# Greenplum.

# Library of Congress

# I/O Considerations in Big Data Analytics

**26 September 2011**

Marshall Presser
Federal Field CTO
EMC, Data Computing Division

# Paradigms in Big Data

Structured (relational) data
  Very Large Databases (100's TB +)
  SQL is the access method
  Can be monolithic or distributed/parallel
  Vendors distribute software only or appliance

Unstructured (non-relational data)
  Hadoop Clusters (100's + nodes)
  Map/Reduce is the access method
  Vendors distribute software only (mostly)

**Greenplum.**

EMC²
where information lives'

2

# Obstacles in Big Data

Both Relational and Non Relational Approaches must deal with I/O issues:

- Latency
- Bandwidth
- Data movement in/out of cluster
- Backup/Recovery
- High Availability

**Greenplum.**

EMC²
where information lives

# MPP (Massively Parallel Processing) Shared-Nothing Architecture

- MPP has extreme scalability on general purpose systems

- Provides automatic parallelization
  - Just load and query like any database
  - Map/Reduce jobs run in parallel

- All nodes can scan and process in parallel
  - Extremely scalable and I/O optimized

- Linear scalability by adding nodes
  - Each adds storage, query, processing and loading performance



SQL
MapReduce

**Master Servers**
Query planning & dispatch

**Network Interconnect**

**Segment Servers** …
Query processing & data storage

**External Sources**
Loading, streaming, etc.

Greenplum.

EMC²
where information lives®

4

# Software and Appliances in Relational Big Data

Greenplum DCA – EMC   (software and appliance)

Neteeza Twin Fin – IBM (appliance only)

Teradata 2580 – Teradata ( appliance only)

Vertica – HP (software and appliance)

> All above use distributed data with conventional I/O
>
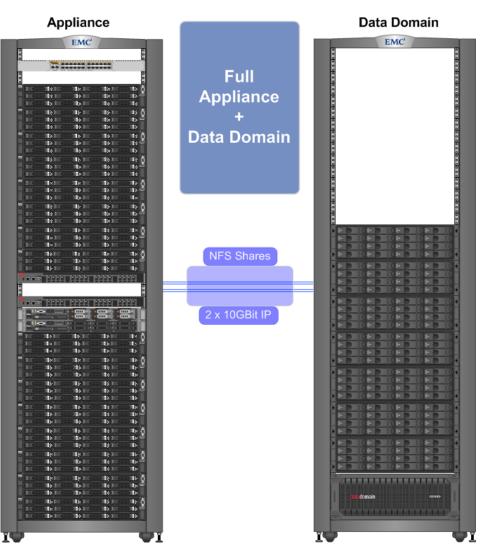> Neteeza and Teradata virtual proprietary network s/w

Exadata – Oracle (appliance only)

> Oracle is the only vendor with a shared disk model
>
> Uses Infiniband to solve latency and bandwidth  issues

**Greenplum.**

EMC²
where information lives'

# Backing up to from an Appliance

- Requirements:
- Parallel backup from all nodes, not just the master
- Incremental or dedup ability via NFS shares or similar
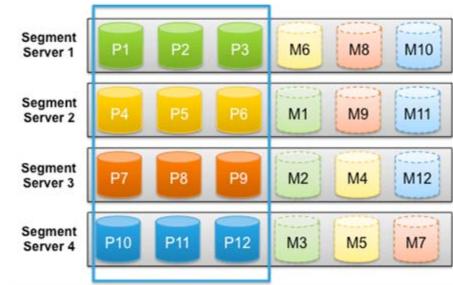- Conneted to private network, not public

**Appliance**

**Data Domain**

Full Appliance + Data Domain

NFS Shares

2 x 10GBit IP

Greenplum.

EMC²
where information lives

# MPP Database Resilience Relies on In-Cluster Mirroring Logic

- Cluster comprises
  - Master servers
  - Multiple Segment servers
- Segment servers support multiple database instances
  - Active primary instances
  - Standby mirror instances
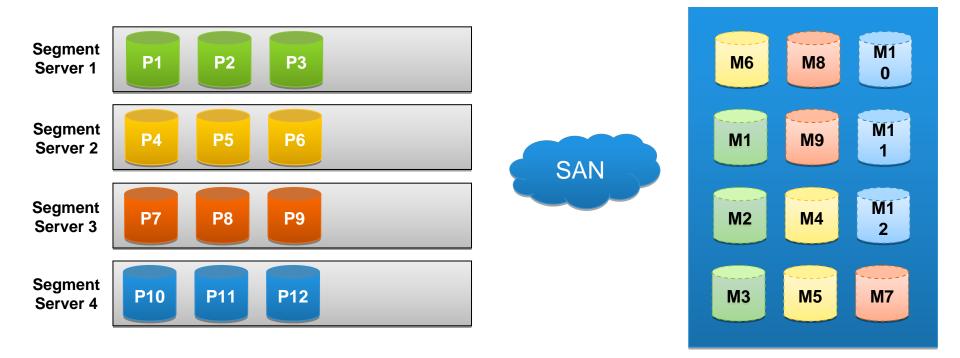- 1:1 mapping between Primary and Mirror instances
- Synchronous mirroring

| Segment Server 1 | P1 | P2 | P3 | M6 | M8 | M10 |
| Segment Server 2 | P4 | P5 | P6 | M1 | M9 | M11 |
| Segment Server 3 | P7 | P8 | P9 | M2 | M4 | M12 |
| Segment Server 4 | P10 | P11 | P12 | M3 | M5 | M7 |

Set of Active Segment Instances

Greenplum.

Data Computing Division

EMC²
where information lives

7

# SAN Mirror Configuration: Mirrors Placed on SAN Storage

**Segment Server 1**

P1  P2  P3

**Segment Server 2**

P4  P5  P6

**Segment Server 3**

P7  P8  P9

**Segment Server 4**

P10  P11  P12

SAN

M6  M8  M10

M1  M9  M11

M2  M4  M12

M3  M5  M7

Doesn't this violate principle of all local storage?  Maybe, maybe not.

**Greenplum.**

**EMC²**
where information lives®
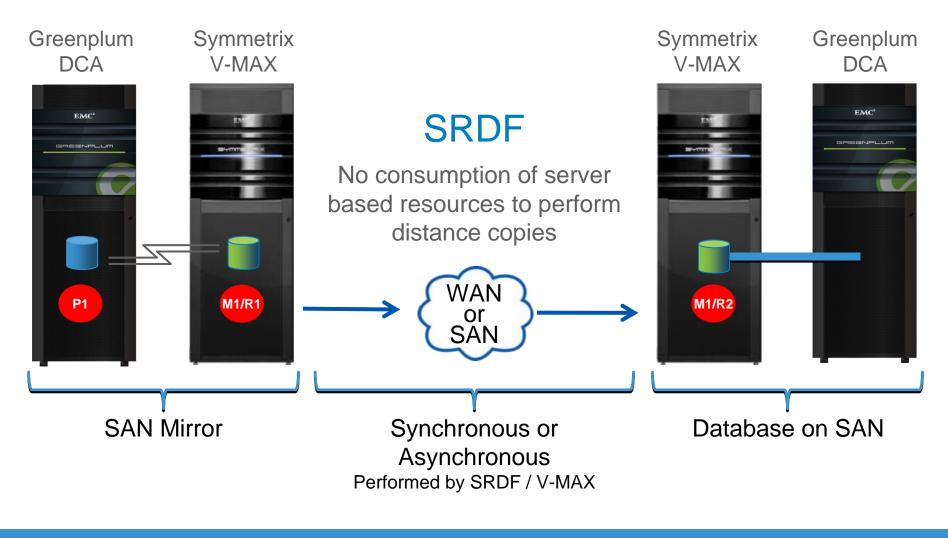
8

# One Example: SAN Mirror to VMAX SAN

- Default DCA configuration has Segment Primaries and Segment Mirrors on internal storage

- SAN Mirror offloads Segment Mirrors to VMAX SAN storage
  - Doubles effective capacity of a DCA
  - Foundation of SAN leverage
  - Seamless off-host backups
  - Data replication

- No performance impact
  - Primaries on internal storage
  - SAN sized for load and failed segment server

**P1**

...

**P96**

**M1**

...

**M96**

Greenplum.

Data Computing Division

EMC²
where information lives®

# One Example: SAN Mirror –With SAN based replication for DR

Greenplum
DCA

Symmetrix
V-MAX

## SRDF

No consumption of server based resources to perform distance copies

Symmetrix
V-MAX

Greenplum
DCA

P1

M1/R1

WAN
or
SAN

M1/R2

SAN Mirror

Synchronous or
Asynchronous
Performed by SRDF / V-MAX

Database on SAN

Greenplum.

EMC²
where information lives"

# What is Hadoop?

Three major components

- An infrastructure for running Map/Reduce jobs
  - Mappers produce name/value pairs
  - Reducers aggregate Mapper Output

- HDFS - A distributed file system for holding input data, output data, and intermediate result

- An ecosystem of higher level  tools overlaid on  MapReduce and HDFS
  - Hive
  - Pig
  - Hbase
  - Zookeeper
  - Mahout
  - Others

Data Computing Division

**Greenplum.**

**EMC²**
where information lives

# Why Hadoop?

- With massive growth of unstructured data Hadoop has quickly become an important new data platform and technology
  - We've seen this first-hand with customers deploying Hadoop alongside relational databases

- A large number of major business/government agency are evaluating Hadoop or have Hadoop in production

- Over 22 organizations running 1+ PB Hadoop Clusters

- Average Hadoop cluster is 30 nodes and growing.

**Greenplum.**

**EMC²**
*where information lives*

# Why Not Hadoop?

Hadoop still a "roll your own" technology

Appliances just appearing Sep/Oct 2011

- Wide scale acceptance requires
  - Better HA features
  - More performant I/O
  - Ease of use and management
- Access to HDFS through a single Name Node
  - Single point of failure
  - Possible contention in large clusters
  - All Name Node data held in memory, limiting number of files in cluster
- Unlike SQL, programming model via Hadoop API a rare skill
- Apache distribution written in Java, good for portability, less good for speed of execution

**Greenplum**

EMC²
where information lives

# Storage Layer Improvements to Apache Hadoop Distribution

- HDFS optimizations
  - Recoded in C, not Java, different I/O philosophy
  - Completely API compatible

- NFS interface for data movement in/out of HDFS

- Distributed Name Node eliminates SPOF for Name-Node

- Remote Mirroring and Snapshots for HA

- Multiple readers/writers – lockless storage

- Built-in transparent compression/encryption

Greenplum.

EMC²
where information lives®

# Thank you

Marshall Presser|   240.401.1750 – cell   |   marshall.presser@emc.com