

Column-Stores vs. Row-Stores

How Different are they Really?

Daniel J. Abadi, Samuel Madden and Nabil Hachem, SIGMOD
2008

Presented by: Mehvish Saleem

Contents

- Structural difference between column and row store
- Contributions
 - Emulation of column-store in row store
 - Column-oriented optimizations esp. Invisible join
 - Column-store performance contributors
- Experiments and Results
- Conclusion

What is Column-Store?



- Row-store stores data in the disk tuple by tuple
- Column-store stores data in the disk column by column

Emulation of column-store in row-store

- Vertical partitioning
- Using index-only plans
- Using materialized views

Vertical Partitioning

- Full Vertical partitioning of each relation
 - Each column = 1 Physical table
 - Can be achieved by adding integer position column to every table
- Join on position for multi-column fetch
- Problems:
 - “Position” - Space and disk bandwidth
 - Header for every tuple – further space wastage

Index-only plans

- Add B+Tree index for every Table column
- Plans never access the actual tuples on disk
- Headers are not stored, so per tuple overhead is less
- Problems:
 - Separate indices may require full index scan, which is slower
 - E.g.: `SELECT AVG(salary) FROM emp WHERE age > 40`

Materialized Views

- Create 'optimal' set of MVs for given query workload
- Objective:
 - Provide just the required data
 - Avoid overheads
 - Performs better
- Expected to perform better than other two approach
- Problems:
 - Require knowledge of query workloads in advance

Column-oriented execution

- Compression
- Late Materialization
- Block Iteration
- Invisible join

Compression

- Perform better on data with low information entropy
- Column-store data super compressible
- Improves performance
 - Less space
 - Less time spent in I/O

Late Materialization

- Delay Tuple Construction
- Might avoid constructing it altogether
- Eg: `SELECT R.a FROM R WHERE R.c = 5 AND R.b = 10`
 - Output of each predicate is a bit string
 - Perform Bitwise AND
 - Use final position list to extract R.a

Block Iteration

- Iterate over blocks of tuples rather than tuple-at-a-time
- Like batch processing
- If column is fixed width, it can be operated as an array
- Exploits potential for parallelism

Invisible Join

Star Schema Benchmark

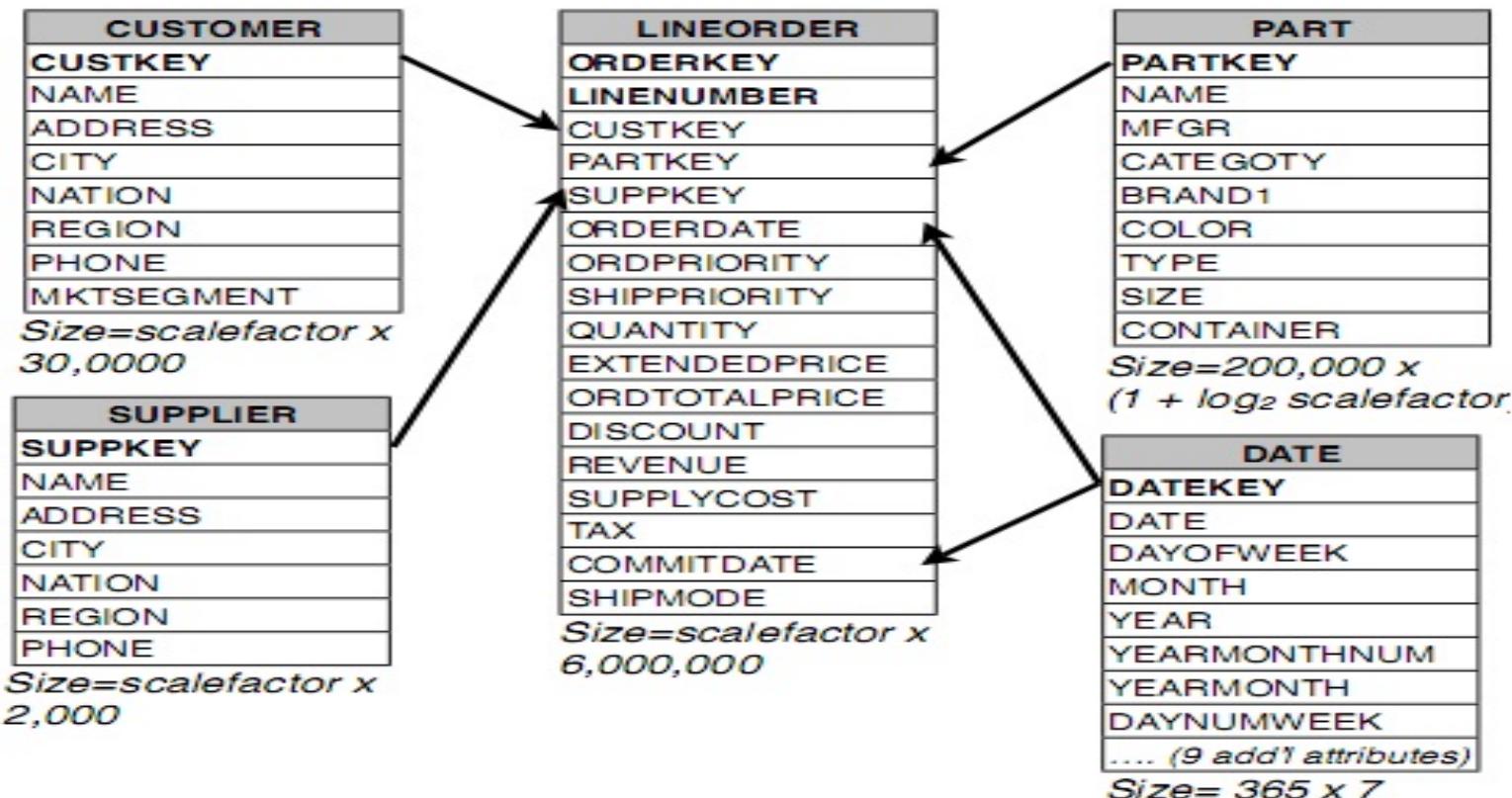


Figure 1: Schema of the SSBM Benchmark

Sample Query

```
SELECT c.nation, s.nation, d.year,  
       sum(lo.revenue) as revenue  
FROM customer AS c, lineorder AS lo,  
      supplier AS s, dwdate AS d  
WHERE lo.custkey = c.custkey  
  AND lo.supkey = s.supkey  
  AND lo.orderdate = d.datekey  
  AND c.region = ASIA  
  AND s.region = ASIA  
  AND d.year >= 1992 and d.year <= 1997  
GROUP BY c.nation, s.nation, d.year  
ORDER BY d.year asc, revenue desc;
```

Find total revenue from Asian customers who purchase a product supplied by an Asian supplier between 1992 and 1997 grouped by nation of the customer, supplier and year of transaction

Phase I

Apply region = 'Asia' on Customer table

custkey	region	nation	...
1	Asia	China	...
2	Europe	France	...
3	Asia	India	...

Hash table
with keys
1 and 3

Apply region = 'Asia' on Supplier table

suppkey	region	nation	...
1	Asia	Russia	...
2	Europe	Spain	...

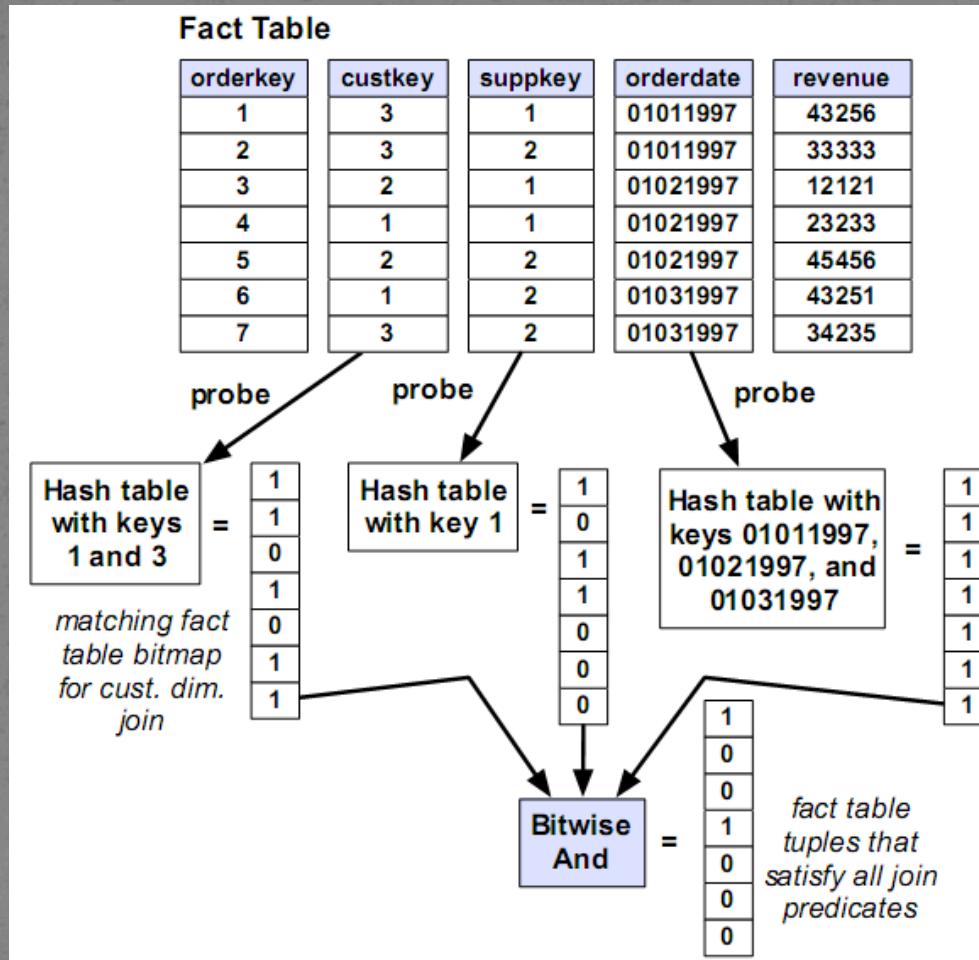
Hash table
with key 1

Apply year in [1992, 1997] on Date table

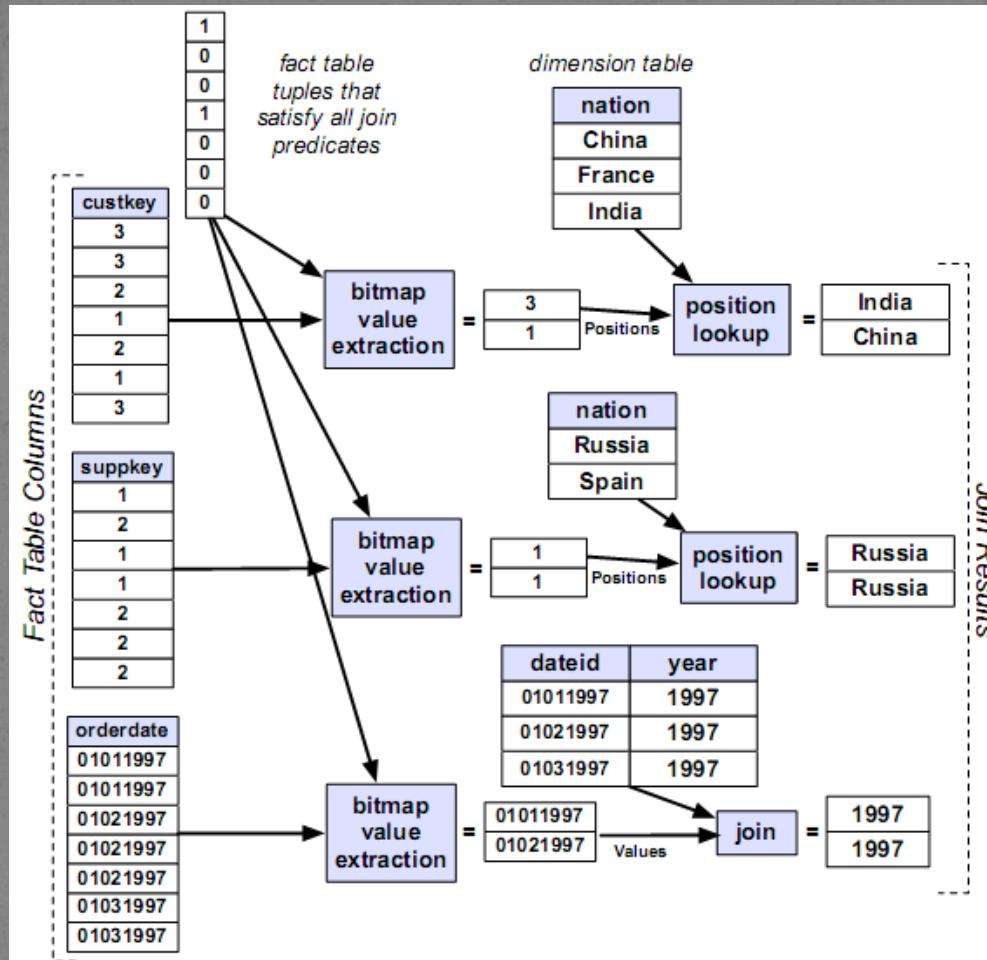
dateid	year	...
01011997	1997	...
01021997	1997	...
01031997	1997	...

Hash table with
keys 01011997,
01021997, and
01031997

Phase II



Phase III

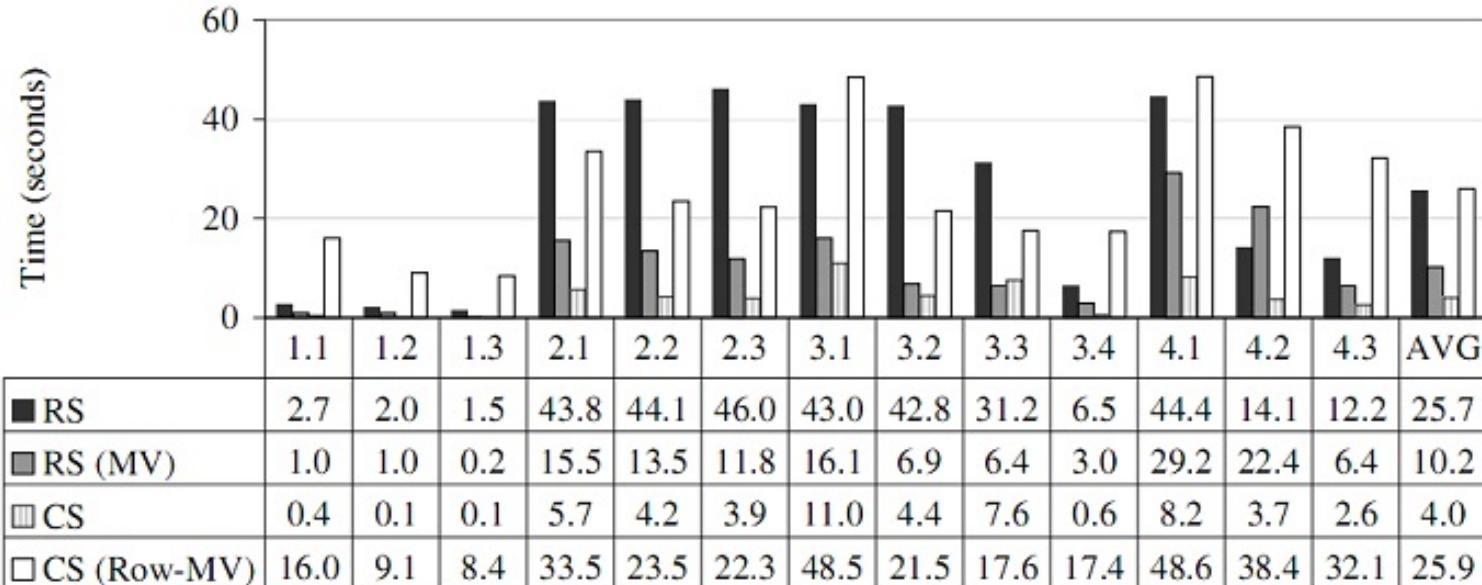


Experiments

Purpose

- Comparing the performance of C-Store with column-store emulation in row-store
- Most significant optimization for column-stores
- Can unmodified row store obtain benefits of column-store?

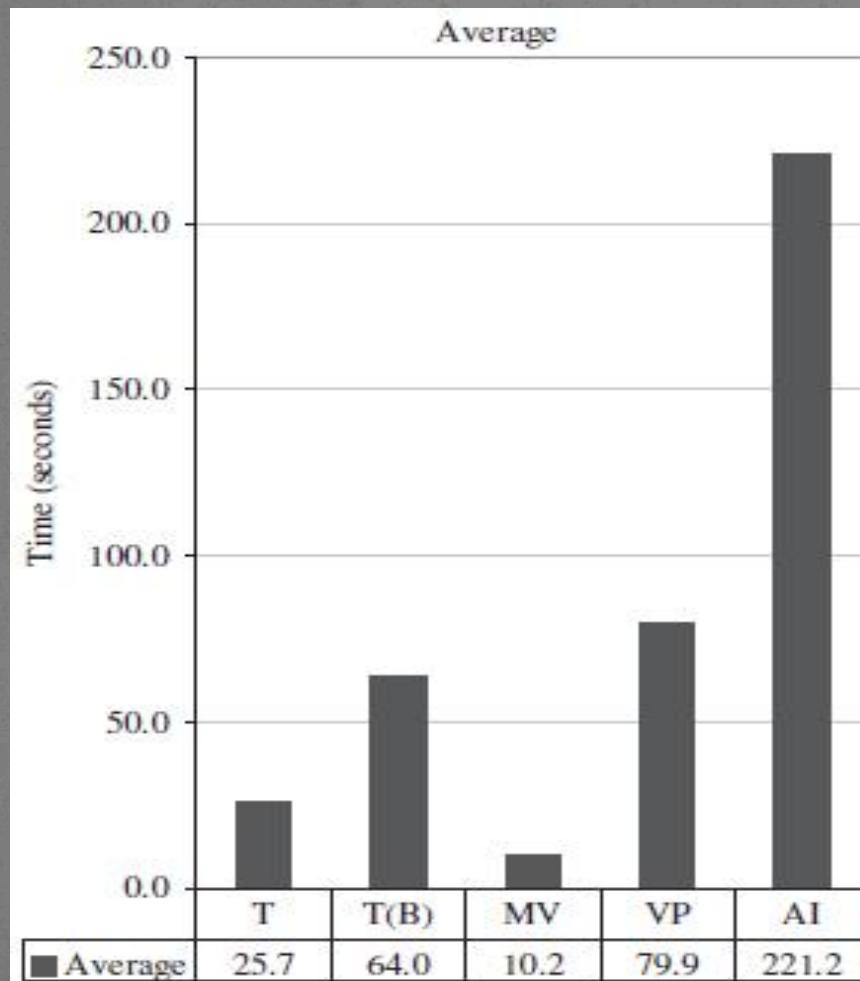
C-Store VS System X



Baseline performance of C-Store “CS” and System X “RS”, compared with materialized view cases on the same systems.

- RS: Base System X
- CS: Base C-Store
- RS (MV): System X with optimal collection of MVs
- CS (Row-MV): Column store constructed from RS(MV)

Different System X Configurations



- T: Traditional
- T(B): Traditional Bitmap
- MV: Materialized View
- VP: Vertical Partitioning
- AI: All Indexes

Results and Analysis

- MV performs best since they read minimal amount of data needed by a query
- Index only plans are the worst:
 - Expensive column joins on fact table
- Vertical partitioning:
 - Tuple overheads and reconstruction
- Column Store perform better than the best case of row store (4.0 sec VS 10.2 sec)

Conclusion

- Emulation of a column-store in a row-store does not yield good results
 - Tuple reconstruction costs
 - A pain to implement
- Most important optimizations for column-store:
Compression and Late Materialization

Thank you!
