

2η σειρά ασκήσεων, Σχεδιασμός βάσεων Δεδομένων

p3170122 | Κων/νος Νικολούτσος

Άσκηση 1

1.

Ο αριθμός των *sorted sublist* μετά την πρώτη φάση είναι:

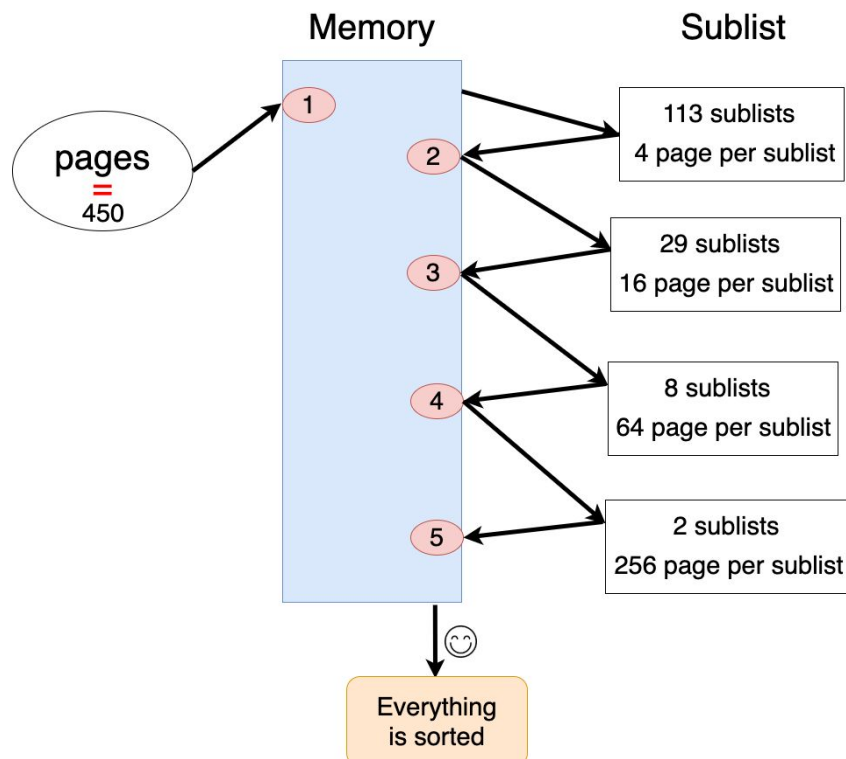
$$\text{Pages} = 45.000/100 = 450$$

$$\text{Sorted_sublists} = \text{Pages}/(\text{buffering_size}) = \text{ceil}(112,5) = 113$$

2.

*(Πέρασμα θεωρώ πόσες φορές φορτώνονται όλα τα δεδομένα στην RAM του DBMS)

Θα χρειαστούν **5 περάσματα**. (Δείτε την παρακάτω εικόνα για μια πιο visual απάντηση)



3.

Είναι κατανοητό και από την παραπάνω εικόνα ότι το κόστος είναι όσα και τα βελάκια μείον ένα:

$$\text{Cost} = 9 * B(R) = 4050 \text{ I/Os}$$

4.

Για να γίνουν μόνο δύο περάσματα αρκεί να ισχύει η παρακάτω συνθήκη (από τα slides του μαθήματος).

$$M \geq \sqrt{B(R)} \Leftrightarrow M > 21.213 \Leftrightarrow \mathbf{M \geq 22}$$

Επιπλέον σε αυτήν την περίπτωση που μεγάλωνε το M το κόστος θα ήταν $3 * B(R)$

Άσκηση 2

1.

Ο πιο αποδοτικός αλγόριθμος για το παράδειγμα είναι ο **Hash-Join** διότι έχει το μικρότερο I/O cost.

NLJ:

Το κόστος είναι: **1290 I/Os**

Επεξήγηση: Χρησιμοποιούμε ως εξωτερική σχέση εκείνη με τις περισσότερες σελίδες. Άρα ως εξωτερική σχέση έχουμε την R_2
 $B(R_2) + \lceil B(R_2)/(M-1) \rceil * B(R_1) = 1140 \text{ I/Os}$

SMJ:

Το κόστος είναι: **1200 I/Os**

Επεξήγηση: Το μέγεθος του ενταμιευτή μνήμης είναι μικρό για να εφαρμοστεί η αποδοτική μορφή του αλγορίθμου. Αυτό συμβαίνει διότι:

$\text{ceil}((B(R1) + B(R2))/M) = 19 > M$ (Δεν χωράνε στην μνήμη)

Οπότε θα χρησιμοποιήσουμε κανονική έκδοση του αλγορίθμου και θα έχουμε κόστος:
 $5(B(R1) + B(R2)) = 1200$ I/Os

Επίσης να τονίσουμε ότι το μέγεθος του ενταμιευτή μνήμης είναι αρκετό για την εφαρμογή της κανονικής έκδοσης του αλγορίθμου διότι ισχύει:
 $M > \text{sqrt}(\max(B(R1), B(R2))) \Leftrightarrow 13 > 12,24$

HJ:

Το κόστος είναι: **720 I/Os**

Επεξήγηση: Το μέγεθος του ενταμιευτή μνήμης είναι μικρό για να εφαρμοστεί η βέλτιστη μορφή του αλγορίθμου (Καμία από τις δύο σχέσεις δεν χωράει ολοκληρή στην μνήμη). $M < \min(B(R1), B(R2))$

Ωστόσο το μέγεθος του ενταμιευτή είναι αρκετό για να εφαρμόσουμε hash-join στην μνήμη. $M > \text{sqrt}(\min(B(R1), B(R2)))$

Οπότε το κόστος θα είναι:
 $3(B(R1) + B(R2)) = 720$ I/Os

2.

Αν αυξήσουμε το M σε 100 σελίδες.

Ακολουθώντας την ίδια λογική προκύπτει ότι:

NLJ:

#pages	M = 13	M = 100
#I/Os	1140	240

Επεξήγηση: Η έκδοση δεν αλλάζει αλλά κερδίζουμε από το γεγονός ότι φορτώνουμε περισσότερα pages από το R1. Το οποίο σημαίνει ότι θα φορτώσουμε λιγότερες φορές σελίδες της σχέσης R2.

SMJ:

#pages	M = 13	M = 100
#I/Os	1200	720

Επεξήγηση: Χρησιμοποιούμε την καλύτερη έκδοση του SMJ διότι:
 $M > \sqrt{B(R1) + B(R2)} \Leftrightarrow 100 > 15.49$

HJ:

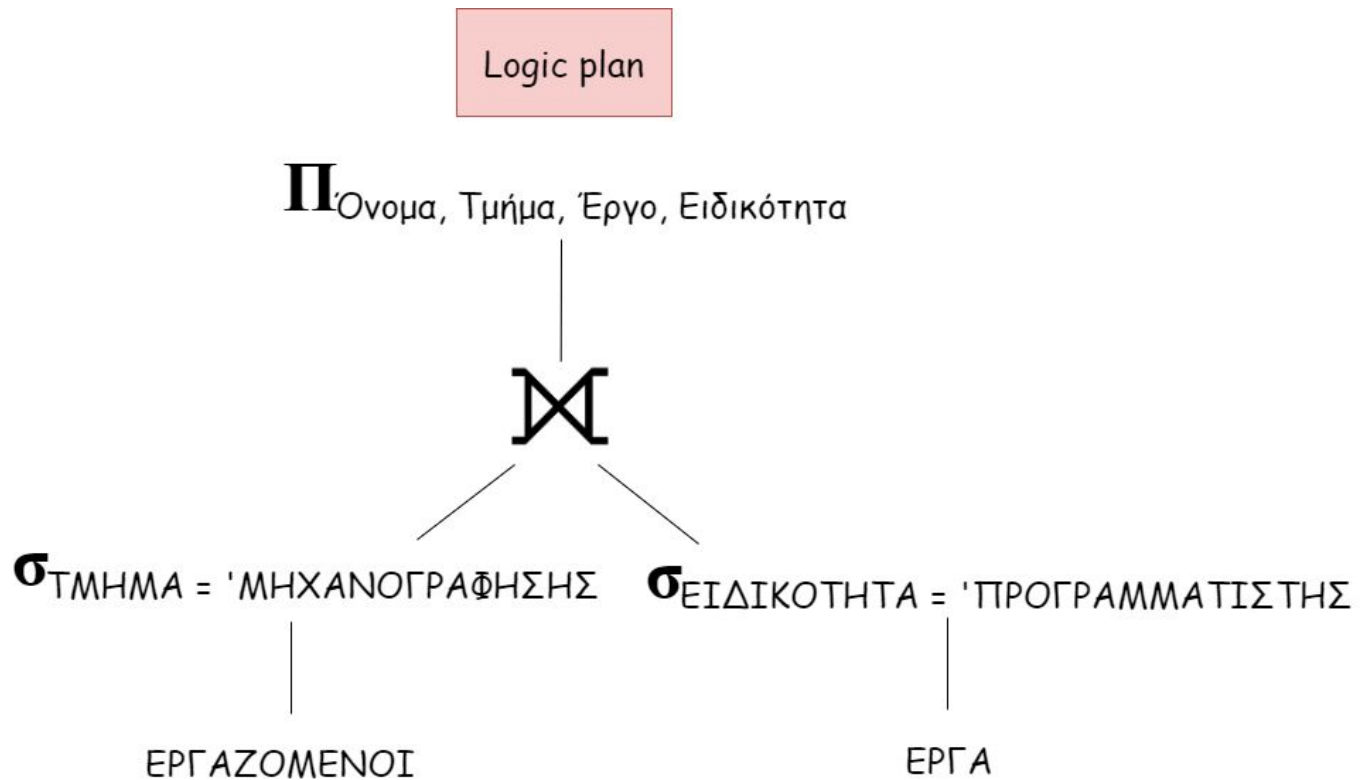
#pages	M = 13	M = 100
#I/Os	720	240

Επεξήγηση: Το μέγεθος του ενταμιευτή μνήμης είναι αρκετό για να χωρέσει η σχέση R2 άρα το κόστος ανέρχεται στο: $B(R1) + B(R2)$

Άσκηση 3

1.

Για την παραγωγή ενός καλού λογικού πλάνου έγινε χρήση της push down selection καλής πρακτικής



2.

α)

*Χρησιμοποιήσαμε οτι τα δεδομένα είναι ομοιόμορφα άρα η ειδικοτητα = 'προγραμματιστής' θα είναι 10 σελίδες.

*Χρησιμοποιήσαμε οτι τα δεδομένα είναι ομοιόμορφα άρα η τμημα = μηχανογραφησης θα είναι 25 σελίδες.

Ας υπολογίσουμε το κόστος ακολουθώντας bottom-up approach.

Πρέπει να γίνουν τα εξής βήματα:

- Θα κάνουμε two-phase-shorting και Θα κανουμε SMJ
- Γραφουμε τα αποτελέσματα

Τα κόστη των παραπάνω είναι:

- $3(B(\text{Εργαζομενοι}) + B(\text{ΕΡΓΑ})) = 105 \text{ I/Os}$
- 0 (Μπορει να μην βγάξει κανενα αποτελεσμα στην καλυτερη περίπτωση)

Αρα το ελάχιστο κόστος (lower bound) του παραπάνου φυσικού πλάνου είναι **105 I/Os**

β) Στην άλλη περίπτωση το θα πρέπει να κάνουμε τα βήματα:

Θα εφαρμόσουμε την καλύτερη μορφή του αλγορίθμου αφήνοντας μονο μια θέση για την εσωτερική σχέση.

$$B(S) + \lceil B(S)/(M-1) \rceil * B(R)$$

Εφαρμόζοντας λοιπον το κόστος απο τα slides του μαθήματος προκύπτει ότι:

Το κόστος είναι: **60 I/Os**

Παρατήρηση: Από την άλλη πλευρά όμως ο SMJ θα καταναλώνει λιγότερη ισχύ στο cpu.

Άσκηση 4

α)

1 κόστος:

Εφόσον γίνεται το table scan το DBMS θα διαβάσει όλα τα blocks. Άρα θα έχει κόστος **100 I/Os**. Όπου υποθέτοντας ομοιομορφία δεδομένων θα κρατήσει 100 εγγραφές δηλαδή θα βρίσκονται σε 5 pages ($5 \times 20 = 100$).

2 κόστος:

Ομοίως, γίνεται το table scan και το DBMS θα διαβάσει όλα τα blocks. Άρα θα έχει κόστος **10 I/Os**. Υποθέτοντας πάλι ομοιομορφία δεδομένων το DBMS θα κρατήσει μόνο 9 εγγραφές, άρα 1 σελίδα.

3 κόστος:

Το κόστος εδώ είναι 0.

Τα προηγούμενα operators φέραν στην κύρια μνήμη $5+1 = 6$ σελίδες που χωράνε. Άρα το κόστος αυτού του τελεστή είναι 0 για το input. Αφού θα κάνει sort και join απευθείας **δίχως κάποιο I/O κόστος αφού αυτά χωράνε στην μνήμη με την μια**. Απο αυτόν τον operator στην καλύτερη περίπτωση δεν θα βγάλει κανένα αποτέλεσμα (0 εγγραφές) το join και στην χειρότερη θα γίνουν όλα με όλα (9 εγγραφές).

Με απλή αναλογία βγαίνει ότι οι σελίδες που θα αποθηκευτούν αυτές οι 9 εγγραφές δεν θα υπερβαίνουν το μέγεθος της μνήμης.

4 κόστος:

Εφόσον λοιπόν τώρα εφαρμόζουμε **indexed** nested loop join για κάθε εγγραφή του προηγούμενου 3ού κόστους θα ψαχνουμε στο index αν υπάρχει. Αν ναι τότε το I/O cost θα αυξανεται κατά 1 μονάδα. Υποθέτοντας τώρα ότι από τις 9 εγγραφές οι μισές κατά μέσο όρο θα υπάρχουν. Έχουμε κόστος 4.5 δηλαδή **5 I/Os**.

Να τονίσουμε εδώ μάλιστα ότι οι σελίδες του 3ου κόστους υπάρχουν ήδη στην μνήμη.

5 κόστος:

Το κόστος εδώ σε Input είναι μηδενικό αφού εφόσον χωράνε όλα τα blocks στην μνήμη αρκεί ο επεξεργαστής να κάνει show μόνο συγκεκριμένα.

Φυσικά τώρα αρκεί να κάνει output τα δεδομένα το οποίο ισούται με τον αριθμό των σελίδων που έμειναν μετά από όλα αυτά τα φιλτραρίσματα στην μνήμη.

6 κόστος:

Μηδενικό κόστος όπως στο 5ο κόστος.

Αρα το συνολικό κόστος θα είναι: **115 I/Os**

Αξίζει να σημειώσουμε εδώ μάλιστα ότι τώρα λόγω ομοιομορφίας οι σελίδες μπορούσαν να φορτώνονται όλες στην κύρια μνήμη. Σε κανονικές συνθήκες όμως αυτό είναι πιθανό να μην συμβαίνει και εκεί θα έχουμε μεγαλύτερο I/O κόστος ώστε να έχουμε το αποτέλεσμα που θέλουμε!