

Parallel Databases

Solutions to Practice Exercises

21.1 If there are few tuples in the queried range, then each query can be processed quickly on a single disk. This allows parallel execution of queries with reduced overhead of initiating queries on multiple disks.

On the other hand, if there are many tuples in the queried range, each query takes a long time to execute as there is no parallelism within its execution. Also, some of the disks can become hot-spots, further increasing response time.

Hybrid range partitioning, in which small ranges (a few blocks each) are partitioned in a round-robin fashion, provides the benefits of range partitioning without its drawbacks.

21.2 a. When there are many small queries, inter-query parallelism gives good throughput. Parallelizing each of these small queries would increase the initiation overhead, without any significant reduction in response time.

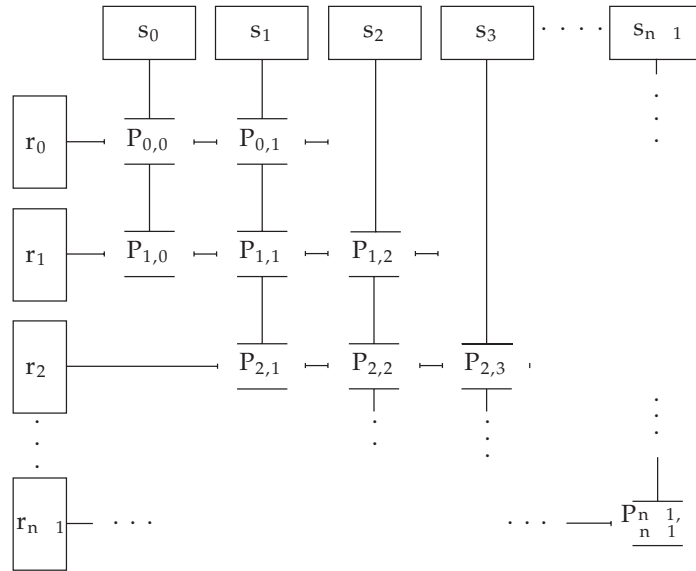
b. With a few large queries, intra-query parallelism is essential to get fast response times. Given that there are large number of processors and disks, only intra-operation parallelism can take advantage of the parallel hardware – for queries typically have few operations, but each one needs to process a large number of tuples.

21.3 a. The speed-up obtained by parallelizing the operations would be offset by the data transfer overhead, as each tuple produced by an operator would have to be transferred to its consumer, which is running on a different processor.

b. In a shared-memory architecture, transferring the tuples is very efficient. So the above argument does not hold to any significant degree.

- c. Even if two operations are independent, it may be that they both supply their outputs to a common third operator. In that case, running all three on the same processor may be better than transferring tuples across processors.

21.4 Relation r is partitioned into n partitions, r_0, r_1, \dots, r_{n-1} , and s is also partitioned into n partitions, s_0, s_1, \dots, s_{n-1} . The partitions are replicated and assigned to processors as shown below.



Each fragment is replicated on 3 processors only, unlike in the general case where it is replicated on n processors. The number of processors required is now approximately $3n$, instead of n^2 in the general case. Therefore given the same number of processors, we can partition the relations into more fragments with this optimization, thus making each local join faster.

- 21.5 a. A partitioning vector which gives 5 partitions with 20 tuples in each partition is: [21, 31, 51, 76]. The 5 partitions obtained are 1 – 20, 21 – 30, 31 – 50, 51 – 75 and 76 – 100. The assumption made in arriving at this partitioning vector is that within a histogram range, each value is equally likely.
- b. Let the histogram ranges be called h_1, h_2, \dots, h_h , and the partitions p_1, p_2, \dots, p_p . Let the frequencies of the histogram ranges be n_1, n_2, \dots, n_h . Each partition should contain N/p tuples, where $N = \sum_{i=1}^h n_i$.

To construct the load balanced partitioning vector, we need to determine the value of the k_1^{th} tuple, the value of the k_2^{th} tuple and so on, where $k_1 = N/p, k_2 = 2N/p$ etc, until k_{p-1} . The partitioning vector will then be $[k_1, k_2, \dots, k_{p-1}]$. The value of the k_i^{th} tuple is determined as follows. First determine the histogram range h_j in which it falls. Assuming all values in

a range are equally likely, the k_i^{th} value will be

$$s_j + (e_j - s_j) * \frac{k_{ij}}{n_j}$$

where

s_j : first value in h_j

e_j : last value in h_j

k_{ij} : $k_i - \sum_{l=1}^{j-1} n_l$

- 21.6 a.** The copies of the data items at a processor should be partitioned across multiple other processors, rather than stored in a single processor, for the following reasons:
- to better distribute the work which should have been done by the failed processor, among the remaining processors.
 - Even when there is no failure, this technique can to some extent deal with hot-spots created by read only transactions.
- b.** RAID level 0 itself stores an extra copy of each data item (mirroring). Thus this is similar to mirroring performed by the database itself, except that the database system does not have to bother about the details of performing the mirroring. It just issues the write to the RAID system, which automatically performs the mirroring.
- RAID level 5 is less expensive than mirroring in terms of disk space requirement, but writes are more expensive, and rebuilding a crashed disk is more expensive.