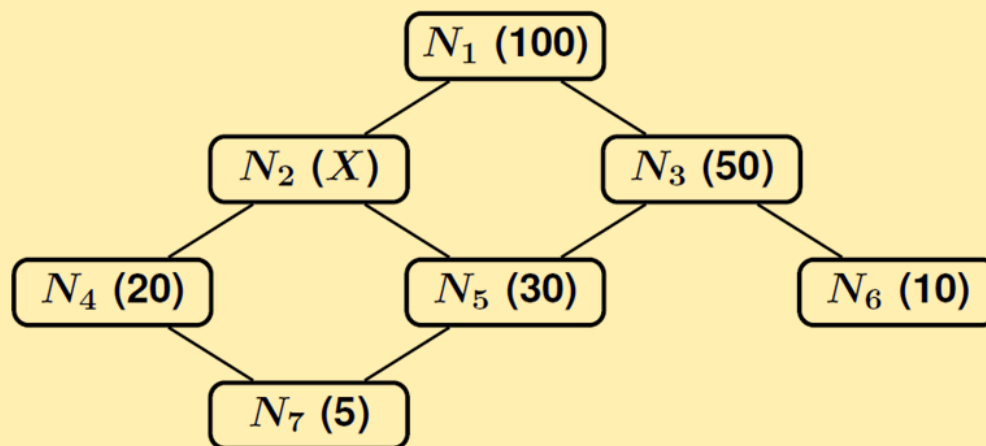


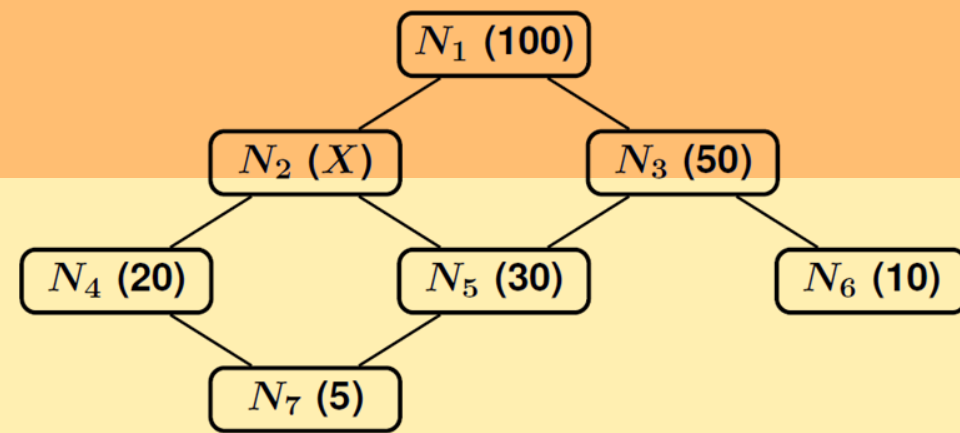
EXERCISE AT HOME

Let us consider the following lattice of possible candidate views to materialize. The numbers associated with the nodes represent the view size, measured in terms of the number of tuples in the view.



- (1 point)** The value X for N_2 is unknown. What is the domain of admissible values for X ?
- (5 points)** Select 2 views to materialize, different from N_1 , with the greedy algorithm HRU. Determine the various possible results of HRU on the basis of the unknown value X for N_2 .

EXERCISE AT HOME

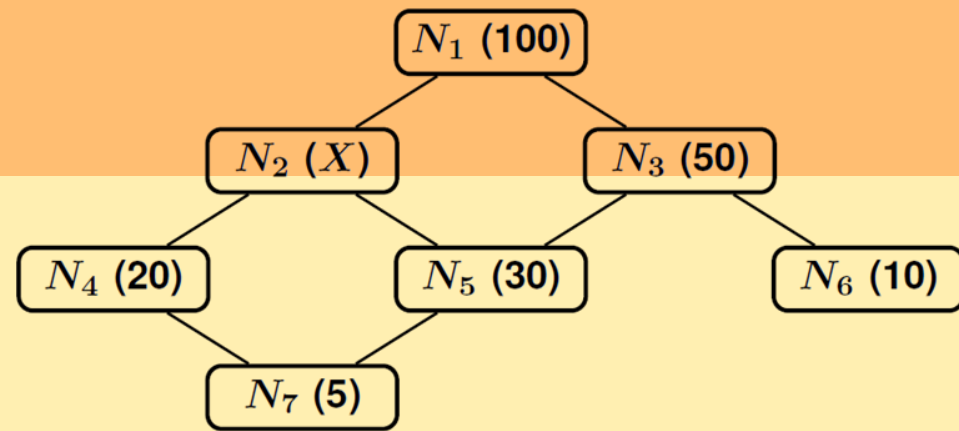


- $\max\{\text{descendants of } N_2\} = 30 \leq X \leq 100 = \min\{\text{ascendants of } N_2\}$

| View | First Choice | | |
|-------|--------------------------|--|--|
| N_2 | $(100-X) \cdot 4$ | | |
| N_3 | $(100-50) \cdot 4 = 200$ | | |
| N_4 | $(100-20) \cdot 2 = 160$ | | |
| N_5 | $(100-30) \cdot 2 = 140$ | | |
| N_6 | $(100-10) = 90$ | | |
| N_7 | $(100-5) = 95$ | | |

- $(100-X) \cdot 4 > 200$ iff $X < 50$

EXERCISE AT HOME



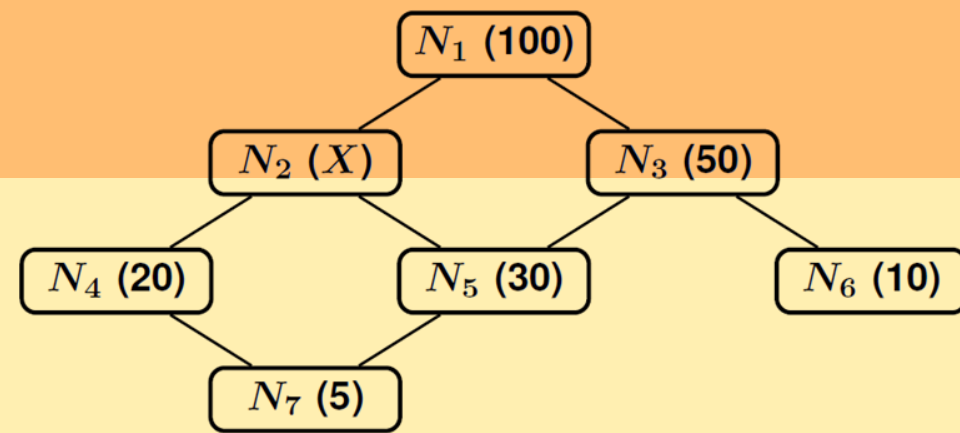
- $\max\{\text{descendents of } N_2\} = 30 \leq X \leq 100 = \min\{\text{ascendents of } N_2\}$

| View | First Choice | Second choice $30 \leq X < 50$ | |
|-------|--------------------------|---|--|
| N_2 | $(100-X) \cdot 4$ | - | |
| N_3 | $(100-50) \cdot 4 = 200$ | $(100-50) \cdot 2 + (X-50) \cdot 2 = 100$ | |
| N_4 | $(100-20) \cdot 2 = 160$ | $(X-20) \cdot 2 \leq 60$ | |
| N_5 | $(100-30) \cdot 2 = 140$ | $(X-30) \cdot 2 \leq 40$ | |
| N_6 | $(100-10) = 90$ | $(100-10) = 90$ | |
| N_7 | $(100-5) = 95$ | $(X-5) \leq 45$ | |

- $(100-X) \cdot 4 > 200$ iff $X < 50$

$$M = \{N_1, N_2, N_3\}$$

EXERCISE AT HOME



- $\max\{\text{descendents of } N_2\} = 30 \leq X \leq 100 = \min\{\text{ascendents of } N_2\}$

| View | First Choice | Second choice $30 \leq X < 50$ | Second choice $50 \leq X \leq 100$ |
|-------|--------------------------|---|---|
| N_2 | $(100-X) \cdot 4$ | - | $(100-X) \cdot 2 + (50-X) \cdot 2 \leq 100$ |
| N_3 | $(100-50) \cdot 4 = 200$ | $(100-50) \cdot 2 + (X-50) \cdot 2 = 100$ | - |
| N_4 | $(100-20) \cdot 2 = 160$ | $(X-20) \cdot 2 \leq 60$ | $(100-20) + (50-20) = 110$ |
| N_5 | $(100-30) \cdot 2 = 140$ | $(X-30) \cdot 2 \leq 40$ | $(50-30) \cdot 2 = 40$ |
| N_6 | $(100-10) = 90$ | $(100-10) = 90$ | $(50-10) = 40$ |
| N_7 | $(100-5) = 95$ | $(X-5) \leq 45$ | $(50-5) = 45$ |

- $(100-X) \cdot 4 > 200$ iff $X < 50$

$$M = \{N_1, N_2, N_3\}$$

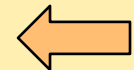
$$M = \{N_1, N_3, N_4\}$$

RELATIONAL DBMS EXTENSIONS FOR DW

- SQL extensions
- Index and storage structures
- Star query physical plans
- Materialized views

NEXT LESSON

- Optimization techniques for star queries with grouping and aggregations



TODAY: functional dependencies and their usage in query optimization

FUNCTIONAL DEPENDENCIES

Functional dependencies

Given a relation schema $R(T)$ and $X, Y \subseteq T$, a **functional dependency (FD)** is a constraint on R of the form $X \rightarrow Y$, i.e. **X functionally determines Y or Y is determined by X** , if

$\forall r$ valid instance of R .

$\forall t1, t2 \in r$. if $t1[X] = t2[X]$ then $t1[Y] = t2[Y]$

Convention: ..., X, Y, Z represent sets of attributes; A, B, C, \dots represent single attributes; a set of attributes $\{A, B, C\}$ is represented just as ABC .

FUNCTIONAL DEPENDENCIES: EXAMPLE

StudentsExams(StudCode, Name, City, Region, BirthYear, Subject, Grade)

| StudCode | Name | City | Region | BirthYear | Subject | Grade |
|----------|------|------|--------|-----------|---------|-------|
| 1234567 | N1 | C1 | R1 | 1995 | DB | 30 |
| 1234567 | N1 | C1 | R1 | 1995 | SE | 28 |
| 1234568 | N2 | C2 | R2 | 1994 | DB | 30 |
| 1234568 | N2 | C2 | R2 | 1994 | SE | 26 |

StudCode \rightarrow Name City Region BirthYear ?

City \rightarrow Region ?

Subject \rightarrow Grade ? **NO**

StudCode Subject \rightarrow Grade ?

Subject \rightarrow Subject ? **trivial**

X \rightarrow {} ? **trivial**

{ } \rightarrow University ? **YES, if University is constant**

REASONING ABOUT FDs: LOGICAL IMPLICATION

Notation:

- $R \langle T, F \rangle$ is a relational schema with **attributes** T and a **set** of functional dependencies F .
Example: $F = \{ X \rightarrow Y, Y \rightarrow Z \}$
- A $FD \in F$ is a constraint on relational instances r of $R \langle T, F \rangle$
 - r is a valid instance if $\forall t1, t2 \in r$. if $t1[X] = t2[X]$ then $t1[Y] = t2[Y]$

Usage: defined by the designer, enforced by the DBMS. In practice, only the functional dependency $K \rightarrow T$ are enforced, when K is a key.

REASONING ABOUT FDs: LOGICAL IMPLICATION

Notation:

- $R \langle T, F \rangle$ is a relational schema with **attributes** T and a **set** of functional dependencies F . Example: $F = \{ X \rightarrow Y, Y \rightarrow Z \}$
- A $FD \in F$ is a constraint on relational instances

Given a set F of FDs, other FDs will generally be ‘implied’ by this set in the following sense:

Definition Given a schema $R \langle T, F \rangle$, we say that F *implies* $X \rightarrow Y$, if every instance r of R that satisfies F also satisfies $X \rightarrow Y$.

Example: $\{ X \rightarrow Y, Y \rightarrow Z \}$ implies $X \rightarrow Z$?

To test if a FD implied by a set F, a set of inference rules can be used with the property of being **sound** and **complete** (F implies FD iff $F \vdash$ FD)

Armstrong axioms:

- If $Y \subseteq X$, then $F \vdash X \rightarrow Y$ (reflexivity **R**)
- If $F \vdash X \rightarrow Y$ and $Z \subseteq T$, then $F \vdash XZ \rightarrow YZ$ (augmentation **A**)
- If $F \vdash X \rightarrow Y$ and $F \vdash Y \rightarrow Z$, then $F \vdash X \rightarrow Z$ (transitivity **T**)

Exercises:

1. $\{ X \rightarrow Y, X \rightarrow Z \} \vdash X \rightarrow YZ$ (union **U**)
2. if $Z \subseteq Y$ then: $X \rightarrow Y \vdash X \rightarrow Z$ (decomposition **D**)
3. $F \vdash X \rightarrow A_1, \dots, A_n$ iff $F \vdash X \rightarrow A_1$ and ... $F \vdash X \rightarrow A_n$

CLOSURE OF A SET OF FDs

The FDs implied by F (the **closure** of F) are defined as:

$$F^+ = \{ X \rightarrow Y \mid F \vdash X \rightarrow Y \}$$

Implication problem: to test whether a FD $X \rightarrow Y \in F^+$ (without computing the whole closure of F)

Exercises:

3. $F \vdash X \rightarrow A_1, \dots, A_n$ iff $F \vdash X \rightarrow A_1$ and ... $F \vdash X \rightarrow A_n$

CLOSURE OF A SET OF ATTRIBUTES

Definition Given a scheme $R \langle T, F \rangle$, and $X \subseteq T$, the closure of X is

$$X^+ = \{ A \in T \mid F \vdash X \rightarrow A \}$$

A procedure to solve the implication problem without computing the whole closure of F follows from the following result.

Theorem $F \vdash X \rightarrow Y$ iff $Y \subseteq X^+$

Proof. We proved it as exercise (3)

SLOW CLOSURE

A simple algorithm to compute X^+ is the following

- faster algorithm exist

Algorithm SLOW CLOSURE

input $R\langle T, F \rangle, X \subseteq T$

output X^+

begin

$X^+ = X$

while (changes to X^+) do

for each $W \rightarrow V$ in F with $W \subseteq X^+$ and $V \notin X^+$

do $X^+ = X^+ \cup V$

end

EXAMPLE

$F = \{DB \rightarrow E, B \rightarrow C, A \rightarrow B\}$.

Is $AD \rightarrow E$ in F^+ ?

$X^+ = AD$

$X^+ = ADB$

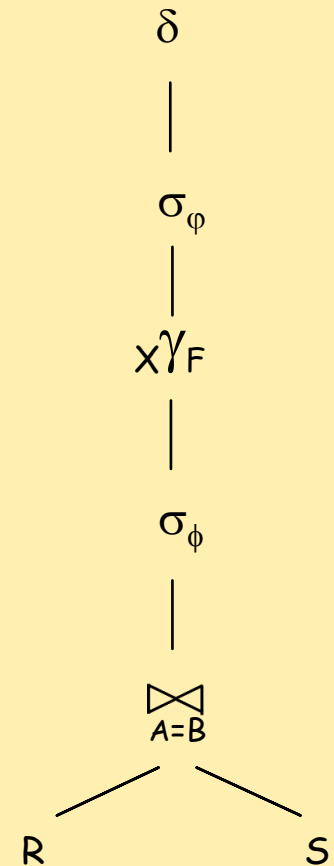
$X^+ = ADBE$



$X^+ = ADBEC$

ASSUMPTIONS

- The tables do not have null values, and **have primary keys**:
 - A key constraint uniquely identifies each record in a table.
 - Tables are then **sets** of tuples
- Table in FROM clause have no attribute with the same name
- Queries are a single *SELECT* with possibly *GROUP BY* and *HAVING* but without subselect and *ORDER BY* clauses.



EXERCISE: FD's lift over cartesian product

If $X \rightarrow Y$ holds in R , is this still the case in $R \times S$?

Hypothesis. If $\forall t1, t2 \in R$. if $t1[X] = t2[X]$ then $t1[Y] = t2[Y]$

Conclusion. $\forall w1, w2 \in R \times S$. if $w1[X] = w2[X]$ then $w1[Y] = w2[Y]$

$$w1 = t1 \circ s1 \Rightarrow w1[X] = t1[X]$$

If X is a key for R and Y a key for S , then is XY a key for $R \times S$?

- XY is a superkey:

$$X \rightarrow R1, \dots, Rn, Y \rightarrow S1, \dots, Sm \mid - XY \rightarrow R1, \dots, Rn, S1, \dots, Sm$$

$$\text{iff } \{R1, \dots, Rn, S1, \dots, Sm\} \subseteq \{X, Y\}^+$$

$$\{X, Y\}^+ = \{X, Y, R1, \dots, Rn, S1, \dots, Sm\}$$

- no subset of XY is a superkey: exercise at home

EXERCISE: FD's lift over selections (and then over joins)

If $X \rightarrow Y$ holds in R , is this still the case in $\sigma_c(R)$?

Hypothesis. If $\forall t1, t2 \in R$. if $t1[X] = t2[X]$ then $t1[Y] = t2[Y]$

Conclusion. $\forall w1, w2 \in \sigma_c(R)$. if $w1[X] = w2[X]$ then $w1[Y] = w2[Y]$

$$w1 = t1 \text{ for some } t1 \in R \Rightarrow w1[X] = t1[X]$$

If X is a key for R , then is X a (super)key for $\sigma_c(R)$?

- X is a superkey:

$$X \rightarrow R1, \dots, Rn \mid - X \rightarrow R1, \dots, Rn$$

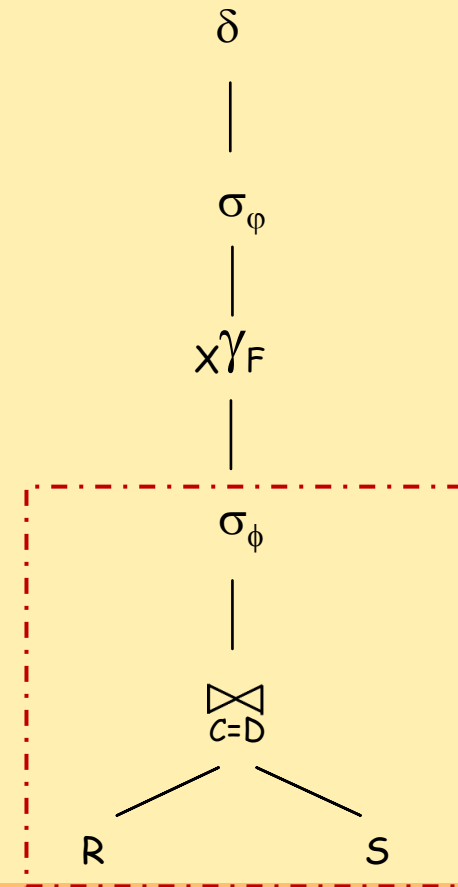
- no subset of X is a superkey: FALSE

for $R(\underline{A}, \underline{B})$, AB is a key

for $\sigma_{B=1}(R)$, A is a key

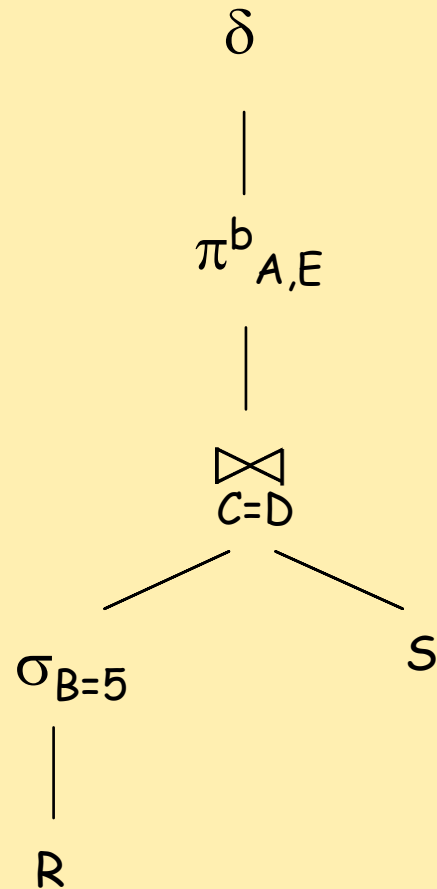
ASSUMPTIONS

- The tables do not have null values, and **have primary keys**:
 - A key constraint uniquely identifies each record in a table.
 - Tables are then **sets** of tuples
- Table in FROM clause have no attribute with the same name
- Queries are a single SELECT with possibly GROUP BY and HAVING but without subselect and ORDER BY clauses.
- Since superkeys are lifted after join and restriction, then $\sigma_C(R \bowtie S)$ is a **set** of tuples



EXAMPLE

$R(\underline{A}, B, C)$ $S(\underline{D}, E)$



$F = \{ A \rightarrow BC, D \rightarrow E, \dots \text{what else?} \dots \}$

$F = \{ A \rightarrow BC, D \rightarrow E, \{\} \rightarrow B, C \rightarrow D, D \rightarrow C \}$

$F \models AE \rightarrow AD ?$

$A \in \{A, E\}^+ \quad D \in \{A, E\}^+ ?$

Apply SLOW CLOSURE

$F \models AE \rightarrow AD \quad \text{YES}$

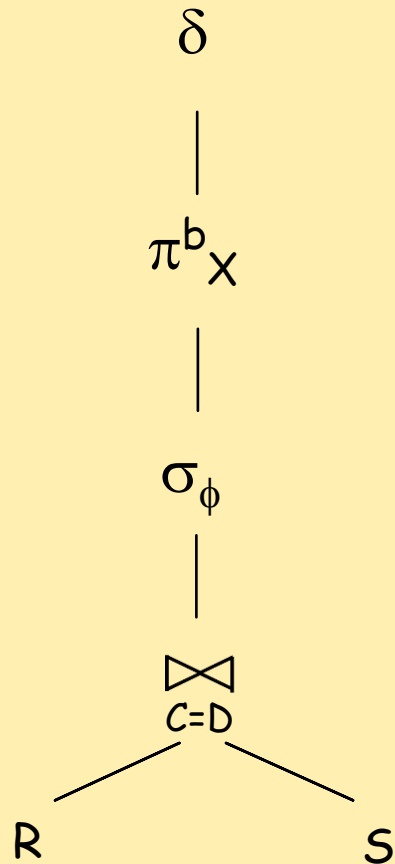
Is δ needed ???

In general: if $X \rightarrow Y$ and Y is a super-key then X is a super-key

Since AD is a superkey after the join, AE is a superkey!

δ is **NOT** needed

GENERALIZATION OF THE EXAMPLE



SELECT DISTINCT X
FROM R, S
WHERE C=D AND ϕ

When **DISTINCT** is useless ?

X is a (super)key in $\sigma_\phi(R \bowtie S)$

e.g., if X determines the keys of R and S

DERIVING FUNCTIONAL DEPENDENCIES IN SQL RESULTS

Which functional dependencies hold in the result of a query

SELECT * FROM-WHERE when all the tables in **FROM** have a key ?

1. Let F the initial set of FDs where their determinants are the keys of every table used in the query.
2. Let C the **WHERE** condition. If a conjunct of C is a predicate $A_i = c$, then F is extended with the functional dependency $\{ \} \rightarrow A_i$.
3. If a conjunct of C is a predicate $A_j = A_k$, e.g. a join condition, F is extended with the functional dependencies $A_j \rightarrow A_k$ and $A_k \rightarrow A_j$.

DERIVING FUNCTIONAL DEPENDENCIES IN SQL RESULTS

An algorithm to compute the closure of an attribute set X in $\sigma_{\phi}(R \bowtie S)$, which works **directly on SQL** without explicitly using functional dependencies.

1. Let $X^+ = X$
2. Add to X^+ all attributes A_i such that $A_i = c$ is a conjunct of the selection.
3. Repeat until X^+ is changed
 - a) Add to X^+ all attributes A_j such that predicate $A_j = A_k$ is a conjunct of the selection, and $A_k \in X^+$.
 - b) Add to X^+ all attributes of a table if X^+ contains a key for that table.

$R(\underline{A}, B, C)$ $S(\underline{D}, E)$

```
SELECT DISTINCT A, E  
FROM R, S  
WHERE C=D AND B=5
```

$\{A, E\}^+ = \{A, E, B, C, D\}$