

CS 640
Principles of Database Management and Use
Winter 2013

Today: Discussion about Parallel DBSs

D. DeWitt and J. Gray: **Parallel Database Systems: The Future of High Performance Database Systems**. *CACM* 36(6), 1992.

These lecture nodes are based on Ramakrishnan and Gehrke's textbook "Database Management Systems"

What are the three main architectures for parallel DBMSs and what are their respective characteristics?

Shared memory:

- Communication overhead is low (memory can be used)
- Larger number of processors: memory contention becomes bottleneck
- Interference

Shared disk:

- Large amounts of data may need to be shipped over the interconnect
- Interference

Shared nothing:

- More complex to implement
- Near-linear scaleup, near-linear speedup

2

What is *pipeline parallelism* and what is *data-partitioned parallelism*?

For which operations are those approaches most useful, respectively?

- Pipelined parallelism most useful for non-blocking operators
- Data-partitioned parallelism most useful for operators that do not need to combine data across partitions (e.g., scan, selection, projection)

3

What are the respective advantages and disadvantages of each of the three basic data partitioning schemes,

- *round-robin*,
- *hashing*, and
- *range partitioning*?

- Round-robin suitable for queries that access entire relation
- Range partitioning superior for queries with range selections
- Hashing keeps data evenly distributed
- Problem for hashing and for range partitioning: data skew

4

How can we achieve data-partitioned parallelism with conventional physical operators?

- New physical operators: split and merge

5

How does the parallel hash join algorithm work?

6