

# Case Study: Query Cost Calculation

---

## 1. Introduction

In a distributed database management system (DDBMS), query processing involves translating a high-level query into an efficient sequence of low-level operations to retrieve desired results. A crucial aspect of query optimization is query cost calculation, which determines the most efficient strategy to execute a query based on the system's resource consumption. The cost of a query is measured in terms of I/O cost, CPU cost, communication cost, and response time.

This case study explores the query cost calculation process, presenting key components, a step-by-step breakdown of cost estimation, and a real-world scenario where query cost affects performance in distributed systems.

## 2. Key Components of Query Cost Calculation

### a. I/O Cost:

- Refers to the cost incurred when reading and writing data blocks from/to disk. It's dependent on the number of blocks read, size of the relations (tables), and whether the data is indexed. I/O cost is a primary concern in databases as reading from disk is usually slower than accessing data from memory.

### b. CPU Cost:

- The cost associated with the CPU's time to process the query (e.g., parsing, compiling, and executing the query). Affected by the number of tuples (rows) involved, operations performed (e.g., joins, aggregations), and predicate evaluation.

### c. Communication Cost:

- Relevant in distributed databases where data needs to be transferred across nodes. Includes data transfer time, which can be significant, especially in geographically dispersed databases.

### d. Memory Usage:

- Cost related to memory usage, including buffer management and temporary storage for intermediate results.

### e. Response Time:

- The time taken by the system to return the query result to the user. This is influenced by I/O and CPU costs, as well as network latency in distributed systems.

### 3. Cost Components in Query Execution Plans

Each query can have multiple execution strategies, and the DBMS uses a query optimizer to select the best plan based on cost estimates. The query execution plan is evaluated using various methods such as:

- Full table scans
- Index scans
- Nested-loop joins
- Merge joins
- Hash joins

Each method has a different cost profile depending on the database size, indexes, and data distribution.

### 4. Real-World Scenario: Distributed Database Query Optimization

Scenario: A large e-commerce company stores customer and transaction data in multiple distributed databases across different geographic locations. Consider a query to retrieve all orders from the 'Orders' table where the total amount exceeds a specified value, with related customer data from the 'Customers' table.

SQL Query Example:

```
SELECT c.customer_name, o.order_id, o.total_amount
FROM Customers c, Orders o
WHERE c.customer_id = o.customer_id
AND o.total_amount > 1000;
```

### 5. Query Cost Estimation Breakdown

Step 1: Parsing the Query

The SQL query is parsed into an intermediate representation (query tree).

Step 2: Logical Plan

The logical plan specifies the relations and joins, and filters required to produce the desired results. No cost is associated with this phase yet, but the plan determines the potential operations needed.

Step 3: Selection of Access Paths

Two main access paths can be used for this query: full table scan or index scan on the 'total\_amount' column (if indexed).

Step 4: Join Strategy Selection

Depending on the size of the data and indexes, various join methods are evaluated. Cost varies based on whether the join is done via nested-loop, hash join, or merge join. A nested-loop join may be more costly for large data sets, while a hash join could be more optimal if hash indexes exist.

#### Step 5: Calculation of I/O and CPU Costs

For the Orders table with 10 million records, a full scan might cost around 10,000 disk reads, while an index scan could reduce this to 500 reads. CPU cost is evaluated based on the number of tuples that pass through filters and need to be joined.

#### Step 6: Distributed System Consideration

The Orders and Customers tables may be distributed across different data centers. Communication cost is a major factor since customer data from one location needs to be joined with orders from another. This requires transferring large datasets between nodes, potentially increasing network latency and overall response time.

## 6. Cost Calculation Example

Assume:

- Orders Table (10 million rows) and Customers Table (1 million rows).
- I/O cost per block read is 1ms, CPU cost per tuple processing is 0.1ms, and network transfer rate is 1000 rows/sec.

Full table scan of Orders:

- 10,000 disk reads = 10,000ms (I/O Cost)
- Processing 10 million rows = 1 million ms (CPU Cost)

Hash Join Strategy:

- Index scan of Orders reduces I/O to 500 reads = 500ms
- Join with Customers using hash join: CPU cost for hash computation and matching ~200ms.

Communication Cost:

- Transferring 1 million rows from Customers node to Orders node: 1000 seconds = 1,000,000ms.

Total Cost Estimate:

- I/O Cost: 500ms
- CPU Cost: 1 million ms
- Communication Cost: 1,000,000ms
- Total Estimated Cost = 2,000,500ms (2 seconds)

## 7. Conclusion

Query cost calculation is a critical process in optimizing the execution of database queries. By considering various components like I/O, CPU, and communication costs, query optimizers can select the most efficient query execution plan. In distributed systems, minimizing communication cost becomes particularly important to enhance query performance and reduce response time. In the given scenario, the combination of index scans and hash joins helped reduce the overall cost, making the query execution more efficient.

## 8. References

- Silberschatz, A., Korth, H.F., & Sudarshan, S. (2020). Database System Concepts.
- Elmasri, R., & Navathe, S.B. (2016). Fundamentals of Database Systems.
- Ramakrishnan, R., & Gehrke, J. (2003). Database Management Systems.