# HPCC Systems

## Pointers on how to calculate the *real* ROI on a Big Data Analytics System

# Original claim from SGI using Hadoop
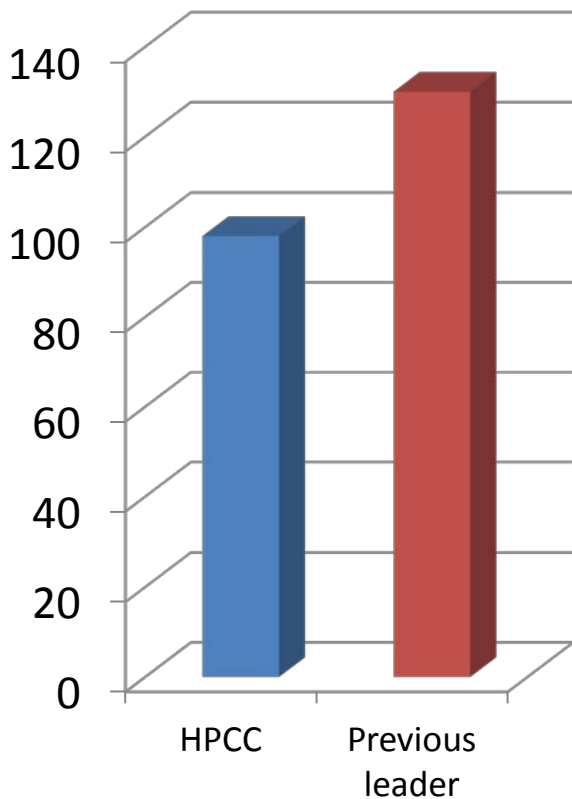
Link to the SGI Press Release:

http://www.sgi.com/company_info/newsroom/press_releases/2011/october/hadoop.html

SGI Press Release summary (October 2011)

- 20 nodes running on SGI® Rackable™ C2005-TY6 half-depth servers
- Each node:
  - 1-2 Intel Xeon E5630 CPU (unspecified)
  - 48GB RAM
  - 1-2 network uplinks (unspecified)
  - (4) 1TB SATA HDD's (configuration unspecified)
- Operating System (unspecified)
- Cloudera distribution of Apache Hadoop (CDH3u0)
- Hadoop configuration (unspecified)

# Results in pictures (less is better)

## Execution Time
### (seconds)



| | |
|---|---|
| 140 | |
| 120 | |
| 100 | |
| 80 | |
| 60 | |
| 40 | |
| 20 | |
| 0 | |
| HPCC | Previous leader |

## Productivity

3 ECL statements

```
// Perform global terasort
rec := record
        string10  key;
        string10  seq;
        string80  fill;
        end;
in := DATASET('nhtest::terasort1',rec,FLAT);
OUTPUT(SORT(in,key,UNSTABLE),,'nhtest::terasort1out',overwrite);
//End
```

700+ Lines of Java MapReduce Code

```
}
abstract int findPartition(Text key);
abstract void print(PrintStream strm) throws IOException;
int getLevel() {
  return level;
}
}

/**
* An inner trie node that contains 256 children based on the next
* character.
*/
static class InnerTrieNode extends TrieNode {
  private TrieNode[] child = new TrieNode[256];

  InnerTrieNode(int level) {
    super(level);
  }
  int findPartition(Text key) {
    int level = getLevel();
    if (key.getLength() <= level) {
      return child[0].findPartition(key);
    }
    return child[key.getBytes()[level] & 0xff].findPartition(key);
```
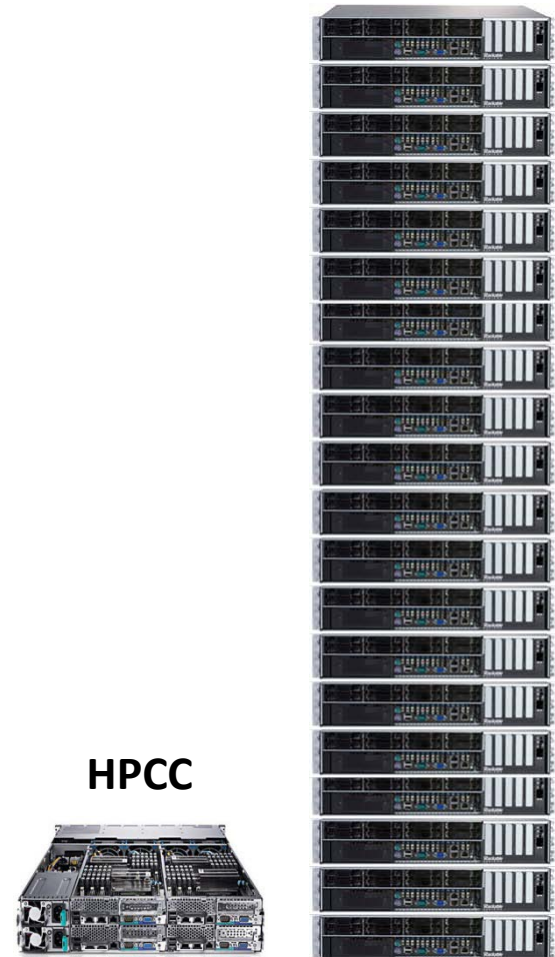
## Space/Cost



**Previous leader**

**HPCC**

# What is really important for the ROI?

**Do you want to:**

- Save on Datacenter Space, Power & Cooling
- Save on Operations personnel
- Save on Server Hardware and ongoing HW maintenance
- Save on Software licenses and Support (less nodes)

-------------------------------------------------------------------------

- Reduce project development/support cost/time
- Run your Solutions much faster and with more capabilities
- Use a Enterprise ready, fully featured, consistent platform
- Get Support from someone with 10+ yrs of Big Data expertise?
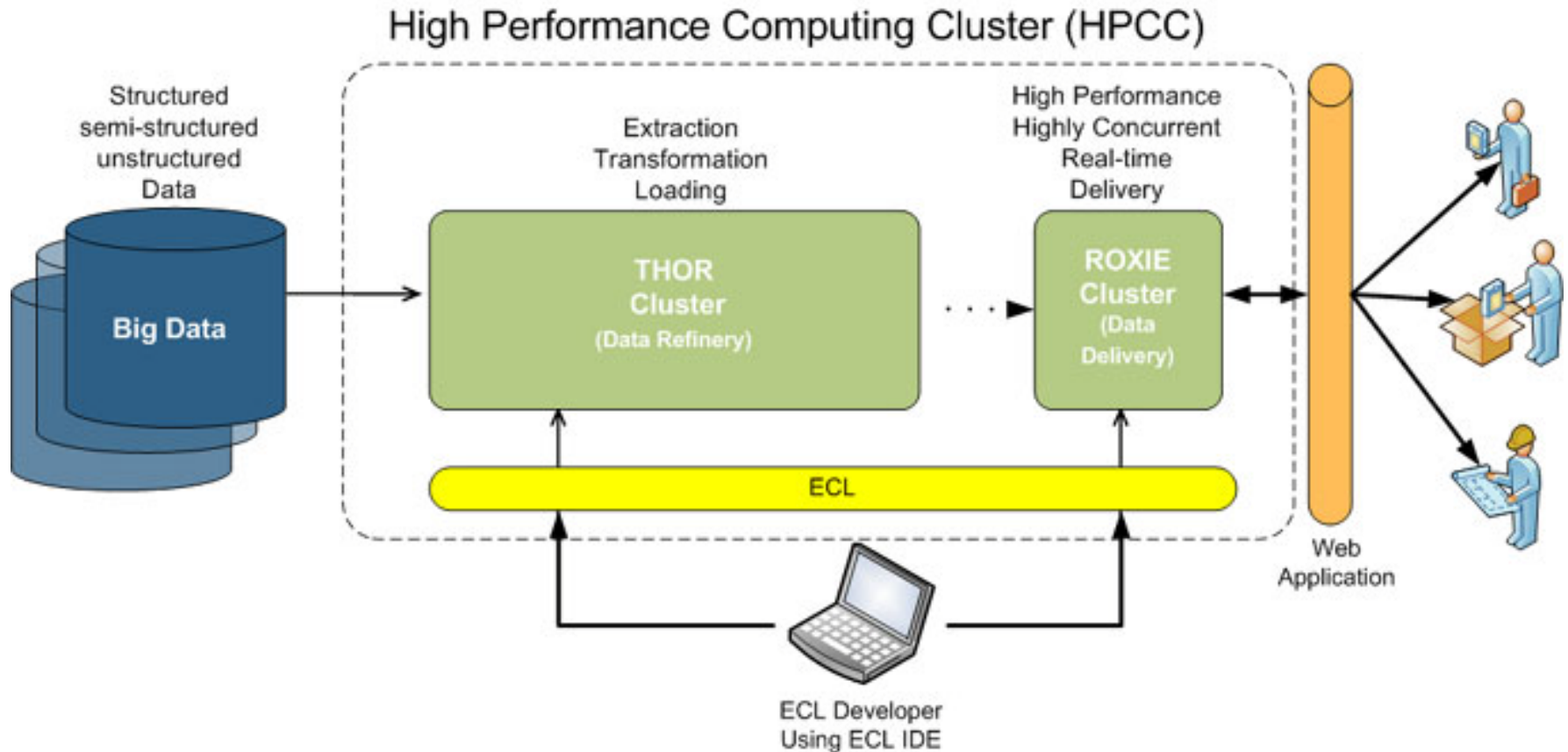
**………Call us, we can help!**

# Appendix

Risk Solutions

# The original claim debunked

- The methodology uses the Terasort specifications (http://sortbenchmark.org/)

- 100 GB are definitively not a Terabyte of data

- But someone threw the challenge, and we needed to respond (knowing how efficient our platform is)

- We are just responding to the SGI/Hadoop claim made on October 2011

High Performance Computing Cluster (HPCC)

# Terasort benchmark process and results

- Data generated following the exact data structure and distribution used in the Terasort benchmark (http://www.ordinal.com/gensort.html)

- 100GB total data size, across 1 billion records
  - 10 bytes as the key
  - 90 bytes as the value

- Flushed files system caches before execution

- Timed the total execution and repeated it 6 times

- Powered all systems down, waited for 15 minutes, powered all systems back up and repeated test

- Verified results

- **Average run time: 98 seconds (vs. 130 seconds previous leader: using Hadoop)**

# Our hardware and software

- 4 nodes running on **one (1)** Dell PowerEdge C6100 2U server
- Each node:
    - Intel Xeon E5675 CPU
    - 48GB RAM
    - (6) SAS Seagate Cheetah HDD's
- Linux CentOS 5.6
    - Deadline scheduler
    - Ext4 filesystem with noatime and nodiratime
- **HPCC Systems Thor**
- Thor configuration:
    - 24 worker threads per physical node
    - 1.5GB RAM allocated per Thor worker thread

# We welcome anyone to contact us to review what we did and how we did it

And yes, a **Sort** is not a complete indication of the performance of a Big Data Analytics solution, but it is an indicator. We know that our platform performs even better on "real jobs" with a more complex set of functions besides a simple sort

HPCC Systems - http://hpccsystems.com