



Scheduler Activations: Effective Kernel Support for the User-Level Management of Parallelism

Thomas E. Anderson, Brian N. Bershad, Edward D.
Lazowska, Henry M. Levy

Presenters: Yexi Jiang, Jorge Arenas, Hao Jiang



The Authors



- Thomas E. Anderson . Robert E. Dinning Professor.
 - constructing robust, secure, and efficient computer systems
- Brian N. Bershad. Professor.
 - computer operating systems, distributed systems and architecture
- Edward D. Lazowska. Bill & Melinda Gates Chair in Computer Science & Engineering.
 - the design, implementation, and analysis of high performance computing and communication systems
- Henry M. Levy. Chairman and Wissner-Slivka Chair
 - operating systems, distributed systems, security, the world-wide web, and computer architecture



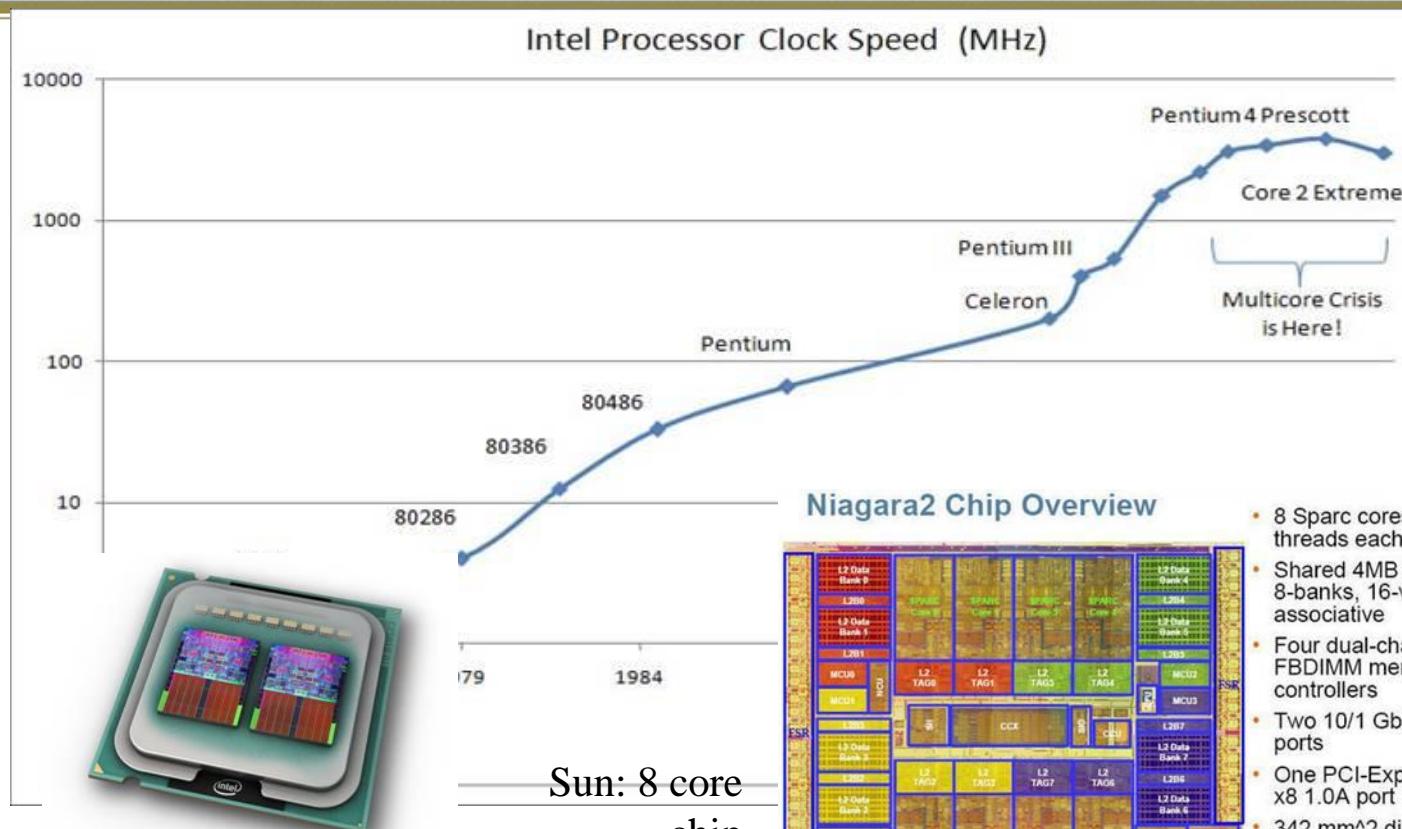
- Background
- Problem Statement
- The Approach
- Related Works
- User-level Thread
- Kernel-level Thread
- Scheduler Activation
- Implementation
- Experimental Evaluation
- Summary



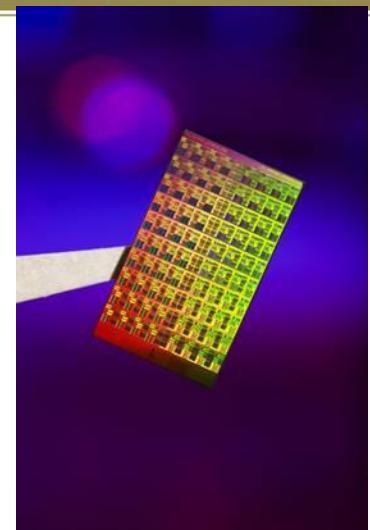
- Technologies Trend.
- Process vs. Thread, and Why Thread?
- User Mode vs. Kernel Mode
- User-Level Thread vs. Kernel-Level Thread



Technologies Driving Multi-Programming



- 8 Sparc cores, 8 threads each
- Shared 4MB L2, 8-banks, 16-way associative
- Four dual-channel FBDIMM memory controllers
- Two 10/1 Gb Enet ports
- One PCI-Express x8 1.0A port
- 342 mm² die size in 65 nm
- 711 signal I/O, 1831 total



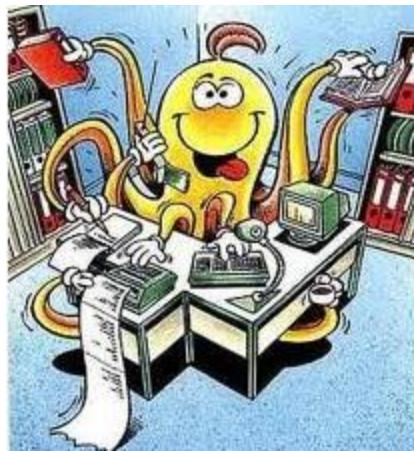
Intel: 80 core experimental system



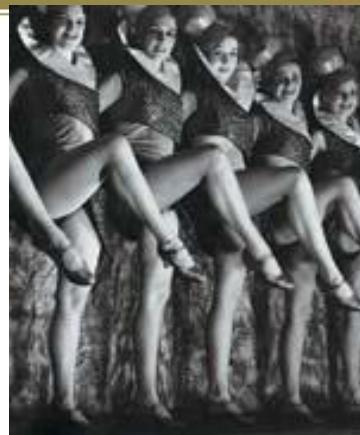
Reasons for Multi-Programming



performance



multitasking



coordination



parallelism



- Process is a program in execution, thread is a light-weight process.

	Process	Thread
Data Structure	PCB	TCB
Communication among processes/threads	Message, Shared Memory	Inherent Shared Memory
Time overhead	High	Low
Protection among processes/threads	Yes	No



- ❖ How to develop an editor that allow typo checking, document printing, and user input simultaneously?

Process based	Thread based
1 process for typo checking, 1 process for printing, 1 process for handling user input. Explicitly leverage inter-process communication. Complicate!	1 thread for typo checking, 1 thread for printing, 1 thread for handling user input. Use shared data structure.



➤ **Kernel Mode**

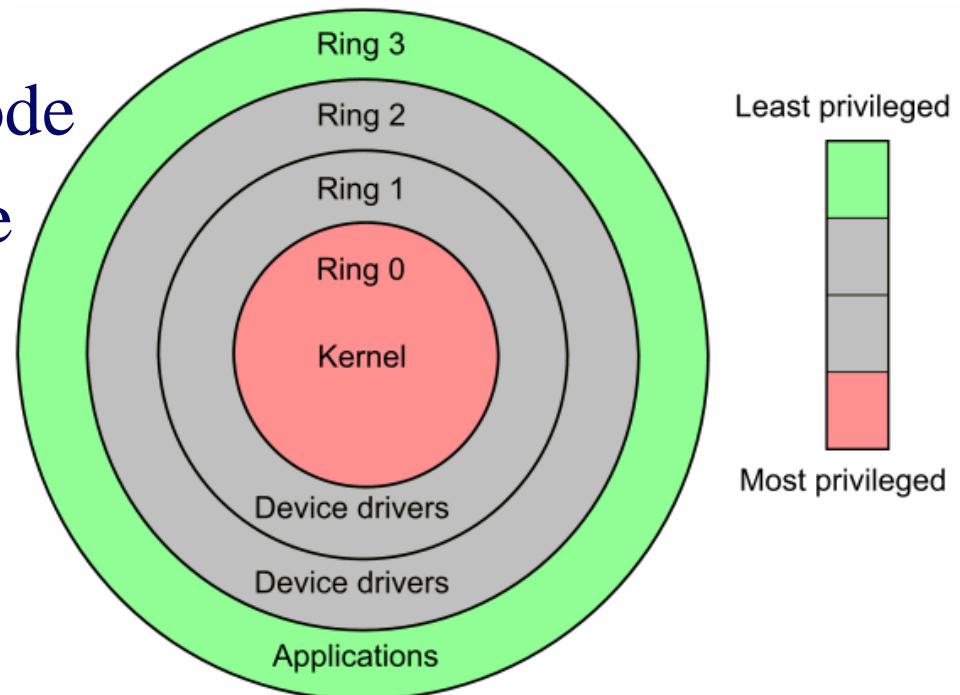
- Unrestricted access to hardware.
- Execute trusted functions.
- Crash is severe.

➤ **User Mode**

- No ability to directly access hardware.
- Crash is OK.



- X86 architecture
 - Ring 0: kernel mode
 - Ring 3: user mode
 - Ring 1,2: driver?



User-Level Threads vs. Kernel-Level Threads

- **User-level threads:** Build without any modifications to the underlying OS. Running on top of the system in user mode.
- **Kernel-level threads:** Build in kernel. Running directly in kernel mode.



User-Level Threads vs. Kernel-Level Threads

	User-level threads	Kernel-level threads
Pros	<ol style="list-style-type: none">1. No modification for kernel2. Low overhead	Full-knowledge of all threads.
Cons	<ol style="list-style-type: none">1. Threads in a process is invisible to the OS.2. Block one thread would block whole process.3. Access to multiple processors is not guaranteed.	<ol style="list-style-type: none">1. High overhead.2. Synchronization is more expensive.



User-Level Threads vs. Kernel-Level Threads

- This paper is one of the pioneers that takes advantages from both user-level and kernel-level threads.
- Models: One-to-one, many-to-one, many-to-many, two-level. (See details later)



- Background
- Problem Statement
- The Approach
- Related Works
- User-level Thread
- Kernel-level Thread
- Scheduler Activation
- Implementation
- Experimental Evaluation
- Summary



Problem Statement

- Existing implementations
 - User-level choices: *PCR*, *FastThreads*
 - Kernel-level choices: threads support in *Mach*, *Topaz*, *V*



- Choose user-level threads
 - High performance.
 - Exhibit poor performance or incorrect behavior in presence of I/O, page faults, preemption.
- Choose kernel-level threads
 - No system integration problem. Threads are directly scheduled by kernel.
 - Low performance.

Neither is good enough!



- Functionality
 - No idled process when ready threads exist.
 - No high-priority threads wait for low-priority threads.
 - Block a thread would not block whole process
- Performance
 - As fast as user-level threads
- Flexibility
 - User-level part enables application-specific customization.



- Background
- Problem Statement
- **The Approach**
- Related Works
- User-level Thread
- Kernel-level Thread
- Scheduler Activation
- Implementation
- Experimental Evaluation
- Summary



- A hybrid threads design.
 - Threads are implemented in both user-level and kernel-level . Only transmit necessary information between user-level and kernel-level.
- Each application:
 - Provided with a virtual multiprocessor.
 - Has full knowledge about the allocated processors.
 - Can control which threads run on which processor.

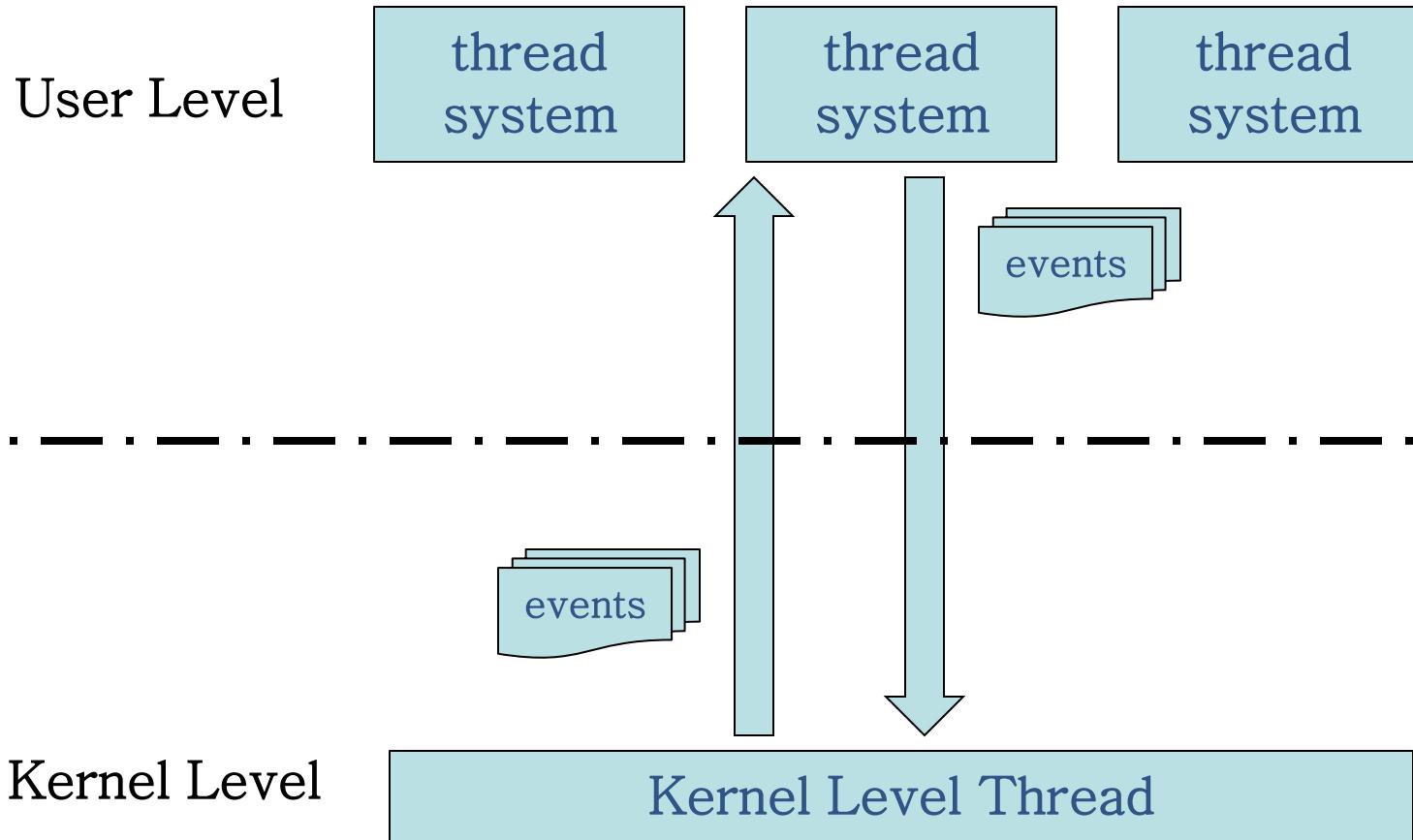


- Communicating upward
 - Improve functionality
- Communicating downward
 - Preserve performance

Scheduler activations: The execution context for an event vectored from the kernel to an address space.



The Approach



The Approach

User Level

thread
system

thread
system

thread
system

Notify the address space
thread scheduler of
every event affecting the
address space

Kernel Level Thread



- Background
- Problem Statement
- The Approach
- Related Works
- User-level Thread
- Kernel-level Thread
- Scheduler Activation
- Implementation
- Experimental Evaluation
- Summary



Related Works

- Psyche and Symunix

	Scheduler Activations	Psyche and Symunix
Platform	Topaz (UNIX)	Psyche: UNIX Symunix: New OS
Communication	Between user-level and kernel-level	Between kernel and application
Synchronization	Low overhead	High overhead, Should be handled between applications and kernel



- Compare with traditional asynchronous I/O system.
 - Uniformly addresses preemption, I/O, and page faults.
 - No need to major change the structure for application and kernel code.



Related Works

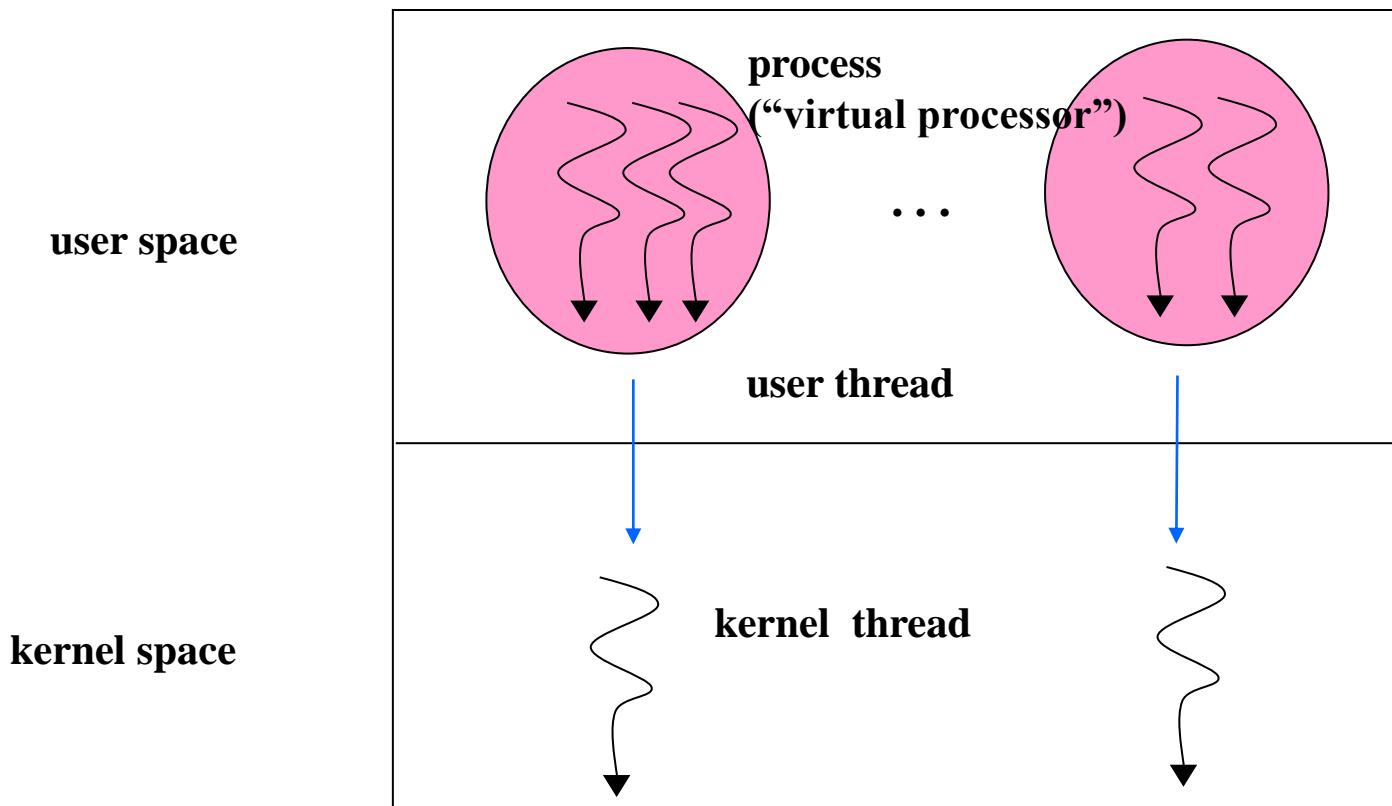
- Compare with Hydra
 - Both implement scheduling policy out of kernel.
 - Hydra: scheduling requires coordinating between scheduler server and kernel. Scheduler Activation: Scheduling is autonomic.



- Background
- Problem Statement
- The Approach
- Related Works
- User-level Thread
- Kernel-level Thread
- Scheduler Activation
- Implementation
- Experimental Evaluation
- Summary



User-level Thread



- **Divide up the process resources between threads.**
 - N:1 arrangement, with N user threads sharing one kernel thread.
 - Managed by runtime library routines.
 - Require no kernel intervention.
- **E.g.: GNU PTH Thread Library, FSU Threads.**
- **The thread package views each process as a virtual processor.**
 - Each virtual processor runs user-level code that pulls threads off the ready queue and runs them.



Advantages of User-level Threads

- **Occupy the user address space of the process.**
 - No overhead of kernel intervention when switching between threads.
 - Not much more expensive than a regular procedure call.
- **Flexible and amenable to specialization.**
 - Only need to meet the requirements of the specific applications that will use them.
 - If priority scheduling is not needed, for example, leave it out.



- **There is a lack of kernel support in existing systems.**
 - A process gets the same amount of time regardless of how many threads it is running.
 - System is not able to make use of multiprocessors.
- **Kernel threads block, resume, and are preempted without the user level threads being notified.**
 - Page faults, I/O can block the whole process.
 - Poor performance or cause errors.



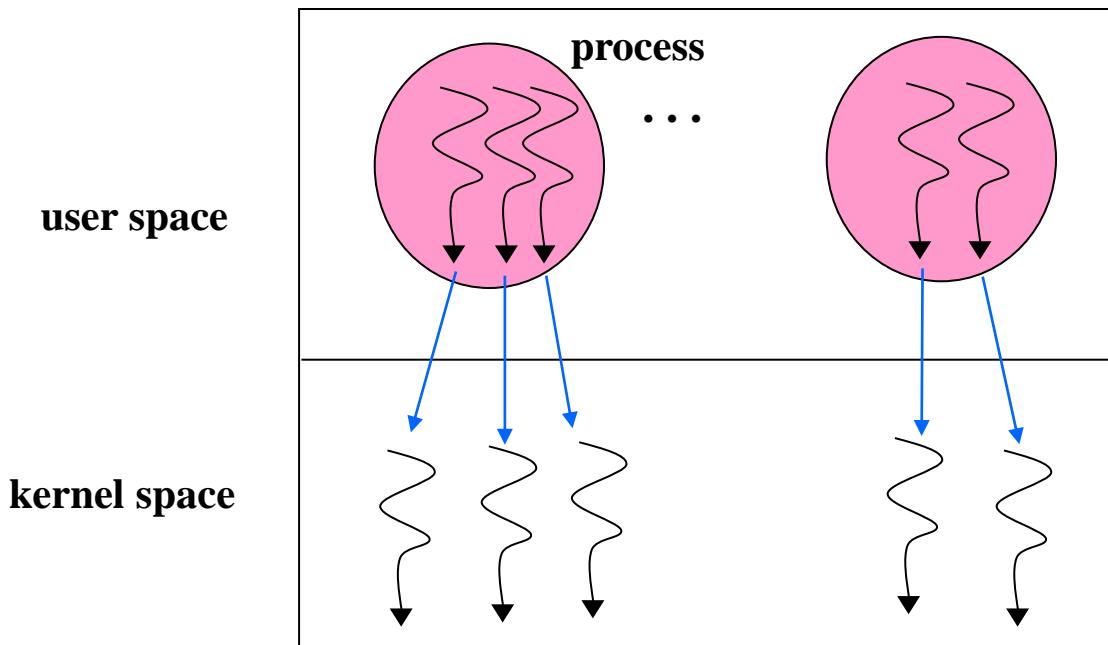
- **One application running five user-level threads.**
 - Performs inexpensive application-level context switches between threads.
- **Thread3 performs a system call for a file read().**
 - Kernel blocks the whole process, including the remaining four threads that are waiting to run.
- **Solution?**
 - Application intercepts system calls and substitutes an asynchronous (nonblocking) call.
 - An asynchronous interface call appears to return control to the caller immediately. Results will arrive later.
 - The remaining threads can now be run while thread3 is blocked.
 - Page faults cannot be predicted and mitigated, however, and will also result in the whole process being blocked.



- Background
- Problem Statement
- The Approach
- Related Works
- User-level Thread
- Kernel-level Thread
- Scheduler Activation
- Implementation
- Experimental Evaluation
- Summary



Kernel-level Thread



- **Kernel threads directly schedule each processes threads onto physical processors.**
 - 1:1 arrangement, each application threads maps to one kernel thread.
 - Enjoy direct OS support, so lack the problems that hinder user-level thread reliability.
- **E.g.: Linux, IRIX, Windows NT.**
- **Too much overhead to use in many parallel programs.**
 - Performance isn't as good as that of user threads.



- **Kernel is aware of the threads.**
 - Scheduler can assign time to processes based on the number of threads they carry.
 - I/O, page faults, and interrupts are not a problem.
 - Supports the needs of any application that runs it.
- **Hundreds of times slower than user-level threads.**
 - Kernel threads use system calls to perform operations.
 - Switching between threads requires a full context switch.
 - Kernel threads are heavyweights, must support every need of an application.
 - Each kernel thread consumes kernel memory.



- Background
- Problem Statement
- The Approach
- Related Works
- User-level Thread
- Kernel-level Thread
- **Scheduler Activation**
- Implementation
- Experimental Evaluation
- Summary



Goal: System with the functionality of kernel threads and the performance and flexibility of user-level threads.

- Maps some no. N of application threads onto a no. M of kernel entities (usually $N > M$).
- N:M thread systems are more complicated and effective
- Implemented by associating groups of threads with single kernel entities.
- Both AIX and Solaris use N:M thread system.

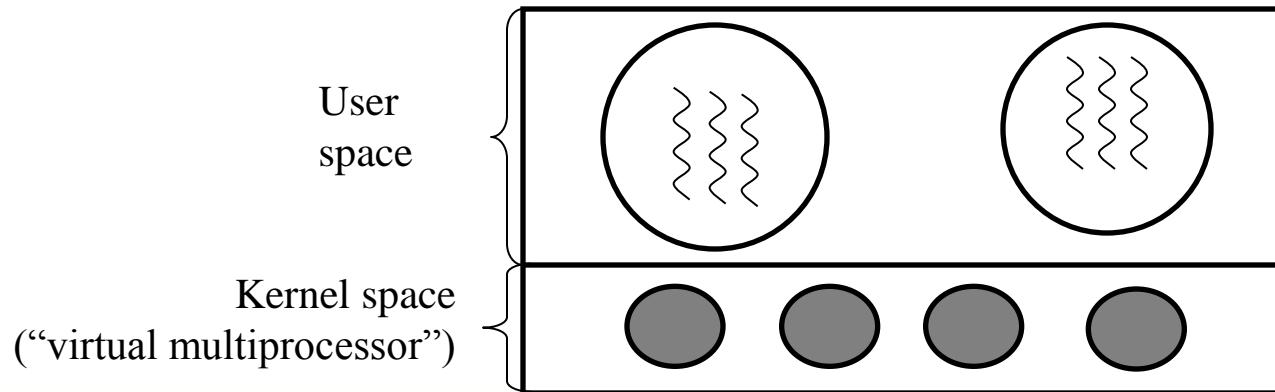


- Application has knowledge of the user-level thread state but has little knowledge of or influence over critical kernel-level events (by design! to achieve the virtual machine abstraction)
- Kernel has inadequate knowledge of user-level thread state to make optimal scheduling decisions

Solution: a mechanism that facilitates exchange of information between user-level and kernel-level mechanisms.

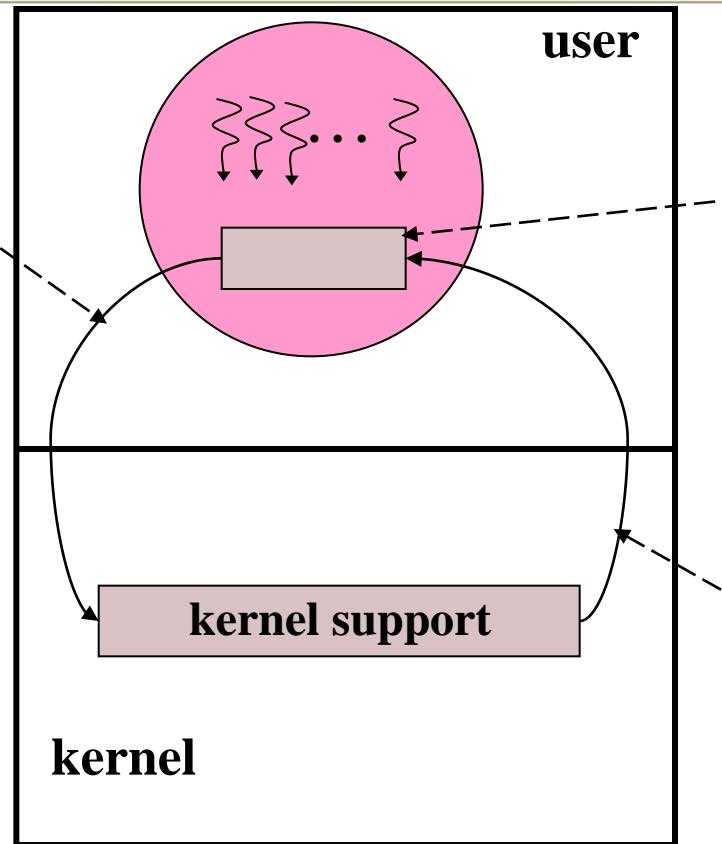


- Application knows how many and which processors allocated to it by kernel.
- Application has complete control over which threads are running on processors.
- Kernel notifies thread scheduler of events affecting address space.
- Thread scheduler notifies kernel regarding processor allocation.



- Change in processor requirements

Scheduler activations



thread
library

- change in processor allocation
- change in thread status



- **Structure to provide communication between the kernel processor and the user-level thread system.**
- **A scheduler activation:**
 - Takes the place of a kernel thread as an execution context for user-level threads.
 - Notifies the user-level thread system of kernel changes.
 - Provides space in the kernel for storing the processor context when the thread is swapped out.

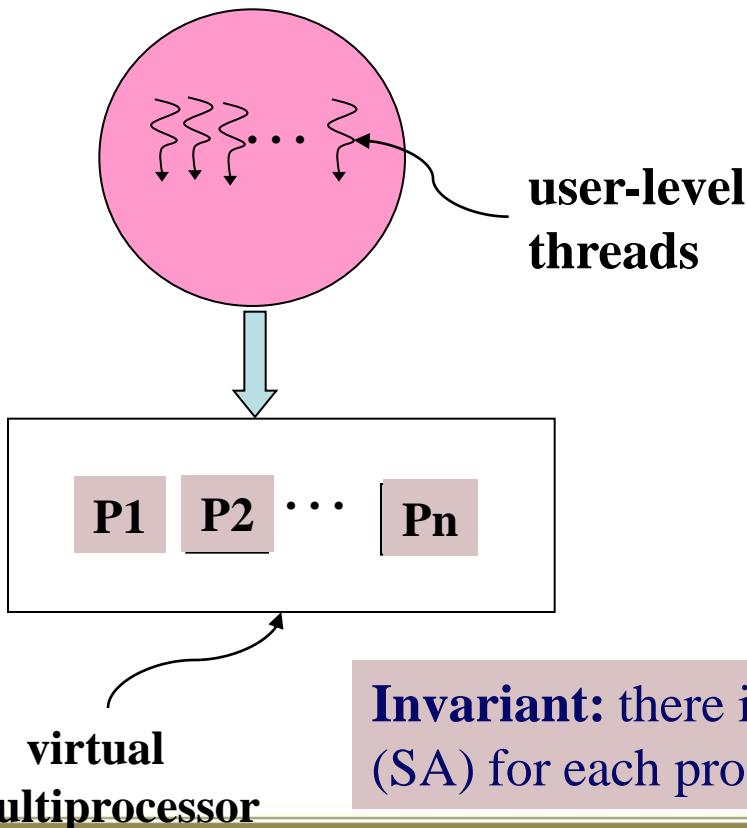


- **The kernel provides each user-level thread system with its own virtual processor.**
 - Kernel has control over how many processors to give.
 - Each user-level thread system chooses which threads to run on it's allocated processors.
 - Kernel notifies the user-level thread system of any changes, through the scheduler activation.
 - processors assigned, I/O interrupts, page faults.
- **One scheduler activation per processor assigned to address space.**

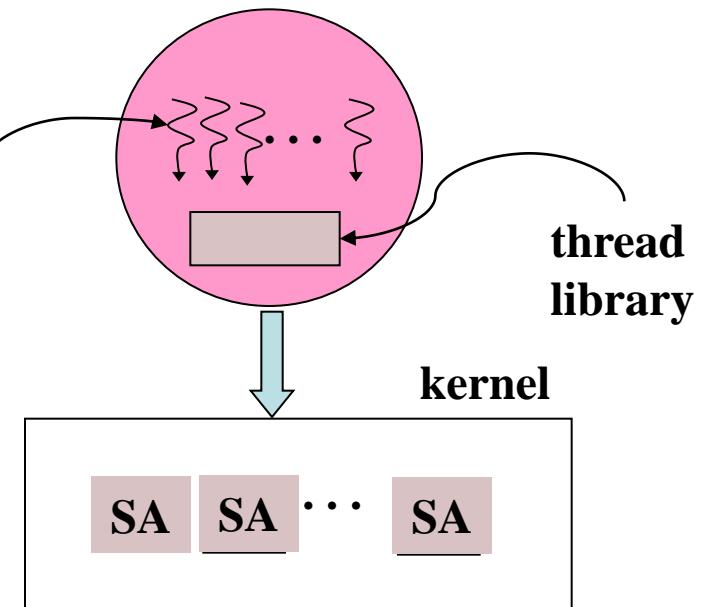


Role of Scheduler Activations

abstraction



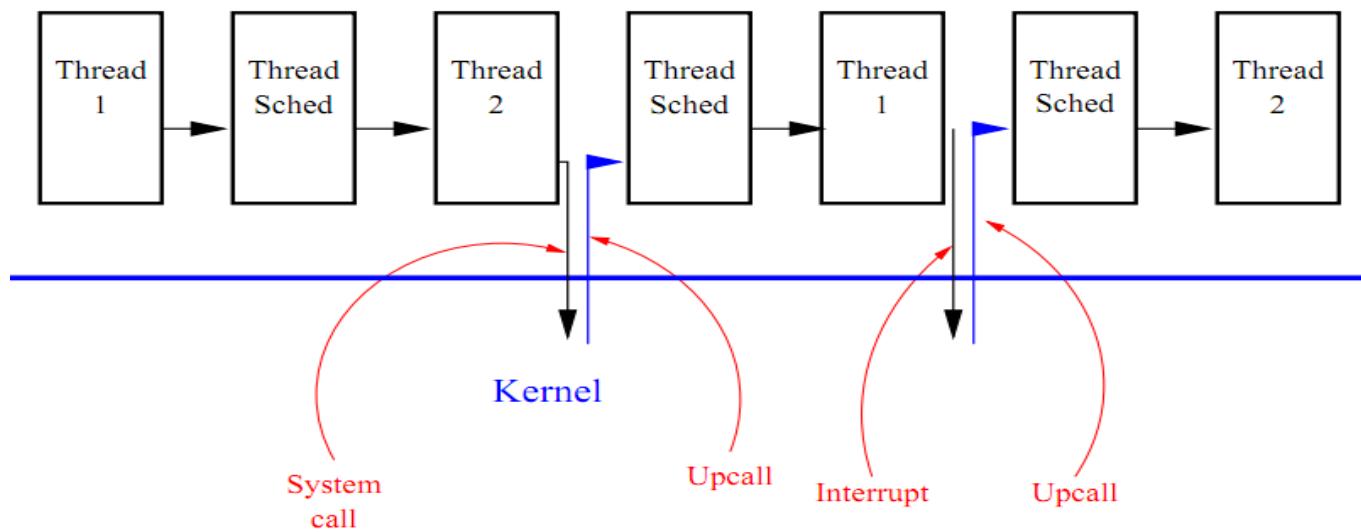
implementation



Invariant: there is one running scheduler activation (SA) for each processor assigned to the user process.



The kernel-level scheduler activation mechanism communicates with the user-level thread library by a set of upcalls:



- **Add this processor** (processor #)
 - *Execute a runnable user-level thread.*
- **Processor has been preempted** (preempted activation # and its machine state)
 - *Return to the ready list the user-level thread that was executing in the context of the preempted scheduler activation.*
- **Scheduler activation has blocked** (blocked activation #)
 - *The blocked scheduler activation is no longer using its processor.*
- **Scheduler activation has unblocked** (unblocked activation # and its machine state)
 - *Return to the ready list the user-level thread that was executing in the context*



Communication : Kernel to User Level

Notice that these four messages correspond neatly with this graph.

- It shows the different states a process or kernel thread can be in: running, ready to run, or blocked.
- Whenever a kernel thread system would change the state of a thread "transparently", scheduler activations contact the user-level scheduler.

It can update its information and rethink.

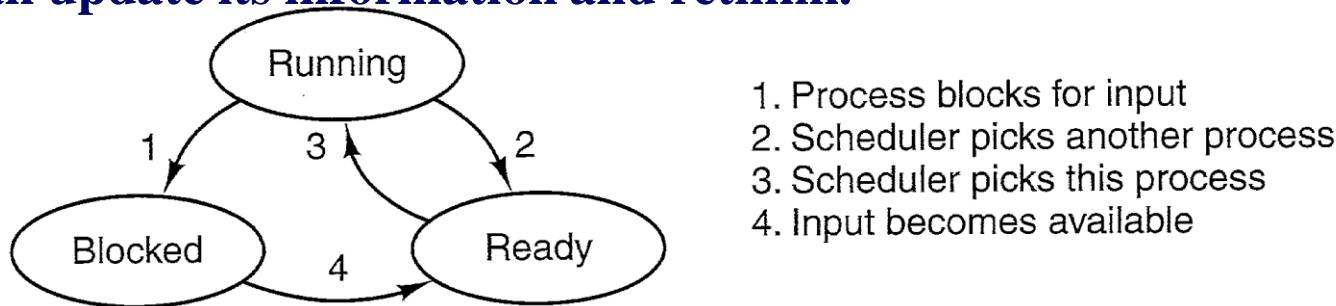
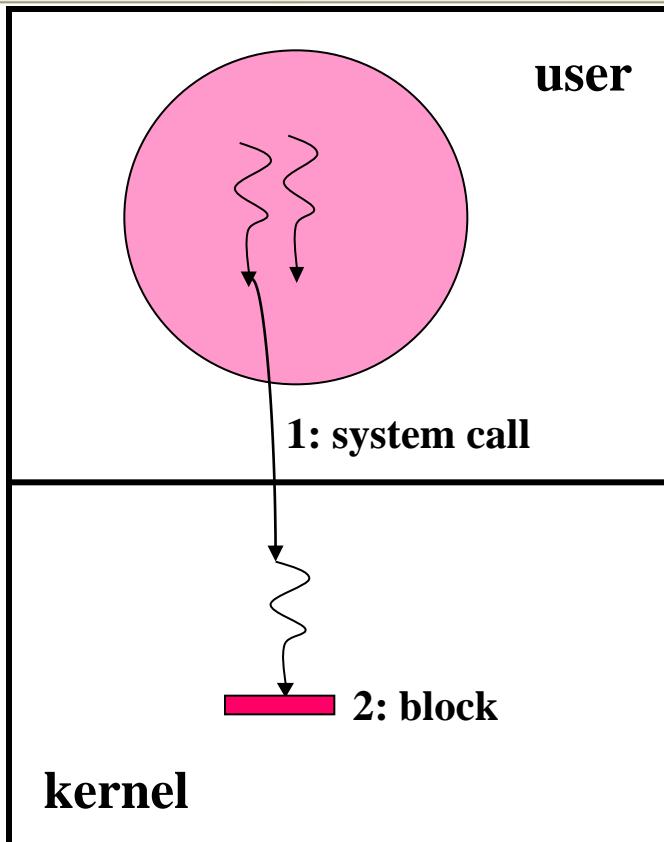


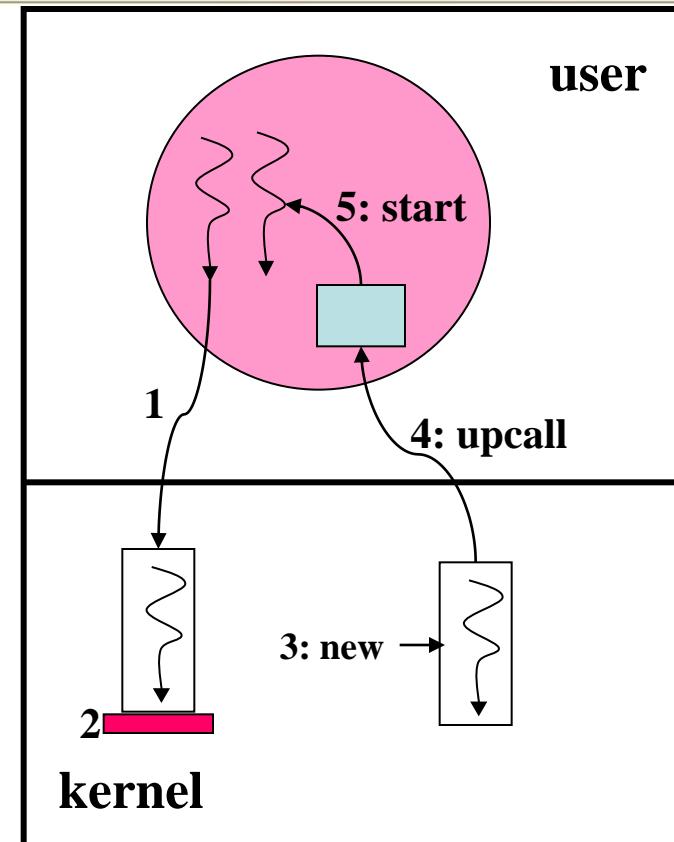
Figure 2-2. A process can be in running, blocked, or ready state. Transitions between these states are as shown.



Avoid Effects of Blocking



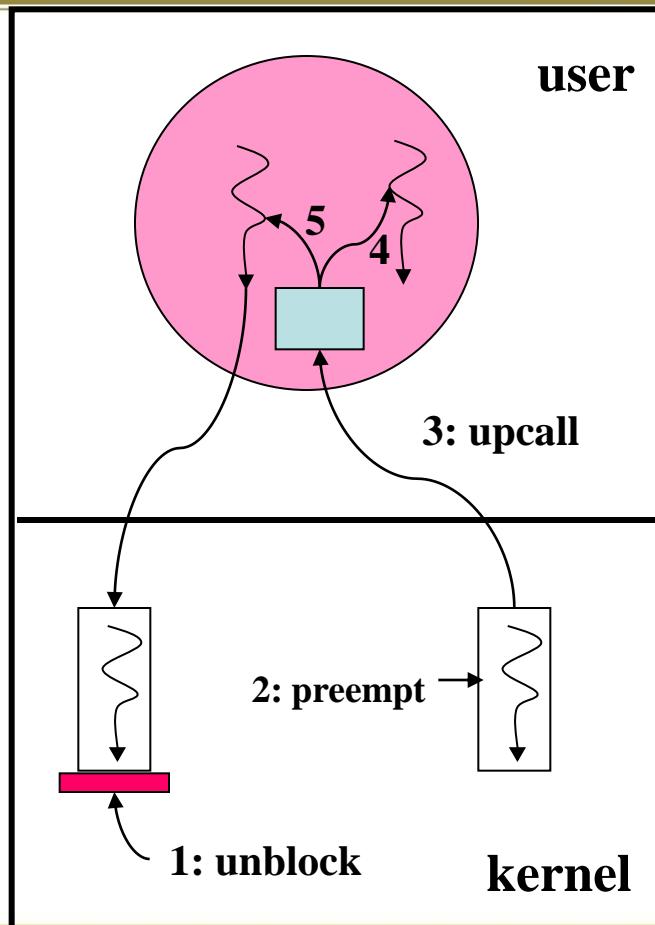
Kernel threads



Scheduler Activations



Resuming Blocked Thread



4: preempt
5: resume



- **#Runnable Threads != # Processors**
- **Add more processors (additional # of processors needed)**
 - Allocate more processors to this address space and start them running scheduler activations.
- **This processor is idle ()**
 - Preempt this processor if another address space needs it.



- **A user-level thread can be preempted while it is executing a critical section.**
 - Could result in deadlock or performance drop.
- **If a thread that has been preempted was executing a critical section, a user-level context switch is necessary.**
 - Thread continues temporarily, finishes its critical section, and then relinquishes control back to the upcall.
- **Implemented by making a copy of each low-level critical section.**
 - In the copy, code is added to yield the processor back to the resuming thread.
 - If a preemption occurs, the kernel resumes the thread in the corresponding place in the copy.



Example: Managing I/O

T1: The application receives two processors to run on from the kernel.

It starts running threads on them.

T2: One of the user-level threads blocks in the kernel.

Kernel uses a fresh scheduler activation to notify the user-level system of this.

T3: The blocking-event completes. Kernel preempts second processor to do the upcall.

Upcall notifies user level of upcall and of event completion.

T4: Upcall takes the thread from the ready list and begins running it.

The first two scheduler activations have now been discarded and replaced.

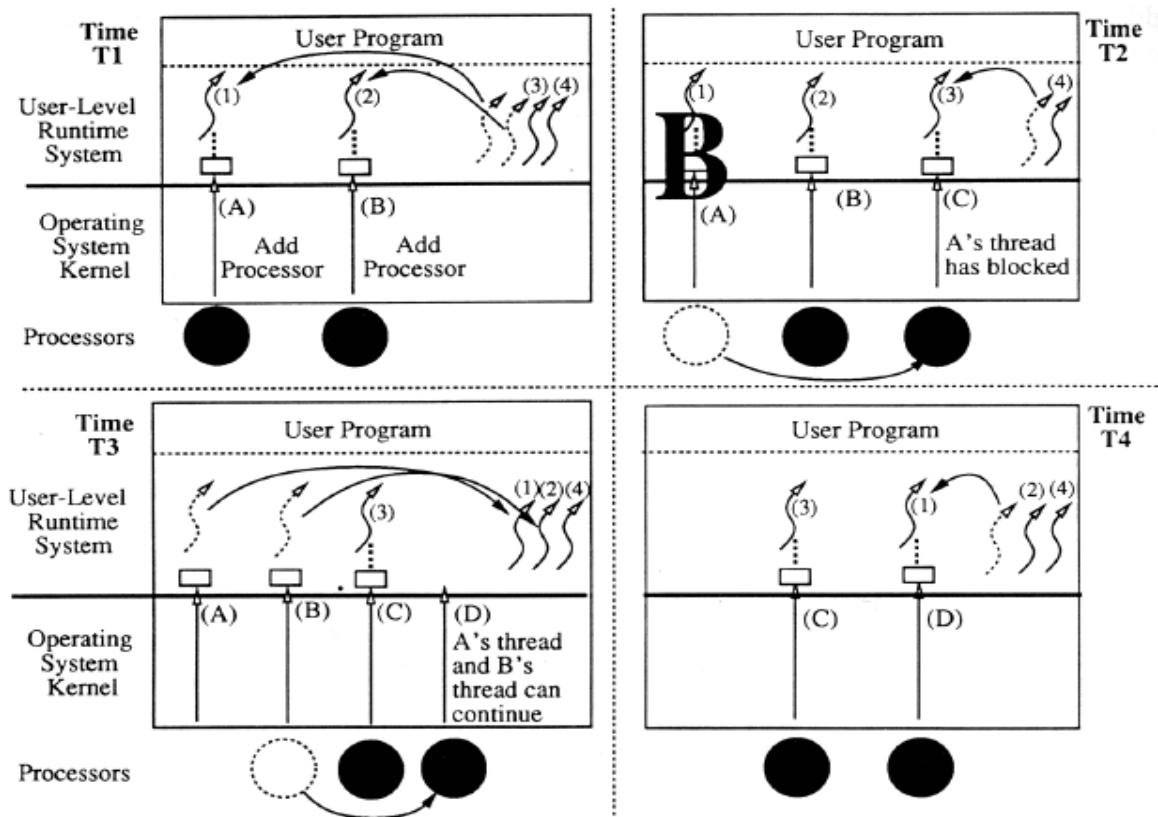


Fig. 1. Example: I/O request/completion.



- Background
- Problem Statement
- The Approach
- Related Works
- User-level Thread
- Kernel-level Thread
- Scheduler Activation
- Implementation
- Experimental Evaluation
- Summary



- Modification of Topaz Operating System
- Topaz is OS for DEC SRC Firefly a multi-processor workstation
- Modification of FastThreads a user level thread package
- Lines of code used



Implementation Sections

- Processor Allocation Policy
- Thread Scheduling Policy
- Performance Enhancement
- Debugging Considerations



- Similar to dynamic policy of Zahorjan and McCann using Space-Shares
- Possible for address space to use kernel threads
- Support for Topaz kernel threads
- No need for static partitioning of processors



- Kernel has not knowledge of an application's concurrency model or scheduling policy
- Applications has freedom to decide how to handle scheduling
- FastThreads use of per processor ready lists with Last In First out order



- Hysteresis used to avoid unnecessary processor reallocation
- Hysteresis Example



- Critical sections
 - Typical solution using flags. Disadvantage: has overhead
 - Solution used, similar to Trellis/Owl garbage collector. Advantage: no overhead in common case



- Management of scheduler activations
 - Unused scheduler activations can be cached for reuse
 - Scheduler activation sent in bulk for reuse



- Firefly Topaz Debugger
 - User level
 - Application



- Background
- Problem Statement
- The Approach
- Related Works
- User-level Thread
- Kernel-level Thread
- Scheduler Activation
- Implementation
- Experimental Evaluation
- Summary



- Three questions
 - What is the cost of user-level thread operations
 - What is the cost of communication between kernel and the user level
 - What is the overall effect on the performance of applications



Experimental Evaluation Sections

- Thread Performance
- Upcalls
- Application Performance



- Cost of user threads, same as FastThreads

Operation	FastThreads on		FastThreads on	
	Topaz threads	Scheduler Activations	Topaz threads	Ultron processes
Null Fork	34	37	948	11300
Signal-Wait	37	42	441	1840

Table 4: Thread Operation Latencies (μ sec.)



- Measure of Upcall performance
 - Two user-level threads to signal and wait



- Parallel Application using Topaz kernel threads
- Setup of application and modifications
- Procedure and methods.
- Comparison to FastThreads



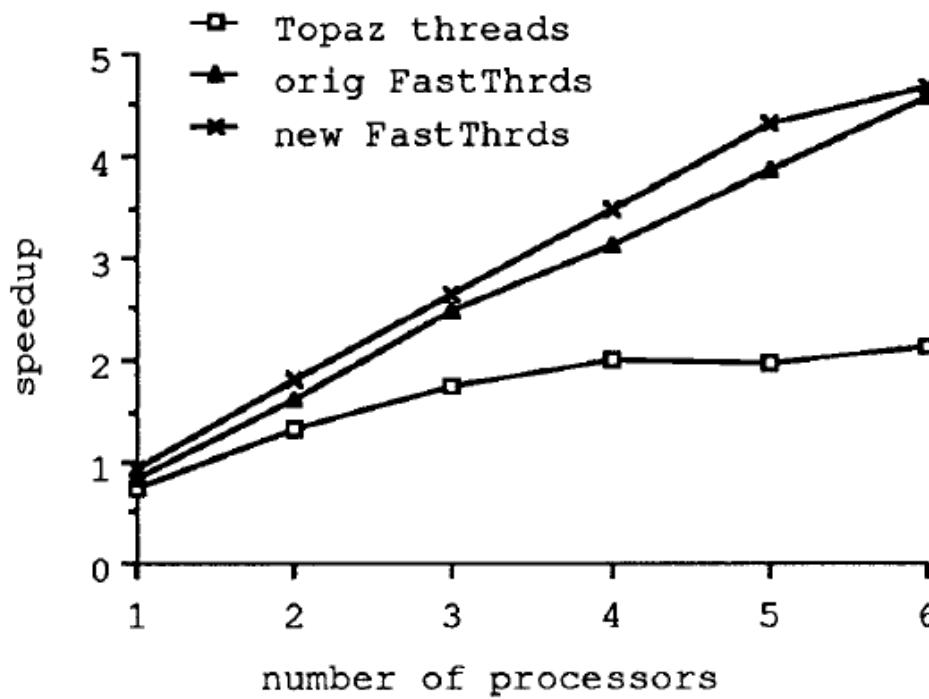


Figure 1: Speedup of N-Body Application vs. Number of Processors, 100% of Memory Available



- Effect of performance application-induced kernel events
- Better performance than FastThreads



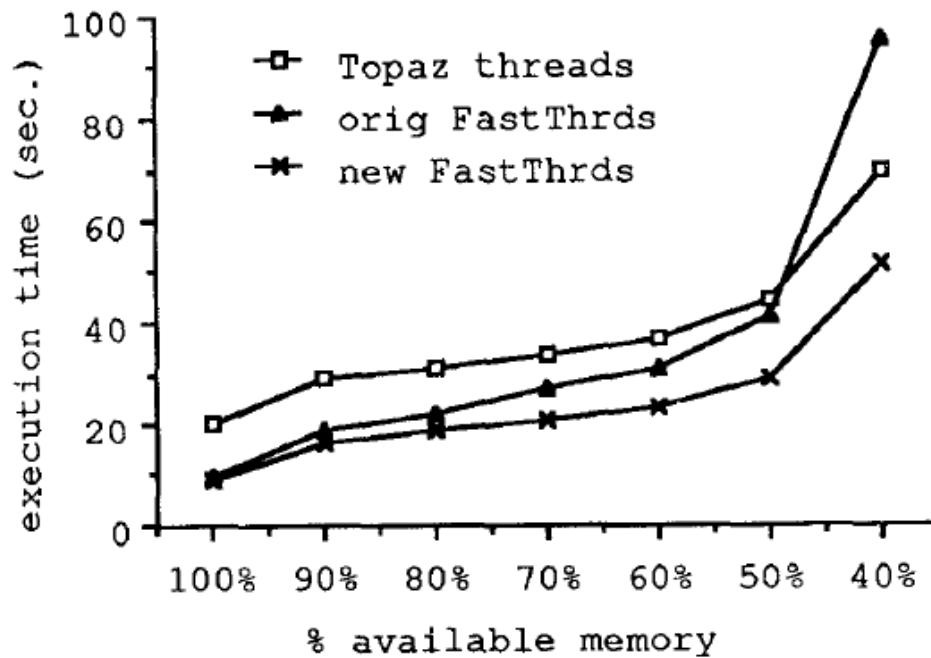


Figure 2: Execution Time of N-Body Application vs.
Amount of Available Memory, 6 Processors



Topaz threads	Original FastThreads	New FastThreads
1.29	1.26	2.45

Table 5: Speedup for N-Body Application, Multiprogramming Level = 2, 6 Processors, 100% of Memory Available



- Background
- Problem Statement
- The Approach
- Related Works
- User-level Thread
- Kernel-level Thread
- Scheduler Activation
- Implementation
- Experimental Evaluation
- **Summary**



- Implement user-level thread package
- Virtual Multiprocessor in application address space
- Divided responsibilities



- Processor allocation
- Thread Scheduling
- Kernel notifications to address space thread scheduler
- Address space notifications to kernel for processor allocations



- Questions?

