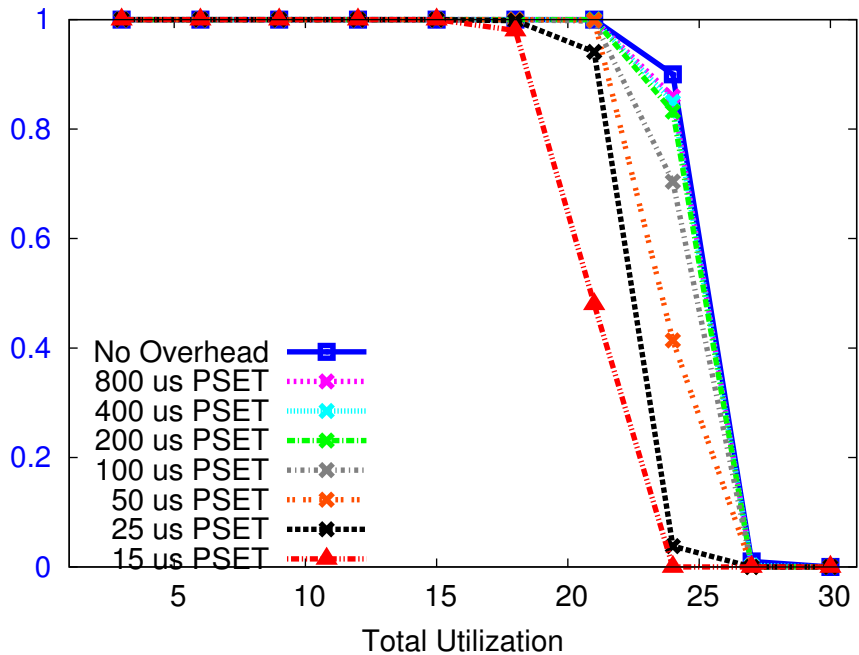
Pair this with:

- a smart assignment algorithm, and
- optimized holistic analysis to analyze schedulability.

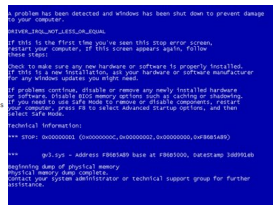
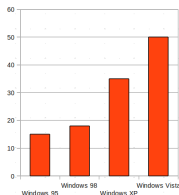
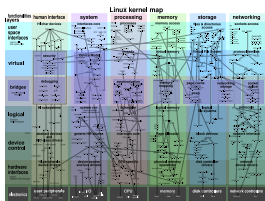
- bin-packing overheads for partitioned systems
- *cut* a task across cores
 - synchronous communication across cores
- specialized mechanisms for cross-core thread activation
 - intra-component: 4x faster than Linux (WC)
 - inter-component: harness non-blocking, async APIs

Schedulability Ratio





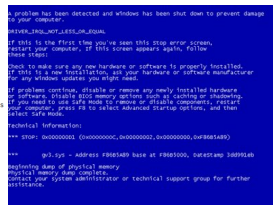
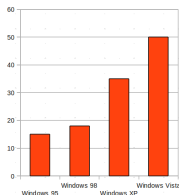
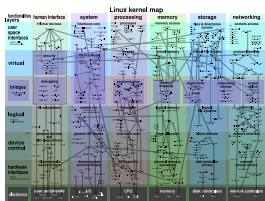
Transparent, System-Provided, Fault Tolerance



Decreasing process sizes

- + faster
- + less power
- + smaller
- increased vulnerability to HW transient faults
- 65% of HW faults corrupt OS state

Transparent, System-Provided, Fault Tolerance



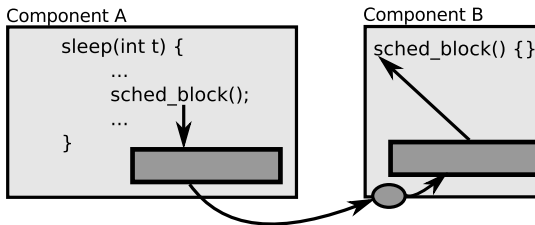
Decreasing process sizes

- + faster
- + less power
- + smaller
- increased vulnerability to HW transient faults
- 65% of HW faults corrupt OS state

Can we provide fault tolerance

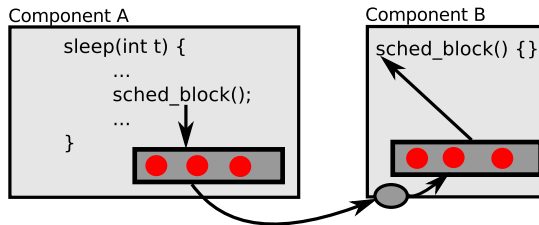
- even for the lowest-level components?
- predictably and efficiently?

Computational Crash Cart: C³



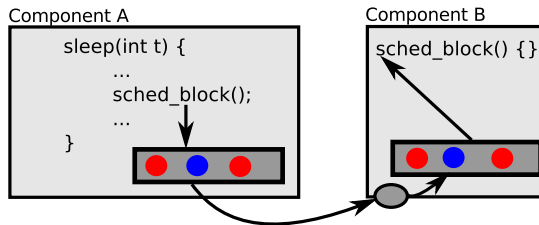
- 1 interpose on communication between components
- 2 track state of each “shared” object
 - file, thread, lock, ...
- 3 fault in server!
- 4 μ -reboot component
- 5 rebuild state via functions in interface

Computational Crash Cart: C³



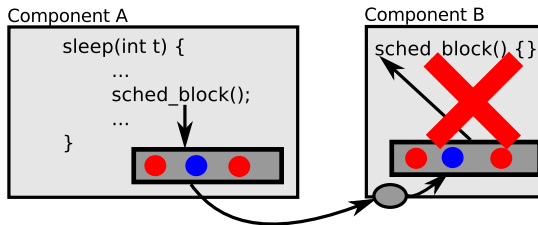
- 1 interpose on communication between components
- 2 track state of each “shared” object
 - file, thread, lock, ...
- 3 fault in server!
- 4 μ -reboot component
- 5 rebuild state via functions in interface

Computational Crash Cart: C³



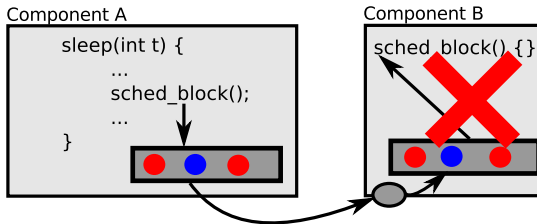
- 1 interpose on communication between components
- 2 track state of each “shared” object
 - file, thread, lock, ...
- 3 fault in server!
- 4 μ -reboot component
- 5 rebuild state via functions in interface

Computational Crash Cart: C³



- 1 interpose on communication between components
- 2 track state of each “shared” object
 - file, thread, lock, ...
- 3 fault in server!
- 4 μ -reboot component
- 5 rebuild state via functions in interface

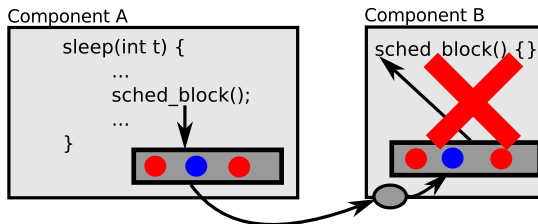
Computational Crash Cart: C³



Recovery affects timing of multiple threads

- performed on-demand by thread using object
- rebuild objects at proper priority
- avoid *recovery inversion*

Computational Crash Cart: C³



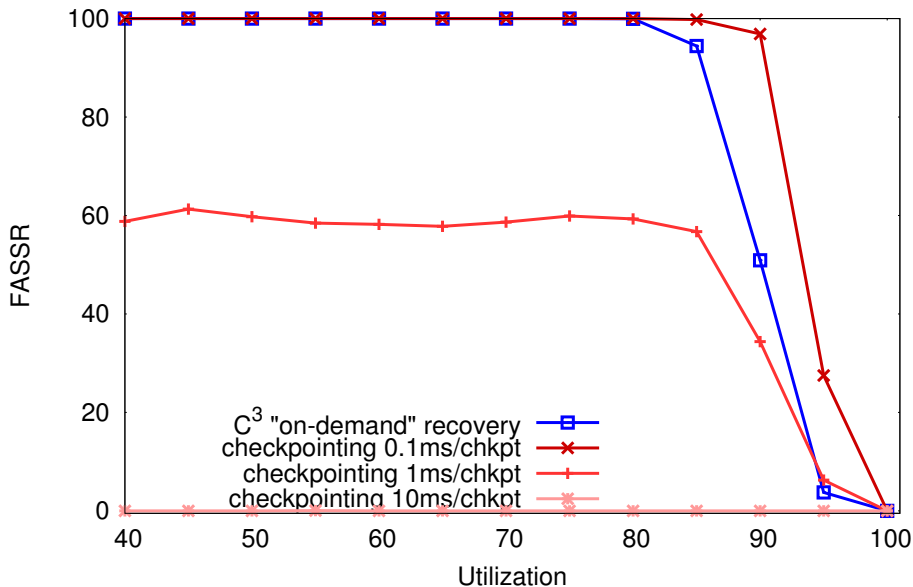
C³: Efficient, system-wide fault tolerance

- recovers 100% injected faults (scheduler, memmgr, fs)
- μ -reboot in $< 20\mu$ -sec
- rebuild object: $< 5\mu$ -sec

Versus checkpointing

- CRIU: 10ms, Xen: 10sec
- C³ : 0.1ms per MB

Fault-Tolerant Systems Schedulability:
Checkpointing and C^3 , 50 tasks, 100ms period



The State of COMPOSITE is...

...in progress.

- MC²: Full-system, predictable parallelism
- C³: Predictable, system-level fault tolerance
- HIEROS: hierarchical paravirtualization
(FreeRTOS done, Linux in-progress)
- ISOLOS: separation kernel support
- SECCOS: fine-grained authentication + monitoring
- ...POSIX support (see Rob Pike's polemic)

COMPOSITE as CBOS:

- configurable to system reqs; as complex as required
- generalizes system structures

COMPOSITE as memory isolation + function call indirection

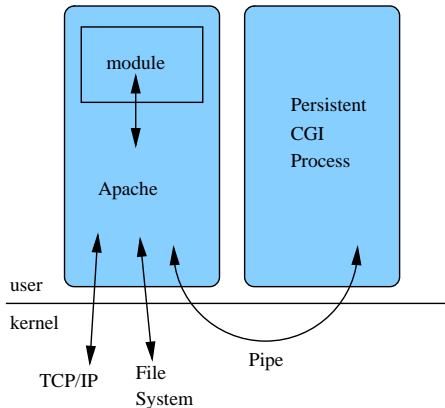
- general, transparent parallelism
- system-level fault tolerance

Thank You!

? || /* */

composite.seas.gwu.edu

Comparison Case: Apache Web-Server, Linux

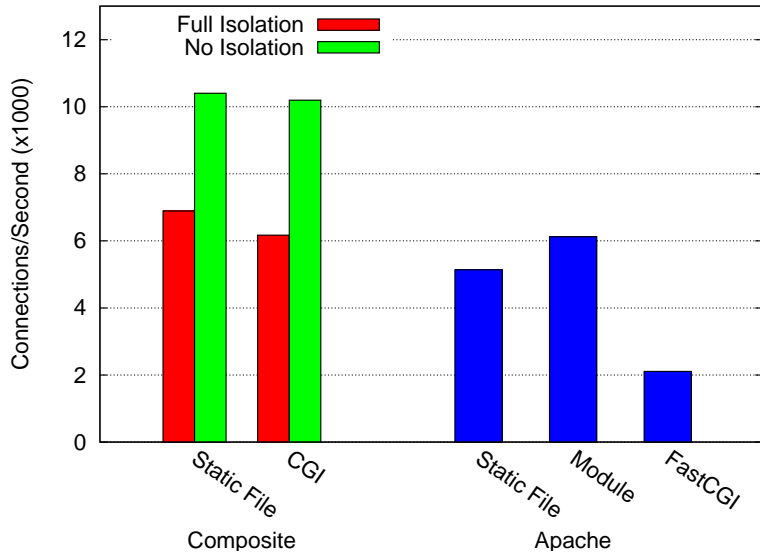


Apache provides multiple content sources

Figures to keep in mind:

- Linux CGI communication (pipe RPC): $6.4 \mu\text{-sec}$
- COMPOSITE component communication: $0.67 \mu\text{-sec}$

Apache, Composite Comparison



Resource Management

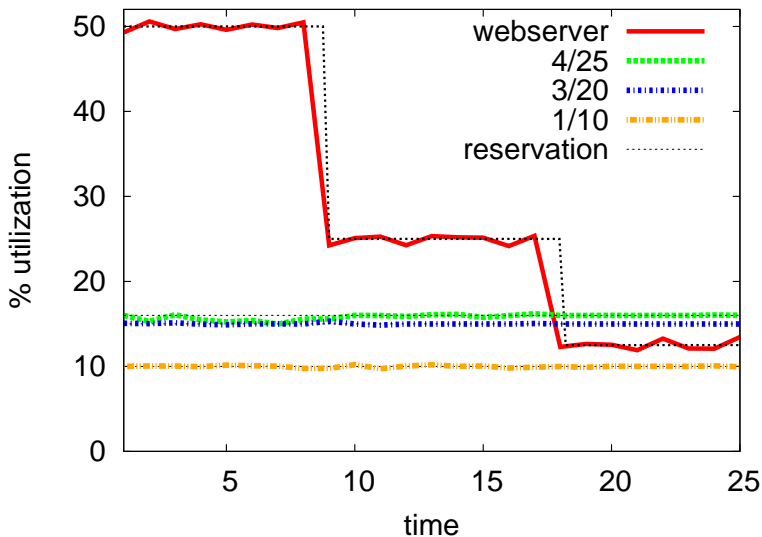
Components configured in the system:

- schedulers
- memory mappers
- I/O managers
- file systems
- networking protocols
- ...

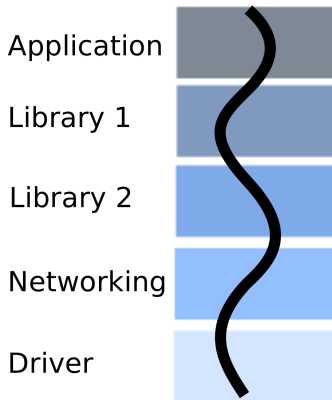
Cost of component resource mgmt? (in μ -seconds)

- Scheduler: thread switch – 0.4 (cos) vs. 0.8 (linux)
- Memory mapping: `mmap` – 2 (cos) vs. 6 (linux)
- I/O: receive packet – 9.69 (cos) vs. 10.3 (linux)

Best Effort Subsystem vs. RT Task Execution



System Management of Parallelism



Traditional model of computation

- thread executes through system layers
- each layer has its own data working set

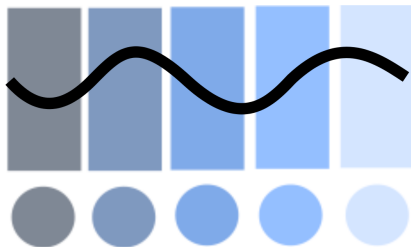
System Management of Parallelism



Traditional model of utilizing parallelism

- thread execute through same layers
- same data working sets in each cache
→ inefficient use of caches!

System Management of Parallelism



COMPOSITE w/ invocations spreading computation across cores

- CPU caches specialize around a specific working set
- controlled cache inefficiency
 - factor of 100 performance difference
- control the parallelism of any one component