

CS 245: Database System Principles

Notes 6: Query Processing

Hector Garcia-Molina

Vyhodnocení dotazu

- Dotaz
- Strom dotazu
- Logický plán
- Úpravy
- Fyzický plán
- Vyhodnocení

Příklad

Select B,D

From R,S

Where $R.A = "c" \wedge S.E = 2 \wedge R.C = S.C$

R	A	B	C	S	C	D	E
a	1	10		10	x	2	
b	1	20		20	y	2	
c	2	10		30	z	2	
d	2	35		40	x	1	
e	3	45		50	y	3	

Výsledek: $\begin{array}{c|c} B & D \\ \hline 2 & x \end{array}$

- Jak vyhodnotit uvedený dotaz?

1. způsob

- kartézský součin
- výběr záznamů
- projekce

RXS

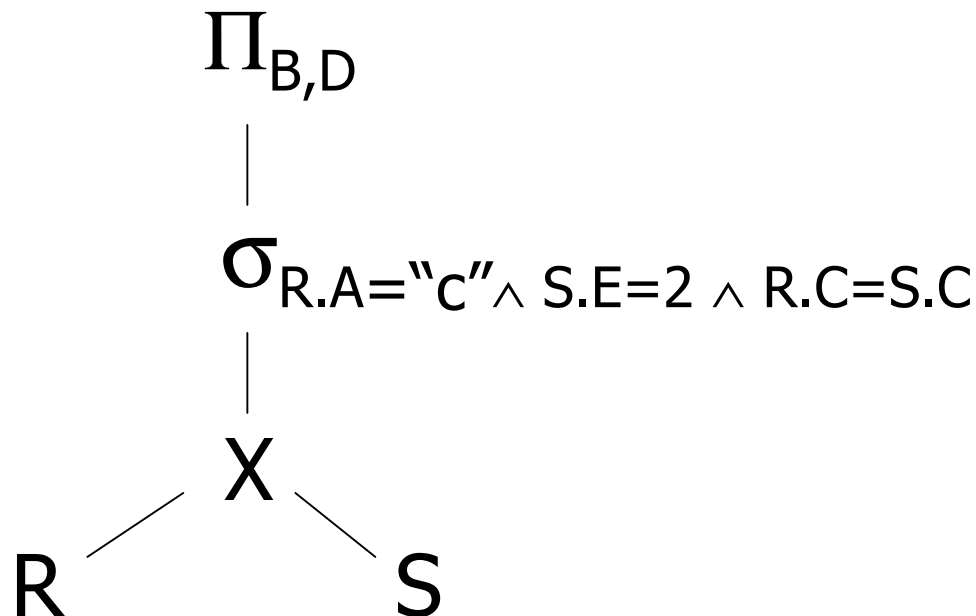
R.A	R.B	R.C	S.C	S.D	S.E
a	1	10	10	x	2
a	1	10	20	y	2
.					
.					
C	2	10	10	x	2
.					
.					

Bingo! →

Jeden máme...

Relační Algebra – použití pro popis plánů

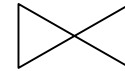
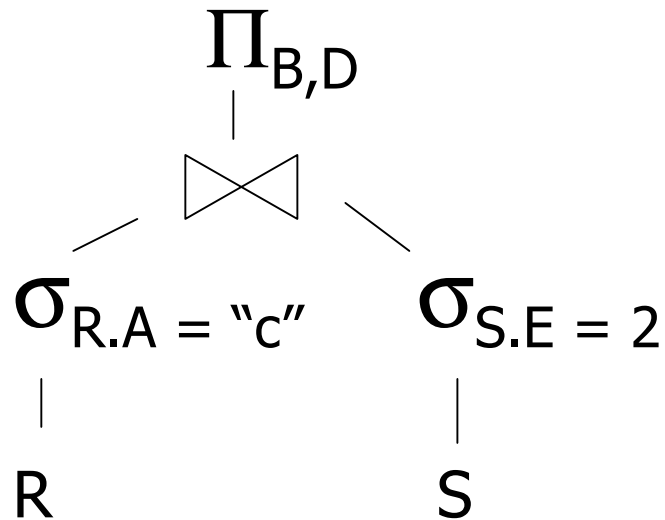
Př.: Plán I



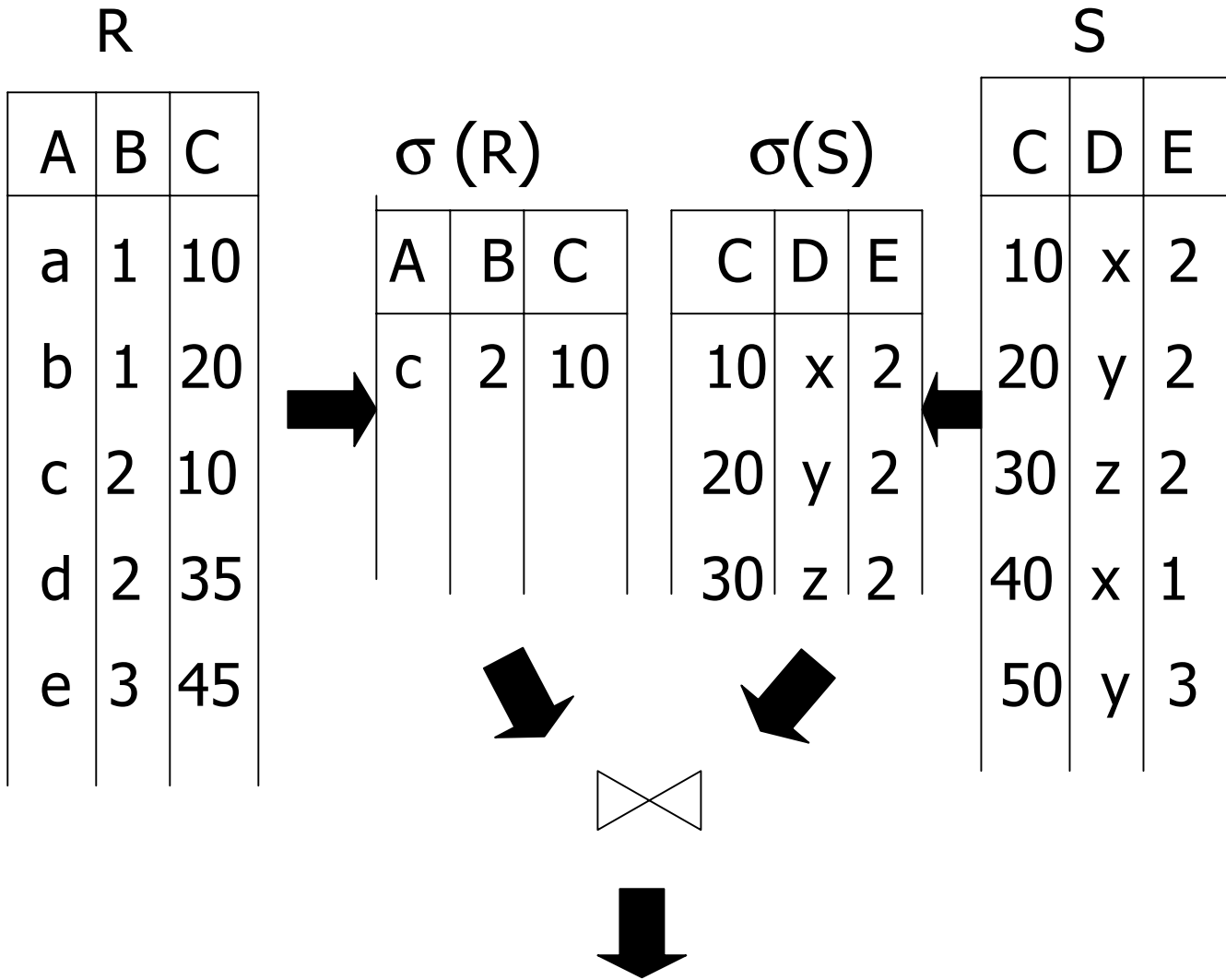
$\Pi_{B,D} [\sigma_{R.A="c" \wedge S.E=2 \wedge R.C=S.C} (RXS)]$

Jiný způsob:

Plán II



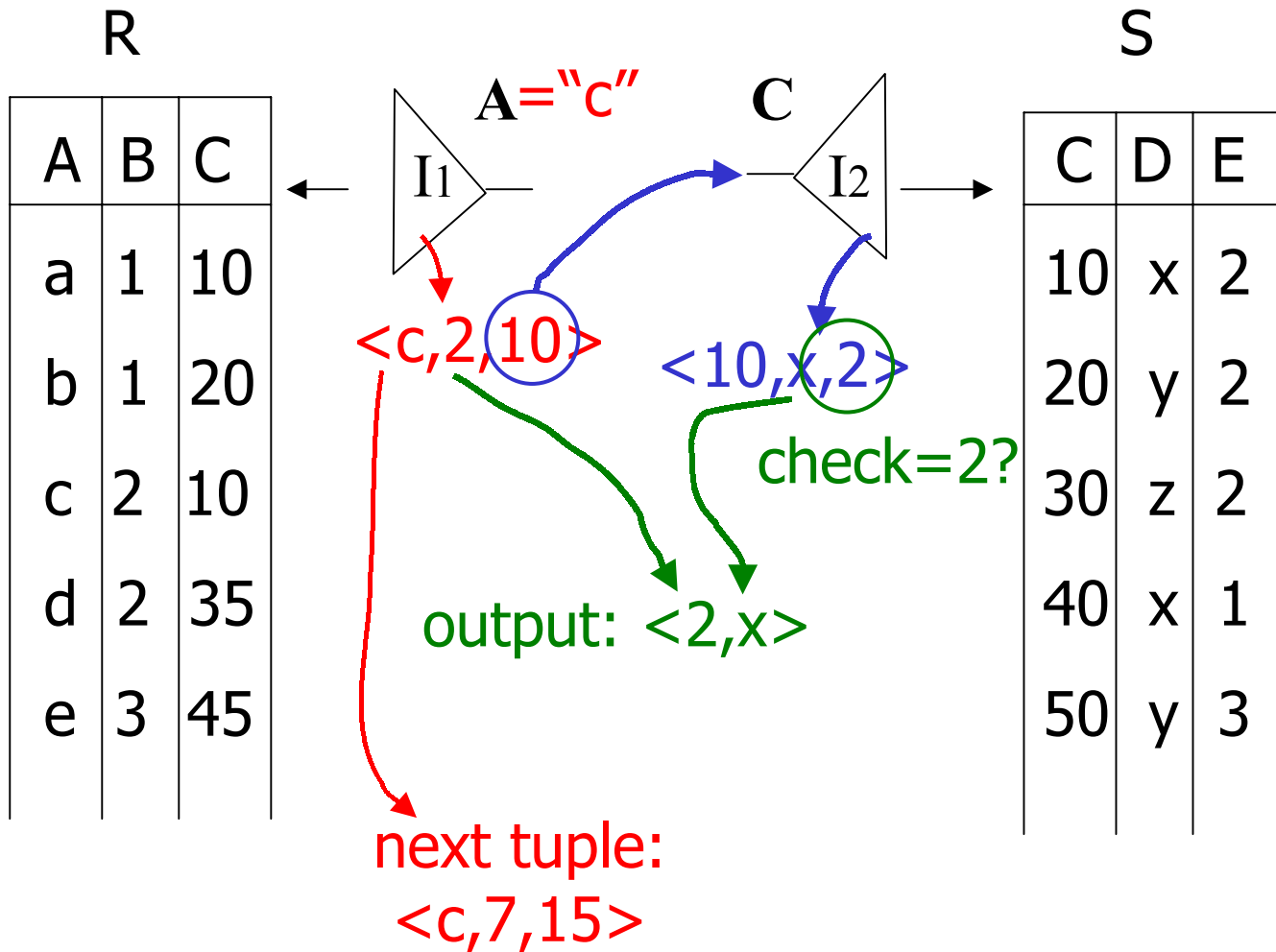
přirozené spojení



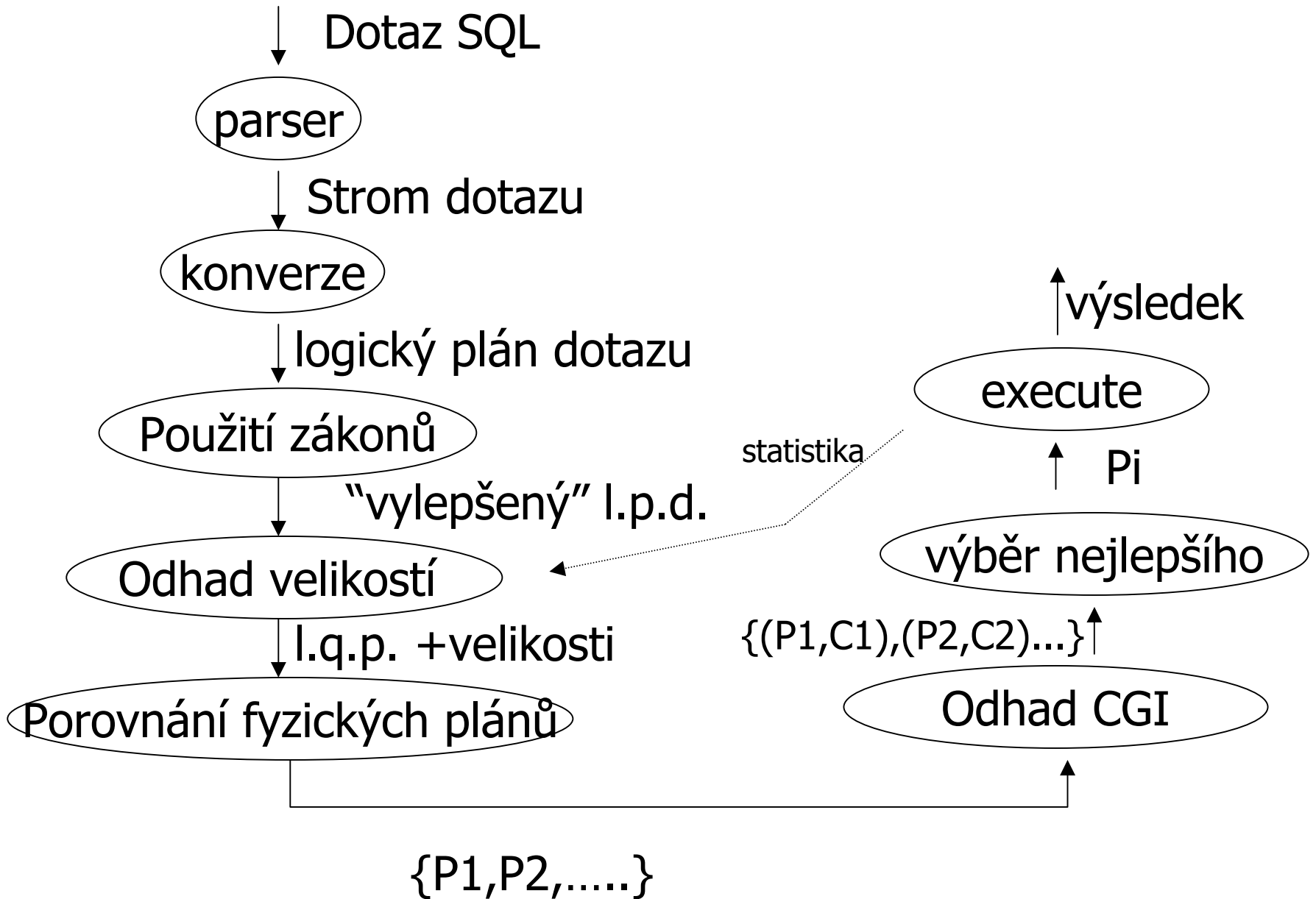
Plán III

Použijeme indexy R.A a S.C

- (1) použijeme index R.A k nalezení záznamů R splňující $R.A = "c"$
- (2) Pro každou nalezenou hodnotu R.C použijeme index S.C pro nalezení odpovídajících záznamů
- (3) Vypustíme záznamy S, kde $S.E \neq 2$
- (4) Spojíme odpovídající záznamy R,S provedeme projekci na atributy B,D



Přehled optimalizace dotazů

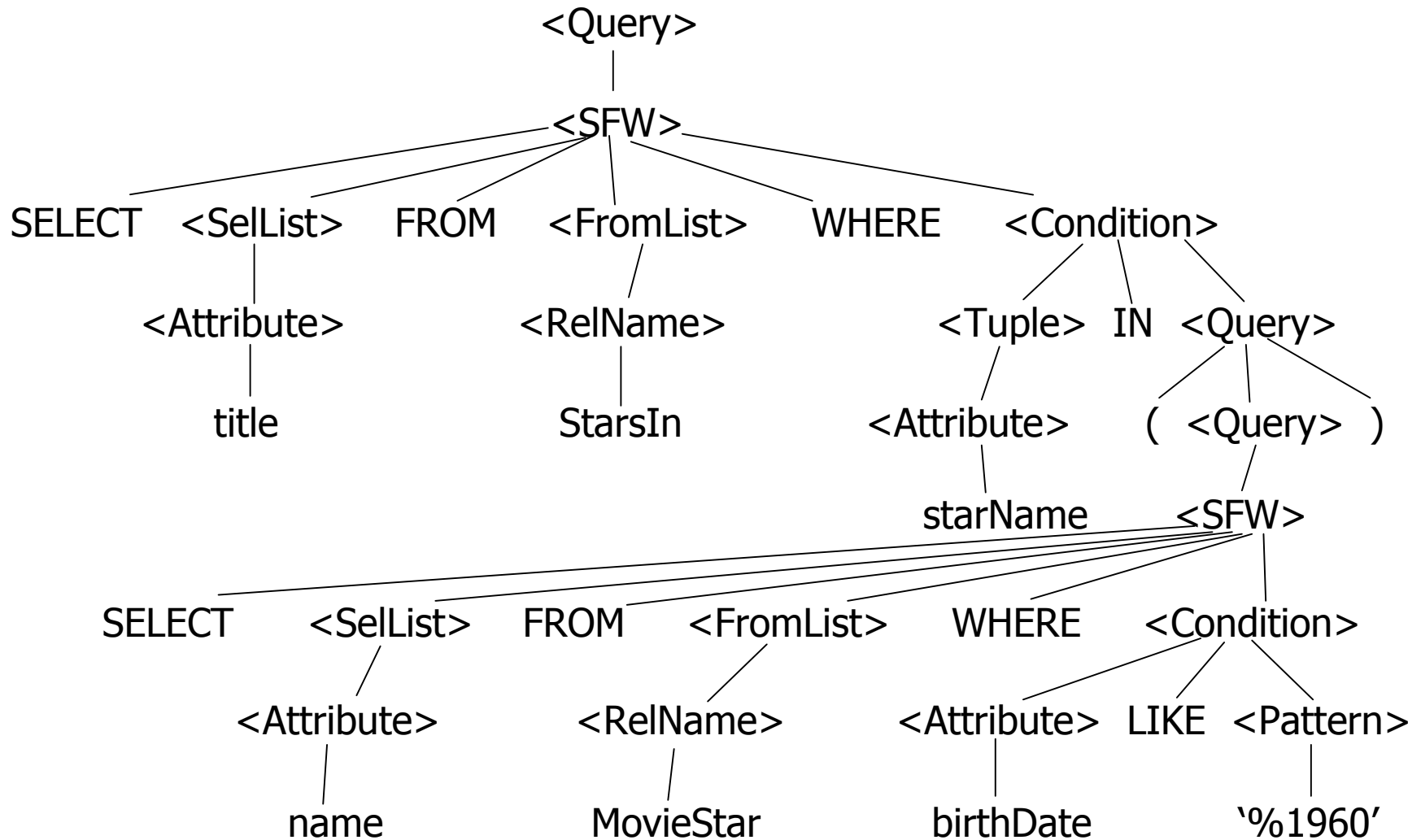


Příklad: dotaz SQL

```
SELECT title
FROM StarsIn
WHERE starName IN (
    SELECT name
    FROM MovieStar
    WHERE birthdate LIKE '%1960'
);
```

(Find the movies with stars born in 1960)

Příklad: Strom dotazu



Příklad: Generování Relační Algebry

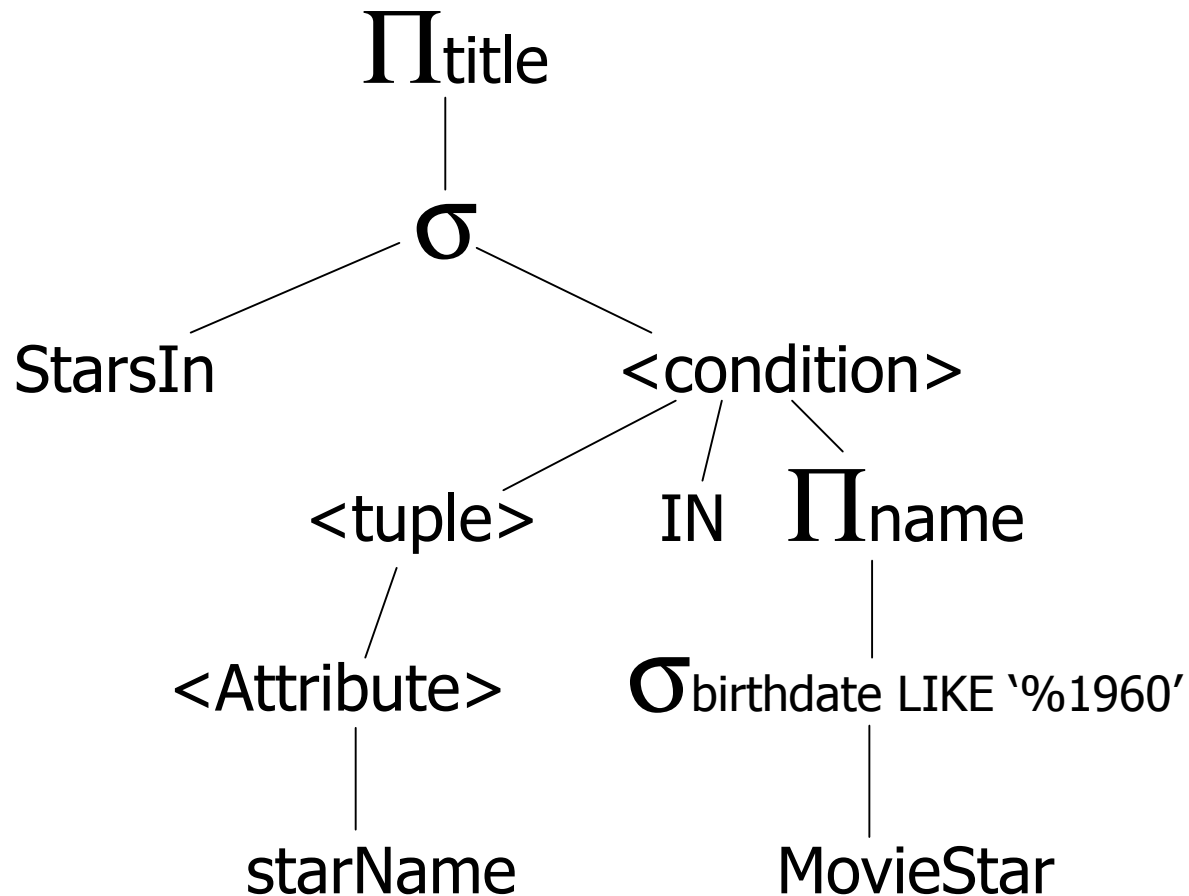


Fig. 7.15: An expression using a two-argument σ , midway between a parse tree and relational algebra

Příklad: Logický plán dotazu

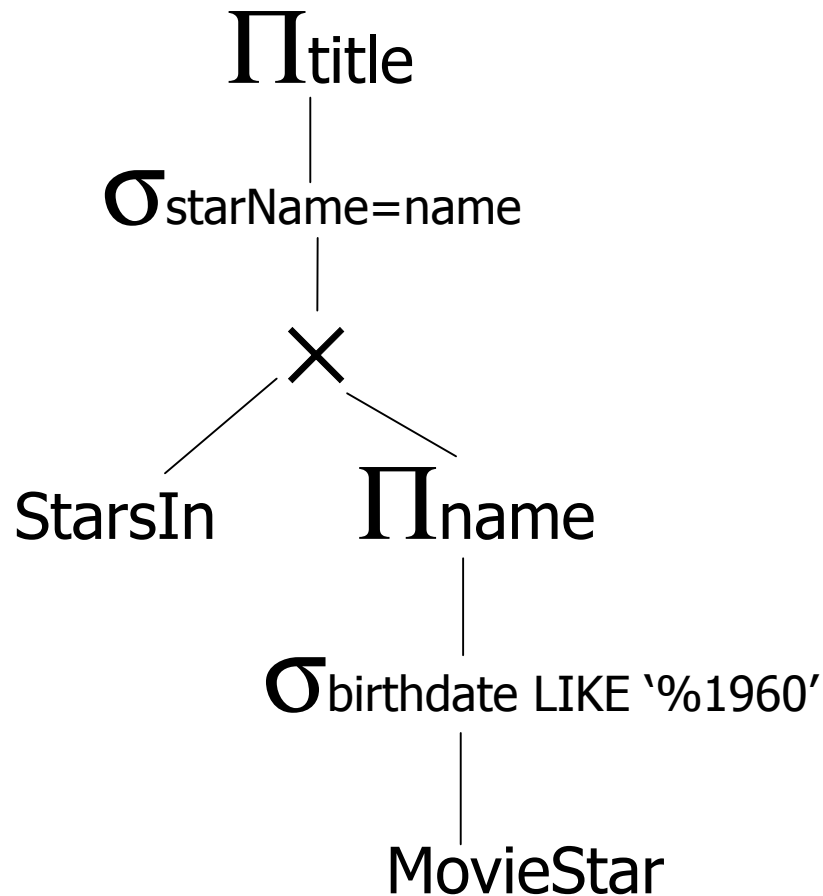
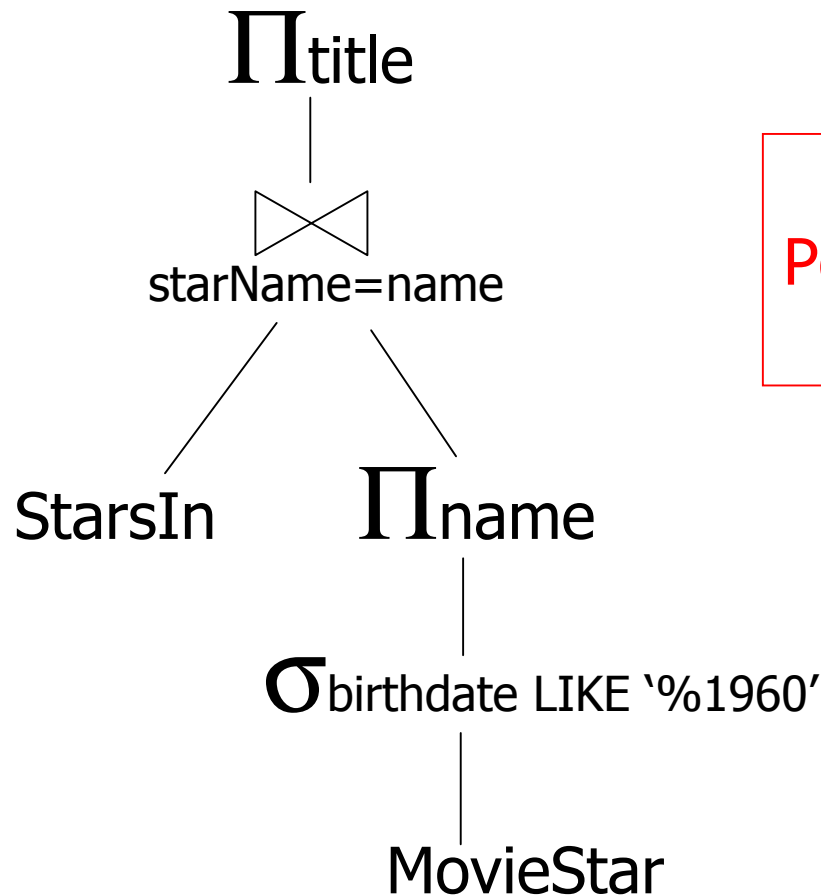


Fig. 7.18: Applying the rule for IN conditions

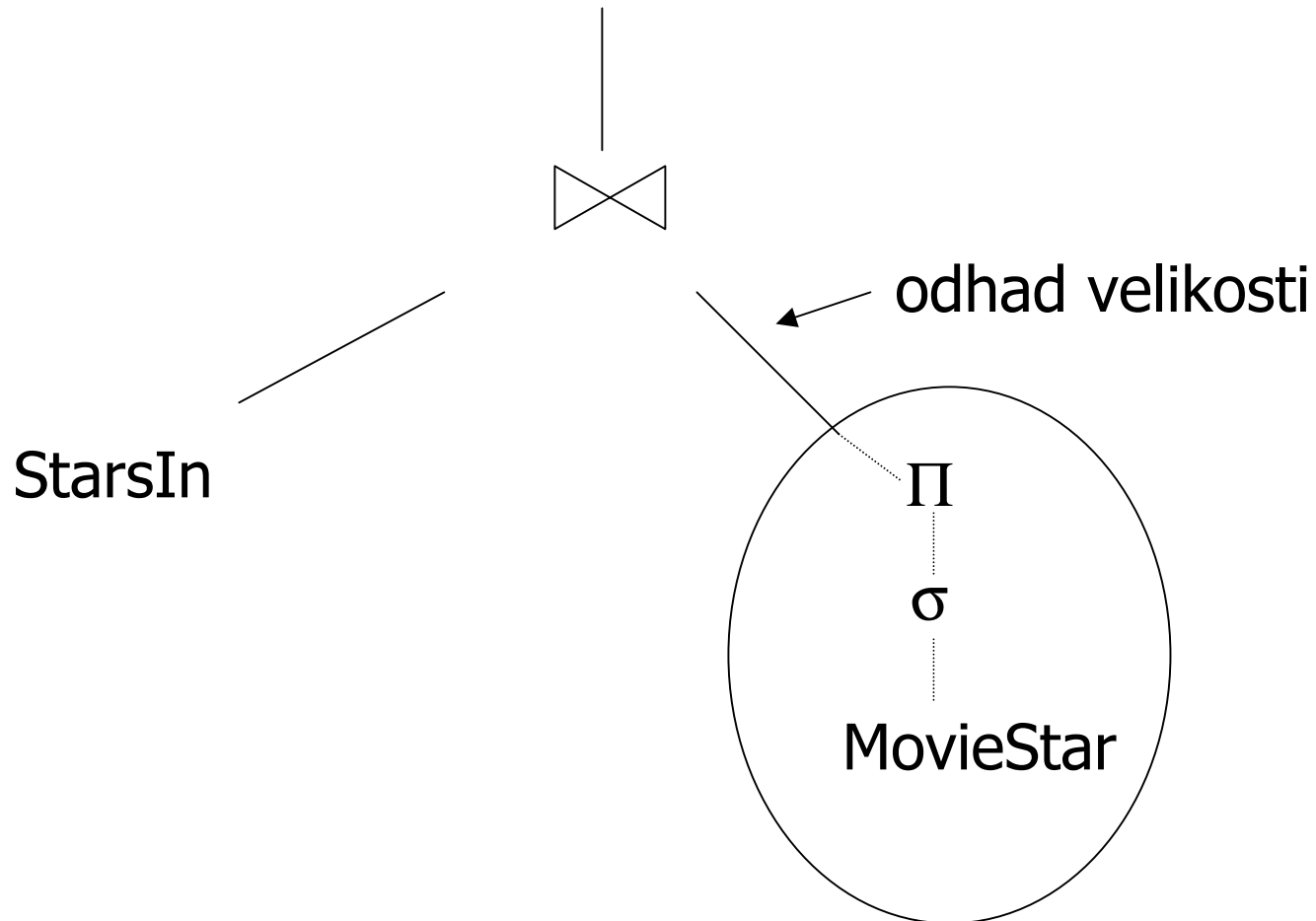
Příklad: Vylepšení logického plánu



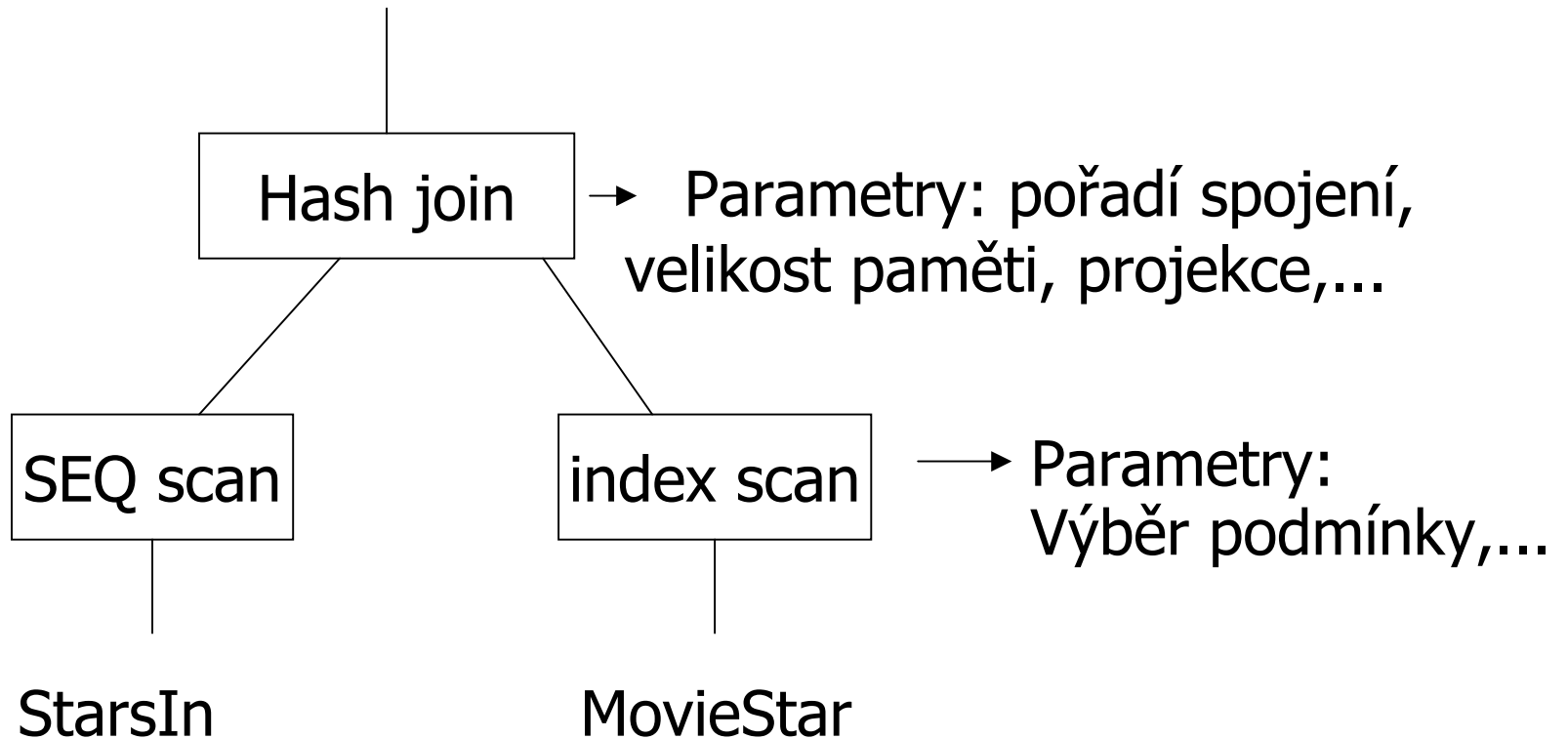
Otázka:
Posunout projekci
na StarsIn?

Fig. 7.20: An improvement on fig. 7.18.

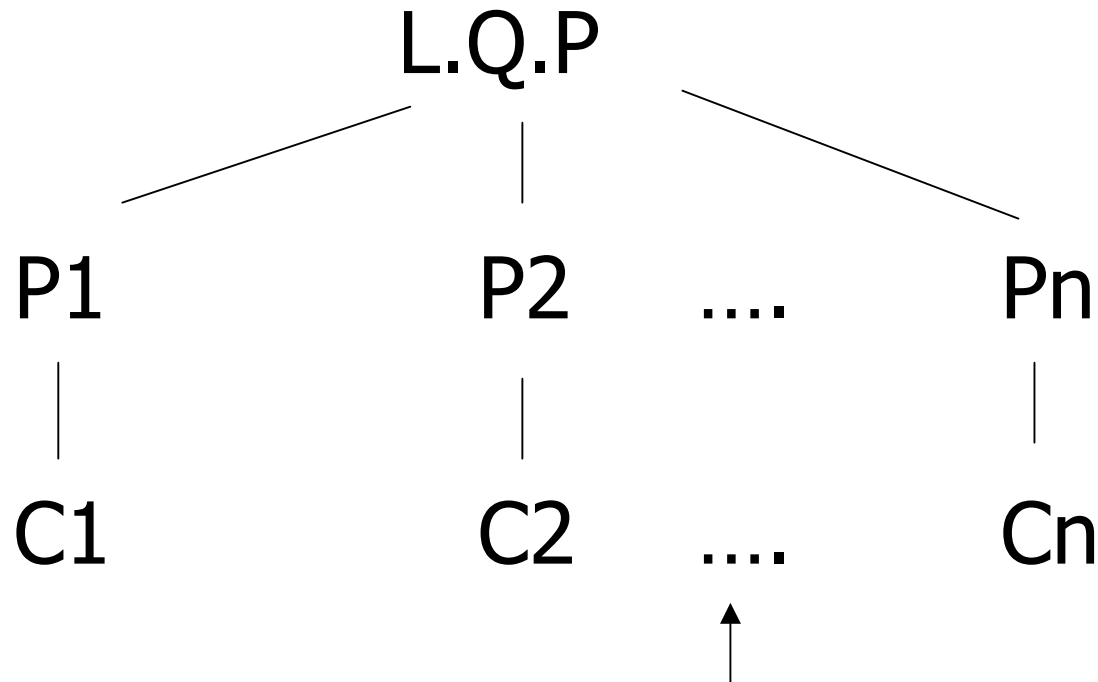
Příklad: Odhad velikostí výsledků



Příklad: jeden fyzický plán



Příklad: Odhad ceny



Vyber nejlepšího!

Optimalizace dotazu

- Úroveň relační algebry
- Úroveň podrobného plánu dotazu
 - Odhad ceny
 - bez indexů
 - s indexy
 - Vytvoření a porovnání plánů

Optimalizace

- Transformační pravidla
(zajištění ekvivalence)
- Jaké transformace?

pravidla: Přir. spojení & součin & sjednocení

$$R \bowtie S = S \bowtie R$$

$$(R \bowtie S) \bowtie T = R \bowtie (S \bowtie T)$$

$$R \times S = S \times R$$

$$(R \times S) \times T = R \times (S \times T)$$

$$R \cup S = S \cup R$$

$$R \cup (S \cup T) = (R \cup S) \cup T$$

pravidla: výběr

$$\sigma_{p_1 \wedge p_2}(R) = \sigma_{p_1} [\sigma_{p_2}(R)]$$

$$\sigma_{p_1 \vee p_2}(R) = [\sigma_{p_1}(R)] \cup [\sigma_{p_2}(R)]$$

Bags vs. Sets

$$R = \{a, a, b, b, b, c\}$$

$$S = \{b, b, c, c, d\}$$

$$R \cup S = ?$$

- Možnost 1 SUM

$$R \cup S = \{a, a, b, b, b, b, b, c, c, c, d\}$$

- Možnost 2 MAX

$$R \cup S = \{a, a, b, b, b, c, c, d\}$$

Možnost 2 (MAX) funguje:

$$\sigma_{p_1 \vee p_2}(R) = \sigma_{p_1}(R) \cup \sigma_{p_2}(R)$$

Příklad: $R = \{a, a, b, b, b, c\}$

P1 satisfied by a,b; P2 satisfied by b,c

$$\sigma_{p_1 \vee p_2}(R) = \{a, a, b, b, b, c\}$$

$$\sigma_{p_1}(R) = \{a, a, b, b, b\}$$

$$\sigma_{p_2}(R) = \{b, b, b, c\}$$

$$\sigma_{p_1}(R) \cup \sigma_{p_2}(R) = \{a, a, b, b, b, c\}$$

Pragmatické rozhodnutí

- > použití "SUM" pro sjednocení multimnožin
- > některá pravidla nemůžeme pro multimnožiny použít

pravidla: kombinace σ + \bowtie

Necht'

p = výraz obsahující pouze atributy R

q = výraz obsahující pouze atributy R

m = výraz obsahující atributy R a S

$$\sigma_p (R \bowtie S) = [\sigma_p (R)] \bowtie S$$

$$\sigma_q (R \bowtie S) = R \bowtie [\sigma_q (S)]$$

pravidla: kombinace $\sigma + \bowtie$

$$\sigma_{p \wedge q} (R \bowtie S) = [\sigma_p (R)] \bowtie [\sigma_q (S)]$$

$$\sigma_{p \wedge q \wedge m} (R \bowtie S) = \sigma_m [(\sigma_p R) \bowtie (\sigma_q S)]$$

$$\sigma_{p \vee q} (R \bowtie S) = [(\sigma_p R) \bowtie S] \cup [R \bowtie (\sigma_q S)]$$

Příklad odvození:

$$\sigma_{p \wedge q} (R \bowtie S) =$$

$$\sigma_p [\sigma_q (R \bowtie S)] =$$

$$\sigma_p [R \bowtie \sigma_q (S)] =$$

$$[\sigma_p (R)] \bowtie [\sigma_q (S)]$$

pravidla: kombinace π, σ

x = podmnožina atributů R

z = atributy ve výrazu P
(podmnožina atributů R)

$$\pi_x[\sigma_p(R)] = \pi_x \left\{ \sigma_p \left[\overset{\pi_{xz}}{\cancel{\pi_x}}(R) \right] \right\}$$

pravidla: kombinace π , \bowtie

x = podmnožina atributů R

y = podmnožina atributů S

z = průnik atributů R, S

$$\pi_{xy} (R \bowtie S) =$$

$$\pi_{xy} \{ [\pi_{xz} (R)] \bowtie [\pi_{yz} (S)] \}$$

$$\pi_{xy} \{ \sigma_p (R \bowtie S) \} =$$

$$\pi_{xy} \{ \sigma_p [\pi_{xz'} (R) \bowtie \pi_{yz'} (S)] \}$$

$$z' = z \cup \{ \text{atributy použité v } P \}$$

pravidla kombinace σ , \cup :

$$\sigma_p(R \cup S) = \sigma_p(R) \cup \sigma_p(S)$$

$$\sigma_p(R - S) = \sigma_p(R) - S = \sigma_p(R) - \sigma_p(S)$$

Jaké transformace použít?

$$\square \sigma_{p_1 \wedge p_2} (R) \rightarrow \sigma_{p_1} [\sigma_{p_2} (R)]$$

$$\square \sigma_p (R \bowtie S) \rightarrow [\sigma_p (R)] \bowtie S$$

$$\square R \bowtie S \rightarrow S \bowtie R$$

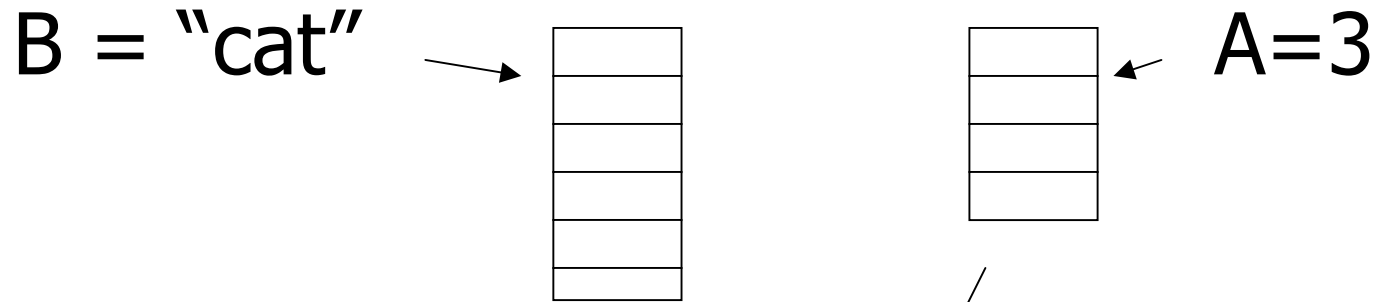
$$\square \pi_x [\sigma_p (R)] \rightarrow \pi_x \{ \sigma_p [\pi_{xz} (R)] \}$$

Obecné rada: projekci co nejdříve

Příklad: $R(A,B,C,D,E)$ $x=\{E\}$
P: $(A=3) \wedge (B=\text{"cat"})$

$\pi_x \{ \sigma_p (R) \}$ vs. $\pi_E \{ \sigma_p \{ \pi_{ABE}(R) \} \}$

□ Co když máme na A, B indexy?



Průnik seznamu ukazatelů dává
výsledek přímo

Obecná pravidla:

- bez transformací neuděláme chybu
- Většinou výhodné: brzký výběr
- eliminace společných podvýrazů

Outline - Query Processing

- Úroveň relační algebry
 - transformační pravidla
 - volba pravidel
- Úroveň podrobného plánu dotazu
 - odhad ceny
 - vytvoření a porovnání plánů

Odhad velikosti výsledku

- Statistika pro relaci R
 - $T(R)$: # záznamů R
 - $S(R)$: # bytů v každém záznamu R
 - $B(R)$: # bloků obsazených R
 - $V(R, A)$: # různých hodnot R
na attributech A

Příklad

R

A	B	C	D
cat	1	10	a
cat	1	20	b
dog	1	30	a
dog	1	40	c
bat	1	50	d

A: 20 byte string

B: 4 byte integer

C: 8 byte date

D: 5 byte string

$$T(R) = 5 \quad S(R) = 37$$

$$V(R,A) = 3$$

$$V(R,C) = 5$$

$$V(R,B) = 1$$

$$V(R,D) = 4$$

Odhad velikostí $W = R1 \times R2$

$$T(W) = T(R1) \times T(R2)$$

$$S(W) = S(R1) + S(R2)$$

Odhad velikostí $W = \sigma_{A=a}(R)$

$$S(W) = S(R)$$

$$T(W) = ?$$

Příklad

R	A	B	C	D
	cat	1	10	a
	cat	1	20	b
	dog	1	30	a
	dog	1	40	c
	bat	1	50	d

$$V(R,A)=3$$

$$V(R,B)=1$$

$$V(R,C)=5$$

$$V(R,D)=4$$

$$W = \sigma_{z=\text{val}}(R) \quad T(W) = \frac{T(R)}{V(R,Z)}$$

Předpoklad:

Hodnoty ve výrazu výběru $Z = \text{val}$ jsou rovnoměrně rozloženy mezi všechny různé hodnoty v R , kterých je $V(R,Z)$.

Alternativní předpoklad:

Hodnoty ve výrazu výběru $Z = \text{val}$
jsou rovnoměrně rozloženy
v oboru hodnot atributu Z : $\text{DOM}(R, Z)$

Příklad

R

A	B	C	D
cat	1	10	a
cat	1	20	b
dog	1	30	a
dog	1	40	c
bat	1	50	d

Alternativní předpoklad

$$V(R,A)=3 \quad \text{DOM}(R,A)=10$$

$$V(R,B)=1 \quad \text{DOM}(R,B)=10$$

$$V(R,C)=5 \quad \text{DOM}(R,C)=10$$

$$V(R,D)=4 \quad \text{DOM}(R,D)=10$$

$$W = \sigma_{z=\text{val}}(R) \quad T(W) = ?$$

$$\begin{aligned} C=\text{val} \Rightarrow T(W) &= (1/10)1 + (1/10)1 + \dots \\ &= (5/10) = 0.5 \end{aligned}$$

$$B=\text{val} \Rightarrow T(W) = (1/10)5 + 0 + 0 = 0.5$$

$$\begin{aligned} A=\text{val} \Rightarrow T(W) &= (1/10)2 + (1/10)2 + (1/10)1 \\ &= 0.5 \end{aligned}$$

Příklad

R

A	B	C	D
cat	1	10	a
cat	1	20	b
dog	1	30	a
dog	1	40	c
bat	1	50	d

Alternativní předpoklad

$$V(R,A)=3 \quad \text{DOM}(R,A)=10$$

$$V(R,B)=1 \quad \text{DOM}(R,B)=10$$

$$V(R,C)=5 \quad \text{DOM}(R,C)=10$$

$$V(R,D)=4 \quad \text{DOM}(R,D)=10$$

$$W = \sigma_{z=\text{val}}(R) \quad T(W) = \frac{T(R)}{\text{DOM}(R,Z)}$$

Velikost výběru

$SC(R,A)$ = průměrný počet záznamů
splňujících podmínku rovnosti na R.A

$$SC(R,A) = \left\{ \begin{array}{l} \frac{T(R)}{V(R,A)} \\ \frac{T(R)}{DOM(R,A)} \end{array} \right.$$

Co když $W = \sigma_{z \geq \text{val}}(R)$?

$$T(W) = ?$$

- Řešení # 1:

$$T(W) = T(R)/2$$

- Řešení # 2:

$$T(W) = T(R)/3$$

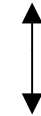
- Řešení # 3: Odhad podle rozsahu

Příklad R

	Z

Min=1

$V(R,Z)=10$



$W = \sigma_{z \geq 15} (R)$

Max=20

$$f = \frac{20-15+1}{20-1+1} = \frac{6}{20} \quad (\text{fraction of range})$$

$$T(W) = f \times T(R)$$

Analogicky:

$f \times V(R,Z)$ = fraction of distinct values

$$T(W) = [f \times V(Z,R)] \frac{T(R)}{V(Z,R)} = f \times T(R)$$

Odhad velikostí $W = R1 \bowtie R2$

$x =$ atributy z R1

$y =$ atributy z R2

Případ 1

$$X \cap Y = \emptyset$$

Stejně jako $R1 \times R2$

Případ 2

$$W = R1 \bowtie R2$$

$$X \cap Y = A$$

R1	A	B	C

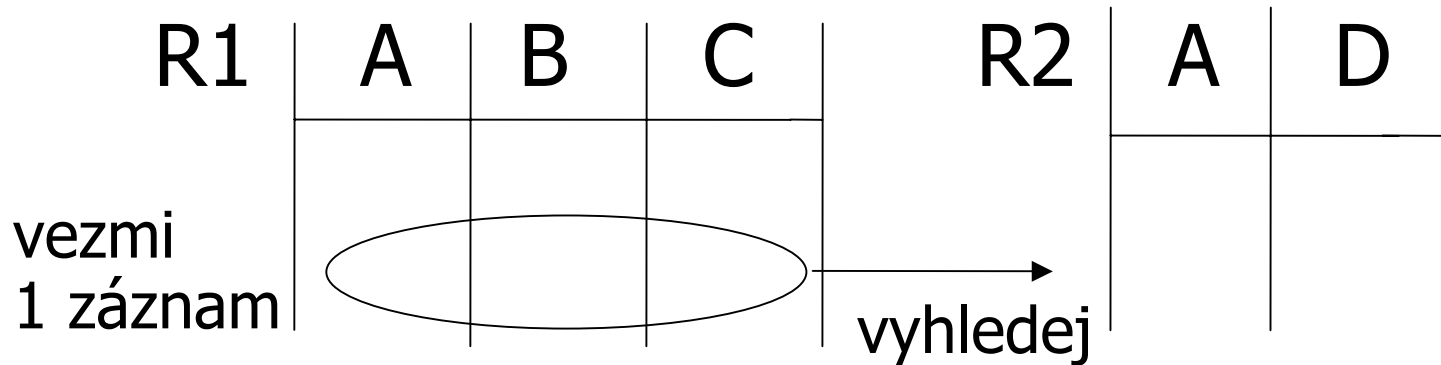
R2	A	D

Předpoklad:

$V(R1,A) \leq V(R2,A) \Rightarrow$ každá hodnota A z R1 je i v R2

$V(R2,A) \leq V(R1,A) \Rightarrow$ každá hodnota A z R2 je i v R1

Výpočet $T(W)$ pro $V(R1,A) \leq V(R2,A)$



1 záznam R1 se spojí s $\frac{T(R2)}{V(R2,A)}$ záznamy

$$\text{tj. } T(W) = \frac{T(R2)}{V(R2, A)} \times T(R1)$$

- $V(R1,A) \leq V(R2,A) \quad T(W) = \frac{T(R2) T(R1)}{V(R2,A)}$

- $V(R2,A) \leq V(R1,A) \quad T(W) = \frac{T(R2) T(R1)}{V(R1,A)}$

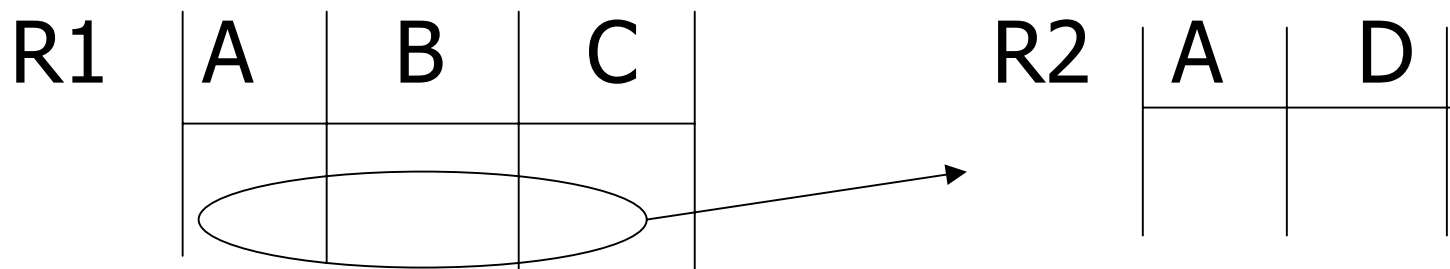
[A je společný atribut]

Obecně $W = R1 \bowtie R2$

$$T(W) = \frac{T(R2) T(R1)}{\max\{ V(R1,A), V(R2,A) \}}$$

Případ 2 Alternativní předpoklad

Hodnoty jsou rovnoměrně rozděleny v doméně



tyto záznamy se spojí s $T(R2)/\text{DOM}(R2,A)$

$$T(W) = \frac{T(R2) T(R1)}{\text{DOM}(R2, A)} = \frac{T(R2) T(R1)}{\text{DOM}(R1, A)}$$

↙ ↘
Předpokládáme stejné

Ve všech případech:

$$S(W) = S(R1) + S(R2) - S(A)$$

↙
velikost atributu A

Poznámka: pro složitější výrazy potřebujeme velikosti mezivýsledků.

$$W = [\underbrace{\sigma_{A=a}(R1)}] \bowtie R2$$

bereme jako relaci U

$$T(U) = T(R1)/V(R1,A) \quad S(U) = S(R1)$$

potřebujeme $V(U, *)$!!

Odhad V

Př.: $U = \sigma_{A=a}(R1)$

R1 má atributy A,B,C,D

$$V(U, A) =$$

$$V(U, B) =$$

$$V(U, C) =$$

$$V(U, D) =$$

Příklad

R1

A	B	C	D
cat	1	10	10
cat	1	20	20
dog	1	30	10
dog	1	40	30
bat	1	50	10

$$V(R1,A)=3$$

$$V(R1,B)=1$$

$$V(R1,C)=5$$

$$V(R1,D)=3$$

$$U = \sigma_{A=a}(R1)$$

$$V(U,A) = 1 \quad V(U,B) = 1 \quad V(U,C) = \frac{T(R1)}{V(R1,A)}$$

$V(D,U)$... něco mezi

Možné odhady $U = \sigma_{A=a}(R)$

$$V(U,A) = 1$$

$$V(U,B) = V(R,B)$$

Pro spojení $U = R1(A,B) \bowtie R2(A,C)$

$$V(U,A) = \min \{ V(R1, A), V(R2, A) \}$$

$$V(U,B) = V(R1, B)$$

$$V(U,C) = V(R2, C)$$

Příklad:

$$Z = R1(A,B) \bowtie R2(B,C) \bowtie R3(C,D)$$

R1	$T(R1) = 1000$	$V(R1,A)=50$	$V(R1,B)=100$
----	----------------	--------------	---------------

R2	$T(R2) = 2000$	$V(R2,B)=200$	$V(R2,C)=300$
----	----------------	---------------	---------------

R3	$T(R3) = 3000$	$V(R3,C)=90$	$V(R3,D)=500$
----	----------------	--------------	---------------

Mezivýsledek: $U = R \bowtie S$

$$T(U) = \frac{1000 \times 2000}{200}$$

$$V(U,A) = 50$$

$$V(U,B) = 100$$

$$V(U,C) = 300$$

$$Z = U \bowtie R3$$

$$T(Z) = \frac{1000 \times 2000 \times 3000}{200 \times 300}$$

$$V(Z,A) = 50$$

$$V(Z,B) = 100$$

$$V(Z,C) = 90$$

$$V(Z,D) = 500$$

Souhrn

- Odhad velikosti výsledků je “umění”
- Nezapomeňte:
 - Pro korektní odhad potřebujeme korektní statistiky – nutnost udržovat při modifikacích tabulky

Přehled

- Odhad ceny plánu dotazu
 - Odhad velikosti výsledku ← Máme
 - Odhad počtu V/V operací ← Příště