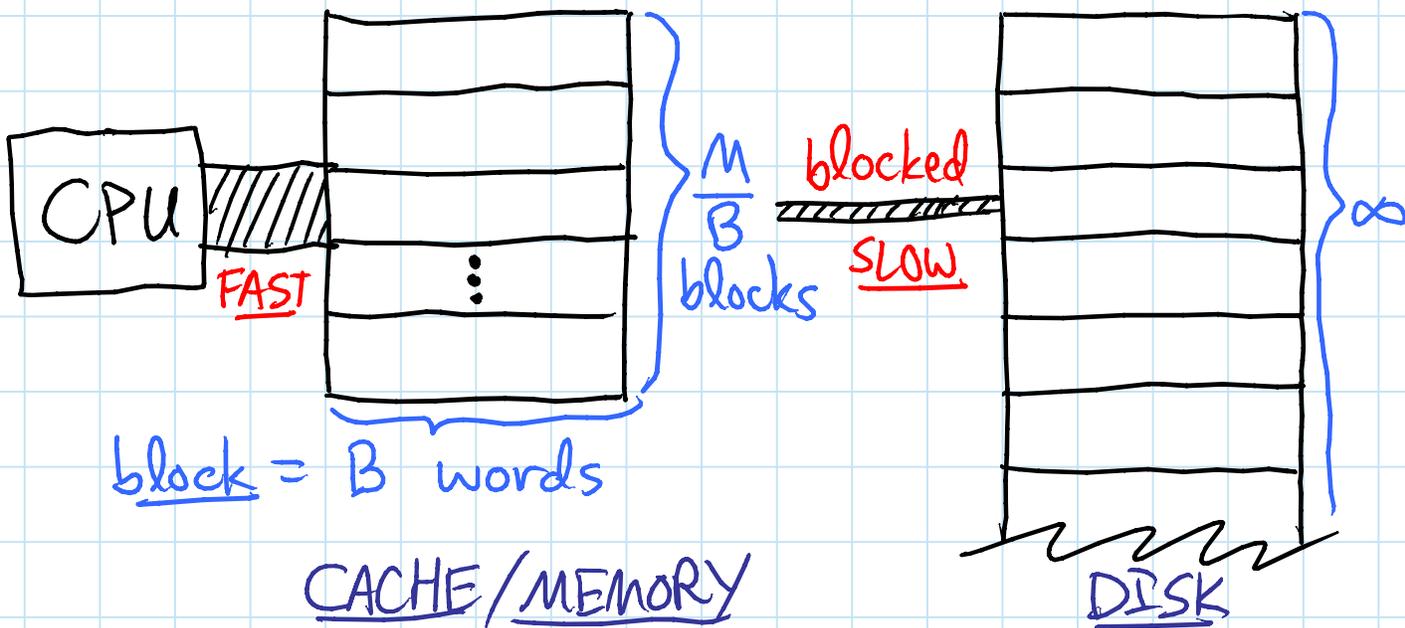


- TODAY: memory hierarchies
- external-memory model
  - cache-oblivious model
  - cache-oblivious B-trees

## External memory / I/O / Disk Access Model:

[Aggarwal & Vitter - CACM 1988]

two-level memory hierarchy



- focus on # memory transfers:  
blocks read/written between cache & disk
- $\leq$  RAM running time
- $\geq \frac{\text{cell-probe LB}}{B}$
- when can we save this factor of  $\geq B$ ?

# Basic results in external memory:

⑧ Scanning:  $O(\lceil \frac{N}{B} \rceil)$  to read/write  $N$  words in order

① Search trees:

- B-trees with branching factor  $\Theta(B)$  support insert, delete, predecessor search in  $O(\log_{B+1} N)$  memory transfers (&  $O(\lg N)$  time, with care, in comparison model)
- $\Omega(\log_{B+1} N)$  for search in comparison model:
  - where query fits among  $N$  items requires  $\lg(N+1)$  bits of information 
  - each block read reveals where query fits among  $B$  items  $\Rightarrow \leq \lg(B+1)$  bits of info.
  - $\Rightarrow$  need  $\geq \frac{\lg(N+1)}{\lg(B+1)}$  memory transfers

- also optimal in "block-probe model" if  $B \geq w$

[Patrascu & Thorup - see L11]

② Sorting:  $O(\frac{N}{B} \log_{MB} \frac{N}{B})$  memory transfers

$\hookrightarrow B \times$  faster than B-tree sort!

$\Omega(\text{ditto})$  in comparison model

③ Permuting:  $O(\min \{ N, \frac{N}{B} \log_{MB} \frac{N}{B} \})$

physical execution

$\Omega(\text{ditto})$  in indivisible model

$\hookrightarrow$  can't pack pieces of input words in words

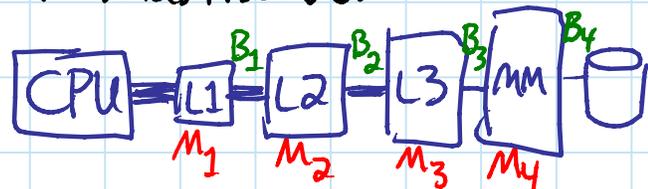
④ Buffer tree:  $O(\frac{1}{B} \log_{MB} \frac{N}{B})$  amortized mem. transf. for delayed queries & batched updates &  $O(\phi)$  delete-min ( $\Rightarrow$  priority queues)

# Cache-oblivious model: [Frigo, Leiserson, 6.046 Prokop, Ramachandran - FOCs 1999; Prokop - MEng 1999]

- like external-memory model
- but algorithm doesn't know  $B$  or  $M$  (!)
- ⇒ must work for all  $B$  &  $M$
- automatic block transfers triggered by word access with offline optimal block replacement
  - FIFO, LRU, or any conservative replacement is 2-competitive given cache of  $2x$  size  
(resource augmentation)
- dropping  $M \Rightarrow M/2$  doesn't affect typical bounds e.g. sorting bound

## Cool:

- clean model: algorithm just like RAM
- adapts to changing  $B$  (disk tracks & cache) &  $M$  (competing processes)
- OPEN: formalize this
- adapts to all levels of multilevel memory hierarchy:



- often possible!

## Basic cache-oblivious results:

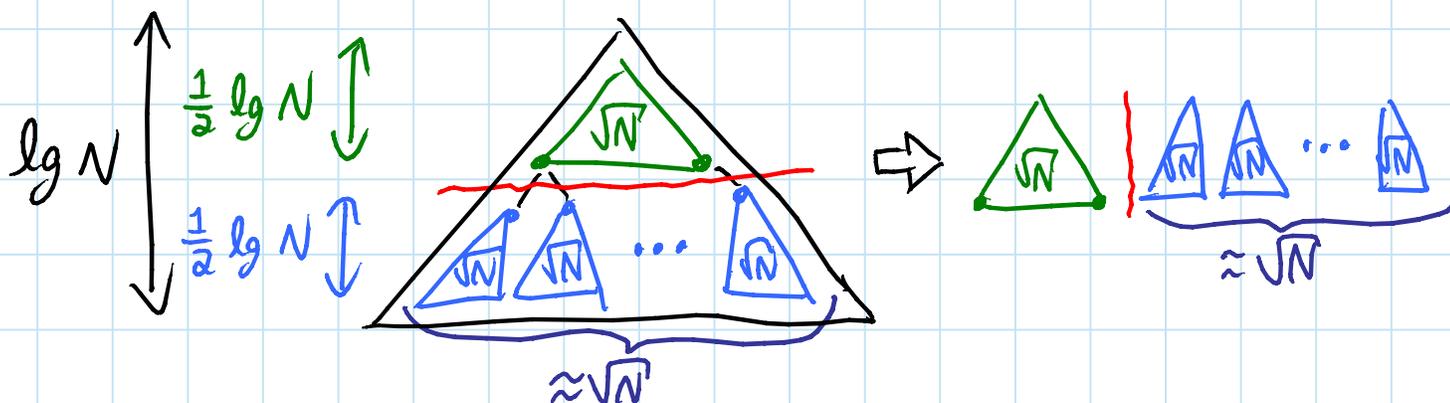
- ① Scanning: same (algorithm & bound)
- \* ① Search trees: insert, delete, & search } TODAY & L21  
in  $O(\log_{B+1} N)$  memory transfers  
[Bender, Demaine, Farach-Colton - FOCs 2000/SICOMP 2005]  
[Bender, Duan, Iacono, Wu - SODA 2002/J. Alg. 2004]  
[Brodal, Fagerberg, Jacob - SODA 2002]  
- best constant is  $\lg e$ , not 1  
[Bender, Brodal, Fagerberg, Ge, He, Hu, Iacono, López-Ortiz - FOCs 2003]
- ② Sorting:  $O(\frac{N}{B} \log_{M/B} \frac{N}{B})$  memory transfers  
[Frigo et al. 1999; Brodal & Fagerberg - ICALP 2002]  
- uses tall-cache assumption:  $M = \Omega(B^{1+\epsilon})$   
- impossible otherwise [Brodal & Fagerberg - STOC 2003]
- ③ Permuting: min impossible [Brodal & Fagerberg - same]
- \* ④ Priority queue:  $O(\frac{1}{B} \log_{M/B} \frac{1}{B})$  amortized mem. transf. } L21  
- uses tall-cache assumption  
[Arge, Bender, Demaine, Holland-Minkley, Munro - STOC 2002/SICOMP 2007; Brodal & Fagerberg - ISAAC 2002]

# Cache-oblivious static search trees:

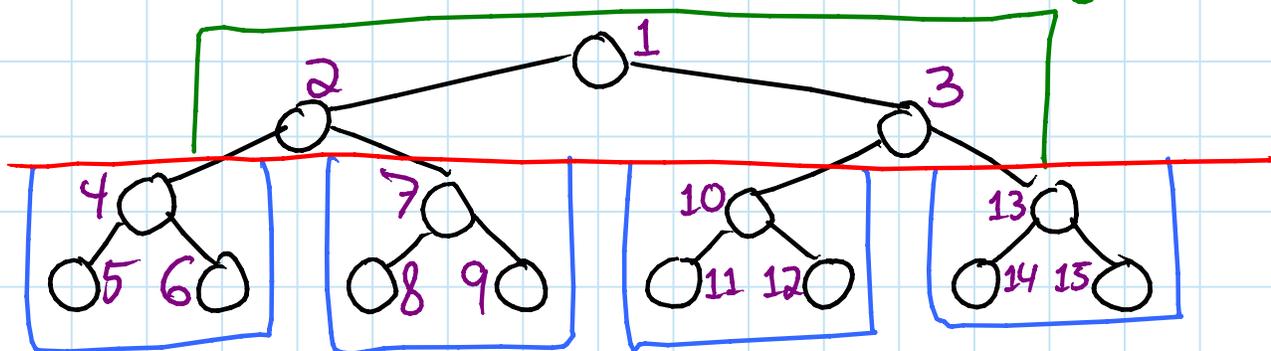
(binary search)

[Prokop-MEng 1999]

- store  $N$  elements in  $N$ -node complete BST
- carve tree at middle level of edges
- $\Rightarrow \approx \sqrt{N} + 1$  pieces of size  $\approx \sqrt{N}$



- recursively lay out pieces & concatenate:  
(in any order)



$\Rightarrow$  order to store nodes

"van Emde Boas layout"

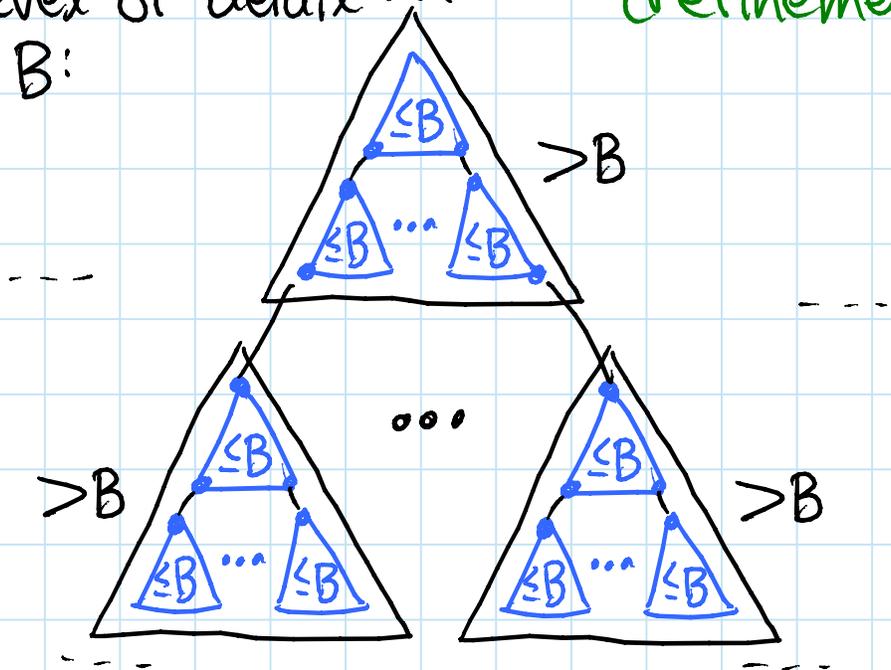
- generalizes to [Bender, Demaine, Farach-Colton 2000]
- height not a power of 2
- node degrees  $\geq 2$  &  $O(1)$

# Analysis of van Emde Boas layout:

- consider level of detail ...

(refinement)

straddling B:



- cutting height in half until  $\leq \lg B \rightarrow$  chunks  
 $\Rightarrow$  chunks have height between  $\frac{1}{2} \lg B$  &  $\lg B$

( $\Rightarrow$  size between  $\sqrt{B}$  &  $B$ )

$\Rightarrow$  #chunks visited on root-to-leaf path  $\leq \frac{\lg N}{\frac{1}{2} \lg B} = 2 \log_B N$

(assume  $B \geq 2 \Rightarrow \lg B \geq 1$ )

- each chunk stores  $\leq B$  words consecutively

$\Rightarrow$  occupies  $\leq 2$  blocks (depending on alignment)

$\Rightarrow$  #memory transfers  $\leq 4 \log_B N$

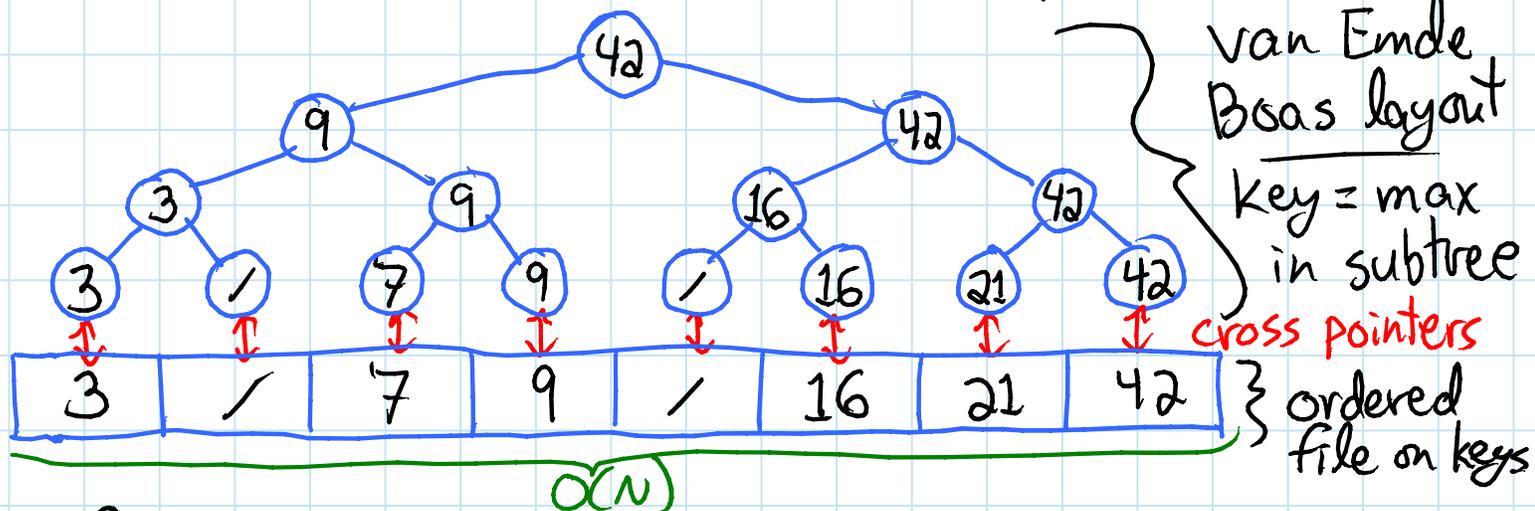
(assuming  $M \geq 2B$ )

# Cache-oblivious B-trees as in [Bender, Duan, Jacons, Wu]

## ① ordered file maintenance: (to do in L21)

- store  $N$  elements in specified order in an array of size  $O(N)$  ~ allow gaps
- updates: delete element:  
insert element between two specified elts.  
by moving elements in array interval of  $O(\lg^2 N)$  amortized via  $O(1)$  interleaved scans

## ② build static search tree on top:

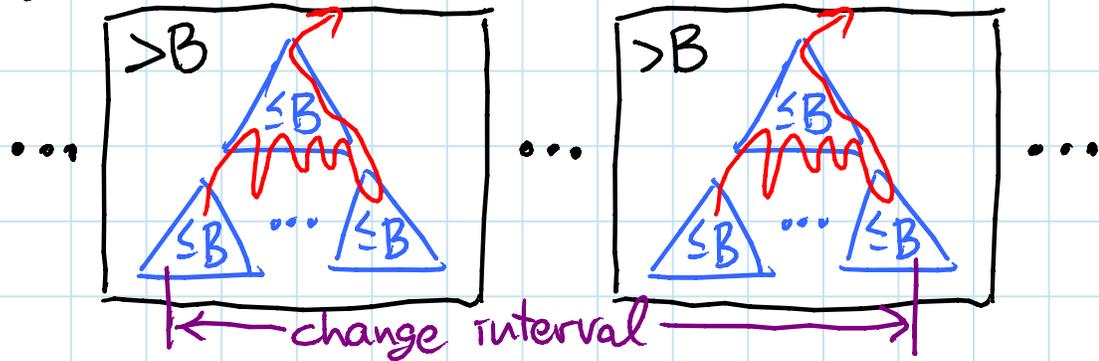


## ③ ops:

- search looks at left child's key to decide L/R  
- still  $O(\log_{B+1} N)$  ← at most double
- insert( $x$ ):
  - search( $x$ ) to find predecessor & successor
  - insert  $x$  in between in ordered file
  - update values in leaves corresp. to changed cells & propagate changes up tree in postorder traversal
- delete( $x$ ) similar

④ update analysis: if  $K$  cells change in ordered file then  $O(\frac{K}{B} + \log_{B+1} N)$  memory transfers

- look at level of detail straddling  $B$
- look at bottom two levels:



- within superchunk of  $>B$ , jumping between  $\leq 2$  chunks of  $\leq B$   $\sim$  assume  $M \geq 2B$

$\Rightarrow O(\lceil \frac{\text{superchunk}}{B} \rceil)$  memory transfers per superchunk  
 ↳ actually visited portion  
 ↳ unnecessary because superchunk  $>B$ , except first & last

$\Rightarrow O(\frac{K}{B} + 1)$  memory transfers in bottom two levels

- # ancestors above these two levels

$$\leq \frac{K}{B} + \lg N$$

ancestors to lca      path to root

cost  $\leq 1$  each

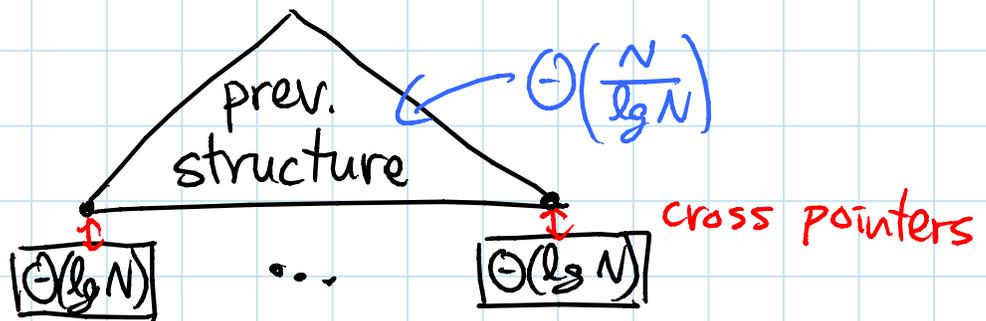
$\log_{B+1} N$  cost as in search

$\Rightarrow O(\frac{K}{B} + \log_{B+1} N)$  total memory transfers

So far: search in  $O(\log_{B+1} N)$   
 update in  $O(\log_{B+1} N + \frac{\lg^2 N}{B})$  amortized  
 suboptimal if  $B = o(\lg N \lg \lg N)$

## ⑤ indirection:

- divide elements into clusters of  $\Theta(\lg N)$
- use previous structure on min(each cluster)



- search: vEB top, scan bottom  
 $\Rightarrow O(\log_{B+1} N + \frac{\lg N}{B}) = O(\log_{B+1} N)$
- update cluster by complete rewrite  $\Rightarrow O(\frac{\lg N}{B})$
- keep clusters between 25% & 100% full
- split / merge & split when necessary (like B-tree)  
 $\Rightarrow \Omega(\lg N)$  updates to charge to
- $\Rightarrow O(1)$  updates in top structure  
only every  $\Omega(\lg N)$  updates
- $\Rightarrow$  amortized update cost  $= O\left(\frac{\log_B N + \frac{\lg^2 N}{B}}{\lg N}\right)$   
 $= O\left(\frac{\lg N}{B}\right)$   
plus search cost  $= O(\log_B N)$

So:  $O(\log_{B+1} N)$  insert, delete, search