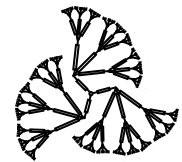# Engineering Cache-oblivious Sorting Algorithms

Master's thesis by

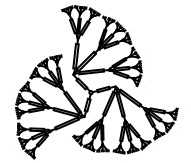Kristoffer Vinther

# Sorting Algorithms

- The Problem

  Given a set of elements, put them in non-decreasing order.

- Motivation

  Very commonly used as a subroutine in other algorithms (such as graph-, geometric-, and scientific algorithms).

  A good sorting implementation is thus important to achieving good implementations of many other algorithms.
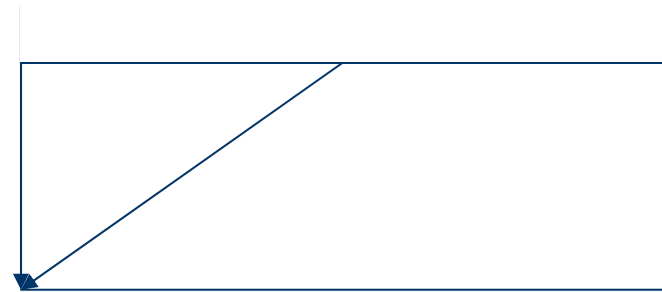
  Performance of sorting algorithms seem greatly influenced by many aspects of modern computers, such as the memory hierarchy and pipelined execution.
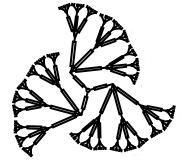
# Sorting Algorithms – Binary mergesort

- Ex. binary mergesort:
  1. Split elements into two halves.
  2. Sort each half recursively.
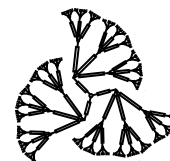  3. Make space for sorted elements and *merge* the sorted halves

# Engineering Cache-oblivious Sorting Algorithms

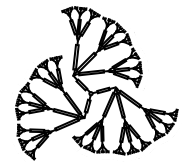Master's thesis by

Kristoffer Vinther

# Cache-oblivious – Motivation

- The presence of a memory hierarchy has become a fact of life.

- Accessing non-local storage may take a very long time.

- Good locality is important to achieving high performance.

|  | 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 | $10^7$ |

# Cache-oblivious Algorithms – Models
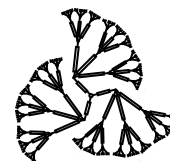
### Random Access Memory
- All basic operations take constant time.
- Complexity is the number of operations executed (instruction count), i.e. the total running time of the algorithm.

### External Memory
- Computation is done in *main memory*.
- Data is brought to and from main memory in *I/Os*, explicitly controlled by the algorithm
- Complexity is the number of I/Os done by the algorithm.

### Cache-oblivious
- Algorithms designed for the RAM model; algorithm does not control the I/Os.
- Algorithms analyzed for the EM model.
- Complexity is *both* the instruction count and the number of I/Os (memory transfers) incurred by the algorithm.
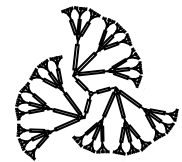
# Cache-oblivious Algorithms – Sorting

## Random Access Memory

- Complexity of binary mergesort: $\mathcal{O}(N \log N)$.

- Complexity of any (comparison-based) sorting algorithm: $\Omega(N \log N)$.

## External Memory

- Complexity of binary mergesort: $\mathcal{O}(N/B \log N/M)$.

- Complexity of any sorting algorithm: $\Omega(N/B \log_{M/B} N/M)$.

---

- Binary mergesort optimal in External Memory only if $M = 2B$.

- What if $M > 2B$? *Multiway mergesort* incurs $\mathcal{O}(N/B \log_{M/B} N/M)$ I/Os, given the right $M$ and $B$.

- Multiway mergesort is suboptimal with the wrong $M$ and $B$.

  - $M$ and $B$ cannot in general be determined.

  - Running the algorithm on a machine different from the one to which it was designed.

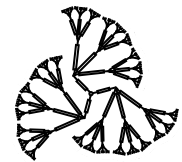- Funnelsort and LOWSCOSA incurs $\mathcal{O}(N/B \log_{M/B} N/M)$ memory transfers, without knowing $M$ and $B$.

# Cache-oblivious Algorithms – Assumptions

- To analyze the cache complexity of an algorithm that is oblivious to caches, some issues need to be settled:

  - How is an I/O initiated?

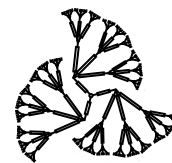  - Where in memory should the block be placed?
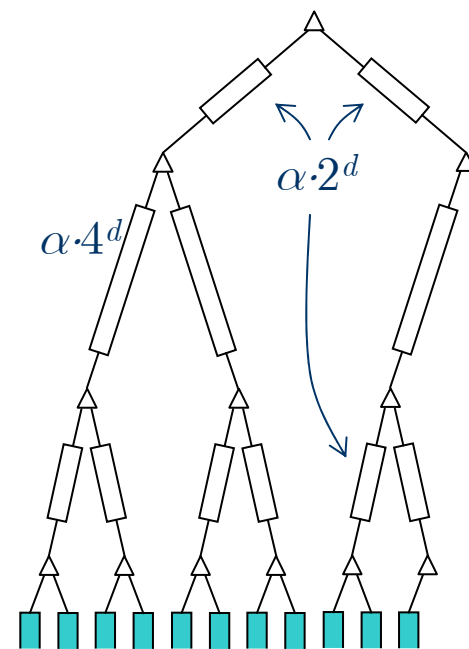
# Cache-oblivious Algorithms – Ideal Cache

- We analyze in the ideal cache model:
  - Automatic replacement
  - Full associativity
  - Optimal replacement strategy: Underlying replacement policy is the optimal offline algorithm.
  - Two levels of memory
  - Tall cache: $M/B \geq cB^{2/(d-1)}$, for some $c > 0$ and $d > 1$.

- Unrealistic assumptions?
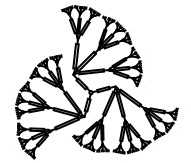
# Cache-oblivious Algorithms – Sorting cont'd.

- Funnelsort and LOWSCOSA achieve optimality by merging with *funnels*.

- A funnel is a tree with buffers on the edges. These buffers are *inputs* and *outputs* of the nodes.

- Buffer capacity is determined by following the van Emde Boas recursion; the capacity of the output buffer of a tree with $k$ inputs is $\alpha k^d$.
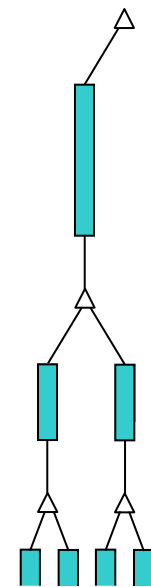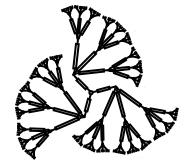
$\alpha \cdot 2^d$

$\alpha \cdot 4^d$

# Merging – Two-phase funnel with refilling

- Elements are merged from the input of a node to the output in a fill() operation.
- In an explicit warm-up phase, fill() is called on all nodes bottom-up. Elements are output from the funnel by then calling fill() on the root.
- When fill() merges at leaf nodes, a custom Refill() function is invoked to signal that elements have been read in from the input of the funnel, so that the space they occupy may be reused.
- fill() merges until either the output is full or one of the inputs is empty. In the latter case, it calls recursively the fill the input. In the first, it is done.
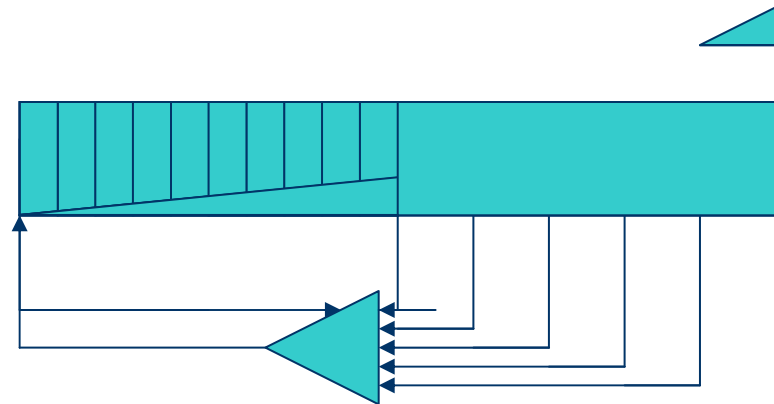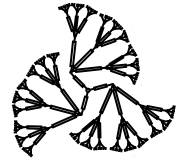
Refill()

# LOWSCOSA

- World's first low-order working space cache-oblivious sorting algorithm.
  1. Partition small elements to the back.
  2. Sort recursively (or by using funnelsort).
  3. Attach refiller that moves elements from the front of the array to newly freed space in the input streams.
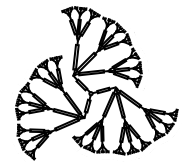  4. Sort right half recursively.

# Engineering Cache-oblivious Sorting Algorithms

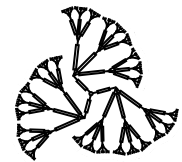Master's thesis by

Kristoffer Vinther
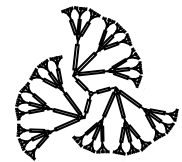
# Algorithm Engineering

It's all about speed!

- ...and
  - Correctness
  - Robustness
  - Flexibility
  - Portability

# Algorithm Engineering – What is speed?

- Theoretician: Asymptotic worst-case running time.

- Algorithm engineer:
  - Good asymptotic performance
  - Low proportionality constants
  - Fast running times on real-world data
  - Robust performance across variety of data
  - Robust performance across variety of platforms
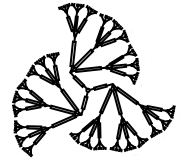
# Algorithm Engineering – How to gain speed?

- Optimize low-level data structures.
- Optimize low-level algorithmic details.
- Optimize low-level coding.
- Optimize memory consumption.
- Maximize locality of reference.

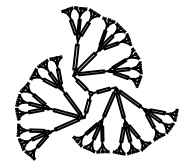A good understanding of the algorithms is extremely important.

# Algorithm Engineering
## – Pencil & paper vs. implementation

- Moret defines algorithm engineering as ”Transforming ”paper-and-pencil” algorithms into efficient and useful implementations.”



- Filling in the details.
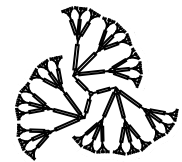
# Experimental Methodology
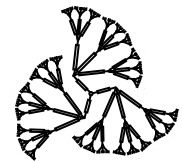
# Methodology
## – Algorithmic details

- How should the funnel be laid out in memory?
- How do we locate nodes and buffers?
- How should we implement merge functionality?
- What is a good value for $z$ and how do we merge multiple streams efficiently?
- How do we reduce the overhead of the sorting algorithm?
- How do we sort at the base of the recursion?
- What are good values for $\alpha$ and $d$?
- How do we handle the output of the funnel?
- How do we best manage memory during sorting?
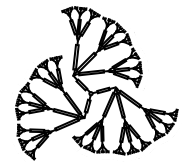- ...

# Methodology – Algorithmic details cont'd.

- Inspired by knowledge of the memory hierarchy and modern processor technology, we develop several solutions to each of these questions.
- All solutions are implemented and benchmarked to locate the best performing combination of approaches.
- It turns out, the simpler the faster (except perhaps memory management).
- Increasing $\alpha$ and $d$ is a cheap way of decreasing the overhead of the funnel.
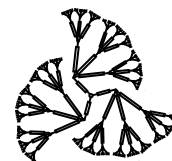
# Methodology – What answers do we seek?

- Are the assumptions of the ideal cache model too unrealistic, i.e. are the algorithms *only* optimal/competitive under ideal conditions?
- Will the better utilization of caches improve running time of our sorting algorithms?
- Will the better utilization of virtual memory improve running time of our sorting algorithms?
- Can our algorithms compete with classic instruction count optimized RAM-based sorting algorithms and memory-tuned cache-aware EM-based sorting algorithms?
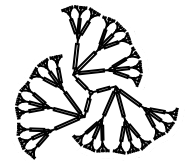
# Methodology – Platforms

- To avoid "accidental optimization," we benchmark on several different architectures:
  - MIPS R10k: Classic RISC; short pipeline, large L2 cache, low clock rate, software TLB miss handling. 64-bit.
  - Pentium 3: Classic CISC; twice as deep a pipeline as MIPS, good branch prediction, many execution units.
  - Pentium 4: Extremely deep pipeline, compensated by very good branch prediction. Very high clock rates.
- Several different operating systems supported: IRIX (64-bit), Linux (32-bit), Windows (32-bit). Benchmarks run on IRIX and Linux.
- Tested with several different compilers: GCC, MSVC, MIPSPRO, ICC.
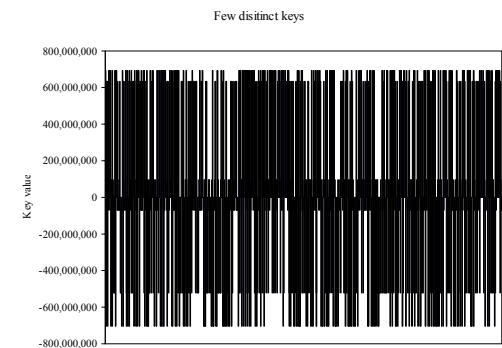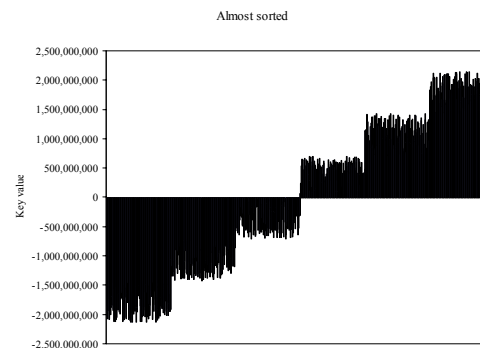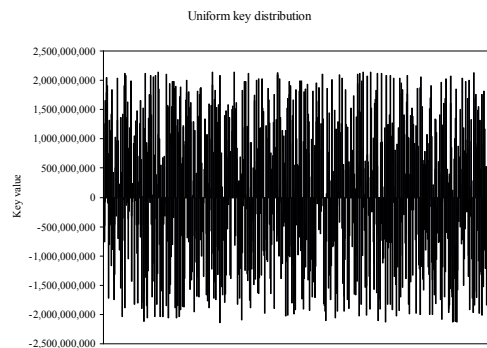
# Methodology – Data types

- To demonstrate robustness, we benchmark several different data types:
  - Key/pointer pairs: **class** { **long** key; **void** *p; }
  - Simple keys: **long**.
  - Records: **class** { **char** record[100]; }
    - Inspired by the official sorting benchmark, Datamation Benchmark.
    - Order determined by strncmp().

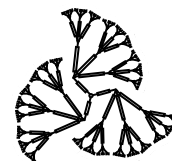# Methodology – Input data

- To demonstrate robustness, we benchmark several different input distributions:
  - Uniformly distributed.
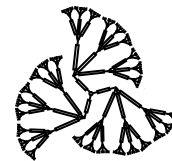  - Almost sorted.
  - Few distinct elements.



Uniform key distribution



Almost sorted



Few disitinct keys
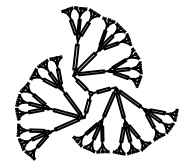
# Methodology – What to measure

- Primarily wall clock time.
- CPU time is no good, since it does not take into account the time spent waiting for page faults to be serviced.
- L2 cache misses.
- TLB misses.
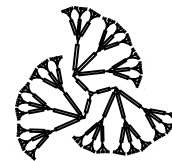- Page faults.

# Methodology – Validity

- For time considerations, we run benchmarks only once.

- Benchmarks are run on such massive datasets that they each take several minuets, even several hours.

- Periodic external disturbances affect all algorithms, are always present, and cannot be eliminated by e.g. averaging.

# Methodology – Competitors

- To answer the question of whether our cache-oblivious algorithms can compete with RAM-based and memory-tuned cache-aware sorting algorithms, we compare them with
  - *Introsort*, developed by SGI as part of STL.
  - *Multiway mergesort*, a part of TPIE, tuned for disk.
  - *Multi-mergesort*, developed by Kubricht *et al.*, tuned for L2 cache.
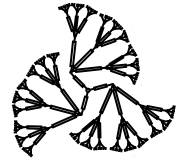  - *Tiled mergesort*, developed by Kubricht *et al.*, tuned for L2 cache.

# Methodology – The problem

- A file stored on a local disk in a native file system contains a number of contiguous elements.
- The problem is solved when there exist a (possibly different) file with the same elements stored contiguously in non-decreasing order.
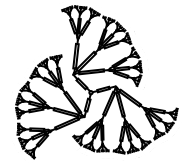- No part of the (original) file is in memory when sorting begins.

Motivation: We don't want to favor any particular initial approach; we believe that real-life applications of sorting doesn't.
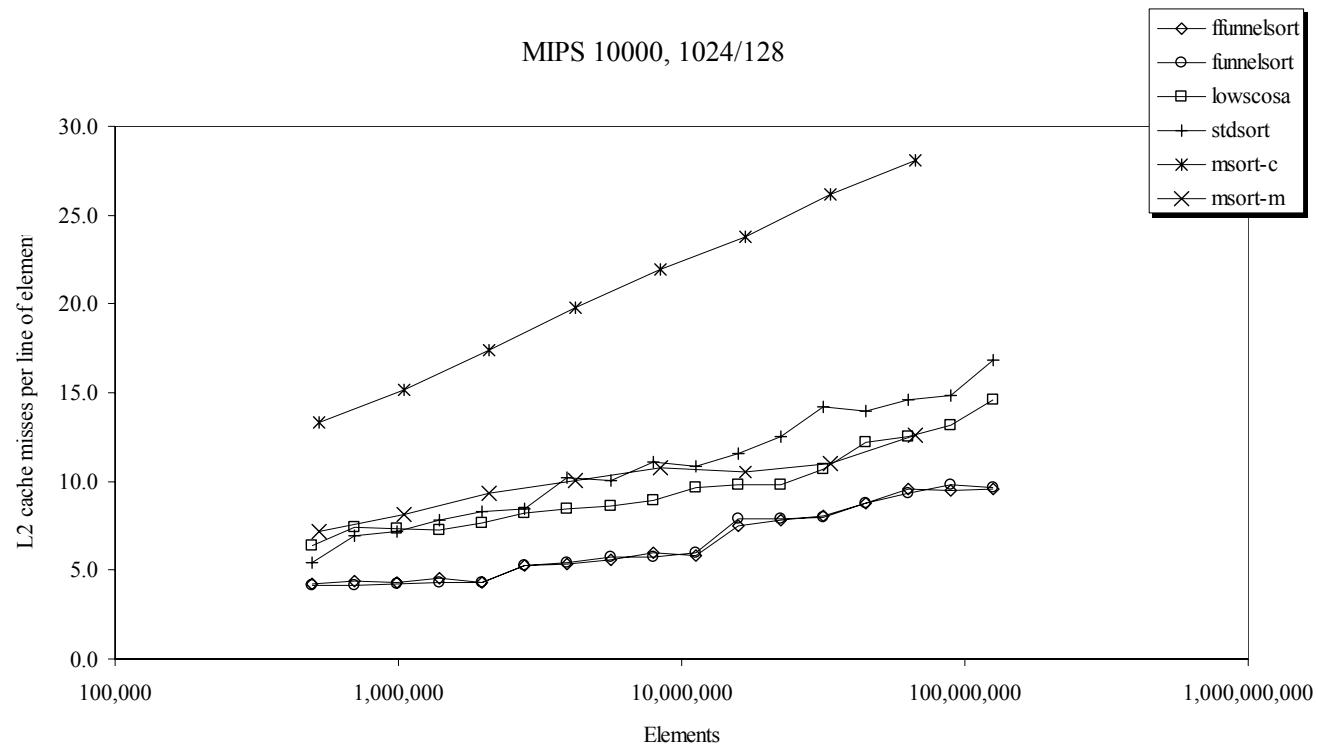
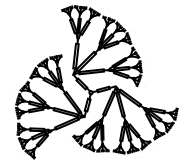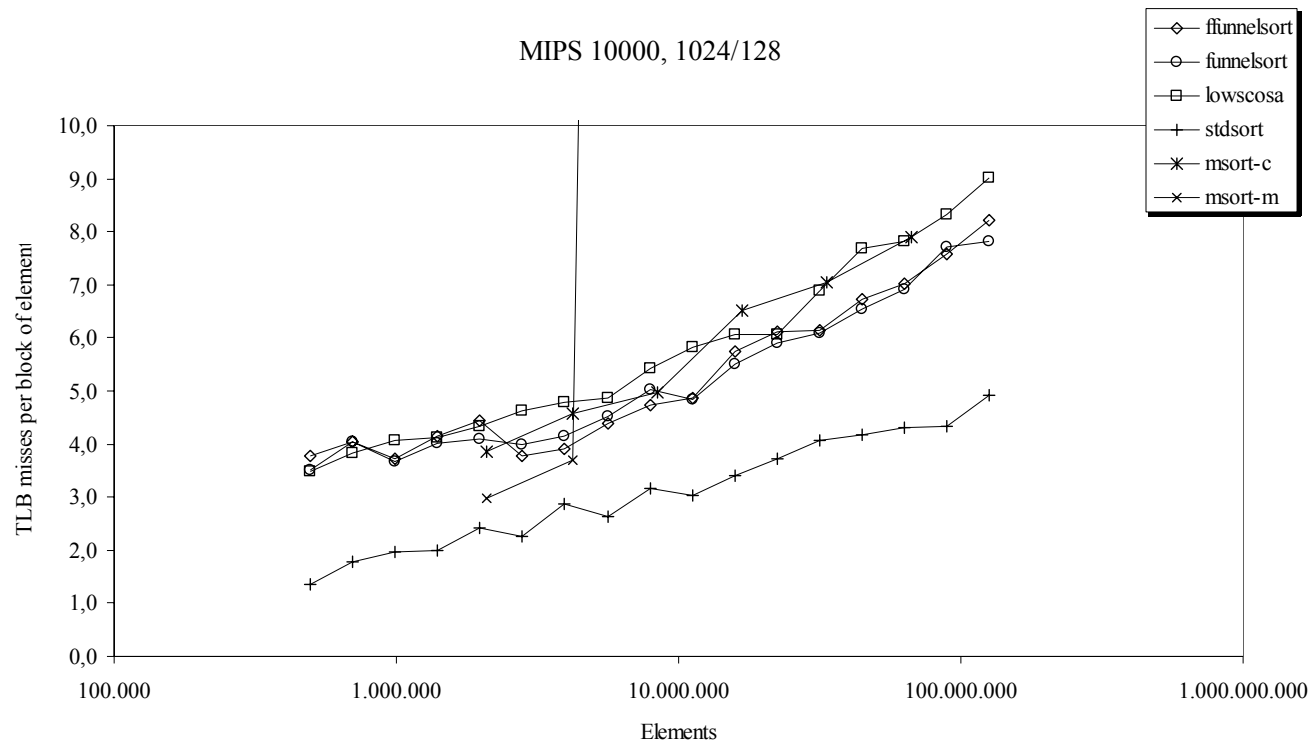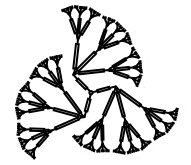Inspiration: Datamation Benchmark.

# Results

# Results
## – L2 cache misses



MIPS 10000, 1024/128

Legend: ffunnelsort, funnelsort, lowscosa, stdsort, msort-c, msort-m

Y-axis: L2 cache misses per line of elemen

X-axis: Elements

# Results – TLB misses



MIPS 10000, 1024/128

Legend:
- ffunnelsort
- funnelsort
- lowscosa
- stdsort
- msort-c
- msort-m

Y-axis: TLB misses per block of element
X-axis: Elements

# Results
# – Page faults



Pentium 3, 256/256

# Results
# – Wall clock time, Pentium 3



Pentium 3, 256/256

Legend:
- ffunnelsort
- funnelsort
- lowscosa
- stdsort
- ami_sort
- msort-c
- msort-m

Y-axis: Wall clock time per element (100.0μs, 10.0μs, 1.0μs, 0.1μs)

X-axis: Elements (1,000,000; 10,000,000; 100,000,000; 1,000,000,000)

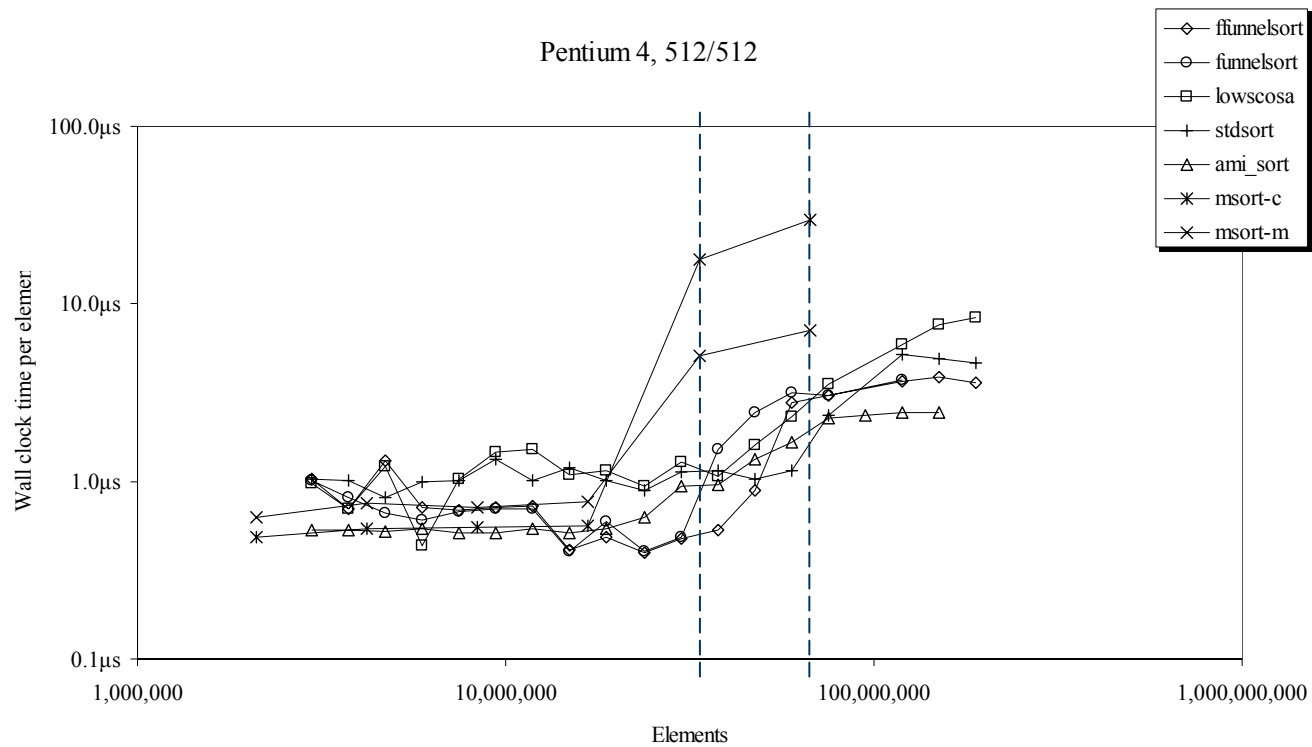# Results
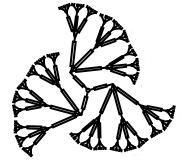# – Wall clock time, Pentium 4



Pentium 4, 512/512

# Results
# – Wall clock time, MIPS R10000
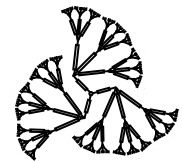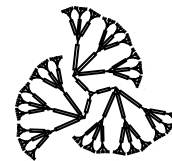


MIPS 10000, 1024/128

# Conclusion

# Conclusion

- Very high performing generic sorting algorithm.
- Unique to our algorithms, performance remains robust
  - across wide range of input sizes.
  - on several different data types.
  - on several different input distributions.
  - across several different processor and operating system architectures.

# References

- Bernard M.E. Moret and Henry D. Shapiro. Algorithms and Experiments: The New (and Old) Methodology, Journal of Universal Computer Science, 7:434-446, 2001.

- Alok Aggarwal and Jeffery S. Vitter. The input/output complexity of sorting and related problems. Communications of the ACM, 1988.

- Matreo Frigo, Charles E. Leiserson, Harald Prokop, and Shidhar Ramachandran. Cache-oblivious algorithms. Proceedings of the 40th Annual Symposium on Foundations of Computer Science, New York, 1999.

- David S. Johnson. A Theoretician's Guide to the Experimental Analysis of Algorithms. Proceedings of the 5th and 6th DIMACS Implementation Challenges. Goldwasser, Johnson, and McGeoch (eds), American Mathematical Society, 2001.

- Datamation Benchmark. Sort Benchmark Home Page, hosted by Microsoft. World Wide Web document, http://research.microsoft.com/barc/SortBenchmark/, 2003.

- Duke University. A Transparent Parallel I/O Environment, World Wide Web Document, http://www.cs.duke.edu/TPIE/, 2002.

- Li Xiao, Xiaodong Zhang, and Stefan A. Kubricht. Improving Memory Performance of Sorting Algorithms, ACM Journal of Experimental Algorithms, Vol. 5, 2000.