**x y r a t e x ·**

**Advancing Digital Storage Innovation**

# Map/Reduce on Lustre
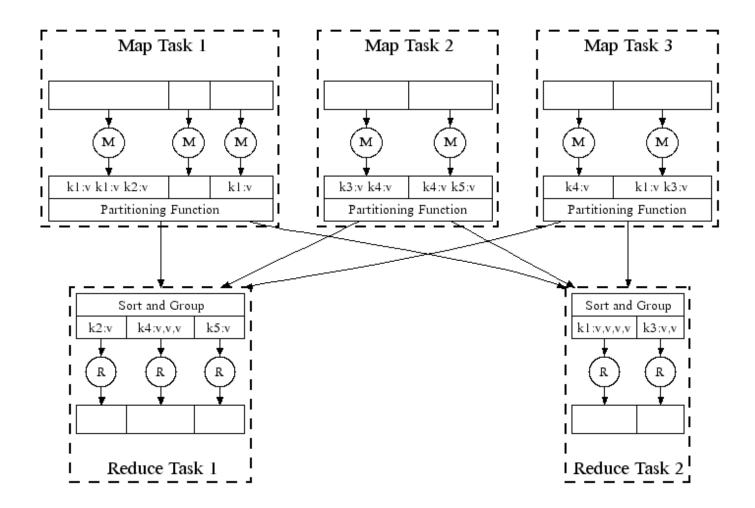## Hadoop Performance in HPC Environments

**Nathan Rutman, Xyratex**

**James B. Hofmann, Naval Research Laboratory**

# Agenda

- Map Reduce Overview
- The Case for Moving Data
- A Combined Lustre / HDFS Cluster
- Theoretical Comparisons
- Benchmark Study
- The Effects of Tuning
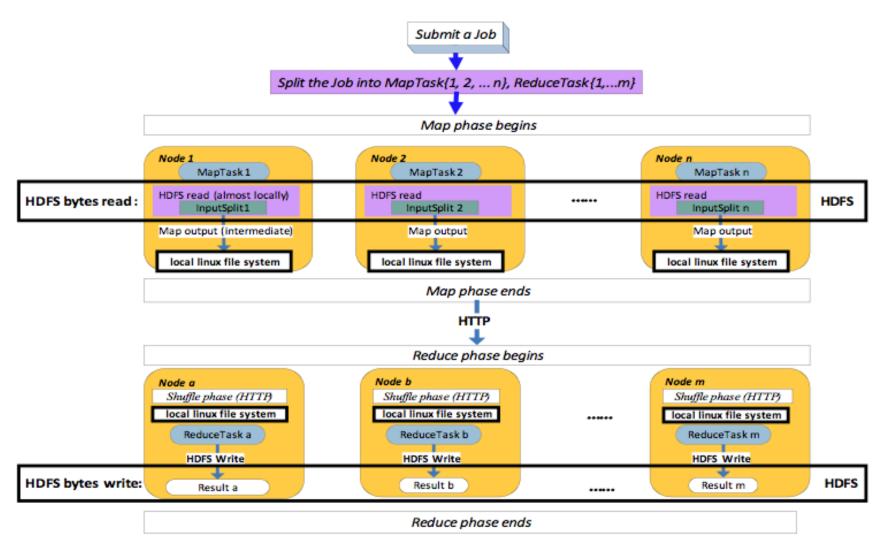- Cost Considerations

**xyratex**

# Map Reduce overview

# Apache Hadoop disk usage

Submit a Job

Split the Job into MapTask{1, 2, ... n}, ReduceTask{1,...m}

Map phase begins

| Node 1 | Node 2 | Node n |
|---|---|---|
| MapTask 1 | MapTask 2 | MapTask n |

**HDFS bytes read :**

| HDFS read (almost locally) | HDFS read | ...... | HDFS read | **HDFS** |
|---|---|---|---|---|
| InputSplit1 | InputSplit 2 | | InputSplit n | |

| Map output (intermediate) | Map output | Map output |
|---|---|---|
| **local linux file system** | **local linux file system** | **local linux file system** |

Map phase ends

**HTTP**

Reduce phase begins

| Node a | Node b | Node m |
|---|---|---|
| *Shuffle phase (HTTP)* | *Shuffle phase (HTTP)* | *Shuffle phase (HTTP)* |
| **local linux file system** | **local linux file system** | **local linux file system** |
| ReduceTask a | ReduceTask b | ReduceTask m |
| **HDFS Write** | **HDFS Write** | **HDFS Write** |

**HDFS bytes write:**

| Result a | ...... | Result b | ...... | Result m | **HDFS** |
|---|---|---|---|---|---|

Reduce phase ends

4

**x y r a t e x**

**Crossing the Chasm: Sneaking a Parallel File System Into Hadoop** , *Carnegie Mellon*

## Grep (64GB, 32 nodes, no replication)

# Other Studies: Hadoop with GPFS
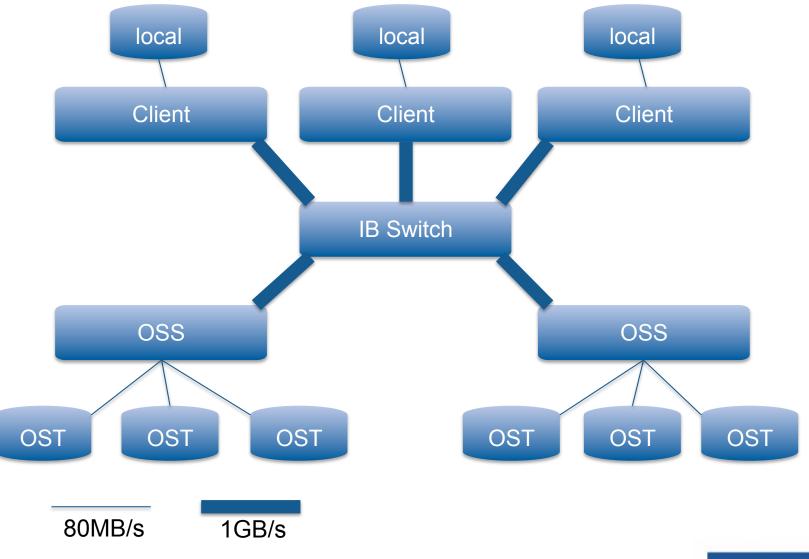
Execution time HDFS and GPFS with metablocks

# A Critical Oversight

- "Moving Computation is Cheaper Than Moving Data"
- The data ALWAYS has to be moved
  - Either from local disk
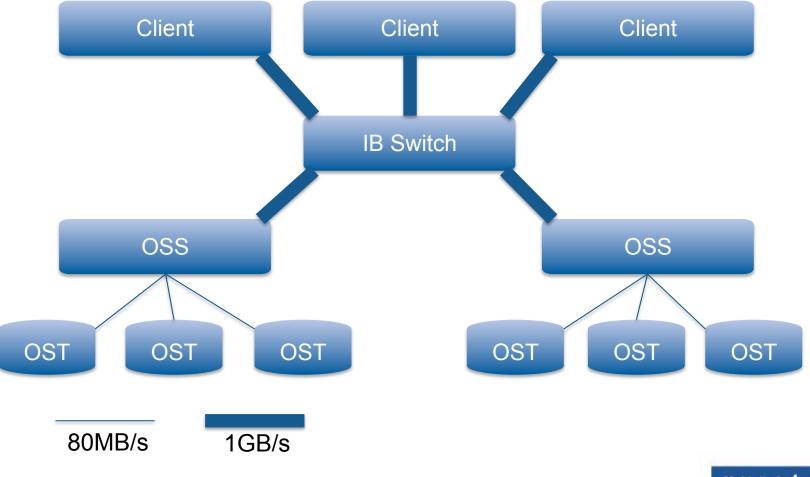  - Or from the network
- And with a good network: the network wins.

**xyratex**

# Cluster Setup: HDFS vs Lustre

- 100 clients, 100 disks, Infiniband
- Disks: 1 TB FATSAS drives (Seagate Barracuda)
  - 80 MB/sec bandwidth with cache off
- Network: 4xSDR Infiniband
  - 1GB/s
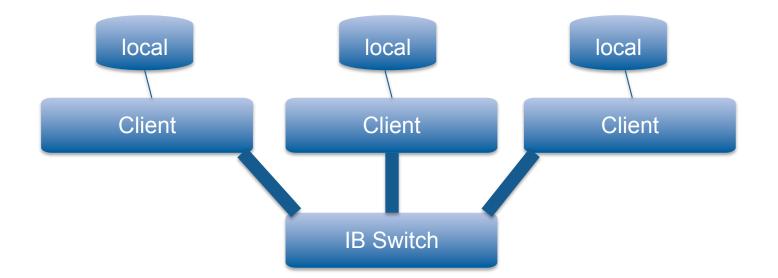- HDFS: 1 drive per client
- Lustre: 10 OSSs with 10 OSTs

# Cluster Setup



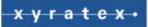| local | local | local |

Client   Client   Client

IB Switch

OSS   OSS

OST   OST   OST   OST   OST   OST

80MB/s   1GB/s

xyratex

# Lustre Setup

Client     Client     Client

IB Switch

OSS                              OSS

OST   OST   OST              OST   OST   OST
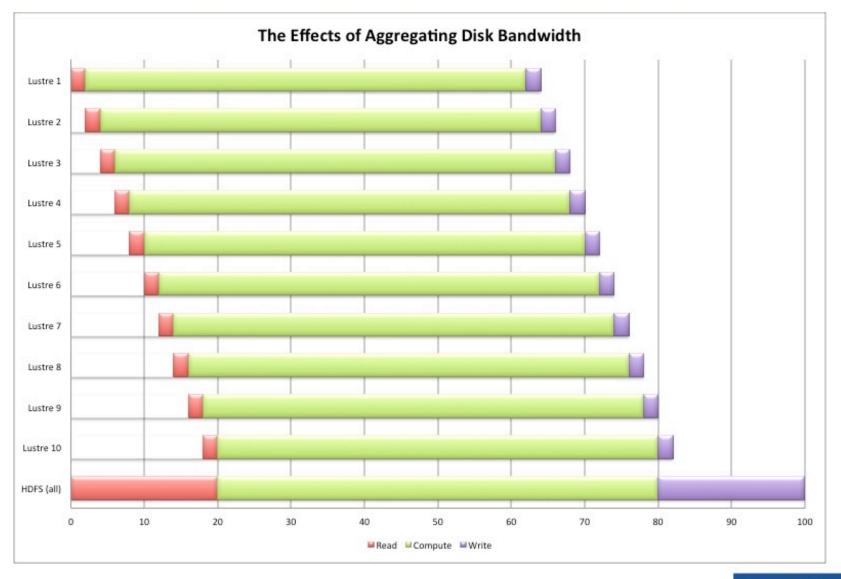
80MB/s            1GB/s

xyratex

# HDFS Setup

# Theoretical Comparison: HDFS vs Lustre

- 100 clients, 100 disks, Infiniband
- HDFS: 1 drive per client
    - Capacity 100 TB
    - Disk bandwidth 8 GB/s aggregate (80MB/s * 100)
- Lustre: Each OSS has
    - Disk bandwidth 800MB/s aggregate (80MB/s * 10)
        - Assuming bus bandwidth to access all drives simultaneously
    - Net bandwidth 1GB/s (IB is point to point)
- With 10 OSSs, we have same the capacity and bandwidth
- Network is not the limiting factor!

**xyratex**

# Striping

- In terms of raw bandwidth, network does not limit data access rate

- Striping the data for each Hadoop data block, we can focus our bandwidth on delivering a single block

- HDFS limit, for any 1 node: 80MB/s

- Lustre limit, for any 1 node: 800MB/s
    - Assuming striping across 10 OSTs
    - Can deliver that to 10 nodes simultaneously

- Typical MR workload is not simultaneous access (after initial job kickoff)
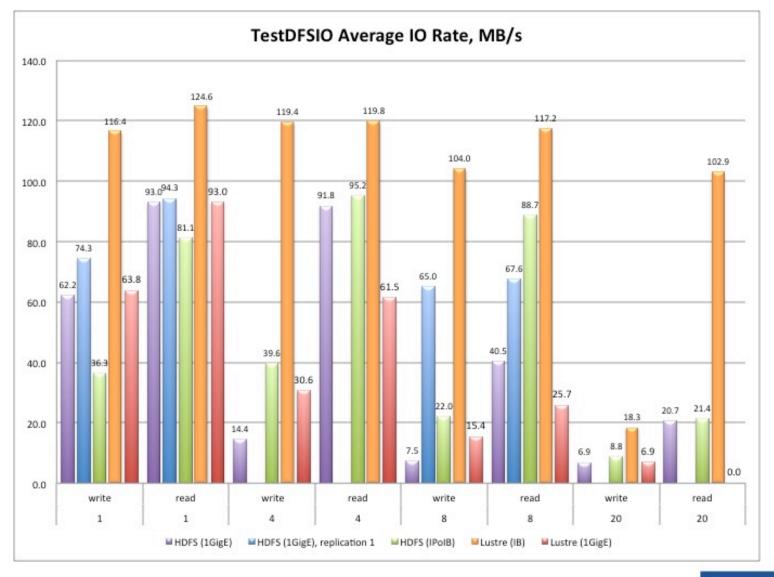
**xyratex**

The Effects of Aggregating Disk Bandwidth

# Replication

- HDFS replicates data 3x by default
- Recently Facebook added HDFS-RAID, which effectively trades off some computation (parity) for capacity
  - Can e.g. bring 3x safety for 2.2x storage cost when used
- Replicas should be done "far away"
- Replicas are synchronous
- HDFS writes are VERY expensive
  - 2 network hops, "far"
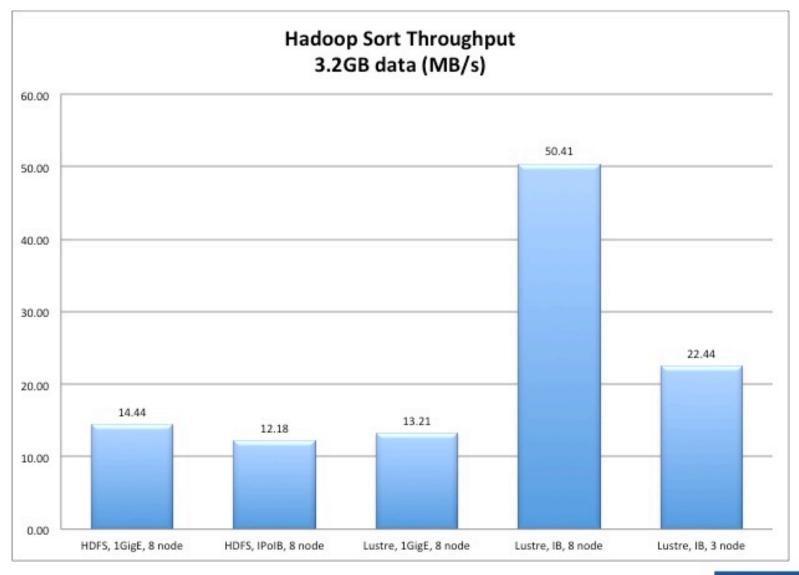  - 3x storage
- Can trade off data safety for some performance

LUG 2011

# Data Locality

- HDFS reads are efficient ONLY on nodes that store data
  - Not network optimized (HTTP, no DIRECTIO, no DMA)
  - No striping = no aggregating drive bandwidth
  - 1GigE = 100MB/s = quick network saturation for non-local reads
  - Reduced replication = reduced node flexibility

- Lustre reads are equally efficient on any client node
  - Flexible number of map tasks
  - Arbitrary choice of mapper nodes
  - Better cluster utilization

- Lustre reads are fast
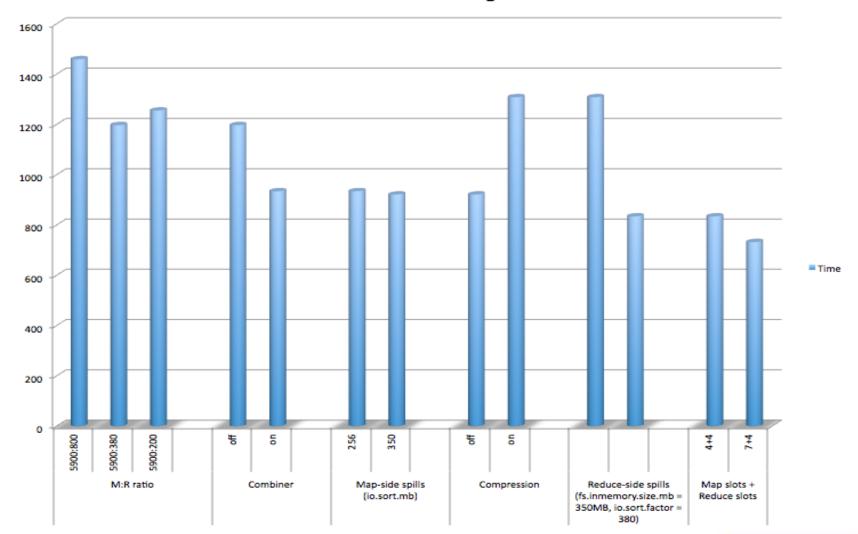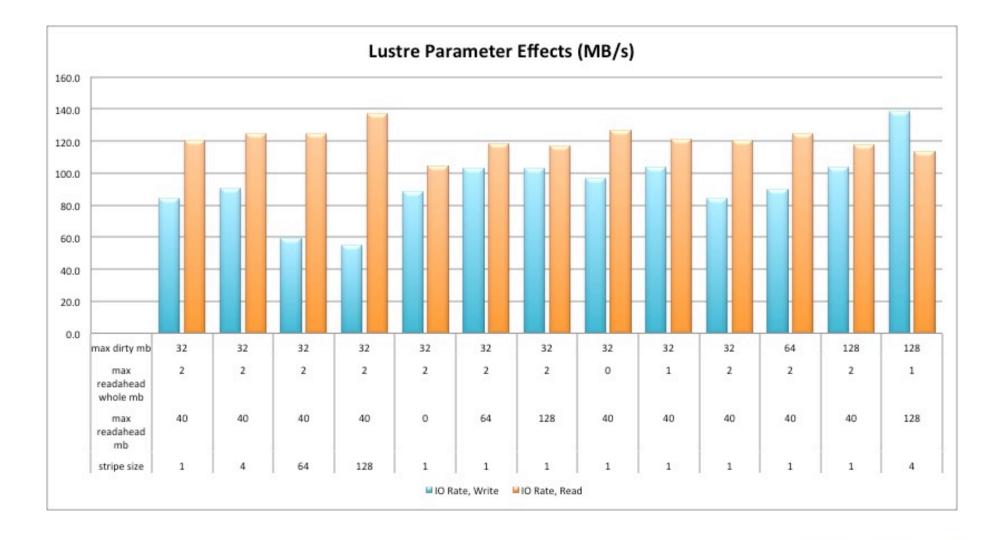  - Striping aggregates disk bandwidth

**xyratex**

# MR I/O Benchmark



TestDFSIO Average IO Rate, MB/s

LUG 2011

# MR Sort Benchmark



**Hadoop Sort Throughput**
**3.2GB data (MB/s)**

LUG 2011

**xyratex**

# MR tuning

Affect of Vaious MR Tuning Parameters

Lustre Parameter Effects (MB/s)

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| max dirty mb | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 64 | 128 | 128 |
| max readahead whole mb | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 0 | 1 | 2 | 2 | 2 | 1 |
| max readahead mb | 40 | 40 | 40 | 40 | 0 | 64 | 128 | 40 | 40 | 40 | 40 | 40 | 128 |
| stripe size | 1 | 4 | 64 | 128 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 4 |

IO Rate, Write    IO Rate, Read

**xyratex·**

# Data Staging: Not a Fair Comparison



Copy 1GB file

LUG 2011

**xyratex**

# Hypothetical Cost Comparison

- Assume Lustre IB has 2x performance of HDFS 1GigE
  - 3x for our sort benchmark
  - Top 500 LINPACK efficiency: 1GigE ~45-50%, 4xQDR ~90-95%

| | Lustre / IB Cluster | | | HDFS / 1 GigE Cluster | | |
|---|---|---|---|---|---|---|
| | Count | Price | Subtotal | Count | Price | Subtotal |
| Nodes | 100 | $7,500 | $750,000 | 200 | $7,500 | $1,500,000 |
| Switches | 9 | $6,500 | $58,500 | 12 | $4,000 | $48,000 |
| Cables | 178 | $100 | $17,800 | 450 | $10 | $4,500 |
| OSS | 2 | $52,000 | $104,000 | 0 | --- | --- |
| Storage | 128TB | --- | --- | 384TB | $100 | $38,400 |
| MDS | 1 | $34,000 | $34,000 | 0 | --- | --- |
| Racks | 4 | $8,000 | $32,000 | 6 | $8,000 | $48,000 |
| **Total** | | | **$996,300** | | | **$1,638,900** |

**xyratex**

# Cost Considerations

- Client node count dominates the overall cost of the cluster
- Doubling size = doubling power, cooling, maintenance costs
- Cluster utilization efficiency
- Data transfer time
- Necessity of maintaining a second cluster

**xyratex**

# Conclusions

- HPC environments have fast networks
- MR should show theoretical performance gains on an appropriately-designed Lustre cluster
- Test results on a small cluster support these propositions
- Performance effects for a particular job may vary widely
- No reason why Hadoop and Lustre can't live happily together
  - Shared storage
  - Shared compute nodes
  - Better performance

x·y·r·a·t·e·x·

# Fini

Thanks!

**xyratex**