# English Edition

# Scheduling (5)

---

## Multiprocessor Thread Scheduling (1)

- Load Sharing
  - no (permanent) assignment of CPUs to processes
- Gang Scheduling
  - a number of threads (that belong together) gets an according number of CPUs by the scheduler simultaneously

---

## Threads

- When handling threads: Treat execution aspect and the other process properties separately
- Application can consist of a number of threads which cooperate and simultaneously run in the same address space
- Threads (of one process) that run on different CPUs can dramatically increase process performance

---

## Multiprocessor Thread Scheduling (2)

- Dedicated CPU allocation
  - threads are a assigned to a specific CPU
- Dynamic scheduling
  - number of threads can change during program execution

# Load Sharing

- distribute compute-load uniformly across all CPUs
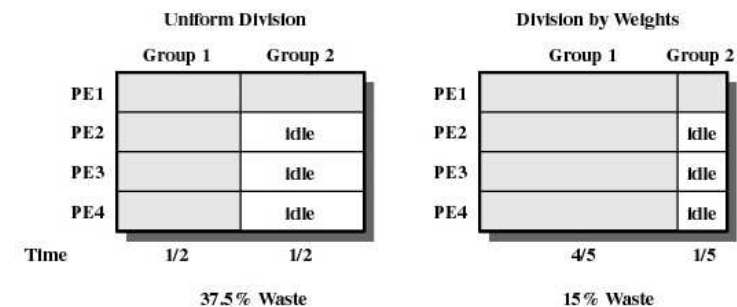- no centralized scheduler needed
- uses global queue(s)

# Gang Scheduling

- simultaneous scheduling of Threads which belong to one process
- useful for applications in which performance suffers immensly when any part (thread) is not active
- threads must often synchronize
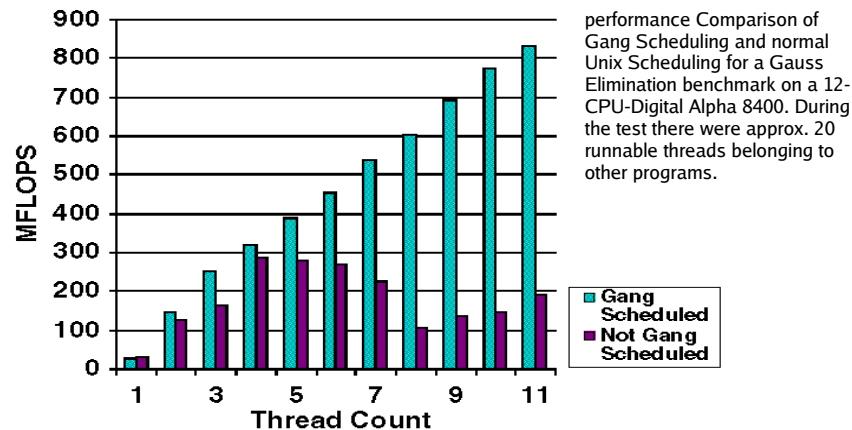
# Disadvantages of Load Sharing

- central queue requires *mutual exclusion*
  - can be a bottleneck when more than one CPU searches for work at the same time
- preemptive threads are likely to change the CPU often
  - cache usage is less efficient
- when all threads are in the global queue, it is unlikely that all threads of one process are scheduled simultaneously

# Scheduling Groups



Scheduling groups with 4 / 1 threads

# Gang Scheduling: Performance



performance Comparison of Gang Scheduling and normal Unix Scheduling for a Gauss Elimination benchmark on a 12-CPU-Digital Alpha 8400. During the test there were approx. 20 runnable threads belonging to other programs.

Picture: http://www.llnl.gov/asci/pse_trilab/sc98.summary.html

---

# Dedicated CPU Allocation (2)



**Situation:**
- 16-CPU-machine
- two compute processes (FFT: Fast Fourier Transformation, Matmul: matrix multiplication) can split their calculations into variably many threads
- As long as they create less than 8 threads (i.e. total # threads < # CPUs), performance increases – afterwards no longer

**Figure 10.3**   Application Speedup as a Function of Number of Processes [TUCK89]

---

# Dedicated CPU Allocation (1)

- when the scheduler picks a process, it schedules each of its threads on one CPU
- some CPUs may be idle
- mo multitasking on the CPUs

---

# Dynamic Scheduling

- number of threads (of one process) can change dynamically: process creates new threads and destroys them – *depending on the resources* that the operating system provides
- operating system adjusts the load in order to optimize CPU utilization
  - idle CPUs are scheduled to run threads
  - new processes can receive a CPU that is currently used by a process having several CPUs
  - ignore request until a CPU becomes available
  - new processes receive a CPU with higher priority than already running applications do

## Example: K42 (1)

- Linux-compatible OS for machines with hundreds of CPUs
- object-oriented; OS system calls are IPC calls
- two scheduler levels
- thread scheduling: completely handled in user mode
- kernel scheduler
  - manages so called dispatchers which in turn manage the threads (often: one dispatcher per process)
  - runs independently on each CPU



Applications
- Application Object
- Application Object

Linux API/ABI
- Linux Libraries/glibc
- Linux Emulation

K42 Operating System Libraries

File Server | Name Server | Linux File Server

Kernel
- Memory Manager
  - File Cache Manager
  - File Cache Manager
- Linux Device Drivers/Internet Protocol Stack

## Example: K42 (2)



Processor 3, Processor 2, Processor 1, Processor 0
P, Q, R, ... a, b, c, f, g
Dispatchers, Process, Resource Domains
S, d, e

## Example: K42 (3)

- Thread Migration: K42 thread scheduler does Load Balancing, by migrating threads from busy dispatchers to idle ones
- kernel may, as an exception, also migrate dispatchers to a different CPU
- kernel scheduler activates resource domains (inside one domain: dispatchers in a circle)
- does Gang Scheduling (see processes P,R)
- dispatcher can cope with single threads blocking (e.g. for page faults or I/O), without also blocking

## Example: K42 (4)

- processes have a choice:
  - real multitasking -> use several dispatchers
  - only programming comfort of threads -> one dispatcher is enough
- miscellaneous thread libraries for programmers, including POSIX threads
- K42 scheduler: http://www.research.ibm.com/K42/white-papers/Scheduling.pdf (2002)
- introduction to K42: http://www.research.ibm.com/journal/sj/442/appavoo.pdf

# Linux O(1) Scheduler

---

## Linux O(1) Scheduler   (2)

### Causes (in kernel 2.4)

- one common queue for all processes on all CPUs; no sorting in that queue

- scheduler must search whole queue in order to find the next process to schedule

- one single lock for the runqueue
  -> one CPU accessing the queue blocks access for all further CPUs

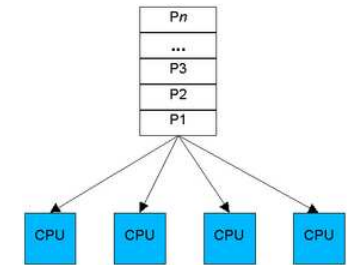- result: schedule action very complex

Bild: Linux Journal,
http://www.linuxjournal.com/node/7178/

---

## Linux O(1) Scheduler   (1)

- change with Linux kernel 2.6:
  new scheduler which remedies some problems of the old 2.4.x scheduler:
  - schedule time was (linearly) dependent on the number of processes, i.e. O(n)
    -> poor performance with very many processes
  - poor performance on SMP machines

---

## Linux O(1) Scheduler   (3)

### Kernel 2.4

- processes not bound to a CPU, assignments random (no processor affinity)
  -> processes change CPUs regularly
  -> poor utilization of CPU caches

# Linux O(1) Scheduler   (4)

**Kernel 2.6: new O(1) Scheduler with the following features:**

- O(1) scheduler: time required for selecting the next process (for one CPU) is constant – independent on the number of processes
- CPUs don't block one another in case of simultaneous scheduling decisions
- load balancer distributes compute-load uniformly across several CPUs