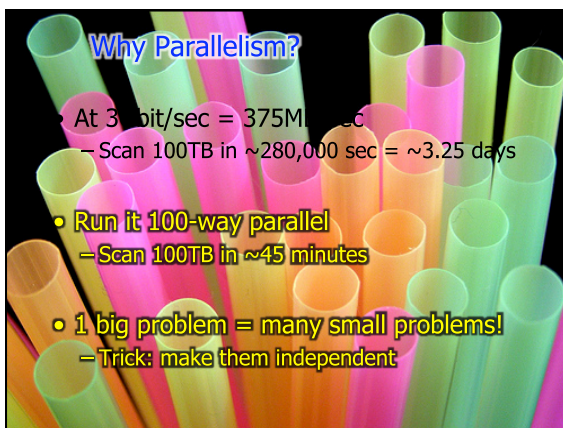


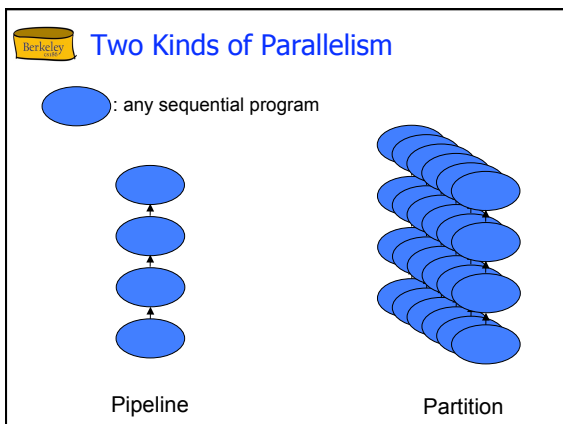
Why Parallelism?

- At 3Gbit/sec = 375MB/sec
 - Scan 100TB in $\sim 280,000$ sec = ~ 3.25 days



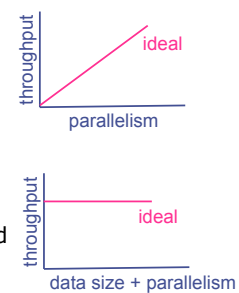
A little history

- Relational revolution
 - declarative set-oriented primitives
 - 1970's
- Parallel relational database systems
 - on commodity hardware
 - 1980's
- Renaissance: MapReduce etc.
 - now



Two Kinds of Benefit

- Speed-up
 - Increase HW
 - Fix workload
- Scale-up
 - Increase HW
 - Increase workload



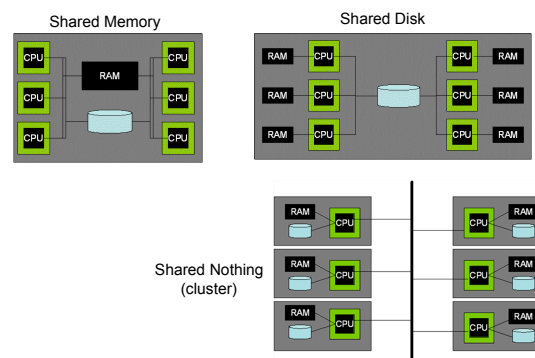


"Big Data" is GREAT for Parallelism!

- Why?
 - Set-oriented languages
 - Batch operations
 - Pre-existing divide-and-conquer algorithms
 - Natural pipelining



Parallel Architectures



Some Early Systems

- Research
 - XPRS (Berkeley, shared-memory)
 - Gamma (Wisconsin, shared-nothing)
 - Volcano (Colorado, shared-nothing)
 - Bubba (MCC, shared-nothing)
- Industry
 - Teradata (shared-nothing)
 - Tandem Non-Stop SQL (shared-nothing)



Uses of Parallelism

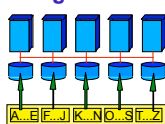
- Inter-query
 - Esp. for Transaction Processing
 - Wait for discussion of Concurrency Control
- Intra-query
 - Inter-operator
 - Tree
 - Pipeline
 - Intra-operator
 - Divide & Conquer
 - Focus here – best bang for the buck



Data Partitioning

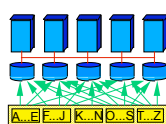
Partitioning a table:

Range



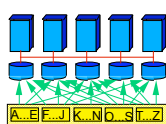
Good for equijoins,
range queries
group-by

Hash



Good for equijoins,
group-by

Round Robin



Good for spreading load

Shared nothing benefits from "good" partitioning
Remind you of something?




Parallel Scans

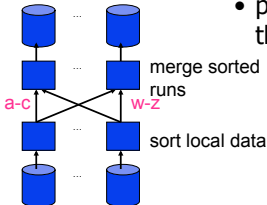
- Scan in parallel, merge (concat) output
- Selection: may skip sites in some cases
 - range or hash partitioning
- Indexes can be built at each partition
- Question: How do indexes differ in the different schemes?
 - Think about both lookups and inserts
 - What about unique indexes (keys)?

Lookup by key

- data partitioned on function of key?
 - great!
- otherwise
 - umm...

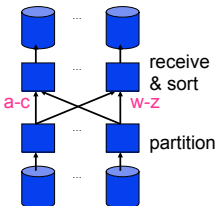


Sort



- problems with this approach?

Sort, improved



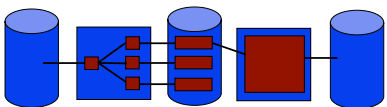
- a key issue: avoiding skew
 - sample to estimate data distribution
 - choose ranges to get uniformity

Sorting Records!

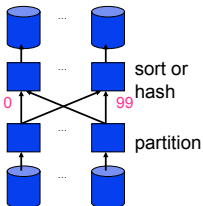
- Sorting has become a blood sport!
 - Parallel sorting is the name of the game ...
- Minute Sort: how many 100-byte records can you sort in a minute?
 - Current World record: 264 **GB**
 - 400 x (dual-3.2GHz processors, 6 disks, 8GB RAM) (2007)
- Penny Sort: how many can you sort for a penny?
 - Current world record: 190 GB
 - 2408 seconds on a \$393 Linux/AMD system (2008)
 - \$393 spread over 3 years worth of seconds = 2408 seconds/penny
- Joule Sort: how many can you sort with a Joule?
- See <http://www.hpl.hp.com/hosted/sortbenchmark/>

Parallel Hashing

- Well ... think about single-site hashing
 - Phase 1: partition input
 - Phase 2: build in-memory hashtable
- Where do you put the communication?
 - And how do you split data across sites?



Grouping



- again, skew is an issue
- approaches:
 - **avoid** (choose partition function carefully)
 - **react** (migrate groups to balance load)

Parallel Aggregates/GroupBy

- For each aggregate function, need a local/global decomposition:
 - $\text{sum}(S) = \Sigma \Sigma(s)$
 - $\text{count} = \Sigma \text{count}(s)$
 - $\text{avg}(S) = (\Sigma \Sigma(s)) / \Sigma \text{count}(s)$
 - etc...
- Group by:
 - Do a local aggregate for each group
 - 3 machines (M^1, \dots, M^3), k "dno" sums each
 - Ship each local agg to an appointed "collector" to compute global agg
 - Hash function maps "dno" to one of 3 collector "buckets" (B_1, \dots, B_3)
 - Sum up each group at exactly 1 collector

SELECT dno, SUM(sales)
FROM dept
GROUP BY dno

Advanced topic: can you think of a summary statistic this doesn't work for?

Parallel Joins

- Nested loop:
 - Compare each outer tuple with each inner tuple that might join.
 - Easy for range/hash partitioning and equijoin
 - Hard otherwise!
- Sort-Merge (or plain Merge-Join):
 - Sorting gives range-partitioning.
 - But what about handling 2 skews?
 - Merging partitioned tables is local.
- Hash-join
 - Hash-partition both inputs
 - build/probe phase is local

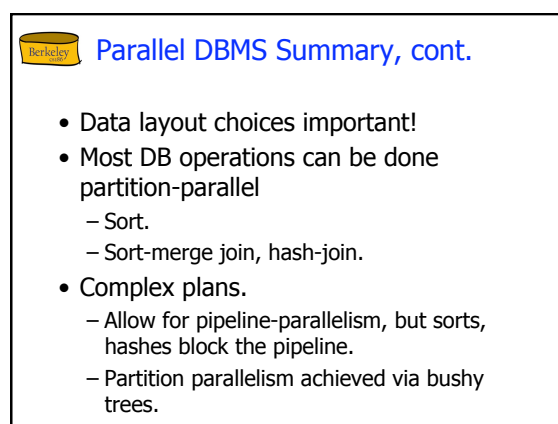
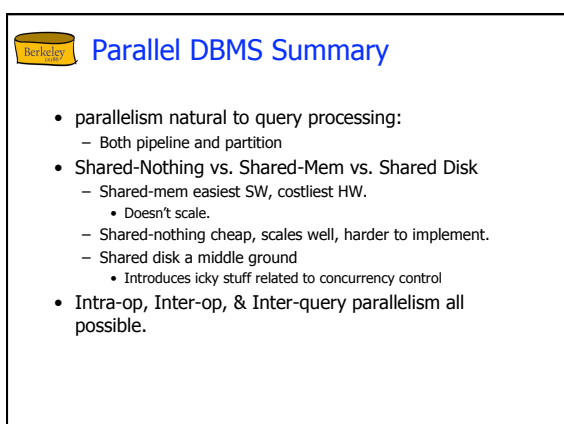
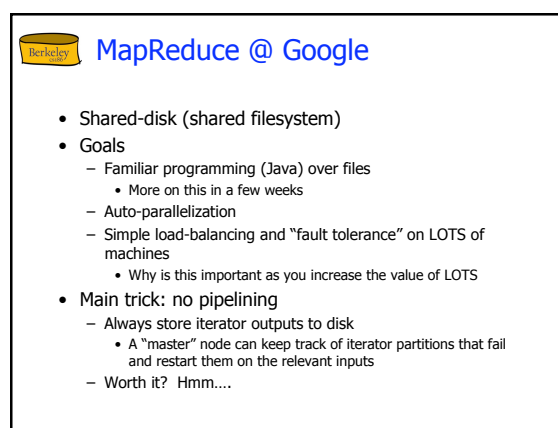
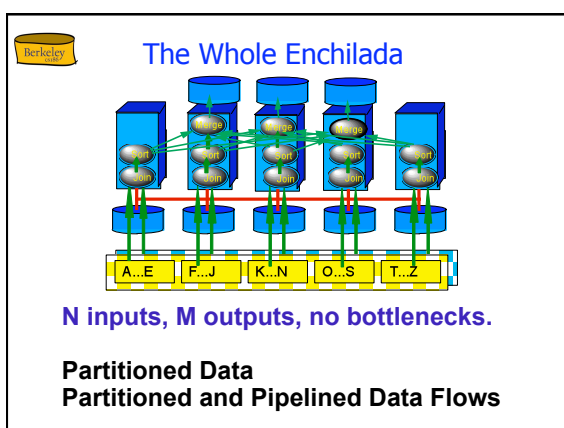
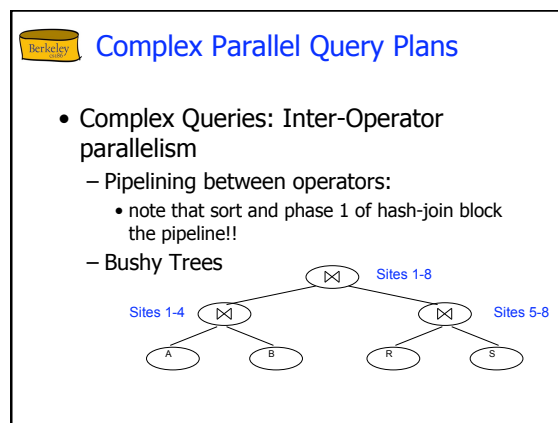
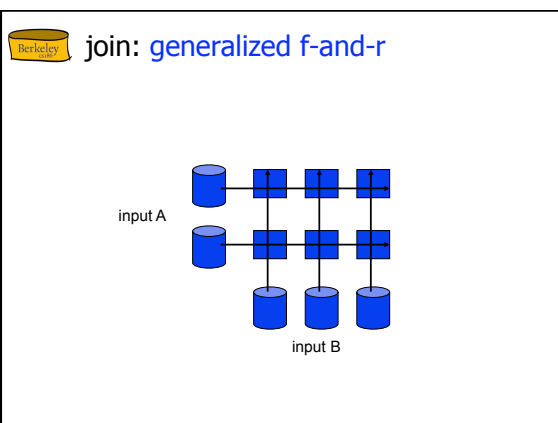
Joins: Bigger picture

- alternatives:
 - symmetric repartitioning
 - what we did so far
 - asymmetric repartitioning
 - fragment and replicate
 - generalized f-and-r

join: symmetric repartitioning

join: asymmetric repartitioning

join: fragment and replicate





Parallel DBMS Summary, cont.

- Hardest part of the equation: query optimization.
 - Wait for it!
- We haven't said anything about Xacts, logging.
 - Familiar in shared-memory architecture.
 - Takes some care in shared-nothing.
 - Yet more tricky in shared-disk