

Adatbázisok elmélete 16. előadás

Katona Gyula Y.
Budapesti Műszaki és Gazdaságtudományi Egyetem
Számítástudományi Tsz.
I. B. 137/b
kiskat@cs.bme.hu
<http://www.cs.bme.hu/~kiskat>

2004

Fizikai végrehajtás

$R(X, Y) \bowtie S(Y, Z)$ végrehajtása:

- **Ha $B(S) < M - 1$, azaz S belefér a memóriába: egymenetes algoritmus**
 1. Beolvassuk S -et és hashtáblát vagy B-fát készítünk, ahol a kulcs Y attribútumai.
 2. Beolvasunk egy blokkot R -ből. Minden sorára kikeressük a passzoló S -beli sorokat. Az eredményt kiírjuk.
I/O műveletigény: $B(S) + B(R)$
- **Ha $B(R) > B(S) > M - 1$: beágyazott ciklusú algoritmus**

Beolvasunk minél több blokkot a memóriába S -ből, utána ugyanazt csináljuk mint fenn.
I/O műveletigény: $B(S)B(R)/M$
- **Ha $B(R), B(S) \leq M^2$: rendezéses algoritmus**

Y kulcs szerint rendezzük R -et és S -et összefésüléses rendezéssel. Vesszük az összes y kulcsú sort a két lista elejéről és kiírjuk az összes párt. (Feltettük, hogy az összes y kulcsú sor elfér a memóriában.)
I/O műveletigény: $5(B(S) + B(R))$

Fizikai végrehajtás

Leginkább az I/O műveletigény érdekes. Ha „túl nagy” a számítási igény az is baj lehet.

Soronkénti, unáris műveletek: Kiválasztás és vetítés. Egyszerre csak egy sort kell vizsgálni, az algoritmus nem függ a memória nagyságától.

Unáris, teljes relációs műveletek: Pl. $\delta(R), \gamma(R)$. Ha nem fér el a reláció a memóriában, akkor mást kell csinálni.

Bináris, teljes relációs műveletek: $\cup, \cap, \setminus, \times, \bowtie$. Sok minden függ a méretektől.

Fizikai végrehajtás

- **Ha $\min(B(R), B(S)) \leq M^2$: hasheléses algoritmus**

Y kulcs szerint vödörös hashelést végzünk R -re és S -re. (Ha közben megtelik egy vödör, azt kiírjuk.) A kapott R_i, S_i vödrökkel egymenetes algoritmust végzünk.
I/O műveletigény: $3(B(S) + B(R))$
- **Ha van index S -re Y szerint: indexet használó algoritmus**

R -et blokkonként olvassuk be, az index alapján keressük ki a hozzá passzoló sorokat.
Átlagos I/O műveletigény: $B(S)B(R)/V(S, Y)$, ahol $V(S, Y)$: Y értékészletének számossága S -ben.

Optimalizálás

Triviális egyszerűsítések (főleg generált lekérdezések esetén hasznos):

- $r \cap r = r$; $r \bowtie r = r$; $r \cup \emptyset = r$; $\sigma_C(\emptyset) = \emptyset$
- $\pi_X(r \cup s) = \pi_X(r) \cup \pi_X(s)$
- $\sigma_{A=B \wedge B=C \wedge A=C}(r) = \sigma_{A=B \wedge B=C}(r)$

Optimalizálás

Milyen sorrendben érdemes kiszámolni $r(A, B) \bowtie s(B, C) \bowtie q(C, D)$ -t?

Szélsőséges esetben lehet, hogy bár r, s, q mindegyikének **1000** sora van, de $r \bowtie s$ -nek csak 1 sora, és $s \bowtie q$ -nak **1000000** sora.
 \Rightarrow Sokkal gyorsabb $(r(A, B) \bowtie s(B, C)) \bowtie q(C, D)$ kiszámolása.

Ezt persze előre nem lehet tudni biztosan. \Rightarrow Statisztikákat vezetünk a relációk attribútumaiban előforduló értékekről
 Ebből lehet becsülni a költségeket + dinamikus programozás vagy mohó algoritmus.

Igazi optimumot nehéz megtalálni:

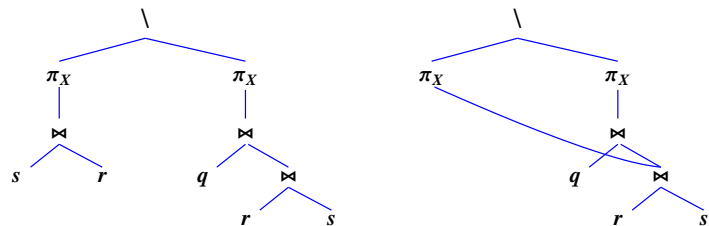
Tétel. *Annak eldöntése NP-teljes, hogy néhány reláció természetes illesztésének van-e legalább egy sora.*

Tétel. *Az optimum megtalálása NP-nehéz probléma.*

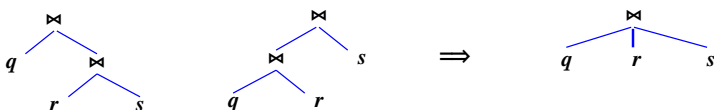
Optimalizálás

Ami többször előfordul, nem biztos, hogy érdemes mindig kiszámolni:

Pl. $\pi_X(s \bowtie r) \setminus \pi_X(q \bowtie r \bowtie s)$



Asszociativitás kihasználható:



Bizonyítás

Tétel. *Annak eldöntése NP-nehéz, hogy néhány reláció természetes illesztésének van-e legalább egy sora.*

Bizonyítás: Visszavezetjük rá a **3-SZÍN** problémát. Adott egy gráf, kérdés színezhető-e 3 színnel.

A gráf minden $e = \{x, y\}$ éléhez vegyünk fel egy-egy relációt:

e	X	Y	e'	X	Z
	piros	kék		piros	kék
	piros	sárga		piros	sárga
	kék	piros		kék	piros
	kék	sárga		kék	sárga
	sárga	piros		sárga	piros
	sárga	kék		sárga	kék

Ha a természetes illesztésnek van sora \Rightarrow egy sor minden csúcshoz rendel egy színt. Mivel az illesztés megfelel az egyes relációknak, egy él két végpontján nem lesz ugyanolyan szín.

Ha van színezés \Rightarrow a színezésben minden élre vegyük ki a megfelelő színpárt. Ezek a sorok összeillenek. \checkmark

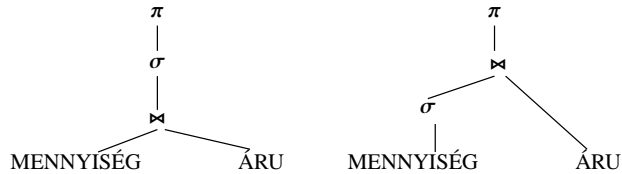
Kiválasztás tologatása

ÁRU(ÁRUKÓD, ÁRUNÉV, EGYSÉGÁR)
MENNYISÉG(DÁTUM, ÁRUKÓD, DB)

Hány darabot adtak el 2002. jan. 15-én az A123 kódú áruból, mi a neve és az ára?

$\pi_{DB, \text{ÁRUNÉV}, \text{EGYSÉGÁR}} (\sigma_{\text{ÁRUKÓD}='A123' \wedge \text{DÁTUM}='2002-01-15'} (\text{MENNYISÉG} \bowtie \text{ÁRU})) \Rightarrow$

$\pi_{DB, \text{ÁRUNÉV}, \text{EGYSÉGÁR}} (\sigma_{\text{ÁRUKÓD}='A123' \wedge \text{DÁTUM}='2002-01-15'} (\text{MENNYISÉG} \bowtie \text{ÁRU}))$



Felhasznált azonosság:

$\sigma_C(\mathbf{R} \bowtie \mathbf{S}) = \sigma_C(\mathbf{R}) \bowtie \mathbf{S}$, ha minden C -beli attribútum szerepel \mathbf{R} -ben.

Hasonló azonosságok:

$\sigma_C(\mathbf{R} \cup \mathbf{S}) = \sigma_C(\mathbf{R}) \cup \sigma_C(\mathbf{S})$,

$\sigma_C(\mathbf{R} \times \mathbf{S}) = \sigma_C(\mathbf{R}) \times \mathbf{S}$, ha minden C -beli attribútum szerepel \mathbf{R} -ben.

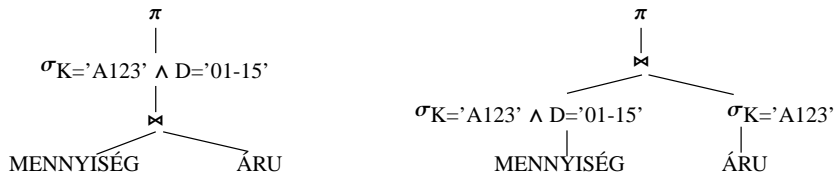
$\sigma_C(\mathbf{R} \bowtie \mathbf{S}) = \sigma_C(\mathbf{R}) \bowtie \sigma_C(\mathbf{S})$, ha minden C -beli attribútum szerepel \mathbf{R} -ben és \mathbf{S} -ben is.

Kiválasztás tologatása

Összetett C szétszedhető:

$\sigma_{C_1 \wedge C_2}(\mathbf{R}) = \sigma_{C_1}(\sigma_{C_2}(\mathbf{R}))$,

$\sigma_{C_1 \vee C_2}(\mathbf{R}) = \sigma_{C_1}(\mathbf{R}) \cup_H \sigma_{C_2}(\mathbf{R})$, ha \mathbf{R} nem multihalmaz.



Lehet, hogy érdemes előbb feltolni, aztán le.

Más műveletekre vonatkozó szabályok

Hasonló szabályok π, δ, γ -ra is vannak, de ezek nem annyira csökkentik a műveletigényt. De csak olyan attribútumot lehet eltüntetni, amire nem hivatkozunk feljebb.

- $\pi_L(\mathbf{R} \bowtie \mathbf{S}) = \pi_L(\pi_M(\mathbf{R}) \bowtie \pi_N(\mathbf{S}))$, ahol M az \mathbf{R} olyan attribútumai, hogy vagy összekapcsolási attribútum, vagy L -beli, N pedig ...
- $\delta(\mathbf{R} \bowtie \mathbf{S}) = \delta(\mathbf{R}) \bowtie \delta(\mathbf{S})$
- $\delta(\sigma_C(\mathbf{R})) = \sigma_C(\delta(\mathbf{R}))$
- $\delta(\gamma_L(\mathbf{R})) = \gamma_L(\mathbf{R})$

De pl. δ nem tolató át \cup_M, π_n .

Még egy fontos kérdés az alkérdések kezelése, de erről most nem szólnak.

Összefoglalás

Sok ilyen szabály alkalmazásával több féle logikai terv előállítható.

Ezeknek megbecsüljük a költségét és választunk egyet. Ehhez készítünk fizikai tervet.

Jobb rendszerekben ez automatikus. Ilyenekben kérdéses, hogy a lekérdezés beírásakor mire kell figyelni.

Inkább olyan egyszerűsítéseket érdemes csak elvégezni, ami a függések következménye, mert ezeket nehezebben lehet automatizálni.

Pl. Oracle-ban van rá mód, hogy megnézzük mi a logikai és fizikai terv és meg lehet adni, hogy pontosan mit csináljon.

Fizikai szervezés, tárkezelés

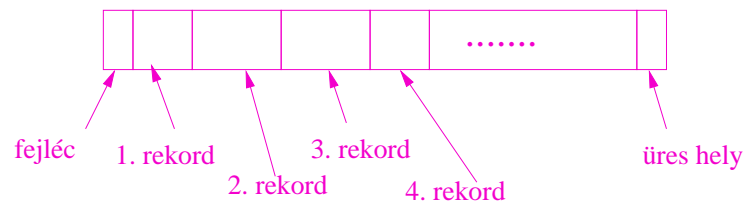
Célja: a rekordokból (egy rekord = a reláció egy sora) álló állomány kezelése úgy, hogy az adatokhoz való hozzáférés gyors legyen.

Fontos jellemzők:

- **külső táras adatkezelés**, mert sok az adat \Rightarrow ha valamivel dolgozni akarunk, akkor be kell hozni a belső memóriába \Rightarrow a költséget a beolvasás/kiírás jelenti \Rightarrow az I/O műveletek számára akarunk optimalizálni
- a műveletek, amiket gyorsan meg kell tudni csinálni: **rekordok beillesztése, törlése, módosítása, keresése**

Blokkokról általában

Tipikus blokk

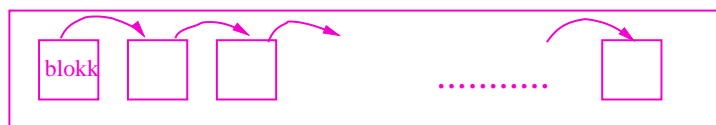


A fejléc tartalmazza a blokkra vonatkozó infókat, (pl: melyik relációhoz tartozik, mennyi a szabad hely benne, hol kezdődik); ezután jönnek a rekordok egymás után, a végén általában marad üres hely.

Fontos feltevés: rekordok blokkhatárt nem lépnek át, ezért általában van üres, fel nem használható hely a blokkok végén. (Ha nagyok a rekordok, pl. képfile-ok, és mégis át kell lépni laphatárt, akkor extra technikák kellene, de ezzel most nem foglalkozunk.)

Az állomány felépítése

Az adatállomány a külső táron van, blokkok (lapok) elérésfolytonos sorozatán.

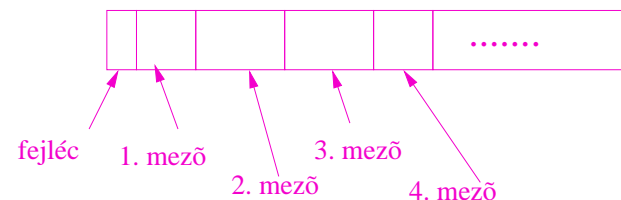


- egyszerre egy blokk írható ki/olvasható be
- blokk méret fix (ált. 2^{10} , 2^{12} byte)
- az operációs rendszer tartja nyilván, hogy melyik reláció rekordjai hol vannak és ő biztosítja az elérésfolytonosságot is

Rekordok típusai

Kötött formátum

Ekkor a mezők száma, mérete, típusa és sorrendje fix



Fejléc:

- a rekord kezelésével kapcsolatos infók: törölt-e, melyik relációhoz tartozik
- a mezők típusa
- időbélyeg (mikor módosult utoljára)

Rekordok típusai

Változó formátum

- Mezők hossza esetleg nem fix (szöveget tartalmazó adatbázisok)
- Ismétlődő mezők lehetnek, és az ismétlések száma nem fix vagy pedig többértékű mezők (szereplők felsorolása filmnél)

Ilyenkor bonyolultabb a fejléc, kevésbé lehet gazdaságosan/előre tervezhetően tárolni a rekordokat, ezért érijük el inkább, hogy ne legyen ez az eset:

Vezessük vissza ezt az esetet a kötöttre, pl. mutatók alkalmazásával a problémás helyeken
⇒ Mostantól feltesszük, hogy a rekordok kötött formátumúak és hogy az egész állományon belül ugyanaz a formátum van.

Fontos fogalmak

1. **mutató:** blokk vagy rekord címét tartalmazó bejegyzés
2. **kötött blokk/rekord:** mutathat rá mutató, ezért nem mozgatható el szabadon. Ez típusszinten adott, azaz ha egy reláció rekordjaira/blokkjaira mutathat mutató, akkor még akkor is kötöttnek számít, ha éppen nem mutat egyre se semmi.
3. **szabad blokk/rekord:** nem mutathat rá mutató
4. **Kulcs, keresési kulcs (néha csak kulcsnak hívjuk):**
 - a rekordok mezőinek egy kitüntetett halmaza (a reláció attribútumainak egy részhalmaza)
 - ez alapján megy a keresés (ezeknek az értékét adjuk meg és azokat a rekordokat (sorokat a relációban) keressük, amiknél pont ezek az értékek szerepelnek)
 - a keresési kulcs nem egyezik meg feltétlenül a reláció egyik kulcsával sem (pl. név a telefonkönyvnel)
 - de az azért elvárás, hogy ne legyen nagyon sok egy-egy értékre illeszkedő rekord