

On a Model of Virtual Address Translation¹

Tomasz Jurkiewicz, Max Planck Institute for Informatics, Saarbrücken, Germany and the Saarbrücken Graduate School of Computer Science
 Kurt Mehlhorn, Max Planck Institute for Informatics, Saarbrücken, Germany

Modern computers are not random access machines (RAMs). They have a memory hierarchy, multiple cores, and a virtual memory. We address the computational cost of the address translation in the virtual memory.

Starting point for our work on virtual memory is the observation that the analysis of some simple algorithms (random scan of an array, binary search, heapsort) in either the RAM model or the EM model (external memory model) does not correctly predict growth rates of actual running times. We propose the VAT model (virtual address translation) to account for the cost of address translations and analyze the algorithms mentioned above and others in the model. The predictions agree with the measurements. We also analyze the VAT-cost of cache-oblivious algorithms.

General Terms: Design, Algorithms, Performance

Additional Key Words and Phrases: External memory, virtual address translation

ACM Reference Format:

Tomasz Jurkiewicz and Kurt Mehlhorn, 2013. Computational Complexity of the Virtual Address Translation. *ACM J. Exp. Algor.* ??, ??, Article ?? (April ??), 29 pages.
 DOI : <http://dx.doi.org/10.1145/0000000.0000000>

1. INTRODUCTION

The role of models of computation in algorithmics is to provide abstractions of real machines for algorithm analysis. Models should be mathematically pleasing and have a predictive value. Both aspects are essential. If the analysis has no predictive value, it is merely a mathematical exercise. If a model is not clean and simple, researchers will not use it. The standard models for algorithm analysis are the RAM (random access machine) model [Shepherdson and Sturgis 1963] and the EM (external memory) model [Aggarwal and Vitter 1988].

The RAM model is by far the most popular model. It is an abstraction of the von Neumann architecture. A computer consists of a control and processing unit and an unbounded memory. Each memory cell can hold a word, and memory access as well as logical and arithmetic operations on words take constant time. The word length is either an explicit parameter or assumed to be logarithmic in the size of the input. The model is very simple and has a predictive value.

Modern machines have virtual memory, multiple processor cores, an extensive memory hierarchy involving several levels of cache memory, main memory, and disks. The external memory model was introduced because the RAM model does not account for

A preliminary version of this article appeared in ALENEX 2013. The article is based on the first author's PhD-thesis.

Author's mail and web addresses: {tojot, mehlhorn}@mpi-inf.mpg.de,

<http://www.mpi-inf.mpg.de/~{tojot, mehlhorn}>

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.

© ?? ACM 1084-6654/??/04-ART?? \$15.00

DOI : <http://dx.doi.org/10.1145/0000000.0000000>

the memory hierarchy, and hence, the RAM model has no predictive value for computations involving disks. We give more details on both models in Section 2.

This research started with a simple experiment. We timed six simple programs for different input sizes, namely, permuting the elements of an array of size n , random scan of an array of size n , n random binary searches in an array of size n , heapsort of n elements, introsort² of n elements, and sequential scan of an array of size n . For some of the programs, e.g., sequential scan through an array and quicksort, the measured running times agree very well with the predictions of the models. *However, the running time of random scan seems to grow as $O(n \log n)$, and the running time of the binary searches seems to grow as $O(n \log^2 n)$, a blatant violation of what either model predicts.* We give the details of the experiments in Section 3.

Why do measured and predicted running times differ? Modern computers have virtual memories. Each process has its own virtual address space $\{0, 1, 2, \dots\}$. Whenever a process accesses memory, the virtual address has to be translated into a physical address. *The translation of virtual addresses into physical addresses incurs cost.* The translation process is usually implemented as a hardware-supported walk in a prefix tree, see Section 4 for details. The tree is stored in the memory hierarchy, and hence, the translation process may incur cache faults. The number of cache faults depends on the locality of memory accesses: the less local, the more cache faults. The depth of the translation tree is logarithmic in the size of an algorithm's address space and hence, in the worst case, every memory access may lead to a logarithmic number of cache faults during the translation process. For random scan and random binary searches, it apparently does.

We propose an extension of the EM model, the VAT (Virtual Address Translation) model, that accounts for the cost of address translation, see Section 5. We show that we may assume that the translation process makes optimal use of the cache memory by relating the cost of optimal use with the cost under the LRU strategy, see Section 5. We analyze a number of programs, including the six mentioned above, in the VAT model and obtain good agreement with the measured running times, see Section 6. We relate the cost of a cache-oblivious algorithm in the EM model to the cost in the VAT model, see Section 7. In particular, cache-oblivious algorithms that do not need a tall-cache assumption incur no or little overhead. In Section 8, we address comments made by reviewers and readers of the paper. We close with some suggestions for further research and consequences for teaching, see Section 9.

Related Work: It is well-known in the architecture and systems community that virtual memory and address translation comes at a cost. Many textbooks on computer organization, e.g., [Hennessy and Patterson 2007], discuss virtual memories. The papers by Drepper [Drepper 2007; 2008] describe computer memories, including virtual translation, in great detail. [Advanced Micro Devices 2010] provides further implementation details.

The cost of address translation has received little attention from the algorithms community. The survey paper by N. Rahman [Rahman 2003] on algorithms for hardware caches and TLB summarizes the work on the subject. She discusses a number of theoretical models for memory. All models discussed in [Rahman 2003] treat address translation atomically, i.e., the translation from virtual to physical addresses is a single operation. However, this is no longer true. In 64-bit systems, the translation process is a tree walk. Our paper is the first that proposes a theoretical model for address translation and analyzes algorithms in this model.

²Introsort is the version of quicksort used in modern versions of the STL. For the purpose of this paper, introsort is a synonym for quicksort.

2. THE RANDOM ACCESS MACHINE AND THE EXTERNAL MEMORY MACHINE

A RAM machine consists of a central processing unit and a memory. The memory consists of cells indexed by nonnegative integers. A cell can hold a bitstring. The CPU has a finite number of registers, in particular an accumulator and an address register. In any one step, a RAM can either perform an operation (simple arithmetic or boolean operations) on its registers or access memory. In a memory access, the content of the memory cell indexed by the content of the address register is either loaded into the accumulator or written from the accumulator. Two timing models are used: in the unit-cost RAM, each operation has cost one, and the length of the bitstrings that can be stored in memory cells and registers is bounded by the logarithm of the size of the input; in the logarithmic-cost RAM, the cost of an operation is equal to the sum of the lengths (in bits) of the operands, and the contents of memory cells and registers are unrestricted.

An EM machine is a RAM with two levels of memory. The levels are referred to as cache and main memory or memory and disk, respectively. We use the terms cache and main memory. The CPU can only operate on data in the cache. Cache and main memory are each divided into blocks of B cells, and data is transported between cache and main memory in blocks. The cache has size M and hence consists of M/B blocks; the main memory is infinite in size. The analysis of algorithms in the EM-model bounds the number of CPU-steps and the number of block transfers. The time required for a block transfer is equal to the time required by $\Theta(B)$ CPU-steps. The hidden constant factor is fairly large, and therefore, the emphasis of the analysis is usually on the number of block transfers.

3. SOME PUZZLING EXPERIMENTS

We used the following seven programs in our experiments. Let A be an array of size n .

- permute: for $j \in [n - 1..0]$ do: $i := \text{random}(0..j)$; $\text{swap}(A[i], A[j])$;
- random scan: $\pi := \text{random permutation}$; for i from 0 to $n - 1$ do: $S := S + A[\pi(i)]$;
- n binary searches for random positions in A ; A is sorted for this experiment.
- heapify
- heapsort
- quicksort
- sequential scan

On a RAM, the first two, the last, and heapify are linear time $\Theta(n)$, and the others are $\Theta(n \log n)$. Figure 1 shows the measured running times for these programs divided by their RAM complexity; we refer to this quantity as *normalized operation time*. More details about our experimental methodology are available in Subsection 3.2. If RAM complexity is a good predictor, the normalized operation times should be approximately constant. We observe that two of the linear time programs show linear behavior, namely, sequential access and heapify, that one of the $\Theta(n \log n)$ programs shows $\Theta(n \log n)$ behavior, namely, quicksort, and that for the other programs (heapsort, repeated binary search, permute, random access), the actual running time grows faster than what the RAM model predicts.

How much faster and why?

Figure 1 also answers the “how much faster” part of the question. Normalized operation time seems to be a piecewise linear in the logarithm of the problem size; observe that we are using a logarithmic scale for the abscissa in this figure. For heapsort and repeated binary search, normalized operation time is almost perfectly piecewise linear, for permute and random scan, the piecewise linear should be taken with a grain of

salt.³ The pieces correspond to the memory hierarchy. *The measurements suggest that the running times of permute and random scan grow like $\Theta(n \log n)$ and the running times of heapsort and repeated binary search grow like $\Theta(n \log^2 n)$.*

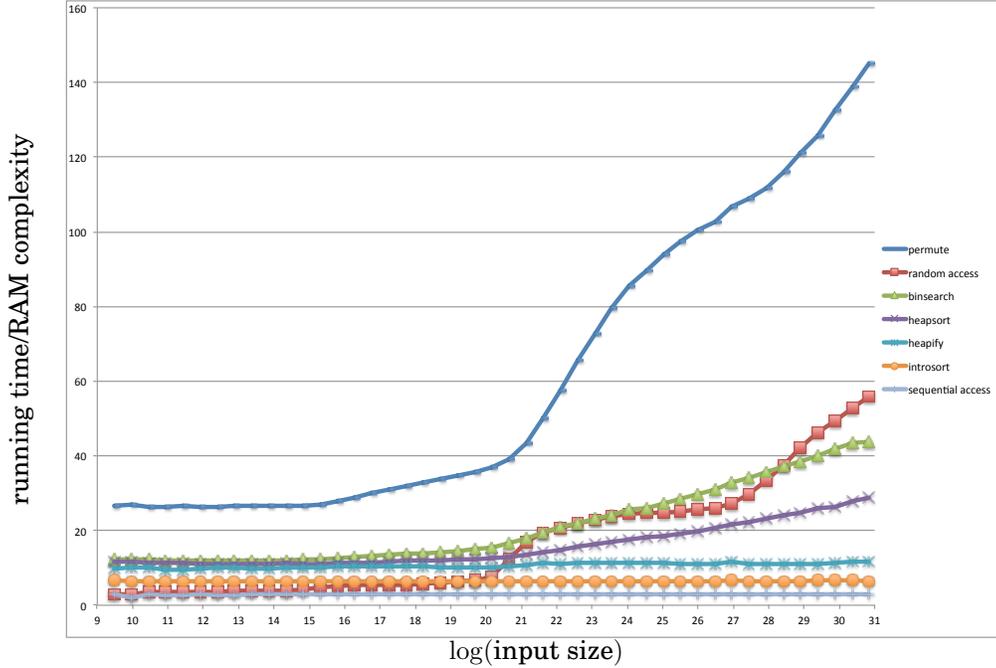


Fig. 1. The abscissa shows the logarithm of the input size. The ordinate shows the measured running time divided by the RAM-complexity (normalized operation time). The normalized operation times of sequential access, quicksort, and heapify are constant, the normalized operation times of the other programs are not.

3.1. Memory Hierarchy Does Not Explain It

We argue in this section that the memory hierarchy does not explain the experimental findings. We give a detailed analysis of the cost of a random scan of an array of size n in a hierarchical memory and relate it to the measured running time. We will see that the prediction by the model and the measured running times differ widely. A simpler argument for a one-level memory hierarchy will be given in Section 5.1.

Let $s_i, i \geq 0$ be the size of the i -th level C_i of the memory hierarchy; $s_{-1} = 0$. We assume $C_i \subset C_{i+1}$ for all i . Let ℓ be such that $s_\ell < n \leq s_{\ell+1}$, i.e., the array fits into level $\ell + 1$ but does not fit into level ℓ . For $i \leq \ell$, a random address is in C_i but not in C_{i-1} , with probability $(s_i - s_{i-1})/n$. Let c_i be the cost of accessing an address that is in C_i but not in C_{i-1} . The expected total cost in the external memory model is equal to

$$T_{EM}(n) := n \cdot \left(\frac{n - s_\ell}{n} c_{\ell+1} + \sum_{0 \leq i \leq \ell} \frac{s_i - s_{i-1}}{n} c_i \right) = n c_{\ell+1} - \sum_{0 \leq i \leq \ell} s_i (c_{i+1} - c_i).$$

³We are leaving it as an open problem to give a satisfactory explanation for the bumpy shape of the graphs for permute and random access.

This is a piecewise linear function whose slope is $c_{\ell+1}$ for $s_\ell < n \leq s_{\ell+1}$. The slopes are increasing but change only when a new level of the memory hierarchy is used. Figure 2 shows the measured running time of random scan divided by EM-complexity as a function of the logarithm of the problem size. Clearly, the figure does not show the graph of a constant function.⁴

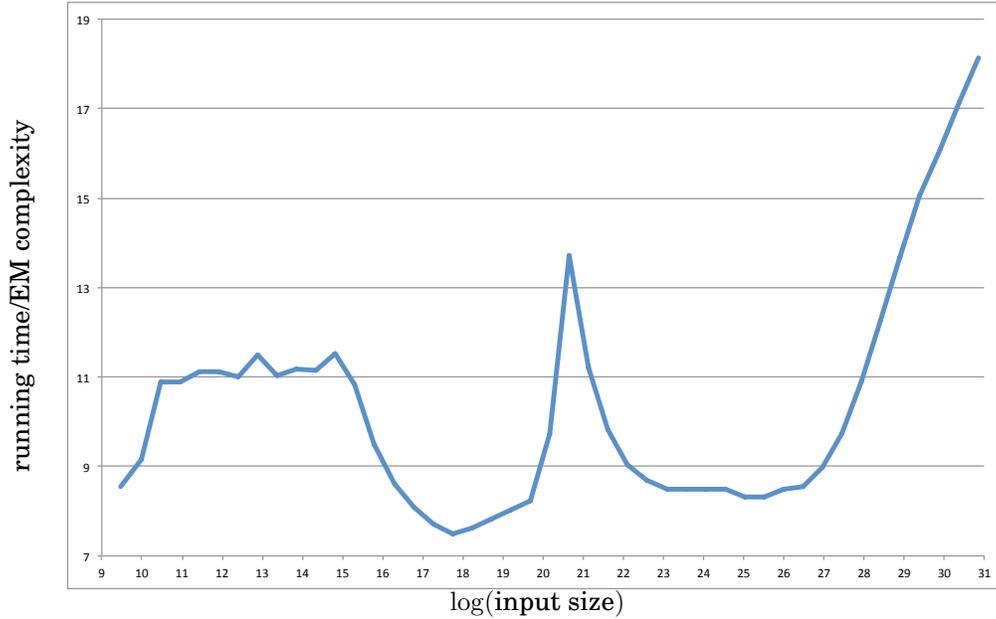


Fig. 2. The running time of random scan divided by the EM-complexity. We used the following parameters for the memory hierarchy: the sizes are taken from the machine specification, and the access times were determined experimentally.

Memory Level	Size	log(maximum number of elements)	Access Time in Picoseconds
L1	32kiB	12	4080
L2	256kiB	15	4575
L3	12MiB	20,58	9937
RAM			38746

3.2. Methodology

Programs used for the preparation of Figure 1 were compiled by gcc in version “Debian 4.4.5-8”, and run on Debian Linux in version 6.0.3, on a machine with an Intel Xeon X5690 processor (3,46 GHz, 12MiB⁵ Smart Cache, 6,4 GT/s QPI). The caption of Figure 2 lists further machine parameters. In each case, we performed multiple repetitions and took the minimum measurement for each considered size of the input data. We chose the minimum because we are estimating the cost that must be incurred. We

⁴The semi-log plot of a function of the form $n \mapsto (n \log(n/a))/(bn - c)$ with $a, b, c > 0$ is convex. Note that $T_{EM}(n) = bn - c$, where b and c depend on the level of the memory hierarchy required for n and that Figure 1 suggests that the actual running time grows like $n \log(n/a)$. The plot may be interpreted as the plot of a piecewise convex function and hence does not contradict the conclusion drawn from Figure 1.

⁵KiB and MiB are modern, non-ambiguous notations for 2^{10} and 2^{20} bytes, respectively. For more details, refer to http://en.wikipedia.org/wiki/Binary_prefix.

also experimented with average or median; moreover, we performed the experiments on other machines and operating systems and obtained consistent results in each case. We grew input sizes by factors of 1.4 to exclude the influence of memory associativity, and we made sure that the largest problem size still fitted in the main memory to eliminate swapping.

For each experiment, we computed its normalized operation time, which we define as the measured execution time divided by the RAM complexity. This way, we eliminate the known factors. The resulting function represents cost of a single RAM-operation in relation to the problem size.

We use semi-log plots for showing normalized operation cost as a function of the logarithm in the input size. In such a plot, linear functions of the logarithm of the input size are easily identified as straight lines.

4. VIRTUAL MEMORY

Virtual addressing was motivated by multi-processing. When several processes are executed concurrently on the same machine, it is convenient and more secure to give each program a linear address space indexed by the nonnegative integers. However, these addresses are now virtual and no longer directly correspond to physical (real) addresses. Rather, it is the task of the operating system to map the virtual addresses of all processes to a single physical memory. The mapping process is hardware supported.

The memory is viewed as a collection of pages of $P = 2^p$ cells (= addressable units). Both virtual and real addresses consist of an *index* and an *offset*. The index selects a page and the offset selects a cell in a page. The index is broken into d segments of length $k = \log K$. For example, for the long addressing mode of the processors of the AMD64 family (see <http://en.wikipedia.org/wiki/X86-64>) the numbers are: $d = 4$, $k = 9$, and $p = 12$; the remaining 16 bits are used for other purposes. The choice of k is not arbitrary. A page consists of 2^{12} bytes. An address consists of 8 bytes and hence a node of the translation tree requires $2^9 \cdot 2^3 = 2^{12}$ bytes. Thus nodes fit exactly into pages.

Logically, the translation process is a walk in a tree with outdegree K ; this tree is usually called the page table [Drepper 2008; Hennessy and Patterson 2007]. The walk starts at the root; the first segment of the index determines the child of the root, the second segment of the index determines the child of the child, and so on. The leaves of the tree store indices of physical pages. The offset then determines the cell in the physical address, i.e., offsets are not translated but taken verbatim. Here quoting [Advanced Micro Devices 2010]:

“Virtual addresses are translated to physical addresses through hierarchical translation tables created and managed by system software. Each table contains a set of entries that point to the next-lower table in the translation hierarchy. A single table at one level of the hierarchy can have hundreds of entries, each of which points to a unique table at the next-lower hierarchical level. Each lower-level table can in turn have hundreds of entries pointing to tables further down the hierarchy. The lowest-level table in the hierarchy points to the translated physical page.

Figure 3 on page 7 shows an overview of the page-translation hierarchy used in long mode. Legacy mode paging uses a subset of this translation hierarchy. As this figure shows, a virtual address is divided into fields, each of which is used as an offset into a translation table. The complete translation chain is made up of all table entries referenced by the virtual-address fields. The lowest-order virtual-address bits are used as the byte offset into the physical page.”

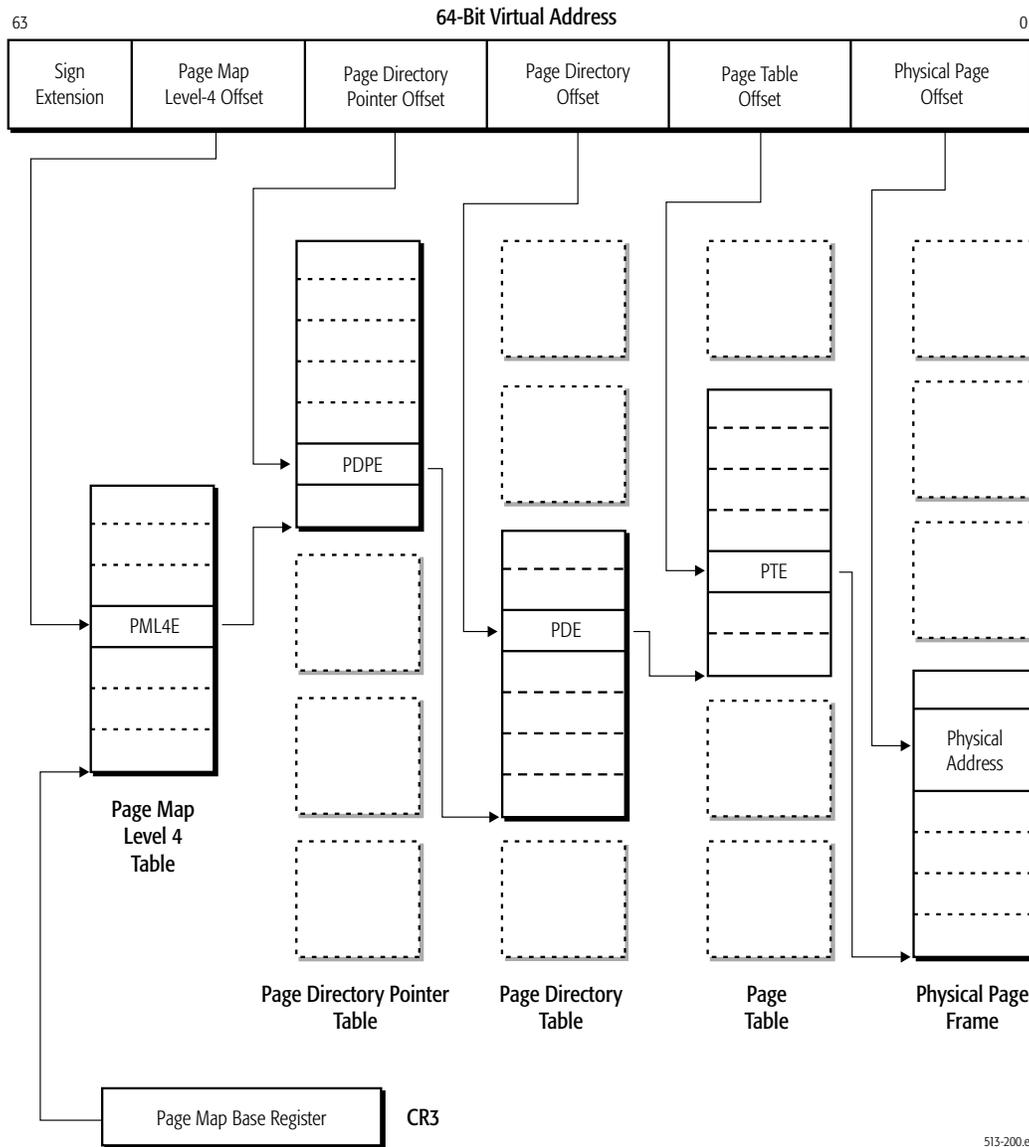


Fig. 3. Virtual to Physical Address Translation in AMD64, figure from [Advanced Micro Devices 2010]. Note that levels of the prefix tree have distinct historical names, as this system was originally not designed to have multiple levels (Page Map Level 4 Table; Page Directory Pointer Table; Page Directory Table; and Page Table).

Due to its size, the page table is stored in the RAM, but nodes accessed during the page table walk have to be brought to faster memory. A small number of recent translations is stored in the translation-lookaside-buffer (TLB). The TLB is a small associative memory that contains physical indices indexed by the virtual ones. This is akin to the first level cache for data. Quoting [Advanced Micro Devices 2010] further:

“Every memory access has its virtual address automatically translated into a physical address using the page-translation hierarchy. *Translation-lookaside buffers* (TLBs), also known as *page-translation caches*, nearly eliminate the performance penalty associated with page translation. TLBs are special on-chip caches that hold the most-recently used virtual-to-physical address translations. Each memory reference (instruction and data) is checked by the TLB. If the translation is present in the TLB, it is immediately provided to the processor, thus avoiding external memory references for accessing page tables.

TLBs take advantage of the *principle of locality*. That is, if a memory address is referenced, it is likely that nearby memory addresses will be referenced in the near future.”

5. VAT, THE VIRTUAL ADDRESS TRANSLATION MODEL

P	page size ($P = 2^p$)
K	arity of translation tree ($K = 2^k$)
d	depth of translation tree ($= \lceil \log_K(\text{max used virtual address}) \rceil$)
W	size of TC cache
τ	cost of a cache fault (number of RAM instructions)

Fig. 4. Notation

VAT machines are RAM machines that use virtual addresses. We concentrate on the virtual memory of a single program. Both real (physical) and virtual addresses are strings in $\{0, K - 1\}^d \{0, \dots, P - 1\}$. 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. It is assumed that $d = \lceil \log_K(\text{maximum used virtual address}/P) \rceil$. The $\{0, \dots, P - 1\}$ part of the address is called *page offset* and P is the page size. 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 are physical pages of the main memory of a RAM machine. In order to translate 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 $(0, \sum_{0 \leq i \leq d-1} x_i K^i)$. The offset y selects the y -th cell in this page.

The translation process uses a translation cache TC that can store W nodes of the translation tree. Note that a node is either an internal node and then points to K other nodes or is a leaf and then stores data. The cache is updated by insertions and evictions. 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. Translating a requires accessing all nodes of the translation path in order. Only nodes in the TC can be accessed. The translation of a ends when v_0 is accessed.

The *number of cache faults incurred by the memory access* is the number of insertions performed during the translation process, and the *cost of the memory access* is τ times the number of cache faults. The number of cache faults is at least the number of nodes of the translation path that are not present in the cache at the beginning of the translation. Figure 4 summarizes the notation.

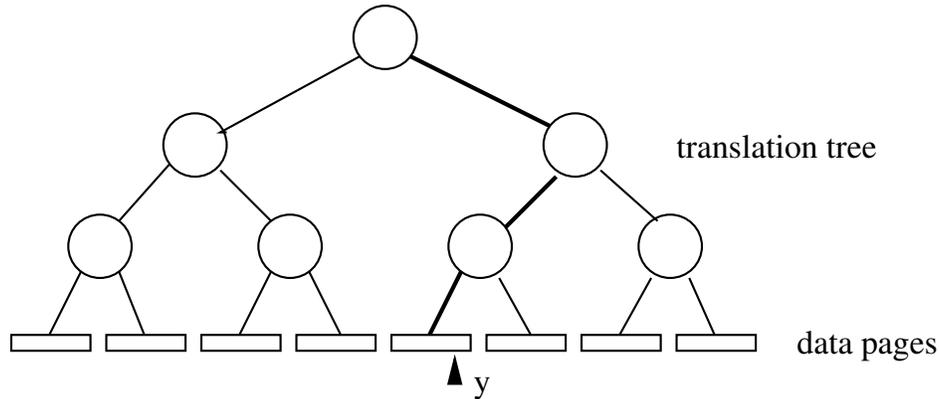


Fig. 5. The pages holding the data are shown at the bottom and the translation tree is shown above the data pages. The translation tree has fan-out K and depth d ; here $K = 2$ and $d = 3$. The translation path for the virtual index 100 is shown. The offset y selects a cell in the physical page with 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 all nodes of the translation tree.

We close the introduction of the model with a trivial, but useful observation.

LEMMA 5.1. *Let $D \geq 0$ and $i \geq 0$ be integers. The translation paths for addresses i and $i + D$ differ in at least the last $\max(0, \lceil \log_K(D/P) \rceil)$ nodes.*

PROOF. Let u_d, \dots, u_0 and v_d, \dots, v_0 be the translation paths for addresses i and $i + D$, respectively. If $u_\ell = v_\ell$, then $D \leq K^\ell P$. As a consequence, if $D > K^\ell P$, the translation paths differ in the last $\ell + 1$ nodes. \square

5.1. Relation to EM-Model

In the EM-model, one counts only cache faults caused by data pages, in the VAT-model one counts cache faults caused by all translation tree nodes.

A comparison of jumping scan and random scan illustrates the difference; this example is due to Jirka Katajainen (personal communication). A jumping scan (with stride P) is a sequential scan of an array with stride P , i.e., cells $0, P, 2P, 3P, \dots$ are accessed in order. In the EM-model each access causes one page fault and hence the EM-cost of jumping scan and random scan are identical.

In the VAT-model, a random scan is much more costly than a jumping scan. By Lemma 6.1, the cost of a random scan of an array of size n is at least $\tau n \log_K(n/(PW))$. However, the cost of a jumping scan is at most $\tau n(1 + 1/(K - 1))$. Observe (see Figure 5) that K subsequent data pages share the last internal of the translation path and hence the number of page faults caused by nodes of the translation tree is bounded by $\sum_{i \geq 1} n/K^i = n/K \cdot K/(K - 1) = n/(K - 1)$.

A second comparison is also informative. Assume that pages are accessed in random order and once a page is accessed all cells of the page are accessed. With respect to cache faults, n such accesses are equivalent to n/P random accesses. In the EM-model the number of cache faults is n/P , which is the same as for a linear scan, and in the VAT-model the number is at least $(n/P) \log_K(n/(P^2W))$.

5.2. TC Replacement Strategies

Since the TC is a special case of a cache in a classic EM machine, the following classic result applies.

LEMMA 5.2 ([SLEATOR AND TARJAN 1985; FRIGO ET AL. 2012]). *An optimal replacement strategy is at most by factor 2 better than LRU⁶ on a cache of double size, assuming both caches start empty.*

This result is useful for upper bounds and lower bounds. LRU is easy to implement. In upper bound arguments, we may use any replacement strategy and then appeal to the Lemma. In lower bound arguments, we may assume the use of LRU. For TC caches, it is natural to assume the initial segment property.

Definition 5.3. An **initial segment** of a rooted tree is an empty tree or a connected subgraph of the tree containing the root. A TC has the **initial segment property (ISP)** if the TC contains an initial segment of the translation tree. A TC replacement strategy has ISP if under this strategy a TC has ISP at all times.

PROPOSITION 5.4. *Strategies with ISP exist only for TCs with $W > d$.*

ISP is important because, as we show later, ISP can be realized at no additional cost for LRU and at little additional cost for the optimal replacement strategy. Therefore, strategies with ISP can significantly simplify proofs for upper and lower bounds. Moreover, strategies with ISP are easier to implement. Any implementation of a caching system requires some way to search the cache. This requires an indexing mechanism. RAM memory is indexed by the memory translation tree. In the case of the TC itself, ISP allows to integrate the indexing structure into the cached content. One only has to store the root of the tree at a fixed position or store the location of the root in a special register, called the Page Map Base Register in Figure 3.

We establish the following relations in this section:

$$\text{MIN}(W) \leq \text{ISMIN}(W) \leq \text{ISLRU}(W) \leq \text{LRU}(W) \leq 2\text{MIN}(W/2)$$

$$\text{MIN}(W) \geq \text{ISMIN}(W + d) \quad \text{ISLRU}(W) \geq \text{LRU}(W + d);$$

here W is the size of the translation cache, d is the index length, LRU is the least-recently-used replacement strategy, MIN is the optimal cache replacement strategy, ISLRU is LRU with the ISP-property, and ISMIN is the optimal replacement strategy with the ISP-property.

5.3. Eager Strategies and the Initial Segment Property

Before we prove an ISP analogue of Lemma 5.2, we need to better understand the behavior of replacement strategies with ISP. For classic caches, premature evictions and insertions do not improve efficiency. We will show that the same holds true for TCs with ISP. This will be useful because we will use early evictions and insertions in some of our arguments.

Definition 5.5. A replacement strategy is **lazy** if it performs an insertion of a missing node, only if the node is accessed directly after this insertion, and performs an eviction only before an insertion for which there would be no free cell otherwise. In the opposite case the strategy is **eager**. Unless stated otherwise, we assume that a strategy being discussed is lazy.

⁶LRU is a strategy that always evicts the Least Recently Used node.

Eager strategies can perform replacements before they are needed and can even insert nodes that are not needed at all. Also, they can insert and re-evict, or evict and re-insert nodes during a single translation. We eliminate this behavior *translation by translation* as follows. Consider a fixed translation and define the sets of **effective evictions and insertions** as follows.

$$EE = \{\text{evict}(a) : \text{there are more } \text{evict}(a) \text{ than } \text{insert}(a) \text{ in the translation.}\}$$

$$EI = \{\text{insert}(a) : \text{there are more } \text{insert}(a) \text{ than } \text{evict}(a) \text{ in the translation.}\}$$

Please note that in this case “there are more” means “there is *one* more” as there cannot be two $\text{evict}(a)$ without an $\text{insert}(a)$ between them, or two $\text{insert}(a)$ without $\text{evict}(a)$.

PROPOSITION 5.6. *The effective evictions and insertions modify the content of the TC in the same way as the original evictions and insertions.*

PROPOSITION 5.7. *During a single translation while a strategy with ISP is in use:*

- (1) *No node from the current translation path is effectively evicted, and all the nodes missing from the current translation path are effectively inserted.*
- (2) *If a node is effectively inserted, no ancestor or descendant of it is effectively deleted. Subject to obeying the size restriction of the TC, we may therefore reorder effective insertions and effective deletions with respect to each other (but not changing the order of the insertions and not changing the order of the evictions).*

LEMMA 5.8. *Any eager replacement strategy with ISP can be transformed into a lazy replacement strategy with ISP with no efficiency loss.*

PROOF. We modify the original evict/insert/access sequence *translation by translation*. Consider the current translation and let EI and EE be the set of effective insertions and evictions. We insert the missing nodes from the current translation path exactly at the moment they are needed. Whenever this implies an insertion into a full cache, we perform one of the lowest effective evictions, where lowest means that no children of the node are in the TC. There must be such an effective eviction because, otherwise, the original sequence would overuse the cache as well. When all nodes of the current translation path are accessed, we schedule all remaining effective evictions and insertions at the beginning of the next translation; first the evictions in descendant-first order and then the insertions in ancestor-first order. The modified sequence is operationally equivalent to the original one, performs no more insertions, and does not exceed cache size. Moreover, the current translation is now lazy. \square

5.4. ISLRU, or LRU with the Initial Segment Property

Even without ISP, LRU has the property below.

LEMMA 5.9. *When the LRU policy is in use, the number of TC misses in a translation is equal to the layer number of the highest missing node on the translation path.*

PROOF. The content of the LRU cache is easy to describe. Concatenate all translation paths and delete all occurrences of each node except the last. The last W nodes of the resulting sequence form the TC. Observe that an occurrence of a node is only deleted if the node is part of a latter translation path. This implies that the TC contains at most two incomplete translation paths, namely, the least recent path that still has nodes in the TC and the current path. The former path is evicted top-down and the latter path is inserted top-down. The claim now easily follows. Let v be the highest missing node on the current translation path. If no descendant of v is contained in the TC, the claim is obvious. Otherwise, the topmost descendant present in the TC is the

first node on the part of the least recent path that is still in the TC. Thus, as the current translation path is loaded into the TC, the least recent path is evicted top-down. Consequently, the gap is never reduced. \square

The proof above also shows that whenever LRU detaches nodes from the initial segment, the detached nodes will never be used again. This suggests a simple (implementable) way of introducing ISP to LRU. If LRU evicts a node that still has descendants in the TC, it also evicts the descendants. The descendants actually form a single path. Next, we use Lemma 5.8 to make this algorithm lazy again. It is easy to see that the resulting algorithm is the ISLRU, as defined next.

Definition 5.10. **ISLRU** (Initial Segment preserving LRU) is the replacement strategy that always evicts the lowest descendant of the least recently used node.

Due to the construction and Lemma 5.8, we have the following.

PROPOSITION 5.11. *ISLRU for TCs with $W > d$ is at least as good as LRU.*

Remark 5.12. In fact, the proposition holds also for $W \leq d$, even though ISLRU no longer has ISP in this case.

5.5. ISMIN: The Optimal Strategy with the Initial Segment Property

Definition 5.13. **ISMIN** (Initial Segment property preserving MIN) is the replacement strategy for TCs with ISP that always evicts from a TC the node that is not used for the longest time into the future among the nodes that are not on the current translation path and have no descendants. Nodes that will never be used again are evicted before the others in arbitrary descendant–first order.

THEOREM 5.14. *ISMIN is an optimal replacement strategy among those with ISP.*

PROOF. Let R be any replacement strategy with ISP, and let t be the first point in time when it departs from ISMIN. We will construct R' with ISP that does not depart from ISMIN, including time t , and has no more TC misses than R . Let v be the node evicted by ISMIN at time t .

We first assume that R evicts v at some later time t' without accessing it in the interval $(t, t']$. Then R' simply evicts v at time t and shifts the other evictions in the interval $[t, t')$ to one later replacement. Postponing evictions to the next replacement does not cause additional insertions and does not break connectivity. It may destroy laziness by moving an eviction of a node directly before its insertion. In this case, R' skips both. Since no descendant of v is in the TC at time t , and v will not be used for the longest time into the future, none of its children will be added by R before time t' ; therefore, the change does not break the connectivity.

We come to the case that R stores v until it is accessed for the next time, say at time t' . Let a be the node evicted by R at time t . R' evicts v instead of a and remembers a as being **special**. We guarantee that the content of the TCs in the strategies R and R' differs only by v and the current special node until time t' and is identical afterwards. To reach this goal, R' replicates the behavior of R except for three situations.

- (1) If R evicts the parent of the special node, R' evicts the special node to preserve ISP and from then on remembers the parent as being special. As long as only Rule 1 is applied, the special node is an ancestor of a .
- (2) If R replaces some node b with the current special node, R' skips the replacement and from then on remembers b as the special node. Since a will be accessed before v , Rule 2 is guaranteed to be applied, and hence, R' is guaranteed to save at least one replacement.
- (3) At time t' , R' replaces the special node with v , performing one extra replacement.

We have shown how to turn an arbitrary replacement strategy with ISP into ISMIN without efficiency loss. This proves the optimality of ISMIN. \square

We can now state an ISP-aware extension of Lemma 5.2.

THEOREM 5.15.

$$\text{MIN}(W) \leq \text{ISMIN}(W) \leq \text{ISLRU}(W) \leq \text{LRU}(W) \leq 2\text{MIN}(W/2),$$

where MIN is an optimal replacement strategy and $A(s)$ denotes a number of insertions performed by replacement strategy A to an initially empty TC of size $s > d$ for an arbitrary but fixed sequence of translations.

PROOF. MIN is an optimal replacement strategy, so it is better than ISMIN. ISMIN is an optimal replacement strategy among those with ISP, so it is better than ISLRU. ISLRU is better than LRU by Proposition 5.11. $\text{LRU}(W) < 2\text{MIN}(W/2)$ holds by Lemma 5.2. \square

5.6. Tighter Relationships

Theorem 5.15 implies $\text{LRU}(W) \leq 2\text{ISLRU}(W/2)$ and $\text{ISMIN}(W) \leq 2\text{MIN}(W/2)$. In this section, we sharpen both inequalities.

LEMMA 5.16. $\text{LRU}(W + d) \leq \text{ISLRU}(W)$.

PROOF. d nodes are sufficient for LRU to store one extra path, hence, from the construction, LRU on a larger cache always stores a superset of nodes stored by ISLRU. Therefore, it causes no more TC misses because it is lazy. \square

THEOREM 5.17. $\text{ISMIN}(W + d) \leq \text{MIN}(W)$.

The proof of Theorem 5.17 is lengthy. We will first derive some properties of Belady's optimal algorithm $\text{MIN}(W)$ and then transform any $\text{MIN}(W)$ -strategy in two steps into a $\text{ISMIN}(W + d)$ -strategy.

Recall that Belady's algorithm MIN , also called the clairvoyant algorithm, is an optimal replacement policy. The algorithm always replaces the node that will not be accessed for the longest time into the future. An elegant optimality proof for this approach is provided in [Michaud 2007]. MIN does not differentiate between nodes that will not be used again. Therefore, without loss of generality, let us from now on consider the descendant-first version of MIN . For any point in time, let us call all the nodes that are still to be accessed in the current translation **the required nodes**. The required nodes are exactly those nodes on the current translation path that are descendants of the last accessed node (or the whole path if the translation is only about to begin).

LEMMA 5.18.

- (1) Let w be in the TC. As long as w has a descendant v in the TC that is not a required node, MIN will not evict w .
- (2) If $W > d$, MIN never evicts the root.
- (3) If $W > d$, MIN never evicts a required node.

PROOF. Ad. 1. If v will be accessed ever again, then w will be used earlier (in the same translation), and so, MIN evicts v before w . If v will never be accessed again, then MIN evicts it before w because it is the descendants-first version. Ad. 2. Either the TC stores the whole current translation path, and no eviction occurs, or there is a cell in the TC that contains a node off the current translation path; hence, the root is not evicted as it has a non-required descendant in the TC. Ad. 3. Either the TC stores the whole current translation path, or there is a cell c in the TC with content that will

not be used before any required node; hence, no required node is the node that will not be needed for the longest time into the future. \square

COROLLARY 5.19. *If $W > d$, MIN inserts the root into the TC as the first insertion during the first translation, and never evicts it.*

LEMMA 5.20. *If $W > d$, MIN evicts only (non-required) nodes with no stored descendants or the node that was just used.*

PROOF. If MIN evicts a node on the current translation path, it cannot be a descendant of the just translated node (Lemma 5.18, claim 3), it also cannot be an ancestor of the just translated node (Lemma 5.18, claim 1). Hence, only the just translated node is admissible. If the algorithm evicts a node off the current translation path, it must have no descendants (Lemma 5.18, claim 1). \square

LEMMA 5.21. *If MIN has evicted the node that was just accessed, it will continue to do so for all the following evictions in the current translation. We will refer to this as **round robin** approach.*

PROOF. If MIN has evicted a node w that was just accessed, it means that all the other nodes stored in the TC will be reused before the evicted node. Moreover, all subsequent nodes traversed after w in the current translation will be reused even later than w if at all. In case of $W > d$, the claim holds by Lemma 5.20. \square

COROLLARY 5.22. *During a single translation, MIN proceeds in the following way:*

- (1) *It starts with the **regular phase** when it inserts missing nodes of a connected path from the root up to some node w , as long as it can evict nodes that will not be reused before just used ones.*
- (2) *It switches to the **round robin phase** for the remaining part of the path.*

It is easy to see that for $W > d$, in the path that was traversed in the round robin fashion, informally speaking, all gaps move up by one. For each gap between stored nodes, the very TC cell that was used to store the node above the gap now stores the last node of the gap. Storage of other nodes does not change. This way, the number of nodes from this path stored in the TC does not change either. However, it reduces numbers of stored nodes on side paths attached to the path.

In order to reach our goal, we will prove the following lemmas by modifying an optimal replacement strategy into intermediate strategies with no additional replacements.

LEMMA 5.23. *There is an eager replacement strategy on a TC of size $W + 1$ that, except for a single special cell, has ISP and causes no more TC misses than an optimal replacement strategy on a TC of size W with no restrictions.*

PROOF. We introduce a replacement strategy RRMIN⁷. We add a special cell rr to the TC, and we refer to the remaining W cells as **regular TC**. We will show that the cell rr allows us to preserve ISP in the regular TC with no additional TC misses. We start with an empty TC, and we run MIN on a separate TC of size W on a side and observe its decisions.

We keep track of a partial bijection⁸ φ_t on the nodes of the translation tree. We put one timestamp t on every TC access and one more between every two accesses in the regular phase of MIN. We position evictions and insertions between the timestamps,

⁷Round Robin MIN

⁸A partial bijection on a set is a bijection between two subsets of the set.

at most one of each between two consecutive accesses. At time t , φ_t maps every node stored by MIN in its TC to a node stored by RRMIN in its regular TC. Function φ_t always maps nodes to (not necessarily proper) ancestors in the memory translation tree. We denote this as $\varphi_t(a) \sqsubseteq a$, and in case of proper ancestors as $\varphi_t(a) \sqsubset a$. We say that a is a witness for $\varphi_t(a)$.

PROPOSITION 5.24. *Since the partial bijection φ_t always maps nodes to ancestors, for every path of the translation tree, RRMIN always stores at least as many nodes as MIN.*

In order to prove the Lemma 5.23, we need to show how to preserve the properties of the bijection φ_t and ISP. In accordance with Corollary 5.22, MIN inserts a number of highest missing nodes in the regular phase and uses the round robin approach on the remaining ones.

Let us first consider the case when MIN has only the regular phase and inserts the complete path. In this case, we substitute evictions and insertions of MIN with these described below.

Let MIN evict a node a . If $\varphi_t(a)$ has no descendants, RRMIN evicts it. In the other case, we find $\varphi_t(b)$ a descendant of $\varphi_t(a)$ with no descendants on his own. RRMIN evicts $\varphi_t(b)$, and we set $\varphi_{t+1}(b) := \varphi_t(a)$. Clearly, we have preserved the properties of φ_{t+1} ⁹, and ISP holds.

Now let MIN insert a new node e . At this point, we know that both RRMIN and MIN store all ancestors of e . If RRMIN has not yet stored e , RRMIN inserts it, and we set $\varphi_{t+1}(e) := e$. If e is already stored, it means it has a witness $\varphi_t^{-1}(e)$ that is a proper descendant of e . We find a sequence $e \sqsupset \varphi_t^{-1}(e) \sqsupset \varphi_t^{-2}(e) \sqsupset \dots \sqsupset \varphi_t^{-k}(e) = g$ that ends with g RRMIN has not stored yet. Such g exists because φ_t^{-1} is an injection on a finite set and is undefined for e . We set $\varphi_{t+1}(h) := h$ for all elements of the sequence except g . RRMIN inserts the highest not stored ancestor f of g , and we set $\varphi_{t+1}(g) := f$. Note that the inserted node f might not be a required node. Properties of φ_t are preserved, and RRMIN did not disconnect the tree it stores. Also, RRMIN performed the same number of evictions and insertions as MIN. Note as well that for all nodes on the translation path, φ_t is identity. Finally, Proposition 5.24 guarantees that all accesses are safe to perform at the time they were scheduled.

Now let us consider the case when MIN has both regular and round robin phase. Assume that the regular phase ends with the visit of node v . At this point, MIN stores the (nonempty for $W > d$ due to Corollary 5.19) initial segment p_v of the current path ending in v ; it does not contain v 's child on the current path, and it contains some number (maybe zero) of required nodes. Starting with v 's child, MIN uses the round robin strategy. Whenever it has to insert a required node, it evicts its parent. Let ℓ_r and ℓ_{rr} be the number of evictions in the regular and round robin phase, respectively.

RRMIN also proceeds in two phases. In the first phase, RRMIN simulates the regular phase as described above. RRMIN also performs ℓ_r evictions in the first phase and φ_t is the identity on p_v at the end of the first phase; this holds because φ_t maps nodes to ancestors and since MIN contains p_v in its entirety at the end of the regular phase. Let d' be the number of nodes on the current path below v ; MIN stores $d' - \ell_{rr}$ of them at the beginning of the round robin phase, which it does not have to insert, and it does not store ℓ_{rr} of them, which it has to insert. Since φ_t is the identity on p_v after phase 1 of the simulation and maps the $d' - \ell_{rr}$ required nodes stored by MIN to ancestors, RRMIN stores at least the next $d' - \ell_{rr}$ required nodes below v in the beginning of phase 2 of the simulation. In the round robin phase, RRMIN inserts the required nodes missing from

⁹ φ_{t+1} is equal to φ_t on all arguments not explicitly specified.

the regular TC one after the other into `rr`, disregarding what `MIN` does. Whenever `MIN` replaces a node a with its child g , in case of $W > d$ we fix φ_t by setting $\varphi_{t+1}(g) := \varphi_t(a)$. By Proposition 5.24, `RRMIN` does no more evictions than `MIN`. Therefore, as it also preserves `ISP` in the regular TC, Lemma 5.23 holds. \square

LEMMA 5.25. *There is a replacement strategy with `ISP` on a TC of size $W + d$, that causes no more TC misses than a general optimal replacement strategy on a TC of size W .*

PROOF. In order to prove the lemma, we will show how to use additional d regular cells in the TC to provide functionality of the special cell `rr` while preserving `ISP` in the whole TC. We run the `RRMIN` algorithm aside on a separate TC of size $W + 1$, and we introduce another replacement strategy, which we call `LIS`¹⁰, on a TC of size $W + d$. `LIS` starts with an empty TC where d cells are marked. `LIS` preserves the following invariants.

- (1) The set of nodes stored in the unmarked cells by `LIS` is equal to the set of nodes stored in the regular TC by `RRMIN`.
- (2) The set of nodes stored in the marked cells by `LIS` contains the node stored in the cell `rr` by `RRMIN`.
- (3) Exactly d cells are marked.
- (4) `LIS` has `ISP`.
- (5) No node is stored twice (once marked, once unmarked).

Whenever `RRMIN` can replicate evictions/insertions of `LIS` without violating the invariants, it does. Otherwise, we consider the following cases.

- (1) Let `RRMIN` in the regular phase evict a node a that has marked descendants in `LIS`. Then, `LIS` marks the cell containing a and unmarks and evicts one of the marked nodes with no descendants that does not store the node stored in `rr`. Such a node exists because the only other case is that the marked cells contain all nodes of some path excluding the root, and the leaf is stored in `rr`. Therefore, a is the root, but the root is never evicted due to `ISP`.
- (2) In the regular phase, `RRMIN` inserts a node c to an empty cell while `LIS` already stores c in a marked cell. In this case, `LIS` unmarks the cell with c and marks the empty cell.
- (3) In the round robin phase, `RRMIN` replaces the content of the cell `rr`, `LIS` (if needed) replaces the content of an arbitrary marked node with no descendants that is not on the current translation path. Since the root is always in the TC and there are d marked cells, such a cell always exists. `ISP` is preserved, as the parent of this node is already in the TC.

At this stage, if we drop notions of φ_t and marked nodes, `LIS` becomes an eager replacement strategy on a standard TC. Therefore, we can use Lemma 5.8 to make it lazy. This concludes the proof of Lemma 5.25. \square

Since `ISMIN` is an optimal strategy with `ISP`, Theorem 5.17 follows from Lemma 5.25.

Remark 5.26. We believe that the requirement for d additional cache size is essentially optimal. Consider the scenario when we access memory cells uniformly at random. Informally speaking, `MIN` will tend to permanently store the first $\log_K(W)$ levels of the translation tree because they are frequently used and will use a single cell to traverse the lower levels. In order to preserve `ISP`, we need $d - \log_K(W) + 1$ additional cells for storing the current path. Thus only little improvement seems to be possible.

¹⁰Lazy strategy preserving the Initial Segments property

CONJECTURE 5.27. *The strategy of storing higher nodes (Lemma 5.23) and using extra d cells to not evict nodes from the current translation path (Lemma 5.25) can be used to add ISP to any replacement strategy without efficiency loss.*

6. ANALYSIS OF ALGORITHMS

In this section, we analyze the translation cost of some algorithms as a function of the problem size n and memory requirement m . For all the algorithms analyzed, $m = \theta(n)$.

In the RAM model, there is a crucial assumption that usually goes unspoken, namely, the size of a machine word is logarithmic in the number of memory cells used. If the words were shorter, one could not address the memory. If the words were longer, one could intelligently pack multiple values in one cell. This technique can be used to solve NPC problems in polynomial time. This effectively puts an upper bound on n , namely, $n < 2^{\text{word length}}$, while asymptotic notations make sense only when n can grow to infinity. However, this is not a limitation of the RAM model, it merely shows that to handle bigger inputs, one needs more powerful machines.

In the VAT model there is also a set of assumptions on the model constants. The assumptions bound n by machine parameters in the same sense as in the RAM model. However, unlike in the RAM model, they can hardly go unspoken. We call them the *asymptotic order relations* between parameters. The assumptions we found necessary for the analysis to be meaningful are as follows:

- (1) $1 \leq \tau d \leq P$; moving a single translation path to the TC costs more than a single instruction, but not more than size-of-a-page many instructions, i.e., if at least one instruction is performed for each cell in a page, the cost of translating the index of the page can be amortized.
- (2) $K \geq 2$, i.e., the fanout of the translation tree is at least two.
- (3) $m/P \leq K^d \leq 2m/P$, i.e., the translation tree suffices to translate all addresses but is not much larger. As a consequence, $\log(m/P) \leq d \log K = dk \leq \log(2m/P) = 1 + \log(m/P)$, and hence, $\log_K(m/P) \leq d \leq \log_K(2m/P) = (1 + \log(m/P))/k$.
- (4) $d \leq W < m^\theta$, for $\theta \in (0, 1)$, i.e., the translation cache can hold at least one translation path, but is significantly smaller than the main memory.

6.1. Sequential Access

We scan an array of size n , i.e., we need to translate addresses $b, b+1, \dots, b+n-1$ in this order, where b is the base address of the array. The translation path stays constant for P consecutive accesses, and hence, at most $2n/P$ indices must be translated for a total cost of at most $\tau d(2+n/P)$. By assumption (1), this is at most $\tau d(n/P+2) \leq n+2P$.

The analysis can be sharpened significantly. We keep the current translation path in the cache, and hence, the first translation incurs at most d faults. The translation path changes after every P -th access and hence changes at most a total of $\lceil n/P \rceil$ times. Of course, whenever the path changes, the last node changes. The next to last node changes after every K -th change of the last node and hence changes at most $\lceil n/(PK) \rceil$ times. In total, we incur

$$d + \sum_{0 \leq i \leq d} \left\lceil \frac{n}{PK^i} \right\rceil < 2d + \frac{K}{K-1} \frac{n}{P}$$

cache faults; of these faults, at most $1 + n/P$ are caused by data pages and the remaining ones are caused by internal nodes of the translation tree. The cost is therefore bounded by $2\tau d + 2\tau n/P \leq 2P + 2n/d$, which is asymptotically smaller than the RAM complexity.

6.2. Random Access

In the worst case, no node of any translation path is in the cache. Thus the total translation cost is bounded by τdn . This is at most $\tau n \log_K(2n/P)$.

We will next argue a lower bound. We may assume that the TC satisfies the initial segment property. The translation path ends in a random leaf of the translation tree. For every leaf, some initial segment of the path ending in this leaf is cached. Let u be an uncached node of the translation tree of minimal depth, and let v be a cached node of maximal depth. If the depth of v is larger by two or more than the depth of u , then it is better to cache u instead of v (because more leaves use u instead of v). Thus, up to one, the same number of nodes is cached on every translation path, and hence, the expected length of the path cached is at most $\log_K W$, and hence, the expected number of faults during a translation is $d - \log_K W$. The total expected cost is therefore at least $\tau n(d - \log_K W) \geq \tau n \log_K n/(PW)$, which is asymptotically larger than the RAM complexity.

LEMMA 6.1. *The memory access cost of a random scan of an array of size n is at least $\tau n \log_K(n/(PW))$ and at most $\tau n \log_K(2n/P)$.*

6.3. Binary Search

We do n binary searches in an array of length n . Each search searches for a random element of the array. For simplicity, we assume that n is a power of two minus one. A binary search in an array is equivalent to a search in a balanced tree where the root is stored in location $n/2$, the children of the root are stored in locations $n/4$ and $3n/4$, and so on. We cache the translation paths of the top ℓ layers of the search tree and the translation path of the current node of the search. The top ℓ layers contain $2^{\ell+1} - 1$ vertices, and hence, we need to store at most $d2^{\ell+1}$ nodes¹¹ of the translation tree. This is feasible if $d2^{\ell+1} \leq W$. For next two paragraphs, let $\ell = \log(W/2d)$.

Any of the remaining $\log n - \ell$ steps of the binary search causes at most d cache faults. Therefore, the total cost per search is bounded by

$$\tau d(\log n - \ell) \leq \tau \log_K(2n/P)(\log n - \ell) = \tau \log_K \frac{2n}{P} \log \frac{2nd}{W}.$$

This analysis may seem unrefined. After all once the search leaves the top ℓ layers of the search tree, addresses of subsequent nodes differ only by $n/2^\ell, n/2^{\ell+1}, \dots, 1$. However, we will next argue that the bound above is essentially sharp for our caching strategy. In a second step, we will extend the bound to all caching strategies. By Lemma 5.1, if two virtual addresses differ by D , their translation paths differ in the last $\lceil \log_K(D/P) \rceil$ nodes. Thus, the scheme above incurs at least

$$\begin{aligned} \sum_{\ell \leq i \leq \log n - p} \left\lceil \frac{1}{k} \log \frac{n}{2^i P} \right\rceil &\geq \sum_{0 \leq j \leq \log n - \ell - p} \frac{1}{k} \log 2^j \geq \\ &\geq \frac{1}{2k} (\log n - \ell - p)^2 = \frac{1}{2} \log_K \frac{2nd}{PW} \log \frac{2nd}{PW} \end{aligned}$$

cache faults. We next show that it essentially holds true for any caching strategy.

By Theorem 5.15, we may assume that ISLRU is used as the cache replacement strategy, i.e., TC contains top nodes of the recent translation paths. Let $\ell = \lceil \log(2W) \rceil$. There are $2^\ell \geq 2W$ vertices of depth ℓ in a binary search tree. Their addresses differ by at least $n/2^\ell$, and hence, for any two such addresses, their translation paths differ in at least the last $z = \lceil \log_K(n/(2^\ell P)) \rceil$ nodes. Call a node at depth ℓ *expensive* if none of the

¹¹We use vertex for the nodes of the search tree and node for the nodes of the translation tree.

last z nodes of its translation path are contained in the TC and *inexpensive* otherwise. There can be at most W inexpensive nodes, and hence, with probability at least $1/2$, a random binary search goes through an expensive node, call it v , at depth ℓ . Since ISLRU is the cache replacement strategy, the last z nodes of the translation path are missing for all descendants of v . Thus, by the argument in the preceding paragraph, the expected number of cache misses per search is at least

$$\frac{1}{2} \sum_{\ell \leq i \leq \log n - p} \left\lceil \frac{1}{k} \log \frac{n}{2^i P} \right\rceil \geq \frac{1}{4} \log_K \frac{2nd}{PW} \log \frac{2nd}{PW}$$

LEMMA 6.2. *The memory access cost of n random binary searches in an array of size n is at most $\tau \log_K \frac{2n}{P} \log \frac{2nd}{W}$ and at least $\frac{\tau}{4} \log_K \frac{2nd}{PW} \log \frac{2nd}{PW}$*

We know from cache-oblivious algorithms that the van-Emde Boas layout of a search tree improves locality. We will show in Section 7 that this improves the translation cost.

6.4. Heapify and Heapsort

We prove a bound on the translation cost of heapify. The following proposition generalizes the analysis of sequential scan.

Definition 6.3. Extremal translation paths of n consecutive addresses are the paths to the first and last address in the range. **Non-extremal nodes** are the nodes on translation paths to addresses in the range that are not on the extremal paths.

PROPOSITION 6.4. *A sequence of memory accesses that gains access to each page in a range causes at least one TC miss for each non-extremal node of the range. If the sequence of pages in the range n is accessed in the decreasing order, this bound is matched by storing the extremal paths and dedicating $\log_K(n/P)$ cells in the TC for the required translations.*

PROPOSITION 6.5. *Let n , ℓ , and x be nonnegative integers. The number of non-extremal nodes in the union of the translation paths of any x out of n consecutive addresses is at most*

$$x\ell + \frac{2n}{PK^\ell}.$$

Moreover, there is a set of $x = \lceil n/(PK^\ell) \rceil$ addresses such that the union of the paths has size at least $x(\ell + 1) + d - \ell$.

PROOF. The union of the translation paths to all n addresses contains at most n/P non-extremal nodes on the leaf level (= level 0) of the translation tree. On level i , $i \geq 0$, from the bottom, it contains at most $n/(PK^i)$ non-extremal nodes.

We overestimate the size of the union of x translation paths by counting one node each on levels 0 to $\ell - 1$ for every translation path and all non-extremal nodes contained in all the n translation paths on the levels above. Thus, the size of the union is bounded by

$$x\ell + \sum_{\ell \leq i \leq d} n/(PK^i) < x\ell + \frac{K}{K-1} \frac{n}{PK^\ell} \leq x\ell + \frac{2n}{PK^\ell}.$$

A node on level ℓ lies on the translation path of $K^\ell P$ consecutive addresses. Consider addresses $z + iPK^\ell$ for $i = 0, 1, \dots, \lceil n/(PK^\ell) \rceil - 1$, where z is the smallest in our set of n addresses. The translation paths to these addresses are disjoint from level ℓ down to level zero and use at least one node on levels $\ell + 1$ to d . Thus, the size of the union is at least $x(\ell + 1) + d - \ell$. \square

An array $A[1..n]$ storing elements from an ordered set is heap-ordered if $A[i] \leq A[2i]$ and $A[i] \leq A[2i + 1]$ for all i with $1 \leq i \leq \lfloor n/2 \rfloor$. An array can be turned into a heap by calling operation $sift(i)$ for $i = \lfloor n/2 \rfloor$ down to 1. $sift(i)$ repeatedly interchanges $z = A[i]$ with the smaller of its two children until the heap property is restored. We use the following translation replacement strategy. Let $z = \min(\log n, \lfloor (W - 2d - 1) / \lfloor \log_K(n/P) \rfloor \rfloor - 1)$. We store the extremal translation paths ($2d - 1$ nodes), non-extremal parts of the translation paths for z addresses a_0, \dots, a_{z-1} , and one additional translation path a_∞ ($\lfloor \log_K(n/P) \rfloor$ nodes for each). The additional translation path is only needed when $z \neq \log n$. During the siftdown of $A[i]$, a_0 is equal to the address of $A[i]$, a_1 is the address of one of the children of i (the one to which $A[i]$ is moved, if it is moved), a_2 is the address of one of the grandchildren of i (the one to which $A[i]$ is moved, if it is moved two levels down), and so on. The additional translation path a_∞ is used for all addresses that are more than z levels below the level containing i .

Let us upper bound the number of the TC misses. Preparing the extremal paths causes up to $2d + 1$ misses. Next, consider the translation cost for a_i , $0 \leq i \leq z - 1$. a_i assumes $n/2^i$ distinct values. Assuming that siblings in the heap always lie in the same page¹², the index (= the part of the address that is being translated) of each a_i decreases over time, and hence, Proposition 6.4 bounds the number of TC misses to the number of the non-extremal nodes in the range. We use Proposition 6.5 to count them. For $i \in \{0, \dots, p\}$, we use the Proposition with $x = n$ and $\ell = 0$ and obtain a bound of

$$\frac{2n}{P} = O\left(\frac{n}{P}\right)$$

TC misses. For i with $p + (\ell - 1)k < i \leq p + \ell k$, where $\ell \geq 1$ and $i \leq z - 1$, we use the proposition with $x = n/2^i$ and obtain a bound of at most

$$\frac{n}{2^i} \cdot \ell + \frac{2n}{PK^\ell} = O\left(\frac{n}{2^i} \cdot \ell + \frac{2n}{2^i}\right) = O\left(\frac{n}{2^i}(\ell + 2)\right) = O\left(n \frac{i}{2^i}\right)$$

TC misses. There are $n/2^z$ siftdowns starting in layers z and above; they use a_∞ . For each such siftdown, we need to translate at most $\log n$ addresses, and each translation causes less than d misses. The total is less than $n(\log n)d/2^z$. Summation yields

$$2d + 1 + (p + 1)O\left(\frac{n}{P}\right) + \sum_{p < i \leq z-1} O\left(n \frac{i}{2^i}\right) + \frac{nd \log n}{2^z} = O\left(d + \frac{np}{P} + \frac{nd \log n}{2^z}\right).$$

For any realistic values of the parameters, the third term is insignificant, hence, the cost is $O\left(\tau\left(d + \frac{np}{P}\right)\right)$.

We next prove the corresponding lower bound under the additional assumption that $W < \frac{1}{2}n/P$. At least one address must be completely translated; hence the cost of $\Omega(\tau d)$. The addresses in $a_0 \dots a_{p-1}$ assume at least one address per page in the subarray $[n/2..n]$ because a_i can never jump by more than 2^{i+1} . First, the addresses are swept by a_0 , then by a_1 , and so on, and no other accesses to the subarray occur in the meantime. Hence, if the LRU strategy is in use and $W < \frac{1}{2}n/P$, there are at least $pn/(2P)$ TC misses to the lowest level of the translation tree. This gives the $\Omega\left(\frac{np}{P}\right)$ part of the lower bound. Hence, the total cost is $\Omega\left(\tau\left(d + \frac{np}{P}\right)\right)$.

In the sorting phase of heapsort, we repeatedly remove the element stored in the root, move the element in the rightmost leaf to the root, and then let this element sift-

¹²This assumption can be easily lifted by allowing an additional constant in the running time or in the TC size.

down to restore the heap-property. The sift-down starts in the root and after accessing address i of the heap moves to address $2i$ or $2i + 1$. For the analysis, we make the additional assumption $W = M$, i.e., the data cache and the TC cache have the same size. We store the top ℓ layers of the heap in the data cache and the translation paths to the vertices to these layers in the TC cache, where $2^{\ell+1} < M$, say $\ell = \log(M/4) = \log(W/4)$. Each of the remaining $\log n - \ell$ sift-down steps may cause d cache misses. The total number of cache faults is therefore bounded by

$$nd(\log n - \ell) \leq n \log_K(2n/P) \log(4n/W).$$

We leave the lower bound as an open problem.

7. CACHE-OBLIVIOUS ALGORITHMS

Algorithms for the EM model are allowed to use the parameters of the memory hierarchy in the program code. For any two adjacent levels of the hierarchy, there are two parameters, the size M of the faster memory and the size B of the blocks in which data is transferred between the faster and the slower memory. Cache-oblivious algorithms are formulated without reference to these parameters, i.e., they are formulated as RAM-algorithms. Only the analysis makes use of the parameters. A transfer of a block of memory is called an IO-operation. For a cache-oblivious algorithm, let $C(M, B, n)$ be the number of IO-operations on an input of size n , where M is the size of the faster memory (also called cache memory) and B is the block size. Of course, $B \leq M$.

Good cache-oblivious algorithms exhibit good locality of reference at all scales, and therefore, one may hope that they also show good behavior in the VAT model. The following theorem gives an upper bound of VAT-complexity in terms of the EM-complexity of an algorithm.

THEOREM 7.1. *Consider a cache-oblivious algorithm with IO-complexity $C(M, B, n)$, where M is the size of the cache, B is the size of a block, and n is the input size. Let $a := \lfloor W/d \rfloor$, where W is the size of the translation cache, and let $P = 2^p$ be the size of a page. Then, the number of cache faults in the VAT-model is at most*

$$\sum_{i=0}^d C(aK^i P, K^i P, n).$$

PROOF. We divide the cache into d parts of size a and reserve one part for each level of the translation tree.

Consider any level i , where the leaves of the translation tree (= data pages) are on level 0. Each node on level i stands for $K^i P$ addresses, and we can store a nodes. Thus, the number of faults on level i in the translation process is the same as the number of faults of the algorithm on blocks of size $K^i P$ and a memory of a blocks (i.e., size $aK^i P$). Therefore, the number of cache faults is at most

$$\sum_{i=0}^d C(aK^i P, K^i P, n).$$

□

Theorem 7.1 allows us to rederive some of the results from Section 6. For example, a linear scan of an array of length n has IO-complexity at most $2 + \lfloor n/B \rfloor$. Thus, the number of cache faults in the VAT-model is at most

$$\sum_{i=0}^d \left(2 + \frac{n}{K^i P} \right) < 2d + \frac{K}{K-1} \frac{n}{P}.$$

It also allows us to derive new results. Quicksort has IO-complexity $O((n/B) \log(n/B))$, and hence, the number of cache faults in the VAT-model is at most

$$\sum_{i=0}^d O\left(\frac{n}{K^i P} \log \frac{n}{K^i P}\right) = O\left(\frac{n}{P} \log \frac{n}{P}\right).$$

Binary search in the van Emde Boas layout has IO-complexity $\log_B n$, and hence, the number of cache faults in the VAT-model is at most

$$\begin{aligned} \sum_{i=0}^d \frac{\log n}{\log(K^i P)} &\leq \frac{\log n}{\log P} + \log n \int_0^d \frac{1}{\log P + x \log K} dx \\ &= \log_P n + \frac{\log n}{\log K} \ln(\log P + x \log K) \Big|_0^d \\ &= \log_P n + (\log_K n)(\ln(\log K^d P) - \ln(\log P)) \\ &= \log_P n + (\log_K n) \ln \log_P(PK^d) \leq \log_P n + \log_K n \ln \log_P 4n, \end{aligned}$$

where the last inequality follows from our assumption that $K^d P$ is at most twice the memory footprint of an algorithm and that the memory footprint of a binary tree with n leaves is bounded by $2n$.

A matrix multiplication with a recursive layout of matrices has IO-complexity $n^3/(M^{1/2}B)$, and hence, the number of cache faults in the VAT-model is at most

$$\sum_{i=0}^d \frac{n^3}{(aK^i P)^{1/2} K^i P} < \frac{K^{3/2}}{K^{3/2} - 1} \frac{n^3}{a^{1/2} P^{3/2}}.$$

Cache-oblivious algorithms that match the performance of the best EM-algorithm for the problem are known for several fundamental algorithmic problems, e.g., sorting, FFT, matrix multiply, and searching [Frigo et al. 2012]. Do all these algorithms automatically have small VAT-complexity via Theorem 7.1? Unfortunately, the answer is *no*. Observe that the theorem refers to the cache misses in a machine with memory size $aK^i P$ and block size $K^i P$, i.e., memory consists of a blocks. However, many of the good cache-oblivious algorithms require a tall-cache assumption $M \geq B^2$; sometimes, the assumption $M \geq B^{1+\epsilon}$ for some positive ϵ suffices. For such algorithms, the theorem above does not give good bounds.

In joint work with Pat Nicholson [Jurkiewicz et al. 2014], we have recently shown that cache-oblivious algorithms requiring a tall-cache assumption also perform well in the VAT-model provided a somewhat more stringent tall cache assumption holds. More precisely, 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. 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 $M \geq 4g(dB)$. 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 [Frigo et al. 2012] 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 $\widetilde{M} \geq \widetilde{B}^2$. As a consequence of our main theorem, we obtain:

THEOREM 7.2. *Funnel sort sorts n items, each of size $a \geq 1$ addressable units, on a VAT-machine with cache size \overline{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 $M = \overline{M}/a$ 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.

8. DISCUSSION

In this section, we discuss additional topics that extend the scope of our research. In particular, we address the comments that we received from the ALENEX13 program committee and other researchers.

8.1. Double Address Translation on Virtual Machines

Nowadays, more and more computation is performed on virtual machines in the clouds. In this environment, address translation must be performed twice, first to the virtual machine addressing space and then to the host. The cost of address translation to host can be as high as $O(\tau \log(\text{size of virtual machine}))$. Moreover, big enough virtual machines may require translation for memory tables in the virtual machine, not just for the data. This is independent of the problem input size and significant in the case of random access, but still negligible in the case of sequential access. To test the impact of the double address translation, we timed permutation and introsort on a virtual machine; results are provided in Figure 6.

Please note that STL introsort takes actually less time than the permutation generator, even for very small data. This is very surprising at first but means that a high VAT cost is especially harmful for programs launched on virtual machines. Since many cloud systems are meant primarily for computing, the discussed phenomenon should be of primal concern for such environments.

8.2. The Model is Too Complicated

While we received comments that the model is too simple, we also received ones saying that the model is too complicated. This impression is probably due to the fact that some of our proofs are somewhat technical. Some arguments simplify if asymptotic notation is used earlier, or if the VAT cost is upper bounded by the RAM cost ahead of time (for sequential access patterns to the memory), or the other way around for the randomized access. However, as this is the first work addressing the subject, we find it appropriate to be more detailed than absolutely necessary. With time, more and more simplifications will appear. Let us briefly discuss a few candidates.

8.2.1. Value of K . There is evidence that for many algorithms, the exact value of K does not matter, and hence, $K = 2$ may be used. In some cases, like repeated binary search, the exact value of K seems to have only a little impact both in theory and practice. In other cases, like permutation, it seems to be the cause of bumps on the chart in Figure 1, but the impact is moderate. A notable exception is matrix transpose and matrix multiplication, where the value of K is blatantly visible. The classic matrix transpose algorithm uses $O(n)$ operations, where n is the input size. However, if the matrix is stored and read row by row, the output matrix must be written one element per row. For a square matrix, this means a jump of \sqrt{n} cells between writes,

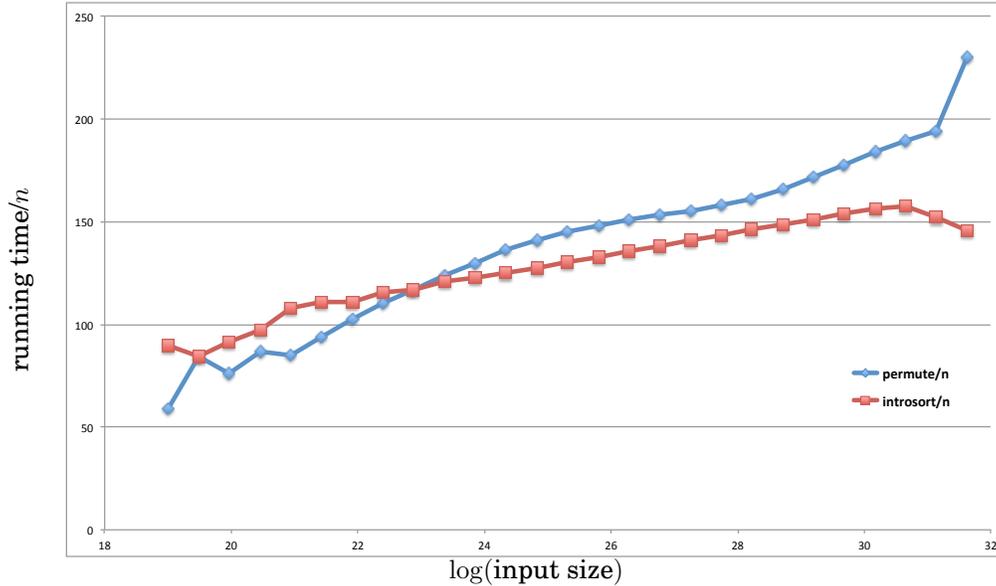


Fig. 6. Execution times divided by n (not the normalized operation time) on a virtual machine with the following specification:
 cpu: Intel(R) Xeon(R) CPU X5660 @ 2.67GHz
 operating system: Windows 7 Enterprise
 compiler: Microsoft Visual Studio 2010

which means \sqrt{n} translations of cost $\Theta(\tau d)$ to produce the first column. As there are \sqrt{n} translations before another element is written to the same row, no translation path can be reused if we consider the LRU algorithm. Therefore, the total VAT cost is $\Theta(\tau nd)$, which is $\Theta(\tau n \log n)$. Figure 7 shows that even though the asymptotic growth is intact, the translation cost grows in jumps rather than in a smooth logarithmic fashion. The distance between the jumps appears to be directly related to the value of K , namely, the jump occurs when the matrix dimension is K times greater than during the previous jump. Note that the EM cost of this algorithm is $\Theta(n)$ for $\sqrt{n} \cdot B > M$, and $\Theta(n/B)$ for $\sqrt{n} \cdot B < M$. In fact, the first cost jump is due to this barrier itself.

8.2.2. *CAT, or Sequence of Consecutive Address Translations.* In our analysis, for many algorithms precisely calculated VAT complexity was much smaller than the RAM complexity. We believe that our approach can bring valuable insight for future research, but some of our results can be obtained in a simpler way. The memory access patterns in the algorithms in question share some common characteristics. There are not too few elements, they are not overspread in the memory, and the accesses are more or less performed in a sequence. We formalize these properties in the following definition.

Definition 8.1. We call a sequence of ℓ memory accesses a **CAT** (sequence of consecutive address translations) if it fulfills the following conditions:

- $\ell = \Omega(\tau d)$.
- On average, $\Omega(\tau)$ elements are accessed per page in the access range.
- The pages are used in the increasing or decreasing order. (But accesses in the page need not follow this rule).

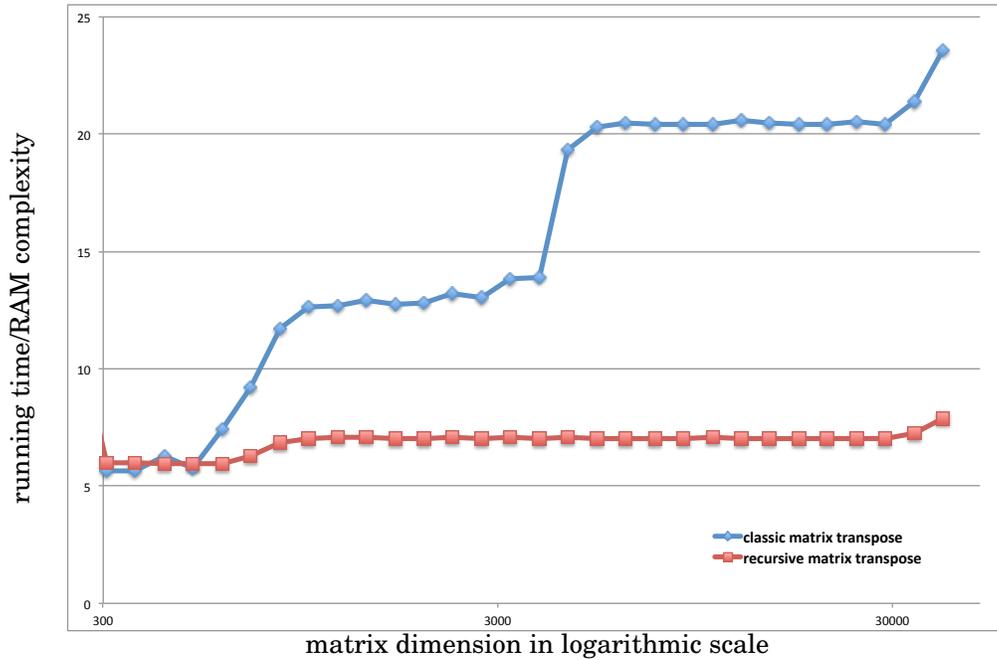


Fig. 7. Running time of the matrix transpose in row by row layout and in the recursive layout

— Memory accesses from the sequence are separated by at most a constant number of other operations.

LEMMA 8.2. *In case of a CAT, the cost of the address translation is dominated by the cost of RAM operations and therefore negligible. Hence, for CATs, it is sufficient to account for them only in the RAM part of the analysis.*

PROOF. We assume the LRU replacement strategy. First, let us assess the cost of translating addresses for all the $O(\ell/\tau)$ pages in the increasing order. The first translation causes d TC misses. Since we allow only a constant number of operations between accesses from the considered sequence, the LRU replacement strategy holds translation path of the last translation when the next one starts. Hence, the addresses to be translated change like in a classic K -nary counter. The amortized cost of an update of a K -nary counter is $O(1)$. Since on average $\Omega(\tau)$ elements are accessed per page, the access range is at most of length $O(\ell/\tau)$, and so the cost of updates is $O(\ell)$. However, we do not start counting from zero, and the potential function in the K -nary counter analysis can reach up to \log (the highest number seen), which in our case can reach d . Hence, we need to add the cost of another d TC misses to our estimation. The cost of all translations is therefore equal to $O(\tau d + \ell) = O(\ell)$.

In the definition of a CAT, we do not assume that every page is used exactly once. However, neither skipping values in the counter, nor reusing them causes extra TC misses.

Since the RAM cost is exactly $\Theta(\ell)$, it dominates the translation cost. \square

8.2.3. RAT, or Sequence of Random Address Translations. Similarly, algorithms with a high VAT cost share common properties.

Definition 8.3. A sequence of memory accesses is called a **RAT** (sequence of random address translations) if:

- There is a memory partitioning such that each part consists of all memory cells with some common virtual address prefix, and parts are of size at least Pm^θ for $\theta \in (0, 1)$.
- For at least a constant fraction of the accesses with at least a constant probability, each access is to a part that was not accessed since W TC misses.

LEMMA 8.4. *The cost of a RAT of length ℓ is $\Theta(\tau\ell d)$. It is the same as the cost of the address translations.*

PROOF. We assume the LRU replacement strategy. Since parts are of size at least Pm^θ for $\theta \in (0, 1)$, a translation of an address from each part uses $\Theta(d)$ translation nodes unique to its translation subtree. Therefore, an access to a part that was not accessed since W TC misses, misses the root of the subtree, and by Lemma 5.9, the access causes $\Theta(d)$ misses. As this happens for at least a constant fraction of the accesses with at least a constant probability, the total cost is $\Theta(\tau\ell d)$. The RAM cost is only $\Theta(\ell)$, which is less than the VAT cost by the order assumption 1. \square

8.3. Larger Page Sizes

The straight-forward method to determine how the Virtual Address Translation affects the running time of programs would be to switch it off and compare the results. Unfortunately, no modern operating system provides such an option. One can approximate the elimination of address translation by increasing the page size. If all the data fits into a single page, address translation is essentially eliminated. If all the data fits into a small number of pages, the number of translations and their cost is reduced. We performed experiments with larger page sizes. However, while hardware architectures support pages sized in gigabytes, operating systems do not. Quoting [Hennessy and Patterson 2007]:

“Relying on the operating systems to change the page size over time.

The Alpha architects had an elaborate plan to grow the architecture over time by growing its page size, even building it into the size of its virtual address. When it came time to grow page sizes with later Alphas, the operating system designers balked and the virtual memory system was revised to grow the address space while maintaining the 8 KB page.

Architects of other computers noticed very high TLB miss rates, and so added multiple, larger page sizes to the TLB. The hope was that operating systems programmers would allocate an object to the largest page that made sense, thereby preserving TLB entries. After a decade of trying, most operating systems use these superpages only for handpicked functions: mapping the display memory or other I/O devices, or using very large pages for the database code.”

There are good reasons why operating systems designer are reluctant to offer larger pages. The main concern is space. Pages must be correctly aligned in memory, so bigger pages lead to a greater waste of memory and limited flexibility while paging to disk. Another problem is that since most processes are small, using bigger pages would lengthen their initialization time. Therefore, current operating systems kernels provide only basic, nontransparent support for bigger pages. The *hugetlbpage* feature of current Linux kernels allows one to use pages of size 2MiB on AMD64 machines. The following links describe the *hugetlbpage*-feature.

- <http://linuxgazette.net/155/krishnakumar.html>
- <https://www.kernel.org/doc/Documentation/vm/hugetlbpage.txt>

— <https://www.kernel.org/doc/Documentation/vm/hugepage-shm.c>
 — <http://man7.org/linux/man-pages/man2/shmget.2.html>

The feature attaches a final real address one level higher in the memory table, i.e., the last layer of nodes is eliminated from the translation trees and pages are now of size 2^{9+12} . This slightly decreases cache usage, decreases the number of nodes needed in each single translation but one, and finally, increases the range of addresses covered by the related entry in the TLB by 512.

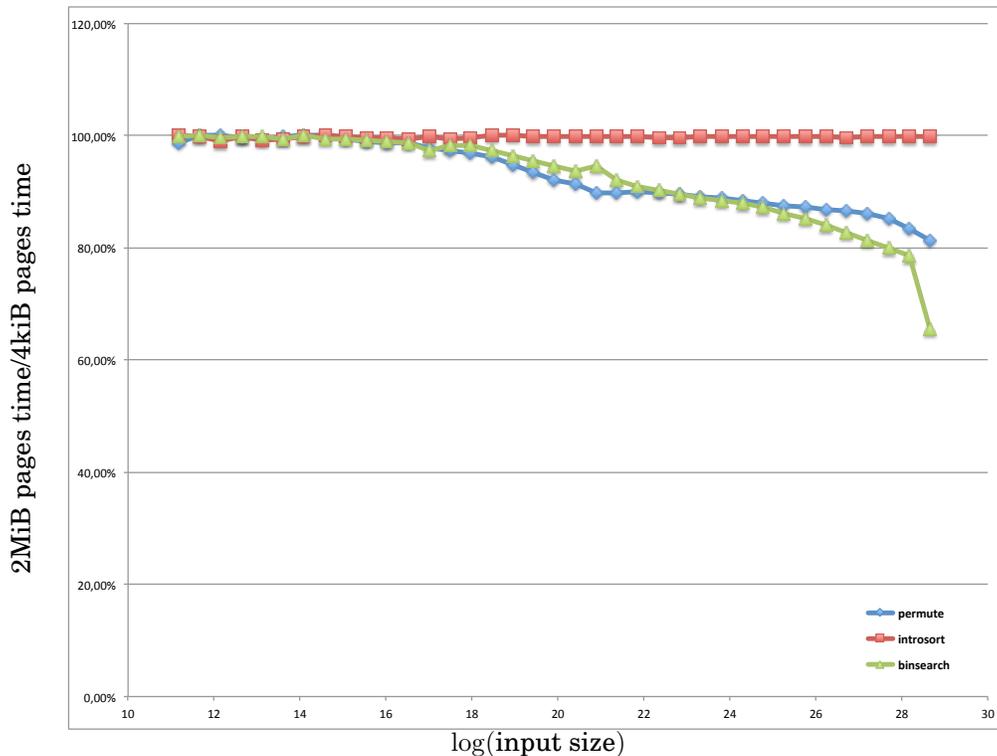


Fig. 8. The abscissa shows the logarithm of the input size. The ordinate shows the measured running time with use of the 2MiB pages as a percentage of the running time with 4kiB pages.

We rerun the *permute*, *introsort* and *binsearch* on the same machine, with and without use of the big pages. Figure 8 clearly shows that use of bigger pages can introduce a speedup. In other words, the cost of virtual address translation can be partially reduced by use of the bigger pages.

8.4. The Translation Tree is Shallow

It is true that the height of the translation tree on today's machines is bounded by 4, and so the translation cost is bounded. However, even though our experiments use only 3 levels, the slowdown appears to be at least as significant in practice as the one caused by a factor of $\log n$ in the operational complexity. Therefore, decreasing VAT complexity has a prominent practical significance. Please note that while 64 bit addresses are sufficient to address any memory that can ever be constructed according to known

physics, there are other practical reasons to consider longer addresses. Therefore, the current bound for the height of the translation tree is not absolute.

8.5. What about Hashing?

We have been asked whether the current virtual address translation system could be replaced with one based on hashing tables to achieve a constant amortized translation time. Let us argue that it is not a good idea. First and foremost, hashing tables sometimes need rehashing, and this would mean the complete blockage of an operating system. Moreover, an adversary can try to increase a number of necessary rehashes. Note that probabilistic guarantees are on the frequencies of the rehashes and the program isolation is insufficient to discard this concern, because an attack can be performed with side-channels like, for example, differential power analysis (see [Tiri 2007]). Finally, a tree walk is simple enough to be supported by hardware to obtain significant speedups; in case of hashing, this would be not so easy.

On the other hand, simple hash tables can be used to implement efficient caches. In fact, associative memory can be seen as a hardware implementation of a hashing table. If we no longer require from the associative memory that it reliably stores all the previous entries, then associative memories of small enough sizes can be well supported by hardware. This is in fact how the TLB is implemented, and one of the reasons why it is so small.

9. CONCLUSIONS

We have introduced the VAT model and have analyzed some fundamental algorithms in this model. We have shown that the predictions made by the model agree well with measured running times. Our work is just the beginning. In follow-up, we show together with Patrick Nicholson [Jurkiewicz et al. 2014] that all cache-oblivious algorithms perform well in the VAT-model provided a tall cache assumption that is somewhat more stringent than for the EM-model. It would be interesting to know, whether this more stringent assumption is necessary.

We believe that every data structure and algorithms course must also discuss algorithm engineering issues. One such issue is that the RAM model ignores essential aspects of modern hardware. The EM model and the VAT model capture additional aspects.

REFERENCES

- Advanced Micro Devices. 2010. AMD64 architecture programmer’s manual volume 2: System programming. (2010).
- Alok Aggarwal and S. Vitter, Jeffrey. 1988. The input/output complexity of sorting and related problems. *Commun. ACM* 31, 9 (Sept. 1988), 1116–1127. DOI: <http://dx.doi.org/10.1145/48529.48535>
- Ulrich Drepper. 2007. What Every Programmer Should Know About Memory. (2007). <http://lwn.net/Articles/250967/>.
- Ulrich Drepper. 2008. The Cost of Virtualization. *ACM Queue* 6, 1 (2008), 28–35.
- M. Frigo, C.E. Leiserson, H. Prokop, and S. Ramachandran. 2012. Cache-Oblivious Algorithms. *ACM Transactions on Algorithms* (2012), 4:1 – 4:22. a preliminary version appeared in FOCS 1999.
- John L. Hennessy and David A. Patterson. 2007. *Computer Architecture: A Quantitative Approach*. Morgan Kaufmann, San Diego.
- Tomasz Jurkiewicz, Kurt Mehlhorn, and Patrick Nicholson. 2014. Cache-Oblivious VAT-Algorithms. (April 2014).
- Pierre Michaud. 2007. (Yet another) proof of optimality for MIN replacement. (October 30 2007). <http://www.irisa.fr/caps/people/michaud/yap.pdf>.
- Naïla Rahman. 2003. Algorithms for Hardware Caches and TLB. In *Algorithms for Memory Hierarchies*, Ulrich Meyer, Peter Sanders, and Jop Sibeyn (Eds.). Lecture Notes in Computer Science, Vol. 2625.

Springer Berlin / Heidelberg, 171–192. http://dx.doi.org/10.1007/3-540-36574-5_8 10.1007/3-540-36574-5_8.

J. C. Shepherdson and H. E. Sturgis. 1963. Computability of Recursive Functions. *J. ACM* 10, 2 (1963), 217–255.

D. Sleator and R.E. Tarjan. 1985. Amortized Efficiency of List Update and Paging Rules. *Commun. ACM (CACM)* 28, 2 (1985), 202–208.

Kris Tiri. 2007. Side-channel attack pitfalls. In *Proceedings of the 44th annual Design Automation Conference (DAC '07)*. ACM, New York, NY, USA, 15–20. DOI: <http://dx.doi.org/10.1145/1278480.1278485>

Received March 2013; revised ??; accepted ??