

Efficient Tree Layout in a Multilevel Memory Hierarchy*

Stephen Alstrup[†] Michael A. Bender[‡] Erik D. Demaine[§]
Martin Farach-Colton[¶] J. Ian Munro^{||} Theis Rauhe[†]
Mikkel Thorup^{**}

November 12, 2002

Abstract

We consider the problem of laying out a tree with fixed parent/child structure in hierarchical memory. The goal is to minimize the expected number of block transfers performed during a search along a root-to-leaf path, subject to a given probability distribution on the leaves. This problem was previously considered by Gil and Itai, who developed optimal algorithms when the block-transfer size B is known. We show how to extend any approximately optimal algorithm to the *cache-oblivious* setting in which the block-transfer size is unknown to the algorithm. The query performance of the cache-oblivious layout is within a constant factor of the query performance of the optimal known-block-size layout. Computing the cache-oblivious layout requires only logarithmically many calls to the layout algorithm with known block size. Finally, we analyze two greedy strategies, and show that they have a performance ratio between $\Omega(\lg B / \lg \lg B)$ and $O(\lg B)$ when compared to the optimal layout.

1 Introduction

The B-tree [3] is the classic optimal search tree for external memory, but it is only optimal when accesses are uniformly distributed. In practice most distributions are nonuniform, e.g., distributions with heavy tails arise almost universally throughout computer science.

*A preliminary version of this paper appeared in ESA 2002 [7].

[†]IT University of Copenhagen, Glentevej 65-67, DK-2400 Copenhagen NV, Denmark; email: {stephen, theis}@diku.dk.

[‡]Department of Computer Science, State University of New York, Stony Brook, NY 11794-4400, USA; email: bender@cs.sunysb.edu. Supported in part by Sandia National Laboratories and the National Science Foundation grants EIA-0112849 and CCR-0208670.

[§]MIT Laboratory for Computer Science, 200 Technology Square, Cambridge, MA 02139, USA; email: edemaine@mit.edu. Supported in part by NSF Grant EIA-0112849.

[¶]Department of Computer Science, Rutgers University, Piscataway, NJ 08855, USA; email: farach@cs.rutgers.edu. Supported in part by NSF Grant CCR-9820879.

^{||}School of Computer Science, University of Waterloo, Waterloo, Ontario N2L 3G1, Canada; email: imunro@uwaterloo.ca.

^{**}AT&T Labs—Research, Shannon Laboratory, Florham Park, NJ 07932, USA; email: mthorup@research.att.com.

Consequently, there is a large body of work on optimizing search trees for nonuniform distributions in a variety of contexts:

1. *Known distribution on a RAM* — optimal binary search trees [1, 18] and variations [15], and Huffman codes [16].
2. *Unknown distribution on a RAM* — splay trees [17, 21].
3. *Known distribution in external memory* — optimal binary search trees in the HMM model [22].
4. *Unknown distribution in external memory* — alternatives to splay trees [17].¹

Fixed Tree Topology. Search trees frequently encode decision trees that cannot be rebalanced because the operations lack associativity. Such trees naturally arise in the context of string or geometric data, where each node represents a character in the string or a geometric predicate on the data. Examples of such structures include tries, suffix trees, Cartesian trees, k-d trees and other BSP trees, quadtrees, etc. Almost always their contents are not uniformly distributed, and often these search trees are unbalanced.

The first question is how to optimize these fixed-topology trees when the access distribution is known. On a RAM there is nothing to optimize because there is nothing to vary. In external memory, however, we can choose the layout of the tree structure in memory, that is, which nodes of the tree are stored in which blocks in memory. This problem was first considered by Gil and Itai [14]. Among other results described below, they presented a dynamic-programming algorithm for optimizing the partition of the N nodes into blocks of size B , given the probability distribution on the leaves. The algorithm runs in $O(NB^2 \lg \Delta)$ time, where Δ is the maximum degree of a node, and uses $O(B \lg N)$ space.

This problem brings up important issues in external memory algorithms because when trees are unbalanced or distributions are skewed, there is even more advantage to a good layout. Whereas uniform distributions lead to B-trees, which save a factor of only $\lg B$ over standard $(\lg N)$ -time balanced binary trees, the savings grow with nonuniformity in the tree. In the extreme case of a linear-height tree or a very skewed distribution we obtain a factor of B savings over a naïve memory layout.

Cache-Oblivious Algorithms. Recently, there has been a surge of interest in data structures for multilevel memory hierarchies. Frigo, Leiserson, Prokop, and Ramachandran [12, 19] introduced the notion of *cache-oblivious algorithms*, which have asymptotically optimal memory performance for all possible values of the memory-hierarchy parameters (block size and memory-level size). As a consequence, such algorithms tune automatically to arbitrary memory hierarchies with an arbitrarily many memory levels. Examples of cache-oblivious *data structures* include cache-oblivious B-trees [6] and its simplifications [8, 9, 20], cache-oblivious persistent trees [5], cache-oblivious priority queues [2], and cache-oblivious linked lists [4]. However, all of these data structures assume a uniform distribution on operations.

¹Although [17] does not explicitly state its results in the external-memory model, its approach easily applies to this scenario.

Our Results. In the conference version of this paper [7], two claims were made: that a simple greedy algorithm is within an additive 1 of optimal for known block size, and that a particular recursive greedy partition of the tree yields a cache-oblivious algorithm within a multiplicative $O(1)$ of optimal. Both these claims turn out to be false, and their corrected versions appear below.

1. We develop a general technique called Split-and-Refine for converting a family of layouts with known block size into a layout with unknown block size, while increasing the expected block cost by at most a constant factor.
2. In addition, we show how to adapt this technique to other objective functions, specifically, minimizing the maximum block cost.
3. In particular, combining Split-and-Refine with known results, we obtain a cache-oblivious layout algorithm whose performance is within a constant factor of optimal.
4. We analyze two natural greedy algorithms for tree layout with known block size. We show that their performance can be as bad as a factor of $\Omega(\lg B / \lg \lg B)$ away from optimal, but is no more than $O(\lg B)$ away from optimal.

Related Work. In addition to the result mentioned above, Gil and Itai [13, 14] consider other algorithmic questions on tree layouts. They prove that minimizing the number of distinct blocks visited by each query is equivalent to minimizing the number of block transfers over several queries; in other words, caching blocks over multiple queries does not change the optimal solutions. Gil and Itai also consider the situation in which the total number of blocks must be minimized (the external store is expensive) and prove that optimizing the tree layout subject to this constraint is NP-hard. In contrast, with the same constraint, it is possible to optimize the expected query cost within an additive $1/2$ in $O(N \lg N)$ time and $O(B)$ space. This algorithm is a variant of their polynomial-time dynamic program for the unconstrained problem.

Clark and Munro [11, 10] consider a worst-case version of the problem in which the goal is to minimize the maximum number of block transfers over all queries, instead of minimizing the expected number of block transfers. They show how the exact optimal layout can be computed in $O(N)$ time for a known block size B . We show how to extend this result to the cache-oblivious setting.

2 Basics

We begin with a few definitions and structural properties.

In the *tree layout* (or *trie layout*) problem, we are given a fixed-topology tree with a known probability distribution on the leaves. When the layout algorithm knows the memory block size B , the goal is to produce a *layout*, which clusters nodes into memory blocks. In the cache-oblivious model, where the layout algorithm does not know B , the goal is to determine a *cache-oblivious layout*, which specifies an order of the nodes in memory. A cache-oblivious

layout must be efficient no matter how the ordering of nodes is partitioned into consecutive blocks of size B , and for all values of B .

The *expected block cost* of a layout is the expected number of blocks along the root-to-leaf path for a randomly selected leaf. The *optimal layout with block size B* minimizes the expected block cost over all possible layouts with block size B .

A simple but useful idea is to propagate the probability distribution on leaves up to internal nodes. Define the *probability of an internal node* to be the probability that the node is on a root-to-leaf path for a randomly chosen leaf; that is, the probability of an internal node is the sum of the probabilities of the leaves in its subtree. These probabilities can be computed and stored at the nodes in linear time.

All of our algorithms are based on the following structural lemmas:

Lemma 1 (Convexity Lemma [14]) *Any fixed-topology tree has an optimal layout that is convex, i.e., in which every block forms a connected subtree of the original tree.*

Lemma 2 (Monotonicity Lemma) *The expected search cost for the optimal tree layout with block size $B - 1$ is no greater than for the optimal tree layout with block size B .*

Proof: The blocks of an optimal layout with block size B need not be full. In particular, each node could store just $B - 1$ elements. \square

Lemma 3 (Smoothness Lemma) *An optimal tree layout with block size $B/2$ has an expected search cost of no more than twice that of an optimal layout with block size B .*

Proof: Consider an optimal convex layout with block size B . Partition each block into two pieces of size $B/2$, where one piece is an arbitrary subtree containing the root of the block, and the other piece consists of the forest of remaining nodes of the block. This new layout has at most twice as many block transfers as the optimal layout with block size B , and so the optimal layout with block size $B/2$ does at least this well. \square

3 Performance of Greedy

In this section we analyze the performance of two greedy heuristics for tree layout with known block size B . Despite the natural appeal of these algorithms, we will show that their performance is far from optimal, by a roughly $\lg B$ factor.

The (most natural) *Weight-Greedy* heuristic incrementally grows a *root block*, that is, the block containing the root of the tree. Initially, the root block contains just the root node; then the heuristic repeatedly adds the maximum-probability node not already in the root block (which is necessarily adjacent to the subtree so far). When the root block fills, i.e., contains B nodes, the heuristic conceptually removes its nodes from the tree, and lays out the remainder recursively.

The *DFS-Greedy* heuristic orders the nodes according to a locally greedy depth-first search, and then partitions the nodes into consecutive blocks of size B . More precisely, the

nodes are ordered according to when they are traversed by a depth-first search that visits the children of a node by decreasing probability. This layout is not necessarily convex.

The conference version of this paper [7] falsely claimed that the layout from the Weight-Greedy heuristic is within an additive 1 of optimal. We rectify this result by (1) demonstrating a tree in which the performance ratio is $\Omega(\lg B / \lg \lg B)$, and (2) proving a complementary near-matching upper bound of $O(\lg B)$.

3.1 Lower Bound of $\Omega(\lg B / \lg \lg B)$

Theorem 1 *The expected block cost of either greedy layout (Weight-Greedy or DFS-Greedy) can be a factor of $\Omega(\lg B / \lg \lg B)$ more than the expected block cost of the optimal layout.*

Proof: We exhibit a tree T with $N \geq B^2$ nodes, for which the expected greedy block cost is $\Theta(\lg B / \lg \lg B)$ and the optimal block cost is $\Theta(1)$. Then we show how to replicate this tree so that the expected greedy block cost is $\Theta(\lg N / \lg B)$ and the optimal block cost is $\Theta(\lg N / \lg \lg B)$.

We build tree T by starting with a complete B -node tree T' with fanout $\Theta(\lg B)$ and height $\Theta(\lg B / \lg \lg B)$, in which all nodes at each level have equal probability. We next augment tree T' by attaching an *escape path* of length B to every node. The probability of the first node on an escape path is slightly higher than the probability of its sibling nodes. (For example, we can let the $\lg B$ children in T' of a parent node have probability $1/(2 + \lg B)$ times the probability of the parent, and then we let the escape path have probability $2/(2 + \lg B)$ times the probability of the parent.) Thus, greedy favors all nodes along the escape paths instead of any of the adjacent children. This construction of T has $\Theta(B^2)$ nodes.

The optimal layout assigns the original tree T' to its own block, and assigns each escape path to its own block. Because there are $\Theta(B^2)$ nodes, there are $\Theta(B)$ blocks. The expected search cost is 2.

Greedy performs worse. Because each escape path has length B , and because greedy favors the escape paths, each node in the tree T' is in its own block. Thus, a search pays a cost of 1 for each level in T' visited before the search chooses an escape path and leaves T' . The probability that a random search reaches the leaves in T' is $(1 - 1/\lg B)^{\Theta(\lg B / \lg \lg B)} \approx (1/e)^{1/\lg \lg B}$, which is at least $1/e$ (and in fact much closer to 1). Thus, at least a constant fraction of searches reach the bottom of the tree, visiting one block for each of $\Theta(\lg B / \lg \lg B)$ levels.

In summary, for a tree T with $\Theta(B^2)$ nodes, optimal has expected block cost of 2, whereas greedy has expected block cost $\Theta(\lg B / \lg \lg B)$.

We can grow tree T to have arbitrary many nodes by splitting each leaf into many leaves all with same parent. This change can only affect the costs of greedy and optimal by at most 1. Thus, the ratios between the costs remains $\Theta(\lg B / \lg \lg B)$, but the costs do not grow with N .

Alternatively, we can replicate the tree T so that the block costs grow with N . Specifically, we attach another copy of T to the end of each escape path, and we repeat this process any desired number of times. Each iteration increases the size of T by roughly a

factor of B . The result is that the optimal expected search cost is $S(N) = S(N/B) + \Theta(1) = \Theta(\lg N / \lg B)$, whereas the greedy expected search cost is $S(N) = S(N/B) + \Theta(\lg B / \lg \lg B) = \Theta(\lg N / \lg \lg B)$. \square

3.2 Upper Bound of $O(\lg B)$

We now prove that no tree is much worse than the class of examples in the previous section:

Theorem 2 *The expected block cost of either greedy layout (Weight-Greedy or DFS-Greedy) is within a factor of $O(\lg B)$ of optimal.*

To simplify boundary cases, we preprocess the tree as follows: to each leaf, we attach $2B$ children called *hair nodes*. The probability of each hair node is $1/(2B)$ of the probability of the leaf. This preprocessing only increases the expected search cost in greedy, and it increases the optimal search cost by at most 1 (because we could put each hair node in its own block, and the optimal solution can only be better).

We partition the tree into *chunks* as follows, and treat each chunk separately. We grow the first chunk starting at the root. The *probability of a chunk* is the probability of entering its root; for the first chunk, this probability is 1. Define the *relative probability* of a node in a chunk to be the probability of entering that node divided by the probability of the chunk. We grow a chunk by repeating the following process: for every node in the chunk that has relative probability more than $1/2B$, we add the children of that node to the chunk. When this process terminates, we have completed a chunk; we conceptually remove those nodes and recurse on the remainder of the tree to complete the partition into chunks.

As a postprocessing step, we demote from chunkhood any chunk that consists solely of a hair node, leaving that hair node separate from all chunks. The reason for the careful handling of hair nodes is to ensure that all of the leaves of a chunk have relative probability at most $1/2B$. Furthermore, the parent of each leaf of a chunk, which we call a *twig node*, has relative probability more than $1/2B$.

We prove two lower bounds on the optimal block partition, and one upper bound on the greedy block partition.

Claim 1 (Lower Bound 1) *Consider the optimal block partition of the tree. The number of blocks along a path from the root of the tree to a leaf of the tree is at least the length of that path divided by B .*

Proof: We need to visit every node along the path from the root of the tree to the leaf of the tree. If there are P nodes along the path, the best we could hope for is that every block contains B of these P nodes along the path. Thus, the number of blocks visited is at least $\lceil P/B \rceil$. This lower bound can only be larger than P/B , which is the claimed bound. \square

The next lower bound first considers each chunk separately, and then combines the chunk-by-chunk bounds to apply to the entire tree. If we counted the number of blocks along a root-to-leaf path separately for each chunk, and then added these counts together, we might double-count blocks because a single block can overlap multiple chunks. Instead, we count

the number of *block transitions* along a path, i.e., the number of traversed edges that straddle two blocks, which is exactly 1 smaller than the number of blocks along that path. Now we can count the number of block transitions separately for the subpath within each chunk, then add these counts together, and the resulting total only underestimates the true number of block transitions.

Claim 2 (Lower Bound 2 Within Chunk) *Consider the optimal block partition of the tree. Within any chunk, the expected number of block transitions from the root of the chunk to a leaf of the chunk is at least $1/2$.*

Proof: The memory block containing the root of the chunk has at most B nodes, so it can contain at most B leaves of the chunk. Each leaf of the chunk has relative probability at most $1/2B$, so B leaves of the chunk have total relative probability at most $1/2$. Thus, the leaves of the chunk that do not fit in the root block of the chunk have a total relative probability of at least $1/2$. For these leaves, the number of block transitions within the chunk is at least 1. Therefore, the expected number of block transitions within the chunk is at least $1/2$. \square

Now we combine the estimates at each chunk to form a bound on the entire tree. Define the *expected chunk count* to be the expected number of chunks along a root-to-leaf path.

Corollary 3 (Lower Bound 2) *Consider the optimal block partition of the tree. The expected number of blocks along a path from the root of the tree to a leaf of the tree is at least half the expected chunk count.*

Proof: Label the chunks from 1 to k . Consider a random root-to-leaf path. Define random variable X_i to be the number of block transitions along the path that occur within chunk i , i.e., the number of block transitions along the subpath entirely contained in chunk i . If the path does not visit chunk i , then X_i is 0. Conditioned upon the path visiting chunk i , X_i is at least $1/2$, by Claim 2.

Define random variable X to be $\sum_{i=1}^k X_i$, which counts all block transitions strictly within chunks, but ignores block transitions that align with chunk boundaries. Thus, X is a lower bound on the number of block transitions along the path. By linearity of expectation, $E[X]$ equals $\sum_{i=1}^k E[X_i]$. By the argument above, $E[X_i]$ is at least $1/2$ times the probability of entering chunk i . Thus, $E[X]$ is at least half the expected chunk count.

This lower bound ignores any additional cost potentially caused by hair nodes that do not belong to chunks, which would only further increase the lower bound. \square

Now we establish an upper bound on either greedy block partition.

Claim 4 (Upper Bound Within Chunk) *Consider either greedy block partition of the tree (Weight-Greedy or DFS-Greedy). Within any chunk, the number of blocks along a path from the root of the chunk to a leaf of the chunk is at most the length of that path divided by B , plus $2 \lg B + 7$.*

Proof: We divide the chunk into *strata* based on the nearest power of two of the relative probability of a node. More precisely, define *stratum* i of the chunk to contain the nodes

with relative probability at most $1/2^i$ and more than $1/2^{i+1}$. We consider strata $0, 1, 2, \dots, \lceil \lg 2B \rceil$. Thus, there are $1 + \lceil \lg 2B \rceil \leq \lg B + 3$ strata. Leaf nodes may have sufficiently small probability to be excluded from all strata. However, the strata partition all nonleaf nodes of the chunk (i.e., down to the twig nodes).

We claim that each stratum is a vertex-disjoint union of paths. Consider any node in stratum i , which by definition has relative probability at most $1/2^i$. At most one child of that node can have relative probability more than $1/2^{i+1}$. Thus, at most one child of each node in stratum i is also in stratum i , so each connected component of stratum i is a path. We call these connected components *subpaths*.

Any path from the root of the chunk to a leaf of the chunk starts in stratum 0, and visits some of the strata in increasing order. The path never revisits a stratum that it already left, because the probabilities of the nodes along the path are monotonically decreasing. Thus, the path passes through at most $\lg B + 3$ strata.

When the greedy heuristic adds a node from a particular stratum to a block, it continues adding all nodes from that stratum to the block, until either the block fills or the subpath within the stratum is exhausted. Thus, after an initial startup of at most B nodes in a subpath, greedy clusters the remainder of the subpath into full blocks of size B . Because of potential roundoff at the top and the bottom of this subpath, the number of blocks along this subpath is at most the length of the subpath divided by B , plus 2. In addition, the leaves of the chunk may not belong to any stratum, so we may visit one additional block after visiting the strata. Summing over all strata, the total number of blocks along the path from the root of the chunk to a leaf of the chunk is at most the length of the path divided by B , plus $2 \lg B + 6$ (2 per stratum), plus 1 (for a leaf). \square

Corollary 5 (Upper Bound) *Consider either greedy block partition of the tree (Weight-Greedy or DFS-Greedy). The number of blocks along a path from the root of the tree to a leaf of the tree is at most the length of that path divided by B , plus $2 \lg B + 8$ times the expected chunk count.*

Proof: We follow the same proof outline as Corollary 3 except for three difference. First, we plug in Claim 4 instead of Claim 2. Second, before we counted the number of block transitions along the path, to avoid over-counting for a lower bound, whereas here we count the number of blocks along a path, to avoid under-counting for an upper bound. Third, we add an additional 1 to the bound because of potential hair nodes separate from all chunks. \square

Finally, we combine the lower bounds in Lemma 1 and Corollary 3, and the upper bound in Corollary 5, to prove Theorem 2:

Proof of Theorem 2: By Corollary 5, the expected number of blocks along a root-to-leaf path is at most the expected path length divided by B , plus $2 \lg B + 8$ times the expected chunk count. By Lemma 1, this expected cost is at most the optimal expected search cost, plus $2 \lg B + 8$ times the expected chunk count. Furthermore, by Corollary 3, the optimal expected search cost is at least half the expected chunk count. Therefore, the ratio of the greedy expected search cost over the optimal expected search cost is at most $1 + 2(2 \lg B + 8)$. That is, greedy performs within a factor of $4 \lg B + 17$ from optimal. \square

4 Cache-Oblivious Layout

In this section, we develop a cache-oblivious layout whose expected block cost is within a constant factor of the optimal layout. This result re-establishes the main claim of the conference version of this paper [7].

Theorem 3 *The Split-and-Refine algorithm produces a cache-oblivious layout whose expected block cost is within a constant multiplicative factor of optimal. The algorithm can use any given black box that lays out a tree within a constant factor of the optimal expected block cost for a known block size B . If the running time of the black box is $T(N, B) = \Omega(N)$, then the Split-and-Refine running time is $O(\sum_{l=0}^{\lceil \lg N \rceil} T(N, 2^l))$.*

In fact, our technique is quite general, and in addition to computing a layout that minimizes the expected block cost, it can compute a layout that minimizes the *maximum* block cost.

Theorem 4 *The Split-and-Refine algorithm produces a cache-oblivious layout whose maximum block cost is within a constant multiplicative factor of optimal. The algorithm can use any given black box that lays out a tree within a constant factor of the optimal maximum block cost for a known block size B . The running times are as in Theorem 3.*

4.1 Split-and-Refine Algorithm

The basic idea of the cache-oblivious layout is to recursively combine optimal layouts for several carefully chosen block sizes. These layouts are computed independently, and the block sizes are chosen so that costs of the layouts grow exponentially. The layouts may be radically different; all we know is their order from coarser (larger B) to finer (smaller B). To create a recursive structure among these layouts, we further partition each layout to be consistent with all coarser layouts. Then we store the tree according to this recursive structure.

More precisely, our cache-oblivious layout algorithm works as follows. For efficiency, we only consider block sizes that are powers of two.² We begin with a block size of the smallest power of two that is at least N , i.e., the *hyperceiling* $\lceil\lceil N \rceil\rceil = 2^{\lceil \lg N \rceil}$. Because this block size is at least N , the optimal partition places all nodes in a single block, and the expected (and worst-case) search cost is 1. We call this trivial partition *level of detail 0*. Now we repeatedly halve the current block size, and at each step we compute the optimal partition. We stop when we reach the coarsest partition whose expected search cost is between 2 and 4; such a partition exists by the Smoothness Lemma (Lemma 3). We call this partition *level of detail 1*. Then we continue halving the current block size, until we reach the coarsest partition whose expected search cost is between 2 and 4 times the expected search cost at level of detail 1. This partition defines *level of detail 2*. We proceed in defining levels of detail until we reach a level of detail ℓ whose block size is 1. In contrast to all other levels,

²Restricting to powers of two speeds up the layout algorithm; if we did not care about speed, we could consider all possible values of B .

the expected search cost at level of detail ℓ may be less than a factor of 2 larger than level of detail $\ell - 1$.

The levels of detail are inconsistent in the sense that a block at one level of detail may not be wholly contained in a block at any coarser level of detail. To define the layout recursively, we require the blocks to form a *recursive structure*: a block at one level of detail should be made up of *subblocks* at the next finer level of detail. To achieve this property, we define the *refined level of detail i* to be the refinement of the partition at level of detail i according to the partitions of all coarser levels of detail $< i$. That is, if two nodes are in different blocks at level of detail i , then we separate them into different blocks at all finer levels of detail $> i$.

The recursive structure allows us to build a recursive layout as follows. Each block at any refined level of detail is stored in a contiguous segment of memory. The subblocks of a block can be stored in any order as long as they are stored contiguously.

4.2 Running Time

We can compute the partition at each level of detail within the claimed running time because we call the black box precisely for block sizes 2^l where $l = 0, 1, \dots, \lceil \lg N \rceil$. Each call to the black box produces a partition on the N nodes, which we represent by arbitrarily assigning each block a unique integer between 1 and N , and labeling each node with the integer assigned to the block containing it.

From these unrefined partitions, it is easy to compute the cache-oblivious layout. To each node we assign a *key* consisting of at most $\lceil \lg N \rceil + 1$ components, assembled from the labels from all levels of detail, where the coarsest level of detail specifies the most significant component of the key. Then we sort the nodes according to their keys in $O(N \lg N)$ time using a radix sort, and lay out the nodes in this order. This layout automatically adheres to the refined levels of detail, without having to compute them explicitly.

4.3 Expected Block Cost

We begin by analyzing the cost of the unrefined levels of detail. Define the random variable X_i to be the number of blocks along a random root-to-leaf path in the partition defined by level of detail i . Thus, $E[X_i]$ is the expected search cost at level of detail i , as above. By construction, $E[X_{i+1}]$ is a factor between 2 and 4 larger than $E[X_i]$.

Let B be the true block size of the cache, not known to the algorithm. We focus on two levels of detail: level of detail $L - 1$ whose block size is at least B , and level of detail L whose block size is at most B . By the Monotonicity Lemma (Lemma 2), the ideal optimal partition with block size B has expected search cost at least $E[X_{L-1}]$, because the block size at level of detail $L - 1$ is only larger than B . By construction, $E[X_L]$ is at most 4 times larger than $E[X_{L-1}]$, and thus is at most 4 times larger than the optimal partition with block size B .

Consider the partition defined by level of detail L , which is designed for blocks smaller than B , but laid out in a memory with block size B . Each block in this partition has size at most B , and hence is contained in at most two memory blocks of size B , depending on alignment. Thus, the expected search cost measured according to block size B is at most 8 times the optimal partition with block size B .

It remains to consider the additional cost introduced by refining level of detail L by all coarser levels of detail. Call an edge of the tree *straddling at level i* if its endpoints lie in different blocks at level of detail i . Connecting to the previous analysis, X_i is 1 plus the number of straddling edges at level i along a random root-to-leaf path. The important property of straddling edges is this: along a root-to-leaf path, the straddling edges at level i count the number of extra memory transfers (block refinements) caused by refining level of detail L according to coarser level of detail i . Thus, an edge spans two different refined blocks in the refined level of detail L precisely when the edge is straddling at some level of detail $\leq L$.

To capture these properties algebraically, define the random variable X to be the number of blocks along a random root-to-leaf path in the partition defined by the refined level of detail L . Because X counts X_L (the cost of the unrefined level of detail) as well as the extra memory transfers caused by straddling edges at levels of detail $< L$, we have the following equation:

$$\begin{aligned} X &= (X_1 - 1) + (X_2 - 1) + \cdots + (X_{L-1} - 1) + X_L \\ &= X_1 + X_2 + \cdots + X_{L-1} + X_L - (L - 1). \end{aligned}$$

Now we want to compute $E[X]$, that is, the expected search cost of the partition defined by the refined level of detail L , ignoring block alignment. By linearity of expectation,

$$\begin{aligned} E[X] &= E[X_1 + X_2 + \cdots + X_L - (L - 1)] \\ &= E[X_1] + E[X_2] + \cdots + E[X_L] - (L - 1) \\ &\leq E[X_1] + E[X_2] + \cdots + E[X_L]. \end{aligned}$$

As mentioned above, each $E[X_i]$ is a factor between 2 and 4 more than $E[X_{i-1}]$, so the series is geometric. Thus,

$$E[X] \leq 2E[X_L].$$

But we already argued that $E[X_L]$ is at most 4 times larger than $E[X_{L-1}]$, and that $E[X_{L-1}]$ is at most the expected search cost of the ideal optimal partition with block size B . Therefore, $E[X]$ is at most 8 times this ideal cost, so the refined level of detail L has an expected number of memory transfers that is at most 16 times the optimal. This concludes the proof of Theorem 3.

4.4 Minimizing the Maximum Block Cost

Clark and Munro [11, 10] consider an analogous tree-layout problem in which the objective is to minimize the maximum block cost (which is independent of any probability distribution of the leaves). They give a simple greedy algorithm that computes the exact optimal layout in $O(N)$ time for a known block size B . The basic idea behind the algorithm is to build a partition bottom-up, starting with the leaves in their own blocks, and merging blocks locally as much as possible while giving priority to the most expensive subtrees (those with highest maximum block cost).

We can use this layout algorithm for known B as an alternative subroutine in the cache-oblivious Split-and-Refine algorithm. The algorithm behaves correctly as before because of

an analogous Smoothness Lemma for the metric of maximum block cost (exactly the same proof of Lemma 3 applies). The analysis also proceeds as before, except that we replace X_i and $E[X_i]$ by Y_i , which is defined as the maximum number of blocks along any root-to-leaf path in the block partition defined by level of detail i . There is no longer any randomization or expectation, but otherwise the proof is identical.

5 Conclusion

In this paper, we developed cache-oblivious layouts of fixed-topology trees, whose performance is within a constant factor of the optimal layout with known block size. The running time of the layout algorithm is dominated by $O(\lg N)$ calls to any given layout algorithm for known block size. We also showed that two natural greedy strategies have a performance ratio between $\Omega(\lg B / \lg \lg B)$ and $O(\lg B)$ when compared to the optimal layout.

The main open problems are to what extent fixed-topology tree layouts can be made dynamic and/or self-adjusting, both with known block size and in the cache-oblivious setting. By *dynamic* we mean the ability to add and delete leaves, and to redistribute probability from leaf to leaf. A *self-adjusting* data structure would adapt (as in splay trees) to an online search distribution without a priori knowledge of the distribution.

Acknowledgments

We thank Stefan Langerman for many helpful discussions.

References

- [1] Alfred V. Aho, John E. Hopcroft, and Jeffrey D. Ullman. *The Design and Analysis of Computer Algorithms*. Addison-Wesley, 1974.
- [2] Lars Arge, Michael A. Bender, Erik D. Demaine, Bryan Holland-Minkley, and J. Ian Munro. Cache-oblivious priority queue and graph algorithm applications. In *Proceedings of the 34th Annual ACM Symposium on Theory of Computing*, pages 268–276, Montréal, Canada, May 2002.
- [3] Rudolf Bayer and Edward M. McCreight. Organization and maintenance of large ordered indexes. *Acta Informatica*, 1(3):173–189, February 1972.
- [4] Michael A. Bender, Richard Cole, Erik D. Demaine, and Martin Farach-Colton. Scanning and traversing: Maintaining data for traversals in a memory hierarchy. In *Proceedings of the 10th Annual European Symposium on Algorithms*, volume 2461 of *Lecture Notes in Computer Science*, pages 139–151, Rome, Italy, September 2002.
- [5] Michael A. Bender, Richard Cole, and Rajeev Raman. Exponential structures for efficient cache-oblivious algorithms. In *Proceedings of the 29th International Colloquium*

- on Automata, Languages and Programming*, volume 2380 of *Lecture Notes in Computer Science*, pages 195–207, Málaga, Spain, July 2002.
- [6] Michael A. Bender, Erik D. Demaine, and Martin Farach-Colton. Cache-oblivious b-trees. In *Proceedings of the 41st Annual Symposium on Foundations of Computer Science*, pages 399–409, Redondo Beach, California, November 2000.
 - [7] Michael A. Bender, Erik D. Demaine, and Martin Farach-Colton. Efficient tree layout in a multilevel memory hierarchy. In *Proceedings of the 10th Annual European Symposium on Algorithms*, volume 2461 of *Lecture Notes in Computer Science*, pages 165–173, Rome, Italy, September 2002.
 - [8] Michael A. Bender, Ziyang Duan, John Iacono, and Jing Wu. A locality-preserving cache-oblivious dynamic dictionary. In *Proceedings of the 13th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 29–38, San Francisco, California, January 2002.
 - [9] Gerth Stølting Brodal, Rolf Fagerberg, and Riko Jacob. Cache oblivious search trees via binary trees of small height. In *Proceedings of the 13th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 39–48, San Francisco, California, January 2002.
 - [10] David Clark. *Compact Pat Trees*. PhD thesis, University of Waterloo, 1996.
 - [11] David R. Clark and J. Ian Munro. Efficient suffix trees on secondary storage. In *Proceedings of the 7th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 383–391, Atlanta, January 1996.
 - [12] Matteo Frigo, Charles E. Leiserson, Harald Prokop, and Sridhar Ramachandran. Cache-oblivious algorithms. In *Proceedings of the 40th Annual Symposium on Foundations of Computer Science*, pages 285–297, New York, October 1999.
 - [13] Joseph Gil and Alon Itai. Packing trees. In *Proceedings of the 3rd Annual European Symposium on Algorithms (ESA)*, pages 113–127, 1995.
 - [14] Joseph Gil and Alon Itai. How to pack trees. *Journal of Algorithms*, 32(2):108–132, 1999.
 - [15] T. C. Hu and A. C. Tucker. Optimal computer search trees and variable-length alphabetic codes. *SIAM Journal on Applied Mathematics*, 21(4):514–532, December 1971.
 - [16] David A. Huffman. A method for the construction of minimum-redundancy codes. *Proceedings of the IRE*, 40(9):1098–1101, 1952.
 - [17] John Iacono. Alternatives to splay trees with $O(\lg n)$ worst-case access times. In *Proceedings of the 12th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 516–522, Washington, D.C., January 2001.
 - [18] Donald E. Knuth. *The Art of Computer Programming*, volume 3 (Sorting and Searching). Addison-Wesley, 1968.

- [19] Harald Prokop. Cache-oblivious algorithms. Master's thesis, Massachusetts Institute of Technology, Cambridge, MA, June 1999.
- [20] Naila Rahman, Richard Cole, and Rajeev Raman. Optimised predecessor data structures for internal memory. In *Proceedings of the 5th International Workshop on Algorithm Engineering*, volume 2141 of *Lecture Notes in Computer Science*, pages 67–78, Aarhus, Denmark, August 2001.
- [21] Daniel Dominic Sleator and Robert Endre Tarjan. Self-adjusting binary search trees. *Journal of the ACM*, 32(3):652–686, July 1985.
- [22] Shripad Thite. Optimum binary search trees on the hierarchical memory model. Master's thesis, Department of Computer Science, University of Illinois at Urbana-Champaign, 2001.