

Maximizing Hadoop Performance with Hardware Compression

Robert Reiner

Director of Marketing

Data Compression and Security

Exar Corporation

What is Big Data?

- “Data sets whose size is beyond the ability of typical data base software tools to capture, store, and analyze”
 - McKinsey Global Institute
- “Data sets so large and complex that it becomes difficult to process using on-hand database management tools. ”
 - Wikipedia
- “Data that’s an order of magnitude bigger than what you’re accustomed to, Grasshopper”
 - eBay

Sources of Big Data

Type of Data	Industry
Web Logs	Social Networking, Online Transactions
RFID	Retail, Manufacturing, Casinos
Smart Grid	Utilities
Sensor	Industrial Equipment, Engines
Telemetry	Video Games
Telematics	Auto Insurance
Text Analysis	Multiple
Time and Location Analysis	Multiple

Big Data Growth

40% projected growth in global data generated/year¹

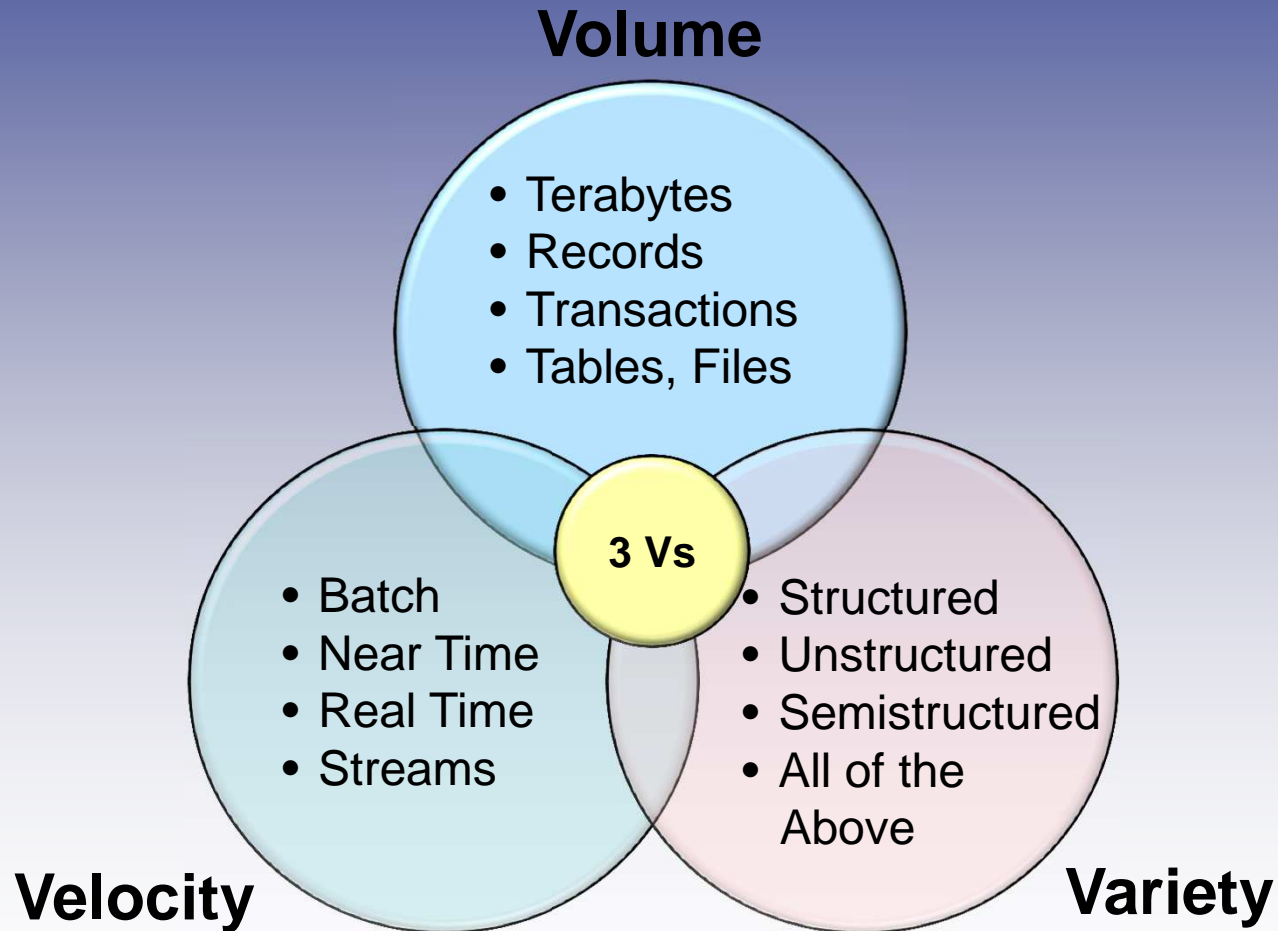
5% growth in global IT spending¹

Twitter generates >7TB of data per day²

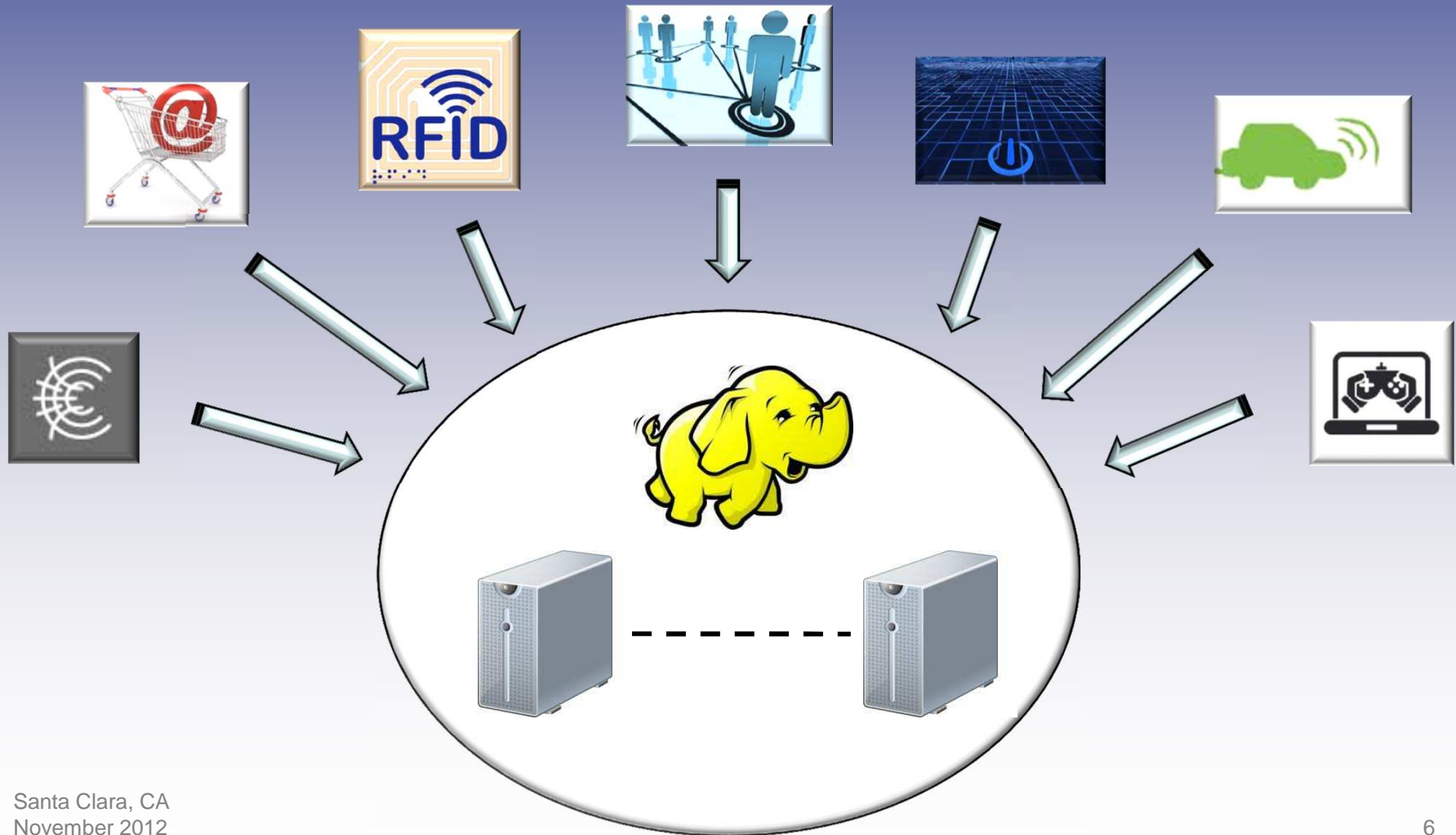
Facebook generates >10TB of data per day²

80% of world's information is unstructured²

Three Vs of Big Data



Hadoop Addresses Big Data Challenges

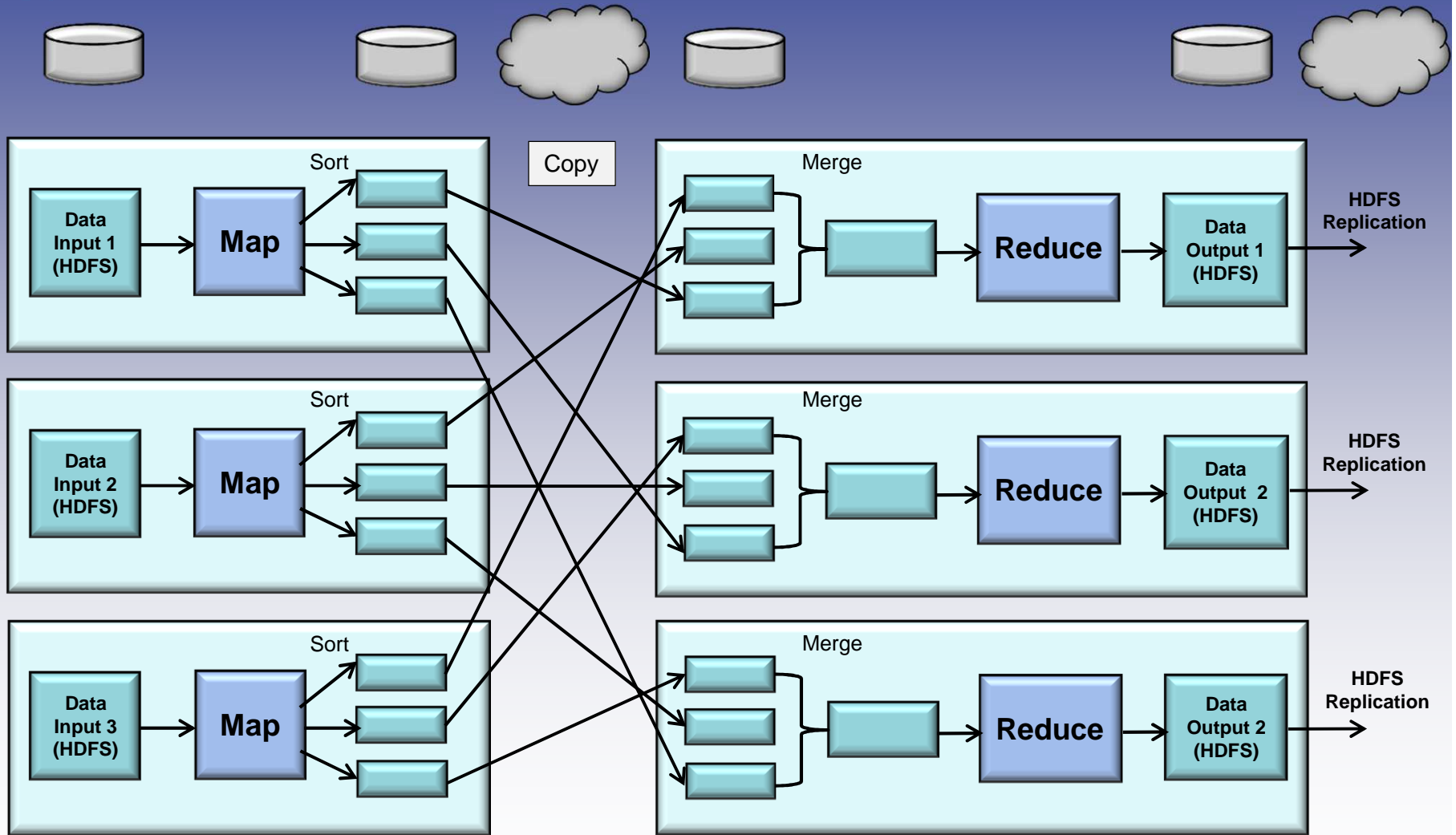


Map Reduce is Core of Hadoop

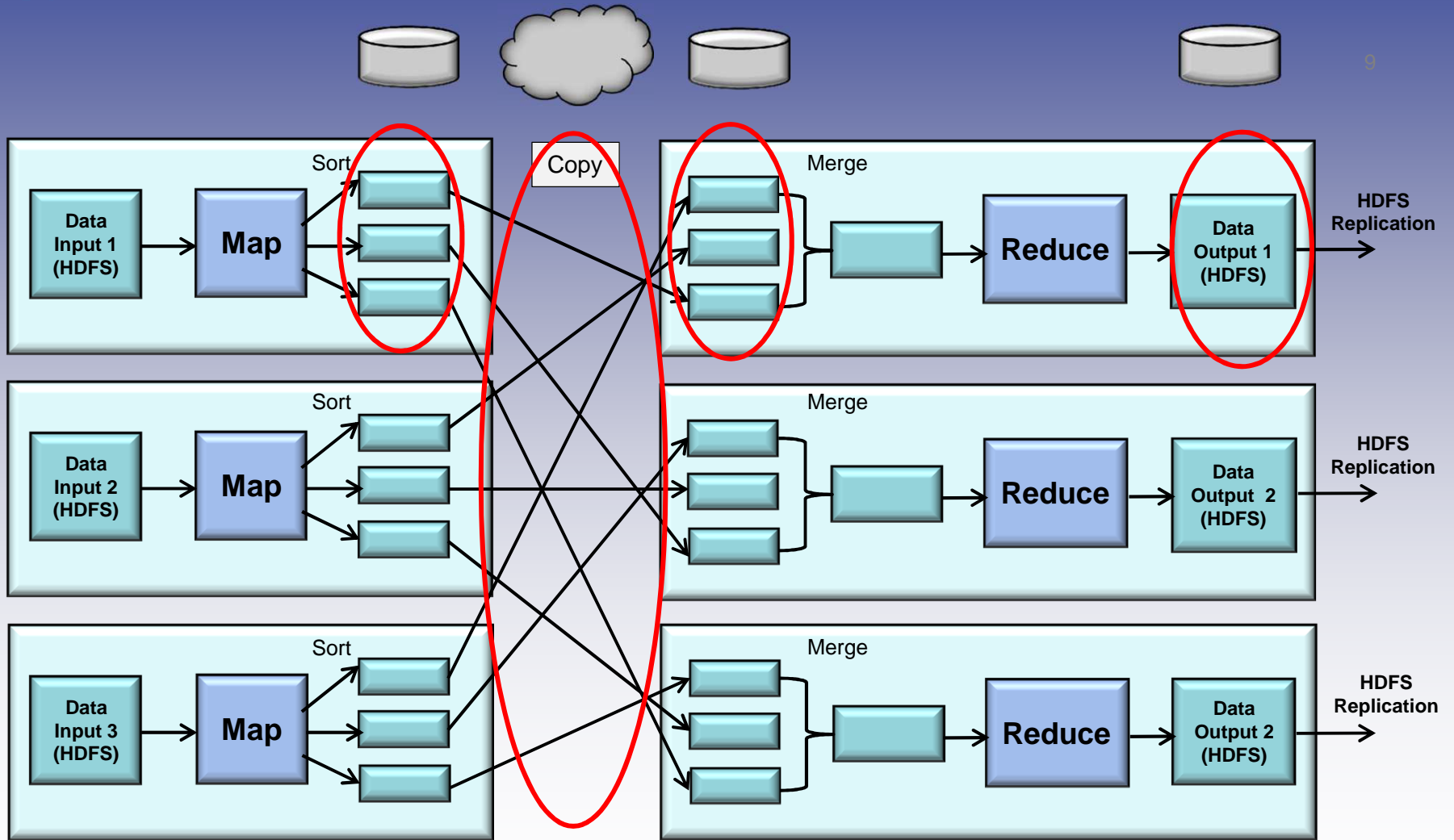
- Parallel Programming Framework
- Simplifies Data Processing Across Massive Data Sets
 - Enables Processing Data in a File System without being Stored into a Data Base
- Ability to Process Unstructured Data
- Excels at Sifting through Huge Masses of Data to Find what is Useful



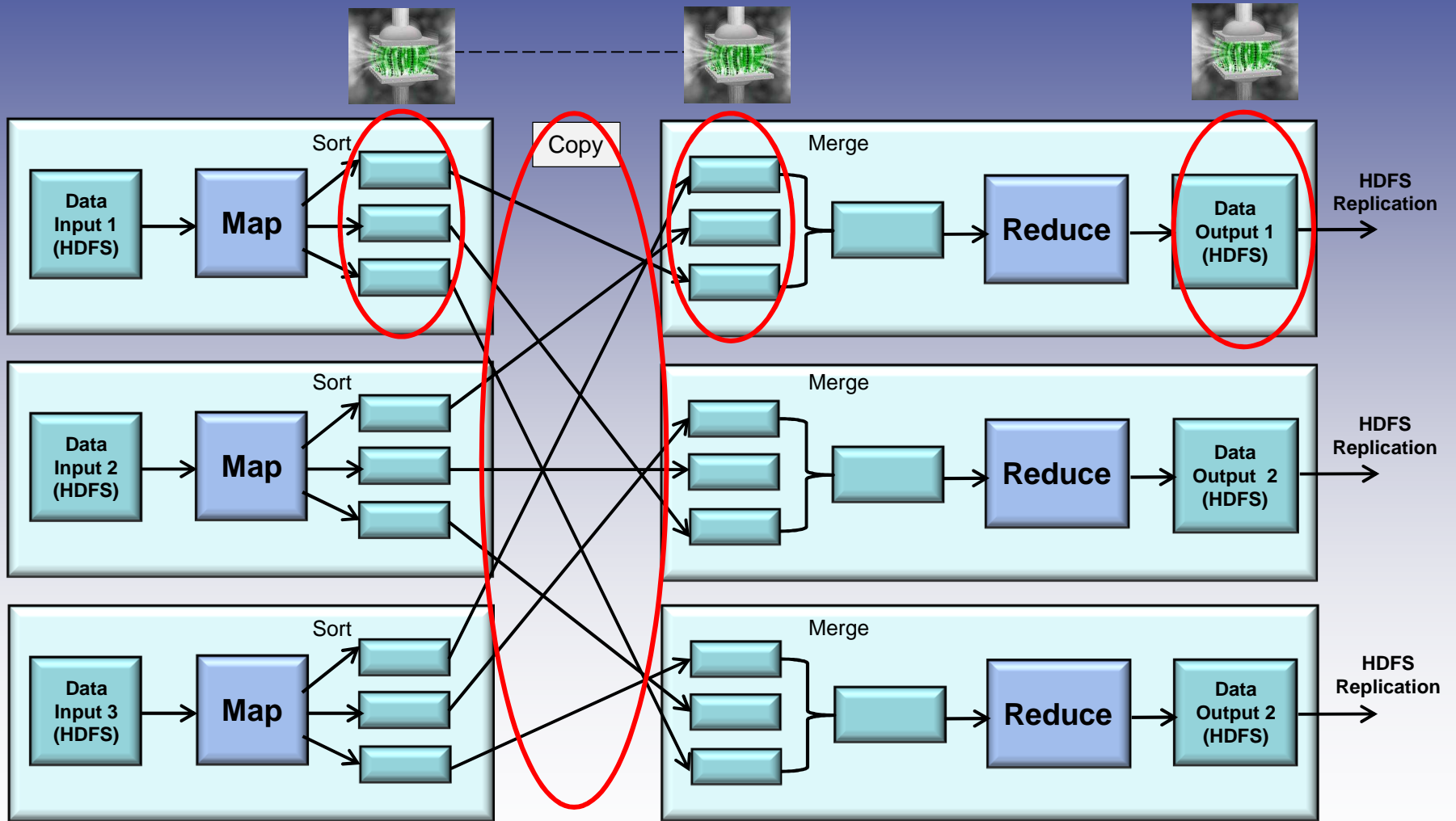
MapReduce Data Flow



MapReduce I/O Bottlenecks



Compression Reduces Networking and Disk I/O – Addresses Bottlenecks



Hadoop Software Compression Codecs

Codec	Compression Performance	Compression Ratio	CPU Overhead
Deflate/ Gzip	Low	High	High
Bzip2	Very Low	Very High	Very High
LZO	Medium	Medium	Medium

Hardware Accelerated Compression

- Processor Intensive Compression Algorithms Executed in Hardware
 - Increases Performance
 - Reduces MapReduce Execution Time
 - Offloads CPU
 - Lowers Energy Consumption



Hardware Accelerated Compression - Exar Solutions

- DX1800 and DX1700 Series PCIe Cards
 - Java codec calls C libraries and invokes card SDK



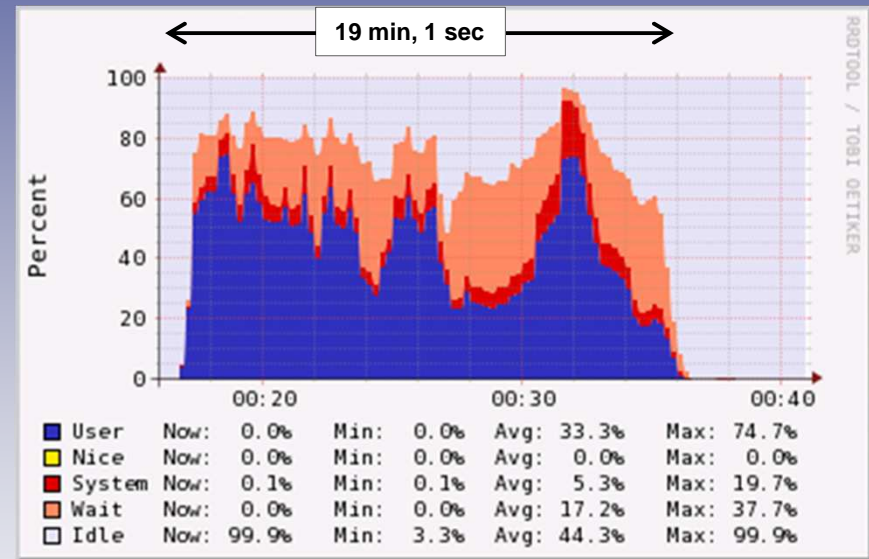
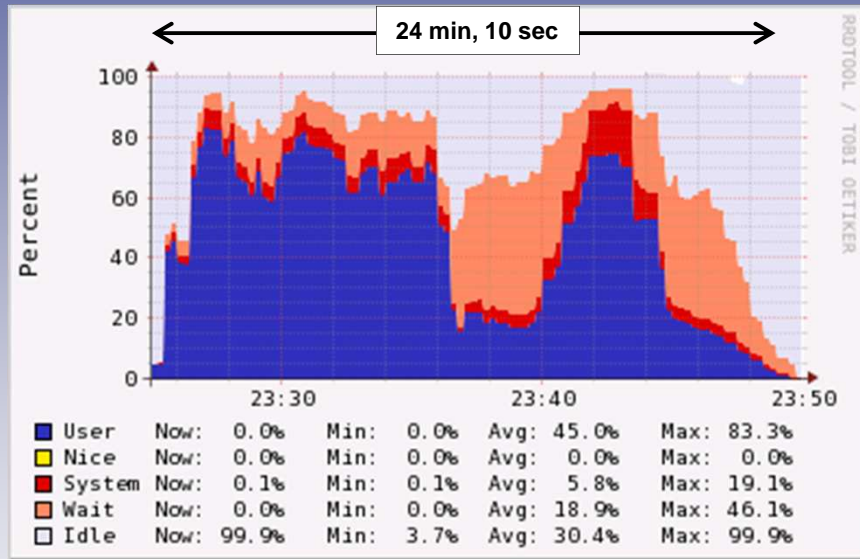
Hardware and Software Compression for Hadoop

Codec	Compression Performance	Compression Ratio	CPU Overhead
HW eLZS	Very High	Medium	Low
HW Deflate/ Gzip	Very High	High	Low
Deflate/ Gzip	Low	High	High
Bzip2	Very Low	Very High	Very High
LZO	Medium	Medium	Medium

Hadoop Codec Benchmarking

- Benchmarked Multiple Codecs using Terasort
- Terasort Input Size: 100GB
- Three Node Hadoop Cluster
 - 1GbE Switch Interconnect
- Hadoop version 1.0.0
- Node Configuration
 - Dual E5620/ node (8 cores, 16 threads)
 - 16 GB DRAM
 - 4 x SATA-300 (108 MB/sec, 500 GB)
 - RHEL 5.4

Benchmarking Results – SW LZO and HW eLZS Codecs



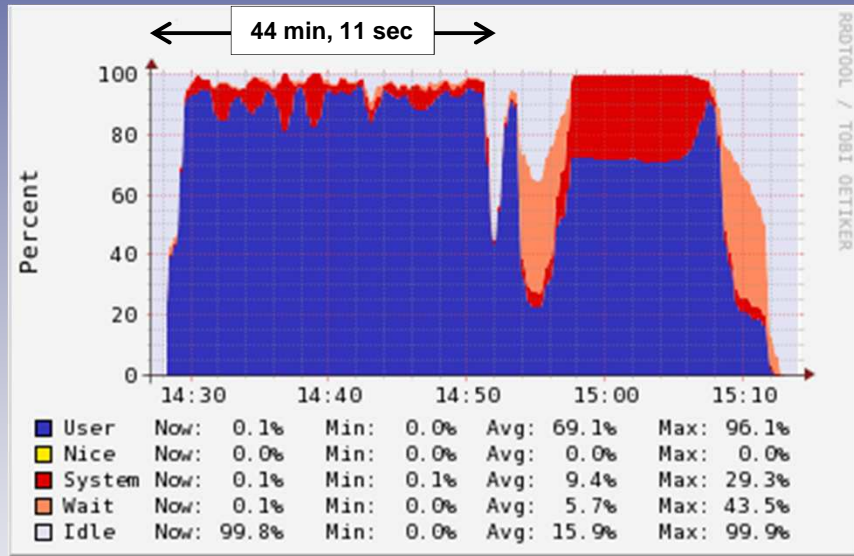
SW LZO

- Total Time: 24 min, 8 sec
- Compression Ratio: 5.126
- .2478 kWh

HW eLZS**

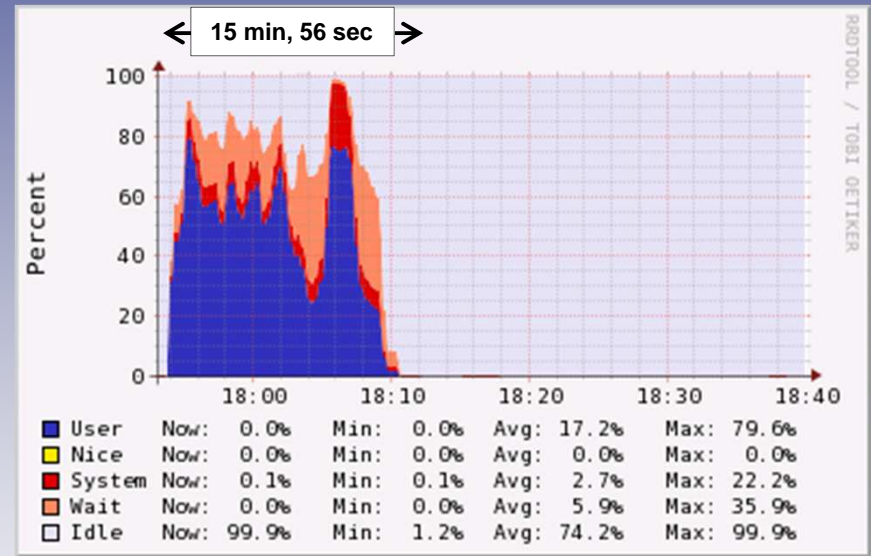
- Total Time: 19 min, 8 sec
- Compression Ratio: 5.136
- .2061 kWh

Benchmarking Results – SW Gzip and HW Gzip Codecs



SW Gzip

- Total Time: 44 min, 11 sec
- Compression Ratio: 7.645



HW Gzip**

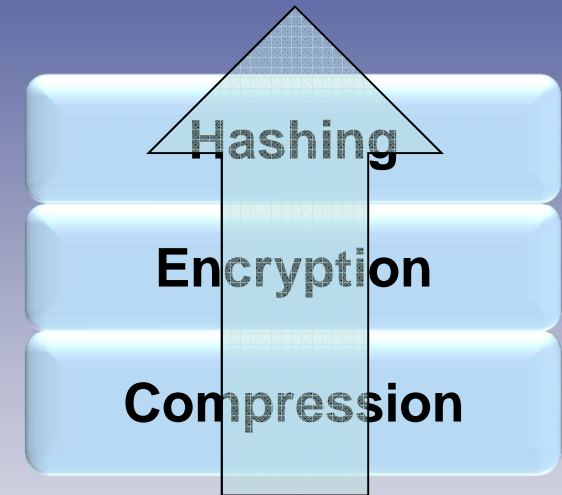
- Total Time: 15 min, 56 sec
- Compression Ratio: 6.329

Terasort Benchmarking Summary

- HW eLZS is 21% faster than SW LZO
- HW Gzip is 64% faster than SW Gzip
- HW Gzip provides the fastest Terasort time of all codecs benchmarked
- HW Accelerated Compression is more energy efficient than SW compression

Future of Hardware Acceleration for Hadoop

- Security
 - Add encryption to Hardware-Accelerated Codec
 - Exar's technology enables single pass compression and encryption
- Enhanced Benchmarks
 - Expand beyond Terasort to include benchmarks that represent additional production workloads



Thank You

Robert Reiner

Director of Marketing

Data Compression and Security

Exar Corporation

rob.reiner@exar.com