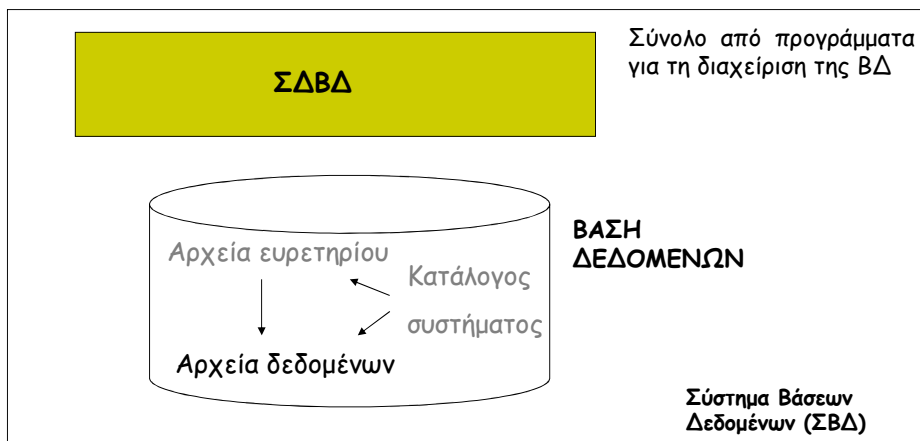




Εισαγωγή στην Επεξεργασία Ερωτήσεων

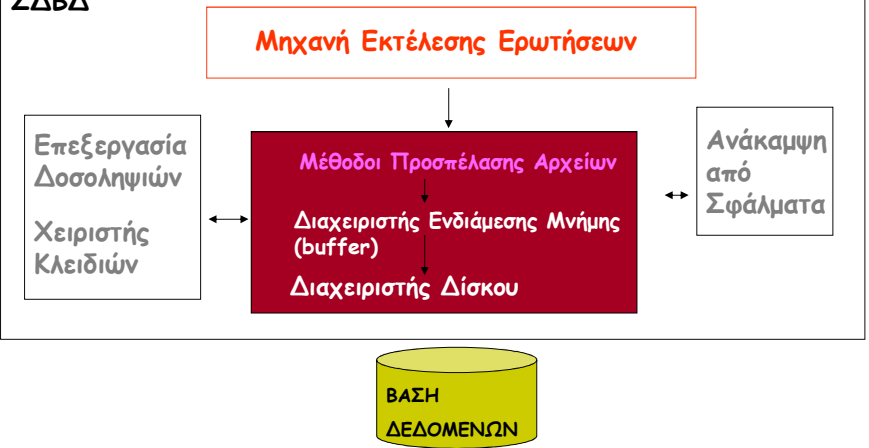


Εισαγωγή

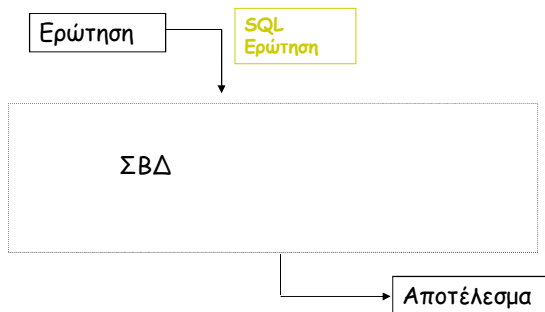


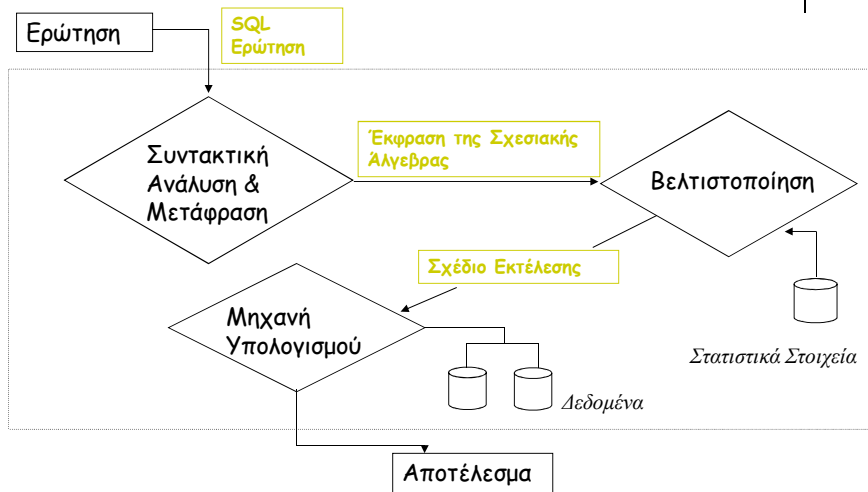


ΣΔΒΔ



Θα δούμε την «πορεία» μιας SQL ερώτησης (πως εκτελείται)





Τα βασικά βήματα στην επεξεργασία μιας ερώτησης είναι

1. Συντακτική Ανάλυση & Μετάφραση
2. Βελτιστοποίηση
3. Υπολογισμός



1. Συντακτική Ανάλυση (Parsing) & Μετάφραση

Η SQL ερώτηση μεταφράζεται σε μια εσωτερική μορφή αφού γίνει ο απαραίτητος συντακτικός και σημασιολογικός έλεγχος (π.χ., τα ονόματα που αναφέρονται είναι ονόματα σχέσεων που υπάρχουν)

Αντικατάσταση των όψεων από τον ορισμό τους

Σε ποια εσωτερική μορφή; Έκφραση της σχεσιακής άλγεβρας

```
select A1, A2, ..., An
from R1, R2, ..., Rm      πA1, A2, ..., An (σP (R1 × R2 × ... × Rm))
where P
```



2. Βελτιστοποίηση

Μια SQL ερώτηση μπορεί να μεταφραστεί σε διαφορετικές (ισοδύναμες) εκφράσεις της σχεσιακής άλγεβρας

```
select balance
from account
where balance < 25000
```

- π_{balance} (σ_{balance < 2500} (account))
- σ_{balance < 2500} (π_{balance} (account))

Με ποιο κριτήριο γίνεται η επιλογή της έκφρασης:

- το πιο «δύσκολο» βήμα



Κάθε πράξη της σχεσιακής άλγεβρας μπορεί να υλοποιηθεί με *διαφορετικούς αλγορίθμους*:

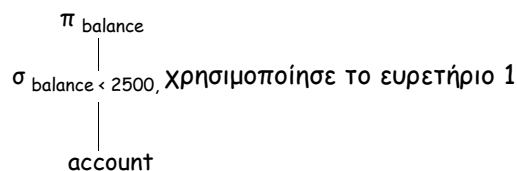
π.χ., για την υλοποίηση της επιλογής μπορεί
 είτε να σαρώσουμε (scan) όλο το αρχείο ελέγχοντας
 κάθε εγγραφή αν ικανοποιεί τη συνθήκη
 είτε αν υπάρχει π.χ., ένα B⁺ ευρετήριο στο γνώρισμα
 balance να χρησιμοποιήσουμε το ευρετήριο

Άρα δεν αρκεί ο προσδιορισμός της πράξης - πρέπει να προσδιορίζεται και ο αλγόριθμος που θα χρησιμοποιηθεί για την υλοποίησή της



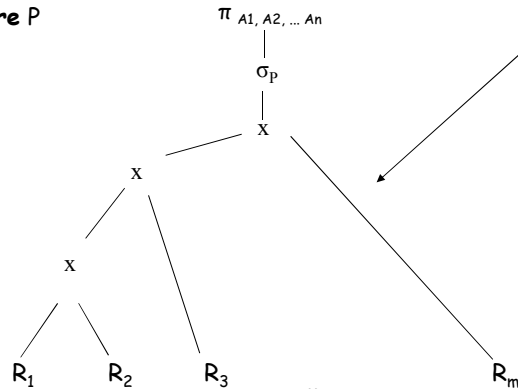
βασικές (primitive) πράξεις: πράξη + αλγόριθμος

Σχέδιο εκτέλεσης (execution plan): μια ακολουθία από βασικές πράξεις





select A_1, A_2, \dots, A_n **Μετάφραση**
from R_1, R_2, \dots, R_m \longrightarrow $\pi_{A_1, A_2, \dots, A_n} (\sigma_P (R_1 \times R_2 \times \dots \times R_m))$
where P



Πλάνο εκτέλεσης (ποιες πράξεις και με ποιον αλγόριθμο)

Φύλλα: σχέσεις

Εσωτερικοί κόμβοι: βασικές πράξεις της σχεσιακής άλγεβρας

Βελτιστοποίηση του πλάνου



- Τα διαφορετικά σχέδια εκτέλεσης έχουν και διαφορετικό κόστος
- **Βελτιστοποίηση:** η διαδικασία επιλογής του σχεδίου εκτέλεσης που έχει το μικρότερο κόστος
- Εκτίμηση του κόστους (συνήθως χρήση στατιστικών στοιχείων)



Μερικοί ευριστικοί κανόνες

Γενική ιδέα: εκτέλεση πρώτα των πράξεων με μικρή επιλεκτικότητα ώστε να περιοριστεί το μέγεθος των ενδιάμεσων αποτελεσμάτων

1. Διάσπαση των πράξεων επιλογής με συζευκτικές συνθήκες σε ακολουθίες πράξεων επιλογής
2. Μετατοπίζουμε την *πράξη επιλογής όσο πιο κάτω* επιτρέπεται από τα γνωρίσματα που περιλαμβάνονται στη συνθήκη
3. Επαναδιευθέτηση των φύλλων ώστε να εκτελούνται πρώτα οι σχέσεις που έχουν τις πιο περιοριστικές πράξεις επιλογής



4. Συνδυασμός μιας πράξης καρτεσιανού γινομένου με μια πράξη επιλογής που ακολουθεί
5. Διάσπαση και *μετακίνηση των λιστών προβολής όσο πιο κάτω* γίνεται στο δέντρο
6. Εντοπισμός υποδέντρων με ομάδες πράξεων που μπορεί να εκτελεστούν με κοινό αλγόριθμο



3. Εκτέλεση

Μηχανή εκτέλεσης που εκτελεί τις βασικές πράξεις



Υπάρχουν υλοποιημένοι μια σειρά από αλγόριθμοι για κάθε βασική πράξη (π.χ., που χρησιμοποιούν ή όχι ευρετήρια κλπ)

Γενικά, το ΣΔΒΔ με βάση κάποια *στατιστικά στοιχεία* κάνει μια *εκτίμηση του κόστους* και *επιλέγει τον αλγόριθμο* για κάθε πράξη με τον μικρότερο (με βάση την εκτίμηση) κόστος

Αποθηκεύονται διάφορα στατιστικά στοιχεία, τα οποία χρησιμοποιούνται για την αποτίμηση του κόστους

Αλγόριθμοι Εκτέλεσης Βασικών Πράξεων



Για να επιλέξουμε ποιόν αλγόριθμο θα χρησιμοποιήσουμε, διατηρούμε στατιστικά στοιχεία

Για ένα *αρχείο δεδομένων* μιας σχέσης R:

- n_R : αριθμός πλειάδων της σχέσης R
- b_R : αριθμός blocks της σχέσης R
- s_R : μέγεθος σε bytes κάθε πλειάδας της σχέσης R
- f_R : παράγοντας ομαδοποίησης (αριθμός εγγραφών ανά block)

αν μη εκτεινόμενη, $f_R = \lfloor B / s_R \rfloor$ και $b_R = \lceil n_R / f_R \rceil$

Ενημέρωση στατιστικών στοιχείων;

Αλγόριθμοι Εκτέλεσης Βασικών Πράξεων



Άλλα στατιστικά στοιχεία;

- $V(A, R)$: αριθμός διαφορετικών τιμών που παίρνει το γνώρισμα A
 $|\pi_A(R)|$ -- αν το A κλειδί;
- $SC(A, R)$: μέσος αριθμός πλειάδων που ικανοποιεί μια συνθήκη (δεδομένου ότι υπάρχει μια τουλάχιστον που την ικανοποιεί)
1 αν κλειδί, αν ομοιόμορφη;

Αλγόριθμοι Εκτέλεσης Βασικών Πράξεων



Στατιστικά στοιχεία επίσης για το *αρχείο ευρετηρίου* (αν υπάρχει)

- f_i : παράγοντας διακλάδωσης,
πολυεπίπεδο f_0 , B^* δέντρο \sim τάξη
- H_i : αριθμός επιπέδων
- LB_i : αριθμός block φύλλων

Με βάση τα στατιστικά επιλέγεται ο αλγόριθμος με το μικρότερο κόστος
I/O Κόστος (Αριθμό blocks που μεταφέρονται)

Επεξεργασία Ερωτήσεων



Αλγόριθμοι εκτέλεσης βασικών πράξεων

Επιλογή

*Διαφορετικοί αλγόριθμοι ανάλογα με το αν το αρχείο είναι ή όχι
διατεταγμένο, αν υπάρχει ή όχι ευρετήριο και από το είδος του
ευρετηρίου*



Πιθανοί αλγόριθμοι εκτέλεσης για την **επιλογή**:

E1: Σειριακή αναζήτηση

E2: Δυαδική αναζήτηση (αν το αρχείο είναι ταξινομημένο)

E3: Χρήση πρωτεύοντος ευρετηρίου/κατακερματισμού (αν υπάρχει)

E4: Χρήση δευτερεύοντος ευρετηρίου/κατακερματισμού (αν υπάρχει)

Αν υπάρχει κάποιο ευρετήριο, λέμε ότι έχουμε μονοπάτι προσπέλασης (access path)



Επιλεκτικότητα επιλογής:

το πλήθος των εγγραφών (πλειάδων) που επιλέγονται (δηλ. ικανοποιούν την συνθήκη)

το πλήθος των εγγραφών (πλειάδων) του αρχείου (σχέσης)

• Έστω $s_i = |\sigma_{\theta_i}(R)|$

επιλεκτικότητα: s_i / n_R

Αν θ_i **συνθήκη ισότητας** σε ένα γνώρισμα υποψήφιο κλειδί $s_i = 1 / n_R$

Αν θ_i **συνθήκη ισότητας** σε ένα γνώρισμα, ομοιόμορφη κατανομή, k διακριτές τιμές, $s_i = k / n_R$



Επιλογή - συνθήκη ισότητας $\sigma_A = \alpha(R)$

E1 Σειριακή αναζήτηση

Διάβασμα (scan) όλου του αρχείου

b_R : αριθμός blocks της σχέσης R

b_R

$b_R/2$ (μέσος όρος) αν το A υποψήφιο κλειδί
(οπότε το αποτέλεσμα έχει μόνο μία πλειάδα, σταματάμε την αναζήτηση μόλις τη βρούμε)

Μπορεί να χρησιμοποιηθεί σε οποιοδήποτε αρχείο



E2 Δυαδική αναζήτηση

b_R : αριθμός blocks της σχέσης R
 $SC(A, R)$: μέσος αριθμός πλειάδων που ικανοποιεί μια συνθήκη
 f_R : παράγοντας ομαδοποίησης

Μπορεί να χρησιμοποιηθεί μόνο αν το αρχείο είναι διατεταγμένο με βάση το A (δηλαδή, το γνώρισμα της επιλογής)

$$\lceil \log(b_R) \rceil \quad \leftarrow \quad \text{Εύρεση της πρώτης}$$

$$+ \quad \lceil SC(A, r)/f_R \rceil - 1 \quad \leftarrow \quad \text{Εύρεση των υπόλοιπων}$$

Αν το A υποψήφιο κλειδί:



E3 Χρήση πρωτεύοντος

Επιλογή: Συνθήκη Ισότητας

Δεντρικού ευρετηρίου

Πρωτεύον ευρετήριο
σημαίνει ταξινομημένο
αρχείο

b_R : αριθμός blocks της σχέσης R
 $SC(A, R)$: μέσος αριθμός πλειάδων που
ικανοποιεί μια συνθήκη
 f_R : παράγοντας ομαδοποίησης
 HT_i : αριθμός επιπέδων

Μπορεί να χρησιμοποιηθεί μόνο αν υπάρχει τέτοιο ευρετήριο στο A

$HT_i + 1$ ← Εύρεση και μεταφορά της
πρώτης

Αν το A δεν είναι υποψήφιο κλειδί -- ευρετήριο συστάδων

$HT_i + \lceil SC(A, R) / f_R \rceil$ ← Εύρεση και των υπόλοιπων

ΣΗΜΕΙΩΣΗ: Πρωτεύον ευρετήριο στο A , σημαίνει ότι οι εγγραφές του αρχείου δεδομένων είναι ταξινομημένες (διατεταγμένες) ως προς A άρα οι υπόλοιπες εγγραφές με την ίδια τιμή (αν υπάρχουν) βρίσκονται σε γειτονικά blocks του αρχείου δεδομένων



Επιλογή: Συνθήκη Ισότητας

E4 Χρήση δευτερεύοντος

Δεντρικού ευρετηρίου

b_R : αριθμός blocks της σχέσης R
 $SC(A, R)$: μέσος αριθμός πλειάδων που
ικανοποιεί μια συνθήκη
 f_R : παράγοντας ομαδοποίησης
 HT_i : αριθμός επιπέδων

Μπορεί να χρησιμοποιηθεί μόνο αν υπάρχει τέτοιο ευρετήριο στο A

Αν το A είναι υποψήφιο κλειδί

$HT_i + 1$ ← Εύρεση και μεταφορά της πρώτης

Αν το A δεν είναι υποψήφιο κλειδί \pm κόστος για την εύρεση των
υπολοίπων

$HT_i + \text{ενδιάμεσο επίπεδο}$

$+SC(A, R)$ ← Εύρεση και των υπόλοιπων

Στη χειρότερη περίπτωση κάθε εγγραφή που ικανοποιεί τη συνθήκη σε διαφορετικό block



Επιλογή - συνθήκη με σύγκριση

$$\sigma_{A \leq u}(\mathbf{R}) \text{ ή } \sigma_{A \geq u}(\mathbf{R})$$



Επιλογή - συνθήκη με σύγκριση

$$\sigma_{A \leq u}(\mathbf{R})$$

Έστω αύξουσα διάταξη

Από το 1^ο block του αρχείου έως την πρώτη εγγραφή με $A > u$

Κόστος?



Επιλογή - συνθήκη με σύγκριση

$$\sigma_{A \leq u}(\mathcal{R})$$

Έστω αρχείου σωρού (δεν υπάρχει διάταξη) και B+ δέντρο

Εύρεση στο B+ δέντρο της τιμής u

Χρήση εγγραφών στο φύλλο για τις υπόλοιπες τιμές

Κόστος?

Τέλος