

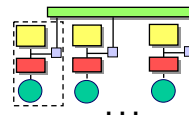
Scheduling

Jesus Labarta

Scheduling

■ Applications submitted to system → Resources x Time

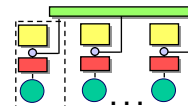
- Resources:
 - ✓ Processors
 - ✓ Memory
 - ✓ ...



■ Objective

- Maximize resource utilization
- Maximize throughput
- Minimize response time

} Not necessarily the same

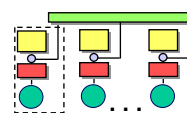
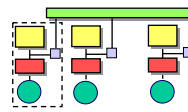




Scheduling

■ Applicable policies depend on

- Application structure
 - ✓ Fixed
 - ✓ Moldable
 - ✓ Malleable
- Resources
 - ✓ Level of sharing support
 - ✓ Level of configurability



Malleability

■ Dynamically adaptable parallelization structure

■ Responsiveness to

- Resource availability: Processors, memory,....
- Application behavior dependence on input data, time,

■ Malleability requires:

- Separation between
 - ✓ Algorithm: Problem logic. Full responsibility / Job of the programmer
 - ✓ Scheduling: Efficiency Responsibility of the system. Programmer hints
- Frequent control points

■ Issue of:

- Programming model support → OpenMP
- Programming practices





Programming practices and malleability

Scheduling decisions: Once for all

Explicit code only related to parallelism

```

C$OMP PARALLEL
WhoAmI=RunTimeCall()
myBlock=f(WhoAmI)
...
Call Compute1(myBlock)
...
DO iters=1, #iters
  Call Compute2(myBlock)
END DO
...
C$OMP END PARALLEL

```

```

C$OMP PARALLEL DO
DO Block=1, #blocks
  Call Compute1(Block)
END DO
C$OMP END PARALLEL
...
DO iter=1, #iters
  C$OMP PARALLEL DO
  DO Block=1, #blocks
    Call Compute1(Block)
  END DO
  C$OMP END PARALLEL
END DO
...
C$OMP END PARALLEL

```

$myBlock \in [1, \#processors] \leftarrow f(resources)$
 $Block \in [1, \#blocks] \leftarrow f(algorithm)$



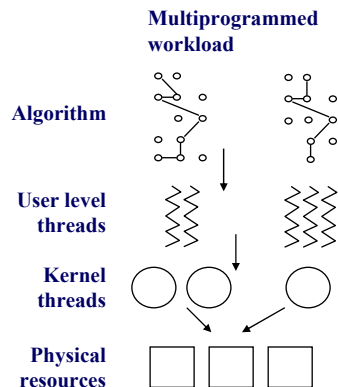
Scheduling

Hierarchy

- Grid
- Jobs scheduling
- Process scheduling
 - ✓ Computation task
 - ✓ Processes
 - ✓ User threads
 - ✓ Kernel threads
 - ✓ Processors

Typical

- Very simple approaches at many levels
 - ✓ 1 to 1
- One most relevant level
- Lack of coordination





Job scheduling

- **Control of Multiprogramming level**
 - Static
 - Dynamic
- **Long term scheduling (minutes – hours – days)**
- **Coarse grain allocation of resources**
- **Basic operation modes**
 - Batch queuing systems
 - Interactive (≡ no job scheduling)



Job scheduling

- **Basic approaches**
 - Fixed Partitions
 - ✓ Different sizes of processors, memory
 - ✓ One queue per partition
 - ✓ Internal fragmentation
 - Variable partitions
 - ✓ Single queue
 - ✓ External fragmentation





Job scheduling

■ Basic policies

- FIFO
- SJF
- Backfilling

■ Requirements

- Estimates of job duration
- Estimates of job performance vs. resources
- User supplied / prediction



Job scheduling

■ Typical setups

- Different queues
 - Memory, duration,...
- User submits to queue that fits estimate of job needs
- A lot of manual supervision by operator, modifying job limits,....





Processor scheduling

■ Mapping of kernel threads to processors

■ Basic approaches

- Time sharing
- Space sharing
- Gang scheduling

■ Types of policies

- Fixed
- Adaptive
- Dynamic



Processor scheduling

■ Time sharing

- Extensions from sequential
- Multiprocessor specific effects/features
 - ✓ Applied to individual kernel threads
 - Fairness issues → fair share schedulers
 - ✓ Load balance vs. locality trade off
 - Local queues
 - Binding / Migration
 - Global queues
 - Contention / Affinity



Processor scheduling

■ Time sharing

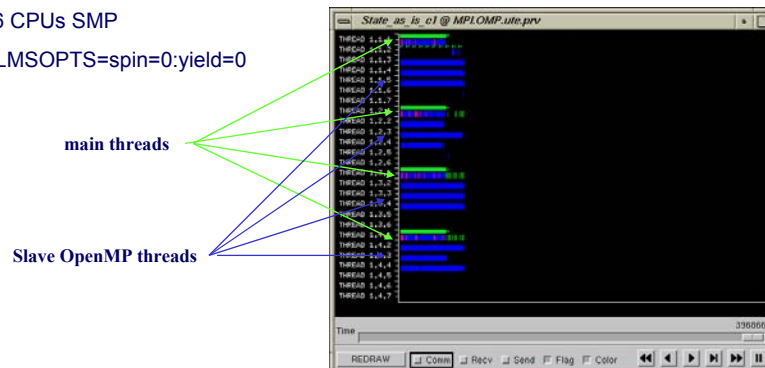
- Multiprocessor specific effects/features
 - ✓ Throughput vs. response time
 - Capability to fill all gaps → more threads than processors
 - Interaction between OS scheduling and application synchronization → high risk of inefficiency → threads == processors



Motivation example

■ Sweep3d 4 MPI tasks with 4 OpenMP threads each

- 16 CPUs SMP
- XLMSOPTS=spin=0:yield=0





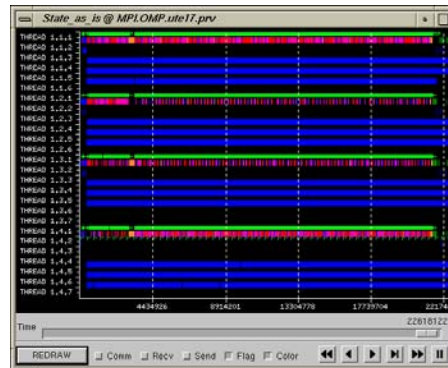
Motivation example

■ Sweep3d 4 MPI tasks with 4 OpenMP threads each

- 16 CPUs SMP
- XLMSOPTS=spin=0:yield=0

■ + 1 sequential program

- Non instrumented



Same scale as previous slide

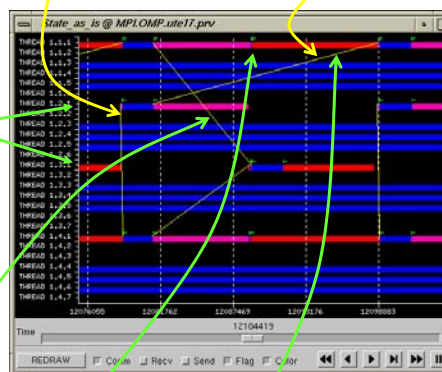


Motivation example

- Quantum == 10 ms
- Intricate precedence and scheduling interferences

2 threads time sharing CPU

Very fast communication Very slow communication



Can not finish because receiver descheduled

Receive called after finishing send

Can not finish because sender descheduled



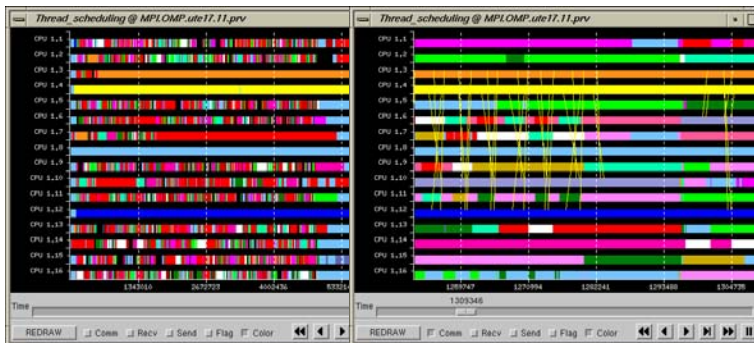
Motivation example

■ Thread scheduling

- XLMSOPTS=spin=1:yield=1

Lots of process migrations of
OpenMP slave threads

CPU bursts \ll Quantum



Jesús Labarta, MP, 2002



Processor scheduling

■ Space sharing

- Partition of processors among applications
 - ✓ Possible fragmentation
- Example policies
 - ✓ Equipartition
 - ✓ Equal efficiency

■ Gang Scheduling

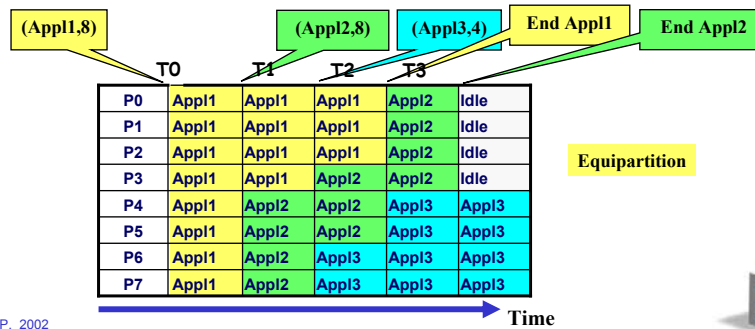
- Coarse grain time sharing

Jesús Labarta, MP, 2002



Equipartition

- Equal allocation among running jobs [McCann93]
 - ✓ Allocation = $P / \text{num_applications}$
- Re-allocation at each job arrival and completion
- Reasonable approach if job performance not known

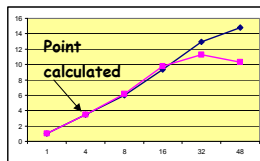


Jesús Labarta, MP, 2002



Equal_efficiency

- Dynamic space-sharing policy [Nguyen96]
- Job performance analysis + Eff. extrapolation
- Assigns processors (one per turn) to those applications that achieve the higher efficiency



$$\text{Efficiency}(p) = \frac{(1 + \beta)}{(p + \beta)}$$

P0	App1	App1	App1	App1	App1	App2	Idle
P1	App1	App1	App1	App1	App1	App2	Idle
P2	App1	App1	App2	App2	App2	App2	Idle
P3	App1	App1	App2	App2	App2	App2	Idle
P4	App1	App2	App2	App2	App2	App2	App3
P5	App1	App2	App2	App3	App2	App2	App3
P6	App1	App2	App2	App3	App3	App3	App3
P7	App1	App2	App2	App3	App3	App3	App3

Time



High Efficiency



Medium Efficiency



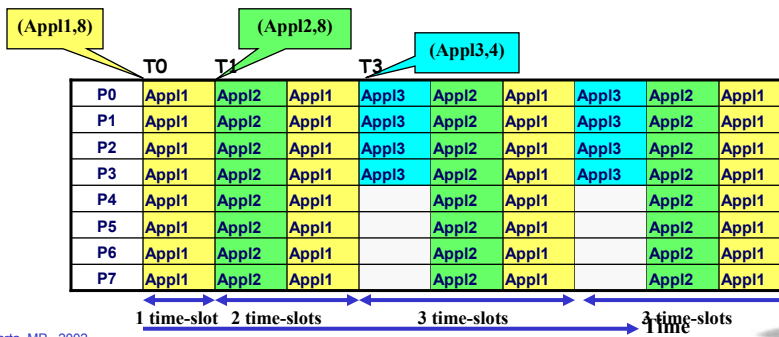
Low Efficiency

Jesús Labarta, MP, 2002



Gang Scheduling

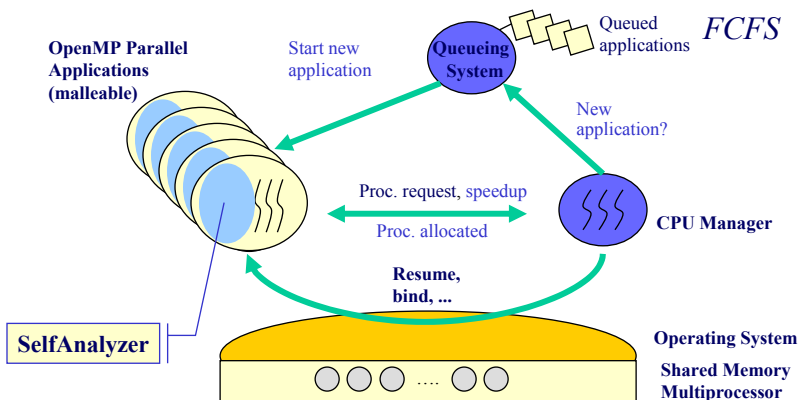
- Threads are grouped into gangs
- Threads in a gang are executed simultaneously
- Time-sharing is used among gangs



Jesús Labarta, MP, 2002



NANOS OS scheduling environment



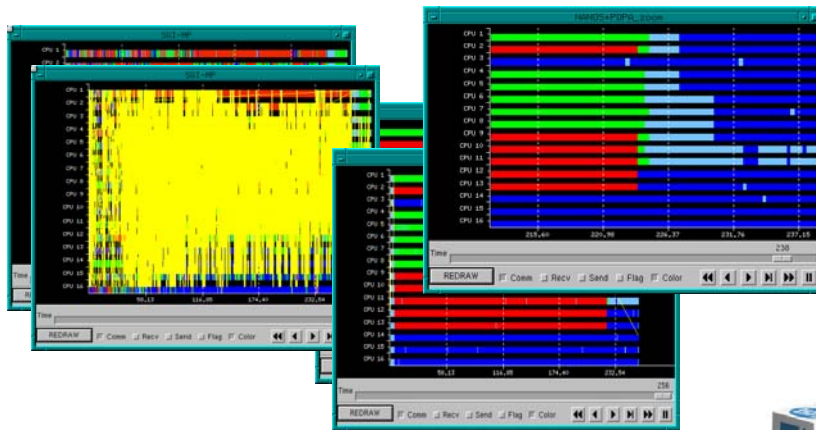
Jesús Labarta, MP, 2002





Native time sharing vs. PDPA

- 3 BTs x 8 processes @ 16 CPUs



Jesús Labarta, MP, 2002



Page placement

- Mapping of pages to nodes
- UMA
 - Not an issue
- NUMA
 - Round Robin - Random
 - First touch
 - Page migration

Jesús Labarta, MP, 2002





Page placement

Should the programmer specify data distribution ?

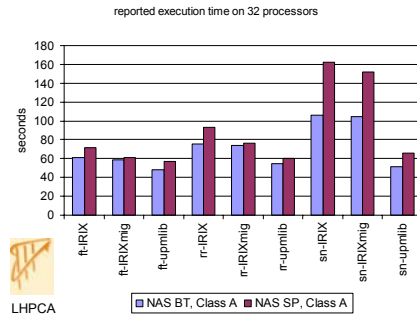
- Difficult problem → Automatic

Approach

- Application run time measures
- Computes “good” distribution
- Asks the system to redistribute

Provides

- Responsiveness
- Local view → Stability



Jesús Labarta, MP, 2002



Page placement & OS scheduling

Native page migration enabled

disabled

