

# RESOURCE MANAGEMENT



# Introduction

2

- Distributed systems contain a set of resources interconnected by a network
- Processes are migrated to fulfill their resource requirements
- Resource manager controls the assignment of resources to processes
- Resources can be logical (shared file) or physical (CPU)
- We consider a resource to be a processor

# Types of process scheduling techniques

3

- ❑ Task assignment approach
  - ❑ User processes are collections of related tasks
  - ❑ Tasks are scheduled to improve performance
- ❑ Load-balancing approach
  - ❑ Tasks are distributed among nodes so as to equalize the workload of nodes of the system
- ❑ Load-sharing approach
  - ❑ Simply attempts to avoid idle nodes while processes wait for being processed

# Desirable features of a scheduling algorithm

4

- No A Priori Knowledge about Processes
  - ▣ Scheduling algorithm should operate without prior knowledge about processes to be executed
  - ▣ User does not want to specify information about characteristics and requirements
- Dynamic in nature
  - ▣ Decision should be based on the changing load of nodes and not on fixed static policy
  - ▣ Scheduling decision should be based on current load
  - ▣ Should support preemptive scheduling

# Desirable features of a scheduling algorithm

5

- Quick decision-making capability
  - ▣ Algorithm must make quick decision about the assignment of task to nodes of system
  - ▣ heuristic methods requiring less computational effort are preferable to find near optimal solution
  
- Balanced system performance and scheduling overhead
  - ▣ Great amount of information gives more intelligent decision, but increases overhead

# Desirable features of a scheduling algorithm

6

## □ Stability

- ▣ Scheduling algorithm is unstable when all processes are migrating without accomplishing any useful work, which is called as processor thrashing
- ▣ It occurs when the nodes turn from lightly-loaded to heavily-loaded state and vice versa

## □ Scalability

- ▣ A scheduling algorithm should be capable of handling small as well as large networks
- ▣ Inquiring workload of all nodes before making scheduling decision can be done in small networks, but not for large one.
- ▣ Select  $m$  of  $N$  nodes for selecting host for process execution

# Desirable features of a scheduling algorithm

7

## □ Fault tolerance

### ▣ Should be capable of working after

- the crash of one or more nodes of the system
- Partitioning into two or more groups due to link failure
- Algorithms with decentralized decision making and consider only available nodes in their decision making have better fault tolerant capability

## □ Fairness of Service

### ▣ How fairly a service is allocated

- ▣ two users initiating equivalent processes expect to receive the same quality of service
- ▣ Load balancing approach is not suitable for fairness of service
- ▣ Load sharing approach is suitable as a node will share its resources as long as its users are not significantly

# Task assignment approach

8

- Main assumptions
  - ▣ Processes have been split into tasks
  - ▣ Computation requirement of tasks and speed of processors are known
  - ▣ Cost of processing tasks on nodes are known
  - ▣ Communication cost between every pair of tasks are known
  - ▣ Resource requirements and available resources on node are known
  - ▣ Reassignment of tasks are not possible



# Task assignment approach

9

- Basic idea: Finding an optimal assignment to achieve goals such as the following:
  - ▣ Minimization of IPC costs
  - ▣ Quick turnaround time of process
  - ▣ High degree of parallelism
  - ▣ Efficient utilization of resources

# Task assignment example

10

- There are two nodes,  $\{n1, n2\}$  and six tasks  $\{t1, t2, t3, t4, t5, t6\}$ . There are two task assignment parameters –
- the task execution cost ( $x_{ab}$  the cost of executing task  $a$  on node  $b$ ) and
- the inter-task communication cost ( $c_{ij}$  the inter-task communication cost between tasks  $i$  and  $j$ ).

## Inter-task communication cost

	t1	t2	t3	t4
t1	0	6	4	0
t2	6	0	8	12
t3	4	8	0	0
t4	0	12	0	0

## Execution costs

	<u>Nodes</u>	
	n1	n2
t1	5	10
t2	2	$\infty$
t3	4	4
t4	6	3

# Task assignment example

11

**Serial assignment**, where tasks t1, t2 are assigned to node n1 and tasks t3, t4 are assigned to node n2:

$$\text{Execution cost, } x = x_{11} + x_{21} + x_{32} + x_{42} = 5 + 2 + 4 + 3 = 14$$

$$\text{Communication cost, } c = c_{13} + c_{14} + c_{23} + c_{24} = 4 + 0 + 8 + 12 = 24.$$

$$\text{Hence total cost} = 14 + 24 = 38.$$

2) **Optimal assignment**, where tasks t1, t3 are assigned to node n2 and task t2 and t4 is assigned to node n1.

$$\text{Execution cost, } x = x_{12} + x_{32} + x_{21} + x_{41} = 10 + 4 + 2 + 6 = 22$$

$$\text{Communication cost, } c = c_{12} + c_{14} + c_{32} + c_{34} = 6 + 0 + 8 + 0 = 14$$

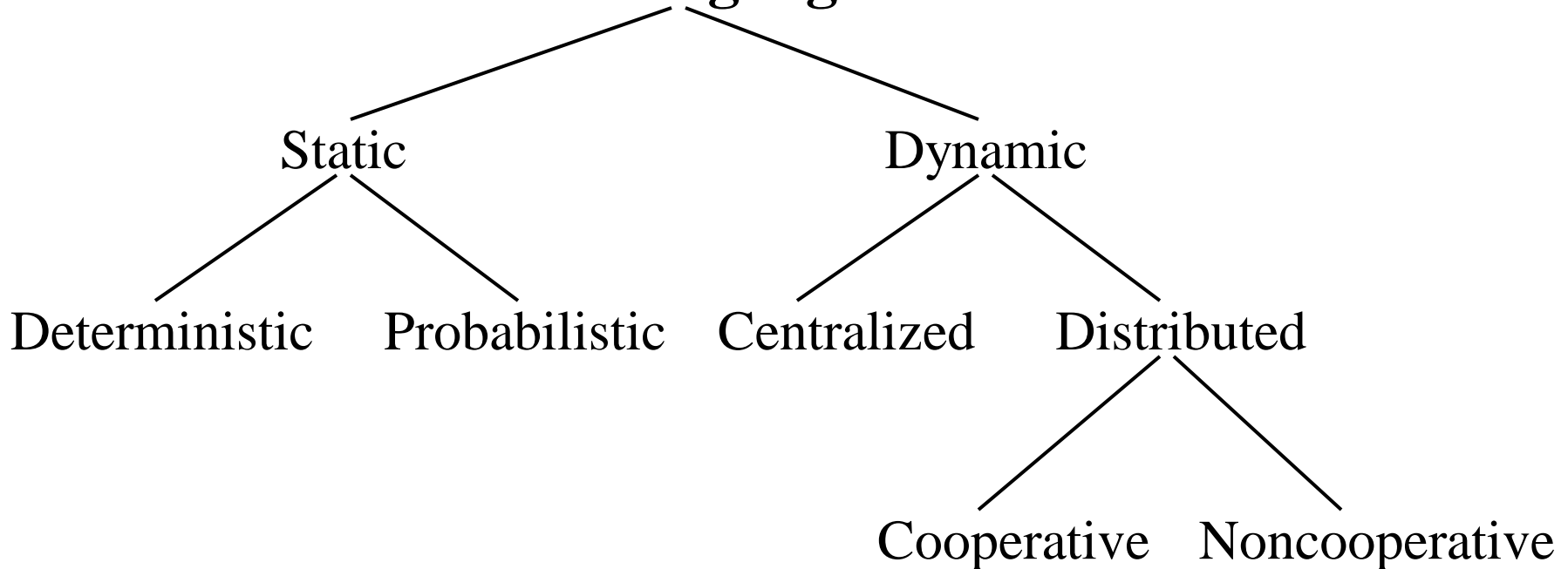
$$\text{Total cost} = 22 + 14 = 36$$

# Load-balancing approach

12

## A Taxonomy of Load-Balancing Algorithms

### **Load-balancing algorithms**



# Load-balancing approach

## Type of load-balancing algorithms

13

### □ **Static versus Dynamic**

- ▣ Static algorithms use only information about the average behavior of the system
- ▣ Static algorithms ignore the current state or load of the nodes in the system
- ▣ Static algorithms are much more simpler
- ▣ Dynamic algorithms collect state information and react to system state if it changed
- ▣ Dynamic algorithms are able to give significantly better performance

# Load-balancing approach

Type of static load-balancing algorithms

14

## □ **Deterministic versus Probabilistic**

- ▣ Deterministic algorithms use the information about the properties of the nodes and the characteristic of processes to be scheduled
- ▣ Probabilistic algorithms use information of static attributes of the system (e.g. number of nodes, processing capability, topology) to formulate simple process placement rules
- ▣ Deterministic approach is difficult to optimize & cost is more
- ▣ Probabilistic approach is easier to implement & has poor performance

# Load-balancing approach

Type of dynamic load-balancing algorithms

15

## □ **Centralized versus Distributed**

- ▣ Centralized approach collects information to server node and makes assignment decision
- ▣ Distributed approach contains entities to make decisions on a predefined set of nodes
- ▣ Centralized algorithms can make efficient decisions, have lower fault-tolerance
- ▣ Distributed algorithms avoid the bottleneck of collecting state information and react faster

# Load-balancing approach

Type of distributed load-balancing algorithms

16

## □ **Cooperative versus Noncooperative**

- In Noncooperative algorithms entities act as autonomous ones and make scheduling decisions independently from other entities
- In Cooperative algorithms distributed entities cooperate with each other
- Cooperative algorithms are more complex and involve larger overhead
- Stability of Cooperative algorithms are better



# Issues in designing Load-balancing algorithms

17

- Load estimation policy
  - ▣ determines how to estimate the workload of a node
- Process transfer policy
  - ▣ determines whether to execute a process locally or remote
- State information exchange policy
  - ▣ determines how to exchange load information among nodes
- Location policy
  - ▣ determines to which node the transferable process should be sent
- Priority assignment policy
  - ▣ determines the priority of execution of local and remote processes
- Migration limiting policy
  - ▣ determines the total number of times a process can migrate

# Load estimation policy I.

for Load-balancing algorithms

18

- To balance the workload on all the nodes of the system, it is necessary to decide how to measure the workload of a particular node
- Some measurable parameters (with time and node dependent factor) can be the following:
  - ▣ Total number of processes on the node
  - ▣ Resource demands of these processes
  - ▣ Instruction mixes of these processes
  - ▣ Architecture and speed of the node's processor
- Several load-balancing algorithms use the total number of processes to achieve big efficiency

# Load estimation policy II.

for Load-balancing algorithms

19

- In some cases the true load could vary widely depending on the remaining service time, which can be measured in several way:
  - ▣ *Memoryless method* assumes that all processes have the same expected remaining service time, independent of the time used so far
  - ▣ *Pastrepeats* assumes that the remaining service time is equal to the time used so far
  - ▣ *Distribution method* states that if the distribution of service times is known, the associated process's remaining service time is the expected remaining time conditioned by the time already used

# Load estimation policy III.

for Load-balancing algorithms

20

- None of the previous methods can be used in modern systems because of periodically running processes and daemons
- An acceptable method for use as the load estimation policy in these systems would be to measure the **CPU utilization of the nodes**
- Central Processing Unit utilization is defined as the number of CPU cycles actually executed per unit of real time
- It can be measured by setting up a timer to periodically check the CPU state (idle/busy)

# Process transfer policy I.

for Load-balancing algorithms

21

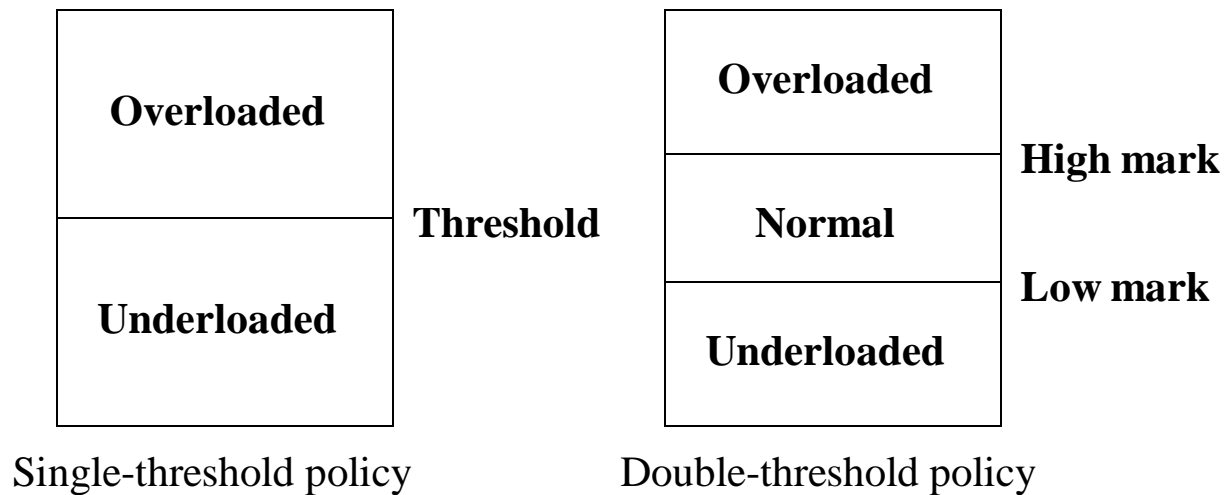
- Most of the algorithms use the *threshold policy* to decide on whether the node is lightly-loaded or heavily-loaded
- Threshold value is a limiting value of the workload of node which can be determined by
  - ▣ Static policy: predefined threshold value for each node depending on processing capability
  - ▣ Dynamic policy: threshold value is calculated from average workload and a predefined constant
- Below threshold value node accepts processes to execute, above threshold value node tries to transfer processes to a lightly-loaded node

# Process transfer policy II.

for Load-balancing algorithms

22

- Single-threshold policy may lead to unstable algorithm because underloaded node could turn to be overloaded right after a process migration



- To reduce instability double-threshold policy has been proposed which is also known as high-low policy

# Process transfer policy III.

for Load-balancing algorithms

23

## □ Double threshold policy

- When node is in overloaded region new local processes are sent to run remotely, requests to accept remote processes are rejected
- When node is in normal region new local processes run locally, requests to accept remote processes are rejected
- When node is in underloaded region new local processes run locally, requests to accept remote processes are accepted

# Location policy I.

for Load-balancing algorithms

24

## □ Threshold method

- ▣ Policy selects a random node, checks whether the node is able to receive the process, then transfers the process. If node rejects, another node is selected randomly. This continues until probe limit is reached.

## □ Shortest method

- ▣ L distinct nodes are chosen at random, each is polled to determine its load. The process is transferred to the node having the minimum value unless its workload value prohibits to accept the process.
- ▣ Simple improvement is to discontinue probing whenever a node with zero load is encountered.



# Location policy II.

for Load-balancing algorithms

25

## □ Bidding method

- ▣ Nodes contain managers (to send processes) and contractors (to receive processes)
- ▣ Managers broadcast a request for bid, contractors respond with bids (prices based on capacity of the contractor node) and manager selects the best offer
- ▣ Winning contractor is notified and asked whether it accepts the process for execution or not
- ▣ Full autonomy for the nodes regarding scheduling
- ▣ Big communication overhead
- ▣ Difficult to decide a good pricing policy

# Location policy III.

for Load-balancing algorithms

26

## □ Pairing

- ▣ Contrary to the former methods the pairing policy is to reduce the variance of load only between pairs
- ▣ Each node asks some randomly chosen node to form a pair with it
- ▣ If it receives a rejection it randomly selects another node and tries to pair again
- ▣ Two nodes that differ greatly in load are temporarily paired with each other and migration starts
- ▣ The pair is broken as soon as the migration is over
- ▣ A node only tries to find a partner if it has at least two processes

# State information exchange policy I. for

## Load-balancing algorithms

27

- Dynamic policies require frequent exchange of state information, but these extra messages arise two opposite impacts:
  - ▣ Increasing the number of messages gives more accurate scheduling decision
  - ▣ Increasing the number of messages raises the queuing time of messages
- State information policies can be the following:
  - ▣ Periodic broadcast
  - ▣ Broadcast when state changes
  - ▣ On-demand exchange
  - ▣ Exchange by polling

# State information exchange policy II.

for Load-balancing algorithms

28

## □ Periodic broadcast

- ▣ Each node broadcasts its state information after the elapse of every  $T$  units of time
- ▣ Problem: heavy traffic, fruitless messages, poor scalability since information exchange is too large for networks having many nodes

## □ Broadcast when state changes

- ▣ Avoids fruitless messages by broadcasting the state only when a process arrives or departures
- ▣ Further improvement is to broadcast only when state switches to another region (double-threshold policy)

# State information exchange policy III.

for Load-balancing algorithms

29

## □ On-demand exchange

- In this method a node broadcast a State-Information-Request message when its state switches from normal to either underloaded or overloaded region.
- On receiving this message other nodes reply with their own state information to the requesting node
- Further improvement can be that **only those nodes reply which are useful to the requesting node**

## □ Exchange by polling

- To avoid poor scalability (coming from broadcast messages) the partner node is searched by polling the other nodes on by one, until poll limit is reached

# Priority assignment policy for Load-balancing algorithms

30

## □ Selfish

- ▣ Local processes are given higher priority than remote processes. Worst response time performance of the three policies.

## □ Altruistic

- ▣ Remote processes are given higher priority than local processes. Best response time performance of the three policies.

## □ Intermediate

- ▣ When the number of local processes is greater or equal to the number of remote processes, local processes are given higher priority than remote processes. Otherwise, remote processes are given higher priority than local processes.

# Migration limiting policy

for Load-balancing algorithms

31

- This policy determines the total number of times a process can migrate
  - ▣ Uncontrolled
    - A remote process arriving at a node is treated just as a process originating at a node, so a process may be migrated any number of times
  - ▣ Controlled
    - Avoids the instability of the uncontrolled policy
    - Use a *migration count* parameter to fix a limit on the number of time a process can migrate
    - Irrevocable migration policy: *migration count* is fixed to 1
    - For long execution processes *migration count* must be greater than 1 to adapt for dynamically changing states

# Load-sharing approach

32

- Drawbacks of Load-balancing approach
  - ▣ Load balancing technique with attempting equalizing the workload on all the nodes is not an appropriate object since big overhead is generated by gathering exact state information
  - ▣ Load balancing is not achievable since number of processes in a node is always fluctuating and temporal unbalance among the nodes exists every moment
- Basic ideas for Load-sharing approach
  - ▣ It is necessary and sufficient to prevent nodes from being idle while some other nodes have more than two processes
  - ▣ Load-sharing is much simpler than load-balancing since it only attempts to ensure that no node is idle when heavily node exists
  - ▣ Priority assignment policy and migration limiting policy are the same as that for the load-balancing algorithms



# Load estimation policies

for Load-sharing algorithms

33

- Since load-sharing algorithms simply attempt to avoid idle nodes, it is sufficient to know whether a node is busy or idle
- Thus these algorithms normally employ the simplest load estimation **policy of counting the total number of processes**
- In modern systems where permanent existence of several processes on an idle node is possible, algorithms measure CPU utilization to estimate the load of a node

# Process transfer policies

for Load-sharing algorithms

34

- Algorithms normally use all-or-nothing strategy
- This strategy uses the threshold value of all the nodes fixed to 1
- Nodes become receiver node when it has no process, and become sender node when it has more than 1 process
- To avoid processing power on nodes having zero process load-sharing algorithms use a threshold value of 2 instead of 1
- When CPU utilization is used as the load estimation policy, the double-threshold policy should be used as the process transfer policy

# Location policies I.

for Load-sharing algorithms

35

- Location policy decides whether the sender node or the receiver node of the process takes the initiative to search for suitable node in the system, and this policy can be the following:
  - ▣ Sender-initiated location policy
    - Sender node decides where to send the process
    - Heavily loaded nodes search for lightly loaded nodes
  - ▣ Receiver-initiated location policy
    - Receiver node decides from where to get the process
    - Lightly loaded nodes search for heavily loaded nodes

# Location policies II.

for Load-sharing algorithms

36

- Sender-initiated location policy
  - ▣ Node becomes overloaded, it either broadcasts or randomly probes the other nodes one by one to find a node that is able to receive remote processes
  - ▣ When broadcasting, suitable node is known as soon as reply arrives
- Receiver-initiated location policy
  - ▣ Nodes becomes underloaded, it either broadcast or randomly probes the other nodes one by one to indicate its willingness to receive remote processes
- Receiver-initiated policy require preemptive process migration facility since scheduling decisions are usually made at process departure epochs

# Location policies III.

for Load-sharing algorithms

37

- Experiences with location policies
  - ▣ Both policies gives substantial performance advantages over the situation in which no load-sharing is attempted
  - ▣ Sender-initiated policy is preferable at light to moderate system loads
  - ▣ Receiver-initiated policy is preferable at high system loads
  - ▣ Sender-initiated policy provide better performance for the case when process transfer cost significantly more at receiver-initiated than at sender-initiated policy due to the preemptive transfer of processes

# State information exchange policies for

## Load-sharing algorithms

38

- In load-sharing algorithms it is not necessary for the nodes to periodically exchange state information, but needs to know the state of other nodes when it is either underloaded or overloaded
- Broadcast when state changes
  - ▣ In sender-initiated/receiver-initiated location policy a node broadcasts State Information Request when it becomes overloaded/underloaded
  - ▣ It is called broadcast-when-idle policy when receiver-initiated policy is used with fixed threshold value value of 1
- Poll when state changes
  - ▣ In large networks polling mechanism is used
  - ▣ Polling mechanism randomly asks different nodes for state information until find an appropriate one or probe limit is reached
  - ▣ It is called poll-when-idle policy when receiver-initiated policy is used with fixed threshold value value of 1

# SUMMARY

39

- Resource manager of a distributed system schedules the processes to optimize combination of resources usage, response time, network congestion, scheduling overhead
- Three different approaches has been discussed
  - ▣ Task assignment approach deals with the assignment of task in order to minimize inter process communication costs and improve turnaround time for the complete process, by taking some constraints into account
  - ▣ In load-balancing approach the process assignment decisions attempt to equalize the avarage workload on all the nodes of the system
  - ▣ In load-sharing approach the process assignment decisions attempt to keep all the nodes busy if there are sufficient processes in the system for all the nodes