

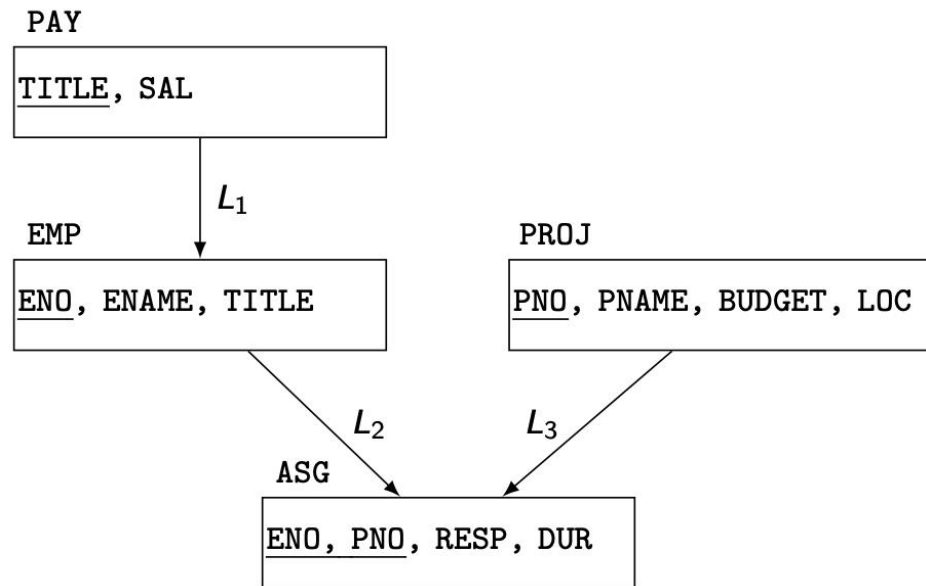
Fragmentation

- Horizontal Fragmentation (HF)
 - ▣ Primary Horizontal Fragmentation (PHF)
 - ▣ Derived Horizontal Fragmentation (DHF)
- Vertical Fragmentation (VF)
- Hybrid Fragmentation (HF)

PHF – Information Requirements

■ Database Information

□ relationship



□ cardinality of each relation: $card(R)$

PHF – Information Requirements

■ Application Information

□ **minterm selectivities:** $sel(m_i)$

- The number of tuples of the relation that would be accessed by a user query which is specified according to a given minterm predicate m_i .

□ **access frequencies:** $acc(q_i)$

- The frequency with which a user application q_i accesses data.
- Access frequency for a minterm predicate can also be defined.

Primary Horizontal Fragmentation

Definition :

$$R_j = \sigma_{F_j}(R), \quad 1 \leq j \leq w$$

where F_j is a selection formula, which is (preferably) a minterm predicate.

Therefore,

A horizontal fragment R_j of relation R consists of all the tuples of R which satisfy a minterm predicate m_j .



Given a set of minterm predicates M , there are as many horizontal fragments of relation R as there are minterm predicates.

Set of horizontal fragments also referred to as **minterm fragments**.

PHF – Algorithm

Given: A relation R , the set of simple predicates Pr

Output: The set of fragments of $R = \{R_1, R_2, \dots, R_w\}$ which obey the fragmentation rules.

Preliminaries :

- ❑ Pr should be *complete*
- ❑ Pr should be *minimal*

Completeness of Simple Predicates

- A set of simple predicates Pr is said to be *complete* if and only if the accesses to the tuples of the minterm fragments defined on Pr requires that two tuples of the same minterm fragment have the same probability of being accessed by any application.
- Example :
 - ❑ Assume PROJ[PNO,PNAME,BUDGET,LOC] has two applications defined on it.
 - ❑ Find the budgets of projects at each location. (1)
 - ❑ Find projects with budgets less than \$200000. (2)

PHF – Example

■ Fragmentation of relation PROJ

□ Applications:

- Find the name and budget of projects given their no.

- Issued at three sites

- Access project information according to budget

- one site accesses ≤ 200000 other accesses > 200000

□ Simple predicates

□ For application (1)

$p_1 : \text{LOC} = \text{"Montreal"}$

$p_2 : \text{LOC} = \text{"New York"}$

$p_3 : \text{LOC} = \text{"Paris"}$

□ For application (2)

$p_4 : \text{BUDGET} \leq 200000$

$p_5 : \text{BUDGET} > 200000$

□ $Pr = Pr' = \{p_1, p_2, p_3, p_4, p_5\}$

PHF – Example

■ Fragmentation of relation PROJ continued

□ Minterm fragments left after elimination

$m_1 : (\text{LOC} = \text{"Montreal"}) \wedge (\text{BUDGET} \leq 200000)$

$m_2 : (\text{LOC} = \text{"Montreal"}) \wedge (\text{BUDGET} > 200000)$

$m_3 : (\text{LOC} = \text{"New York"}) \wedge (\text{BUDGET} \leq 200000)$

$m_4 : (\text{LOC} = \text{"New York"}) \wedge (\text{BUDGET} > 200000)$

$m_5 : (\text{LOC} = \text{"Paris"}) \wedge (\text{BUDGET} \leq 200000)$

$m_6 : (\text{LOC} = \text{"Paris"}) \wedge (\text{BUDGET} > 200000)$

PHF – Example

PROJ₁

PNO	PNAME	BUDGET	LOC
P1	Instrumentation	150000	Montreal

PROJ₃

PNO	PNAME	BUDGET	LOC
P2	Database Develop.	135000	New York

PROJ₄

PNO	PNAME	BUDGET	LOC
P3	CAD/CAM	255000	New York

PROJ₆

PNO	PNAME	BUDGET	LOC
P4	Maintenance	310000	Paris

PHF – Correctness

■ Completeness

- Since Pr' is complete and minimal, the selection predicates are complete

■ Reconstruction

- If relation R is fragmented into $F_R = \{R_1, R_2, \dots, R_r\}$

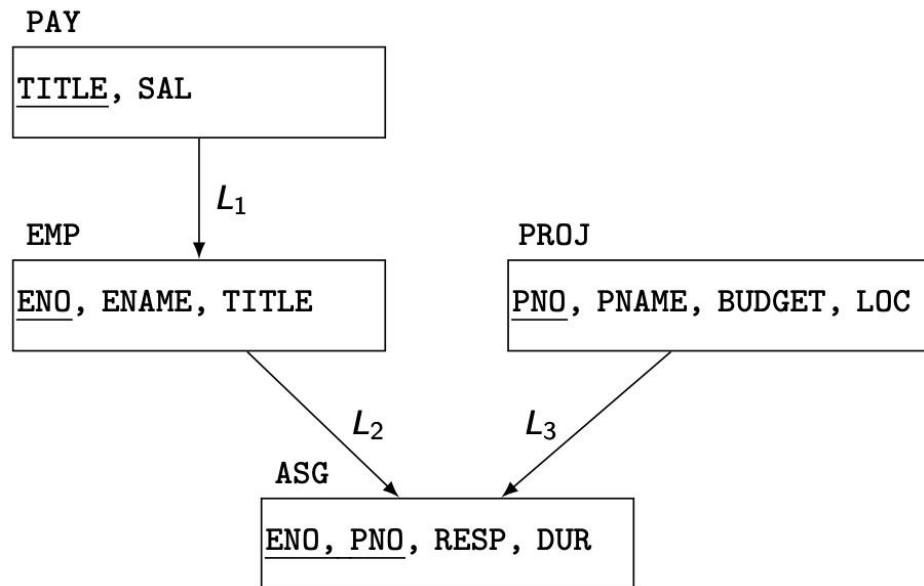
$$R = \bigcup_{\forall R_i \in F_R} R_i$$

■ Disjointness

- Minterm predicates that form the basis of fragmentation should be mutually exclusive.

Derived Horizontal Fragmentation

- Defined on a member relation of a link according to a selection operation specified on its owner.
 - ❑ Each link is an equijoin.
 - ❑ Equijoin can be implemented by means of semiioins.



DHF – Definition

Given a link L where $owner(L)=S$ and $member(L)=R$, the derived horizontal fragments of R are defined as

$$R_i = R \bowtie_F S_i, 1 \leq i \leq w$$

where w is the maximum number of fragments that will be defined on R and

$$S_i = \sigma_{F_i}(S)$$

where F_i is the formula according to which the primary horizontal fragment S_i is defined.

DHF – Example

Given link L_1 where $\text{owner}(L_1)=\text{SKILL}$ and $\text{member}(L_1)=\text{EMP}$

$$\text{EMP}_1 = \text{EMP} \times \text{SKILL}_1$$

$$\text{EMP}_2 = \text{EMP} \times \text{SKILL}_2$$

where

$$\text{SKILL}_1 = \sigma_{\text{SAL} \leq 30000}(\text{SKILL})$$

$$\text{SKILL}_2 = \sigma_{\text{SAL} > 30000}(\text{SKILL})$$

ASG₁

ENO	PNO	RESP	DUR
E3	P3	Consultant	10
E3	P4	Engineer	48
E4	P2	Programmer	18
E7	P3	Engineer	36

ASG₂

ENO	PNO	RESP	DUR
E1	P1	Manager	12
E2	P1	Analyst	24
E2	P2	Analyst	6
E5	P2	Manager	24
E6	P4	Manager	48
E8	P3	Manager	40

DHF – Correctness

■ Completeness

- Referential integrity
- Let R be the member relation of a link whose owner is relation S which is fragmented as $F_S = \{S_1, S_2, \dots, S_n\}$. Furthermore, let A be the join attribute between R and S . Then, for each tuple t of R , there should be a tuple t' of S such that

$$t[A] = t'[A]$$

■ Reconstruction

- Same as primary horizontal fragmentation.

■ Disjointness

- Simple join graphs between the owner and the member fragments.

Vertical Fragmentation

- Has been studied within the centralized context
 - design methodology
 - physical clustering
- 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	q_2 : SELECT	PNAME, BUDGET
FROM	PROJ	FROM	PROJ
WHERE	PNO=Value		
q_3 : SELECT	PNAME	q_4 : SELECT	SUM (BUDGET)
FROM	PROJ	FROM	PROJ
WHERE	LOC=Value	WHERE	LOC=Value

	PNO	PNAME	BUDGET	LOC
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 – Algorithm

Two problems :

❑ Cluster forming in the middle of the Clustered Affinity 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 $O(2^m)$

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

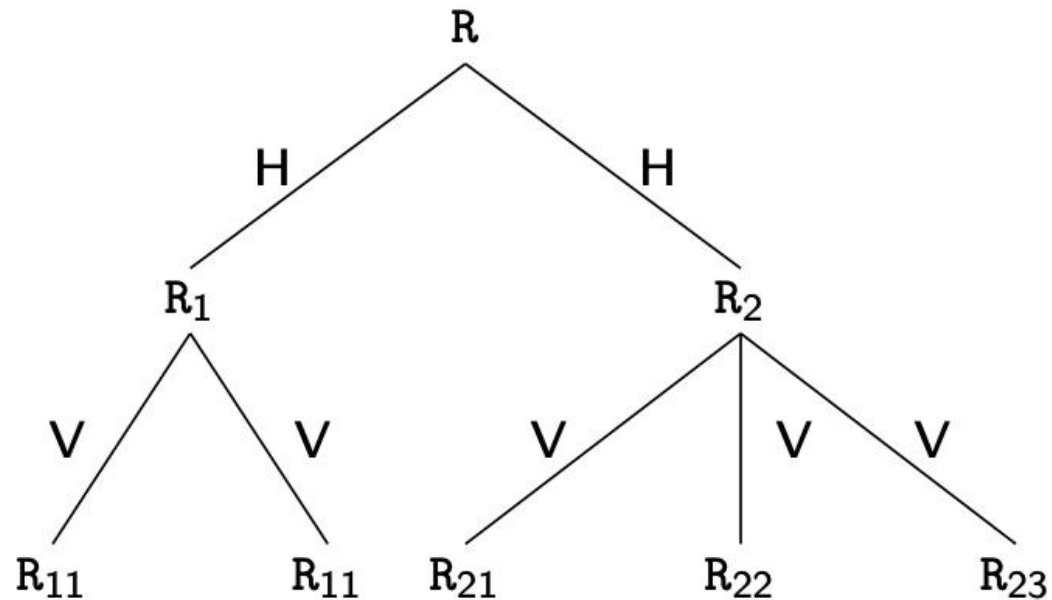
- Reconstruction can be achieved by

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

■ Disjointness

- Tuple Identifiers (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



Reconstruction of Hybrid Fragmentation

