On a Model of Indexability and Its Bounds for Range Queries

JOSEPH M. HELLERSTEIN

University of California, Berkeley, Berkeley, California

ELIAS KOUTSOUPIAS

University of California, Los Angeles, Los Angeles, California

DANIEL P. MIRANKER

University of Texas at Austin, Austin, Texas

CHRISTOS H. PAPADIMITRIOU

University of California, Berkeley, Berkeley, California

AND

VASILIS SAMOLADAS

University of Texas at Austin, Austin, Texas

Abstract. We develop a theoretical framework to characterize the hardness of indexing data sets on block-access memory devices like hard disks. We define an indexing workload by a data set and a set of potential queries. For a workload, we can construct an indexing scheme, which is a collection of fixed-sized subsets of the data. We identify two measures of efficiency for an indexing scheme on a workload: storage redundancy, r (how many times each item in the data set is stored), and access overhead, A (how many times more blocks than necessary does a query retrieve).

The work of J. M. Hellerstein was supported by National Science Foundation (NSF) grant IRI-9703972, NASA grant FDNAGW-5198, and a Sloan Foundation Fellowship.

The work of E. Koutsoupias was supported by NSF grant CCR-9521606 and CCR-0105752.

The work of C. H. Papadimitriou was supported by NSF grant CCR-9820897 and by an IBM Faculty Development Award.

Authors' addresses: J. M. Hellerstein and C. H. Papadimitriou, Computer Science Division, EECS Department, 387 Soda Hall #1776, University of California, Berkeley, Berkeley, CA 94720-1776, e-mail: jmh@cs.berkeley.edu and christos@cs.berkeley.edu; E. Koutsoupias, Computer Science Department, 3731J Boelter Hall, University of California, Los Angeles, Los Angeles, CA 90095-1596, e-mail: elias@cs.ucla.edu; D. P. Miranker and V. Samoladas, Department of Computer Sciences, 2.124 Taylor Hall, University of Texas at Austin, Austin, TX 78712, e-mail: miranker@cs.utexas.edu and vsam@cs.utexas.edu.

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 direct 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 worked 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., 1515 Broadway, New York, NY 10036 USA, fax +1 (212) 869-0481, or permissions@acm.org. © 2002 ACM 0004-5411/02/0100-0035 \$5.00

Journal of the ACM, Vol. 49, No. 1, January 2002, pp. 35-55.

For many interesting families of workloads, there exists a trade-off between storage redundancy and access overhead. Given a desired access overhead A, there is a minimum redundancy that any indexing scheme must exhibit. We prove a lower-bound theorem for deriving the minimum redundancy. By applying this theorem, we show interesting upper and lower bounds and trade-offs between A and r in the case of multidimensional range queries and set queries.

Categories and Subject Descriptors: H.2.2 [Database Management]: Physical Design-access methods

General Terms: Theory

Additional Key Words and Phrases: Database, index, indexability, lower bounds, multidimensional, query, redundancy

1. Introduction

Upon its definition, the B-tree promptly proved to be an effective access method for the primary applications of relational databases [Bayer and McCreight 1972]. The success and ubiquity of the relational data model arguably owes much to the timely definition of the B-tree. Since then, a major thrust of database research has been to extend the relational model and relational systems to manage more complex types and more expressive query languages. The B-tree is widely recognized to be an inadequate data structure in many of the novel contexts, and no single, general-purpose successor has emerged to enable the diversity of applications and requirements for contemporary information systems. Therefore, it is important to develop general methodologies and tools for the design of new indexing methods, as well as mathematical tools, to evaluate their performance and identify their limitations, a priori.

A systems approach to the "generalized indexing" problem has been proposed and implemented [Hellerstein et al. 1995; Kornacker et al. 1997; Aoki 1998; Kornacker 1999]. The results highlighted the need for theoretical tools to rigorously analyze indexing problems. To aid developers of new indexes in this general framework, a kind of *theory of indexability* is required: a mathematical model that allows the performance and scalability of an indexing scheme to be evaluated much as complexity theory is used to evaluate algorithms. Where complexity theory considers in-memory data structures, a theory of indexability must consider the impact of disk-based *secondary storage*.

Our pragmatic results focus on the multidimensional range search problem, a common workload for many advanced applications. An enormous amount of experimental research has been devoted to this problem: a recent survey cites over 50 different multidimensional data structures [Gaede and Günther 1998]. Many commercial vendors of Object-Relational Database Systems and Geographic Information Systems use one of these structures, typically some variant of the R-tree [Guttman 1984], the Grid File [Nievergelt et al. 1984], or disk-resident adaptations of the quad-tree [Samet 1989]. This research is primarily experimental. Analytic research on these structures has concentrated on probabilistic and empirical studies of their average-case performance, under various data and query distributions.

At the same period of time that heuristic disk-based indices such as the R-Tree were introduced, the computational geometry community was studying mainmemory data structures for range searching, paying little attention to secondary memory. In contrast to most multidimensional indexing research by the database community, the work in computational geometry is mostly theoretically oriented,

with an emphasis on worst-case asymptotic performance. We believe that the striking contrast between these two approaches to the same problem arose from a fundamental fact: for two-dimensional range searching (and more so for higher dimensions), optimal query cost cannot be achieved with space proportional to the data set, but instead requires significant storage redundancy, typically by a mul*tiplicative factor* at least logarithmic to the size of the data set [Chazelle 1990a]. Of course, redundancy has often been used in databases to accelerate performance: index structures are themselves typically redundant to the data sets they index, and the addition of logarithmic space is standard for upper levels in search trees. However, the space cost of redundancy in databases has rarely been as high as a logarithmic *multiple* to the size of the data set. This is only reasonable: databases usually store very big data sets, on top of which a logarithmic factor of redundancy makes the solution considerably more expensive in space. Also, high redundancy increases the I/O cost of online updates, at least proportionally to the redundancy. For these reasons, low-redundancy access methods are typically used in practice [Kornacker 1999; Kanth et al. 1999].¹

Thus, database research concentrated on data structures with low redundancy, with very bad worst-case behavior, but with the hope of reasonable average-case behavior for real workloads. Lowering the observed average-case cost has typically been achieved through problem-specific heuristics, which take into account the particularities of various data sets and query workloads. For example, access methods for two-dimensional geographic data have been differentiated from access methods for temporal queries over one-dimensional data, by a choice of heuristics appropriate to the expected distributions of respective data sets and typical queries, despite the fact that from a conceptual point of view the two problems are equivalent.

The two approaches presented above (ad hoc and application-dependent indexing schemes versus highly redundant computational geometry data structures) are in a sense extreme: one penalizes worst-case query performance by keeping space linear, the other strongly favors query performance without regard to the storage cost becoming prohibitive. The results in this paper strive to reconcile these two approaches by exposing and studying the fundamental trade-off between I/O time and space for these problems, and investigating techniques that are parameterized by the total space or the desired worst-case query cost.

1.1. INDEXING WORKLOADS AND INDEXING SCHEMES. Database access methods must be evaluated in the context of a particular *workload*. A workload consists of an *instance* of a database (a finite subset of some domain), together with a set of queries (a given set of subsets of the instance). In one-dimensional indices such as B-trees, for example, the instance is some totally ordered finite set, and the most common queries considered are *range queries*, that is to say, intervals of this order. Other common workloads include multi-dimensional point sets with range queries, sets of intervals with stabbing queries, and powersets with intersection or inclusion queries.

¹A notable exception to this rule is the inverted index technique widely used for text retrieval (see, e.g., Witten et al. [1999]), in which each document identifier is replicated in the index about as many times as terms in the document. This replication means that online, concurrent updates to text indexes are not widely supported in practice in text retrieval systems.

In what we term indexability theory, the workload plays a role similar to the role a partially recursive language plays in complexity or decidability theory: it is the unit whose complexity must be characterized.² For each workload, we have a space of possible indexing schemes; the analog of algorithms that partially decide the language. Such an indexing scheme is a collection of *B*-subsets of the instance, which we call blocks. The *block size B* is assumed fixed and very large (usually in the hundreds). The union of the blocks exhausts the instance. Each query is answered by retrieving a set of blocks, whose union is a superset of the query.

Our approach suppresses important aspects of indexing, such as the algorithms for determining the partition of the instance into blocks (possibly with repetitions), as well as the algorithms for determining, given a query, the blocks in the indexing scheme that cover it (e.g., the cost of traversing a tree to its leaf level). Furthermore, we also ignore the storage and retrieval costs necessary to support such algorithms, for example, auxiliary information such as "directories" or "internal nodes." These omissions are justified in three ways: first, we are mostly interested in lower bounds, and therefore we are free to disregard aspects of the complexity of the problem. Second, in practice, these aspects do not appear to be the source of design difficulties or of complexity—it appears that good assignment of data items to blocks tends to suggest efficient traversal algorithms, and to have low storage overhead. Third, secondary storage techniques such as buffer management mask and absorb many of these auxiliary cost components. However, our model also ignores the dynamic aspect of the problem, that is, the cost of insertion and deletion. Its consideration *could* be a source of added complexity, and in a more general model the source of more powerful lower bounds.

In this article, we propose a model for the indexing of a data set with respect to a given workload, and explore in its light the fundamental properties and tradeoffs of indexing, with an emphasis on lower bounds. In particular, we introduce a lower-bound theorem that is applicable to arbitrary workloads—although it is not guaranteed to always yield tight bounds. We also analyze within this model a number of interesting families of workloads, including multidimensional point sets with range queries, and powersets with subsumption queries. Besides revealing some interesting laws, our results indicate positive prospects for the use of limited redundancy. For example, for two-dimensional range queries, even a small amount of redundancy can significantly decrease the worst-case query cost.

2. Related Work

This work was initially motivated by the work of Hellerstein, Naughton and Pfeffer on the Generalized Search Tree (or GiST) [Hellerstein et al. 1995]. The GiST is an extensible template indexing structure, organized as a balanced search tree. In their discussion of indexing issues, the authors stated the need for a "theory of indexability," a formal framework that would "describe whether or not trying to index a given data set is practical for a given set of queries."

The research into external data structures has largely been experimental. Theoretical work on the B-tree and its variants, as well as on external hashing, concentrated

 $^{^{2}}$ More accurately, the analog of a language is a *family of workloads*, one for each cardinality of the instance. Such growing families of workloads allow us to focus on asymptotic analysis and ignore additive constants.

mainly on probabilistic analysis of performance, under various distributions of the indexed data. For these problems, the worst-case asymptotic performance has been known for a long time.

Previous work on index data structures concentrated on the study of specialized problems. In the area of multidimensional indexing, data structures are often classified into two categories: those that partition the data set, such as R-trees and their variants, and those that partition the search space, such as quad-trees and their variants. In both categories, most of the proposed algorithms are based on heuristics, and all have relatively bad worst-case asymptotic performance. It is not clear whether this classification has any definitive bearing on performance, and no clear winner has emerged among the many proposals, even for well-understood families of workloads. A comprehensive exposition of the relevant work can be found in the survey of spatial access methods of Gaede and Günther [1998] and the survey of temporal access methods of Salzberg and Tsotras [1999].

This situation has been changing in the past few years, mostly due to the work of Kanellakis, Vitter, and their collaborators. In Kanellakis et al. [1993], it was shown that multidimensional range search generalizes indexing problems in new database paradigms such as constraint databases and class hierarchies. In subsequent publications [Ramaswamy and Subramanian 1994; Ramaswamy and Kanellakis 1995; Subramanian and Ramaswamy 1995; Vengroff and Vitter 1996], asymptotically efficient dynamic algorithms are presented for two-sided and threesided range queries, and for interval stabbing queries. An optimal solution to the interval management problem has recently been found by Arge and Vitter [1996]. Most of this work involves upper bounds and is therefore mainly concerned with the analysis of the searching aspect of the problem. There are two exceptions: First, in Kanellakis et al. [1996], there is an argument (proof of Lemma 2.7) that anticipates our Theorem 4.1, namely, that the access overhead must be \sqrt{B} in the special case in which the blocks are restricted to be rectangular. Second, in the last section of Subramanian and Ramaswamy [1995], there is an interesting lower bound, where it is shown (by extending a result by Chazelle [1990a] to the case of block accesses) that storage redundancy $\Omega(\log n / \log \log_R n)$ is necessary if *additive* (as opposed to our multiplicative) access overhead is to remain polynomial in $\log_{B} n$. Also related are the results in Nodine et al. [1993], who use cost metrics similar to ours, to characterize the locality in external graph searching.

The question of lower bounds in multidimensional searching has been addressed in Mehlhorn [1984], without, however, our emphasis on block accesses. Similar work is presented in Smid and Overmars [1990], where lower bounds are derived in a model involving binary trees with certain further restrictions; the block size is considered in that paper as a function of n, the number of points. Finally, in the database literature there has been extensive analysis (worst case, expected case, or empirical/experimental) of many access methods for multidimensional searching (see, e.g., Pagel et al. [1993], Faloutsos and Kamel [1994], and Belussi and Faloutsos [1995]). More recently, the ideas presented here have been used in a more rigorous framework for empirically analyzing and tuning indexing performance [Shah et al. 1999].

The concept of a space/time tradeoff in main-memory range searching has been studied thoroughly [Fredman 1980, 1981; Yao 1982; Vaidya 1989; Chazelle 1990a, 1990b, 1995]. All these works consider variants of either the RAM machine, or the

pointer machine. These memory models are fundamentally different from blockstructured secondary memory.

The *cell probe model*, originally introduced by Yao [1981], is a general framework for dealing with data structure problems, especially valuable for proving lower bounds, and space-time trade-offs in particular. Let f be any mapping from query $Q \in \{0, 1\}^q$ and dataset $d \in \{0, 1\}^n$ to the answer f(Q, d) of query Q over d. The cell probe model assumes the existence of a memory of s cells, each cell of b bits. Let t be the maximum, over all Q and d, of the least number of cells that must be accessed in order to compute f(Q, d). We are interested in trade-offs between s and t, with b a parameter of the model. Miltersen et al. [1995] proposed some general lower bounds techniques, employing asymmetric communication complexity, and applied them to certain data structure problems related to set membership.

The cell probe model is more general than the indexability model in this paper, because in it memory can be organized in an arbitrary way, whereas in indexability we assume that the memory contains explicit representations of the records. The cell probe model has been used in the past to derive lower bounds in geometric problems; for example, Chakrabarti et al. [1999] and Barkol and Rabani [2000] applied this model to the *nearest neighbor* problem, a pure search problem for which indexability yields trivial results. However, to date the cell probe model has not been applied to range *reporting* problems, which is the class of problems with which indexability is concerned. By "reporting problems" we mean, informally, data structure problems in which the output of the algorithm must be a set of records (think of them as strings or pointers), and the algorithm is not allowed to look inside these records. Reporting problems are an appropriate framework for database storage problems, as they reflect the data independence present in databases. For reporting problems, it makes sense to restrict the data structure solutions so that each memory location holds a record, as we do in the indexability framework, forfeiting the generality of the cell probe model. As one of the referees pointed out to us, by restricting the records stored to be single bits, the cell probe model can be adapted to prove certain lower bounds similar to our bounds for set workload reporting problems (Section 7), starting from the communication complexity results of Miltersen et al. [1995]. These bounds are quantitatively weaker than ours, but hold in a more general model (albeit a model the generality of which is inaccessible to the database problems of interest here). It would be interesting to find such cell probe lower bounds for range search workloads (our main concern in this article).

3. Definitions

In this section, we set out a simple framework for defining an indexing problem, and for measuring the efficiency of a particular indexing scheme for the problem.

3.1. INDEXING WORKLOADS. Indexing schemes must be evaluated in the context of a particular *workload*, consisting of a finite subset of some domain together with a set of queries. More formally, we have the following definition:

Definition 3.1. A workload W is a tuple W = (D, I, Q), where D is a nonempty set (the *domain*), $I \subseteq D$ is a nonempty finite set (the *instance*), and Q is a set of subsets of I (the *query set*). A workload we consider extensively is the set of two-dimensional range queries. This workload consists of the domain \mathbb{R}^2 , the instance $I = \{(i, j) : 1 \le i, j, \le n\}$, and the family of "range queries" $Q[a, b, c, d] = \{(i, j) : a \le i \le b, c \le j \le d\}$, one for each quadruple (a, b, c, d) with $1 \le a \le b \le n, 1 \le c \le d \le n$. Notice that this is a *family of workloads*, with instances of increasing cardinality, one for each $n \ge 0$. Another family of workloads (the set inclusion queries) has as its domain, for each n, all subsets of $\{1, 2, ..., n\}$, and for each subset I of the domain, the set of queries $Q = \{Q_S : S \subseteq \{1, 2, ..., n\}\}$, where $Q_S = \{T \in I : T \subseteq S\}$.

In the terminology of combinatorics, W is a simple hypergraph, where I is the vertex set, and Q is the edge set. The hypergraph abstraction has been used in related work to measure the quality of existing indexing schemes on particular workloads [Shah et al. 1999]. We do not use this terminology here, choosing instead to define terms more natural for databases. There is no analog of the domain Din hypergraphs. We could have dropped it from our definition, but it is suggestive of a parameterization of workloads. For example, all two-dimensional range-query workloads have the same domain.

3.2. INDEXING SCHEMES

Definition 3.2. An indexing scheme S = (W, B) consists of a workload W = (D, I, Q), and for some positive integer B a set B of B-subsets of I, such that B covers I.

We refer to the elements of \mathcal{B} as blocks, and to \mathcal{B} as the set of blocks. We refer to B as the block size, and K stands for the total number of blocks $|\mathcal{B}|$. Notice that an indexing scheme is a simple, B-regular hypergraph with vertex set I.

As a convention in this article, we use lower-case letters from the end of the alphabet, x, y, z to represent elements of I, letter Q, possibly with subscripts, to denote queries, and letter b, possibly with subscripts, to denote blocks. Also, we typically use U to represent sets of blocks.

3.3. PERFORMANCE MEASURES. We now define two performance measures on indexing schemes, *redundancy* and *access overhead*, evaluating the performance of the scheme in terms of space and I/O, respectively. In particular, redundancy measures the amount of space needed by the indexing scheme, while access overhead measures the amount of I/O required by queries. In both cases, the measures are normalized by the ideal performance (linear space and size of the query, respectively). In the following definitions, let S = (W, B) be an indexing scheme of block size *B* on workload W = (D, I, Q), and let N = |I|.

3.3.1. Storage Redundancy

Definition 3.3. The redundancy r(x) of $x \in I$ is the number of blocks that contain x:

$$r(x) = |\{b \in \mathcal{B} : x \in b\}|.$$

The redundancy r of S is then defined as the average of r(x) over all objects:

$$r = \frac{1}{N} \sum_{x \in I} r(x).$$

It is easy to see that the number of blocks is K = rN/B.

We also define the maximum redundancy \hat{r} in S, as $\hat{r} = \max_{x \in I} r(x)$.

3.3.2. Access Overhead

Definition 3.4. A set of blocks, $U \subseteq \mathcal{B}$, covers a query $Q \in \mathcal{Q}$, iff $Q \subseteq \bigcup_{b \in U} b$.

Definition 3.5. A cover set, $C_Q \subseteq \mathcal{B}$, for query $Q \in \mathcal{Q}$ is a minimum-size set of blocks that covers Q.

Notice that a query may have multiple cover sets.

Definition 3.6. The access overhead A(Q) of query $Q \in Q$ is defined as

$$A(Q) = \frac{|C_Q|}{\lceil |Q|/B \rceil}$$

where $C_0 \subseteq \mathcal{B}$ is a cover set for Q.

It is easy to see that $1 \le A(Q) \le B$, since any query Q will be covered by at most |Q| blocks.

Informally, A(Q) models the observed cost of query Q normalized by its ideal cost, in terms of block accesses. For a given query Q, $\lceil |Q|/B \rceil$ is the minimum number of blocks required. A(Q) is the multiplicative overhead associated with Q for a particular indexing scheme.

We now define the access overhead A of indexing scheme S, to be the maximum of A(Q) over all queries.

Definition 3.7. The access overhead A for indexing scheme S is

$$A = \max_{Q \in \mathcal{Q}} A(Q)$$

Notice that, although the redundancy is defined as an average (over all data items), the access overhead is a maximum (over all queries). This is less arbitrary than it may seem at first. By averaging over all data items we capture the true (worst-case) space performance of the indexing scheme, while averaging I/O performance over all queries would be much less defensible since queries are generally not equiprobable, and guarantees, and thus worst-case analysis, are desirable in the context of query response time.

3.4. SOME TRIVIAL BOUNDS AND TRADE-OFFS. Based on standard properties of databases and disks, we assume that the number of objects N is always much greater than the block size B, although B is not limited in any concrete way.

For some indexing scheme S, the minimum possible redundancy is 1, when \mathcal{B} is a partition of I, and the maximum redundancy is $\binom{N-1}{B-1}$, when $\mathcal{B} = \binom{I}{B}$.³ For S having maximum redundancy, A is exactly 1, which is minimum; in that case, every query Q can be covered by a set of disjoint blocks whose union contains Q.

42

³For set *S* and $n \ge 0$, $\binom{S}{n}$ denotes the set of all *n*-subsets of *S*.

Also, for r = 1 it is easy to devise a problem where A = B, which is maximum (e.g., $Q = {I \choose B}$).

4. Trade-Offs for a Two-Dimensional Workload

Given this framework for indexability, we proceed to examine some families of workloads that have received significant attention in the indexing literature. Our main goal is to expose the trade-offs in lower bounds of r and A for these workloads, delimiting the potential efficiency of indexes for these workloads. We start with some simple positive results (upper bounds) that are useful in two ways. First, they illustrate the framework of indexing schemes. And second, they allow us to conclude later that the lower bounds of this article are tight.

Our main lower bound results are driven by the *Redundancy Theorem* that we develop in Section 5. However, we do not need the Redundancy Theorem to obtain our first interesting lower bound, which is presented in the second part of section.

4.1. TWO-DIMENSIONAL QUERIES. We shall consider here workloads over the two-dimensional domain \mathbb{R}^2 , with $I = \{(i, j) : 1 \le i, j, \le n\}$, and 2-d range queries over this instance. We are interested in determining the minimum possible access overhead when the redundancy r is fixed.

PROPOSITION 1. For each integer r, there is an indexing scheme S_r with redundancy r and access overhead $2B^{1/2r} + 2$.

PROOF. The main idea for the indexing scheme S_r is that each query Q of $x \times y$ points will be covered by disjoint blocks of S_r that have "almost" the same aspect ratio y/x with Q. The ideal situation is to have blocks with aspect ratio y/x, so that the query Q is tiled nicely by these blocks; compare this with the worst case when the query Q is "long and narrow" and it is covered by "short and wide" blocks. Because of the restriction on the redundancy r of the indexing scheme S_r , it is not possible to have blocks for each aspect ratio. However, we can choose blocks so that any aspect ratio can be approximated.

More precisely, for each i = 1, 2, ..., r, our indexing scheme S_r contains all $B^{(2i-1)/2r} \times B^{(2r-2i+1)/2r}$ blocks that partition *I*. The aspect ratios $B^{(r-2i+1)/r}$, for i = 1, 2, ..., r, of these blocks are evenly distributed. It is immediate that S_r has redundancy *r* (maximum as well as average). To see that S_r has access overhead at most $2B^{1/2r} + 2$, consider the set of $B^{j/r} \times B^{(r-j)/r}$ queries, j = 0, 1, ..., r. Clearly, the best coverage of such a query is by blocks that have almost the same aspect ratio, that is, blocks of size $B^{(2j-1)/2r} \times B^{(2r-2j+1)/2r}$ or blocks of size $B^{(2j+1)/2r} \times B^{(2r-2j-1)/2r}$. In both cases, when the query is "aligned" with the blocks, it requires $B^{1/2r}$ blocks (either one row of $B^{1/2r}$ blocks or one column of $B^{1/2r}$ blocks). For non-aligned queries the ratio can be as high as $2B^{1/2r} + 2$; to see this, consider the case where an aligned query is satisfied by a row of $B^{1/2r}$ blocks. If we shift this query out of horizontal and vertical alignment, we need two rows of blocks instead of one, and at one of the ends we need an additional column of two blocks as well. On the other hand, it is not difficult to see that these are the worst queries for this indexing scheme. \Box

If the access ratio is A, the above scheme has both average and maximum redundancy $r = \Omega(\log B/\log A)$. We show that this is the best possible relation between

r and A. Indeed, in the remainder of this section, we prove that this is the case when the maximum redundancy is one. We defer the study of the general case after we introduce our lower-bound theorem.

4.2. A LOWER BOUND FOR REDUNDANCY r = 1. We show that up to a constant factor the above indexing scheme is optimal when r = 1. In Kanellakis et al. [1993], the result below was shown for the special case when the blocks are restricted to be rectangular.

THEOREM 4.1. Any indexing scheme for 2-dimensional range queries with redundancy r = 1 has access overhead at least $A = B^{1/2}$. For the d-dimensional case, the lower bound is $A = B^{1-(1/d)}$.

PROOF. We consider first the 2-dimensional case, the general case being a straightforward generalization. For simplicity, we assume that n is a multiple of B.

For the lower bound, we consider only queries of size $1 \times B$ and $B \times 1$. The queries of size $1 \times B$ partition the instance and so do the queries of size $B \times 1$. The total number of queries is $2n^2/B$.

Now fix a block $b \in \mathcal{B}$ that intersects x_1 horizontal lines and x_2 vertical lines (by a "line" we mean a set of data points of the form $\{(1, j), (2, j), \ldots, (n, j)\}$ or $\{(i, 1), (i, 2), \ldots, (i, n)\}$). Since every block has *B* points, we must have $x_1x_2 \ge B$; hence $x_1 + x_2$ is at least $2B^{1/2}$. Therefore, every block intersects at least $2B^{1/2}$ of the above queries. Taking into account that the number of queries is twice the number of blocks, we can conclude that, on the average, every query of the above collection is intersected by $B^{1/2}$ blocks at least. (To see this in detail, consider the number of pairs of intersecting blocks and queries; it is no less than $2B^{1/2}$ times the total number of blocks, which is $2B^{1/2}n^2/B$; since there are $2n^2/B$ queries in total in the collection, the average number of intersecting blocks per query is $B^{1/2}$.) When the redundancy is r = 1, all these blocks are needed to cover the query.

Notice that we showed not only that there exists a query with access overhead $B^{1/2}$, but that this is the expected access overhead for a random query from the above set.

The generalization to the *d*-dimensional case is straightforward (for example we now have $x_1 + \cdots + x_d \ge dB^{1/d}$ which gives access overhead at least $B^{1-(1/d)}$).

5. The Redundancy Theorem

We now turn our attention to a workload-independent analysis of the indexability model that culminates with the Redundancy Theorem.

We first state and prove a set-theoretic result that is of central importance to our work. Note that this theorem is not specific to indexing schemes; it arises in extremal set theory. The reader is warned that the notation does not correspond to indexing schemes.

THEOREM 5.1. Let S_1, S_2, \ldots, S_a $(a \ge 1)$ be nonempty finite sets, $S = S_1 \cup S_2 \cup \cdots \cup S_a$ be their union, and $L \le |S|$ be a positive integer. Let k denote the maximum integer such that there exist k pair-wise disjoint sets P_1, P_2, \ldots, P_k , so that for all $i, 1 \le i \le k$,

- (1) $|P_i| = L$, and
- (2) $P_i \subseteq S_j$ for some $j, 1 \leq j \leq a$.

or k = 0 if no such sets exist. Then,

$$k \ge \frac{|S|}{L} - a. \tag{1}$$

PROOF. Let P_1, \ldots, P_k be sets that satisfy the properties of the theorem and let P be their union, $P = P_1 \cup P_2 \cup \cdots \cup P_k$. The maximality of k guarantees that P contains all but at most L elements from every $S_i, i = 1, \ldots, a$. That is, $|S_i \setminus P| < L$ (otherwise, we can add any subset of L elements of $S_i \setminus P$ to the collection of P_i 's). We can now estimate $|S \setminus P| = |(\bigcup S_i) \setminus P| = |\bigcup (S_i \setminus P)| \le \sum_{i=1}^{a} |S_i \setminus P| < aL$. Since every P_j has cardinality L we conclude that kL = |P| > |S| - aL which implies the desired k > |S|/L - a. \Box

To apply the above theorem to the domain of indexing schemes, we define a convenient concept, *flakes*, to capture the overlap of queries and blocks.

Definition 5.2. Let S = (W, B) be an indexing scheme on workload W = (D, I, Q). A flake is any set of objects $F \subseteq I$ such that for some query Q and some block $b, F \subseteq Q \cap b$.

Note that a flake is a subset (potentially a proper subset) of the intersection of a block and a query. The flexibility to deal with proper subsets will allow us to consider flakes of a fixed size, allowing us to apply certain combinatorial results below.

We now have the following lemma on flakes:

LEMMA 5.1 (FLAKING LEMMA). Let S be an indexing scheme, A be its access overhead, and ϑ be a real number in the interval [2, B/A] such that $B/\vartheta A$ is an integer. Then, any query Q with $|Q| \ge B/2$ will contain at least $(\vartheta - 2)A|Q|/B$ pair-wise disjoint flakes of size $B/\vartheta A$.

PROOF. The parameter ϑ exists only to guarantee that $B/\vartheta A$ is integer.

Choose a cover set for Q, say $C_Q = \{b_1, \dots, b_a\}$, of size a. Let S_1, \dots, S_a be flakes defined by $S_i = Q \cap b_i$ for $1 \le i \le a$. We have

$$a = A(Q) \left\lceil \frac{|Q|}{B} \right\rceil \le A \left\lceil \frac{|Q|}{B} \right\rceil \le 2A \frac{|Q|}{B}$$

(because $|Q| \ge B/2$). We apply Theorem 5.1 on S_i for $L = B/\vartheta A$, and conclude that the number k of flakes of size $B/\vartheta A$ is at least

$$k \geq \frac{|Q|}{B/\vartheta A} - a$$

$$\geq \vartheta A \frac{|Q|}{B} - 2A \frac{|Q|}{B}$$

$$= (\vartheta - 2)A \frac{|Q|}{B}.$$

We proceed to prove a second technical tool from extremal set theory. In coding theory, under a slightly different statement, this result is known as Johnson's bound [Johnson 1962]. Again, the notation does not correspond to indexing schemes.

THEOREM 5.3 (JOHNSON'S BOUND). Let S be a finite set, and S_1, S_2, \ldots, S_k be subsets of S, each of size at least $\alpha |S|$, such that the intersection of any two of them is of size at most $\beta |S|$. If $\beta < \alpha^2/(2 - \alpha)$ the number of subsets k is at most α/β .

PROOF. Since $S_1, S_2, \ldots, S_t, t \le k$, are subsets of S, their union $S_1 \cup S_2 \cup \cdots \cup S_t$ is also a subset of S and therefore

$$\left|\bigcup_{j=1}^{t} S_{j}\right| \leq |S|.$$

It follows that

$$\sum_{j=1}^{t} |S_j| - \sum_{j=1}^{t} \sum_{l=j+1}^{t} |S_j \cap S_l| \le |S|.$$

By the assumptions about the sizes of the subsets and their pairwise intersection, the last inequality implies that

$$t\alpha|S| - \binom{t}{2}\beta|S| \le |S|.$$

Therefore, every $t \le k$ must satisfy the inequality $\alpha t - \beta {t \choose 2} - 1 \le 0$. It immediately follows that if a positive integer t does not satisfy this inequality, then the number k of subsets must be less than t. So, in order to upper bound the number k of subsets, we need to guarantee that the above inequality is not satisfied by at least one positive integer. This can be easily done if we require that the two roots of the polynomial $\alpha t - \beta {t \choose 2} - 1$ differ by more than 1. Since the roots of the polynomial are

$$\frac{\alpha+\beta/2\pm\sqrt{(\alpha+\beta/2)^2-2\beta}}{\beta},$$

it is easy to verify that they differ by more than 1 when $\beta < \alpha^2/(2-\alpha)$.

But then, the number of subsets is at most equal to the minimum root of the above polynomial. Thus

$$k \leq \frac{\alpha + \beta/2 - \sqrt{(\alpha + \beta/2)^2 - 2\beta}}{\beta}.$$

This last inequality implies that $k \leq \alpha/\beta$. \Box

Note that the hypotheses of the above lemma cannot be improved by a factor of more than 2, because when $\beta \ge \alpha^2$, the number of possible subsets is unbounded, that is, it is an increasing function of |S|.

We are now ready to state and prove our main result.

THEOREM 5.4. Let S be an indexing scheme, and let Q_1, Q_2, \ldots, Q_M be queries, such that for every $i, 1 \le i \le M$:

- (1) $|Q_i| \ge B/2$, and
- (2) $|Q_i \cap Q_j| \leq B/2(\vartheta A)^2$ for all $j \neq i, 1 \leq j \leq M$.

Then, the redundancy is bounded by

$$r \ge \frac{\vartheta - 2}{2\vartheta} \frac{1}{N} \sum_{i=1}^{M} |Q_i|,$$

where ϑ is any real number in the interval [2, B/A] such that $B/\vartheta A$ is integer.

PROOF. We prove the lower bound in two steps. First, we compute the *minimum* number of flakes contained in queries Q_1, Q_2, \ldots, Q_M . Let this number be f_1 . Then we will compute the maximum number of flakes contained in each block. Let this number be f_2 . Clearly, there will be at least f_1/f_2 blocks in \mathcal{B} .

Step 1. Consider any query Q_i . By the flaking lemma, this query contains at least $(\vartheta - 2)A|Q_i|/B$ disjoint flakes of size $B/\vartheta A$. Let F be such a flake. F cannot be contained in some other query Q_j , $j \neq i$, because if it were, then it would be a subset of Q_j as well as of Q_i , and thus $|Q_i \cap Q_j| \ge B/\vartheta A > B/2(\vartheta A)^2$. We conclude that

$$f_1 = \sum_{i=1}^{M} (\vartheta - 2) A \frac{|Q_i|}{B} = (\vartheta - 2) A \sum_{i=1}^{M} \frac{|Q_i|}{B}.$$

Step 2. Consider any block *b*, and let F_1, F_2, \ldots, F_k be the flakes contained in this block. Since all these flakes are subsets of *b*, we upper bound the number of flakes *k*, using Johnson's bound. Each flake F_i is of size $B/\vartheta A$. Also, for two distinct flakes F_i and $F_j, i \neq j, |F_i \cap F_j| \leq B/2(\vartheta A)^2$, by the following argument: If the flakes are contained in the same query, then they are disjoint. If the flakes are contained in different queries, then their intersection is bounded by the intersection of these queries. Thus, Johnson's bound is applicable with $\alpha = 1/\vartheta A$, and $\beta = 1/2(\vartheta A)^2$. It can easily be checked that $\beta < \alpha^2/(2 - \alpha)$. Thus, we conclude that

$$f_2 = \frac{\alpha}{\beta} = 2\vartheta A.$$

Substituting, we get

$$\frac{f_1}{f_2} = \frac{\vartheta - 2}{2\vartheta} \sum_{i=1}^M \frac{|Q_i|}{B}$$

The proof is complete, by the inequality $K = rN/B \ge f_1/f_2$ which simplifies to

$$r \ge \frac{\vartheta - 2}{2\vartheta} \frac{1}{N} \sum_{i=1}^{M} |Q_i|.$$

Notice that the theorem is useful only for access overhead $A = O(\sqrt{B})$: either all queries are disjoint (implying access overhead 1), or for nondisjoint queries Q_i and Q_j , the second premise in the statement of the theorem implies that $B/2(\vartheta A)^2 \ge 1$.

We now simplify the theorem by removing the parameter ϑ .

THEOREM 5.5 (REDUNDANCY THEOREM). Let S be an indexing scheme with access overhead $A \leq \sqrt{B}/4$, and let Q_1, Q_2, \ldots, Q_M be queries, such that for

every $i, 1 \leq i \leq M$:

(1) $|Q_i| \ge B/2$, and (2) $|Q_i \cap Q_j| \le B/16A^2$ for all $j \ne i, 1 \le j \le M$.

Then, the redundancy is bounded by

$$r \ge \frac{1}{12N} \sum_{i=1}^{M} |Q_i|.$$

PROOF. Let $\vartheta_1 = 12/5$ and $\vartheta_2 = 2\sqrt{2}$. We first show that there exists $\vartheta \in [\vartheta_1, \vartheta_2]$ such that $B/\vartheta A$ is integer. This follows from

$$\frac{B}{\vartheta_1 A} - \frac{B}{\vartheta_2 A} = \left(\frac{1}{\vartheta_1} - \frac{1}{\vartheta_2}\right) \frac{B}{A} \ge \left(\frac{1}{\vartheta_1} - \frac{1}{\vartheta_2}\right) \frac{B}{A^2} \ge \left(\frac{1}{\vartheta_1} - \frac{1}{\vartheta_2}\right) 16 > 1.$$

Using such a ϑ in Theorem 5.4, the second premise becomes

$$|Q_i \cap Q_j| \le \frac{B}{16A^2} = \frac{B}{2(\vartheta_2 A)^2} \le \frac{B}{2(\vartheta A)^2}$$

and the factor $(\vartheta - 2)/2\vartheta$ of the conclusion becomes

$$\frac{\vartheta - 2}{2\vartheta} \ge \frac{\vartheta_1 - 2}{2\vartheta_1} = \frac{1}{12}.$$

Observe that given any set of queries $\mathcal{M} = \{Q_1, \dots, Q_M\}$, we can construct blocks for each query independently, for a total of

$$T_{\mathcal{M}} = \sum_{i=1}^{M} \left\lceil \frac{|Q_i|}{B} \right\rceil$$

blocks, achieving a perfect access overhead of one, with redundancy $r = T_{\mathcal{M}} \frac{B}{N} \ge \sum_{i=1}^{M} |Q_i|/N$. The Redundancy Theorem states that when the queries intersect pairwise in at most $B/16A^2$ elements for some A, increasing the access overhead to A does not yield an improvement in space by more than a constant factor of $T_{\mathcal{M}}$.

6. Lower Bounds for Multidimensional Range Queries

We now apply the Redundancy Theorem to d-dimensional range queries. First, we examine the case for 2-dimensional range queries, and then we generalize to d dimensions.

For any $d \ge 1$, we define the *d*-dimensional range query workload, \mathcal{R}_n^d , whose domain is \mathbb{R}^d , with instance $I = [1:n]^d$ and query set

$$\mathcal{Q} = \{[a_1:b_1] \times \cdots \times [a_d:b_d] \mid 1 \le a_i \le b_i \le n\}.$$

For this workload, $N = n^d$.

6.1. 2-D RANGE QUERIES. In order to apply the Redundancy Theorem, we must identify queries Q_1, Q_2, \ldots, Q_M , each of size at least B/2, and with pairwise intersections at most $B/16A^2$. We consider only queries of size $c^j \times B/c^j$, for $j = 0, 1, \ldots, \log_c B$. For each aspect ratio, we partition the $n \times n$ grid, obtaining a

48



FIG. 1. Two rectangles of sizes $c^j \times B/c^j$ and $c^{j'} \times B/c^{j'}$, j < j', intersecting in at most $B/c^{j'-j}$ points.

total of $M = n^2/B(1 + \log_c B)$ queries of size B each. Before we apply the theorem, we compute the parameter c.

Let j and j' be integers $0 \le j < j' \le \log_c B$, and Q_j and $Q_{j'}$ be queries of dimensions $c^j \times B/c^j$ and $c^{j'} \times B/c^{j'}$, respectively. Figure 1 depicts the setup. It is easy to see that for any j and j', $|Q_j \cap Q_{j'}| \le B/c^{j'-j} \le B/c$. Thus, we take $c = 16A^2$.

We are now ready to apply the Redundancy Theorem. From the theorem,

$$r \geq \frac{1}{12} \frac{MB}{n^2}$$

= $\frac{1}{12} \frac{1}{n^2} \left(B \frac{n^2}{B} (1 + \log_c B) \right)$
= $\frac{1}{12} (1 + \log_c B)$
 $\geq \frac{1}{12} \log_c B$
= $\frac{1}{12} \frac{\log B}{\log(16A^2)}$

and thus we have

$$r = \Omega\left(\frac{\log B}{\log A}\right).$$

6.2. *d*-DIMENSIONAL QUERIES. We can generalize the above technique to *d*-dimensional queries. We consider queries of size *B*, with dimensions $c^{j_1} \times c^{j_2} \times \cdots \times c^{j_d}$, for all nonnegative integer j_1, j_2, \ldots, j_d , such that $\sum_{k=1}^d j_k = \log_c B$. For each sequence j_1, j_2, \ldots, j_d , we partition the *d*-dimensional cube into n^d/B (hyper)rectangles, of dimensions $c^{j_1} \times c^{j_2} \times \cdots \times c^{j_d}$.

In order to select the appropriate value for c, we consider the size of pairwise intersections of rectangles with different dimensions. It is easy to see that $c = 16A^2$ is applicable in this case also, guaranteeing that the intersection of any two rectangles will have size at most $B/16A^2$.

We also use the well-known fact that the number of distinct sequences of d nonnegative integers, whose sum is n, is given by

$$\binom{n+d-1}{d-1}$$

(cf. Bose-Einstein distribution).

Thus, the total number of queries (each of size B) will be

$$M = \frac{n^d}{B} \left(\frac{\log_c B + d - 1}{d - 1} \right) = \frac{n^d}{B} \left(\frac{\frac{\log B}{\log(16A^2)} + d - 1}{d - 1} \right)$$

and for the redundancy we have

$$r \ge \frac{1}{12} \left(\frac{\log B}{\log(16A^2)} + d - 1 \\ d - 1 \right).$$

For d a constant, the above quantity is a polynomial of degree d - 1. Thus, we have shown the following theorem:

THEOREM 6.1. For workload \mathcal{R}_n^d , the storage redundancy is bound by

$$r = \left(\frac{\Omega(\frac{\log B}{\log A}) + d - 1}{d - 1} \right) = \Omega\left(\left(\frac{\log B}{\log A} \right)^{d - 1} \right).$$

6.3. FIBONACCI WORKLOAD. So far, our trade-offs have depended only on the block size B, but not on the size of the instance. Unfortunately, this is not always the case. In this section, we study a family of workloads for two-dimensional range queries that exhibits much worse performance.

Using our framework of indexing schemes, it was shown in Koutsoupias and Taylor [1998] that there exist simple 2-dimensional workloads with trade-offs that depend on the instance size. In particular, they studied range queries of the Fibonacci lattice (to be defined shortly) and showed that any indexing scheme with redundancy less than $\Theta(\log n)$ has the worst possible overhead A = B. The bound $\Theta(\log n)$ is tight up to a constant factor. They later extended the results to random sets of points and higher dimensions Koutsoupias and Taylor [1999].

Here we illustrate the power of the Redundancy Theorem by extending the results for the Fibonacci lattice when the access overhead is small, $A = O(\sqrt{B})$. Furthermore, we give the precise trade-off between redundancy and access overhead.

We now define the Fibonacci lattice, which is the regular lattice rotated appropriately. Let $n = f_k$ be the *k*th Fibonacci number. The Fibonacci lattice F_n is the

set of points defined by:

$$F_n = \{(i, if_{k-1} \mod n) : i = 0, 1, \dots, n-1\}$$
 for $n = f_k$

The Fibonacci workload over domain \mathbb{R}^2 is defined by taking the Fibonacci lattice as the instance *I*, and all rectangular queries as \mathcal{Q} .

We only need the following property of the Fibonacci lattice, from Fiat and Shamir [1989]:

PROPOSITION 2. For the Fibonacci lattice F_n of n points, and for $t \ge 0$, any rectangle with area $t \cdot n$ contains between $\lfloor t/c_1 \rfloor$ and $\lceil t/c_2 \rceil$ points, where $c_1 \approx 1.9$ and $c_2 \approx 0.45$.

Now we apply the Redundancy Theorem to the Fibonacci workload. We have to define an appropriate set of queries Q_1, \ldots, Q_M , each of cardinality at least B/2.

We consider rectangles of area $a = c_1 Bn/2$. By Proposition 2, each such rectangle will contain at least B/2 points. Let c be a parameter to be specified later. We consider rectangles of dimensions $c^i \times a/c^i$, for appropriate values of i. For each such aspect ratio, we partition the Fibonacci lattice into non-overlapping rectangles, in a tiling fashion. Each of these rectangles will define a query.

Because no rectangle can have a side longer than n, we must constrain i to obey

$$c^i \le n$$
 and $\frac{a}{c^i} \le n$.

From these, we compute that *i* must range between $\log_c (c_1 B/2)$ and $\log_c n$, that is, approximately $\log_c 2n/c_1 B$ aspect ratios. Since, for each *i*, we cover the whole set of points, the Redundancy Theorem gives

$$r \ge \frac{1}{12}\log_c \frac{2n}{c_1B} = \Omega\left(\frac{\log(n/B)}{\log c}\right).$$

Now we specify an appropriate value of parameter c that satisfies the second premise of the Redundancy Theorem—which states that no two queries can intersect by more than $B/16A^2$ points. We observe that rectangles of the same aspect ratio do not intersect, and rectangles of different aspect ratios have intersections of area at most a/c. Again by Proposition 2, it suffices to have

$$\left\lceil \frac{a/c}{c_2 n} \right\rceil \le \frac{B}{16A^2},$$

which is satisfied by

$$c \approx 8 \frac{c_1}{c_2} A^2.$$

Thus, we have the following theorem:

THEOREM 6.2. For the Fibonacci workload, any indexing scheme with the access overhead $A \le \sqrt{B}/4$ must have redundancy

$$r = \Omega\left(\frac{\log(n/B)}{\log A}\right).$$

The Fibonacci lattice is only one of many low-discrepancy [Matousek 1999], planar point sets we could have used. For example, we could have used the point

set used by Chazelle [1990a], in his proof of a lower bound for range search in the pointer machine model. Matousek [1999] discusses the discrepancy properties of the Fibonacci lattice, and many other point sets. However, none of these will improve the trade-off of Theorem 6.2 by more than a small constant factor.

7. Set Workloads

We now turn our attention to the problem of indexing for arbitrary sets. An interesting workload is the λ -set workload $\mathcal{K}_{n,\lambda}$, whose instance is the set $\{1, \ldots, n\}$ and whose query set is the set of all λ -subsets of the instance. We show that these workloads are far worse than 2-dimensional queries.

Our Redundancy Theorem is applicable only when $\lambda > B/2$. In practice, we are also interested in workloads with small values for λ . To analyze these workloads, we prove a corollary of the following famous theorem by Turán [Turán 1941; van Lint and Wilson 1992]:

THEOREM 7.1 (TURÁN'S THEOREM). If a simple graph of n vertices has more than

$$\frac{(p-2)n^2}{2(p-1)} - \frac{r(p-1-r)}{2(p-1)} \quad (r = n \bmod p)$$

edges, then it contains a complete graph of p vertices (a p-clique).

For a given graph, an *independent set* is a subset of its vertices such that there is no edge between any pair of these vertices.

COROLLARY 1. In a simple graph G(V, E), with |V| = n, if

$$|E| \le \frac{n^2 - n(p-1)}{2(p-1)},$$

then G has an independent set of size p.

PROOF. Let $\tilde{G}(V, \tilde{E})$ be the graph with

$$\tilde{E} = \left\{ (v_1, v_2) \in \binom{V}{2} \mid (v_1, v_2) \notin E \right\}.$$

Then,

$$\tilde{E}| = \binom{n}{2} - |E| > \frac{(p-2)n^2}{2(p-1)},$$

and thus by Turán's Theorem \tilde{G} has a *p*-clique. The vertices of the clique form an independent set in *G*. \Box

We now show a lower bound for set workloads.

THEOREM 7.2. For workload $\mathcal{K}_{n,\lambda}(I, \mathcal{Q}), B \geq \lambda$, any indexing scheme with redundancy

$$r < \frac{n - \lambda + 1}{(\lambda - 1)(B - 1)}$$

has the worst possible access overhead $A = \lambda$.

PROOF. Construct a graph G(I, E) where $(x_1, x_2) \in E$ iff there exists a block containing both x_1 and x_2 . This graph will have at most

$$r\frac{n}{B}\binom{B}{2} < \frac{n^2 - n(\lambda - 1)}{2(\lambda - 1)}$$

edges. By Corollary 1, it has an independent set of size λ . This set, taken as a query, will require exactly λ distinct blocks to be covered (by the construction of *G*).

The last theorem states that $\mathcal{K}_{n,\lambda}$ requires space at least *quadratic* in n/B to avoid the worst possible access overhead. We show that within a factor of 2, the bound of the theorem is tight.

THEOREM 7.3. For workload $\mathcal{K}_{n,\lambda}$ and $B \ge \lambda$, there exists an indexing scheme of access overhead $A = \lambda - 1$ and redundancy

$$r = \frac{2n}{(\lambda - 1)B} - 1.$$

PROOF. We arbitrarily partition the instance into $\lambda - 1$ sets of roughly equal size, $S_1, \ldots, S_{\lambda-1}$. For each set S_i , we construct suitable blocks so that for any $x, x' \in S_i$ there is a single block containing both. Then, for every query Q, some elements x_1 and x_2 will belong to the same set S_i , and thus will be covered by a single block, and so $A(Q) \leq \lambda - 1$.

To construct blocks for set S_i , we arbitrarily partition the set S_i into $k = 2n/(\lambda - 1)B$ sets t_j , j = 1, ..., k of size B/2 each. For each pair of these sets, we construct a block containing their union. Thus, for any pair of elements of S_i , there exists a block containing both.

For each of the $\lambda - 1$ sets S_i , we constructed $\binom{k}{2}$ blocks. The total number of blocks constructed thus is

$$(\lambda - 1) \left(\frac{\frac{2n}{(\lambda - 1)B}}{2} \right) = \frac{n}{B} \left(\frac{2n}{(\lambda - 1)B} - 1 \right)$$

which yields the required redundancy. \Box

8. Conclusions

We have presented a new framework for the modeling and study of indexing in external memory. Our cost model is minimalistic, in that it ignores important parameters of external memory and indexing. This is not by accident but rather by design. There exist more precise (and more complex) cost models that are more accurate in predicting space and/or I/O indexing costs, for example, models that include the search aspects of indexing, or models that describe hard disk performance more accurately. In our view, however, a successful model is not one that represents reality faithfully, but rather one that manages to capture the essence of a facet of the real world in a way that allows for deeper study and understanding of this facet.

Having argued in favor of the minimalistic aspects of indexability, we should stress that we expect indexability results to often carry over to more detailed models straight-forwardly, and also to the implementation domain. Recent results by Arge et al. [1999] indicate that it may be possible to employ indexability techniques as subroutines in external data structures, as part of a systematic approach to the "externalization" of main memory data structures. Shah et al. [1999] have developed an index analysis tool called **amdb**; among its features is a test for unit redundancy indexability, which serves as a concrete performance target for index developers.

REFERENCES

- AOKI, P. M. 1998. Generalizing "search" in generalized search trees (extended abstract). In *Proceedings* of the 14th International Conference on Data Engineering (Orlando, Fla., Feb. 23–27). pp. 380–389.
- ARGE, L., SAMOLADAS, V., AND VITTER, J. S. 1999. On two-dimensional indexability and optimal range search indexing. In *Proceedings of the 18th Annual ACM Symposium on the Principle of Database Systems*. ACM, New York, pp. 346–357.
- ARGE, L., AND VITTER, J. S. 1996. Optimal dynamic interval management in external memory. In Proceedings of the 37th Annual Symposium on Foundations of Computer Science (Oct.). IEEE Computer Society Press, Los Alamitos, Calif., pp. 560–569.
- BARKOL, O., AND RABANI, Y. 2000. Tighter bounds for nearest neighbor search and related problems in the cell probe model. In *Proceedings of the 32nd Annual ACM Symposium on the Theory of Computing*. ACM, New York, pp. 388–396.
- BAYER, R., AND MCCREIGHT, E. 1972. Organization and maintenance of large ordered indexes. Acta Inf. 1, 173–189.
- BELUSSI, A., AND FALOUTSOS, C. 1995. Estimating the selectivity of spatial queries using the 'correlation' fractal dimension. In *Proceedings of the 21st International Conference on Very Large Databases*. pp. 299–310.
- CHAKRABARTI, A., CHAZELLE, B., GUM, B., AND LVOV, A. 1999. A lower bound on the complexity of approximate nearest-neighbor searching on the hamming cube. In *Proceedings of the 31st Annual ACM Symposium on the Theory of Computing*. ACM, New York, pp. 305–311.
- CHAZELLE, B. 1990a. Lower bounds for orthogonal range searching. I: The reporting case. J. ACM 37, 2 (Apr.), 200–212.
- CHAZELLE, B. 1990b. Lower bounds for orthogonal range searching. II: The arithmetic model. *J. ACM* 37, 3 (June), 439–463.
- CHAZELLE, B. 1995. Lower bounds for off-line range searching. In *Proceedings of the 27th Annual ACM Symposium on the Theory of Computing*. ACM, New York, pp. 733–740.
- FALOUTSOS, C., AND KAMEL, I. 1994. Beyond uniformity and independence: Analysis of R-trees using the concept of fractal dimension. In *Proceedings of the 13th Annual ACM Symposium on the Principles* of Database Systems. ACM, New York, pp. 4–13.
- FIAT, A., AND SHAMIR, A. 1989. How to find a battleship. Networks 19, 361-371.
- FREDMAN, M. L. 1980. The inherent complexity of dynamic data structures which accomodate range queries. In *Proceedings of the IEEE Symposium on Foundations of Computer Science*. IEEE Computer Society Press, Los Alamitos, Calif., pp. 191–199.
- FREDMAN, M. L. 1981. Lower bounds on the complexity of some optimal data structures. SIAM J. Comput. 10, 1–10.
- GAEDE, V., AND GÜNTHER, O. 1998. Multidimensional access methods. *ACM Comput. Surv. 30*, 2 (June), 170–231.
- GUTTMAN, A. 1984. R-trees: A dynamic index structure for spatial searching. In *Proceedings of the ACM SIGMOD International Conference on the Management of Data*. ACM, New York, pp. 47–57.
- HELLERSTEIN, J. M., NAUGHTON, J. E., AND PFEFFER, A. 1995. Generalized search trees for database systems. In Proceedings of the 21st International Conference on Very Large Databases. pp. 562–573.
- JOHNSON, S. M. 1962. A new upper bound for error-correcting codes. *IEEE Trans. Inf. Theory* 8, 203– 207.
- KANELLAKIS, P. C., RAMASWAMY, S., VENGROFF, D. E., AND VITTER, J. S. 1996. Indexing for data models with constraints and classes. J. Comput. Syst. Sci. 52, 3, 589–612.
- KANELLAKIS, P. C., RAMASWAMY, S., VENGROFF, D. E., AND VITTER, J. S. 1993. Indexing for data models with constraints and classes. In *Proceedings of the 12th Annual ACM Symposium on Principles* of Database Systems. ACM, New York, pp. 233–243.
- KANTH, K. V. R., RAVADA, S., SHARMA, J., AND BANERJEE, J. 1999. Indexing medium-dimensionality data in Oracle. In *Proceedings of the ACM SIGMOD International Conference on Management of Data* (Philadelphia, Pa., June 1–3). ACM, New York, pp. 521–522.

- KORNACKER, M. 1999. High-performance extensible indexing. In Proceedings of 25th International Conference on Very Large Data Bases (Edinburgh, Scotland, UK, Sept. 7–10). ACM, New York, pp. 699–708.
- KORNACKER, M., MOHAN, C., AND HELLERSTEIN, J. M. 1997. Concurrency and recovery in generalized search trees. In *Proceedings of the ACM SIGMOD International Conference on Management of Data* (Tucson, Az., May 13–15). ACM, New York, pp. 62–72.
- KOUTSOUPIAS, E., AND TAYLOR, D. S. 1998. Tight bounds for 2-dimensional indexing schemes. In Proceedings of the 17th Annual ACM Symposium on Principles of Database Systems. ACM, New York, pp. 52–58.
- KOUTSOUPIAS, E., AND TAYLOR, D. S. 1999. Indexing schemes for random points. In *Proceedings of the* 10th Annual ACM-SIAM Symposium on Discrete Algorithms. ACM, New York, pp. 596–602.

MATOUSEK, J. 1999. Geometric Discrepancy. Springer-Verlag, New York.

- MEHLHORN, K. 1984. Data Structures and Algorithms 3: Multi-dimensional Searching and Computational Geometry. EATCS Monographs on Theoretical Computer Science, Springer-Verlag, New York.
- MILTERSEN, P. B., NISAN, N., SAFRA, S., AND WIDGERSON, A. 1995. On data structures and asymmetric communication complexity. In *Proceedings of the 27th Annual ACM Symposium on the Theory of Computing*. ACM, New York, pp. 103–111.
- NIEVERGELT, J., HINTERBERGER, H., AND SEVCIK, K. C. 1984. The grid file: An adaptable, symmetric multikey file structure. ACM Trans. Datab. Syst. 9, 1 (Mar.), 38–71.
- NODINE, M. H., GOODRICH, M. T., AND VITTER, J. S. 1993. Blocking for external graph searching. In Proceedings of the 12th Annual ACM Symposium on Principles of Database Systems. ACM, New York, pp. 222–232.
- PAGEL, B.-U., SIX, H.-W., TOBEN, H., AND WIDMAYER, P. 1993. Towards an analysis of range query performance in spatial data structures. In *Proceedings of the 12th Annual ACM Symposium on Principles* of Database Systems. ACM, New York, pp. 214–221.
- RAMASWAMY, S., AND KANELLAKIS, P. C. 1995. OODB indexing by class-division. In *Proceedings of* the ACM SIGMOD International Conference on Management of Data. ACM, New York, pp. 139–150.
- RAMASWAMY, S., AND SUBRAMANIAN, S. 1994. Path caching: A technique for optimal external searching. In Proceedings of the 13th Annual ACM Symposium on the Principles of Database Systems. ACM, New York, pp. 25–35.
- SALZBERG, B., AND TSOTRAS, V. 1999. Comparison of access methods for time-evolving data. ACM Comput. Surv. 31, 2, 158–221.

SAMET, H. 1989. The Design and Analysis of Spatial Data Structures. Addison-Wesley, Reading, Mass.

SHAH, M. A., KORNACKER, M., AND HELLERSTEIN, J. M. 1999. AMDB: A visual access method development tool. In User Interfaces to Data Intensive Systems (UIDIS). pp. 130–140.

- SMID, M., AND OVERMARS, M. 1990. Maintaining range trees in secondary memory. Part II: Lower bounds. Acta Inf. 27, 453–480.
- SUBRAMANIAN, S., AND RAMASWAMY, S. 1995. The P-range tree: A new data structure for range searching in secondary memory. In *Proceedings of the 6th ACM-SIAM Symposium on Discrete Algorithms*. ACM, New York, pp. 378–387.
- TURÁN, P. 1941. An extrenal problem in graph theory (in Hungarian). Mat. Fiz. Lapok. 48, 435-452.
- VAIDYA, P. M. 1989. Space-time trade-offs for orthogonal range queries. *SIAM J. Comput.* 18, 4 (Aug.), 748–758.
- VAN LINT, J. H., AND WILSON, R. M. 1992. A Course in Combinatorics. Cambridge University Press, Cambridge, Mass.
- VENGROFF, D. E., AND VITTER, J. S. 1996. Efficient 3-D range searching in external memory. In Proceedings of the 28th Annual ACM Symposium on the Theory of Computing. ACM, New York, pp. 192–201.

WITTEN, I. H., MOFFAT, A., AND BELL, T. C. 1999. Managing Gigabytes: Compressing and Indexing Documents and Images. Morgan-Kaufmann, San Francisco, Calif.

YAO, A. C. 1981. Should tables be sorted? J. ACM 28, 3, 615-628.

YAO, A. C. 1982. Space-time tradeoff for answering range queries. In Proceedings of the 14th Annual ACM Symposium on the Theory of Computing. ACM, New York, pp. 128–136.

RECEIVED OCTOBER 2000; REVISED OCTOBER 2001; ACCEPTED OCTOBER 2001

Journal of the ACM, Vol. 49, No. 1, January 2002.