

Aim of Scheduling

- Assign processes to be executed by the processor(s)
 - Response time
 - Throughput
 - Processor utilization
 - Tardiness etc.

1

Scheduling Environments

- Single vs. multiple processors
- Static vs. dynamic process arrival
- Preemptive vs. nonpreemptive
- Independent vs. dependent tasks
- etc.

2

Table 9.1 Types of Scheduling

Long-term scheduling	The decision to add to the pool of processes to be executed
Medium-term scheduling	The decision to add to the number of processes that are partially or fully in main memory
Short-term scheduling	The decision as to which available process will be executed by the processor
I/O scheduling	The decision as to which process's pending I/O request shall be handled by an available I/O device

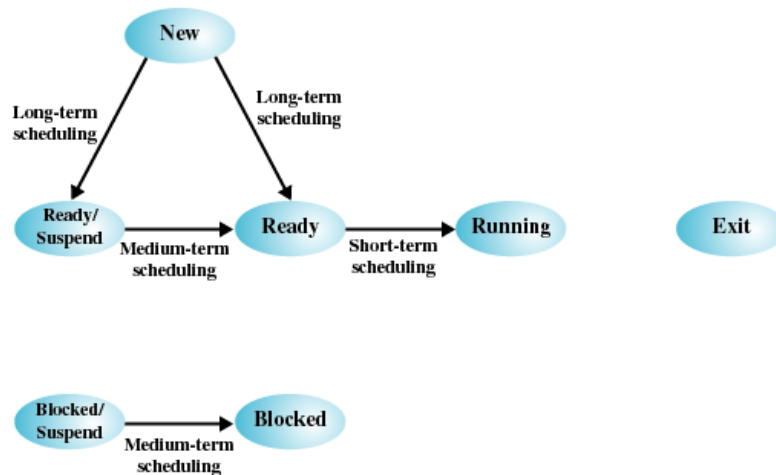


Figure 9.1 Scheduling and Process State Transitions

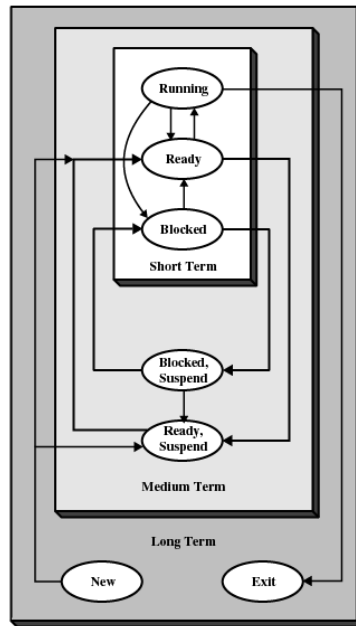


Figure 9.2 Levels of Scheduling

5

Long-Term Scheduling

- Determines which programs are admitted to the system for processing
- Controls the degree of multiprogramming
- More processes, smaller percentage of time each process is executed

6

Medium-Term Scheduling

- Part of the swapping function
- Based on the need to manage the degree of multiprogramming

7

Short-Term Scheduling

- Known as the dispatcher
- Executes most frequently
- Invoked when an event occurs
 - Clock interrupts
 - I/O interrupts
 - Operating system calls
 - Signals

8

Short-Term Scheduling Criteria

- User-oriented
 - Response Time
 - Elapsed time between the submission of a request until there is output.
- System-oriented
 - Effective and efficient utilization of the processor

9

Short-Term Scheduling Criteria

- Performance-related
 - Quantitative
 - Measurable such as response time and throughput

10

Table 9.2 Scheduling Criteria

User Oriented, Performance Related

Turnaround time This is the interval of time between the submission of a process and its completion. Includes actual execution time plus time spent waiting for resources, including the processor. This is an appropriate measure for a batch job.

Response time For an interactive process, this is the time from the submission of a request until the response begins to be received. Often a process can begin producing some output to the user while continuing to process the request. Thus, this is a better measure than turnaround time from the user's point of view. The scheduling discipline should attempt to achieve low response time and to maximize the number of interactive users receiving acceptable response time.

Deadlines When process completion deadlines can be specified, the scheduling discipline should subordinate other goals to that of maximizing the percentage of deadlines met.

User Oriented, Other

Predictability A given job should run in about the same amount of time and at about the same cost regardless of the load on the system. A wide variation in response time or turnaround time is distracting to users. It may signal a wide swing in system workloads or the need for system tuning to cure instabilities.

11

System Oriented, Performance Related

Throughput The scheduling policy should attempt to maximize the number of processes completed per unit of time. This is a measure of how much work is being performed. This clearly depends on the average length of a process but is also influenced by the scheduling policy, which may affect utilization.

Processor utilization This is the percentage of time that the processor is busy. For an expensive shared system, this is a significant criterion. In single-user systems and in some other systems, such as real-time systems, this criterion is less important than some of the others.

System Oriented, Other

Fairness In the absence of guidance from the user or other system-supplied guidance, processes should be treated the same, and no process should suffer starvation.

Enforcing priorities When processes are assigned priorities, the scheduling policy should favor higher-priority processes.

Balancing resources The scheduling policy should keep the resources of the system busy. Processes that will underutilize stressed resources should be favored. This criterion also involves medium-term and long-term scheduling.

12

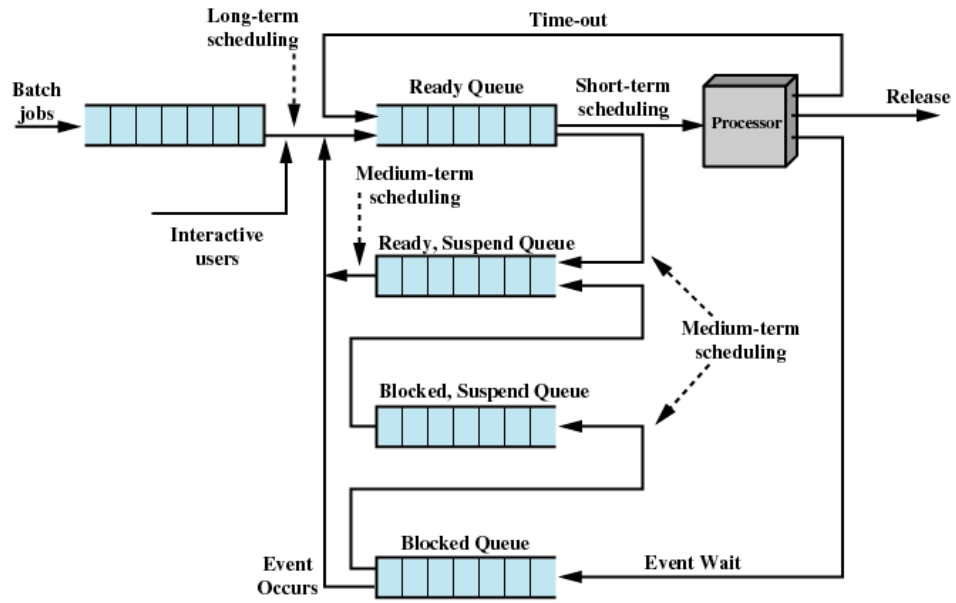


Figure 9.3 Queuing Diagram for Scheduling