

Cache-Oblivious VAT-Algorithms

Tomasz Jurkiewicz* Kurt Mehlhorn† Patrick Nicholson‡

April 15, 2014

Abstract

The VAT-model (virtual address translation model) extends the EM-model (external memory model) and takes the cost of address translation in virtual memories into account. In this model, the cost of a single memory access may be logarithmic in the largest address used. We show that the VAT-cost of cache-oblivious algorithms is only a constant factor larger than their EM-cost; this requires a somewhat more stringent tall cache assumption than for the EM-model.

1 Introduction

Modern processors have a memory hierarchy and use virtual memory. We concentrate on two-levels of the hierarchy and refer to the faster memory as the cache. Data is moved between the fast and the slow memory in blocks of contiguous memory cells, and only data residing in the fast memory can be accessed directly. Whenever data in the slow memory is accessed, a cache fault occurs and the block containing the data must be moved to the fast memory. In the EM-model of computation, the complexity of an algorithm is defined as the number of cache faults.

In general, many processes are running concurrently on the same machine. Each running process has its own linear address space $0, 1, 2, \dots$, and the operating system ensures that the distinct linear address spaces of the distinct processes are mapped injectively to the physical address space of the processor. To this effect, the operating system maintains a translation tree for each process. The translation from virtual addresses to physical addresses is implemented as a tree walk and incurs cost; see Section 3 for details. The depth of the tree for a particular process is $d = \log_K(m/P)$ where m is the maximum address used by the process, K is the arity of the tree and P is the page size. Typically, $K \approx P \approx 2^{10}$, and K is chosen such that the space requirement of a translation tree node is equal to the size of a page. The translation process accesses d pages and only pages residing in fast memory can be accessed directly. Any node visited during the translation process must be brought into fast memory if not already there, and hence a single memory access may cause up to d cache faults.

The cost of the translation process is clearly noticeable in some cases. Jurkiewicz and Mehlhorn [JM13] timed some simple programs and observed that for some of them the quotient

$$\text{measured running time for input size } n / \text{RAM-running time for input size } n$$

seems to grow logarithmically in n ; see Figure 1. Jurkiewicz and Mehlhorn introduced the VAT-model as an extension of the EM-model to account for the cost of address translation; see Section 3 for a definition of their model. They showed that the growth rates of the measured running times of the programs mentioned in Figure 1 are correctly predicted by the model.

The EM-model penalizes non-locality and the VAT-model penalizes it even more. Cache-oblivious algorithms [FLPR12] show good data locality for all memory sizes. Jurkiewicz and Mehlhorn [JM13] showed that some cache-oblivious algorithms, namely those that do not need a tall cache assumption, also

*Google Zürich, email: tomasz.tojot.jurkiewicz@gmail.com

†MPI for Informatics, email: mehlhorn@mpi-inf.mpg.de

‡MPI for Informatics, email: pnichols@mpi-inf.mpg.de

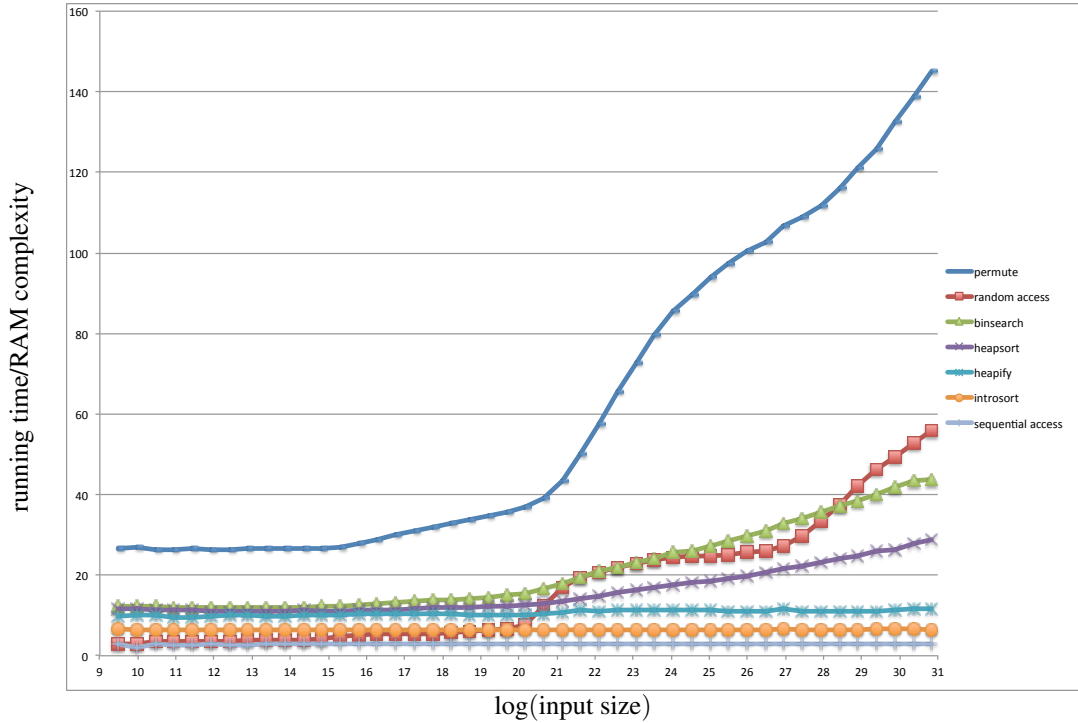


Figure 1: The abscissa shows the logarithm of the input size. The ordinate shows the measured running time divided by the RAM-complexity (normalized operation time). RAM complexity of the programs shown is either cn or $cn \log n$; [JM13] ignores lower order terms. The constants c were chosen such that the different plots fit nicely into the same figure. The normalized operation times of sequential access, quicksort (introsort), and heapify are constant, the normalized operation times of permute, random scan, repeated binary search, heapsort grow as functions of the problem size. Note that a straight-line corresponds to a linear function in $\log(\text{input size})$.

perform well in the VAT-model. Their paper poses the question of whether a similar statement can be made for the much larger class of cache-oblivious algorithms that require a tall cache assumption. We answer their question in the affirmative.

Our main result is as follows. Consider a cache-oblivious algorithm that incurs $C(\tilde{M}, \tilde{B}, n)$ cache faults, when run on a machine with cache size \tilde{M} and block size \tilde{B} , provided that $\tilde{M} \geq g(\tilde{B})$. Here $g: \mathbb{N} \mapsto \mathbb{N}$ is a function that captures the “tallness” requirement on the cache [FLPR12]. We consider the execution of the algorithm on a VAT-machine with cache size \bar{M} and page size P and show that the number of cache faults is bounded by $4dC(\bar{M}/4, dB, n)$ provided that $\bar{M} \geq 4g(dB)$. Here $M = \bar{M}/a$, $B = P/a$ and $a \geq 1$ is the size (in addressable units) of the items handled by the algorithm.

Funnel sort [FLPR12] is an optimal cache-oblivious sorting algorithm. On an EM-machine with cache size \tilde{M} and block size \tilde{B} , it sorts n items with

$$C(\tilde{M}, \tilde{B}, n) = O\left(\frac{n}{\tilde{B}} \left\lceil \frac{\log n / \tilde{M}}{\log \tilde{M} / \tilde{B}} \right\rceil\right)$$

cache faults provided that $\tilde{M} \geq \tilde{B}^2$: thus g is quadratic¹ for this algorithm. As a consequence of our main theorem, we obtain:

Theorem 1 *Funnel sort sorts n items, each of size $a \geq 1$ addressable units, on a VAT-machine with cache*

¹This constraint can be reduced to $\tilde{M} \geq \tilde{B}^{1+\epsilon}$ for any constant $\epsilon > 0$ [BF02], however we do not wish to introduce additional notation.

size \bar{M} and page size P , with at most

$$O\left(\frac{4n}{B} \left\lceil \frac{\log 4n/M}{\log M/(4dB)} \right\rceil\right)$$

cache faults, where $M = \bar{M}/a$ and $B = P/a$. This assumes $(B \log_K(2n/P))^2 \leq M/4$.

Since $M/(4dB) \geq (M/B)^{1/2}$ for realistic values of M, B, K , and n , this implies funnel-sort is essentially optimal also in the VAT-model.

2 The EM-Model

An EM-machine has a fast memory of size \tilde{M} and a data is moved between fast and slow memory in blocks of size \tilde{B} . Algorithms for EM-machines may use \tilde{M} and \tilde{B} in the program code; algorithms are not written for specific values of \tilde{M} and \tilde{B} , but work for any values of \tilde{M} and \tilde{B} satisfying certain constraints, e.g., that the fast memory can hold a certain number of blocks. We capture these constraints by a function $g: \mathbb{N} \mapsto \mathbb{N}$ and the requirement $\tilde{M} \geq g(\tilde{B})$.

Cache-oblivious algorithms are algorithms that do not refer to the parameters \tilde{M} and \tilde{B} in the code. Only the analysis is done in terms of \tilde{M} and \tilde{B} . Frequently, the analysis only holds for \tilde{M} and \tilde{B} satisfying certain constraints; e.g., that the cache is tall and satisfies $\tilde{M} \geq \tilde{B}^2$. Again, this can be captured by an appropriate function g .

It is customary in the EM-literature that the size of the fast memory and the size of a block are expressed in terms of number of items handled by the algorithm. For example, for a sorting algorithm the items are the objects to be sorted.

3 The VAT-Model [JM13, Section XXX]

VAT-machines are EM-machines that use virtual addresses. We concentrate on the virtual memory of a single program. Both real (physical) and virtual addresses are strings in $\{0, K-1\}^d \{0, P-1\}$. Any such string corresponds to a number in the interval $[0, K^d P - 1]$ in a natural way. The $\{0, K-1\}^d$ part of the address is called an *index*, and its length d is an execution parameter fixed prior to the execution. We assume $d = \lceil \log_K(\text{last used address}/P) \rceil$. The $\{0, P-1\}$ part of the address is called page offset and P is the page size. A page contains P addressable units, usually bytes.² The translation process is a tree walk. We have a K -ary tree T of height d . The nodes of the tree are pairs (ℓ, i) with $\ell \geq 0$ and $i \geq 0$. We refer to ℓ as the layer of the node and to i as the number of the node. The leaves of the tree are on layer zero and a node (ℓ, i) on layer $\ell \geq 1$ has K children on layer $\ell - 1$, namely the nodes $(\ell - 1, Ki + a)$, for $a = 0 \dots K - 1$. In particular, node $(d, 0)$, the root, has children $(d - 1, 0), \dots, (d - 1, K - 1)$. The leaves of the tree correspond to physical pages of the main memory of a RAM machine. In order to translate a virtual address $x_{d-1} \dots x_0 y$, we start in the root of T , and then follow the path described by $x_{d-1} \dots x_0$. We refer to this path as the *translation path* for the address. The path ends in the leaf $(0, \sum_{0 \leq i \leq d-1} x_i K^i)$. Then the offset y selects the y -th cell in this page.

A VAT-machine has a fast memory (cache) of size \bar{M} and data is moved between fast and slow memory in units of P cells. We assume that a node of the translation tree fits into a page. Only nodes and data pages in the cache can be accessed directly, nodes and data pages currently in slow memory have to be brought into fast memory before they can be accessed. More precisely, let a be a virtual address, and let v_d, v_{d-1}, \dots, v_0 be its translation path. Here, v_d is the root of the translation tree, v_d to v_1 are internal nodes of the translation tree, and v_0 is a data page. We assume that the root node is always in cache. Translating a requires accessing all nodes of the translation path in order. Only nodes in the cache can be accessed. If a node needs to be accessed, but is not in the cache, it needs to be added to the cache, and some other

²In actual systems K is chosen such that a node fits exactly into a page. For example, for the 64-bit addressing mode of the processors of the AMD64 family (see <http://en.wikipedia.org/wiki/X86-64>), the addressable units are bytes and $P = 2^{12}$. Since an address consists of 2^3 bytes, $K = 2^9$.

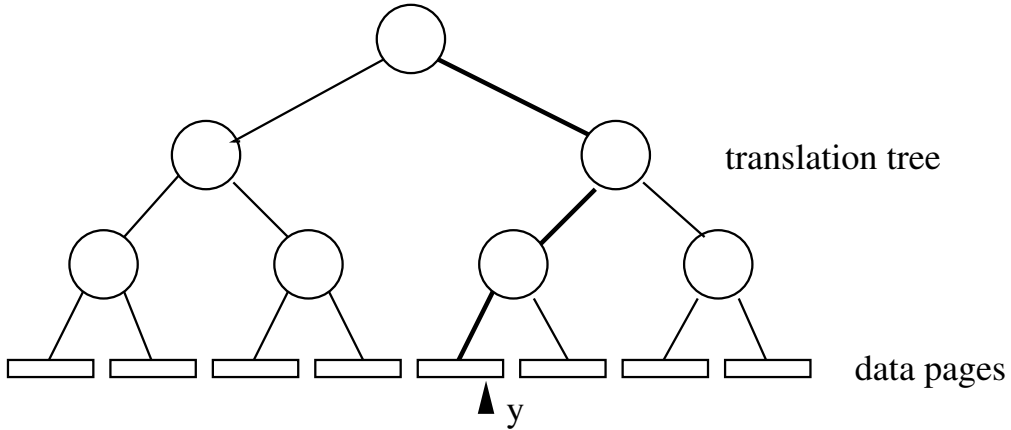


Figure 2: The pages holding the data are shown at the bottom and the translation tree is shown above the data pages. The translation tree has fan-out K and depth d ; here $K = 2$ and $d = 3$. The translation path for the virtual index 100 is shown. The offset y selects a cell in the physical page with virtual index 100. The nodes of the translation tree and the data pages are stored in memory. Only nodes and data pages in fast memory (cache memory) can be accessed directly, nodes and data pages currently in slow memory have to be brought into fast memory before they can be accessed. Each such move is a cache fault. In the EM-model only cache faults for data pages are counted, in the VAT-model, we count cache faults for nodes of the translation tree and for data pages.

page has to be evicted. The translation of a ends when v_0 is accessed. The cost of the memory access is the number of page faults incurred during the translation process.

EM- and VAT-machines move cells in contiguous blocks. If the items handled by an EM-machine comprise $a \geq 1$ addressable units, we have $M = \bar{M}/a$ and $B = P/a$. In the EM cost model only cache faults for data pages are charged and in the VAT cost model all cache faults arising in a memory access are charged.

4 EM-Algorithms as VAT-Algorithms

In the worst case, a memory access causes one cache fault in the EM-model and d cache faults in the VAT-model. In order to amortize the cost of a cache fault in the EM-model, it suffices to access a significant fraction of the memory cells in each page brought to fast memory. Therefore, an EM-algorithm that is efficient for block size dP/a should also be efficient in the VAT-model. The following discussion captures this intuition.

For an EM-algorithm, let $C(\tilde{M}, \tilde{B}, n)$ be the number of IO-operations on an input of size n , where \tilde{M} is the size of the faster memory (also called cache memory), \tilde{B} is the block size, and $\tilde{M} \geq g(\tilde{B})$. The following theorem shows that any EM-algorithm that is aware of \tilde{M} and \tilde{B} implies the existence of a VAT-algorithm that is aware of \bar{M} and P .

Theorem 2 *Let $g : \mathbb{N} \rightarrow \mathbb{N}$. Consider an EM-algorithm with IO-complexity $C(\tilde{M}, \tilde{B}, n)$, where \tilde{M} is the size of the cache, \tilde{B} is the size of a block, and n is the input size, provided that $\tilde{M} \geq g(\tilde{B})$. Let $d = \log_K(n/P)$ and assume $g(dP/a) \leq \bar{M}/(4a)$, where $a \geq 1$ is the item size in number of addressable units. The program can be made to run on a VAT-machine with cache size \bar{M} and page size P with at most $4dC(\bar{M}/(4a), dP/a, n)$ cache faults. With the notation $M = \bar{M}/a$ and $B = P/a$, the upper bound can be stated as $4dC(M/4, dB, n)$ and the tallness requirement becomes $M \geq 4g(dB)$.*

Proof: We run the EM-algorithm with a cache size of $\tilde{M} = \bar{M}/(4a)$ and a block size $\tilde{B} = dP/a$ and show how to execute it efficiently in a VAT-machine with page size P and cache size \bar{M} . We use one-fourth of

the cache for data and three-fourth for nodes of the translation tree. Since

$$\tilde{M} = \bar{M}/(4a) \geq g(dP/a) = g(\tilde{B}),$$

the number of cache faults incurred by the EM-algorithm is at most $C(\tilde{M}, \tilde{B}, n)$.

Whenever, the EM-algorithm moves a block containing \tilde{B} items to its data cache, the VAT-machine moves the corresponding d pages to the data cache and also moves all internal nodes of the translation paths to these d pages to then translation cache. The number of internal nodes is bounded by $2d + \sum_{i \geq 1} d/K^i \leq 2d + d/(K-1) \leq 3d$. Thus a translation cache of size $3\bar{M}/4$ suffices to store the translation paths to all pages in the data cache. For every cache fault of the EM-model, the VAT-machine incurs $4d$ cache faults. The theorem follows. \blacksquare

We apply the theorem to multi-way mergesort, an optimal sorting algorithm in the EM-model. It first creates n/\tilde{M} sorted runs of size \tilde{M} each by sorting chunks of size \tilde{M} in internal memory. It then performs multi-way merge sort on these runs. For the merge, it keeps one block from each input run, one block of the output run, and a heap containing the first elements of each run in fast memory. If $\tilde{M} \geq 5\tilde{B}$, we can merge two sequences since the space for the heap is certainly no more than the space for the input runs. The scheme results in a merge factor of $\Theta(\tilde{M}/\tilde{B})$. We assume for simplicity that the factor is exactly \tilde{M}/\tilde{B} . Thus

$$C(\tilde{M}, \tilde{B}, n) = \frac{n}{\tilde{B}} \left(1 + \left\lceil \frac{\log n/\tilde{M}}{\log \tilde{M}/\tilde{B}} \right\rceil \right).$$

By Theorem 2 and with a cache size of \bar{M} and page size P , the number of cache faults in the VAT-model is at most

$$4dC(\bar{M}/(4a), dP/a, n) = \frac{4n}{P/a} \left(1 + \left\lceil \frac{\log 4an/\bar{M}}{\log \bar{M}/(4dP)} \right\rceil \right) = \frac{4n}{B} \left(1 + \left\lceil \frac{\log 4n/M}{\log M/(4dB)} \right\rceil \right)$$

Here a is the item size in number of addressable units, $M = \bar{M}/a$ and $B = P/a$. This assumes $5dP \leq \bar{M}/4$. If $M/(4dB) \geq (M/B)^{1/2}$, the asymptotic number of cache faults is the same in both models. For realistic values of M , B , and n , this will be the case.

5 Cache-Oblivious VAT-Algorithms

We next extend the theorem to cache-oblivious algorithms. Unlike the previous theorem, the following theorem indicates that *any* algorithm, regardless of whether it is aware of the cache and block sizes, implies the existence of a VAT-algorithm that is similarly oblivious to \bar{M} and P .

Theorem 3 *Let $g : \mathbb{N} \mapsto \mathbb{N}$. Let A be an algorithm that incurs $C(\tilde{M}, \tilde{B}, n)$ cache faults in the EM-model with cache size \tilde{M} and block size \tilde{B} on an input of size n , provided that $\tilde{M} \geq g(\tilde{B})$. Let $d = \log_K(n/P)$ and assume $g(dP/a) \leq \bar{M}/(4a)$, where $a \geq 1$ is the item size in number of addressable units. On a VAT-machine with cache size \bar{M} and page size P and optimal use of the cache, the program incurs at most $4dC(\bar{M}/(4a), dP/a, n)$. With the notation $M = \bar{M}/a$ and $B = P/a$, the upper bound can be stated as $4dC(M/4, dB, n)$ and the tallness requirement becomes $M/4 \geq g(dB)$.*

Proof: We can almost literally reuse the proof of Theorem 2. The optimal execution of A on a VAT-machine with cache size \bar{M} and page size P cannot incur more cache faults than the particular execution that we describe next.

Consider an execution of A on an EM-machine with cache size $\tilde{M} = \bar{M}/(4a)$ and block size $\tilde{B} = dP/a$. Since

$$\tilde{M} = \bar{M}/(4a) \geq g(dP/a) = g(\tilde{B}),$$

the number of cache faults incurred by the EM-algorithm is at most $C(\tilde{M}, \tilde{B}, n)$. We execute the program on the VAT-machine as in the proof of Theorem 2.

We use one-fourth of the cache for data and three-fourth for nodes of the translation tree. Whenever, the EM-machine moves a block (of size \tilde{B} items) to its data cache, the VAT-machine moves the corresponding d pages to the data cache and also moves all internal nodes of the translation paths to these d pages to then translation cache. The number of internal nodes is bounded by $2d + \sum_{i \geq 1} d/K^i \leq 2d + d/(K-1) \leq 3d$. Thus a translation cache of size $3\tilde{M}/4$ suffices to store the translation paths to all pages in the data cache. For every cache fault of the EM-machine, the VAT-machine incurs $4d$ cache faults. The theorem follows. ■

The optimal cache replacement strategy may be replaced by LRU at the cost of doubling the cache size and doubling the number of cache faults [JM13].

Recall that a cache-oblivious algorithm has no knowledge of the memory and the block size. It does well for any choice of \tilde{M} and \tilde{B} as long as $\tilde{M} \geq g(\tilde{B})$. The theorem tells us that it also does well in the VAT-model as long as the more stringent requirement $M \geq 4ag(dB)$ is satisfied.

6 Conclusion

We have shown that performance of cache-oblivious algorithms in the VAT-model matches their performance in the EM-model provided a somewhat more stringent tall cache assumption holds.

References

- [BF02] Gerth Stølting Brodal and Rolf Fagerberg. Cache oblivious distribution sweeping. In *ICALP*, volume 2380 of *LNCS*, pages 426–438, 2002.
- [FLPR12] M. Frigo, C.E. Leiserson, H. Prokop, and S. Ramachandran. Cache-oblivious algorithms. *ACM Transactions on Algorithms*, pages 4:1 – 4:22, 2012. a preliminary version appeared in FOCS 1999.
- [JM13] Tomasz Jurkiewicz and Kurt Mehlhorn. The Cost of Address Translation. In *ALLENEX*, pages 148–162, 2013. Full version available at <http://arxiv.org/abs/1212.0703>.