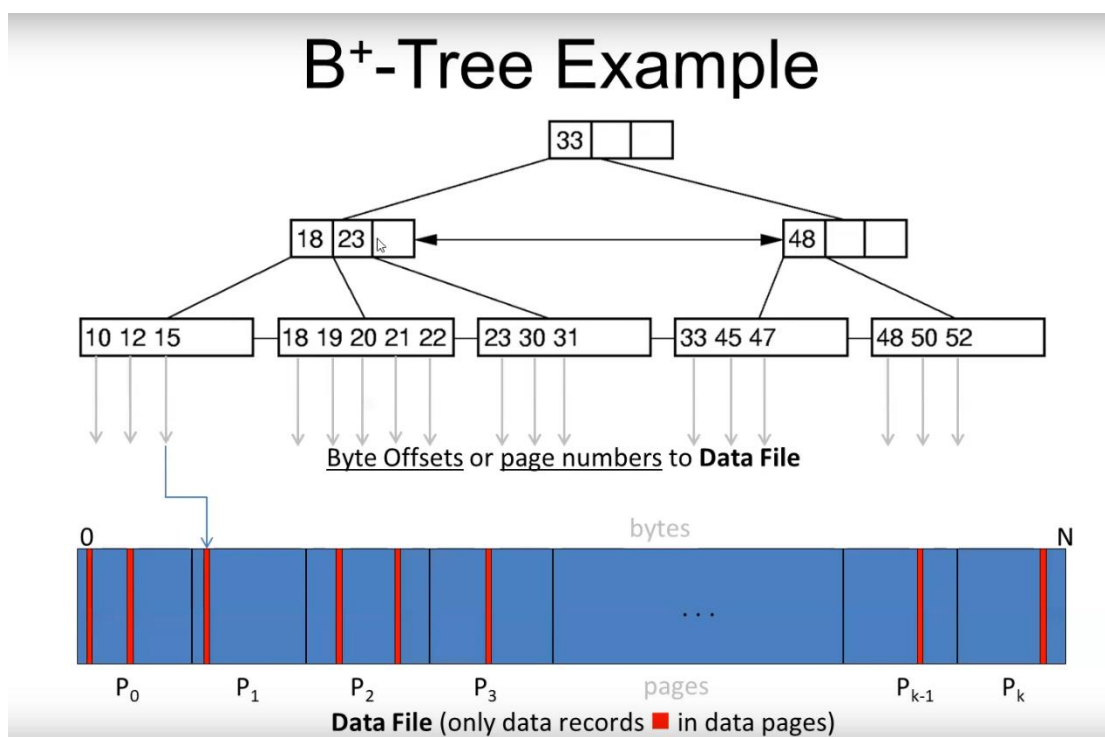


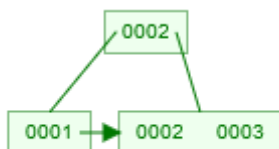
Μέσος αριθμός προσβάσεων δίσκου εισαγωγή ανά	Μέσος αριθμός προσβάσεων δίσκου ανά αναζήτηση αυτομάτη	Μέσος αριθμός προσβάσεων δίσκου για διάστημα K τιμών (K=10)	Μέσος αριθμός προσβάσεων δίσκου για διάστημα K τιμών (K=100)
15.3	6.0	16.85	1077.7



Σε αυτή την άσκηση κληθήκαμε να υλοποιήσουμε το B+tree στον δίσκο. Το B+tree είναι ένα ευρετήριο που συχνά χρησιμοποιείται στις βάσεις δεδομένων για να βελτιωθεί ο χρόνος εκτέλεσης των ερωτημάτων. Η ιδέα που καθιστά το B+tree χρήσιμο εργαλείο είναι η δυνατότητα του να δείχνει απευθείας στην σελίδα του δίσκου που είναι αποθηκευμένες οι τιμές ενός κλειδιού σε $\log_d N$ I/Os, όπου N ο συνολικός αριθμός των σελίδων (data file) στον δίσκο και d ο βαθμός του δέντρου. Χωρίς την χρήση του δέντρου θα έπρεπε κανείς να αναζητήσει σειριακά σε όλες τις σελίδες του δίσκου την τιμή ενός κλειδιού, κάτι που κοστίζει $O(N)$.

Σε αυτή την άσκηση υλοποιούμε την εισαγωγή στοιχείων στον δέντρο, την αναζήτηση τυχαίες τιμές καθώς και εύρος τιμών. Η εισαγωγή στοιχείων στον δέντρο περιλαμβάνει

κάποια σημαντικά βήματα, τα οποία συντελούν στο να παραμείνει στο δέντρο ισορροπημένο, ακόμα και αν τα στοιχεία εισάγονται ταξινομημένα. Το σημαντικότερο από αυτά τα βήματα είναι η διάσπαση, κατά την οποία ένας κόμβος που έχει γεμίσει διασπάται σε δύο κομμάτια. Το κάθε κομμάτι περιλαμβάνει τα μισά στοιχεία του αρχικού κόμβου και το ενδιάμεσο (median) στοιχείο του κόμβου μεταβαίνει στο πατέρα-κόμβο. Για παράδειγμα η παρακάτω εικόνα δείχνει ένα δέντρο βαθμού 3 μετά από την διάσπαση.



Με αυτόν το τρόπο, όλοι οι κόμβοι είναι ομοιόμορφα φορτωμένοι με κλειδιά και το δέντρο ποτέ δεν εκφυλίζεται. Επίσης φαίνεται στην παραπάνω εικόνα ότι οι κόμβοι φύλλα είναι συνδεδεμένοι με ένα pointer. Αυτή η προσθήκη μειώνει δραματικά τον χρόνο και τα I/Os που θα χρειαστούν για να γίνουν τα ερωτήματα εύρους τιμών αφού οι κόμβοι στα φύλλα είναι μια ταξινομημένη συνδεδεμένη λίστα. Συνεπώς, αφού κανείς βρει την θέση έναρξης (low) του ερωτήματος εύρους τιμών στα φύλλα, τότε αρκεί να συνεχίζει να διασχίζει την λίστα των φύλλων σειριακά μέχρι το ανώτερο όριο του εύρους τιμών να μην ικανοποιείται.

Πώς όμως τα φύλλα του δέντρου μας οδηγούν με ένα I/O στην σελίδα του δίσκου που είναι αποθηκευμένη η τιμή του κλειδιού; Τα φύλλα έχουν ένα pointer για κάθε κλειδί, που δείχνει στην σελίδα του δίσκου που μας ενδιαφέρει, παραλείποντας σελίδες που δεν περιέχουν πληροφορία χρήσιμη για το ερώτημα μας. Συνεπώς, αναμένουμε θεωρητικά μέσο αριθμό προσβάσεων στον δίσκο ανα αναζήτηση $\log_d N$ προσβάσεις στο δίσκο μέχρι την εύρεση του φύλλου που περιέχει το κλειδί που αναζητούμε + 1 I/O για ανάγνωση της τιμής του κλειδιού από τον δίσκο. Αυτή η υπόθεση επαληθεύεται πειραματικά αφού για $N = 10^5$ κλειδιά/σελίδες χρειαζόμαστε περίπου 6 προσβάσεις στον δίσκο ανα αναζήτηση.

Για ερωτήματα εύρους κλειδιών- K αναμένουμε $\log_d N$ προσβάσεις στο δίσκο μέχρι την εύρεση του φύλλου που περιέχει το ελάχιστο κλειδί μέσα στο ζητούμενου εύρους + K αναγνώσεις από το data file + K/d αναγνώσεις κόμβων φύλλα. Η παραπάνω θεωρητικός αριθμός από I/Os επιβεβαιώνεται στην πράξη αφού για $K = 10$ χρειαζόμαστε ~16 προσβάσεις στον δίσκο και για $K = 1000$ χρειαζόμαστε ~1077 I/O's, αντίστοιχα.

Ο μέσος όρος προσβάσεων στον δίσκο για εισαγωγές εξαρτάται από τον βαθμό του δέντρου, συνεπώς από το μέγεθος της σελίδας που είναι διαθέσιμο για την αποθήκευση κόμβων. Είναι γεγονός ότι όσο μεγαλύτερος είναι αυτός ο αριθμός, τόσο λιγότερες διασπάσεις συμβαίνουν, συνεπώς μειώνεται ο μέσος αριθμός προσβάσεων. Για σελίδα 256 bytes, ο βαθμός του κόμβου είναι $d = 29$. Αυτό προκύπτει υπολογίζοντας τον χώρο που μένει διαθέσιμος για κλειδιά στην σελίδα αφού αφαιρεθεί ο χώρος που χρειάζονται τα υπόλοιπα στοιχεία του κόμβου (child pointer, right and left pointer, etch).

References

- <https://www.cs.usfca.edu/~galles/visualization/BPlusTree.html>
- <https://www.youtube.com/watch?v=gV5G3UXdwS0>