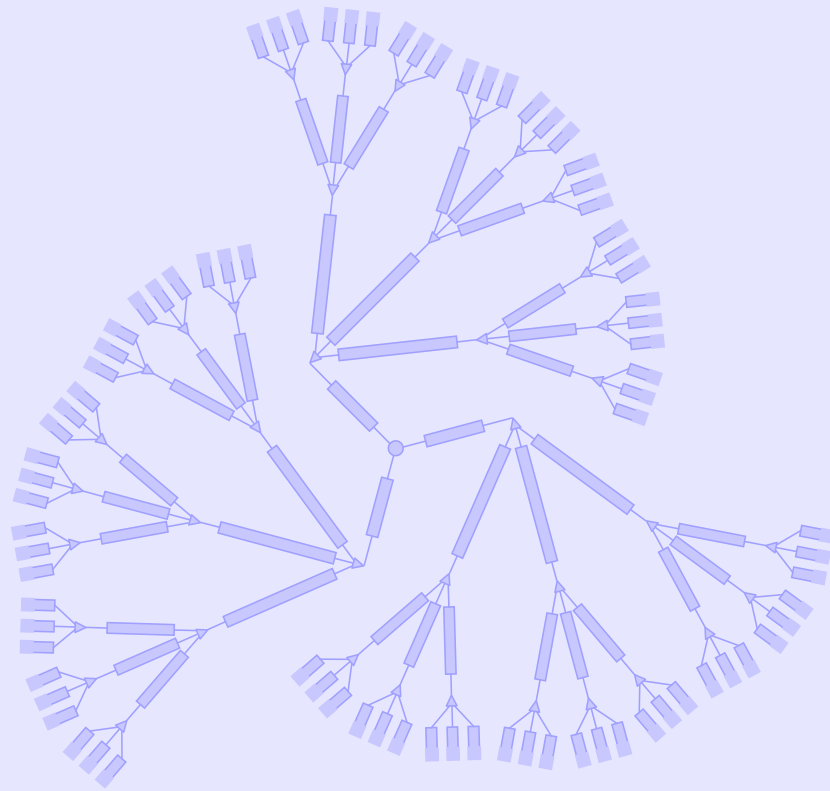


Engineering a Cache-Oblivious Sorting Algorithm



Gerth Stølting Brodal
University of Aarhus

Rolf Fagerberg
University of Aarhus

Kristoffer Vinther
Systematic Software Engineering

Motivation

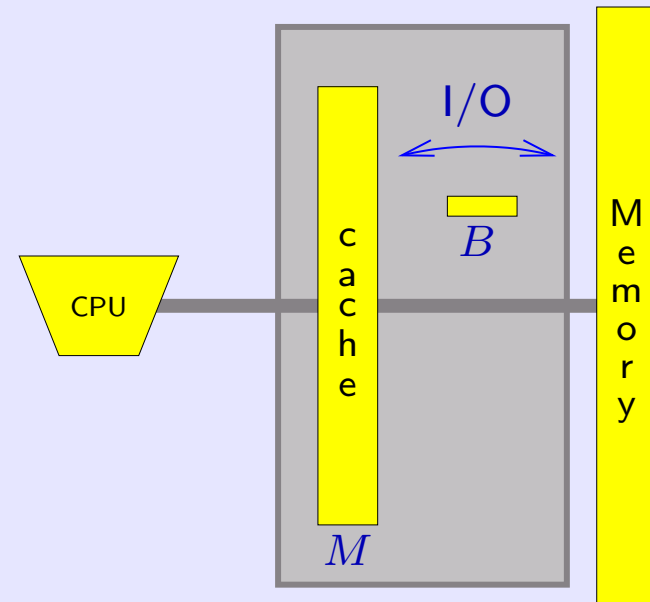
- Memory hierarchy has become a fact of life
- Accessing non-local storage may take a very long time
- Good locality is important to achieving high performance

	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 Model

Frigo, Leiserson, Prokop, Ramachandran, FOCS'99

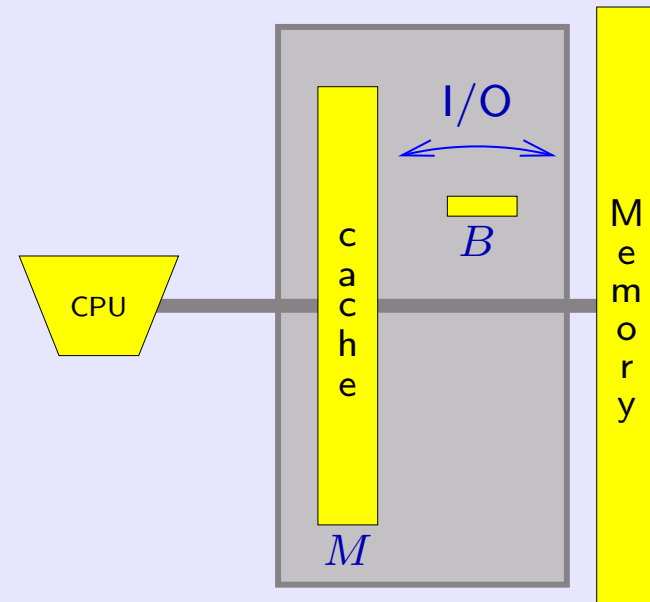
- Program in the RAM model
- Analyze in the I/O model for arbitrary B and M
- Optimal off-line cache replacement strategy



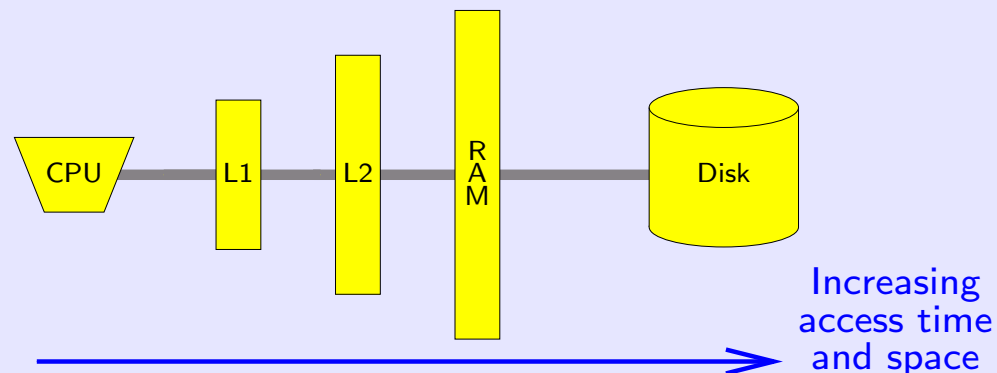
Cache-Oblivious Model

Frigo, Leiserson, Prokop, Ramachandran, FOCS'99

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



- Optimal on arbitrary level \Rightarrow optimal on **all levels**
- Portability



Sorting — I/O Bounds

QuickSort ★

Binary MergeSort ★

$\Theta\left(\frac{M}{B}\right)$ -way MergeSort ★

(Lazy) Funnelsort ★ ★

$$O\left(\frac{N}{B} \cdot \log_2 \frac{N}{M}\right)$$

$$O(\text{Sort}_{M,B}(N))$$

Aggarwal and Vitter 1988

$$O(\text{Sort}_{M,B}(N))$$

Frigo, Leiserson, Prokop and Ramachandran 1999

Brodal and Fagerberg 2002

★ cache-aware

★ cache-oblivious

★ requires $M \geq B^{1+\varepsilon}$

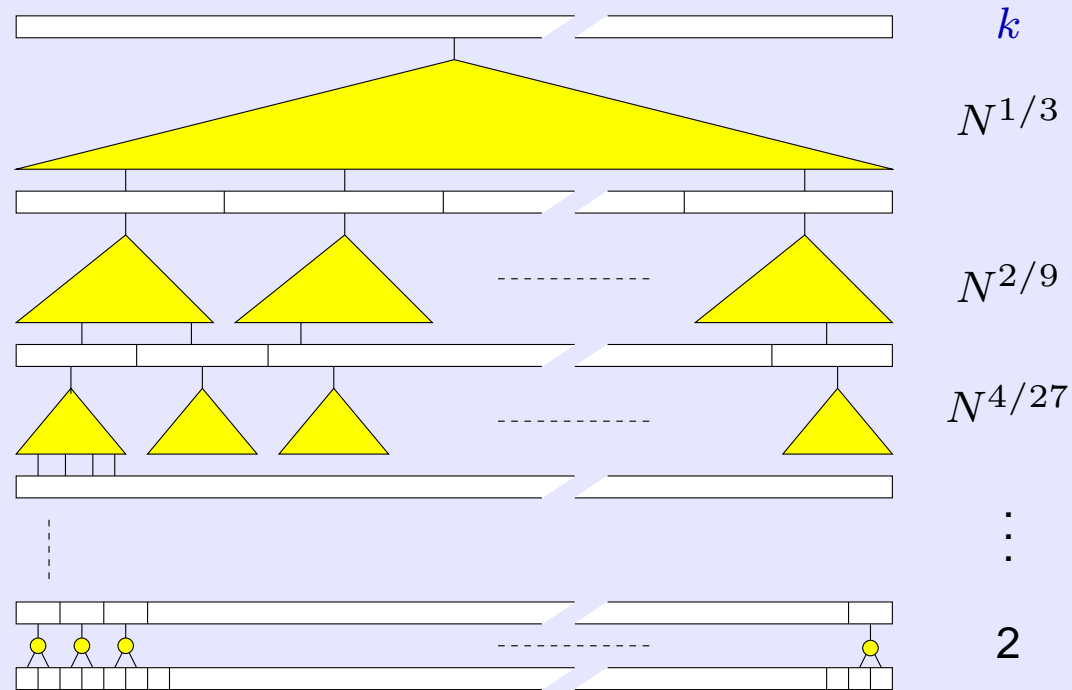
$$\text{Sort}_{M,B}(N) = \frac{N}{B} \cdot \log_{M/B} \frac{N}{B}$$

Lazy Funnel Sort

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

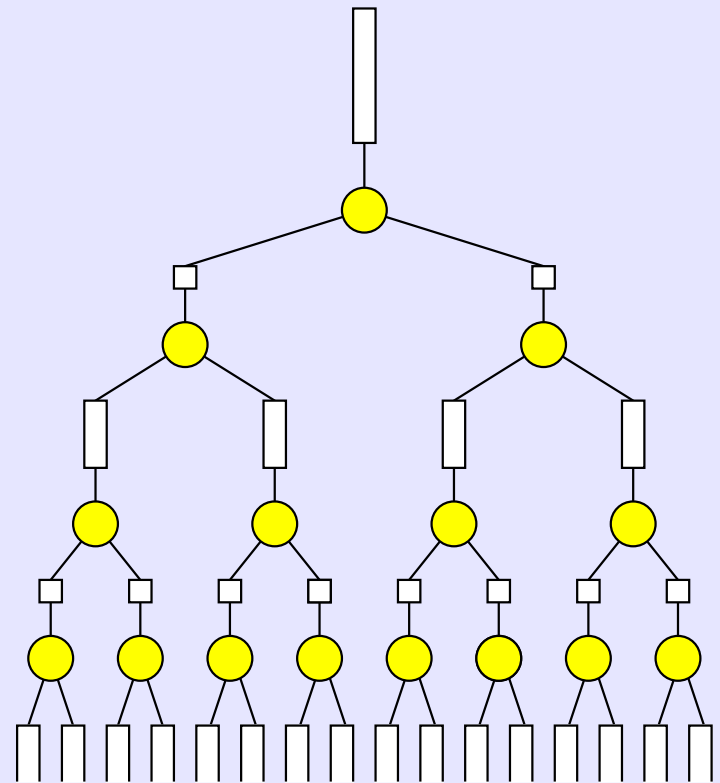
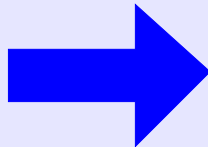
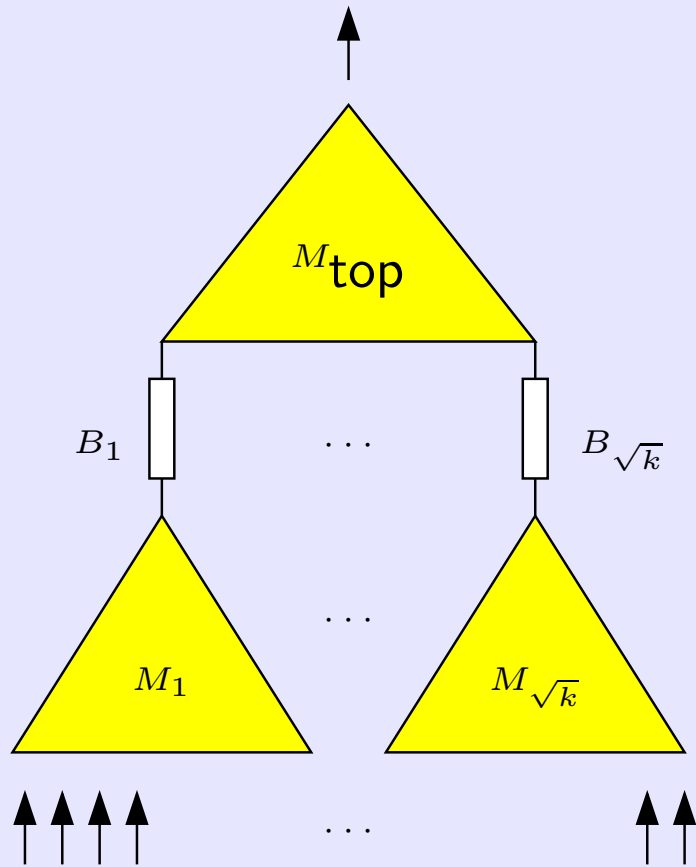
Recursively **Funnel Sort** each segment

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



Brodal and Fagerberg 2002

Lazy k -merger



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

Brodal and Fagerberg 2002

Engineering Lazy Funnel Sort

- Recursive implementation beats iterative implementation
- 4- or 5-way merge beats 2-way merge
- Standard memory allocator beats hand-coded allocator
- Pointer based van Emde Boas layout beats implicit layouts
- Nodes and buffers stored separately beats one layout
- Straightforward beats hand-coded branch elimination in core loop
- $d = 2.5$ and $\alpha = 16$
- Reuse merger data structures
- For k -mergers of height < 2 switch to Quicksort

Evaluating Funnelsort

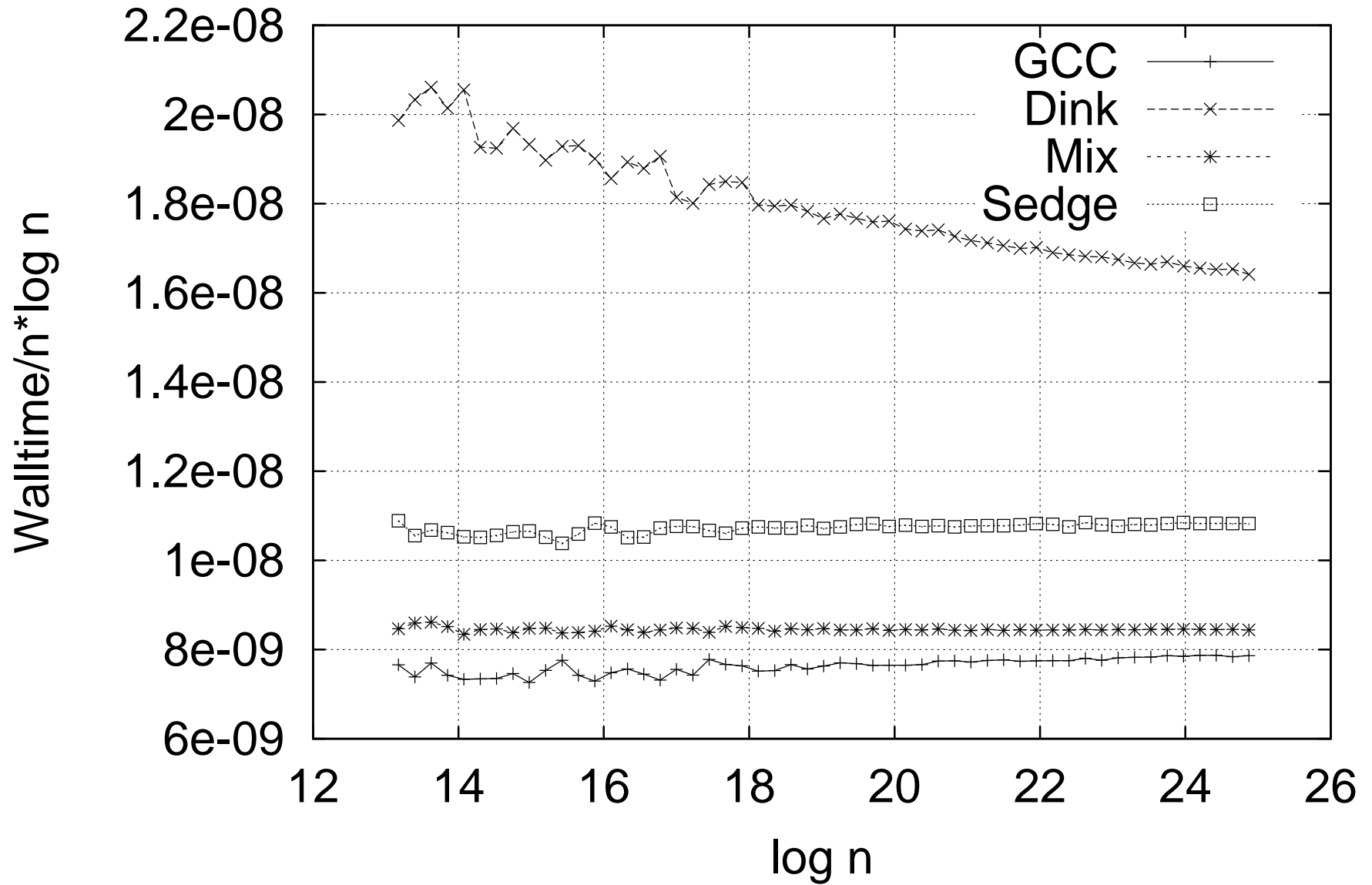
- 2-way and 4-way Funnelsort
- Quicksort
 - STL GCC
 - STL Intel C++
 - Sedgewick
 - Bentley & Maclroy with pivot tuning
- TPIE - optimized for external memory
- Rmerge - optimized for registers
- Mergesort by LaMarca and Ladner

Hardware Specifications

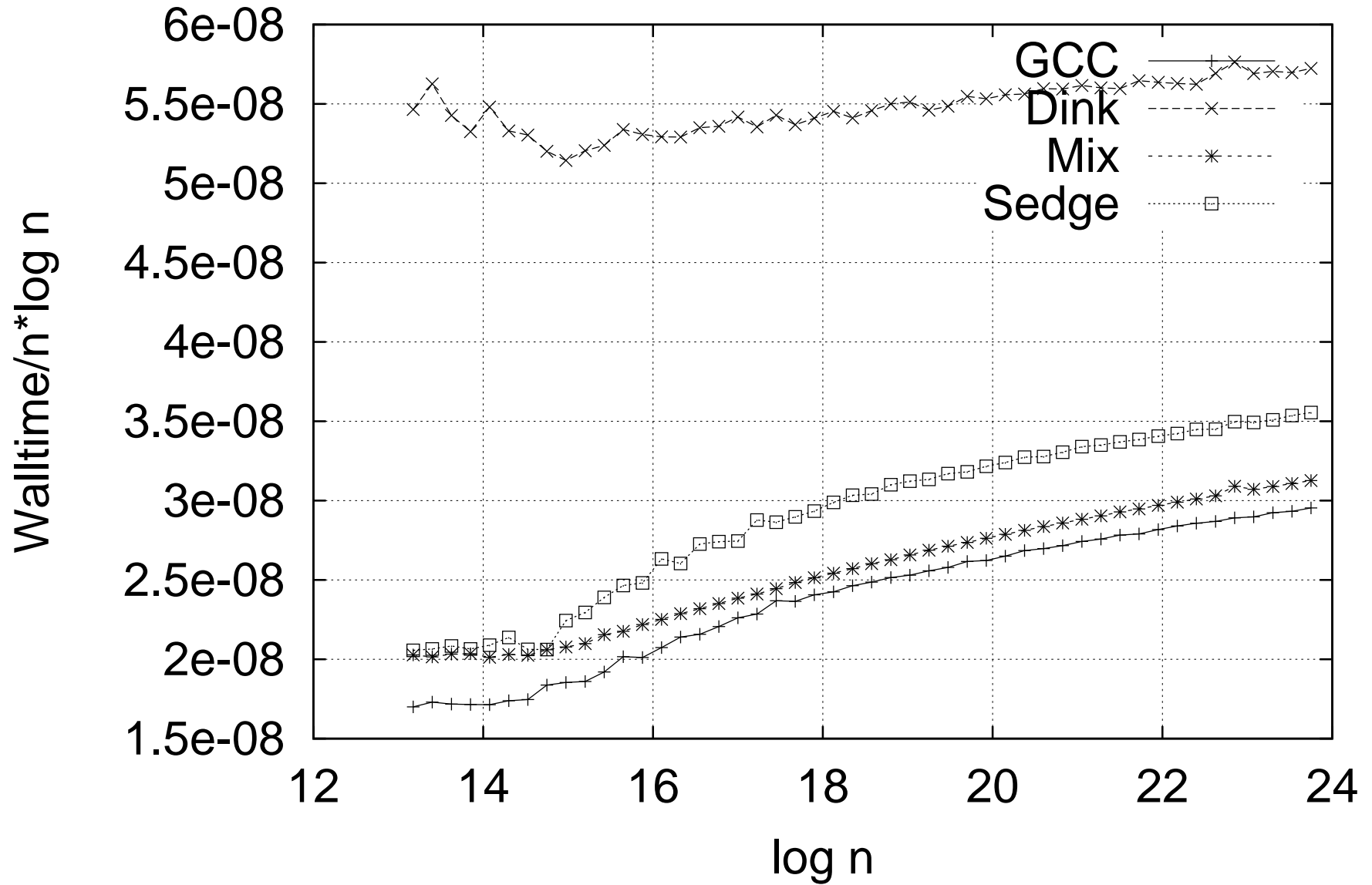
	<i>Pentium 4</i>	<i>Pentium III</i>	<i>MIPS 10000</i>	<i>AMD Athlon</i>	<i>Itanium 2</i>
Architecture type	Modern CISC	Classic CISC	RISC	Modern CISC	EPIC
Operation system	Linux v. 2.4.18	Linux v. 2.4.18	IRIX v. 6.5	Linux 2.4.18	Linux 2.4.18
Clock rate	2400MHz	800MHz	175MHz	1333 MHz	1137 MHz
Address space	32 bit	32 bit	64 bit	32 bit	64 bit
Pipeline stages	20	12	6	10	8
L1 data cache size	8 KB	16 KB	32 KB	128 KB	32 KB
L1 line size	128 B	32 B	32 B	64 B	64 B
L1 associativity	4-way	4-way	2-way	2-way	4-way
L2 cache size	512 KB	256 KB	1024 KB	256 KB	256 KB
L2 line size	128 B	32 B	32 B	64 B	128 B
L2 associativity	8-way	4-way	2-way	8-way	8-way
TLB entries	128	64	64	40	128
TLB associativity	full	4-way	64-way	4-way	full
TLB miss handling	hardware	hardware	software	hardware	?
RAM size	512 MB	256 MB	128 MB	512 MB	3072 MB

Comparison of Quicksort Implementations

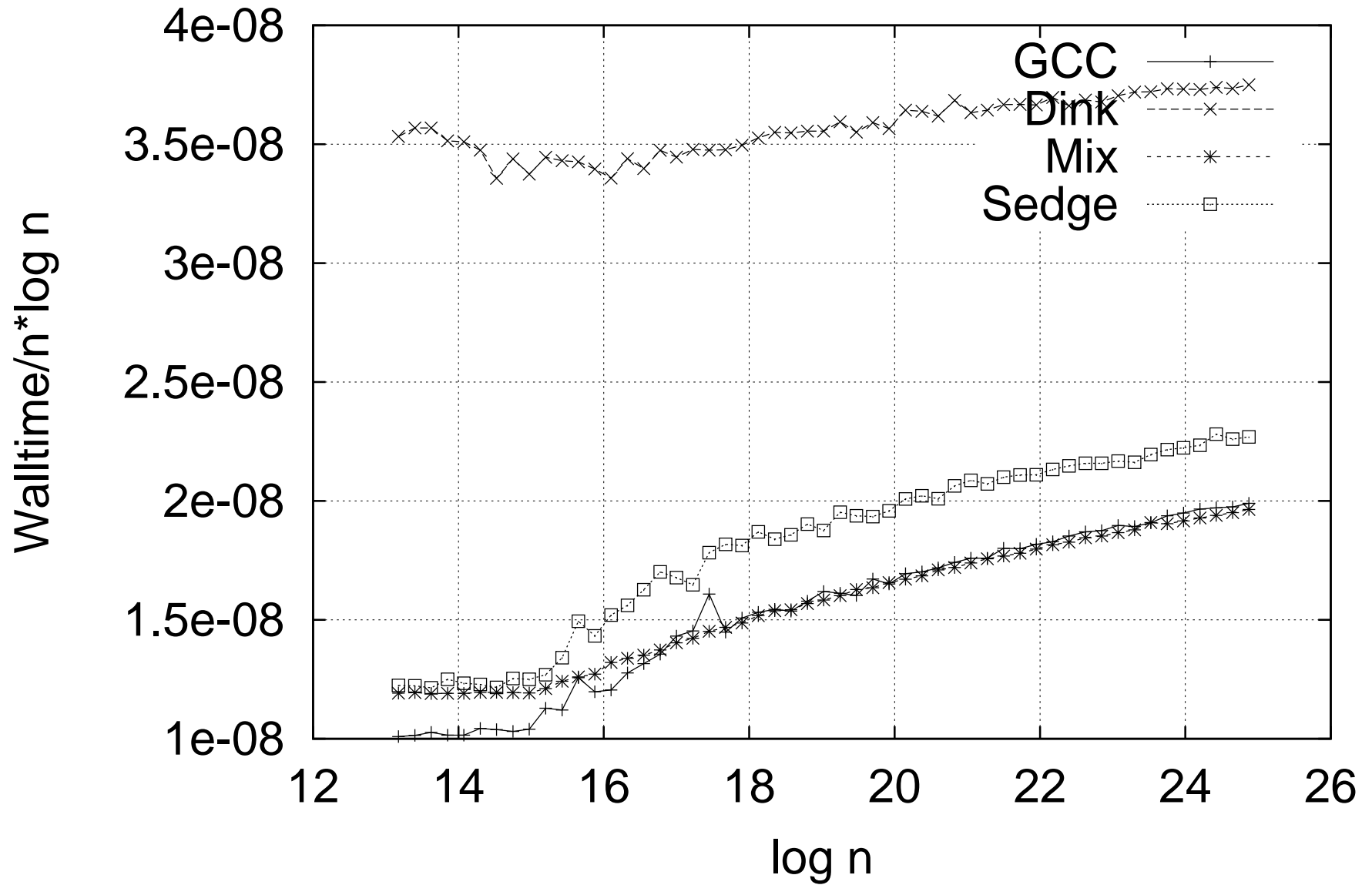
Uniform pairs - Pentium 4



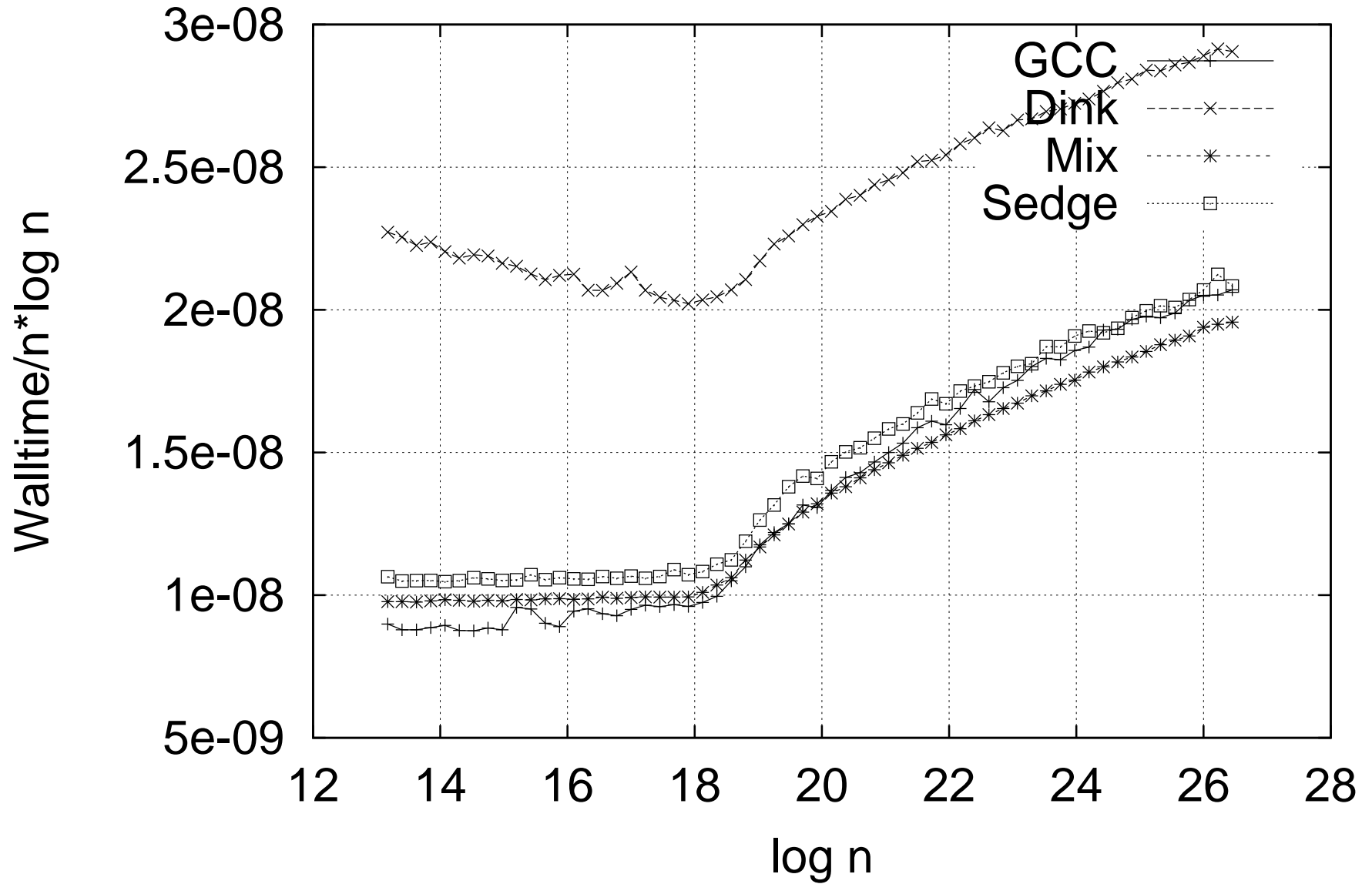
Uniform pairs - Pentium III



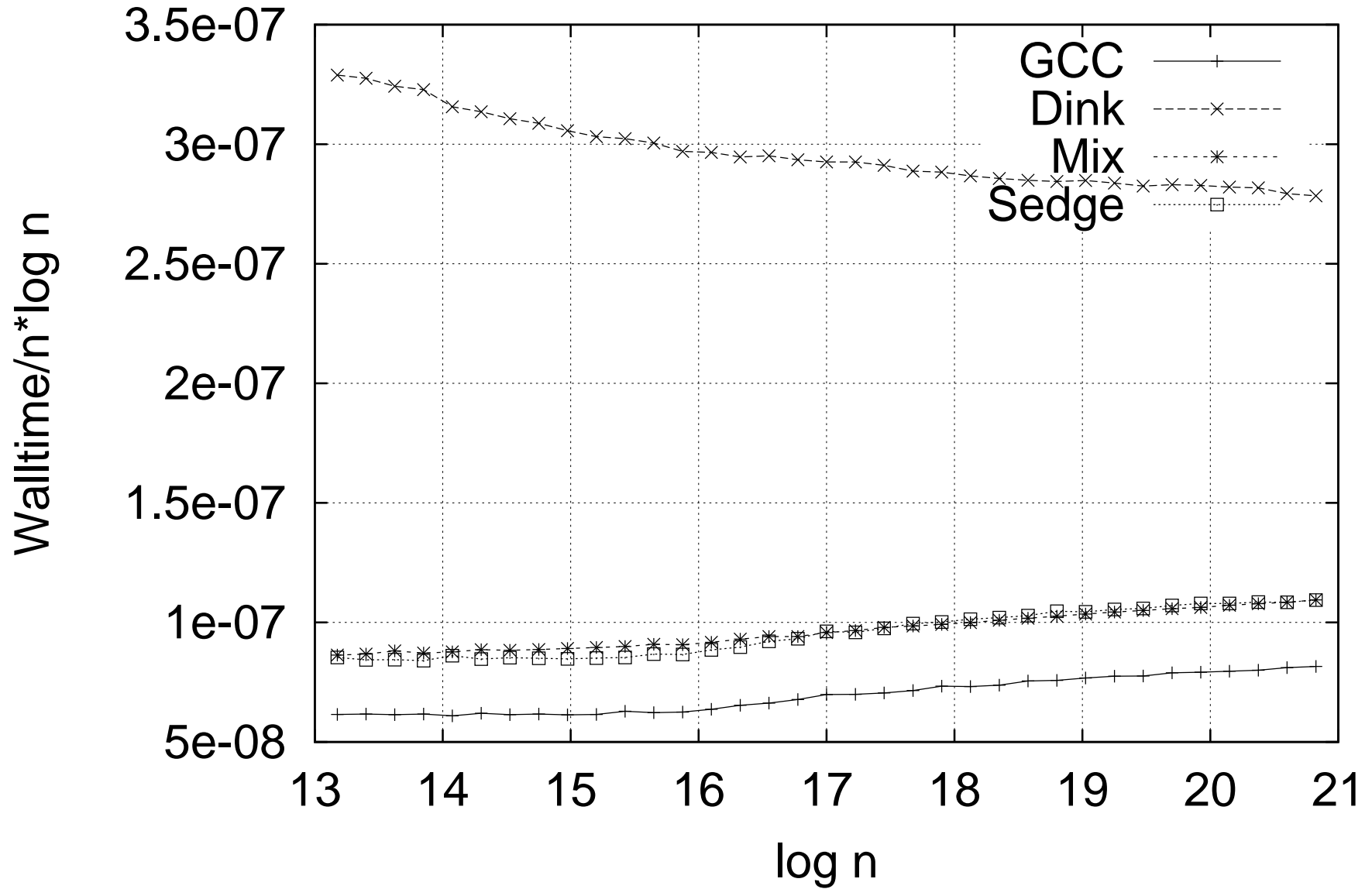
Uniform pairs - AMD Athlon



Uniform pairs - Itanium 2

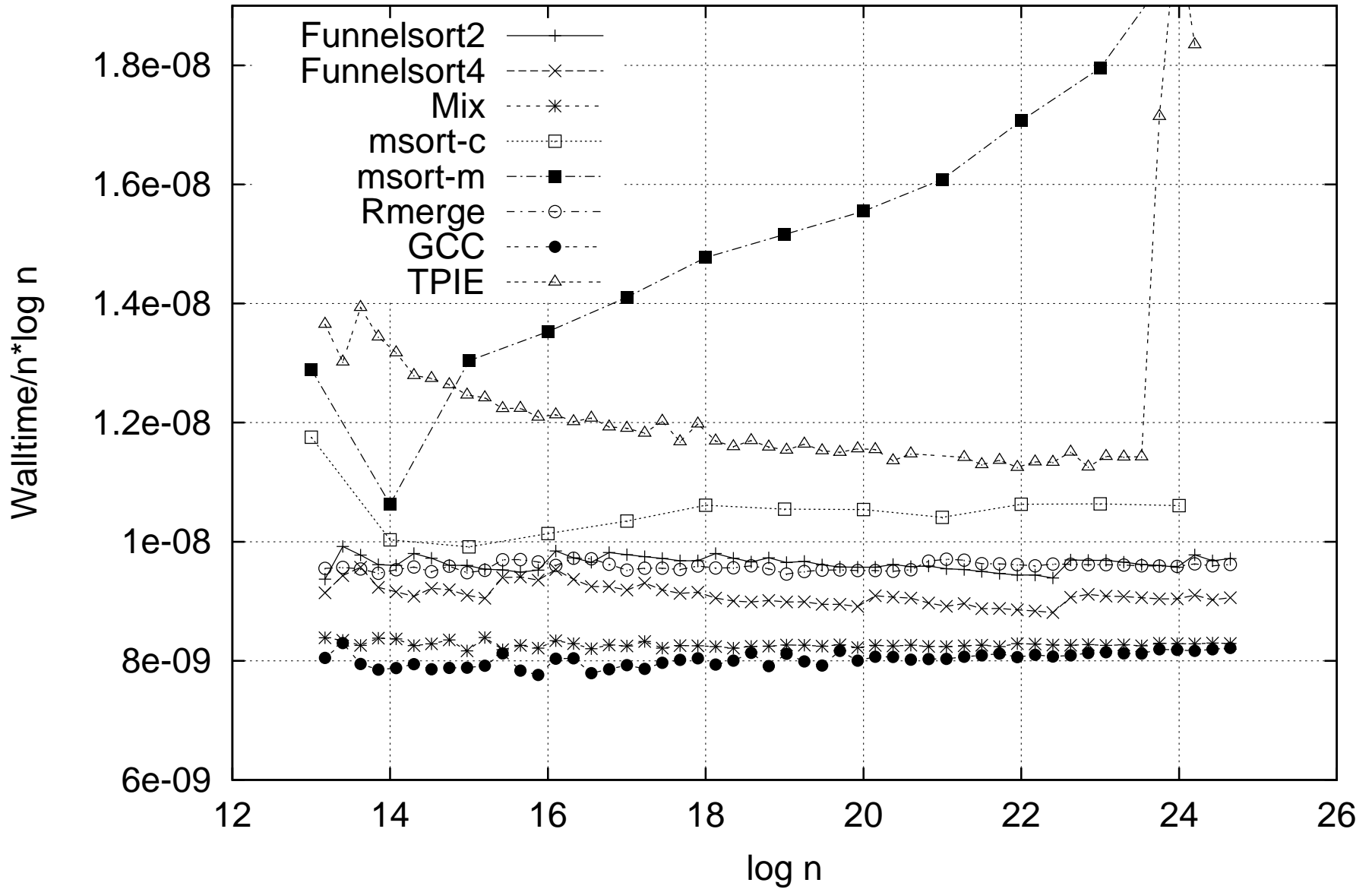


Uniform pairs - MIPS 10000

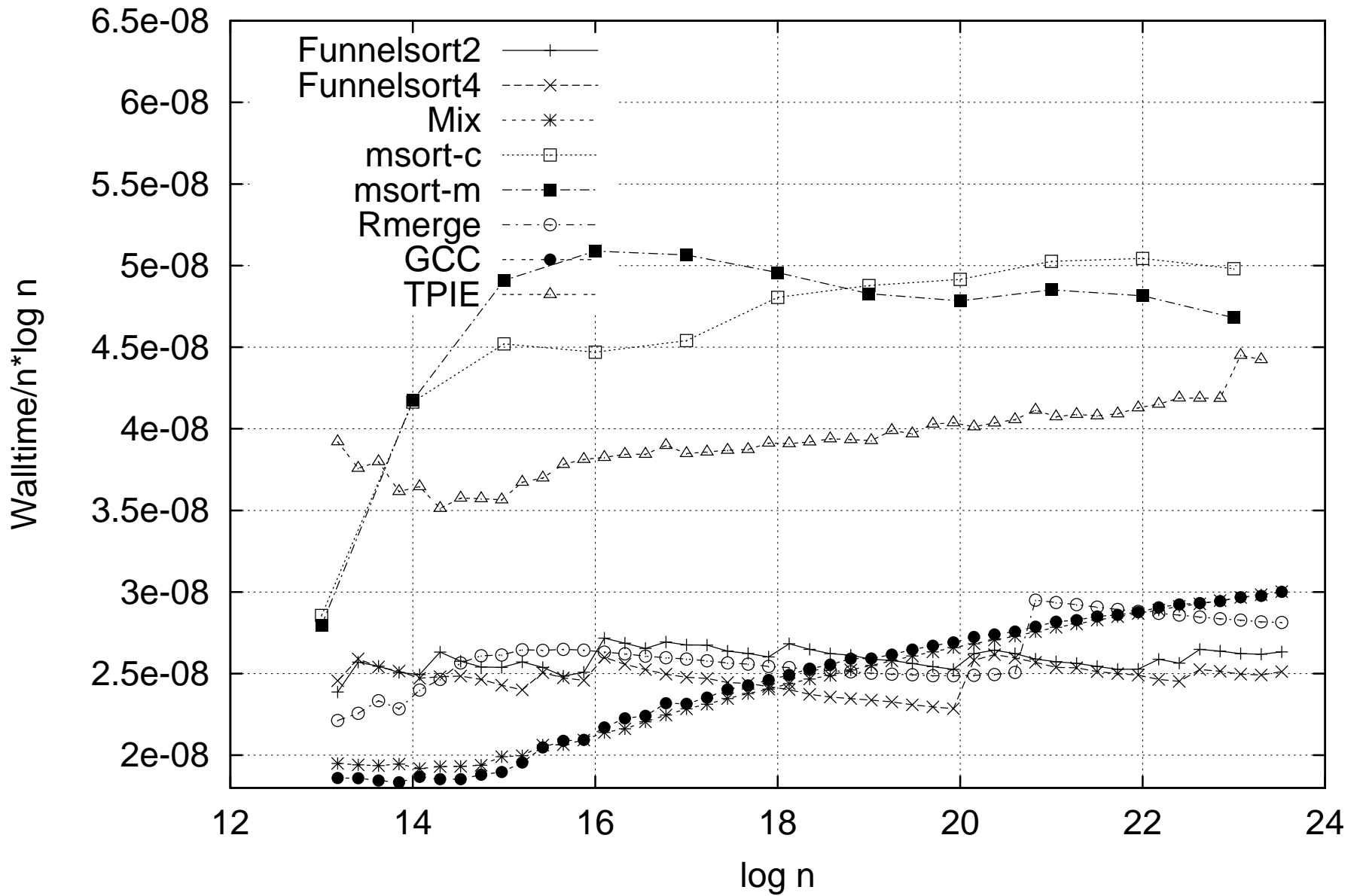


Results for Inputs in RAM

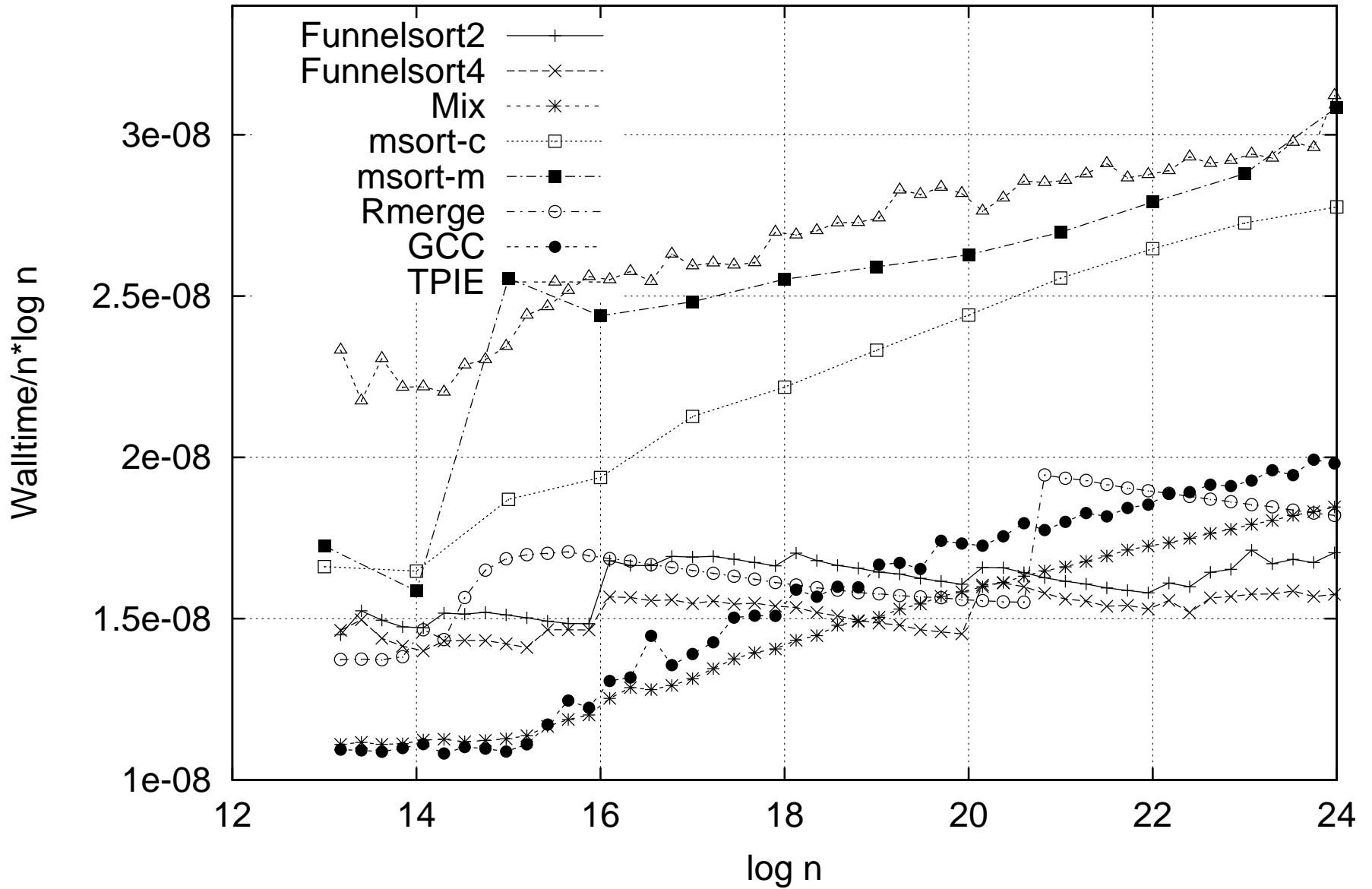
Uniform pairs - Pentium 4



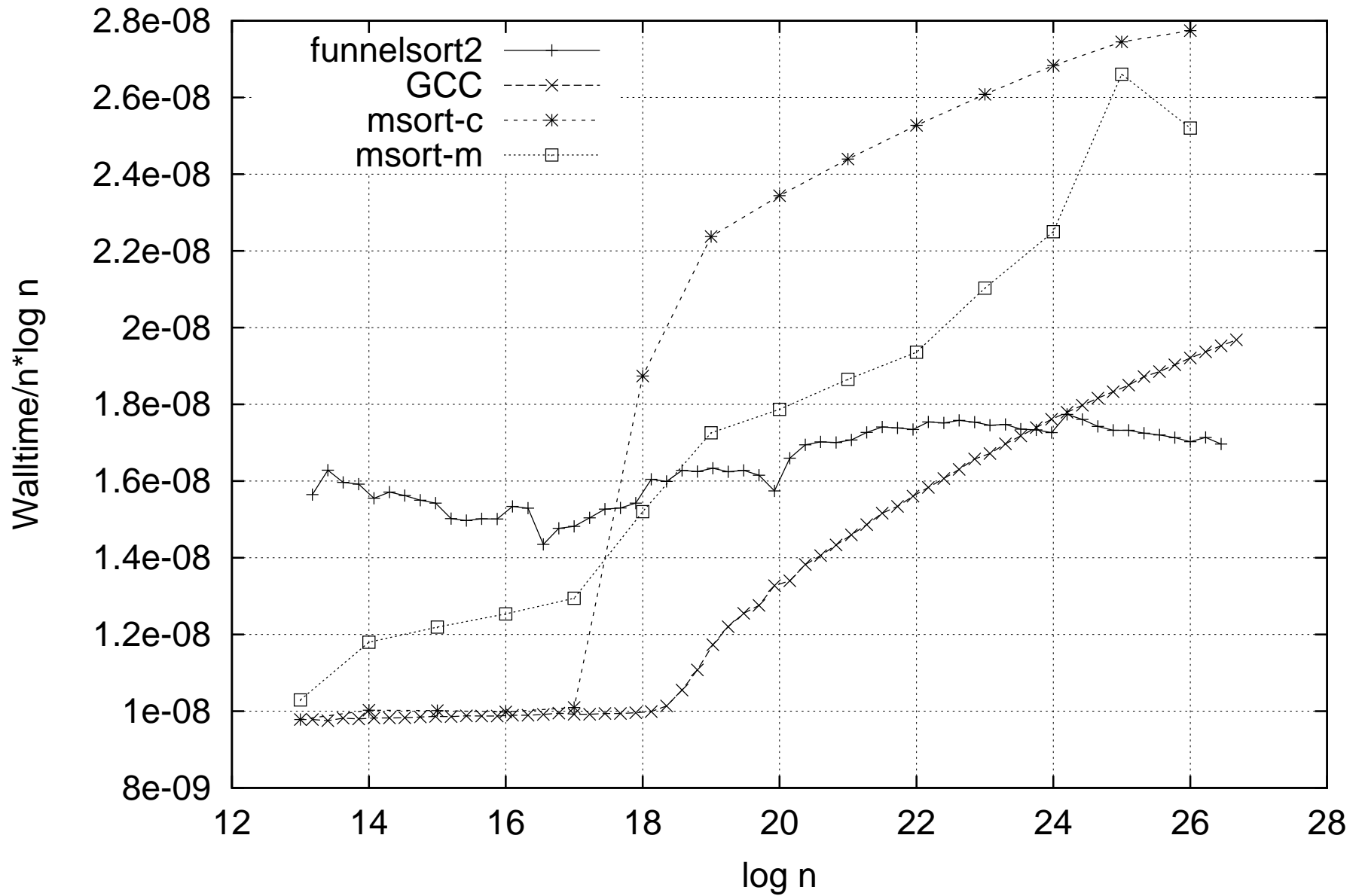
Uniform pairs - Pentium III



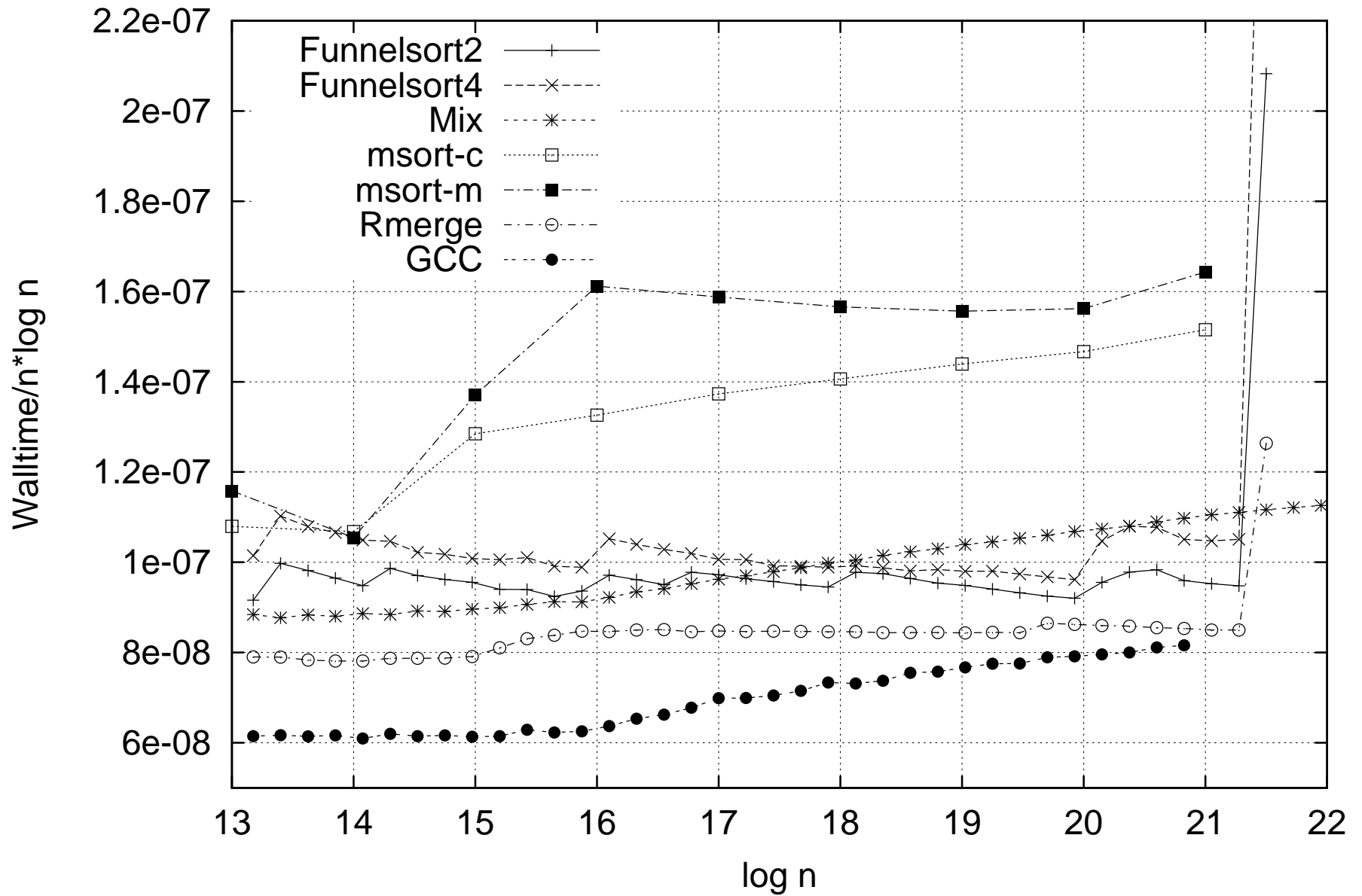
Uniform pairs - AMD Athlon



Uniform pairs - Itanium 2

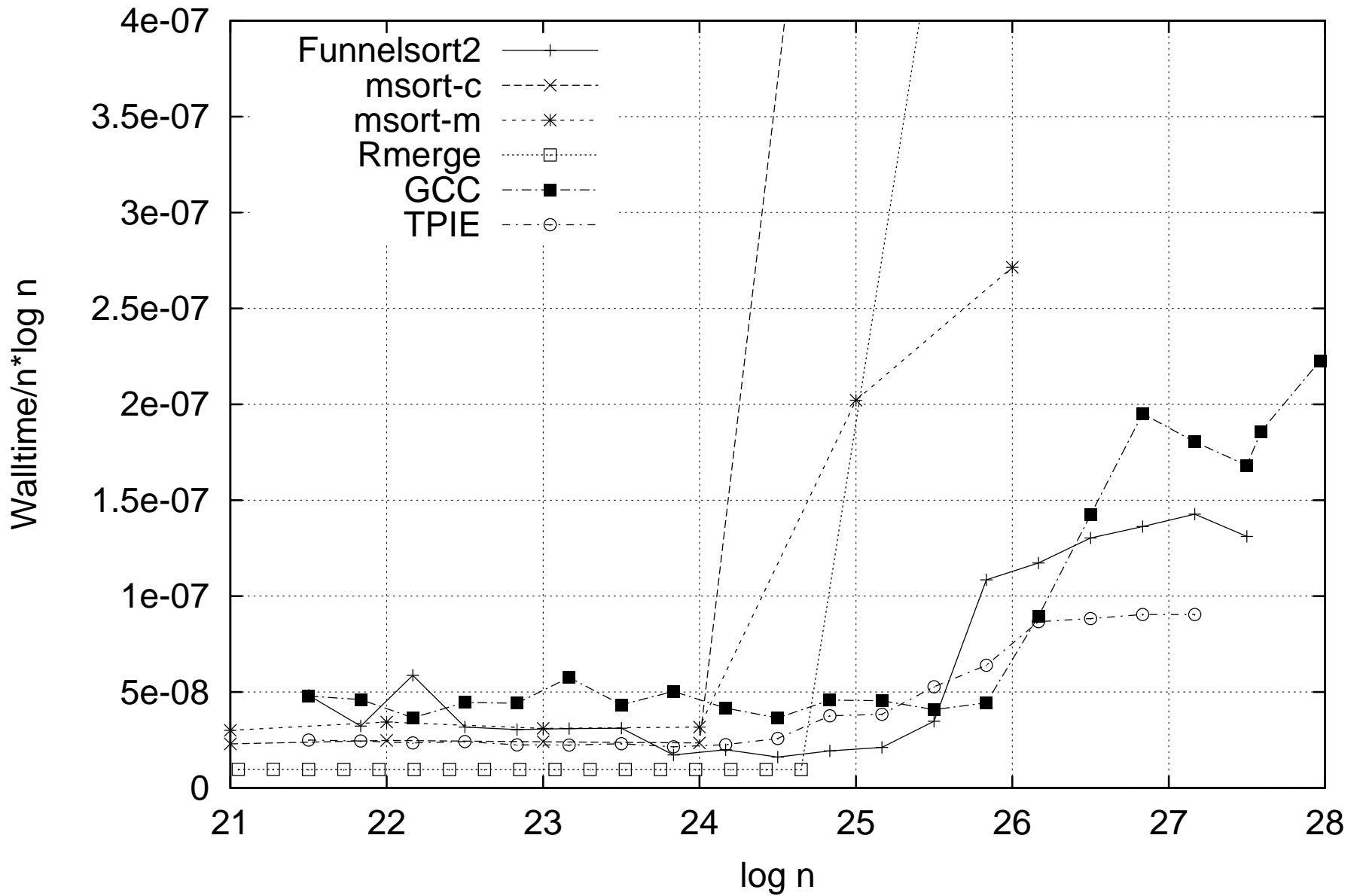


Uniform pairs - MIPS 10000

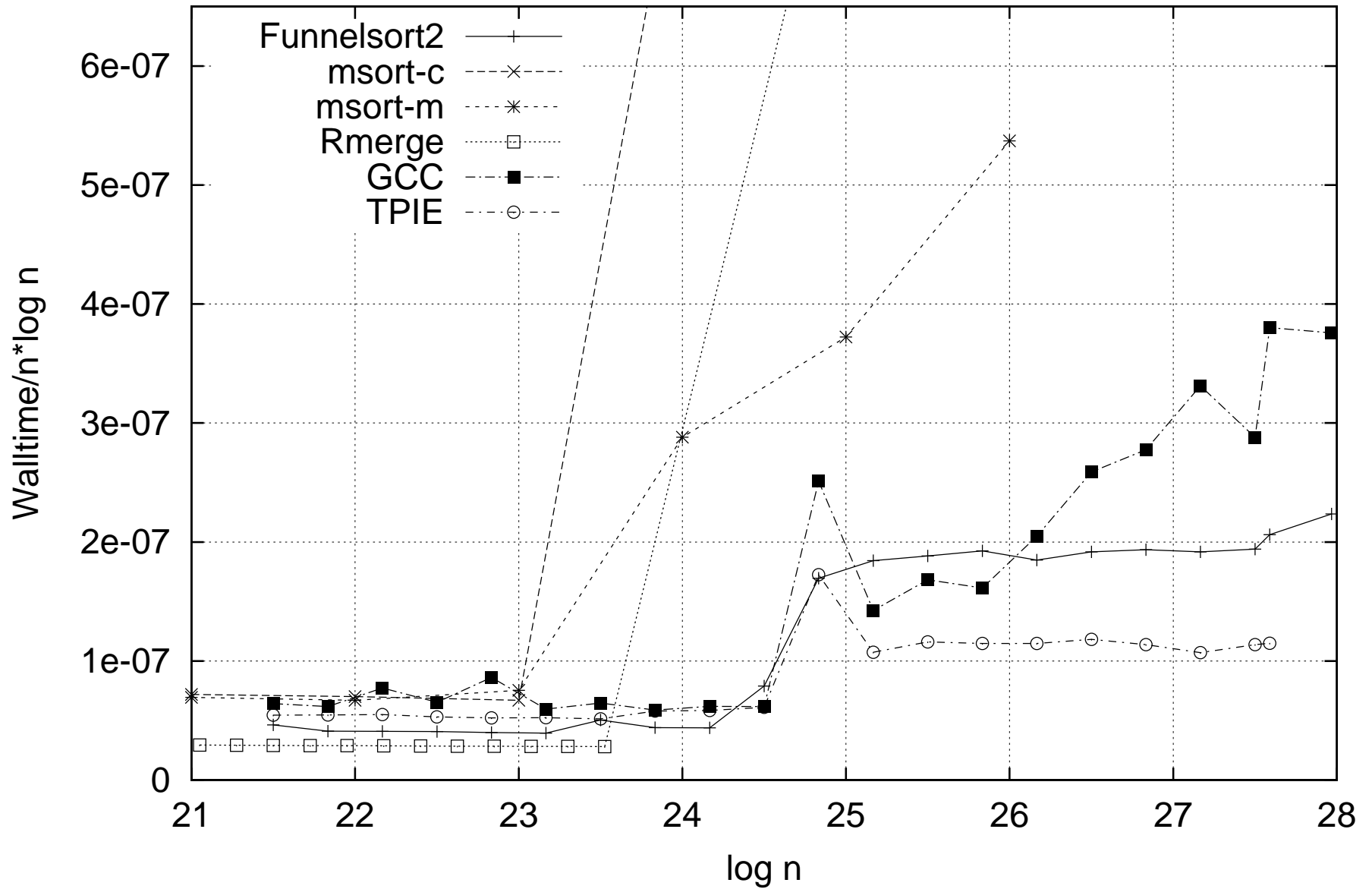


Results for Inputs on Disk

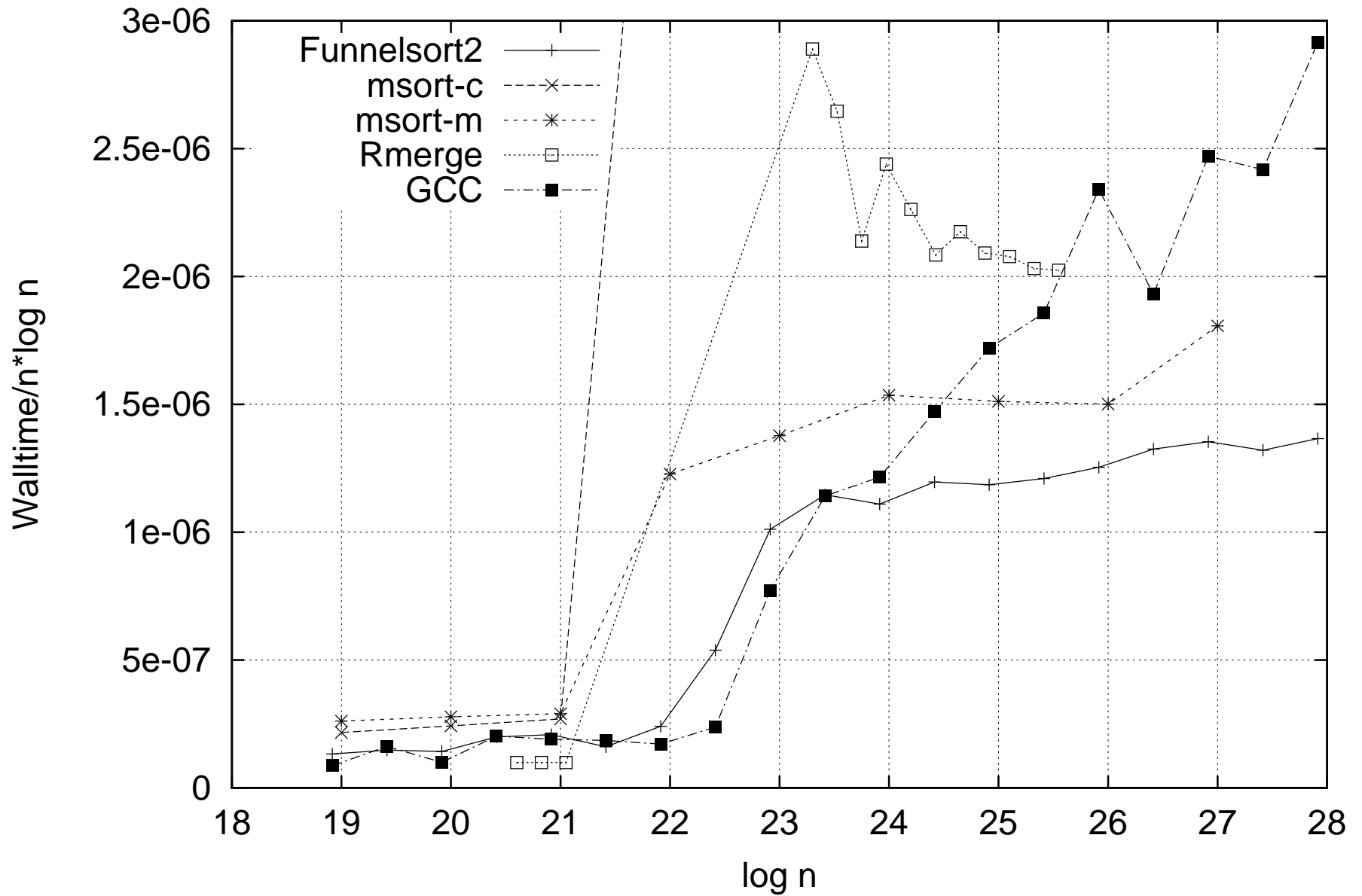
Uniform pairs - Pentium 4



Uniform pairs - Pentium III



Uniform pairs - MIPS 10000



Conclusion

- Very high performing generic sorting algorithm
- Performance remains robust across wide range of input sizes
- Across several different processor and operating system architectures
- On several different data types
- On several different input distributions
- *Overhead involved in being cache-oblivious can be small enough for the nice theoretical properties to actually transfer into practical advantages*