



***The Modern Mainframe...
At the Heart of Your Business***

A Mainframe Primer - Mainframe Clustering



© 2006 IBM Corporation

Superior Qualities of Service

- How does the mainframe deliver superior qualities of service?
 - ▶ Unmatched scale-up
 - ▶ Continuous operation
 - ▶ Systematic disaster recovery

- Mainframe clustering technology hardware and software are optimized to provide these qualities of service
 - ▶ Unique Parallel Sysplex design is better than anything else

Mainframe Clustering is Superior

■ System z

- ▶ Specialized hardware for clustering
- ▶ Dedicated high speed fiber interconnect
 - Low latency
- ▶ Integrated exploitation by operating system and all software subsystems

■ Distributed

- ▶ No special hardware
- ▶ No special networking
 - Full software path length
- ▶ Each subsystem (database, application server) is designed to run on commodity servers



1. **Very low overhead yields ultimate scalability (up to 32 mainframe systems in a cluster)**
2. **Highest of high availability**

04 - Mainframe Clustering v1.3.ppt

4

A Primer on Mainframe Clustering

■ Coupling Facility

- ▶ Dedicated processor with specialized microcode to coordinate shared resources
- ▶ Supported by machine instruction set
- ▶ Large amounts of fast memory
- ▶ High speed inter-connect to clustered systems
- ▶ Timing facilities to maintain logical execution-order across coupled systems
- ▶ Highly Fault-Tolerant

■ Parallel Sysplex

- ▶ Multiple z/OS images clustered using the coupling facility for coordination

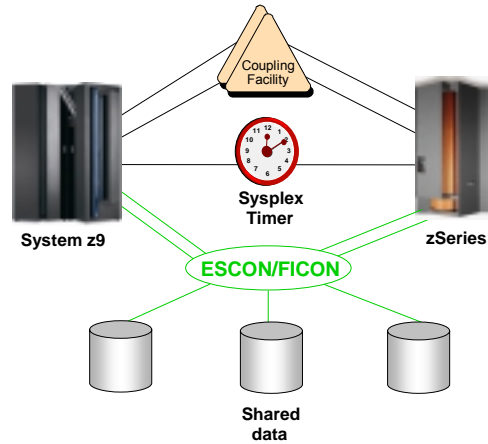
This presentation will use the word "image" to refer to a node in a sysplex cluster, "LPAR" may also be used to describe this

04 - Mainframe Clustering v1.3.ppt

5

Parallel Sysplex – What is it ?

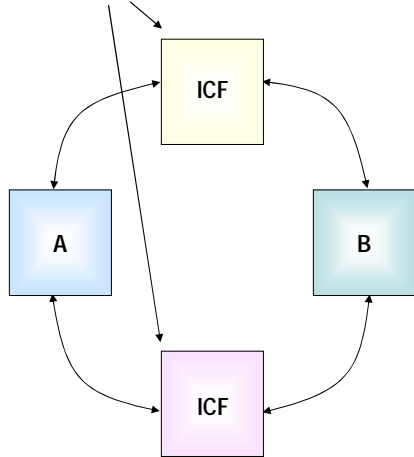
- Hardware
 - ▶ Redundant coupling facilities providing processing and true shared storage
 - ▶ Timing facilities
 - Sysplex timers (Hardware)
 - STP protocol (Software)
 - ▶ High speed interconnections (up to 16 Gigabits/sec, up to 10 meters)
 - ▶ Fiber switch provides access to data
- Micro-code + Software
 - ▶ CFCC (coupling facility control code)
 - High throughput, low latency, micro code control program for the coupling facility
- Clustering service APIs within z/OS
 - ▶ XES APIs support program connectivity
 - ▶ XCF connectivity configuration
- Workload Management
 - ▶ WLM (workload manager within a z/OS instance)
 - ▶ IRD (intelligent resource director across LPAR's)
 - ▶ Both manage workload across the sysplex



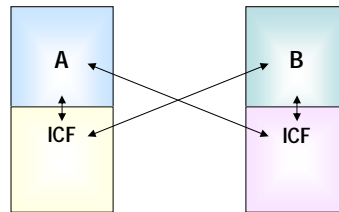
...A Key IBM Unique Differentiator in the IT Industry

Implementation of Coupling Facility

Standalone Hardware Dual Coupling Facilities



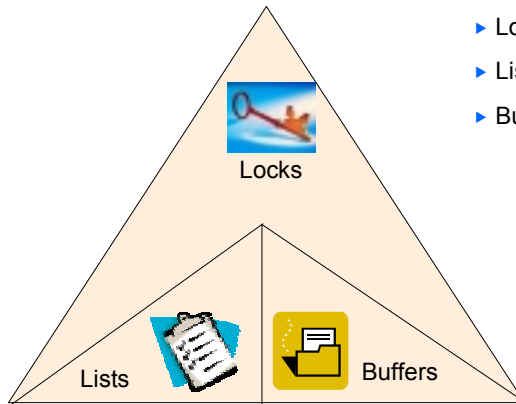
ICF within System z LPARs



Coupling Facility is a Flexible Effective Hardware Technology for Managing Clusters

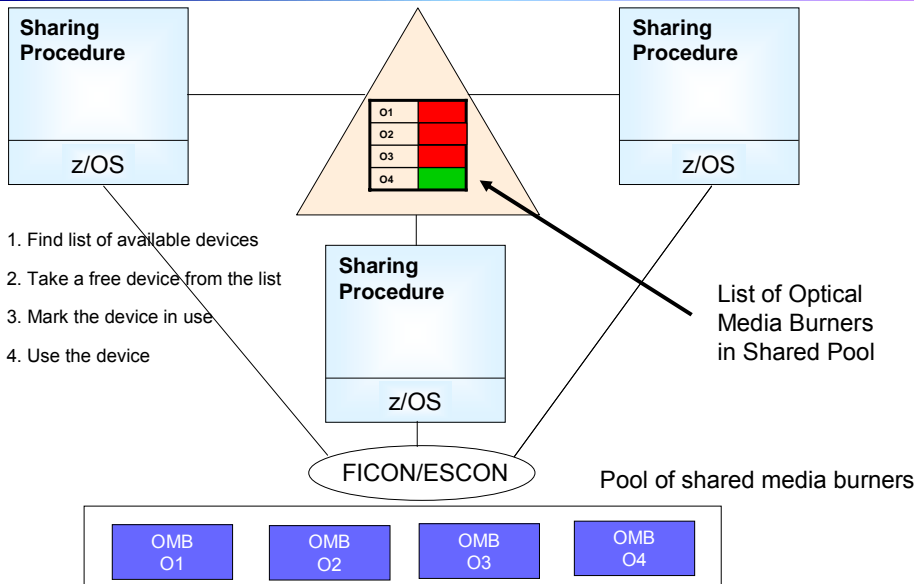
The Coupling Facility implements

- ▶ Locks for synchronizing data
- ▶ Lists for sharing data
- ▶ Buffers for database consistency

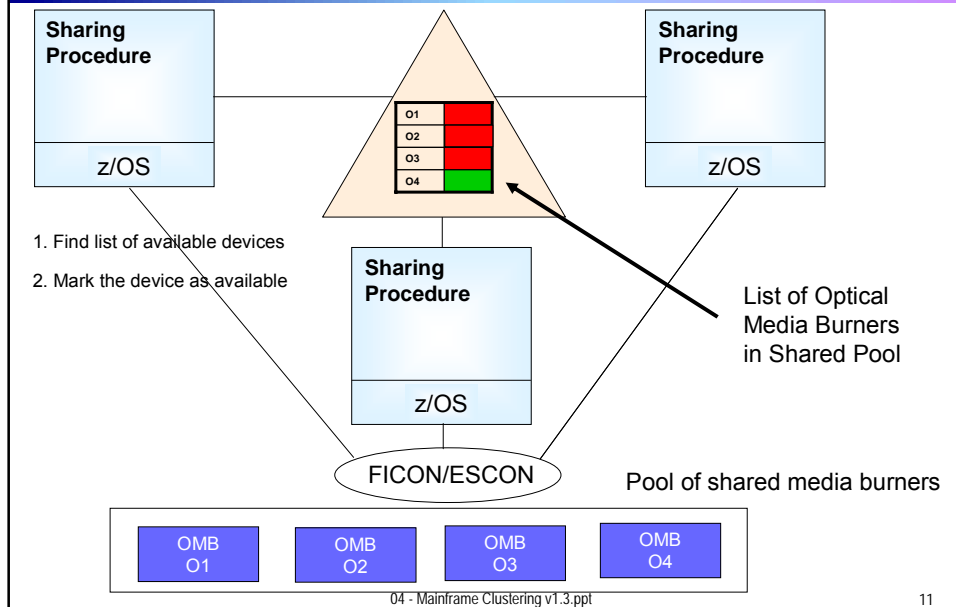


Let's look at examples of how each of these are used in a cluster

Using the List Capability for Sharing Devices Getting a Device for Making a Backup Copy



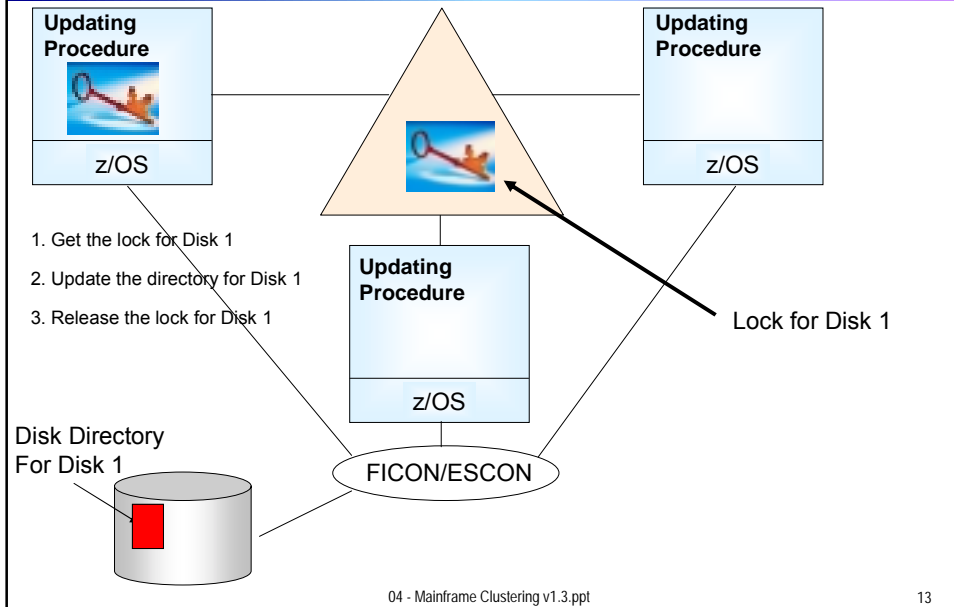
Using the List Capability for Sharing Devices Releasing a Device for Use by Others



Other System Uses of Lists

- Shared Resources
 - ▶ Tapes
 - ▶ Files
 - ▶ Consoles
 - ▶ Etc
- Sysplex-wide information
 - ▶ Workload-balancing information
 - ▶ Status of each system in the sysplex
- Subsystem information
 - ▶ Logfiles for recovery
 - ▶ Configuration and Restart Data

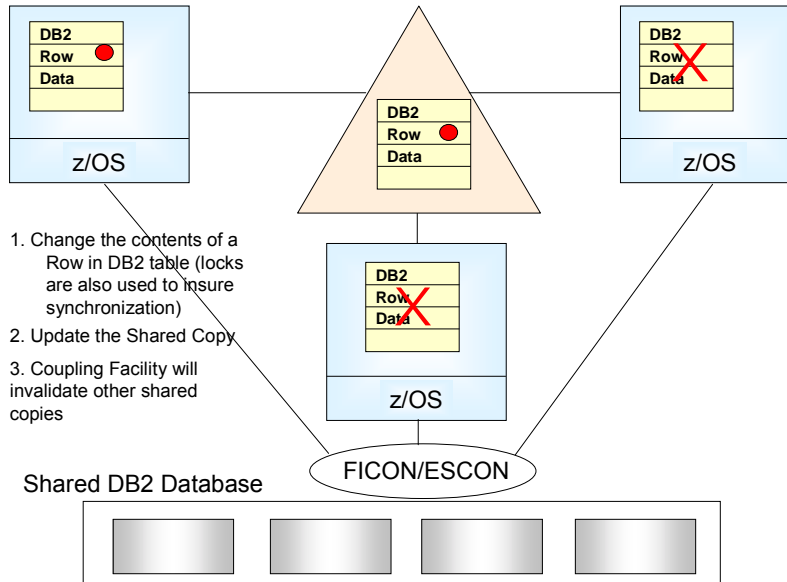
Use of the Lock Capability for Updating Information



Other System Uses of Locks

- Any synchronization of shared information
 - ▶ Files
 - ▶ Databases
 - ▶ System-wide resources

Using the Buffer Capability for DB2 Data Consistency



04 - Mainframe Clustering v1.3.ppt

15

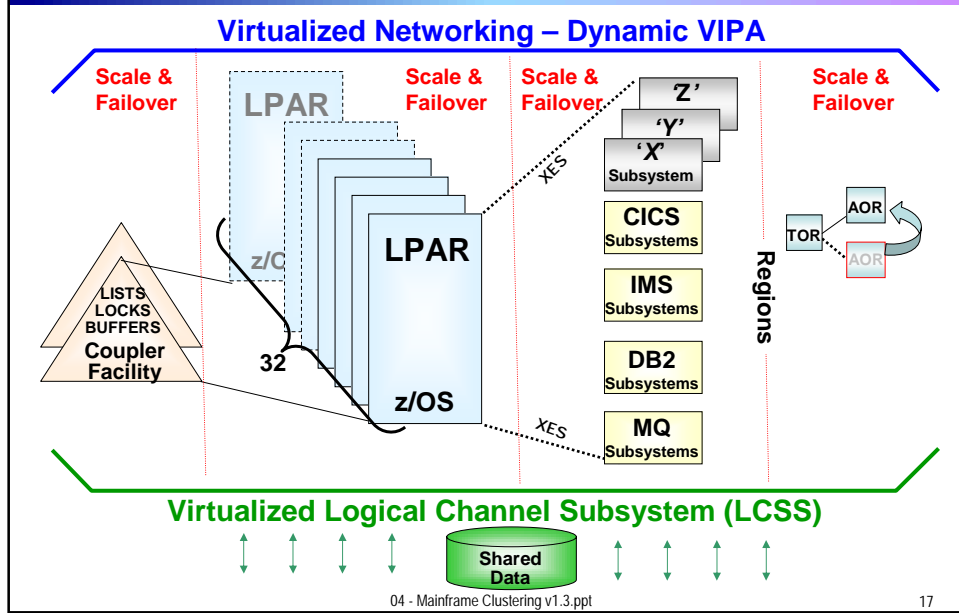
Other Uses of Buffers for Data Consistency

- DB2 for System z
- IMS
- VSAM
- Computer Associates IDMS
- Computer Associates Datacom

04 - Mainframe Clustering v1.3.ppt

16

Multi-layered Approach Results in High Availability and Scalability



We'll discuss how exploitation of the parallel sysplex helps DB2 beat Oracle RAC later.

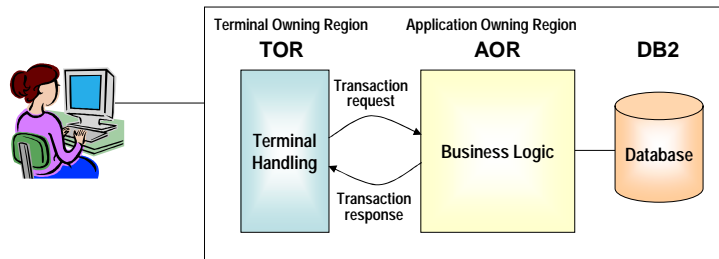
Let's take a quick look at how CICS benefits from the parallel sysplex and these multiple layers



IBM

Multi-layer Benefits for CICS Layer 1 – Regions

- CICS takes a transaction request from an end user, accesses a database, performs business logic and returns a response (similar to J2EE)



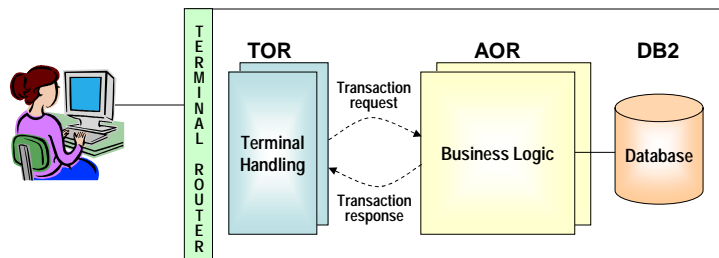
- Each CICS region (TOR and AOR) provides a single thread of execution
- Regions provide transaction isolation

04 - Mainframe Clustering v1.3.ppt

19

Multi-layer Benefits for CICS Layer 2 – LPARs

- Uses multiple TOR and AOR



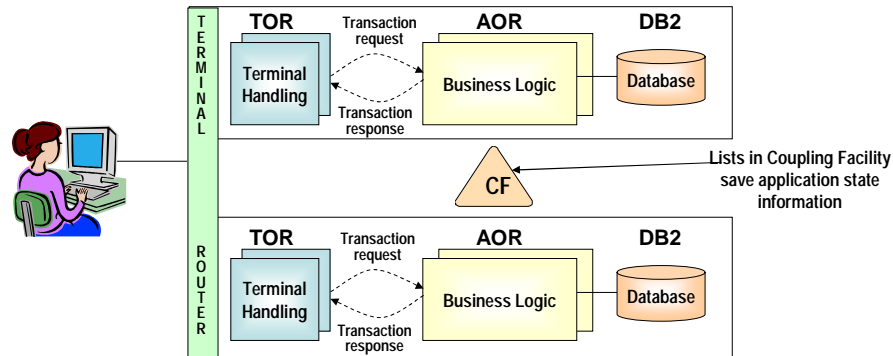
- Terminal router routes transaction to appropriate TOR
- Multiple TORs & AORs scale by adding system resources (threads, memory, etc)
- Multiple TORs & AORs provide availability
 - ▶ A software failure could bring down an AOR or TOR (e.g. programmer error)
 - ▶ Current in flight transactions are rolled back
 - ▶ New transactions are routed to other TOR or AOR

04 - Mainframe Clustering v1.3.ppt

20

Multi-layer Benefits for CICS Layer 3 – Sysplex

- Multiple TOR and AOR on multiple machines in parallel sysplex



- Scalability is enhanced In that processing resources from up to 32 LPARS in the sysplex can be utilized
- The work of a failed TOR or AOR can be taken over by any other TOR or AOR in the sysplex
- Protects against hardware failure of an entire machine or operating system image

04 - Mainframe Clustering v1.3.ppt

21

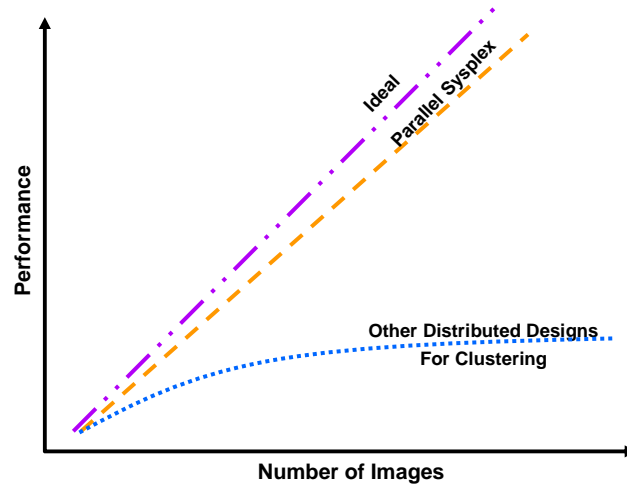
Parallel Sysplex Performance

- High performance interconnect and low latency in coupling facility causes minimal overhead.
- Typical overhead
 - ▶ Multisystem Management - 3%
 - ▶ Resource Sharing - 3%
 - ▶ Application data sharing - <10%
 - ▶ Incremental cost of adding an image - 1/2%
- Result
 - ▶ Near-linear scalability as more systems are added
 - ▶ Better efficiency than other clustering schemes

04 - Mainframe Clustering v1.3.ppt

22

Mainframe Clustering Delivers Near-Linear Scalability



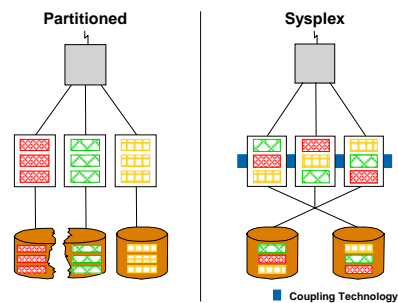
04 - Mainframe Clustering v1.3.ppt

23

To Achieve Higher Scale, Most Distributed Cluster Designs Must Resort to Partitioning the Data

- Data is partitioned, so that each processor is the only one that can access that data
 - ▶ Requires application-level design to accomplish
 - ▶ Growth of processors or data requires re-partitioning
 - ▶ No ability to workload balance some partitions may be busy while others idle
 - ▶ Failover requires re-partitioning the data to the remaining processors

- Result – harder to build, manage, and grow lower availability



04 - Mainframe Clustering v1.3.ppt

24

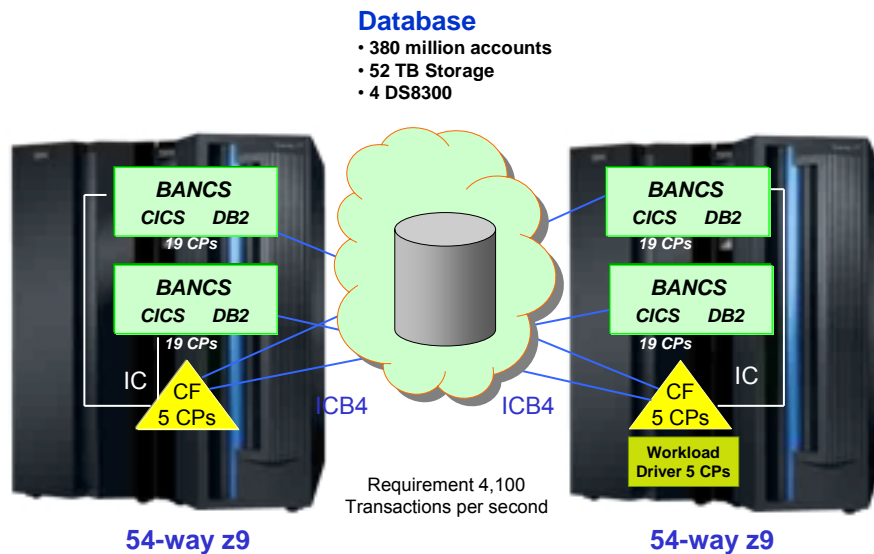
Imagine the Scale...

- A single 54-way* System z delivers 17,801 MIPs and huge I/O bandwidth
 - ▶ This is roughly 6 times the processing capacity of the largest HP Itanium Superdome with 768 processor cores**
- Up to 32 of these systems can be clustered in a parallel sysplex, single system image

* Using z/OS V1.9 shipping September 2007 (previous maximum was 32 processors out of 54)

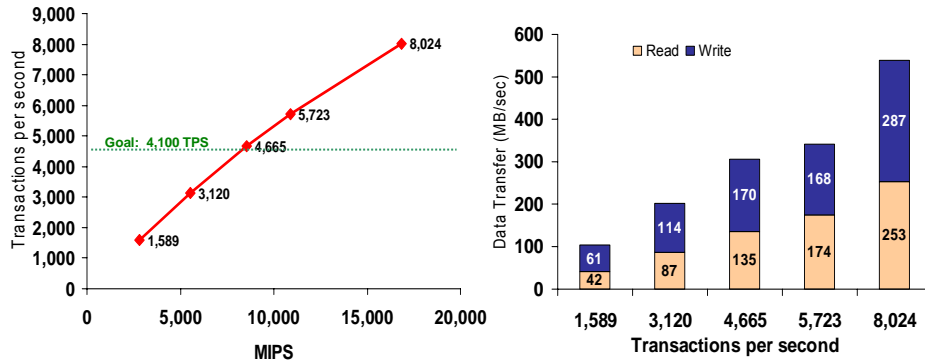
** Based on equivalence factor of 1 MIP = 122 RPE's from HP presentations

The Largest System z Benchmark Known... Bank of China Benchmark System Configuration



Bank of China System z Benchmark

Near-Linear Scalability on a Parallel Sysplex running CICS and DB2 in a single system image with No Partitioning Required



Huge scale up, requires huge I/O bandwidth capacity

04 - Mainframe Clustering v1.3.ppt

27

Mainframe Parallel Sysplex Summary

- Very low overhead to create single system image
- Ultimate scalability
 - ▶ Sysplex up to 32 systems each with 32 processors
- Highest of high availability
 - ▶ Hardware and software
- Foundation for a systematic disaster recovery capability

04 - Mainframe Clustering v1.3.ppt

28