

Coping with the Memory Hierarchy the Cache-Oblivious Way

Rolf Fagerberg

University of Aarhus

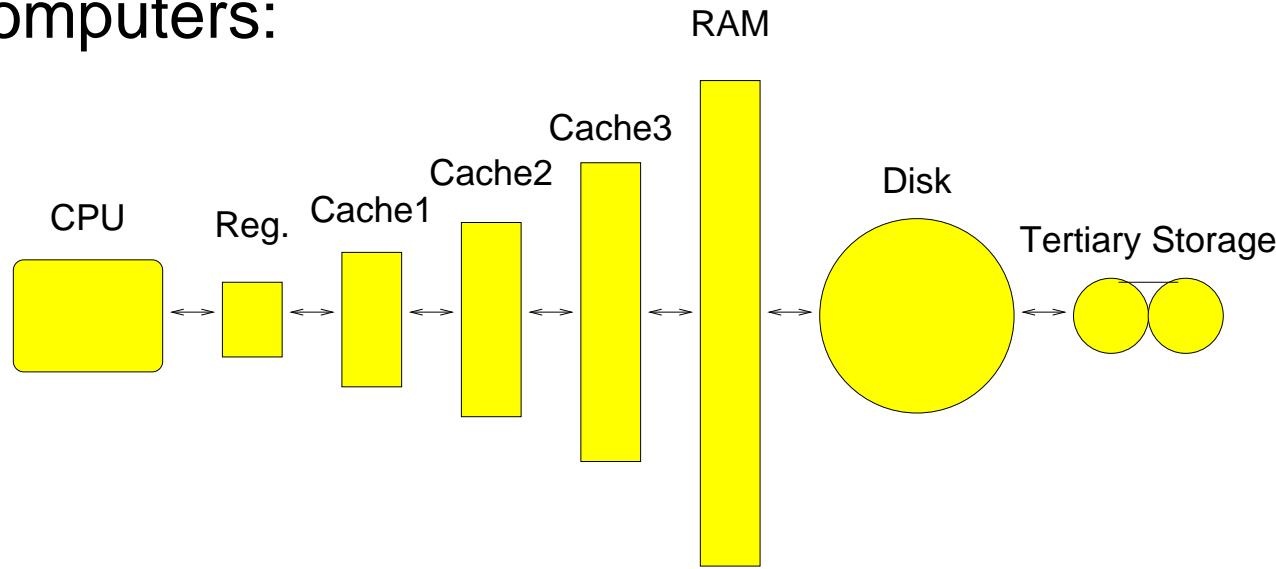
Imada, SDU, February 18, 2004

Overview

- The memory hierarchy
- The I/O-model
- The cache-oblivious model
- Examples of cache-oblivious algorithms
 - Double for-loop (with applications)
 - Searching
 - Sorting
- Theoretical limits of cache-obliviousness

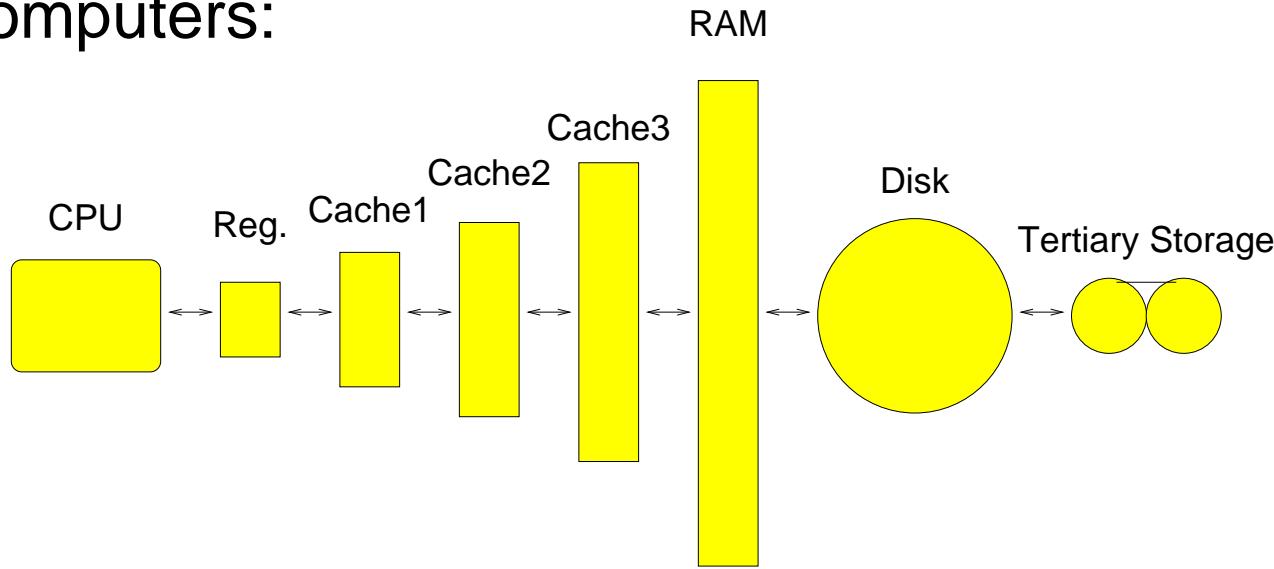
The Memory Hierarchy

Modern computers:



The Memory Hierarchy

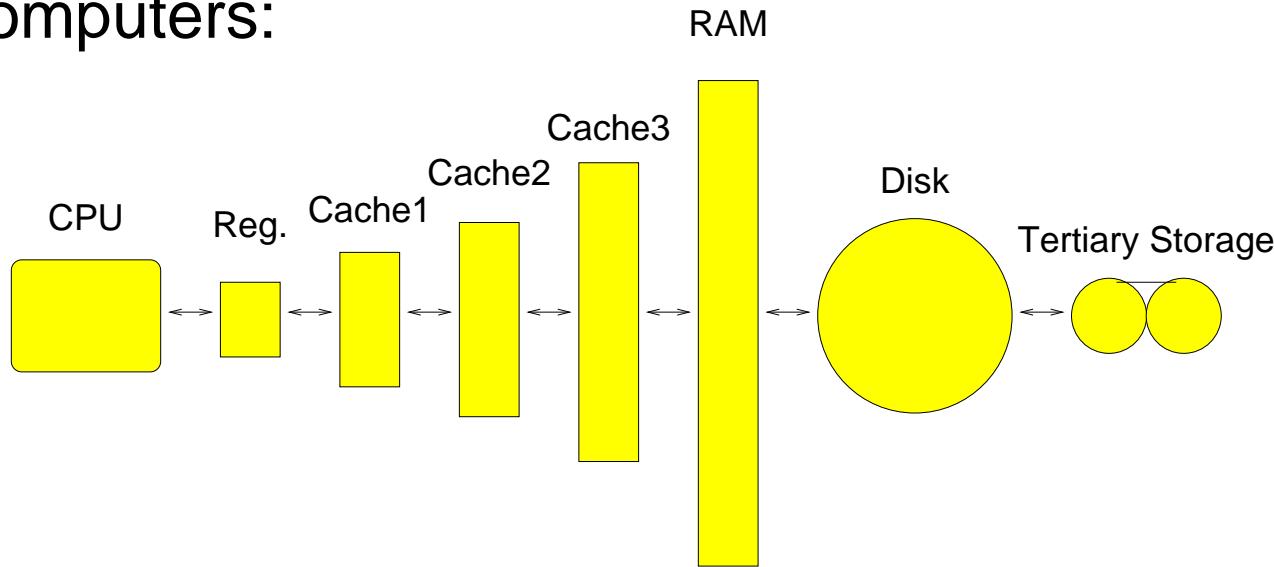
Modern computers:



	<i>Access time</i>	<i>Volume</i>
Registers	1 cycle	1 Kb
Cache	10 cycles	512 Kb
RAM	100 cycles	512 Mb
Disk	20,000,000 cycles	80 Gb

The Memory Hierarchy

Modern computers:



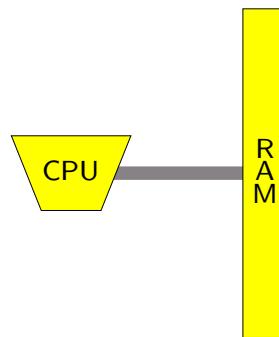
Gap increases over time.

	<i>Access time</i>	<i>Volume</i>
Registers	1 cycle	1 Kb
Cache	10 cycles	512 Kb
RAM	100 cycles	512 Mb
Disk	20,000,000 cycles	80 Gb

Real problems of **Gigabyte**, **Terabyte**, and even **Petabyte** size: Databases (finance, phone companies, banks, weather, geology, geography, astronomy), WWW, GIS systems, computer graphics.

Classic RAM Model

The RAM model:



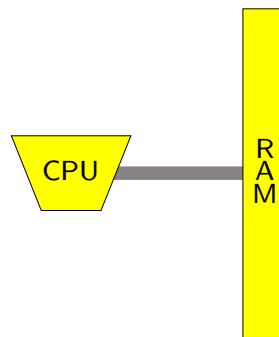
Add: $O(1)$

Branch: $O(1)$

Mem access: $O(1)$

Classic RAM Model

The RAM model:



Add: $O(1)$

Branch: $O(1)$

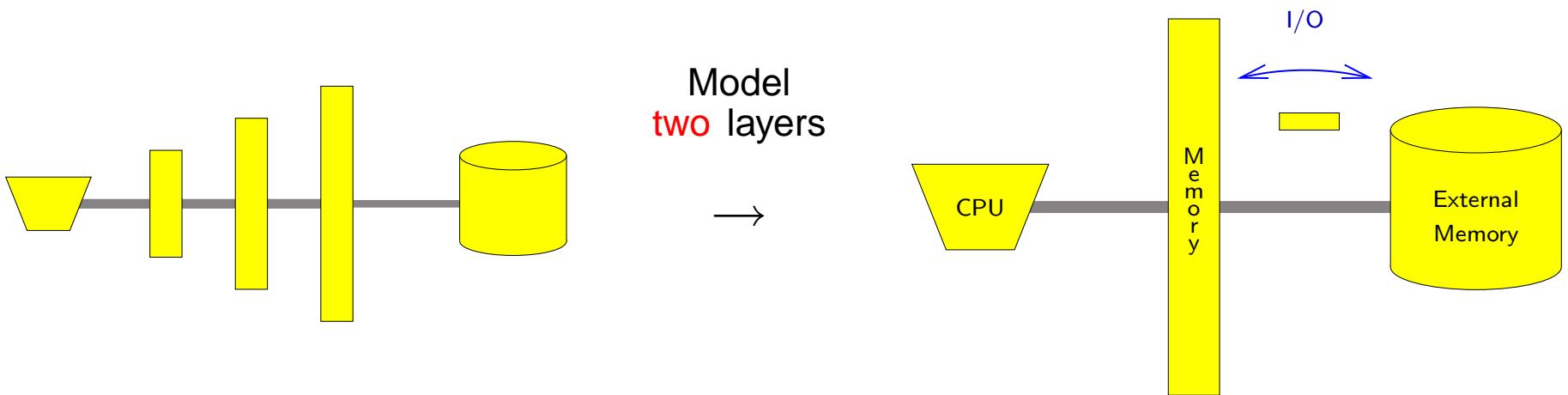
~~Mem access: $O(1)$~~

Increasingly inadequate

Overview

- ✓ The memory hierarchy
- The I/O-model
- The cache-oblivious model
- Examples of cache-oblivious algorithms
 - Double for-loop (with applications)
 - Searching
 - Sorting
- Theoretical limits of cache-obliviousness

I/O Model



N = problem size

M = memory size

B = I/O block size

Aggarwal and Vitter 1988

Cost: number of I/Os.

Example

	CPU time	Inplace	Worstcase
Heapsort	$N \log N$	✓	✓
Quicksort	$N \log N$	✓	
Mergesort	$N \log N$		✓

Example

	CPU time	Inplace	Worstcase	I/O
Heapsort	$N \log N$	✓	✓	$N \log N$
Quicksort	$N \log N$	✓		$(N \log N)/B$
Mergesort	$N \log N$		✓	$(N \log N)/B$

Random memory access \Rightarrow page fault at every access.
Sequential memory access \Rightarrow page fault every B accesses.

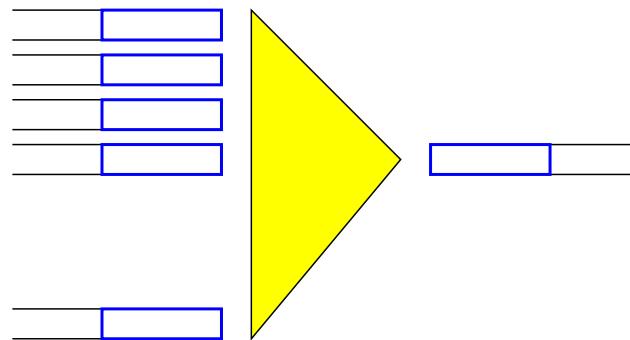
Typically, $B \sim 10^3$

I/O-Optimal Sorting

Binary Mergesort:

$$\frac{N}{B} \log_2 N \text{ I/Os}$$

Multi-Way Merging:



Maximal
merge degree
 $\approx M/B$

Multi-Way Mergesort:

$$\frac{N}{B} \log_{M/B} \frac{N}{M} \text{ I/Os}$$

I/O Model Facts

- Scanning: $\Theta(N/B)$ I/Os.
- Searching: $\Theta(\log_B N)$ I/Os by B -trees.
- Sorting: $\Theta\left(\frac{N}{B} \log_{M/B} \frac{N}{M}\right)$ I/Os by $\frac{M}{B}$ -way merge-sort.
- Permuting: $\Theta\left(\min\left\{N, \frac{N}{B} \log_{M/B} \frac{N}{M}\right\}\right)$ by direct move or sorting

1988-2004:

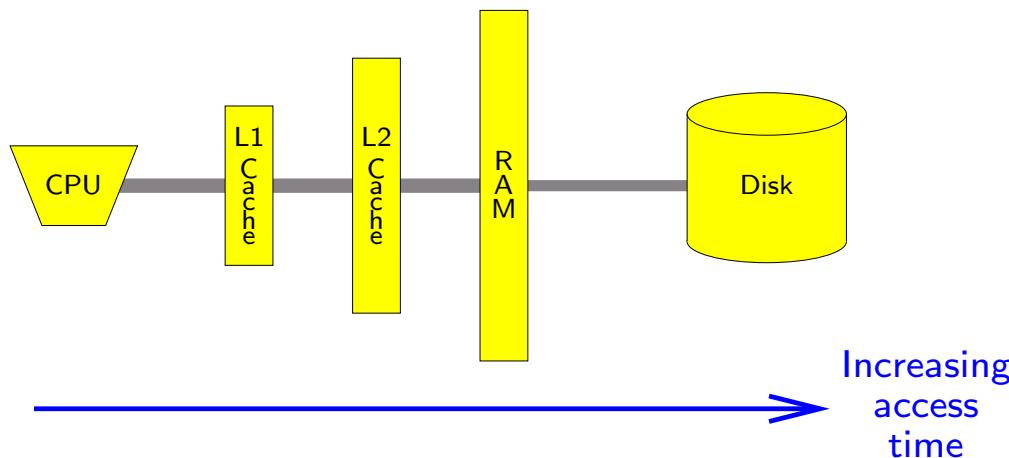
Many algorithms and data structures for problems from computational geometry, graphs, strings, . . .

Overview

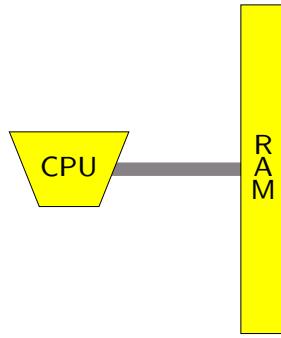
- ✓ The memory hierarchy
- ✓ The I/O-model
- The cache-oblivious model
- Examples of cache-oblivious algorithms
 - Double for-loop (with applications)
 - Searching
 - Sorting
- Theoretical limits of cache-obliviousness

Computer Models

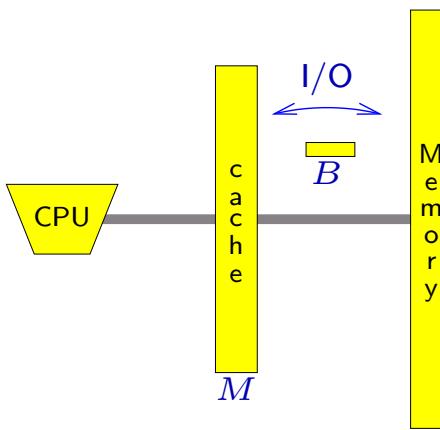
Reality:



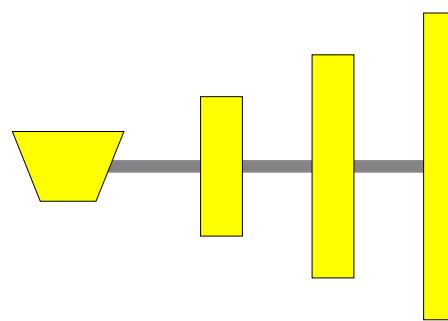
Models:



RAM model



I/O model



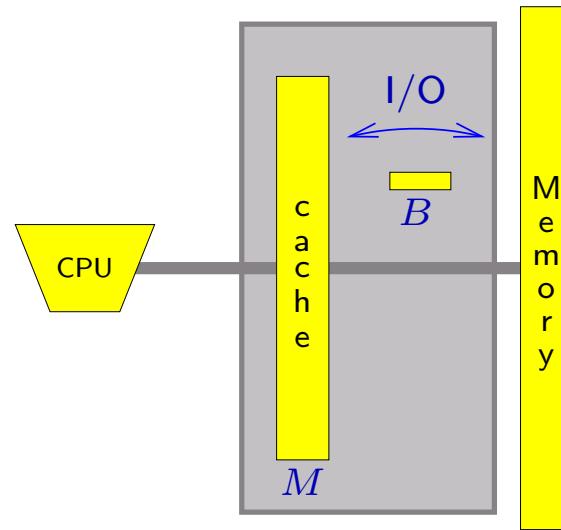
Multi-level
models

Cache-
Oblivious-
ness

New Model

Cache-Oblivious Model

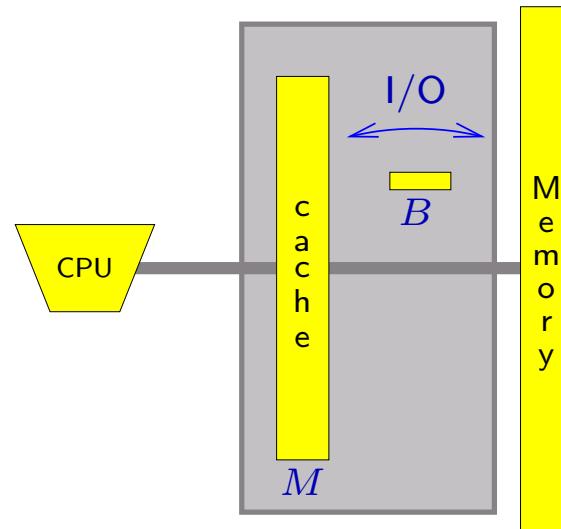
- Program in the RAM model
- Analyze in the I/O model for arbitrary B and M
- Optimal off-line cache replacement strategy



Frigo, Leiserson, Prokop, Ramachandran, FOCS'99

Cache-Oblivious Model

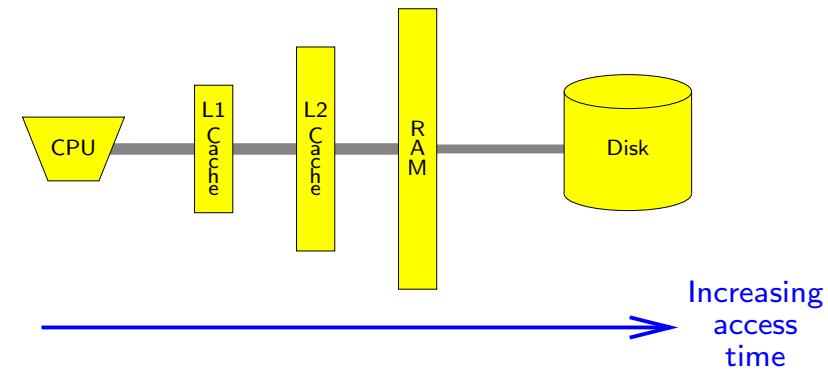
- Program in the RAM model
- Analyze in the I/O model for arbitrary B and M
- Optimal off-line cache replacement strategy



Advantages:

Frigo, Leiserson, Prokop, Ramachandran, FOCS'99

- Optimal on arbitrary level \Rightarrow optimal on all levels
- Portability
- Simplicity of model.



Cache-Oblivious Results

Scanning ⇒ stack, queue, selection, . . .

Cache-Oblivious Results

Scanning ⇒ stack, queue, selection,

Matrix multiplication, FFT: FOCS'99

Sorting: FOCS'99, ICALP'02, ALENEX'04

Search trees: Prokop 99, FOCS'00,
WAE'01, SODA'02 × 2,
ESA'02, FOCS'03

Priority queues: STOC'02, ISAAC'02

Graph algorithms: STOC'02, BRICS-04-2

Computational geometry: 2 × ICALP'02 , SCG'03

Scanning dynamic sets: ESA'02

Power of cache-obliviousness: STOC'03

Cache-Oblivious Results

Scanning ⇒ stack, queue, selection,

Matrix multiplication, FFT: FOCS'99

Sorting: FOCS'99, ICALP'02, ALENEX'04

Search trees: Prokop 99, FOCS'00,
WAE'01, SODA'02 × 2,
ESA'02, FOCS'03

Priority queues: STOC'02, ISAAC'02

Graph algorithms: STOC'02, BRICS-04-2

Computational geometry: 2 × ICALP'02 , SCG'03

Scanning dynamic sets: ESA'02

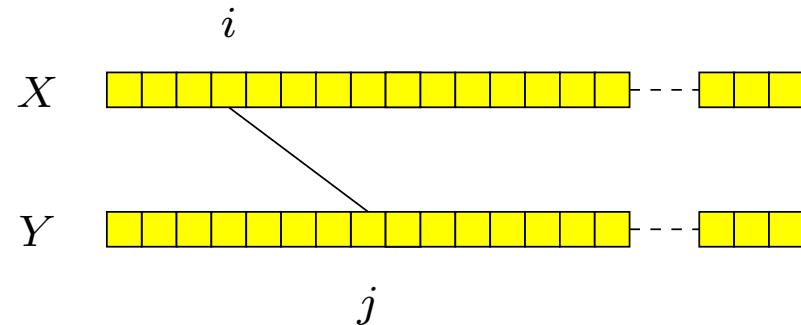
Power of cache-obliviousness: STOC'03

Overview

- ✓ The memory hierarchy
- ✓ The I/O-model
- ✓ The cache-oblivious model
- Examples of cache-oblivious algorithms
 - Double for-loop (with applications)
 - Searching
 - Sorting
- Theoretical limits of cache-obliviousness

Double for-loop

X, Y arrays of length n :



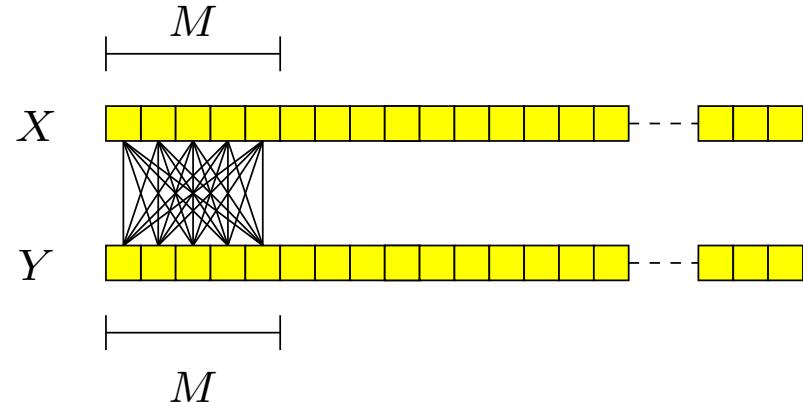
```
for(i=0; i<n; i++)
    for(j=0; j<n; j++)
        f(X[i], Y[j])
```

I/O complexity:

$$n \times \frac{n}{B} = \frac{n^2}{B}$$

Double for-loop

More efficient version
in the I/O-model:



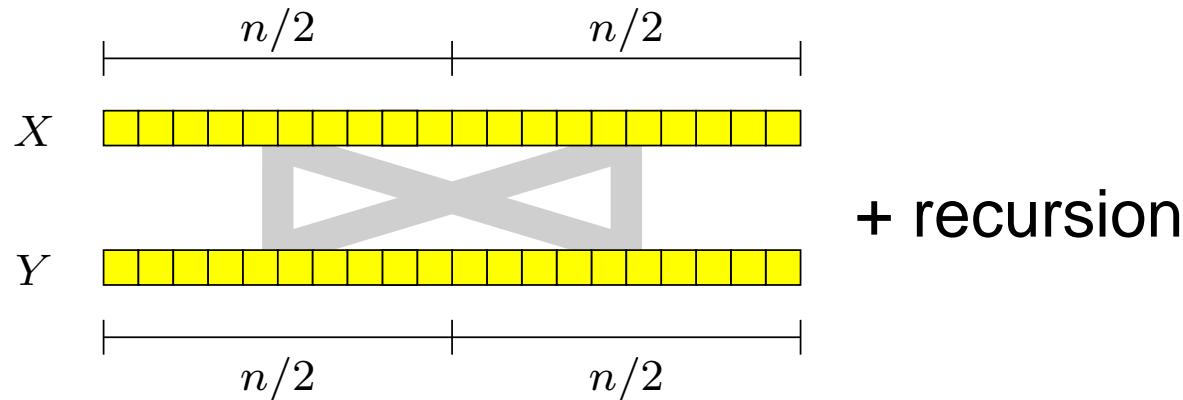
I/O complexity:

$$\frac{n}{M} \times \frac{n}{M} \times \frac{M}{B} = \frac{n^2}{MB}$$

```
for(i=0; i<n; i=i+M)
    for(j=0; j<n; j++)
        for(k=i; k<i+M; k++)
            f(X[i], Y[j])
```

Double for-loop

Cache-oblivious version:



I/O complexity:

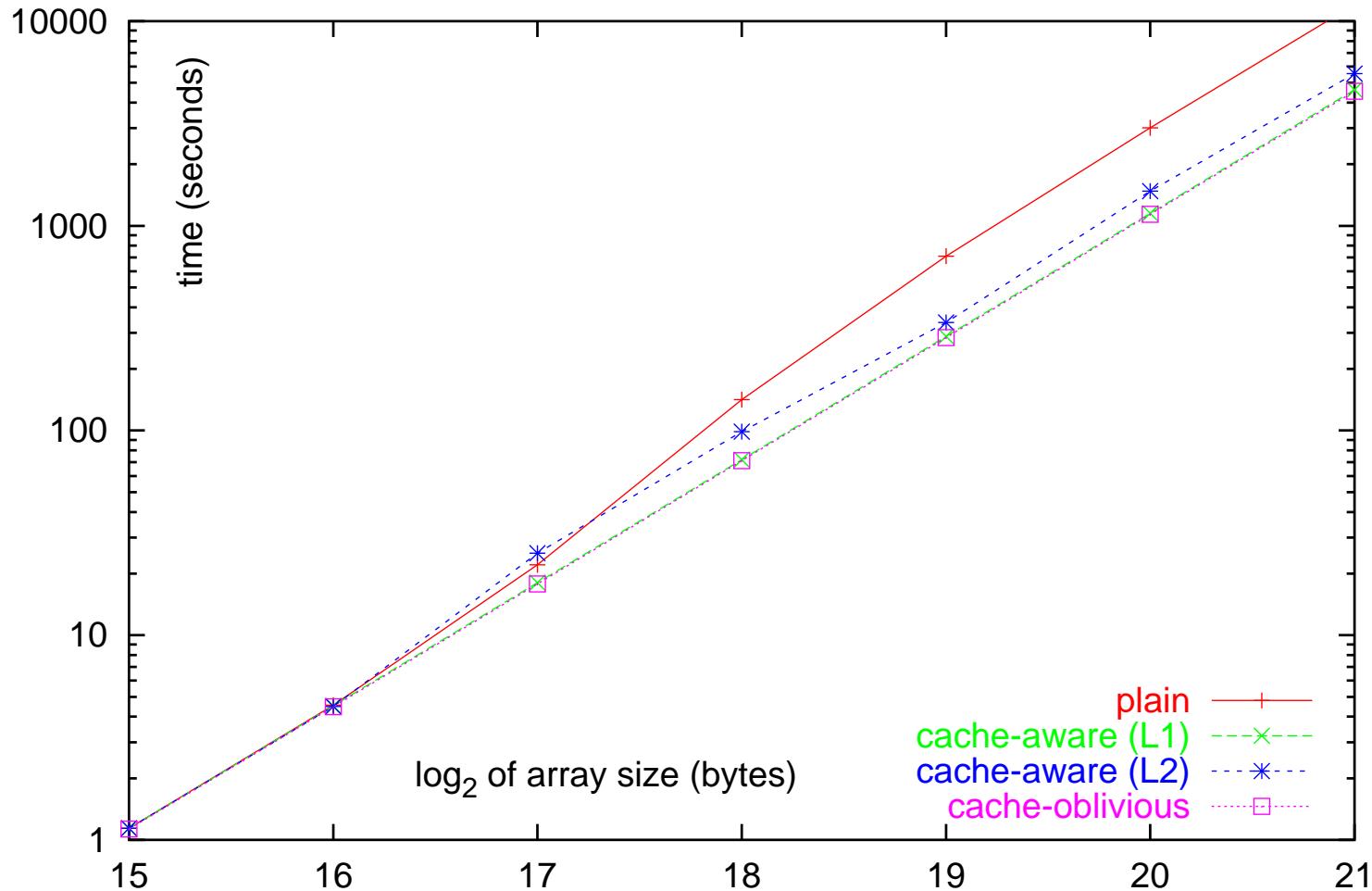
Again $\frac{n^2}{MB}$

Double for-loop

Cache-oblivious version

```
doubleLoop(i, j, length)
    if (length == 1)
        f(X[i], Y[j])
    else
        doubleLoop(i, j, length/2)
        doubleLoop(i, j+length/2, length/2)
        doubleLoop(i+length/2, j, length/2)
        doubleLoop(i+length/2, j+length/2, length/2)
```

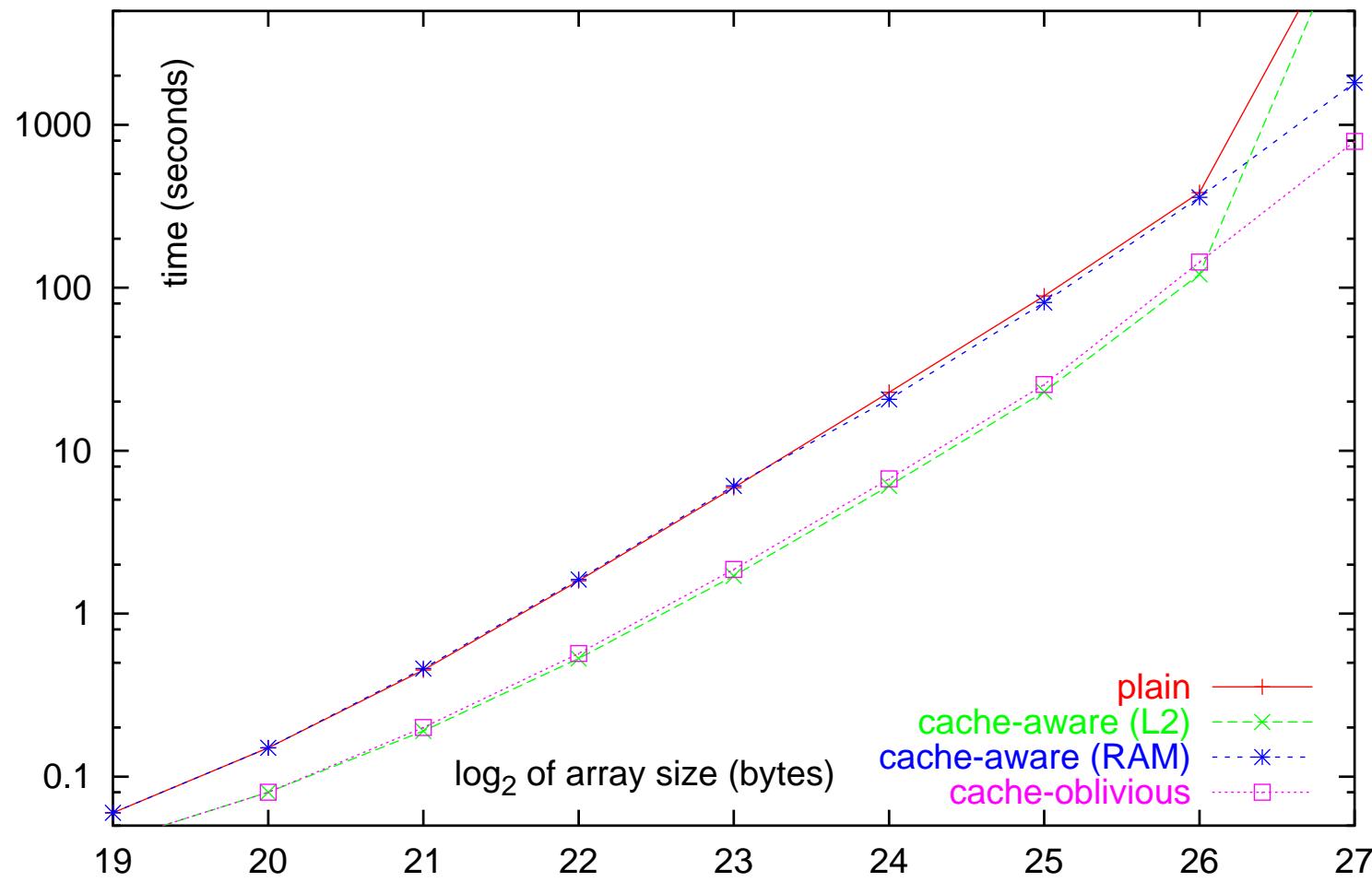
Experiments



Sizes within RAM (element size 4 bytes)

366 MHz Pentium II, 128 MB RAM, 256 KB Cache, gcc -O3, Linux

Experiments



Sizes exceeding RAM (element size 1 KB)

366 MHz Pentium II, 128 MB RAM, 256 KB Cache, gcc -O3, Linux

For-loop Applications

Join in databases

Dynamic programming
(bioinformatics)

Matrix multiplication
(scientific computing)

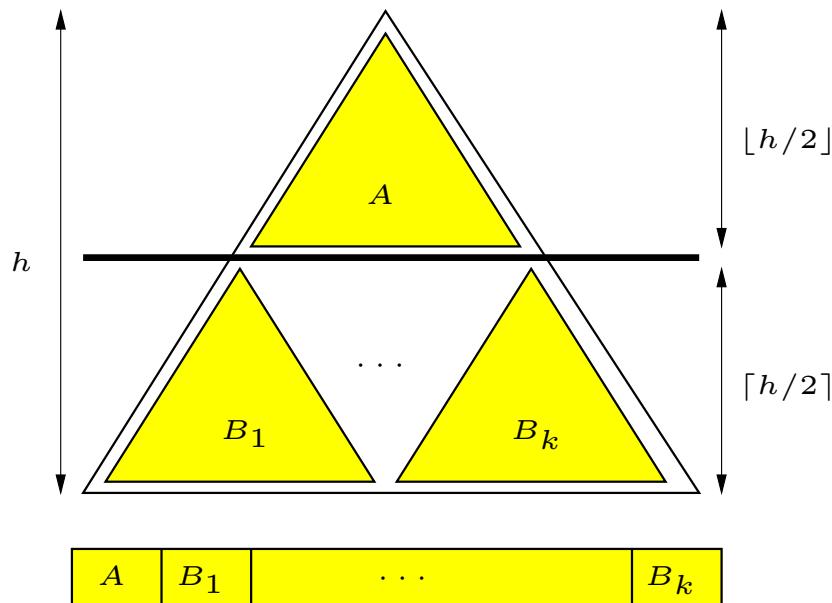
Overview

- ✓ The memory hierarchy
- ✓ The I/O-model
- ✓ The cache-oblivious model
- Examples of cache-oblivious algorithms
 - ✓ Double for-loop (with applications)
 - Searching
 - Sorting
- Theoretical limits of cache-obliviousness

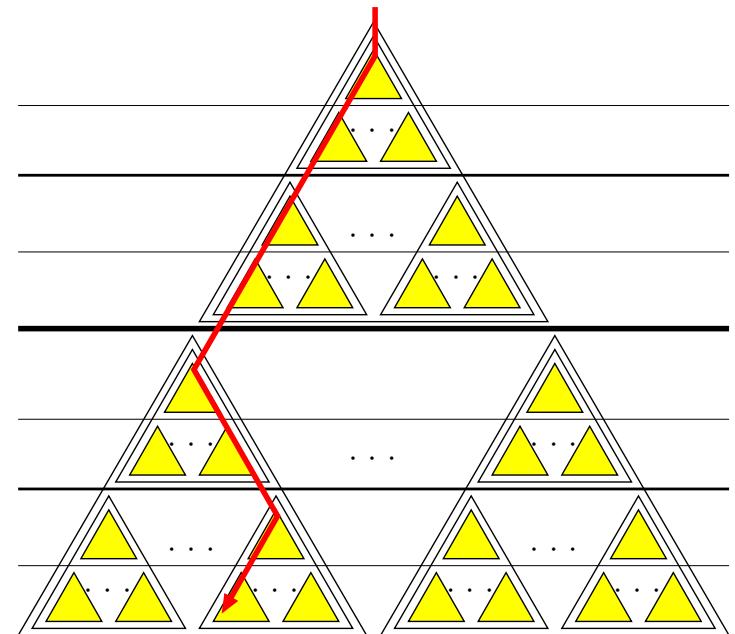
Static Cache-Oblivious Trees

Recursive memory layout (van Emde Boas layout)

Prokop 1999



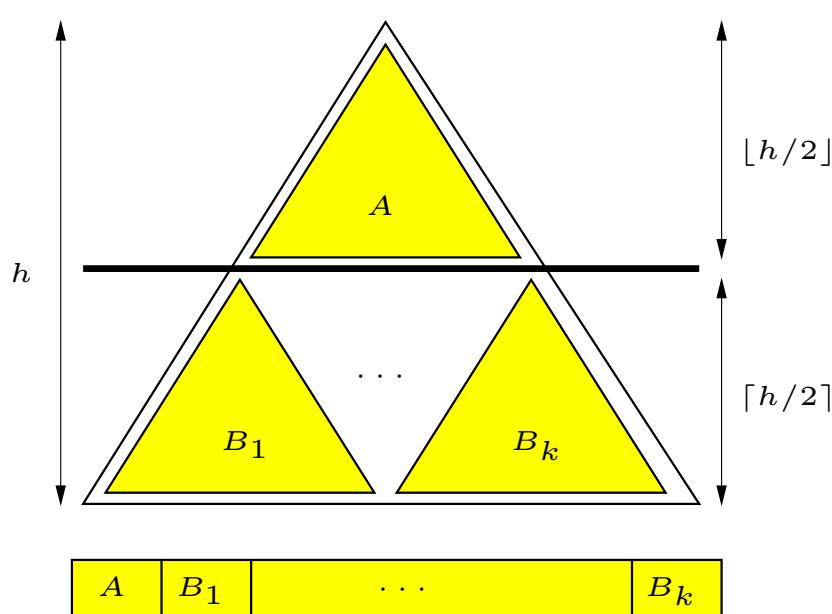
Binary tree



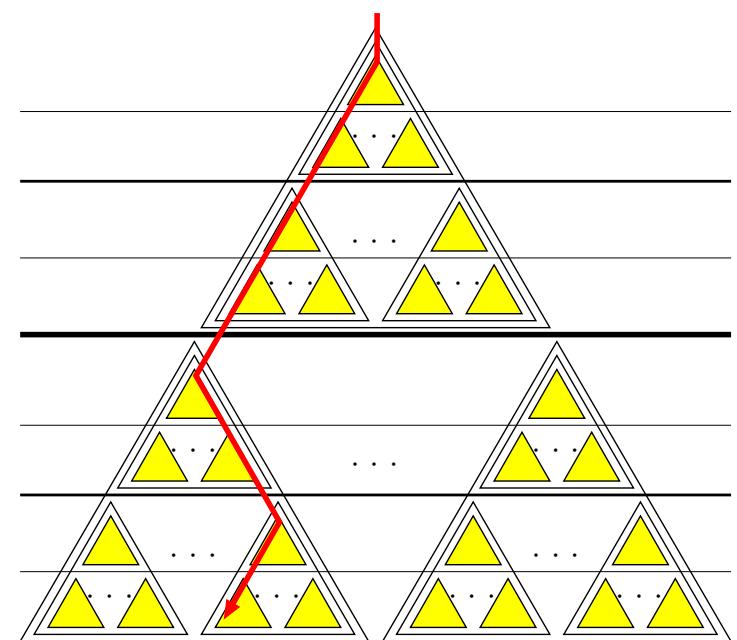
Searches use $O(\log_B N)$ I/Os

Static Cache-Oblivious Trees

Recursive memory layout (van Emde Boas layout)



Binary tree

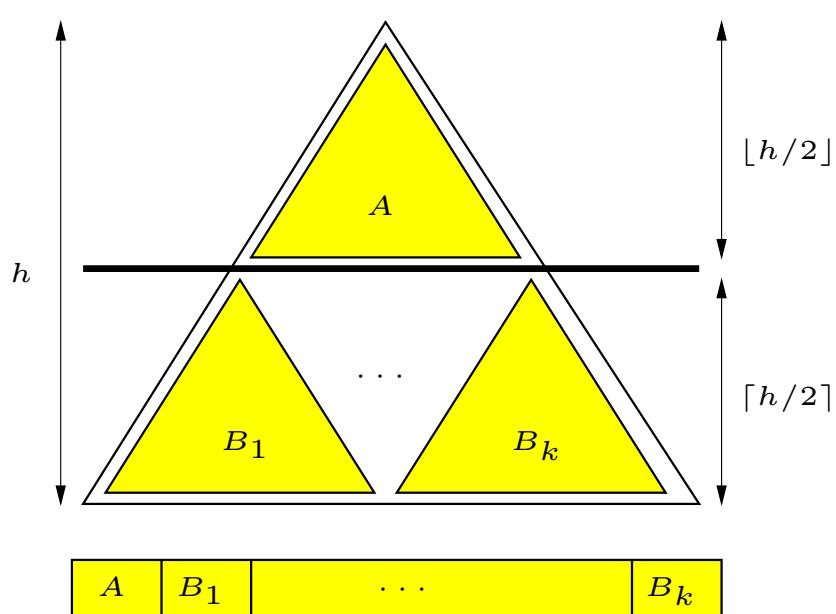


Searches use $O(\log_B N)$ I/Os

Dynamization?

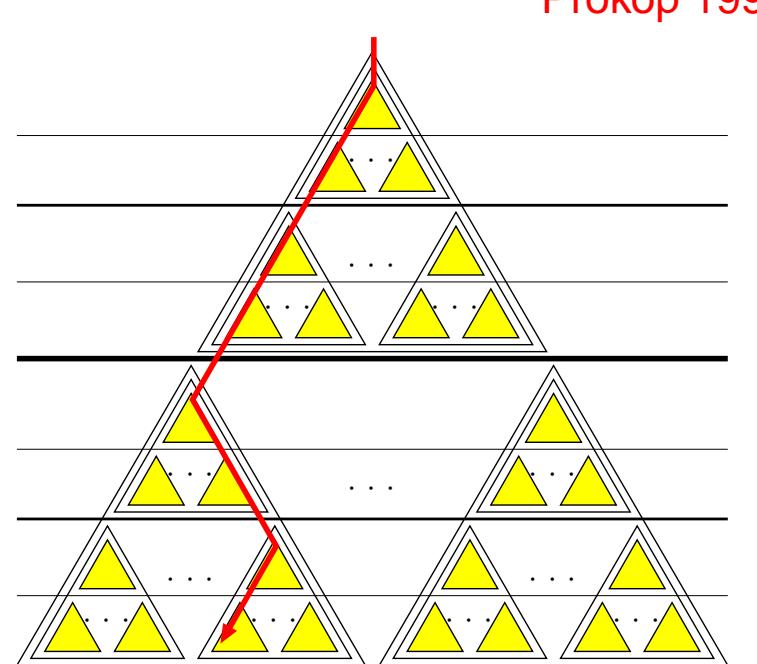
Static Cache-Oblivious Trees

Recursive memory layout (van Emde Boas layout)



Binary tree

Dynamization?



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

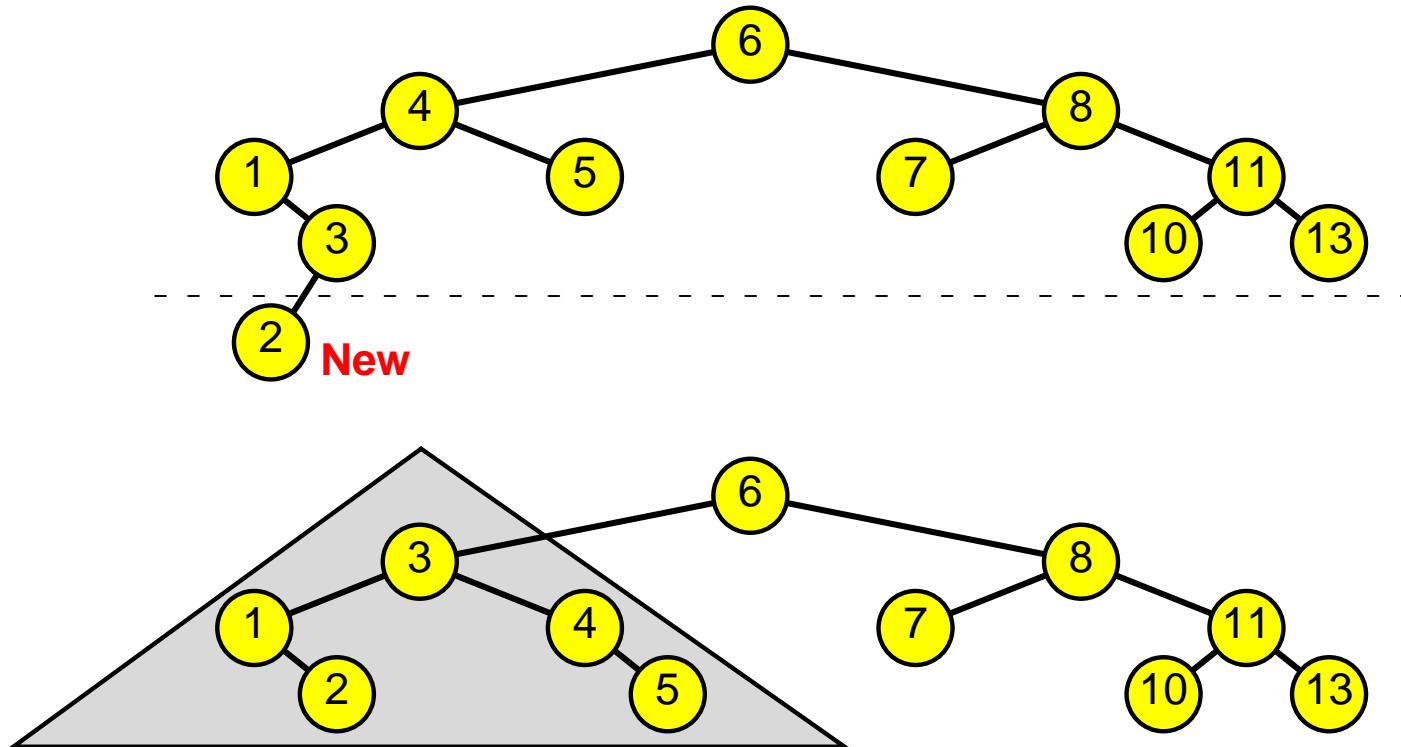
Bender, Demaine, Farach-Colton, FOCS'00

Rahman, Cole, Raman, WAE'01

Bender, Duan, Iacono, Wu, SODA 02

Brodal, Fagerberg, Jacob, SODA'02

Binary Trees of Height $\log_2(n) + O(1)$



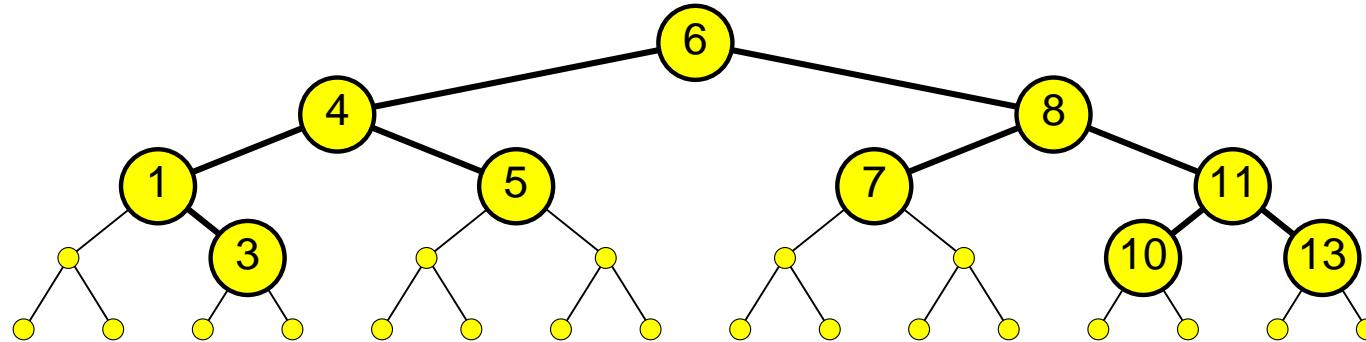
- If an insertion causes non-small height then **rebuild** subtree at nearest ancestor with sufficient few descendants
- Insertions require amortized time $O(\log^2 N)$

Itai, Konheim, Rodeh, 1981

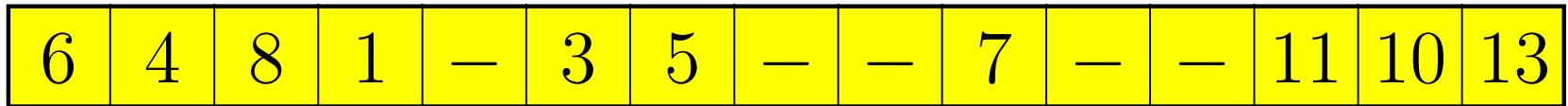
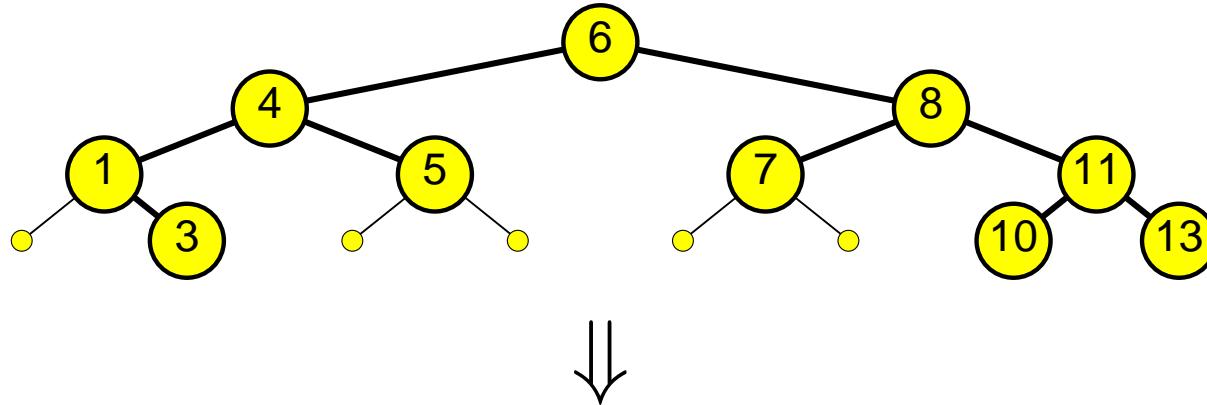
Andersson, Lai, 1990

Simple Dynamic Cache-Oblivious Trees

- Embed a dynamic tree of height $\log_2(n) + O(1)$ into a complete tree
- Static van Emde Boas layout of the complete tree in array.



Example



Search

$O(\log_B N)$

Range Reporting

$O\left(\log_B N + \frac{k}{B}\right)$

Updates

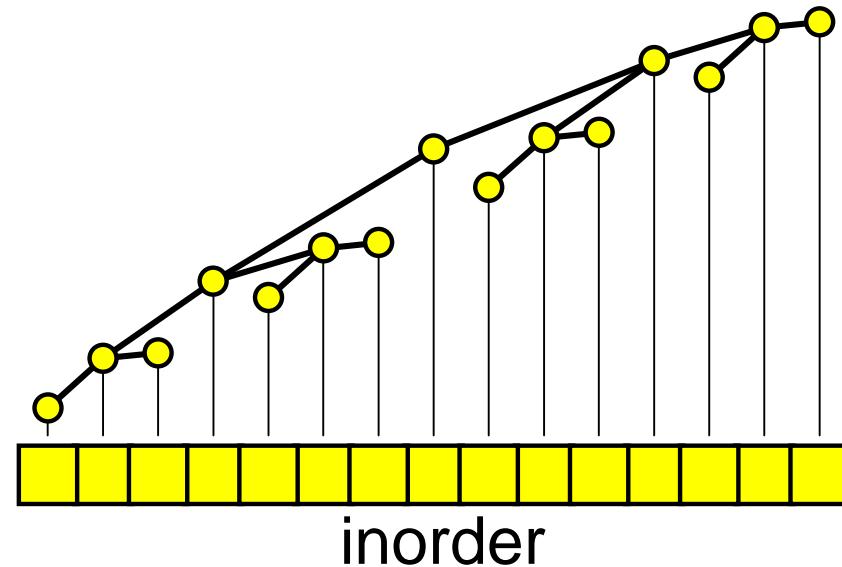
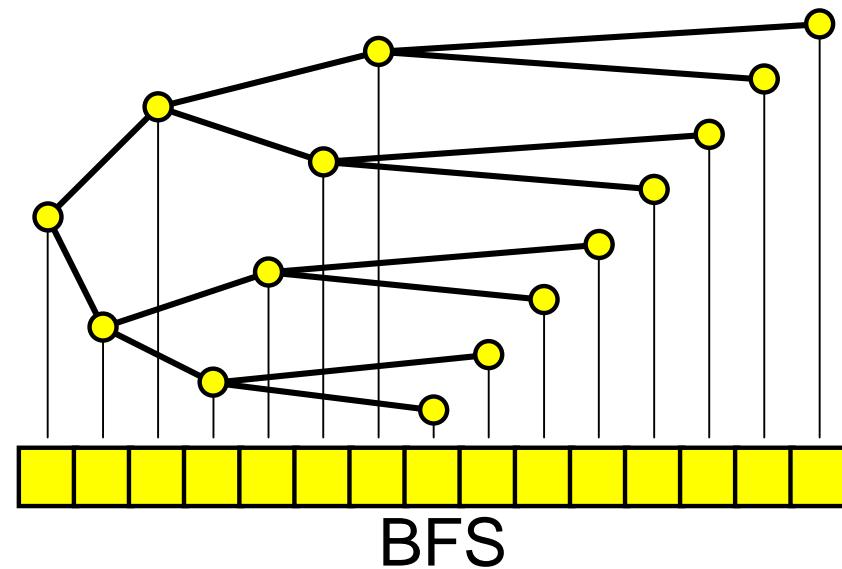
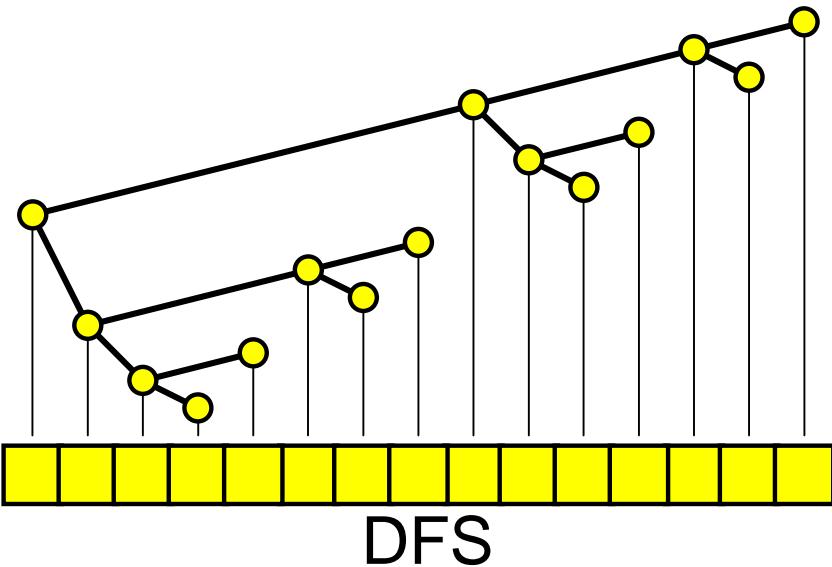
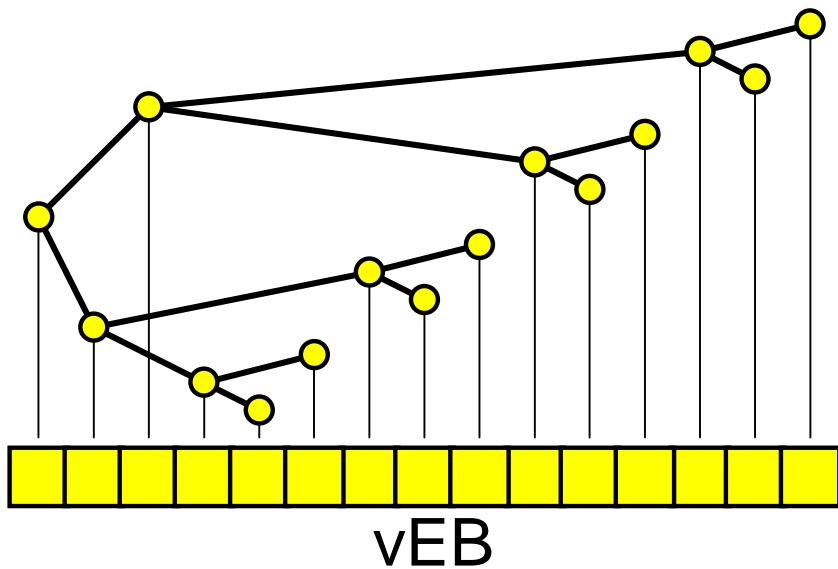
$O\left(\log_B N + \frac{\log^2 N}{B}\right)$

Experiments

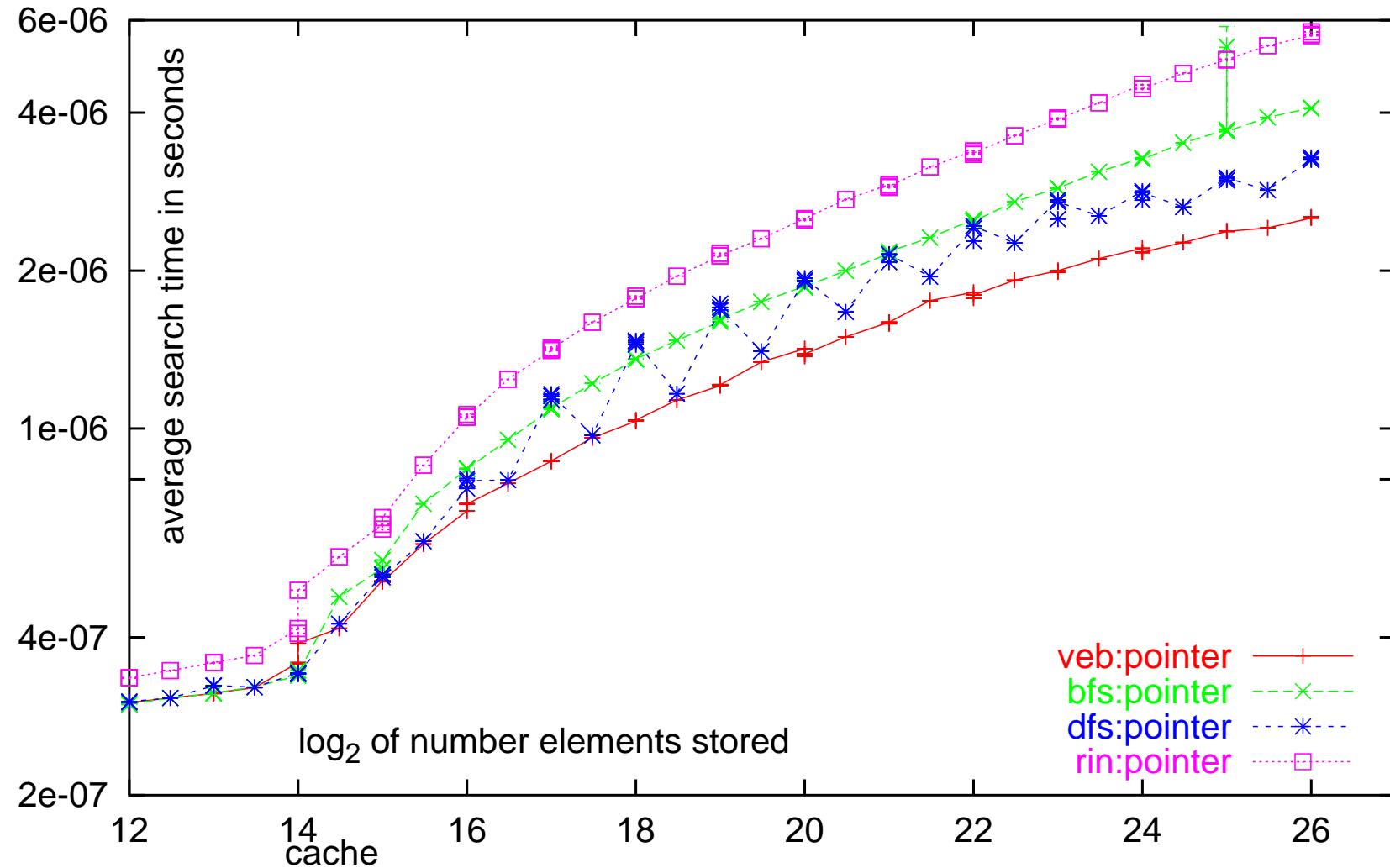
Brodal, Fagerberg, Jacob, SODA'02

1. Study search time in static tree layouts.
 - Classic layouts: BFS, DFS, inorder, randomly built trees
 - Cache-aware multi-way trees
 - Cache-oblivious vEB layout
2. Study pointer-based vs. implicit representation.

Different Memory Layouts



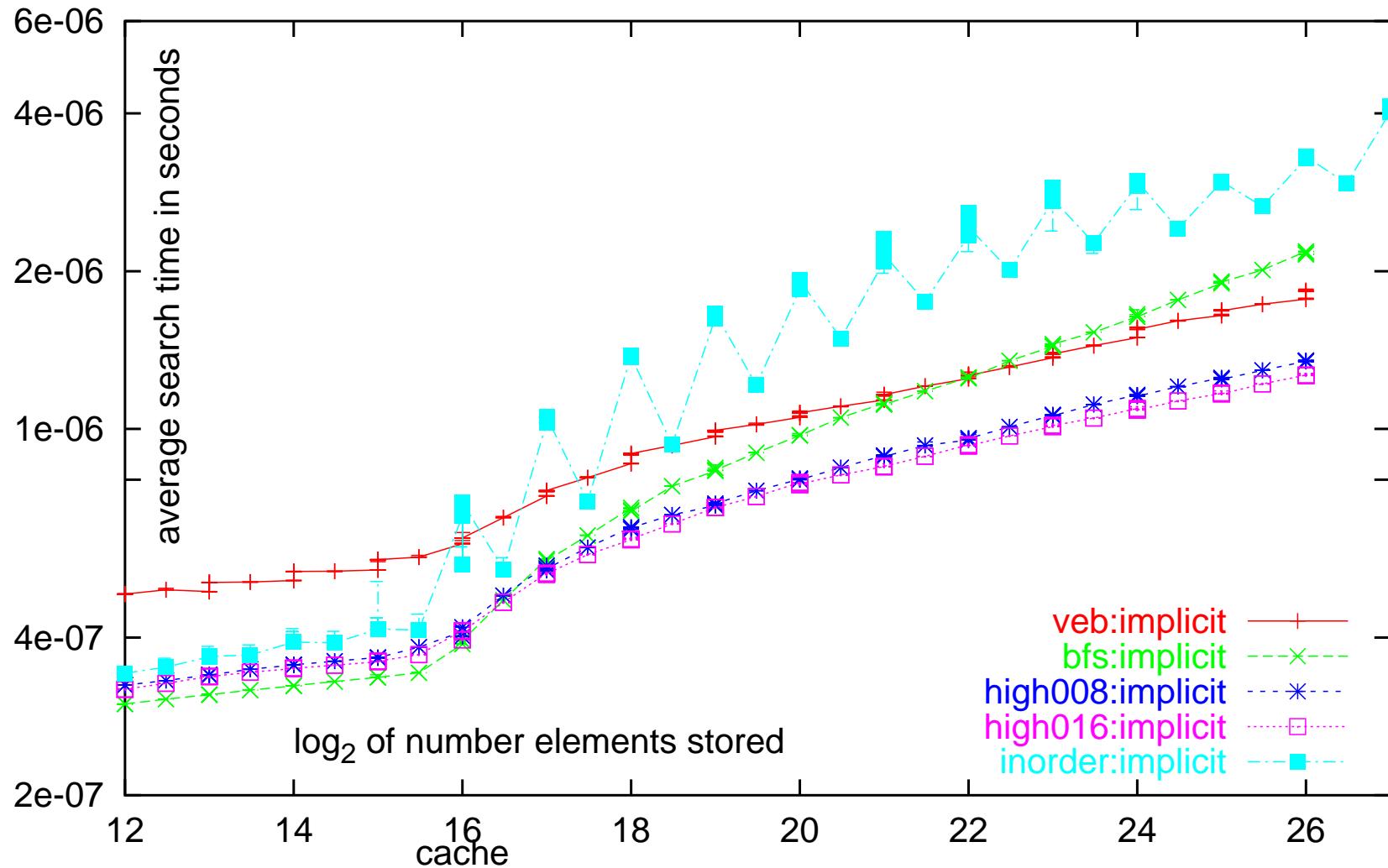
Search in Pointer Based Tree



Pointer based vEB, BFS, DFS, and randomly built.

1 GHz Pentium III, 1GB RAM, 256 KB Cache, gcc -O3 / linux

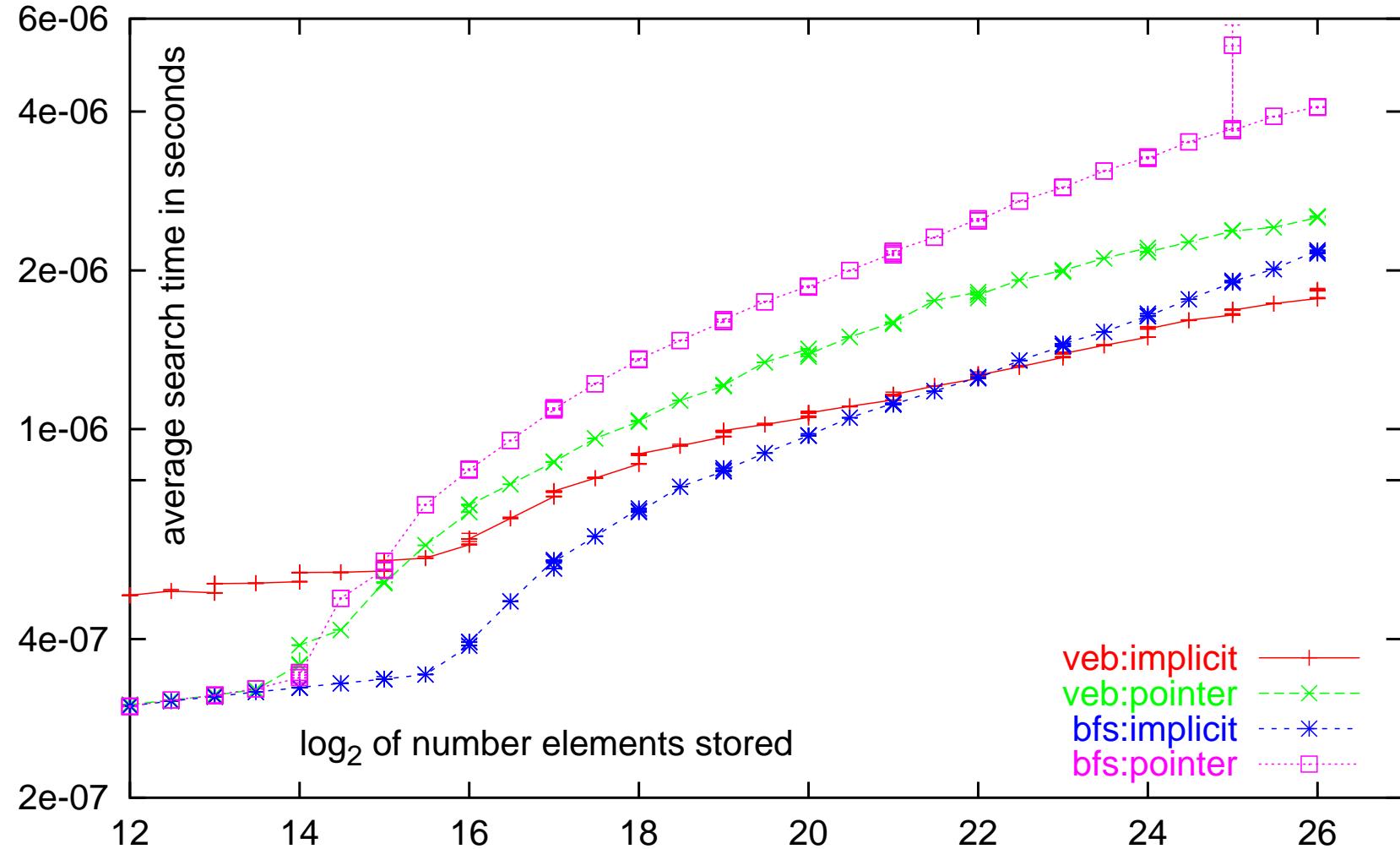
Search in Implicit Tree



Implicit vEB, BFS, inorder, multi-way.

1 GHz Pentium III, 1GB RAM, 256 KB Cache, gcc -O3 / linux

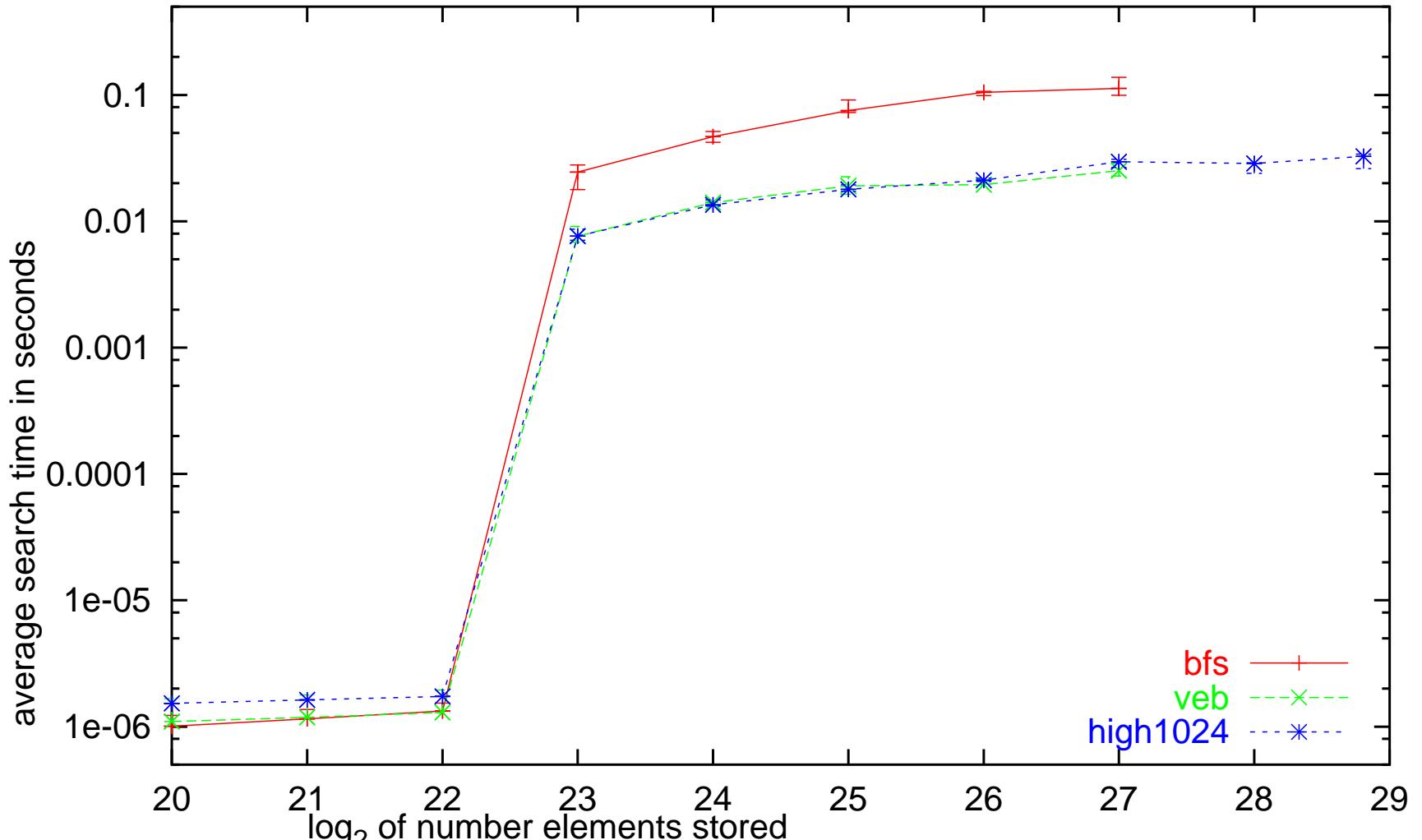
Pointer versus Implicit



static van Emde Boas layout and bfs-layout

1 GHz Pentium III, 1GB RAM, 256 KB Cache, gcc -O3 / linux

Beyond Main Memory



Multiway-tree, BFS, vEB implicit layout

1 GHz Pentium III, 32MB RAM, 256 KB Cache, gcc -O3 / linux

Overview

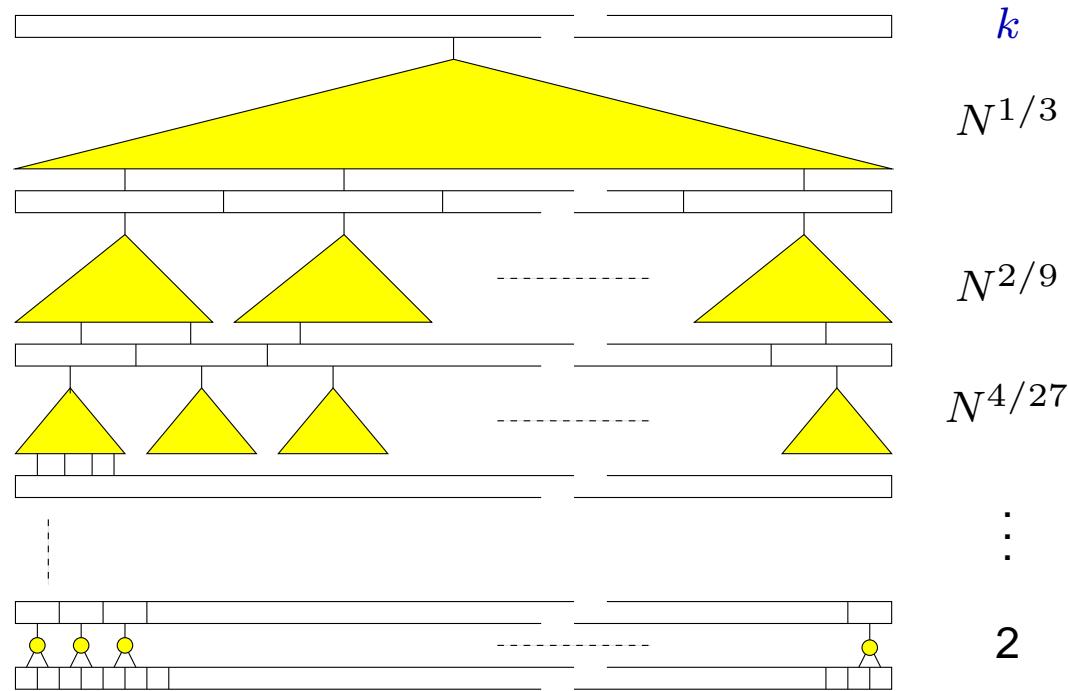
- ✓ The memory hierarchy
- ✓ The I/O-model
- ✓ The cache-oblivious model
- Examples of cache-oblivious algorithms
 - ✓ Double for-loop (with applications)
 - ✓ Searching
 - Sorting
- Theoretical limits of cache-obliviousness

Funnelsort

Divide input in $N^{1/3}$ segments of size $N^{2/3}$

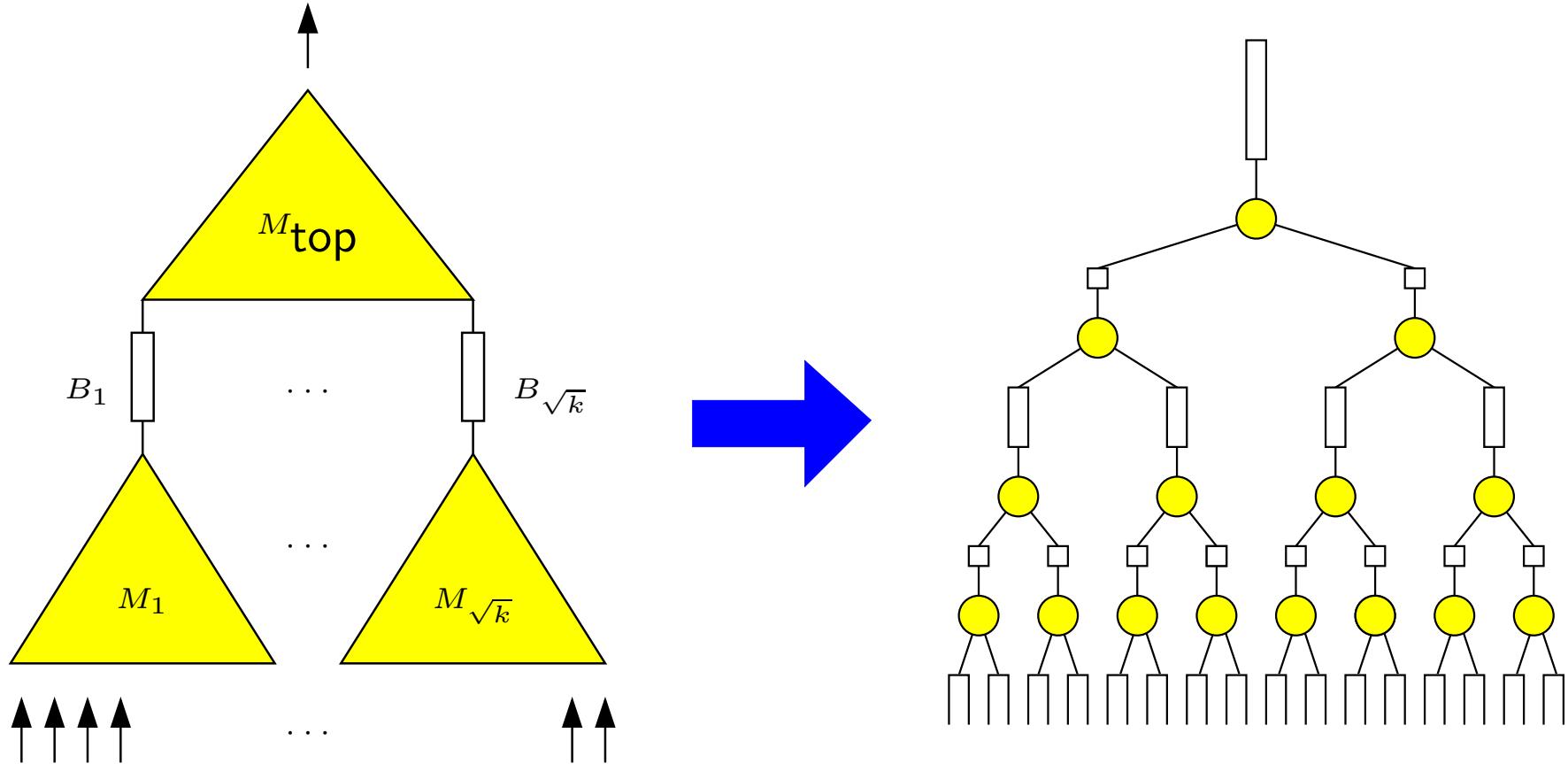
Recursively **Funnelsort** each segment

Merge sorted segments by an $N^{1/3}$ -merger



Frigo, Leiserson, Prokop, Ramachandran, 1999

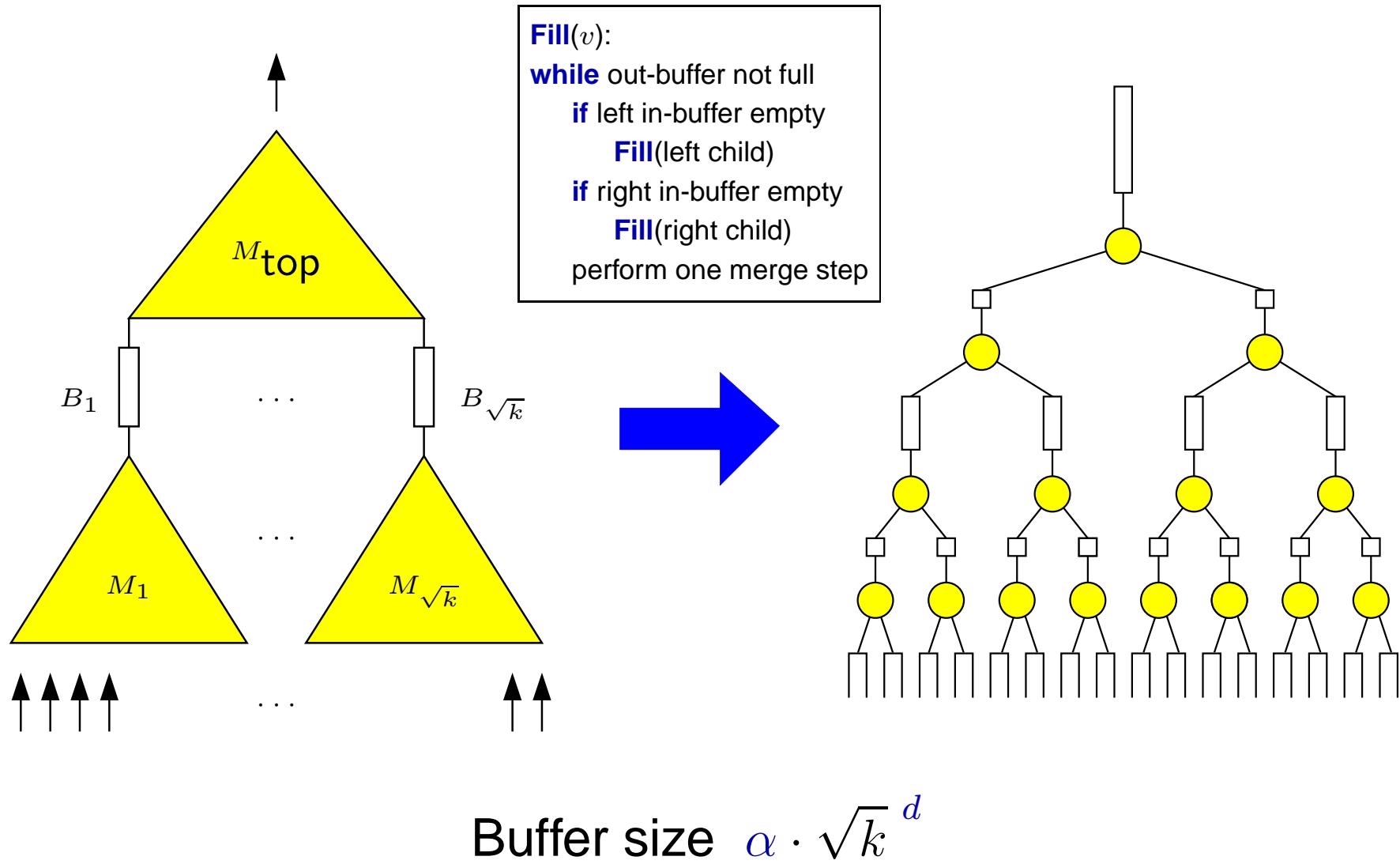
k -merger



Buffer size $\alpha \cdot \sqrt{k}^d$

Brodal, Fagerberg 2002

k -merger



Brodal, Fagerberg 2002

Experiments

Engineer: Try a large number of possibilities for parameter choices, design choices, code optimizations (human vs. compiler+CPU), memory layouts,...

Choose a hard pack of competitors:

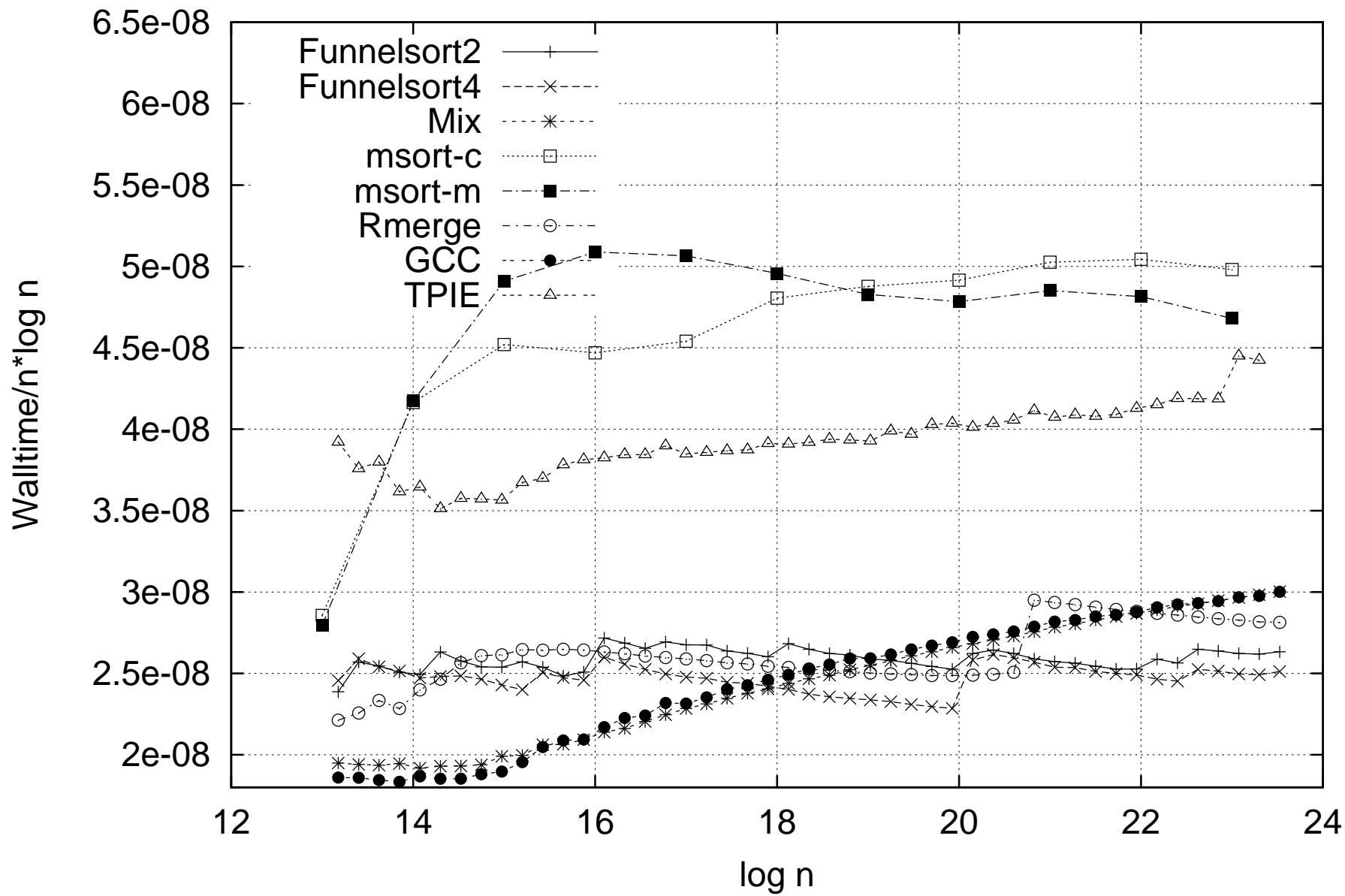
- 2-way and 4-way Funnelsort
- Best library Quicksort we can find
- Recent cache-aware proposals (tuned for RAM or for disk)

Run on a number of different machines: Pentium 4, Pentium III, MIPS 10000, AMD Athlon, Itanium 2.

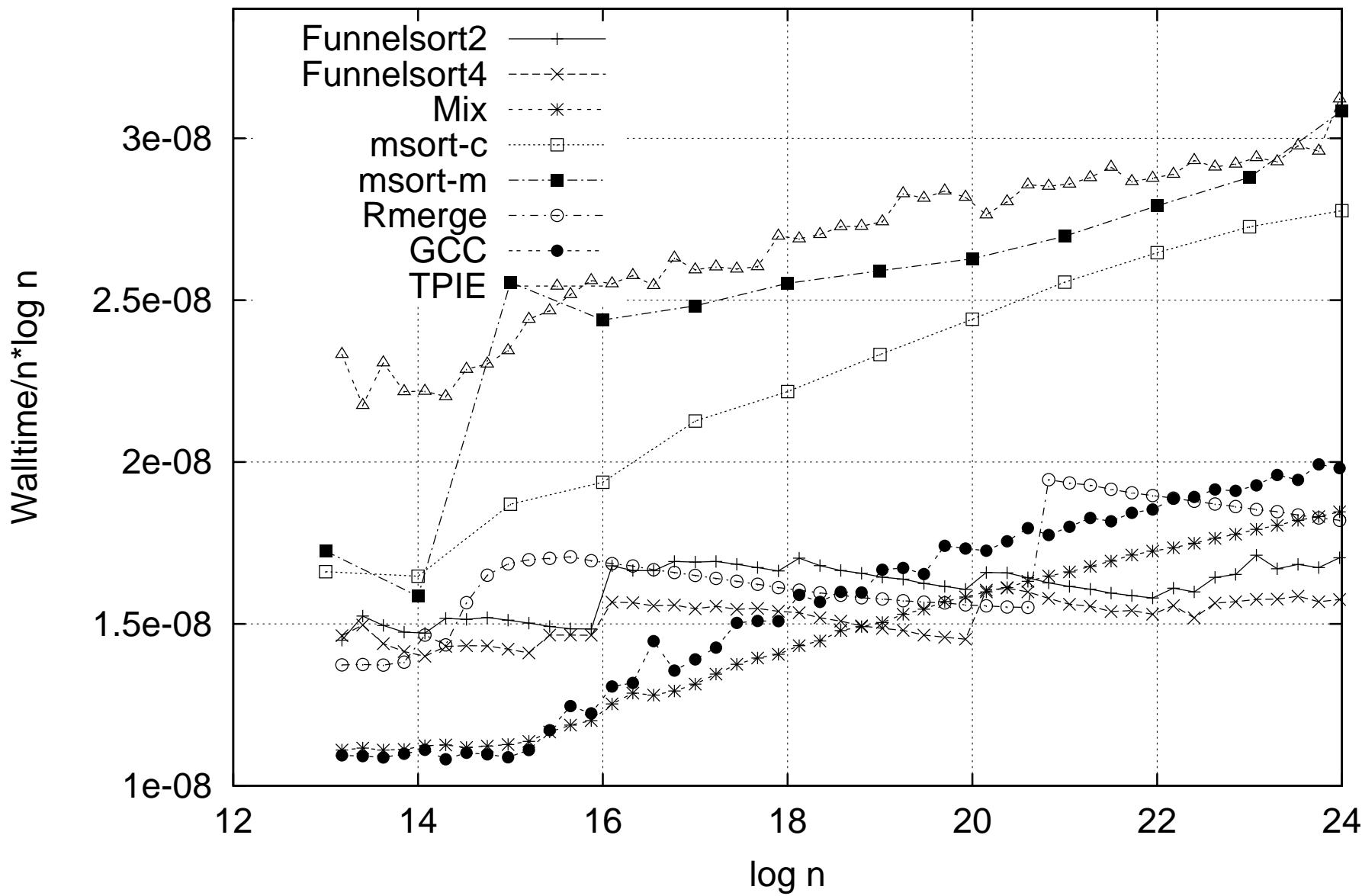
Brodal, Fagerberg, Vinther, ALENEX'04

Results for Inputs in RAM

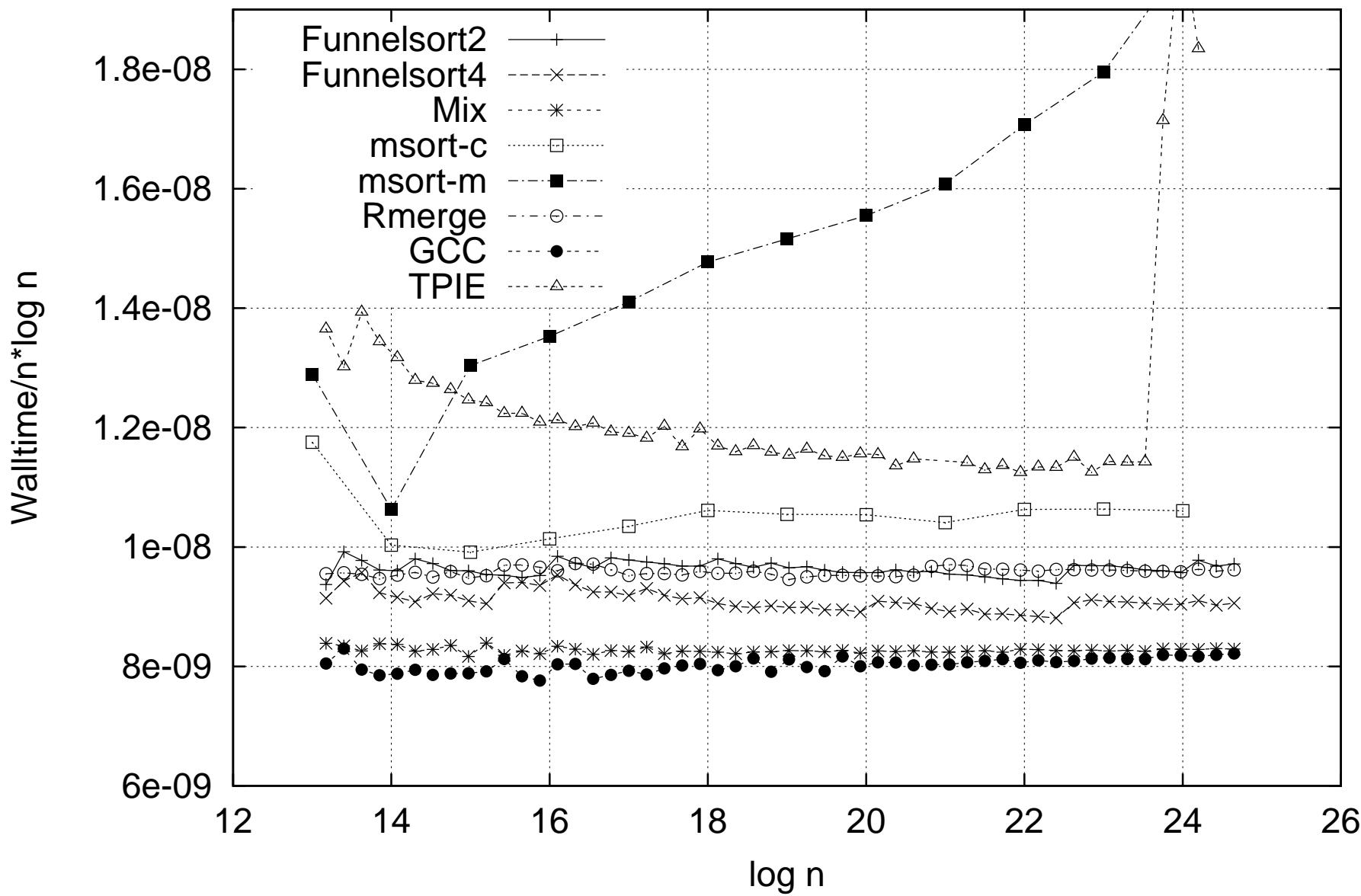
Uniform pairs - Pentium III



Uniform pairs - AMD Athlon

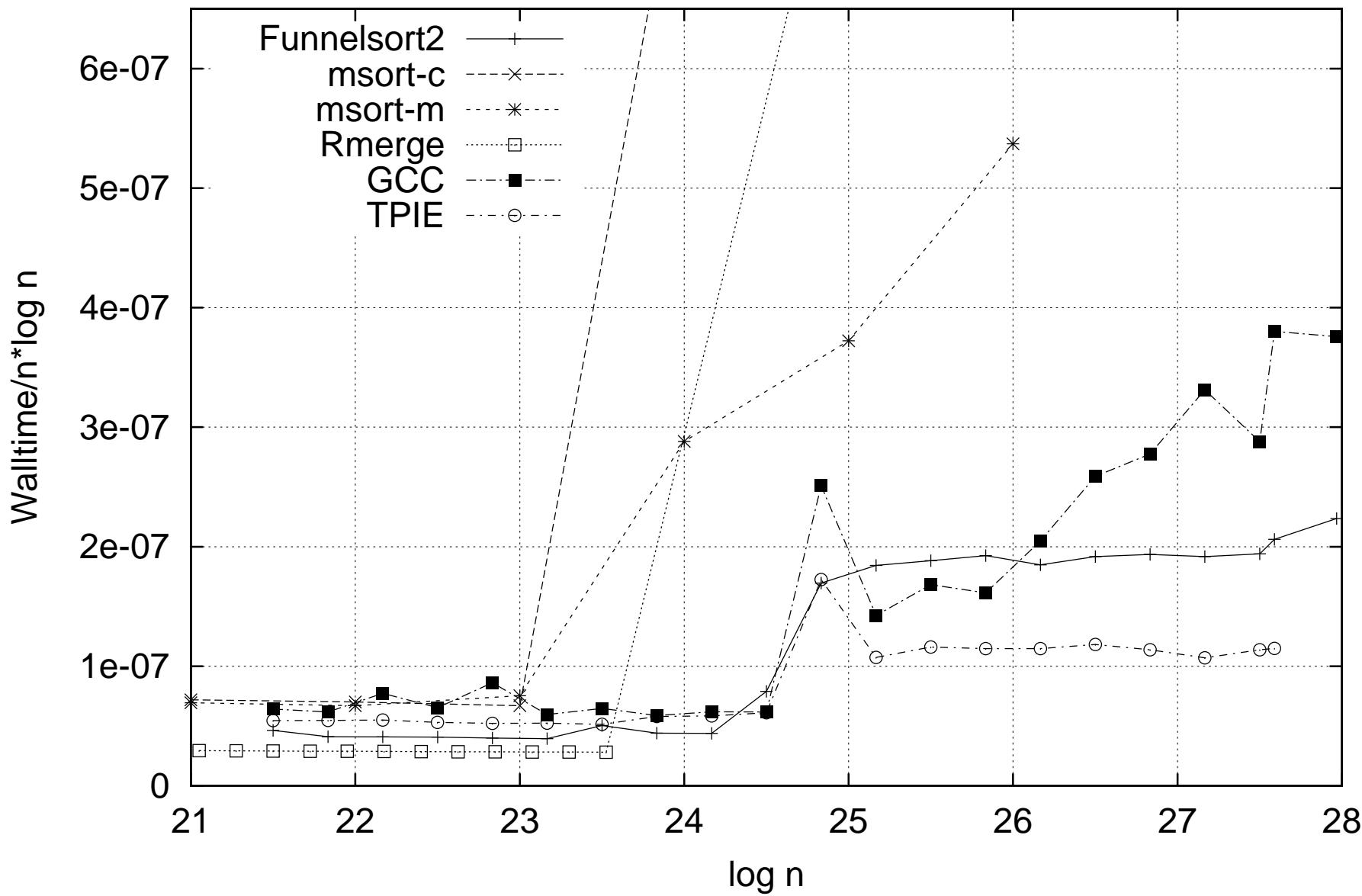


Uniform pairs - Pentium 4

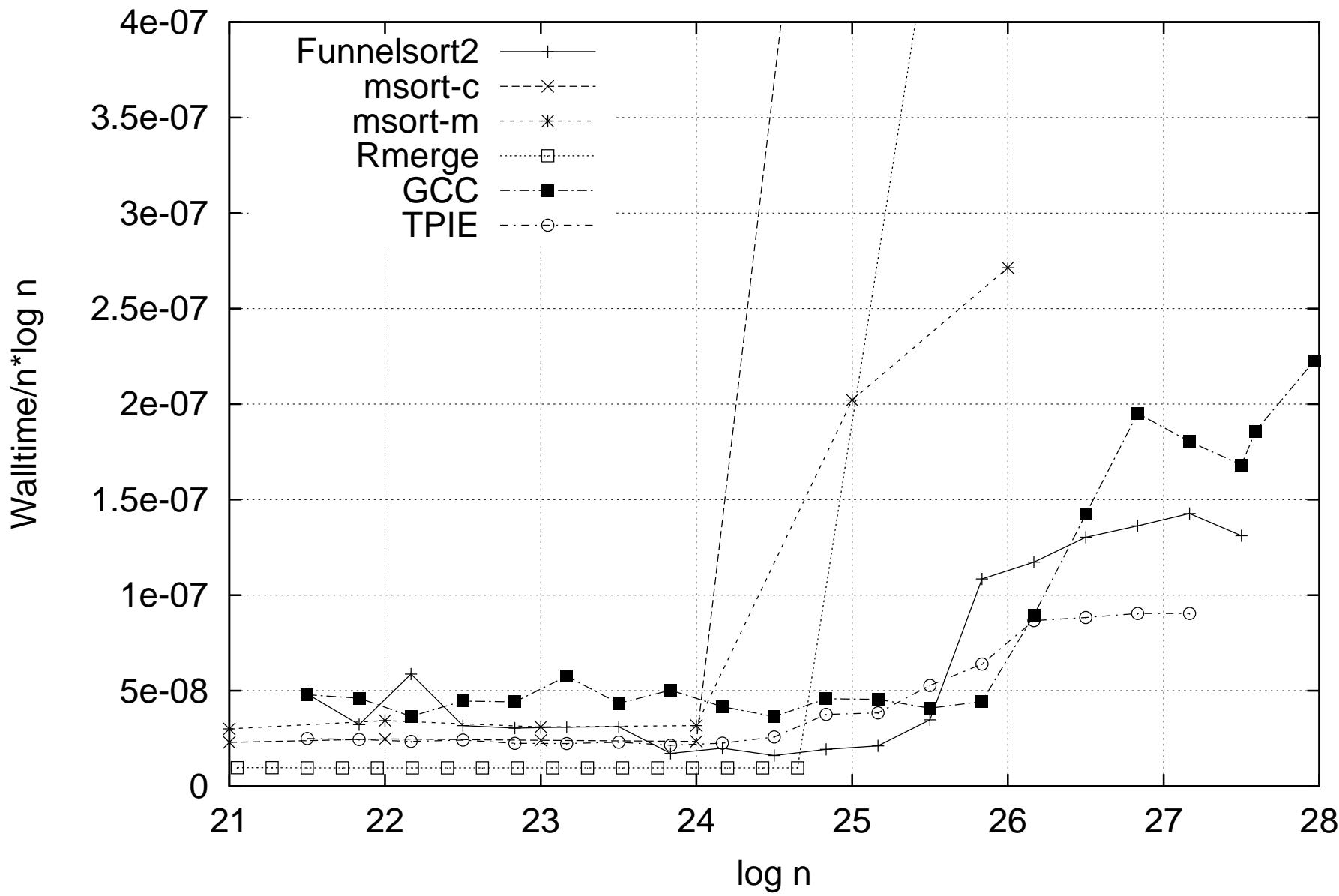


Results for Inputs on Disk

Uniform pairs - Pentium III



Uniform pairs - Pentium 4



Practical Conclusion

For a number of basic algorithmic problems there exist solutions which

- Are theoretically I/O-efficient
- Are simple
- Are robust - adapt automatically to the specifics of the memory hierarchy
- Compete well with explicit (cache-aware) I/O algorithms.

Practical Conclusion

For a number of basic algorithmic problems there exist solutions which

- Are theoretically I/O-efficient
- Are simple
- Are robust - adapt automatically to the specifics of the memory hierarchy
- Compete well with explicit (cache-aware) I/O algorithms.

Consider putting cache-obliviousness in
your toolbox of algorithmic techniques.

Overview

- ✓ The memory hierarchy
- ✓ The I/O-model
- ✓ The cache-oblivious model
- Examples of cache-oblivious algorithms
 - ✓ Double for-loop (with applications)
 - ✓ Searching
 - ✓ Sorting
- Theoretical limits of cache-obliviousness

Basic Question

In terms of power:

Cache-oblivious algorithms = I/O-algorithms?

Basic Question

In terms of power:

Cache-oblivious algorithms = I/O-algorithms?

Or is there a separation: A problem for which no cache-oblivious algorithm can match the best I/O-algorithm?

Frigo, Leiserson, Prokop, Ramachandran, FOCS'99

Recent Answers

Brodal, Fagerberg, STOC'03

Bender, Brodal, Fagerberg, Ge, He, Hu, Iacono, López-Ortiz, FOCS'03

Sorting:

$$\text{Sort}_{B,M}(N) \text{ I/Os} \quad ?$$

Permuting:

$$\text{Perm}_{B,M}(N) \text{ I/Os} \quad ?$$

Searching:

$$c \cdot \log_B \frac{N}{M} \text{ I/Os} \quad ?$$

Recent Answers

Brodal, Fagerberg, STOC'03

Bender, Brodal, Fagerberg, Ge, He, Hu, Iacono, López-Ortiz, FOCS'03

Sorting:

Sort_{B,M}(N) I/Os ?

Not possible
without
tall cache
assumption

Permuting:

Perm_{B,M}(N) I/Os ?

Not possible
even with
tall cache
assumption

Searching:

c · log_B $\frac{N}{M}$ I/Os ?

$c \geq \log(e) \approx 1.443$

Recent Answers

Brodal, Fagerberg, STOC'03

Bender, Brodal, Fagerberg, Ge, He, Hu, Iacono, López-Ortiz, FOCS'03

Sorting:

$$\text{Sort}_{B,M}(N) \text{ I/Os} \quad ?$$

Not possible
without
tall cache
assumption

Permuting:

$$\text{Perm}_{B,M}(N) \text{ I/Os} \quad ?$$

Not possible
even with
tall cache
assumption

Searching:

$$c \cdot \log_B \frac{N}{M} \text{ I/Os} \quad ?$$

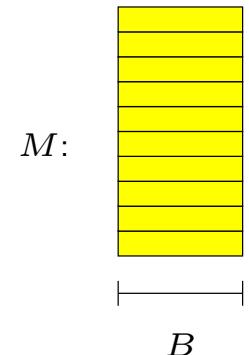
$$c \geq \log(e) \approx 1.443$$

Matching upper bounds for sorting and searching

Cache-Oblivious Sorting

Previous cache-oblivious sorting results use a

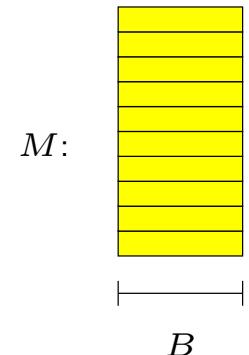
Tall Cache Assumption $B \leq M^{1/2}$



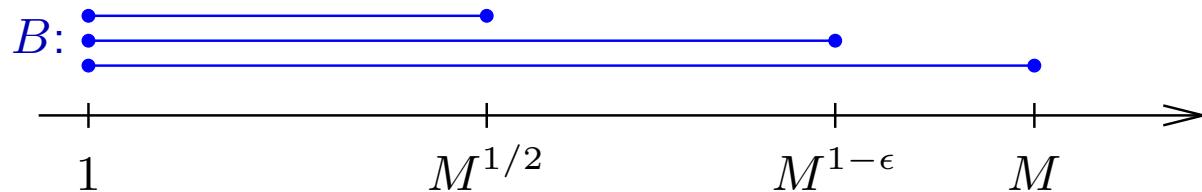
Cache-Oblivious Sorting

Previous cache-oblivious sorting results use a

Tall Cache Assumption $B \leq M^{1/2}$

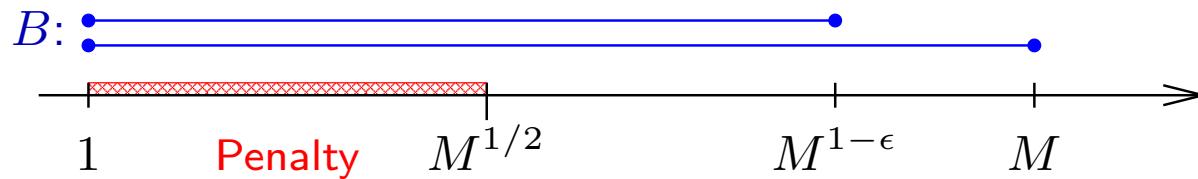


Algorithm	Assumption	I/O bound
Funnelsort [FLPR99]	$B \leq M^{1/2}$	$\text{Sort}_{B,M}(N)$
Lazy Funnelsort [BF02]	$B \leq M^{1-\epsilon}$	$\frac{N}{\epsilon B} \log_M \frac{N}{M}$
Binary Mergesort	$(B \leq M)$	$\frac{N}{B} \log_2 \frac{N}{M}$



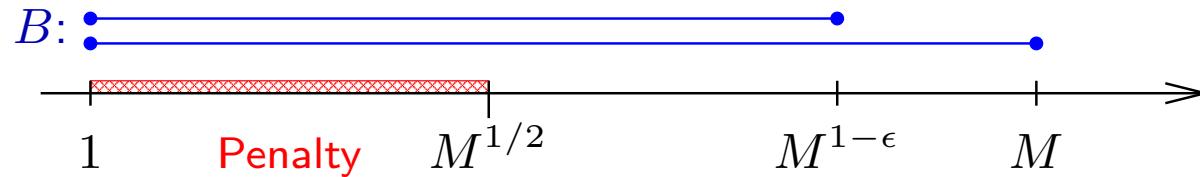
Result for Sorting

	<i>Assumption</i>	<i>I/Os</i>
Lazy Funnelsort	$B \leq M^{1-\epsilon}$	$B = M^{1-\epsilon} : \text{Sort}_{B,M}(N)$ $B \leq M^{1/2} : \text{Sort}_{B,M}(N) \cdot \boxed{\frac{1}{\epsilon}}$
Binary Mergesort	$(B \leq M)$	$B = M : \text{Sort}_{B,M}(N)$ $B \leq M^{1/2} : \text{Sort}_{B,M}(N) \cdot \boxed{\log M}$



Result for Sorting

	<i>Assumption</i>	<i>I/Os</i>
Lazy Funnelsort	$B \leq M^{1-\epsilon}$	$(a) B = M^{1-\epsilon} : \text{Sort}_{B,M}(N)$ $(b) B \leq M^{1/2} : \text{Sort}_{B,M}(N) \cdot \boxed{\frac{1}{\epsilon}}$
Binary Mergesort	$(B \leq M)$	$(a) B = M : \text{Sort}_{B,M}(N)$ $(b) B \leq M^{1/2} : \text{Sort}_{B,M}(N) \cdot \boxed{\log M}$



Theorem This is tight. For any cache-oblivious comparison based sorting algorithm:

$$(a) \Rightarrow (b)$$

The end

Proof

One algorithm, two machines ($B_1 \leq B_2$):

	Block Size	Memory	I/Os
Machine 1	B_1	M	t_1
Machine 2	B_2	M	t_2

Main result:

$$8t_1B_1 + 3t_1B_1 \log \frac{8Mt_2}{t_1B_1} \geq N \log \frac{N}{M} - 1.45N \quad (*)$$

Theorem 1:

Inserting (a) in (*) leads to (b)

Fake Proof

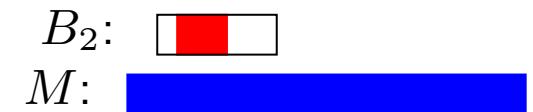
Goal:

$$8t_1B_1 + 3t_1B_1 \log \frac{8Mt_2}{t_1B_1} \geq N \log \frac{N}{M} - 1.45N \quad (*)$$

Merging sorted lists X and Y takes $|X| \log \frac{|Y|}{|X|}$ comparisons.

In total t_1B_1 elements touched $\Rightarrow t_1B_1/t_2$ elements touched on average per B_2 -I/O \Rightarrow effective B_2 is t_1B_1/t_2 .

Comparisons gained per B_2 -I/O:



$$t_1B_1/t_2 \cdot \log \frac{M}{t_1B_1/t_2} .$$

Hence:

$$t_1B_1 \cdot \log \frac{Mt_2}{t_1B_1} \geq N \log N - 1.45N .$$

□

Fake Proof

Goal:

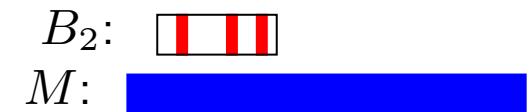
$$8t_1B_1 + 3t_1B_1 \log \frac{8Mt_2}{t_1B_1} \geq N \log \frac{N}{M} - 1.45N \quad (*)$$

Merging sorted lists X and Y takes $|X| \log \frac{|Y|}{|X|}$ comparisons.

In total t_1B_1 elements touched $\Rightarrow t_1B_1/t_2$ elements touched on average per B_2 -I/O \Rightarrow effective B_2 is t_1B_1/t_2 .

Comparisons gained per B_2 -I/O:

$$t_1B_1/t_2 \cdot \log \frac{M}{t_1B_1/t_2} .$$



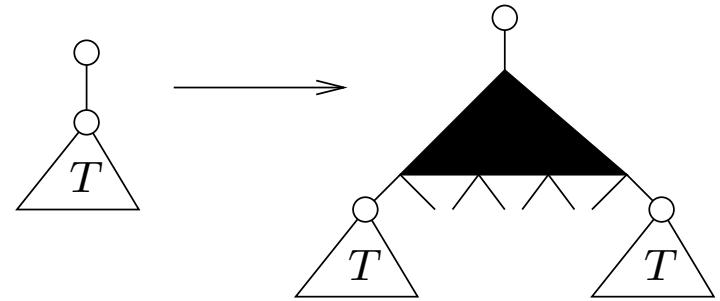
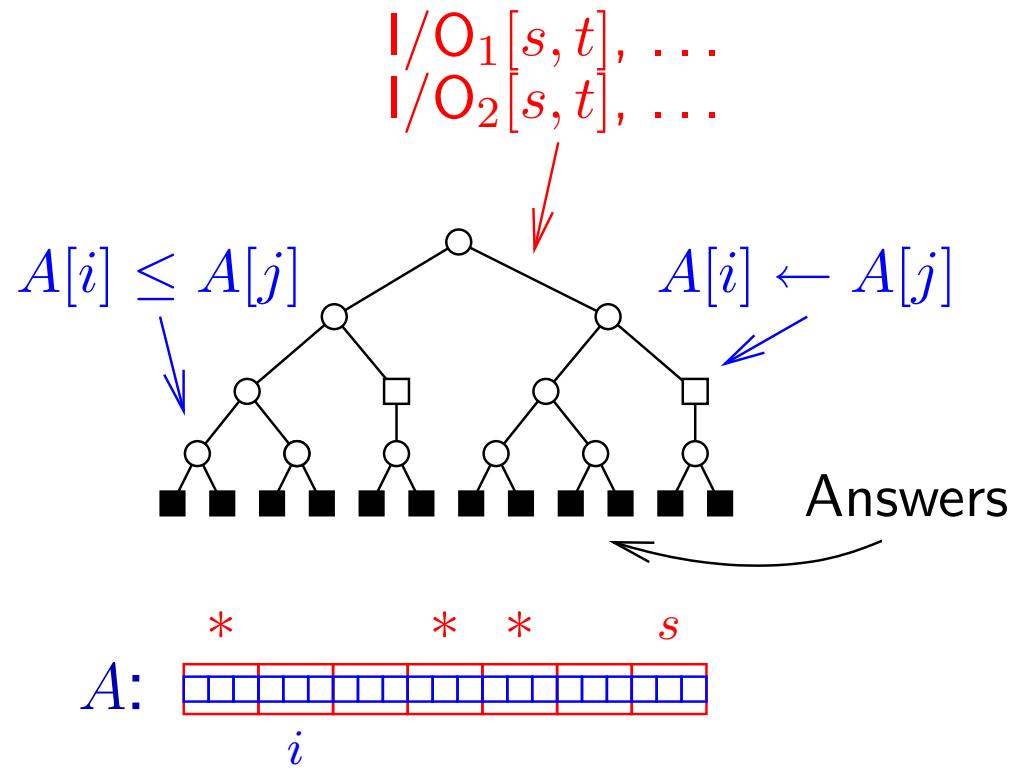
One problem:
Online choice

Hence:

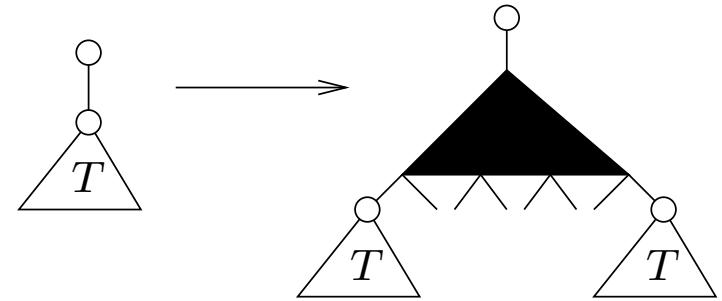
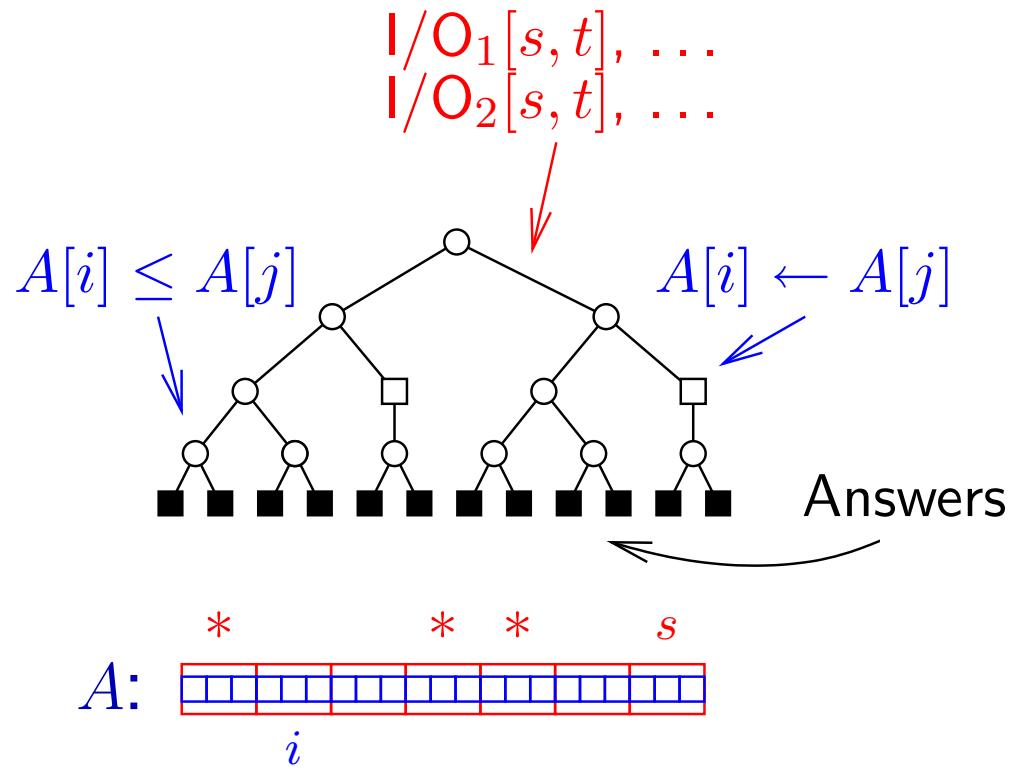
$$t_1B_1 \cdot \log \frac{Mt_2}{t_1B_1} \geq N \log N - 1.45N .$$

□

Ideas from Real Proof



Ideas from Real Proof



$$8t_1B_1 + 3t_1B_1 \log \frac{8Mt_2}{B_1t_1} \geq \text{height} \geq N \log \frac{N}{M} - 1.45N . \quad (*)$$

The end