

Solution 22.8

Answer 1. In the following, when we refer to the "highest paid employee/manager", note there could be duplicates, and the queries will return all such tuples that qualify. Also, we assume there is no cost involved in shipping a query from the site it is posted at, or the site(s) where it is evaluated.

- (a) We need to do a *total exchange* of the Departments partitions. In other words, each site must ship its partition to *each* of the other nine sites. Each of the Departments partition is stored in 500 pages. For each site, the cost of shipping a partition to the other nine sites is $500 \times 9t_s$. Hence, the cost of the *total exchange* is $45,000t_s$. The cost of a local natural join (using Sort-Merge Join) is $3(M + N)t_d$, where $M = 10,000$ Departments records, and $N = 2 \times 10^6$ Employees records. (Remember the Employees relation is also partitioned, so each site has 10,000 pages, each with 200 records.) Hence the cost of the natural join across all the sites is $603 \times 10^5 t_d$. Finally, the result fragments have to be shipped to the query site. Each site will have a result fragment of 500 pages (the same size as the local Departments fragment. Why?), and the total cost of shipping these fragments from nine sites to the query site is $4500t_s$. Adding all the costs together, the total cost of the plan is $49,500t_s + 60,300,000t_d$.
- (b) Clearly, the highest paid employee will have a salary in the bracket $900,000 < sal \leq 1,000,000$. (We are assured that there is at least one employee in this salary range because the *sal* field has values uniformly distributed in the range 0 to 1,000,000, and there are 100,000 pages of Employees tuples). So we ship the query to the tenth site, evaluate the highest paid employee at that site, and ship the result back to the query site. Finding the employee(s) with the maximum salary at the tenth site requires a scan of the Employees fragment, which costs $10,000t_d$. Since the salary is uniformly distributed, we may assume there are $\frac{200 \times 10,000}{100,000} = 20$ result tuples, which fit on a single page. Thus total cost of the query is $10,000t_d + t_s$.
- (c) The best plan for this query results from the observation that the partition $sal \leq 100,000$ is replicated at every site. So we just need to query the partition at the *query site* for the highest paid employee in its partition. This does not involve any shipping costs. Thus, the cost of the query is the cost of scanning the Employees partition, which is $10,000t_d$.
- (d) Again, we will find the highest paid employee in the salary bracket $400,000 < sal \leq 500,000$ in the fifth site. So we ship the query to the fifth site, evaluate the query there, and ship the result to the query site. The cost is the same as in part (b) above.
- (e) Employees in the salary bracket $450,000 < sal \leq 550,000$ are split across the fifth and the sixth sites. But clearly, by the partitioning of the Employees relation, the employees in the sixth site earn more than those at the
-

fifth site. So we need to query only the partition in the sixth site, but we need to qualify our query by selecting only those employees whose salary is less than 550,000, and then selecting the one with the highest salary. The cost is again the same as in part (b) above.

- (f) We need to do a partial join of the two relations. The Departments fragment at the query site is shipped across to each of the other nine sites, and the fragment is joined with the Employees fragment at the site. (Why would SemiJoin or BloomJoin at the query site not buy us much in this case? Hint: The *eid* field is a key for Employees and the join is on this field.) At each site, the result is scanned to select the tuple with the maximum salary, and this tuple is shipped to the query site, where a final scan of the ten tuples contributed by each site yields the department (amongst the ones stored at the query site) with the highest paid manager. The cost of the initial shipping is $4500t_s$, the cost of ten local joins is $603 \times 10^5 t_d$ (refer to part (a)), the cost of the scans is , the cost of shipping the result tuples is $9t_s$, and the cost of the final scan is $10t_d$.
 - (g) We do the above operation for each fragment of the Departments relation. Thus, the initial shipping costs, and the local join costs increase ten-fold. Everything else remains the same. Thus, the cost of the query is .
2. (a) We use the **read-only write-all** policy for synchronous replication here. Since every site has a copy of the fragment, we can raise the salary of all those with salary less than 100,000 without having to ship any part of the relation. So each site simply raises the salary of all employees in the partition by 10 have a salary greater than 100,000, and we just reset their salary to exactly 100,000. In order to do this, each site needs to obtain an exclusive lock on its local copy of the fragment, and the query is executed at each site.
- (b) We again use the **read-only write-all** policy, and we need to obtain an exclusive lock on all the copies of the fragment with employees whose salary is less than 100,000. The remaining fragments of the relation are not replicated, and hence, just one exclusive lock for each fragment is required, at the site where it is stored. Once the update is complete, we need to do some shipping to satisfy the original partitioning of Employees. Each site will ship to the "next" site employees whose salaries have been raised to the next bracket of the original partitioning. In other words, the first site will ship all records will salaries greater than 100,000 to the second site, the second site will ship all records will salaries greater than 200,000 to the third site, and so on. To do this, each site needs to get an exclusive lock on the fragment they store. Once each site has modified its relation, shipped the appropriate records, and has received its set of additional records, it obtains

a lock on its fragment to append the new records. The only tricky fragment to handle is the one that has the employees with the lowest salaries, since it is replicated at all the sites. Since we are using synchronous replication, the changes made at the first site must be executed at all the sites. Hence, each site must obtain an exclusive lock on its copy of the fragment, and run the query against it.

3. (a) Using *primary site* asynchronous replication, only the site with the primary copy of the fragment whose employees earn less than 100,000 is visited. None of the other copies need to be locked. The primary site is predetermined, and in our case, it is possibly the first site. The changes to the primary copy will be captured at a latter time, and applied to the secondary copies.
- (b) We again use *primary site* asynchronous replication. We follow a procedure very similar to the one described when synchronous replication is used. The only difference is in the handling of the fragment which is replicated at all the sites. Since we are using asynchronous replication here, we can just lock and change the primary copy, and copy the changes to the other sites at some later time.