



Fault Tolerant Design in the EOS Archive

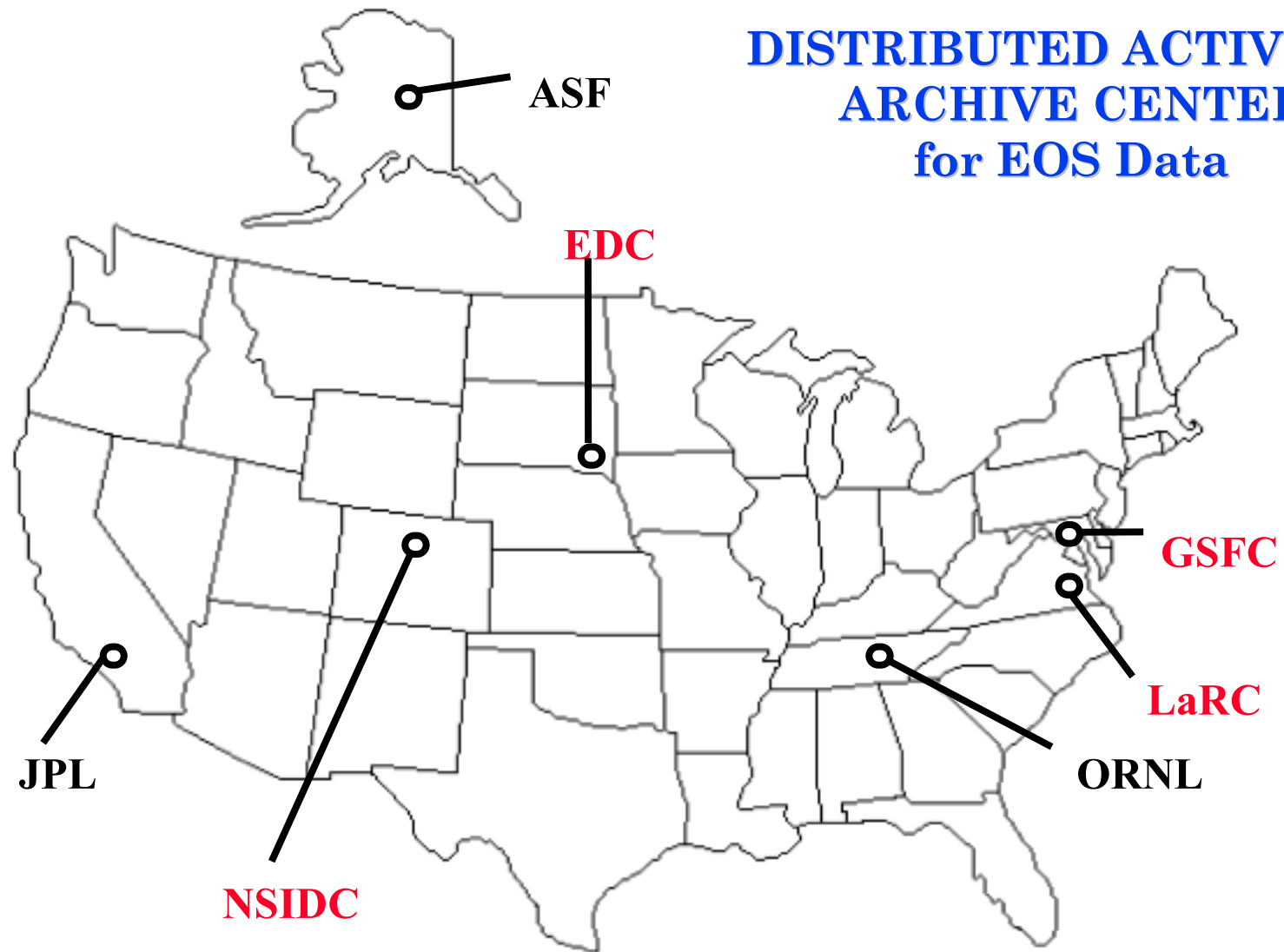
Alla Lake, Jonathan Crawford,
Raymond Simanowith, Bradley Koenig
Lockheed Martin Corporation

IEEE Mass Storage & Technology
Conference

March 29, 2000



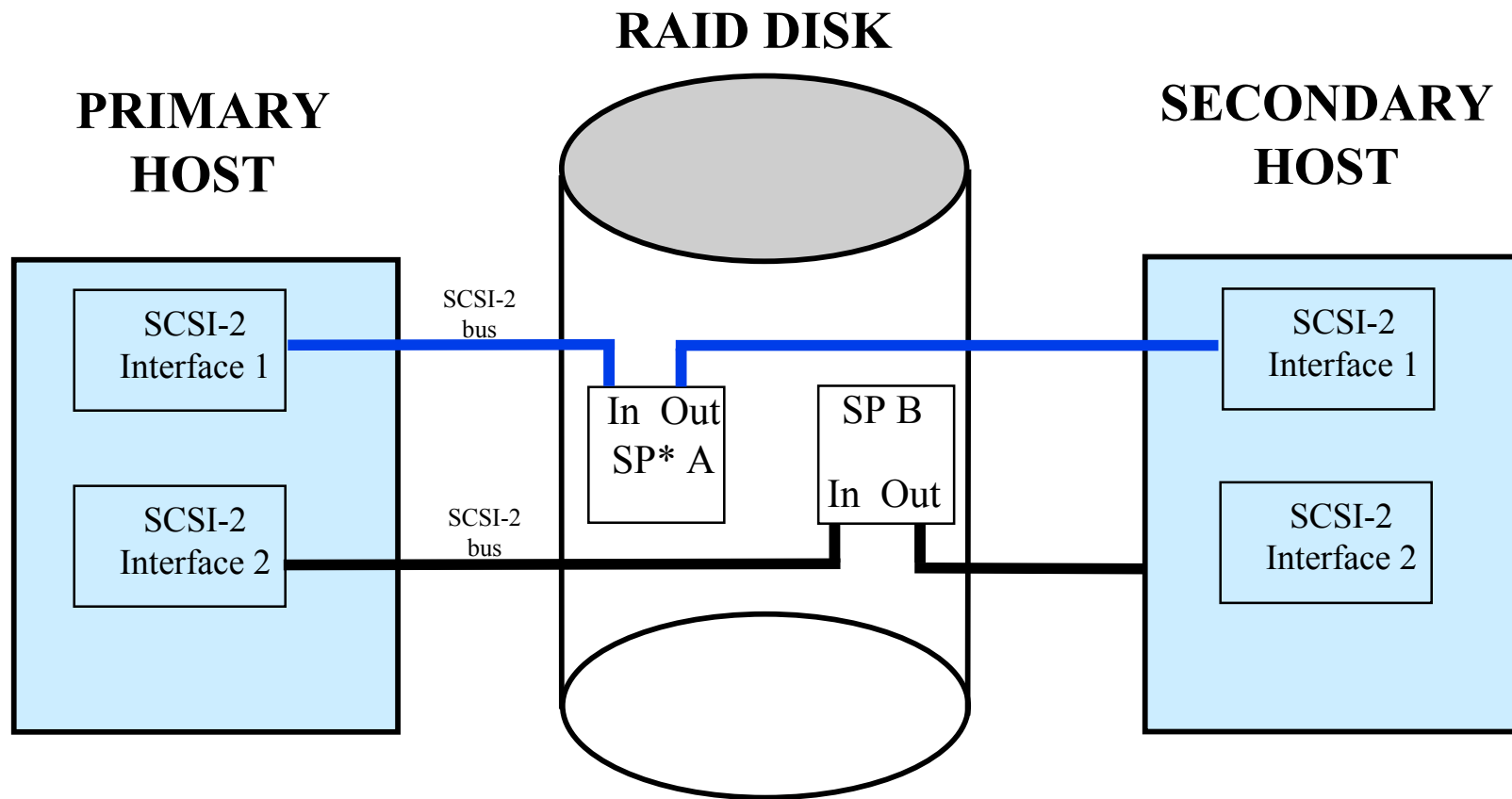
DISTRIBUTED ACTIVE ARCHIVE CENTERS for EOS Data





Server Failover Pair

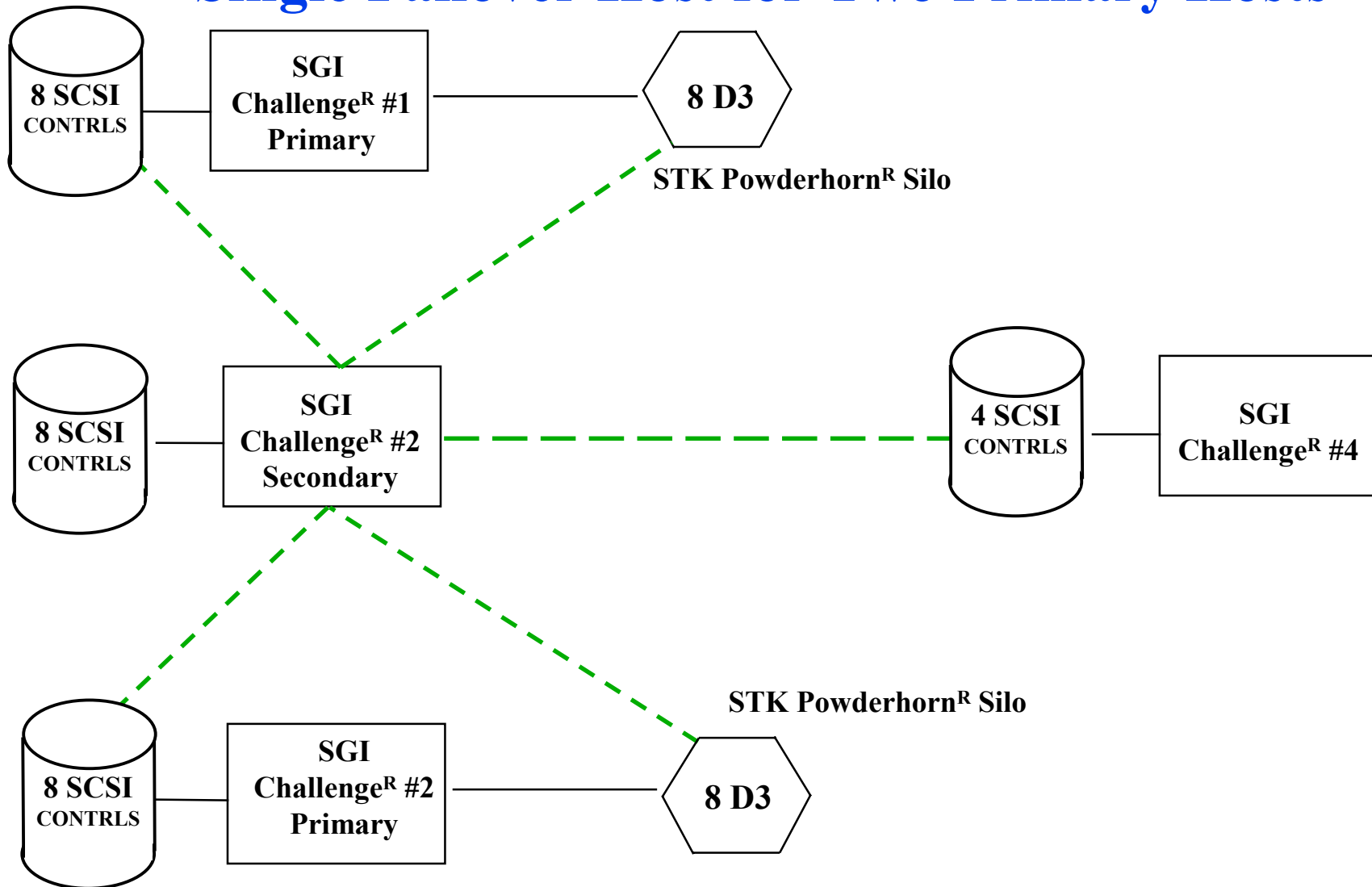
(Dual-Bus/Dual-Initiator Configuration)



*SP - Storage Processor

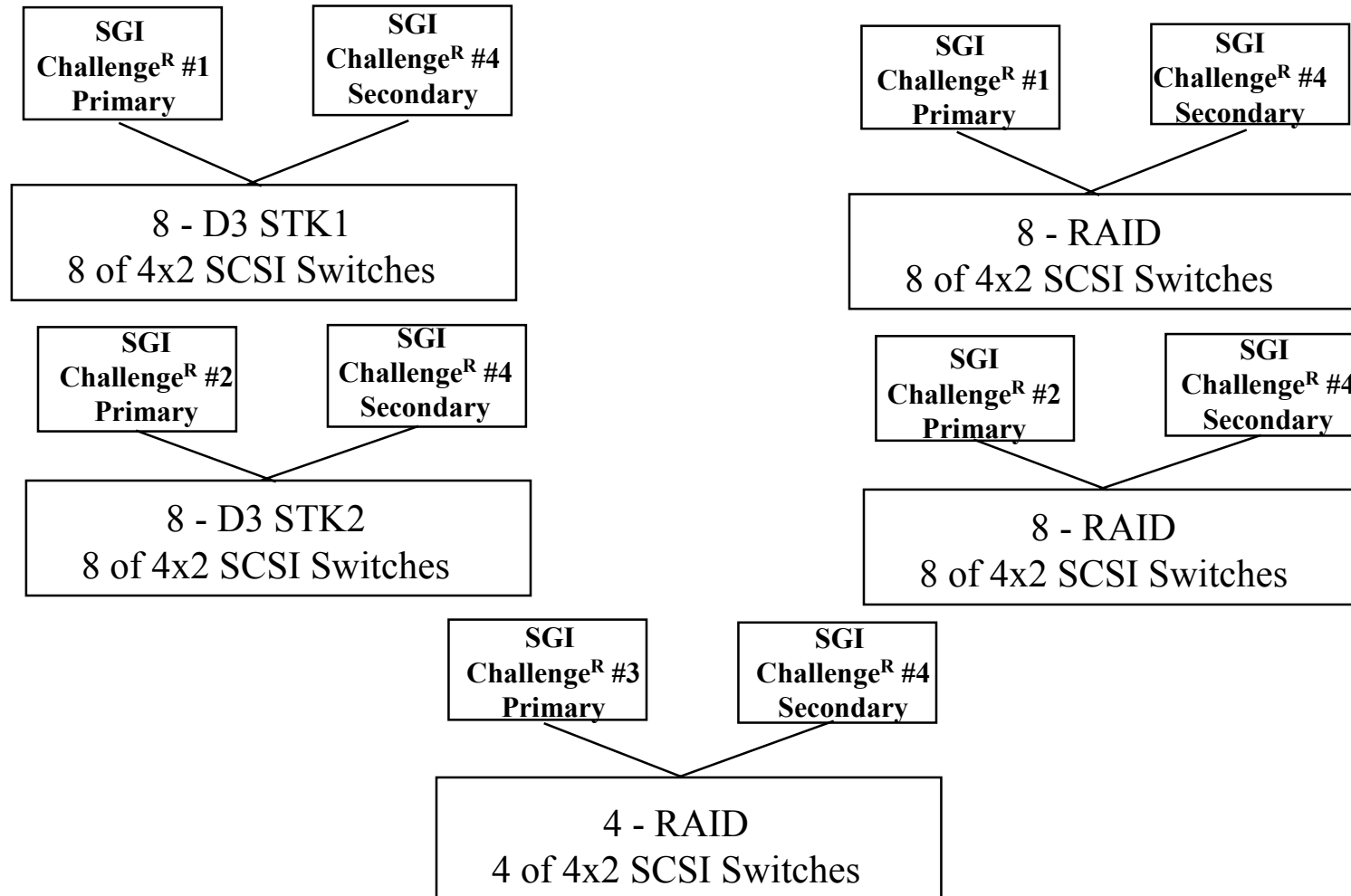


Archive Failover Switching at GSFC - Single Failover Host for Two Primary Hosts





Archive Failover Switching at GSFC - 36 SCSI switches



SCSI Connection Switching is automated using TCL scripts



Hardware Configuration Constraints

- Only one host is actively addressing RAID at any one time
- One host is always Normally Primary to simplify administration
- Both hardware hosts are identically configured
- All Applications reside on the shared RAID to maintain version and configuration synchronization between the primary and secondary hosts
- Hardware Fault Detection and Failover Activation are both fully manual, failover process itself is scripted



Hardware Configuration Constraints

- Presence of SCSI and HiPPI connections greatly complicates the physical aspects hardware failure recovery
- Expect significant improvements in that area with Fibre attachments



Network Failover Implementation

- Both hosts have their own Primary IP address, but also share a single Virtual IP address
- All references (mounts) to the Primary Host are via the Virtual IP address
- The Primary Host has IP aliasing enabled and the Secondary Host Disabled
- The FDDI routers must be manually flushed (ARP Cache cleared) to enable route failover
- requires an administrator to initiate



Sybase Implementation - Fully Scripted

- Database tables are on *xlv* RAID 1 partition
- To enable switchover, the Sybase Open Client Interface File must be changed at Failover to a pre-fabricated “failover” version

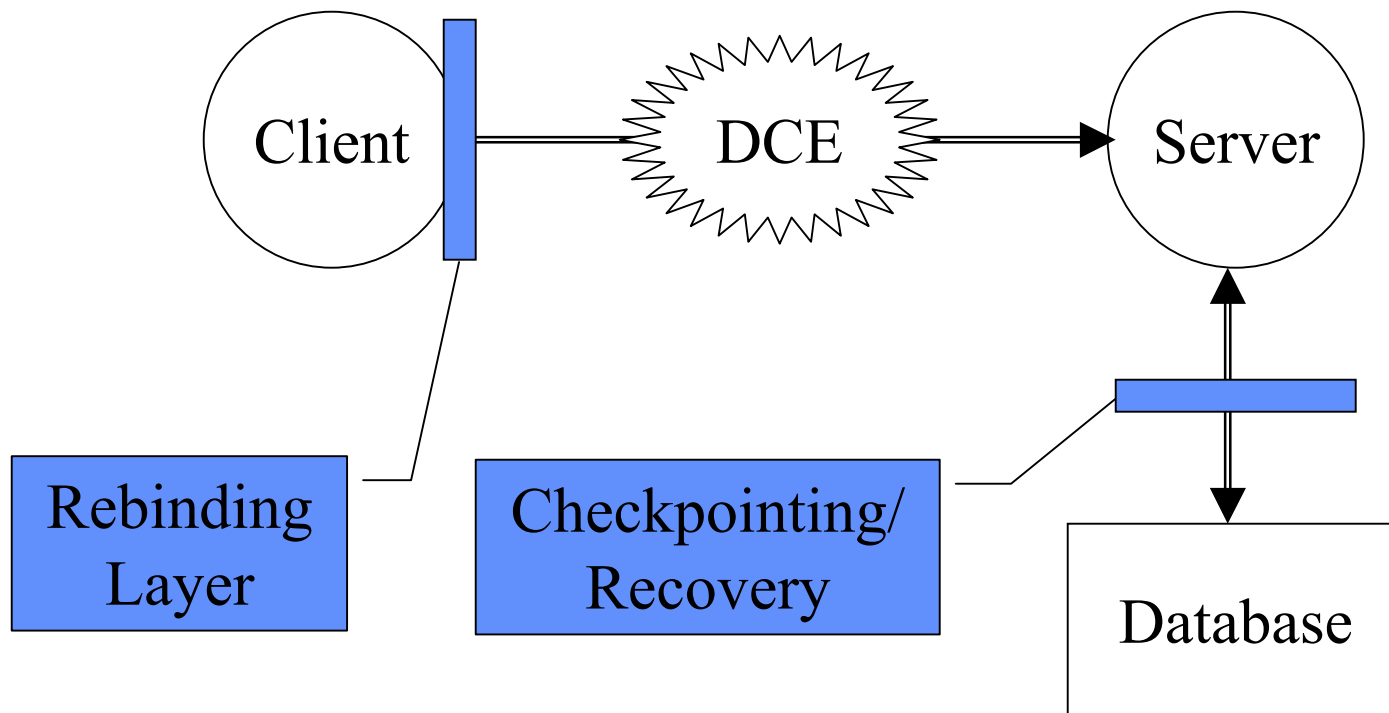


HiPPI Switchover Implementation - Manual

- Manually disable HiPPI interface on the Failover Host, if necessary (i.e. if the the Host is not fully shut down)
- Manually setup “listen” mode on the HiPPI Switch
- Manually enable HiPPI interface on the new Primary Host by invoking a pre-fabricated script



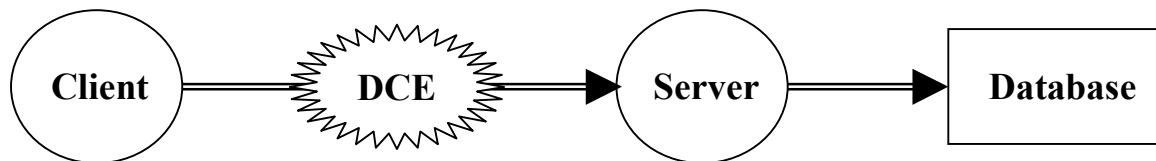
Software Fault Recovery Mechanisms



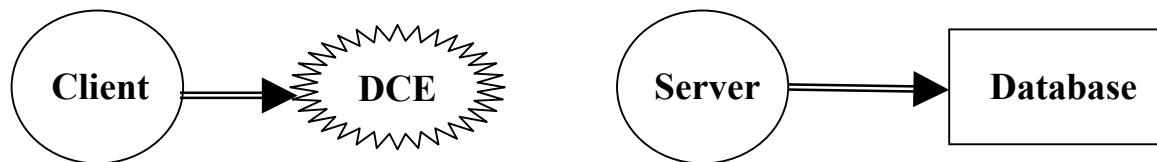


Network Disruption Recovery

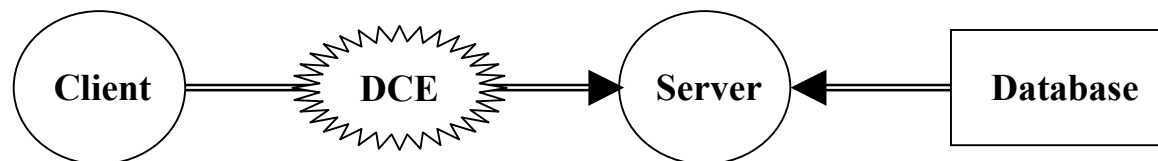
- Client submits request
- Server checkpoints request to database



- Client attempts to rebind to server
- Server continues processing, checkpointing results to database



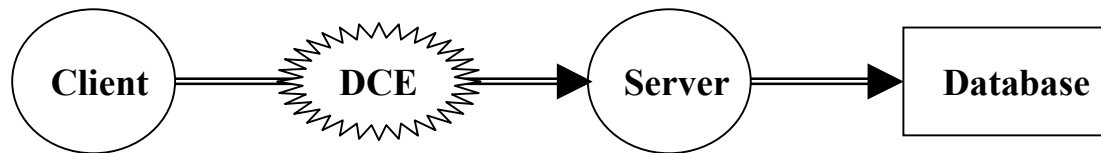
- Upon reconnection, server restores results from database



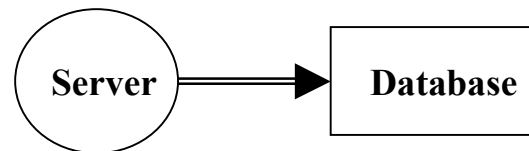


Client Crash Recovery

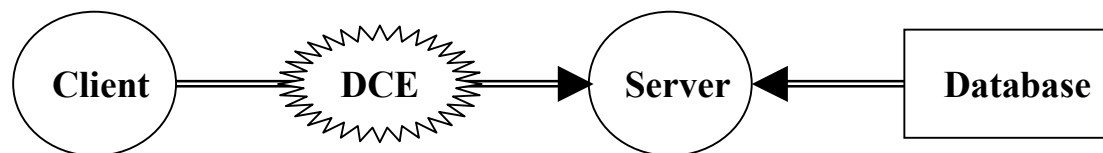
- Client submits request
- Server checkpoints request to database



- Server continues processing, checkpointing results to database



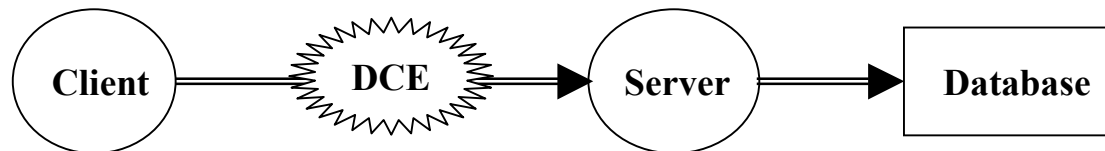
- Client restarts and resubmits request
- Server restores request results from database



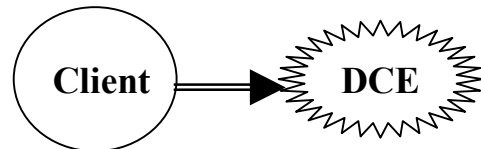


Server Crash Recovery

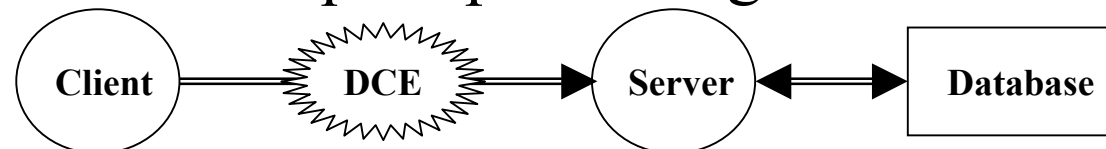
- Client submits request
- Server checkpoints request to database



- Client attempts to rebind to server



- Upon reconnection, client resubmits request
- Server restores from last checkpointed state and resumes request processing





Start Temperatures

- **Warm restart**
 - Typical restart; all request processing resumed from last checkpointed state
- **Cold start**
 - All outstanding requests are cancelled and removed from the queue
 - If a request failed by cold start is resubmitted, it will be processed as a new request
 - Typically used to have a client flush their requests out of the system
- **Cold restart**
 - Requests are cancelled but left in queue
 - All attempts to reprocess request are failed immediately
 - Rarely used; allows a downstream server to “backflush” requests regardless of origin



Contact Information

- alake@eos.hitc.com
- bkoenig@eos.hitc.com
- jcrawfor@eos.hitc.com
- rsimanow@eos.hitc.com