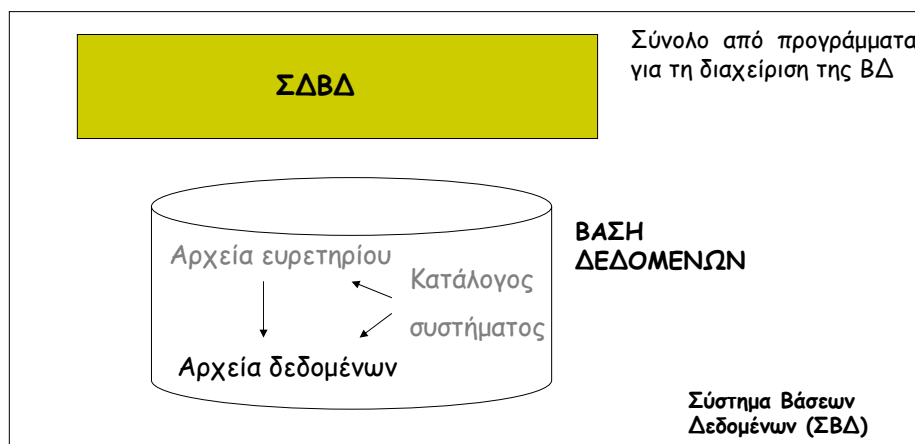
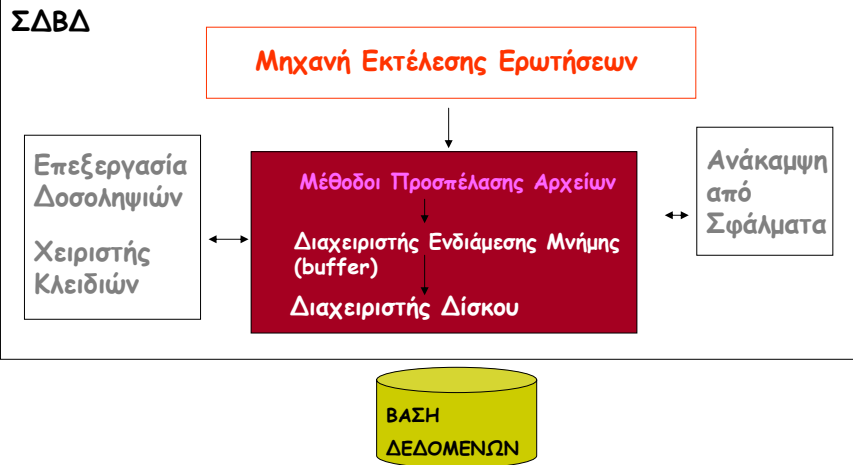




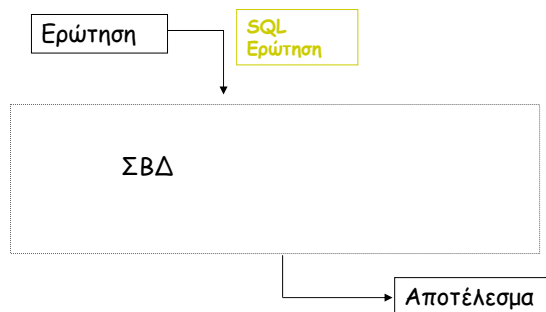
Εισαγωγή στην Επεξεργασία Ερωτήσεων

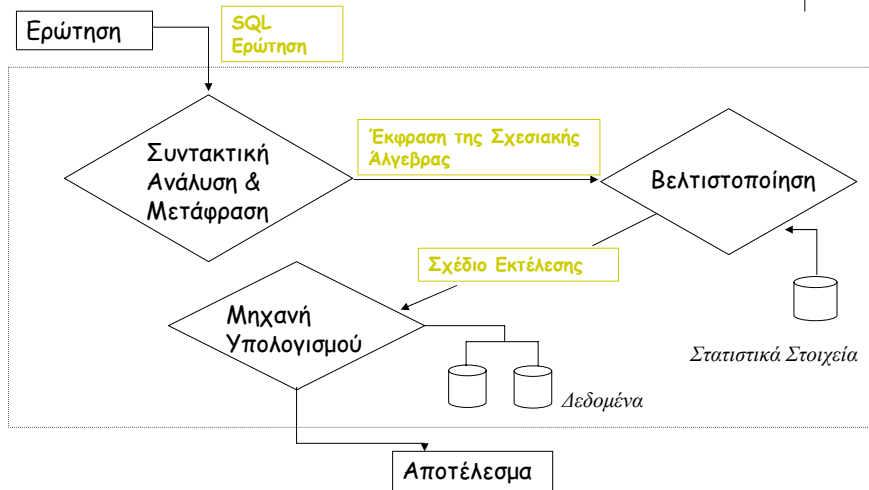
Εισαγωγή





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;



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

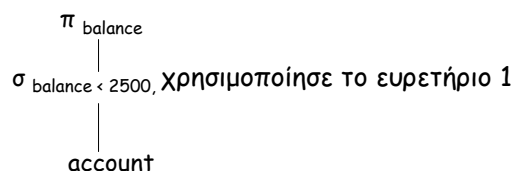
π.χ., για την υλοποίηση της επιλογής μπορεί
 είτε να σαρώσουμε (scan) όλο το αρχείο ελέγχοντας
 κάθε εγγραφή αν ικανοποιεί τη συνθήκη
 είτε αν υπάρχει π.χ., ένα B⁺ ευρετήριο στο γνώρισμα
 balance να χρησιμοποιήσουμε το ευρετήριο

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



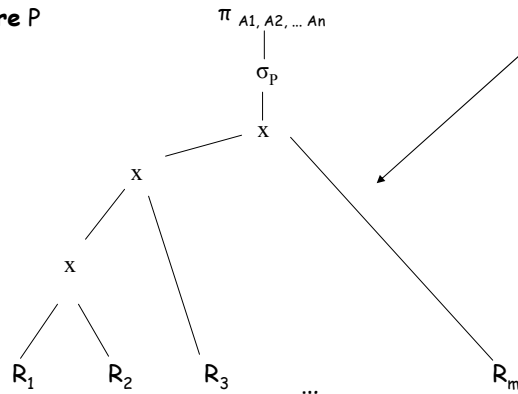
βασικές (primitive) πράξεις: πράξη + αλγόριθμος

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 διαφορεικό 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



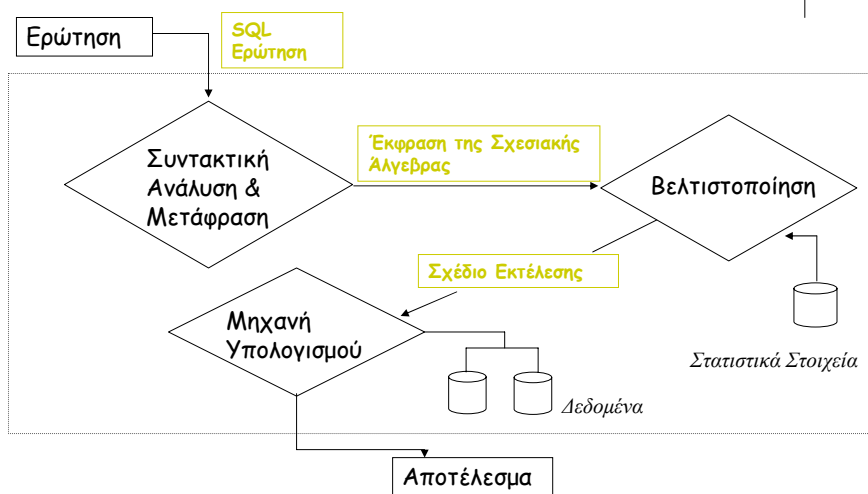
Υπάρχουν υλοποιημένοι μια σειρά από αλγόριθμοι για κάθε βασική πράξη (π.χ., που χρησιμοποιούν ή όχι ευρετήρια κλπ)

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



Αλγόριθμους εκτέλεσης βασικών πράξεων Επιλογή

Επεξεργασία Ερωτήσεων (ανακεφαλαίωση)



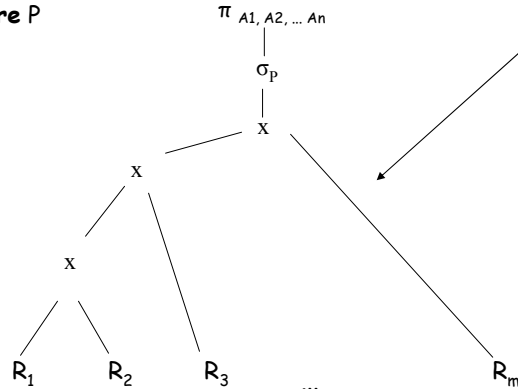
Επεξεργασία Ερωτήσεων (ανακεφαλαίωση)



select A_1, A_2, \dots, A_n
 from R_1, R_2, \dots, R_m
 where P

Μετάφραση

$\pi_{A_1, A_2, \dots, A_n} (\sigma_P (R_1 \times R_2 \times \dots \times R_m))$



Πλάνο εκτέλεσης (ποιες πράξεις και με ποιον αλγόριθμο)

Βελτιστοποίηση

Εκτέλεση

Επιλογή κατάλληλου αλγορίθμου για κάθε βασική πράξη της σχεσιακής άλγεβρας

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;

Π.χ., για μια πράξη επιλογής στο γνώρισμα A

- $V(A, R)$: αριθμός διαφορετικών τιμών του 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 ~ 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)



Επιλογή

Θα εξετάσουμε:

- Επιλογή με συνθήκη ισότητας ($\sigma_{A = a} (R)$)
- Επιλογή με συνθήκη σύγκρισης - διαστήματος/περιοχής (range query) ($\sigma_{A \leq u} (R)$) ή ($\sigma_{A \geq u} (R)$)
- Επιλογή με σύζευξη ($\sigma_{\theta_1 \text{ AND } \theta_2 \dots \text{ AND } \theta_n} (R)$)
- Επιλογή με διάζευξη ($\sigma_{\theta_1 \text{ OR } \theta_2 \dots \text{ OR } \theta_n} (R)$)

Θα δούμε ποιος είναι ο καλύτερος (με το μικρότερο κόστος σε blocks) αλγόριθμος για την εκτέλεση της πράξης



Πιθανοί αλγόριθμοι εκτέλεσης για την **επιλογή**:

E1: Σειριακή αναζήτηση

E2: Δυαδική αναζήτηση

E3: Χρήση πρωτεύοντος ευρετηρίου/κατακερματισμού

E4: Χρήση δευτερεύοντος ευρετηρίου/κατακερματισμού

Για τον E2 πρέπει το αρχείο να είναι ταξινομημένο

Για τους E3 και E4 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

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

Μπορεί να χρησιμοποιηθεί σε οποιοδήποτε αρχείο



Επιλογή: Συνθήκη Ισότητας

E2 Δυναδική 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 γνώρισμα της επιλογής

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

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

Αν το A υποψήφιο κλειδί;



Επιλογή: Συνθήκη Ισότητας

E3 Χρήση πρωτεύοντος (πολυεπίπεδου) 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

$$HT_i + 1 \quad \longleftarrow \quad \text{Εύρεση και μεταφορά της πρώτης}$$

Αν το A δεν είναι υποψήφιο κλειδί -- ευρετήριο συστάδων

$$HT_i + \lceil SC(A, R)/f_R \rceil \quad \longleftarrow \quad \text{Εύρεση και των υπόλοιπων}$$

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 (πολυεπίπεδου) 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

Επιλογή - συνθήκη με σύγκριση

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

Έστω ότι c πλειάδες ικανοποιούν τη συνθήκη

Γενικά $c = n_R/2$ (δηλαδή, οι μισές)

Έστω \min , \max (μικρότερη, μεγαλύτερη τιμή του A), αν ομοιόμορφη κατανομή και $\sigma_{A \leq u}(R)$

$$c = \begin{cases} 0 & \text{αν } u < \min \\ n_R & \text{αν } u \geq \max \\ n_R * [(u - \min) / (\max - \min)] & \end{cases}$$



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

Θα δούμε δύο αλγορίθμους:

E5 Χρήση πρωτεύοντος πολυ-επίπεδου ευρετηρίου

E6 Χρήση δευτερεύοντος πολυ-επίπεδου ευρετηρίου



E5 Χρήση πρωτεύοντος (πολυεπίπεδου) ευρετηρίου

*Πρωτεύον, σημαίνει ταξινομημένο αρχείο,
έστω σε αύξουσα διάταξη*

$$A \geq u$$

1. Χρήση ευρετηρίου για την εύρεση της πρώτης εγγραφής $A \geq u$
2. Σάρωση όλου του αρχείου ξεκινώντας από αυτήν την εγγραφή

$$HT_i + \lceil c / f_R \rceil$$

$$A \leq u$$

Δε χρειάζεται ευρετήριο, γιατί:

*c: επιλεξιμότητα (πλειάδες που
ικανοποιούν την συνθήκη)
f_R: παράγοντας ομαδοποίησης
HT_i: αριθμός επιπέδων*



Επιλογή: Συνθήκη με Σύγκριση

Ε6 Χρήση δευτερεύοντος (πολυεπίπεδου) ευρετηρίου

- Εύρεση του πρώτου φύλλου του ευρετηρίου *Μη ταξινομημένο αρχείο*
- Για κάθε block (φύλλο) του ευρετηρίου διάβασε το αντίστοιχο block δεδομένων *(σημείωση, αν B+ δέντρο χρησιμοποιούμε το δείκτη ανάμεσα στα φύλλα)*

Σάρωση των φύλλων του δέντρου

$A \leq u$ από την αρχή έως το u

$A \geq u$ από το u έως το τέλος

c : επιλεξιμότητα (πλειάδες που ικανοποιούν την συνθήκη)
 n_R : αριθμός εγγραφών
 LB_i : αριθμός block φύλλων
 HT_i : αριθμός επιπέδων

Π.χ., αν $c = n_R/2$, και A κλειδί τότε (αν κάθε εγγραφή σε διαφορετικό block)

$$HT_i + LB_i/2 + n_R/2$$



Επιλογή: Συνθήκη Σύζευξης

Επιλογή - συνθήκη σύζευξης

$$\sigma_{\theta_1 \text{ AND } \theta_2 \dots \text{ AND } \theta_n (R)}$$

Επιλεκτικότητα μιας συνθήκης:

το πλήθος των εγγραφών (πλειάδων) που την ικανοποιούν
 το πλήθος των εγγραφών (πλειάδων) του αρχείου (σχέσης)

- Αν οι συνθήκες είναι ανεξάρτητες, το μέγεθος του αποτελέσματος:

$$\frac{n_R * s_1 * s_2 * \dots * s_n}{n_R^n}$$

s_i επιλεκτικότητα της θ_i



E7 Συζευκτική επιλογή με χρήση ενός απλού ευρετηρίου

Υπάρχει διαδρομή προσπέλασης για ένα από τα γνωρίσματα που εμφανίζονται σε οποιαδήποτε απλή συνθήκη

Επιλογή του γνωρίσματος στην απλή συνθήκη με τη **μικρότερη** επιλεκτικότητα (γιατί:)

Χρήση μιας από τις προηγούμενες μεθόδους για την ανάκτηση των εγγραφών που ικανοποιούν αυτήν την συνθήκη και

έλεγχος για κάθε επιλεγμένη εγγραφή αν ικανοποιεί και τις υπόλοιπες συνθήκες



E8 Συζευκτική επιλογή με χρήση σύνθετου ευρετηρίου

Αν υπάρχει ευρετήριο στο συνδυασμό δύο ή περισσότερων γνωρισμάτων που εμφανίζονται σε οποιαδήποτε απλές συνθήκες



E9 Συζευκτική επιλογή με τομή δεικτών

Αν υπάρχουν ευρετήρια σε περισσότερα από ένα από τα γνωρίσματα

Τότε διαβάζουμε τα blocks του αρχείου δεδομένων που δίνονται από όλα τα ευρετήρια



Επιλογή - συνθήκη διάζευξης

$$\sigma_{\theta_1 \text{ OR } \theta_2 \dots \text{ OR } \theta_n} (R)$$

Αν έστω και μία από τις συνθήκες δεν έχει διαδρομή προσπέλασης -> σάρωση όλου του αρχείου



Τέλος!