# Summary

| | |
|---|---|
| 👤 | Architecture of Database Systems |
| 🔧 | Transaction Management |
| 🔗 | Modern Database Technology |
| 🗄 | Data Warehouses and OLAP |
| ⛏ | **Data Mining** |
| 📊 | Big Data Analytics |

## Time Series Forecasting Models

Domain Specific Extensions — White-Box
Base Forecast Models

EGRV

(Auto)Regression — Gray-Box
(S)AR(I)MA(X)
VAR

Exponential Smoothing

Machine Learning — Black-Box

$$\begin{pmatrix} \hat{x}_{1,t+1} \\ \vdots \\ \hat{x}_{n,t+1} \end{pmatrix} = \begin{pmatrix} \alpha_{1,1} & \cdots & \alpha_{n,1} \\ \vdots & \ddots & \vdots \\ \alpha_{1,n} & \cdots & \alpha_{n,n} \end{pmatrix} \cdot \begin{pmatrix} x_{1,t} \\ \vdots \\ x_{n,t} \end{pmatrix}$$

Auto Regression (AR) → ARMA ← Moving Average (MA)

Exogenous Influence → ARIMAX
Seasonality → SARIMA

UK Energy Demand – Month 11

## Forecasting Process & Time Series Storage

| Method | Error ε | AIC |
|---|---|---|
| ARIMA | 23.47 | 25.72 |
| HoltWinters | 18.38 | 20.43 |
| VAR | | |
| MARS | | |

Optimize, Predict, Evaluate

### Relational Storage

| ID | rec_date | sales |
|---|---|---|
| 1 | 20.11.2021 | 3 |
| 1 | 21.11.2021 | 7 |
| 1 | 22.11.2021 | 4 |

### Array Storage

| ID | start | interv | sales |
|---|---|---|---|
| 1 | 20.11.2021 | 1 day | {3,7,4,5,3,2,8,7,4,6} |
| 2 | 20.11.2021 | 1 day | {7,5,8,9,1,5,5,8,9,6} |

## Dependency Mining

Support

$$s(X \rightarrow Y) = s(X \cup Y) = \frac{|\{T_k | X \cup Y \subseteq T_k\}|}{|D|} = p(XY)$$

Confidence

$$c(X \rightarrow Y) = \frac{s(X \cup Y)}{s(X)} = \frac{|\{T_k | X \cup Y \subseteq T_k\}|}{|\{T_k | X \subseteq T_k\}|} = p(Y|X)$$

**Strong rule: High Support + High Confidence**

| $I_k^1$ | # | $s(I_k^1)$ |
|---|---|---|
| {chips} | 3 | 0.75 |
| {cola} | 2 | 0.5 |
| {peanuts} | 1 | 0.25 |
| {beer} | 3 | 0.75 |
| {cigarettes} | 3 | 0.75 |

| $I_k^2$ | # | $s(I_k^2)$ |
|---|---|---|
| {chips, cola} | 2 | 0.5 |
| {beer, chips} | 2 | 0.5 |
| {chips, cigarettes} | 2 | 0.5 |
| {beer, cola} | 1 | 0.25 |
| {cigarettes, cola} | 1 | 0.25 |
| {beer, cigarettes} | 3 | 0.75 |

| $I_k^2$ | # | $s(I_k^2)$ |
|---|---|---|
| {chips, cola} | 2 | 0.5 |
| {beer, chips} | 2 | 0.5 |
| {chips, cigarettes} | 2 | 0.5 |
| {beer, cola} | 1 | 0.25 |
| {cigarettes, cola} | 1 | 0.25 |
| {beer, cigarettes} | 3 | 0.75 |

| $I_k^3$ | # | $s(I_k^3)$ |
|---|---|---|
| {beer, chips, cigar.} | 2 | 0.5 |
| {chips, cigar., cola} | 1 | 0.25 |

*Prune because of non-frequent subset*

# Course Outline

| | |
|---|---|
| | Architecture of Database Systems |
| | Transaction Management |
| | Modern Database Technology |
| | Data Warehouses and OLAP |
| | Data Mining |
| | Big Data Analytics |

**The Big Data Vs**

Volume
- Terabytes
- Records
- Transactions
- Tables, Files

Variety
- Structured
- Unstructured
- Multi-factor
- Probabilistic

Velocity
- Batch
- Real/near-time
- Processes
- Streams

→ *see lectures 10-15*

→ **New Data Management Methods and Systems (ACID→BASE, NoSQL,…)**

→ **Parallelization to meet performance requirements** — Today

→ **Complexities of scaling hidden by data processing frameworks** — Wednesday

→ *see lecture 10*

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

Big Data Analytics

4

**More Vs:**
- *Veracity*: Is your data (source) trustworthy/ meaningful?
- *Visualization*: How to communicate? insights& knowledge?
- *Value*: How to use data for (machine) learning, optimization, …?
- (Volatility, Vulnerability, Validity, …)

# Recap Query Processing

**SQL Query (What!)**

```
SELECT A_Name, SUM(S_Qty)
FROM Article, Sales
WHERE A_Anr = S_Anr AND
      S_Date >= '2021-01-01'
GROUP BY A_Name
```

YES, but HOW do we get there (efficiently)?

### Result

| A_Name | SUM |
|--------|-----|
| Article A | 12 |
| Acticle B | 1 |

$\gamma_{SUM(S\_Qty); A\_Name}$

$\sigma_{S\_Date >= '2021-01-01'}$

$\bowtie_{S\_Anr=A\_Anr}$

Sales     Article

*Query Plan (How!)*

### Sales

| ... | S_Anr | S_Date | S_Qty |
|-----|-------|--------|-------|
| ... | 1 | 2020-09-20 | 7 |
| ... | 1 | 2021-01-17 | 2 |
| ... | 1 | 2021-02-21 | 5 |
| ... | 2 | 2021-02-22 | 1 |
| ... | 1 | 2021-03-07 | 5 |

### Article

| A_Anr | A_Name | ... |
|-------|--------|-----|
| 1 | Article A | ... |
| 2 | Acticle B | ... |

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

Big Data Analytics

# Query Processing

# Logical Query Plan

## Mono-Block Query

**SELECT** C_NAME, C_ADDRESS
**FROM** TPCD.CUSTOMER, TPCD.SUPPLIER
**WHERE** C_NAME = S_NAME
**AND** C_MKTSEGMENT = 'MACHINERY';

$\pi_{C\_NAME, C\_ADDRESS}$ — Projection

$\sigma_{C\_MKTSEGMENT='MACHINERY'}$ — Selection

$\underset{C\_NAME = S\_NAME}{\bowtie}$ — Join Operation

CUSTOMER   SUPPLIER — Source Relations

## Star Query

**SELECT**
P_BRAND, O_SHIPPRIORITY,
SUM(L_QUANTITY*L_EXTENDEDPRICE)
AS TURNOVER
**FROM**
TPCD.LINEITEM, TPCD.ORDERS, TPCD.PART
**WHERE** L_ORDERKEY = O_ORDERKEY
**AND** L_PARTKEY = P_PARTKEY
**AND** O_ORDERSTATUS = 'F'
**AND** P_CONTAINER = 'LG_BAG'
**GROUP BY**
P_BRAND, O_SHIPPRIORITY
**HAVING**
AVG(L_QUANTITY) > 250;

$\pi_{P\_BRAND, O\_SHIPPRIORITY, TURNOVER}$

$\sigma_{\$X > 250}$

$\gamma_{P\_BRAND, O\_SHIPPRIORITY:}$
TURNOVER = SUM(L_QUANTITY*L_EXTENDEDPRICE),
$\$X$ = AVG(L_QUANTITY)

$\sigma_{O\_ORDERSTATUS='F'\ AND\ P\_CONTAINER='LG\_BAG'}$

$\underset{L\_PARTKEY = P\_PARTKEY}{\bowtie}$

$\underset{L\_ORDERKEY = O\_ORDERKEY}{\bowtie}$

ORDERS   LINEITEM   PART

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

Big Data Analytics

7

# Query Optimization – Example

**Logical Query Plan**

$\gamma_{SUM(S\_Qty); A\_Name}$
|
$\sigma_{S\_Date >= \text{‘2021-01-01’}}$
|
$\bowtie_{S\_Anr=A\_Anr}$

Sales    Article

**Logical Plan Rewrite**

Optimized Logical Query Plan

$\gamma_{SUM(S\_Qty); A\_Name}$
|
$\bowtie_{S\_Anr=A\_Anr}$

$\sigma_{S\_Date >= \text{‘2021-01-01’}}$    Article
|
Sales

**Plan Transformation**

Query Execution Plan

$\text{HSGRPBY}_{SUM(S\_Qty); A\_Name}$
|
$\text{NLJOIN}_{S\_Anr=A\_Anr}$
(index nested loop join)

FETCH    IX:Article

RIDSCN    TBL:Sales
|
$\text{SORT}_{RID<p,i>}$
|
$\text{IXSCAN}_{S\_Date >= \text{‘2021-01-01’}}$
|
IX:Sales

Logical Operation
Physical Operation

Big Data Analytics

8

# Parallel Execution Opportunities

| Query Perspective | Operator Perspective |
|---|---|
| **Inter-query parallelism → Important for OLTP** | **Inter-operator parallelism** |
| <ul><li>Parallel execution of multiple queries</li><li>Goal: increase throughput (e.g., 100.000 queries / second better than 10.000 queries / second)</li></ul> | <ul><li>Parallel execution of multiple operators</li><li>Goal: increase operator throughput, but limited through operator dependencies<br>*→ See lectures 7 and 8*</li></ul> |
| **Intra-query parallelism → Important for OLAP** | **Intra—operator parallelism** |
| <ul><li>Parallel execution of a single query</li><li>Goal: decrease response time (10s better than 100s runtime)</li></ul> | <ul><li>Parallel execution of a single operator</li><li>Goal: decrease operator run time</li></ul> |

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

Big Data Analytics

# Scale-Out Architecture and Parallelization

→ *see lecture 8*

**Shared Nothing**
Each node has autonomy over a subset of the data



Partition-parallelism can be used for many relational operators:

- sort
- aggregation
- join
- …

Some are trivial to parallelize: selection, projection

```
SELECT *
FROM orders
WHERE amount>50
```

Compile

U

$\sigma_{amount>50}$          $\sigma_{amount>50}$

**Orders$_1$ on Node$_1$**       **Orders$_2$ on Node$_2$**

| 6/1/17 | 424252 | Couch | $570 |
| 6/1/17 | 256823 | Car | $1123 |
| 6/2/17 | 636353 | Bike | $86 |
| 6/5/17 | 662113 | Chair | $10 |

| 6/7/17 | 121113 | Lamp | $19 |
| 6/9/17 | 887734 | Bike | $59 |
| 6/11/17 | 252111 | Scooter | $18 |
| 6/11/17 | 116458 | Hammer | $8000 |

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

10

Big Data Analytics

**Parallel Operators**

| Selection | Sort | Aggregation |
|---|---|---|

Selection:
σ_p  σ_p
P_1  P_2
partitioned data (on two nodes)
→ Projection works similarly as selection

Sort:
sort  sort — merge sorted runs
Re-partition on sort-key  0-40  41-99
sort  sort — sort local
P_1  P_2
partitioned data

Aggregation:
γ
U — post aggregate
γ  γ — local aggregate
P_1  P_2
partitioned data

11

Big Data Analytics

**Selection**
- Execute selection $\sigma_p$ in parallel on each fragment
  → Filter out data based on some predicate p (e.g., age=30)
- Pruning can be used for parallel selection
  - Query asks for customers with age>40
  - Customer table is range partitioned into two fragments "age<30 and age>=30"
  - Query can prune fragment with age<30

**Sort**
1. Sort locally (e.g., on age)
2. Re-partition on sort-key (aka "shuffling") using range-partitioning
3. Sort re-partitioned data (can use merge-sort!)

**Aggregation**
Main idea
- Aggregate locally
- Then re-partition using N:1 shuffle and post aggregate

- Pictured plan works for SUM, MIN, MAX
  - COUNT to be re-written to use SUM for post-aggregation
  - AVG needs to be re-written to use SUM / COUNT locally
  - COUNT DISTINCT can only eliminate duplicates locally

# Example: Parallel Sort

**Example: Sort-by Age Ascending**

| Name | Age |
|------|-----|
| Adler | 22 |
| Muller | 38 |

Re-partition on sort-key

| Name | Age |
|------|-----|
| Muller | 38 |
| Bishop | 57 |

| Name | Age |
|------|-----|
| Bishop | 57 |
| Muller | 38 |

0-40          41-99

merge sorted runs

sort local

P₁   P₂

partitioned data

| Name | Age |
|------|-----|
| Zeller | 47 |
| Bishop | 57 |

| Name | Age |
|------|-----|
| Adler | 22 |
| Zeller | 47 |

| Name | Age |
|------|-----|
| Zeller | 47 |
| Adler | 22 |

**Data Shuffling (AKA Re-Partitioning)**

Re-partition on sort-key

0-40    41-99

P₁    P₂

partitioned data

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

Big Data Analytics

- Data shuffling is implemented as a separate operator (send-receive)
- Sort operators are un-aware of parallel (distributed) execution

# Data Shuffling: Details

Data Shuffling can use different partitioning strategies (range vs. hash-partitioning, N:M vs. N:1) to re-partition data during query execution

**Range-based N:M**   **Hash-based N:M**   **N:1 (no part. function)**

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

Big Data Analytics

Not only parallel sort uses data shuffling but also join, aggregation, ...

# Example: Parallel Aggregation

**Example: SELECT SUM(Salary) FROM Census GROUP BY Gender**

| Gender | Salary |
|--------|--------|
| m      | 107k   |

re-partition on group-key

| Gender | Salary |
|--------|--------|
| m      | 95k    |

| Gender | Salary |
|--------|--------|
| m      | 38k    |
| m      | 57k    |

γ — post aggregate

γ — local aggregate

P₁   P₂

partitioned data

| Gender | Salary |
|--------|--------|
| f      | 47k    |

| Gender | Salary |
|--------|--------|
| f      | 47k    |
| m      | 22k    |

| Gender | Salary |
|--------|--------|
| f      | 47k    |
| m      | 22k    |

- Joins are the most expensive distributed operation
    - might need to shuffle large intermediate results

- Different Strategies for Parallel Execution
    - symmetric repartitioning
      (synonym: (fully) re-partitioned join)
    - asymmetric re-partitioning
      (synonym: one-way re-partitioned join, directed join)
    - fragment and replicate
      (synonym: broadcast join)
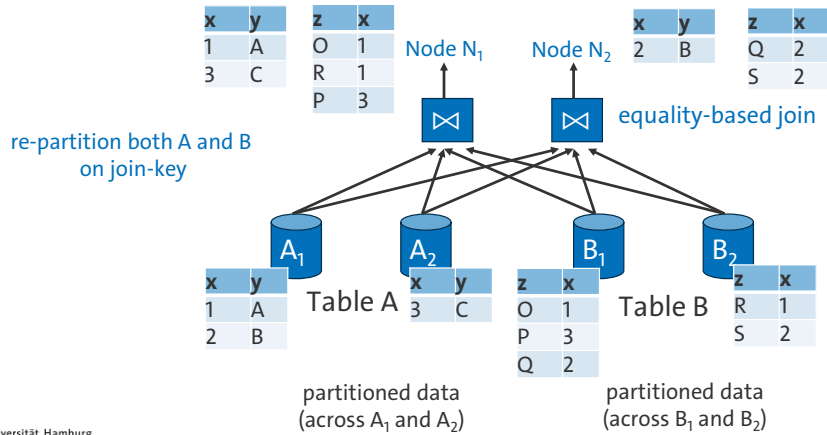
# Join: Symmetric Repartitioning

**„Repartitioned join"**

- Both join partners are repartitioned after the join attribute
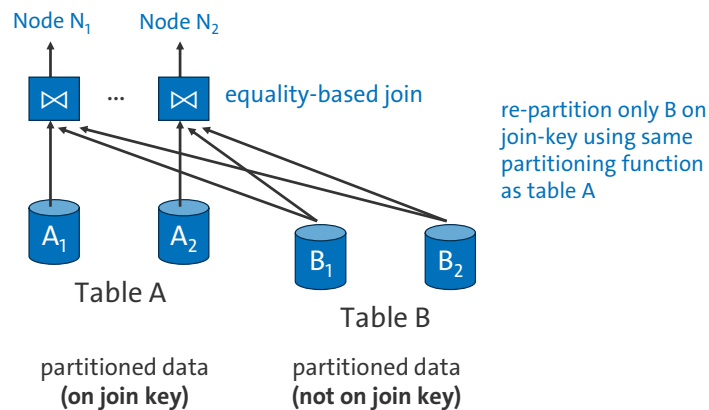
- High communication costs

    → Avoid!

Node $N_1$    Node $N_2$

re-partition both
A and B
on join-key

equality-based join

$A_1$   $A_2$        $B_1$   $B_2$

Table A           Table B

partitioned data
(across two nodes
$N_1$ and $N_2$)

partitioned data
(across two nodes
$N_1$ and $N_2$)

Big Data Analytics

# Join: Symmetric Repartitioning

**Example:** $A \bowtie_{A.x=B.x} B$



re-partition both A and B on join-key

equality-based join

Table A

Table B

partitioned data (across $A_1$ and $A_2$)

partitioned data (across $B_1$ and $B_2$)

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

18

## Join: Asymmetric Repartitioning

Node $N_1$   Node $N_2$

⋈ ... ⋈   equality-based join

re-partition only B on join-key using same partitioning function as table A

$A_1$   $A_2$   $B_1$   $B_2$

Table A

Table B

partitioned data **(on join key)**   partitioned data **(not on join key)**

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

Big Data Analytics

**One-way redistribution join "directed join"**
- One of the two join partners is partitioned after the join attribute
- Partitions of the other join partner are partitioned newly at runtime after the join attribute
- Example:
    - Order relation is partitioned according to the customer key
    - Repartitioning by the attribute O_ORDERKEY

19

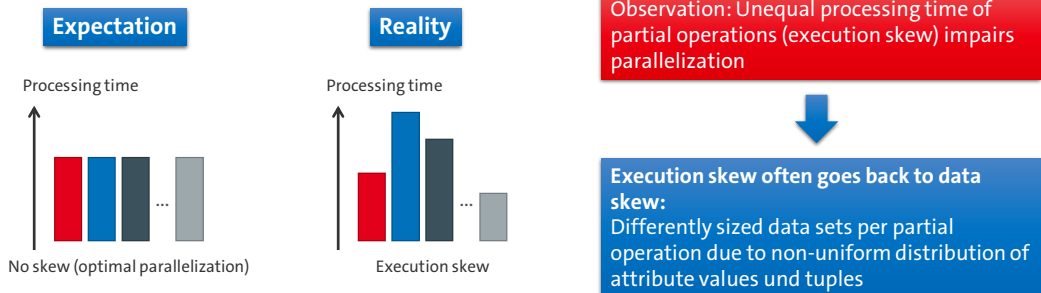# Example: Asymmetric Repartitioning

**Example:** $A \bowtie_{A.x=B.x} B$



re-partition B on join-key

Table A

Table B

Partitioned data (on join key "x")

Partitioned data

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

Big Data Analytics

# Join: Fragment and Replicate

partitioned data
(across two nodes
$A_1$ and $A_2$)

Replicate A

Replicate (all fragments of) smaller table
(here: table A) to all nodes

Table A

partitioned data
(across two nodes $B_1$ and $B_2$)

Table B

**Broadcast Join**
- Partitioning with small relations is not worth it. ...
- Assign a copy of the smaller join partner to the partitions of the larger join partner
     → Advantage: no relation must be partitioned after the join attribute

**Data Distribution Skew (tuple placement skew)**
- Different partition sizes
- Uneven duration of scanning operations
- Treatment: Best knowledge of the distribution of values for distribution attributes
    - Histograms
    - Sampling
    - Computed during sorting on sort merge joins

**Redistribution Skew**
- Distribution function leads to different fragment sizes
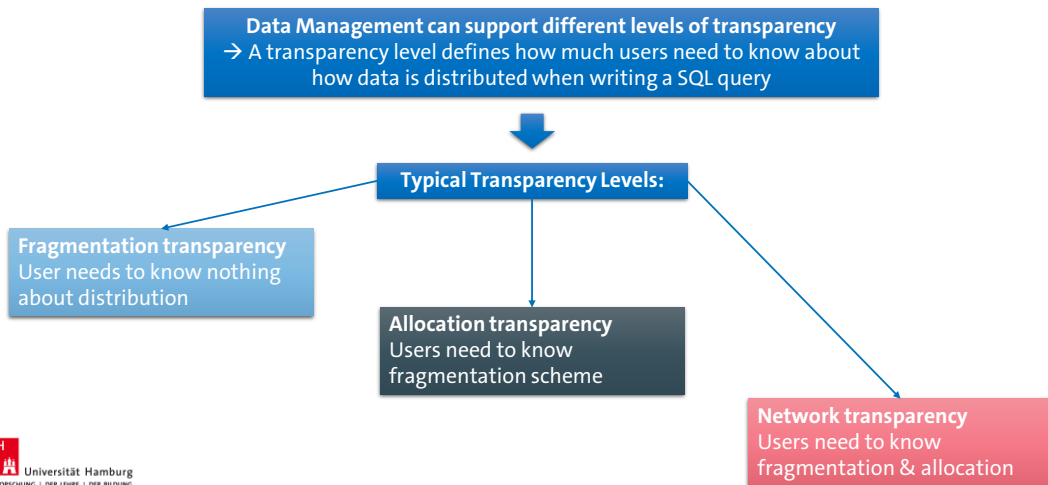- Treatment: Same as data distribution skew

**Selectivity Skew**
- Different hit rates per partition, e.g. range queries regarding distribution attribute for range partitioning
- Treatment: Hardly treatable, as determined by request and data transfer
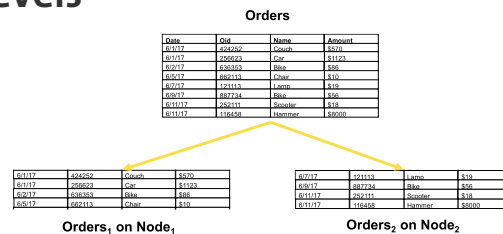
**Join Product Skew**
- Different join selectivity per node
- Treatment:
    - Estimation of the total size of the join result as well as the resulting value distribution for the join attribute
    - Determine area partitioning, which provides a roughly equal partial result for each of the p join processors

**Transparency for Query Formulation**

**Data Management can support different levels of transparency**
→ A transparency level defines how much users need to know about how data is distributed when writing a SQL query

**Typical Transparency Levels:**

**Fragmentation transparency**
User needs to know nothing about distribution

**Allocation transparency**
Users need to know fragmentation scheme

**Network transparency**
Users need to know fragmentation & allocation

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

Big Data Analytics

More examples for transparency levels which might or might not be available in your chosen system: location transparency (like network transparency), naming transparency (no equal names must exist on two different sites),...

**Transparency levels**

Orders table with columns: Date, Oid, Name, Amount — split into Orders₁ on Node₁ and Orders₂ on Node₂

Fragmentation transparency:
```
select COUNT(*)
  from Orders
 where Amount>=5
```

Allocation transparency:
```
select COUNT(*)
  from (Orders₁ union Orders₂)
 where Amount>=5
```

Network transparency:
```
select COUNT(*)
  from (Orders₁@Node₁
        union
         Orders₂@Node₂)
 where Amount>=5
```

**In a distributed DBMS with fragmentation transparency, users only need to know the database schema (tables & attributes)**
- Fragmentation and allocation are "hidden" from users
- Behaves like a non-distributed database system

**Allocation transparency: Users need to know fragmentation information (i.e., number and names of partitions, as well as partitioning function used)**
- Allocation is "hidden" from users
- Users have to be aware of the partitioning scheme

**Network transparency: Users need to know fragmentation and allocation**
- However, users do not need to know physical network addresses (e.g., IPs) of machines

**Which level of transparency?**
- Compromise: Ease of use versus ability to control query processing by user
- Choice often depends on application
    - Full transparency often poor choice for geographically distributed DBMSs;  Users want to control which operations are executed remotely
    - However, parallel DBMSs often implement full transparency