

# Anfrageoptimierung in Datenbanken — Theorie und Praxis

Vorlesung im WS 2017/18, HU Berlin,  
Prof. J.C. Freytag & Prof. N. Schweikardt

Theorie Teil  
Prof. N. Schweikardt

# Kapitel 1: Grundbegriffe

- $\mathbb{N} := \{0, 1, 2, 3, \dots\}$
- $\mathbb{N}_{\geq 1} := \mathbb{N} \setminus \{0\}$

Für  $m, n \in \mathbb{N}$  ist

- $[m, n] := \{i \in \mathbb{N} : m \leq i \leq n\}$
- $[n] := [1, n]$

Ist  $M$  eine Menge, so schreiben wir  $X \subseteq_e M$ , um auszudrücken, dass  $X$  eine endliche Teilmenge von  $M$  ist.

Wir schreiben  $\mathcal{P}(M)$  oder  $2^M$ , um die Potenzmenge von  $M$  zu bezeichnen.

Für  $k \in \mathbb{N}$  ist

$$\binom{M}{k} := \{X \subseteq_e M : |X| = k\}$$

die Menge aller  $k$ -elementigen Teilmengen von  $M$ ,

und 
$$\binom{M}{\leq k} := \bigcup_{i=0}^k \binom{M}{i}$$

Für ein Tupel  $t = (t_1, \dots, t_k)$  und Zahlen  $i_1, \dots, i_k \in [k]$  ist  $\pi_{i_1, \dots, i_k}(t) := (t_{i_1}, \dots, t_{i_k})$ .

Abkürzungen:

- "f.a." steht für "für alle"
- "ex." steht für "es existiert" bzw. "es gibt"
- "s.d." steht für "so dass"
- "DB" steht für "Datenbank"

## 1.1 Datenbanken

Für den gesamten Theorieteil sei  $\text{dom}$  eine fest gewählte, abzählbar unendliche Menge. Elemente in  $\text{dom}$  werden auch Konstanten genannt; wir bezeichnen sie oft mit Kleinbuchstaben wie  $a, b, c$ .

Ein (Datenbank-) Schema  $\sigma$  ist eine endliche Menge von Relationssymbolen (auch: Relationsnamen), wobei jedem  $R \in \sigma$  eine feste Stelligkeit  $\text{ar}(R) \in \mathbb{N}$  zugeordnet ist.

Eine Datenbank vom Schema  $\sigma$  (kurz:  $\sigma$ -DB) besteht aus einer Relation  $R^D \subseteq \text{dom}^{\text{ar}(R)}$  für jedes  $R \in \sigma$ .

Für  $\sigma = \{R_1, \dots, R_s\}$  schreiben wir auch  $D = (R_1^D, \dots, R_s^D)$ , um eine  $\sigma$ -DB zu bezeichnen.

Der active domain  $\text{adom}(D)$  von  $D$  ist die kleinste Menge  $A \subseteq \text{dom}$ , für die gilt:  
 $R^D \subseteq A^{\text{ar}(R)}$  f.a.  $R \in \sigma$ .

## 1.2 Konjunktive Anfragen

Sei  $\text{var}$  eine fest gewählte, abzählbar unendliche Menge mit  $\text{var} \cap \text{dom} = \emptyset$ .

Elemente in  $\text{var}$  nennen wir Variablen; wir bezeichnen sie oft mit Großbuchstaben wie  $A, B, C, X, Y, \dots$ .

Ein Term ist ein Element aus  $\text{var} \cup \text{dom}$ .

Ein freies Tupel der Stelligkeit  $k \in \mathbb{N}$  ist ein Element aus  $(\text{var} \cup \text{dom})^k$ .

Sei  $\sigma$  ein Schema.

Ein  $\sigma$ -Atom  $\alpha$  ist von der Form  $R(u_1, \dots, u_r)$  mit  $R \in \sigma$ ,  $r = \text{ar}(R)$  und  $u_1, \dots, u_r \in \text{var} \cup \text{dom}$ .

Wir setzen  $\text{vars}(\alpha) := \{u_1, \dots, u_r\} \cap \text{var}$

und  $\text{cons}(\alpha) := \{u_1, \dots, u_r\} \cap \text{dom}$ .

Eine (regelbasierte) konjunktive Anfrage  
(kurz: CQ, für "conjunctive query")  $Q$  vom  
Schema  $\sigma$  ist von der Form

$$Q(v_1, \dots, v_k) \leftarrow \alpha_1, \dots, \alpha_m \quad (*)$$

- mit
- $k \geq 0$ ,  $m \geq 1$ ,  $v_1, \dots, v_k \in \text{var } \sigma \text{ dom}$ ,
  - für jedes  $i \in [m]$  ist  $\alpha_i$  ein  $\sigma$ -Atom und
  - $\{v_1, \dots, v_k\} \cap \text{var} \subseteq \bigcup_{i \in [m]} \text{vars}(\alpha_i)$ .

Eine Anfrage der Form  $(*)$  wird auch Regel  
genannt; das "Atom"  $Q(v_1, \dots, v_k)$  heißt  
Kopf der Regel; " $\alpha_1, \dots, \alpha_m$ " heißt Rumpf der Regel;  
die Zahl  $k$  ist die Stelligkeit der Anfrage.

Mit  $\text{vars}(Q)$  (bzw.  $\text{cons}(Q)$ ) bezeichnen wir die  
Menge aller Variablen (bzw. Konstanten), die in  $Q$   
vorkommen. D.h.:

$$\text{vars}(Q) = \bigcup_{i \in [m]} \text{vars}(\alpha_i) = (\{v_1, \dots, v_k\} \cap \text{var}) \cup \bigcup_{i \in [m]} \text{vars}(\alpha_i)$$

$$\text{cons}(Q) = (\{v_1, \dots, v_k\} \cap \text{dom}) \cup \bigcup_{i \in [m]} \text{cons}(\alpha_i).$$

5  
Eine Belegung für  $Q$  ist eine Abbildung

$$\beta: \text{vars}(Q) \cup \text{dom} \rightarrow \text{dom} \quad \text{mit}$$

$$\beta(a) = a \quad \text{f. a. } a \in \text{dom}.$$

Eine Belegung  $\beta$  für  $Q$  ist ein Homomorphismus von  $Q$  auf eine  $\sigma$ -DB  $D$

(kurz:  $\beta: Q \rightarrow D$ ), falls für jedes Atom  $\alpha$  der Form  $R(u_1, \dots, u_r)$  im Rumpf von  $Q$  gilt:  $(\beta(u_1), \dots, \beta(u_r)) \in D^R$ .

Das Anfrageergebnis (kurz: Ergebnis, Resultat) von  $Q$  auf  $D$  ist die Menge

$$Q(D) := \left\{ (\beta(v_1), \dots, \beta(v_k)) : \beta \text{ ist ein Homomorphismus von } Q \text{ auf } D \right\}$$

Beobachtung 1.1

$$Q(D) \subseteq \underbrace{(\text{adom}(D) \cup \text{cons}(Q))}_{{=: \text{adom}(Q, D)}}^k$$

## Beispiel 1.2

Sei  $\sigma_E := \{E\}$  mit  $\text{ar}(E) = 2$ .

Eine  $\sigma_E$ -DB  $D$  besteht aus einer 2-stelligen Relation  $E^D \subseteq_e \text{dom} \times \text{dom}$ .

Wir können sie z.B. als "soziales Netzwerk" interpretieren, bei dem ein Tupel  $(a,b) \in E^D$  besagt "Person  $a$  hat die Nachricht von Person  $b$  abonniert" — kurz: Person  $a$  ist ein "Follower" von Person  $b$ . Wir können uns  $D$  also vorstellen als Datenbank, die aus einer Tabelle  $E^D$  der Form

Follower	Author
...	...

(b) Gesucht: Eine Anfrage  $Q_\Delta$ , s.d. auf jeder  $\sigma_E$ -DB  $D$  gilt:

$Q_\Delta(D)$  besteht aus genau denselbigen Tripeln  $(a,b,c)$ , für die gilt:

- $a$  ist Follower von  $b$ ,
- $b$  ist Follower von  $c$  und
- $c$  ist Follower von  $a$ .

In SQL:

```
SELECT DISTINCT E1.Follower, E2.Follower, E3.Follower
FROM E AS E1, E AS E2, E AS E3
WHERE E1.Author = E2.Follower AND
      E2.Author = E3.Follower AND
      E3.Author = E1.Follower
```

Als CQ:

$Q_A(A, B, C) \leftarrow E(A, B), E(B, C), E(C, A)$

(a) Gesucht: Eine Anfrage  $Q_{(b)}$ , s.d. auf jeder  $\sigma_E$ -DB  $D$  gilt:  $Q_{(b)}(D)$  besteht aus genau denjenigen 1-Tupeln  $(a)$ , für die gilt:  
 $a$  ist Follower von Sascha Lobo und von Mesut Özil

Als CQ:

$Q_{(b)}(A) \leftarrow E(A, \text{Sascha Lobo}), E(A, \text{Mesut Özil})$

In SQL:

```
SELECT DISTINCT E1.Follower
FROM E AS E1, E AS E2
WHERE E1.Author = 'Sascha Lobo' AND
      E2.Author = 'Mesut Özil' AND
      E1.Follower = E2.Follower
```



## 1.3 Das Auswertungsproblem

Das Auswertungsproblem für CQ ist das Berechnungsproblem mit

Eingabe: Ein Schema  $\sigma$ ,  
eine  $\sigma$ -DB  $D$  und  
eine CQ  $Q$  von Schema  $\sigma$

Ausgabe:  $Q(D)$

### Beobachtung 1.3

Das Auswertungsproblem für CQ lässt sich mit dem folgenden einfachen, aber nicht gerade effizienten Algorithmus lösen:

Bei Eingabe von  $\sigma, D, Q$ :

- 1) Berechne  $s := |\text{vars}(Q)|$  und  $\{x_1, \dots, x_s\} = \text{vars}(Q)$
- 2) Berechne  $N := \text{atom}(Q, D)$  und  $\{a_1, \dots, a_N\} = \text{atom}(Q, D)$
- 3) Ergebnis :=  $\emptyset$
- 4) Für jedes  $(b_1, \dots, b_s) \in \{a_1, \dots, a_N\}^s$  tue Folgendes:  
Sei  $\beta$  die Belegung für  $Q$  mit  $\beta(x_i) = b_i$  für  $i \in [s]$ .  
Teste, ob für jedes Atom  $\alpha$  der Form  $R(x_1, \dots, x_r)$   
im Rumpf von  $Q$  gilt:  $(\beta(x_1), \dots, \beta(x_r)) \in R$ .  
Falls ja, so Ergebnis := Ergebnis  $\cup \{(\beta(x_1), \dots, \beta(x_r))\}$   
(wobei der Kopf von  $Q$  von der Form  $Q(x_1, \dots, x_s)$  ist).
- 5) Gib Ergebnis aus.

Unter Verwendung von Beobachtung 1.1 kann man sich leicht davon überzeugen, dass dieser Algorithmus tatsächlich  $Q(D)$  ausführt.

Die Laufzeit wird im Wesentlichen dominiert durch die Anzahl  $N^s$  der Schleifendurchläufe in Zeile 4).

Satz 1.4 (Satz von Chandra und Merlin, 1977)

Das Auswertungsproblem für  $Q$ -stellige konjunktive Anfragen ist NP-vollständig.

Beweisidee:

Zugehörigkeit zu NP:

Bei Eingabe von  $\sigma, D, Q$ , führe die Schritte 1), 2), 3) aus Beobachtung 1.3 durch.

An Stelle von Schritt 4) tue Folgendes:

4) Rate ein Tupel  $(b_1, \dots, b_s) \in \{a_1, \dots, a_n\}^s$ .

Betrachte die Belegung  $\beta$  mit  $\beta(x_i) = b_i$  f.a.  $i \in \{s\}$ .  
Teste, ob für jedes Atom  $\alpha$  der Form  $R(u_1, \dots, u_r)$  im Rumpf von  $Q$  gilt:  $(\beta(u_1), \dots, \beta(u_r)) \in R^D$ .

Falls ja, so Ergebnis :=  $\{()\}$ .

5) Gib Ergebnis aus

Man kann sich leicht davon überzeugen, dass dies ein nichtdeterministischer Polynomialzeit-Algorithmus ist, der das Auswertungsproblem für  $Q$  löst.

## NP-Härte:

wir reduzieren das NP-vollständige Problem

### CLIQUE

Eingabe: Ein ungerichteter endlicher Graph  
 $G = (V(G), E(G))$  und  
 eine Zahl  $k \in \mathbb{N}$ .

Frage: Besitzt  $G$  eine Clique der Größe  $k$ ,  
 d.h. gibt es Knoten  $v_1, \dots, v_k \in V(G)$ , s.d.  
 f.a.  $i, j \in [k]$  mit  $i \neq j$  gilt:  $\{v_i, v_j\} \in E(G)$ ?

auf's Auswertungsproblem für CQ:

Für eine Eingabe  $(G, k)$  für's CLIQUE-Problem  
 konstruieren wir eine Eingabe  $(\sigma, D, Q)$  für's  
 Auswertungsproblem für CQ, s.d. gilt:

$G$  besitzt eine Clique der Größe  $k \iff Q(D) \ni \{()\}$ .

$\Updownarrow$  Def. Semantik von CQs

es gibt einen Homomorphismus  
 $\beta: Q \rightarrow D$

Dann wählen wir  $\sigma := \sigma_E = \{E\}$ ,

$D := (E^D)$  mit  $E^D := \{(a, b) : \{a, b\} \in E(G)\}$  und

$Q := Q_{\min\{k, |V(G)|+1\}}$ , wobei für jedes  $k' \geq 1$  die  
 Anfrage  $Q_{k'}$  die CQ mit Kopf  $Q_{k'}()$  ist, deren  
 Rumpf aus allen Atomen  $E(x_i, x_j)$  mit  $1 \leq i < j \leq k'$   
 besteht, wobei  $x_1, \dots, x_{k'}$   $k'$  verschiedene Variablen sind.

Beispiel:  $Q_3() \leftarrow E(x_1, x_2), E(x_2, x_3), E(x_1, x_2)$

Man sieht leicht, dass gilt

$Q_k(D) \ni \{1\} \Leftrightarrow G$  besitzt eine Clique der Größe  $k$ .

Somit gilt auch:

$Q(D) \ni \{1\} \Leftrightarrow G$  besitzt eine Clique der Größe  $k$ .

Außerdem kann  $(G, D, Q)$  bei Eingabe von  $(G, k)$  deterministisch in Zeit polynomiell in der Größe von  $G$  und (der Binärdarstellung von)  $k$  erzeugt werden.

Somit haben wir eine Polynomialzeit-Reduktion vom CLIQUE-Problem auf's Antwortungsproblem für  $Q$  konstruiert.

□

### Folgerung 1.5

Falls  $P \neq NP$  ist, so gibt es keinen Algorithmus, der das Antwortungsproblem für  $Q$ -Anfragen deterministisch in Zeit  $(k+n)^{O(1)}$  löst, wobei  $k$  die Größe der Anfrage und  $n$  die Größe der Datenbank ist.

## Bemerkung 1.6

Unter Verwendung einer stärkeren Annahme aus der parametrisierten Komplexitätstheorie lässt sich sogar Folgendes zeigen:

### Satz von Papadimitriou und Yannakakis, 1997

Falls  $FPT \neq W[1]$ , so gibt es keinen Algorithmus, der das Auswertungsproblem für  $Q$ -stellige CQ-Anfragen deterministisch in Zeit in Zeit  $f(k) \cdot n^c$  löst, wobei  $c$  irgendeine natürliche Zahl, und  $f$  irgendeine berechenbare Funktion ist, und  $k$  die Größe der Anfrage und  $n$  die Größe der Datenbank bezeichnet.

$FPT$  und  $W[1]$  sind Komplexitätsklassen, die in der parametrisierten Komplexitätstheorie Rollen spielen, die mit den Rollen von  $P$  und  $NP$  in der klassischen Komplexitätstheorie vergleichbar sind.