

18. Relational Database Design

The relational datamodel was introduced by Codd in 1970. It is the most widely used datamodel—extended with the possibilities of the World Wide Web—, because of its simplicity and flexibility. The main idea of the relational model is that data is organised in relational tables, where rows correspond to individual *records* and columns to *attributes*. A *relational schema* consists of one or more relations and their attribute sets. In the present chapter only schemata consisting of one relation are considered for the sake of simplicity. In contrast to the mathematical concept of relations, in the relational schema the order of the attributes is not important, always *sets* of attributes are considered instead of *lists*. Every attribute has an associated *domain* that is a set of elementary values that the attribute can take values from. As an example, consider the following schema.

Employee(Name,Mother's name,Social Security Number,Post,Salary)

The domain of attributes *Name* and *Mother's name* is the set of finite character strings (more precisely its subset containing all possible names). The domain of *Social Security Number* is the set of integers satisfying certain formal and parity check requirements. The attribute *Post* can take values from the set {Director,Section chief,System integrator,Programmer,Receptionist,Janitor,Handyman}. An *instance* of a schema R is a relation r if its columns correspond to the attributes of R and its rows contain values from the domains of attributes at the attributes' positions. A typical row of a relation of the Employee schema could be

(John Brown,Camille Parker,184-83-2010,Programmer,\$172,000)

There can be dependencies between different data of a relation. For example, in an instance of the Employee schema the value of Social Security Number determines all other values of a row. Similarly, the pair (Name,Mother's name) is a unique identifier. Naturally, it may occur that some set of attributes do not determine all attributes of a record uniquely, just some of its subsets.

A relational schema has several *integrity constraints* attached. The most important kind of these is *functional dependency*. Let U and V be two sets of attributes. V *functionally depends* on U , $U \rightarrow V$ in notation, means that whenever two records are identical in the attributes belonging to U , then they must agree in the attribute belonging to V , as well. Throughout this chapter the attribute set $\{A_1, A_2, \dots, A_k\}$ is denoted by $A_1A_2 \dots A_k$ for the sake of convenience.

Example 18.1 *Functional dependencies* Consider the schema

$$R(\mathbf{P}\text{rofessor}, \mathbf{S}\text{ubject}, \mathbf{R}\text{oom}, \mathbf{S}\text{tudent}, \mathbf{G}\text{rade}, \mathbf{T}\text{ime}) .$$

The meaning of an individual record is that a given student got a given grade of a given subject that was taught by a given professor at a given time slot. The following functional dependencies are satisfied.

$\mathbf{S}\mathbf{u} \rightarrow \mathbf{P}$: One subject is taught by one professor.

$\mathbf{P}\mathbf{T} \rightarrow \mathbf{R}$: A professor teaches in one room at a time.

$\mathbf{S}\mathbf{t}\mathbf{T} \rightarrow \mathbf{R}$: A student attends a lecture in one room at a time.

$\mathbf{S}\mathbf{t}\mathbf{T} \rightarrow \mathbf{S}\mathbf{u}$: A student attends a lecture of one subject at a time.

$\mathbf{S}\mathbf{u}\mathbf{S}\mathbf{t} \rightarrow \mathbf{G}$: A student receives a unique final grade of a subject.

In Example 18.1 the attribute set $\mathbf{S}\mathbf{t}\mathbf{T}$ uniquely determines the values of all other attributes, furthermore it is minimal such set with respect to containment. This kind attribute sets are called *keys*. If all attributes are functionally dependent on a set of attributes X , then X is called a *superkey*. It is clear that every superkey contains a key and that any set of attributes containing a superkey is also a superkey.

18.1. Functional dependencies

Some functional dependencies valid for a given relational schema are known already in the design phase, others are consequences of these. The $\mathbf{S}\mathbf{t}\mathbf{T} \rightarrow \mathbf{P}$ dependency is implied by the $\mathbf{S}\mathbf{t}\mathbf{T} \rightarrow \mathbf{S}\mathbf{u}$ and $\mathbf{S}\mathbf{u} \rightarrow \mathbf{P}$ dependencies in Example 18.1. Indeed, if two records agree on attributes $\mathbf{S}\mathbf{t}$ and \mathbf{T} , then they must have the same value in attribute $\mathbf{S}\mathbf{u}$. Agreeing in $\mathbf{S}\mathbf{u}$ and $\mathbf{S}\mathbf{u} \rightarrow \mathbf{P}$ implies that the two records agree in \mathbf{P} , as well, thus $\mathbf{S}\mathbf{t}\mathbf{T} \rightarrow \mathbf{P}$ holds.

Definition 18.1 Let R be a relational schema, F be a set of functional dependencies over R . The functional dependency $U \rightarrow V$ is *logically implied* by F , in notation $F \models U \rightarrow V$, if each instance of R that satisfies all dependencies of F also satisfies $U \rightarrow V$. The *closure* of a set F of functional dependencies is the set F^+ given by

$$F^+ = \{U \rightarrow V : F \models U \rightarrow V\} .$$

18.1.1. Armstrong-axioms

In order to determine keys, or to understand logical implication between functional dependencies, it is necessary to know the closure F^+ of a set F of functional dependencies, or for a given $X \rightarrow Z$ dependency the question whether it belongs to F^+ must be decidable. For this, *inference rules* are needed that tell that from a set of functional dependencies what others follow. The *Armstrong-axioms* form a system of *sound* and *complete* inference rules. A system of rules is sound if only valid functional dependencies can be derived using it. It is complete, if every dependency $X \rightarrow Z$ that is logically implied by the set F is derivable from F using the inference rules.

ARMSTRONG-AXIOMS

- (A1) **Reflexivity** $Y \subseteq X \subseteq R$ implies $X \rightarrow Y$.
- (A2) **Augmentation** If $X \rightarrow Y$, then for arbitrary $Z \subseteq R$, $XZ \rightarrow YZ$ holds.
- (A3) **Transitivity** If $X \rightarrow Y$ and $Y \rightarrow Z$ hold, then $X \rightarrow Z$ holds, as well.

Example 18.2 *Derivation by the Armstrong-axioms.* Let $R = ABCD$ and $F = \{A \rightarrow C, B \rightarrow D\}$, then AB is a key:

1. $A \rightarrow C$ is given.
2. $AB \rightarrow ABC$ 1. is augmented by (A2) with AB .
3. $B \rightarrow D$ is given.
4. $ABC \rightarrow ABCD$ 3. is augmented by (A2) with ABC .
5. $AB \rightarrow ABCD$ transitivity (A3) is applied to 2. and 4..

Thus it is shown that AB is superkey. That it is really a key, follows from algorithm $\text{CLOSURE}(R, F, X)$.

There are other valid inference rules besides (A1)–(A3). The next lemma lists some, the proof is left to the Reader (Exercise 18.1-5).

Lemma 18.2

1. **Union rule** $\{X \rightarrow Y, X \rightarrow Z\} \models X \rightarrow YZ$.
2. **Pseudo transitivity** $\{X \rightarrow Y, WY \rightarrow Z\} \models XW \rightarrow YZ$.
3. **Decomposition** If $X \rightarrow Y$ holds and $Z \subseteq Y$, then $X \rightarrow Z$ holds, as well.

The soundness of system (A1)–(A3) can be proven by easy induction on the length of the derivation. The completeness will follow from the proof of correctness of algorithm $\text{CLOSURE}(R, F, X)$ by the following lemma. Let X^+ denote the **closure** of the set of attributes $X \subseteq R$ with respect to the family of functional dependencies F , that is $X^+ = \{A \in R: X \rightarrow A \text{ follows from } F \text{ by the Armstrong-axioms}\}$.

Lemma 18.3 *The functional dependency $X \rightarrow Y$ follows from the family of functional dependencies F by the Armstrong-axioms iff $Y \subseteq X^+$.*

Proof Let $Y = A_1A_2 \dots A_n$ where A_i 's are attributes, and assume that $Y \subseteq X^+$. $X \rightarrow A_i$ follows by the Armstrong-axioms for all i by the definition of X^+ . Applying the union rule of Lemma 18.2 $X \rightarrow Y$ follows. On the other hand, assume that $X \rightarrow Y$ can be derived by the Armstrong-axioms. By the decomposition rule of Lemma 18.2 $X \rightarrow A_i$ follows by (A1)–(A3) for all i . Thus, $Y \subseteq X^+$. ■

18.1.2. Closures

Calculation of closures is important in testing equivalence or logical implication between systems of functional dependencies. The first idea could be that for a given

family F of functional dependencies in order to decide whether $F \models \{X \rightarrow Y\}$, it is enough to calculate F^+ and check whether $\{X \rightarrow Y\} \in F^+$ holds. However, the size of F^+ could be exponential in the size of input. Consider the family F of functional dependencies given by

$$F = \{A \rightarrow B_1, A \rightarrow B_2, \dots, A \rightarrow B_n\}.$$

F^+ consists of all functional dependencies of the form $A \rightarrow Y$, where $Y \subseteq \{B_1, B_2, \dots, B_n\}$, thus $|F^+| = 2^n$. Nevertheless, the closure X^+ of an attribute set X with respect to F can be determined in linear time of the total length of functional dependencies in F . The following is an algorithm that calculates the closure X^+ of an attribute set X with respect to F . The input consists of the schema R , that is a finite set of attributes, a set F of functional dependencies defined over R , and an attribute set $X \subseteq R$.

CLOSURE(R, F, X)

```

1   $X^{(0)} \leftarrow X$ 
2   $i \leftarrow 0$ 
3   $G \leftarrow F$  ▷ Functional dependencies not used yet.
4  repeat
5      $X^{(i+1)} \leftarrow X^{(i)}$ 
6     for all  $Y \rightarrow Z$  in  $G$ 
7         do if  $Y \subseteq X^{(i)}$ 
8             then  $X^{(i+1)} \leftarrow X^{(i+1)} \cup Z$ 
9                  $G \leftarrow G \setminus \{Y \rightarrow Z\}$ 
10     $i \leftarrow i + 1$ 
11 until  $X^{(i-1)} = X^{(i)}$ 

```

It is easy to see that the attributes that are put into any of the $X^{(j)}$'s by CLOSURE(R, F, X) really belong to X^+ . The harder part of the correctness proof of this algorithm is to show that each attribute belonging to X^+ will be put into some of the $X^{(j)}$'s.

Theorem 18.4 CLOSURE(R, F, X) correctly calculates X^+ .

Proof First we prove by induction that if an attribute A is put into an $X^{(j)}$ during CLOSURE(R, F, X), then A really belongs to X^+ .

Base case: $j = 0$. In this case $A \in X$ and by reflexivity (A1) $A \in X^+$.

Induction step: Let $j > 0$ and assume that $X^{(j-1)} \subseteq X^+$. A is put into $X^{(j)}$, because there is a functional dependency $Y \rightarrow Z$ in F , where $Y \subseteq X^{(j-1)}$ and $A \in Z$. By induction, $Y \subseteq X^+$ holds, which implies using Lemma 18.3 that $X \rightarrow Y$ holds, as well. By transitivity (A3) $X \rightarrow Y$ and $Y \rightarrow Z$ implies $X \rightarrow Z$. By reflexivity (A1) and $A \in Z$, $Z \rightarrow A$ holds. Applying transitivity again, $X \rightarrow A$ is obtained, that is $A \in X^+$.

On the other hand, we show that if $A \in X^+$, then A is contained in the result of CLOSURE(R, F, X). Suppose in contrary that $A \in X^+$, but $A \notin X^{(i)}$, where $X^{(i)}$ is the result of CLOSURE(R, F, X). By the stop condition in line 9 this means

$X^{(i)} = X^{(i+1)}$. An instance r of the schema R is constructed that satisfies every functional dependency of F , but $X \rightarrow A$ does not hold in r if $A \notin X^{(i)}$. Let r be the following two-rowed relation:

Attributes of $X^{(i)}$				Other attributes			
1	1	...	1	1	1	...	1
1	1	...	1	0	0	...	0

Let us suppose that the above r violates a $U \rightarrow V$ functional dependency of F , that is $U \subseteq X^{(i)}$, but V is not a subset of $X^{(i)}$. However, in this case $\text{CLOSURE}(R, F, X)$ could not have stopped yet, since $X^{(i)} \text{NEX}^{(i+1)}$.

$A \in X^+$ implies using Lemma 18.3 that $X \rightarrow A$ follows from F by the Armstrong-axioms. (A1)–(A3) is a sound system of inference rules, hence in every instance that satisfies F , $X \rightarrow A$ must hold. However, the only way this could happen in instance r is if $A \in X^{(i)}$. ■

Let us observe that the relation instance r given in the proof above provides the completeness proof for the Armstrong-axioms, as well. Indeed, the closure X^+ calculated by $\text{CLOSURE}(R, F, X)$ is the set of those attributes for which $X \rightarrow A$ follows from F by the Armstrong-axioms. Meanwhile, for every other attribute B , there exist two rows of r that agree on X , but differ in B , that is $F \models X \rightarrow B$ *does not* hold.

The running time of $\text{CLOSURE}(R, F, X)$ is $O(n^2)$, where n is the length of the input. Indeed, in the **repeat – until** loop of lines 4–11 every not yet used dependency is checked, and the body of the loop is executed at most $|R \setminus X| + 1$ times, since it is started again only if $X^{(i-1)} \text{NEX}^{(i)}$, that is a new attribute is added to the closure of X . However, the running time can be reduced to linear with appropriate bookkeeping.

1. For every yet unused $W \rightarrow Z$ dependency of F it is kept track of how many attributes of W are not yet included in the closure ($i[W, Z]$).
2. For every attribute A those yet unused dependencies are kept in a doubly linked list L_A whose left side contains A .
3. Those not yet used dependencies $W \rightarrow Z$ are kept in a linked list J , whose left side W 's every attribute is contained in the closure already, that is for which $i[W, Z] = 0$.

It is assumed that the family of functional dependencies F is given as a set of attribute pairs (W, Z) , representing $W \rightarrow Z$. The $\text{LINEAR-CLOSURE}(R, F, X)$ algorithm is a modification of $\text{CLOSURE}(R, F, X)$ using the above bookkeeping, whose running time is linear. R is the schema, F is the given family of functional dependencies, and we are to determine the closure of attribute set X .

Algorithm $\text{LINEAR-CLOSURE}(R, F, X)$ consists of two parts. In the initialisation phase (lines 1–13) the lists are initialised. The loops of lines 2–5 and 6–8, respectively,

take $O(\sum_{(W,Z) \in F} |W|)$ time. The loop in lines 9–11 means $O(|F|)$ steps. If the length of the input is denoted by n , then this is $O(n)$ steps altogether.

During the execution of lines 14–23, every functional dependency (W, Z) is examined at most once, when it is taken off from list J . Thus, lines 15–16 and 23 take at most $|F|$ steps. The running time of the loops in line 17–22 can be estimated by observing that the sum $\sum i[W, Z]$ is decreased by one in each execution, hence it takes $O(\sum i_0[W, Z])$ steps, where $i_0[W, Z]$ is the $i[W, Z]$ value obtained in the initialisation phase. However, $\sum i_0[W, Z] \leq \sum_{(W,Z) \in F} |W|$, thus lines 14–23 also take $O(n)$ time in total.

LINEAR-CLOSURE(R, F, X)

```

1                                     ▷ Initialisation phase.
2 for all  $(W, Z) \in F$ 
3     do for all  $A \in W$ 
4         do add  $(W, Z)$  to list  $L_A$ 
5      $i[W, Z] \leftarrow 0$ 
6 for all  $A \in R \setminus X$ 
7     do for all  $(W, Z)$  of list  $L_A$ 
8         do  $i[W, Z] \leftarrow i[W, Z] + 1$ 
9 for all  $(W, Z) \in F$ 
10    do if  $i[W, Z] = 0$ 
11        then add  $(W, Z)$  to list  $J$ 
12  $X^+ \leftarrow X$ 
13                                     ▷ End of initialisation phase.
14 while  $J$  is nonempty
15    do  $(W, Z) \leftarrow \text{head}(J)$ 
16        delete  $(W, Z)$  from list  $J$ 
17    for all  $A \in Z \setminus X^+$ 
18        do for all  $(W, Z)$  of list  $L_A$ 
19            do  $i[W, Z] \leftarrow i[W, Z] - 1$ 
20            if  $i[W, Z] = 0$ 
21                then add  $(W, Z)$  to list  $J$ 
22                delete  $(W, Z)$  from list  $L_A$ 
23     $X^+ \leftarrow X^+ \cup Z$ 
24 return  $X^+$ 

```

18.1.3. Minimal cover

Algorithm LINEAR-CLOSURE(R, F, X) can be used to test equivalence of systems of dependencies. Let F and G be two families of functional dependencies. F and G are said to be *equivalent*, if exactly the same functional dependencies follow from both, that is $F^+ = G^+$. It is clear that it is enough to check for all functional dependencies $X \rightarrow Y$ in F whether it belongs to G^+ , and vice versa, for all $W \rightarrow Z$ in G , whether it is in F^+ . Indeed, if some of these is not satisfied, say $X \rightarrow Y$ is not in G^+ , then surely $F^+ \text{NEG}^+$. On the other hand, if all $X \rightarrow Y$ are in G^+ , then a proof

of a functional dependency $U \rightarrow V$ from F^+ can be obtained from dependencies in G in such a way that to the derivation of the dependencies $X \rightarrow Y$ of F from G , the derivation of $U \rightarrow V$ from F is concatenated. In order to decide that a dependency $X \rightarrow Y$ from F is in G^+ , it is enough to construct the closure $X^+(G)$ of attribute set X with respect to G using $\text{LINEAR-CLOSURE}(R, G, X)$, then check whether $Y \subseteq X^+(G)$ holds. The following special functional dependency system equivalent with F is useful.

Definition 18.5 *The system of functional dependencies G is a **minimal cover** of the family of functional dependencies F iff G is equivalent with F , and*

1. *functional dependencies of G are in the form $X \rightarrow A$, where A is an attribute and $A \notin X$,*
2. *no functional dependency can be dropped from G , i.e., $(G - \{X \rightarrow A\})^+ \subsetneq G^+$,*
3. *the left sides of dependencies in G are minimal, that is $X \rightarrow A \in G$, $Y \subsetneq X \implies ((G - \{X \rightarrow A\}) \cup \{Y \rightarrow A\})^+ \neq G^+$.*

Every set of functional dependencies have a minimal cover, namely algorithm $\text{MINIMAL-COVER}(R, F)$ constructs one.

$\text{MINIMAL-COVER}(R, F)$

```

1  $G \leftarrow \emptyset$ 
2 for all  $X \rightarrow Y \in F$ 
3   do for all  $A \in Y - X$ 
4     do  $G \leftarrow G \cup X \rightarrow A$ 
5                                      $\triangleright$  Each right hand side consists of a single attribute.
6 for all  $X \rightarrow A \in G$ 
7   do while there exists  $B \in X: A \in (X - B)^+(G)$ 
8      $X \leftarrow X - B$ 
9                                      $\triangleright$  Each left hand side is minimal.
10 for all  $X \rightarrow A \in G$ 
11   do if  $A \in X^+(G - \{X \rightarrow A\})$ 
12     then  $G \leftarrow G - \{X \rightarrow A\}$ 
13                                      $\triangleright$  No redundant dependency exists.
```

After executing the loop of lines 2–4, the right hand side of each dependency in G consists of a single attribute. The equality $G^+ = F^+$ follows from the union rule of Lemma 18.2 and the reflexivity axiom. Lines 6–8 minimise the left hand sides. In line 11 it is checked whether a given functional dependency of G can be removed without changing the closure. $X^+(G - \{X \rightarrow A\})$ is the closure of attribute set X with respect to the family of functional dependencies $G - \{X \rightarrow A\}$.

Claim 18.6 $\text{MINIMAL-COVER}(R, F)$ *calculates a minimal cover of F .*

Proof It is enough to show that during execution of the loop in lines 10–12, no functional dependency $X \rightarrow A$ is generated whose left hand side could be decreased.

Indeed, if a $X \rightarrow A$ dependency would exist, such that for some $Y \subsetneq X$ $Y \rightarrow A \in G^+$ held, then $Y \rightarrow A \in G'^+$ would also hold, where G' is the set of dependencies considered when $X \rightarrow A$ is checked in lines 6–8. $G \subseteq G'$, which implies $G^+ \subseteq G'^+$ (see Exercise 18.1-1). Thus, X should have been decreased already during execution of the loop in lines 6–8. ■

18.1.4. Keys

In database design it is important to identify those attribute sets that uniquely determine the data in individual records.

Definition 18.7 Let (R, F) be a relational schema. The set of attributes $X \subseteq R$ is called a **superkey**, if $X \rightarrow R \in F^+$. A superkey X is called a **key**, if it is minimal with respect to containment, that is no proper subset $Y \subsetneq X$ is key.

The question is how the keys can be determined from (R, F) ? What makes this problem hard is that the number of keys could be super exponential function of the size of (R, F) . In particular, Yu and Johnson constructed such relational schema, where $|F| = k$, but the number of keys is $k!$. Békéssy and Demetrovics gave a beautiful and simple proof of the fact that starting from k functional dependencies, at most $k!$ key can be obtained. (This was independently proved by Osborne and Tompa.)

The proof of Békéssy and Demetrovics is based on the operation $*$ they introduced, which is defined for functional dependencies.

Definition 18.8 Let $e_1 = U \rightarrow V$ and $e_2 = X \rightarrow Y$ be two functional dependencies. The binary operation $*$ is defined by

$$e_1 * e_2 = U \cup ((R - V) \cap X) \rightarrow V \cup Y.$$

Some properties of operation $*$ is listed, the proof is left to the Reader (Exercise 18.1-3). Operation $*$ is associative, furthermore it is idempotent in the sense that if $e = e_1 * e_2 * \dots * e_k$ and $e' = e * e_i$ for some $1 \leq i \leq k$, then $e' = e$.

Claim 18.9 (Békéssy and Demetrovics). Let (R, F) be a relational schema and let $F = \{e_1, e_2, \dots, e_k\}$ be a listing of the functional dependencies. If X is a key, then $X \rightarrow R = e_{\pi_1} * e_{\pi_2} * \dots * e_{\pi_s} * d$, where $(\pi_1, \pi_2, \dots, \pi_s)$ is an ordered subset of the index set $\{1, 2, \dots, k\}$, and d is a trivial dependency in the form $D \rightarrow D$.

Proposition 18.9 bounds in some sense the possible sets of attributes in the search for keys. The next proposition gives lower and upper bounds for the keys.

Claim 18.10 Let (R, F) be a relational schema and let $F = \{U_i \rightarrow V_i : 1 \leq i \leq k\}$. Let us assume without loss of generality that $U_i \cap V_i = \emptyset$. Let $\mathcal{U} = \bigcup_{i=1}^k U_i$ and $\mathcal{V} = \bigcup_{i=1}^k V_i$. If K is a key in the schema (R, F) , then

$$\mathcal{H}_L = R - \mathcal{V} \subseteq K \subseteq (R - \mathcal{V}) \cup \mathcal{U} = \mathcal{H}_U.$$

The proof is not too hard, it is left as an exercise for the Reader (Exercise 18.1-4). The algorithm LIST-KEYS(R, F) that lists the keys of the schema (R, F) is based on the bounds of Proposition 18.10. The running time can be bounded by $O(n!)$, but one cannot expect any better, since to list the output needs that much time in worst case.

LIST-KEYS(R, F)

```

1                                     ▷ Let  $\mathcal{U}$  and  $\mathcal{V}$  be as defined in Proposition 18.10
2 if  $\mathcal{U} \cap \mathcal{V} = \emptyset$ 
3   then return  $R - \mathcal{V}$ 
4                                     ▷  $R - \mathcal{V}$  is the only key.
5 if  $(R - \mathcal{V})^+ = R$ 
6   then return  $R - \mathcal{V}$ 
7                                     ▷  $R - \mathcal{V}$  is the only key.
8  $\mathcal{K} \leftarrow \emptyset$ 
9 for all permutations  $A_1, A_2, \dots, A_h$  of the attributes of  $\mathcal{U} \cap \mathcal{V}$ 
10  do  $K \leftarrow (R - \mathcal{V}) \cup \mathcal{U}$ 
11    for  $i \leftarrow 1$  to  $h$ 
12      do  $Z \leftarrow K - A_i$ 
13        if  $Z^+ = R$ 
14          then  $K \leftarrow Z$ 
15     $\mathcal{K} \leftarrow \mathcal{K} \cup \{K\}$ 
16 return  $\mathcal{K}$ 

```

Exercises

18.1-1 Let R be a relational schema and let F and G be families of functional dependencies over R . Show that

- $F \subseteq F^+$.
- $(F^+)^+ = F^+$.
- If $F \subseteq G$, then $F^+ \subseteq G^+$.

Formulate and prove similar properties of the closure X^+ – with respect to F – of an attribute set X .

18.1-2 Derive the functional dependency $AB \rightarrow F$ from the set of dependencies $G = \{AB \rightarrow C, A \rightarrow D, CD \rightarrow EF\}$ using Armstrong-axioms (A1)–(A3).

18.1-3 Show that operation $*$ is associative, furthermore if for functional dependencies e_1, e_2, \dots, e_k we have $e = e_1 * e_2 * \dots * e_k$ and $e' = e * e_i$ for some $1 \leq i \leq k$, then $e' = e$.

18.1-4 Prove Proposition 18.10.

18.1-5 Prove the union, pseudo transitivity and decomposition rules of Lemma 18.2.

18.2. Decomposition of relational schemata

A *decomposition* of a relational schema $R = \{A_1, A_2, \dots, A_n\}$ is a collection $\rho = \{R_1, R_2, \dots, R_k\}$ of subsets of R such that

$$R = R_1 \cup R_2 \cup \dots \cup R_k .$$

The R_i 's need not be disjoint, in fact in most application they must not be. One important motivation of decompositions is to avoid *anomalies*.

Example 18.3 *Anomalies* Consider the following schema

SUPPLIER-INFO(SNAME,ADDRESS,ITEM,PRICE)

This schema encompasses the following problems:

1. *Redundancy*. The address of a supplier is recorded with every item it supplies.
2. *Possible inconsistency (update anomaly)*. As a consequence of redundancy, the address of a supplier might be updated in some records and might not be in some others, hence the supplier would not have a unique address, even though it is expected to have.
3. *Insertion anomaly*. The address of a supplier cannot be recorded if it does not supply anything at the moment. One could try to use NULL values in attributes ITEM and PRICE, but would it be remembered that it must be deleted, when a supplied item is entered for that supplier? More serious problem that SNAME and ITEM together form a key of the schema, and the NULL values could make it impossible to search by an index based on that key.
4. *Deletion anomaly* This is the opposite of the above. If all items supplied by a supplier are deleted, then as a side effect the address of the supplier is also lost.

All problems mentioned above are eliminated if schema SUPPLIER-INFO is replaced by two sub-schemata:

SUPPLIER(SNAME,ADDRESS),
SUPPLIES(SNAME,ITEM,PRICE).

In this case each suppliers address is recorded only once, and it is not necessary that the supplier supplies a item in order its address to be recorded. For the sake of convenience the attributes are denoted by single characters S (SNAME), A (ADDRESS), I (ITEM), P (PRICE).

Question is that is it correct to replace the schema $SAIP$ by SA and SIP ? Let r be an instance of schema $SAIP$. It is natural to require that if SA and SIP is used, then the relations belonging to them are obtained projecting r to SA and SIP , respectively, that is $r_{SA} = \pi_{SA}(r)$ and $r_{SIP} = \pi_{SIP}(r)$. r_{SA} and r_{SIP} contains the same information as r , if r can be reconstructed using only r_{SA} and r_{SIP} . The calculation of r from r_{SA} and r_{SIP} can be done by the *natural join* operator.

Definition 18.11 The *natural join* of relations r_i of schemata R_i ($i = 1, 2, \dots, n$) is the relation s belonging to the schema $\cup_{i=1}^n R_i$, which consists of all rows μ that for all i there exists a row ν_i of relation r_i such that $\mu[R_i] = \nu_i[R_i]$. In notation $s = \bowtie_{i=1}^n r_i$.

Example 18.4 Let $R_1 = AB$, $R_2 = BC$, $r_1 = \{ab, a'b', ab''\}$ and $r_2 = \{bc, bc', b'c''\}$. The natural join of r_1 and r_2 belongs to the schema $R = ABC$, and it is the relation $r_1 \bowtie r_2 = \{abc, abc', a'b'c''\}$.

If s is the natural join of r_{SA} and r_{SIP} , that is $s = r_{SA} \bowtie r_{SIP}$, then $\pi_{SA}(s) = r_{SA}$ és $\pi_{SIP}(s) = r_{SIP}$ by Lemma 18.12. If r NEs, then the original relation could not be reconstructed knowing only r_{SA} and r_{SIP} .

18.2.1. Lossless join

Let $\rho = \{R_1, R_2, \dots, R_k\}$ be a decomposition of schema R , furthermore let F be a family of functional dependencies over R . The decomposition ρ is said to have **lossless join property** (with respect to F), if every instance r of R that satisfies F also satisfies

$$r = \pi_{R_1}(r) \bowtie \pi_{R_2}(r) \bowtie \dots \bowtie \pi_{R_k}(r).$$

That is, relation r is the natural join of its projections to attribute sets R_i , $i = 1, 2, \dots, k$. For a decomposition $\rho = \{R_1, R_2, \dots, R_k\}$, let m_ρ denote the mapping which assigns to relation r the relation $m_\rho(r) = \bowtie_{i=1}^k \pi_{R_i}(r)$. Thus, the lossless join property with respect to a family of functional dependencies means that $r = m_\rho(r)$ for all instances r that satisfy F .

Lemma 18.12 Let $\rho = \{R_1, R_2, \dots, R_k\}$ be a decomposition of schema R , and let r be an arbitrary instance of R . Furthermore, let $r_i = \pi_{R_i}(r)$. Then

1. $r \subseteq m_\rho(r)$.
2. If $s = m_\rho(r)$, then $\pi_{R_i}(s) = r_i$.
3. $m_\rho(m_\rho(r)) = m_\rho(r)$.

The proof is left to the Reader (Exercise 18.2-7).

18.2.2. Checking the lossless join property

It is relatively not hard to check that a decomposition $\rho = \{R_1, R_2, \dots, R_k\}$ of schema R has the lossless join property. The essence of algorithm JOIN-TEST(R, F, ρ) is the following.

A $k \times n$ array T is constructed, whose column j corresponds to attribute A_j , while row i corresponds to schema R_i . $T[i, j] = 0$ if $A_j \in R_i$, otherwise $T[i, j] = i$.

The following step is repeated until there is no more possible change in the array. Consider a functional dependency $X \rightarrow Y$ from F . If a pair of rows i and j agree in all attributes of X , then their values in attributes of Y are made equal. More precisely, if one of the values in an attribute of Y is 0, then the other one is set to 0, as well, otherwise it is arbitrary which of the two values is set to be equal to the other one. If a symbol is changed, then **each** of its occurrences in that column must be changed accordingly. If at the end of this process there is an all 0 row in T , then the decomposition has the lossless join property, otherwise, it is lossy.

JOIN-TEST(R, F, ρ)

```

1                                     ▷ Initialisation phase.
2 for  $i \leftarrow 1$  to  $|\rho|$ 
3   do for  $j \leftarrow 1$  to  $|R|$ 
4     do if  $A_j \in R_i$ 
5       then  $T[i, j] \leftarrow 0$ 
6       else  $T[i, j] \leftarrow i$ 
7                                     ▷ End of initialisation phase.
8  $S \leftarrow T$ 
9 repeat
10    $T \leftarrow S$ 
11   for all  $\{X \rightarrow Y\} \in F$ 
12     do for  $i \leftarrow 1$  to  $|\rho| - 1$ 
13       do for  $j \leftarrow i + 1$  to  $|R|$ 
14         do if for all  $A_h$  in  $X$  ( $S[i, h] = S[j, h]$ )
15           then EQUATE( $i, j, S, Y$ )
16 until  $S = T$ 
17 if there exist an all 0 row in  $S$ 
18   then return “Lossless join”
19   else return “Lossy join”

```

Procedure EQUATE(i, j, S, Y) makes the appropriate symbols equal.

EQUATE(i, j, S, Y)

```

1 for  $A_l \in Y$ 
2   do if  $S[i, l] \cdot S[j, l] = 0$ 
3     then
4       for  $d \leftarrow 1$  to  $k$ 
5         do if  $S[d, l] = S[i, l] \vee S[d, l] = S[j, l]$ 
6           then  $S[d, l] \leftarrow 0$ 
7     else
8       for  $d \leftarrow 1$  to  $k$ 
9         do if  $S[d, l] = S[j, l]$ 
10          then  $S[d, l] \leftarrow S[i, l]$ 

```

Example 18.5 *Checking lossless join property* Let $R = ABCDE$, $R_1 = AD$, $R_2 = AB$, $R_3 = BE$, $R_4 = CDE$, $R_5 = AE$, furthermore let the functional dependencies be $\{A \rightarrow C, B \rightarrow C, C \rightarrow D, DE \rightarrow C, CE \rightarrow A\}$. The initial array is shown on Figure 18.1(a). Using $A \rightarrow C$ values 1,2,5 in column C can be equated to 1. Then applying $B \rightarrow C$ value 3 of column C again can be changed to 1. The result is shown on Figure 18.1(b). Now $C \rightarrow D$ can be used to change values 2,3,5 of column D to 0. Then applying $DE \rightarrow C$ (the only nonzero) value 1 of column C can be set to 0. Finally, $CE \rightarrow A$ makes it possible to change values 3 and 4 in column A to be changed to 0. The final result is shown on Figure 18.1(c). The third row consists of only zeroes, thus the decomposition has the lossless join property.

A	B	C	D	E
0	1	1	0	1
0	0	2	2	2
3	0	3	3	0
4	4	0	0	0
0	5	5	5	0

(a)

A	B	C	D	E
0	1	1	0	1
0	0	1	2	2
3	0	1	3	0
4	4	0	0	0
0	5	1	5	0

(b)

A	B	C	D	E
0	1	0	0	1
0	0	0	2	2
0	0	0	0	0
0	4	0	0	0
0	5	0	0	0

(c)

Figure 18.1 Application of $\text{JOIN-TEST}(R, F, \rho)$.

It is clear that the running time of algorithm $\text{JOIN-TEST}(R, F, \rho)$ is polynomial in the length of the input. The important thing is that it uses only the schema, not the instance r belonging to the schema. Since the size of an instance is larger than the size of the schema by many orders of magnitude, the running time of an algorithm using the schema only is negligible with respect to the time required by an algorithm processing the data stored.

Theorem 18.13 *Procedure $\text{JOIN-TEST}(R, F, \rho)$ correctly determines whether a given decomposition has the lossless join property.*

Proof Let us assume first that the resulting array T contains no all zero row. T itself can be considered as a relational instance over the schema R . This relation satisfies all functional dependencies from F , because the algorithm finished since there was no more change in the table during checking the functional dependencies. It is true for the starting table that its projections to every R_i 's contain an all zero row, and this property does not change during the running of the algorithm, since a 0 is never changed to another symbol. It follows, that the natural join $m_\rho(T)$ contains the all zero row, that is $\text{TNE}m_\rho(T)$. Thus the decomposition is lossy. The proof of the other direction is only sketched.

Logic, domain calculus is used. The necessary definitions can be found in the books of Abiteboul, Hull and Vianu, or Ullman, respectively. Imagine that variable a_j is written in place of zeroes, and b_{ij} is written in place of i 's in column j , and

JOIN-TEST(R, F, ρ) is run in this setting. The resulting table contains row $a_1 a_2 \dots a_n$, which corresponds to the all zero row. Every table can be viewed as a shorthand notation for the following domain calculus expression

$$\{a_1 a_2 \dots a_n \mid (\exists b_{11}) \dots (\exists b_{kn}) (R(w_1) \wedge \dots \wedge R(w_k))\}, \quad (18.1)$$

where w_i is the i th row of T . If T is the starting table, then formula (18.1) defines m_ρ exactly. As a justification note that for a relation r , $m_\rho(r)$ contains the row $a_1 a_2 \dots a_n$ iff r contains for all i a row whose j th coordinate is a_j if A_j is an attribute of R_i , and arbitrary values represented by variables b_{il} in the other attributes.

Consider an arbitrary relation r belonging to schema R that satisfies the dependencies of F . The modifications (equating symbols) of the table done by JOIN-TEST(R, F, ρ) do not change the set of rows obtained from r by (18.1), if the modifications are done in the formula, as well. Intuitively it can be seen from the fact that only such symbols are equated in (18.1), that can only take equal values in a relation satisfying functional dependencies of F . The exact proof is omitted, since it is quiet tedious.

Since in the result table of JOIN-TEST(R, F, ρ) the all a 's row occurs, the domain calculus formula that belongs to this table is of the following form:

$$\{a_1 a_2 \dots a_n \mid (\exists b_{11}) \dots (\exists b_{kn}) (R(a_1 a_2 \dots a_n) \wedge \dots)\}. \quad (18.2)$$

It is obvious that if (18.2) is applied to relation r belonging to schema R , then the result will be a subset of r . However, if r satisfies the dependencies of F , then (18.2) calculates $m_\rho(r)$. According to Lemma 18.12, $r \subseteq m_\rho(r)$ holds, thus if r satisfies F , then (18.2) gives back r exactly, so $r = m_\rho(r)$, that is the decomposition has the lossless join property. ■

Procedure JOIN-TEST(R, F, ρ) can be used independently of the number of parts occurring in the decomposition. The price of this generality is paid in the running time requirement. However, if R is to be decomposed only into *two* parts, then CLOSURE(R, F, X) or LINEAR-CLOSURE(R, F, X) can be used to obtain the same result faster, according to the next theorem.

Theorem 18.14 *Let $\rho = (R_1, R_2)$ be a decomposition of R , furthermore let F be a set of functional dependencies. Decomposition ρ has the lossless join property with respect to F iff*

$$(R_1 \cap R_2) \rightarrow (R_1 - R_2) \text{ or } (R_1 \cap R_2) \rightarrow (R_2 - R_1).$$

These dependencies need not be in F , it is enough if they are in F^+ .

Proof The starting table in procedure JOIN-TEST(R, F, ρ) is the following:

	$R_1 \cap R_2$	$R_1 - R_2$	$R_2 - R_1$	
row of R_1	00...0	00...0	11...1	(18.3)
row of R_2	00...0	22...2	00...0	

It is not hard to see using induction on the number of steps done by JOIN-TEST(R, F, ρ) that if the algorithm changes both values of the column of an attribute

A to 0, then $A \in (R_1 \cap R_2)^+$. This is obviously true at the start. If at some time values of column A must be equated, then by lines 11–14 of the algorithm, there exists $\{X \rightarrow Y\} \in F$, such that the two rows of the table agree on X , and $A \in Y$. By the induction assumption $X \subseteq (R_1 \cap R_2)^+$ holds. Applying Armstrong-axioms (transitivity and reflexivity), $A \in (R_1 \cap R_2)^+$ follows.

On the other hand, let us assume that $A \in (R_1 \cap R_2)^+$, that is $(R_1 \cap R_2) \rightarrow A$. Then this functional dependency can be derived from F using Armstrong-axioms. By induction on the length of this derivation it can be seen that procedure JOIN-TEST(R, F, ρ) will equate the two values of column A , that is set them to 0. Thus, the row of R_1 will be all 0 iff $(R_1 \cap R_2) \rightarrow (R_2 - R_1)$, similarly, the row of R_2 will be all 0 iff $(R_1 \cap R_2) \rightarrow (R_1 - R_2)$. ■

18.2.3. Dependency preserving decompositions

The lossless join property is important so that a relation can be recovered from its projections. In practice, usually not the relation r belonging to the underlying schema R is stored, but relations $r_i = r[R_i]$ for an appropriate decomposition $\rho = (R_1, R_2, \dots, R_k)$, in order to avoid anomalies. The functional dependencies F of schema R are *integrity constraints* of the database, relation r is consistent if it satisfies all prescribed functional dependencies. When during the life time of the database updates are executed, that is rows are inserted into or deleted from the projection relations, then it may happen that the natural join of the new projections does not satisfy the functional dependencies of F . It would be too costly to join the projected relations – and then project them again – after each update to check the integrity constraints. However, the *projection* of the family of functional dependencies F to an attribute set Z can be defined: $\pi_Z(F)$ consists of those functional dependencies $\{X \rightarrow Y\} \in F^+$, where $XY \subseteq Z$. After an update, if relation r_i is changed, then it is relatively easy to check whether $\pi_{R_i}(F)$ still holds. Thus, it would be desired if family F would be logical implication of the families of functional dependencies $\pi_{R_i}(F)$ $i = 1, 2, \dots, k$. Let $\pi_\rho(F) = \bigcup_{i=1}^k \pi_{R_i}(F)$.

Definition 18.15 *The decomposition ρ is said to be **dependency preserving**, if*

$$\pi_\rho(F)^+ = F^+.$$

Note that $\pi_\rho(F) \subseteq F^+$, hence $\pi_\rho(F)^+ \subseteq F^+$ always holds. Consider the following example.

Example 18.6 Let $R = (\text{City}, \text{Street}, \text{Zip code})$ be the underlying schema, furthermore let $F = \{CS \rightarrow Z, Z \rightarrow C\}$ be the functional dependencies. Let the decomposition ρ be $\rho = (CZ, SZ)$. This has the lossless join property by Theorem 18.14. $\pi_\rho(F)$ consists of $Z \rightarrow C$ besides the trivial dependencies. Let $R_1 = CZ$ and $R_2 = SZ$. Two rows are inserted into each of the projections belonging to schemata R_1 and R_2 , respectively, so that functional dependencies of the projections are satisfied:

R_1	C	Z	R_2	S	Z
	Fort Wayne	46805		Coliseum Blvd	46805
	Fort Wayne	46815		Coliseum Blvd	46815

In this case R_1 and R_2 satisfy the dependencies prescribed for them separately, however in $R_1 \bowtie R_2$ the dependency $CS \rightarrow Z$ does not hold.

It is true as well, that none of the decompositions of this schema preserves the dependency $CS \rightarrow Z$. Indeed, this is the only dependency that contains Z on the right hand side, thus if it is to be preserved, then there has to be a subschema that contains C, S, Z , but then the decomposition would not be proper. This will be considered again when decomposition into normal forms is treated.

Note that it may happen that decomposition ρ preserves functional dependencies, but does not have the lossless join property. Indeed, let $R = ABCD$, $F = \{A \rightarrow B, C \rightarrow D\}$, and let the decomposition be $\rho = (AB, CD)$.

Theoretically it is very simple to check whether a decomposition $\rho = (R_1, R_2, \dots, R_k)$ is dependency preserving. Just F^+ needs to be calculated, then projections need to be taken, finally one should check whether the union of the projections is equivalent with F . The main problem with this approach is that even calculating F^+ may need exponential time.

Nevertheless, the problem can be solved without explicitly determining F^+ . Let $G = \pi_\rho(F)$. G will not be calculated, only its equivalence with F will be checked. For this end, it needs to be decidable for all functional dependencies $\{X \rightarrow Y\} \in F^+$ that if X^+ is taken with respect to G , whether it contains Y . The trick is that X^+ is determined *without* full knowledge of G by repeatedly taking the effect to the closure of the projections of F onto the individual R_i 's. That is, the concept of S -operation on an attribute set Z is introduced, where S is another set of attributes: Z is replaced by $Z \cup ((Z \cap S)^+ \cap S)$, where the closure is taken with respect to F . Thus, the closure of the part of Z that lies in S is taken with respect to F , then from the resulting attributes those are added to Z , which also belong to S .

It is clear that the running time of algorithm $\text{PRESERVE}(\rho, F)$ is polynomial in the length of the input. More precisely, the outermost **for** loop is executed at most once for each dependency in F (it may happen that it turns out earlier that some dependency is not preserved). The body of the **repeat-until** loop in lines 3–7. requires linear number of steps, it is executed at most $|R|$ times. Thus, the body of the **for** loop needs quadratic time, so the total running time can be bounded by the cube of the input length.

$\text{PRESERVE}(\rho, F)$

```

1  for all  $(X \rightarrow Y) \in F$ 
2      do  $Z \leftarrow X$ 
3      repeat
4           $W \leftarrow Z$ 
5          for  $i \leftarrow 1$  to  $k$ 
6              do  $Z \leftarrow Z \cup (\text{LINEAR-CLOSURE}(R, F, Z \cap R_i) \cap R_i)$ 
7          until  $Z = W$ 
8      if  $Y \not\subseteq Z$ 
9          then return "Not dependency preserving"
10 return "Dependency preserving"
```


Example 18.7 Consider the schema $R = ABCD$, let the decomposition be $\rho = \{AB, BC, CD\}$, and dependencies be $F = \{A \rightarrow B, B \rightarrow C, C \rightarrow D, D \rightarrow A\}$. That is, by the visible cycle of the dependencies, every attribute determines all others. Since D and A do not occur together in the decomposition one might think that the dependency $D \rightarrow A$ is not preserved, however this intuition is wrong. The reason is that during the projection to AB , not only the dependency $A \rightarrow B$ is obtained, but $B \rightarrow A$, as well, since not F , but F^+ is projected. Similarly, $C \rightarrow B$ and $D \rightarrow C$ are obtained, as well, but $D \rightarrow A$ is a logical implication of these by the transitivity of the Armstrong axioms. Thus it is expected that $\text{PRESERVE}(\rho, F)$ claims that $D \rightarrow A$ is preserved.

Start from the attribute set $Y = \{D\}$. There are three possible operations, the AB -operation, the BC -operation and the CD -operation. The first two obviously does not add anything to $\{D\}^+$, since $\{D\} \cap \{A, B\} = \{D\} \cap \{B, C\} = \emptyset$, that is the closure of the empty set should be taken, which is empty (in the present example). However, using the CD -operation:

$$\begin{aligned} Z &= \{D\} \cup ((\{D\} \cap \{C, D\})^+ \cap \{C, D\}) \\ &= \{D\} \cup (\{D\}^+ \cap \{C, D\}) \\ &= \{D\} \cup (\{A, B, C, D\} \cap \{C, D\}) \\ &= \{C, D\}. \end{aligned}$$

In the next round using the BC -operation the actual $Z = \{C, D\}$ is changed to $Z = \{B, C, D\}$, finally applying the AB -operation on this, $Z = \{A, B, C, D\}$ is obtained. This cannot change, so procedure $\text{PRESERVE}(\rho, F)$ stops. Thus, with respect to the family of functional dependencies

$$G = \pi_{AB}(F) \cup \pi_{BC}(F) \cup \pi_{CD}(F) ,$$

$\{D\}^+ = \{A, B, C, D\}$ holds, that is $G \models D \rightarrow A$. It can be checked similarly that the other dependencies of F are in G^+ (as a fact in G).

Theorem 18.16 *The procedure $\text{PRESERVE}(\rho, F)$ determines correctly whether the decomposition ρ is dependency preserving.*

Proof It is enough to check for a single functional dependency $X \rightarrow Y$ whether whether the procedure decides correctly if it is in G^+ . When an attribute is added to Z in lines 3-7, then Functional dependencies from G are used, thus by the soundness of the Armstrong-axioms if $\text{PRESERVE}(\rho, F)$ claims that $X \rightarrow Y \in G^+$, then it is indeed so.

On the other hand, if $X \rightarrow Y \in G^+$, then $\text{LINEAR-CLOSURE}(R, F, X)$ (run by G as input) adds the attributes of Y one-by-one to X . In every step when an attribute is added, some functional dependency $U \rightarrow V$ of G is used. This dependency is in one of $\pi_{R_i}(F)$'s, since G is the union of these. An easy induction on the number of functional dependencies used in procedure $\text{LINEAR-CLOSURE}(R, F, X)$ shows that sooner or later Z becomes a subset of U , then applying the R_i -operation all attributes of V are added to Z . ■

18.2.4. Normal forms

The goal of transforming (decomposing) relational schemata into *normal forms* is to avoid the anomalies described in the previous section. Normal forms of many different strengths were introduced in the course of evolution of database theory, here only the *Boyce–Codd* normal form (BCNF) and the *third*, furthermore *fourth* normal form (3NF and 4NF) are treated in detail, since these are the most important ones from practical point of view.

Boyce-Codd normal form

Definition 18.17 Let R be relational schema, F be a family of functional dependencies over R . (R, F) is said to be in *Boyce-Codd normal form* if $X \rightarrow A \in F^+$ and $A \not\subseteq X$ implies that A is a superkey.

The most important property of BCNF is that it eliminates redundancy. This is based on the following theorem whose proof is left to the Reader as an exercise (Exercise 18.2-8).

Theorem 18.18 Schema (R, F) is in BCNF iff for arbitrary attribute $A \in R$ and key $X \subset R$ there exists no $Y \subseteq R$, for which $X \rightarrow Y \in F^+$; $Y \rightarrow X \notin F^+$; $Y \rightarrow A \in F^+$ and $A \notin Y$.

In other words, Theorem 18.18 states that “BCNF \iff There is no transitive dependence on keys”. Let us assume that a given schema is not in BCNF, for example $C \rightarrow B$ and $B \rightarrow A$ hold, but $B \rightarrow C$ does not, then the same B value could occur besides many different C values, but at each occasion the same A value would be stored with it, which is redundant. Formulating somewhat differently, the meaning of BCNF is that (only) using functional dependencies an attribute value in a row cannot be predicted from other attribute values. Indeed, assume that there exists a schema R , in which the value of an attribute can be determined using a functional dependency by comparison of two rows. That is, there exists two rows that agree on an attribute set X , differ on the set Y and the value of the remaining (unique) attribute A can be determined in one of the rows from the value taken in the other row.

X	Y	A
x	y_1	a
x	y_2	?

If the value ? can be determined by a functional dependency, then this value can only be a , the dependency is $Z \rightarrow A$, where Z is an appropriate subset of X . However, Z cannot be a superkey, since the two rows are distinct, thus R is not in BCNF.

3NF Although BCNF helps eliminating anomalies, it is not true that every schema can be decomposed into subschemata in BCNF so that the decomposition is dependency preserving. As it was shown in Example 18.6, no proper decomposition of schema CSZ preserves the $CS \rightarrow Z$ dependency. At the same time, the schema is clearly not in BCNF, because of the $Z \rightarrow C$ dependency.

Since dependency preserving is important because of consistency checking of

a database, it is practical to introduce a normal form that every schema has dependency preserving decomposition into that form, and it allows minimum possible redundancy. An attribute is called *prime attribute*, if it occurs in a key.

Definition 18.19 *The schema (R, F) is in **third normal form**, if whenever $X \rightarrow A \in F^+$, then either X is a superkey, or A is a prime attribute.*

The schema *SAIP* of Example 18.3 with the dependencies $SI \rightarrow P$ and $S \rightarrow A$ is not in 3NF, since SI is the only key and so A is not a prime attribute. Thus, functional dependency $S \rightarrow A$ violates the 3NF property.

3NF is clearly weaker condition than BCNF, since “or A is a prime attribute” occurs in the definition. The schema *CSZ* in Example 18.6 is trivially in 3NF, because every attribute is prime, but it was already shown that it is not in BCNF.

Testing normal forms Theoretically every functional dependency in F^+ should be checked whether it violates the conditions of BCNF or 3NF, and it is known that F^+ can be exponentially large in the size of F . Nevertheless, it can be shown that if the functional dependencies in F are of the form that the right hand side is a single attribute always, then it is enough to check violation of BCNF, or 3NF respectively, for dependencies of F . Indeed, let $X \rightarrow A \in F^+$ be a dependency that violates the appropriate condition, that is X is not a superkey and in case of 3NF, A is not prime. $X \rightarrow A \in F^+ \iff A \in X^+$. In the step when $\text{CLOSURE}(R, F, X)$ puts A into X^+ (line 8) it uses a functional dependency $Y \rightarrow A$ from F that $Y \subset X^+$ and $A \notin Y$. This dependency is non-trivial and A is (still) not prime. Furthermore, if Y were a superkey, then by $R = Y^+ \subseteq (X^+)^+ = X^+$, X would also be a superkey. Thus, the functional dependency $Y \rightarrow A$ from F violates the condition of the normal form. The functional dependencies easily can be checked in polynomial time, since it is enough to calculate the closure of the left hand side of each dependency. This finishes checking for BCNF, because if the closure of each left hand side is R , then the schema is in BCNF, otherwise a dependency is found that violates the condition. In order to test 3NF it may be necessary to decide about an attribute whether it is prime or not. However this problem is NP-complete, see Problem 18-4.

Lossless join decomposition into BCNF Let (R, F) be a relational schema (where F is the set of functional dependencies). The schema is to be decomposed into union of subschemata R_1, R_2, \dots, R_k , such that the decomposition has the lossless join property, furthermore each R_i endowed with the set of functional dependencies $\pi_{R_i}(F)$ is in BCNF. The basic idea of the decomposition is simple:

- If (R, F) is in BCNF, then ready.
- If not, it is decomposed into two proper parts (R_1, R_2) , whose join is lossless.
- Repeat the above for R_1 and R_2 .

In order to see that this works one has to show two things:

- If (R, F) is not in BCNF, then it has a lossless join decomposition into smaller parts.
- If a part of a lossless join decomposition is further decomposed, then the new decomposition has the lossless join property, as well.

Lemma 18.20 *Let (R, F) be a relational schema (where F is the set of functional dependencies), $\rho = (R_1, R_2, \dots, R_k)$ be a lossless join decomposition of R . Furthermore, let $\sigma = (S_1, S_2)$ be a lossless join decomposition of R_1 with respect to $\pi_{R_1}(F)$. Then $(S_1, S_2, R_2, \dots, R_k)$ is a lossless join decomposition of R .*

The proof of Lemma 18.20 is based on the associativity of natural join. The details are left to the Reader (Exercise 18.2-9).

This can be applied for a simple, but unfortunately exponential time algorithm that decomposes a schema into subschemata of BCNF property. The projections in lines 4–5 of NAïV-BCNF(S, G) may be of exponential size in the length of the input. In order to decompose schema (R, F) , the procedure must be called with parameters R, F . Procedure NAïV-BCNF(S, G) is recursive, S is the actual schema with set of functional dependencies G . It is assumed that the dependencies in G are of the form $X \rightarrow A$, where A is a single attribute.

NAïV-BCNF(S, G)

```

1 while there exists  $\{X \rightarrow A\} \in G$ , that violates BCNF
2     do  $S_1 \leftarrow \{XA\}$ 
3          $S_2 \leftarrow S - A$ 
4          $G_1 \leftarrow \pi_{S_1}(G)$ 
5          $G_2 \leftarrow \pi_{S_2}(G)$ 
6     return (NAïV-BCNF( $S_1, G_1$ ), NAïV-BCNF( $S_2, G_2$ ))
7 return  $S$ 
```

However, if the algorithm is allowed overdoing things, that is to decompose a schema even if it is already in BCNF, then there is no need for projecting the dependencies. The procedure is based on the following two lemmatae.

Lemma 18.21

1. *A schema of only two attributes is in BCNF.*
2. *If R is not in BCNF, then there exists two attributes A and B in R , such that $(R - AB) \rightarrow A$ holds.*

Proof If the schema consists of two attributes, $R = AB$, then there are at most two possible non-trivial dependencies, $A \rightarrow B$ and $B \rightarrow A$. It is clear, that if some of them holds, then the left hand side of the dependency is a key, so the dependency does not violate the BCNF property. However, if none of the two holds, then BCNF is trivially satisfied.

On the other hand, let us assume that the dependency $X \rightarrow A$ violates the BCNF property. Then there must exist an attribute $B \in R - (XA)$, since otherwise X would be a superkey. For this B , $(R - AB) \rightarrow A$ holds. ■

Let us note, that the converse of the second statement of Lemma 18.21 is not true. It may happen that a schema R is in BCNF, but there are still two attributes $\{A, B\}$ that satisfy $(R - AB) \rightarrow A$. Indeed, let $R = ABC$, $F = \{C \rightarrow A, C \rightarrow B\}$. This schema is obviously in BCNF, nevertheless $(R - AB) = C \rightarrow A$.

The main contribution of Lemma 18.21 is that the projections of functional dependencies need not be calculated in order to check whether a schema obtained during the procedure is in BCNF. It is enough to calculate $(R - AB)^+$ for pairs $\{A, B\}$ of attributes, which can be done by $\text{LINEAR-CLOSURE}(R, F, X)$ in linear time, so the whole checking is polynomial (cubic) time. However, this requires a way of calculating $(R - AB)^+$ without actually projecting down the dependencies. The next lemma is useful for this task.

Lemma 18.22 *Let $R_2 \subset R_1 \subset R$ and let F be the set of functional dependencies of scheme R . Then*

$$\pi_{R_2}(\pi_{R_1}(F)) = \pi_{R_2}(F) .$$

The proof is left for the Reader (Exercise 18.2-10). The method of lossless join BCNF decomposition is as follows. Schema R is decomposed into two subschemata. One is XA that is in BCNF, satisfying $X \rightarrow A$. The other subschema is $R - A$, hence by Theorem 18.14 the decomposition has the lossless join property. This is applied recursively to $R - A$, until such a schema is obtained that satisfies property 2 of Lemma 18.21. The lossless join property of this recursively generated decomposition is guaranteed by Lemma 18.20.

POLYNOMIAL-BCNF(R, F)

```

1   $Z \leftarrow R$ 
2       $\triangleright Z$  is the schema that is not known to be in BCNF during the procedure.
3   $\rho \leftarrow \emptyset$ 
4  while there exist  $A, B$  in  $Z$ , such that  $A \in (Z - AB)^+$  and  $|Z| > 2$ 
5      do Let  $A$  and  $B$  be such a pair
6           $E \leftarrow A$ 
7           $Y \leftarrow Z - B$ 
8          while there exist  $C, D$  in  $Y$ , such that  $C \in (Z - CD)^+$ 
9              do  $Y \leftarrow Y - D$ 
10              $E \leftarrow C$ 
11              $\rho \leftarrow \rho \cup \{Y\}$ 
12              $Z \leftarrow Z - E$ 
13  $\rho \leftarrow \rho \cup \{Z\}$ 
14 return  $\rho$ 

```

The running time of $\text{POLYNOMIAL-BCNF}(R, F)$ is polynomial, in fact it can be bounded by $O(n^5)$, as follows. During each execution of the loop in lines 4–12 the size of Z is decreased by at least one, so the loop body is executed at most n times. $(Z - AB)^+$ is calculated in line 4 for at most $O(n^2)$ pairs that can be done in linear time using LINEAR-CLOSURE that results in $O(n^3)$ steps for each execution of the loop body. In lines 8–10 the size of Y is decreased in each iteration, so during each execution of lines 3–12, they give at most n iteration. The condition of the

command **while** of line 8 is checked for $O(n^2)$ pairs of attributes, each checking is done in linear time. The running time of the algorithm is dominated by the time required by lines 8–10 that take $n \cdot n \cdot O(n^2) \cdot O(n) = O(n^5)$ steps altogether.

Dependency preserving decomposition into 3NF We have seen already that it is not always possible to decompose a schema into subschemata in BCNF so that the decomposition is dependency preserving. Nevertheless, if only 3NF is required then a decomposition can be given using $\text{MINIMAL-COVER}(R, F)$. Let R be a relational schema and F be the set of functional dependencies. Using $\text{MINIMAL-COVER}(R, F)$ a minimal cover G of F is constructed. Let $G = \{X_1 \rightarrow A_1, X_2 \rightarrow A_2, \dots, X_k \rightarrow A_k\}$.

Theorem 18.23 *The decomposition $\rho = (X_1A_1, X_2A_2, \dots, X_kA_k)$ is dependency preserving decomposition of R into subschemata in 3NF.*

Proof Since $G^+ = F^+$ and the functional dependency $X_i \rightarrow A_i$ is in $\pi_{R_i}(F)$, the decomposition preserves every dependency of F . Let us suppose indirectly, that the schema $R_i = X_iA_i$ is not in 3NF, that is there exists a dependency $U \rightarrow B$ that violates the conditions of 3NF. This means that the dependency is non-trivial and U is not a superkey in R_i and B is not a prime attribute of R_i . There are two cases possible. If $B = A_i$, then using that U is not a superkey $U \subsetneq X_i$ follows. In this case the functional dependency $U \rightarrow A_i$ contradicts to that $X_i \rightarrow A_i$ was a member of minimal cover, since its left hand side could be decreased. In the case when $BNEA_i$, $B \in X_i$ holds. B is not prime in R_i , thus X_i is not a key, only a superkey. However, then X_i would contain a key Y such that $Y \subsetneq X_i$. Furthermore, $Y \rightarrow A_i$ would hold, as well, that contradicts to the minimality of G since the left hand side of $X_i \rightarrow A_i$ could be decreased. ■

If the decomposition needs to have the lossless join property besides being dependency preserving, then ρ given in Theorem 18.23 is to be extended by a key X of R . Although it was seen before that it is not possible to list **all** keys in polynomial time, **one** can be obtained in a simple greedy way, the details are left to the Reader (Exercise 18.2-11).

Theorem 18.24 *Let (R, F) be a relational schema, and let $G = \{X_1 \rightarrow A_1, X_2 \rightarrow A_2, \dots, X_k \rightarrow A_k\}$ be a minimal cover of F . Furthermore, let X be a key in (R, F) . Then the decomposition $\tau = (X, X_1A_1, X_2A_2, \dots, X_kA_k)$ is a lossless join and dependency preserving decomposition of R into subschemata in 3NF.*

Proof It was shown during the proof of Theorem 18.23 that the subschemata $R_i = X_iA_i$ are in 3NF for $i = 1, 2, \dots, k$. There cannot be a non-trivial dependency in the subschema $R_0 = X$, because if it were, then X would not be a key, only a superkey.

The lossless join property of τ is shown by the use of $\text{JOIN-TEST}(R, G, \rho)$ procedure. Note that it is enough to consider the minimal cover G of F . More precisely, we show that the row corresponding to X in the table will be all 0 after running $\text{JOIN-TEST}(R, G, \rho)$. Let A_1, A_2, \dots, A_m be the order of the attributes of $R - X$ as $\text{CLOSURE}(R, G, X)$ inserts them into X^+ . Since X is a key, every attribute of $R - X$ is taken during $\text{CLOSURE}(R, G, X)$. It will be shown by induction on i that

the element in row of X and column of A_i is 0 after running $\text{JOIN-TEST}(R, G, \rho)$.

The base case of $i = 0$ is obvious. Let us suppose that the statement is true for $i - 1$ and consider when and why A_i is inserted into X^+ . In lines 6–8 of $\text{CLOSURE}(R, G, X)$ such a functional dependency $Y \rightarrow A_i$ is used where $Y \subseteq X \cup \{A_1, A_2, \dots, A_{i-1}\}$. Then $Y \rightarrow A_i \in G$, $YA_i = R_j$ for some j . The rows corresponding to X and $YA_i = R_j$ agree in columns of X (all 0 by the induction hypothesis), thus the entries in column of A_i are equated by $\text{JOIN-TEST}(R, G, \rho)$. This value is 0 in the row corresponding to $YA_i = R_j$, thus it becomes 0 in the row of X , as well. ■

It is interesting to note that although an arbitrary schema can be decomposed into subschemata in 3NF in polynomial time, nevertheless it is NP-complete to decide whether a given schema (R, F) is in 3NF, see Problem 18-4. However, the BCNF property can be decided in polynomial time. This difference is caused by that in order to decide 3NF property one needs to decide about an attribute whether it is prime. This latter problem requires the listing of all keys of a schema.

18.2.5. Multivalued dependencies

Example 18.8 Besides functional dependencies, some other dependencies hold in Example 18.1, as well. There can be several lectures of a subject in different times and rooms. Part of an instance of the schema could be the following.

Professor	Subject	Room	Student	Grade	Time
Caroline Doubtfire	Analysis	MA223	John Smith	A ⁻	Monday 8–10
Caroline Doubtfire	Analysis	CS456	John Smith	A ⁻	Wednesday 12–2
Caroline Doubtfire	Analysis	MA223	Ching Lee	A ⁺	Monday 8–10
Caroline Doubtfire	Analysis	CS456	Ching Lee	A ⁺	Wednesday 12–2

A set of values of Time and Room attributes, respectively, belong to each given value of Subject, and all other attribute values are repeated with these. Sets of attributes SR and StG are independent, that is their values occur in each combination.

The set of attributes Y is said to be **multivalued dependent** on set of attributes X , in notation $X \twoheadrightarrow Y$, if for every value on X , there exists a set of values on Y that is not dependent in any way on the values taken in $R - X - Y$. The precise definition is as follows.

Definition 18.25 *The relational schema R satisfies the **multivalued dependency** $X \twoheadrightarrow Y$, if for every relation r of schema R and arbitrary tuples t_1, t_2 of r that satisfy $t_1[X] = t_2[X]$, there exists tuples $t_3, t_4 \in r$ such that*

- $t_3[XY] = t_1[XY]$
- $t_3[R - XY] = t_2[R - XY]$
- $t_4[XY] = t_2[XY]$
- $t_4[R - XY] = t_1[R - XY]$

*holds.*¹

¹It would be enough to require the existence of t_3 , since the existence of t_4 would follow. However, the symmetry of multivalued dependency is more apparent in this way.

In Example 18.8 $S \rightarrow TR$ holds.

Remark 18.26 Functional dependency is **equality generating** dependency, that is from the equality of two objects it deduces the equality of other other two objects. On the other hand, multivalued dependency is **tuple generating dependency**, that is the existence of two rows that agree somewhere implies the existence of some other rows.

There exists a sound and complete axiomatisation of multivalued dependencies similar to the Armstrong-axioms of functional dependencies. Logical implication and inference can be defined analogously. The multivalued dependency $X \twoheadrightarrow Y$ is **logically implied** by the set M of multivalued dependencies, in notation $M \models X \twoheadrightarrow Y$, if every relation that satisfies all dependencies of M also satisfies $X \twoheadrightarrow Y$.

Note, that $X \rightarrow Y$ implies $X \twoheadrightarrow Y$. The rows t_3 and t_4 of Definition 18.25 can be chosen as $t_3 = t_2$ and $t_4 = t_1$, respectively. Thus, functional dependencies and multivalued dependencies admit a common axiomatisation. Besides Armstrong-axioms (A1)–(A3), five other are needed. Let R be a relational schema.

(A4) **Complementation:** $\{X \rightarrow Y\} \models X \twoheadrightarrow (R - X - Y)$.

(A5) **Extension:** If $X \twoheadrightarrow Y$ holds, and $V \subseteq W$, then $WX \twoheadrightarrow VY$.

(A6) **Transitivity:** $\{X \twoheadrightarrow Y, Y \twoheadrightarrow Z\} \models X \twoheadrightarrow (Z - Y)$.

(A7) $\{X \rightarrow Y\} \models X \twoheadrightarrow Y$.

(A8) If $X \twoheadrightarrow Y$ holds, $Z \subseteq Y$, furthermore for some W disjoint from Y $W \rightarrow Z$ holds, then $X \rightarrow Z$ is true, as well.

Beeri, Fagin and Howard proved that (A1)–(A8) is sound and complete system of axioms for functional and Multivalued dependencies together. Proof of soundness is left for the Reader (Exercise 18.2-12), the proof of the completeness exceeds the level of this book. The rules of Lemma 18.2 are valid in exactly the same way as when only functional dependencies were considered. Some further rules are listed in the next Proposition.

Claim 18.27 *The followings are true for multivalued dependencies.*

1. **Union rule:** $\{X \twoheadrightarrow Y, X \twoheadrightarrow Z\} \models X \twoheadrightarrow YZ$.
2. **Pseudotransitivity:** $\{X \twoheadrightarrow Y, WY \twoheadrightarrow Z\} \models WX \twoheadrightarrow (Z - WY)$.
3. **Mixed pseudotransitivity:** $\{X \twoheadrightarrow Y, XY \twoheadrightarrow Z\} \models X \twoheadrightarrow (Z - Y)$.
4. **Decomposition rule** for multivalued dependencies: *has $X \twoheadrightarrow Y$ and $X \twoheadrightarrow Z$ holds, then $X \twoheadrightarrow (Y \cap Z)$, $X \twoheadrightarrow (Y - Z)$ and $X \twoheadrightarrow (Z - Y)$ holds, as well.*

The proof of Proposition 18.27 is left for the Reader (Exercise 18.2-13).

Dependency basis Important difference between functional dependencies and multivalued dependencies is that $X \rightarrow Y$ immediately implies $X \rightarrow A$ for all A in Y , however $X \twoheadrightarrow A$ is deduced by the decomposition rule for multivalued dependencies from $X \twoheadrightarrow Y$ only if there exists a set of attributes Z such that $X \twoheadrightarrow Z$ and $Z \cap Y = A$, or $Y - Z = A$. Nevertheless, the following theorem is true.

Theorem 18.28 *Let R be a relational schema, $X \subseteq R$ be a set of attributes. Then there exists a partition Y_1, Y_2, \dots, Y_k of the set of attributes $R - X$ such that for $Z \subseteq R - X$ the multivalued dependency $X \twoheadrightarrow Z$ holds if and only if Z is the union of some Y_i 's.*

Proof We start from the one-element partition $W_1 = R - X$. This will be refined successively, while the property that $X \twoheadrightarrow W_i$ holds for all W_i in the actual decomposition, is kept. If $X \twoheadrightarrow Z$ and Z is not a union of some of the W_i 's, then replace every W_i such that neither $W_i \cap Z$ nor $W_i - Z$ is empty by $W_i \cap Z$ and $W_i - Z$. According to the decomposition rule of Proposition 18.27, both $X \twoheadrightarrow (W_i \cap Z)$ and $X \twoheadrightarrow (W_i - Z)$ hold. Since $R - X$ is finite, the refinement process terminates after a finite number of steps, that is for all Z such that $X \twoheadrightarrow Z$ holds, Z is the union of some blocks of the partition. In order to complete the proof one needs to observe only that by the union rule of Proposition 18.27, the union of some blocks of the partition depends on X in multivalued way. ■

Definition 18.29 *The partition Y_1, Y_2, \dots, Y_k constructed in Theorem 18.28 from a set D of functional and multivalued dependencies is called the **dependency basis** of X (with respect to D).*

Example 18.9 Consider the familiar schema

$$R(\text{Professor}, \text{Subject}, \text{Room}, \text{Student}, \text{Grade}, \text{Time})$$

of Examples 18.1 and 18.8. $\mathbf{Su} \twoheadrightarrow \mathbf{RT}$ was shown in Example 18.8. By the complementation rule $\mathbf{Su} \twoheadrightarrow \mathbf{PStG}$ follows. $\mathbf{Su} \twoheadrightarrow \mathbf{P}$ is also known. This implies by axiom (A7) that $\mathbf{Su} \twoheadrightarrow \mathbf{P}$. By the decomposition rule $\mathbf{Su} \twoheadrightarrow \mathbf{Stg}$ follows. It is easy to see that no other one-element attribute set is determined by \mathbf{Su} via multivalued dependency. Thus, the dependency basis of \mathbf{Su} is the partition $\{\mathbf{P}, \mathbf{RT}, \mathbf{StG}\}$.

We would like to compute the set D^+ of logical consequences of a given set D of functional and multivalued dependencies. One possibility is to apply axioms (A1)–(A8) to extend the set of dependencies repeatedly, until no more extension is possible. However, this could be an exponential time process in the size of D . One cannot expect any better, since it was shown before that even D^+ can be exponentially larger than D . Nevertheless, in many applications it is not needed to compute the whole set D^+ , one only needs to decide whether a given functional dependency $X \rightarrow Y$ or multivalued dependency $X \twoheadrightarrow Y$ belongs to D^+ or not. In order to decide about a multivalued dependency $X \twoheadrightarrow Y$, it is enough to compute the dependency basis of X , then to check whether $Z - X$ can be written as a union of some blocks of the partition. The following is true.

Theorem 18.30 (Beeri). *In order to compute the dependency basis of a set of attributes X with respect to a set of dependencies D , it is enough to consider the following set M of multivalued dependencies:*

1. All multivalued dependencies of D and

2. for every $X \rightarrow Y$ in D the set of multivalued dependencies $X \twoheadrightarrow A_1, X \twoheadrightarrow A_2, \dots, X \twoheadrightarrow A_k$, where $Y = A_1 A_2 \dots A_k$, and the A_i 's are single attributes.

The only thing left is to decide about functional dependencies based on the dependency basis. $\text{CLOSURE}(R, F, X)$ works correctly only if multivalued dependencies are not considered. The next theorem helps in this case.

Theorem 18.31 (Beeri). *Let us assume that $A \notin X$ and the dependency basis of X with respect to the set M of multivalued dependencies obtained in Theorem 18.30 is known. $X \rightarrow A$ holds if and only if*

1. A forms a single element block in the partition of the dependency basis, and
2. There exists a set Y of attributes that does not contain A , $Y \rightarrow Z$ is an element of the originally given set of dependencies D , furthermore $A \in Z$.

Based on the observations above, the following polynomial time algorithm can be given to compute the dependency basis of a set of attributes X .

DEPENDENCY-BASIS(R, M, X)

```

1  $\mathcal{S} \leftarrow \{R - X\}$  ▷ The collection of sets in the dependency basis is  $\mathcal{S}$ .
2 repeat
3   for all  $V \twoheadrightarrow W \in M$ 
4     do if there exists  $Y \in \mathcal{S}$  such that  $Y \cap W \neq \emptyset \wedge Y \cap V = \emptyset$ 
5       then  $\mathcal{S} \leftarrow \mathcal{S} - \{\{Y\}\} \cup \{\{Y \cap W\}, \{Y - W\}\}$ 
6 until  $\mathcal{S}$  does not change
7 return  $\mathcal{S}$ 
```

It is immediate that if \mathcal{S} changes in lines 3–5. of $\text{DEPENDENCY-BASIS}(R, M, X)$, then some block of the partition is cut by the algorithm. This implies that the running time is a polynomial function of the sizes of M and R . In particular, by careful implementation one can make this polynomial to $O(|M| \cdot |R|^3)$, see Problem 18-5.

Fourth normal form 4NF The Boyce-Codd normal form can be generalised to the case where multivalued dependencies are also considered besides functional dependencies, and one needs to get rid of the redundancy caused by them.

Definition 18.32 *Let R be a relational schema, D be a set of functional and multivalued dependencies over R . R is in **fourth normal form (4NF)**, if for arbitrary multivalued dependency $X \twoheadrightarrow Y \in D^+$ for which $Y \not\subseteq X$ and $R \text{ NEXY}$, holds that X is superkey in R .*

Observe that $4\text{NF} \implies \text{BCNF}$. Indeed, if $X \rightarrow A$ violated the BCNF condition, then $A \notin X$, furthermore XA could not contain all attributes of R , because that would imply that X is a superkey. However, $X \rightarrow A$ implies $X \twoheadrightarrow A$ by (A8), which in turn would violate the 4NF condition.

Schema R together with set of functional and multivalued dependencies D can be decomposed into $\rho = (R_1, R_2, \dots, R_k)$, where each R_i is in 4NF and the decomposition has the lossless join property. The method follows the same idea as the decomposition into BCNF subschemata. If schema S is not in 4NF, then there exists a multivalued dependency $X \twoheadrightarrow Y$ in the projection of D onto S that violates the 4NF condition. That is, X is not a superkey in S , Y neither is empty, nor is a subset of X , furthermore the union of X and Y is not S . It can be assumed without loss of generality that X and Y are disjoint, since $X \twoheadrightarrow (Y - X)$ is implied by $X \twoheadrightarrow Y$ using (A1), (A7) and the decomposition rule. In this case S can be replaced by subschemata $S_1 = XY$ and $S_2 = S - Y$, each having a smaller number of attributes than S itself, thus the process terminates in finite time.

Two things has to be dealt with in order to see that the process above is correct.

- Decomposition S_1, S_2 has the lossless join property.
- How can the projected dependency set $\pi_S(D)$ be computed?

The first problem is answered by the following theorem.

Theorem 18.33 *The decomposition $\rho = (R_1, R_2)$ of schema R has the lossless join property with respect to a set of functional and multivalued dependencies D iff*

$$(R_1 \cap R_2) \twoheadrightarrow (R_1 - R_2).$$

Proof The decomposition $\rho = (R_1, R_2)$ of schema R has the lossless join property iff for any relation r over the schema R that satisfies all dependencies from D holds that if μ and ν are two tuples of r , then there exists a tuple φ satisfying $\varphi[R_1] = \mu[R_1]$ and $\varphi[R_2] = \nu[R_2]$, then it is contained in r . More precisely, φ is the natural join of the projections of μ on R_1 and of ν on R_2 , respectively, which exist iff $\mu[R_1 \cap R_2] = \nu[R_1 \cap R_2]$. Thus the fact that φ is always contained in r is equivalent with that $(R_1 \cap R_2) \twoheadrightarrow (R_1 - R_2)$. ■

To compute the projection $\pi_S(D)$ of the dependency set D one can use the following theorem of Aho, Beeri and Ullman. $\pi_S(D)$ is the set of multivalued dependencies that are logical implications of D and use attributes of S only.

Theorem 18.34 (Aho, Beeri és Ullman). $\pi_S(D)$ consists of the following dependencies:

- For all $X \twoheadrightarrow Y \in D^+$, if $X \subseteq S$, then $X \twoheadrightarrow (Y \cap S) \in \pi_S(D)$.
- For all $X \twoheadrightarrow Y \in D^+$, if $X \subseteq S$, then $X \twoheadrightarrow (Y \cap S) \in \pi_S(D)$.

Other dependencies cannot be derived from the fact that D holds in R .

Unfortunately this theorem does not help in computing the projected dependencies in polynomial time, since even computing D^+ could take exponential time. Thus, the algorithm of 4NF decomposition is not polynomial either, because the 4NF condition must be checked with respect to the projected dependencies in the subschemata. This is in deep contrast with the case of BCNF decomposition. The reason is, that to check BCNF condition one does not need to compute the projected dependencies,

only closures of attribute sets need to be considered according to Lemma 18.21.

Exercises

18.2-1 Are the following inference rules sound?

- If $XW \rightarrow Y$ and $XY \rightarrow Z$, then $X \rightarrow (Z - W)$.
- If $X \rightarrow Y$ and $Y \rightarrow Z$, then $X \rightarrow Z$.
- If $X \rightarrow Y$ and $XY \rightarrow Z$, then $X \rightarrow Z$.

18.2-2 Prove Theorem 18.30, that is show the following. Let D be a set of functional and multivalued dependencies, and let $m(D) = \{X \twoheadrightarrow Y : X \twoheadrightarrow Y \in D\} \cup \{X \twoheadrightarrow A : A \in Y \text{ for some } X \rightarrow Y \in D\}$. Then

- $D \models X \rightarrow Y \implies m(D) \models X \twoheadrightarrow Y$, and
- $D \models X \twoheadrightarrow Y \iff m(D) \models X \twoheadrightarrow Y$.

Hint. Use induction on the inference rules to prove b.

18.2-3 Consider the database of an investment firm, whose attributes are as follows: B (stockbroker), O (office of stockbroker), I (investor), S (stock), A (amount of stocks of the investor), D (dividend of the stock). The following functional dependencies are valid: $S \rightarrow D$, $I \rightarrow B$, $IS \rightarrow A$, $B \rightarrow O$.

- Determine a key of schema $R = BOISAD$.
- How many keys are in schema R ?
- Give a lossless join decomposition of R into subschemata in BCNF.
- Give a dependency preserving and lossless join decomposition of R into subschemata in 3NF.

18.2-4 The schema R of Exercise 18.2-3 is decomposed into subschemata SD , IB , ISA and BO . Does this decomposition have the lossless join property?

18.2-5 Assume that schema R of Exercise 18.2-3 is represented by ISA , IB , SD and ISO subschemata. Give a minimal cover of the projections of dependencies given in Exercise 18.2-3. Exhibit a minimal cover for the union of the sets of projected dependencies. Is this decomposition dependency preserving?

18.2-6 Let the functional dependency $S \rightarrow D$ of Exercise 18.2-3 be replaced by the multivalued dependency $S \twoheadrightarrow D$. That is, D represents the stock's dividend "history".

- Compute the dependency basis of I .
- Compute the dependency basis of BS .
- Give a decomposition of R into subschemata in 4NF.

18.2-7 Consider the decomposition $\rho = \{R_1, R_2, \dots, R_k\}$ of schema R . Let $r_i = \pi_{R_i}(r)$, furthermore $m_\rho(r) = \bowtie_{i=1}^k \pi_{R_i}(r)$. Prove:

- $r \subseteq m_\rho(r)$.
- If $s = m_\rho(r)$, then $\pi_{R_i}(s) = r_i$.
- $m_\rho(m_\rho(r)) = m_\rho(r)$.

18.2-8 Prove that schema (R, F) is in BCNF iff for arbitrary $A \in R$ and key $X \subset R$, it holds that there exists no $Y \subseteq R$, for which $X \rightarrow Y \in F^+$; $Y \rightarrow X \notin F^+$; $Y \rightarrow A \in F^+$ and $A \notin Y$.

18.2-9 Prove Lemma 18.20.

18.2-10 Let us assume that $R_2 \subset R_1 \subset R$ and the set of functional dependencies of schema R is F . Prove that $\pi_{R_2}(\pi_{R_1}(F)) = \pi_{R_2}(F)$.

18.2-11 Give a $O(n^2)$ running time algorithm to find a key of the relational schema (R, F) . *Hint.* Use that R is superkey and each superkey contains a key. Try to drop attributes from R one-by-one and check whether the remaining set is still a key.

18.2-12 Prove that axioms (A1)–(A8) are sound for functional and multivalued dependencies.

18.2-13 Derive the four inference rules of Proposition 18.27 from axioms (A1)–(A8).

18.3. Generalised dependencies

Two such dependencies will be discussed in this section that are generalizations of the previous ones, however cannot be axiomatised with axioms similar to (A1)–(A8).

18.3.1. Join dependencies

Theorem 18.33 states that multivalued dependency is equivalent with that some decomposition the schema into two parts has the lossless join property. Its generalisation is the *join dependency*.

Definition 18.35 Let R be a relational schema and let $R = \bigcup_{i=1}^k X_i$. The relation r belonging to R is said to satisfy the *join dependency*

$$\bowtie[X_1, X_2, \dots, X_k]$$

if

$$r = \bowtie_{i=1}^k \pi_{X_i}(r) .$$

In this setting r satisfies multivalued dependency $X \twoheadrightarrow Y$ iff it satisfies the join dependency $\bowtie[XY, X(R - Y)]$. The join dependency $\bowtie[X_1, X_2, \dots, X_k]$ expresses that the decomposition $\rho = (X_1, X_2, \dots, X_k)$ has the lossless join property. One can define the *fifth normal form, 5NF*.

Definition 18.36 The relational schema R is in *fifth normal form*, if it is in 4NF and has no non-trivial join dependency.

The fifth normal form has theoretical significance primarily. The schemata used in practice usually have *primary keys*. Using that the schema could be decomposed into subschemata of two attributes each, where one of the attributes is a superkey in every subschema.

Example 18.10 Consider the database of clients of a bank (Client-number, Name, Address, accountBalance). Here C is unique identifier, thus the schema could be decomposed into (CN, CA, CB), which has the lossless join property. However, it is not worth doing so, since no storage place can be saved, furthermore no anomalies are avoided with it.

There exists an *axiomatisation* of a dependency system if there is a finite set of inference rules that is sound and complete, i.e. logical implication coincides with being derivable by using the inference rules. For example, the Armstrong-axioms give an axiomatisation of functional dependencies, while the set of rules (A1)–(A8)

is the same for functional and multivalued dependencies considered together. Unfortunately, the following negative result is true.

Theorem 18.37 *The family of join dependencies has no finite axiomatisation.*

In contrary to the above, Abiteboul, Hull and Vianu show in their book that the logical implication problem can be decided by an algorithm for the family of functional and join dependencies taken together. The complexity of the problem is as follows.

Theorem 18.38

- *It is NP-complete to decide whether a given join dependency is implied by another given join dependency and a functional dependency.*
- *It is NP-hard to decide whether a given join dependency is implied by given set of multivalued dependencies.*

18.3.2. Branching dependencies

A generalisation of functional dependencies is the family of **branching dependencies**. Let us assume that $A, B \subset R$ and there exists no $q + 1$ rows in relation r over schema R , such that they contain at most p distinct values in columns of A , but all $q + 1$ values are pairwise distinct in some column of B . Then B is said to be **(p, q) -dependent** on A , in notation $A \xrightarrow{p,q} B$. In particular, $A \xrightarrow{1,1} B$ holds if and only if functional dependency $A \rightarrow B$ holds.

Example 18.11 Consider the database of the trips of an international transport truck.

- One trip: four distinct countries.
- One country has at most five neighbours.
- There are 30 countries to be considered.

Let x_1, x_2, x_3, x_4 be the attributes of the countries reached in a trip. In this case $x_i \xrightarrow{1,1} x_{i+1}$ does not hold, however another dependency is valid:

$$x_i \xrightarrow{1,5} x_{i+1} .$$

The storage space requirement of the database can be significantly reduced using these dependencies. The range of each element of the original table consists of 30 values, names of countries or some codes of them (5 bits each, at least). Let us store a little table ($30 \times 5 \times 5 = 750$ bits) that contains a numbering of the neighbours of each country, which assigns to them the numbers 0,1,2,3,4 in some order. Now we can replace attribute x_2 by these numbers (x_2^*), because the value of x_1 gives the starting country and the value of x_2^* determines the second country with the help of the little table. The same holds for the attribute x_3 , but we can decrease the number of possible values even further, if we give a table of numbering the possible third countries for each x_1, x_2 pair. In this case, the attribute x_3^* can take only 4 different values. The same holds for x_4 , too. That is, while each element of the original table could be encoded by 5 bits, now for the cost of two little auxiliary tables we could decrease the length of the elements in the second column to 3 bits, and that of the elements in the third and fourth columns to 2 bits.

The (p, q) -closure of an attribute set $X \subseteq R$ can be defined:

$$C_{p,q}(X) = \{A \in R : X \xrightarrow{p,q} A\}.$$

In particular, $C_{1,1}(X) = X^+$. In case of branching dependencies even such basic questions are hard as whether there exists an **Armstrong-relation** for a given family of dependencies.

Definition 18.39 Let R be a relational schema, F be a set of dependencies of some dependency family \mathcal{F} defined on R . A relation r over schema R is **Armstrong-relation** for F , if the set of dependencies from \mathcal{F} that r satisfies is exactly F , that is $F = \{\sigma \in \mathcal{F} : r \models \sigma\}$.

Armstrong proved that for an arbitrary set of functional dependencies F there exists Armstrong-relation for F^+ . The proof is based on the three properties of closures of attributes sets with respect to F , listed in Exercise 18.1-1 For branching dependencies only the first two holds in general.

Lemma 18.40 Let $0 < p \leq q$, furthermore let R be a relational schema. For $X, Y \subseteq R$ one has

1. $X \subseteq C_{p,q}(X)$ and
2. $X \subseteq Y \implies C_{p,q}(X) \subseteq C_{p,q}(Y)$.

There exists such $C: 2^R \rightarrow 2^R$ mapping and natural numbers p, q that there exists no Armstrong-relation for C in the family if (p, q) -dependencies.

Grant Minker investigated **numerical dependencies** that are similar to branching dependencies. For attribute sets $X, Y \subseteq R$ the dependency $X \xrightarrow{k} Y$ holds in a relation r over schema R if for every tuple value taken on the set of attributes X , there exists at most k distinct tuple values taken on Y . This condition is stronger than that of $X \xrightarrow{1,k} Y$, since the latter only requires that in each column of Y there are at most k values, independently of each other. That allows $k^{|Y-X|}$ different Y projections. Numerical dependencies were axiomatised in some special cases, based on that Katona showed that branching dependencies have no finite axiomatisation. It is still an open problem whether logical implication is algorithmically decidable amongst branching dependencies.

Exercises

18.3-1 Prove Theorem 18.38.

18.3-2 Prove Lemma 18.40.

18.3-3 Prove that if $p = q$, then $C_{p,p}(C_{p,p}(X)) = C_{p,p}(X)$ holds besides the two properties of Lemma 18.40.

18.3-4 A $C: 2^R \rightarrow 2^R$ mapping is called a **closure**, if it satisfies the two properties of Lemma 18.40 and and the third one of Exercise 18.3-3. Prove that if $C: 2^R \rightarrow 2^R$ is a closure, and F is the family of dependencies defined by $X \rightarrow Y \iff Y \subseteq C(X)$, then there exists an Armstrong-relation for F in the family of $(1, 1)$ -dependencies (functional dependencies) and in the family of $(2, 2)$ -dependencies, respectively.

18.3-5 Let C be the closure defined by

$$C(X) = \begin{cases} X, & \text{if } |X| < 2 \\ R & \text{otherwise} \end{cases} .$$

Prove that there exists no Armstrong-relation for C in the family of (n, n) -dependencies, if $n > 2$.

Problems

18-1 External attributes

Maier calls attribute A an **external attribute** in the functional dependency $X \rightarrow Y$ with respect to the family of dependencies F over schema R , if the following two conditions hold:

1. $(F - \{X \rightarrow Y\}) \cup \{X \rightarrow (Y - A)\} \models X \rightarrow Y$, or
2. $(F - \{X \rightarrow Y\}) \cup \{(X - A) \rightarrow Y\} \models X \rightarrow Y$.

Design an $O(n^2)$ running time algorithm, whose input is schema (R, F) and output is a set of dependencies G equivalent with F that has no external attributes.

18-2 The order of the elimination steps in the construction of minimal cover is important

In the procedure $\text{MINIMAL-COVER}(R, F)$ the set of functional dependencies was changed in two ways: either by dropping redundant dependencies, or by dropping redundant attributes from the left hand sides of the dependencies. If the latter method is used first, until there is no more attribute that can be dropped from some left hand side, then the first method, this way a minimal cover is obtained really, according to Proposition 18.6. Prove that if the first method applied first and then the second, until there is no more possible applications, respectively, then the obtained set of dependencies is not necessarily a minimal cover of F .

18-3 BCNF subschema

Prove that the following problem is coNP-complete: Given a relational schema R with set of functional dependencies F , furthermore $S \subset R$, decide whether $(S, \pi_S(F))$ is in BCNF.

18-4 3NF is hard to recognise

Let (R, F) be a relational schema, where F is the system of functional dependencies.

The **k size key problem** is the following: given a natural number k , determine whether there exists a key of size at most k .

The **prime attribute problem** is the following: for a given $A \in R$, determine whether it is a prime attribute.

- a. Prove that the k size key problem is NP-complete. *Hint.* Reduce the **vertex cover** problem to the prime attribute problem.
- b. Prove that the prime attribute problem is NP-complete by reducing the k size key problem to it.

- c. Prove that determining whether the relational schema (R, F) is in 3NF is NP-complete. *Hint.* Reduce the prime attribute problem to it.

18-5 Running time of Dependency-basis

Give an implementation of procedure DEPENDENCY-BASIS, whose running time is $O(|M| \cdot |R|^3)$.

Chapter Notes

The relational data model was introduced by Codd [7] in 1970. Functional dependencies were treated in his paper of 1972 [?], their axiomatisation was completed by Armstrong [?]. The logical implication problem for functional dependencies were investigated by Beeri and Bernstein [4], furthermore Maier [19]. Maier also treats the possible definitions of minimal covers, their connections and the complexity of their computations in that paper. Maier, Mendelzon and Sagiv found method to decide logical implications among general dependencies [20]. Beeri, Fagin and Howard proved that axiom system (A1)–(A8) is sound and complete for functional and multivalued dependencies taken together [?]. Yu and Johnson [?] constructed such relational schema, where $|F| = k$ and the number of keys is $k!$. Békéssy and Demetrovics [6] gave a simple and beautiful proof for the statement, that from k functional dependencies at most $k!$ keys can be obtained, thus Yu and Johnson's construction is extremal.

Armstrong-relations were introduced and studied by Fagin [?, 13], furthermore by Beeri, Fagin, Dowd and Statman [5].

Multivalued dependencies were independently discovered by Zaniolo [?], Fagin [12] and Delobel [8].

The necessity of normal forms was recognised by Codd while studying update anomalies [?, ?]. The Boyce–Codd normal form was introduced in [?]. The definition of the third normal form used in this chapter was given by Zaniolo [28]. Complexity of decomposition into subschemata in certain normal forms was studied by Lucchesi and Osborne [18], Beeri and Bernstein [4], furthermore Tsou and Fischer [26].

Theorems 18.30 and 18.31 are results of Beeri [3]. Theorem 18.34 is from a paper of Aho, Beeri és Ullman [2].

Theorems 18.37 and 18.38 are from the book of Abiteboul, Hull and Vianu [1], the non-existence of finite axiomatisation of join dependencies is Petrov's result [21].

Branching dependencies were introduced by Demetrovics, Katona and Sali, they studied existence of Armstrong-relations and the size of minimal Armstrong-relations [9, 10, 11, 23]. Katona showed that there exists no finite axiomatisation of branching dependencies in (ICDT'92 Berlin, invited talk) but never published.

Possibilities of axiomatisation of numerical dependencies were investigated by Grant and Minker [15, 16].

Good introduction of the concepts of this chapter can be found in the books of Abiteboul, Hull and Vianu [1], Ullman [27] furthermore Thalheim [24], respectively.

19. Query Rewriting in Relational Databases

In chapter “Relational database design” basic concepts of relational databases were introduced, such as relational schema, relation, instance. Databases were studied from the designer point of view, the main question was how to avoid redundant data storage, or various anomalies arising during the use of the database.

In the present chapter the schema is considered to be given and focus is on fast and efficient ways of answering user queries. First, basic (theoretical) query languages and their connections are reviewed in Section 19.1.

In the second part of this chapter (Section 19.2) views are considered. Informally, a view is nothing else, but result of a query. Use of views in query efficiency, providing physical data independence and data integration is explained.

Finally, the third part of the present chapter (Section 19.3) introduces query rewriting.

19.1. Queries

Consider the database of cinemas in Budapest. Assume that the schema consists of three relations:

$$\mathbf{CinePest} = \{Film, Theater, Show\} . \quad (19.1)$$

The schemata of individual relations are as follows:

$$\begin{aligned} Film &= \{\mathbf{T}itle, \mathbf{D}irector, \mathbf{A}ctor\} , \\ Theater &= \{\mathbf{T}heater, \mathbf{A}ddress, \mathbf{P}hone\} , \\ Show &= \{\mathbf{T}heater, \mathbf{T}itle, \mathbf{T}ime\} . \end{aligned} \quad (19.2)$$

Possible values of instances of each relation are shown on Figure 19.1.

Typical user queries could be:

19.1 Who is the director of “Control”?

19.2 List the names address of those theatres where Kurosawa films are played.

19.3 Give the names of directors who played part in some of their films.

Film

Title	Director	Actor
Control	Antal, Nimród	Csányi, Sándor
Control	Antal, Nimród	Mucsi, Zoltán
Control	Antal, Nimród	Pindroch, Csaba
⋮	⋮	⋮
Rashomon	Akira Kurosawa	Toshiro Mifune
Rashomon	Akira Kurosawa	Machiko Kyo
Rashomon	Akira Kurosawa	Mori Masayuki

Theatre

Theater	Address	Phone
Bem	II., Margit Blvd. 5/b.	316-8708
Corvin	VIII., Corvin alley 1.	459-5050
Európa	VII., Rákóczi st. 82.	322-5419
Művész	VI., Teréz blvd. 30.	332-6726
⋮	⋮	⋮
Uránia	VIII., Rákóczi st. 21.	486-3413
Vörösmarty	VIII., Üllői st. 4.	317-4542

Show

Theater	Title	Time
Bem	Rashomon	19:00
Bem	Rashomon	21:30
Uránia	Control	18:15
Művész	Rashomon	16:30
Művész	Control	17:00
⋮	⋮	⋮
Corvin	Control	10:15

Figure 19.1 The database **CinePest**.

These queries define a mapping from the relations of database schema **CinePest** to some other schema (in the present case to schemata of single relations). Formally, *query* and *query mapping* should be distinguished. The former is a syntactic concept, the latter is a mapping from the set of instances over the input schema to the set of instances over the output schema, that is determined by the query according to some semantic interpretation. However, for both concepts the word “query” is used for the sake of convenience, which one is meant will be clear from

the context.

Definition 19.1 Queries q_1 and q_2 over schema R are said to be **equivalent**, in notation $q_1 \equiv q_2$, if they have the same output schema and for every instance \mathcal{I} over schema R $q_1(\mathcal{I}) = q_2(\mathcal{I})$ holds.

In the remaining of this chapter the most important query languages are reviewed. The expressive powers of query languages need to be compared.

Definition 19.2 Let \mathcal{Q}_1 and \mathcal{Q}_2 be query languages (with appropriate semantics). \mathcal{Q}_2 is **dominated by** \mathcal{Q}_1 (\mathcal{Q}_1 is **weaker**, than \mathcal{Q}_2), in notation $\mathcal{Q}_1 \sqsubseteq \mathcal{Q}_2$, if for every query q_1 of \mathcal{Q}_1 there exists a query $q_2 \in \mathcal{Q}_2$, such that $q_1 \equiv q_2$. \mathcal{Q}_1 and \mathcal{Q}_2 are **equivalent**, if $\mathcal{Q}_1 \sqsubseteq \mathcal{Q}_2$ and $\mathcal{Q}_1 \supseteq \mathcal{Q}_2$.

Example 19.1 *Query.* Consider Question 19.2. As a first try the next solution is obtained:

if there exist in relations *Film*, *Theater* and *Show* tuples $(x_T, \text{"Akira Kurosawa"}, x_A)$, (x_{Th}, x_{Ad}, x_P) and (x_{Th}, x_T, x_{Ti})
then put the tuple (*Theater* : x_{Th} , *Address* : x_A) into the output relation.

$x_T, x_A, x_{Th}, x_{Ad}, x_P, x_{Ti}$ denote different variables that take their values from the domains of the corresponding attributes, respectively. Using the same variables implicitly marked where should stand identical values in different tuples.

19.1.1. Conjunctive queries

Conjunctive queries are the simplest kind of queries, but they are the easiest to handle and have the most good properties. Three equivalent forms will be studied, two of them based on logic, the third one is of algebraic nature. The name comes from first order logic expressions that contain only existential quantors (\exists), furthermore consist of atomic expressions connected with logical “and”, that is conjunction.

Datalog – rule based queries The tuple (x_1, x_2, \dots, x_m) is called **free tuple** if the x_i 's are variables or constants. This is a generalisation of a tuple of a relational instance. For example, the tuple $(x_T, \text{"Akira Kurosawa"}, x_A)$ in Example 19.1 is a free tuple.

Definition 19.3 Let \mathbf{R} be a relational database schema. **Rule based conjunctive query** is an expression of the following form

$$ans(u) \leftarrow R_1(u_1), R_2(u_2), \dots, R_n(u_n), \quad (19.3)$$

where $n \geq 0$, R_1, R_2, \dots, R_n are relation names from \mathbf{R} , ans is a relation name not in \mathbf{R} , u, u_1, u_2, \dots, u_n are free tuples. Every variable occurring in u must occur in one of u_1, u_2, \dots, u_n , as well.

The rule based conjunctive query is also called a **rule** for the sake of simplicity. $ans(u)$ is the **head** of the rule, $R_1(u_1), R_2(u_2), \dots, R_n(u_n)$ is the **body** of the rule,

$R_i(u_i)$ is called a (*relational*) *atom*. It is assumed that each variable of the head also occurs in some atom of the body.

A rule can be considered as some tool that tells how can we deduce newer and newer *facts*, that is tuples, to include in the output relation. If the variables of the rule can be assigned such values that each atom $R_i(u_i)$ is true (that is the appropriate tuple is contained in the relation R_i), then tuple u is added to the relation *ans*. Since all variables of the head occur in some atoms of the body, one never has to consider *infinite* domains, since the variables can take values from the actual instance queried. Formally, let \mathcal{I} be an instance over relational schema \mathbf{R} , furthermore let q be a the query given by rule (19.3). Let $var(q)$ denote the set of variables occurring in q , and let $dom(\mathcal{I})$ denote the set of constants that occur in \mathcal{I} . The *image of \mathcal{I} under q* is given by

$$q(\mathcal{I}) = \{\nu(u) \mid \nu: var(q) \rightarrow dom(\mathcal{I}) \text{ and } \nu(u_i) \in R_i \ i = 1, 2, \dots, n\}. \quad (19.4)$$

An immediate way of calculating $q(\mathcal{I})$ is to consider all possible valuations ν in some order. There are more efficient algorithms, either by equivalent rewriting of the query, or by using some indices.

An important difference between atoms of the body and the head is that relations R_1, R_2, \dots, R_n are considered given, (physically) stored, while relation *ans* is not, it is thought to be calculated by the help of the rule. This justifies the names: R_i 's are *extensional relations* and *ans* is *intensional relation*.

Query q over schema \mathbf{R} is *monotone*, if for instances \mathcal{I} and \mathcal{J} over \mathbf{R} , $\mathcal{I} \subseteq \mathcal{J}$ implies $q(\mathcal{I}) \subseteq q(\mathcal{J})$. q is *satisfiable*, if there exists an instance \mathcal{I} , such that $q(\mathcal{I}) \neq \emptyset$. The proof of the next simple observation is left for the Reader (Exercise 19.1-1).

Claim 19.4 *Rule based queries are monotone and satisfiable.*

Proposition 19.4 shows the limitations of rule based queries. For example, the query *Which theatres play only Kurosawa films?* is obviously not monotone, hence cannot be expressed by rules of form (19.3).

Tableau queries. If the difference between variables and constants is not considered, then the body of a rule can be viewed as an instance over the schema. This leads to a tabular form of conjunctive queries that is most similar to the visual queries (QBE: Query By Example) of database management system Microsoft Access.

Definition 19.5 *A **tableau** over the schema \mathbf{R} is a generalisation of an instance over \mathbf{R} , in the sense that variables may occur in the tuples besides constants. The pair (\mathbf{T}, u) is a **tableau query** if \mathbf{T} is a tableau and u is a free tuple such that all variables of u occur in \mathbf{T} , as well. The free tuple u is the **summary**.*

The summary row u of tableau query (\mathbf{T}, u) shows which tuples form the result of the query. The essence of the procedure is that the pattern given by tableau \mathbf{T} is searched for in the database, and if found then the tuple corresponding to is included in the output relation. More precisely, the mapping $\nu: var(\mathbf{T}) \rightarrow dom(\mathcal{I})$ is an *embedding* of tableau (\mathbf{T}, u) into instance \mathcal{I} , if $\nu(\mathbf{T}) \subseteq \mathcal{I}$. The output relation

of tableau query (\mathbf{T}, u) consists of all tuples $\nu(u)$ that ν is an embedding of tableau (\mathbf{T}, u) into instance \mathcal{I} .

Example 19.2 *Tableau query* Let \mathbf{T} be the following tableau.

<i>Film</i>	<i>Title</i>	<i>Director</i>	<i>Actor</i>
	x_T	“Akira Kurosawa”	x_A
<i>Theater</i>	<i>Theater</i>	<i>Address</i>	<i>Phone</i>
	x_{Th}	x_{Ad}	x_P
<i>Show</i>	<i>Theater</i>	<i>Title</i>	<i>Time</i>
	x_{Th}	x_T	x_{Ti}

The tableau query $(\mathbf{T}, \langle Theater: x_{Th}, Address: x_{Ad} \rangle)$ answers question 19.2. of the introduction.

The syntax of tableau queries is similar to that of rule based queries. It will be useful later that conditions for one query to contain another one can be easily formulated in the language of tableau queries.

Relational algebra*. A database consists of relations, and a relation is a set of tuples. The result of a query is also a relation with a given attribute set. It is a natural idea that output of a query could be expressed by algebraic and other operations on relations. The *relational algebra** consists of the following operations.¹

Selection: It is of form either $\sigma_{A=c}$ or $\sigma_{A=B}$, where A and B are attributes while c is a constant. The operation can be applied to all such relations R that has attribute A (and B), and its result is relation ans that has the same set of attributes as R has, and consists of all tuples that satisfy the *selection condition*.

Projection: The form of the operation is $\pi_{A_1, A_2, \dots, A_n}$, $n \geq 0$, where A_i 's are distinct attributes. It can be applied to all such relations whose attribute set includes each A_i and its result is the relation ans that has attribute set $\{A_1, A_2, \dots, A_n\}$,

$$val = \{t[A_1, A_2, \dots, A_n] | t \in R\},$$

that is it consists of the restrictions of tuples in R to the attribute set $\{A_1, A_2, \dots, A_n\}$.

Natural join: This operation has been defined earlier in chapter “Relational database design”. Its notation is \bowtie , its input consists of two (or more) relations R_1, R_2 , with attribute sets V_1, V_2 , respectively. The attribute set of the output relation is $V_1 \cup V_2$.

$$R_1 \bowtie R_2 = \{t \text{ tuple over } V_1 \cup V_2 | \exists v \in R_1, \exists w \in R_2, t[V_1] = v \text{ \&es } t[V_2] = w\}.$$

Renaming: Attribute renaming is nothing else, but an injective mapping from a finite set of attributes U into the set of all attributes. Attribute renaming f can

¹ The relational algebra* is the *monotone* part of the (full) relational algebra introduced later.

be given by the list of pairs $(A, f(A))$, where $ANEf(A)$, which is written usually in the form $A_1A_2 \dots A_n \rightarrow B_1B_2 \dots B_n$. The **renaming operator** δ_f maps from inputs over U to outputs over $f[U]$. If R is a relation over U , then

$$\delta_f(R) = \{v \text{ over } f[U] \mid \exists u \in R, v(f(A)) = u(A), \forall A \in U\} .$$

Relational algebra* queries are obtained by finitely many applications of the operations above from **relational algebra base queries**, which are

Input relation: R .

Single constant: $\{ \langle A : a \rangle \}$, where a is a constant, A is an attribute name.

Example 19.3 *Relational algebra* query.* The question 19.2. of the introduction can be expressed with relational algebra operations as follows.

$$\pi_{Theater, Address} ((\sigma_{Director="Akira Kurosawa"}(Film) \bowtie Show) \bowtie Theater) .$$

The mapping given by a relational algebra* query can be easily defined via induction on the operation tree. It is easy to see (Exercise 19.1-2) that non-satisfiable queries can be given using relational algebra*. There exist no rule based or tableau query equivalent with such a non-satisfiable query. Nevertheless, the following is true.

Theorem 19.6 *Rule based queries, tableau queries and satisfiable relational algebra* are equivalent query languages.*

The proof of Theorem 19.6 consists of three main parts:

1. Rule based \equiv Tableau
2. Satisfiable relational algebra* \sqsubseteq Rule based
3. Rule based \sqsubseteq Satisfiable relational algebra*

The first (easy) step is left to the Reader (Exercise 19.1-3). For the second step, it has to be seen first, that the language of rule based queries is closed under composition. More precisely, let $\mathbf{R} = \{R_1, R_2, \dots, R_n\}$ be a database, q be a query over \mathbf{R} . If the output relation of q is S_1 , then in a subsequent query S_1 can be used in the same way as any extensional relation of \mathbf{R} . Thus relation S_2 can be defined, then with its help relation S_3 can be defined, and so on. Relations S_i are **intensional** relations. The **conjunctive query program** P is a list of rules

$$\begin{aligned} S_1(u_1) &\leftarrow body_1 \\ S_2(u_2) &\leftarrow body_2 \\ &\vdots \\ S_m(u_m) &\leftarrow body_m , \end{aligned} \tag{19.5}$$

where S_i 's are pairwise distinct and not contained in \mathbf{R} . In rule body $body_i$ only relations R_1, R_2, \dots, R_n and S_1, S_2, \dots, S_{i-1} can occur. S_m is considered to be the output relation of P , its evaluation is done by computing the results of the rules

one-by-one in order. It is not hard to see that with appropriate renaming the variables P can be substituted by a single rule, as it is shown in the following example.

Example 19.4 *Conjunctive query program.* Let $\mathbf{R} = \{Q, R\}$, and consider the following conjunctive query program

$$\begin{aligned} S_1(x, z) &\leftarrow Q(x, y), R(y, z, w) \\ S_2(x, y, z) &\leftarrow S_1(x, w), R(w, y, v), S_1(v, z) \\ S_3(x, z) &\leftarrow S_2(x, u, v), Q(v, z) . \end{aligned} \quad (19.6)$$

S_2 can be written using Q and R only by the first two rules of (19.6)

$$S_2(x, y, z) \leftarrow Q(x, y_1), R(y_1, w, w_1), R(w, y, v), Q(v, y_2), R(y_2, z, w_2) . \quad (19.7)$$

It is apparent that some variables had to be renamed to avoid unwanted interaction of rule bodies. Substituting expression (19.7) into the third rule of (19.6) in place of S_2 , and appropriately renaming the variables

$$S_3(x, z) \leftarrow Q(x, y_1), R(y_1, w, w_1), R(w, u, v_1), Q(v_1, y_2), R(y_2, v, w_2), Q(v, z) . \quad (19.8)$$

is obtained.

Thus it is enough to realise each single relational algebra* operation by an appropriate rule.

$P \bowtie Q$: Let \vec{x} denote the list of variables (and constants) corresponding to the common attributes of P and Q , let \vec{y} denote the variables (and constants) corresponding to the attributes occurring only in P , while \vec{z} denotes those of corresponding to Q 's own attributes. Then rule $ans(\vec{x}, \vec{y}, \vec{z}) \leftarrow P(\vec{x}, \vec{y}), Q(\vec{x}, \vec{z})$ gives exactly relation $P \bowtie Q$.

$\sigma_F(R)$: Assume that $R = R(A_1, A_2, \dots, A_n)$ and the selection condition F is of form either $A_i = a$ or $A_i = A_j$, where A_i, A_j are attributes a is constant. Then

$$ans(x_1, \dots, x_{i-1}, a, x_{i+1}, \dots, x_n) \leftarrow R(x_1, \dots, x_{i-1}, a, x_{i+1}, \dots, x_n) ,$$

respectively,

$$ans(x_1, \dots, x_{i-1}, y, x_{i+1}, \dots, x_{j-1}, y, x_{j+1}, \dots, x_n) \leftarrow R(x_1, \dots, x_{i-1}, y, x_{i+1}, \dots, x_{j-1}, y, x_{j+1}, \dots, x_n)$$

are the rules sought. The satisfiability of relational algebra* query is used here. Indeed, during composition of operations we never obtain an expression where two distinct constants should be equated.

$\pi_{A_{i_1}, A_{i_2}, \dots, A_{i_m}}(R)$: If $R = R(A_1, A_2, \dots, A_n)$, then

$$ans(x_{i_1}, x_{i_2}, \dots, x_{i_m}) \leftarrow R(x_1, x_2, \dots, x_n)$$

works.

$A_1 A_2 \dots A_n \rightarrow B_1 B_2 \dots B_n$: The renaming operation of relational algebra* can be achieved by renaming the appropriate variables, as it was shown in Example 19.4.

For the proof of the third step let us consider rule

$$ans(\vec{x}) \leftarrow R_1(\vec{x}_1), R_2(\vec{x}_2), \dots, R_n(\vec{x}_n) . \quad (19.9)$$

By renaming the attributes of relations R_i 's, we may assume without loss of generality that all attribute names are distinct. Then $R = R_1 \bowtie R_2 \bowtie \dots \bowtie R_n$ can be constructed that is really a direct product, since the attribute names are distinct. The constants and multiple occurrences of variables of rule (19.9) can be simulated by appropriate selection operators. The final result is obtained by projecting to the set of attributes corresponding to the variables of relation ans .

19.1.2. Extensions

Conjunctive queries are a class of query languages that has many good properties. However, the set of expressible questions are rather narrow. Consider the following.

19.4. List those pairs where one member directed the other member in a film, and vice versa, the other member also directed the first in a film.

19.5. Which theatres show “La Dolce Vita” or “Rashomon”?

19.6. Which are those films of Hitchcock that Hitchcock did not play a part in?

19.7. List those films whose every actor played in some film of Fellini.

19.8. Let us recall the game “Chain-of-Actors”. The first player names an actor/actress, the next another one who played in some film together with the first named. This is continued like that, always a new actor/actress has to be named who played together with the previous one. The winner is that player who could continue the chain last time. List those actors/actresses who could be reached by “Chain-of-Actors” starting with “Marcello Mastroianni”.

Equality atoms. Question **19.4.** can be easily answered if equalities are also allowed rule bodies, besides relational atoms:

$$ans(y_1, y_2) \leftarrow Film(x_1, y_1, z_1), Film(x_2, y_2, z_2), y_1 = z_2, y_2 = z_1. \quad (19.10)$$

Allowing equalities raises two problems. First, the result of the query could become infinite. For example, the rule based query

$$ans(x, y) \leftarrow R(x), y = z \quad (19.11)$$

results in an infinite number of tuples, since variables y and z are not bounded by relation R , thus there can be an infinite number of evaluations that satisfy the rule body. Hence, the concept of *domain restricted* query is introduced. Rule based query q is *domain restricted*, if all variables that occur in the rule body also occur in some relational atom.

The second problem is that equality atoms may cause the body of a rule become

unsatisfiable, in contrast to Proposition 19.4. For example, query

$$ans(x) \leftarrow R(x), x = a, x = b \quad (19.12)$$

is domain restricted, however if a and b are distinct constants, then the answer will be empty. It is easy to check whether a rule based query with equality atoms is satisfiable.

SATISFIABLE(q)

- 1 Compute the transitive closure of equalities of the body of q .
- 2 **if** Two distinct constants should be equal by transitivity
- 3 **then return** “Not satisfiable.”
- 4 **else return** “Satisfiable.”

It is also true (Exercise 19.1-4) that if a rule based query q that contains equality atoms is satisfiable, then there exists a another rule based query q' without equalities that is equivalent with q .

Disjunction – union. The question 19.5. cannot be expressed with conjunctive queries. However, if the union operator is added to relational algebra, then 19.5. can be expressed in that extended relational algebra:

$$\pi_{Theater}(\sigma_{Title="La Dolce Vita"}(Show) \cup \sigma_{Title="Rashomon"}(Show)) . \quad (19.13)$$

Rule based queries are also capable of expressing question 19.5. if it is allowed that the same relation is in the head of many distinct rules:

$$\begin{aligned} ans(x_M) &\leftarrow Show(x_{Th}, "La Dolce Vita", x_{Ti}) , \\ ans(x_M) &\leftarrow Show(x_{Th}, "Rashomon", x_{Ti}) . \end{aligned} \quad (19.14)$$

Non-recursive datalog program is a generalisation of this.

Definition 19.7 A **non-recursive datalog program** over schema \mathbf{R} is a set of rules

$$\begin{aligned} S_1(u_1) &\leftarrow body_1 \\ S_2(u_2) &\leftarrow body_2 \\ &\vdots \\ S_m(u_m) &\leftarrow body_m , \end{aligned} \quad (19.15)$$

where no relation of \mathbf{R} occurs in a head, the same relation may occur in the head of several rules, furthermore there exists an ordering r_1, r_2, \dots, r_m of the rules such that the relation in the head of r_i does not occur in the body of any rule r_j for $j \leq i$.

The semantics of the non-recursive datalog program (19.15) is similar to the conjunctive query program (19.5). The rules are evaluated in the order r_1, r_2, \dots, r_m of Definition 19.7, and if a relation occurs in more than one head then the union of the sets of tuples given by those rules is taken.

The union of tableau queries (\mathbf{T}_i, u) $i = 1, 2, \dots, n$ is denoted by $(\{\mathbf{T}_1, \mathbf{T}_2, \dots, \mathbf{T}_n\}, u)$. It is evaluated by individually computing the result of each tableau query (\mathbf{T}_i, u) , then the union of them is taken. The following holds.

Theorem 19.8 *The language of non-recursive datalog programs with unique output relation and the relational algebra extended with the union operator are equivalent.*

The proof of Theorem 19.8 is similar to that of Theorem 19.6 and it is left to the Reader (Exercise 19.1-5). Let us note that the expressive power of the union of tableau queries is weaker. This is caused by the requirement having the same summary row for each tableau. For example, the non-recursive datalog program query

$$\begin{aligned} ans(a) &\leftarrow \\ ans(b) &\leftarrow \end{aligned} \tag{19.16}$$

cannot be realised as union of tableau queries.

Negation. The query 19.6. is obviously not monotone. Indeed, suppose that in relation *Film* there exist tuples about Hitchcock’s film *Psycho*, for example (“Psycho”, “A. Hitchcock”, “A. Perkins”), (“Psycho”, “A. Hitchcock”, “J. Leigh”), . . . , however, the tuple (“Psycho”, “A. Hitchcock”, “A. Hitchcock”) is not included. Then the tuple (“Psycho”) occurs in the output of query 19.6. With some effort one can realize however, that Hitchcock appears in the film *Psycho*, as “a man in cowboy hat”. If the tuple (“Psycho”, “A. Hitchcock”, “A. Hitchcock”) is added to relation *Film* as a consequence, then the instance of schema **CinePest** gets larger, but the output of query 19.6. becomes smaller.

It is not too hard to see that the query languages discussed so far are monotone, hence query 19.6. cannot be formulated with non-recursive datalog program or with some of its equivalents. Nevertheless, if the difference (−) operator is also added to relation algebra, then it becomes capable of expressing queries of type 19.6. For example,

$$\pi_{Title}(\sigma_{Director="A. Hitchcock"}(Film)) - \pi_{Title}(\sigma_{Actor="A. Hitchcock"}(Film)) \tag{19.17}$$

realises exactly query 19.6. Hence, the (full) relational algebra consists of operations $\{\sigma, \pi, \bowtie, \delta, \cup, -\}$. The importance of the relational algebra is shown by the fact, that Codd calls a query language \mathcal{Q} **relationally complete**, exactly if for all relational algebra query q there exists $q' \in \mathcal{Q}$, such that $q \equiv q'$.

If **negative literals**, that is atoms of the form $\neg R(u)$ are also allowed in rule bodies, then the obtained **non-recursive datalog with negation**, in notation **nr-datalog[−]** is relationally complete.

Definition 19.9 *A non-recursive datalog[−] (nr-datalog[−]) rule is of form*

$$q : S(u) \leftarrow L_1, L_2, \dots, L_n, \tag{19.18}$$

where S is a relation, u is a free tuple, L_i ’s are **literals**, that is expression of form $R(v)$ or $\neg R(v)$, such that v is a free tuple for $i = 1, 2, \dots, n$. S does not occur in the body of the rule. The rule is **domain restricted**, if each variable x that occurs in the rule also occurs in a **positive literal** (expression of the form $R(v)$) of the body. Every nr-datalog[−] rule is considered domain restricted, unless it is specified otherwise.

The semantics of rule (19.18) is as follows. Let \mathbf{R} be a relational schema that contains all relations occurring in the body of q , furthermore, let \mathcal{I} be an instance over \mathbf{R} . The *image of \mathcal{I} under q* is

$$q(\mathcal{I}) = \{\nu(u) \mid \nu \text{ is an valuation of the variables and for } i = 1, 2, \dots, n \\ \nu(u_i) \in \mathcal{I}(R_i), \text{ if } L_i = R_i(u_i) \text{ and} \\ \nu(u_i) \notin \mathcal{I}(R_i), \text{ if } L_i = \neg R_i(u_i)\}. \quad (19.19)$$

A *nr-datalog $^\neg$ program* over schema \mathbf{R} is a collection of nr-datalog $^\neg$ rules

$$\begin{aligned} S_1(u_1) &\leftarrow \text{body}_1 \\ S_2(u_2) &\leftarrow \text{body}_2 \\ &\vdots \\ S_m(u_m) &\leftarrow \text{body}_m, \end{aligned} \quad (19.20)$$

where relations of schema \mathbf{R} do not occur in heads of rules, the same relation may appear in more than one rule head, furthermore there exists an ordering r_1, r_2, \dots, r_m of the rules such that the relation of the head of rule r_i does not occur in the head of any rule r_j if $j \leq i$.

The computation of the result of nr-datalog $^\neg$ program (19.20) applied to instance \mathcal{I} over schema \mathbf{R} can be done in the same way as the evaluation of non-recursive datalog program (19.15), with the difference that the individual nr-datalog $^\neg$ rules should be interpreted according to (19.19).

Example 19.5 *Nr-datalog $^\neg$ program.* Let us assume that all films that are included in relation *Film* have only one director. (It is not always true in real life!) The nr-datalog $^\neg$ rule

$$\text{ans}(x) \leftarrow \text{Film}(x, \text{"A. Hitchcock"}, z), \neg \text{Film}(x, \text{"A. Hitchcock"}, \text{"A. Hitchcock"}) \quad (19.21)$$

expresses query **19.6**. Query **19.7**. is realised by the nr-datalog $^\neg$ program

$$\begin{aligned} \text{Fellini-actor}(z) &\leftarrow \text{Film}(x, \text{"F. Fellini"}, z) \\ \text{Not-the-answer}(x) &\leftarrow \text{Film}(x, y, z), \neg \text{Fellini-actor}(z) \\ \text{Answer}(x) &\leftarrow \text{Film}(x, y, z), \neg \text{Not-the-answer}(x). \end{aligned} \quad (19.22)$$

One has to be careful in writing nr-datalog $^\neg$ programs. If the first two rules of program (19.22) were to be merged like in Example 19.4

$$\begin{aligned} \text{Bad-not-ans}(x) &\leftarrow \text{Film}(x, y, z), \neg \text{Film}(x', \text{"F. Fellini"}, z), \text{Film}(x', \text{"F. Fellini"}, z') \\ \text{Answer}(x) &\leftarrow \text{Film}(x, y, z), \neg \text{Bad-not-ans}(x), \end{aligned} \quad (19.23)$$

then (19.23) answers the following query (assuming that all films have unique director)

19.9. List all those films whose every actor played in *each* film of Fellini, instead of query **19.7**.

It is easy to see that every satisfiable nr-datalog $^\neg$ program that contains equality atoms can be replaced by one without equalities. Furthermore the following proposition is true, as well.

Claim 19.10 *The satisfiable (full) relational algebra and the nr-datalog⁻ programs with single output relation are equivalent query languages.*

Recursion. Query 19.8. cannot be formulated using the query languages introduced so far. Some *a priori* information would be needed about how long a *chain-of-actors* could be formed starting with a given actor/actress. Let us assume that the maximum length of a chain-of-actors starting from “Marcello Mastroianni” is 117. (It would be interesting to know the *real* value!) Then, the following non-recursive datalog program gives the answer.

$$\begin{aligned}
 \text{Film-partner}(z_1, z_2) &\leftarrow \text{Film}(x, y, z_1), \text{Film}(x, y, z_2), z_1 < z_2^2 \\
 \text{Partial-answer}_1(z) &\leftarrow \text{Film-partner}(z, \text{“Marcello Mastroianni”}) \\
 \text{Partial-answer}_1(z) &\leftarrow \text{Film-partner}(\text{“Marcello Mastroianni”}, z) \\
 \text{Partial-answer}_2(z) &\leftarrow \text{Film-partner}(z, y), \text{Partial-answer}_1(y) \\
 \text{Partial-answer}_2(z) &\leftarrow \text{Film-partner}(y, z), \text{Partial-answer}_1(y) \\
 &\vdots \\
 \text{Partial-answer}_{117}(z) &\leftarrow \text{Film-partner}(z, y), \text{Partial-answer}_{116}(y) \\
 \text{Partial-answer}_{117}(z) &\leftarrow \text{Film-partner}(y, z), \text{Partial-answer}_{116}(y) \\
 \text{Mastroianni-chain}(z) &\leftarrow \text{Partial-answer}_1(z) \\
 \text{Mastroianni-chain}(z) &\leftarrow \text{Partial-answer}_2(z) \\
 &\vdots \\
 \text{Mastroianni-chain}(z) &\leftarrow \text{Partial-answer}_{117}(z)
 \end{aligned} \tag{19.24}$$

It is much easier to express query 19.8. using *recursion*. In fact, the *transitive closure* of the graph *Film-partner* needs to be calculated. For the sake of simplicity the definition of *Film-partner* is changed a little (thus approximately doubling the storage requirement).

$$\begin{aligned}
 \text{Film-partner}(z_1, z_2) &\leftarrow \text{Film}(x, y, z_1), \text{Film}(x, y, z_2) \\
 \text{Chain-partner}(x, y) &\leftarrow \text{Film-partner}(x, y) \\
 \text{Chain-partner}(x, y) &\leftarrow \text{Film-partner}(x, z), \text{Chain-partner}(z, y) .
 \end{aligned} \tag{19.25}$$

The datalog program (19.25) is *recursive*, since the definition of relation *Chain-partner* uses the relation itself. Let us suppose for a moment that this is meaningful, then query 19.8. is answered by rule

$$\text{Mastroianni-chain}(x) \leftarrow \text{Chain-partner}(x, \text{“Marcello Mastroianni”}) \tag{19.26}$$

Definition 19.11 *The expression*

$$R_1(u_1) \leftarrow R_2(u_2), R_3(u_3), \dots, R_n(u_n) \tag{19.27}$$

is a datalog rule, if $n \geq 1$, the R_i 's are relation names, the u_i 's are free tuples of

²Arbitrary comparison atoms can be used, as well, similarly to equality atoms. Here $z_1 < z_2$ makes it sure that all pairs occur at most once in the list.

appropriate length. Every variable of u_1 has to occur in one of u_2, \dots, u_n , as well. The head of the rule is $R_1(u_1)$, the body of the rule is $R_2(u_2), R_3(u_3), \dots, R_n(u_n)$. A **datalog program** is a finite collection of rules of type (19.27). Let P be a datalog program. The relation R occurring in P is **extensional** if it occurs in only rule bodies, and it is **intensional** if it occurs in the head of some rule.

If ν is a valuation of the variables of rule (19.27), then $R_1(\nu(u_1)) \leftarrow R_2(\nu(u_2)), R_3(\nu(u_3)), \dots, R_n(\nu(u_n))$ is a **realisation** of rule (19.27). The **extensional (database) schema** of P consists of the extensional relations of P , in notation $edb(P)$. The **intensional schema** of P , in notation $idb(P)$ is defined similarly as consisting of the intensional relations of P . Let $sch(P) = edb(P) \cup idb(P)$. The semantics of datalog program P is a mapping from the set of instances over $edb(P)$ to the set of instances over $idb(P)$. This can be defined proof theoretically, model theoretically or as a fixpoint of some operator. This latter one is equivalent with the first two, so to save space only the fixpoint theoretical definition is discussed.

There are no negative literals used in Definition 19.11. The main reason of this is that recursion and negation together may be meaningless, or contradictory. Nevertheless, sometimes negative atoms might be necessary. In those cases the semantics of the program will be defined specially.

Fixpoint semantics. Let P be a datalog program, \mathcal{K} be an instance over $sch(P)$.

Fact A , that is a tuple consisting of constants is an **immediate consequence** of \mathcal{K} and P , if either $A \in \mathcal{K}(R)$ for some relation $R \in sch(P)$, or $A \leftarrow A_1, A_2, \dots, A_n$ is a realisation of a rule in P and each A_i is in \mathcal{K} . The **immediate consequence operator** T_P is a mapping from the set of instances over $sch(P)$ to itself. $T_P(\mathcal{K})$ consists of all immediate consequences of \mathcal{K} and P .

Claim 19.12 *The immediate consequence operator T_P is monotone.*

Proof Let \mathcal{I} and \mathcal{J} be instances over $sch(P)$, such that $\mathcal{I} \subseteq \mathcal{J}$. Let A be a fact of $T_P(\mathcal{I})$. If $A \in \mathcal{I}(R)$ for some relation $R \in sch(P)$, then $A \in \mathcal{J}(R)$ is implied by $\mathcal{I} \subseteq \mathcal{J}$. on the other hand, if $A \leftarrow A_1, A_2, \dots, A_n$ is a realisation of a rule in P and each A_i is in \mathcal{I} , then $A_i \in \mathcal{J}$ also holds. ■

The definition of T_P implies that $\mathcal{K} \subseteq T_P(\mathcal{K})$. Using Proposition 19.12 it follows that

$$\mathcal{K} \subseteq T_P(\mathcal{K}) \subseteq T_P(T_P(\mathcal{K})) \subseteq \dots \quad (19.28)$$

Theorem 19.13 *For every instance \mathcal{I} over schema $sch(P)$ there exists a unique minimal instance $\mathcal{K} \subseteq \mathcal{I}$ that is a **fixpoint** of T_P , i.e. $\mathcal{K} = T_P(\mathcal{K})$.*

Proof Let $T_P^i(\mathcal{I})$ denote the consecutive application of operator T_P i -times, and let

$$\mathcal{K} = \bigcup_{i=0}^{\infty} T_P^i(\mathcal{I}) \quad (19.29)$$

By the monotonicity of T_P and (19.29) we have

$$T_P(\mathcal{K}) = \bigcup_{i=1}^{\infty} T_P^i(\mathcal{I}) \subseteq \bigcup_{i=0}^{\infty} T_P^i(\mathcal{I}) = \mathcal{K} \subseteq T_P(\mathcal{K}), \quad (19.30)$$

that is \mathcal{K} is a fixpoint. It is easy to see that every fixpoint that contains \mathcal{I} , also contains $T_P^i(\mathcal{I})$ for all $i = 1, 2, \dots$, that is it contains \mathcal{K} , as well. ■

Definition 19.14 *The **result** of datalog program P on instance \mathcal{I} over $edb(P)$ is the unique minimal fixpoint of T_P containing \mathcal{I} , in notation $P(\mathcal{I})$.*

It can be seen, see Exercise 19.1-6, that the chain in (19.28) is finite, that is there exists an n , such that $T_P(T_P^n(\mathcal{I})) = T_P^n(\mathcal{I})$. The naive evaluation of the result of the datalog program is based on this observation.

NAIV-DATALOG(P, \mathcal{I})

```

1  $\mathcal{K} \leftarrow \mathcal{I}$ 
2 while  $T_P(\mathcal{K}) \neq \mathcal{K}$ 
3     do  $\mathcal{K} \leftarrow T_P(\mathcal{K})$ 
4 return  $\mathcal{K}$ 

```

Procedure NAIV-DATALOG is not optimal, of course, since every fact that becomes included in \mathcal{K} is calculated again and again at every further execution of the **while** loop.

The idea of SEMI-NAIV-DATALOG is that it tries to use only recently calculated new facts in the **while** loop, as much as it is possible, thus avoiding recomputation of known facts. Consider datalog program P with $edb(P) = \mathbf{R}$, and $idb(P) = \mathbf{T}$. For a rule

$$S(u) \leftarrow R_1(v_1), \dots, R_n(v_n), T_1(w_1), \dots, T_m(w_m) \quad (19.31)$$

of P where $R_k \in \mathbf{R}$ and $T_j \in \mathbf{T}$, the following rules are constructed for $j = 1, 2, \dots, m$ and $i \geq 1$

$$temp_S^{i+1}(u) \leftarrow R_1(v_1), \dots, R_n(v_n), \\ T_1^i(w_1), \dots, T_{j-1}^i(w_{j-1}), \Delta_{T_j}^i(w_j), T_{j+1}^{i-1}(w_{j+1}), \dots, T_m^{i-1}(w_m). \quad (19.32)$$

Relation $\Delta_{T_j}^i$ denotes the change of T_j in iteration i . The union of rules corresponding to S in layer i is denoted by P_S^i , that is rules of form (19.32) for $temp_S^{i+1}$, $j = 1, 2, \dots, m$. Assume that the list of idb relations occurring in rules defining the idb relation S is T_1, T_2, \dots, T_ℓ . Let

$$P_S^i(\mathcal{I}, T_1^{i-1}, \dots, T_\ell^{i-1}, T_1^i, \dots, T_\ell^i, \Delta_{T_1}^i, \dots, \Delta_{T_\ell}^i) \quad (19.33)$$

denote the set of facts (tuples) obtained by applying rules (19.32) to input instance \mathcal{I} and to idb relations $T_j^{i-1}, T_j^i, \Delta_{T_j}^i$. The input instance \mathcal{I} is the actual value of the edb relations of P .

SEMI-NAIV-DATALOG(P, \mathcal{I})

```

1   $P' \leftarrow$  those rules of  $P$  whose body does not contain  $idb$  relation
2  for  $S \in idb(P)$ 
3      do  $S^0 \leftarrow \emptyset$ 
4           $\Delta_S^1 \leftarrow P'(\mathcal{I})(S)$ 
5   $i \leftarrow 1$ 
6  repeat
7      for  $S \in idb(P)$ 
8           $\triangleright T_1, \dots, T_\ell$  are the  $idb$  relations of the rules defining  $S$ .
9          do  $S^i \leftarrow S^{i-1} \cup \Delta_S^i$ 
10              $\Delta_S^{i+1} \leftarrow P_S^i(\mathcal{I}, T_1^{i-1}, \dots, T_\ell^{i-1}, T_1^i, \dots, T_\ell^i, \Delta_{T_1}^i, \dots, \Delta_{T_\ell}^i) - S^i$ 
11              $i \leftarrow i + 1$ 
12 until  $\Delta_S^i = \emptyset$  for all  $S \in idb(P)$ 
13 for  $S \in idb(P)$ 
14     do  $S \leftarrow S^i$ 
15 return  $\mathcal{S}$ 

```

Theorem 19.15 *Procedure SEMI-NAIV-DATALOG correctly computes the result of program P on instance \mathcal{I} .*

Proof We will show by induction on i that after execution of the loop of lines 6–12 i^{th} times the value of S^i is $T_P^i(\mathcal{I})(S)$, while Δ_S^{i+1} is equal to $T_P^{i+1}(\mathcal{I})(S) - T_P^i(\mathcal{I})(S)$ for arbitrary $S \in idb(P)$. $T_P^i(\mathcal{I})(S)$ is the result obtained for S starting from \mathcal{I} and applying the immediate consequence operator T_P i -times.

For $i = 0$, line 4 calculates exactly $T_P(\mathcal{I})(S)$ for all $S \in idb(P)$. In order to prove the induction step, one only needs to see that $P_S^i(\mathcal{I}, T_1^{i-1}, \dots, T_\ell^{i-1}, T_1^i, \dots, T_\ell^i, \Delta_{T_1}^i, \dots, \Delta_{T_\ell}^i) \cup S^i$ is exactly equal to $T_P^{i+1}(\mathcal{I})(S)$, since in lines 9–10 procedure SEMI-NAIV-DATALOG constructs S^i -t and Δ_S^{i+1} using that. The value of S^i is $T_P^i(\mathcal{I})(S)$, by the induction hypothesis. Additional new tuples are obtained only if that for some idb relation defining S such tuples are considered that are constructed at the last application of T_P , and these are in relations $\Delta_{T_1}^i, \dots, \Delta_{T_\ell}^i$, also by the induction hypothesis.

The halting condition of line 12 means exactly that all relations $S \in idb(P)$ are unchanged during the application of the immediate consequence operator T_P , thus the algorithm has found its minimal fixpoint. This latter one is exactly the result of datalog program P on input instance \mathcal{I} according to Definition 19.14. ■

Procedure SEMI-NAIV-DATALOG eliminates a large amount of unnecessary calculations, nevertheless it is not optimal on some datalog programs (Exercise gy:snaiv). However, analysis of the datalog program and computation based on that can save most of the unnecessary calculations.

Definition 19.16 *Let P be a datalog program. The **precedence graph** of P is the directed graph G_P defined as follows. Its vertex set consists of the relations of $idb(P)$, and (R, R') is an arc for $R, R' \in idb(P)$ if there exists a rule in P whose*

head is R' and whose body contains R . P is **recursive**, if G_P contains a directed cycle. Relations R and R' are **mutually recursive** if they belong to the same strongly connected component of G_P .

Being mutually recursive is an equivalence relation on the set of relations $idb(P)$. The main idea of procedure IMPROVED-SEMI-NAIV-DATALOG is that for a relation $R \in idb(P)$ only those relations have to be computed “simultaneously” with R that are in its equivalence class, all other relations defining R can be calculated “in advance” and can be considered as *edb* relations.

IMPROVED-SEMI-NAIV-DATALOG(P, \mathcal{I})

- 1 Determine the equivalence classes of $idb(P)$ under mutual recursivity.
- 2 List the equivalence classes $[R_1], [R_2], \dots, [R_n]$
according to a topological order of G_P .
- 3 \triangleright There exists no directed path from R_j to R_i in G_P for all $i < j$.
- 4 **for** $i \leftarrow 1$ **to** n
- 5 **do** Use SEMI-NAIV-DATALOG to compute relations of $[R_i]$
taking relations of $[R_j]$ as *edb* relations for $j < i$.

Lines 1–2 can be executed in time $O(v_{G_P} + e_{G_P})$ using depth first search, where v_{G_P} and e_{G_P} denote the number of vertices and edges of graph G_P , respectively. Proof of correctness of the procedure is left to the Reader (Exercise 19.1-8).

19.1.3. Complexity of query containment

In the present section we return to conjunctive queries. The costliest task in computing result of a query is to generate the natural join of relations. In particular, if there are no indexes available on the common attributes, hence only procedure FULL-TUPLEWISE-JOIN is applicable.

FULL-TUPLEWISE-JOIN(R_1, R_2)

- 1 $S \leftarrow \emptyset$
- 2 **for** all $u \in R_1$
- 3 **do for** all $v \in R_2$
- 4 **do if** u and v can be joined
- 5 **then** $S \leftarrow S \cup \{u \bowtie v\}$
- 6 **return** S

It is clear that the running time of FULL-TUPLEWISE-JOIN is $O(|R_1| \times |R_2|)$. Thus, it is important that in what order is a query executed, since during computation natural joins of relations of various sizes must be calculated. In case of tableau queries the **Homomorphism Theorem** gives a possibility of a query rewriting that uses less joins than the original.

Let q_1, q_2 be queries over schema \mathbf{R} . q_2 **contains** q_1 , in notation $q_1 \sqsubseteq q_2$, if for all instances \mathcal{I} over schema \mathbf{R} $q_1(\mathcal{I}) \subseteq q_2(\mathcal{I})$ holds. $q_1 \equiv q_2$ according to Definition 19.1 iff $q_1 \sqsubseteq q_2$ and $q_1 \supseteq q_2$. A generalisation of valuations will be needed. **Substitution**

is a mapping from the set of variables to the union of sets of variables and sets of constants that is extended to constants as identity. Extension of substitution to free tuples and tableaux can be defined naturally.

Definition 19.17 Let $q = (\mathbf{T}, u)$ and $q' = (\mathbf{T}', u')$ be two tableau queries over schema \mathbf{R} . Substitution θ is a **homomorphism** from q' to q , if $\theta(\mathbf{T}') = \mathbf{T}$ and $\theta(u') = u$.

Theorem 19.18 (Homomorphism Theorem). Let $q = (\mathbf{T}, u)$ and $q' = (\mathbf{T}', u')$ be two tableau queries over schema \mathbf{R} . $q \sqsubseteq q'$ if and only if, there exists a homomorphism from q' to q .

Proof Assume first, that θ is a homomorphism from q' to q , and let \mathcal{I} be an instance over schema \mathbf{R} . Let $w \in q(\mathcal{I})$. This holds exactly if there exists a valuation ν that maps tableau \mathbf{T} into \mathcal{I} and $\nu(u) = w$. It is easy to see that $\theta \circ \nu$ maps tableau \mathbf{T}' into \mathcal{I} and $\theta \circ \nu(u') = w$, that is $w \in q'(\mathcal{I})$. Hence, $w \in q(\mathcal{I}) \implies w \in q'(\mathcal{I})$, which in turn is equivalent with $q \sqsubseteq q'$.

On the other hand, let us assume that $q \sqsubseteq q'$. The idea of the proof is that both, q and q' are applied to the “instance” \mathbf{T} . The output of q is free tuple u , hence the output of q' also contains u , that is there exists a θ embedding of \mathbf{T}' into \mathbf{T} that maps u' to u . To formalise this argument the instance $\mathcal{I}_{\mathbf{T}}$ isomorphic to \mathbf{T} is constructed.

Let V be the set of variables occurring in \mathbf{T} . For all $x \in V$ let a_x be constant that differs from all constants appearing in \mathbf{T} or \mathbf{T}' , furthermore $xNEa_x' \implies a_xNEa_x'$. Let μ be the valuation that maps $x \in V$ to a_x , furthermore let $\mathcal{I}_{\mathbf{T}} = \mu(\mathbf{T})$. μ is a bijection from V to $\mu(V)$ and there are no constants of \mathbf{T} appearing in $\mu(V)$, hence μ^{-1} well defined on the constants occurring in $\mathcal{I}_{\mathbf{T}}$.

It is clear that $\mu(u) \in q(\mathcal{I}_{\mathbf{T}})$, thus using $q \sqsubseteq q'$ $\mu(u) \in q'(\mathcal{I}_{\mathbf{T}})$ is obtained. That is, there exists a valuation ν that embeds tableau \mathbf{T}' into $\mathcal{I}_{\mathbf{T}}$, such that $\nu(u') = \mu(u)$. It is not hard to see that $\nu \circ \mu^{-1}$ is a homomorphism from q' to q . ■

Query optimisation by tableau minimisation. According to Theorem 19.6 tableau queries and satisfiable relational algebra (without subtraction) are equivalent. The proof shows that the relational algebra expression equivalent with a tableau query is of form $\pi_{\overline{\mathbb{E}}}(\sigma_F(R_1 \bowtie R_2 \bowtie \dots \bowtie R_k))$, where k is the number of rows of the tableau. It implies that if the number of joins is to be minimised, then the number of rows of an equivalent tableau must be minimised.

The tableau query (\mathbf{T}, u) is **minimal**, if there exists no tableau query (\mathbf{S}, v) that is equivalent with (\mathbf{T}, u) and $|\mathbf{S}| < |\mathbf{T}|$, that is S has fewer rows. It may be surprising, but it is true, that a minimal tableau query equivalent with (\mathbf{T}, u) can be obtained by simply dropping some rows from \mathbf{T} .

Theorem 19.19 Let $q = (\mathbf{T}, u)$ be a tableau query. There exists a subset \mathbf{T}' of \mathbf{T} , such that query $q' = (\mathbf{T}', u)$ is minimal and equivalent with $q = (\mathbf{T}, u)$.

Proof Let (\mathbf{S}, v) be a minimal query equivalent with q . According to the Homomorphism Theorem there exist homomorphisms θ from q to (\mathbf{S}, v) , and λ from (\mathbf{S}, v) to

q . Let $\mathbf{T}' = \theta \circ \lambda(\mathbf{T})$. It is easy to check that $(\mathbf{T}', u) \equiv q$ and $|\mathbf{T}'| \leq |\mathbf{S}|$. But (\mathbf{S}, v) is minimal, hence (\mathbf{T}', u) is minimal, as well. ■

Example 19.6 *Application of tableau minimisation* Consider the relational algebra expression

$$q = \pi_{AB}(\sigma_{B=5}(R)) \bowtie \pi_{BC}(\pi_{AB}(R) \bowtie \pi_{AC}(\sigma_{B=5}(R))) \tag{19.34}$$

over the schema \mathbf{R} of attribute set $\{A, B, C\}$. The tableau query corresponding to q is the following tableau \mathbf{T} :

R	A	B	C	(19.35)
x	5	z_1		
x_1	5	z_2		
x_1	5	z		
u	x	5	z	

Such a homomorphism is sought that maps some rows of \mathbf{T} to some other rows of \mathbf{T} , thus sort of “folding” the tableau. The first row cannot be eliminated, since the homomorphism is an identity on free tuple u , thus x must be mapped to itself. The situation is similar with the third row, as the image of z is itself under any homomorphism. However the second row can be eliminated by mapping x_1 to x and z_2 to z , respectively. Thus, the minimal tableau equivalent with \mathbf{T} consists of the first and third rows of \mathbf{T} . Transforming back to relational algebra expression,

$$\pi_{AB}(\sigma_{B=5}(R)) \bowtie \pi_{BC}(\sigma_{B=5}(R)) \tag{19.36}$$

is obtained. Query (19.36) contains one less join operator than query (19.34).

The next theorem states that the question of tableau containment and equivalence is NP-complete, hence tableau minimisation is an NP-hard problem.

Theorem 19.20 *For given tableau queries q and q' the following decision problems are NP-complete:*

19.10. $q \sqsubseteq q'$?

19.11. $q \equiv q'$?

19.12. Assume that q' is obtained from q by removing some free tuples. Is it true then that $q \equiv q'$?

Proof The EXACT-COVER problem will be reduced to the various tableau problems. The input of EXACT-COVER problem is a finite set $X = \{x_1, x_2, \dots, x_n\}$, and a collection of its subsets $\mathcal{S} = \{S_1, S_2, \dots, S_m\}$. It has to be determined whether there exists $\mathcal{S}' \sqsubseteq \mathcal{S}$, such that subsets in \mathcal{S}' cover X exactly (that is, for all $x \in X$ there exists exactly one $S \in \mathcal{S}'$ such that $x \in S$). EXACT-COVER is known to be an NP-complete problem.

Let $\mathcal{E} = (X, \mathcal{S})$ be an input of EXACT COVER. A construction is sketched that produces a pair $q_{\mathcal{E}}, q'_{\mathcal{E}}$ of tableau queries to \mathcal{E} in polynomial time. This construction can be used to prove the various NP-completeness results.

Let the schema \mathbf{R} consist of the pairwise distinct attributes $A_1, A_2, \dots, A_n, B_1, B_2, \dots, B_m$. $q_{\mathcal{E}} = (\mathbf{T}_{\mathcal{E}}, t)$ and $q'_{\mathcal{E}} = (\mathbf{T}'_{\mathcal{E}}, t)$ are tableau

queries over schema \mathbf{R} such that the summary row of both of them is free tuple $t = \langle A_1 : a_1, A_2 : a_2, \dots, A_n : a_n \rangle$, where a_1, a_2, \dots, a_n are pairwise distinct variables.

Let $b_1, b_2, \dots, b_m, c_1, c_2, \dots, c_m$ be another set of pairwise distinct variables. Tableau $\mathbf{T}_\mathcal{E}$ consists of n rows, for each element of X corresponds one. a_i stands in column of attribute A_i in the row of x_i , while b_j stands in column of attribute B_j for all such j that $x_i \in S_j$ holds. In other positions of tableau $\mathbf{T}_\mathcal{E}$ pairwise distinct new variables stand.

Similarly, $\mathbf{T}'_\mathcal{E}$ consists of m rows, one corresponding to each element of \mathcal{S} . a_i stands in column of attribute A_i in the row of S_j for all such i that $x_i \in S_j$, furthermore $c_{j'}$ stands in the column of attribute $B_{j'}$, for all $j' \in \mathbf{NE}j$. In other positions of tableau $\mathbf{T}'_\mathcal{E}$ pairwise distinct new variables stand.

The NP-completeness of problem 19.10. follows from that X has an exact cover using sets of \mathcal{S} if and only if $q'_\mathcal{E} \sqsubseteq q_\mathcal{E}$ holds. The proof, and the proof of the NP-completeness of problems 19.11. and 19.12. are left to the Reader (Exercise 19.1-9). ■

Exercises

19.1-1 Prove Proposition 19.4, that is every rule based query q is monotone and satisfiable. *Hint.* For the proof of satisfiability let K be the set of constants occurring in query q , and let $a \notin K$ be another constant. For every relation schema R_i in rule (19.3) construct all tuples (a_1, a_2, \dots, a_r) , where $a_i \in K \cup \{a\}$, and r is the arity of R_i . Let \mathcal{I} be the instance obtained so. Prove that $q(\mathcal{I})$ is nonempty.

19.1-2 Give a relational schema \mathbf{R} and a relational algebra query q over schema \mathbf{R} , whose result is empty to arbitrary instance over \mathbf{R} .

19.1-3 Prove that the languages of rule based queries and tableau queries are equivalent.

19.1-4 Prove that every rule based query q with equality atoms is either equivalent with the empty query q^0 , or there exists a rule based query q' without equality atoms such that $q \equiv q'$. Give a polynomial time algorithm that determines for a rule based query q with equality atoms whether $q \equiv q^0$ holds, and if not, then constructs a rule based query q' without equality atoms, such that $q \equiv q'$.

19.1-5 Prove Theorem 19.8 by generalising the proof idea of Theorem 19.6.

19.1-6 Let P be a datalog program, \mathcal{I} be an instance over $edb(P)$, $inC(P, \mathcal{I})$ be the (finite) set of constants occurring in \mathcal{I} and P . Let $\mathbf{B}(P, \mathcal{I})$ be the following instance over $sch(P)$:

1. For every relation R of $edb(P)$ the fact $R(u)$ is in $\mathbf{B}(P, \mathcal{I})$ iff it is in \mathcal{I} , furthermore
2. for every relation R of $idb(P)$ every $R(u)$ fact constructed from constants of $C(P, \mathcal{I})$ is in $\mathbf{B}(P, \mathcal{I})$.

Prove that

$$\mathcal{I} \subseteq T_P(\mathcal{I}) \subseteq T_P^2(\mathcal{K}) \subseteq T_P^3(\mathcal{K}) \subseteq \dots \subseteq \mathbf{B}(P, \mathcal{I}). \quad (19.37)$$

19.1-7 Give an example of an input, that is a datalog program P and instance \mathcal{I} over $edb(P)$, such that the same tuple is produced by distinct runs of the loop of SEMI-NAIV-DATALOG.

19.1-8 Prove that procedure IMPROVED-SEMI-NAIV-DATALOG stops in finite time for all inputs, and gives correct result. Give an example of an input on which IMPROVED-SEMI-NAIV-DATALOG calculates less number of rows multiple times than SEMI-NAIV-DATALOG.

19.1-9

1. Prove that for tableau queries $q_{\mathcal{E}} = (\mathbf{T}_{\mathcal{E}}, t)$ and $q'_{\mathcal{E}} = (\mathbf{T}'_{\mathcal{E}}, t)$ of the proof of Theorem 19.20 there exists a homomorphism from $(\mathbf{T}_{\mathcal{E}}, t)$ to $(\mathbf{T}'_{\mathcal{E}}, t)$ if and only if the EXACT-COVER problem $\mathcal{E} = (X, \mathcal{S})$ has solution.
2. Prove that the decision problems 19.11. and 19.12. are NP-complete.

19.2. Views

A database system architecture has three main levels:

- physical layer;
- logical layer;
- outer (user) layer.

The goal of the separation of levels is to reach data independence and the convenience of users. The three views on Figure 19.2 show possible user interfaces: multirelational, universal relation and graphical interface.

The *physical layer* consists of the actually stored data files and the dense and sparse indices built over them.

The separation of the *logical layer* from the physical layer makes it possible for the user to concentrate on the logical dependencies of the data, which approximates the image of the reality to be modelled better. The logical layer consists of the database schema description together with the various integrity constraints, dependencies. This the layer where the database administrators work with the system. The connection between the physical layer and the logical layer is maintained by the database engine.

The goal of the separation of the logical layer and the *outer layer* is that the endusers can see the database according to their (narrow) needs and requirements. For example, a very simple view of the outer layer of a bank database could be the automatic teller machine, or a much more complex view could be the credit history of a client for loan approval.

19.2.1. View as a result of a query

The question is that how can the views of different layers be given. If a query given by relational algebra expression is considered as a formula that will be applied to relational instances, then the *view* is obtained. Datalog rules show the difference

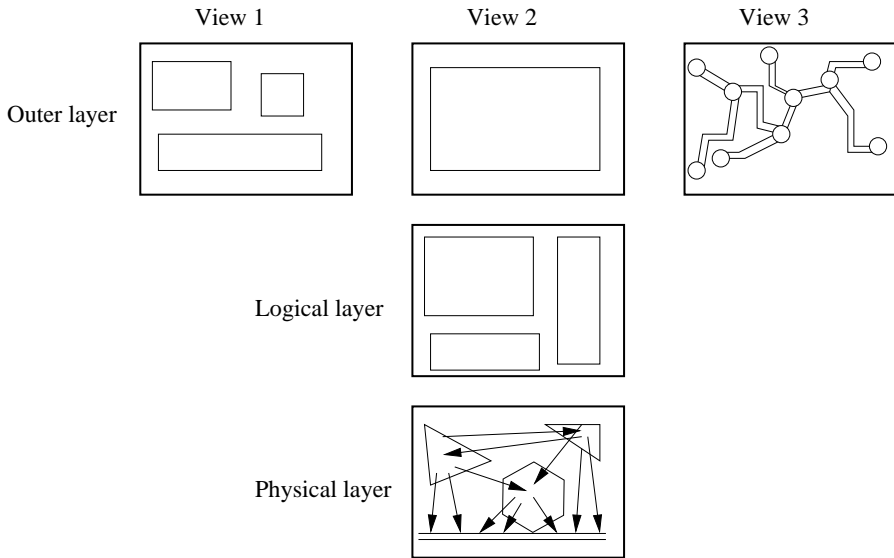


Figure 19.2 The three levels of database architecture.

between views and relations, well. The relations defined by rules are called *intensional*, because these are the relations that do not have to exist on external storage devices, that is to exist extensionally, in contrast to the *extensional* relations.

Definition 19.21 *The \mathcal{V} expression given in some query language \mathcal{Q} over schema \mathbf{R} is called a *view*.*

Similarly to intensional relations, views can be used in definition of queries or other views, as well.

Example 19.7 *SQL view.* Views in database manipulation language SQL can be given in the following way. Suppose that the only interesting data for us from schema **CinePest** is where and when are Kurosawa's film shown. The view **KurosawaTimes** is given by the SQL command

KUROSAWATIMES

```

1 create view KurosawaTimes as
2   select Theater, Time
3   from Film, Show
4   where Film.Title=Show.Title and Film.Director="Akira Kurosawa"
```

Written in relational algebra is as follows.

$$KurosawaTimes(Theater, Time) = \pi_{Theater, Time}(Theater \bowtie \sigma_{Director="Akira Kurosawa"}(Film)) \quad (19.38)$$

Finally, the same by datalog rule is:

$$KurosawaTimes(x_{Th}, x_{Ti}) \leftarrow Theater(x_{Th}, x_T, x_{Ti}), Film(x_T, "Akira Kurosawa", x_A) . \quad (19.39)$$

Line 2 of KUROSAWA TIMES marks the selection operator used, line 3 marks that which two relations are to be joined, finally the condition of line 4 shows that it is a natural join, not a direct product.

Having defined view \mathcal{V} , it can be used in further queries or view definitions like any other (extensional) relation.

Advantages of using views

- Automatic data hiding: Such data that is not part of the view used, is not shown to the user, thus the user cannot read or modify them without having proper access rights to them. So by providing access to the database through views, a simple, but effective security mechanism is created.
- Views provide simple “macro capabilities”. Using the view *KurosawaTimes* defined in Example 19.7 it is easy to find those theatres where Kurosawa films are shown in the morning:

$$KurosawaMorning(Theater) \leftarrow KurosawaTimes(Theater, x_{Ti}), x_{Ti} < 12 . \quad (19.40)$$

Of course the user could include the definition of *KurosawaTimes* in the code directly, however convenience considerations are first here, in close similarity with macros.

- Views make it possible that the same data could be seen in different ways by different users at the same time.
- Views provide **logical data independence**. The essence of logical data independence is that users and their programs are protected from the structural changes of the database schema. It can be achieved by defining the relations of the schema before the structural change as views in the new schema.
- Views make controlled data input possible. The **with check option** clause of command **create view** is to do this in SQL.

Materialised view. Some view could be used in several different queries. It could be useful in these cases that if the tuples of the relation(s) defined by the view need not be calculated again and again, but the output of the query defining the view is stored, and only read in at further uses. Such stored output is called a **materialised view**.

Exercises

19.2-1 Consider the following schema:

FilmStar(Name, Address, Gender, BirthDate)
FilmMogul(Name, Address, Certificate#, Assets)
Studio(Name, Address, PresidentCert#) .

Relation *FilmMogul* contains data of the big people in film business (studio presidents, producers, etc.). The attribute names speak for themselves, *Certificate#* is the number of the certificate of the filmmogul, *PresidentCert#* is the certificate number of the president of the studio. Give the definitions of the following views using datalog rules, relational algebra expressions, furthermore SQL:

1. *RichMogul*: Lists the names, addresses, certificate numbers and assets of those filmmoguls, whose asset value is over 1 million dollars.
2. *StudioPresident*: Lists the names, addresses and certificate numbers of those filmmoguls, who are studio presidents, as well.
3. *MogulStar*: Lists the names, addresses, certificate numbers and assets of those people who are filmstars and filmmoguls at the same time.

19.2-2 Give the definitions of the following views over schema **CinePest** using datalog rules, relational algebra expressions, furthermore SQL:

1. *Marilyn(Title)*: Lists the titles of Marilyn Monroe's films.
2. *CorvinInfo(Title,Time,Phone)*: List the titles and show times of films shown in theatre *Corvin*, together with the phone number of the theatre.

19.3. Query rewriting

Answering queries using views, in other words query rewriting using views has become a problem attracting much attention and interest recently. The reason is its applicability in a wide range of data manipulation problems: query optimisation, providing physical data independence, data and information integration, furthermore planning data warehouses.

The problem is the following. Assume that query Q is given over schema \mathbf{R} , together with views V_1, V_2, \dots, V_n . Can one answer Q using only the results of views V_1, V_2, \dots, V_n ? Or, which is the largest set of tuples that can be determined knowing the views? If the views and the relations of the base schema can be accessed both, what is the cheapest way to compute the result of Q ?

19.3.1. Motivation

Before query rewriting algorithms are studied in detail, some motivating examples of applications are given. The following university database is used throughout this section.

University = $\{Professor, Course, Teach, Registered, Major, Affiliation, Supervisor\}$.
(19.41)

The schemata of the individual relations are as follows:

$$\begin{aligned}
 \textit{Professor} &= \{PName, Area\} \\
 \textit{Course} &= \{C\text{-Number}, Title\} \\
 \textit{Teaches} &= \{PName, C\text{-Number}, Semester, Evaluation\} \\
 \textit{Registered} &= \{Student, C\text{-Number}, Semester\} \\
 \textit{Major} &= \{Student, Department\} \\
 \textit{Affiliation} &= \{PName, Department\} \\
 \textit{Advisor} &= \{PName, Student\} .
 \end{aligned} \tag{19.42}$$

It is supposed that professors, students and departments are uniquely identified by their names. Tuples of relation *Registered* show that which student took which course in what semester, while *Major* shows which department a student choose in majoring (for the sake of convenience it is assumed that one department has one subject as possible major).

Query optimisation. If the computation necessary to answer a query was performed in advance and the results are stored in some materialised view, then it can be used to speed up the query answering.

Consider the query that looks for pairs (*Student*, *Title*), where the student registered for the given Ph.D.-level course, the course is taught by a professor of the *Database* area (the C-number of graduate courses is at least 400, and the Ph.D.-level courses are those with C-number at least 500).

$$\textit{val}(x_S, x_T) \leftarrow \textit{Teach}(x_P, x_C, x_{S_e}, y_1), \textit{Professor}(x_P, \text{“database”}), \textit{Registered}(x_S, x_C, x_{S_e}), \textit{Course}(x_C, x_T), x_C \geq 500 . \tag{19.43}$$

Suppose that the following materialised view is available that contains the registration data of graduate courses.

$$\textit{Graduate}(x_S, x_T, x_C, x_{S_e}) \leftarrow \textit{Registered}(x_S, x_C, x_{S_e}), \textit{Course}(x_C, x_T), x_C \geq 400 . \tag{19.44}$$

View *Graduate* can be used to answer query (19.43).

$$\textit{val}(x_S, x_T) \leftarrow \textit{Teaches}(x_P, x_C, x_{S_e}, y_1), \textit{Professor}(x_P, \text{“database”}), (x_S, x_T, x_C, x_{S_e}), x_C \geq 500 . \tag{19.45}$$

It will be faster to compute (19.45) than to compute (19.43), because the natural join of relations *Registered* and *Course* has already be done by view *Graduate*, furthermore it shelled off the undergraduate courses (that make up for the bulk of registration data at most universities). It worth noting that view *Graduate* could be used event hough **syntactically** did not agree with any part of query (19.43).

On the other hand, it may happen that the the original query can be answered faster. If relations *Registered* and *Course* have an index on attribute *C-Number*, but there exists no index built for *Graduate*, then it could be faster to answer query (19.43) directly from the database relations. Thus, the real challenge is not only that to decide about a materialised view whether it could be used to answer some query logically, but a thorough cost analysis must be done when is it worth using the existing views.

Physical data independence. One of the principles underlying modern database systems is the separation between the logical view of data and the physical view of data. With the exception of horizontal or vertical partitioning of relations into multiple files, relational database systems are still largely based on a one-to-one correspondence between relations in the schema and the files in which they are stored. In object-oriented systems, maintaining the separation is necessary because the logical schema contains significant redundancy, and does not correspond to a good physical layout. Maintaining physical data independence becomes more crucial in applications where the logical model is introduced as an intermediate level after the physical representation has already been determined. This is common in storage of XML data in relational databases, and in data integration. In fact, the Stored System stores XML documents in a relational database, and uses views to describe the mapping from XML into relations in the database.

To maintain physical data independence, a widely accepted method is to use views over the logical schema as mechanism for describing the physical storage of data. In particular, Tsatalos, Solomon and Ioannidis use GMAPs (*Generalised Multi-level Access Paths*) to describe data storage.

A GMAP describes the physical organisation and the indexes of the storage structure. The first clause of the GMAP (**as**) describes the actual data structure used to store a set of tuples (e.g., a B⁺-tree, hash index, etc.) The remaining clauses describe the content of the structure, much like a view definition. The **given** and **select** clauses describe the available attributes, where the **given** clause describes the attributes on which the structure is indexed. The definition of the view, given in the **where** clause uses infix notation over the conceptual model.

In the example shown in Figure 19.3, the GMAP G1 stores a set of pairs containing students and the departments in which they major, and these pairs are indexed by a B⁺-tree on attribute *Student.name*. The GMAP G2 stores an index from names of students to the numbers of courses in which they are registered. The GMAP G3 stores an index from course numbers to departments whose majors are enrolled in the course.

Given that the data is stored in structures described by the GMAPs, the question that arises is how to use these structures to answer queries. Since the logical content of the GMAPs are described by views, answering a query amounts to finding a way of rewriting the query using these views. If there are multiple ways of answering the query using the views, we would like to find the cheapest one. Note that in contrast to the query optimisation context, we **must** use the views to answer the given query, because all the data is stored in the GMAPs.

Consider the following query, which asks for names of students registered for Ph.D.-level courses and the departments in which these students are majoring.

$$\begin{aligned} ans(Student, Department) \leftarrow & Registered(Student, C-number, y), \\ & Major(Student, Department), \\ & C-number \geq 500 . \end{aligned} \tag{19.46}$$

The query can be answered in two ways. First, since *Student.name* uniquely identifies a student, we can take the join of G1 and G2, and then apply selection operator $Course.c-number \geq 500$, finally a projection eliminates the unnecessary attributes.

```

def.gmap G1 as b+-tree by
  given Student.name
  select Department
  where Student major Department

def.gmap G2 as b+-tree by
  given Student.name
  select Course.c-number
  where Student registered Course

def.gmap G3 as b+-tree by
  given Course.c-number
  select Department
  where Student registered Course and Student major Department

```

Figure 19.3 GMAPs for the university domain.

The other execution plan could be to join G3 with G2 and select *Course.c-number* \geq 500. In fact, this solution may even be more efficient because G3 has an index on the course number and therefore the intermediate joins may be much smaller.

Data integration. A *data integration system* (also known as *mediator system*) provides a *uniform* query interface to a multitude of autonomous heterogeneous data sources. Prime examples of data integration applications include enterprise integration, querying multiple sources on the World-Wide Web, and integration of data from distributed scientific experiments.

To provide a uniform interface, a data integration system exposes to the user a *mediated schema*. A mediated schema is a set of *virtual* relations, in the sense that they are not stored anywhere. The mediated schema is designed manually for a particular data integration application. To be able to answer queries, the system must also contain a set of *source descriptions*. A description of a data source specifies the contents of the source, the attributes that can be found in the source, and the (integrity) constraints on the content of the source. A widely adopted approach for specifying source descriptions is to describe the contents of a data source as a *view* over the mediated schema. This approach facilitates the addition of new data sources and the specification of constraints on the contents of sources.

In order to answer a query, a data integration system needs to translate a query formulated on the mediated schema into one that refers directly to the schemata of the data sources. Since the contents of the data sources are described as views, the translation problem amounts finding a way to answer a query using a set of views.

Consider as an example the case where the mediated schema exposed to the user is schema **University**, except that the relations *Teaches* and *Course* have an additional attribute identifying the university at which the course is being taught:

$$\begin{aligned}
 \textit{Course} &= \{C\text{-number}, Title, Univ\} \\
 \textit{Teaches} &= \{PName, C\text{-number}, Semester, Evaluation, Univ.\}
 \end{aligned}
 \tag{19.47}$$

Suppose we have the following two data sources. The first source provides a listing of all the courses entitled "Database Systems" taught anywhere and their instructors. This source can be described by the following view definition:

$$\begin{aligned} DBcourse(Title, PName, C-number, Univ) \leftarrow & Course(C-number, Title, Univ), \\ & Teaches(PName, C-number, Semester, \\ & \quad Evaluation, Univ), \\ & Title = \text{"Database Systems"} . \end{aligned} \tag{19.48}$$

The second source lists Ph.D.-level courses being taught at The Ohio State University (OSU), and is described by the following view definition:

$$\begin{aligned} OSUPhD(Title, PName, C-number, Univ) \leftarrow & Course(C-number, Title, Univ), \\ & Teaches(PName, C-number, Semester, \\ & \quad Evaluation, Univ), \\ & Univ = \text{"OSU"}, C-number \geq 500. \end{aligned} \tag{19.49}$$

If we were to ask the data integration system who teaches courses titled "Database Systems" at OSU, it would be able to answer the query by applying a selection on the source *DB-courses*:

$$ans(PName) \leftarrow DBcourse(Title, PName, C-number, Univ), Univ = \text{"OSU"} . \tag{19.50}$$

On the other hand, suppose we ask for all the graduate-level courses (not just in databases) being offered at OSU. Given that only these two sources are available, the data integration system cannot find **all** tuples in the answer to the query. Instead, the system can attempt to find the maximal set of tuples in the answer available from the sources. In particular, the system can obtain graduate **database** courses at OSU from the *DB-course* source, and the Ph.D.-level courses at OSU from the *OSUPhD* source. Hence, the following non-recursive datalog program gives the maximal set of answers that can be obtained from the two sources:

$$\begin{aligned} ans(Title, C-number) \leftarrow & DBcourse(Title, PName, C-number, Univ), \\ & Univ = \text{"OSU"}, C-number \geq 400 \\ ans(Title, C-number) \leftarrow & OSUPhD(Title, PName, C-number, Univ) . \end{aligned} \tag{19.51}$$

Note that courses that are not Ph.D.-level courses or database courses will not be returned as answers. Whereas in the contexts of query optimisation and maintaining physical data independence the focus is on finding a query expression that is **equivalent** with the original query, here finding a query expression that provides the **maximal answers** from the views is attempted.

Semantic data caching. If the database is accessed via client-server architecture, the data cached at the client can be semantically modelled as results of certain queries, rather than at the physical level as a set of data pages or tuples. Hence, deciding which data needs to be shipped from the server in order to answer a given query requires an analysis of which parts of the query can be answered by the cached views.

19.3.2. Complexity problems of query rewriting

In this section the *theoretical* complexity of query rewriting is studied. Mostly conjunctive queries are considered. *Minimal*, and *complete* rewriting will be distinguished. It will be shown that if the query is conjunctive, furthermore the materialised views are also given as results of conjunctive queries, then the rewriting problem is *NP-complete*, assuming that neither the query nor the view definitions contain comparison atoms. Conjunctive queries will be considered in rule-based form for the sake of convenience.

Assume that query Q is given over schema \mathbf{R} .

Definition 19.22 *The conjunctive query Q' is a **rewriting** of query Q using views $\mathcal{V} = V_1, V_2, \dots, V_m$, if*

- Q and Q' are equivalent, and
- Q' contains one or more literals from \mathcal{V} .

Q' is said to be **locally minimal** if no literal can be removed from Q' without violating the equivalence. The rewriting is **globally minimal**, if there exists no rewriting using a smaller number of literals. (The comparison atoms $=, \neq, \leq, <$ are not counted in the number of literals.)

Example 19.8 *Query rewriting.* Consider the following query Q and view V .

$$\begin{aligned} Q: q(X, U) &\leftarrow p(X, Y), p_0(Y, Z), p_1(X, W), p_2(W, U) \\ V: v(A, B) &\leftarrow p(A, C), p_0(C, B), p_1(A, D). \end{aligned} \quad (19.52)$$

Q can be rewritten using V :

$$Q': q(X, U) \leftarrow v(X, Z), p_1(X, W), p_2(W, U). \quad (19.53)$$

View V replaces the first two literals of query Q . Note that the view certainly satisfies the third literal of the query, as well. However, it cannot be removed from the rewriting, since variable D does not occur in the head of V , thus if literal p_1 were to be removed, too, then the natural join of p_1 and p_2 would not be enforced anymore.

Since in some of the applications the database relations are inaccessible, only the views can be accessed, for example in case of data integration or data warehouses, the concept of **complete rewriting** is introduced.

Definition 19.23 *A rewriting Q' of query Q using views $\mathcal{V} = V_1, V_2, \dots, V_m$ is called a **complete rewriting**, if Q' contains only literals of \mathcal{V} and comparison atoms.*

Example 19.9 *Complete rewriting.* Assume that besides view V of Example 19.8 the following view is given, as well:

$$V_2: v_2(A, B) \leftarrow p_1(A, C), p_2(C, B), p_0(D, E) \quad (19.54)$$

A complete rewriting of query Q is:

$$Q'': q(X, U) \leftarrow v(X, Z), v_2(X, U). \quad (19.55)$$

It is important to see that this rewriting cannot be obtained *step-by-step*, first using only V , then trying to incorporate V_2 , (or just in the opposite order) since relation p_0 of V_2 does not occur in Q' . Thus, in order to find the complete rewriting, use of the two view must be considered *parallel*, at the same time.

There is a close connection between finding a rewriting and the problem of query containment. This latter one was discussed for tableau queries in section 19.1.3. Homomorphism between tableau queries can be defined for rule based queries, as well. The only difference is that it is not required in this section that a homomorphism maps the head of one rule to the head of the other. (The head of a rule corresponds to the summary row of a tableau.) According to Theorem 19.20 it is NP-complete to decide whether conjunctive query Q_1 contains another conjunctive query Q_2 . This remains true in the case when Q_2 may contain comparison atoms, as well. However, if both, Q_1 and Q_2 may contain comparison atoms, then the existence of a homomorphism from Q_1 to Q_2 is only a sufficient but not necessary condition for the containment of queries, which is a Π_2^p -complete problem in that case. The discussion of this latter complexity class is beyond the scope of this chapter, thus it is omitted. The next proposition gives a necessary and sufficient condition whether there exists a rewriting of query Q using view V .

Claim 19.24 *Let Q and V be conjunctive queries that may contain comparison atoms. There exists a rewriting of query Q using view V if and only if $\pi_\emptyset(Q) \subseteq \pi_\emptyset(V)$, that is the projection of V to the empty attribute set contains that of Q .*

Proof Observe that $\pi_\emptyset(Q) \subseteq \pi_\emptyset(V)$ is equivalent with the following proposition: If the output of V is empty for some instance, then the same holds for the output of Q , as well.

Assume first that there exists a rewriting, that is a rule equivalent with Q that contains V in its body. If r is such an instance, that the result of V is empty on it, then every rule that includes V in its body results in empty set over r , too.

In order to prove the other direction, assume that if the output of V is empty for some instance, then the same holds for the output of Q , as well. Let

$$\begin{aligned} Q: q(\tilde{x}) &\leftarrow q_1(\tilde{x}), q_2(\tilde{x}), \dots, q_m(\tilde{x}) \\ V: v(\tilde{a}) &\leftarrow v_1(\tilde{a}), v_2(\tilde{a}), \dots, v_n(\tilde{a}) . \end{aligned} \quad (19.56)$$

Let \tilde{y} be a list of variables disjoint from variables of \tilde{x} . Then the query Q' defined by

$$Q': q'(\tilde{x}) \leftarrow q_1(\tilde{x}), q_2(\tilde{x}), \dots, q_m(\tilde{x}), v_1(\tilde{y}), v_2(\tilde{y}), \dots, v_n(\tilde{y}) \quad (19.57)$$

satisfies $Q \equiv Q'$. It is clear that $Q' \subseteq Q$. On the other hand, if there exists a valuation of the variables of \tilde{y} that satisfies the body of V over some instance r , then fixing it, for arbitrary valuation of variables in \tilde{x} a tuple is obtained in the output of Q , whenever a tuple is obtained in the output of Q' together with the previously fixed valuation of variables of \tilde{y} . ■

As a corollary of Theorem 19.20 and Proposition 19.24 the following theorem is obtained.

Theorem 19.25 *Let Q be a conjunctive query that may contain comparison atoms, and let \mathcal{V} be a set of views. If the views in \mathcal{V} are given by conjunctive queries that do not contain comparison atoms, then it is NP-complete to decide whether there exists a rewriting of Q using \mathcal{V} .*

The proof of Theorem 19.25 is left for the Reader (Exercise 19.3-1).

The proof of Proposition 19.24 uses new variables. However, according to the next lemma, this is not necessary. Another important observation is that it is enough to consider a subset of database relations occurring in the original query when locally minimal rewriting is sought, new database relations need not be introduced.

Lemma 19.26 *Let Q be a conjunctive query that does not contain comparison atoms*

$$Q: q(\tilde{X}) \leftarrow p_1(\tilde{U}_1), p_2(\tilde{U}_2), \dots, p_n(\tilde{U}_n), \quad (19.58)$$

furthermore let \mathcal{V} be a set of views that do not contain comparison atoms either.

1. If Q' is a locally minimal rewriting of Q using \mathcal{V} , then the set of database literals in Q' is isomorphic to a subset of database literals occurring in Q .
2. If

$$q(\tilde{X}) \leftarrow p_1(\tilde{U}_1), p_2(\tilde{U}_2), \dots, p_n(\tilde{U}_n), v_1(\tilde{Y}_1), v_2(\tilde{Y}_2), \dots, v_k(\tilde{Y}_k) \quad (19.59)$$

is a rewriting of Q using the views, then there exists a rewriting

$$q(\tilde{X}) \leftarrow p_1(\tilde{U}_1), p_2(\tilde{U}_2), \dots, p_n(\tilde{U}_n), v_1(\tilde{Y}'_1), v_2(\tilde{Y}'_2), \dots, v_k(\tilde{Y}'_k) \quad (19.60)$$

such that $\{\tilde{Y}'_1 \cup \dots \cup \tilde{Y}'_k\} \subseteq \{\tilde{U}_1 \cup \dots \cup \tilde{U}_n\}$, that is the rewriting does not introduce new variables.

The details of the proof of Lemma 19.26 are left for the Reader (Exercise 19.3-2). The next lemma is of fundamental importance: A minimal rewriting of Q using \mathcal{V} cannot **increase** the number of literals.

Lemma 19.27 *Let Q be conjunctive query, \mathcal{V} be set of views given by conjunctive queries, both without comparison atoms. If the body of Q contains p literals and Q' is a locally minimal rewriting of Q using \mathcal{V} , then Q' contains at most p literals.*

Proof Replacing the view literals of Q' by their definitions query Q'' is obtained. Let φ be a homomorphism from the body of Q to Q'' . The existence of φ follows from $Q \equiv Q''$ by the Homomorphism Theorem (Theorem 19.18). Each of the literals l_1, l_2, \dots, l_p of the body of Q is mapped to at most one literal obtained from the expansion of view definitions. If Q' contains more than p view literals, then the expansion of some view literals in the body of Q'' is disjoint from the image of φ . These view literals can be removed from the body of Q' without changing the equivalence. ■

Based on Lemma 19.27 the following theorem can be stated about complexity of minimal rewritings.

Theorem 19.28 *Let Q be conjunctive query, \mathcal{V} be set of views given by conjunctive queries, both without comparison atoms. Let the body of Q contain p literals.*

1. *It is NP-complete to decide whether there exists a rewriting Q' of Q using \mathcal{V} that uses at most k ($\leq p$) literals.*
2. *It is NP-complete to decide whether there exists a rewriting Q' of Q using \mathcal{V} that uses at most k ($\leq p$) database literals.*
3. *It is NP-complete to decide whether there exists a complete rewriting of Q using \mathcal{V} .*

Proof The first statement is proved, the proof of the other two is similar. According to Lemmas 19.27 and 19.26, only such rewritings need to be considered that have at most as many literals as the query itself, contain a subset of the literals of the query and do not introduce new variables. Such a rewriting and the homomorphisms proving the equivalence can be tested in polynomial time, hence the problem is in NP. In order to prove that it is NP-hard, Theorem 19.25 is used. For a given query Q and view V let V' be the view, whose head is same as the head of V , but whose body is the conjunction of the bodies of Q and V . It is easy to see that there exists a rewriting using V' with a single literal if and only if there exists a rewriting (with no restriction) using V . ■

19.3.3. Practical algorithms

In this section only complete rewritings are studied. This does not mean real restriction, since if database relations are also to be used, then views mirroring the database relations one-to-one can be introduced. The concept of *equivalent* rewriting introduced in Definition 19.22 is appropriate if the goal of the rewriting is query optimisation or providing physical data independence. However, in the context of data integration on data warehouses equivalent rewritings cannot be sought, since all necessary data may not be available. Thus, the concept of maximally contained rewriting is introduced that depends on the query language used, in contrast to equivalent rewritings.

Definition 19.29 *Let Q be a query, \mathcal{V} be a set of views, \mathcal{L} be a query language. Q' is a **maximally contained rewriting** of Q with respect to \mathcal{L} , if*

1. *Q' is a query of language \mathcal{L} using only views from \mathcal{V} ,*
2. *Q contains Q' ,*
3. *if query $Q_1 \in \mathcal{L}$ satisfies $Q' \sqsubseteq Q_1 \sqsubseteq Q$, then $Q' \equiv Q_1$.*

Query optimisation using materialised views. Before discussing how can a traditional optimiser be modified in order to use materialised views instead of database relations, it is necessary to survey when can view be used to answer a given query. Essentially, view V can be used to answer query Q , if the intersection of the sets of database relations in the body of V and in the body of Q is non-empty, furthermore some of the attributes are selected by V are also selected by Q . Besides this, in case of equivalent rewriting, if V contains comparison atoms for such attributes that are also occurring in Q , then the view must apply logically equivalent, or weaker condition, than the query. If logically stronger condition is applied in the view, then it can be part of a (maximally) contained rewriting. This can be shown best via an example. Consider the query Q over schema **University** that list those professor, student, semester triplets, where the advisor of the student is the professor and in the given semester the student registered for some course taught by the professor.

$$\begin{aligned}
 Q: q(Pname, Student, Semester) \leftarrow & Registered(Student, C-number, Semester), \\
 & Advisor(Pname, Student), \\
 & Teaches(Pname, C-number, Semester, x_E), \\
 & Semester \geq \text{“Fall2000”} .
 \end{aligned}
 \tag{19.61}$$

View V_1 below can be used to answer Q , since it uses the same join condition for relations *Registered* and *Teaches* as Q , as it is shown by variables of the same name. Furthermore, V_1 selects attributes *Student*, *PName*, *Semester*, that are necessary in order to properly join with relation *Advisor*, and for select clause of the query. Finally, the predicate $Semester > \text{“Fall1999”}$ is weaker than the predicate $Semester \geq \text{“Fall2000”}$ of the query.

$$\begin{aligned}
 V_1: v_1(Student, PName, Semester) \leftarrow & Teaches(PName, C-number, Semester, x_E), \\
 & Registered(Student, C-number, Semester), \\
 & Semester > \text{“Fall1999”} .
 \end{aligned}
 \tag{19.62}$$

The following four views illustrate how minor modifications to V_1 change the usability in answering the query.

$$\begin{aligned}
 V_2: v_2(Student, Semester) \leftarrow & Teaches(x_P, C-number, Semester, x_E), \\
 & Registered(Student, C-number, Semester), \\
 & Semester > \text{“Fall1999”} .
 \end{aligned}
 \tag{19.63}$$

$$\begin{aligned}
 V_3: v_3(Student, PName, Semester) \leftarrow & Teaches(PName, C-number, x_S, x_E), \\
 & Registered(Student, C-number, Semester), \\
 & Semester > \text{“Fall1999”} .
 \end{aligned}
 \tag{19.64}$$

$$\begin{aligned}
 V_4: v_4(Student, PName, Semester) \leftarrow & Teaches(PName, C-number, Semester, x_E), \\
 & Adviser(PName, x_{St}), Professor(PName, x_A), \\
 & Registered(Student, C-number, Semester), \\
 & Semester > \text{“Fall1999”} .
 \end{aligned}
 \tag{19.65}$$

$$V_5: v_5(Student, PName, Semester) \leftarrow Teaches(PName, C-number, Semester, x_E), \\ Registered(Student, C-number, Semester), \\ Semester > \text{“Fall2001”} . \quad (19.66)$$

View V_2 is similar to V_1 , except that it does not select the attribute $PName$ from relation $Teaches$, which is needed for the join with the relation $Adviser$ and for the selection of the query. Hence, to use V_2 in the rewriting, it has to be joined with relation $Teaches$ again. Still, if the join of relations $Registered$ and $Teaches$ is very selective, then employing V_2 may actually result in a more efficient query execution plan.

In view V_3 the join of relations $Registered$ and $Teaches$ is over only attribute $C-number$, the equality of variables $Semester$ and x_S is not required. Since attribute x_S is not selected by V_3 , the join predicate cannot be applied in the rewriting, and therefore there is little gain by using V_3 .

View V_4 considers only the professors who have at least one area of research. Hence, the view applies an additional condition that does not exist in the query, and cannot be used in an equivalent rewriting unless union and negation are allowed in the rewriting language. However, if there is an integrity constraint stating that every professor has at least one area of research, then an optimiser should be able to realise that V_4 is usable.

Finally, view V_5 applies a stronger predicate than the query, and is therefore usable for a contained rewriting, but not for an equivalent rewriting of the query.

System-R style optimisation Before discussing the changes to traditional optimisation, first the principles underlying the *System-R style optimiser* is recalled briefly. System-R takes a bottom-up approach to building query execution plans. In the first phase, it concentrates on plans of size 1, i.e., chooses the best access paths to every table mentioned in the query. In phase n , the algorithm considers plans of size n , by combining plans obtained in the previous phases (sizes of k and $n - k$). The algorithm terminates after constructing plans that cover all the relations in the query. The efficiency of System-R stems from the fact that it partitions query execution plans into *equivalence classes*, and only considers a single execution plan for every equivalence class. Two plans are in the same equivalence class if they

- cover the same set of relations in the query (and therefore are also of the same size), and
- produce the answer in the same interesting order.

In our context, the query optimiser builds query execution plans by accessing a set of views, rather than a set of database relations. Hence, in addition to the meta-data that the query optimiser has about the materialised views (e.g., statistics, indexes) the optimiser is also given as input the query expressions defining the views. The additional issues that the optimiser needs to consider in the presence of materialised views are as follows.

- A. In the first iteration the algorithm needs to decide which views are *relevant* to the query according to the conditions illustrated above. The corresponding

step is trivial in a traditional optimiser: a relation is relevant to the query if it is in the body of the query expression.

B. Since the query execution plans involve joins over views, rather than joins over database relations, plans can no longer be neatly partitioned into equivalence classes which can be explored in increasing size. This observation implies several changes to the traditional algorithm:

1. **Termination testing:** the algorithm needs to distinguish *partial query execution plans* of the query from *complete query execution plans*. The enumeration of the possible join orders terminates when there are no more unexplored partial plans. In contrast, in the traditional setting the algorithm terminates after considering the equivalence classes that include all the relations of the query.
2. **Pruning of plans:** a traditional optimiser compares between pairs of plans *within* one equivalence class and saves only the cheapest one for each class. In our context, the query optimiser needs to compare between *any pair* of plans generated thus far. A plan p is pruned if there is another plan p' that
 - (a) is cheaper than p , and
 - (b) has greater or equal contribution to the query than p . Informally, a plan p' contributes more to the query than plan p if it covers more of the relations in the query and selects more of the necessary attributes.
3. **Combining partial plans:** in the traditional setting, when two partial plans are combined, the join predicates that involve both plans are explicit in the query, and the enumeration algorithm need only consider the most efficient way to apply these predicates. However, in our case, it may not be obvious a priori which join predicate will yield a correct rewriting of the query, since views are joined rather than database relations directly. Hence, the enumeration algorithm needs to consider several alternative join predicates. Fortunately, in practice, the number of join predicates that need to be considered can be significantly pruned using meta-data about the schema. For example, there is no point in trying to join a string attribute with a numeric one. Furthermore, in some cases knowledge of integrity constraints and the structure of the query can be used to reduce the number of join predicates to be considered. Finally, after considering all the possible join predicates, the optimiser also needs to check whether the resulting plan is still a partial solution to the query.

The following table summarises the comparison of the traditional optimiser versus one that exploits materialised views.

Conventional optimiser

Iteration 1

- a) Find all possible access paths.
- b) Compare their cost and keep the least expensive.
- c) If the query has one relation, **stop**.

Iteration 2

For each query join:

- a) Consider joining the relevant access paths found in the previous iteration using all possible join methods.
- b) Compare the cost of the resulting join plans and keep the least expensive.
- c) If the query has two relations, **stop**.

Iteration 3

⋮

Optimiser using views

Iteration 1

- a1) Find all views that are relevant to the query.
- a2) Distinguish between partial and complete solutions to the query.
- b) Compare all pairs of views. If one has neither greater contribution nor a lower cost than the other, prune it.
- c) If there are no partial solutions, **stop**.

Iteration 2

- a1) Consider joining all partial solutions found in the previous iteration using all possible equi-join methods and trying all possible subsets of join predicates.
- a2) Distinguish between partial and complete solutions to the query.
- b) If any newly generated solution is either not relevant to the query, or dominated by another, prune it.
- c) If there are no partial solutions, **stop**.

Iteration 3

⋮

Another method of equivalent rewriting is using transformation rules. The common theme in the works of that area is that replacing some part of the query with a view is considered as another transformation available to the optimiser. These methods are not discussed in detail here.

The query optimisers discussed above were designed to handle cases where the number of views is relatively small (i.e., comparable to the size of the database schema), and cases where equivalent rewriting is required. In contrast, the context of data integration requires consideration of large number of views, since each data source is being described by one or more views. In addition, the view definitions may contain many complex predicates, whose goal is to express fine-grained distinction between the contents of different data sources. Furthermore, in the context of data integration it is often assumed that the views are not complete, i.e., they may only contain a subset of the tuples satisfying their definition. In the foregoing, some algorithms for answering queries using views are described that were developed specifically for the context of data integration.

The Bucket Algorithm. The goal of the Bucket Algorithm is to reformulate a user query that is posed on a mediated (virtual) schema into a query that refers directly to the available sources. Both the query and the sources are described by conjunctive queries that may include atoms of arithmetic comparison predicates. The set of comparison atoms of query Q is denoted by $C(Q)$.

Since the number of possible rewritings may be exponential in the size of the query, the main idea underlying the bucket algorithm is that the number of query rewritings that need to be considered can be drastically reduced if we first consider each *subgoal* – the relational atoms of the query – is considered in isolation, and determine which views may be relevant to each subgoal.

The algorithm proceeds as follows. First, a *bucket* is created for each subgoal in the query Q that is not in $C(Q)$, containing the views that are relevant to answering the particular subgoal. In the second step, all such conjunctive query rewritings are considered that include one conjunct (view) from each bucket. For each rewriting V obtained it is checked that whether it is semantically correct, that is $V \sqsubseteq Q$ holds, or whether it can be made semantically correct by adding comparison atoms. Finally the remaining plans are minimised by pruning redundant subgoals. Algorithm CREATE-BUCKET executes the first step described above. Its input is a set of source descriptions \mathcal{V} and a conjunctive query Q in the form

$$Q: Q(\tilde{X}) \leftarrow R_1(\tilde{X}_1), R_2(\tilde{X}_2), \dots, R_m(\tilde{X}_m), C(Q). \quad (19.67)$$

CREATE-BUCKET(Q, \mathcal{V})

```

1  for  $i \leftarrow 1$  to  $m$ 
2    do  $Bucket[i] \leftarrow \emptyset$ 
3    for all  $V \in \mathcal{V}$ 
4      do  $\triangleright V$  is of form  $V: V(\tilde{Y}) \leftarrow S_1(\tilde{Y}_1), \dots, S_n(\tilde{Y}_n), C(V)$ .
5      do for  $j \leftarrow 1$  to  $n$ 
6        if  $R_i = S_j$ 
7          then let  $\phi$  be a mapping defined on the variables
8             of  $V$  as follows:
9             if  $y$  is the  $k^{\text{th}}$  variable of  $\tilde{Y}_j$  and  $y \in \tilde{Y}$ 
10              then  $\phi(y) = x_k$ , where  $x_k$  is the  $k^{\text{th}}$  variable of  $\tilde{X}_i$ 
11              else  $\phi(y)$  is a new variable that
12                 does not appear in  $Q$  or  $V$ .
13               $Q'() \leftarrow R_1(\tilde{X}_1), R_m(\tilde{X}_m), C(Q), S_1(\phi(\tilde{Y}_1)), \dots,$ 
14                  $S_n(\phi(\tilde{Y}_n)), \phi(C(V))$ 
15              if SATISFIABLE $^{\geq}(Q')$ 
16                then add  $\phi(V)$  to  $Bucket[i]$ .
17  return  $Bucket$ 

```

Procedure SATISFIABLE $^{\geq}$ is the extension of SATISFIABLE described in section 19.1.2 to the case when comparison atoms may occur besides equality atoms. The necessary change is only that for all variable y occurring in comparison atoms it must be checked whether all predicates involving y are satisfiable simultaneously.

CREATE-BUCKET running time is polynomial function of the sizes of Q and \mathcal{V} . Indeed, the kernel of the nested loops of lines 3 and 5 runs $n \sum_{V \in \mathcal{V}} |V|$ times. The commands of lines 6–13 require constant time, except for line 12. The condition of command **if** in line 12 can be checked in polynomial time.

In order to prove the correctness of procedure CREATE-BUCKET, one should check under what condition is a view V put in $Bucket[i]$. In line 6 it is checked whether relation R_i appears as a subgoal in V . If not, then obviously V cannot give usable information for subgoal R_i of Q . If R_i is a subgoal of V , then in lines 9–10 a mapping is constructed that applied to the variables allows the correspondence between subgoals S_j and R_i , in accordance with relations occurring in the heads of Q and V , respectively. Finally, in line 12 it is checked whether the comparison atoms contradict with the correspondence constructed.

In the second step, having constructed the buckets using CREATE-BUCKET, the bucket algorithm finds a set of *conjunctive query rewritings*, each of them being a conjunctive query that includes one conjunct from every bucket. Each of these conjunctive query rewritings represents one way of obtaining part of the answer to Q from the views. The result of the bucket algorithm is defined to be the union of the conjunctive query rewritings (since each of the rewritings may contribute different tuples). A given conjunctive query Q' is a *conjunctive query rewriting*, if

1. $Q' \sqsubseteq Q$, or
2. Q' can be extended with comparison atoms, so that the previous property holds.

Example 19.10 *Bucket algorithm.* Consider the following query Q that lists those articles x that there exists another article y of the same area such that x and y mutually cites each other. There are three views (sources) available, V_1, V_2, V_3 .

$$\begin{array}{ll}
 Q(x) & \leftarrow \text{cite}(x, y), \text{cite}(y, x), \text{sameArea}(x, y) \\
 V_1(a) & \leftarrow \text{cite}(a, b), \text{cite}(b, a) \\
 V_2(c, d) & \leftarrow \text{sameArea}(c, d) \\
 V_3(f, h) & \leftarrow \text{cite}(f, g), \text{cite}(g, h), \text{sameArea}(f, g) .
 \end{array} \tag{19.68}$$

In the first step, applying CREATE-BUCKET, the following buckets are constructed.

$$\begin{array}{c|c|c}
 \text{cite}(x, y) & \text{cite}(y, x) & \text{sameArea}(x, y) \\
 \hline
 V_1(x) & V_1(x) & V_2(x) \\
 V_3(x) & V_3(x) & V_3(x)
 \end{array} \tag{19.69}$$

In the second step the algorithm constructs a conjunctive query Q' from each element of the Cartesian product of the buckets, and checks whether Q' is contained in Q . If yes, it is given to the answer.

In our case, it tries to match V_1 with the other views, however no correct answer is obtained so. The reason is that b does not appear in the head of V_1 , hence the join condition of Q – variables x and y occur in relation sameArea , as well – cannot be applied. Then rewritings containing V_3 are considered, recognising that equating the variables in the head of V_3 a contained rewriting is obtained. Finally, the algorithm finds that combining V_3 and V_2 rewriting is obtained, as well. This latter is redundant, as it is obtained by simple

checking, that is V_2 can be pruned. Thus, the result of the bucket algorithm for query (19.68) is the following (actually equivalent) rewriting

$$Q'(x) \leftarrow V_3(x, x). \quad (19.70)$$

The strength of the bucket algorithm is that it exploits the predicates in the query to prune significantly the number of candidate conjunctive rewritings that need to be considered. Checking whether a view should belong to a bucket can be done in time polynomial in the size of the query and view definition when the predicates involved are arithmetic comparisons. Hence, if the data sources (i.e., the views) are indeed distinguished by having different comparison predicates, then the resulting buckets will be relatively small.

The main disadvantage of the bucket algorithm is that the Cartesian product of the buckets may still be large. Furthermore, the second step of the algorithm needs to perform a query containment test for every candidate rewriting, which is NP-complete even when no comparison predicates are involved.

Inverse-rules algorithm. The Inverse-rules algorithm is a procedure that can be applied more generally than the Bucket algorithm. It finds a maximally contained rewriting for any query given by arbitrary recursive datalog program that does not contain negation, in polynomial time.

The first question is that for given datalog program \mathcal{P} and set of conjunctive queries \mathcal{V} , whether there exists a datalog program \mathcal{P}_v equivalent with \mathcal{P} , whose *edb* relations are relations v_1, v_2, \dots, v_n of \mathcal{V} . Unfortunately, this is algorithmically undecidable. Surprisingly, the best, maximally contained rewriting can be constructed. In the case, when there exists a datalog program \mathcal{P}_v equivalent with \mathcal{P} , the algorithm finds that, since a maximally contained rewriting contains \mathcal{P}_v , as well. This seemingly contradicts to the fact that the existence of equivalent rewriting is algorithmically undecidable, however it is undecidable about the result of the inverse-rules algorithm, whether it is really equivalent to the original query.

Example 19.11 *Equivalent rewriting.* Consider the following datalog program \mathcal{P} , where *edb* relations *edge* and *black* contain the edges and vertices coloured black of a graph G .

$$\begin{aligned} \mathcal{P}: \quad q(X, Y) &\leftarrow \text{edge}(X, Z), \text{edge}(Z, Y), \text{black}(Z) \\ q(X, Y) &\leftarrow \text{edge}(X, Z), \text{black}(Z), q(Z, Y). \end{aligned} \quad (19.71)$$

It is easy to check that \mathcal{P} lists the endpoints of such paths (more precisely walks) of graph G whose inner points are all black. Assume that only the following two views can be accessed.

$$\begin{aligned} v_1(X, Y) &\leftarrow \text{edge}(X, Y), \text{black}(X) \\ v_2(X, Y) &\leftarrow \text{edge}(X, Y), \text{black}(Y) \end{aligned} \quad (19.72)$$

v_1 stores edges whose tail is black, while v_2 stores those, whose head is black. There exists an equivalent rewriting \mathcal{P}_v of datalog program \mathcal{P} that uses only views v_1 and v_2 as *edb* relations:

$$\begin{aligned} \mathcal{P}_v: \quad q(X, Y) &\leftarrow v_2(X, Z), v_1(Z, Y) \\ q(X, Y) &\leftarrow v_2(X, Z), q(Z, Y) \end{aligned} \quad (19.73)$$

However, if only v_1 , or v_2 is accessible alone, then equivalent rewriting is not possible, since

only such paths are obtainable whose starting, or respectively, ending vertex is black.

In order to describe the Inverse-rules Algorithm, it is necessary to introduce the *Horn rule*, which is a generalisation of *datalog program*, and *datalog rule*. If *function symbols* are also allowed in the free tuple u_i of rule (19.27) in Definition 19.11, besides variables and constants, then *Horn rule* is obtained. A *logic program* is a collection of Horn rules. In this sense a logic program without function symbols is a datalog program. The concepts of *edb* and *idb* can be defined for logic programs in the same way as for datalog programs.

The Inverse-rules Algorithm consists of two steps. First, a logic program is constructed that may contain function symbols. However, these will not occur in recursive rules, thus in the second step they can be eliminated and the logic program can be transformed into a datalog program.

Definition 19.30 The *inverse* v^{-1} of view v given by

$$v(X_1, \dots, X_m) \leftarrow v_1(\tilde{Y}_1), \dots, v_n(\tilde{Y}_n) \quad (19.74)$$

is the following collection of Horn rules. A rule corresponds to every subgoal $v_i(\tilde{Y}_i)$, whose body is the literal $v(X_1, \dots, X_m)$. The head of the rule is $v_i(\tilde{Z}_i)$, where \tilde{Z}_i is obtained from \tilde{Y}_i by preserving variables appearing in the head of rule (19.74), while function symbol $f_Y(X_1, \dots, X_m)$ is written in place of every variable Y not appearing the head. Distinct function symbols correspond to distinct variables. The inverse of a set \mathcal{V} of views is the set $\{v^{-1} : v \in \mathcal{V}\}$, where distinct function symbols occur in the inverses of distinct rules.

The idea of the definition of inverses is that if a tuple (x_1, \dots, x_m) appears in a view v , for some constants x_1, \dots, x_m , then there is a valuation of every variable y appearing in the head that makes the body of the rule true. This “unknown” valuation is denoted by the function symbol $f_Y(X_1, \dots, X_m)$.

Example 19.12 *Inverse of views.* Let \mathcal{V} be the following collection of views.

$$\begin{aligned} v_1(X, Y) &\leftarrow \text{edge}(X, Z), \text{edge}(Z, W), \text{edge}(W, Y) \\ v_2(X) &\leftarrow \text{edge}(X, Z). \end{aligned} \quad (19.75)$$

Then \mathcal{V}^{-1} consists of the following rules.

$$\begin{aligned} \text{edge}(X, f_{1,Z}(X, Y)) &\leftarrow v_1(X, Y) \\ \text{edge}(f_{1,Z}(X, Y), f_{1,W}(X, Y)) &\leftarrow v_1(X, Y) \\ \text{edge}(f_{1,W}(X, Y), Y) &\leftarrow v_1(X, Y) \\ \text{edge}(X, f_{2,Z}(X)) &\leftarrow v_2(X). \end{aligned} \quad (19.76)$$

Now, the maximally contained rewriting of datalog program \mathcal{P} using views \mathcal{V} can easily be constructed for given \mathcal{P} and \mathcal{V} .

First, those rules are deleted from \mathcal{P} that contain such *edb* relation that do not appear in the definition any view from \mathcal{V} . The rules of \mathcal{V}^{-1} are added the datalog program \mathcal{P}^- obtained so, thus forming logic program $(\mathcal{P}^-, \mathcal{V}^{-1})$. Note, that the remaining *edb* relations of \mathcal{P} are *idb* relations in logic program $(\mathcal{P}^-, \mathcal{V}^{-1})$, since they

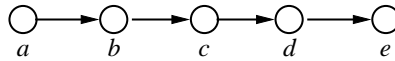


Figure 19.4 The graph G .

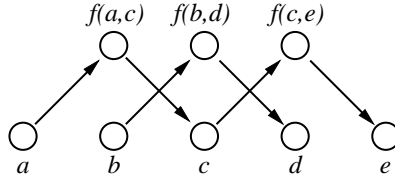


Figure 19.5 The graph G' .

appear in the heads of the rules of \mathcal{V}^{-1} . The names of *idb* relations are arbitrary, so they can be renamed so that their names do not coincide with the names of *edb* relations of \mathcal{P} . However, this is not done in the following example, for the sake of better understanding.

Example 19.13 *Logic program.* Consider the following datalog program that calculates the transitive closure of relation *edge*.

$$\begin{aligned} \mathcal{P}: \quad q(X, Y) &\leftarrow \text{edge}(X, Y) \\ q(X, Y) &\leftarrow \text{edge}(X, Z), q(Z, Y) \end{aligned} \tag{19.77}$$

Assume that only the following materialised view is accessible, that stores the endpoints of paths of length two. If only this view is usable, then the most that can be expected is listing the endpoints of paths of even length. Since the unique *edb* relation of datalog program \mathcal{P} is *edge*, that also appears in the definition of v , the logic program $(\mathcal{P}^-, \mathcal{V}^{-1})$ is obtained by adding the rules of \mathcal{V}^{-1} to \mathcal{P} .

$$\begin{aligned} (\mathcal{P}^-, \mathcal{V}^{-1}): \quad q(X, Y) &\leftarrow \text{edge}(X, Y) \\ q(X, Y) &\leftarrow \text{edge}(X, Z), q(Z, Y) \\ \text{edge}(X, f(X, Y)) &\leftarrow v(X, Y) \\ \text{edge}(f(X, Y), Y) &\leftarrow v(X, Y) . \end{aligned} \tag{19.78}$$

Let the instance of the *edb* relation *edge* of datalog program \mathcal{P} be the graph G shown on Figure 19.4. Then $(\mathcal{P}^-, \mathcal{V}^{-1})$ introduces three new constants, $f(a, c)$, $f(b, d)$ és $f(c, e)$. The *idb* relation *edge* of logic program \mathcal{V}^{-1} is graph G' shown on Figure 19.5. \mathcal{P}^- computes the transitive closure of graph G' . Note that those pairs in the transitive closure that do not contain any of the new constants are exactly the endpoints of even paths of G .

The result of logic program $(\mathcal{P}^-, \mathcal{V}^{-1})$ in Example 19.13 can be calculated by procedure NAÏV-DATALOG, for example. However, it is not true for logic programs in general, that the algorithm terminates. Indeed, consider the logic program

$$\begin{aligned} q(X) &\leftarrow p(X) \\ q(f(X)) &\leftarrow q(X) . \end{aligned} \tag{19.79}$$

If the *edb* relation p contains the constant a , then the output of the program is the infinite sequence $a, f(a), f(f(a)), f(f(f(a))), \dots$. In contrary to this, the output of the logic program $(\mathcal{P}^-, \mathcal{V}^{-1})$ given by the Inverse-rules Algorithm is guaranteed to be finite, thus the computation terminates in finite time.

Theorem 19.31 *For arbitrary datalog program \mathcal{P} and set of conjunctive views \mathcal{V} , and for finite instances of the views, there exist a unique minimal fixpoint of the logic program $(\mathcal{P}^-, \mathcal{V}^{-1})$, furthermore procedures NAIV-DATALOG and SEMI-NAIV-DATALOG give this minimal fixpoint as output.*

The essence of the proof of Theorem 19.31 is that function symbols are only introduced by inverse rules, that are in turn not recursive, thus terms containing nested functions symbols are not produced. The details of the proof are left for the Reader (Exercise 19.3-3).

Even if started from the *edb* relations of a database, the output of a logic program may contain tuples that have function symbols. Thus, a filter is introduced that eliminates the unnecessary tuples. Let database D be the instance of the *edb*relations of datalog program \mathcal{P} . $\mathcal{P}(D)\downarrow$ denotes the set of those tuples from $\mathcal{P}(D)$ that do not contain function symbols. Let $\mathcal{P}\downarrow$ denote that program, which computes $\mathcal{P}(D)\downarrow$ for a given instance D . The proof of the following theorem, exceeds the limitations of the present chapter.

Theorem 19.32 *For arbitrary datalog program \mathcal{P} and set of conjunctive views \mathcal{V} , the logic program $(\mathcal{P}^-, \mathcal{V}^{-1})\downarrow$ is a maximally contained rewriting of \mathcal{P} using \mathcal{V} . Furthermore, $(\mathcal{P}^-, \mathcal{V}^{-1})$ can be constructed in polynomial time of the sizes of \mathcal{P} and \mathcal{V} .*

The meaning of Theorem 19.32 is that the simple procedure of adding the inverses of view definitions to a datalog program results in a logic program that uses the views as much as possible. It is easy to see that $(\mathcal{P}^-, \mathcal{V}^{-1})$ can be constructed in polynomial time of the sizes of \mathcal{P} and \mathcal{V} , since for every subgoal $v_i \in \mathcal{V}$ a unique inverse rule must be constructed.

In order to completely solve the rewriting problem however, a datalog program needs to be produced that is equivalent with the logic program $(\mathcal{P}^-, \mathcal{V}^{-1})\downarrow$. The key to this is the observation that $(\mathcal{P}^-, \mathcal{V}^{-1})\downarrow$ contains only finitely many function symbols, furthermore during a bottom-up evaluation like NAIV-DATALOG and its versions, nested function symbols are not produced. With proper book keeping the appearance of function symbols can be kept track, without actually producing those tuples that contain them.

The transformation is done bottom-up like in procedure NAIV-DATALOG. The function symbol $f(X_1, \dots, X_k)$ appearing in the *idb* relation of \mathcal{V}^{-1} is replaced by the list of variables X_1, \dots, X_k . At same time the name of the *idb* relation needs to be marked, to remember that the list X_1, \dots, X_k belongs to function symbol $f(X_1, \dots, X_k)$. Thus, new “temporary” relation names are introduced. Consider the the rule

$$\text{edge}(X, f(X, Y)) \leftarrow v(X, Y) \quad (19.80)$$

of the logic program (19.78) in Example 19.13. It is replaced by rule

$$\text{edge}^{(1, f(2,3))}(X, X, Y) \leftarrow v(X, Y) \quad (19.81)$$

Notation $\langle 1, f(2, 3) \rangle$ means that the first argument of $edge^{\langle 1, f(2, 3) \rangle}$ is the same as the first argument of $edge$, while the second and third arguments of $edge^{\langle 1, f(2, 3) \rangle}$ together with function symbol f give the second argument of $edge$. If a function symbol would become an argument of an *idb* relation of \mathcal{P}^- during the bottom-up evaluation of $(\mathcal{P}^-, \mathcal{V}^{-1})$, then a new rule is added to the program with appropriately marked relation names.

Example 19.14 *Transformation of logic program into datalog program.* The logic program Example 19.13 is transformed to the following datalog program by the procedure sketched above. The different phases of the bottom-up execution of NAIV-DATALOG are separated by lines.

$$\begin{array}{ll}
 edge^{\langle 1, f(2, 3) \rangle}(X, X, Y) & \leftarrow v(X, Y) \\
 edge^{\langle f(1, 2), 3 \rangle}(X, Y, Y) & \leftarrow v(X, Y) \\
 \hline
 q^{\langle 1, f(2, 3) \rangle}(X, Y_1, Y_2) & \leftarrow edge^{\langle 1, f(2, 3) \rangle}(X, Y_1, Y_2) \\
 q^{\langle f(1, 2), 3 \rangle}(X_1, X_2, Y) & \leftarrow edge^{\langle f(1, 2), 3 \rangle}(X_1, X_2, Y) \\
 \hline
 q(X, Y) & \leftarrow edge^{\langle 1, f(2, 3) \rangle}(X, Z_1, Z_2), q^{\langle f(1, 2), 3 \rangle}(Z_1, Z_2, Y) \\
 q^{\langle f(1, 2), f(3, 4) \rangle}(X_1, X_2, Y_1, Y_2) & \leftarrow edge^{\langle f(1, 2), 3 \rangle}(X_1, X_2, Z), q^{\langle 1, f(2, 3) \rangle}(Z, Y_1, Y_2) \\
 \hline
 q^{\langle f(1, 2), 3 \rangle}(X_1, X_2, Y) & \leftarrow edge^{\langle f(1, 2), 3 \rangle}(X_1, X_2, Z), q(Z, Y) \\
 q^{\langle 1, f(2, 3) \rangle}(X, Y_1, Y_2) & \leftarrow edge^{\langle 1, f(2, 3) \rangle}(X, Z_1, Z_2), q^{\langle f(1, 2), f(3, 4) \rangle}(Z_1, Z_2, Y_1, Y_2)
 \end{array} \tag{19.82}$$

The datalog program obtained shows clearly that which arguments could involve function symbols in the original logic program. However, some rows containing function symbols never give tuples not containing function symbols during the evaluation of the output of the program.

Relation p is called *significant*, if in the precedence graph of Definition 19.16³ there exists oriented path from p to the output relation of q . If p is not significant, then the tuples of p are not needed to compute the output of the program, thus p can be eliminated from the program.

Example 19.15 *Eliminating non-significant relations.* There exists no directed path in the precedence graph of the datalog program obtained in Example 19.14, from relations $q^{\langle 1, f(2, 3) \rangle}$ and $q^{\langle f(1, 2), f(3, 4) \rangle}$ to the output relation q of the program, thus they are not significant, i.e., they can be eliminated together with the rules that involve them. The following datalog program is obtained:

$$\begin{array}{ll}
 edge^{\langle 1, f(2, 3) \rangle}(X, X, Y) & \leftarrow v(X, Y) \\
 edge^{\langle f(1, 2), 3 \rangle}(X, Y, Y) & \leftarrow v(X, Y) \\
 q^{\langle f(1, 2), 3 \rangle}(X_1, X_2, Y) & \leftarrow edge^{\langle f(1, 2), 3 \rangle}(X_1, X_2, Y) \\
 q^{\langle f(1, 2), 3 \rangle}(X_1, X_2, Y) & \leftarrow edge^{\langle f(1, 2), 3 \rangle}(X_1, X_2, Z), q(Z, Y) \\
 q(X, Y) & \leftarrow edge^{\langle 1, f(2, 3) \rangle}(X, Z_1, Z_2), q^{\langle f(1, 2), 3 \rangle}(Z_1, Z_2, Y) .
 \end{array} \tag{19.83}$$

One more simplification step can be performed, which does not decrease the number of necessary derivations during computation of the output, however avoids redundant data copying. If p is such a relation in the datalog program that is defined

³ Here the definition of precedence graph needs to be extended for the *edb* relations of the datalog program, as well.

by a single rule, which in turn contains a single relation in its body, then p can be removed from the program and replaced by the relation of the body of the rule defining p , having equated the variables accordingly.

Example 19.16 *Avoiding unnecessary data copying.* In Example 19.14 the relations $edge^{(1,f(2,3))}$ and $edge^{(f(1,2),3)}$ are defined by a single rule, respectively, furthermore these two rules both have a single relation in their bodies. Hence, program (19.83) can be simplified further.

$$\begin{aligned} q^{(f(1,2),3)}(X, Y, Y) &\leftarrow v(X, Y) \\ q^{(f(1,2),3)}(X, Z, Y) &\leftarrow v(X, Z), q(Z, Y) \\ q(X, Y) &\leftarrow v(X, Z), q^{(f(1,2),3)}(X, Z, Y). \end{aligned} \quad (19.84)$$

The datalog program obtained in the two simplification steps above is denoted by $(\mathcal{P}^-, \mathcal{V}^{-1})^{datalog}$. It is clear that there exists a one-to-one correspondence between the bottom-up evaluations of $(\mathcal{P}^-, \mathcal{V}^{-1})$ and $(\mathcal{P}^-, \mathcal{V}^{-1})^{datalog}$. Since the function symbols in $(\mathcal{P}^-, \mathcal{V}^{-1})^{datalog}$ are kept track, it is sure that the output instance obtained is in fact the subset of tuples of the output of $(\mathcal{P}^-, \mathcal{V}^{-1})$ that do not contain function symbols.

Theorem 19.33 *For arbitrary datalog program \mathcal{P} that does not contain negations, and set of conjunctive views \mathcal{V} , the logic program $(\mathcal{P}^-, \mathcal{V}^{-1})\downarrow$ is equivalent with the datalog program $(\mathcal{P}^-, \mathcal{V}^{-1})^{datalog}$.*

MiniCon. The main disadvantage of the Bucket Algorithm is that it considers each of the subgoals in isolation, therefore does not observe the most of the interactions between the subgoals of the views. Thus, the buckets may contain many unusable views, and the second phase of the algorithm may become very expensive.

The advantage of the Inverse-rules Algorithm is its conceptual simplicity and modularity. The inverses of the views must be computed only once, then they can be applied to arbitrary queries given by datalog programs. On the other hand, much of the computational advantage of exploiting the materialised views can be lost. Using the resulting rewriting produced by the algorithm for actually evaluating queries from the views has significant drawback, since it insists on recomputing the extensions of the database relations.

The MINICON algorithm addresses the limitations of the previous two algorithms. The key idea underlying the algorithm is a change of perspective: instead of building rewritings for each of the query *subgoals*, it is considered how each of the *variables* in the query can interact with the available views. The result is that the second phase of MINICON needs to consider drastically fewer combinations of views. In the following we return to conjunctive queries, and for the sake of easier understanding only such views are considered that do not contain constants.

The MINICON algorithm starts out like the Bucket Algorithm, considering which views contain subgoals that correspond to subgoals in the query. However, once the algorithm finds a partial variable mapping from a subgoal g in the query to a subgoal g_1 in a view V , it changes perspective and looks at the variables in the query. The algorithm considers the join predicates in the query – which are specified by multiple

occurrences of the same variable – and finds the minimal additional set of subgoals that *must* be mapped to subgoals in V , given that g will be mapped to g_1 . This set of subgoals and mapping information is called a *MiniCon Description (MCD)*. In the second phase the MCDs are combined to produce query rewritings. The construction of the MCDs makes the most expensive part of the Bucket Algorithm obsolete, that is the checking of containment between the rewritings and the query, because the generating rule of MCDs makes it sure that their join gives correct result.

For a given mapping $\tau: Var(Q) \rightarrow Var(V)$ subgoal g_1 of view V is said to *cover* a subgoal g of query Q , if $\tau(g) = g_1$. $Var(Q)$, and respectively $Var(V)$ denotes the set of variables of the query, respectively of that of the view. In order to prove that a rewriting gives only tuples that belong to the output of the query, a homomorphism must be exhibited from the query onto the rewriting. An MCD can be considered as a part of such a homomorphism, hence, these parts will be put together easily.

The rewriting of query Q is a union of conjunctive queries using the views. Some of the variables may be equated in the heads of some of the views as in the equivalent rewriting (19.70) of Example 19.10. Thus, it is useful to introduce the concept of *head homomorphism*. The mapping $h: Var(V) \rightarrow Var(V)$ is a *head homomorphism*, if it is an identity on variables that do not occur in the head of V , but it can equate variables of the head. For every variable x of the head of V , $h(x)$ also appear in the head of V , furthermore $h(x) = h(h(x))$. Now, the exact definition of MCD can be given.

Definition 19.34 *The quadruple $C = (h_C, V(\tilde{Y})_C, \varphi_C, G_C)$ is a *MiniCon Description (MCD)* for query Q over view V , where*

- h_C is a head homomorphism over V ,
- $V(\tilde{Y})_C$ is obtained from V by applying h_C , that is $\tilde{Y} = h_C(\tilde{A})$, where \tilde{A} is the set of variables appearing in the head of V ,
- φ_C is a partial mapping from $Var(Q)$ to $h_C(Var(V))$,
- G_C is a set of subgoals of Q that are covered by some subgoal of $H_C(V)$ using the mapping φ_C (note: not all such subgoals are necessarily included in G_C).

The procedure constructing MCDs is based on the following proposition.

Claim 19.35 *Let C be a MiniCon Description over view V for query Q . C can be used for a non-redundant rewriting of Q if the following conditions hold*

C1. *for every variable x that is in the head of Q and is in the domain of φ_C , as well, $\varphi_C(x)$ appears in the head of $h_C(V)$, furthermore*

C2. *if $\varphi_C(y)$ does not appear in the head of $h_C(V)$, then for all such subgoals of Q that contain y holds that*

1. every variable of g appears in the domain of φ_C and
2. $\varphi_C(g) \in h_C(V)$.

Clause **C1** is the same as in the Bucket Algorithm. Clause **C2** means that if a variable x is part of a join predicate which is not enforced by the view, then x must be in the head of the view so the join predicate can be applied by another subgoal in the rewriting. The procedure FORM-MCDS gives the usable MiniCon Descriptions for a conjunctive query Q and set of conjunctive views \mathcal{V} .

FORM-MCDS(Q, \mathcal{V})

```

1   $\mathcal{C} \leftarrow \emptyset$ 
2  for each subgoal  $g$  of  $Q$ 
3    do for  $V \in \mathcal{V}$ 
4      do for every subgoal  $v \in V$ 
5        do Let  $h$  be the least restrictive head homomorphism on  $V$ ,
           such that there exists a mapping  $\varphi$  with  $\varphi(g) = h(v)$ .
6          if  $\varphi$  and  $h$  exist
7            then Add to  $\mathcal{C}$  any new MCD  $C$ ,
           that can be constructed where:
8              (a)  $\varphi_C$  (respectively,  $h_C$ ) is
                 an extension of  $\varphi$  (respectively,  $h$ ),
9              (b)  $G_C$  is the minimal subset of subgoals of  $Q$  such that
                  $G_C, \varphi_C$  and  $h_C$  satisfy Proposition 19.35, and
10             (c) It is not possible to extend  $\varphi$  and  $h$  to  $\varphi'_C$ 
                 and  $h'_C$  such that (b) is satisfied,
                 and  $G'_C$  as defined in (b), is a subset of  $G_C$ .
11 return  $\mathcal{C}$ 

```

Consider again query (19.68) and the views of Example 19.10. Procedure FORM-MCDS considers subgoal $cite(x, y)$ of the query first. It does not create an MCD for view V_1 , because clause **C2** of Proposition 19.35 would be violated. Indeed, the condition would require that subgoal $sameArea(x, y)$ be also covered by V_1 using the mapping $\varphi(x) = a, \varphi(y) = b$, since is not in the head of V_1 .⁴ For the same reason, no MCD will be created for V_1 even when the other subgoals of the query are considered. In a sense, the MiniCon Algorithm shifts some of the work done by the combination step of the Bucket Algorithm to the phase of creating the MCDs by using FORM-MCDS. The following table shows the output of procedure FORM-MCDS.

$V(\tilde{Y})$	h	φ	G	(19.85)
$V_2(c, d)$	$c \rightarrow c, d \rightarrow d$	$x \rightarrow c, y \rightarrow d$	3	
$V_3(f, f)$	$f \rightarrow f, h \rightarrow f$	$x \rightarrow f, y \rightarrow f$	1, 2, 3	

Procedure FORM-MCDS includes in G_C only the *minimal* set of subgoals that are necessary in order to satisfy Proposition 19.35. This makes it possible that in the second phase of the MiniCon Algorithm needs only to consider combinations of MCDs that cover *pairwise disjoint subsets* of subgoals of the query.

Claim 19.36 *Given a query Q , a set of views \mathcal{V} , and the set of MCDs \mathcal{C} for Q over the views \mathcal{V} , the only combinations of MCDs that can result in non-redundant*

⁴ The case of $\varphi(x) = b, \varphi(y) = a$ is similar.

rewritings of Q are of the form C_1, \dots, C_l , where

- C3.** $G_{C_1} \cup \dots \cup G_{C_l}$ contains all the subgoals of Q , and
C4. for every $i \text{NEj } j$ $G_{C_i} \cap G_{C_j} = \emptyset$.

The fact that only such sets of MCDs need to be considered that provide partitions of the subgoals in the query reduces the search space of the algorithm drastically. In order to formulate procedure COMBINE-MCDs, another notation needs to be introduced. The φ_C mapping of MCD C may map a set of variables of Q onto the same variable of $h_C(V)$. One arbitrarily chosen representative of this set is chosen, with the only restriction that if there exists variables in this set from the head of Q , then one of those is the chosen one. Let $EC_{\varphi_C}(x)$ denote the representative variable of the set containing x . The MiniCon Description C is considered extended with $EC_{\varphi_C}(x)$ in the following as a quintet $(h_C, V(\tilde{Y}), \varphi_C, G_C, EC_{\varphi_C})$. If the MCDs C_1, \dots, C_k are to be combined, and for some $i \text{NEj } j$ $EC_{\varphi_{C_i}}(x) = EC_{\varphi_{C_i}}(y)$ and $EC_{\varphi_{C_j}}(y) = EC_{\varphi_{C_j}}(z)$ holds, then in the conjunctive rewriting obtained by the join x, y and z will be mapped to the same variable. Let S_C denote the equivalence relation determined on the variables of Q by two variables being equivalent if they are mapped onto the same variable by φ_C , that is, $xS_Cy \iff EC_{\varphi_C}(x) = EC_{\varphi_C}(y)$. Let \mathcal{C} be the set of MCDs obtained as the output of FORM-MCDs.

COMBINE-MCDs(\mathcal{C})

- 1 *Answer* $\leftarrow \emptyset$
- 2 **for** $\{C_1, \dots, C_n\} \subseteq \mathcal{C}$ such that G_{C_1}, \dots, G_{C_n} is a partition of the subgoals of Q
- 3 **do** Define a mapping Ψ_i on \tilde{Y}_i as follows:
 - 4 **if** there exists a variable x in Q such that $\varphi_i(x) = y$
 - 5 **then** $\Psi_i(y) = x$
 - 6 **else** $\Psi_i(y)$ is a fresh copy of y
- 7 Let S be the transitive closure of $S_{C_1} \cup \dots \cup S_{C_n}$
- 8 $\triangleright S$ is an equivalence relation of variables of Q .
- 9 Choose a representative for each equivalence class of S .
- 10 Define mapping EC as follows:
 - 11 **if** $x \in \text{Var}(Q)$
 - 12 **then** $EC(x)$ is the representative of the equivalence class of x under S
 - 13 **else** $EC(x) = x$
- 14 Let Q' be given as $Q'(EC(\tilde{X})) \leftarrow V_{C_1}(EC(\Psi_1(\tilde{Y}_1))), \dots, V_{C_n}(EC(\Psi_n(\tilde{Y}_n)))$
- 15 *Answer* $\leftarrow \text{Answer} \cup \{Q'\}$
- 16 **return** *Answer*

The following theorem summarises the properties of the MiniCon Algorithm.

Theorem 19.37 *Given a conjunctive query Q and conjunctive views \mathcal{V} , both without comparison predicates and constants, the MiniCon Algorithm produces the union of conjunctive queries that is a maximally contained rewriting of Q using \mathcal{V} .*

The complete proof of Theorem 19.37 exceeds the limitations of the present chapter. However, in Problem 19-1 the reader is asked to prove that union of the conjunctive

queries obtained as output of COMBINE-MCDS is contained in Q .

It must be noted that the running times of the Bucket Algorithm, the Inverse-rules Algorithm and the MiniCon Algorithm are the same in the worst case: $O(nmM^n)$, where n is the number of subgoals in the query, m is the maximal number of subgoals in a view, and M is the number of views. However, practical test runs show that in case of large number of views (3–400 views) the MiniCon Algorithm is significantly faster than the other two.

Exercises

19.3-1 Prove Theorem 19.25 using Proposition 19.24 and Theorem 19.20.

19.3-2 Prove the two statements of Lemma 19.26. *Hint.* For the first statement, write in their definitions in place of views $v_i(\tilde{Y}_i)$ into Q' . Minimise the obtained query Q'' using Theorem 19.19. For the second statement use Proposition 19.24 to prove that there exists a homomorphism h_i from the body of the conjunctive query defining view $v_i(\tilde{Y}_i)$ to the body of Q . Show that $\tilde{Y}'_i = h_i(\tilde{Y}_i)$ is a good choice.

19.3-3 Prove Theorem 19.31 using that datalog programs have unique minimal fixpoint.

Problems

19-1 MiniCon is correct

Prove that the output of the MiniCon Algorithm is correct. *Hint.* It is enough to show that for each conjunctive query Q' given in line 14 of COMBINE-MCDS $Q' \sqsubseteq Q$ holds. For the latter, construct a homomorphism from Q to Q' .

19-2 $(\mathcal{P}^-, \mathcal{V}^{-1})\downarrow$ is correct

Prove that each tuple produced by logic program $(\mathcal{P}^-, \mathcal{V}^{-1})\downarrow$ is contained in the output of \mathcal{P} (part of the proof of Theorem 19.32). *Hint.* Let t be a tuple in the output of $(\mathcal{P}^-, \mathcal{V}^{-1})\downarrow$ that does not contain function symbols. Consider the derivation tree of t . Its leaves are literals, since they are extensional relations of program $(\mathcal{P}^-, \mathcal{V}^{-1})$. If these leaves are removed from the tree, then the leaves of the remaining tree are *edb* relations of \mathcal{P} . Prove that the tree obtained is the derivation tree of t in datalog program \mathcal{P} .

19-3 Datalog views

This problem tries to justify why only conjunctive views were considered. Let \mathcal{V} be a set of views, Q be a query. For a given instance \mathcal{I} of the views the tuple t is a **certain answer** of query Q , if for any database instance \mathcal{D} such that $\mathcal{I} \subseteq \mathcal{V}(\mathcal{D})$, $t \in Q(\mathcal{D})$ holds, as well.

- a. Prove that if the views of \mathcal{V} are given by datalog programs, query Q is conjunctive and may contain non-equality (NE) predicates, then the question whether for a given instance \mathcal{I} of the views tuple t is a certain answer of Q is algorithmically undecidable. *Hint.* Reduce to this question the **Post Correspondence**

Problem, which is the following: Given two sets of words $\{w_1, w_2, \dots, w_n\}$ and $\{w'_1, w'_2, \dots, w'_n\}$ over the alphabet $\{a, b\}$. The question is whether there exists a sequence of indices i_1, i_2, \dots, i_k (repetition allowed) such that

$$w_{i_1}w_{i_2} \cdots w_{i_k} = w'_{i_1}w'_{i_2} \cdots w'_{i_k} . \tag{19.86}$$

The Post Correspondence Problem is well known algorithmically undecidable problem. Let the view V be given by the following datalog program:

$$\begin{aligned} V(0,0) &\leftarrow S(e, e, e) \\ V(X, Y) &\leftarrow V(X_0, Y_0), S(X_0, X_1, \alpha_1), \dots, S(X_{g-1}, Y, \alpha_g), \\ &\quad S(Y_0, Y_1, \beta_1), \dots, S(Y_{h-1}, Y, \beta_h) \\ &\quad \text{where } w_i = \alpha_1 \dots \alpha_g \text{ and } w'_i = \beta_1 \dots \beta_h \\ &\quad \text{is a rule for all } i \in \{1, 2, \dots, n\} \\ S(X, Y, Z) &\leftarrow P(X, X, Y), P(X, Y, Z) . \end{aligned} \tag{19.87}$$

Furthermore, let Q be the following conjunctive query.

$$Q(c) \leftarrow P(X, Y, Z), P(X, Y, Z'), ZNEZ' . \tag{19.88}$$

Show that for the instance \mathcal{I} of V that is given by $\mathcal{I}(V) = \{\langle e, e \rangle\}$ and $\mathcal{I}(S) = \{\}$, the tuple $\langle c \rangle$ is a certain answer of query Q if and only if the Post Correspondence Problem with sets $\{w_1, w_2, \dots, w_n\}$ and $\{w'_1, w'_2, \dots, w'_n\}$ has *no* solution.

- b. In contrast to the undecidability result of a., if \mathcal{V} is a set of conjunctive views and query Q is given by datalog program \mathcal{P} , then it is easy to decide about an arbitrary tuple t whether it is a certain answer of Q for a given view instance \mathcal{I} . Prove that the datalog program $(\mathcal{P}^-, \mathcal{V}^{-1})^{datalog}$ gives exactly the tuples of the certain answer of Q as output.

Chapter Notes

There are several dimensions along which the treatments of the problem “answering queries using views” can be classified. Figure 19.6 shows the taxonomy of the work.

The most significant distinction between the different works is whether their goal is data integration or whether it is query optimisation and maintenance of physical data independence. The key difference between these two classes of works is the output of the the algorithm for answering queries using views. In the former case, given a query Q and a set of views \mathcal{V} , the goal of the algorithm is to produce an expression Q' that references the views and is either equivalent to or contained in Q . In the latter case, the algorithm must go further and produce a (hopefully optimal) query execution plan for answering Q using the views (and possibly the database relations). Here the rewriting must be an equivalent to Q in order to ensure the correctness of the plan.

The similarity between these two bodies of work is that they are concerned with the core issue of whether a rewriting of a query is equivalent or contained in the query. However, while logical correctness suffices for the data integration context, it does

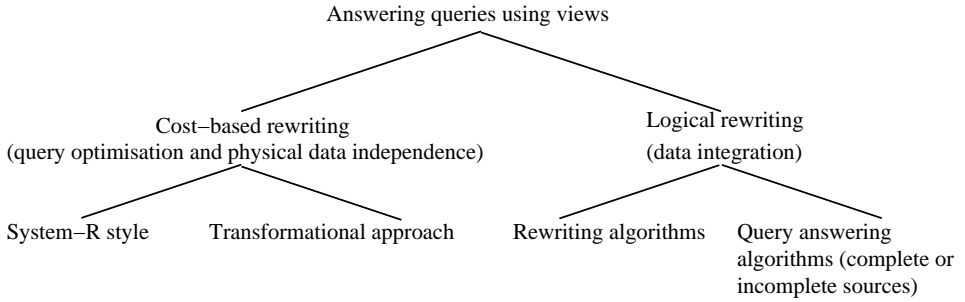


Figure 19.6 A taxonomy of work on answering queries using views.

not in the query optimisation context where we also need to find the *cheapest* plan using the views. The complication arises because the optimisation algorithms need to consider views that do not contribute to the *logical* correctness of the rewriting, but do reduce the cost of the resulting plan. Hence, while the reasoning underlying the algorithms in the data integration context is mostly logical, in the query optimisation case it is both logical and cost-based. On the other hand, an aspect stressed in data integration context is the importance of dealing with a large number of views, which correspond to data sources. In the context of query optimisation it is generally assumed (not always!) that the number of views is roughly comparable to the size of the schema.

The works on query optimisation can be classified further into System-R style optimisers and transformational optimisers. Examples of the former are works of Chaudhuri, Krishnamurty, Potomianos and Shim [?]; Tsatalos, Solomon, and Ioannidis [25]. Papers of Florescu, Raschid, and Valduriez [?]; Bello et. al. [?]; Deutsch, Popa and Tannen [?], Zaharioudakis et. al. [?], furthermore Goldstein és Larson[?] belong to the latter.

Rewriting algorithms in the data integration context are studied in works of Yang and Larson [?]; Levy, Mendelzon, Sagiv and Srivastava [?]; Qian [?]; furthermore Lambrecht, Kambhampati and Gnanaprakasam [?]. The Bucket Algorithm was introduced by Levy, Rajaraman and Ordille [?]. The Inverse-rules Algorithm is invented by Duschka and Genesereth [?, ?]. The MiniCon Algorithm was developed by Pottinger and Halevy [?, 22].

Query answering algorithms and the complexity of the problem is studied in papers of Abiteboul and Duschka [?]; Grahne and Mendelzon [?]; furthermore Calvanese, De Giacomo, Lenzerini and Vardi [?].

The STORED system was developed by Deutsch, Fernandez and Suciu [?]. Semantic caching is discussed in the paper of Yang, Karlapalem and Li [?]. Extensions of the rewriting problem are studied in [?, ?, ?, ?, ?].

Surveys of the area can be found in works of Abiteboul [?], Florescu, Levy and Mendelzon [14], Halevy [?, 17], furthermore Ullman[?].

Research of the authors was (partially) supported by Hungarian National Research Fund (OTKA) grants Nos. T034702, T037846T and T042706.

Bibliography

- [1] S. [Abiteboul](#), V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995. [894](#)
- [2] A. Aho, C. Beeri, J. D. Ullman. The theory of joins in relational databases. *ACM Transactions on Database Systems*, 4(3):297–314, 1979. [894](#)
- [3] C. Beeri. On the membership problem for functional and multivalued dependencies in relational databases. *ACM Transactions on Database Systems*, 5:241–259, 1980. [894](#)
- [4] C. Beeri, P. Bernstein. Computational problems related to the design of normal form relational schemas. *ACM Transactions on Database Systems*, 4(1):30–59, 1979. [894](#)
- [5] C. Beeri, M. Dowd. On the structure of armstrong relations for functional dependencies. *Journal of ACM*, 31(1):30–46, 1984. [894](#)
- [6] A. Békéssy, J. [Demetrovics](#). Contribution to the theory of data base relations. *Discrete Mathematics*, 27(1):1–10, 1979. [894](#)
- [7] E. F. [Codd](#). A relational model of large shared data banks. *Communications of the ACM*, 13(6):377–387, 1970. [894](#)
- [8] C. Delobel. Normalization and hierarchical dependencies in the relational data model. *ACM Transactions on Database Systems*, 3(3):201–222, 1978. [894](#)
- [9] J. [Demetrovics](#), Gy. O. H. [Katona](#), A. [Sali](#). Minimal representations of branching dependencies. *Discrete Applied Mathematics*, 40:139–153, 1992. [894](#)
- [10] J. [Demetrovics](#), Gy. O. H. [Katona](#), A. [Sali](#). Minimal representations of branching dependencies. *Acta Scientiarum Mathematicorum (Szeged)*, 60:213–223, 1995. [894](#)
- [11] J. [Demetrovics](#), Gy. O. H. [Katona](#), A. [Sali](#). Design type problems motivated by database theory. *Journal of Statistical Planning and Inference*, 72:149–164, 1998. [894](#)
- [12] R. Fagin. Multivalued dependencies and a new normal form for relational databases. *ACM Transactions on Database Systems*, 2:262–278, 1977. [894](#)
- [13] R. Fagin. Horn clauses and database dependencies. *Journal of ACM*, 29(4):952–985, 1982. [894](#)
- [14] D. [Florescu](#), A. [Halevy](#), A. O. Mendelzon. Database techniques for the world-wide web: a survey. *SIGMOD Record*, 27(3):59–74, 1998. [943](#)
- [15] J. Grant, J. Minker. Inferences for numerical dependencies. *Theoretical Computer Science*, 41:271–287, 1985. [894](#)
- [16] J. Grant, J. Minker. Normalization and axiomatization for numerical dependencies. *Information and Control*, 65:1–17, 1985. [894](#)
- [17] A. [Halevy](#). Answering queries using views: A survey. *The VLDB Journal*, 10:270–294, 2001. [943](#)
- [18] C. Lucchesi. Candidate keys for relations. *Journal of Computer and System Sciences*, 17(2):270–279, 1978. [894](#)
- [19] D. Maier. Minimum covers in the relational database model. *Journal of the ACM*, 27(4):664–674, 1980. [894](#)
- [20] D. Maier, A. O. Mendelzon, Y. Sagiv. Testing implications of data dependencies. *ACM Transactions on Database Systems*, 4(4):455–469, 1979. [894](#)

- [21] S. Petrov. Finite axiomatization of languages for representation of system properties. *Information Sciences*, 47:339–372, 1989. [894](#)
- [22] R. Pottinger. MinCon: A scalable algorithm for answering queries using views. *The VLDB Journal*, 10(2):182–198, 2001. [943](#)
- [23] A. Sali, Sr., A. [Sali](#). Generalized dependencies in relational databases. *Acta Cybernetica*, 13:431–438, 1998. [894](#)
- [24] B. Thalheim. *Dependencies in Relational Databases*. B. G. [Teubner](#), 1991. [894](#)
- [25] O. G. Tsatalos, M. C. Solomon, Y. Ioannidis. The GMAP: a versatile tool for physical data independence. *The VLDB Journal*, 5(2):101–118, 1996. [943](#)
- [26] D. M. Tsou, P. C. Fischer. Decomposition of a relation scheme into Boyce–Codd normal form. *SIGACT News*, 14(3):23–29, 1982. [894](#)
- [27] J. D. [Ullman](#). *Principles of Database and Knowledge Base Systems. Vol. 1*. Computer Science Press, 1989 (2nd edition). [894](#)
- [28] C. Zaniolo. A new normal form for the design of relational database schemata. *ACM Transactions on Database Systems*, 7:489–499, 1982. [894](#)

This bibliography is made by HBibTeX. First key of the sorting is the name of the authors (first author, second author etc.), second key is the year of publication, third key is the title of the document.

Underlying shows that the electronic version of the bibliography on the homepage of the book contains a link to the corresponding address.

Index

This index uses the following conventions. Numbers are alphabetised as if spelled out; for example, “2-3-4-tree” is indexed as if were “two-three-four-tree”. When an entry refers to a place other than the main text, the page number is followed by a tag: *exe* for exercise, *exa* for example, *fig* for figure, *pr* for problem and *fn* for footnote.

The numbers of pages containing a definition are printed in *italic* font, e.g.

time complexity, *583*.

A

anomaly, [871](#), [879](#)
deletion, [871](#)
insertion, [871](#)
redundancy, [871](#)
update, [871](#)
Armstrong-axioms, [863](#), [870](#)*exe*, [878](#)
Armstrong-relation, [892](#)
atom
relational, [898](#)
attribute, [862](#)
external, [893](#)*pr*
prime, [880](#), [893](#)*pr*
axiomatisation, [890](#)

B

bucket, [930](#)
Bucket Algorithm, [930](#), [937](#)

C

certain answer, [941](#)*pr*
CLOSURE, [865](#)
CLOSURE, [865](#), [883](#), [892](#)*exe*
of a set of attributes, [870](#)*exe*
of a set of functional dependencies, [863](#),
[864](#), [870](#)*exe*
of set of attributes, [864](#), [865](#)
COMBINE-MCDs, [940](#)
CREATE-BUCKET, [930](#)

D

database architecture
layer
logical, [914](#)
outer, [914](#)
physical, [914](#)
data independence
logical, [916](#)

data integration system, [920](#)
datalog
non-recursive, [903](#)
with negation, [904](#)
program, [907](#), [933](#)
precedence graph, [909](#)
recursive, [910](#)
rule, [906](#), [933](#)
decomposition
dependency preserving, [876](#)
dependency preserving into 3NF, [883](#)
lossless join, [872](#)
into BCNF, [880](#)
dependency
branching, [891](#)
equality generating, [885](#)
functional, [862](#)
equivalent families, [867](#)
minimal cover of a family of, [868](#)
join, [890](#)
multivalued, [884](#)
numerical, [892](#)
tuple generating, [885](#)
dependency basis, [886](#)
DEPENDENCY-BASIS, [887](#)
DEPENDENCY-BASIS, [894](#)*pr*
depth first search, [910](#)
domain, [862](#)
domain calculus, [874](#)
domain restricted, [902](#)

E

EQUATE, [873](#)
EXACT-COVER, [912](#)

F

fact, [907](#)
fixpoint, [907](#)
FORM-MCDs, [939](#)

free tuple, [897](#), [904](#)

FULL-TUPLEWISE-JOIN, [910](#)

G

Generalised Multi-level Access Path, [919](#)

GMAP, [919](#)

H

head homomorphism, [938](#)

homomorphism theorem, [910](#)

Horn rule, [933](#)

I

image under a query, [898](#)

immediate consequence, [907](#)

operator, [907](#)

IMPROVED-SEMI-NAIV-DATALOG, [910](#), [914](#)*exe*

inference rules, [863](#)

complete, [863](#)

sound, [863](#)

infinite domain, [898](#)

instance, [862](#), [895](#)

integrity constraint, [862](#), [876](#), [914](#)

Inverse-rules Algorithm, [932](#), [937](#)

J

join

natural, [899](#)

JOIN-TEST, [873](#)

JOIN-TEST, [874](#)*fig*, [883](#)

K

key, [863](#), [869](#), [880](#)

primary, [890](#)

L

LINEAR-CLOSURE, [867](#)

LINEAR-CLOSURE, [877](#), [882](#)

LINEAR-CLOSURE, [878](#), [882](#)

List-Keys, [870](#)

literal, [904](#)

negative, [904](#)

positive, [904](#)

logical implication, [863](#), [885](#)

logic program, [933](#)

lossless join, [872](#)

M

MCD, [938](#)

mediator system, [920](#)

Microsoft

Access, [898](#)

MiniCon Description, [938](#)

MiniCon Description (MCD), [938](#)

MINIMAL-COVER, [868](#)

MINIMAL-COVER, [883](#)

MINIMAL-COVER, [893](#)*pr*

N

NAIV-BCNF, [881](#)

NAIV-DATALOG, [908](#), [935](#)

natural join, [871](#), [881](#), [899](#)

normal form, [879](#)

BCNF, [893](#)*pr*

BCNF, [879](#)

Boyce-Codd, [879](#)

Boyce-Codd, [879](#)

5NF, [890](#)

4NF, [879](#), [887](#)

3NF, [879](#), [880](#), [883](#), [893](#)*pr*

nr-datalog⁺ program, [905](#)

P

POLYNOMIAL-BCNF, [882](#)

Post Correspondence Problem, [942](#)*pr*

precedence graph, [909](#), [936](#)

PRESERVE, [877](#)

projection, [899](#)

Q

QBE, [898](#)

query, [896](#)

conjunctive, [910](#)

domain restricted, [902](#)

program, [900](#)

rule based, [897](#)

subgoal, [930](#)

empty, [913](#)*exe*

equivalent, [897](#)

homomorphism, [911](#), [923](#)

language, [895](#)

equivalent, [897](#)

mapping, [896](#)

monotone, [898](#), [913](#)*exe*

relational algebra, [913](#)*exe*

rewriting, [922](#)

complete, [922](#)

conjunctive, [931](#)

equivalent, [922](#), [925](#)

globally minimal, [922](#)

locally minimal, [922](#)

maximally contained, [925](#)

minimal, [922](#)

rule based, [913](#)*exe*

satisfiable, [898](#), [913](#)*exe*

subgoal, [937](#)

tableau, [898](#), [913](#)*exe*

minimal, [911](#)

summary, [898](#)

variables of, [937](#)

query language

relationally complete, [904](#)

query rewriting, [895](#)–[943](#)

R

record, [862](#)

recursion, [906](#)

relation, [895](#)

extensional, [898](#), [907](#), [915](#)

instance, [895](#), [896](#)*fig*

intensional, [898](#), [900](#), [907](#), [915](#)

mutually recursive, [910](#)

virtual, [920](#)

relational

schema, [862](#)

- decomposition of, [871](#)
- table, [862](#)
- relational algebra*, [899](#)
- relational data bases,
- relational schema, [895](#)
- renaming, [899](#)
- rule, [897](#)
 - body, [897](#)
 - domain restricted, [904](#)
 - head, [897](#)
 - realisation, [907](#)

S

SATISFIABLE, [903](#)

schema

- extensional, [907](#)
- intensional, [907](#)
- mediated, [920](#)

selection, [899](#)

- condition, [899](#)

SEMI-NAIV-DATALOG, [935](#)

SEMI-NAIV-DATALOG, [909](#), [914](#)*exe*

- source description, [920](#)
- SQL, [915](#)
- strongly connected component, [910](#)
- subgoal, [930](#)
- substitution, [910](#)
- superkey, [863](#), [869](#), [879](#)
- System-R style optimiser, [927](#)

T

- transitive closure, [906](#)
- tuple
 - free, [897](#)

V

- view, [895](#), [914](#)
 - inverse, [933](#)
 - materialised, [916](#)

X

- XML, [919](#)

Name Index

This index uses the following conventions. If we know the full name of a cited person, then we print it. If the cited person is not living, and we know the correct data, then we print also the year of her/his birth and death.

A

Abiteboul, Serge, [874](#), [891](#), [894](#), [943](#), [944](#)
Aho, Alfred V., [888](#), [894](#), [944](#)
Armstrong, William Ward, [878](#), [892](#), [894](#)

B

Beeri, Catriel, [885](#), [886](#), [888](#), [894](#), [944](#)
Békéssy, András, [869](#), [894](#), [944](#)
Bello, Randall G., [943](#)
Bernstein, P. A., [894](#), [944](#)
Boyce, Raymond F., [879](#), [894](#)
Buneman, Peter, [943](#)

C

Calvanese, Diego, [943](#)
Chaudhuri, Surajit, [943](#)
Cochrane, Roberta, [943](#)
Codd, Edgar F. (1923–2003), [862](#), [879](#), [894](#),
[904](#), [944](#)

D

De Giacomo, Giuseppe, [943](#)
Delobel, C., [894](#), [944](#)
Demetrovics, János, [862](#), [869](#), [894](#), [895](#), [944](#)
Deutsch, Alin, [943](#)
Dias, Karl, [943](#)
Dowd, M., [894](#), [944](#)
Downing, Alan, [943](#)
Duschka, Oliver M., [943](#)

F

Fagin, R., [885](#), [894](#), [944](#)
Feenan, James J., [943](#)
Fernandez, Mary, [943](#)
Finnerty, James L., [943](#)
Fischer, P. C., [894](#), [945](#)
Florescu, Daniela D., [943](#), [944](#)
Friedman, Marc, [943](#)

G

Genesereth, Michael R., [943](#)

Gnanaprakasam, Senthil, [943](#)

Goldstein, Jonathan, [943](#)
Grahne, Gösta, [943](#)
Grant, John, [892](#), [894](#), [944](#)

H

Halevy, Alon Y., [943](#), [944](#)
Howard, J. H., [885](#), [894](#)
Hull, Richard, [874](#), [891](#), [894](#)

I

Ioannidis, Yannis E., [919](#), [943](#), [945](#)

J

Johnson, D. T., [869](#), [894](#)

K

Kambhampati, Subbarao, [943](#)
Karlalalem, Kamalakar, [943](#)
Katona, Gyula O. H., [892](#), [894](#), [944](#)
Krishnamurty, Ravi, [943](#)
Kwok, Cody T., [943](#)

L

Lambrecht, Eric, [943](#)
Lapis, George, [943](#)
Larson, Per-Åke, [943](#)
Lenzerini, Maurizio, [943](#)
Levy, Alon Y., [943](#)
Li, Qing, [943](#)
Lucchesi, C. L., [894](#), [944](#)

M

Maier, David, [893](#), [894](#), [944](#)
Mendelzon, Alberto O., [894](#), [943](#), [944](#)
Minker, Jack, [892](#), [894](#), [944](#)

N

Norcott, William D., [943](#)

O

Ordille, Joann J., [943](#)
Osborne, Sylvia L., [869](#)

P

Petrov, S. V., [894](#), [944](#)
Pirahesh, Hamid, [943](#)
Popa, Lucian, [943](#)
Potomianos, Spyros, [943](#)
Pottinger, Rachel, [943](#), [945](#)

Q

Qian, Xiaolei, [943](#)

R

Rajaraman, Anand, [943](#)
Raschid, Louiqa, [943](#)

S

Sagiv, Yehoshua, [894](#), [943](#), [944](#)
Sali, Attila, [862](#), [894](#), [895](#), [944](#), [945](#)
Sali, Attila, Sr., [945](#)
Shim, Kyuseok, [943](#)
Solomon, Marvin H., [919](#), [943](#), [945](#)
Srivastava, Divesh, [943](#)
Statman, R., [894](#)
Suciu, Dan, [943](#)
Sun, Harry, [943](#)

T

Tannen, Val, [943](#)
Thalheim, Bernhardt, [894](#), [945](#)
Tomp, Frank Wm., [869](#)
Tsatalos, Odysseas G., [919](#), [943](#), [945](#)
Tsou, D. M., [894](#), [945](#)

U

Ullman, Jeffrey David, [874](#), [888](#), [894](#),
[943-945](#)
Urata, Monica, [943](#)

V

Valduriez, Patrick, [943](#)
Vardi, Moshe Y., [943](#)
Vianu, Victor, [874](#), [891](#), [894](#), [944](#)

W

Weld, Daniel S., [943](#)
Witkowski, Andrew, [943](#)

Y

Yang, H. Z., [943](#)
Yang, Jian, [943](#)
Yu, C. T., [869](#), [894](#)

Z

Zaharioudakis, Markos, [943](#)
Zaniolo, C., [894](#), [945](#)
Ziauddin, Mohamed, [943](#)