

# Einführung in Datenbanken

## Übung 7: Duplikate in SQL

Prof. Dr. Stefan Brass

PD Dr. Alexander Hinneburg

Martin-Luther-Universität Halle-Wittenberg

Wintersemester 2020/21

<http://www.informatik.uni-halle.de/~brass/db20/>

















































# Erstes Beispiel zu Duplikaten (4)

- Damit hat man für diesen Zustand zwei Variablenbelegungen,
  - die beide die `WHERE`-Bedingung erfüllen,
  - beide die gleiche Ausgabe (unter `SELECT`) erzeugen.
- Man kann das auch mit dem naiven Auswertungsalgorithmus (Schleife über allen Tabellenzeilen) verstehen:

```

for S in STUDENTEN do
    if like(S.VORNAME, 'Alex%') then
        print S.NACHNAME
  
```

Mit der Schleife werden gerade die möglichen Variablenbelegungen durchgegangen.

- Im Beispiel wären die Duplikate wohl erwünscht.



# Algorithmus für mögliche Duplikate (1)

- In der Vorlesung ist ein Algorithmus angegeben, der eine hinreichende Bedingung für ein überflüssiges **DISTINCT** enthält:
  - Wenn der Algorithmus ausgibt, dass die Anfrage (ohne **DISTINCT**) keine Duplikate liefern kann, ist das sicher.
  - Im anderen Fall weiss man nicht, ob die Anfrage wirklich Duplikate liefern kann, oder der Algorithmus zu schwach ist.
    - Da die Frage unentscheidbar ist, kann es keinen Algorithmus geben, der immer eine definitive (und korrekte) Antwort gibt. Man kann das Duplikatsproblem auf das Konsistenzproblem zurückführen und umgekehrt.
- Der Algorithmus sollte für die Klausuraufgaben ausreichen.
  - Sie können die Frage nach Duplikaten natürlich auch weniger formal klären. Stellen Sie sich die (geschachtelten) Schleifen über den Tabellen vor (eine für jede Tupelvariable).

# Algorithmus für mögliche Duplikate (2)

- Der Algorithmus verwendet eine Menge  $\mathcal{K}$ 
  - von Attributreferenzen, also Einträgen der Form  $\langle \text{Tupelvariable} \rangle . \langle \text{Attribut} \rangle$ ,
  - von denen bekannt ist, dass sie für eine gegebene Ausgabezeile nur einen Wert haben können.
- Das gilt insbesondere für Attribut-Referenzen, die direkt unter **SELECT** stehen, im Beispiel also **S.Nachname**.
- Falls nun die Menge  $\mathcal{K}$  einen Schlüssel von jeder Tupelvariable enthält, kann es zu einer gegebenen Ausgabezeile nicht zwei verschiedene Variablenbelegungen geben.

Die Ausgabezeile bestimmt dann ja die Schlüsselwerte eindeutig. Zu einer gegebenen Ausgabezeile kann es dann nur jeweils eine Tabellenzeile geben, die man der Tupelvariablen zuweisen kann, um diese Ausgabezeile zu erzeugen.

# Algorithmus für mögliche Duplikate (3)

- Im Beispiel ist:  $\mathcal{K} = \{S.Nachname\}$ .
- **Nachname** ist kein Schlüssel der Tabelle, über die **S** läuft (**STUDENTEN**).
- Daher sagt der Algorithmus nichts aus.
  - Bzw.: „Die Anfrage könnte eventuell Duplikate liefern“.
- Um ein konkretes Beispiel zu finden, versuche man zwei Variablenbelegungen für die Tupelvariablen zu konstruieren,
  - die die gleichen Werte für alle Elemente von  $\mathcal{K}$  haben, und
  - die **WHERE**-Bedingung erfüllen.

Im allgemeinen greift hier die Unentscheidbarkeit des Konsistenzproblems. Für einfache WHERE-Bedingungen wäre das natürlich automatisch möglich. D.h. man könnte den Algorithmus so ausbauen, dass er zumindest manchmal ein klares „Duplikate sind möglich“ ausgibt.

# Zweites Beispiel zu Duplikaten

## Aufgabe:

- Gegeben sei die Tabelle

```
STUDENTEN(SID, VORNAME, NACHNAME, EMAILo)
```

- Die Kombination VORNAME, NACHNAME sei auch Schlüssel.

Neben dem Primärschlüssel SID.

- Kann folgende Anfrage Duplikate liefern?

```
SELECT S.NACHNAME  
FROM   STUDENTEN S  
WHERE  S.VORNAME = 'Alex'
```

Ja/Nein-Abstimmung.

# Zweites Beispiel zum Duplikat-Algorithmus

- Die Menge  $\mathcal{K}$  wird wie eben mit Attributreferenzen initialisiert, die direkt unter SELECT als Ausgabespalte stehen:

$$\mathcal{K} = \{S.NACHNAME\}.$$

Um eventuelle Alternativen mit OR auszuschließen, wandelt man die Bedingung zunächst in konjunktive Normalform um (eine mit AND-verknüpfte Liste), und schaut dann, ob ein Element dieser Liste die Form  $X.A = c$  hat.

- Dann werden alle Attribut-Referenzen hinzugefügt, die aufgrund der WHERE-Klausel einen festen Wert haben:

$$\mathcal{K} = \{S.NACHNAME, S.VORNAME\}.$$

- Da **VORNAME** und **NACHNAME** zusammen einen Schlüssel von **STUDENTEN** bilden (über dieser Tabelle läuft **S**), garantiert der Algorithmus, dass diese Anfrage keine Duplikate liefern kann.

# Einfache Aufgaben zu Duplikaten (1)

a)

- Kann diese Anfrage Duplikate liefern?

```
SELECT SID
FROM STUDENTEN
WHERE NACHNAME LIKE 'M%'
```

- Wie eben sind die Schlüssel von **STUDENTEN**:
  - **SID** (Primärschlüssel)
  - **VORNAME, NACHNAME** (Sekundärschlüssel)
- Ja/Nein-Abstimmung.
- Was ist **K**?

# Einfache Aufgaben zu Duplikaten (2)

a)

- Nein (kann keine Duplikate liefern).
- $\mathcal{K} = \{\text{STUDENTEN.SID}\}$ .

b)

- Kann diese Anfrage Duplikate liefern?

```
SELECT VORNAME
FROM   STUDENTEN
WHERE  SID = 153
```

- Schlüssel wie gehabt.  
 {SID} und {VORNAME, NACHNAME}.
- Was ist  $\mathcal{K}$ ?

# Einfache Aufgaben zu Duplikaten (3)

b)

- Nein (kann keine Duplikate liefern).
- $\mathcal{K} = \{\text{STUDENTEN.VORNAME, STUDENTEN.SID}\}$ .

SID ist Schlüssel von STUDENTEN.

Es ist hier unwesentlich, dass VORNAME auch in  $\mathcal{K}$  enthalten ist.

c)

- Kann diese Anfrage Duplikate liefern?

```
SELECT VORNAME
FROM   STUDENTEN
WHERE  VORNAME = 'Dorothea'
```

- Was ist  $\mathcal{K}$ ?



# Einfache Aufgaben zu Duplikaten (4)

c)

- Ja (kann Duplikate liefern).

Es kann mehrere Studentinnen geben, die Dorothea heißen.

- $\mathcal{K} = \{\text{STUDENTEN.VORNAME}\}$ .

VORNAME allein ist nicht Schlüssel von Studenten.

d)

- Kann diese Anfrage Duplikate liefern?

```
SELECT SID
FROM   STUDENTEN
WHERE  VORNAME = 'Lisa' OR NACHNAME = 'Weiss'
```

# Einfache Aufgaben zu Duplikaten (5)

d)

- Nein (keine Duplikate).

Obwohl es „Lisa Weiss“ gibt, die beide Teile der Disjunktion erfüllt, wird sie nicht doppelt geliefert: Jede Variablenbelegung wird nur ein Mal getestet.

- $\mathcal{K} = \{\text{STUDENTEN.SID}\}$ .

SID.VORNAME und SID.NACHNAME sind nicht dabei, weil die WHERE-Bedingung nicht eine einzelne Konstante erzwingt (durch das OR).

e)

- Kann diese Anfrage Duplikate liefern?

```
SELECT VORNAME
FROM   STUDENTEN
WHERE  SID <= 101 AND SID >= 101
OR     SID > 200 AND SID <= 200
```

# Einfache Aufgaben zu Duplikaten (6)

e)

- Nein, die Bedingung ist äquivalent zu  $SID = 101$ .

Die Anfrage liefert also nur höchstens eine Zeile.

Der Teil nach dem OR ist inkonsistent.

- Der Algorithmus liefert  $\mathcal{K} = \{STUDENTEN.VORNAME\}$  und deswegen „ich weiss nicht/Duplikate eventuell möglich“.

Äquivalenz ist unentscheidbar. Der Algorithmus findet  $X.A = c$  nur, wenn es syntaktisch offensichtlich ist.

f)

- Kann diese Anfrage Duplikate liefern?

```
SELECT X.SID
FROM STUDENTEN X, STUDENTEN Y
```

# Einfache Aufgaben zu Duplikaten (7)

f)

- Ja, Duplikate sind möglich: Die Tupelvariable **Y** erzeugt eine zusätzliche Schleife über den STUDENTEN.

```
for X in STUDENTEN do
    for Y in STUDENTEN do
        print X.SID
```

- Wenn STUDENTEN  $n$  Zeilen hat, werden  $n^2$  Ergebniszeilen ausgegeben, und zwar die  $n$  verschiedenen SIDs jeweils  $n$  mal.  
D.h. im Fall  $n = 0$  und  $n = 1$  gibt es keine Duplikate. Die Frage ist aber immer, ob es mindestens einen DB-Zustand gibt, in dem die Anfrage Duplikate erzeugt. Das gilt hier natürlich.
- $\mathcal{K} = \{X.SID\}$ .  
Dies enthält zwar einen Schlüssel von **X**, aber nicht von **Y**.

# Algorithmus bei mehreren Tupelvariablen

- Am Anfang wird  $\mathcal{K}$  wie gehabt initialisiert mit
  - Attributreferenzen  $X.A$  als Ausgabespalten unter **SELECT**
  - Attributreferenzen, die aufgrund der **WHERE**-Bedingung nur einen Wert annehmen können:  $X.A = c$ .
- Anschließend wird  $\mathcal{K}$  mit den folgenden beiden Regeln erweitert (iteriert, bis sich nichts mehr ändert):
  - Wenn die **WHERE**-Bedingung eine Gleichung  $X.A = Y.B$  enthält (als Teil einer **AND**-Liste), und  $X.A$  ist schon in  $\mathcal{K}$  enthalten, wird  $Y.B$  hinzugefügt (und umgekehrt).
  - Wenn ein Schlüssel einer Tupelvariable enthalten ist, werden alle anderen Attribute dieser Tupelvariable auch hinzugefügt.

# Beispiel mit mehreren Tupelvariablen (1)

- Kann diese Anfrage Duplikate liefern?

```
SELECT S.VORNAME, S.NACHNAME
FROM   STUDENTEN S, BEWERTUNGEN B
WHERE  S.SID = B.SID
AND    B.ATYP = 'H' AND B.ANR = 1
```

- $\mathcal{K} := \{S.VORNAME, S.NACHNAME\}$

Diese Attributreferenzen stehen explizit unter SELECT.

- $\mathcal{K} := \mathcal{K} \cup \{B.ATYP, B.ANR\}$

Wegen der WHERE-Bedingungen  $B.ATYP = 'H'$  und  $B.ANR = 1$ , die für alle Variablenbelegungen gelten müssen (d.h. nicht durch OR ausgehebelt).

- $\mathcal{K} := \mathcal{K} \cup \{S.SID, S.EMAIL\}$

$\{S.VORNAME, S.NACHNAME\} \subseteq \mathcal{K}$  bilden einen Schlüssel von S.

## Beispiel mit mehreren Tupelvariablen (2)

- $\mathcal{K} := \mathcal{K} \cup \{B.SID\}$

Wegen der WHERE-Bedingung  $S.SID = B.SID$  und  $S.SID \in \mathcal{K}$ .

- $\mathcal{K} := \mathcal{K} \cup \{B.PUNKTE\}$

$\{B.SID, B.ATYP, B.ANR\} \subseteq \mathcal{K}$  bilden einen Schlüssel von B.

- Nun enthält  $\mathcal{K}$  alle Attribute von allen Tupelvariablen.

- Damit ist garantiert, dass die Anfrage keine Duplikate liefern kann.

- Der Algorithmus ist nicht prüfungsrelevant.

Der Algorithmus ist (vermutlich) von Brass/Goldberg.

- Die Fähigkeit, festzustellen, ob eine Anfrage Duplikate liefern kann, schon.

# Aufgabe zu Duplikaten (1)

- Kann diese Anfrage Duplikate liefern?

```
SELECT S.VORNAME, S.NACHNAME
FROM   STUDENTEN S, BEWERTUNGEN B
WHERE  S.SID = B.SID
AND    B.PUNKTE = 10
AND    B.ATYP = 'H'
```

- Ja/Nein-Abstimmung.
- Falls Nein: Beweisen Sie die Aussage, z.B. mit dem obigen Algorithmus.
- Falls Ja: Geben Sie einen Datenbankzustand an, in dem die Anfrage ein Duplikat liefert.







# Beispiel-Datenbank

## STUDENTEN

<u>SID</u>	<u>VORNAME</u>	<u>NACHNAME</u>	<u>EMAIL</u>
101	Lisa	Weiss	...
102	Michael	Grau	NULL
103	Daniel	Sommer	...
104	Iris	Winter	...

## AUFGABEN

<u>ATYP</u>	<u>ANR</u>	<u>THEMA</u>	<u>MAXPT</u>
H	1	ER	10
H	2	SQL	10
Z	1	SQL	14

## BEWERTUNGEN

<u>SID</u>	<u>ATYP</u>	<u>ANR</u>	<u>PUNKTE</u>
101	H	1	10
101	H	2	8
101	Z	1	12
102	H	1	9
102	H	2	9
102	Z	1	10
103	H	1	5
103	Z	1	7

# Logische Analyse von Anfragen (5)

1. 

```
SELECT B.SID, A.ANR AS NR
FROM   BEWERTUNGEN B, AUFGABEN A
WHERE  B.ATYP = A.ATYP AND A.ATYP = 'H'
AND    A.ANR = B.ANR AND B.PUNKTE = A.MAXPT
```
2. 

```
SELECT Y.SID, X.ANR AS NR
FROM   AUFGABEN X, BEWERTUNGEN Y
WHERE  X.MAXPT = Y.PUNKTE
AND    X.ATYP='H' AND Y.ATYP='H' AND X.ANR=Y.ANR
```

Kreuzen Sie die erste zutreffende Antwort an:

- A. Die Anfragen liefern immer die gleiche Antwort (äquivalent).
- B. Bis auf Duplikate die gleiche Antwort.
- C. Anfrage 1 liefert immer  $\emptyset$  (inkonsistent).
- D. Anfrage 2 liefert immer  $\emptyset$  (inkonsistent).
- E. Keine der Aussagen trifft zu.











# Präsenzaufgabe: Duplikate (2)

- Kann die folgende Anfrage Duplikate liefern?

```

SELECT S.VORNAME, S.NACHNAME
FROM   STUDENTEN LISA, STUDENTEN S,
       BEWERTUNGEN H1L, BEWERTUNGEN H1S
WHERE  LISA.VORNAME = 'Lisa'
AND    LISA.NACHNAME = 'Weiss'
AND    H1L.SID = LISA.SID
AND    H1S.SID = S.SID
AND    H1L.ATYP = 'H' AND H1S.ATYP = 'H'
AND    H1L.ANR = H1S.ANR
AND    H1S.PUNKTE > H1L.PUNKTE

```

Übliche Schlüssel: {SID} und {VORNAME, NACHNAME} von STUDENTEN,  
{SID, ATYP, ANR} von BEWERTUNGEN.

- Ja/Nein plus gute Begründung (2 Punkte)

