# Informix Cheetah

IBM        /

zhangjb@cn.ibm.com
        IBM

# Informix Cheetah

## Capitalize on its strengths

### *On line Transaction Processing*

V.11 - new features increase the performance of IDS for OLTP. Multi-Active node Cluster with High Availability

### *Integrated Solutions*

IDS is IBM's leading data server for industrial-strength integrated solutions.

### *Low DBA effort*

Low DBA requirements are a strong selling point. Express and Workgroup Editions available.

### *Key Industries*

IDS has significant market presence within key industries including Telco, Government, Retail & Banking

### IDS is a "lead-with" data server for ANY OLTP oppty

# IDS                           -Parallel Everything

**Modern parallel-everywhere base RDBMS**

–
–
–
–
–

–

–

–

I/O

# Availability….

- **?**

  The proportion of time that an application can be used for productive work, measured against the time that it must be functional.

-

  – The reliability of the components that comprise the application: namely, how often they fail.

  – How long it takes for the application to be restored once a failure has occurred.

**Availability:** The accessibility of a system resource in a timely manner; for example, the measurement of a system's uptime. Availability can be measured relative to "100% operational" in terms of the number of "9s" that the system is available.

| Availability Level Class | Uptime | Downtime Limit Per Year |
|---|---|---|
| 2 | 99% | 4 Days |
| 3 | 99.9% | 9 Hours |
| 4 | 99.99% | 1 Hour |
| 5 | 99.999% | 5 Minutes |
| 6 | 99.9999% | 30 Seconds |
| 7 | 99.99999% | 3 Seconds |

Source: "High Availability: A Perspective," Jane Wright, Ann Katan (Gartner Group), 11/24/2004

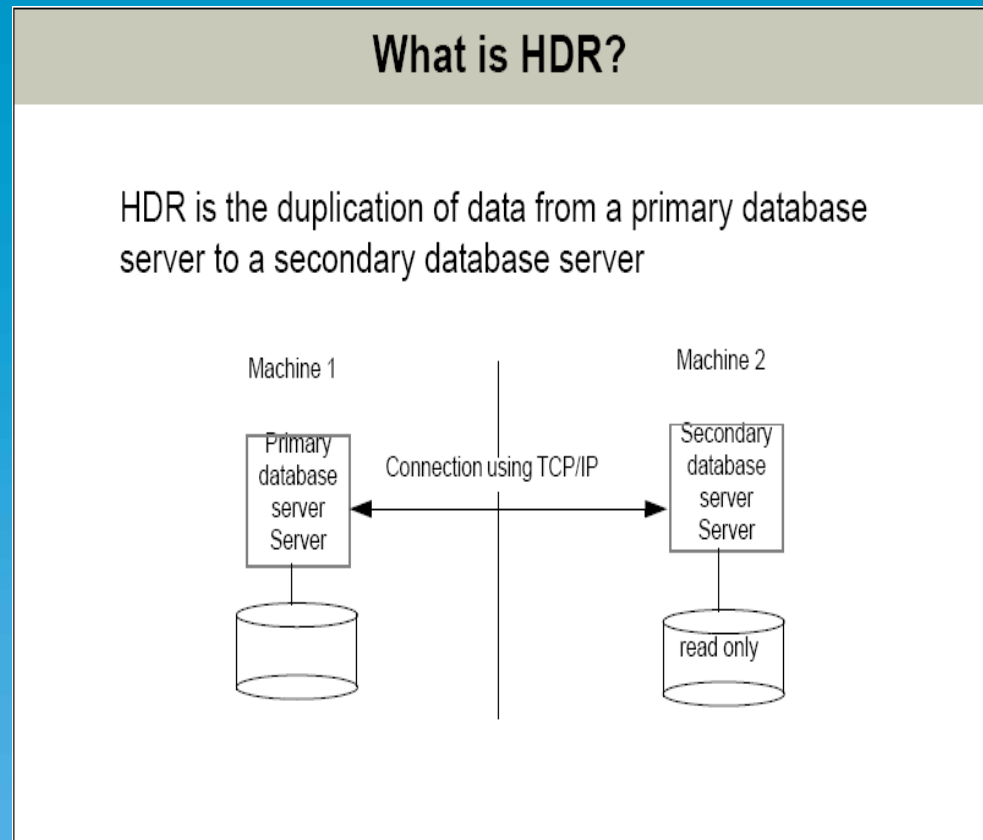Informix
--High Availability Data Replication  (HDR)

Primary

Secondary

# What is HDR?

- Two identical servers on two identical machines
  - Primary server
  - Secondary server
- Primary server
  - Fully functional server
  - **<u>All database activity</u>** – insert/update/deletes, are performed on this instance
  - Sends logs to secondary server



What is HDR?

HDR is the duplication of data from a primary database server to a secondary database server

Machine 1 | Machine 2

Primary database server Server — Connection using TCP/IP — Secondary database server Server
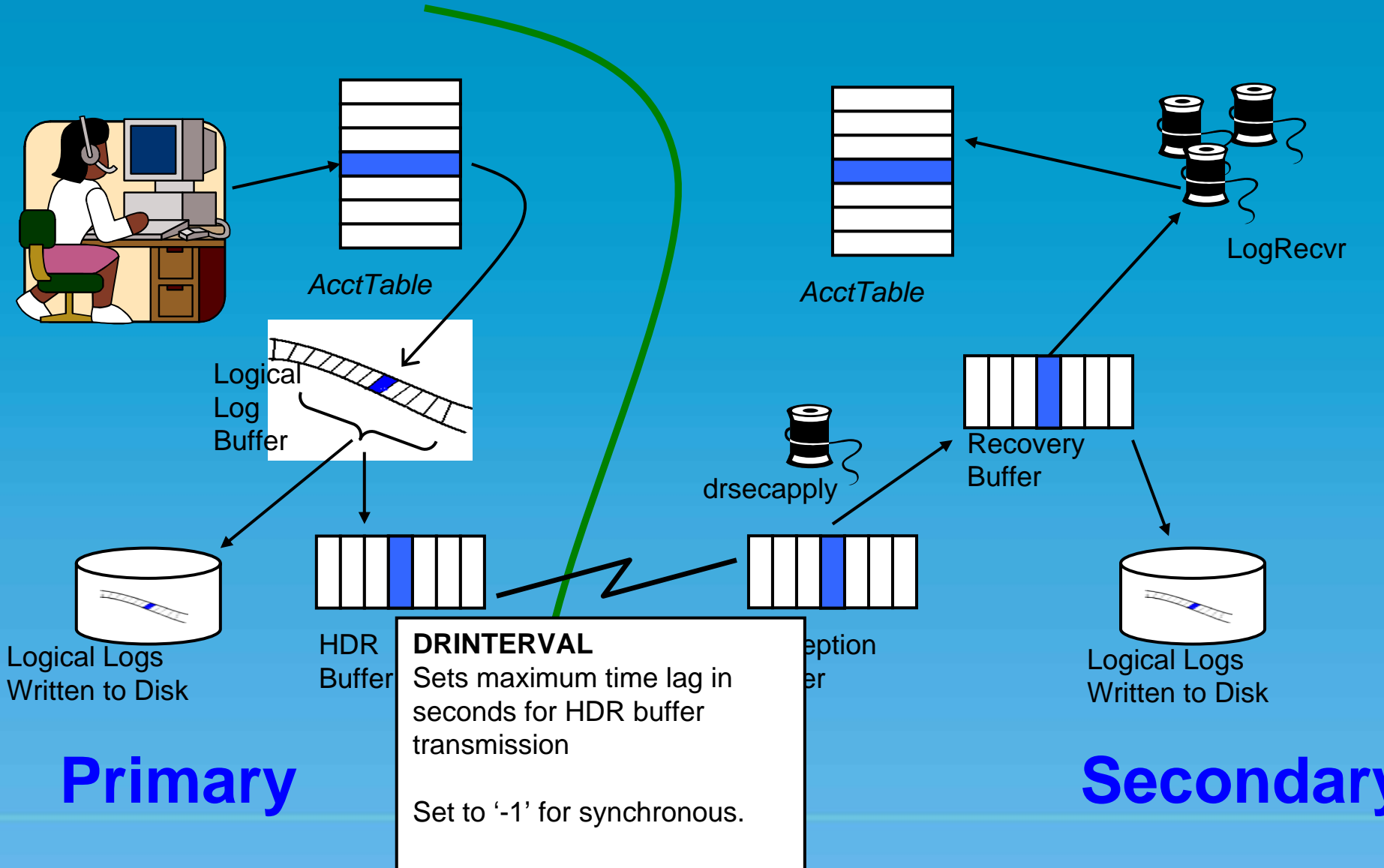
read only

# HDR (cont…)

- Secondary server
  - **Read only server** : allows read only query (Pre cheetah)
  - Always in recovery mode
  - Receives logs from primary and replay them to keep in sync with primary
  - **Cheetah allow read/write operations**
- When Primary server goes down, secondary server takes over as Standard server
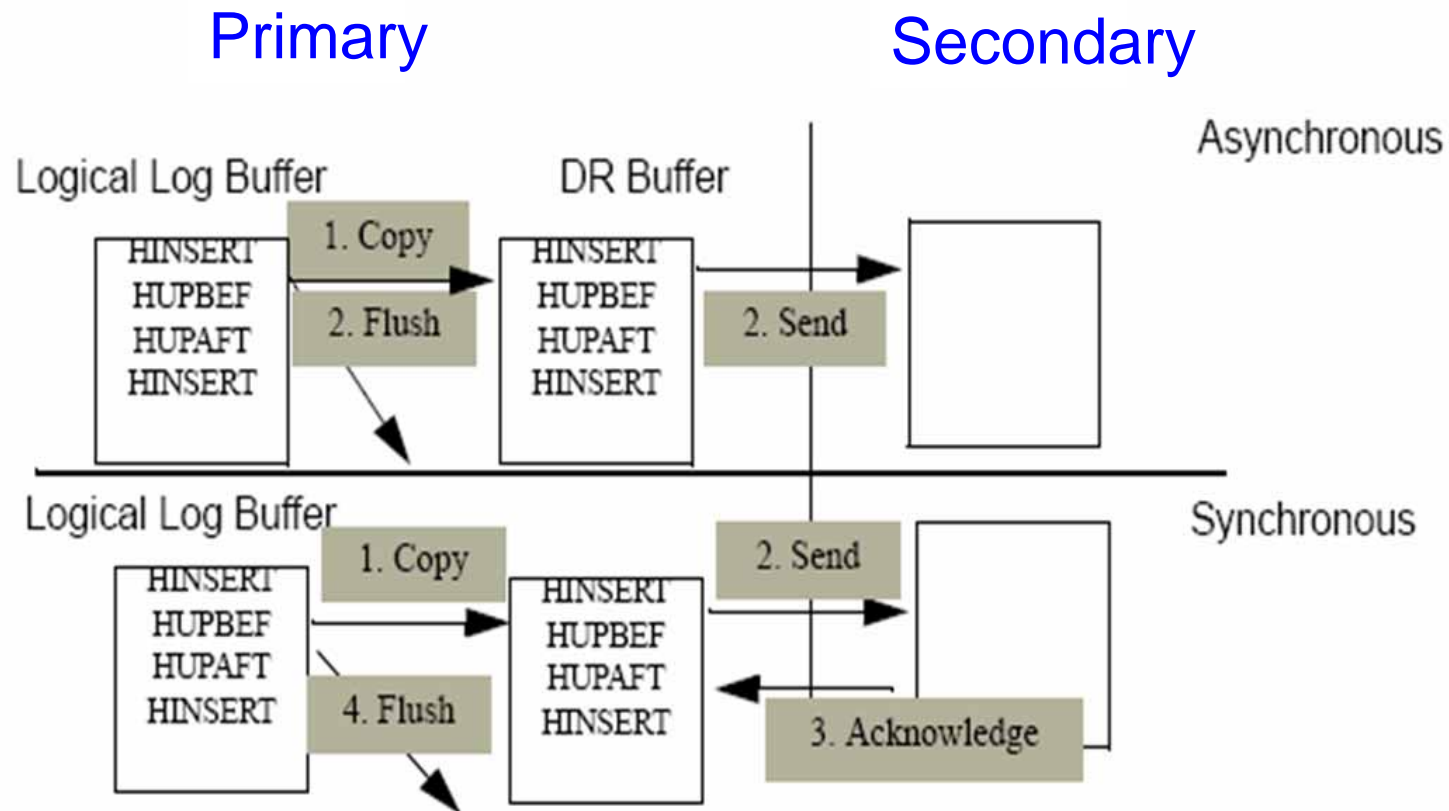
# HDR – How it works

**Primary**

**Secondary**

*AcctTable*

*AcctTable*

LogRecvr

Logical Log Buffer

drsecapply

Recovery Buffer

Logical Logs Written to Disk

HDR Buffer

Logical Logs Written to Disk

**DRINTERVAL**
Sets maximum time lag in seconds for HDR buffer transmission

Set to '-1' for synchronous.

# HDR -- Asynchronous vs Synchronous

# HDR Availability - Failover

Down

onmode –d standard

Standard

A

B

Users reading and updating database

**(Not Any More!!)**

Read-only

Running Reports

# What If Network Is Down: DRAUTO Case

Primary

A

ping timeout

Standard

B

Users reading and updating database
**(Not connected any More!!)**

Running Reports

Users reading and updating database

# HDR Reference Sites

- Walmart                                    (HPUX)
- PZU                                        (AIX)
- Wells Fargo Bank                           (Sun)
- Rakuten                                    (Sun)
- German Border Patrol                       (Sun)
- Phonehouse.DE                              (Sun)
- Huawei                                     (AIX,HPUX,Sun)

Availability with IDS is platform independent

# Factors that Affect HDR Availability

- Failure Detection Time
  - Network speed + OS tuning
- Failover Time
  - HDR takes 3 to 5 seconds
- Recovery Time for Open Transactions
  - When failover occurs, open transactions must be rolled back
  - Time for rollback is dependent on length of transactions
  - OLTP transactions should be short
  - Application may need to be changed to reduce recovery time
  - IDS Cheetah's new Recovery Time Objective feature

Frequency of failures is the main determining factor

# Cheetah – RTO Configuration Parameter

- RTO_SERVER_RESTART allows users to specify a target amount of time the server is allowed for fast recovery.

- RTO_SERVER_RESTART values:
  - 0 – off, uses CKPTINTVL to trigger checkpoints (<= 10.0 functionality).
  - 60 to 1800 seconds (1 – 30 minutes).

- Server will automatically monitor current workload and adjust checkpoint frequency to meet RTO policy.

- Server will fine tune with each fast recovery to improve predictability.

- Dynamically updatable with onmode –wm and –wf.

# HDR vs Hardware HA Solutions

- AIX HACMP
- HP ServiceGuard
- Sun Cluster PDB

Choice of hardware vs. software solution
- Hardware solution is dependent on hardware vendor
- Disk Array is a single point of failure.
- Additional time needed for database initialization.
- Secondary database is not started until failover.
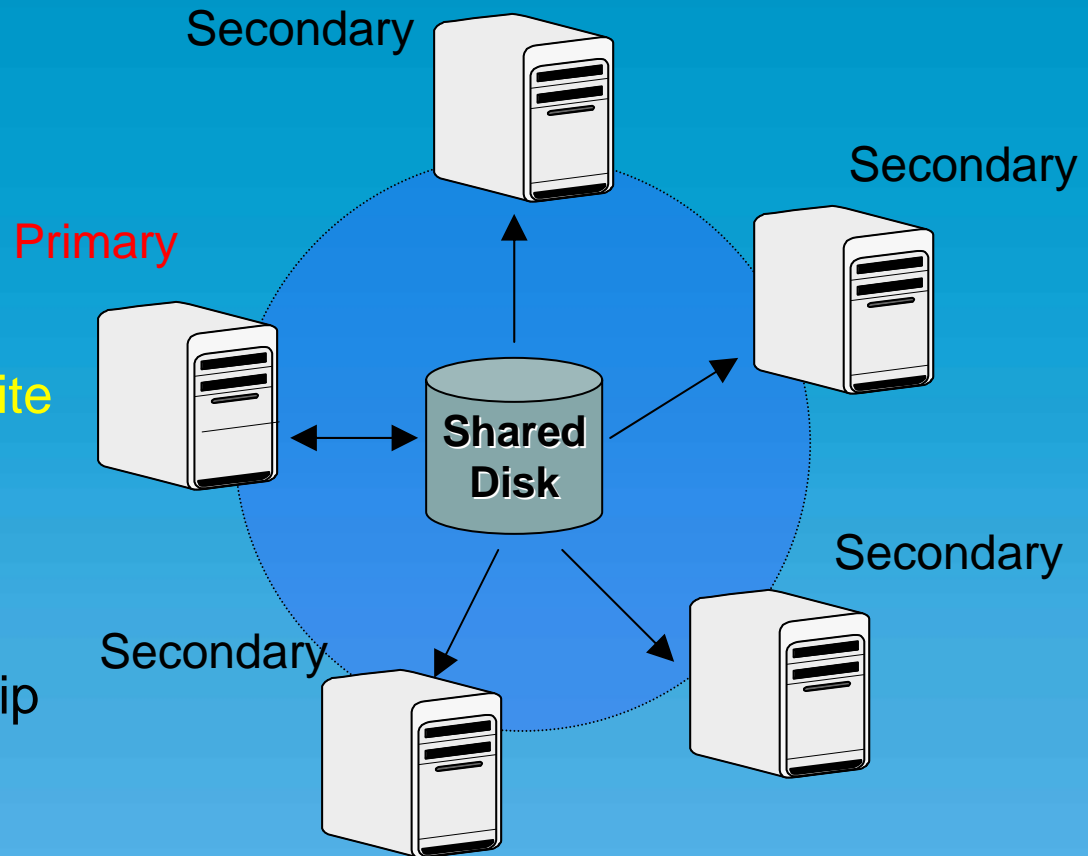- Ultimate choice is cost dependent

# Informix Cheetah's new "MACH-11"

NEW!

- Multi-instance Active Cluster for High Availability.

  – Extends HDR to support more than a single primary with a single secondary instance.

  – Three new types of secondary instances:

    - Shared Disk Secondary (SDS)

    - Remote Standalone Secondary  (RSS)

    - Continuous Log Restore (CLR) or "near-line" standby**

  – HDR, RSS & SDS technology can be used in any combination.

- "MACH-11" is not:

  – Just 1-to-N HDR.

  – "MACH-11" is the treating of all three forms of secondary as a multi-tiered availability solution.
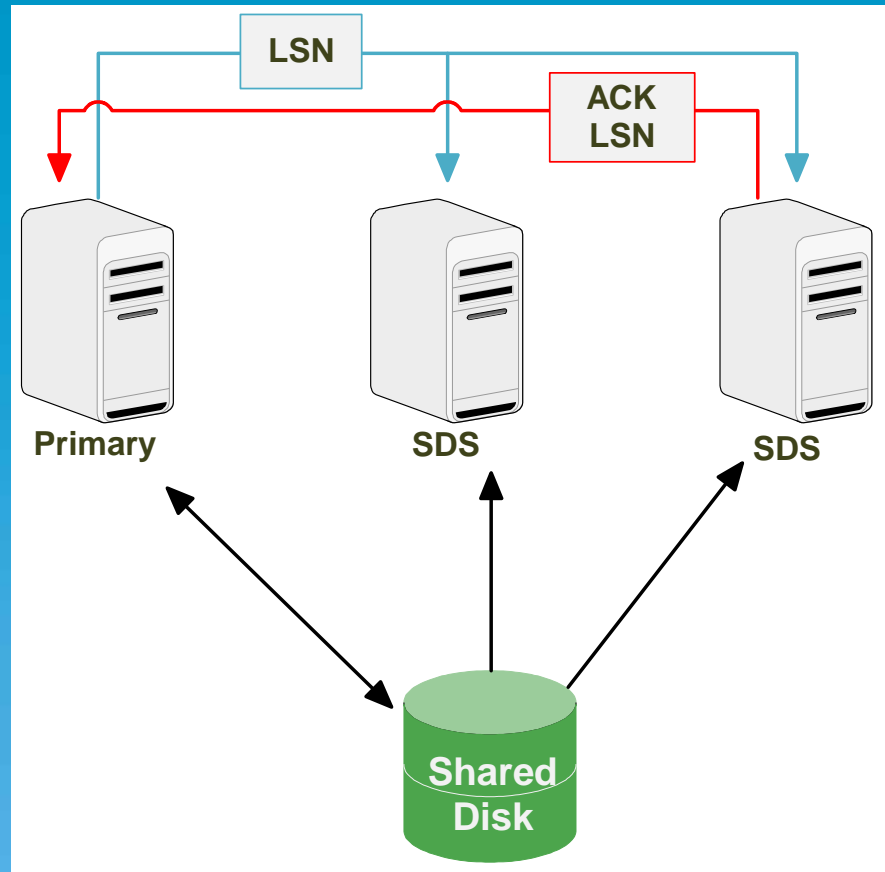
# Shared Disk Secondary

- Share data
  - Single source of data
  - Multiple servers
  - Save data storage
- Optimize Capacity
  - Distributed read/write
  - Flexible
  - Easy scalability
- Multiple Redundancy
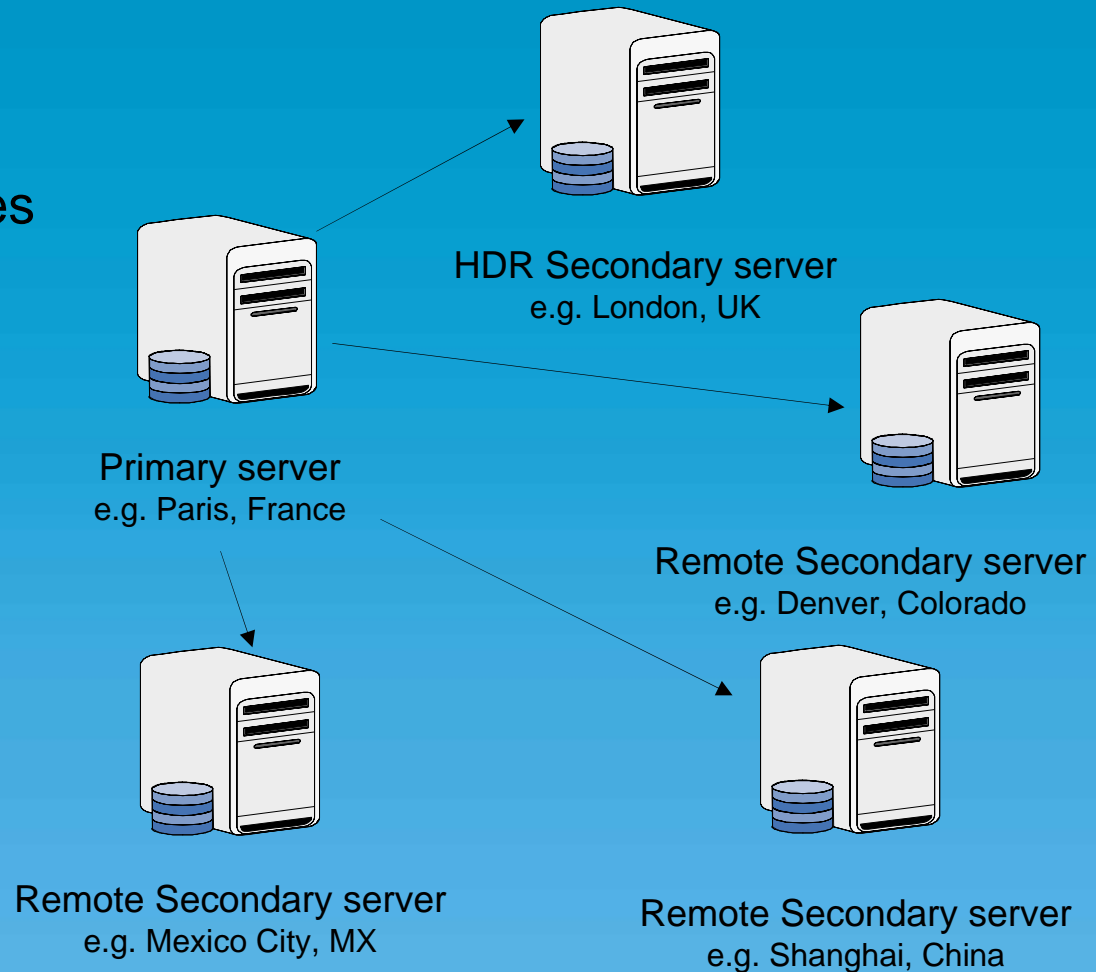- Lower cost of ownership

# Shared Disk Secondary instances

- Primary transmits the current Log Sequence Number (LSN) as it is flushing logs.

- SDS instance(s) receives the LSN from the primary and reads the logs from the shared disks.

- SDS instance(s) applies log changes to its buffer cache.

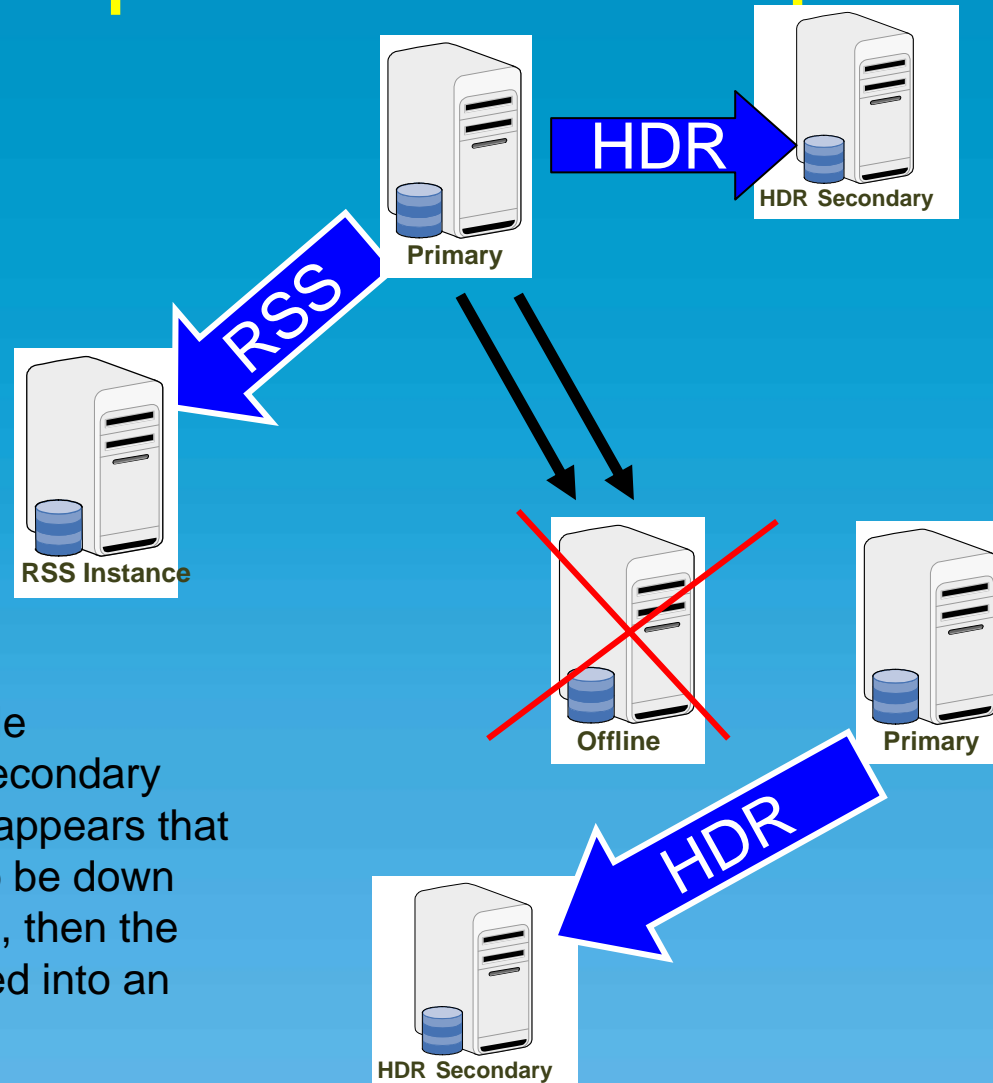- SDS instance(s) resynch processed LSN to primary.

# Remote Standalone Secondary

## *Evolution of HDR*

- Multiple hot backups
  - Additional Secondaries
  - Global
  - Extends HDR

- Optimize Capacity
  - Distribute workload
  - Read/write locally
  - Improve performance
  - Increase capacity
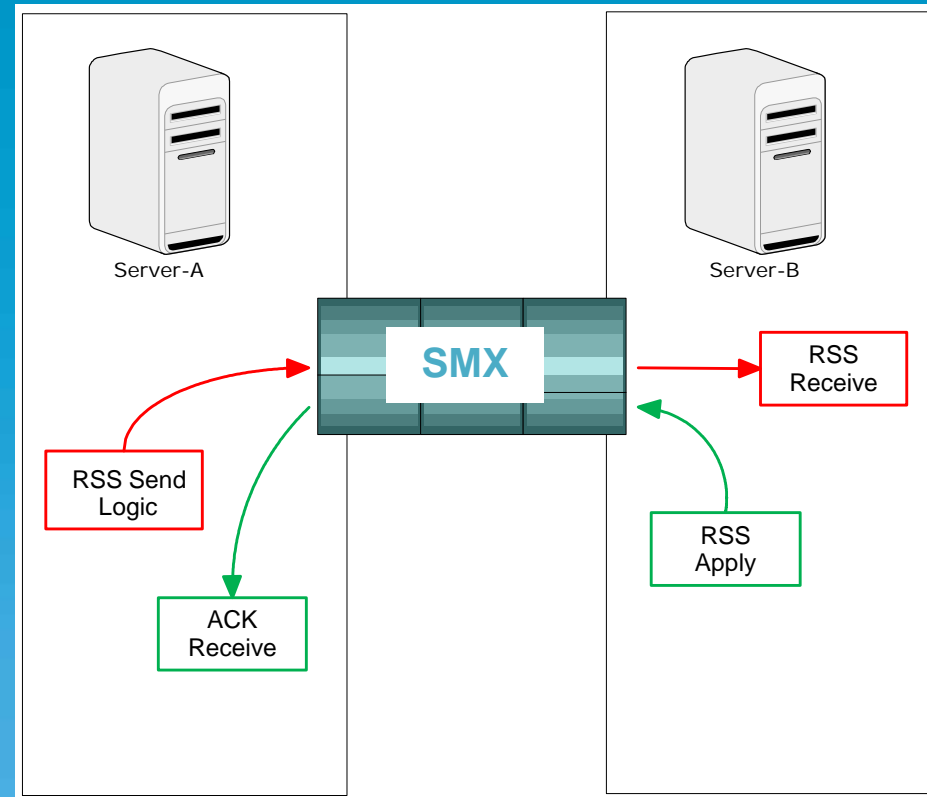
HDR Secondary server
e.g. London, UK

Primary server
e.g. Paris, France

Remote Secondary server
e.g. Denver, Colorado

Remote Secondary server
e.g. Mexico City, MX

Remote Secondary server
e.g. Shanghai, China

# RSS − Backup to the Backup



If the primary fails, it is possible to convert the existing HDR secondary into the primary instance. If it appears that the original primary is going to be down for an extended period of time, then the RSS instance can be converted into an HDR secondary instance.

# Server Multiplexer (SMX)

- Multiplexed network connection.

- Uses full duplex protocol
  - Sends packets without waiting for "ack" back:
  - HDR uses half-duplex so primary knows secondary received transaction before committing.

- Supports encryption.

- Automatically activated.

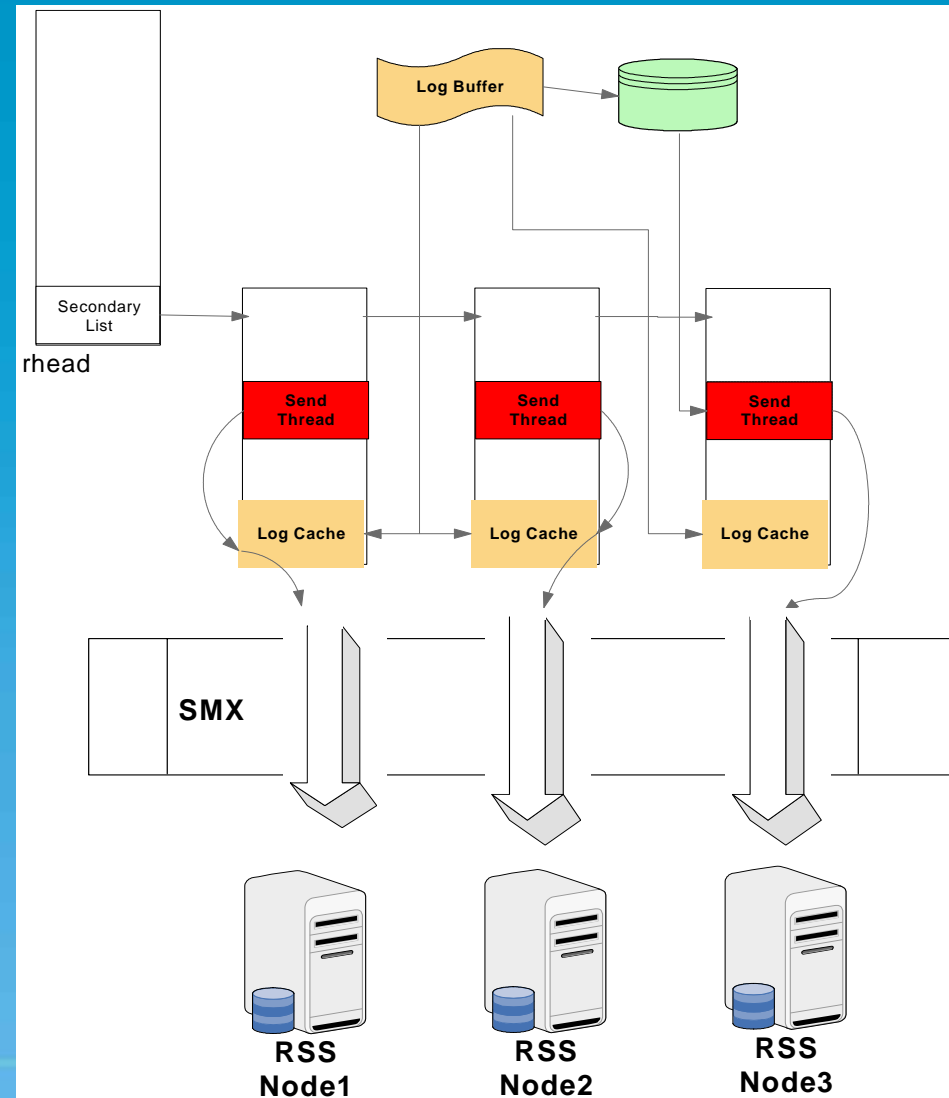- Requires no configuration other than encryption.

Server-A        SMX        Server-B

RSS Send Logic

ACK Receive

RSS Receive

RSS Apply

# RSS: Primary – How It Works

In the current IDS engine, when we flush a log buffer to disk, we also copy that log page to the HDR log buffer and to the ER buffer cache. The primary now also checks to see if there are RSS secondary nodes and will copy that page to a log cache which is used to send that page to the remote node. If the send thread is currently sleeping, it will wake up the RSS_Send thread which will transmit that log page to the remote server.

It is possible that the next page which needs to be sent is not in the log cache. In that case, the RSS_Send thread will read the log pages directly from disk.
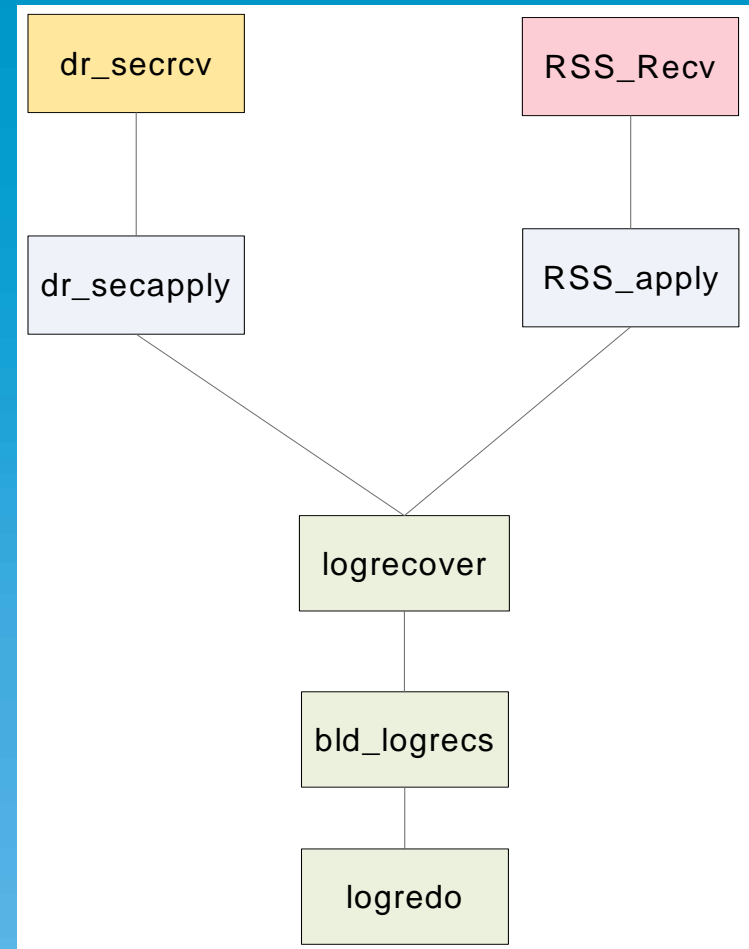
The RSS_Send thread works in a fully duplexed network model. That means that the thread does not wait for an ACK before it sends the next buffer. Flow control is implemented in that up to 32 buffer transmissions will be sent before it requires an acknowledgement from the secondary node. If that limit is reached, then the send thread will block and wait for the RSS_Recv thread to receive an ACK from the peer node.
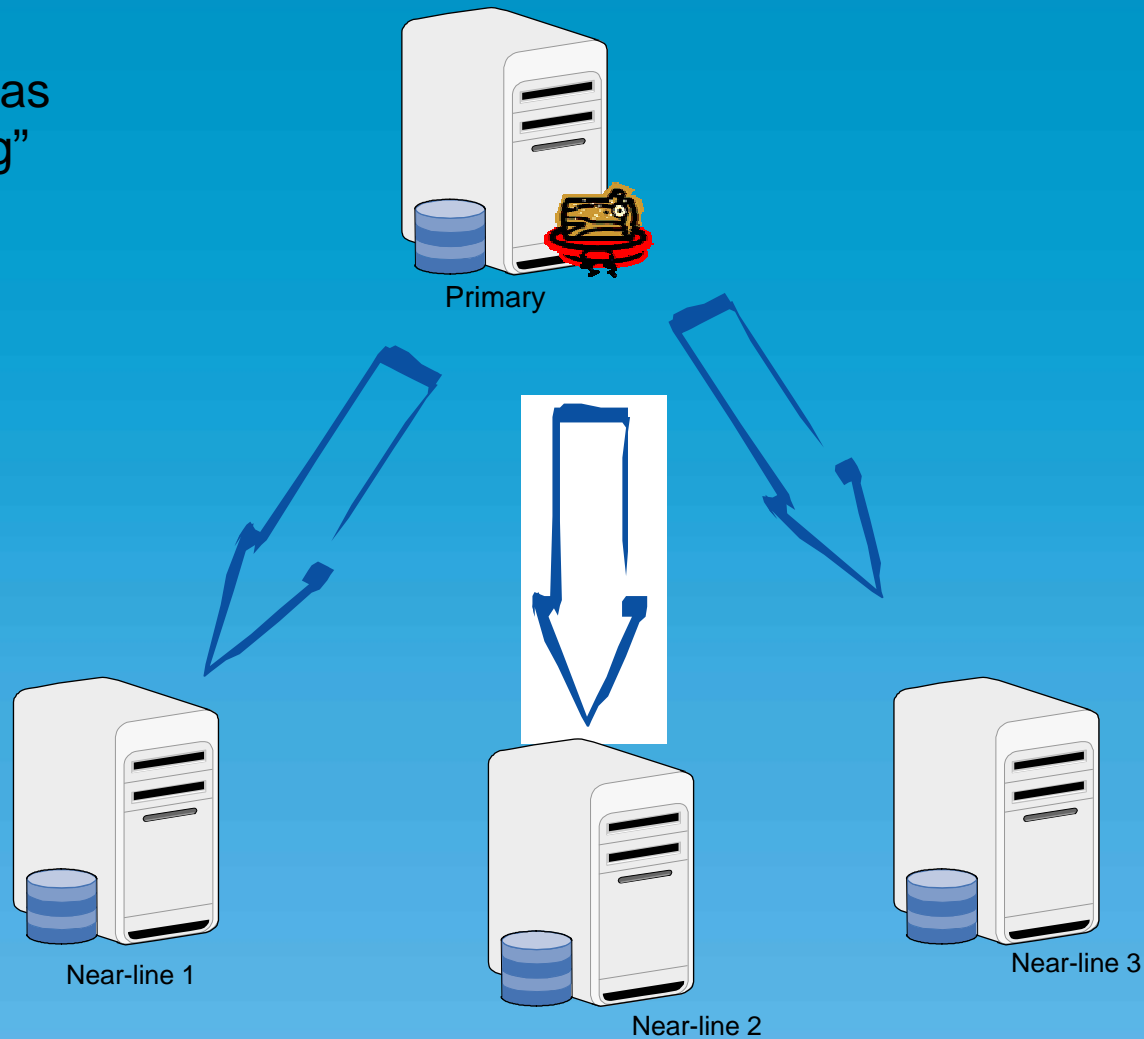
# RSS: Secondary – How it Works

Existing HDR has a dr_secrcv thread which is responsible for receiving log pages from the HDR primary and interfaces with the HDR buffer pool on the secondary node.  This is replaced by RSS_Recv which interfaces with SMX to receive the log pages from the primary node. RSS_apply is a minor modification of dr_secapply and is 98% common code with the HDR code which is used by dr_secapply. The rest of the recovery logic is common with existing HDR code.
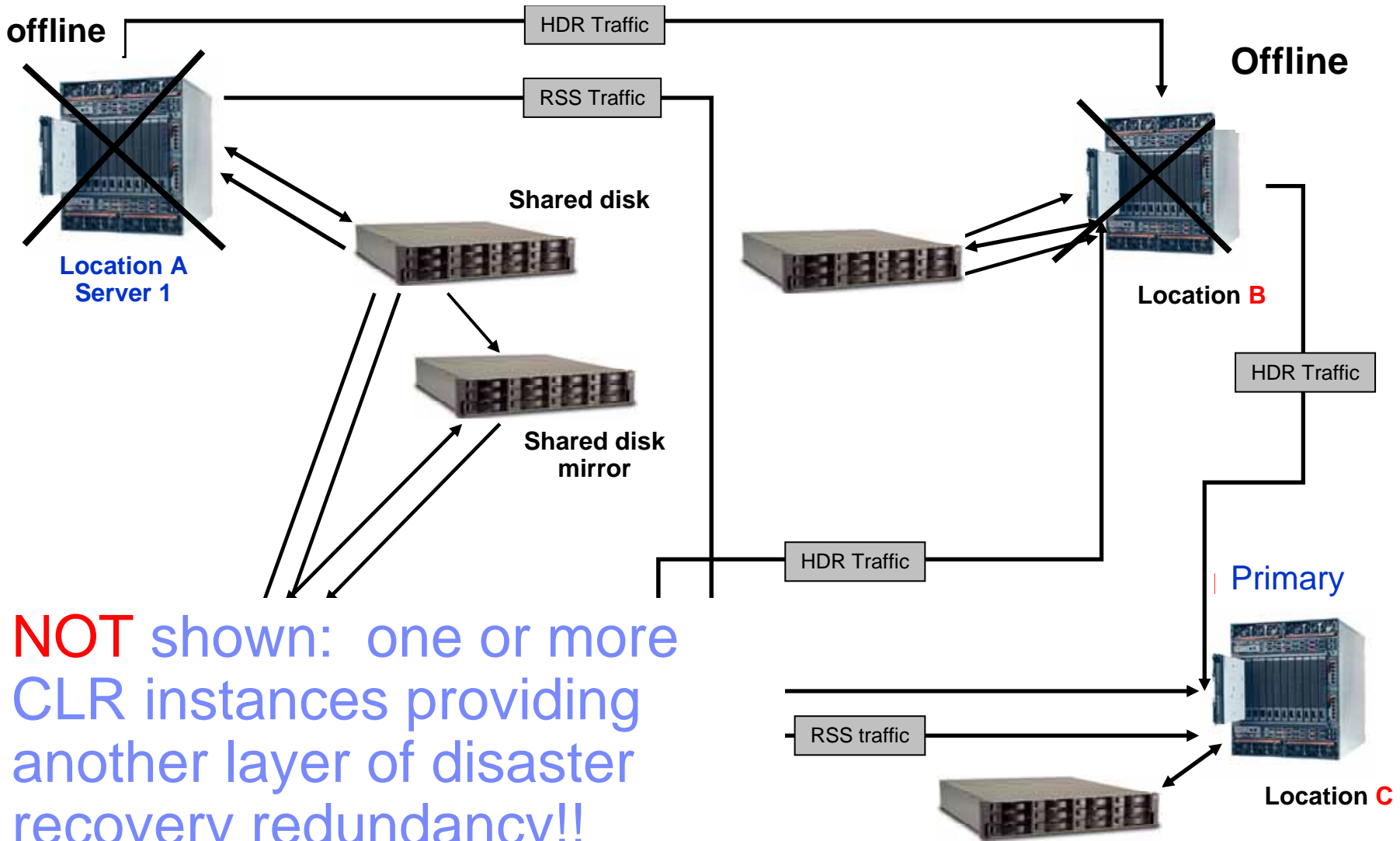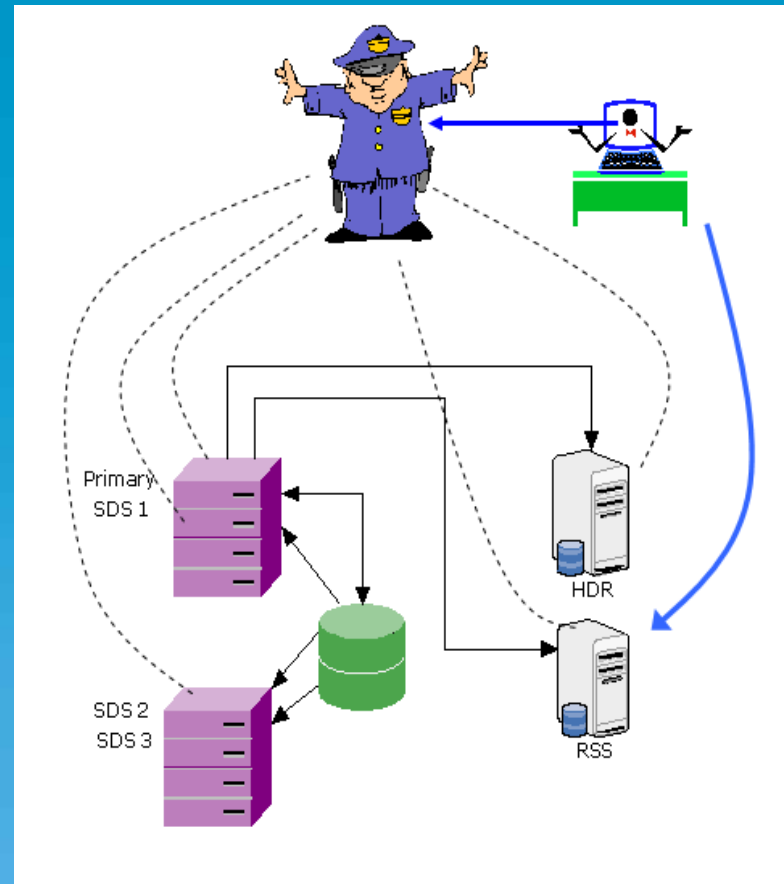
# Near-line Standby

Also known as
"log shipping"

Primary

Near-line 1

Near-line 2

Near-line 3

# The server at location B fails

**offline**

HDR Traffic

**Offline**

RSS Traffic

**Shared disk**

**Location A
Server 1**

**Location B**

**Shared disk
mirror**

HDR Traffic

HDR Traffic

Primary

NOT shown: one or more CLR instances providing another layer of disaster recovery redundancy!!

RSS traffic

**Location C**

# Connection Manager

- Maintains knowledge of all nodes within the cluster

- Records adding/removal of nodes

- Monitors type of node

- Monitors workload of nodes

- Routes the client application to target node

- Automatic failover

- Works on the concept of class of service by resolving the following requirements

•Connect to the best possible secondary
•Connect to the current primary
•Connect to the SDS node or primary with the most free CPU cycles
•Connect to either the HDR primary or the HDR secondary
•Connect to an SDS node or HDR secondary, if any are currently active, otherwise connect to the primary
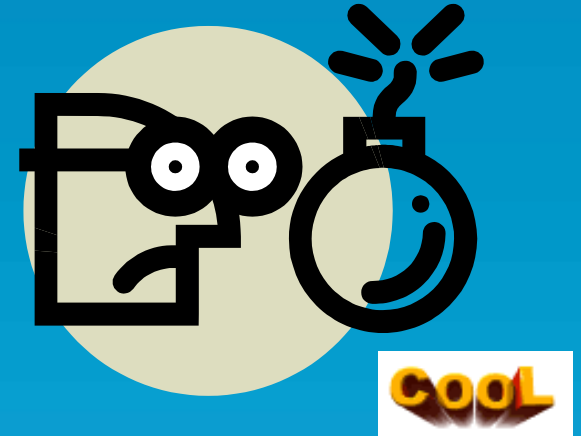
# Arbitrator

- Purpose is to quickly detect the failure of the current primary.
- Part of the connection manager
- Gets confirmation through alternate paths before performing a failover
- Performs a failover using admin API interface
- There is failover for the arbitrator just as there is failover for the connection manager.

# Failover

- SDS can be used with HDR and RSS.

- SDS can be promoted directly to primary:

    onmode –d set primary SDS <sds_instance_to_promote>

    When the primary is changed, other instances 'follow the change'.

- Order of failover should be:
    1. To an SDS instance.
    2. To the HDR secondary.
    3. To an RSS instance.

# MACH11 – Database Cluster

- Redirected Writes
  - Allow HDR,SDS,RSS to handle queries which modify data
    - Update processing spread between SDS node and Primary
- Network Services
  - Runs as a separate middle layer outside IDS
  - Clients connect to this middle layer rather than directly to an SDS or Primary node
  - Detects when an RSS or SDS node fails
    - Provides application redirection
  - Provides automatic failover to RSS or SDS nodes
- Administration via IDSAdmin
  - Configure and start new nodes
  - Better graphical representation of clusters

# IDSAdmin Tool

# IDS MACH11 LifeDemo

➢ Informix MACH11 Database cluster
➢ Informix MACH11 High Availability



Mach11. wmf