



Data Computing Appliance

Presented by Emin ÇALIKLI
Sr. Technology Consultant
emin.calikli@emc.com



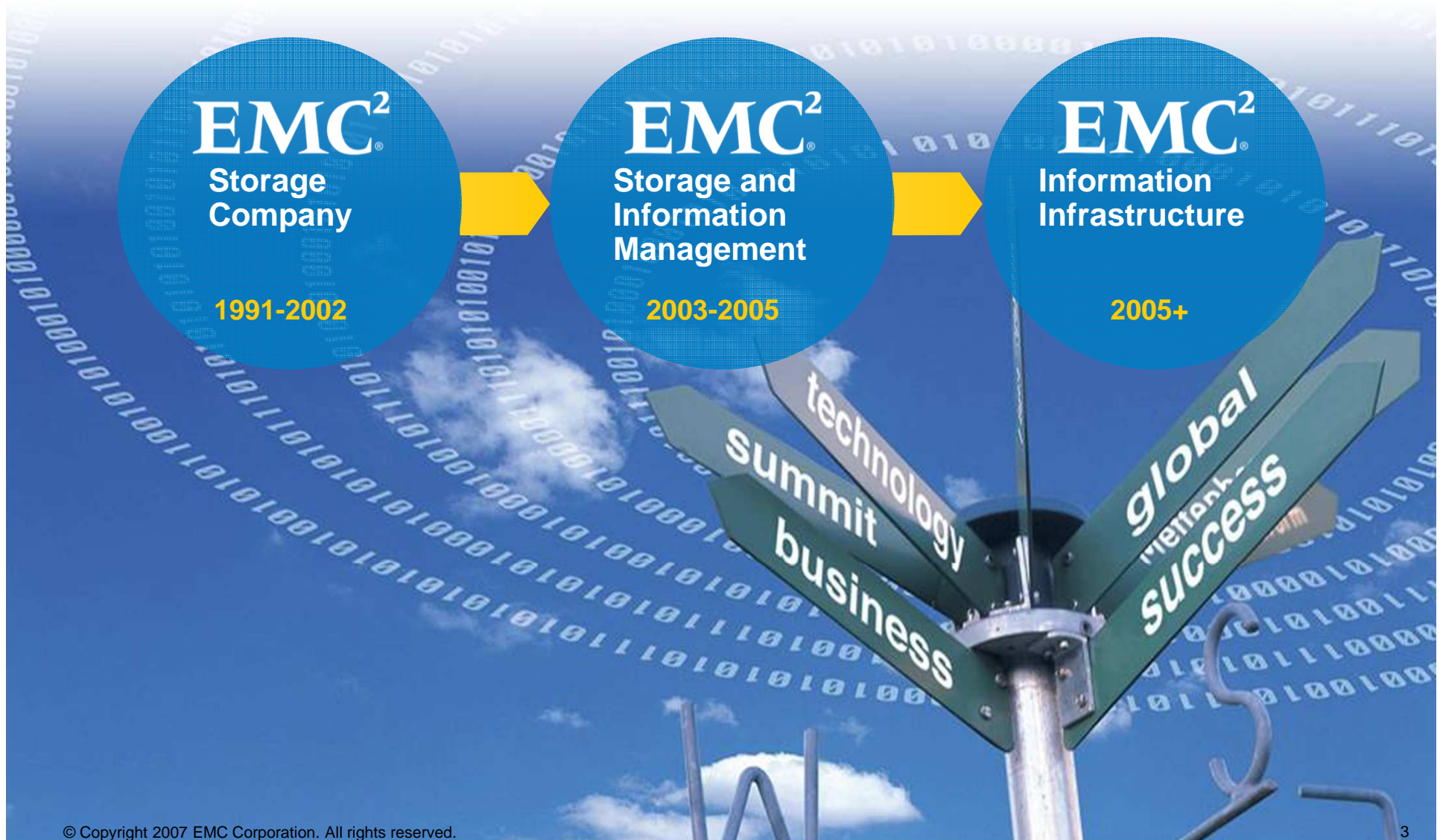
Data Computing Division

EMC²
where information lives™

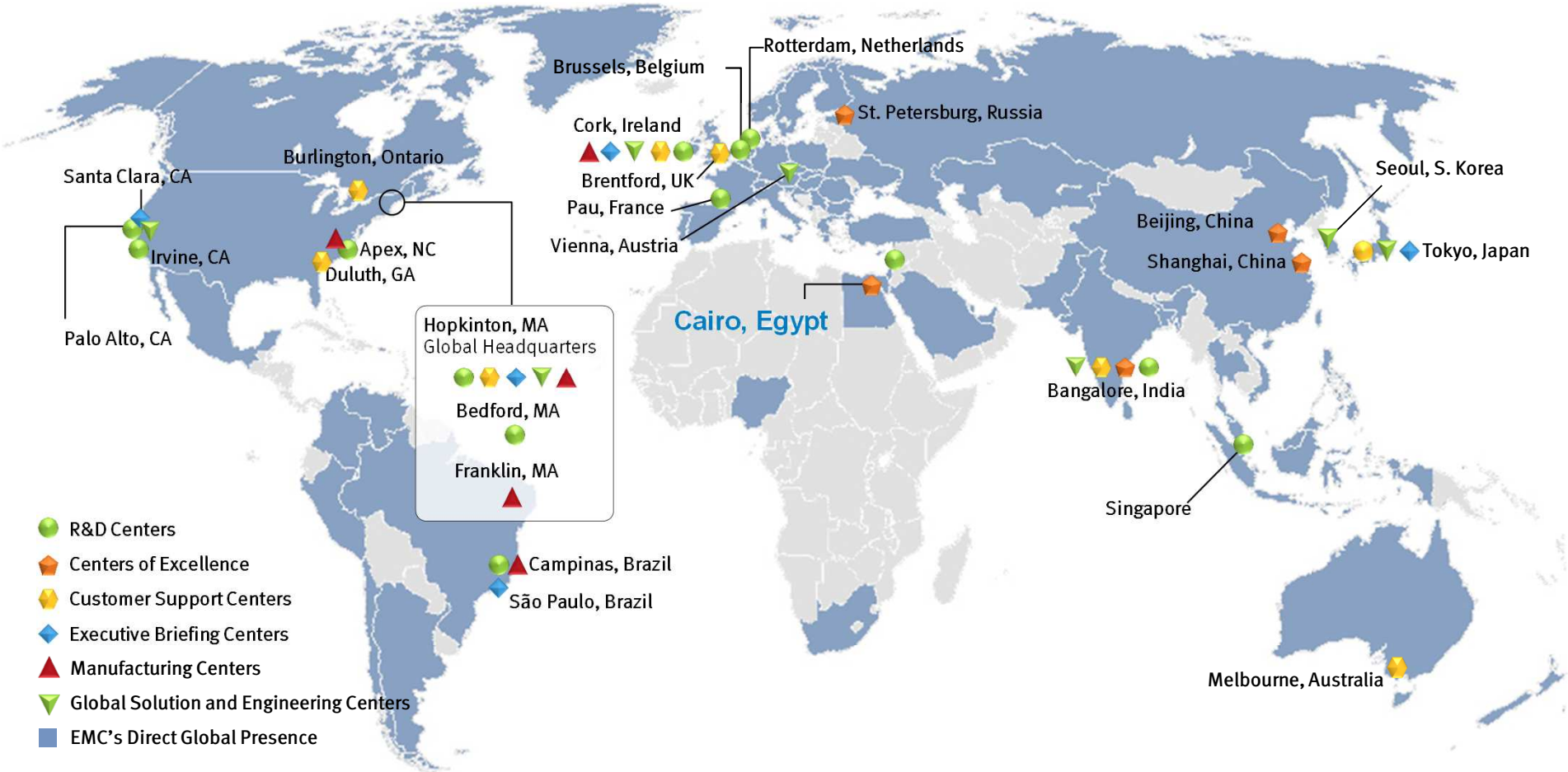
Agenda

- Who is EMC?
- Greenplum Innovative Technology
- Greenplum Data Computing Appliance

EMC's Evolution



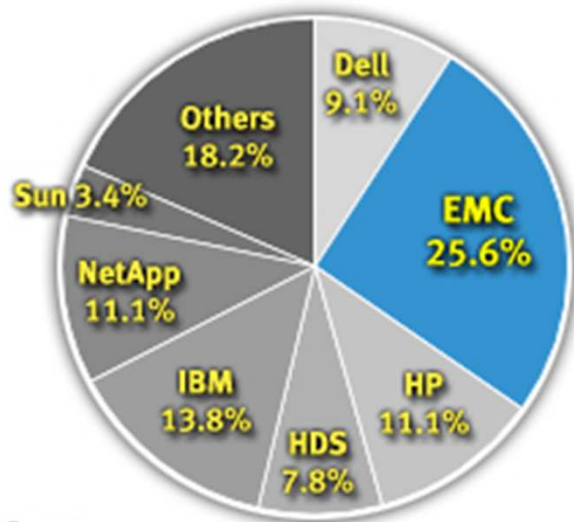
400+ Offices with Operations in 60+ Countries 43,000 Employees



Q4 '10 Highlights

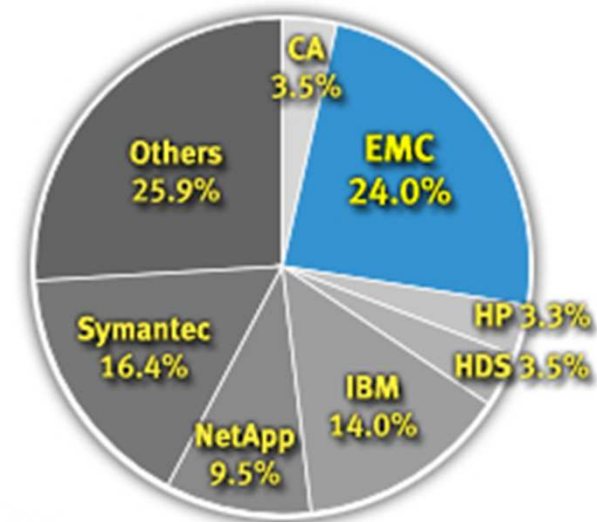
Thriving in a Challenging Global Economy

W/W External Disk Storage Systems, 2010



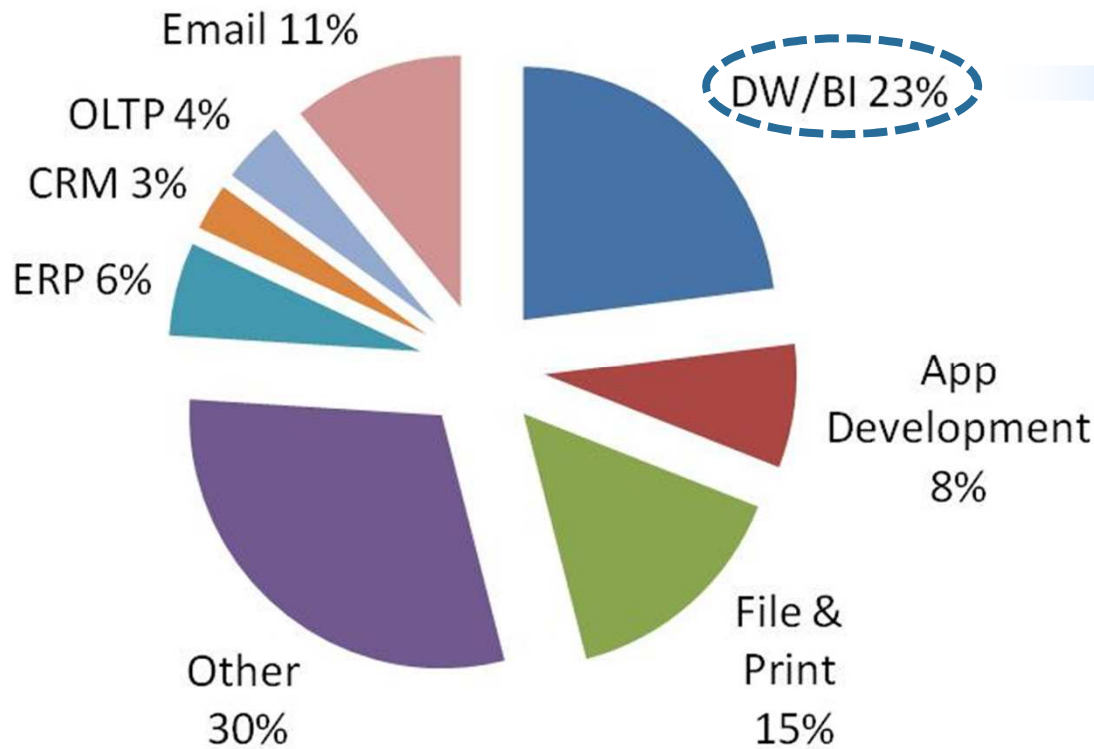
Source: IDC W/W Quarterly Disk Storage Systems Tracker, 3/11.

W/W Storage Software, 2010



Source: IDC W/W Quarterly Storage Software Tracker, 3/11.

DW/BI Storage Volumes



DW/BI is 23% of external storage workloads

Source: IDC, Sep 2009



THE ANSWER MACHINE

DATA IN. DECISIONS OUT.

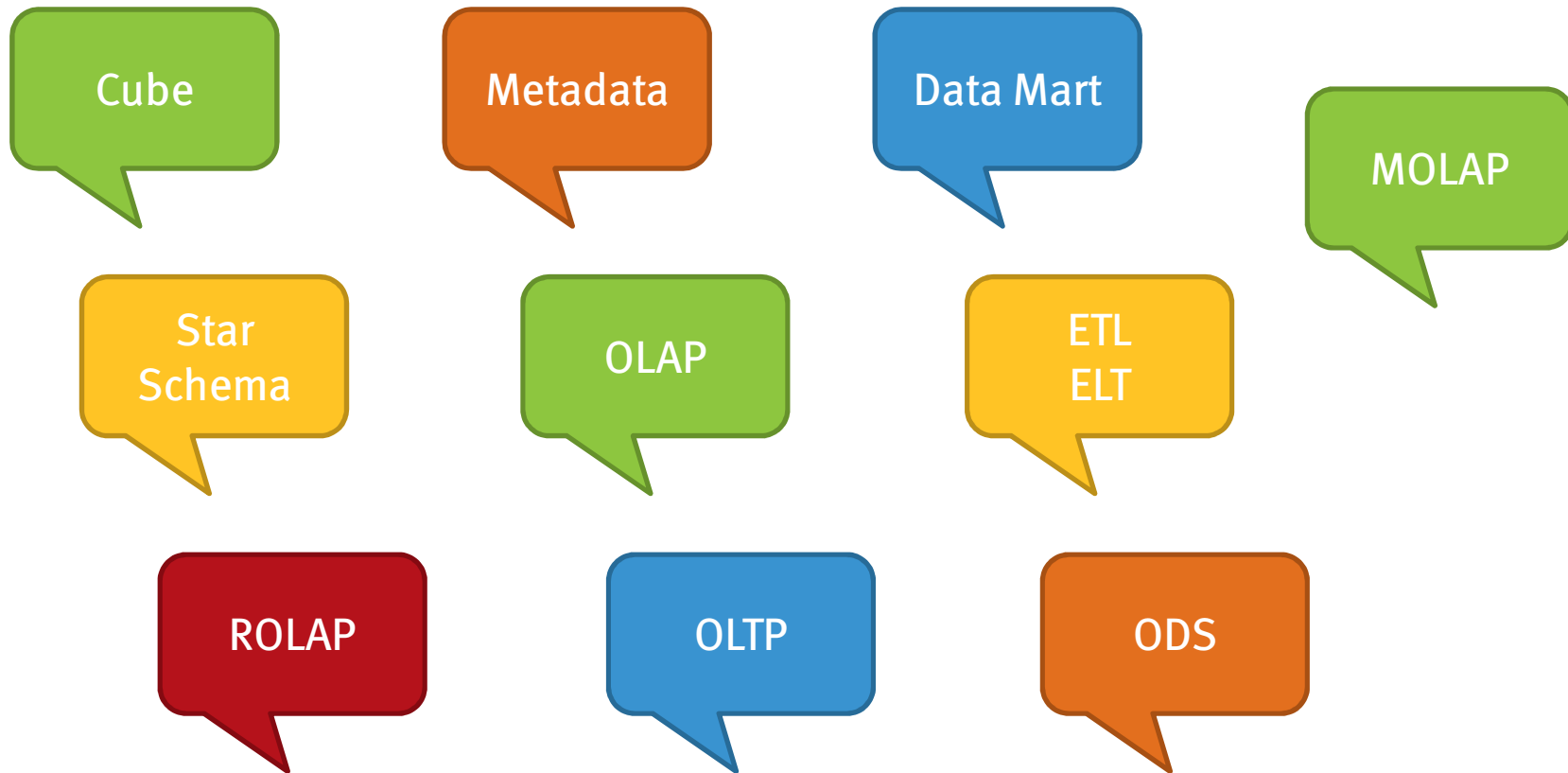
Introducing the
Greenplum Data Computing Appliance



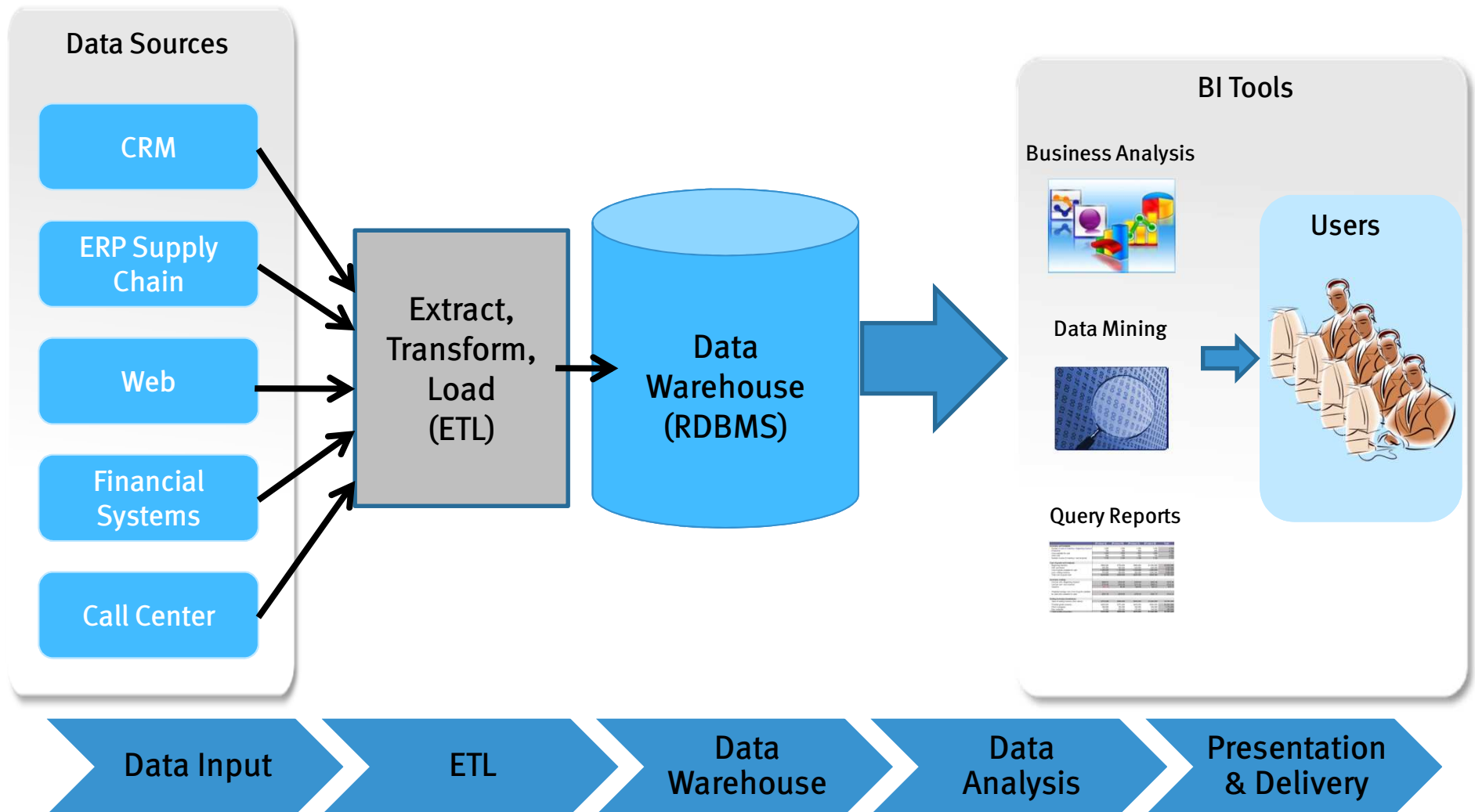
Data Computing Division



Data Warehousing?



Information Flow in a DW/BI Environment



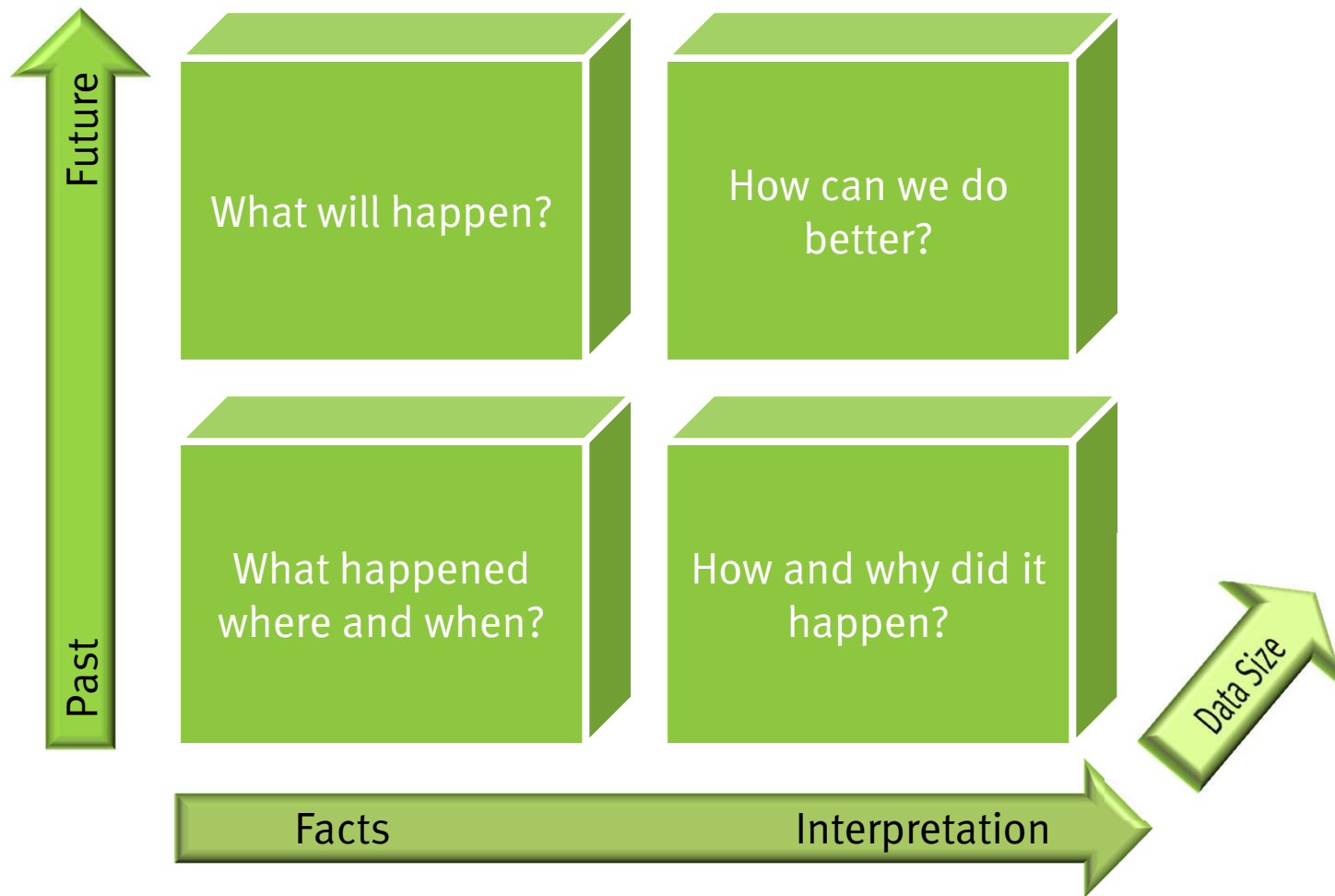
Greenplum Architecture

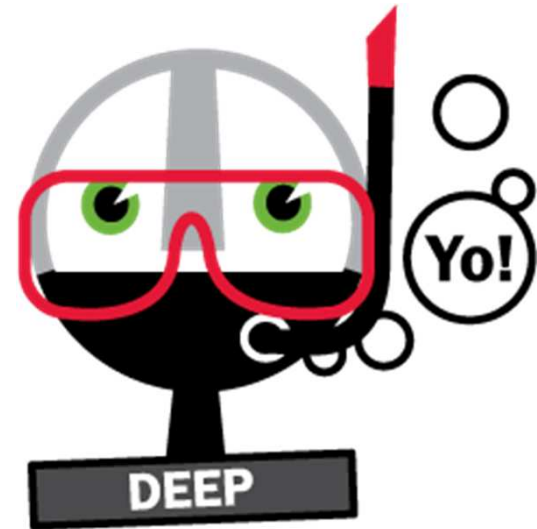
- **MPP = Massively Parallel Processing**
 - Two or more Servers (with own CPU/RAM/Disk) working on the same task
 - Multiple units of parallelism working together
 - Parallel Database Operations
 - Parallel CPU Processing
 - **Segments** = Greenplum Units of Parallelism (one Postgres database)
- **'Shared Nothing' Architecture**
 - Each Segment is a separate Postgres Database
 - Segments only operate on their portion of the data
 - Segments are self-sufficient
 - Dedicated CPU Processes
 - Dedicated storage that is only accessible by the Segment

Things are changing

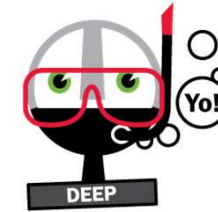
Yesterday's Data Warehouse and Analytic Infrastructure	The Greenplum Future
Proprietary	Commodity
Expensive	Cost-Effective
Centralized, Monolithic	Distributed
Process-Heavy	Self-Service
Batch	Real-Time
Summarized	Deep
Slow	Agile

Deeper Analysis – Bigger Business Value





MAD Skills



- **Business Intelligence**

- MicroStrategy
- Business Objects
- Cognos
- SAS
- Informatica
- ...

- **The Analytics Library**

- Logistic regression
- Mann-Whitney U Test
- Chi-Square Test
- Sparse vectors, matrices
- NLTK
- ...

- **Built-in Analytics**

- Multiple linear regression
- Naïve Bayes
- Matrix operations
- Window functions, OLAP
- PL/R
- ...

- **Methods**

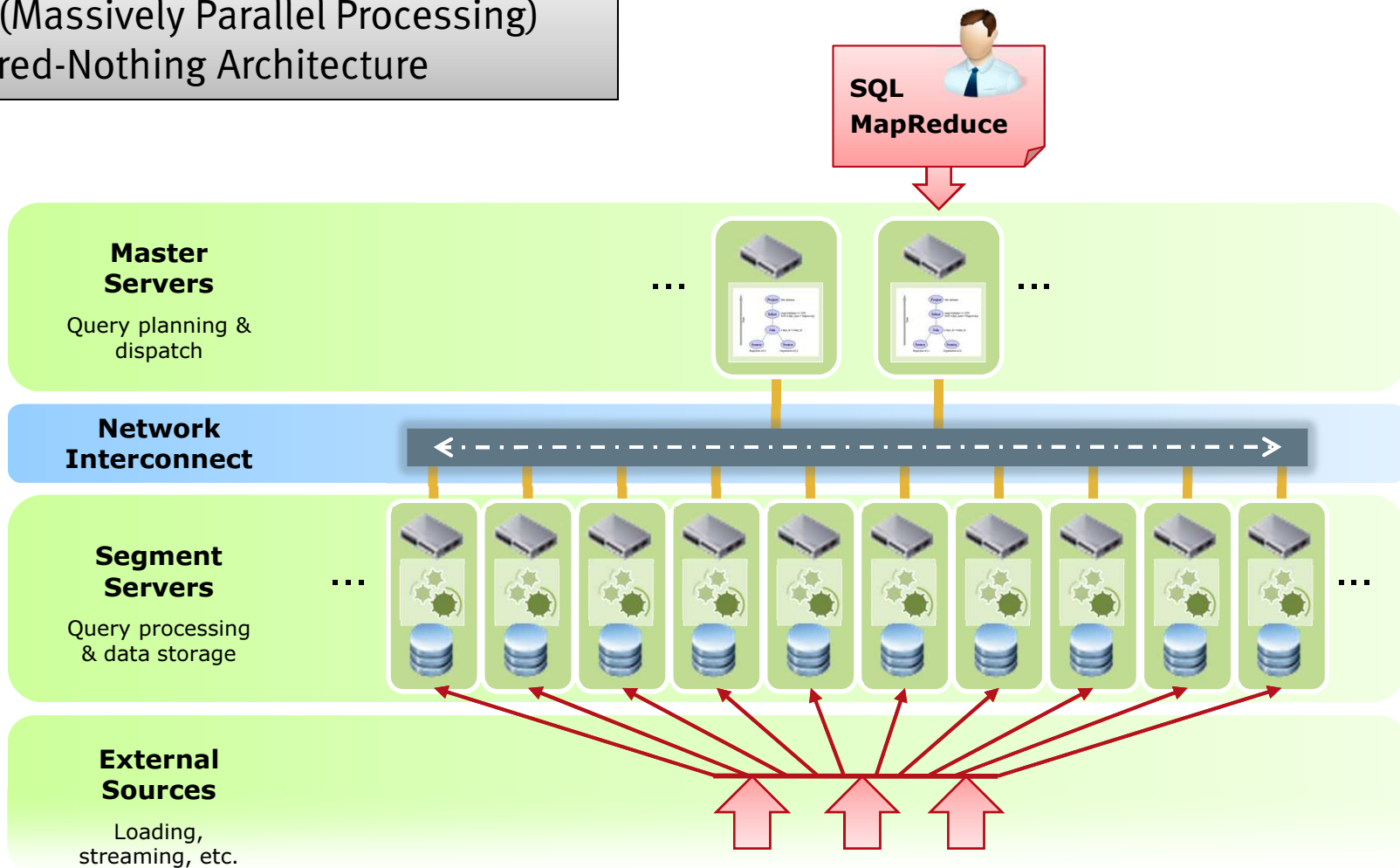
- Log likelihood
- Conjugate gradient
- Re-sampling
- K-means clustering
- Association rules
- ...

Our products

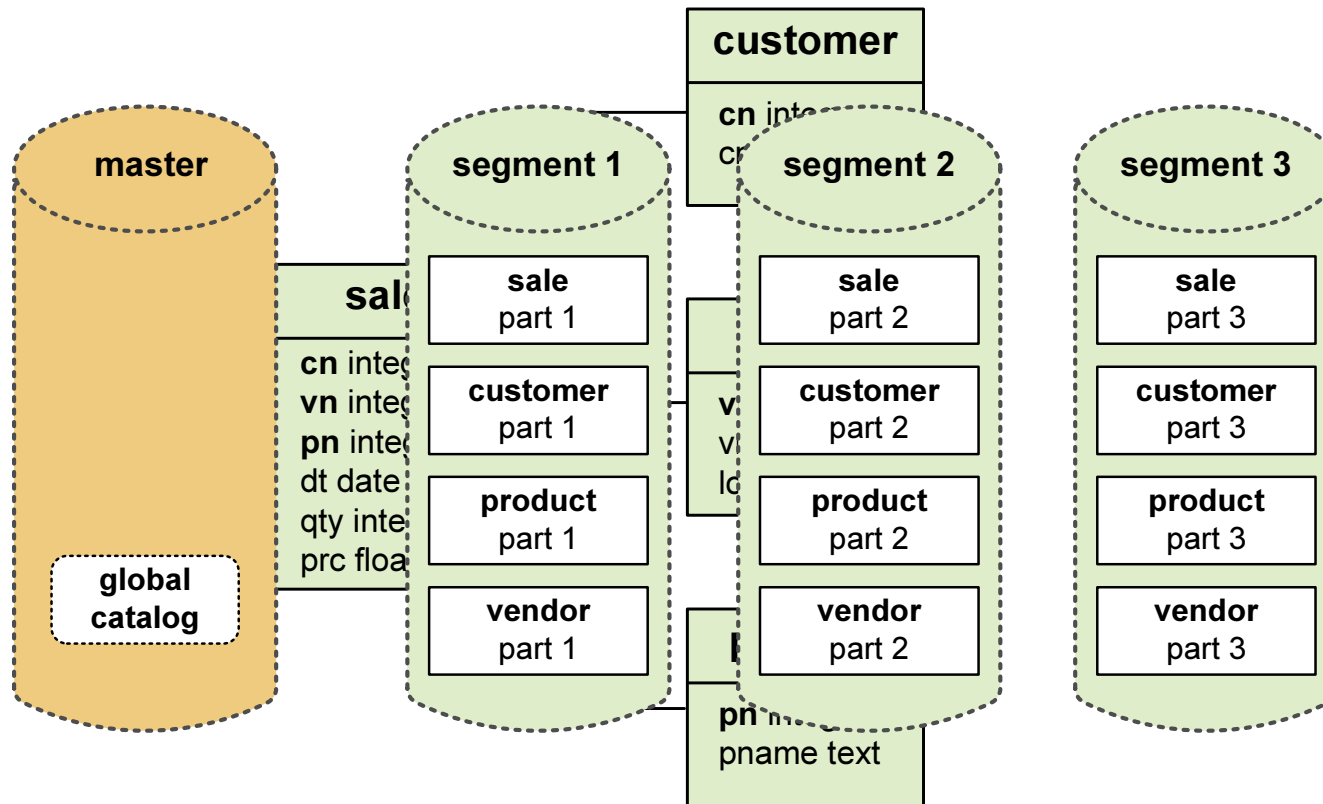
- **Greenplum Enterprise Data Cloud - Chorus**
 - Breakthrough data virtualization and collaboration software platform
- **Greenplum MPP Database**
 - Industry's leading MPP database software for large-scale data warehousing and analytics
- **Greenplum Single Node Edition**
 - Free analytic database software, highly optimized for single server deployments
- **Greenplum DCA – Data Computing Appliance**
 - DCA10, DCA100, DCA1000, DCA10000

Greenplum Database Architecture

MPP (Massively Parallel Processing)
Shared-Nothing Architecture



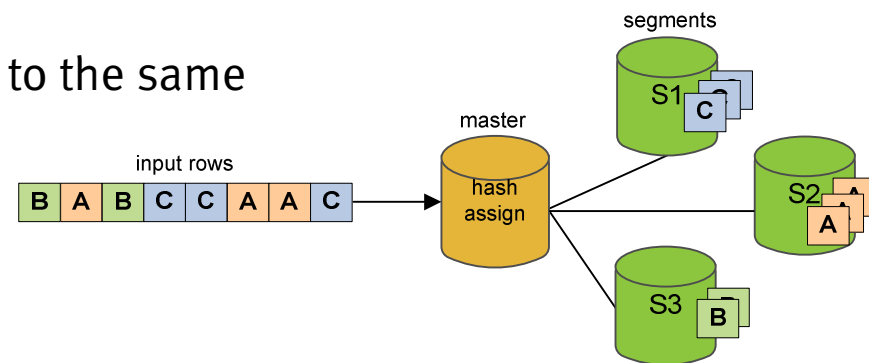
Distributed Tables



Distribution Policies

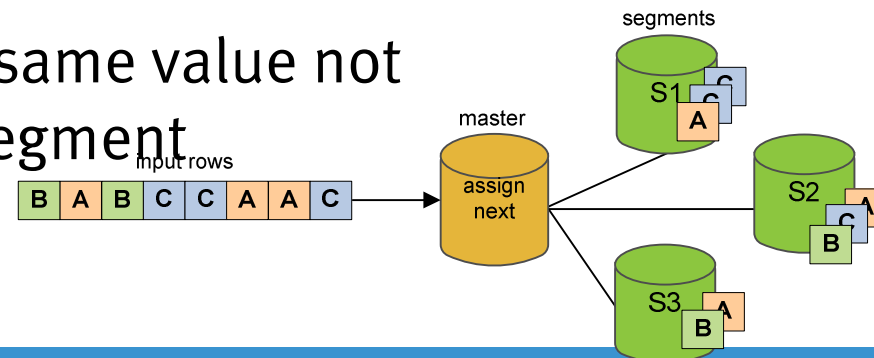
• Hash Distribution

- CREATE TABLE ... DISTRIBUTED BY (column [, ...])
- Keys of the same value always sent to the same segments



• Round-Robin Distribution

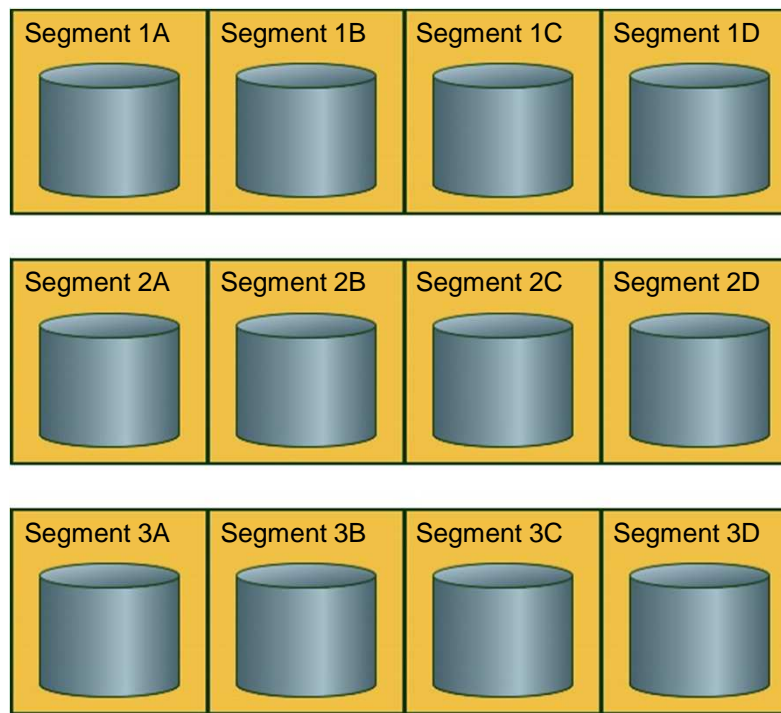
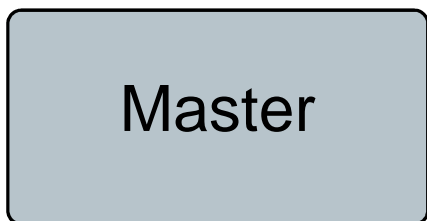
- CREATE TABLE ... DISTRIBUTED RANDOMLY
- Rows with columns of the same value not necessarily on the same segment



Parallel Data Scans

```
SELECT COUNT(*)  
FROM orders  
WHERE order_date >= 'Oct 20 2005'  
AND order_date < 'Oct 27 2005'
```

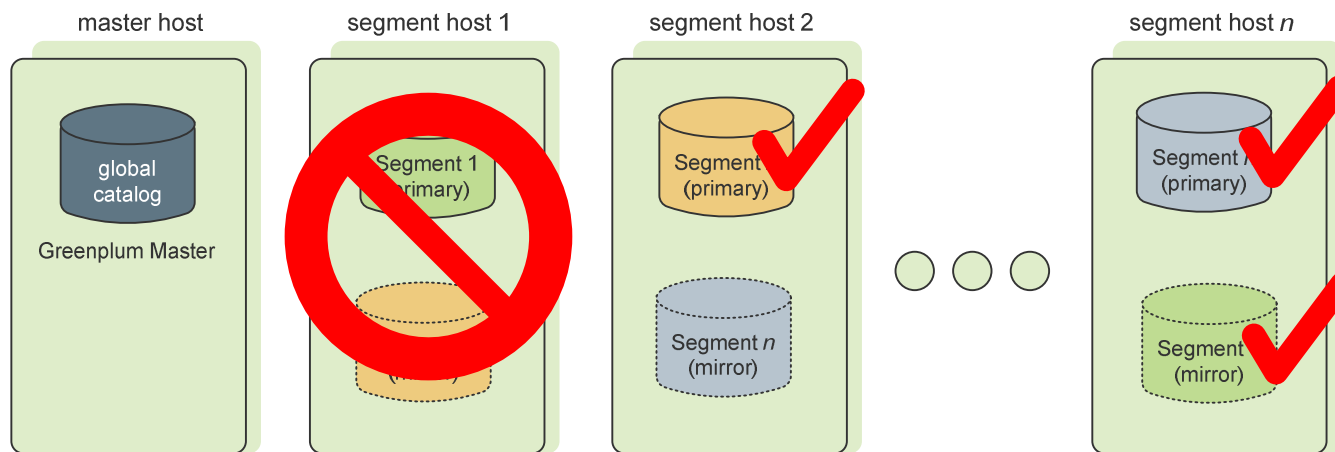
4,423,323



- Develop Query Plan
- Send Plan to Segments
- Segments Return Results
- Return Results

Each Segment Scans Data Simultaneously

Data Redundancy - Segment Mirroring



Configuration Specifications



Expansion bus

Eight segment servers
full rack

Interconnect bus

Two master servers

Eight segment servers

Specifications	GP100: Half Rack	GP1000: Full Rack
Master servers	2	2
Segment servers	8	16
Memory per server	48 GB	48 GB
Total memory	384 GB	768 GB
Segment HDDs (SAS)	96	192
Usable capacity (uncompressed)	18 TB	36 TB
Usable capacity (compressed)	72 TB	144 TB
Scan rate	12 GB/s	24 GB/s
Data load rate	5 TB/hour	10 TB/hour

THE FASTEST

Integrated Data Domain Backup

EMC Greenplum
Data Computing
Appliance



EMC
Data Domain
DD880



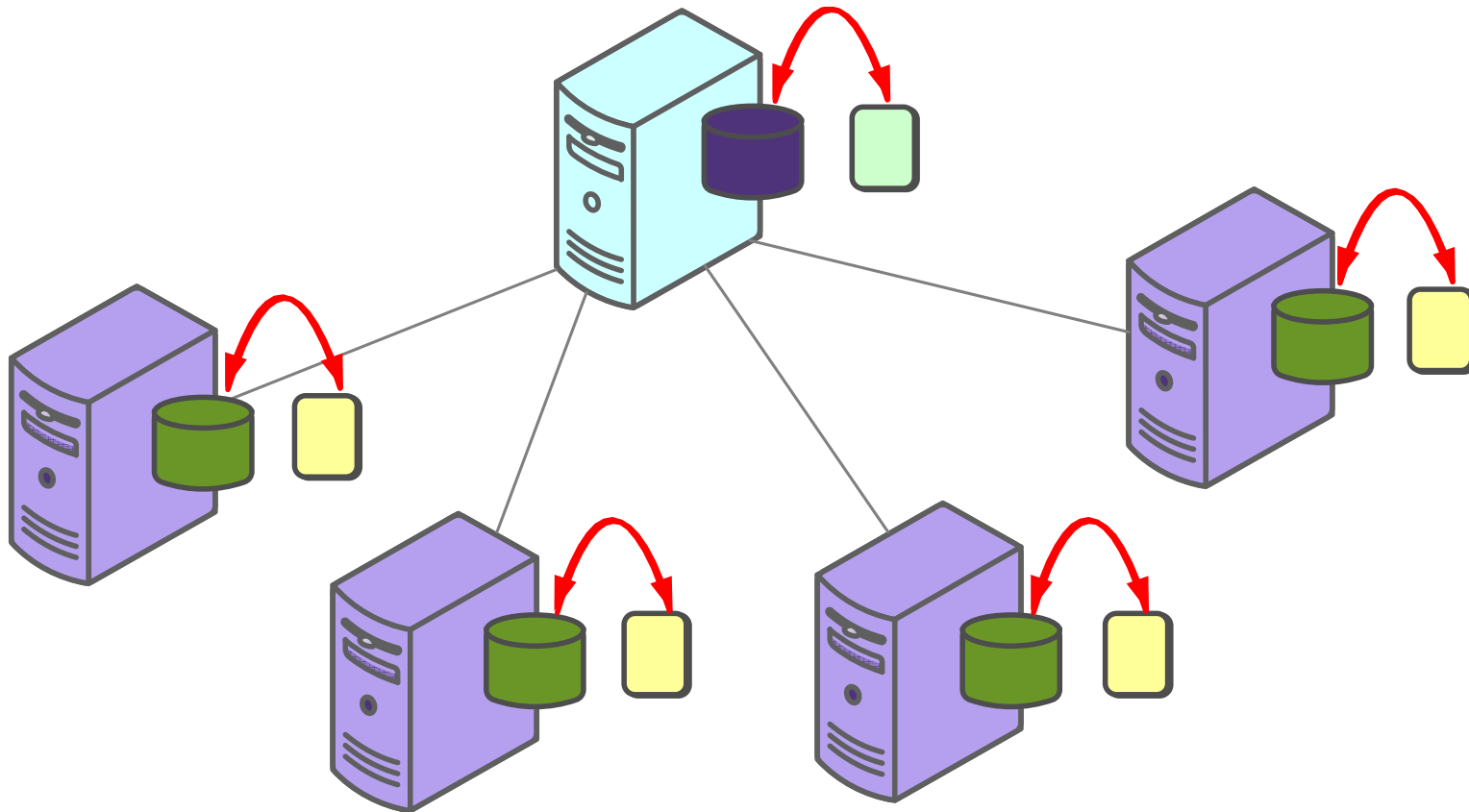
Segment Servers
NFS Shares →

Twin-Ax / FC
Cables

2 x 10GBit
IP Links

- Backup and Recovery
 - With EMC Data Domain / Greenplum native utility
- Reduces storage backup requirements
 - Deduplicates data
- Fast, reliable data recovery
 - Reduced recovery time
- Flexible and efficient
 - Designate intervals to backup
 - Point-in-time copies

Parallel *Everything*, including Backups and Restores



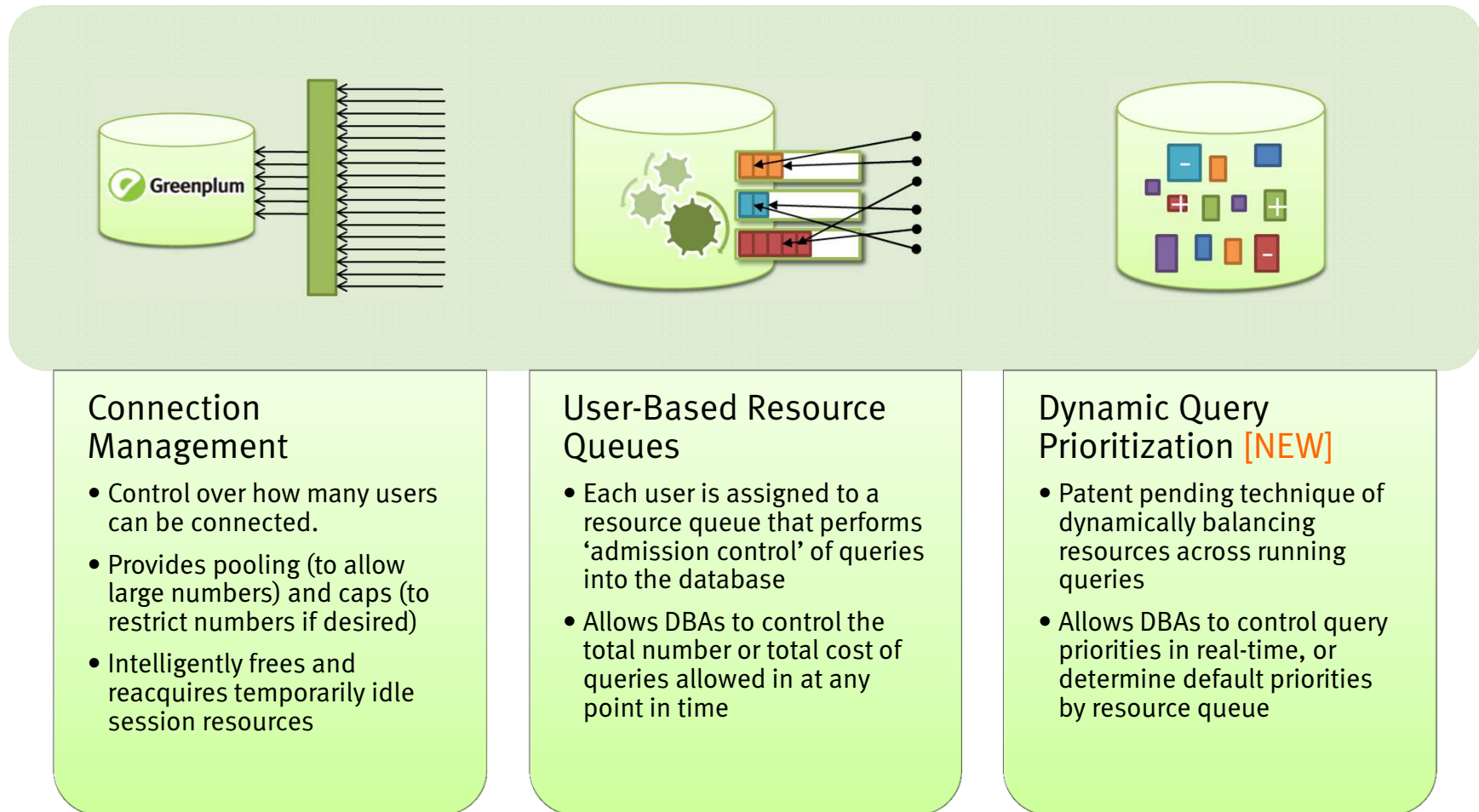
Greenplum Polymorphic Storage™

Flexible Row or Column Oriented Processing



- **Rather than take a side, we give customers the flexibility of both**
 - Results consistent with industry/academic findings of row vs column benefits
- **Row-orientation typically better for general purpose DW**
 - Avoid reassembly overhead that dominates in typical workloads
- **Column-orientation typically better for an important set of use cases**
 - Accessing small # of cols from a wide table – e.g. certain Data Mining use cases
- **Table Orientation**
 - Just specify 'orientation=row' or 'orientation=column' when creating a table
 - Gzip and LZ compression algorithms available with either orientation
- **Gives customers the choice of processing model for any table**
- **Efficient pre-projection, parallel execution in either case**

Workload Management Overview



Connection Management

- Control over how many users can be connected.
- Provides pooling (to allow large numbers) and caps (to restrict numbers if desired)
- Intelligently frees and reacquires temporarily idle session resources

User-Based Resource Queues

- Each user is assigned to a resource queue that performs 'admission control' of queries into the database
- Allows DBAs to control the total number or total cost of queries allowed in at any point in time

Dynamic Query Prioritization [NEW]

- Patent pending technique of dynamically balancing resources across running queries
- Allows DBAs to control query priorities in real-time, or determine default priorities by resource queue

Enhanced Database Monitoring - Dashboard



Enhanced Database Monitoring – System Metrics



Enhanced Database Monitoring – Query Monitor

The screenshot displays the Greenplum Query Monitor interface. At the top, it shows the Greenplum logo, the GPDB version, and user information (Welcome gpadmin, Logout, About). Below this is a navigation bar with tabs for DASHBOARD, SYSTEM METRICS, and QUERY MONITOR. The QUERY MONITOR tab is active, showing a search filter for 'Query search' and search results for 10 finished queries and 0 aborted queries.

A table lists several queries with columns: Watch, Query ID, Submit Time, Username, Database, Status, Wait time, Run time, Rows out, CPU %, CPU total, CPU velocity, CPU skew, Row skew, and Details. The query with ID 1229458791-13482-58 is highlighted.

The 'Query Details' panel for query 1229458791-13482-58 is expanded, showing a query plan diagram and a table of statistics for each node type. The query text is also visible in a separate pane.

Query Text (1229458791-13482-58)

```

select
  s_name,
  s_address
from
  supplier,
  nation
where
  s_suppkey in(
  select
    ps_suppkey
  from
    partsupp,
    (
      select
        sum(i)
      from
        lineitem
      where
        i_ship
        and i_
        group by i_
      where
        g_i_partkey = ps_
        and...
    )
  )
  
```

Query Plan Statistics:

Node type	HashAggregate
Last updated	2008-12-17 15:06:40
Node status	Finished
Start time	2008-12-17 15:06:34
Duration	0.00 (avg)
PMem Size	2,070,970 (avg)
PMem Max	16,595,218 (avg)
Mem Size	302,658,901 (avg)
Mem Resident	63,363,754 (avg)
Shared mem	48,345,088 (avg)
CPU total	413.38 (avg)
CPU %	0.19 (avg)
Phase	
Rows out	0 (sum)
Rows estimate	12 (sum)
CPU skew	86.61
Row skew	0.00

Additional attributes:

- Rows in: 0 Rows (avg) [estimated 0]
- Aggregate Total Spill Tuples: 0 Tuples (avg) [estimated 0]
- Aggregate Total Spill Bytes: 0 Bytes (avg) [estimated 0]
- Aggregate Total Spill Batches: 0 Batches (avg) [estimated 0]
- Aggregate Total Spill Pass: 0 Passes (avg) [estimated 0]
- Aggregate Current Spill Pass Read Tuples: 0 Tuples (avg) [estimated 0]
- Aggregate Current Spill Pass Read Bytes: 0 Bytes (avg) [estimated 0]
- Aggregate Current Spill Pass Tuples: 0 Tuples (avg) [estimated 0]
- Aggregate Current Spill Pass Bytes: 0 Bytes (avg) [estimated 0]
- Aggregate Current Spill Pass Batches: 0 Batches (avg) [estimated 0]

where information lives

EMC Consulting Services for DW/BI

BENCHMARKING AND PROOF OF CONCEPT

- Client use cases implemented on appliance and exercised at scale
- Standard performance metrics captured for current state and competitive comparison
- TCO and ROI calculations for decision makers
- Platform training for client IT staff



PLATFORM MIGRATION AND INTEGRATION

- Database Replatforming and Consolidation
- Application code porting, including Oracle, Teradata, DB2 and SQL Server
- High speed historic data migration, consolidation and harmonization
- Real/near real time and batch system integration via ETL and ELT



PERFORMANCE DESIGN AND DATA MODELING

- Optimal Greenplum configuration, implementation and data model design aligned with business use cases
- MPP Distribution
- Workload Management
- Polymorphic Storage including:
 - Partitioning
 - Compression
 - Columnar

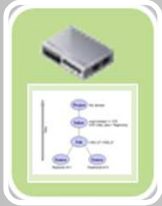


INFORMATION MANAGEMENT AND ANALYTICS

- Near real time predictive analytics across multi-Terabyte data sets
- Business Intelligence and Reporting
- Self provisioning analytic marts and business user collaboration
- Data Governance, Stewardship and Continuous Quality Improvement



Key Technical Innovations



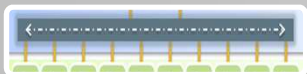
Scatter-Gather Data Streaming

- Industry leading data loading capabilities



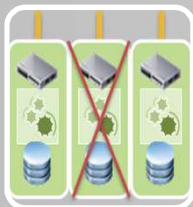
Online Expansion

- Dynamically provision new servers with no downtime



Map-Reduce Support

- Parallel programming on data for advanced analytics



Polymorphic Storage

- Support for both row and column-oriented storage

Greenplum Data Computing Appliance

Data In. Decisions Out.

- Industry leading price-performance
- Most powerful data loading machine on the planet (+10TB per hour)
- First best step to private-cloud, virtualized DW and analytic infrastructure



EMC² where information lives[®]