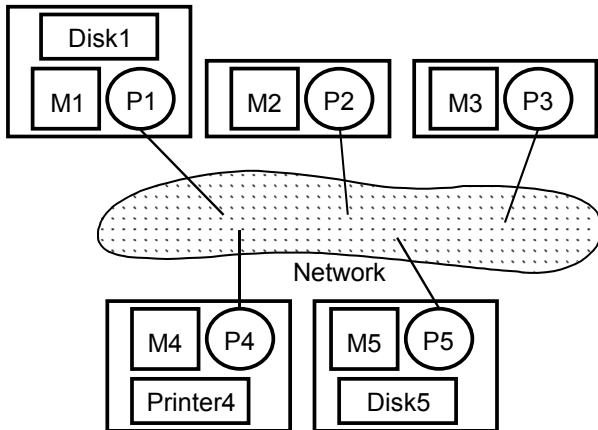


What is a distributed system (again)

- “True” Distributed Operating System
 - Loosely-coupled hardware
 - No shared memory, but provides the “feel” of a single memory
 - Tightly-coupled software
 - One single OS, or at least the feel of one
 - Machines are somewhat, but not completely, autonomous



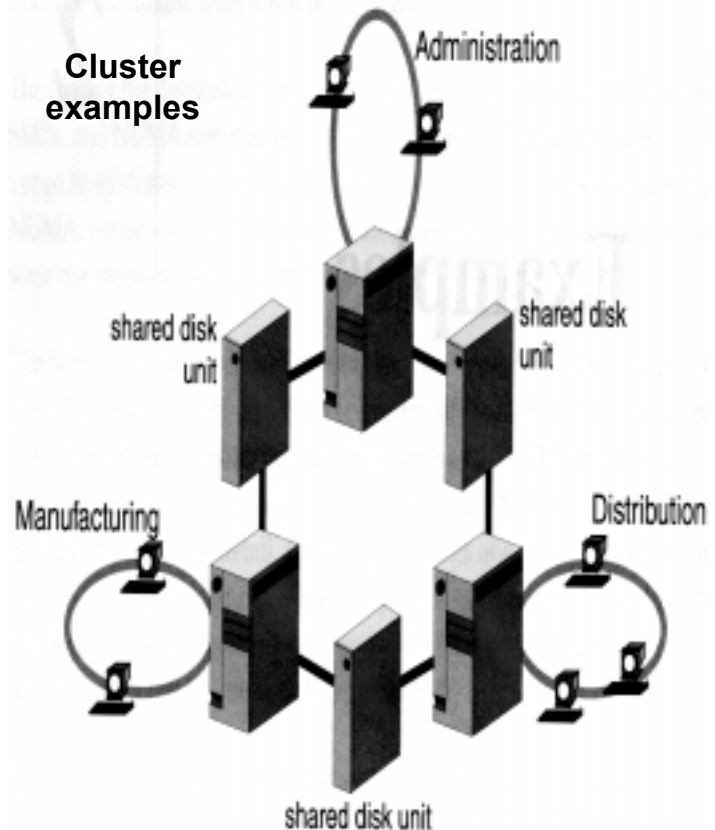
Clusters

- A subclass of distributed systems
- a small scale (mostly) homogeneous (the same hardware and OS) array of computers (located usually in one site) dedicated to small number of well defined tasks in solving of which the cluster acts as one single whole.
- typical tasks for “classic” distributed systems:
 - file services from/to distributed machines over (college) campus
 - distributing workload to all machine on campus
- typical tasks for a cluster:
 - high-availability web-service/file service, other high-availability applications
 - computing “farms”.

Clusters (C) vs. Distributed systems (D)

- structure
 - [C] - homogeneous - purchased to perform a certain task
 - [D] - heterogeneous - built from available hardware
- scale
 - [C] - small scale - setup doesn't have to scale
 - [D] - medium/large - have to span (potentially) large number of machines
- task
 - [C] - specialized - small set of well-defined tasks
 - [D] - general - general-user computing environments
- price
 - [C] - (relatively) cheap [D] - free(?)/expensive
- reliability
 - [C] - as good as it needs to be [D] - high/low?
- security
 - [C] - nodes trust each-other [D] - they don't

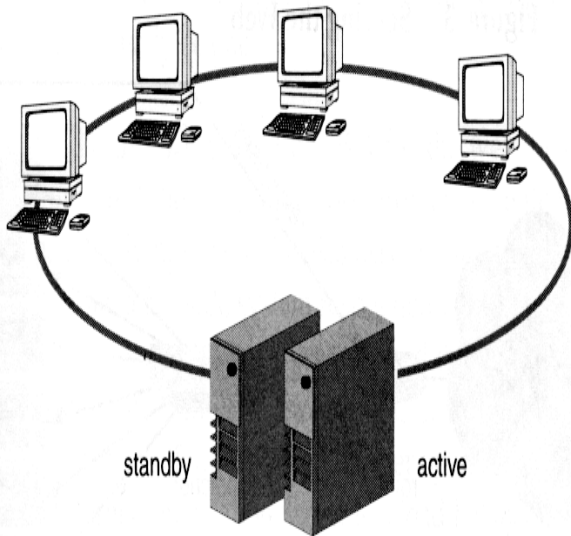
Cluster examples



pictures taken from “In Search of Clusters”, G.F. Pfister, 1998

- branches get access to shared information even if one of the links or computers fails

Cluster examples (cont.)



- active machine - serves files to the network of computers
- standby machine - listens to network and updates it's own copy of files
- in case of machine failure - standby machine takes over file service *transparent* to users

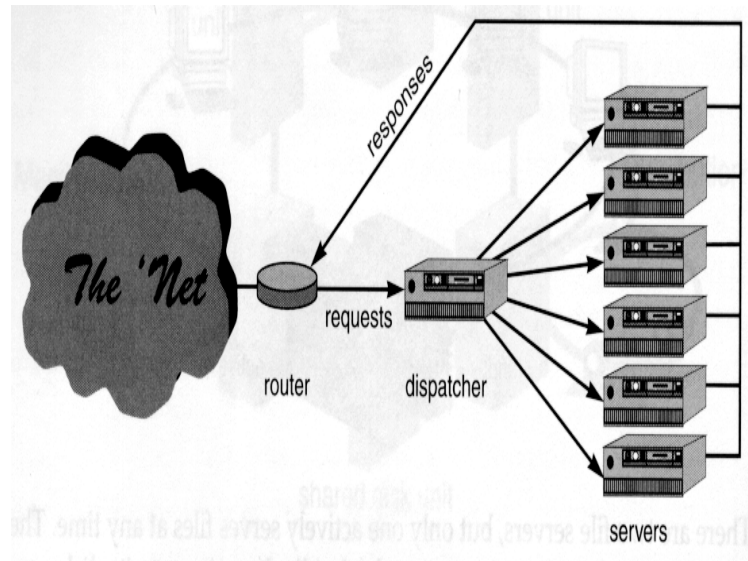
5

Notes by Prof. Mikhail Nesterenko

Spring 2000, Lecture 24

Cluster examples (cont.)

- dispatcher machine - sends the web requests to server machines and makes sure that the servers are evenly loaded
- web service continues even if a server fails



6

Notes by Prof. Mikhail Nesterenko

Spring 2000, Lecture 24

Classification of clusters

- By architecture:
 - with hardware additions - OpenVMS, Tandem Himalaya, Parallel Sysplex
 - pure software - Beowulf, ...
- By task. There is no dividing line between clusters and true distributed systems - as we add features the clusters start to resemble D.S.
 - availability
 - batch processing
 - database
 - generic (scientific) computation
 - full clusters (distributed systems) - single system image

7

Notes by Prof. Mikhail Nesterenko

Spring 2000, Lecture 24

Dependability concepts

- two aspects of dependability
 - reliability - probability of continuous correct operation operation
example: airline navigation system
 - availability - probability that the system operates correctly at given point in time
example: telephone switching system
- error - invalid state of the system
- fault - cause of the error. There are two types:
 - transient - electromagnetic interference, wrong command given by a (human) operator
 - permanent - electric circuit failure, software bug
- fault-tolerance - ability of the system to detect and/or withstand faults. Usually implemented as specialized hardware modules: modular redundancy, inter-module comparators, reliable voting logic
- high-availability - ability of the system to be in operational state with a specified probability.

8

Notes by Prof. Mikhail Nesterenko

Spring 2000, Lecture 24

High-availability

availability	total accumulated outage per year	class (#of 9s)
90%	more than a month	0/1
99%	under 4 days	1/2
99.9%	under 9 hours	2/3
99.99%	about 1 hour	3/4
99.999%	over 5 minutes	4/5
99.9999%	about half a minute	5/6
99.99999%	about 3 seconds	6

- The system is classified by the amount of downtime it allows
 - 1 - campus networks
 - 2 - usual non-clustered commodity stand-alone machines
 - 3 - usual cluster (4 possible)
 - 5 - telephone switches
 - 6 - in-flight aircraft computers

Types of outages, failover

- Two types of outages
 - unplanned - caused by faults
 - planned - need for maintenance of the system (backups, OS upgrades, upgrades, etc.)
- Certain systems should work reliably only part of the time - stock-exchange computers, in-flight computers
- if the system should be available round the clock the objective is to minimize both types of outages
- Simplest high availability cluster: backup server with failover
 - failover - the process of transferring control from failed server to the backup server
 - failback - the process of transferring control from backup server to primary server
- cluster with failover helps avoid planned as well as unplanned outages

Watchdogs

- watchdog is a mechanism of notification (and possible correction) of a failure.
- simplest (software) watchdog - a process monitoring application processes. If the monitored process fails watchdog may take recovery action.
 - watchdog can run on the same machine as the application program - may not be very useful if the machine crashes
 - on different machine - how is communication carried out?
- application process may be programmed to cooperate with the watchdog. Three ways cooperation:
 - heartbeat - periodic notification sent to the watchdog by the application process to confirm its correct execution. Alternate heartbeat paths - network, RS-232, SCSI
 - application initiated
 - watchdog initiated
 - idle notification - application informs watchdog that it is idle
 - error notification - application notifies that it encountered an error it cannot correct

Replication and Switchover

- Two types of cluster failover organization:
 - replication (shared-nothing cluster) - backup server keeps its own copy of data
 - switchover (shared-data cluster) - backup has access to the storage devices used by primary

Replication	Switchover
+ easier to add to an existing single machine	- harder to add - must modify existing cabling
+ easier to configure	- harder to configure
+ can use any old I/O adapters and controllers	- requires specialized I/O devices
+ can use simple storage units	- must use hardened storage like RAID
- 1-to-many backup is hard	+ 1-to-many backup possible as long as interconnect allows
- requires another copy of storage	+ only one copy of storage used
- CPU overhead in normal operation - synchronization needed	+ no overhead in normal operation
- failback requires additional copying	+ no copying on failback

Disaster recovery

- Disaster - failure that affects the large portions or the whole site - fire, flood, storm-damage
- usual recovery technique - resume operations on the system outside the scope of the disaster

tier	description
0	no disaster recovery
1	backups are periodically taken and stored off premises
2	backups are taken to a "hot-site" where they can be loaded on a secondary system if necessary
3	electronic vaulting - network connects primary site and secondary site, back-ups are transferred by network
4	active secondary - data send over the wire, the data is kept loaded and ready to run on secondary
5	secondary is kept completely up-to-date