

CSE 332

Data Structures & Parallelism

B Trees

Melissa Winstanley
Spring 2024

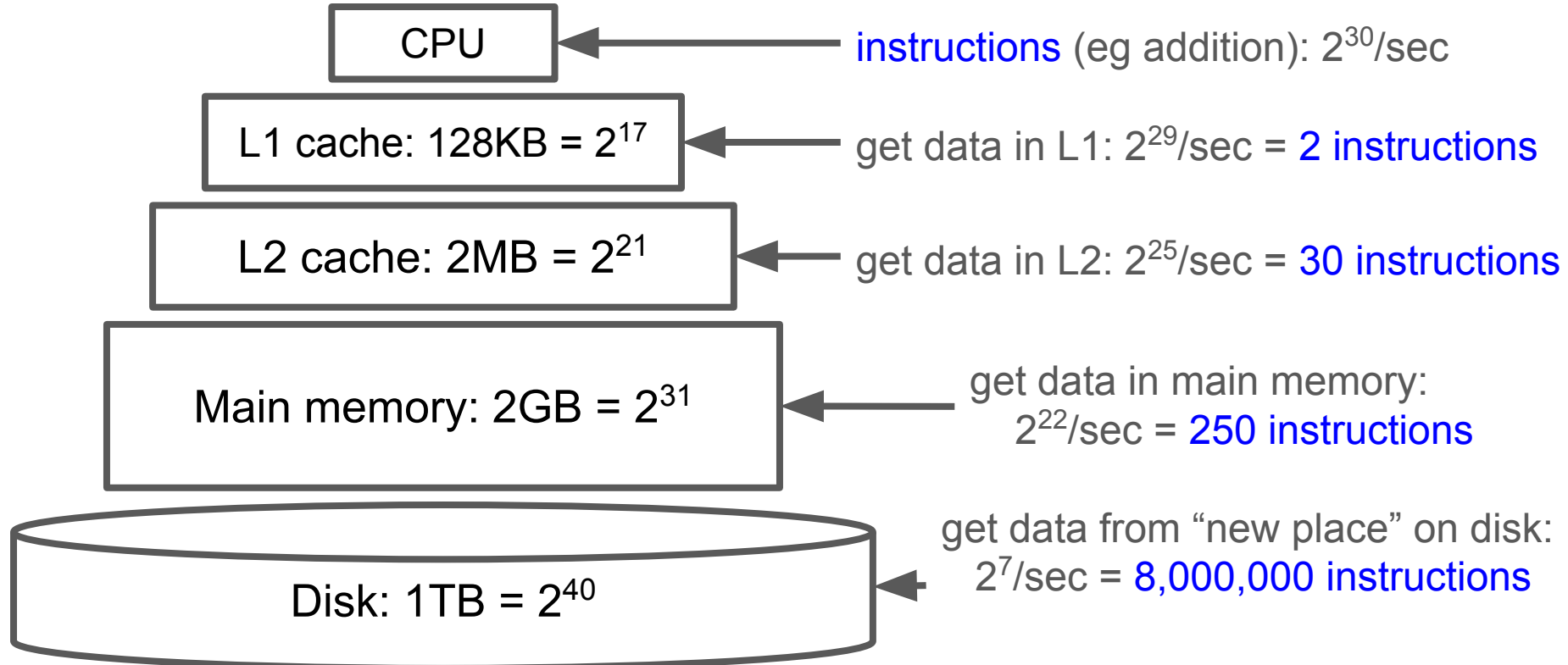
Question Time

One of the assumptions that Big-Oh makes is that all operations take the same amount of time.

Is that really true?

A Typical Memory Hierarchy

“Every desktop/laptop/server is different” but here’s a plausible configuration



“Fuggedaboutit”, usually

- The hardware **automatically** moves data into the caches from main memory for you
 - Replacing items already there
 - So algorithms much faster if “data fits in cache” (often does)
- Disk accesses are done by software (e.g., ask operating system to open a file or database to access some data)
- So most code **“just runs”** but sometimes it’s worth designing algorithms / data structures with knowledge of memory hierarchy
 - And when you do, you often need to know one more thing...

How does data move up the hierarchy?

- Moving data up the memory hierarchy is slow because of latency (think distance-to-travel)
 - Since we're making the trip anyway, may as well carpool
 - Get a block of data in the same time it would take to get a byte
 - Sends **nearby memory** because:
 - It's easy
 - And likely to be asked for soon
- Side note: Once a value is in cache, may as well keep it around for awhile; accessed once, a particular value is more likely to be accessed again in the **near future** (more likely than some random other value)

Spatial locality



Temporal locality



Locality

Temporal Locality (locality in **time**) – If an address is referenced, it will tend to be referenced again soon.

Spatial Locality (locality in **space**) – If an address is referenced, addresses that are close by will tend to be referenced soon.

Arrays vs Linked lists

Which has the potential to best take advantage of spatial locality?

Block/line size

- The amount of data moved from **disk** into **memory** is called the “**block**” size or the “**page**” size
 - Not under program control
- The amount of data moved from **memory** into **cache** is called the cache “**line**” size
 - Not under program control

BSTs?

- Looking things up in balanced binary search trees is $O(\log n)$, so even for $n = 2^{39}$ (512GB) we need not worry about minutes or hours
- Still, number of disk accesses matters:
 - Pretend for a minute we had an AVL tree of height 55
 - The total number of nodes could be? _____
 - Most of the nodes will be on disk: the tree is shallow, but it is still many gigabytes big so the entire *tree* cannot fit in memory
 - Even if memory holds the first 25 nodes on our path, we still potentially need 30 disk accesses if we are traversing the entire height of the tree.

Note about numbers

- **Note:** All the numbers in this lecture are “ballpark” “back of the envelope” figures
- **Moral:** Even if they are off by, say, a factor of 5, the moral is the same:

**If your data structure is mostly on disk,
you want to minimize disk accesses**

- A better data structure in this setting would exploit the block size and relatively fast memory access to avoid disk accesses...

Trees as Dictionaries

(N = 10 million) [Example from Weiss]

In worst case, each node access is a disk access, number of accesses:

Worst case big-O

Disk accesses

- BST
- AVL
- B Tree

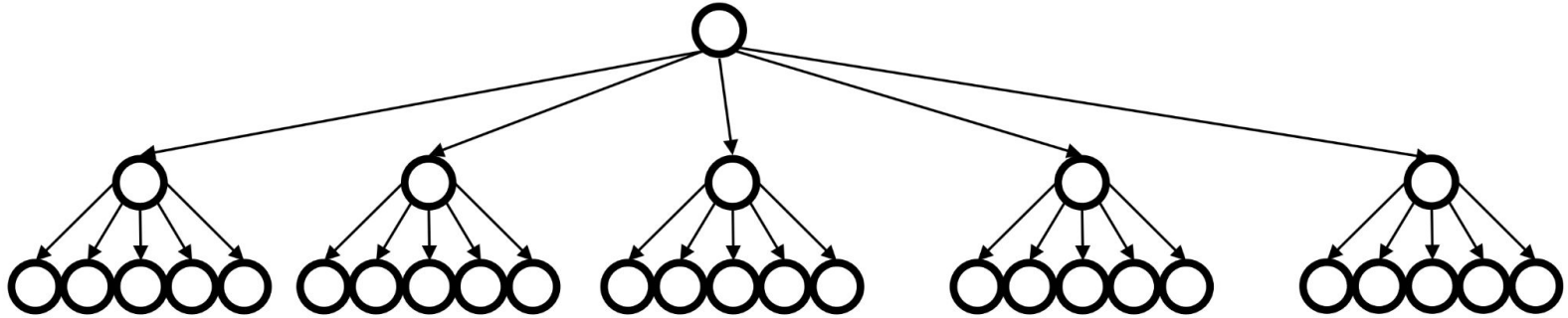
Our goal

- **Problem:** A dictionary with so much data *most of it is on disk*
- **Desire:** A balanced tree (logarithmic height) that is even shallower than AVL trees so that we can minimize disk accesses and exploit disk-block size
- **A key idea:** Increase the branching factor of our tree

M-ary Search Tree

Build some sort of search tree with branching factor M :

- Have an array of sorted children (`Node[]`)
- Choose M to fit snugly into a disk block (1 access for array)

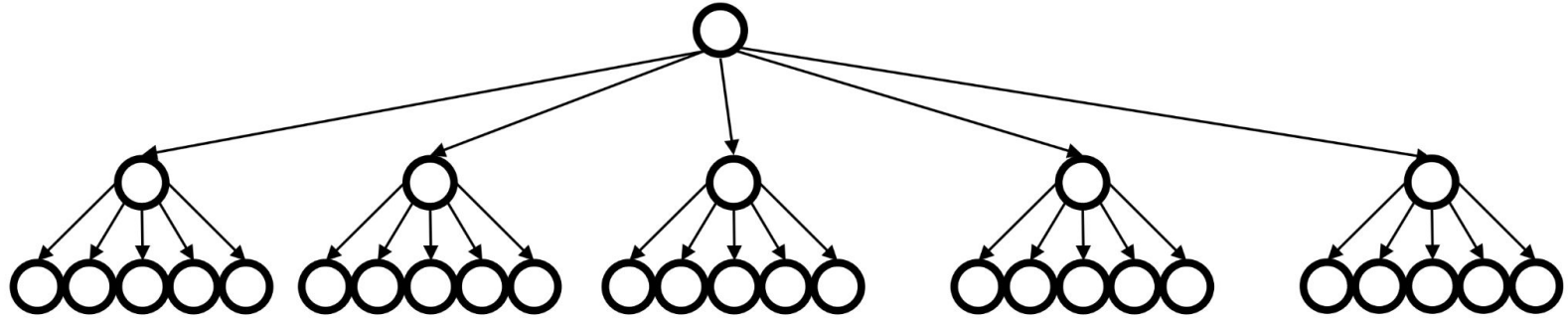


Perfect tree of height h has $(M^{h+1}-1)/(M-1)$ nodes (textbook, page 4)

What is the **height** of this tree?

What is the worst case running time of **find**?

Complexity of Find in M-ary Search Tree



How many **hops**?

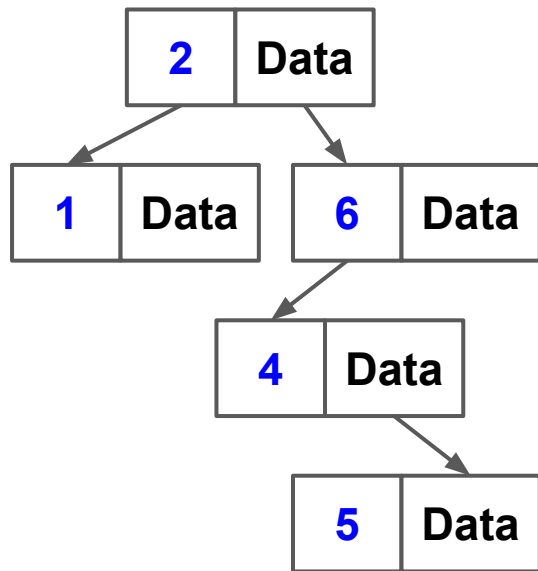
How much **work** at each level?

(find which child to take)

Overall complexity?

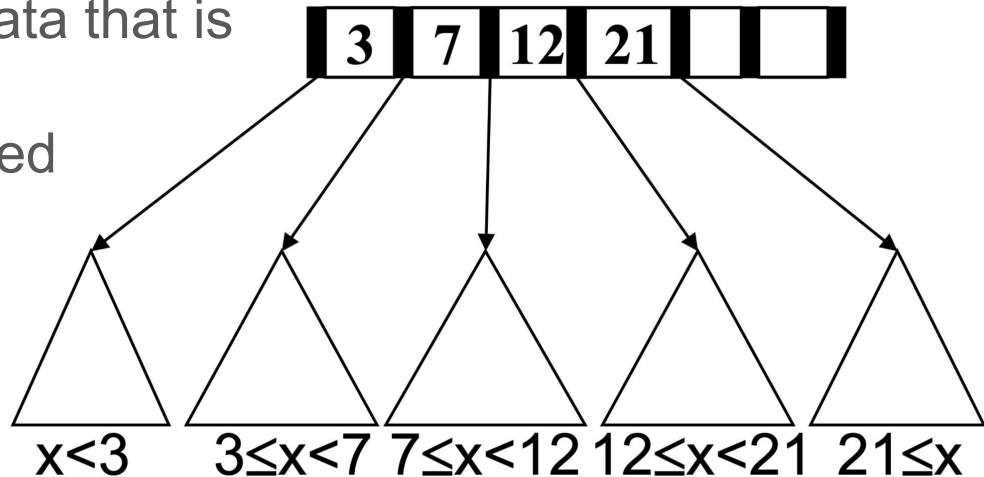
Questions about M-ary search trees

- What should the **order** property be?
- How would you **rebalance** (ideally without more disk accesses)?
- Storing **real data** at inner-nodes (like we do in a BST) seems kind of wasteful...
 - To access the node, will have to load the **data** from disk, even though most of the time we won't use it!!
 - Usually we are just “passing through” a node on the way to the value we are actually looking for.
- So let's use the branching-factor idea, but for a **different kind of balanced tree**:
 - **Not** a *binary search tree*
 - But still logarithmic height for any $M > 2$



B+ Trees (we and the book say “B Trees”)

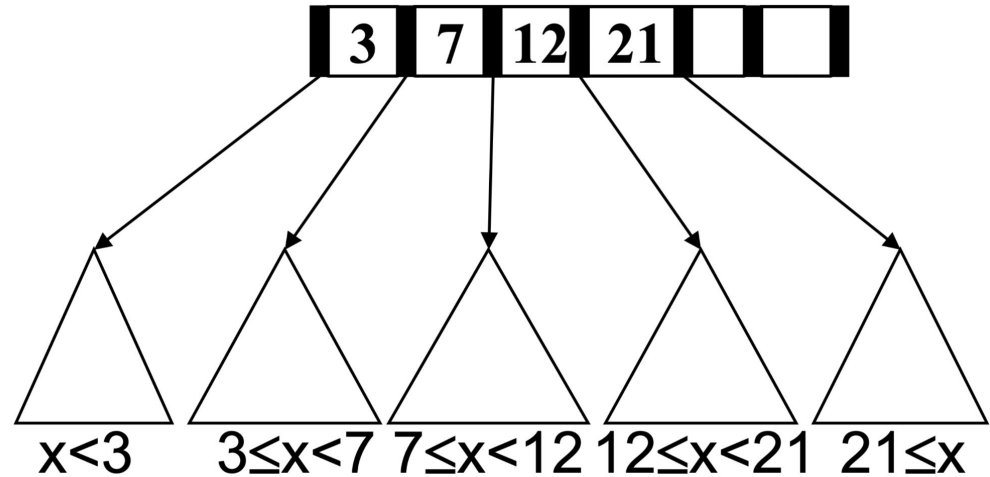
- Two types of nodes: **internal nodes** & **leaves**
- Each **internal node** has room for up to $M-1$ keys and M children
 - No other data; **all data at the leaves!**
- **Order property:** Subtree **between** keys **a** and **b** contains only data that is $\geq a$ and $< b$ (notice the \geq)
- **Leaf** nodes have up to L sorted data items
- As usual, we'll ignore the “along for the ride” data



B Trees: Leaves vs Internal Nodes

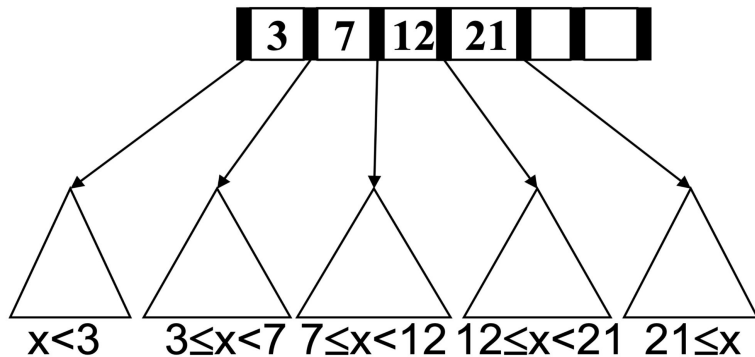
Remember:

- **Leaves** store data
- **Internal nodes** are 'signposts'
- There are different ways to implement these - pay attention to how we talk about them!



Find

- Different from BST in that we don't store data at internal nodes
- But **find** is still an easy root-to-leaf recursive algorithm
 - At each internal node do binary search on (up to) M-1 keys to find the branch to take
 - At the leaf do binary search on the (up to) L data items
- But to get logarithmic running time, we need a balance condition...



Structure Properties

- **Root** (special case)
 - If tree has $\leq L$ items, root is a leaf (occurs when starting up, otherwise unusual)
 - Else has between 2 and M children
- **Internal nodes**
 - Have between $\lceil M/2 \rceil$ and M children, i.e., **at least half full**
- **Leaf nodes**
 - **All leaves at the same depth**
 - Have between $\lceil L/2 \rceil$ and L data items, i.e., **at least half full**
- Any $M > 2$ and L will work, but:
We pick M and L **based on disk-block size**