# Lecture 19

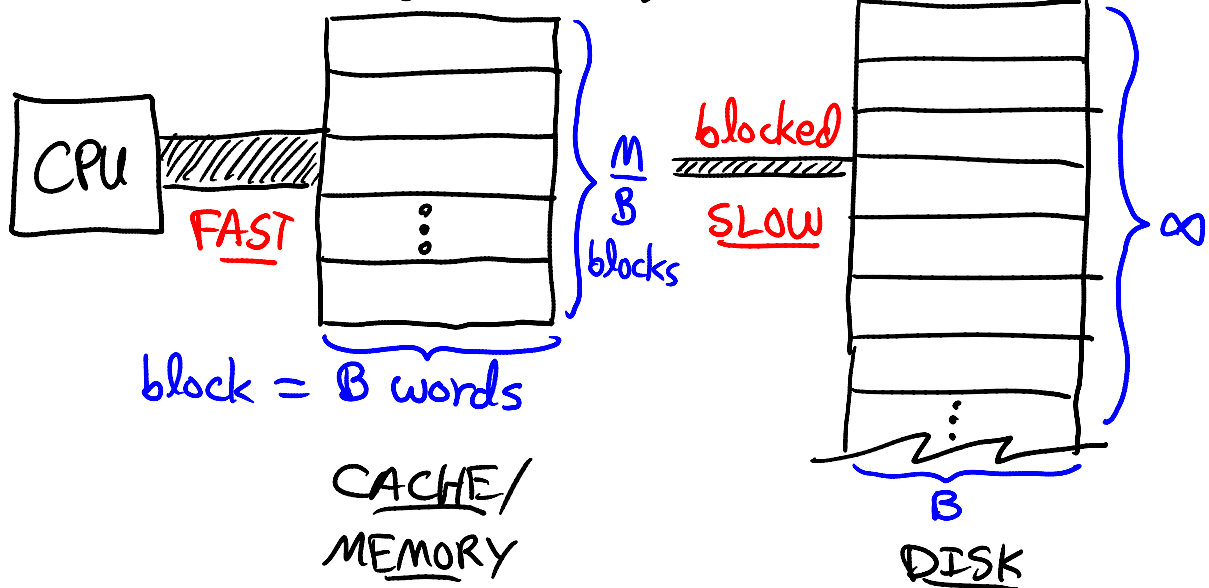## External memory/I/O/Disk Access Model [Aggarwal & Vitter—CACM 1988]
— two-level memory hierarchy:



- charge mainly for <u>memory transfers</u>:
  blocks read/written between cache & disk
- any algorithm with running time $T(N)$
  uses $\leq T(N)$ memory transfers
- when can we use fewer? (& still $T(N)$ time)

# Basic results in external memory:

⓪ <u>Scanning</u>: $O(\lceil \frac{N}{B} \rceil)$ to read/write $N$ words in order

① <u>Search trees</u>:
- B-trees with branching factor $\Theta(B)$ support insert, delete, (predecessor) search in $O(\log_{B+1} N)$ memory transfers
  <span style="color:green">(& $O(\lg N)$ time in comparison model)</span>
- this is optimal 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 information
  $\Rightarrow$ need $\geq \frac{\lg(N+1)}{\lg(B+1)}$ memory transfers

- also optimal in "block-probe model" if $B \geq w$ <span style="color:purple">[Lecture 15]</span>

② <u>Sorting</u>: $O(\frac{N}{B} \log_{M/B} \frac{N}{B})$ memory transfers
  $\hookrightarrow$ <span style="color:green">$\geq B \times$ faster than B-tree sort!</span>
  $\Omega(\text{ditto})$ in comparison model

③ <u>Permutation</u>: $O(\min\{N, \frac{N}{B} \log_{M/B} \frac{N}{B}\})$
  <span style="color:blue">physical execution</span>$\downarrow$  $\Omega(\text{ditto})$ in "indivisible model" —
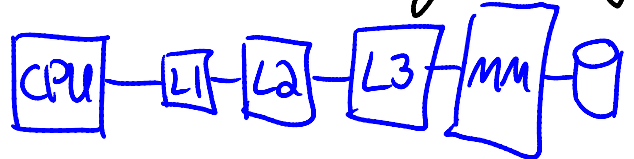  can't pack pieces of input words into word

④ <u>Buffer tree</u>: $O(\frac{1}{B} \log_{M/B} \frac{N}{B})$ amortized memory transfers for delayed queries / batched updates & instant delete-min $\Rightarrow$ priority queues

# Cache-oblivious model [Frigo, Leiserson, Prokop, Ramachandran— FOCS 1999; Prokop—M.Eng. 1999]

— like external-memory model
— but algorithm doesn't know $B$ or $M$ (!)
— automatic block transfers triggered by word access
  with <u>offline optimal</u> block replacement
   — FIFO, LRU, or any conservative replace. strategy
     is 2-competitive given cache of $2\times$ size
   — dropping $M \searrow M/2$ doesn't affect bounds like

## Cool:
   — clean model: algorithm just like RAM alg.
   — <u>adapts</u> to changing $B$ (disk tracks)
      informally...        $\& M$ (competing processes)
   — adapts to all levels of multilevel memory hierarchy
     each with own $B \& M$

     $\boxed{CPU} - \boxed{L1} \boxed{L2} - \boxed{L3} \boxed{MM} - \boxed{\phantom{o}}$

   — often possible!

# Basic cache-oblivious results:

⓪ <u>Scanning</u>: same algorithm & bound

* ① <u>"B-tree"</u>: insert, delete, & search

in $O(\log_{B+1} N)$ memory transfers <span style="color:red">} TODAY</span>
<span style="color:red">} & L20</span>

[Bender, Demaine, Farach-Colton — FOCS 2000/SICOMP 2005;

Bender, Duan, Iacono, Wu — SODA 2002;

Brodal, Fagerberg, Jacob — SODA 2002]

— best constant is $\lg e$, not 1  [Bender, Brodal, Fagerberg, Ge, He, Hu, Iacono, López-Ortiz — FOCS 2003]

② <u>Sorting</u>: $O(\frac{N}{B} \log_{M/B} \frac{N}{B})$ memory transfers

[Frigo et al. 1999; Brodal & Fagerberg — ICALP 2002]

— uses <u>tall-cache assumption</u>: $M = \Omega(B^{1+\varepsilon})$

— impossible otherwise  [Brodal & Fagerberg — STOC 2003]

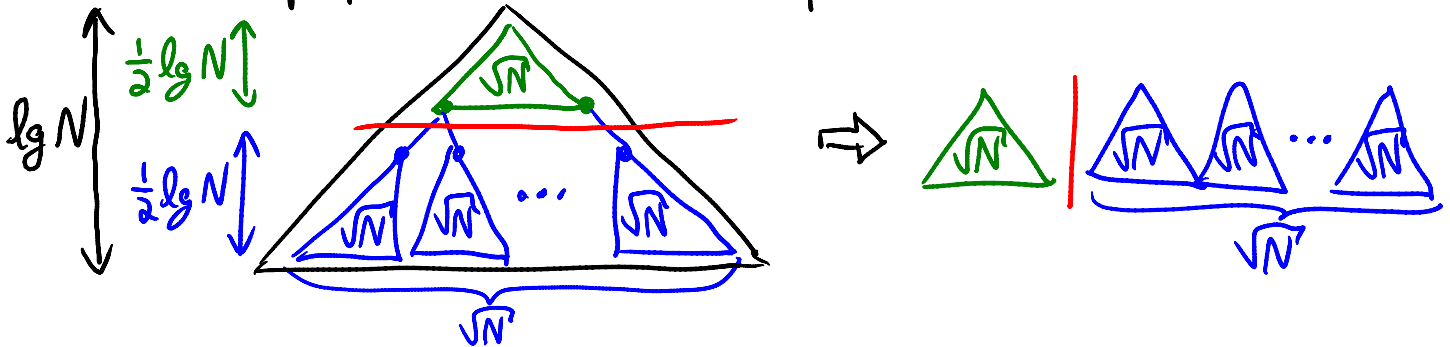③ <u>Permuting</u>: min is impossible  [Brodal & Fagerberg — same]

④ <u>Priority queue</u>: $O(\frac{1}{B} \log_{M/B} \frac{N}{B})$ amortized mem. transfers

(also uses tall-cache assumption)

[Arge, Bender, Demaine, Holland-Minkley, Munro — STOC/SICOMP 2002/2007;

Brodal & Fagerberg — ISAAC 2002]

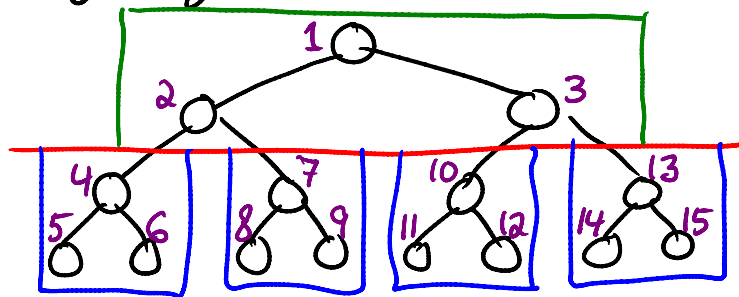# Cache-oblivious binary search / static search trees:
## van Emde Boas layout [Prokop-MEng 1999; Bender, Demaine, Farach-Colton - FOCS 2000]

- store $N$ elements in order in $N$-node complete BST
- carve tree at middle level of edges
$\Rightarrow$ one top piece, $\approx \sqrt{N}$ bottom pieces, each size $\approx \sqrt{N}$
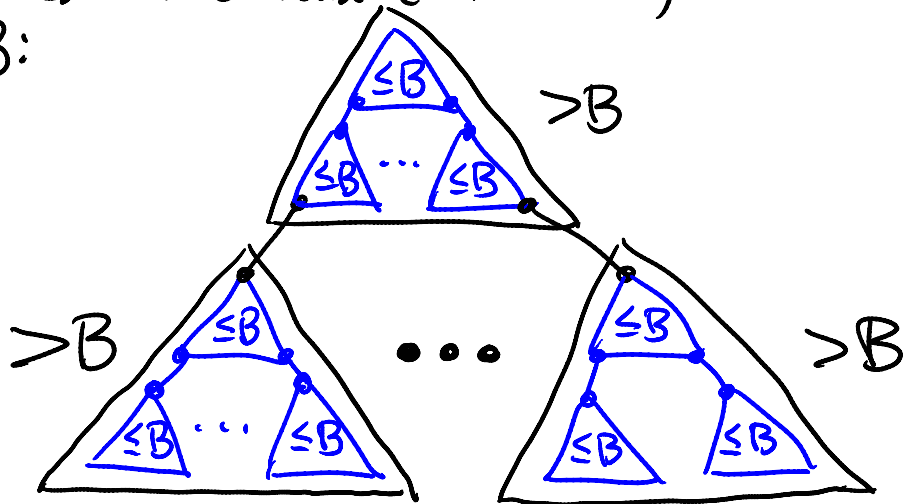


- recursively lay out pieces & concatenate
  e.g:

# Analysis of van Emde Boas layout:

— consider level of detail (refinement) straddling B:



— cutting height in half until $\leq \lg B \to$ pieces

$\Rightarrow$ pieces have height between $\frac{1}{2}\lg B$ & $\lg B$

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

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

  (sloppy on B+1 issue)

— each piece stores $\leq B$ elements consecutively

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

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

  (assuming $M \geq 2B$)
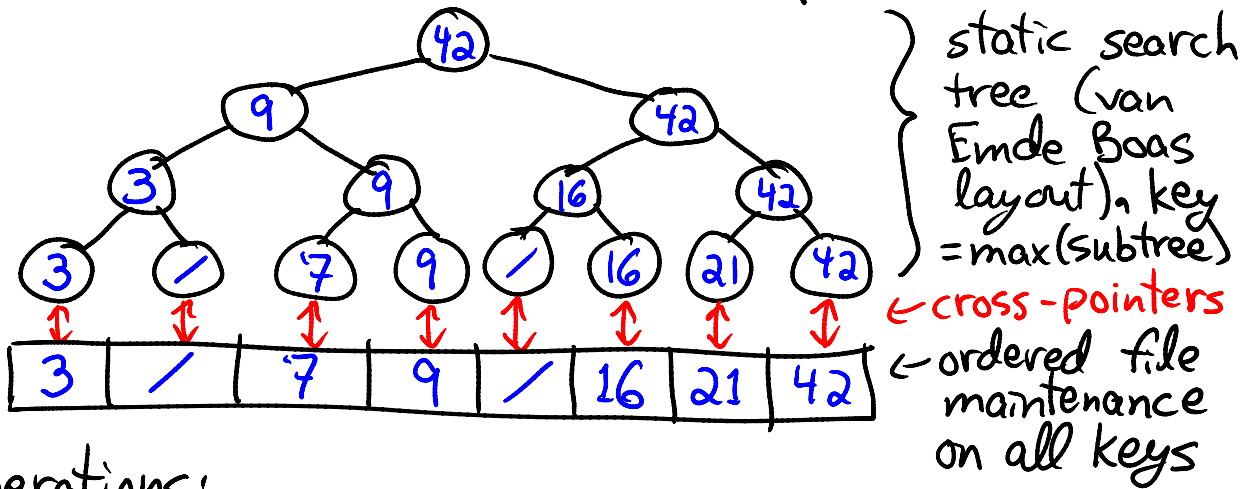

# Generalizations:

— height not a power of 2

— node degrees $\geq 2$ & $O(1)$

# Cache-oblivious B-trees, as in [Bender, Duan, Iacono, Wu — SODA 2002]

① ordered file maintenance:
   - store N elements in specified order in array of size $O(N)$   (allow gaps)
   - updates: insert element between two specified / delete element by moving elements in array interval of $O(\lg^2 N)$ amortized
   - black box to be filled [Lecture 20]

② build static search tree on top:



static search tree (van Emde Boas layout), key = max(subtree)

← cross-pointers

← ordered file maintenance on all keys
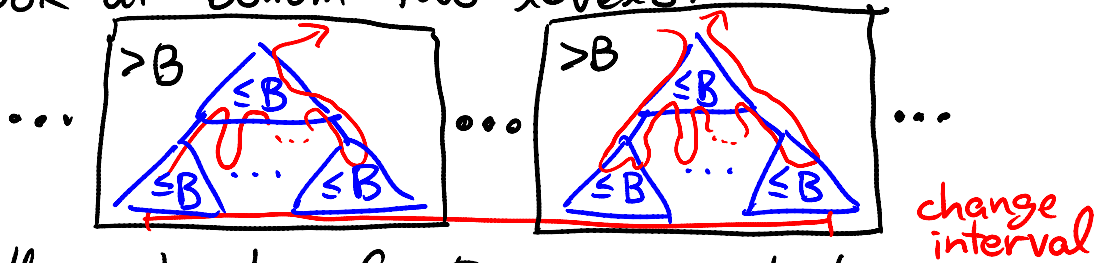
③ operations:
   - search looks at left child's key to decide direction
   - insert(x):
      - search(x) finds predecessor/successor
      - insert x in between in ordered file
      - update values in leaves corresp. to changed cells & propagate changes up tree, in postorder traversal
   - delete(x) similar

④ <u>update analysis</u>: if $K$ cells change in ordered file
then $O\left(\frac{K}{B} + \log_B N\right)$ mem. transf.

— look at level of detail straddling $B$

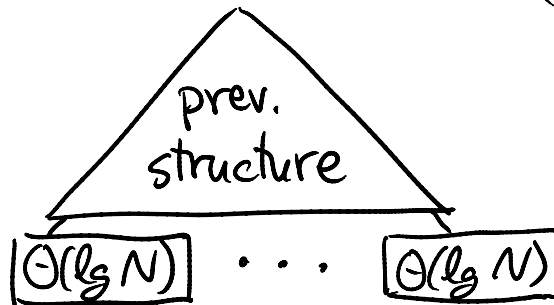— look at bottom two levels:



change interval

— within <u>chunk</u> of $>B$, jumping between
  $\leq 2$ <u>subchunks</u> of $\leq B$  ~ assume $M \geq 2B$

$\Rightarrow O(\text{chunk}/B)$ memory transfers per chunk

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

— # nodes above these two levels $\leq \frac{K}{B} + \lg N$

ancestors to lca || path to root
cost $\leq 1$ each || $\log_B N$ cost
              as before

$\Rightarrow O\left(\frac{K}{B} + \log_B N\right)$ total memory transfers


<u>So far</u>:  search in $O(\log_{B+1} N)$
          update in $O\left(\log_{B+1} N + \frac{\lg^2 N}{B}\right)$ amortized

suboptimal if
$B = o(\lg N \lg \lg N)$

⑤ indirection:
   — cluster elements into $\Theta\left(\frac{N}{\lg N}\right)$ groups of $\Theta(\lg N)$
   — use previous structure on min(each cluster)



   — update cluster by complete rewrite $\sim O\left(\frac{\lg N}{B}\right)$
   — keep cluster between 25% & 100% full
   — split / merge & split when necessary

$\Omega(\lg N)$ updates to charge to

$\Rightarrow$ update in top structure
   only every $\Omega(\lg N)$ updates
$\Rightarrow$ amortized update cost $= O\left(\frac{\lg N}{B}\right)$
   plus search cost

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