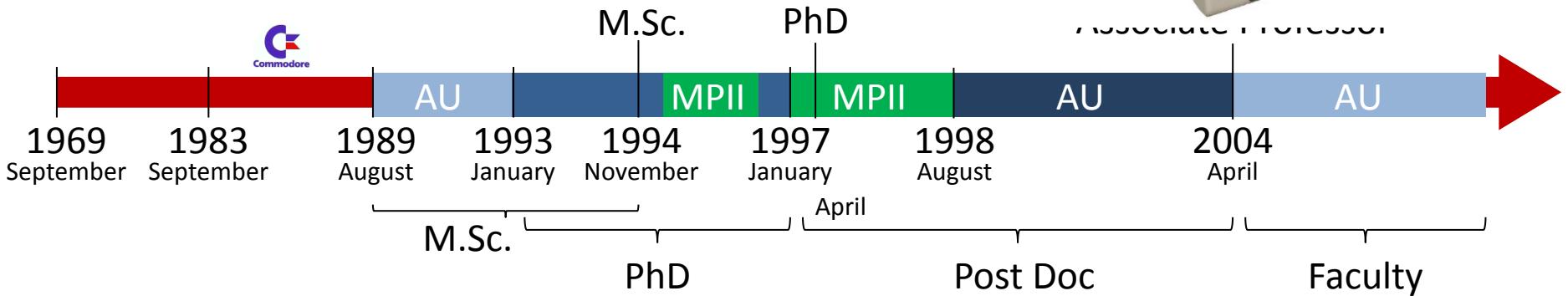


Cache-Oblivious Dynamic Dictionaries with Optimal Update/Query Tradeoff Indeksstrukturer til Ekstern Hukommelse

Gerth Stølting Brodal
Aarhus University



Gerth Stølting Brodal





madalgo

CENTER FOR MASSIVE DATA ALGORITHMIC





Danmarks
Grundforskningsfond
Danish National
Research Foundation

- Center of
- **Lars Arge**, Professor, Centerleader
- Gerth S. Brodal, Associate Professor
- 5 Post Docs, 10 PhD students, 4 TAP
- Total budget for 5 years ca. 60 million DKR

AU



Arge



Brodal

MIT



Demaine



Indyk

MPII



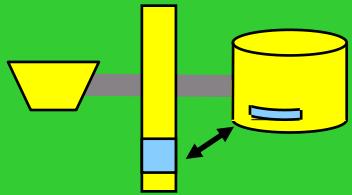
Mehlhorn

Frankfurt

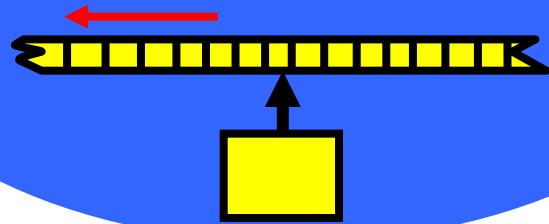


Meyer

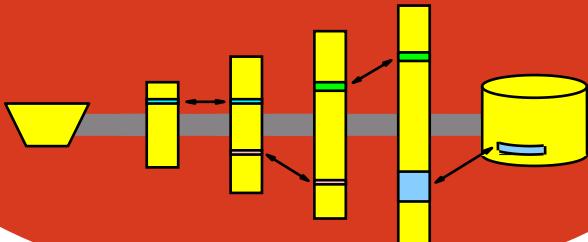
I/O Efficient Algorithms



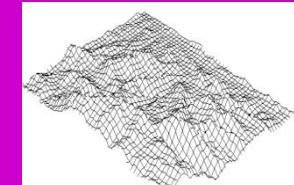
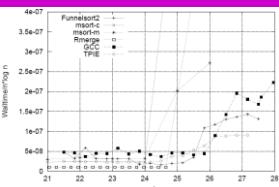
Streaming Algorithms

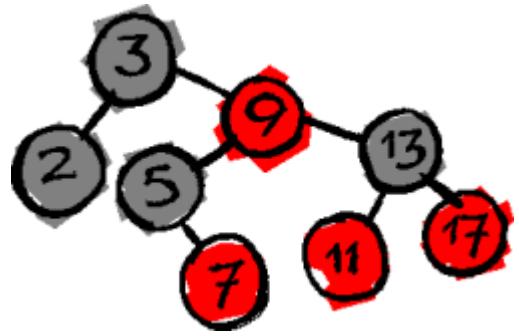


Cache Oblivious Algorithms

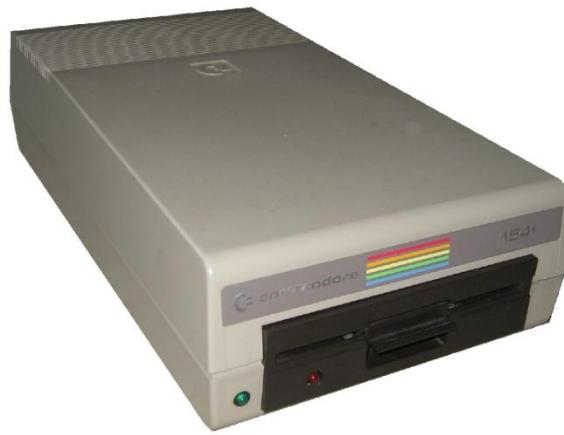


Algorithm Engineering

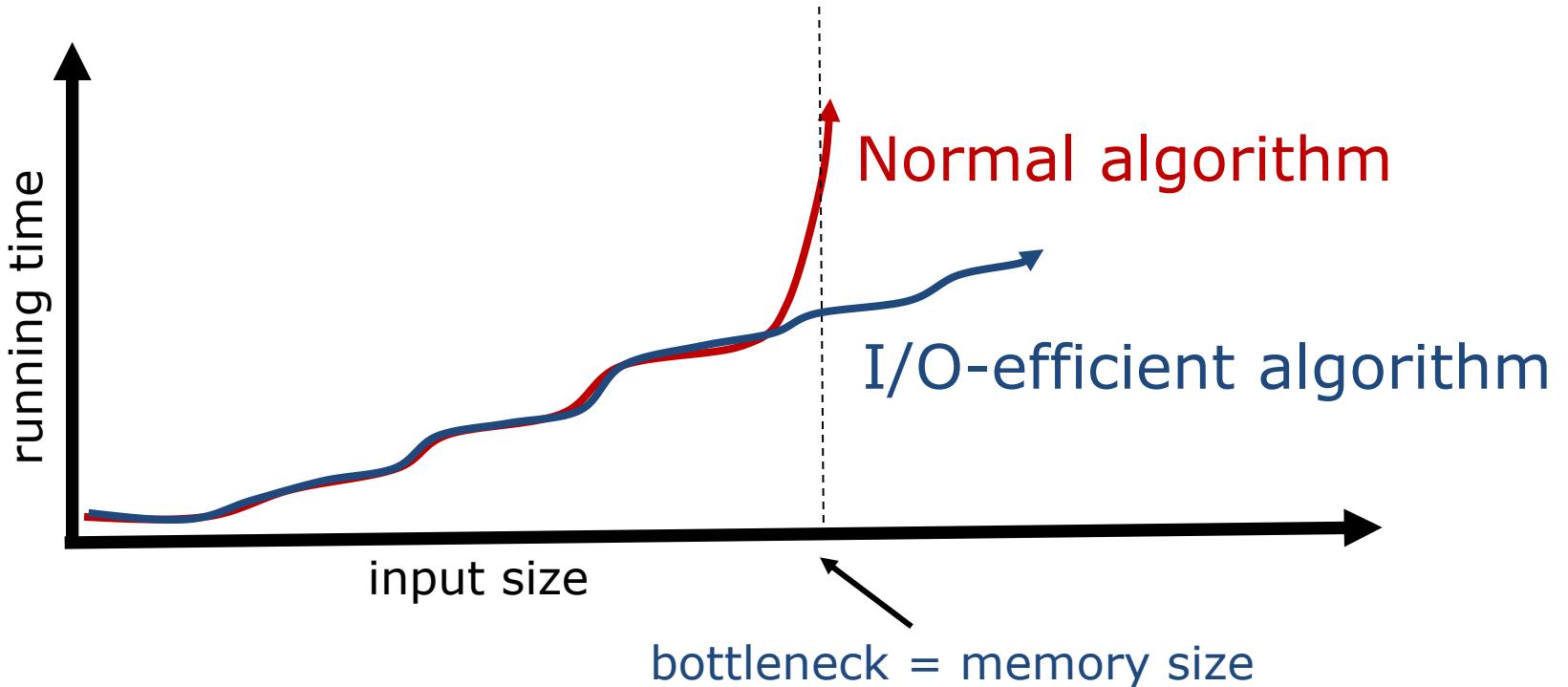




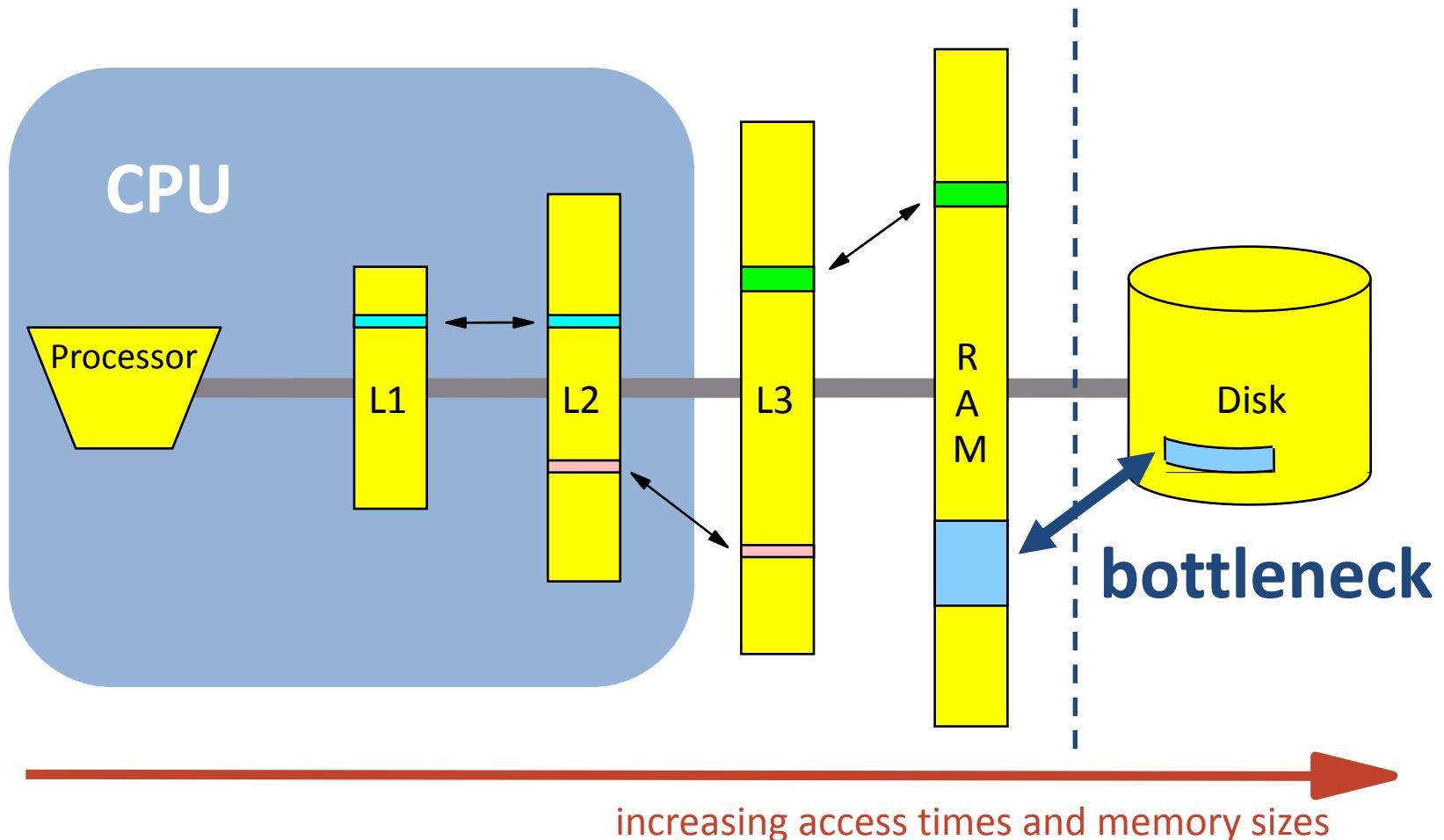
+



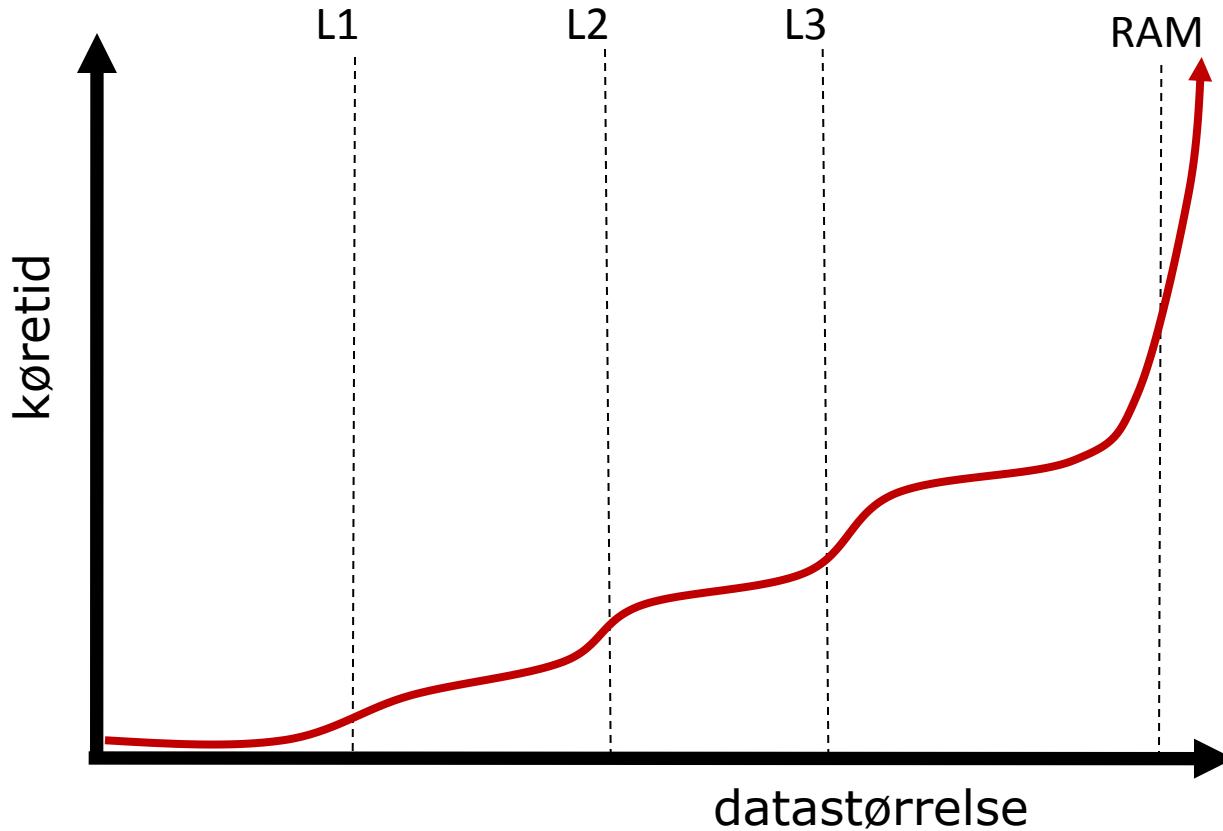
The problem...



Memory Hierarchies



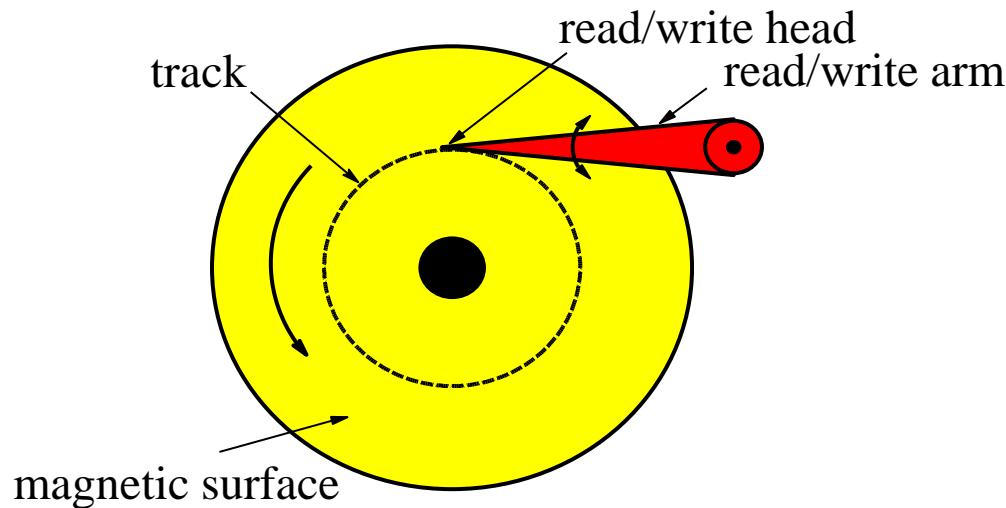
Memory Hierarchies vs. Running Time



Memory Access Times

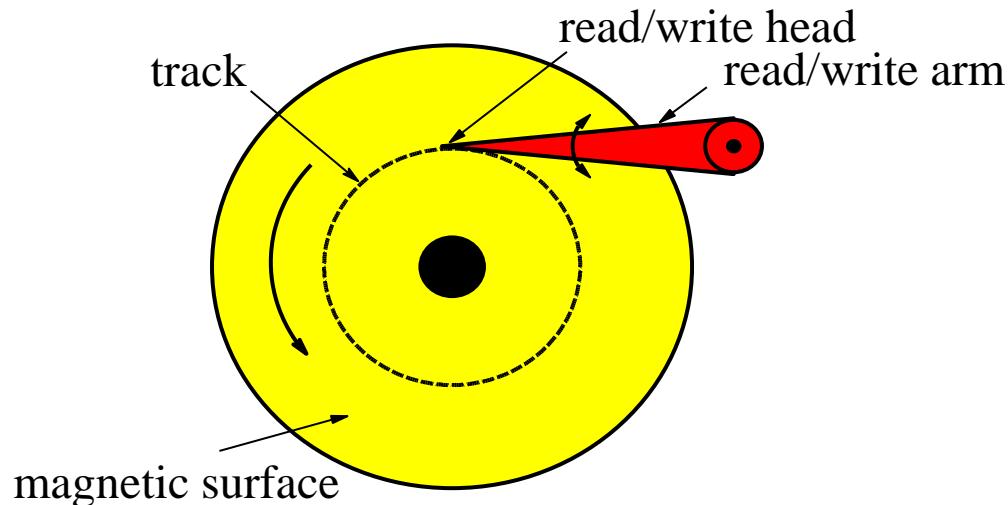
	Latency	Relative to CPU
Register	0.5 ns	1
L1 cache	0.5 ns	1-2
L2 cache	3 ns	2-7
DRAM	150 ns	80-200
TLB	500+ ns	200-2000
Disk	10 ms	10^7

Disk Mechanics



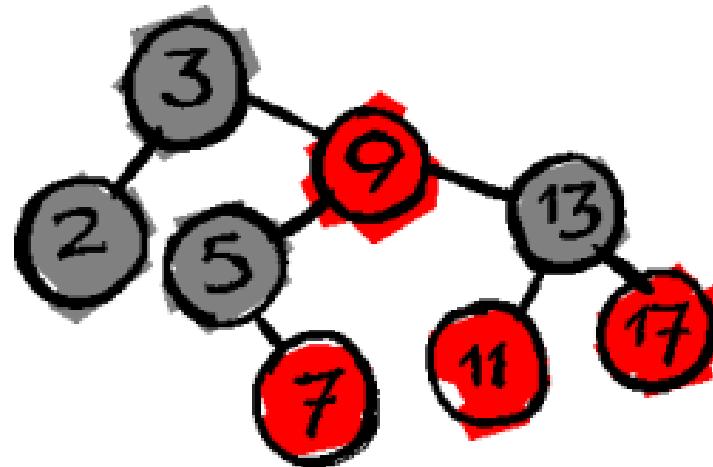
"The difference in speed between modern CPU and disk technologies is analogous to the difference in speed in sharpening a pencil using a sharpener on one's desk or by taking an airplane to the other side of the world and using a sharpener on someone else's desk." (D. Comer)

Disk Mechanics

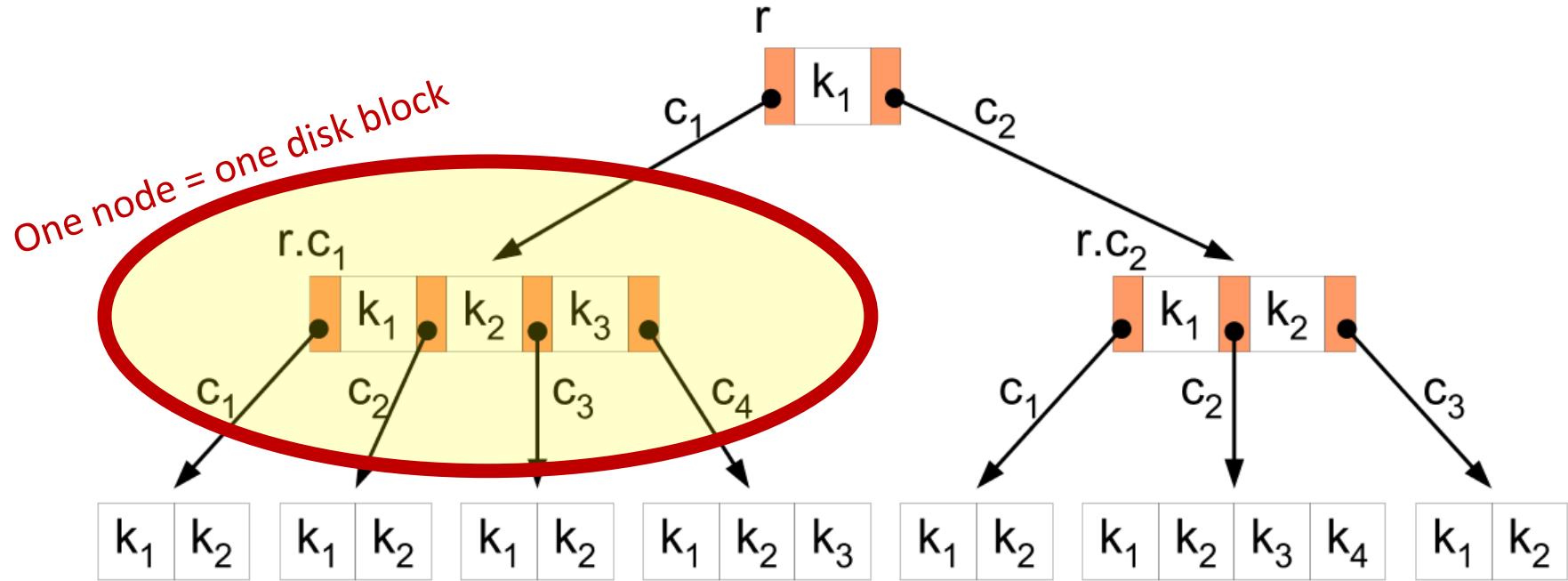


- I/O is often bottleneck when handling massive datasets
- Disk access is **10⁷ times slower** than main memory access!
- Disk systems try to amortize large access time transferring **large contiguous blocks** of data
- Need to store and access data to **take advantage of blocks** !

Internal Memory Index

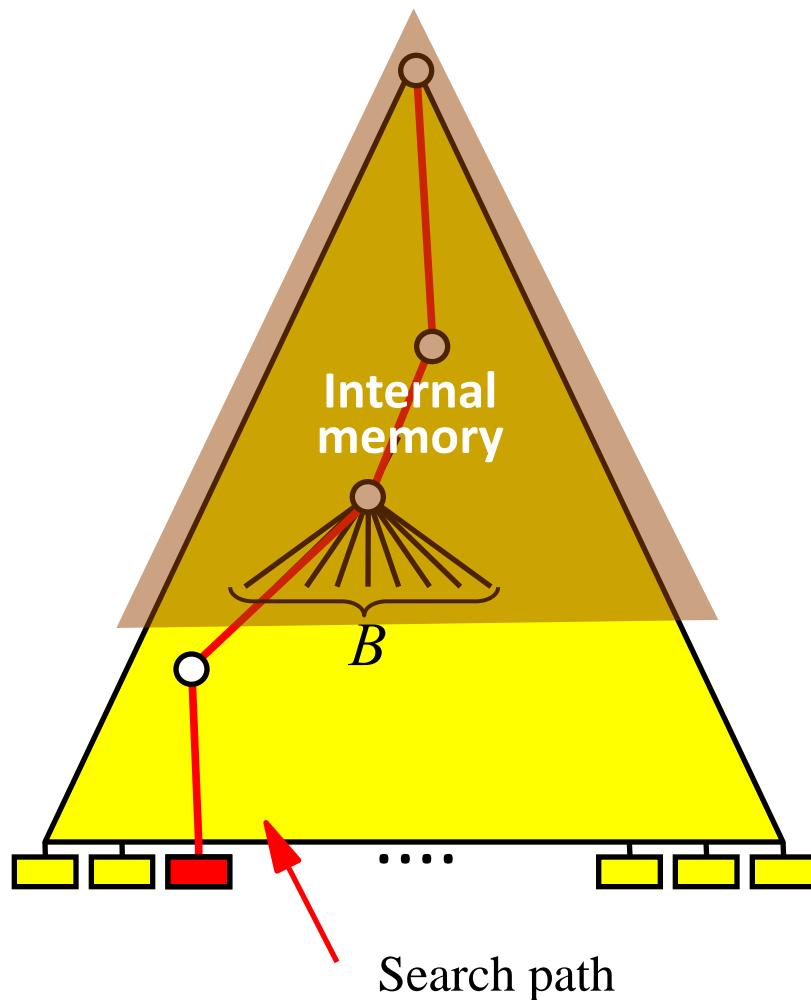


B(oing)-Trees



Height $O(\log_B N)$

B-trees - The Basic Searching Structure



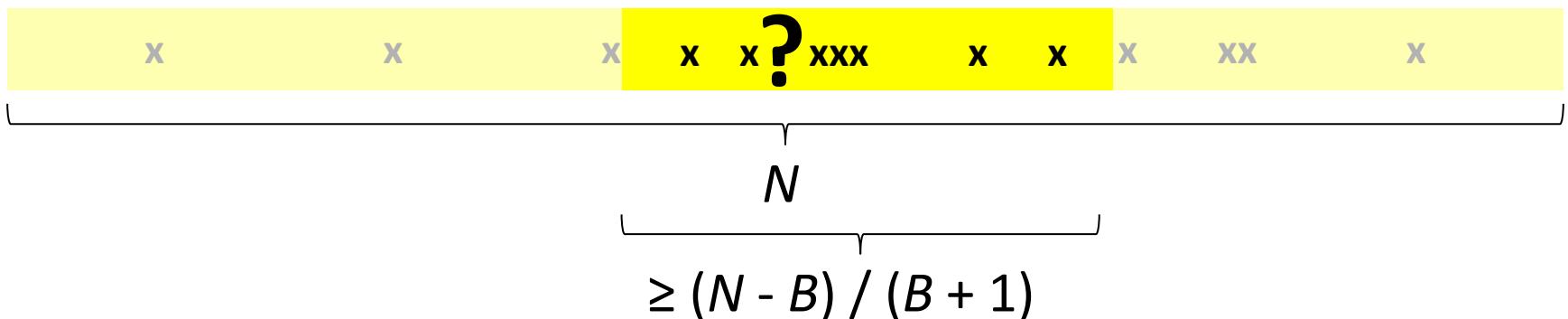
- Searches
Practice: 4-5 I/Os
- Repeated searching
Practice: 1-2 I/Os

The Bad News...

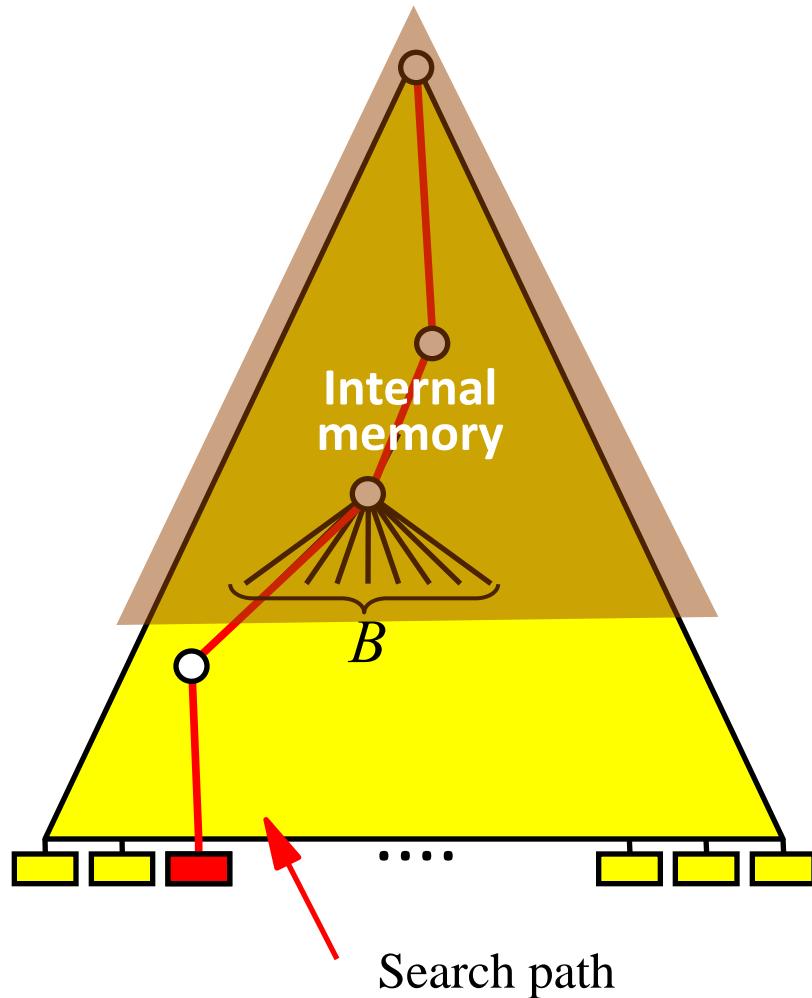
Searching any external memory dictionary
(incl. B-trees) requires worst-case

$$\Omega(\log_B N) \text{ I/Os}$$

Proof idea:



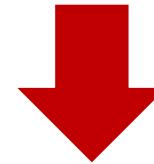
B-trees



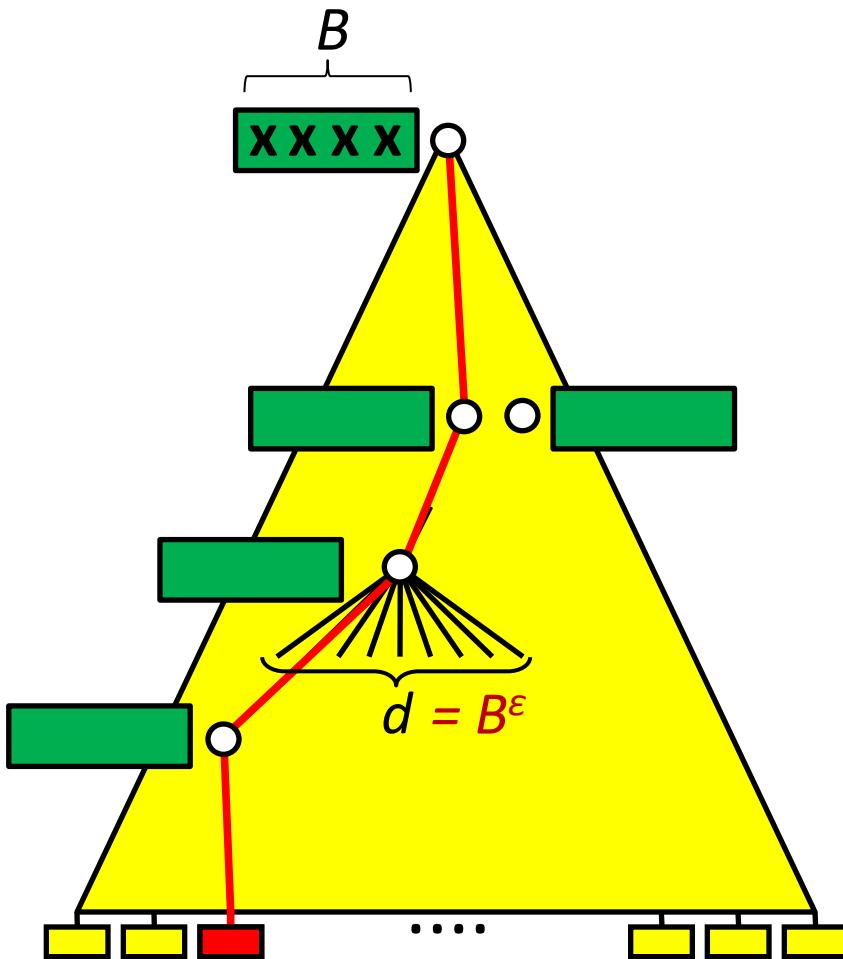
- Searches $O(\log_B N)$ I/Os

- Updates $O(\log_B N)$ I/Os

Best possible



B-trees with Buffered Updates



- Searches cost

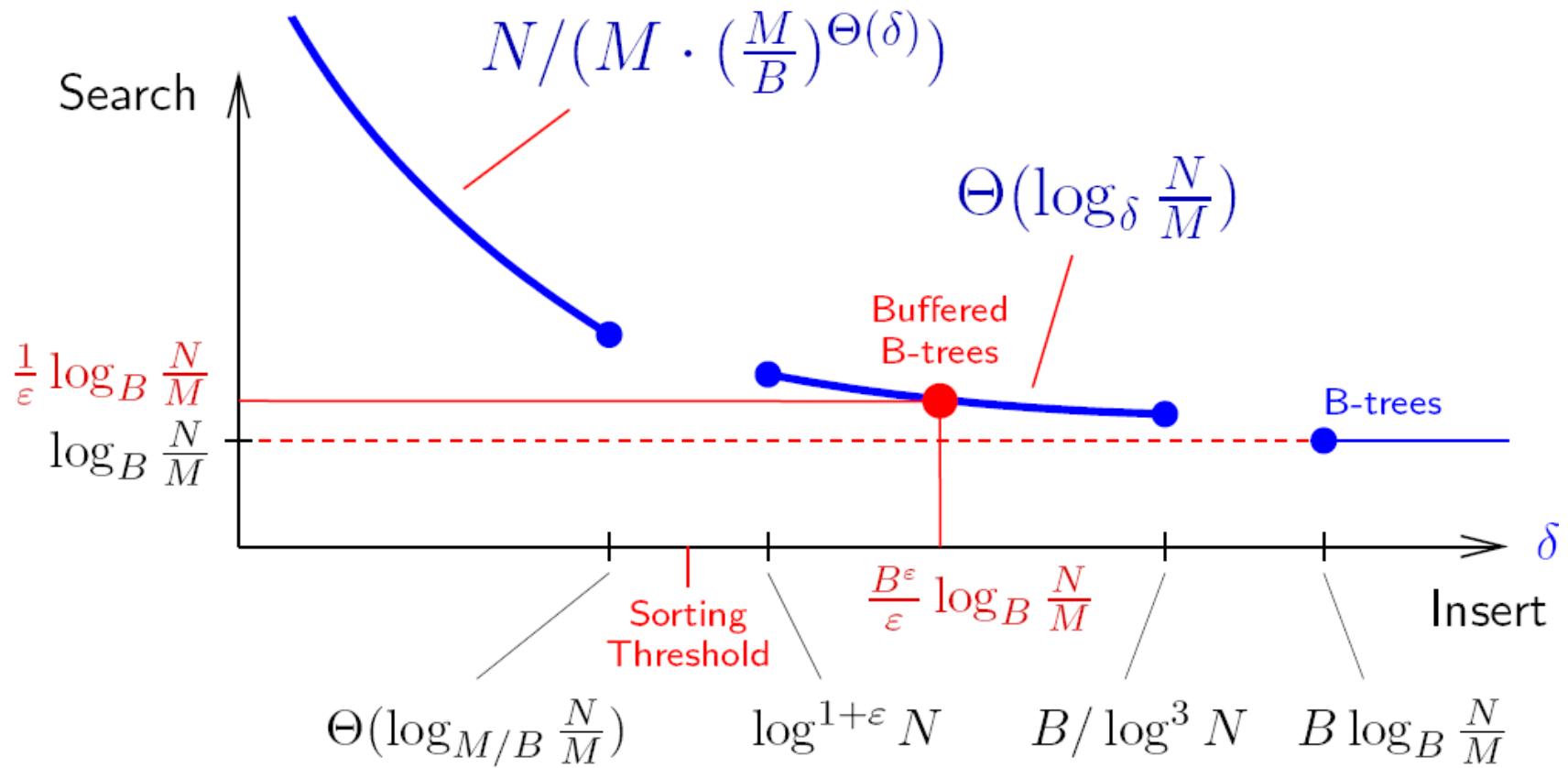
$$\begin{aligned}
 & O(\log_d N) \text{ I/Os} \\
 & = O(\log_B N \cdot 1/\varepsilon) \text{ I/Os} \\
 & \qquad \qquad \qquad \underbrace{2}_{\text{2}}
 \end{aligned}$$

- N updates cost

$$\begin{aligned}
 & O(N \cdot \log_d N \cdot d / B) \text{ I/Os} \\
 & = O(N \cdot \log_B N \cdot 1/\varepsilon B^{1-\varepsilon}) \text{ I/Os} \\
 & \qquad \qquad \qquad \underbrace{2 / \sqrt{B}}_{\text{2 / VB}}
 \end{aligned}$$

Trade-off between search and update times – optimal !

B-trees with Buffered Updates



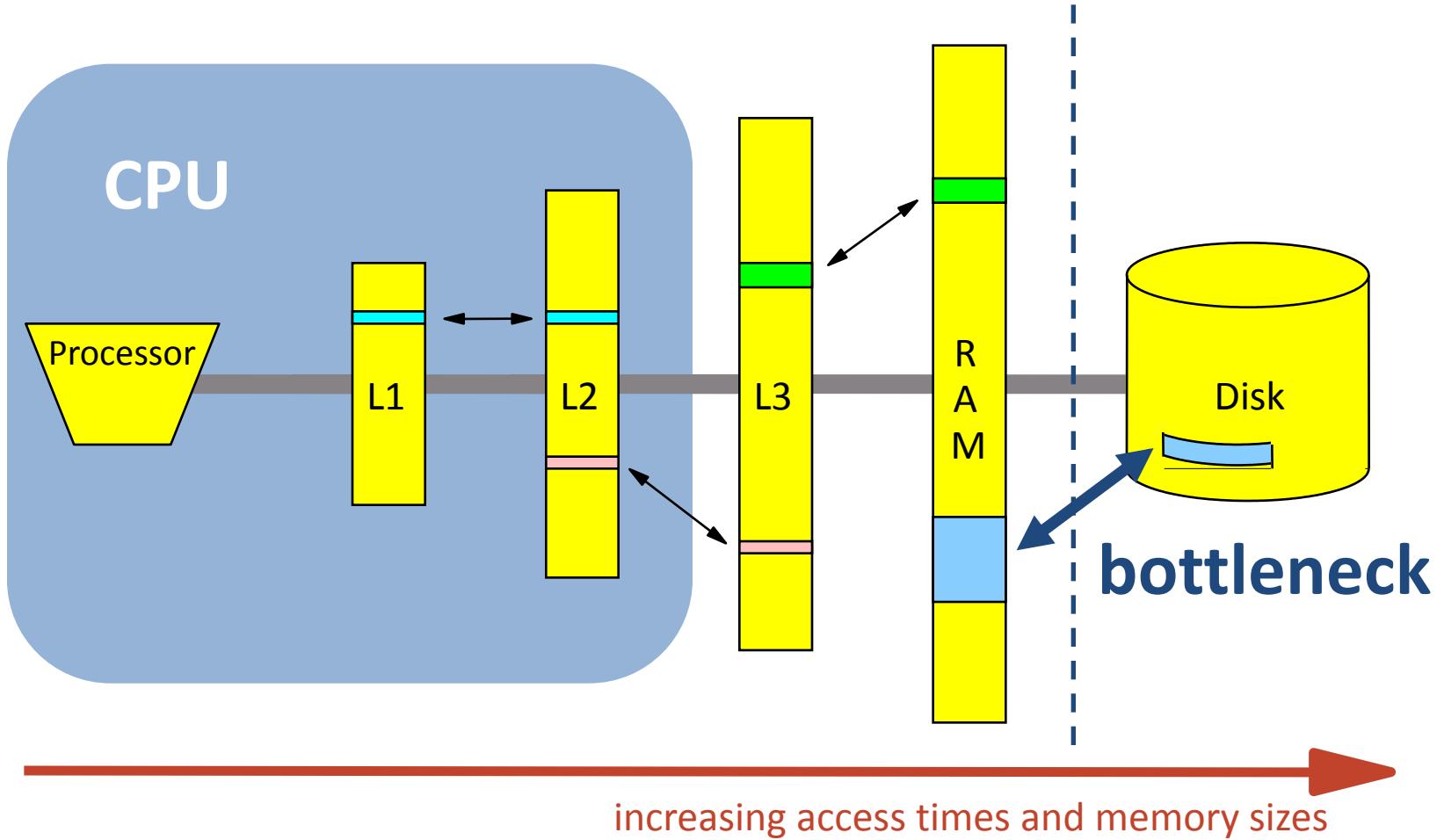
B-trees with Buffered Updates

Experimental Study

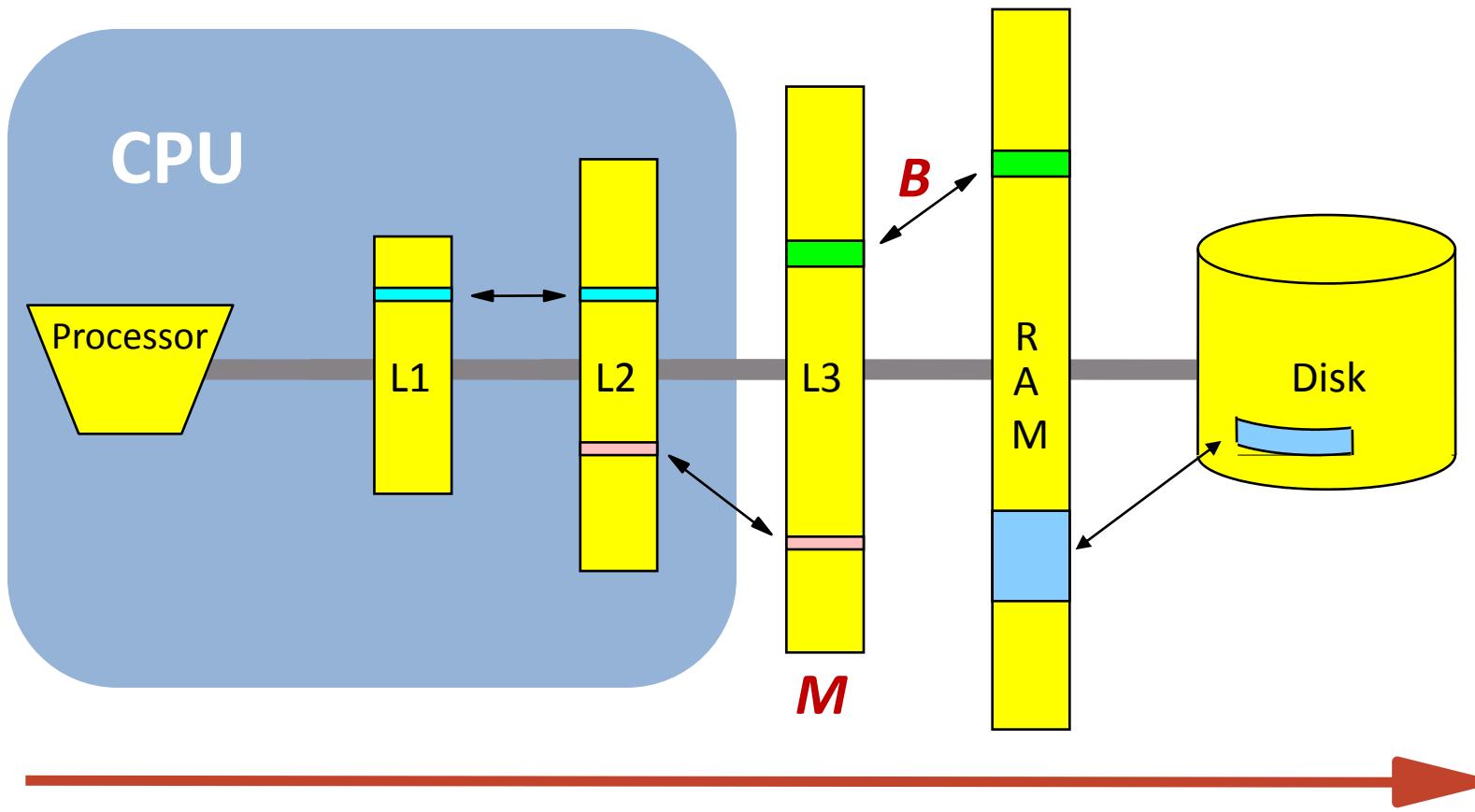
- 100.000.000 elements
- Search time basically unchanged with buffers
- 100 times faster updates

Assumptions until now :

B is known + one level is the bottleneck



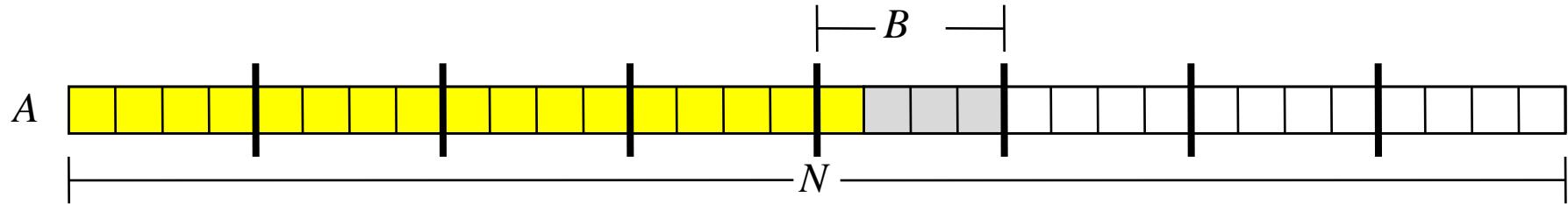
Cache-Obliviousnes Model



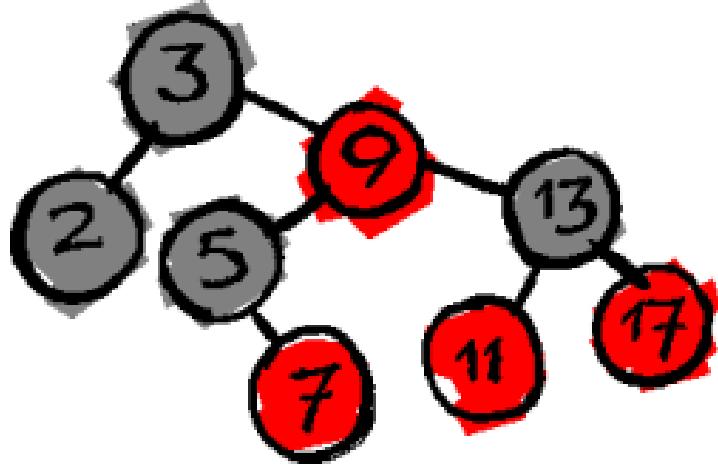
Algorithm does not know B and M
(assume optimal offline cache-replacement strategy)

I/O Efficient Scanning

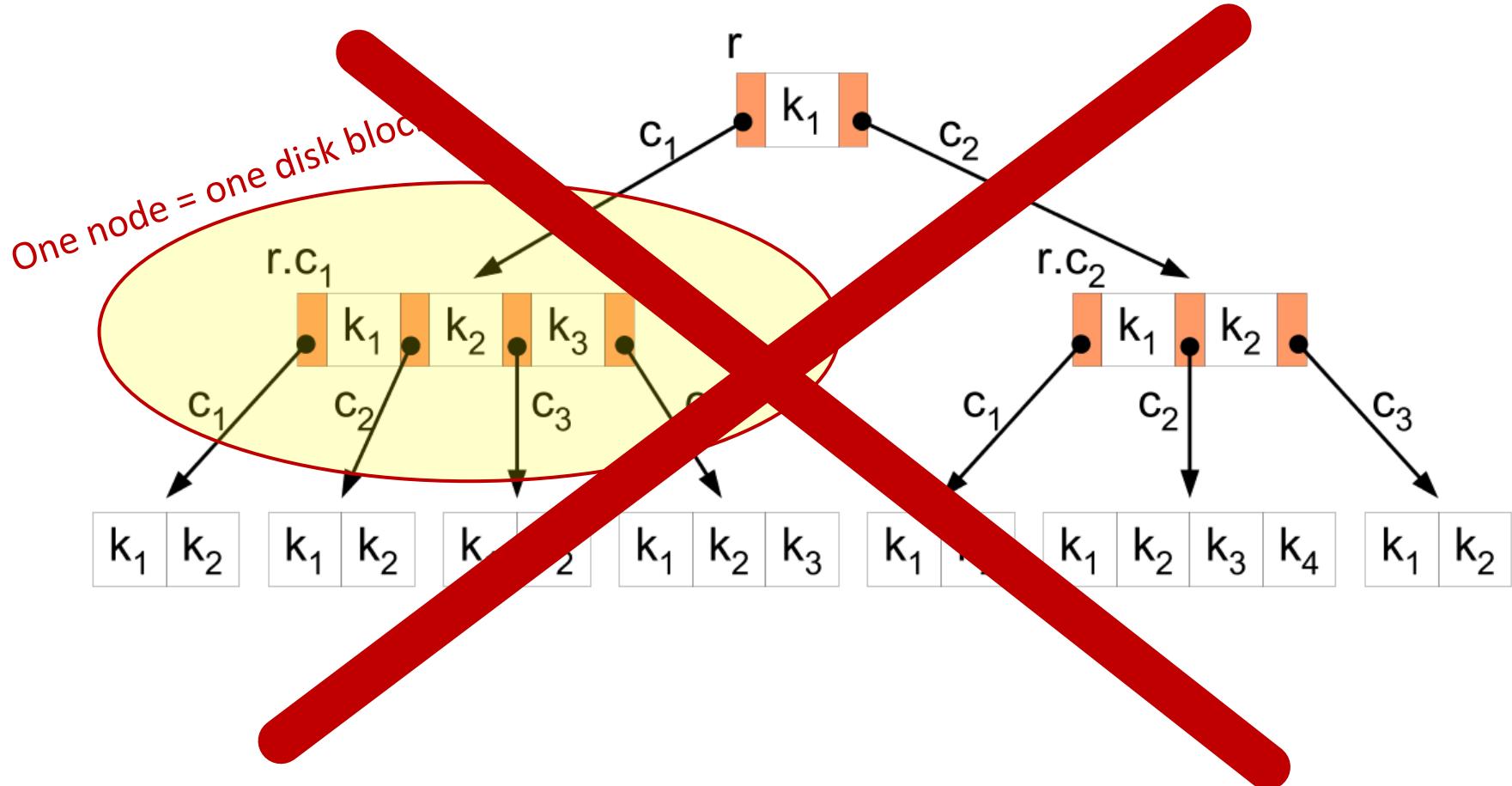
```
sum = 0  
for i = 1 to N do sum = sum + A[i]
```



$O(N/B)$ I/Os

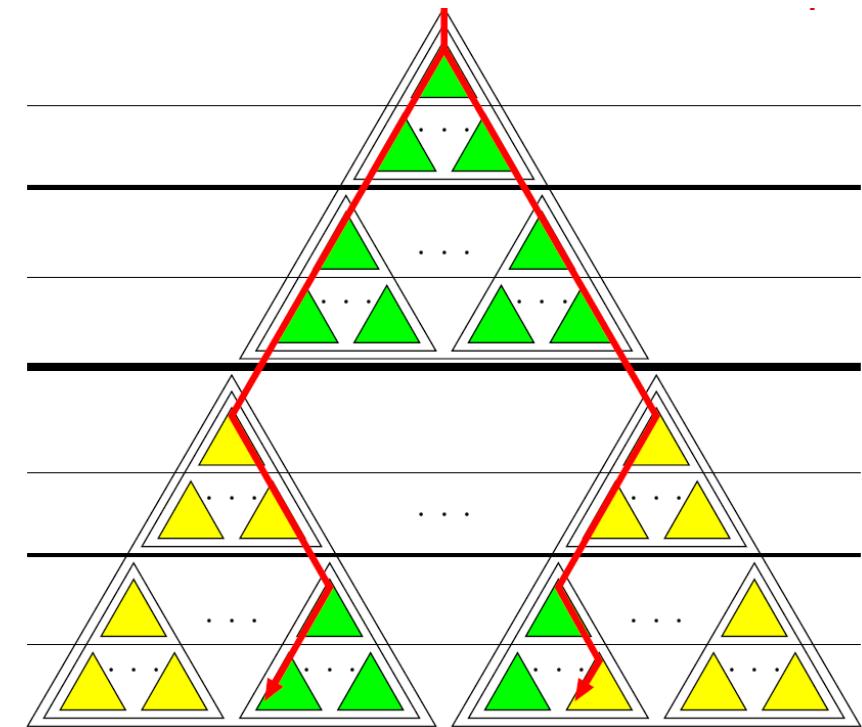
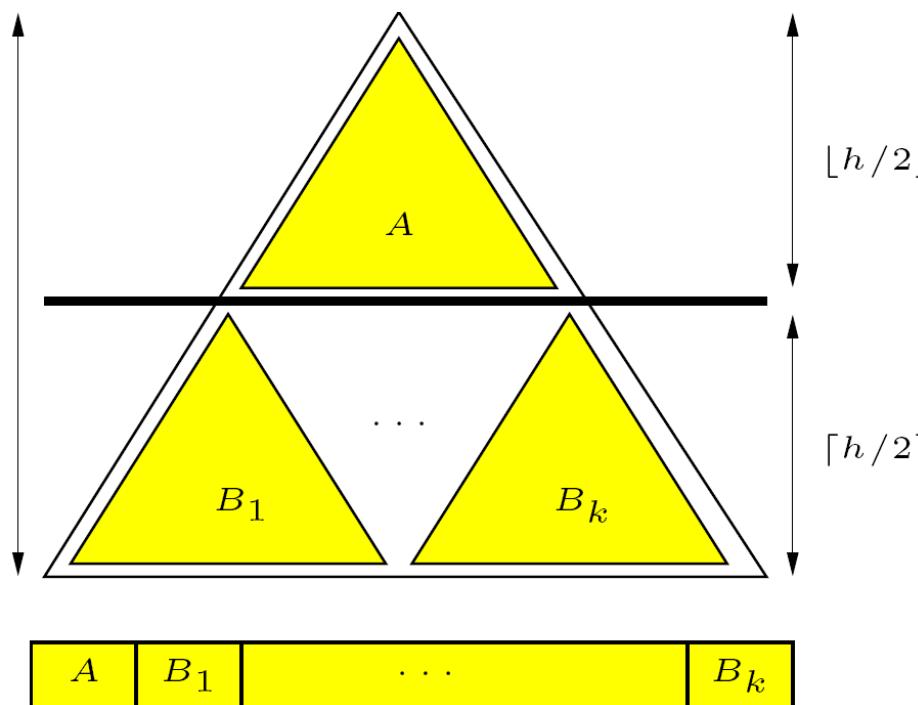


B(oing)-Trees



Height $O(\log_B N)$

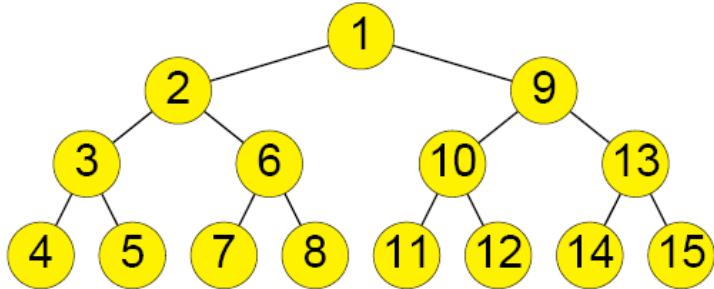
Recursive Search Tree Layout (Cache-Obliviousness)



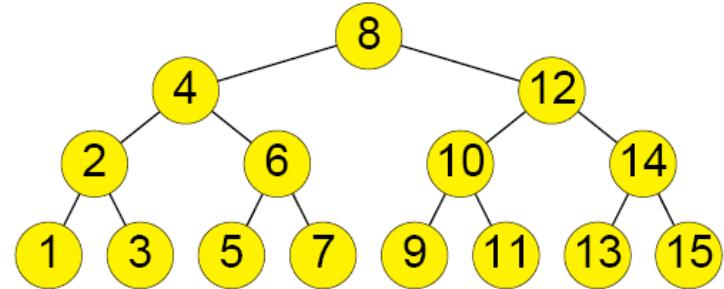
Searches require $O(\log_B N)$ I/Os

(small subtrees have $\sqrt{B} \leq$ size $\leq B$, a path traverses at most $2 \cdot \log_B N$ trees)

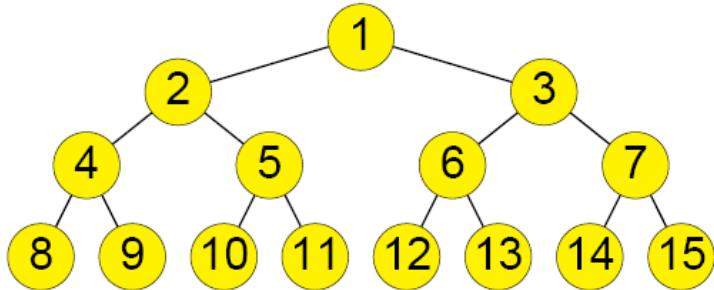
Experiments : Binary Tree Layout



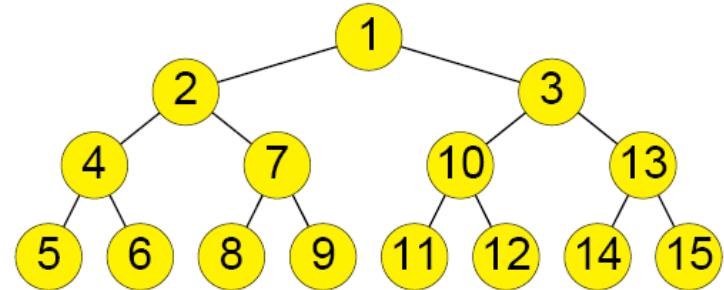
DFS



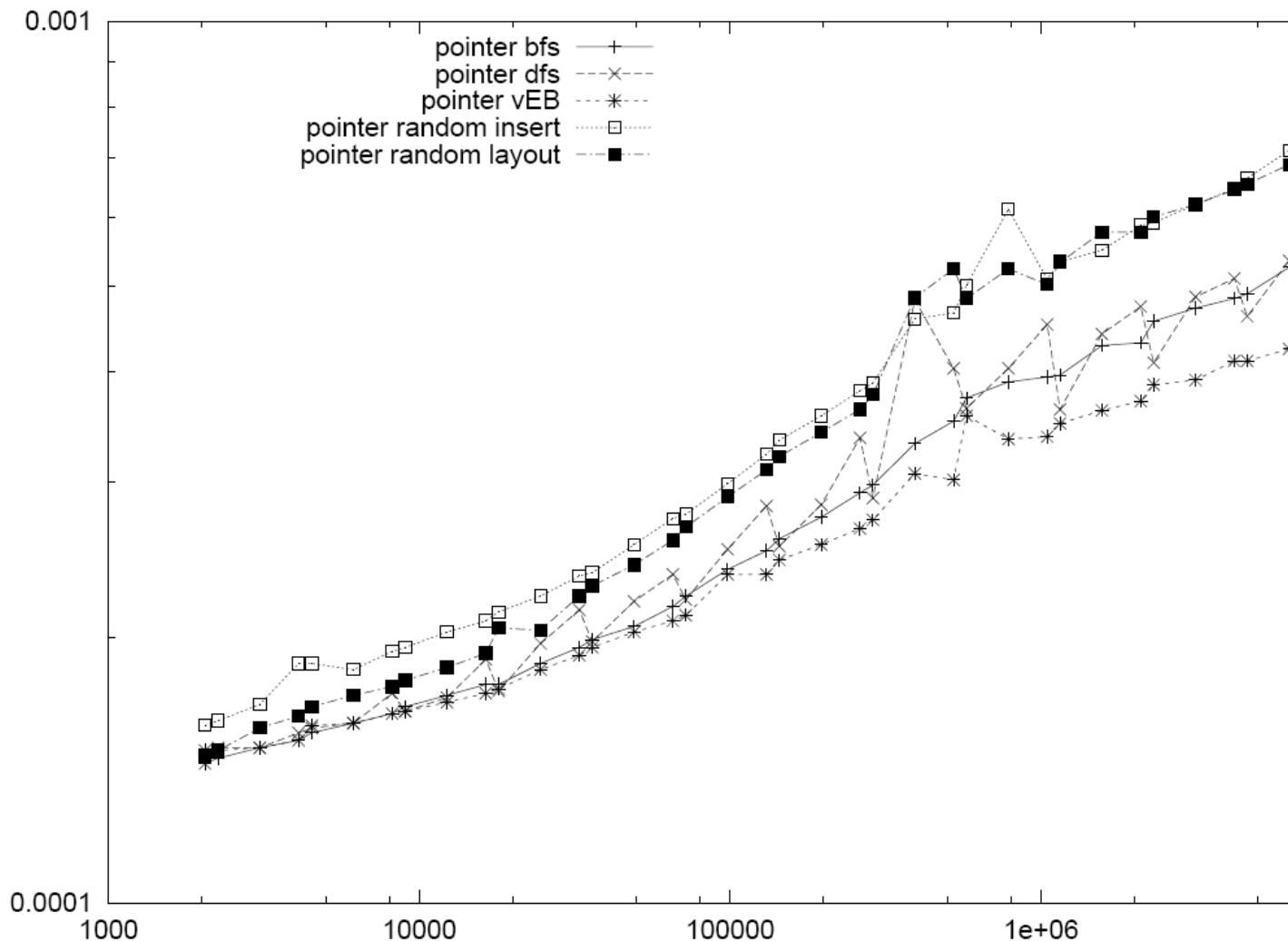
inorder



BFS

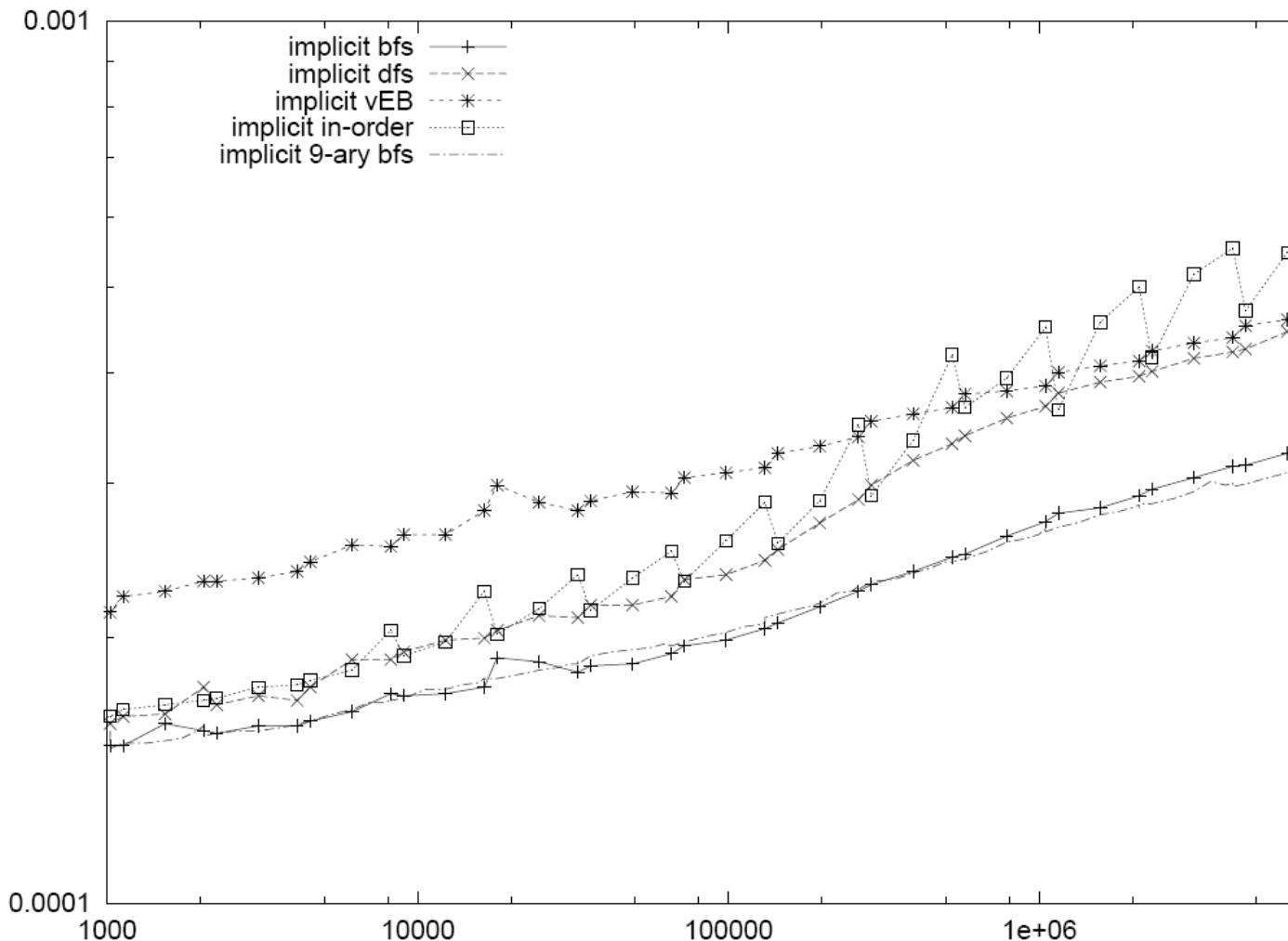
van Emde Boas
(in theory best)

Searches with Pointer Layout



- van Emde Boas layout wins, followed by the BFS layout van Emde Boas layout

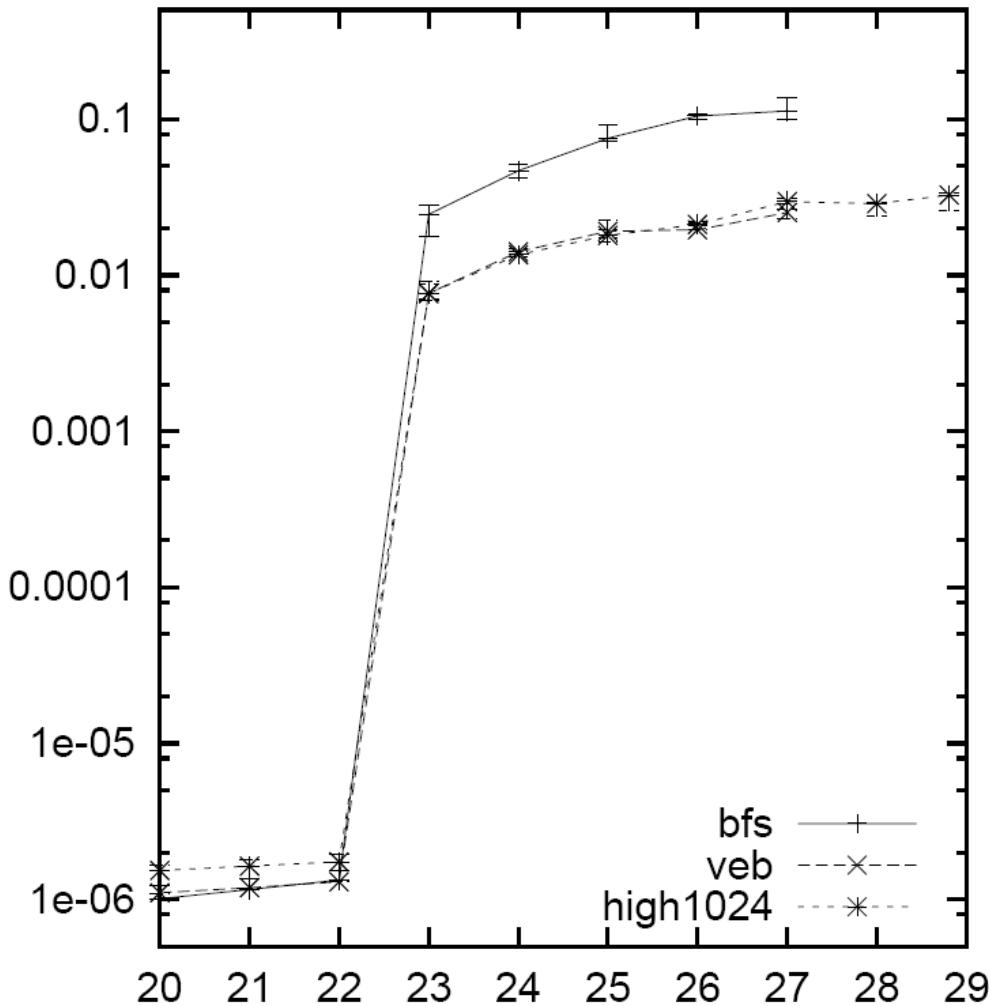
Searches with Implicit Layout



- BFS layout wins due to simplicity and caching of topmost levels
- van Emde Boas layout requires quite complex index computations

Searches with Pointer Layout

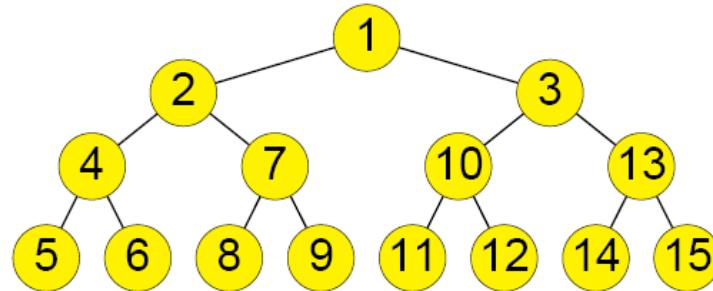
- Beyond Main Memory



- van Emde Boas layout wins, followed by the BFS layout

Optimal Cache-Oblivious Static Index

= van Emde Boas Layout of
Complete Binary Tree

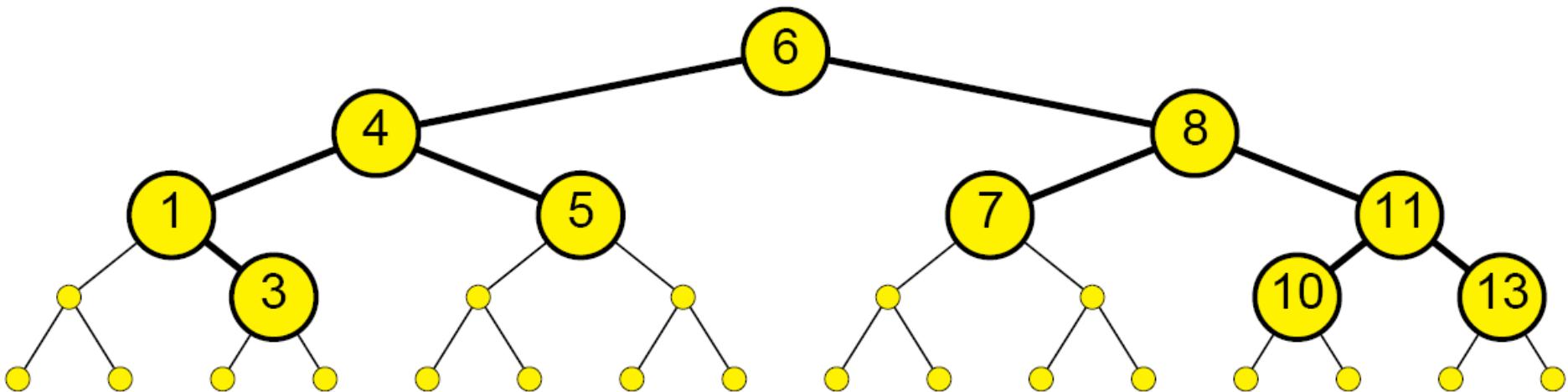


van Emde Boas

Cache-Oblivious Dynamic Index ?

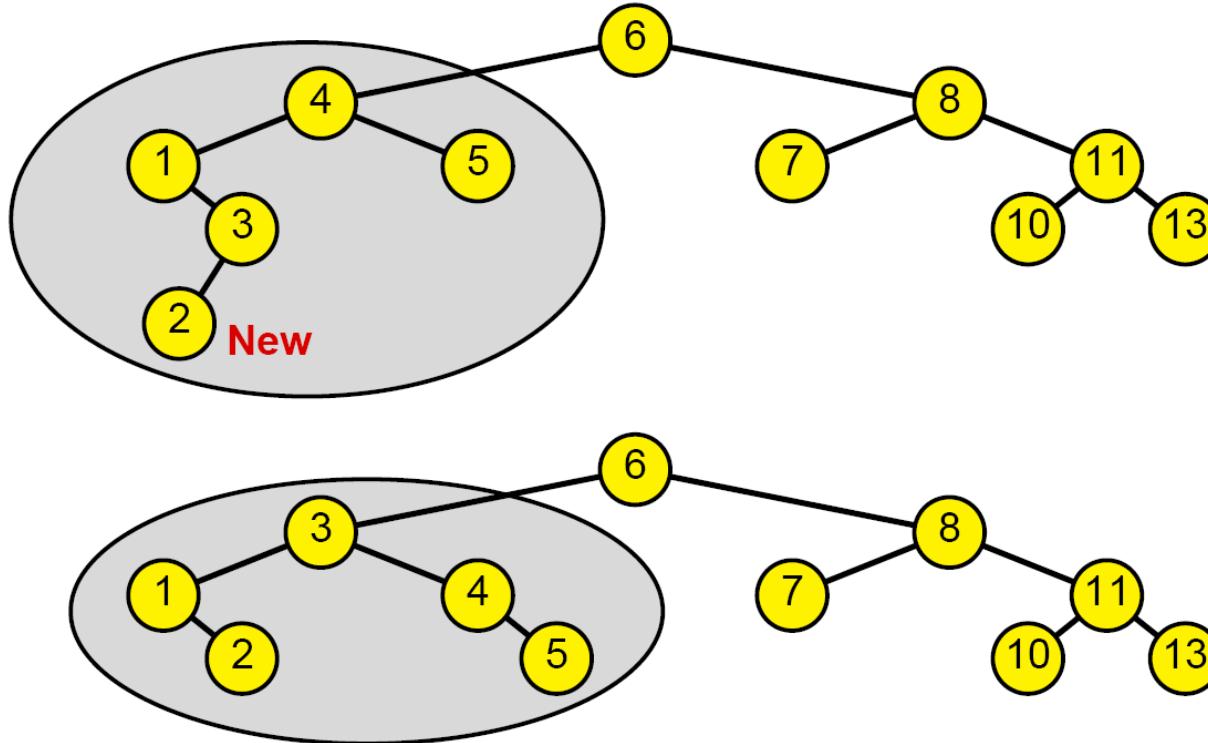
Binary Search Trees

Dynamic



- **Embed** a dynamic tree of small height into a complete tree
- Static van Emde Boas layout
- Rebuild data structure whenever N doubles or halves

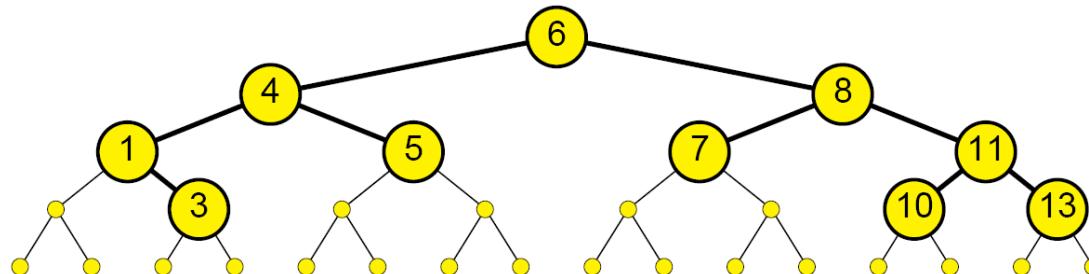
Rebalancing Subtree during Insertion



- If an insertion causes non-small height then **rebuild** subtree at nearest ancestor with sufficient few descendants
- Insertions require amortized time $O(\log^2 N)$

Optimal Cache-Oblivious Dynamic Index

- Search $O(\log_B N)$ I/Os (optimal)
- Updates $O(\log_B N + (\log^2 N) / B)$ I/Os

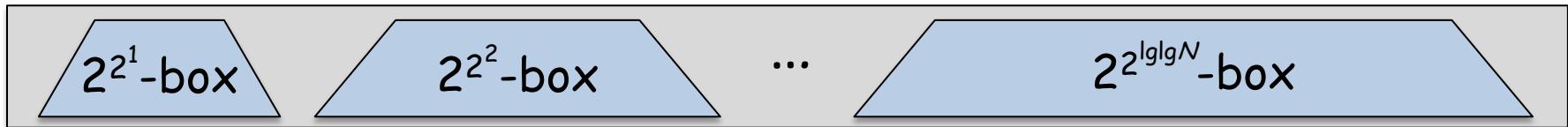


Cache-Oblivious Index with Query-Update Trade-off ?

Cache-Oblivious Dynamic Dictionaries with Query-Update Trade-off

- Solution matching the I/O complexity of Buffered B-trees
- Searches $O(\log_B N \cdot 1/\varepsilon)$ I/Os
- N updates $O(N \cdot \log_B N \cdot 1/\varepsilon B^{1-\varepsilon})$ I/Os

xDict



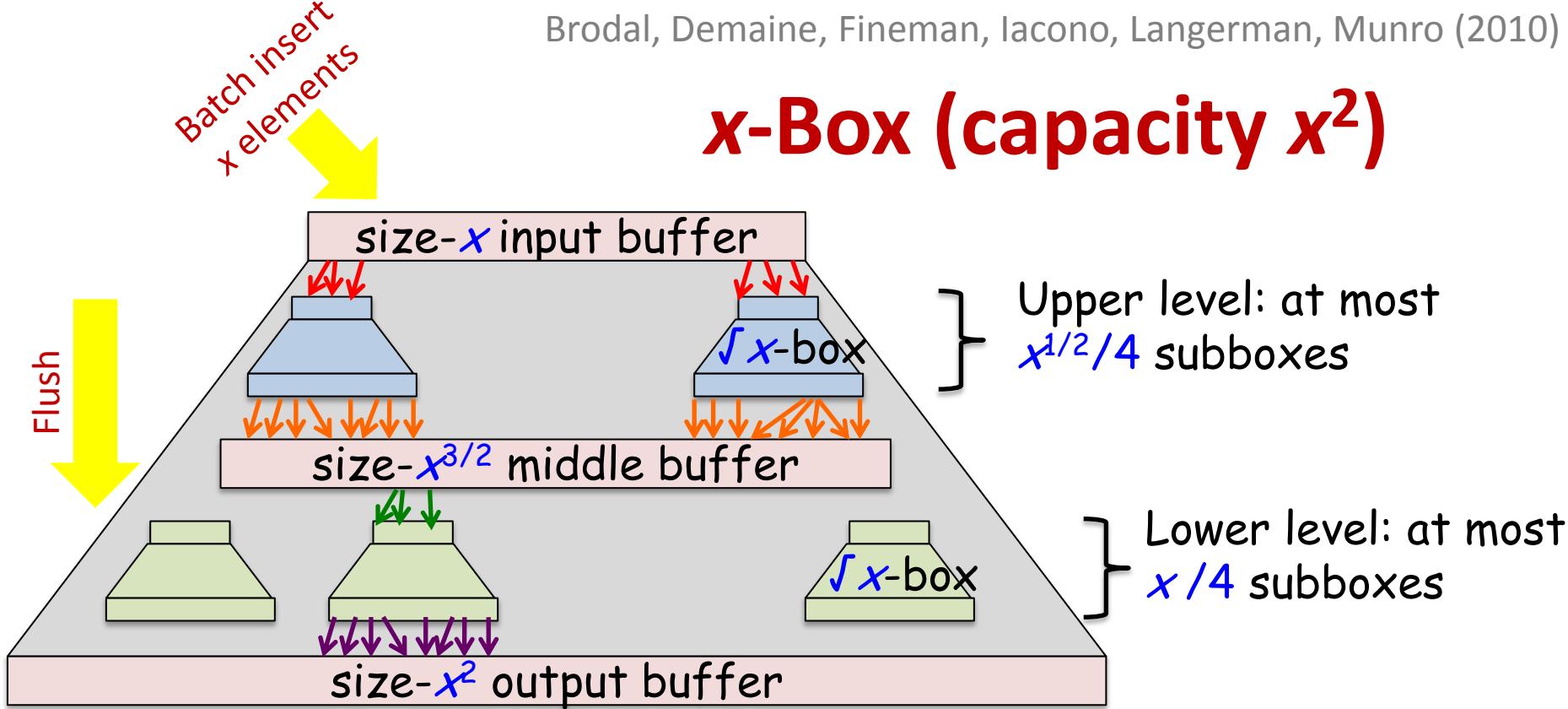
Insert

- Insert into smallest box
- When a box reaches capacity, **Flush** it and **Batch-Insert** into the next box

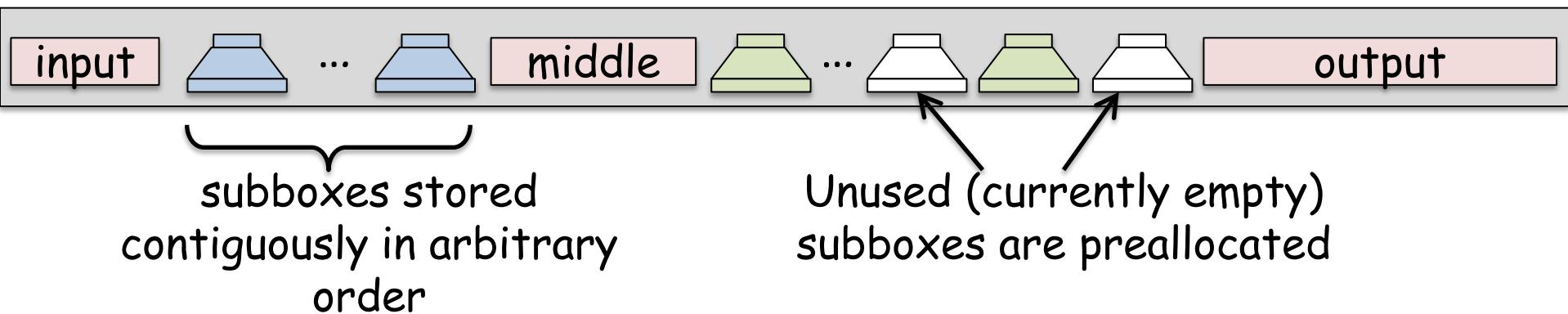
Search

- Search in each x-box
- $O(\log_B x)$ cost is dominated by largest box $O(\log_B N)$

x -Box (capacity x^2)



Memory Layout



External Memory Index

- Searches $O(\log_B N \cdot 1/\varepsilon)$ I/Os
- N updates $O(N \cdot \log_B N \cdot 1/\varepsilon B^{1-\varepsilon})$ I/Os

