

# Vertical Fragmentation

# Vertical Fragmentation

- More difficult than horizontal, because more alternatives exist.
- Two approaches :
  - Grouping
    - ♦ attributes to fragments
  - Splitting
    - ♦ relation to fragments

# Vertical Fragmentation

- Overlapping fragments
  - ↳ grouping
- Non-overlapping fragments
  - ↳ splitting
- We do not consider the replicated key attributes to be overlapping.
- Advantage:
  - Easier to enforce functional dependencies (for integrity checking etc.)

# VF – Information Requirements

- Application Information

- Attribute affinities

- ♦ a measure that indicates how closely related the attributes are
    - ♦ This is obtained from more primitive usage data

- Attribute usage values

- ♦ Given a set of queries  $Q = \{q_1, q_2, \dots, q_q\}$  that will run on the relation  $R[A_1, A_2, \dots, A_n]$ ,

$$use(q_i, A_j) = \begin{cases} 1 & \text{if attribute } A_j \text{ is referenced by query } q_i \\ 0 & \text{otherwise} \end{cases}$$

$use(q_i, \bullet)$  can be defined accordingly

# VF – Definition of $use(q_i, A_j)$

Consider the following 4 queries for relation PROJ

$q_1$ : **SELECT** BUDGET  
FROM PROJ  
WHERE PNO=Value

$q_2$ : **SELECT** PNAME, BUDGET  
FROM PROJ

$q_3$ : **SELECT** PNAME  
FROM PROJ  
WHERE LOC=Value

$q_4$ : **SELECT** SUM(BUDGET)  
FROM PROJ  
WHERE LOC=Value

Let  $A_1 = \text{PNO}$ ,  $A_2 = \text{PNAME}$ ,  $A_3 = \text{BUDGET}$ ,  $A_4 = \text{LOC}$

	$A_1$	$A_2$	$A_3$	$A_4$
$q_1$	1	0	1	0
$q_2$	0	1	1	0
$q_3$	0	1	0	1
$q_4$	0	0	1	1

# VF – Affinity Measure $aff(A_i, A_j)$

The **attribute affinity measure** between two attributes  $A_i$  and  $A_j$  of a relation  $R[A_1, A_2, \dots, A_n]$  with respect to the set of applications  $Q = (q_1, q_2, \dots, q_q)$  is defined as follows :

$$aff(A_i, A_j) = \sum_{\text{all queries that access } A_i \text{ and } A_j} (\text{query access})$$

$$\text{query access} = \sum_{\text{all sites}} \text{access frequency of a query} * \frac{\text{access}}{\text{execution}}$$

# VF – Calculation of $aff(A_i, A_j)$

	$A_1$	$A_2$	$A_3$	$A_4$
$q_1$	1	0	1	0
$q_2$	0	1	1	0
$q_3$	0	1	0	1
$q_4$	0	0	1	1

$$aff(A_1, A_1) = 15 \cdot 1 + 20 \cdot 1 + 10 \cdot 1 = 45$$

$$aff(A_1, A_2) = 0$$

$$aff(A_1, A_3) = 15 \cdot 1 + 20 \cdot 1 + 10 \cdot 1 = 45$$

$$aff(A_1, A_4) = 0$$

	$S_1$	$S_2$	$S_3$
$q_1$	15	20	10
$q_2$	5	0	0
$q_3$	25	25	25
$q_4$	3	0	0

$$aff(A_2, A_2) = 5 + 75 = 80$$

$$aff(A_2, A_3) = 5$$

$$aff(A_2, A_4) = 25 \cdot 1 + 25 \cdot 1 + 25 \cdot 1 = 75$$

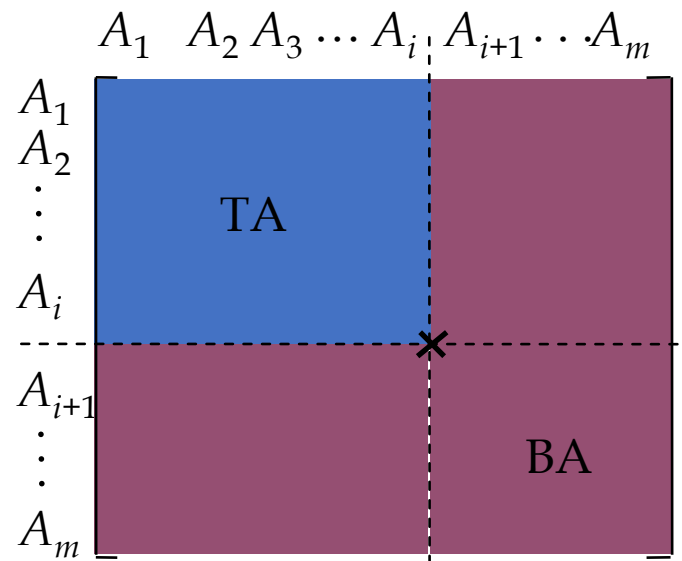
	$A_1$	$A_2$	$A_3$	$A_4$
$A_1$	45	0	45	0
$A_2$	0	80	5	75
$A_3$	45	5	53	3
$A_4$	0	75	3	78

$$aff(A_3, A_4) = 3$$

$$aff(A_4, A_4) = 75 + 3 = 78$$

# VF – Algorithm

How can you divide a set of clustered attributes  $\{A_1, A_2, \dots, A_n\}$  into two (or more) sets  $\{A_1, A_2, \dots, A_i\}$  and  $\{A_i, \dots, A_n\}$  such that there are no (or minimal) applications that access both (or more than one) of the sets.





# VF – Clustering Algorithm

- Take the attribute affinity matrix  $AA$  and reorganize the attribute orders to form clusters where the attributes in each cluster demonstrate high affinity to one another.
- Bond Energy Algorithm (BEA) has been used for clustering of entities. BEA finds an ordering of entities (in our case attributes) such that the global affinity measure is maximized.

$$AM = \sum_i \sum_j \text{(affinity of } A_i \text{ and } A_j \text{ with their neighbors)}$$

# VF – Algorithm

Define

$TQ$  = set of applications that access only  $TA$

$BQ$  = set of applications that access only  $BA$

$OQ$  = set of applications that access both  $TA$  and  $BA$

and

$CTQ$  = total number of accesses to attributes by applications that access only  $TA$

$CBQ$  = total number of accesses to attributes by applications that access only  $BA$

$COQ$  = total number of accesses to attributes by applications that access both  $TA$  and  $BA$

Then find the point along the diagonal that maximizes

$$CTQ * CBQ - COQ^2$$

# VF – Algorithm

Two problems :

- Cluster forming in the middle of the CA matrix
  - Shift a row up and a column left and apply the algorithm to find the “best” partitioning point
  - Do this for all possible shifts
  - Cost  $O(m^2)$
- More than two clusters
  - $m$ -way partitioning
  - try 1, 2, ...,  $m-1$  split points along diagonal and try to find the best point for each of these
  - Cost  $O(2^m)$

# Bond Energy Algorithm

**Input:** The AA matrix

**Output:** The clustered affinity matrix CA which is a perturbation of AA

- ① *Initialization:* Place and fix one of the columns of AA in CA.
- ② *Iteration:* Place the remaining  $n-i$  columns in the remaining  $i+1$  positions in the CA matrix. For each column, choose the placement that makes the most contribution to the global affinity measure.
- ③ *Row order:* Order the rows according to the column ordering.

# Bond Energy Algorithm

“Best” placement? Define contribution of a placement:

$$cont(A_i, A_k, A_j) = 2bond(A_i, A_k) + 2bond(A_k, A_l) - 2bond(A_i, A_j)$$

where

$$bond(A_x, A_y) = \sum_{z=1}^n aff(A_z, A_x) aff(A_z, A_y)$$

# BEA – Example

Consider the following AA matrix and the corresponding CA matrix where  $A_1$  and  $A_2$  have been placed. Place  $A_3$ :

$$AA = \begin{array}{c} A_1 \\ A_2 \\ A_3 \\ A_4 \end{array} \begin{array}{c} A_1 \\ A_2 \\ A_3 \\ A_4 \end{array} \begin{bmatrix} 45 & 0 & 5 & 0 \\ 0 & 80 & 5 & 75 \\ 45 & 5 & 53 & 3 \\ 0 & 75 & 3 & 78 \end{bmatrix} \quad CA = \begin{array}{c} A_1 \\ A_2 \end{array} \begin{array}{c} A_1 \\ A_2 \end{array} \begin{bmatrix} 45 & 0 \\ 0 & 80 \\ 45 & 5 \\ 0 & 75 \end{bmatrix}$$

$$2bond(A_1, A_3) = 45 * 5 + 0 * 5 + 45 * 53 + 0 * 3 = 4410$$

$$\begin{aligned} 2bond(A_2, A_3) \\ = 0 * 5 + 80 * 5 + 5 * 53 + 75 * 3 = 890 \end{aligned}$$

# BEA – Example

Consider the following AA matrix and the corresponding CA matrix where  $A_1$  and  $A_2$  have been placed. Place  $A_3$ :

$$AA = \begin{array}{c} \\ A_1 \\ A_2 \\ A_3 \\ A_4 \end{array} \begin{array}{c} A_1 \\ A_2 \\ A_3 \\ A_4 \end{array} \begin{bmatrix} 45 & 0 & 5 & 0 \\ 0 & 80 & 5 & 75 \\ 45 & 5 & 53 & 3 \\ 0 & 75 & 3 & 78 \end{bmatrix} \quad CA = \begin{array}{c} \\ A_1 \\ A_2 \end{array} \begin{array}{c} A_1 \\ A_2 \end{array} \begin{bmatrix} 45 & 0 \\ 0 & 80 \\ 45 & 5 \\ 0 & 75 \end{bmatrix}$$

Ordering (0-3-1) :

$$\begin{aligned} cont(A_0, A_3, A_1) &= 2bond(A_0, A_3) + 2bond(A_3, A_1) - 2bond(A_0, A_1) \\ &= 2 * 0 + 2 * 4410 - 2 * 0 = 8820 \end{aligned}$$

Ordering (1-3-2) :

$$\begin{aligned} cont(A_1, A_3, A_2) &= 2bond(A_1, A_3) + 2bond(A_3, A_2) - 2bond(A_1, A_2) \\ &= 2 * 4410 + 2 * 890 - 2 * 225 = 10150 \end{aligned}$$

Ordering (2-3-4) :

$$cont(A_2, A_3, A_4) = 1780$$

# BEA – Example

- Therefore, the CA matrix has the form

$$\begin{array}{c} A_1 \quad A_3 \quad A_2 \\ \left[ \begin{array}{ccc} 45 & 5 & 0 \\ 0 & 5 & 80 \\ 45 & 53 & 5 \\ 0 & 3 & 75 \end{array} \right] \end{array}$$

- When  $A_4$  is placed, the final form of the CA matrix (after row organization) is

$$\begin{array}{c} A_1 \quad A_3 \quad A_2 \quad A_4 \\ \begin{array}{c} A_1 \\ A_3 \\ A_2 \\ A_4 \end{array} \left[ \begin{array}{cccc} 45 & 5 & 0 & 0 \\ 45 & 53 & 5 & 3 \\ 0 & 5 & 80 & 75 \\ 0 & 3 & 75 & 78 \end{array} \right] \end{array}$$



# BEA – Example

- After adding A3, CA matrix has the form

$$\begin{array}{c}
 A_1 \\
 A_2 \\
 A_3 \\
 A_4
 \end{array}
 \begin{array}{c}
 A_1 \\
 A_3 \\
 A_2 \\
 A_4
 \end{array}
 \begin{bmatrix}
 45 & 5 & 0 & 0 \\
 0 & 5 & 80 & 75 \\
 45 & 53 & 5 & 3 \\
 0 & 3 & 75 & 78
 \end{bmatrix}$$

Now arranging the rows:

$$\begin{array}{c}
 A_1 \\
 A_2 \\
 A_3 \\
 A_4
 \end{array}
 \begin{array}{c}
 A_1 \\
 A_3 \\
 A_2 \\
 A_4
 \end{array}
 \begin{bmatrix}
 45 & 5 & 0 & 0 \\
 45 & 53 & 5 & 3 \\
 0 & 5 & 80 & 75 \\
 0 & 3 & 75 & 78
 \end{bmatrix}$$

$$\begin{array}{c}
 A_1 \\
 A_2 \\
 A_3 \\
 A_4
 \end{array}
 \begin{array}{c}
 A_1 \\
 A_2 \\
 A_3 \\
 A_4
 \end{array}
 \begin{bmatrix}
 45 & 0 & 5 & 0 \\
 0 & 80 & 5 & 75 \\
 45 & 5 & 53 & 3 \\
 0 & 75 & 3 & 78
 \end{bmatrix}
 \quad
 \begin{array}{c}
 A_1 \\
 A_2
 \end{array}
 \begin{array}{c}
 A_1 \\
 A_2
 \end{array}
 \begin{bmatrix}
 45 & 0 \\
 0 & 80 \\
 45 & 5 \\
 0 & 75
 \end{bmatrix}$$

# BEA – Example

- The final step is to divide the set of attributes into two sets such that the queries that access both sets are minimized.
- The appropriate split point on the diagonal must be figured out.
- This is an optimization problem. The optimal points on the diagonal must be determined.

	$A_1$	$A_3$	$A_2$	$A_4$
$A_1$	45	5	0	0
$A_2$	45	53	5	3
$A_3$	0	5	80	75
$A_4$	0	3	75	78

# VF – Correctness

A relation  $R$ , defined over attribute set  $A$  and key  $K$ , generates the vertical partitioning  $F_R = \{R_1, R_2, \dots, R_r\}$ .

- Completeness

- The following should be true for  $A$ :

$$A = \bigcup A_{R_i}$$

- Reconstruction

- Reconstruction can be achieved by

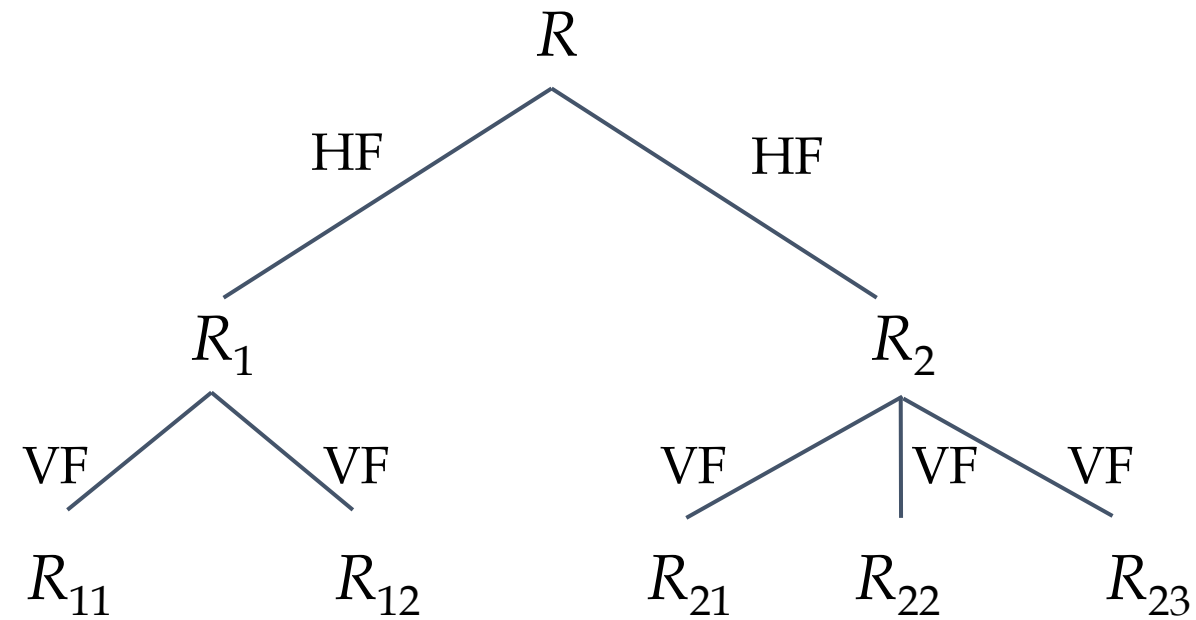
$$R = \bowtie_K R_i \quad \forall R_i \in F_R$$

- Disjointness

- TID's are not considered to be overlapping since they are maintained by the system

- Duplicated keys are not considered to be overlapping

# Hybrid Fragmentation



# Fragment Allocation

- Problem Statement

Given

$F = \{F_1, F_2, \dots, F_n\}$  fragments

$S = \{S_1, S_2, \dots, S_m\}$  network sites

$Q = \{q_1, q_2, \dots, q_q\}$  applications

Find the "optimal" distribution of  $F$  to  $S$ .

- Optimality

→ Minimal cost

- ♦ Communication + storage + processing (read & update)
- ♦ Cost in terms of time (usually)

→ Performance

Response time and/or throughput

→ Constraints

- ♦ Per site constraints (storage & processing)

# Reference

M. Tamer Ozsü and Patrick Valduriez, "Principles of Distributed Database Systems," Prentice Hall