

CSE 373: B-Trees

Pete Morcos
University of Washington
4/14/00

<http://www.cs.washington.edu/education/courses/cse373/00sp>

A new style of tree

- B-trees are unlike the other trees we've seen
 - Data only stored in leaves; interior nodes just for searching
 - Each node has many children (often hundreds)
 - Thus, tree is extremely shallow
- Very important in database systems
- Designed for high performance when managing enormous amounts of data

Disks

- Many databases hold many gigabytes or terabytes (10^{12}) of info
 - Too much to fit in memory
- Disk access time is measured in ms, memory time in ns—about a million times slower
- When disk data is accessed, you read a whole *page* (512 bytes to a few K), not just one byte
 - We'll assume 1000 byte pages for simplicity
- The all-important goal is to **reduce the number of disk accesses!**

Why not just use binary trees?

- Recall that in binary search trees, each node visited reduces search space by about one-half
 - Thus, $\log_2 N$ node accesses needed per search
- $\log_2 1,000,000,000,000$ is about 40
 - 40 disk accesses for each piece of data is unacceptable
- Since we are going to get a whole page of data per disk read anyway, make nodes as big as possible
 - Each node has M children
 - Suppose search keys are strings up to 36 bytes
 - $M = 1000 / (36 + 4 \text{ [for a child pointer]}) = 250$
 - $\log_{250} 1,000,000,000,000 = 5$ (as opposed to 40)
- Often, top 2 levels fit in memory, so medium size B-trees only have to hit the disk once (for the actual data node)

Our simplifying assumptions

- Databases can't use pointers in B-trees since some nodes will be in memory, some on disk
- We'll assume everything fits in memory
- Some of our written problems will ask you about 1000 byte disk pages, as on previous slide
- But most examples and homework will use what I'll call a mini B-tree, with $M = 4$

Example B-tree

- B-tree rules:
 - All leaves at same depth, and hold sorted array of data
 - Root has 2 to M children (labeled P_0 to P_{M-1})
 - Other non-leaf nodes have $\lceil M/2 \rceil$ to M children
- Each non-leaf has up to M-1 values (k_1 to k_{M-1}); child P_i holds values $\geq k_i$; P_0 holds stuff $\leq k_1$
- Note that values in non-leaves are not actual data!



