

Solution of CIS 670 Assignment Two

1. *Explain the difference between data replication in a distributed system to the maintenance of a remote backup site? Compare the advantages and disadvantages of the data replication and data fragmentation approaches in a distributed system. When is it useful to have replication or fragmentation of data?*

Answer:

Data replication provides greater availability and also parallelism as multiple sites are able to cater to more transactions. But there is a reduction in data transfer as more systems have a local copy. Also, there is an update overhead as all replicas need to remain consistent.

Remote Backup systems carry out actions such as concurrency control and recovery at a single site. Also remote backup systems help avoid two phase commits and all the overhead resulting from it. Since transactions contact only one site, the overhead of running transactions at multiple sites are avoided. Therefore remote backup sites provide a lower cost approach.

Replication is useful when the data is required at multiple sites in a distributed system. Fragmentation is useful, as the data is made available only at sites where it is needed and will be useful, thus reduces redundancy. Both replication and fragmentation facilitate parallelism.

2. *Consider the relations:*

employee (name, address, salary, plant number)

machine (machine, number, type, plant number)

Assume that the employee relation is fragmented horizontally by plant number, and that each fragment is stored locally at its corresponding plant site. Assume that the machine relation is stored in its entirety at the Armonk site. Describe a good strategy for processing each of the following queries.

a. Find all employees at the plant that contains machine number 1130.

b. Find all employees at plants that contain machines whose type is “milling machine.”

c. Find all machines at the Almaden plant.

d. Find employee \bowtie machine.

Answer:

- a.
 - i. Perform $\Pi_{\text{plant number}} (\sigma_{\text{machine number}=1130} (\text{machine}))$ at Armonk.
 - ii. Send the query $\Pi_{\text{name}} (\text{employee})$ to all site(s) which are in the result of the previous query.
 - iii. Those sites compute the answers.
 - iv. Union the answers at the destination site.
- b. This strategy is the same as a), except the first step should be to perform $\Pi_{\text{plant number}} (\sigma_{\text{type}=\text{"milling machine"}} (\text{machine}))$ at Armonk.
- c.
 - i. Perform $\sigma_{\text{plant number} = x} (\text{machine})$ at Armonk, where x is the plant number for Almaden.
 - ii. Send the answers to the destination site.
- d. Strategy 1:
 - i. Group *machine* at Armonk by plant number.
 - ii. Send the groups to the sites with the corresponding plant number.
 - iii. Perform a local join between the local data and the received data.
 - iv. Union the results at the destination site.

Strategy 2:

Send the machine relation at Armonk, and all the fragments of the employee relation to the destination site. Then perform the join at the destination site.

There is parallelism in the join computation according to the first strategy but not in the second. Nevertheless, in a WAN the amount of data to be shipped is the main cost factor. We expect that each plant will have more than one machine, hence the result of the local join at each site will be a cross-product of the employee tuples and machines at that plant. This cross-product's size is greater than the size of the employee fragment at that site. As a result the second strategy will result in less data shipping, and will be more efficient.

3. a) Depends on the site which the query was entered. b) Depends on the site which the result was desired.

Answer: a) The cost of transferring the query itself is much less expensive than cost of transferring the data. The site at which the query was entered only needs to transfer the query to each plant_number site, which is not very expensive. Hence, our strategy does not depend on the site at which query was entered.

b) For the first query, we compute the plant_number at Armonk site and then compute the employee tuple at each site locally which are then transferred to the destination site. Our strategy is relatively independent of the site at which the results is desired.

For the second query, we again compute the plant_number at Armonk site and then compute employee tuple at each site locally (based on plant_number result) which are then transferred to destination site. Again, this strategy is relatively independent of the site at which the result is desired.

For the third query, the query has to be computed at Armonk site and then result is send to destination site, which means that this strategy is independent of the site at which the result is desired. Only if Armonk is the site the result is desired, we don't need to transfer the results.

For the fourth query, both possible strategies involve the data (or result) being migrated from each local site to the destination site regardless of the choice of destination. Hence, again this strategy is independent of the site at which the result is desired.

4. *What is the difference and relationship between distributed database and cloud database?*

Answer:

Distributed databases store information separately in each of local databases. When the query is issued by a user, the query is distributed into local databases and the answers are combined. Cloud database is based on cloud infrastructure which is supported by the large companies (e.g., Amazon, Google) or other organizations as vendors. It is naturally that the cloud database use NoSQL database (e.g., key-value stores) to store and manage the data while distributed database use relational/SQL database.

Therefore, cloud database is maintained by the cloud vendors, the availability and scalability are guaranteed. However, the customers do not have full control for the data management. The organization which built the distributed database maintain the database and have full control on the database, but the availability and scalability are not guaranteed as it relies on each distributed site and the SQL database.

