



# ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΑΤΡΩΝ - ΤΜΗΥΠ

## ΒΑΣΕΙΣ ΔΕΔΟΜΕΝΩΝ ΙΙ

---

### *B. Μεγαλοοικονόμου*

## Επεξεργασία Ερωτημάτων/Βελτιστοποίηση

(παρουσίαση βασισμένη εν μέρη σε σημειώσεις των Silberchatz, Korth και Sudarshan και του C. Faloutsos)



# Γενικά – σχεσιακό μοντέλο

---

- Σχεσιακό μοντέλο - SQL
- Συναρτησιακές εξαρτήσεις & Κανονικοποίηση
- Φυσικός σχεδιασμός, Δεικτοδότηση
- Επεξεργασία ερωτημάτων/βελτιστοποίηση
- Επεξεργασία Δοσοληψιών
- Προχωρημένα ζητήματα
  - Κατανεμημένες Βάσεις Δεδομένων
  - OO- και OR-DBMSs

# Επισκόπηση ενός ΣΔΒΔ

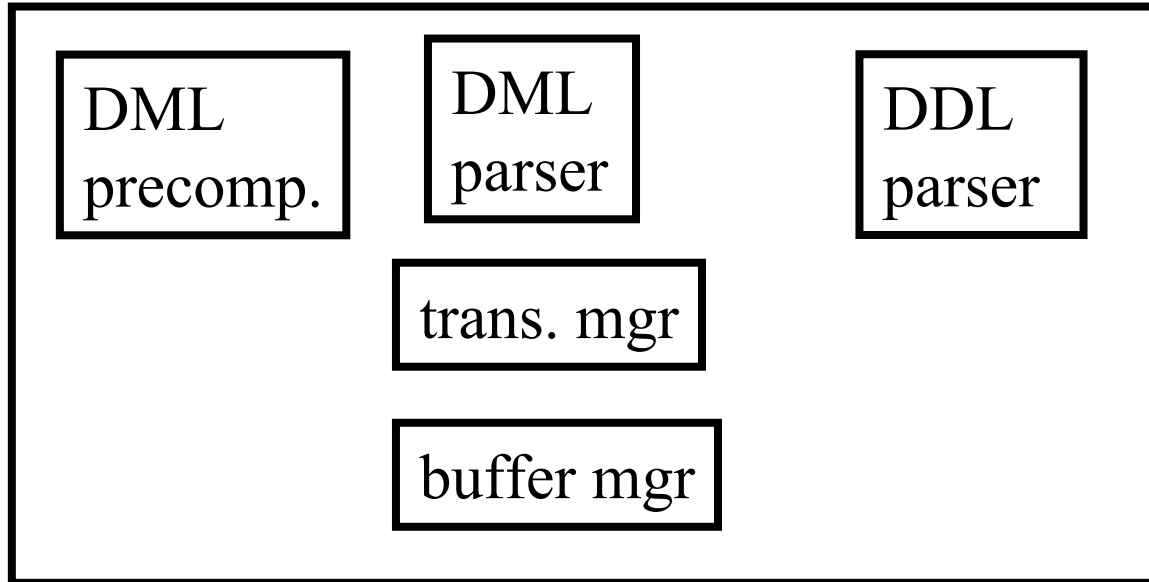
Naïve  
user



casual  
user



DBA



Αρχεία Δεδομένων κατάλογος



# Επισκόπηση- λεπτομερώς

---

- Κίνητρο- Γιατί βελτιστ. Ερωτήματος;
- Ισοδυναμία εκφράσεων
- Εκτίμηση κόστους
- Κόστος ευρετηρίων
- Στρατηγικές συνένωσης (join)



# Γιατί βελτιστ. Ερωτήματος;

---

- SQL: ~δηλωτική
- Καλή Βελτ. Ερωτ. → Μεγάλη διαφορά
  - π.χ., ακολουθιακή σάρωση vs B-tree ευρετήριο, σε  $P=1,000$  σελίδες



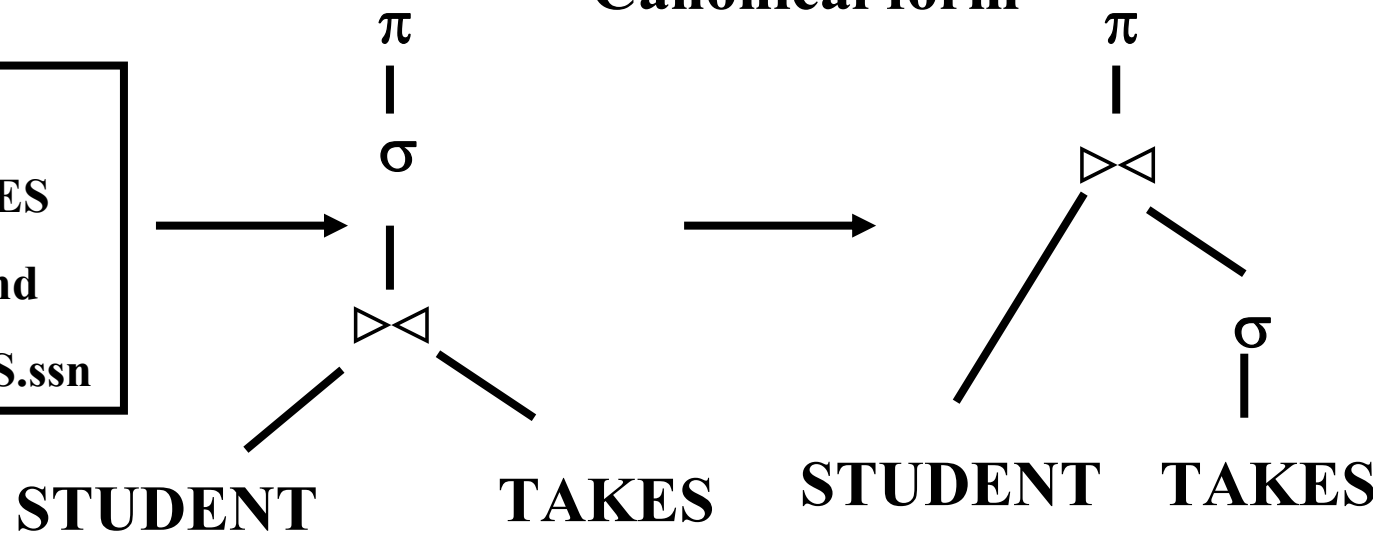
## Βήματα Βελτ.Ερωτ.

---

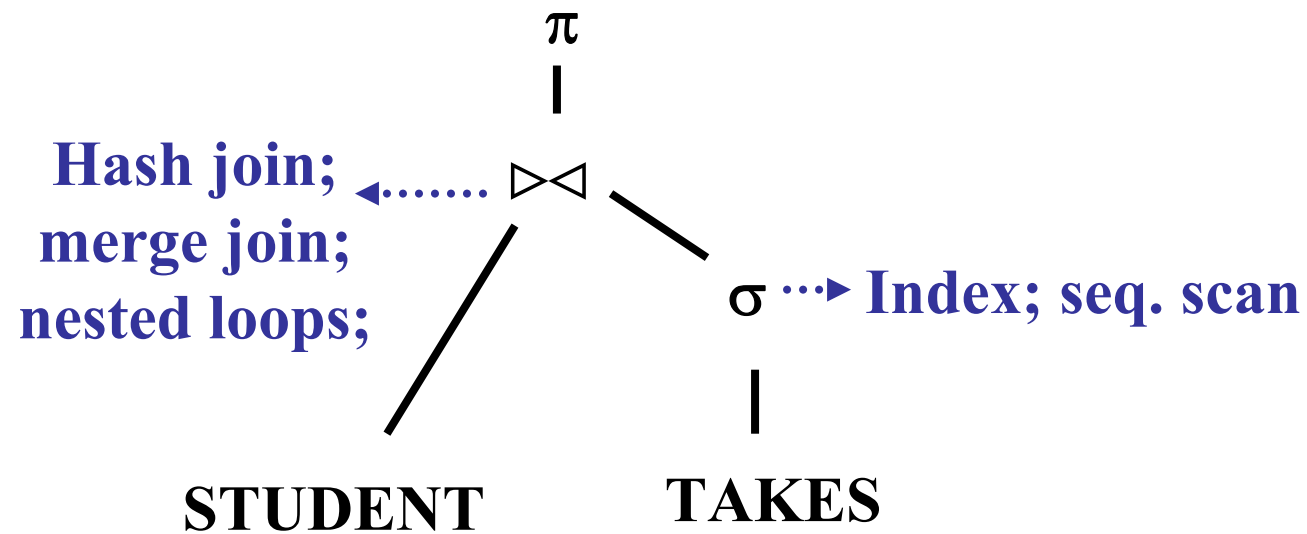
- φέρτε το ερώτημα στην εσωτερική μορφή (π.χ., parse tree)
- ... στην 'κανονική μορφή' (συντακτική Βελτ. Ερωτ.)
- δημιουργήστε εναλλακτικά σχέδια
- εκτιμήστε το κόστος, επιλέξτε το καλύτερο

# Βελτ.Ερωτ. - παράδειγμα

```
select name
from STUDENT, TAKES
where c-id='CIS331' and
STUDENT.ssn=TAKES.ssn
```



# Βελτ.Ερωτ. - παράδειγμα







# Επισκόπηση- λεπτομερώς

---

- Γιατί Βελτ.Ερωτ.;
- Ισοδυναμία εκφράσεων
- Εκτίμηση κόστους
- Κόστος ευρετηρίων
- Στρατηγικές συνένωσης (join)



# Ισοδυναμία εκφράσεων

---

- ... ή συντακτική Βελτ. Ερωτ.
- Εν συντομία: εφάρμοσε επιλογές και προβολές νωρίς
- Περισσότερες λεπτομέρειες:  
    δείτε κανόνες μετατροπής στο βιβλίο



# Ισοδυναμία εκφράσεων

- Ε: Πώς να αποδείξετε έναν κανόνα μετασχηματισμού;

$$\sigma_P(R1 \triangleright\triangleleft R2) \stackrel{?}{=} \sigma_P(R1) \triangleright\triangleleft \sigma_P(R2)$$

- Α: Χρησιμοποιείτε το σχεσιακό λογισμό πλειάδων (TRC) για να δείξετε ότι LHS = RHS, π.χ.:

$$\sigma_P(R1 \cup R2) \stackrel{?}{=} \sigma_P(R1) \cup \sigma_P(R2)$$



# Ισοδυναμία εκφράσεων

---

$$\sigma_P(R1 \cup R2) \stackrel{?}{=} \sigma_P(R1) \cup \sigma_P(R2)$$

$$t \in LHS \Leftrightarrow$$

$$t \in (R1 \cup R2) \wedge P(t) \Leftrightarrow$$

$$(t \in R1 \vee t \in R2) \wedge P(t) \Leftrightarrow$$

$$(t \in R1 \wedge P(t)) \vee (t \in R2) \wedge P(t) \Leftrightarrow$$



# Ισοδυναμία εκφράσεων

---

$$\sigma_P(R1 \cup R2) \stackrel{?}{=} \sigma_P(R1) \cup \sigma_P(R2)$$

...

$$(t \in R1 \wedge P(t)) \vee (t \in R2) \wedge P(t) \iff$$

$$(t \in \sigma_P(R1)) \vee (t \in \sigma_P(R2)) \iff$$

$$t \in \sigma_P(R1) \cup \sigma_P(R2) \iff$$

$t \in RHS$

*QED*



# Ισοδυναμία εκφράσεων

---

- Επιλογή (selection)

- πραγματοποιήστε τις νωρίς
- σπάστε ένα πολύπλοκο κατηγορημα και πιέστε

$$\sigma_{p1 \wedge p2 \wedge \dots \wedge pn}(R) = \sigma_{p1}(\sigma_{p2}(\dots \sigma_{pn}(R))\dots)$$

- απλοποιήστε ένα πολύπλοκο κατηγορημα
  - ('X=Y and Y=3') -> 'X=3 and Y=3'



# Ισοδυναμία εκφράσεων

---

- Προβολές (projections)
  - πραγματοποιήστε τις νωρίς (αλλά προσεκτικά...)
    - Μικρότερες πλειάδες
    - Λιγότερες πλειάδες (αν εξαλείφουν οι διπλότυπες)
  - Εξαιρέστε όλα τα γνωρίσματα εκτός από αυτά που ζητήθηκαν ή είναι απαραίτητα (π.χ., γνωρίσματα συνένωσης)



# Ισοδυναμία εκφράσεων

---

- ΣΥΝΕΝΩΣΕΙΣ

- Αντιστροφή, συσχέτιση

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

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

- E: n-way συνένωση – Πόσες διαφορετικές ακολουθίες; ... Διεξοδική αρίθμηση είναι πολύ αργή...





## Βήματα Βελτ.Ερωτ.

---

- φέρτε τα ερωτήματα στην εσωτερική μορφή (π.χ., parse tree)
- ... στην 'κανονική μορφή' (συντακτική Βελτ.Ερωτ.)
- δημιουργήστε εναλλακτικά σχέδια
- **εκτιμήστε το κόστος;** επιλέξτε το καλύτερο



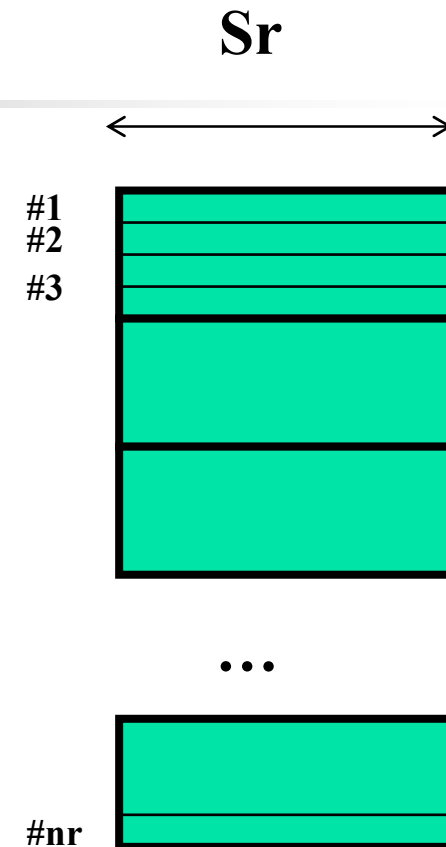
# Εκτίμηση κόστους

---

- Π.χ., βρείτε το Α.Μ. των φοιτητών με βαθμό 10 στις «Βάσεις Δεδομένων» (χρησιμοποιώντας ακολουθιακή αναζήτηση)
- Πόση ώρα θα πάρει ένα ερώτημα?
  - CPU (αλλά: μικρό κόστος, μειώνεται, δύσκολο να εκτιμηθεί)
  - Δίσκος (κυρίως, # μεταφορών blocks)
- Πόσες πλειάδες θα χαρακτηριστούν;
- (Ποια στατιστικά χρειαζόμαστε να κρατάμε;)

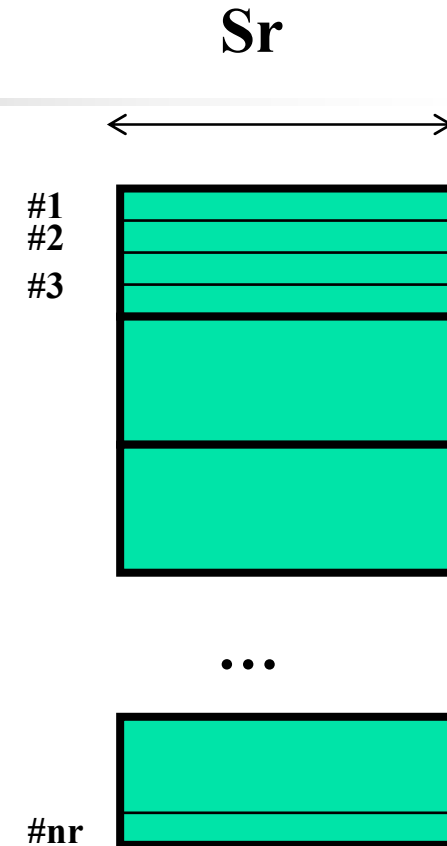
# Εκτίμηση κόστους

- Στατιστικά: για κάθε σχέση 'r' κρατάμε
  - nr : # πλειάδων;
  - Sr : μέγεθος πλειάδων σε bytes



# Εκτίμηση κόστους

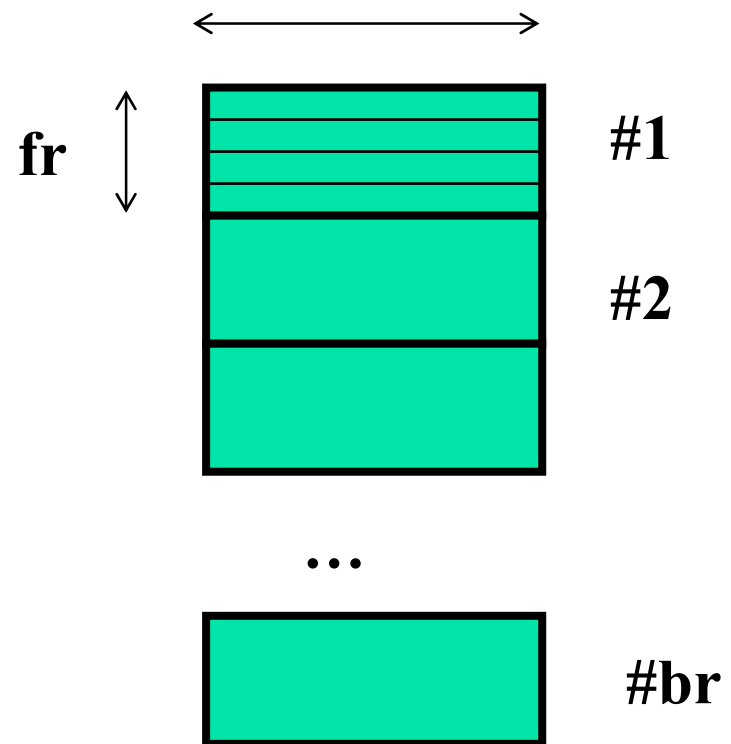
- Στατιστικά: για κάθε σχέση 'r' κρατάμε
  - ...
  - $V(A,r)$ : αριθμό διακριτών τιμών του γνωρίσματος 'A'
  - (πρόσφατα, και ιστογράμματα)



# Παράγωγα στατιστικά

$S_r$

- $fr$ : παράγοντας ομαδοποίησης (blocking factor) =  $\max\#$  εγγραφών/block (=?? )
- $br$ : # blocks (=?? )
- $SC(A,r) =$  επιλεκτικότητα =  $\text{avg}\#$  εγγραφών με  $A = \text{δοθέν}$  (=?? )





# Παράγωγα στατιστικά

---

- $fr$ : παράγοντας ομαδοποίησης =  $\max\#$  εγγραφών/block (=  $B/Sr$  ;  $B$ : μέγεθος block σε bytes)
- $br$ : # blocks (=  $nr / fr$  )



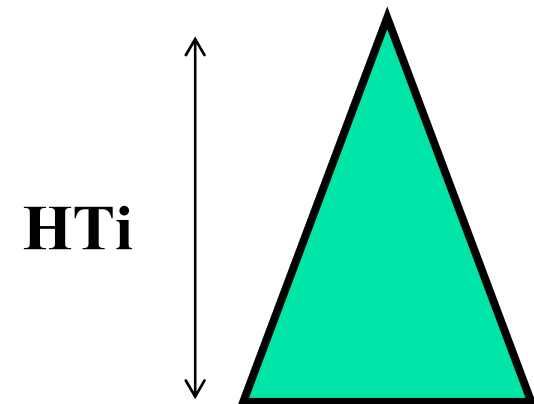
# Παράγωγα στατιστικά

---

- $SC(A,r)$  = επιλεκτικότητα = avg# εγγραφών με  $A$ =δοθέν ( $= nr / V(A,r)$ ) (προϋποθέτει ομοιομορφία...) – πχ: 30,000 φοιτητές, 10 σχολές – πόσοι φοιτητές στην Πολυτεχνική Σχολή;

# Επιπλέον μετρικές που χρειάζονται:

- Για ευρετήριο 'i':
  - $f_i$ : μέσος παράγοντας διακλάδωσης (fanout – degree) ( $\sim 50-100$ )
  - $H_i$ : # επίπεδα του ευρετηρίου 'i' ( $\sim 2-3$ )
    - $\sim \log(\#entries)/\log(f_i)$
  - $L_i$ : # μπλοκς στο επίπεδο των φύλλων







# ΣΤΑΤΙΣΤΙΚΑ

---

- Που τα αποθηκεύουμε;
- Πόσο συχνά τα ανανεώνουμε;



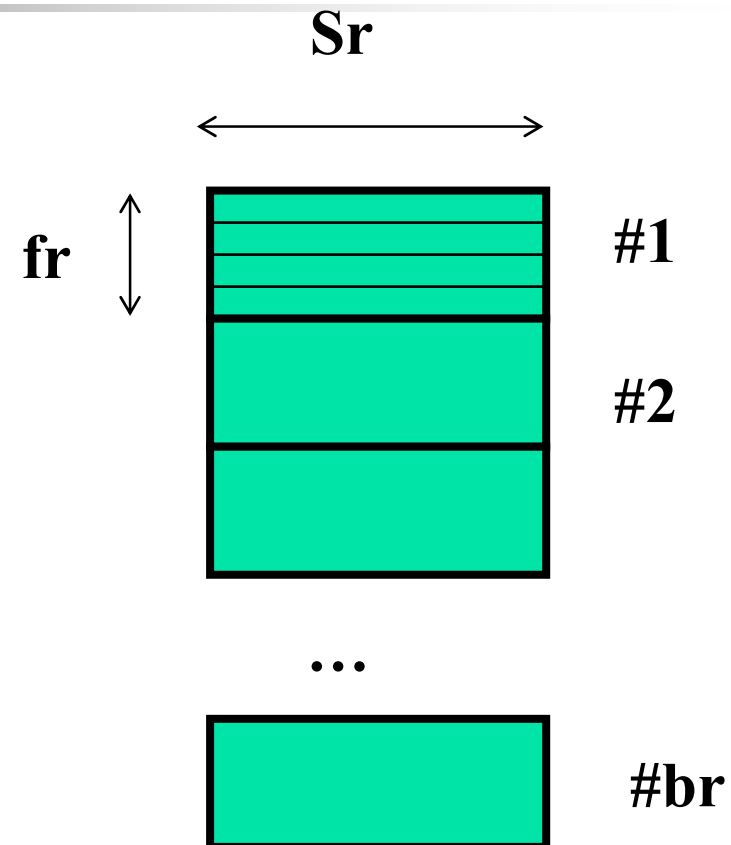
# Βήματα Βελτ.Ερωτ.

---

- Φέρτε το ερώτημα στην εσωτερική μορφή (π.χ., parse tree)
- ... στην 'κανονική μορφή' (συντακτική Βελτ.Ερωτ.)
- δημιουργήστε εναλλακτικά πλάνα
  - επιλογές; ταξινόμηση; προβολές
  - συνενώσεις
- εκτιμήστε το κόστος, επιλέξτε το καλύτερο

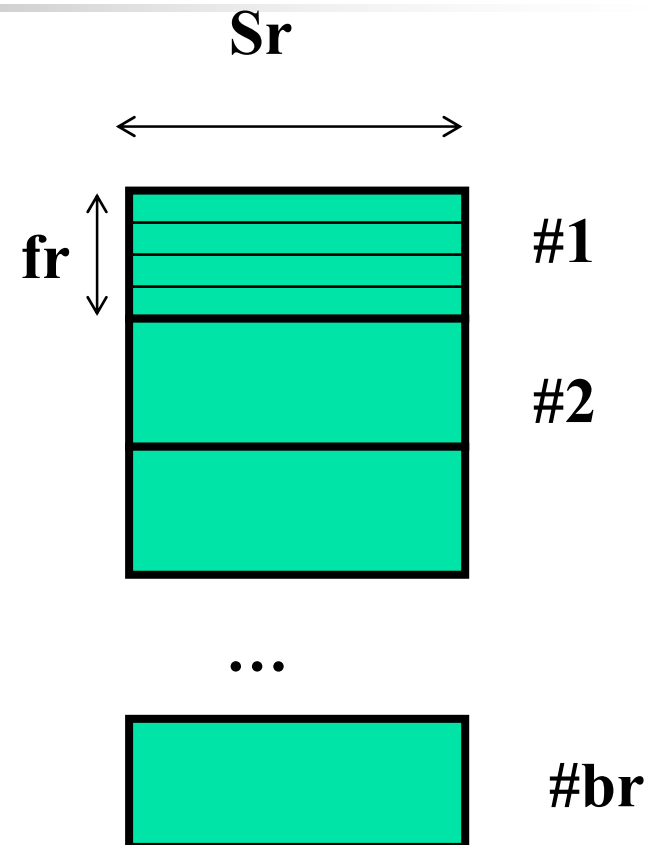
# Εκτίμηση κόστους & δημιουργία πλάνου

- Επιλογές – π.χ.,  
**select** \*  
**from** TAKES  
**where** grade = 'A'
- Πλάνα;



# Εκτίμηση κόστους & δημιουργία πλάνου

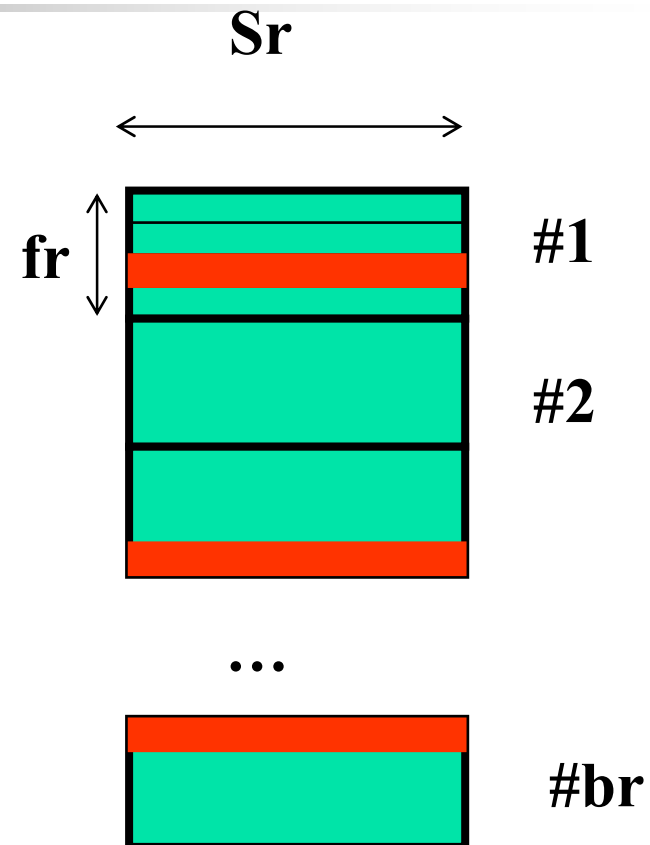
- Πλάνα;
  - Σειριακή αναζήτηση
  - Δυαδική αναζήτηση
    - (αν είναι ταξινομημένο κατά αύξουσα σειρά)
  - Αναζήτηση ευρετηρίου
    - αν υπάρχει ένα ευρετήριο



# Εκτίμηση κόστους & δημιουργία πλάνου

## Διαδοχ. Αναζ.– κόστος;

- $br$  (χειρότερη περίπτωση)
- $br/2$  (μέση περίπτωση, αν αναζητάμε με βάση το πρωτεύων κλειδί)

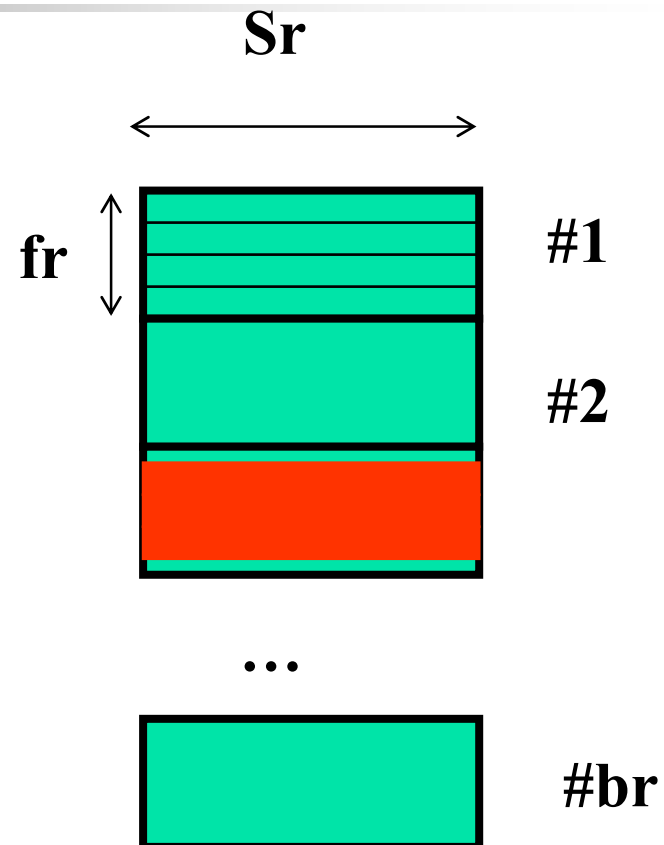


# Εκτίμηση κόστους & δημιουργία πλάνου

## δυναμική αναζ.- κόστος;

αν είναι ταξινομημένο κατά αύξουσα σειρά:

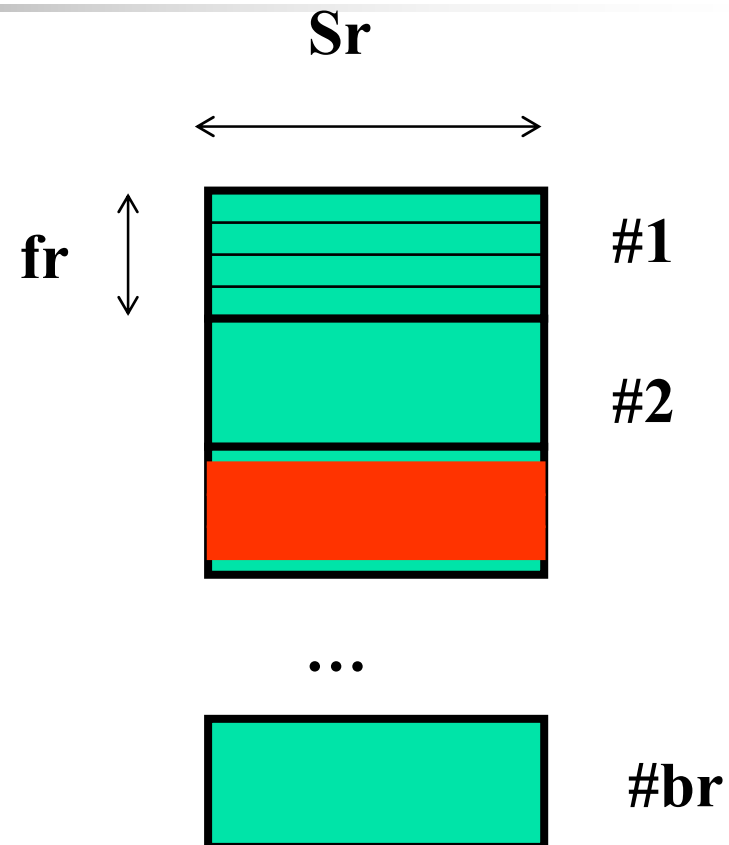
- $\sim \log(br) + SC(A,r)/fr$  (= #μπλοκς καταλυμένων από αρμόδιες πλειάδες)  
-1



# Εκτίμηση κόστους & δημιουργία πλάνου

εκτίμηση της επιλογής  
πρωτευόντων  
 $SC(A,r)$ :

**μη τετριμμένη**



# Εκτίμηση κόστους & δημιουργία πλάνου

## μέθοδος#3:

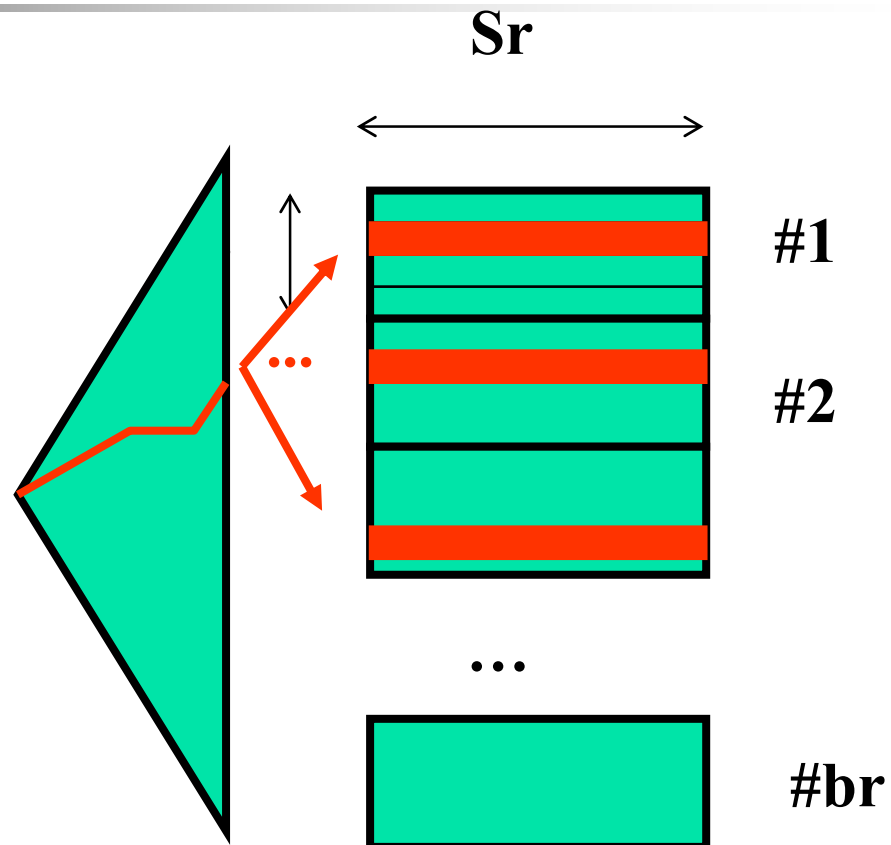
### ευρετήριο – κόστος?

- επίπεδα ευρετηρίου +
- μπλοκς με αρμόδιες πλειάδες

περίπτωση#1: πρωτεύων κλειδί

περίπτωση#2: δευτ. κλειδί – συσταδοποιημένο ευρετήριο

περίπτωση#3: δευτ. κλειδί – μη συσταδοποιημένο ευρετήριο





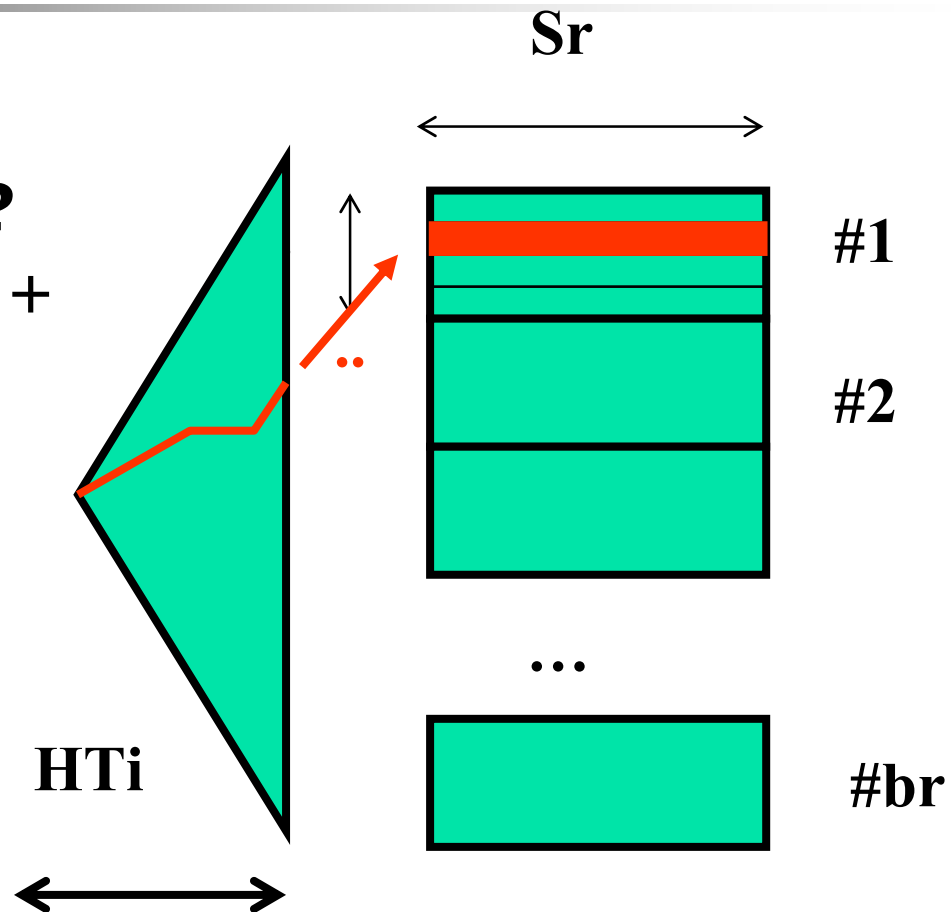
# Εκτίμηση κόστους & δημιουργία πλάνου

## μέθοδος#3: ευρετήριο – κόστος?

- επίπεδα ευρετηρίου +
- μπλοκς με αρμόδιες πλειάδες

περίπτωση#1: πρωτεύων  
κλειδί – κόστος:

$H_{Ti} + 1$



# Εκτίμηση κόστους & δημιουργία πλάνου

## μέθοδος#3:

### ευρετήριο – κόστος?

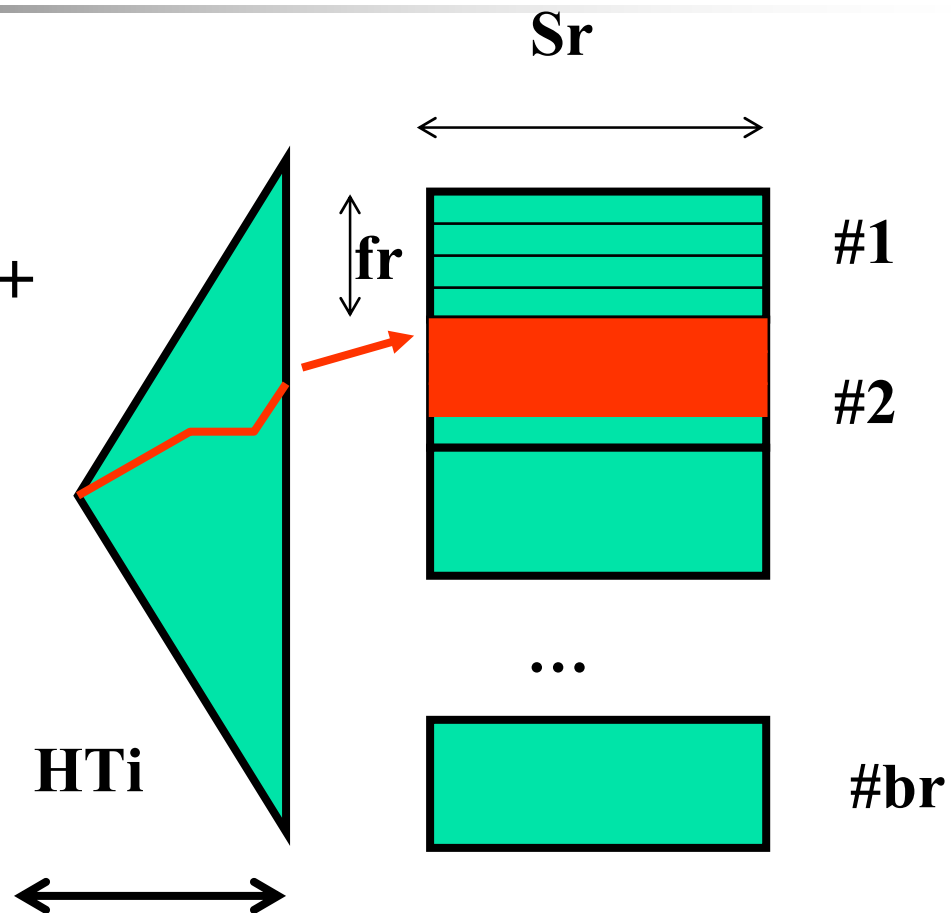
- επίπεδα ευρετηρίου +
- μπλοκς με αρμόδιες πλειάδες

περίπτωση#2: δευτ. κλειδί – συσταδοποιημένο ευρετήριο

Ἡ πρωτ. ευρετ.σε μη-κλειδί

...ανάκληση πολλαπλών εγγραφών

$$HT_i + SC(A,r)/fr$$



# Εκτίμηση κόστους & δημιουργία πλάνου

## μέθοδος#3:

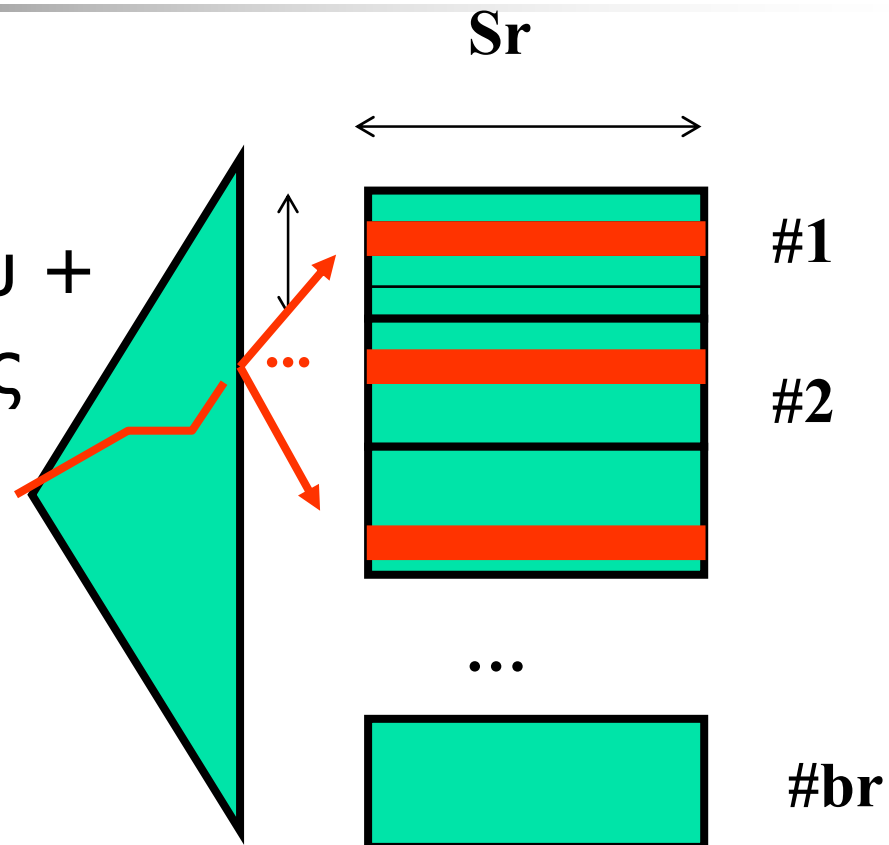
### ευρετήριο – κόστος?

- επίπεδα ευρετηρίου +
- μπλοκς με αρμόδιες πλειάδες

περίπτωση#3: δευτ. κλειδί  
– μη συσταδοποιημένο  
ευρετήριο

$HT_i + SC(A,r)$

(actually, pessimistic...)





# Εκτίμηση κόστους – αριθμητικό παράδειγμα

---

- βρείτε τους λογαριασμούς με *όνομα-υποκαταστήματος* = 'Πάτρα'
- Λογαριασμός (*όνομα-υποκαταστήματος, υπόλοιπο, ...*)



## Αριθμητικό παράδειγμα (συνεχ.)

- $n$ -account = 10,000 πλειάδες
- $f$ -account = 20 πλειάδες/μπλοκ
- $V(\text{υπόλοιπο, λογαριασμός}) = 500$   
διακριτές τιμές
- $V(\text{όνομα-υποκαταστήματος, λογαριασμός}) = 50$  διακριτές τιμές
- για  $\text{branch-index}$ : παράγοντας διακλάδωσης (fanout)  $f_i = 20$



# Αριθμητικό παράδειγμα

---

- E1: κόστος σειριακής αναζήτησης;
- A1: 500 προσπελάσεις δίσκου
- E2: υποθέστε ένα ευρετήριο συστάδων στο όνομα-υποκαταστήματος – κόστος;

# Εκτίμηση κόστους & δημιουργία πλάνου

μέθοδος#3: ευρετήριο

– κόστος;

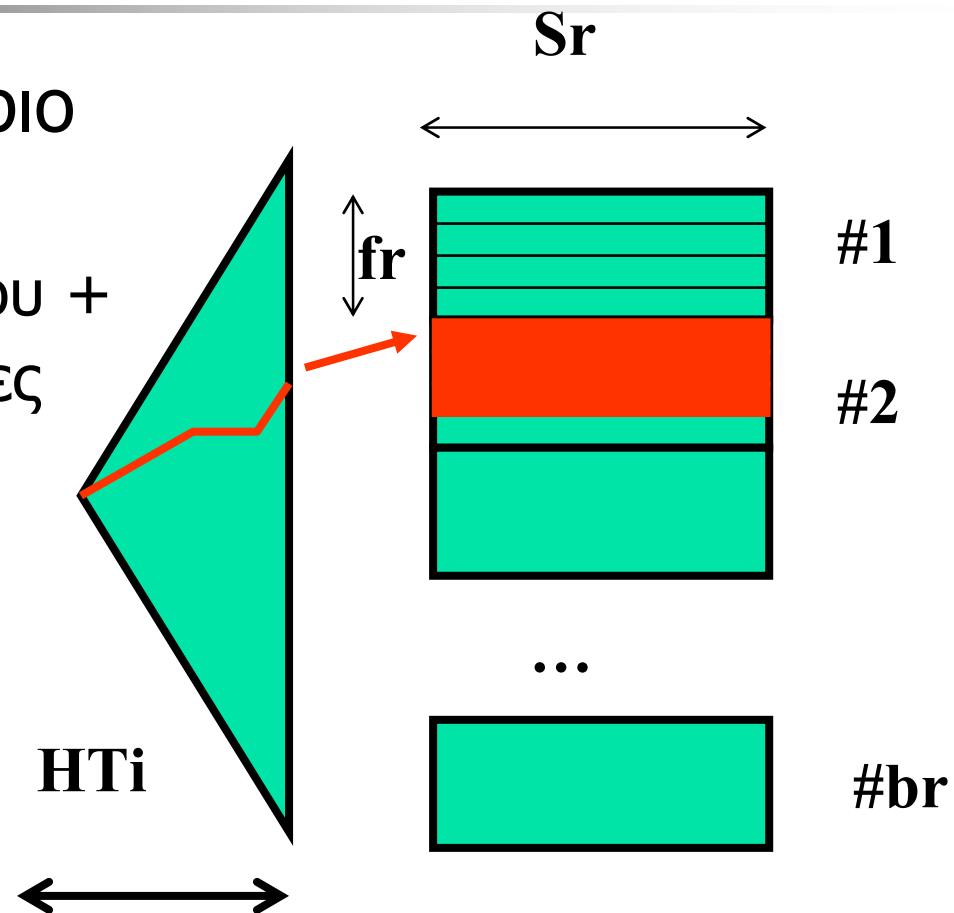
- Επίπεδα ευρετηρίου +
- μπλοκς με αρμόδιες πλειάδες

περίπτωση#2: δευτ. κλειδί

– συσταδοποιημένο

ευρετήριο

$HT_i + SC(A,r)/fr$





# Αριθμητικά παραδείγματα

---

- A2:  
HTi +  
SC(όνομα-υποκαταστήματος,  
λογαριασμός)/f-account
- HTi: 50 τιμές, με παράγοντας  
διακλάδωσης 20 -> HT=2 επίπεδα  
( $\log(50)/\log(20) = 1+$ )
- SC(..) = # αρμόδιες εγγραφές =
  - $nr/V(A,r) = 10,000/50 = 200$  πλειάδες
  - SC/f: έκταση 200/20 μπλοκς = 10 μπλοκς





# Αριθμητικά παραδείγματα

---

- Α2 τελική απάντηση:  
 $2+10=12$  προσπελάσεις μπλοκς
- (αντί 500 προσπελάσεις μπλοκς  
σειριακής αναζήτησης)
- σημείωση: in all fairness
  - σειριακές αναζητήσεις δίσκου:  $\sim 2\text{msec}$  (ή και λιγότερο)
  - Τυχαίες προσπελάσεις μπλοκς:  $\sim 10\text{msec}$



# Επισκόπηση- λεπτομερώς

---

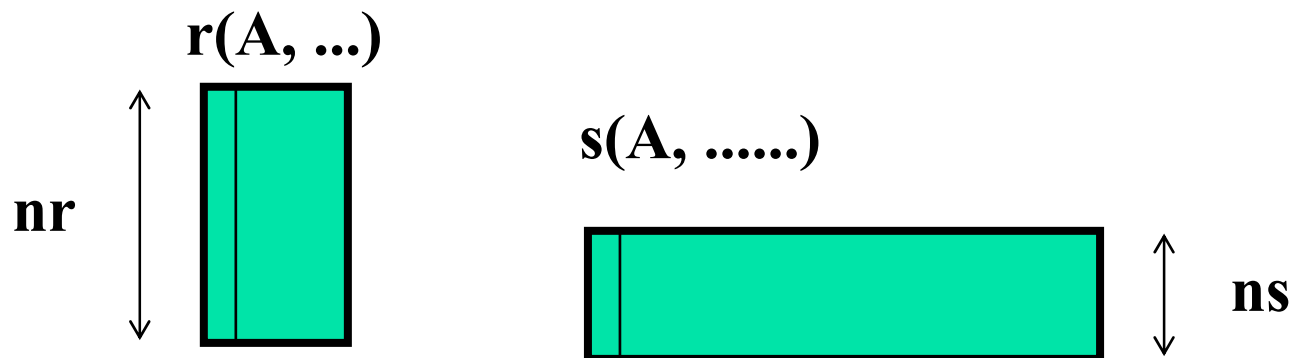
- Κίνητρο- Γιατί βελτιστ. Ερωτήματος;
- Ισοδυναμία εκφράσεων
- Εκτίμηση κόστους
- Κόστος ευρετηρίων
- Στρατηγικές συνένωσης (join)



## 2-way joins

---

- αλγόριθμος(-οι) για  $r$  JOIN  $s$ ?
- $n_r$ ,  $n_s$  πλειάδες κάθε σχέση



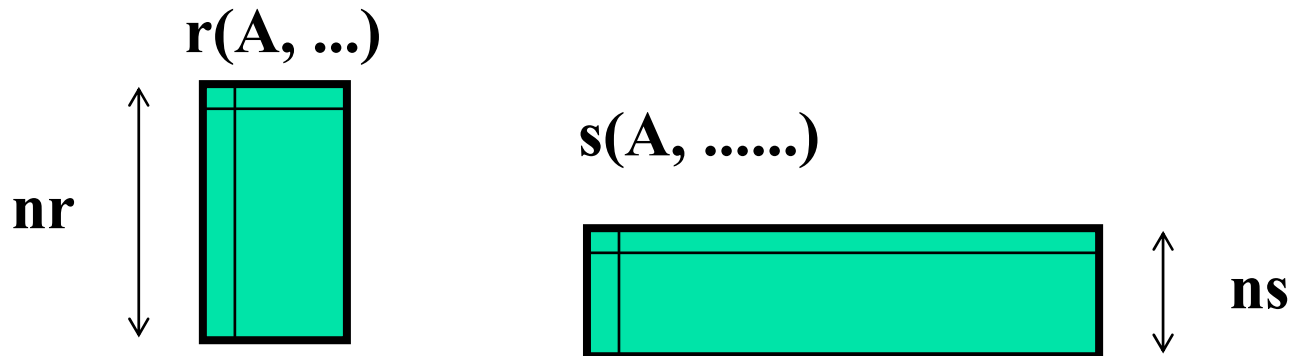
# 2-way joins

- Αλγόριθμος #0: (αφελής) εμφωλευμένες επαναλήψεις (**ΑΡΓΟ!**)

for each tuple  $tr$  of  $r$

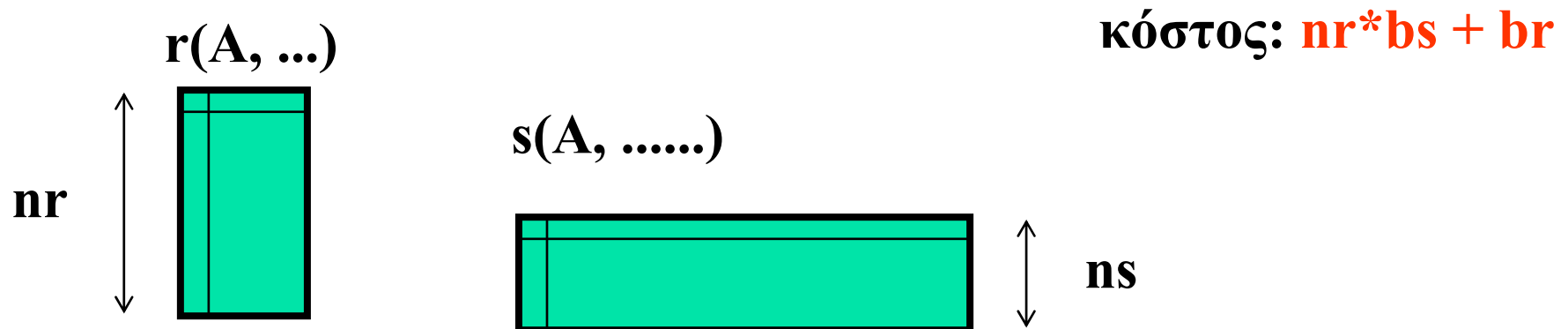
for each tuple  $ts$  of  $s$

print, if they match



## 2-way joins

- Αλγόριθμος #0: γιατί είναι κακή?
- πόσες προσπελάσεις δίσκου ('br' και 'bs' είναι ο αριθμός των μπλοκς για 'r' και 's')?



# 2-way joins

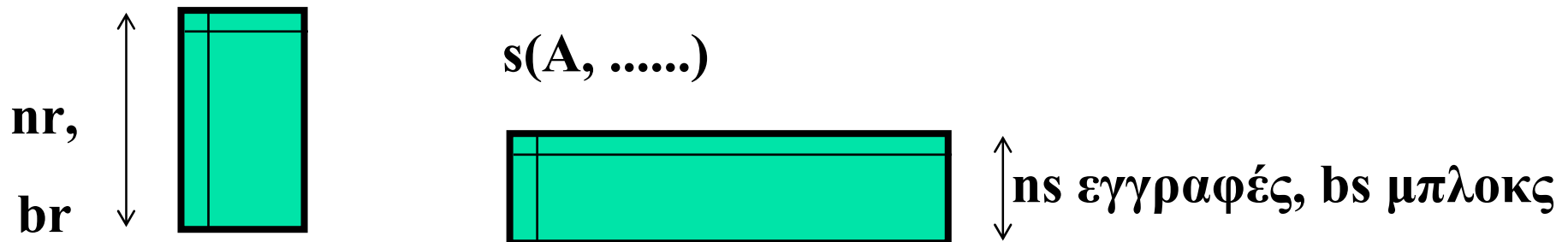
- Αλγόριθμος #1: Blocked εμφωλευμένων επαναλήψεων συνένωση

- read in a block of r

κόστος:  $br + br * bs$

- read in a block of s

- print matching tuples

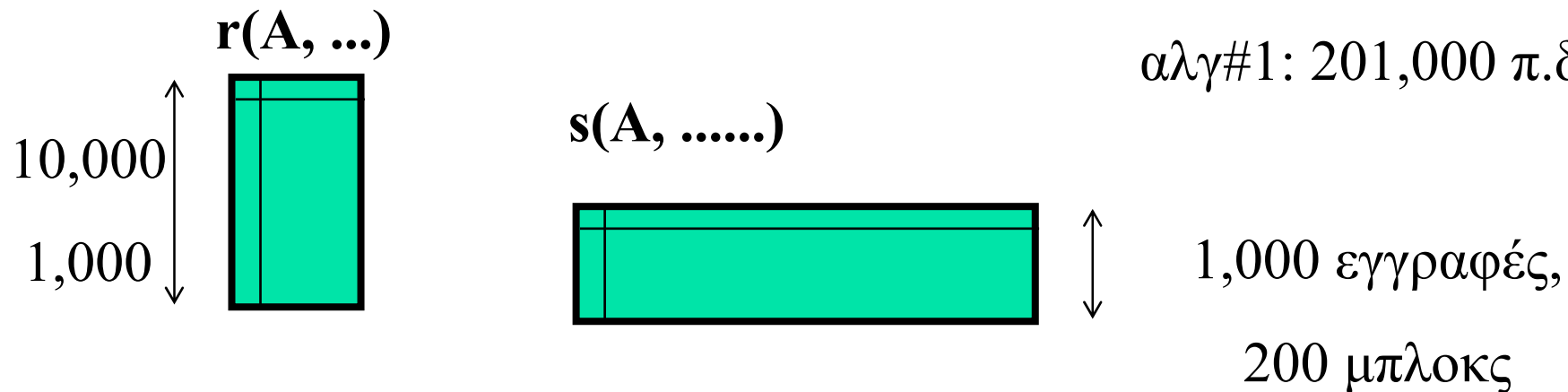


# 2-way joins

- Αριθμητικό παράδειγμα:
  - $nr = 10,000$  πλειάδες,  $br = 1,000$  μπλοκς
  - $ns = 1,000$  πλειάδες,  $bs = 200$  μπλοκς

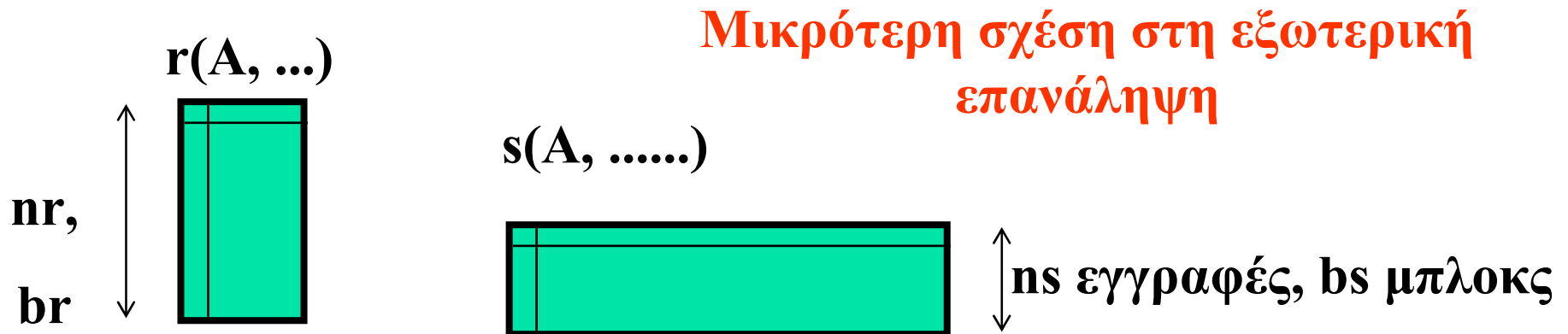
αλγ#0: 2,001,000 π.δ.

αλγ#1: 201,000 π.δ.



# 2-way joins

- Παρατήρηση 1: Αλγ#1: ασύμετρος:
  - κόστος:  $br + br * bs$  – ανάστροφοι ρόλοι:
  - κόστος =  $bs + bs * br$
- Καλύτερη επιλογή?



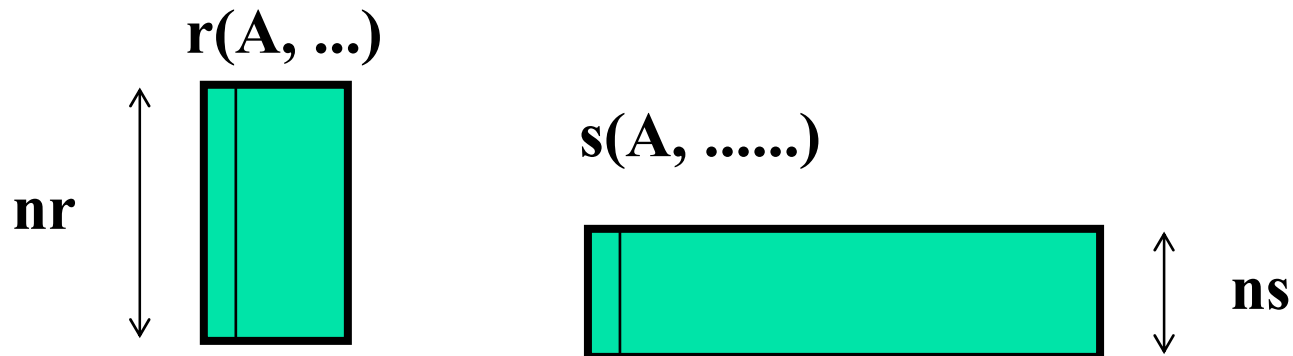


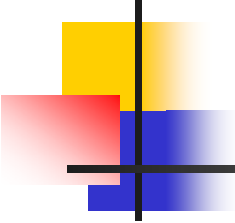


## 2-way joins

---

- Άλλοι αλγόριθμοι για  $r \text{ JOIN } s$ ;
- $n_r$ ,  $n_s$  πλειάδες το καθένα



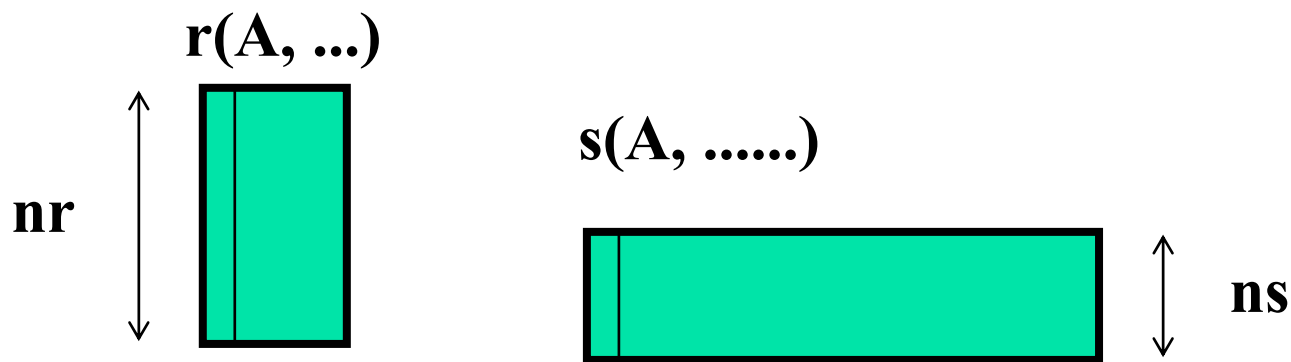


# 2-way joins - other algo's

---

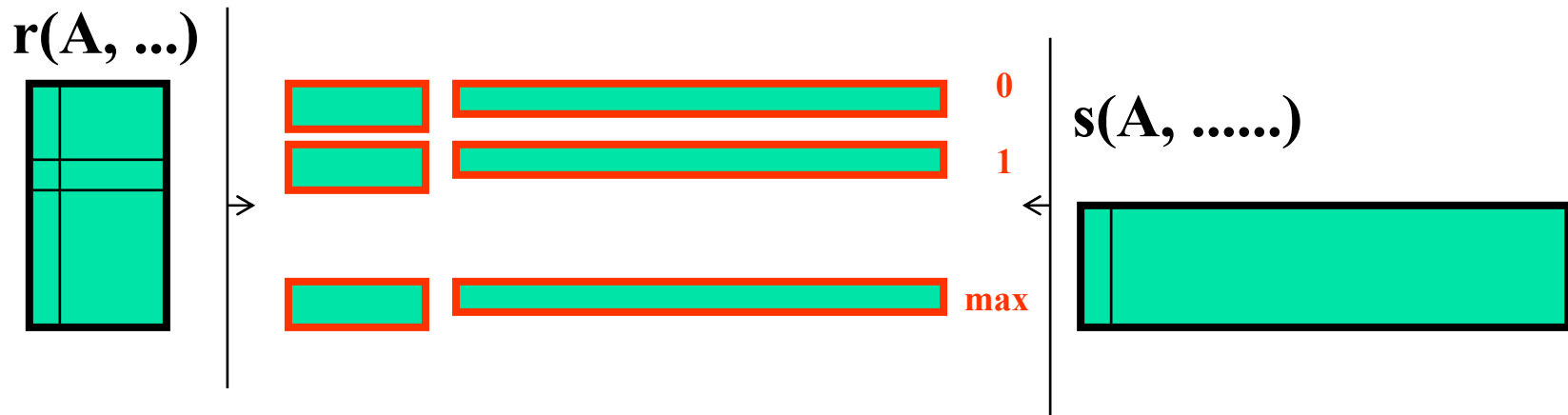
## ■ **sort-merge**

- Ταξινόμησε 'r'; Ταξινόμησε 's'; Merge ταξινομημένες εκδόσεις των r και s.
- Καλό, αν η μια ή και οι δύο σχέσεις είναι
- ήδη ταξινομημένες



# 2-way joins - other algo's

- hash join:
  - hash 'r' into (0, 1, ..., 'max') buckets
  - hash 's' into buckets (same hash function)
  - join each pair of matching buckets





# Δομή των βελτιστοποιητών επερωτήσεων:

---

Περισσότερα heuristics από Oracle,  
Sybase and Starburst (-> DB2) : στο  
βιβλίο

Γενικά: Βελτιστοποίηση Επερωτήσεων:  
Είναι πολύ σημαντική για μεγάλες  
βάσεις δεδομένων

(**explain select** <sql-statement>  
δείχνει το πλάνο)



# Συμπεράσματα – βήματα βελτ.Επερ.

---

- φέρτε το ερώτημα στην εσωτερική μορφή (πχ., parse tree)
- ... στην 'κανονική μορφή' (συντακτική Βελτ.Ερωτ.)
- δημιουργήστε εναλλακτικά σχέδια
  - επιλογές (απλές; σύνθετα κατηγορήματα)
  - ταξινόμηση; προβολές
  - συνενώσεις
- εκτιμήστε το κόστος, επιλέξτε το καλύτερο