

## Kapitel 2: Die AGM-Schranke und andere obere Schranken für die Größe von Anfrageergebnissen

Der Begriff "AGM-Schranke" ist benannt nach den Autoren Atserias, Grohe und Marx des Artikels "Size bounds and query plans for relational joins" in Proc. FOCS 2008, pp. 739-748.

### 2.1 Warm-Up: Die $\Delta$ -Anfrage

Betrachte die Anfrage  $Q_\Delta$  aus Beispiel 1.2

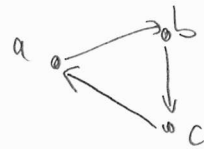
$$Q_\Delta(A, B, C) \leftarrow E(A, B), E(B, C), E(C, A).$$

Sei  $\sigma := \sigma_E = \{E\}$ . Jede  $\sigma$ -DB  $D$  entspricht einem gerichteten Graphen  $G^D = (V^D, E^D)$  mit  $V^D := \text{atom}(D)$ .

Klar: Für jede  $\sigma$ -DB  $D$  ist

$$Q_\Delta(D) = \{ (a, b, c) : (a, b) \in E^D \text{ und } (b, c) \in E^D \text{ und } (c, a) \in E^D \}$$

Skizze:



D.h.:  $Q_\Delta(D)$  besteht aus allen gerichteten Dreiecken in  $G^D$ .

Ziel: Gib eine obere Schranke für die Anzahl  $|Q_\Delta(D)|$  von Tupeln in  $Q_\Delta(D)$  an und finde einen möglichst effizienten Algorithmus, der  $Q_\Delta(D)$  berechnet.

Notation:  $n^D := |\text{dom}(D)|$   
 $N^D := |E^D|$

Offensichtliche obere Schranken für  $|Q_\Delta(D)|$ :

1)  $|Q_\Delta(D)| \leq (n^D)^3$

2)  $|Q_\Delta(D)| \leq (N^D)^2$

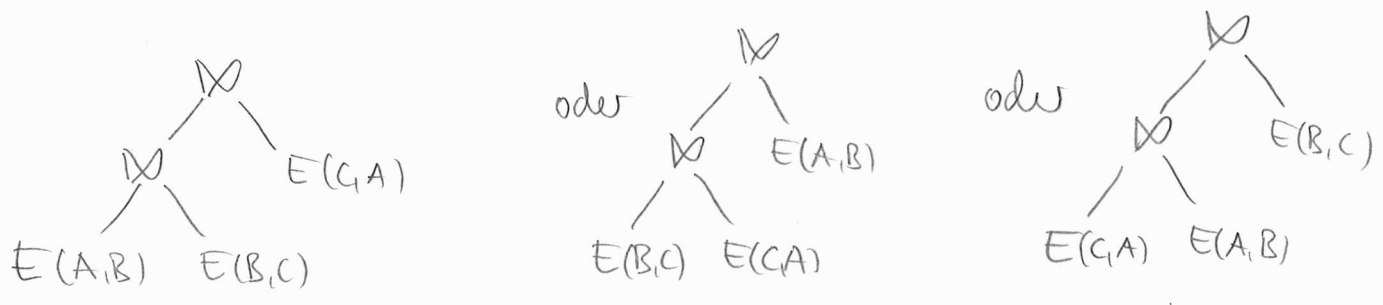
Frage: Geht das besser?

Antwort, die der nächste Satz (Satz  $\Delta$ ) gibt:

Ja:  $|Q_\Delta(D)| \leq 2 \cdot (N^D)^{1,5}$  — und  
 es gibt auch einen Algorithmus, der  $Q_\Delta(D)$   
 mit Laufzeit  $\approx (N^D)^{1,5}$  löst.

Bevor wir uns den Satz (linker Beweis) anschauen,  
 diskutieren wir aber erstmal, wie ein  
 herkömmliches Datenbanksystem die Anfrage  $Q_\Delta$   
 auswertet (siehe Beispiel 1.1 für eine Formulierung  
 der Anfrage in SQL).

Um die Anfrage  $Q_\Delta$  auszuwerten wird ein Datenbanksystem i.d.R. QEPs (query evaluation plans) der folgenden Art betrachtet



Bei jedem dieser QEPs wird als Zwischenergebnis  $Q_2(D)$  berechnet, für

$$Q_2(x,y,z) \leftarrow E(x,y), E(y,z)$$

und dann wird

$$Q_\Delta(x,y,z) \leftarrow Q_2(x,y,z), E(z,x)$$

auf  $D$  ausgewertet.

## Beispiel 2.1

Für jede Zahl  $m \geq 1$  geben wir eine  $\sigma$ -DB  $D_m$  mit  
 $N^{D_m} = \Theta(m)$ ,  $|Q_\Delta(D_m)| = \Theta(m)$ , aber

$$|Q_2(D_m)| = \Omega(m^2) \quad \text{für} \quad Q_2(X, Y, Z) \leftarrow E(X, Y), E(Y, Z).$$

Die "üblichen" QEPs (query evaluation plans), die von Datenbanksystemen zur Auswertung der Anfrage  $Q_\Delta$  erzeugt werden, benötigen zum Auswerten von  $Q_\Delta$  auf  $D_m$  daher Zeit  $\Omega((N^{D_m})^2)$ , während der Algorithmus aus dem folgenden Satz  $\Delta$  nur Zeit  $O((N^{D_m})^{1.5})$  benötigt.

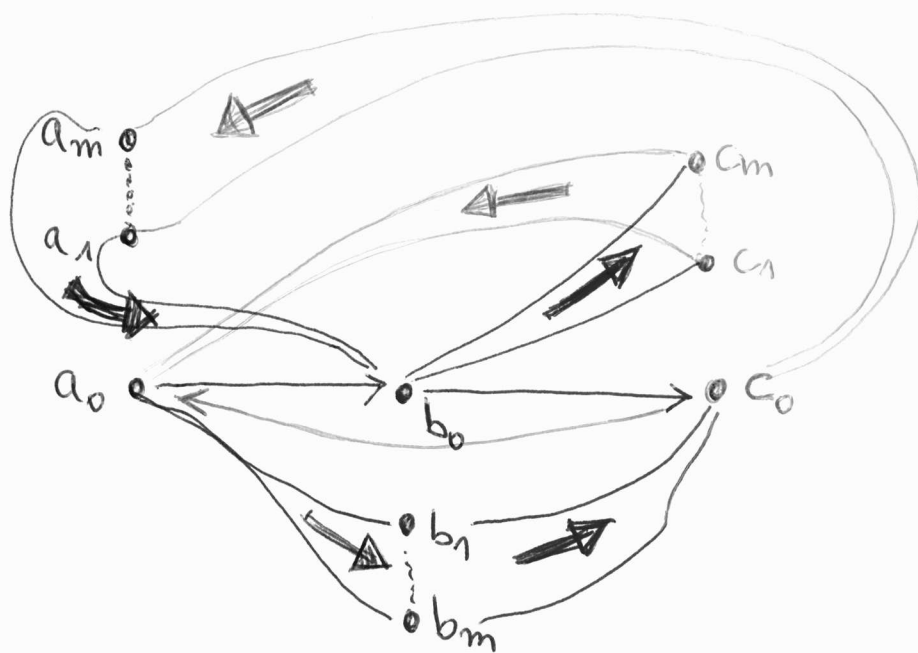
Konstruktion von  $D_m$ :

Seien  $a_0, a_1, \dots, a_m, b_0, b_1, \dots, b_m, c_0, c_1, \dots, c_m$   
 $3(m+1)$  verschiedene Elemente aus  $\text{dom}$ . Sei

$$E^{D_m} := \left\{ \begin{array}{l} \{ (a_0, b_i) : i \in [0, m] \} \cup \\ \{ (a_i, b_0) : i \in [1, m] \} \cup \\ \{ (b_0, c_i) : i \in [0, m] \} \cup \\ \{ (b_i, c_0) : i \in [1, m] \} \cup \\ \{ (c_0, a_i) : i \in [0, m] \} \cup \\ \{ (c_i, a_0) : i \in [1, m] \} \end{array} \right.$$

klar:  $N^{D_m} := |E^{D_m}| = 3 \cdot (m+1+m) = 6m+3$

Skizze:  $D_m$



Man kann sich leicht davon überzeugen, dass gilt:

$$Q_{\Delta}(D_m) = X \cup \pi_{2,3,1}(X) \cup \pi_{3,1,2}(X) \quad \text{für}$$

$$X = \left\{ (a_0, b_0, c_i) : i \in [0, m] \right\} \cup \\ \left\{ (a_0, b_i, c_0) : i \in [1, m] \right\} \cup \\ \left\{ (a_i, b_0, c_0) : i \in [1, m] \right\}.$$

Es gilt:  $|X| = m+1 + 2m = 3m+1$  und  $|Q_{\Delta}(D_m)| = 3|X| = 9m+3$

Außerdem ist  $(a_i, b_0, c_j) \in Q_2(D_m)$ , f.a.  $i, j \in [0, m]$ ,

also  $|Q_2(D_m)| > m^2$

Satz Δ: Für jede  $\sigma$ -DB  $D$  gilt:

$$|Q_{\Delta}(D)| \leq 2 \cdot (N^D)^{3/2}$$

und es gibt einen Algorithmus, der bei Eingabe von  $D$  die Menge  $Q_{\Delta}(D)$  in Zeit

$O((N^D)^{3/2} \cdot t)$  berechnet, wobei  $t$  die Zeit ist, die wir benötigen, um bei Eingabe von  $(v,w)$  zu testen, ob  $(v,w) \in E^D$  ist.

Beweis: Sei  $D$  eine beliebige  $\sigma$ -DB.

Für jedes  $v \in \text{atom}(D)$  sei

$$\begin{aligned} \text{aus-Grad}^D(v) &:= |\{w : (v,w) \in E^D\}| \\ &= |E^D(v,*)| \end{aligned}$$

$$\text{für } E^D(v,*) := \{w : (v,w) \in E^D\}.$$

Ein Knoten  $v \in \text{atom}(D)$  heißt

- schwer, wenn  $\text{aus-Grad}^D(v) \geq \sqrt{N^D}$  ist
- leicht, wenn  $v$  nicht schwer ist.

Ein Tupel  $t = (a,b,c) \in Q_{\Delta}(D)$  heißt

schwer, wenn  $a$  schwer ist, und es heißt  
 leicht, wenn  $a$  leicht ist.

Für  $x \in \{\text{leicht}, \text{schwer}\}$  sei

$$Q_{\Delta, x}(D) := \{t \in Q_{\Delta}(D) : t \text{ ist } x\}$$

klar:  $Q_{\Delta}(D) = Q_{\Delta, \text{leicht}}(D) \cup Q_{\Delta, \text{schwer}}(D)$

und  $|Q_{\Delta}(D)| = |Q_{\Delta, \text{leicht}}(D)| + |Q_{\Delta, \text{schwer}}(D)|$ .

Behauptung 1:  $|Q_{\Delta, \text{leicht}}(D)| \leq (N^D)^{3/2}$

Beweis: Für jeden leichten Knoten  $a \in \text{adom}(D)$  ist  $\text{aus-Grad}^D(a) < \sqrt{N^D}$ , und es gilt:

$$\begin{aligned} & |\{(b,c) : (a,b,c) \in Q_{\Delta}(D)\}| \\ & \leq |\{(b,c) : (a,b) \in E^D \text{ und } (c,a) \in E^D\}| \\ & \leq \text{aus-Grad}^D(a) \cdot \text{ein-Grad}^D(a), \text{ wobei} \end{aligned}$$

$\text{ein-Grad}^D(a) := |\{w : (w,a) \in E^D\}| = |E^D(*,a)|$   
für  $E^D(*,a) := \{w : (w,a) \in E^D\}$ .

Insgesamt gilt:

$$\begin{aligned} |Q_{\Delta, \text{leicht}}(D)| & \leq \sum_{\substack{a \in \text{adom}(D), \\ a \text{ leicht}}} |\{(b,c) : (a,b,c) \in Q_{\Delta}(D)\}| \\ & \leq \sum_{\substack{a \in \text{adom}(D), \\ a \text{ leicht}}} \text{aus-Grad}^D(a) \cdot \text{ein-Grad}^D(a) \\ & \leq \sqrt{N^D} \cdot \sum_{a \in \text{adom}(D)} \text{ein-Grad}^D(a) \\ & \leq \sqrt{N^D} \cdot |E^D| \\ & = \sqrt{N^D} \cdot N^D = (N^D)^{3/2} \end{aligned}$$

□ Beh 1.

Sei  $S^D$  die Menge aller schweren Knoten von  $D$ ,  
d.h.  $S^D := \{v \in \text{dom}(D) : v \text{ ist schwer}\}$ .

20

Behauptung 2:  $|Q_{\Delta, \text{schwer}}(D)| \leq |S^D| \cdot N^D$

Beweis:  $Q_{\Delta, \text{schwer}}(D)$   
 $= \{(a, b, c) : a \in S^D \text{ und } (a, b, c) \in Q_{\Delta}(D)\}$   
 $\subseteq \{(a, b, c) : a \in S^D \text{ und } (b, c) \in E^D\}$   
 $= S^D \times E^D$

Somit:  $|Q_{\Delta, \text{schwer}}(D)| \leq |S^D \times E^D| = |S^D| \cdot |E^D| = |S^D| \cdot N^D$

□ Beh 2

Behauptung 3:  $|S^D| \leq \sqrt{N^D}$

Beweis: Für jedes  $v \in S^D$  ist  $|E^D(v, *)| = \text{aus-Grad}^D(v) \geq \sqrt{N^D}$ .

Es gilt:

$$\begin{aligned} N^D = |E^D| &\geq \left| \bigcup_{v \in S^D} \{(v, w) : w \in E^D(v, *)\} \right| \\ &= \sum_{v \in S^D} |\{(v, w) : w \in E^D(v, *)\}| \\ &= \sum_{v \in S^D} |E^D(v, *)| \\ &\geq \sum_{v \in S^D} \sqrt{N^D} \\ &= |S^D| \cdot \sqrt{N^D} \end{aligned}$$

Somit ist  $|S^D| \leq \frac{N^D}{\sqrt{N^D}} = \sqrt{N^D}$ .

□ Beh 3



Insgesamt gilt:

$$\begin{aligned}
|Q_{\Delta}(D)| &= |Q_{\Delta, \text{leicht}}(D)| + |Q_{\Delta, \text{schwer}}(D)| \\
&\stackrel{\text{Beh 1\&2}}{\leq} (N^D)^{3/2} + |S^D| \cdot N^D \\
&\stackrel{\text{Beh 3}}{\leq} (N^D)^{3/2} + \sqrt{N^D} \cdot N^D \\
&= (N^D)^{3/2} + (N^D)^{3/2} \\
&= 2 \cdot (N^D)^{3/2}
\end{aligned}$$

Unser Algorithmus zur Berechnung von  $Q_{\Delta}(D)$  geht bei

Eingabe von  $D$  wie folgt vor:

0) Initialisiere  $Q_{\Delta}(D) := \emptyset$

1) Berechne für jedes  $v \in \text{atom}(D)$  folgendes:

die Mengen  $E^D(v, *) = \{w : (v, w) \in E^D\}$  und  
 $E^D(*, v) = \{w : (w, v) \in E^D\}$  und

die Zahl  $\text{aus-Grad}^D(v)$ .

Durch geeignetes Sortieren von  $E^D$  geht das in  
Zeit  $O(N^D \cdot \log N^D) \leq O((N^D)^{3/2})$

2) Berechne die Menge  $S^D$  aller schweren  $v \in \text{atom}(D)$  und  
die Menge  $L^D := \text{atom}(D) \setminus S^D$  aller leichten  $v \in \text{atom}(D)$ .  
Unter Verwendung der in 1) gesammelten Infos geht das  
in Zeit  $O(N^D)$ .

3) Betrachte jedes leichte  $a \in L^D$

Betrachte jedes  $b \in E^D(a, *)$

Betrachte jedes  $c \in E^D(*, a)$

und teste, ob  $(b, c) \in E^D$ .

Wenn ja, füge  $(a, (b, c))$  in  $Q_{\Delta}(D)$  ein



Auf die gleiche Art wie im Beweis von Beh 1 erhalten wir, dass  
das in Zeit  $O((N^D)^{3/2} \cdot t)$  geht.

4) Betrachte jedes schwere  $a \in S^D$   
 Betrachte jedes  $(b, c) \in E^D$   
 und teste, ob  $(a, b) \in E^D$   
 und  $(c, a) \in E^D$ .



Wenn ja, füge  $(a, b, c)$  in  $Q_\Delta(D)$  ein.

Das geht in Zeit  $|S^D| \cdot N^D \cdot 2t = O((N^D)^{3/2} \cdot t)$ .

Insgesamt berechnen wir so  $Q_\Delta(D)$  in Zeit  $O((N^D)^{3/2} \cdot t)$ .

□ Satz  $\Delta$

Bemerkung  $\Delta$ :

Der Algorithmus aus Satz  $\Delta$  ist im folgenden Sinn Worst-case-optimal:

Es gibt eine Folge von  $\sigma$ -Datenbanken  $D_1, D_2, D_3, \dots$   
 so dass für  $N_i := |E^{D_i}|$  für  $i \in \mathbb{N}_{\geq 1}$  gilt:

$$N_1 < N_2 < N_3 < \dots \text{ und}$$

$$|Q_\Delta(D_i)| = (N_i)^{3/2}.$$

Dazu wähle für  $i \geq 1$  die DB  $D_i$  mit  $E^{D_i} = [1, i] \times [1, i]$ .

Dann ist  $N_i = i^2$  und  $Q_\Delta(D_i) = [1, i]^3$ , also

$$|Q_\Delta(D_i)| = i^3 = (N_i)^{3/2}.$$

## Literaturhinweis:

Der hier präsentierte Beweis von Satz  $\Delta$  ist aus der folgenden Arbeit entnommen:

"Worst-case optimal join algorithms" von  
 Ngo, Pocar, Ré, Rudra,  
 In Proc. PODS 2012, pp. 37-48

Die Beweisidee wird dort Loomis und Whitney zugeschrieben ("An inequality related to the isoperimetric inequality" von Loomis, Whitney. In Bull. Amer. Math. Soc., 55: 361-362, 1949).