

Vyhodnocování dotazů

slajdy k přednášce NDBI001

Jaroslav Pokorný
MFF UK, Praha
pokorny@ksi.mff.cuni.cz

Časová a prostorová složitost

- Jako dlouho trvá dotaz?
 - CPU (cena je malá; snižuje se; těžko odhadnutelná)
 - Disk (hlavní složka ceny, # I/O operací)
- Kolik n-tic je třeba přenést?
- Jaké statistiky je třeba udržovat?

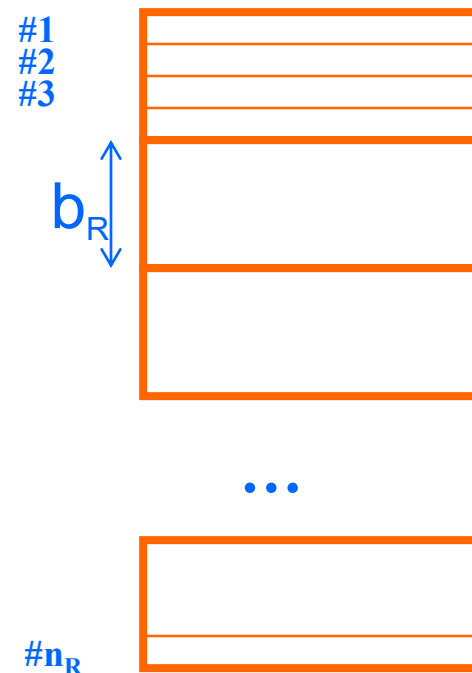
Statistiky

Statistiky pro každou relaci:

n_R	počet n-tic relace R
$V(A,R)$	počet prvků $R[A]$
p_R	počet stránek k uložení R
b_R	blokovací faktor
M	počet stránek volného prostoru v RAM
$l(A,R)$	počet úrovní indexového souboru pro A v R

Značení:

$buffer_R$	celistvý počet stránek pro R v RAM (neuvažujeme caching)
------------	--



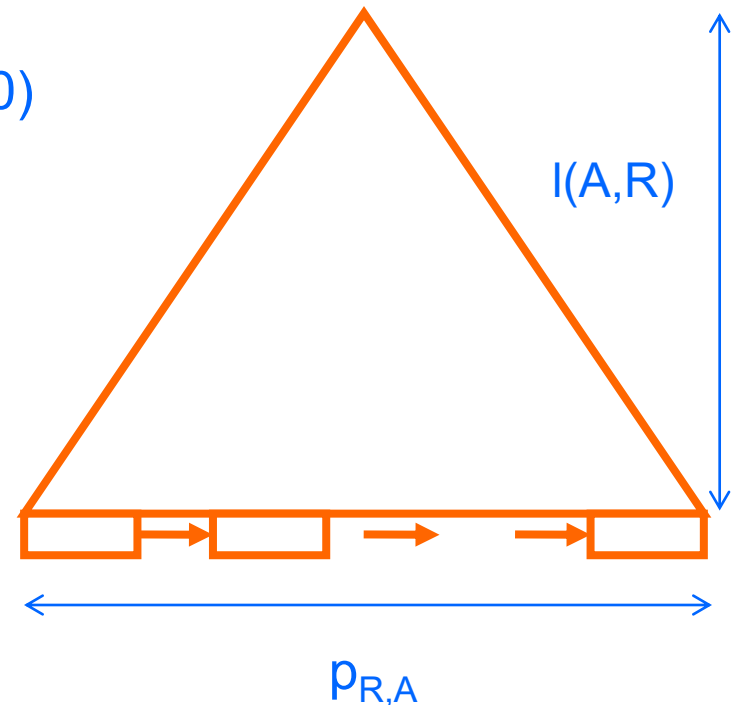
Indexace B⁺-stromy:

Vhodné: existuje-li na $\text{dom}(A)$ uspořádání.

Pro atribut A relace R :

- $f_{A,R}$: průměrný počet následníků
ve vnitřním uzlu ($\sim 50-100$)
- $I(A,R)$: # úrovní indexu ($\sim 2-3$)
– $\sim \log(V(A,R))/\log(f_{A,R})$
- $p_{R,A}$: # listových stránek

Značení: A místo $R.A$



Metody pro výpočet selekce

```
select *  
from R  
where A = 'a'
```

Případy:

- A je primární klíč,
- A je sekundární (alternativní) klíč
- k A existuje index - obyčejný nebo typu CLUSTER
- A je hašovaný klíč

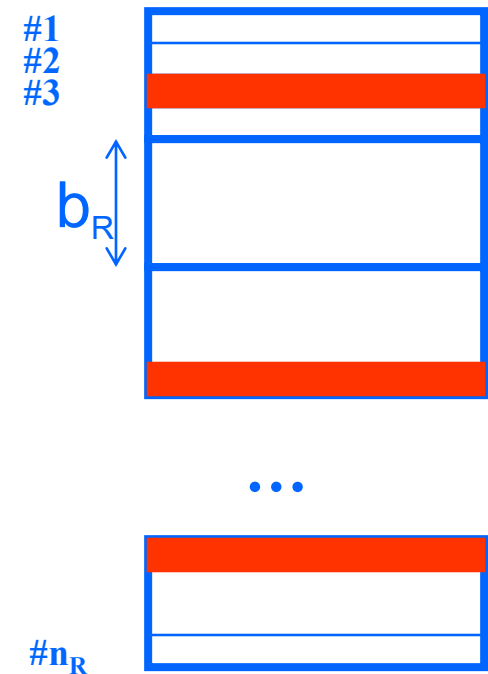
Předpoklady: rovnoměrné rozložení hodnot A v R[A]

$$\Rightarrow n_{R(A=a)} = n_R / V(A,R)$$

Metody pro výpočet selekce

Sekvenční vyhledávání

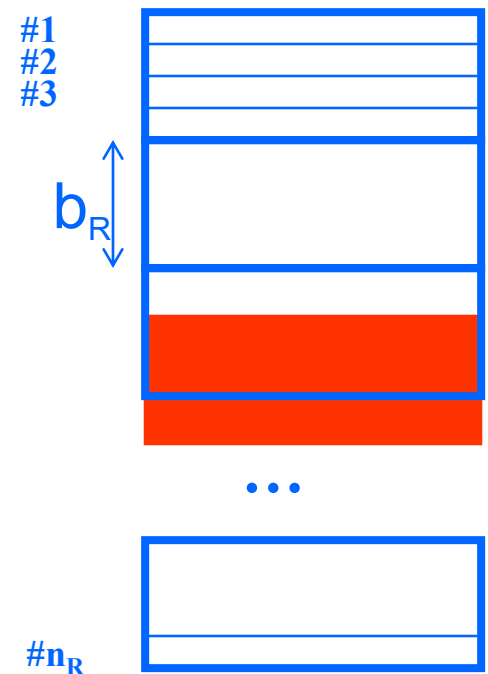
- p_R /*nejhorší případ*/
- $p_R/2$ /* průměrně, je-li A primární klíč*/



Metody pro výpočet selekce

Binární vyhledávání, je-li R
uspořádaná podle A

- $\log_2(p_R)$ /*je-li A primární klíč */
- $\log_2(p_R) + n_{R(A=a)} / b_R$ /*je-li A libovolný atribut */



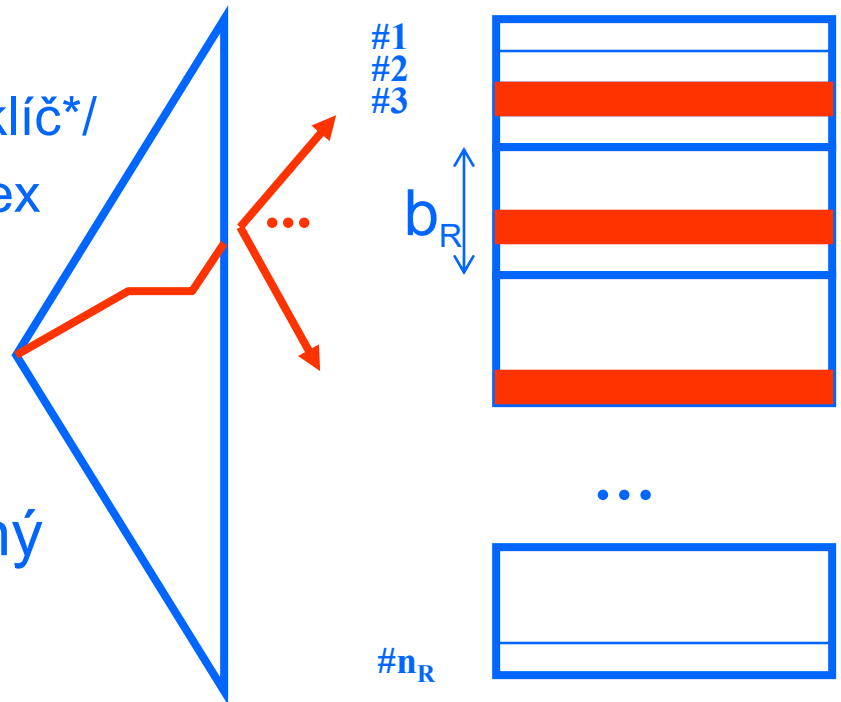
Metody pro výpočet selekce

Vyhledávání, existuje-li pro A index

- $I(A) + 1$ /*je-li A primární klíč*/
- $I(A) + n_{R(A=a)}/b_R$ /*je-li index pro A typu CLUSTER*/
- $I(A) + n_{R(A=a)}$ /*není-li index pro A typu CLUSTER*/

Vyhledávání, je-li A hašovaný atribut

- ≈ 1 přístup



Metody pro výpočet selekce

```
select *  
from R  
where A < 'a'
```

Sekvenční vyhledávání

- p_R
- $p_R(a - \min_A) / (\max_A - \min_A)$

/ nejhorší případ*/*

*/*je-li R setříděná dle A*/*

Vyhledávání, existuje-li index

- $I(A) + p_R / 2$
- $I(A) + p_{R,A} / 2 + n_R / 2$

*/*je-li R setříděná dle A*/*

/ je-li index pro A, A je sekundární klíč*/*

Příklad

Rezervace(č_zák, č_letu, datum, pozn)

$n_{\text{Rezervace}} = 10\ 000$

$b_{\text{Rezervace}} = 20$

$V(\text{č_zák}, \text{Rezervace}) = 500$

$V(\text{č_letu}, \text{Rezervace}) = 50$

$f_{\text{č_letu},R} = 20$

Dotaz: Najdi zákazníky s číslem letu = '77'

Příklad

Sekvenční hledání:

⇒ cena dotazu: 500 diskových přístupů

Klastrovaný index pro č_letu:

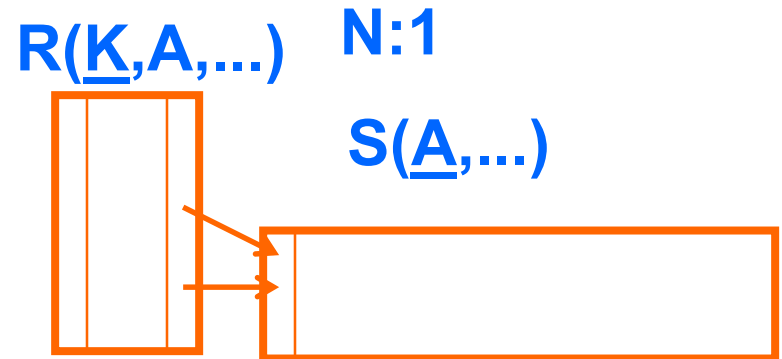
cena dotazu = $I(\text{č_letu}) + n_{R(\text{č_letu}=a)}/b_R$

- $I(A)$: 50 hodnot $f_A = 20 \Rightarrow I(A)=2$
Zdůvodnění: $(\log(50)/\log(20)) \approx 2$
- $n_{R(A=a)} = n_R/V(A,r) = 10,000/50 = 200$ n-tic
 $n_{R(A=a)}/b_R = 200/20 = 10$ stránek
⇒ cena dotazu = $2+10= 12$

Metody pro výpočet spojení

Dva typy implementace:

- nezávislé relace
- s ukazateli
(Starbrust, Winbase,...)



Základní metody:

- hnížděné cykly (varianty s indexováním, vyhledáváním)
- setřídění-slévání
- hašované spojení

Předpoklady: spojovací atribut A , $p_R \leq p_S$,
u varianty s ukazateli $R \rightarrow S$

Pz.: speciální případ - kartézský součin

Hnízděné cykly - binárně

- naivní algoritmus

for each $r \in R$

for each $s \in S$

if $\Theta(r,s)$ then begin $u := r [\Theta] s$; WRITE(u) end

...

- po stránkách

menší relaci jako vnější!

$M=3 \Rightarrow p_R + p_R p_S$ čtení

$(n_R n_S / V(A,S)) / b_{RS}$ zápisů (zdůvodni !)

Vylepšení: - vnitřní relace se čte \longleftrightarrow

ušetří se 1 čtení na začátku (konci)

Hnízděné cykly - binárně

Varianty:

- M velké, pak rozdělení: $M-2$, 1 , 1
vnější vnitřní výsledek
 $\Rightarrow p_R + p_S p_R / (M-2)$ čtení
- R ve vnitřní paměti
 $\Rightarrow p_R + p_S$ čtení
- s ukazateli, $M=3$
 $\Rightarrow p_R + n_R$ čtení

Hnízděné cykly - binárně

- index na S.A (B⁺-strom)

Předpoklady: R setříděná v R.A, S.A primární

$$\Rightarrow p_R + I(A,S) + p_{S,A} + V(A,R) \quad \text{čtení}$$

- S hašovaná dle S.A

Předpoklady: R setříděná v R.A, S.A primární

$$\Rightarrow p_R + V(A,R) \quad \text{čtení}$$

- se selekcí (pomocí vyhledávání),

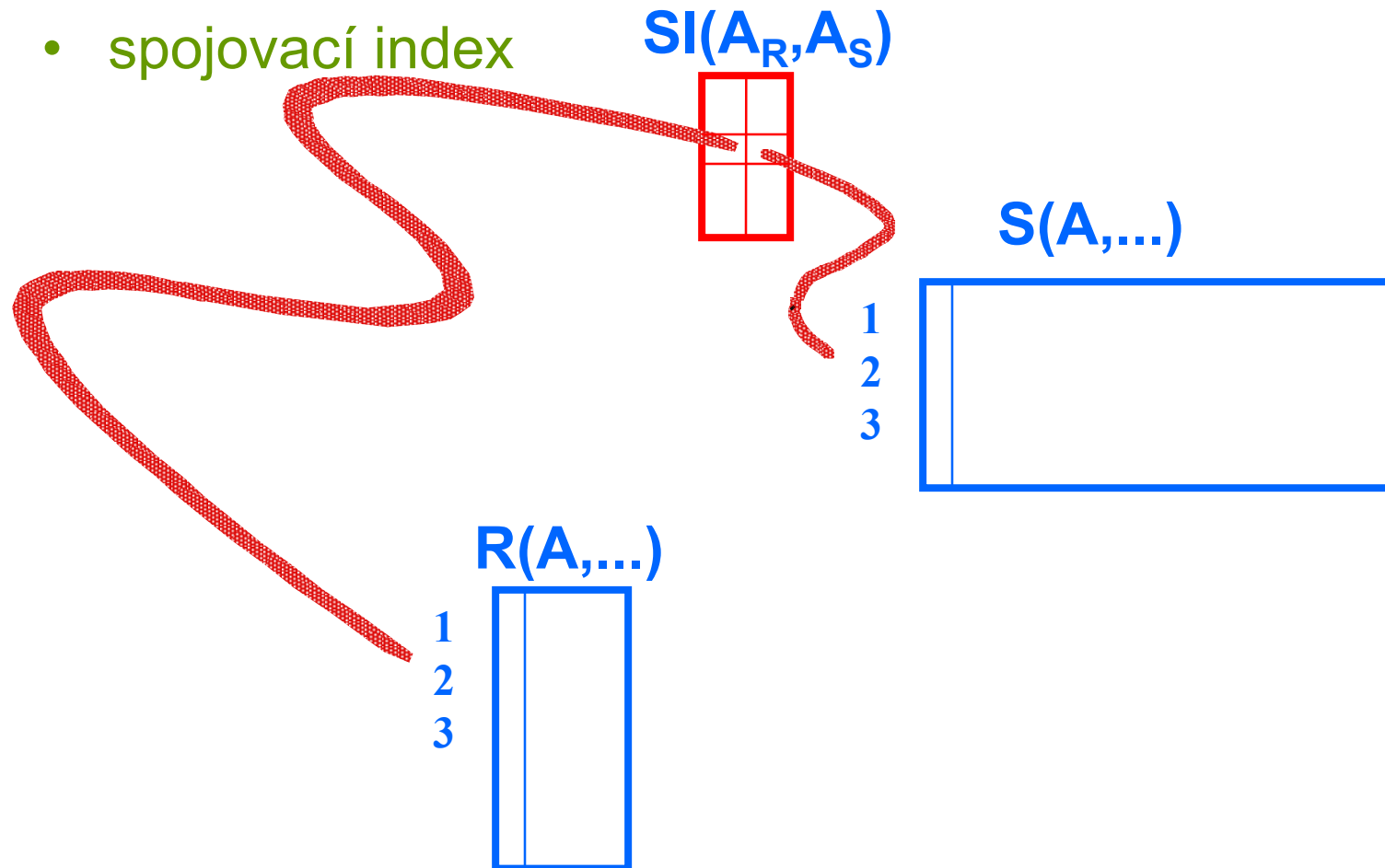
Př.: `SELECT * FROM R,S`
`WHERE R.A=S.A AND R.B=12`

Předpoklady: R.B primární klíč (indexovaný), S.A sekundární klíč (klastr. Index, n-tice s S.A=a v jedné stránce)

$$\Rightarrow I(A,S) + I(B,R) + 2 \quad \text{čtení}$$

Hnízděné cykly - binárně

- spojovací index



Hnízděné cykly - více relací

$$M' = M_1 + M_2 + \dots + M_n < M$$

R_i rozděleny na l_i podrelací velikostí M_i , tj. $l_i = \lceil p_i/M_i \rceil$,
($1 \leq i \leq n$)

cenová funkce [Kim84]

$$C = p_1 + [M_2 + (p_2 - M_2) \lceil p_1/M_1 \rceil] + \dots + [M_n + (p_n - M_n) \lceil p_1/M_1 \rceil \dots \lceil p_{n-1}/M_{n-1} \rceil]$$

\Rightarrow problém hledání celočíselných M_i , aby C minimální

heuristika:

- (1) Srovnej n relací do cyklů algoritmu tak, že $p_1 \leq p_2 \leq \dots \leq p_n$;
- (2) Pro R_n alokuj 1 stránku, $M' - 1$ rozděl rovnoměrně;
- (3) $(M' - 1)/(n-1)$ necelé **then** přiděl větší M_i menším relacím;

Hnízděné cykly - více relací

Struktura základního algoritmu (zde pro tři relace):

```
for j:=1 to L1 do
  begin přečti R1j do bufferM1;
  for k:=1 to L2 do
    begin přečti R2k do bufferM2;
    for s:=1 to L3 do
      begin přečti R3s do bufferM3;
      vytvoř spojení bufferMi, 1 ≤ i ≤ 3;
      zapiš výsledek
    end
  end
end
```

Hnízděné cykly - více relací

Př.:

a) $p_1 = 7, p_2 = 14, p_3 = 21, M' = 11$

\Rightarrow dělení $M' = \langle 5, 5, 1 \rangle$

b) $p_1 \leq \dots \leq p_5, M' = 16$

\Rightarrow dělení $M' = \langle 4, 4, 4, 3, 1 \rangle$

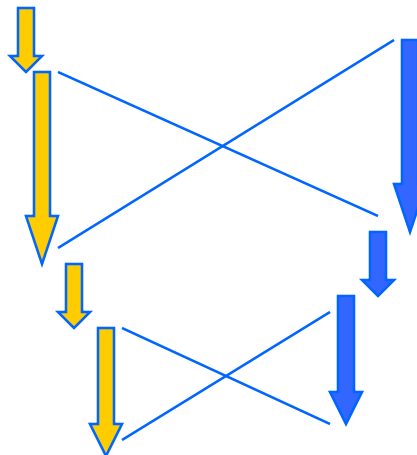
Setřídění-slévání

Idea: třídění, slévání (dvoufázový algoritmus)

Vhodnost: jsou-li R a S setříděné

R(...,A,...)

0001
0002
0002
0002
0004
0005
0005



S(...,A,...)

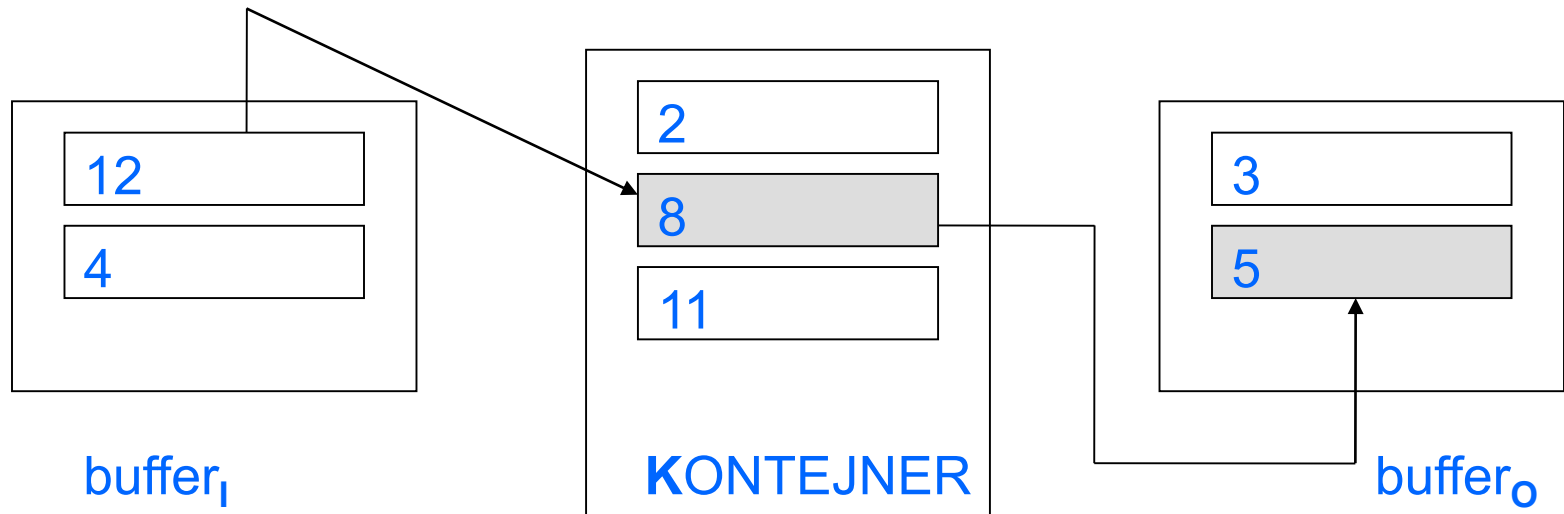
0002
0002
0002
0003
0005
0005

- min. $M = 3$, 2.fáze vyžaduje $p_R + p_S$ čtení
- potřebuje pomocný prostor pro třídění
- **výsledek je setříděný**

Setřídění-slévání

- $M = 3$ (s použitím externího třídění)
 - $\Rightarrow \sim 2p_R \log(p_R) + 2p_S \log(p_S) + p_R + p_S$ bez zápisu výsledku
- $M \geq \sqrt{p_S}$ (dvoufázový algoritmus)
 - (1) Vytvářejí se setříděné běhy velikosti $2M$ stránek (pomocí prioritní fronty) a ukládají na disk;
 - \Rightarrow velikost běhu je $\geq 2\sqrt{p_S}$
 - \Rightarrow pro S jich je nejvýše $p_S/2\sqrt{p_S}$, pro R také ne více než $p_S/2\sqrt{p_S}$
 - \Rightarrow celkem nejvýše $\sqrt{p_S}$
 - (2) Pro každý běh se alokuje v paměti stránka a souběžně se slévá;
 - $\Rightarrow 3(p_R + p_S)$ bez zápisu výsledku

Princip prioritní fronty



1. Naplní se K a $buffer_1$ n-ticemi z R.
2. Z K se vybírají n-tice u takové, že $u.A \geq v.A$, $\forall v$ ve $buffer_0$ a řadí se ve vzestupném pořadí podle hodnoty A . (1)
3. Uvolněné místo v K se zaplní novou n-ticí ze $buffer_1$. Je-li $buffer_1 = \emptyset$, načte se nová stránka R. Je-li $buffer_0$ plný, na vnější paměti se prodlouží daný běh. Jestliže žádná n-tice z kontejneru nevyhovuje (1), je současný stav $buffer_0$ poslední stránkou vytvářeného běhu.

Tímto způsobem lze vytvořit běhy dlouhé v průměru $2M$ stránek.

Setřídění-slévání

- varianta s ukazateli

R se setřídí podle ukazatelů

S se čte pouze jednou, nemusí být setříděna

$\Rightarrow 3p_R + p_S$ bez zápisu výsledku

Porovnání:

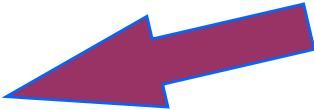
- $|p_R - p_S|$ velké \Rightarrow lepší jsou hnížděné cykly
- $|p_R - p_S|$ malé \Rightarrow lepší je setřídění-slévání
- omezující selekce \Rightarrow lepší pomocí vyhledávání

Hašovaná spojení

Vhodnost:

nejsou-li dostupné indexy pro R.A a S.A

nemusí-li být výsledek setříděn dle A

- klasické hašování
- GRACE algoritmus 
- jednoduché hašování
- rekurzivní rozdělení relací
- hybridní hašování

Spojení klasickým hašováním

Předpoklad: R se vejde do M stránek

$M = p_R * F + 1 + 1$, kde F je koeficient větší než 1

(1) Zahašuj R do vnitřní paměti;

(2) Čti S sekvenčně;

Hašuj s.A a přímým přístupem najdi $r \in R$;

(3) if s.A = r.A then begin u:= r * s; WRITE(u) end

$\Rightarrow p_R + p_S$ čtení

Spojení s rozdělováním relací

Předpoklad: R se nevejde do M stránek

Idea: R a S se rozdělí na disjunktní podmnožiny tak, že se spojují *JEN* ty korespondující.

Dvoufázový algoritmus:

(1) Rozděl R a S;

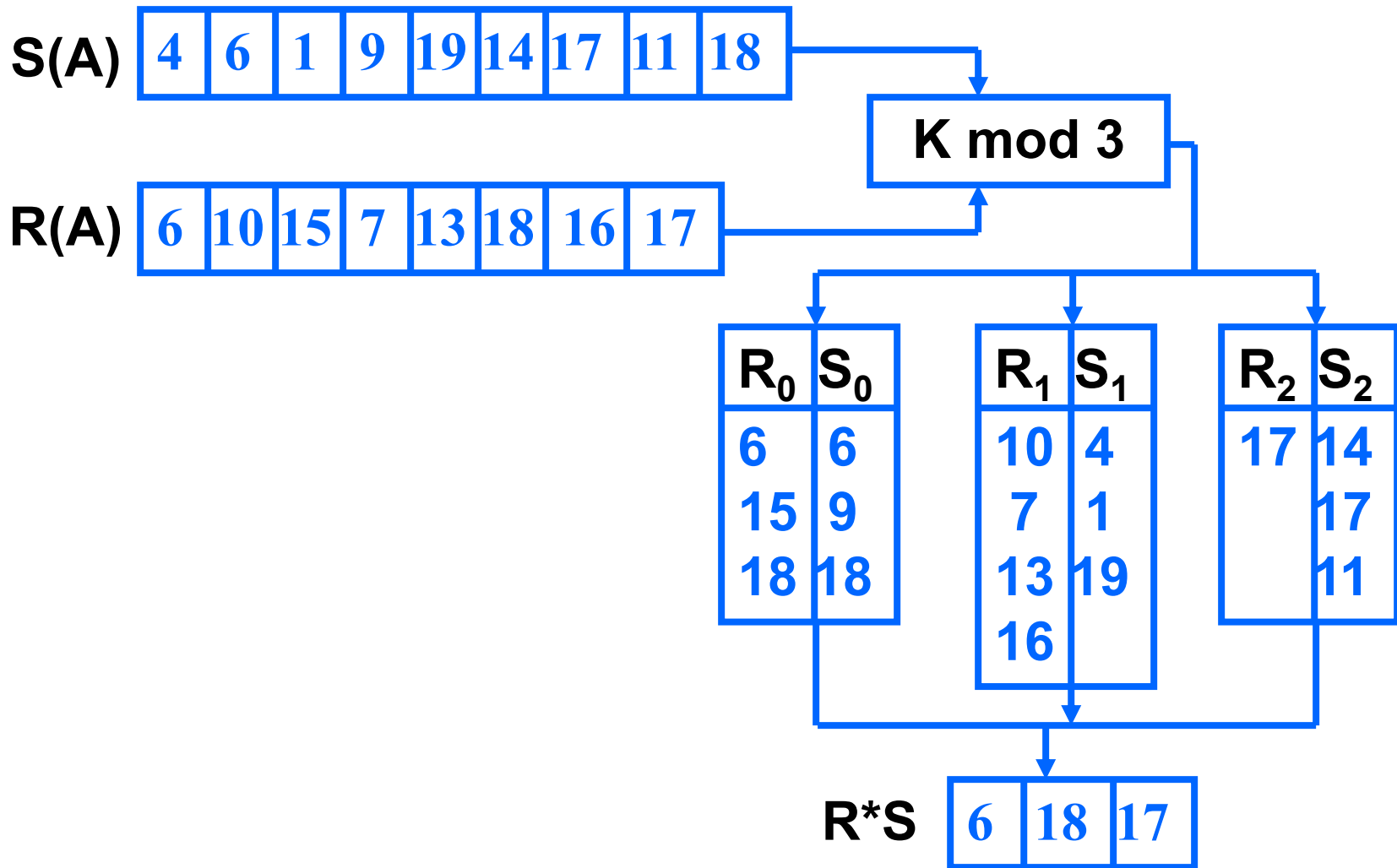
(2) Zahašuj část R (části R) do M-2 stránek;

Čti odpovídající část S;

Hašuj s.A a přímým přístupem najdi n-tice $r \in R$;

Generuj výsledek;

Příklad:



GRACE algoritmus

- „školní“ verze

Datové struktury: n -tice R a S , kapsy ukazatelů HR_i , HS_i ,
 $i \in \{0, 1, \dots, m-1\}$

hašovací funkce $h: \text{dom}(A) \rightarrow \langle 0, m-1 \rangle$

Algoritmus:

for $k:=1$ to n_R do begin $i := h(R[k].A)$; $HR_i := HR_i \cup \{k\}$ end

for $k:=1$ to n_S do begin $i := h(S[k].A)$; $HS_i := HS_i \cup \{k\}$ end

for $i:=0$ to $m-1$ do

begin $POM_R := \emptyset$; $POM_S := \emptyset$;

foreach $j \in HR_i$ do begin $r:=R[j]$; $POM_R := POM_R \cup \{r\}$ end;

foreach $j \in HS_i$ do begin $s:=S[j]$; $POM_S := POM_S \cup \{s\}$ end;

GRACE algoritmus

```
foreach s ∈ POMS do                                     /* ve vnitřní paměti */
  begin foreach r ∈ POMR do
    begin if s.A = r.A then begin u:= r * s;
                                          WRITE(u)
                                          end
    end
  end
end
end
```

$\Rightarrow p_R + p_S + n_R + n_S$ čtení

vhodnost: když $p_R/m + p_S/m < M$

GRACE s ukládáním rozdělených relací

- $M \geq \sqrt{(p_R * F)}$
 - (1) Zvol h tak, že R lze rozdělit do $m = \sqrt{(p_R * F)}$ částí;
 - (2) Čti R a hašuj do (výstupních) $buffer_i$, ($0 \leq i \leq m-1$);
if $buffer_i$ je plný then $WRITE(buffer_i)$;
 - (3) Proveď (2) pro S ;
 - (4) for $i := 0$ to $m-1$ do begin
 - (4.1) Čti R_i a hašuj do paměti velikosti $\sqrt{(p_R * F)}$;
 - (4.2) Čti $s \in S_i$ a hašuj $s.A$.
Existuje-li $r \in R_i$ a $s.A = r.A$, generuj výsledek.

GRACE s ukládáním rozdělených relací

Zdůvodnění 4.1: předpoklad - $R_i \approx$ stejné velké

$$p_R/m = p_R / \sqrt{(p_R * F)} = \sqrt{(p_R/F)}$$

$$R_i \text{ vyžaduje prostor } F \sqrt{(p_R/F)} = \sqrt{(p_R * F)}$$

$$\Rightarrow 3(p_R + p_S) \quad \text{I/O operací}$$

vhodnost: když $p_R/m + p_S/m < M$

poznámky:

- S_i mohou být libovolně velké. Vyžadují 1 stránku paměti;
- problém, když $V(A,R)$ je malý;
- vhodné v situacích, když $R(\underline{K}, \dots)$, $S(\underline{K}_1, K, \dots)$;
- nevejde-li se R_i resp. S_i do $M-2$ stránek \Rightarrow rekurze
tj. R_i se rozdělí na $R_{i0}, R_{i1}, \dots, R_{i(k-1)}$ stránek;

Jednoduché hašování

Předpoklad: $p_R * F > M-2$, A je UNIQUE

Idea: speciální případ GRACE, kdy $R \rightarrow R_1, R_2$

Algoritmus:

```
repeat begin zvol h; Čti R a hašuj r.A; M-2 bufferů  
          tvoří  $R_1$ , ostatní n-tice do  $R_2$  na disk;
```

```
          Čti S a hašuj s.A;
```

```
          if  $h(s.A)$  padne do prostoru  $R_1$ 
```

```
          then begin if  $s.A = r.A$  then generuj výsledek end
```

```
          else ulož s do  $S_2$  na disk;
```

```
           $R := R_2$ ;  $S := S_2$  end
```

```
until  $R_2 \neq \emptyset$ ;
```


Hybridní hašování

Idea: kombinace GRACE a jednoduchého hašování,
R se dělí na části $R_1, R_2, \dots, R_k, R_0$ tak, že R_0 se vejde do vnitřní paměti.

Rozdělení M-2 stránek: $|buffer_i| = 1$ ($1 \leq i \leq k$), $|buffer_0| = p_{R_0}$

Algoritmus:

(1) Zvol h ;

(2) Čti R a hašuj $r.A$; tvoř R_i ($0 \leq i \leq k$); /* R_0 je v $buffer_0$ */

(3) Čti S a hašuj $s.A$; tvoř S_i ($1 \leq i \leq k$);

if $h(s.A)$ padne do prostoru S_0 then realizuj spojení;

(4) for $i:=1$ to k do realizuj spojení podle GRACE;

Porovnání algoritmů

Předpoklady:

$M > \sqrt{p_S}$ pro setřídění-slévání

$M > \sqrt{p_R}$ pro hašování

Značení: $\text{alg1} \gg \text{alg2} \Leftrightarrow \text{alg1}$ je lepší než alg2

	Setřídění-slévání	GRACE	jednoduché hašování	hybridní hašování
GRACE	\gg		\gg (pro menší M)	
jednoduché hašování	\gg	\gg (pro větší M)		
hybridní hašování	\gg	\gg	\gg	

Dělení

Df.: R a S se schématy Ω_1 resp. $\Omega_2 \subset \Omega_1$

$$T = R \div S = R[\Omega_1 - \Omega_2] - ((R[\Omega_1 - \Omega_2] \times S) - R)[\Omega_1 - \Omega_2]$$

Př.:

R	A	B
	3	6
	3	2
	8	2
	1	2
	1	3
	3	4
	3	3

S	B
	2
	4
	3

R	A	B
	1	2
	1	3
	3	4
	3	3
	3	6
	3	2
	8	2

T	A
	3

po setřídění

Dělení-hašováním

Idea: vytvoří se kapsy HS_i pro hodnoty z $S.B$ a do nich se ukládají hodnoty z $R.A$. Hodnoty z $\cap HS_i$ přispívají do výsledku.

Algoritmus: (prvky hašovací tabulky jsou např. typu array nebo set, reprezentují kapsy)

(1) Čti S , spočti $h(s.B)$ a označ prostor (kapsu) $HS_{s.B}$

foreach $s.B$ do $HS_{s.B} := \emptyset$;

(2) for $j:=1$ to n_R do begin $r:=R[j]$;

if existuje kapsa pro $h(r.B)$

then $HS_{r.B} := HS_{r.B} \cup \{r.A\}$ end

(3) foreach $HS_{s.B}$ do $sort(HS_{s.B})$; /*není nutné*/

(4) Vytvoř $\cap HS_i$ a generuj T ;

Další operace

GROUP BY

- index pro A - přes index se získají skupiny
- setřídění podle R.A
- hašováním (jako u dělení)

foreach $a \in R[A]$ do vytvoř kapsu + proměnná pro
výpočet agregační funkce;

DISTINCT

také pomocí hašování