

MapReduce: Simplified Data Processing on Large Clusters

These are slides with a history. I found them on the web... They are apparently based on Dan Weld's class at U. Washington, (who in turn based his slides on those by Jeff Dean, Sanjay Ghemawat, Google, Inc.)

Motivation

- Large-Scale Data Processing
 - Want to use 1000s of CPUs
 - But don't want hassle of *managing* things
- MapReduce provides
 - Automatic parallelization & distribution
 - Fault tolerance
 - I/O scheduling
 - Monitoring & status updates

Map/Reduce

- Map/Reduce
 - Programming model from Lisp
 - (and other functional languages)
- Many problems can be phrased this way
- Easy to distribute across nodes
- Nice retry/failure semantics

Map in Lisp (Scheme)

- $(\text{map } f \text{ list} [list_1, list_2, \dots])$
- $(\text{map square } '(1 2 3 4))$
 - $(1 4 9 16)$
- $(\text{reduce } + \text{ '(1 4 9 16)})$
 - 30
- $(\text{reduce } + (\text{map square } (\text{map } - l_1 l_2))))$

Unary operator

Binary operator

Map/Reduce ala Google

- `map(key, val)` is run on each item in set
 - emits new-key / new-val pairs
- `reduce(key, vals)` is run for each unique key emitted by `map()`
 - emits final output
- Often, one application will need to run map/reduce many times in succession

count words in docs

- Input consists of (url, contents) pairs
- map(key=url, val=contents):
 - For each word w in contents, emit $(w, 1)$
- reduce(key=word, values=uniq_counts):
 - Sum all “1”s in values list
 - Emit result “(word, sum)”

Count, Illustrated

map(key=url, val=contents):

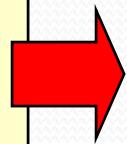
For each word w in contents, emit (w, "1")

reduce(key=word, values=uniq_counts):

Sum all "1"s in values list

Emit result "(word, sum)"

see bob throw
see spot run



see	1	bob	1
bob	1	run	1
run	1	see	2
see	1	spot	1
spot	1	throw	1
throw	1		

Grep

- Input consists of (url+offset, single line)
- map(key=url+offset, val=line):
 - If contents matches regexp, emit (line, “1”)
- reduce(key=line, values=uniq_counts):
 - Don’t do anything; just emit line

Reverse Web-Link Graph

- Map
 - For each URL linking to target, ...
 - Output $\langle \text{target}, \text{source} \rangle$ pairs
- Reduce
 - Concatenate list of all source URLs
 - Outputs: $\langle \text{target}, \text{list } (\text{source}) \rangle$ pairs

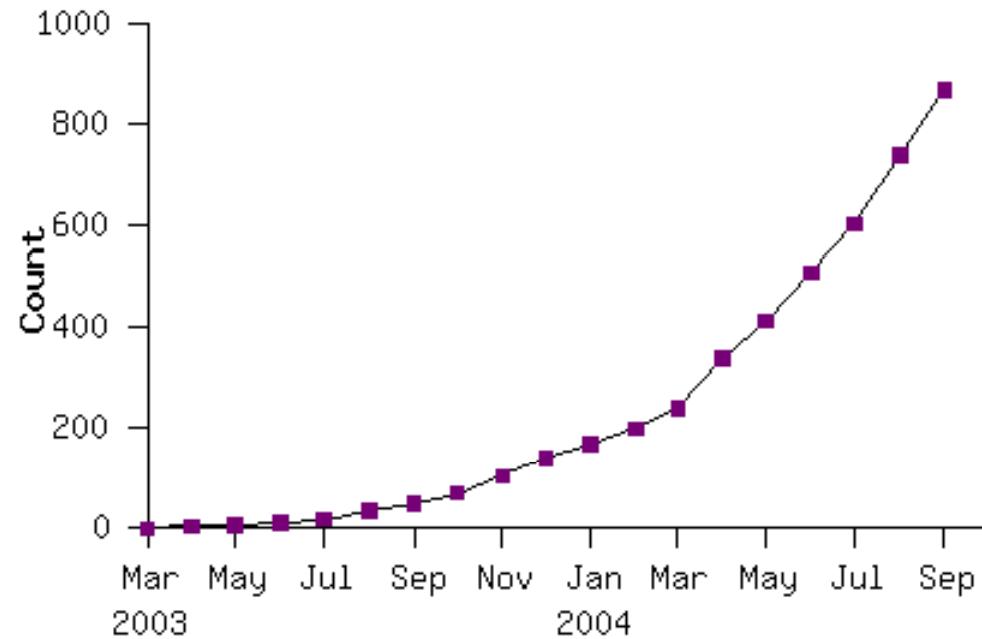
Index maps words to files

Compute an *Inverted Index*

- Map
 - For each file f and each word in the file w
 - Output(f, w) pairs
- Reduce
 - Merge, eliminating duplicates

Model is Widely Applicable

MapReduce Programs In Google Source Tree



Example uses:

distributed grep

term-vector / host

document clustering

...

distributed sort

web access log stats

machine learning

...

web link-graph reversal

inverted index construction

statistical machine
translation

...

Implementation Overview

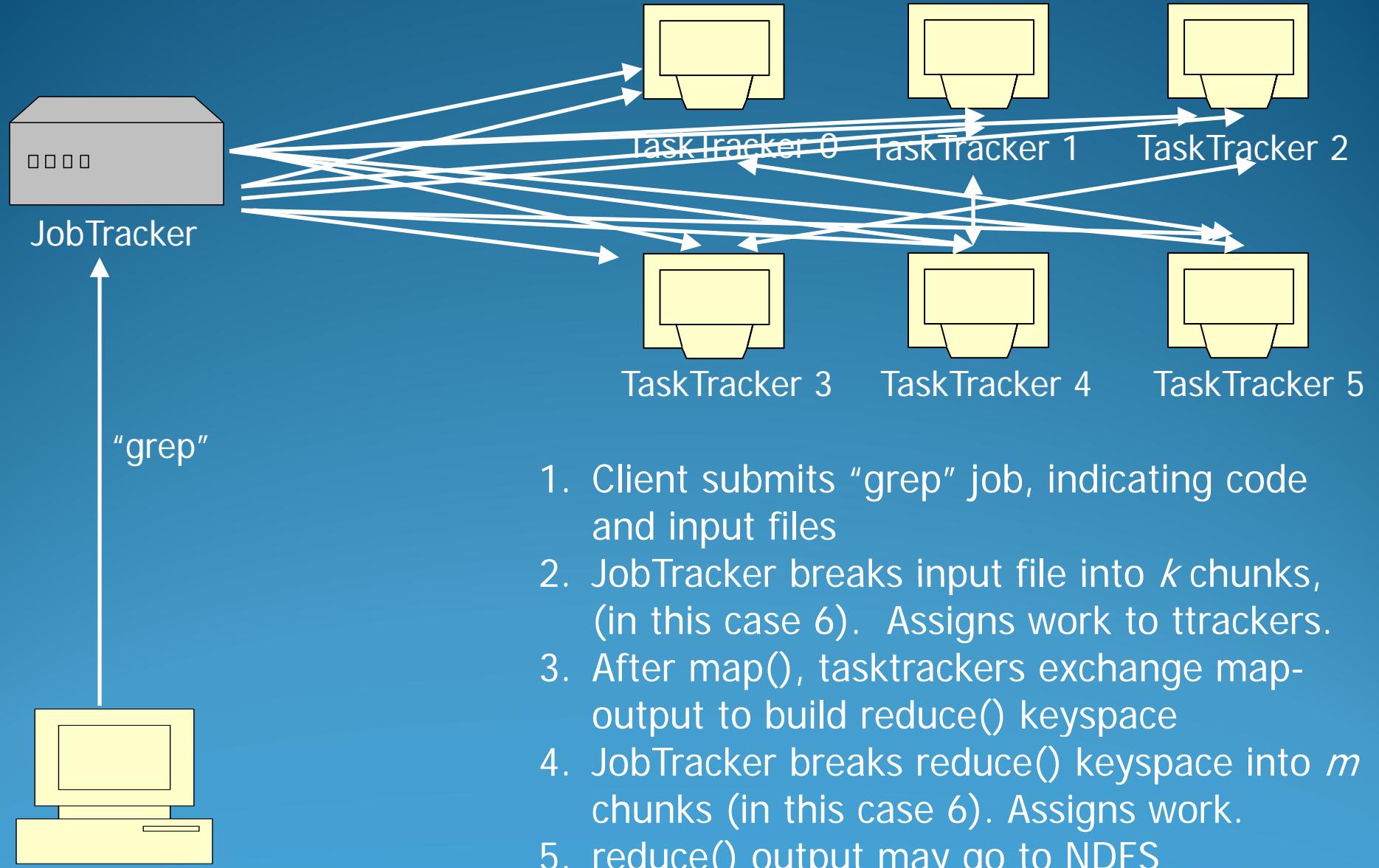
- 100s/1000s of 2-CPU x86 machines, 2-4 GB of memory
- Limited bisection bandwidth
- Storage is on local IDE disks
- GFS: distributed file system manages data (SOSP'03)
- Job scheduling system: jobs made up of tasks, scheduler assigns tasks to machines

Implementation is a C++ library linked into user programs

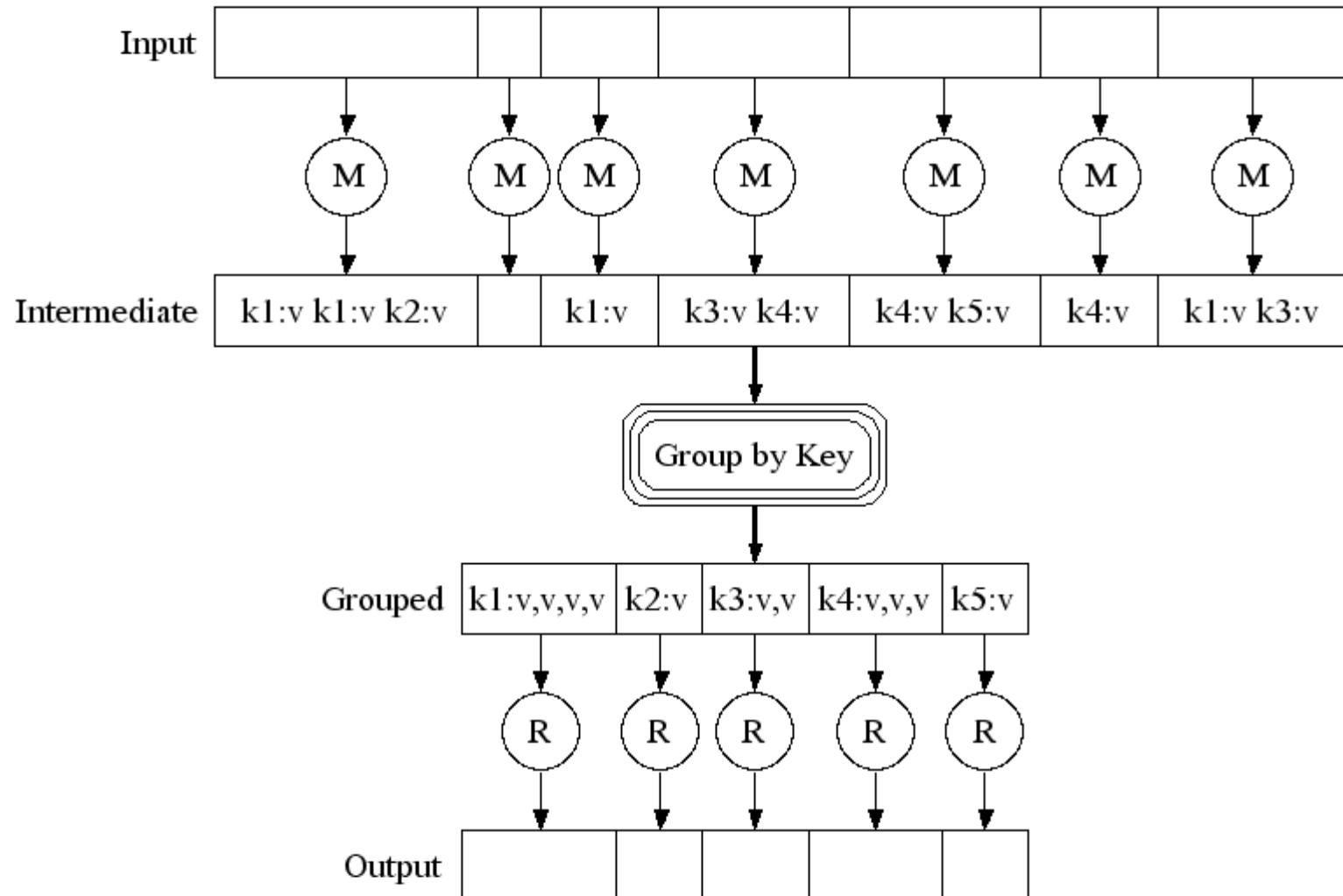
Execution

- How is this distributed?
 1. Partition input key/value pairs into chunks, run map() tasks in parallel
 2. After all map()s are complete, consolidate all emitted values for each unique emitted key
 3. Now partition space of output map keys, and run reduce() in parallel
- If map() or reduce() fails, reexecute!

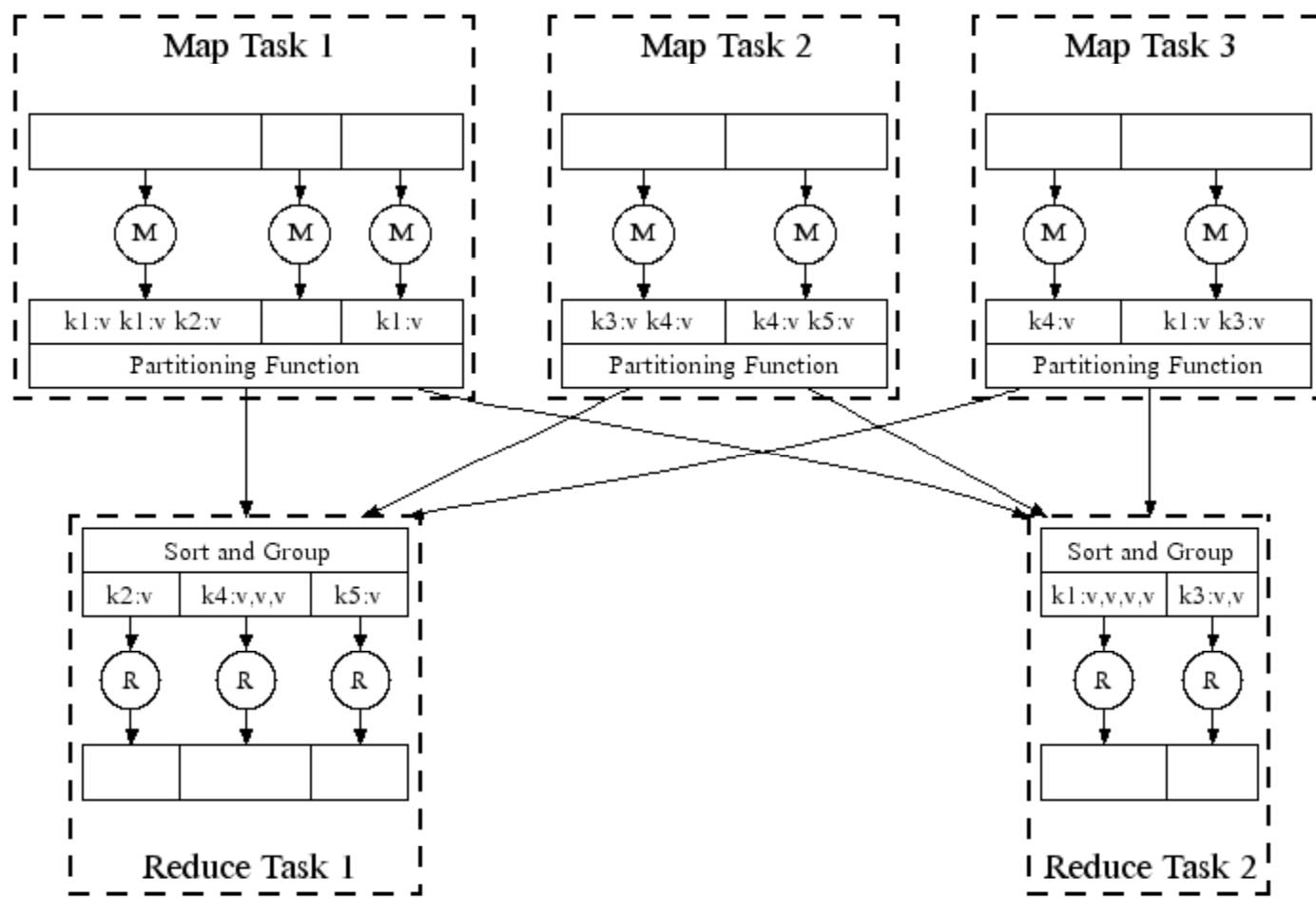
Job Processing



Execution

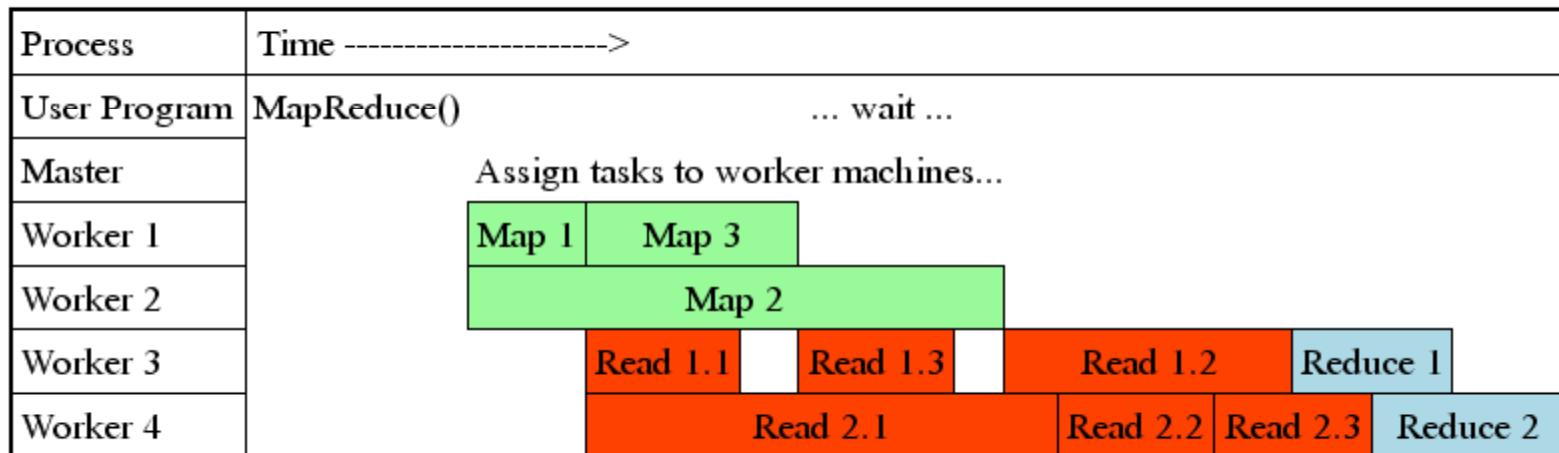


Parallel Execution



Task Granularity & Pipelining

- Fine granularity tasks: map tasks >> machines
 - Minimizes time for fault recovery
 - Can pipeline shuffling with map execution
 - Better dynamic load balancing
- Often use 200,000 map & 5000 reduce tasks, running on 2000 machines

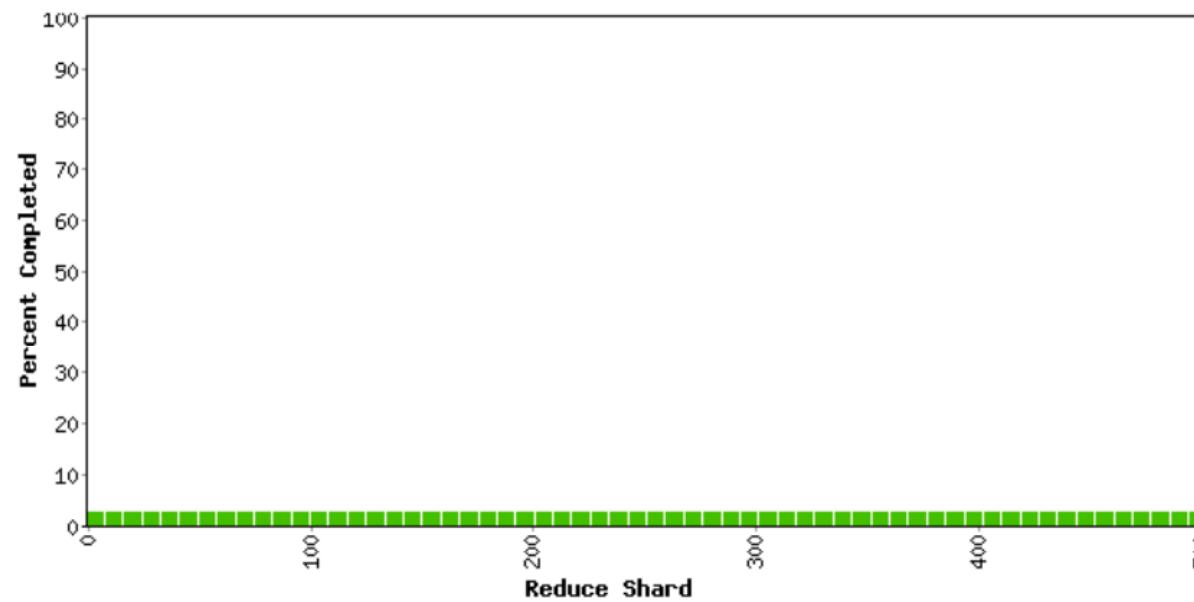


MapReduce status: MR_Indexer-beta6-large-2003_10_28_00_03

Started: Fri Nov 7 09:51:07 2003 -- up 0 hr 00 min 18 sec

323 workers; 0 deaths

Type	Shards	Done	Active	Input(MB)	Done(MB)	Output(MB)
Map	13853	0	323	878934.6	1314.4	717.0
Shuffle	500	0	323	717.0	0.0	0.0
Reduce	500	0	0	0.0	0.0	0.0



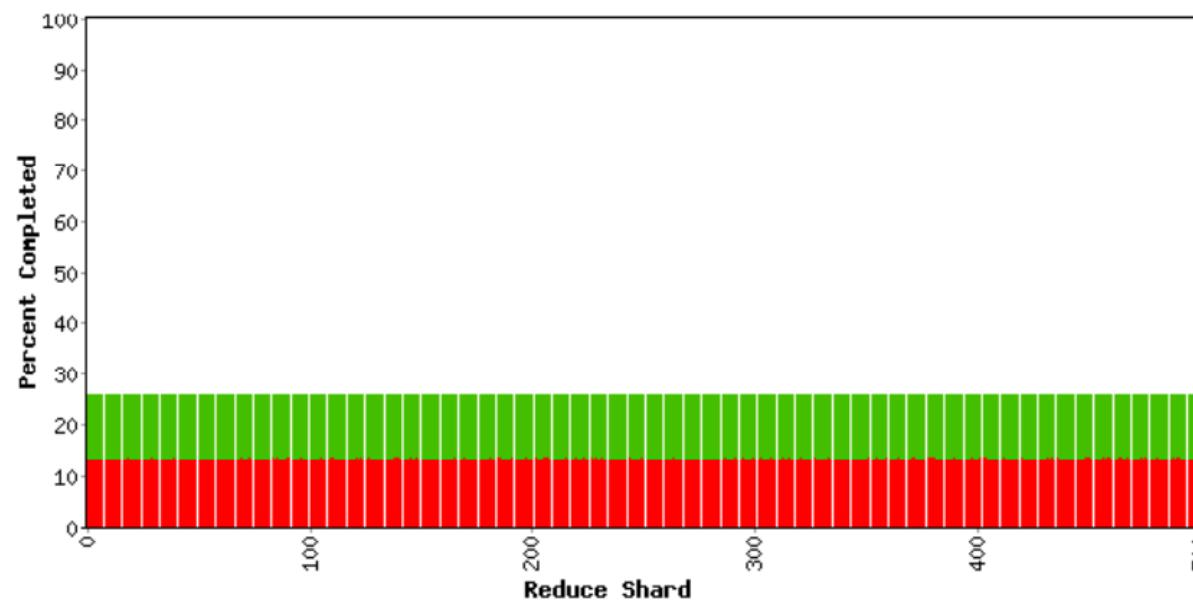
Counters	
Variable	Minute
Mapped (MB/s)	72.5
Shuffle (MB/s)	0.0
Output (MB/s)	0.0
doc-index-hits	145825686
docs-indexed	506631
dups-in-index-merge	0
mr-operator-calls	508192
mr-operator-commits	506631

MapReduce status: MR_Indexer-beta6-large-2003_10_28_00_03

Started: Fri Nov 7 09:51:07 2003 -- up 0 hr 05 min 07 sec

1707 workers; 1 deaths

Type	Shards	Done	Active	Input(MB)	Done(MB)	Output(MB)
Map	13853	1857	1707	878934.6	191995.8	113936.6
Shuffle	500	0	500	113936.6	57113.7	57113.7
Reduce	500	0	0	57113.7	0.0	0.0



Counters

Variable	Minute
Mapped (MB/s)	699.1
Shuffle (MB/s)	349.5
Output (MB/s)	0.0
doc-index-hits	5004411944
docs-indexed	17290135
dups-in-index-merge	0
mr-operator-calls	17331371
mr-operator-outouts	17290135

MapReduce status: MR_Indexer-beta6-large-2003_10_28_00_03

Started: Fri Nov 7 09:51:07 2003 -- up 0 hr 10 min 18 sec

1707 workers; 1 deaths

Type	Shards	Done	Active	Input(MB)	Done(MB)	Output(MB)
Map	13853	5354	1707	878934.6	406020.1	241058.2
Shuffle	500	0	500	241058.2	196362.5	196362.5
Reduce	500	0	0	196362.5	0.0	0.0



Counters

Variable	Minute
Mapped (MB/s)	704.4
Shuffle (MB/s)	371.9
Output (MB/s)	0.0
doc-index-hits	5000364228
docs-indexed	17300709
dups-in-index-merge	0
mr-operator-calls	17342493
mr-operator-outputs	17300709

MapReduce status: MR_Indexer-beta6-large-2003_10_28_00_03

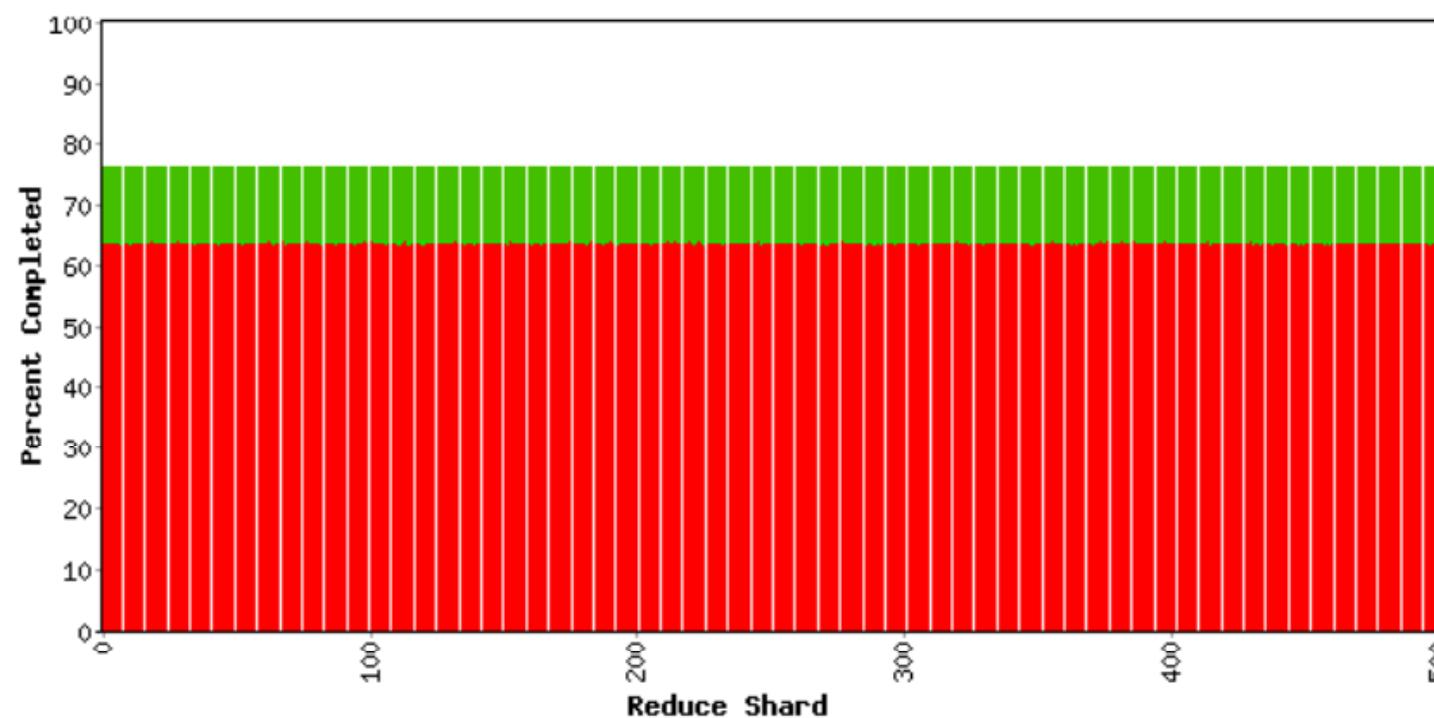
Started: Fri Nov 7 09:51:07 2003 -- up 0 hr 15 min 31 sec

1707 workers; 1 deaths

Type	Shards	Done	Active	Input(MB)	Done(MB)	Output(MB)
Map	13853	8841	1707	878934.6	621608.5	369459.8
Shuffle	500	0	500	369459.8	326986.8	326986.8
Reduce	500	0	0	326986.8	0.0	0.0

Counters

Variable	Minute
Mapped (MB/s)	706.5
Shuffle (MB/s)	419.2
Output (MB/s)	0.0
doc-index-hits	4982870667
docs-indexed	17229926
dups-in-index-merge	0
mr-operator-calls	17272056
mr-operator-outputs	17229926

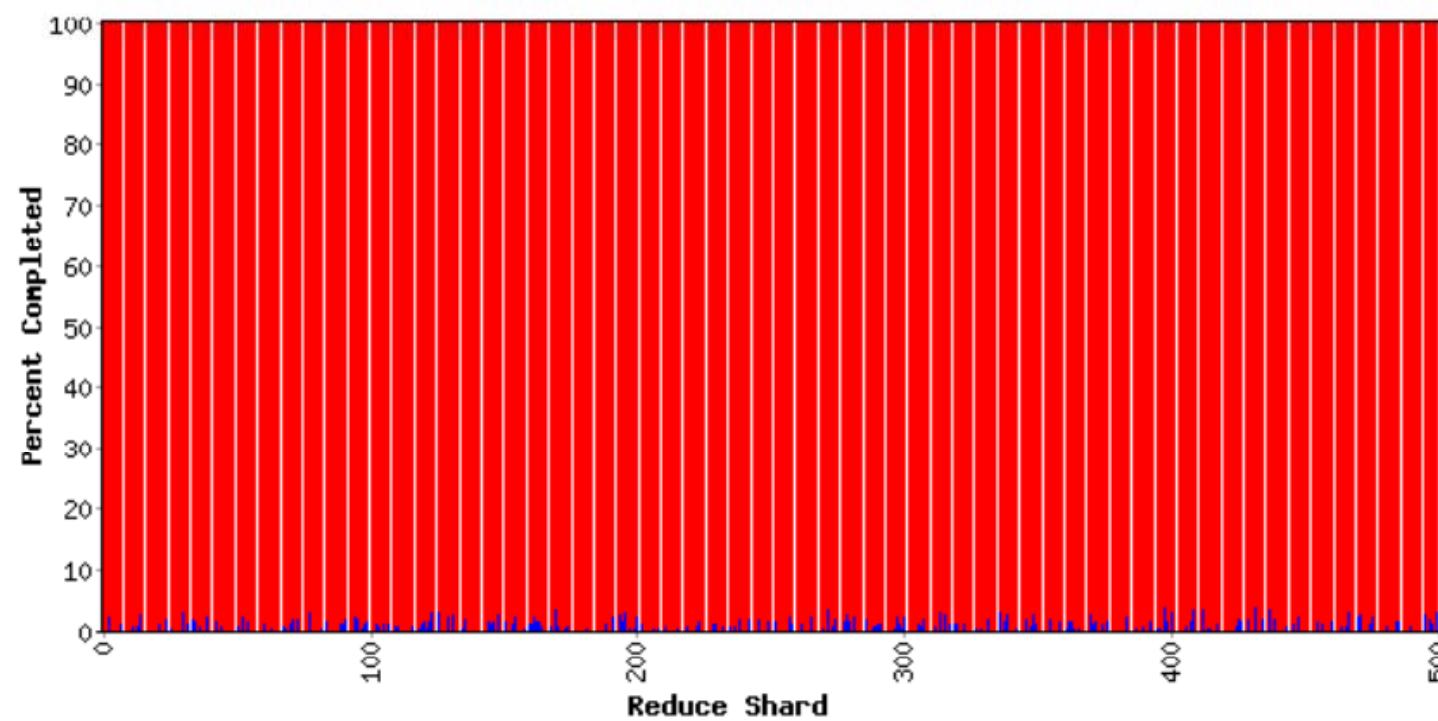


MapReduce status: MR_Indexer-beta6-large-2003_10_28_00_03

Started: Fri Nov 7 09:51:07 2003 -- up 0 hr 29 min 45 sec

1707 workers; 1 deaths

Type	Shards	Done	Active	Input(MB)	Done(MB)	Output(MB)
Map	13853	13853	0	878934.6	878934.6	523499.2
Shuffle	500	195	305	523499.2	523389.6	523389.6
Reduce	500	0	195	523389.6	2685.2	2742.6



Counters

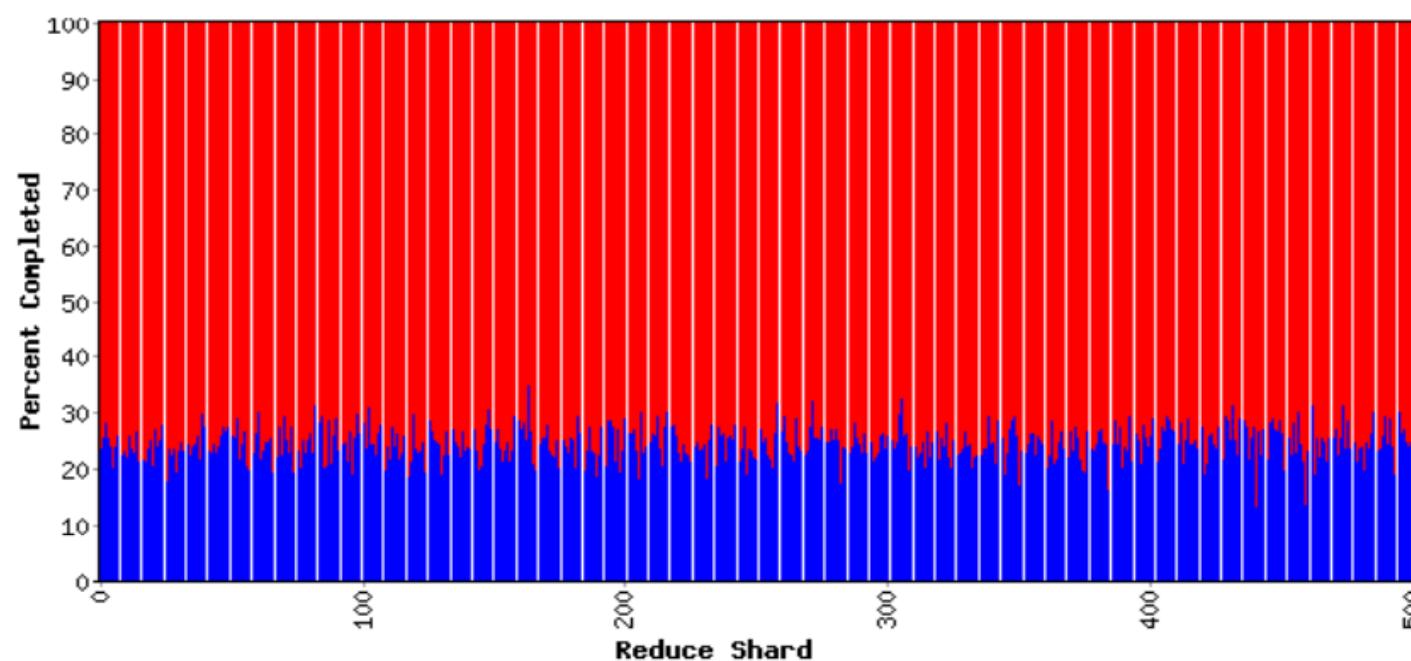
Variable	Minute	
Mapped (MB/s)	0.3	
Shuffle (MB/s)	0.5	
Output (MB/s)	45.7	
doc-index-hits	2313178	10 ⁶
docs-indexed	7936	
dups-in-index-merge	0	
mr-merge-calls	1954105	
mr-merge-outputs	1954105	

MapReduce status: MR_Indexer-beta6-large-2003_10_28_00_03

Started: Fri Nov 7 09:51:07 2003 -- up 0 hr 31 min 34 sec

1707 workers; 1 deaths

Type	Shards	Done	Active	Input(MB)	Done(MB)	Output(MB)
Map	13853	13853	0	878934.6	878934.6	523499.2
Shuffle	500	500	0	523499.2	523499.5	523499.5
Reduce	500	0	500	523499.5	133837.8	136929.6



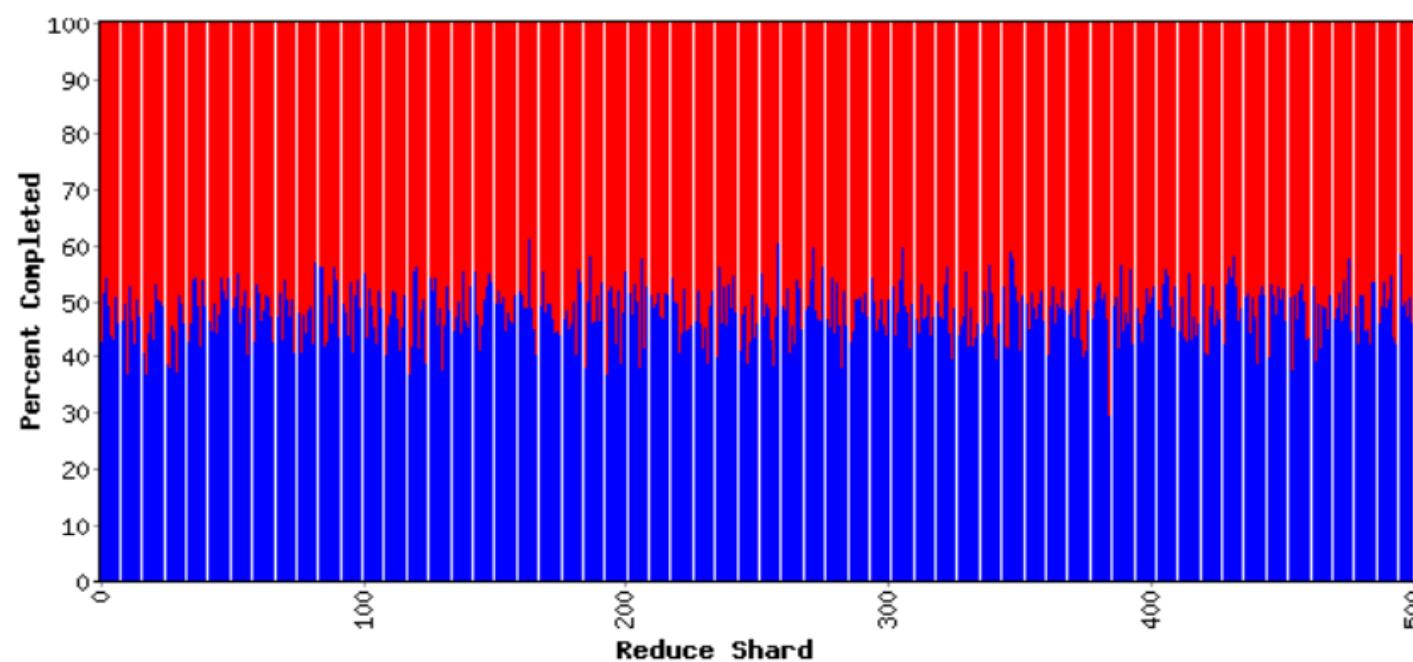
Counters	
Variable	Minute
Mapped (MB/s)	0.0
Shuffle (MB/s)	0.1
Output (MB/s)	1238.8
doc-index-hits	0.10
docs-indexed	0
dups-in-index-merge	0
mr-merge-calls	51738599
mr-merge-outputs	51738599

MapReduce status: MR_Indexer-beta6-large-2003_10_28_00_03

Started: Fri Nov 7 09:51:07 2003 -- up 0 hr 33 min 22 sec

1707 workers; 1 deaths

Type	Shards	Done	Active	Input(MB)	Done(MB)	Output(MB)
Map	13853	13853	0	878934.6	878934.6	523499.2
Shuffle	500	500	0	523499.2	523499.5	523499.5
Reduce	500	0	500	523499.5	263283.3	269351.2



Counters	
Variable	Minute
Mapped (MB/s)	0.0
Shuffle (MB/s)	0.0
Output (MB/s)	1225.1
doc-index-hits	0 10
docs-indexed	0
dups-in-index-merge	0
mr-merge-calls	51842100
mr-merge-outputs	51842100

MapReduce status: MR_Indexer-beta6-large-2003_10_28_00_03

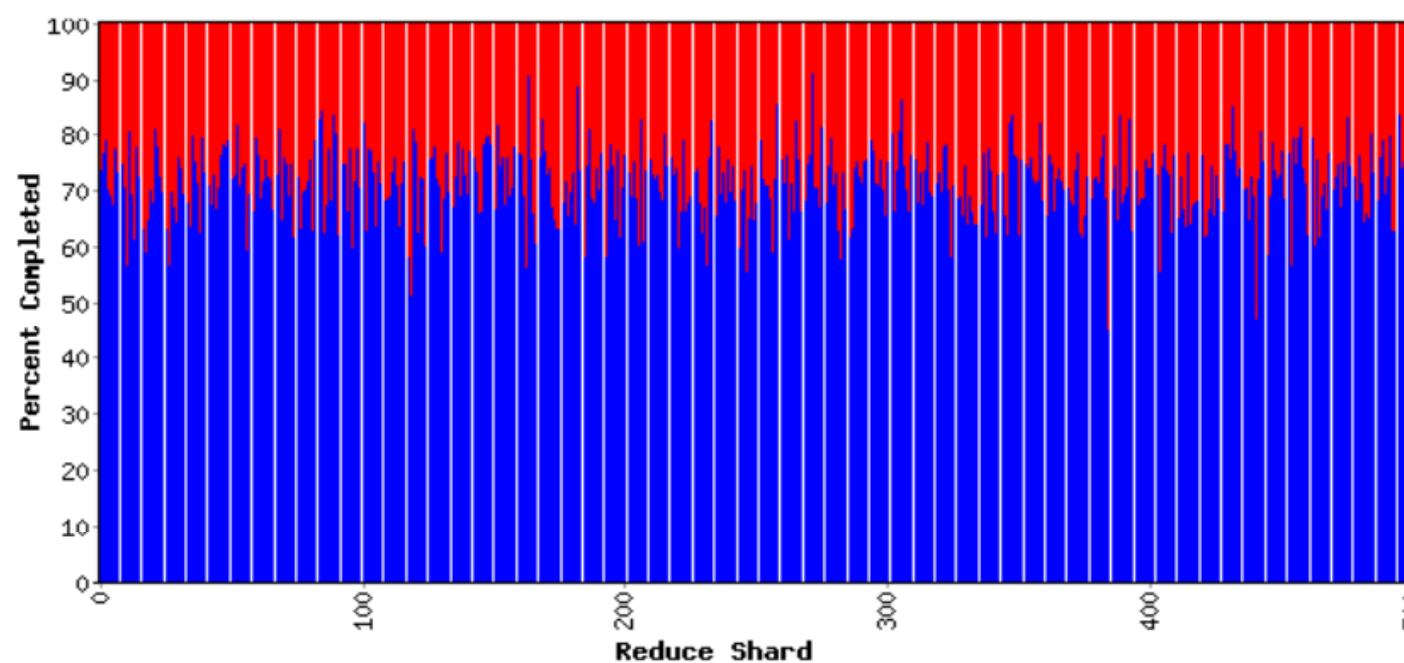
Started: Fri Nov 7 09:51:07 2003 -- up 0 hr 35 min 08 sec

1707 workers; 1 deaths

Type	Shards	Done	Active	Input(MB)	Done(MB)	Output(MB)
Map	13853	13853	0	878934.6	878934.6	523499.2
Shuffle	500	500	0	523499.2	523499.5	523499.5
Reduce	500	0	500	523499.5	390447.6	399457.2

Counters

Variable	Minute
Mapped (MB/s)	0.0
Shuffle (MB/s)	0.0
Output (MB/s)	1222.0
doc-index-hits	0 100
docs-indexed	0
dups-in-index-merge	0
mr-merge-calls	51640600
mr-merge-outputs	51640600

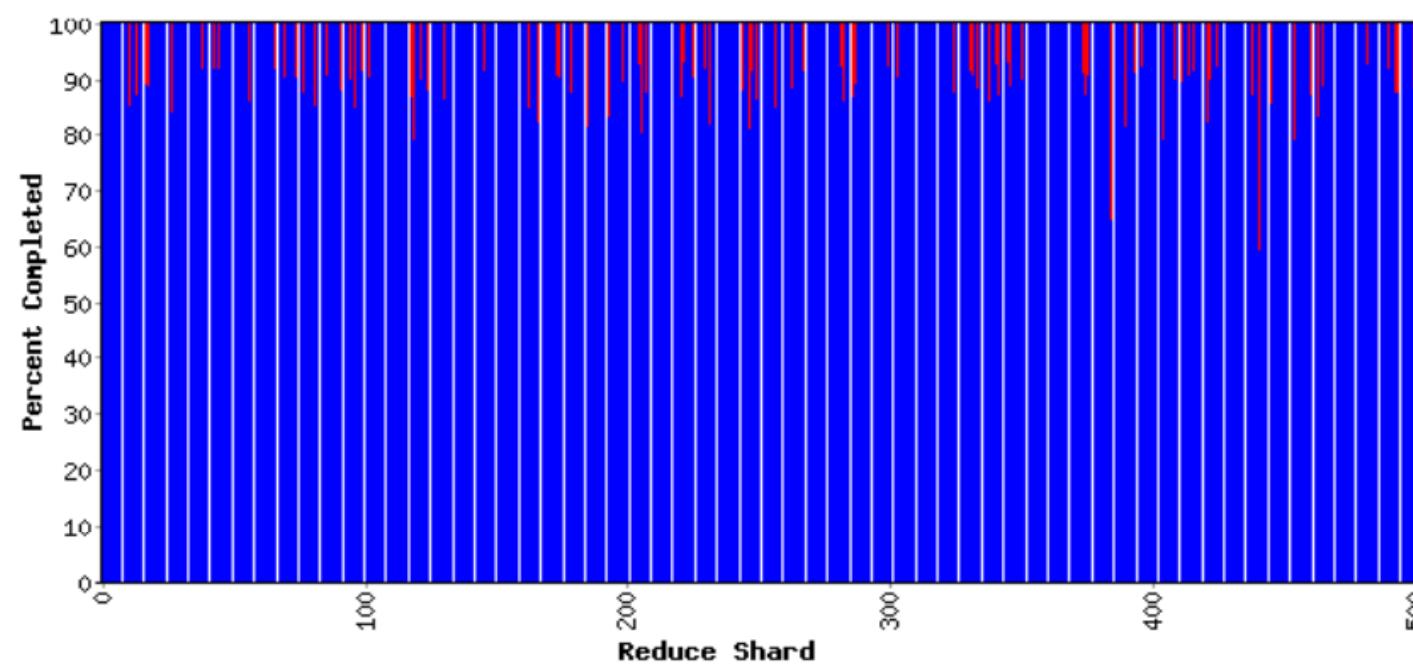


MapReduce status: MR_Indexer-beta6-large-2003_10_28_00_03

Started: Fri Nov 7 09:51:07 2003 -- up 0 hr 37 min 01 sec

1707 workers; 1 deaths

Type	Shards	Done	Active	Input(MB)	Done(MB)	Output(MB)
Map	13853	13853	0	878934.6	878934.6	523499.2
Shuffle	500	500	0	523499.2	520468.6	520468.6
Reduce	500	406	94	520468.6	512265.2	514373.3



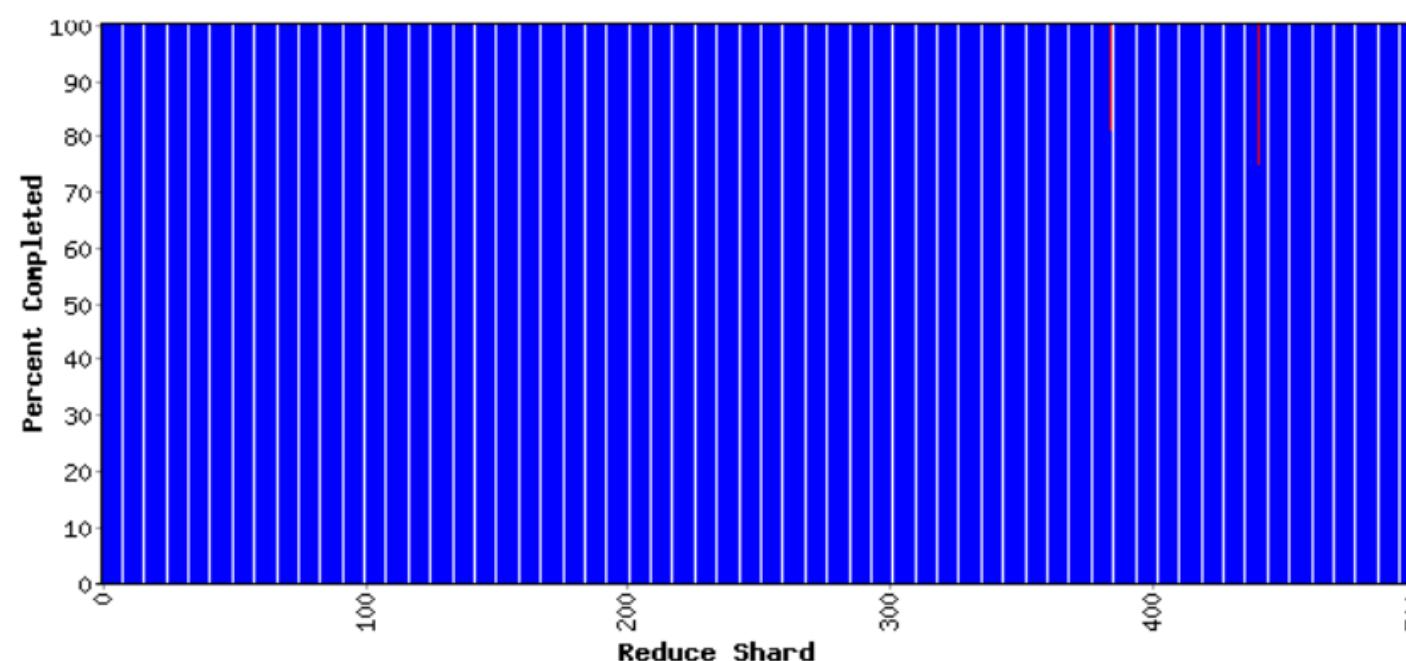
Counters	
Variable	Minute
Mapped (MB/s)	0.0
Shuffle (MB/s)	0.0
Output (MB/s)	849.5
doc-index-hits	0 100
docs-indexed	0
dups-in-index-merge	0
mr-merge-calls	35083350
mr-merge-outputs	35083350

MapReduce status: MR_Indexer-beta6-large-2003_10_28_00_03

Started: Fri Nov 7 09:51:07 2003 -- up 0 hr 38 min 56 sec

1707 workers; 1 deaths

Type	Shards	Done	Active	Input(MB)	Done(MB)	Output(MB)
Map	13853	13853	0	878934.6	878934.6	523499.2
Shuffle	500	500	0	523499.2	519781.8	519781.8
Reduce	500	498	2	519781.8	519394.7	519440.7



Counters

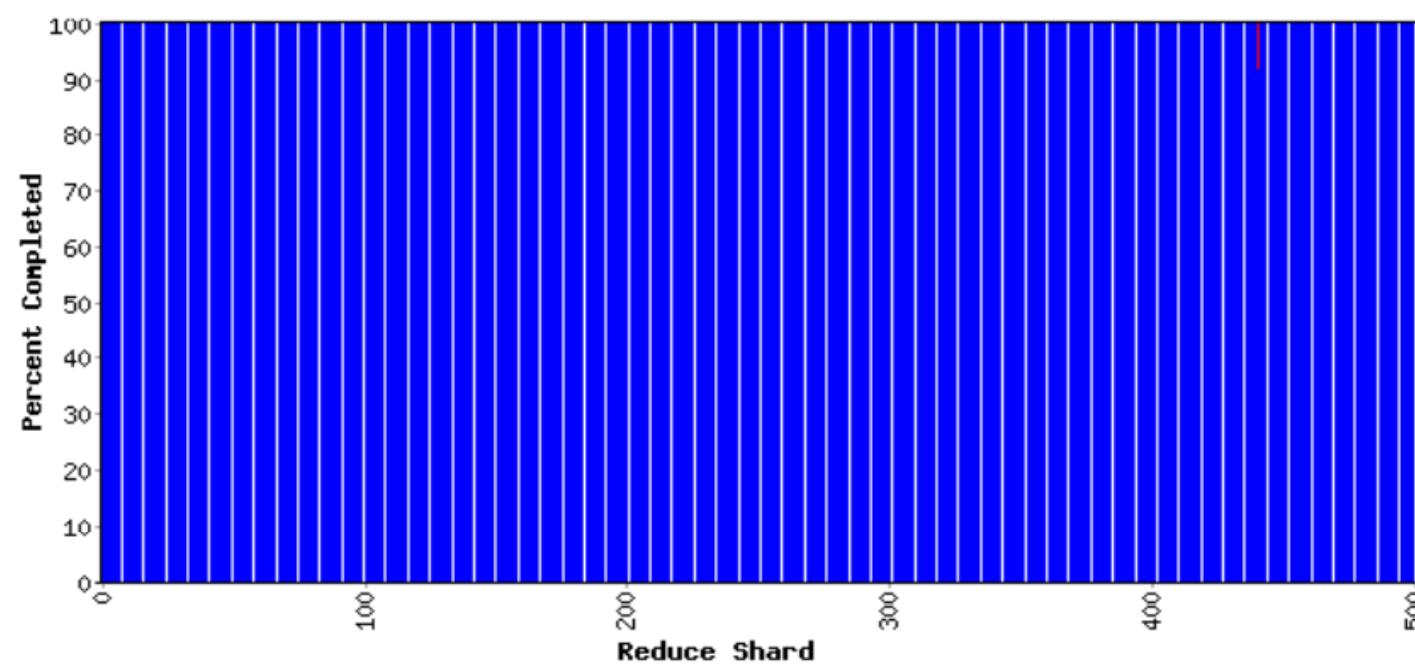
Variable	Minute	Second
Mapped (MB/s)	0.0	
Shuffle (MB/s)	0.0	
Output (MB/s)	9.4	
doc-index-hits	0	1056
docs-indexed	0	
dups-in-index-merge	0	
mr-merge-calls	394792	
mr-merge-outputs	394792	

MapReduce status: MR_Indexer-beta6-large-2003_10_28_00_03

Started: Fri Nov 7 09:51:07 2003 -- up 0 hr 40 min 43 sec

1707 workers; 1 deaths

Type	Shards	Done	Active	Input(MB)	Done(MB)	Output(MB)
Map	13853	13853	0	878934.6	878934.6	523499.2
Shuffle	500	500	0	523499.2	519774.3	519774.3
Reduce	500	499	1	519774.3	519735.2	519764.0



Counters

Variable	Minute	Value
Mapped (MB/s)	0.0	0
Shuffle (MB/s)	0.0	0
Output (MB/s)	1.9	0
doc-index-hits	0	105
docs-indexed	0	0
dups-in-index-merge	0	0
mr-merge-calls	73442	0
mr-merge-outputs	73442	0

Fault Tolerance / Workers

Handled via re-execution

- Detect failure via periodic heartbeats
- Re-execute completed + in-progress *map* tasks
- Re-execute in progress *reduce* tasks
- Task completion committed through master

Robust: lost 1600/1800 machines once → finished ok

Semantics in presence of failures: “at least once”

Master Failure

- Could handle, presumably using the kind of replication mechanisms we'll be studying in near future
- But don't yet
 - (runs are short enough so that master failure is unlikely)

Refinement: Redundant Execution

Slow workers significantly delay completion time

- Other jobs consuming resources on machine
- Bad disks w/ soft errors transfer data slowly
- Weird things: processor caches disabled (!!)

Solution: Near end of phase, spawn backup tasks

- Whichever one finishes first "wins"

Dramatically shortens job completion time

Refinement: Locality Optimization

- Master scheduling policy:
 - Asks GFS for locations of replicas of input file blocks
 - Map tasks typically split into 64MB (GFS block size)
 - Map tasks scheduled so GFS input block replica are on same machine or same rack
- Effect
 - Thousands of machines read input at local disk speed
 - Without this, rack switches limit read rate

Refinement

Skipping Bad Records

- Map/Reduce functions might fail for some inputs
 - Best solution is to debug & fix
 - Not always possible ~ third-party source libraries
 - On segmentation fault:
 - Send UDP packet to master from signal handler
 - Include sequence number of record being processed
 - If master sees two failures for same record:
 - Next worker is told to skip the record

Other Refinements

- Sorting guarantees
 - within each reduce partition
- Compression of intermediate data
- Combiner
 - Useful for saving network bandwidth
- Local execution for debugging/testing
- User-defined counters

Performance

Tests run on cluster of 1800 machines:

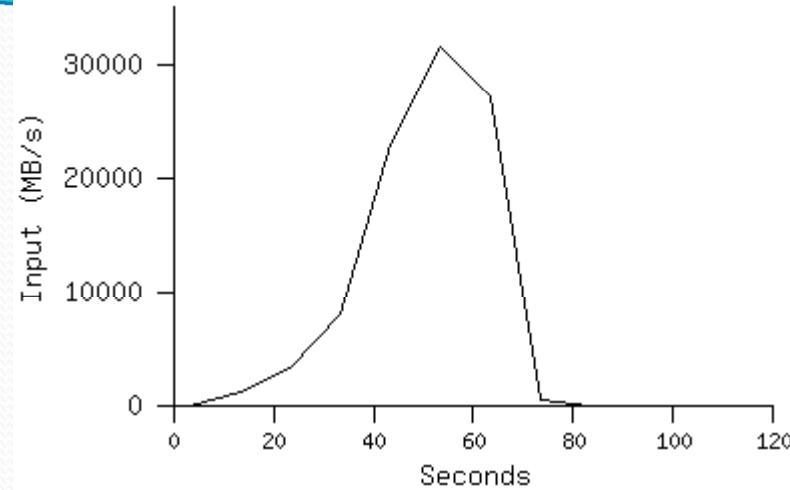
- 4 GB of memory
- Dual-processor 2 GHz Xeons with Hyperthreading
- Dual 160 GB IDE disks
- Gigabit Ethernet per machine
- Bisection bandwidth approximately 100 Gbps

Two benchmarks:

MR_GrepScan 1010 100-byte records to extract records
 matching a rare pattern (92K matching records)

MR_SortSort 1010 100-byte records (modeled after TeraSort
 benchmark)

MR_Grep



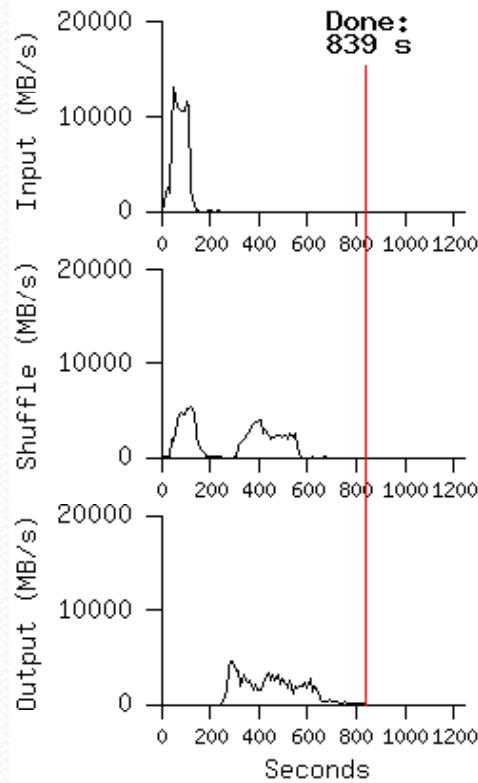
Locality optimization helps:

- 1800 machines read 1 TB at peak ~31 GB/s
- W/out this, rack switches would limit to 10 GB/s

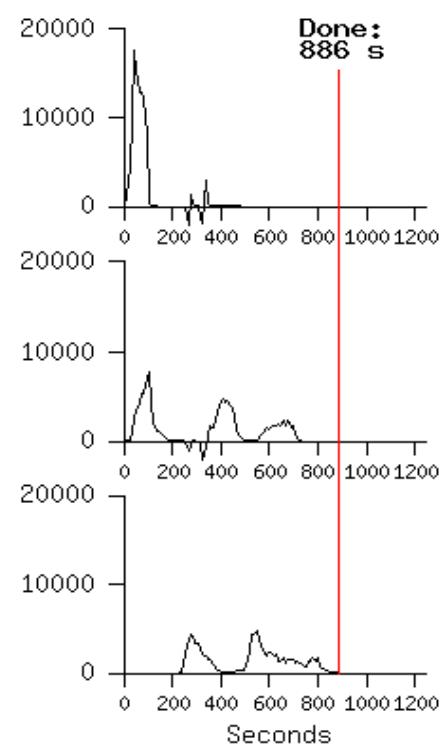
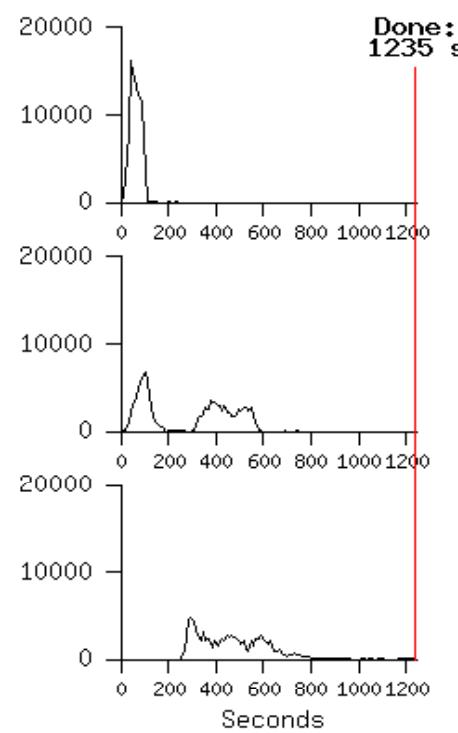
Startup overhead is significant for short jobs

MR_Sort

Normal



No backup tasks 200 processes killed



- Backup tasks reduce job completion time a lot!
- System deals well with failures

Experience

Rewrote Google's production indexing
System using MapReduce

- Set of 10, 14, 17, 21, 24 MapReduce operations
- New code is simpler, easier to understand
 - 3800 lines C++ → 700
- MapReduce handles failures, slow machines
- Easy to make indexing faster
 - add more machines

Usage in Aug 2004

Number of jobs	29,423
Average job completion time	634 secs
Machine days used	79,186 days
Input data read	3,288 TB
Intermediate data produced	758 TB
Output data written	193 TB
Average worker machines per job	157
Average worker deaths per job	1.2
Average map tasks per job	3,351
Average reduce tasks per job	55
Unique <i>map</i> implementations	395
Unique <i>reduce</i> implementations	269
Unique <i>map/reduce</i> combinations	426

Underlying technologies used

- Implementation of Map/Reduce made use of other cloud computing services available at Google
 - System management tools track the available nodes, configurations, current loads
 - Chubby “locking” tool for synchronization
 - Google file system (GFS) provides convenient storage, makes it easy to gather the inputs needed for Reduce (write locally, anywhere, and read anywhere)
 - Big Table: a table-structured database, runs over GFS

Related Work

- Programming model inspired by functional language primitives
- Partitioning/shuffling similar to many large-scale sorting systems
 - NOW-Sort ['97]
- Re-execution for fault tolerance
 - BAD-FS ['04] and TACC ['97]
- Locality optimization has parallels with Active Disks/Diamond work
 - Active Disks ['01], Diamond ['04]
- Backup tasks similar to Eager Scheduling in Charlotte system
 - Charlotte ['96]
- Dynamic load balancing solves similar problem as River's distributed queues
 - River ['99]



Conclusions

- MapReduce proven to be useful abstraction
- Greatly simplifies large-scale computations
- Easy to use:
 - focus on problem,
 - let library deal w/ messy details