

Lower Bounds for High Dimensional Nearest Neighbor Search and Related Problems

Allan Borodin
Rafail Ostrovsky
Yuval Rabani

Abstract

In spite of extensive and continuing research, for various geometric search problems (such as nearest neighbor search), the best algorithms known have performance that degrades exponentially in the dimension. This phenomenon is sometimes called the *curse of dimensionality*. Recent results [38, 37, 40] show that in some sense it is possible to avoid the curse of dimensionality for the approximate nearest neighbor search problem. But must the exact nearest neighbor search problem suffer this curse? We provide some evidence in support of the curse. Specifically we investigate the exact nearest neighbor search problem and the related problem of exact partial match within the asymmetric communication model first used by Miltersen [43] to study data structure problems. We derive non-trivial asymptotic lower bounds for the exact problem that stand in contrast to known algorithms for approximate nearest neighbor search.

1 Introduction

Background. One of the most intriguing problems concerning search structures in computational geometry is the following: Given n points in a d -dimensional Euclidean space (the *database*), pre-process the points so that queries of the form “*find in the database the closest point to location x* ” can be answered quickly. More generally, we can define this *nearest neighbor search* (NNS) problem in any vector space, and with any metric (or even with a non-metric distance function). Recently, theoretical research into this problem gained some momentum, inspired in part by applications to multimedia information retrieval and data mining [49, 20, 47, 27, 34, 50, 10, 26].

Trivially, one solution to the problem is to store the raw data, and in response to a query, to compute the distance from the query to each of the n points. Typically, as in the Euclidean case, this would take $O(nd)$ storage and $O(nd)$ search time. On the other extreme, if the set of possible queries is finite (e.g., in the Hamming cube), one can store the answer to each possible query in a dictionary keyed by the query. In the Hamming cube (as a typical case), this requires $\exp(d)$ storage, but merely $O(d)$ search time. Thus, the challenge is to find a solution enjoying the good aspects of both trivial solutions; that is, storage polynomial (preferably linear) in nd (the size of the data set), and search time polynomial (again preferably linear) in d (the size of the query). This problem is non-trivial in the range¹ $\log n \ll d \ll n^\kappa$, for all $\kappa > 0$. (In non-discrete spaces, such as the Euclidean case, even smaller values of d are challenging and for small d we might rephrase the challenge to allow search time polynomial, preferably linear, in $d + \log n$.)

To date, no such solution is known for arbitrary n and d in any reasonable setting (such as Euclidean space, or the Hamming cube). The common “wisdom” among researchers is that simultaneously getting $\text{poly}(nd)$ storage and $\text{poly}(d)$ search time is impossible. Moreover, it has been conjectured that either storage or search time must grow exponentially in d (at least for certain values of n). This conjecture is known as *the curse of dimensionality* [19]. Consequently, much of the present research emphasizes some restriction of the problem, such as considering “typical inputs”, or approximate solutions.

¹We use $f(n) \ll g(n)$ to denote that f is asymptotically smaller than g .

Our results. This paper aims at providing some more persuasive evidence for the curse of dimensionality in a combinatorial setting. We examine NNS in the context of the cell probe model [52]. In the cell probe model, the database is stored in a data structure consisting of m memory cells, each containing b bits. A query is answered by probing in sequence t cells (the address of each cell may be a function of the query and of the contents of previously probed cells). We exploit the connection between asymmetric communication complexity and the cell probe model [44] to derive tradeoffs between the size m of the data structure and the number $t \cdot b$ of bits retrieved.

Specifically, consider the d -dimensional Hamming cube $C_d = \{0, 1\}^d$. We analyze the communication game between Alice, who gets a query $q \in C_d$, and Bob, who gets a database $D \in C_d^n$ (think of D as a set of n points in C_d). They have to output one iff the minimum Hamming distance from q to a point in D is at most λ , and otherwise they have to output zero. The threshold $\lambda \in \{0, 1, 2, \dots, d\}$ is a fixed parameter of the communication game. Notice that the function that Alice and Bob have to compute is a decision version of NNS. We refer to it as the λ -neighbor problem. Clearly, the NNS problem is at least as hard as the λ -neighbor problem. We show that if there is a randomized two-sided error protocol² to compute the λ -neighbor problem where Alice sends a bits and Bob sends b bits, then either $a = \Omega(\log n \log d)$ or $b = \Omega(n^{1-\epsilon})$, for every $\epsilon > 0$. We derive these bounds using the general *richness* technique developed in [46]. In fact, using this richness technique we can prove much stronger lower bounds for randomized one-sided error protocols (and hence for deterministic protocols) for the *complement* of the λ -neighbor problem. Namely, for every ϵ , there is a δ such that either Alice sends $(1 - \epsilon)d$ bits or Bob sends at least δnd bits. However, as we will show, a “direct application” of this method to the λ -neighbor problem appears to be impossible.³ Our main conceptual contribution then is a way to restrict the instances of the λ -neighbor problem to a subset with a nicely structured communication matrix so as to be able to apply the richness technique to the λ -neighbor problem.

We achieve the restriction by considering a different well studied search problem of significant importance in its own right. In the *exact partial match problem*, the database consists of n points in the Hamming cube C_d . A query has 0-1 values assigned to some of the coordinates. The other coordinates are *don't cares*. In reply, we must check if the query matches one of the points in the database, comparing positions with assigned 0-1 values only. (Notice that the implicit distance function here is not a metric.) We show an easy reduction from the exact partial match problem to the λ -neighbor problem, for some value of λ . The reduction produces restricted instances of the λ -neighbor problem. We then proceed to show lower bounds for the communication complexity of exact partial match. We have to further restrict the instances to allow for the application of the richness technique.

The bounds on communication complexity imply lower bounds on exact partial match and on NNS in the Hamming cube. By a simple reduction, we also get similar results for instances of points in $\ell_p^d(\mathbb{R}^d)$ with distances measured by the L_p norm), for every $1 \leq p < \infty$. We show that if a nearest neighbor query is answered in a constant number of probes t , then we either need super-polynomial storage ($2^{\Omega(\log n \log d)}$ cells), or else we need to retrieve nearly-linear $\Omega(n^{1-\epsilon})$ bits from the database. These lower bounds stand in contrast with recent positive results on approximate search (see below). Alternatively, if we restrict ourselves to cells of size $\text{poly}(d)$, then a polynomial number of cells is unattainable unless the number of rounds $t = \Omega(\log d)$. This improves upon the $\Omega(\sqrt{\log d})$ randomized two-sided error lower bound claimed in [46] for the “notoriously difficult” partial match problem. (Our definition of the exact partial match problem differs from the partial match problem defined in [46]. However, a d -dimensional instance of our problem can be easily embedded into a $2d$ dimensional instance of their problem so that our $\Omega(\log d)$ lower bound does apply to their problem as well.) We note in passing that deriving strong lower bound tradeoffs in the unrestricted cell probe model is a well-recognized and fundamental open problem (see

²For definiteness we say that a two-sided error protocol returns the correct answer with probability at least $2/3$. A one-sided error protocol for a decision problem never incorrectly answers YES (=1) and incorrectly answers NO (=0) with probability at most $1/3$.

³As discussed in the related work mentioned below, the known upper bounds for an approximate version of the λ -neighbor problem use randomization with (depending on the implementation) 1-sided or 2-sided error. A 1-sided error lower bound for the non λ -neighbor problem does not imply a randomized lower bound for the λ -neighbor problem (although it does imply a *deterministic* lower bound). Hence our lower bound for the non λ -neighbor problem does not serve to distinguish between the exact and approximate versions of the nearest neighbor problem. Of course, a 2-sided error lower bound for non λ -neighbor implies the same bound for λ -neighbor.

below).

Related work. There is an extensive body of research concerning nearest neighbor problems for small dimensional (e.g. 2 and 3) Euclidean space (see for example the text by de Berg *et al* [13]). Dobkin and Lipton’s seminal paper [21] marks the beginning of work on the Euclidean case of arbitrary dimension. They achieve a discretization of the problem, so that (super-) exponential storage can be used to answer queries quickly. Dobkin and Lipton use $O\left(n^{2^{d+1}}\right)$ storage to allow $O(2^d \log n)$ search time. Clarkson [18] improves the storage requirement to $O\left(n^{(1+\delta)\lceil d/2 \rceil}\right)$, paying $d^{O(d)} \log n$ search time. Improvements by Yao and Yao [53], Matoušek [41], and Agarwal and Matoušek [1] still give exponential in d storage and search time. Finally, Meiser [42], gives the best result to date (in terms of search time) — $O(d^5 \log n)$ search time using $O\left(n^{d+1+\delta}\right)$ storage.⁴

In the *approximate nearest neighbor search* (approximate NNS) problem, a query is answered by finding a database point whose distance from the query is within a factor of $(1 + \varepsilon)$ of the distance to the closest database point. Usually, the parameter ε is fixed in the pre-processing phase. (To avoid confusion, we refer to the version of the problem requiring an exact answer as the *exact* NNS problem.) For approximate NNS in Euclidean space, Arya and Mount [6] give $O(1/\varepsilon)^d \log n$ search time using $O(1/\varepsilon)^d n$ storage. Clarkson [19] improves the dependence on ε to $(1/\varepsilon)^{(d-1)/2}$. Arya, Mount, Netanyahu, Silverman, and Wu [7] give $O((d/\varepsilon)^d \log n)$ search time using $O(n \log n)$ storage (the pre-processing does not depend on ε). In comparison with exact search, these results have better storage requirements. Still, their search time is better than the trivial $O(nd)$ for small d only (as are the results for exact search excluding Meiser’s).

Kleinberg [38] gives $O(n + d \log^3 n)/\varepsilon^2$ search time using $nd \log^{O(1)} n/\varepsilon^2$ storage, thus providing asymptotic improvement over the trivial search time for all (non constant) d using polynomial storage. Another algorithm of Kleinberg (in the same paper) gives $d^2 \log^{O(1)} n/\varepsilon^2$ search time using $O(n \log d/\varepsilon^2)^{2d}$ storage. This is better than Meiser’s exact search time, using similar storage. Indyk and Motwani [37] give improved bounds of $d \log^{O(1)} n$ search time using $O(1/\varepsilon)^d n \log^{O(1)} n$ storage. Their result extends to any L_p norm. They also show that for $1 \leq p \leq 2$, polynomial $(nd)^{O(1/\varepsilon^2)}$ storage can be used to answer correctly *most* queries in time $d \log^{O(1)} n/\varepsilon^2$. Independent of [37], Kushilevitz, Ostrovsky, and Rabani [40] give $d^2 \log^{O(1)} n/\varepsilon^2$ search time (for *all* queries) using $(nd)^{O(1/\varepsilon^2)}$ storage. Their result holds for the L_1 norm too (in particular for the cube, with somewhat better bounds). With the exception of the first algorithm of Indyk and Motwani [37] (which is still exponential in d for small ε), all of these approximate NNS algorithms are randomized algorithms. Of related interest are results on approximate NNS for a large ε by Bern [14] and Chan [16] (for Euclidean space), and by Indyk [36] (for the L_∞ norm).

Thus, approximate NNS does not suffer from the curse of dimensionality,⁵ at least not from the point of view of randomized algorithms and asymptotic bounds for fixed ε . In fact, most of the results mentioned above can be stated in terms of the cell probe model, using a small number of probes. For example, a λ -neighbor version of the algorithm of Kushilevitz et al. [40] for the cube can be implemented as a randomized two-sided error one round cell probe algorithm with $n^{O(1/\varepsilon^2)}$ cells, each containing one bit.⁶ Using several such structures (for different distances) an approximate NNS implementation takes $O(\log \log d)$ rounds (using binary search on a geometric progression of λ ’s). Therefore, the lower bounds in this paper provide some evidence that in high dimension, exact NNS is indeed far more difficult than approximate NNS.

⁴Meiser considers the more general problem of point location in arrangements of hyperplanes and obtains the storage bound $O\left(n^{d+\delta}\right)$ where n is the number of hyperplanes. The nearest neighbor problem with n data points is easily transformed into this general problem by considering the $\binom{n}{2}$ hyperplanes constructed by bisecting each pair of data points. To obtain the stated space bound, the given problem is transformed to a point location query among n hyperplanes in $(d+1)$ -space which is obtained by replacing each point $(a_1, a_2, \dots, a_d)a_d$ by the hyperplane $x_{d+1} = (x_1 - a_1)a_1 + \dots + (x_d - a_d)$ and replacing the query point (q_1, \dots, q_d) by $(q_1, \dots, q_d, \infty)$. The hyperplane directly below the transformed query point corresponds to the original query point’s nearest neighbor.

⁵However, as one of the referees has pointed out, since the stated complexity bounds are now exponential in $1/\varepsilon^2$, we might say that the curse of dimensionality has been replaced by a “curse of exactness”.

⁶Using d bits per cell, it is possible to derive a one-sided error implementation.

It is easy to see that the exact full match problem (equivalent to 0-neighbor) also does not suffer the curse of dimensionality. Indeed, one can pre-process the database by using a perfect hashing function thereby permitting a query search to be performed in one table lookup. Alternatively, one can arrange the database in a d -depth, n -leaf tree and then perform the query search by a simple search of the tree in time $O(d)$. We do not know of any analogous results for the partial match problem.

However, to the best of our knowledge, lower bounds for exact NNS (or the partial match query problem) in high dimensions do not seem sufficiently “convincing” to justify the curse of dimensionality conjecture. That is, either the models with respect to which lower bounds have been established seem quite restricted or the bounds are quite weak. One nice example of a well structured model (for both dynamic and static data structure problems) is Fredman’s [28, 29, 30] semi-group model. The model is designed for searching problems (e.g., range queries) in which a semi-group value is associated with each data point and one wants to retrieve the semi-group sum of all data points in some specified set (e.g., satisfying a partial match or more generally satisfying a range query). The static model allows pre-processing of sums of arbitrary subsets of database points and Fredman derives strong lower bound tradeoffs between memory (the number of pre-computed subsets) and search time⁷ (the number of semi-group additions on these pre-computed subsets). While this model can be used for decision problems (by letting the semi-group be the Boolean values under the operation of logical OR) the model is clearly restrictive in the limited way information can be obtained from the data structure.

In the setting of real Euclidean space (i.e., \mathbb{R}^d with the usual L_2 metric), an appropriate model is the algebraic computation tree (or the closely related algebraic decision tree). Simply stated for real valued inputs, one can compute and test the signs of polynomials of these inputs. The complexity of such an algebraic computation tree is usually taken as the number of multiplications and tests. The algebraic computation model was first introduced by Ben-Or [12] (having been preceded by the algebraic decision tree model where one only counts tests but then restricts the degree of the polynomials allowed). Following a substantial chain of papers, Grigoriev and Karpinski [32] prove a randomized lower bound to determine membership in a polyhedron, the lower bound being (roughly speaking) the logarithm of the number of faces. Since there are instances of n data points x_1, \dots, x_n in \mathbb{R}^d which give rise to a Voronoi diagram with $\Theta(n^{\lceil d/2 \rceil})$ faces, we can apply the polyhedron lower bound to derive an $\Omega(d \log n)$ randomized algebraic computation tree lower bound for deciding if a given query $q \in \mathbb{R}^d$ is closest to a given data point x_i . Computation tree models do not reflect storage usage and indeed this lower bound holds independent of the storage allowed and hence the bound provides a somewhat matching counterpart to the $d^5 \log n$ upper bound of Meiser. In the context of our paper, we note that the model does not allow modular operations (thereby precluding certain hashing functions), and the analysis used for the Grigoriev and Karpinski lower bound is not applicable for the case of finite domains such as the Hamming cube.

A model which does capture hashing and more combinatorial settings is introduced in Rivest [48] and further developed in Dolev, Harari and Parnas [22] and Dolev, Harari, Linial, Nisan and Parnas [23]. Rivest studies the *all* partial match problem and Dolev *et al.* study the *all* λ -neighbor problem where for each problem all database points satisfying the query must be found. In this model, each database point is hashed to a bucket and interesting tradeoffs are established between the number of buckets containing database points satisfying the query and the maximum size of a bucket. Another lower bound is by Indyk [36]. He shows a lower bound for approximate NNS under the L_∞ norm in the *indexing model* of Hellerstein, Koutsoupias, and Papadimitriou [35]. The indexing model tries to capture the cost of using external memory devices for large data sets, and appears to be computationally more restricted than the cell probe model. Indyk shows that the *superset query* problem of Hellerstein et al. reduces to $(1 + \epsilon)$ -approximate NNS under the L_∞ norm, for any $\epsilon < 1$. The lower bound that Hellerstein et al. give for superset query is weak, unless the storage redundancy is quite small. This does not seem to pose a serious theoretical restriction on a solution, though it may address important considerations in practice.

The cell probe model was formulated by Yao [52]. It is considered as the most general data structure model for proving lower bounds. Ajtai [2] and Xiao [51] obtain further lower bounds in this model.

⁷Associated with each possible query is a straight line program whose operations are either the semi-group addition $v_i = v_j + v_k$ or the scalar multiplication of a semi-group value $v_i = c \cdot v_j = v_j + v_j + \dots + v_j$ by a positive integer c .

Miltersen [43, 44, 45] pioneered the connection between the cell probe model and asymmetric communication complexity. Miltersen, Nisan, Safra, and Wigderson [46] provide general methods for proving lower bounds for asymmetric communication complexity, including the richness technique used here. For more information on communication complexity, see the book by Kushilevitz and Nisan [39]. As Miltersen et al. observe, communication complexity cannot prove strong lower bounds for the cell probe model without restrictions (say on the number of rounds). Furthermore, they observe that lower bounds for the general cell probe model imply time-space tradeoffs for branching programs, one of the notoriously difficult problems in computational complexity. For the best lower bound on branching programs to date, and additional references, see Beame, Saks, and Thathachar [9], and the two recent papers of Ajtai [3, 4].

As mentioned above, Rivest [48] analyses hashing based algorithms for partial match. (See Bentley and Sedgewick [11] for a more recent historical account.) Rivest conjectures that *any* $O(nd)$ -sized data structure would require $\Omega(n^{1-s/d})$ time to search, where s is the number of exposed coordinates.⁸ Also as mentioned above, using a *round elimination* technique, Miltersen et al. [46] claim a lower bound of $\Omega(\sqrt{\log d})$ on the number of probes required to find a partial match in the cell probe model with $(nd)^{O(1)}$ cells, each containing $\text{poly}(d)$ bits. We improve this bound to $\Omega(\log d)$. Moreover, for the corresponding communication problem our lower bounds are on the total number of bits transmitted by either side, irrespective of the number of rounds.

Independent of (and complimentary to) our work, Chakrabarti, Chazelle, Gum and Lvov [15] have established an $\Omega(\log \log d / \log \log \log d)$ deterministic lower bound in the cell probe model for finding an approximate nearest neighbor. The approximation factor here can be as large as $2^{\lfloor (\log d)^{1-\epsilon} \rfloor}$ for any $\epsilon > 0$.

Subsequent to our conference publication, Barkol and Rabani [8] established a significantly improved lower bound for the λ -neighbor problem. They are able to derive a lower bound for *two-sided error* randomized asymmetric communication complexity protocols that is quite comparable to our one-sided error bound in Theorem 2.2. Namely, again assuming some necessary reasonable bounds on the range of the dimension d , they show:

Let $\lambda = \frac{d}{2} - \frac{2}{\ln 2} \sqrt{d \log n}$. For every $0 < \epsilon < 1$ there exist $\delta > 0$, such that in every two-sided error protocol for the λ -neighbor problem in C_d , when n is sufficiently large then either the query side sends at least ϵd bits, or the database side sends at least n^δ bits.

2 A Communication Complexity Lower Bound for the non λ -neighbor problem

We use the Miltersen et al. [46] *richness technique* to analyze the asymmetric communication game between Alice, who gets a query $q \in C_d$, and Bob, who gets a database $D \in C_d^n$.⁹ For the *non λ -neighbor problem*, the players have to output one iff there does *not* exist a data base point $x \in D$ whose Hamming distance to the query q is at most λ . An $[a, b]$ protocol for a communication problem is a protocol in which Alice sends at most a bits and Bob sends at most b bits.

For the sake of completeness, we review the Miltersen et al. [46] richness technique (for one-sided error protocols). We associate a communication matrix M_f with any communication problem f . Namely we index the rows by the possible inputs for Alice and the columns by the possible inputs for Bob and the x, y entry of M_f is the value $f(x, y)$. A communication problem f is $[u, v]$ -rich iff its communication matrix M_f has at least v columns each containing at least u ones. The richness technique is captured by the following *richness lemma*:

⁸The Rivest conjecture was stated without any conditions on s but it seems reasonable to believe that the conjecture assumes that s is not very small (e.g. s is a constant) or very large (e.g. $d - c$, c a constant).

⁹That is, in order to simplify the analysis, we allow repetitions in the database. The results that follow would not change in any significant way if a database was defined to be n distinct vectors.

Lemma 2.1 (Miltersen et al. [46]). *Let f be $[u, v]$ -rich. If f has a randomized one-sided error $[a, b]$ protocol, then M_f contains a submatrix of dimension at least $u/2^{a+2} \times v/2^{a+b+2}$ containing only 1-entries.*

The richness technique then is to show that a given communication problem is sufficiently rich yet does not contain large submatrices containing only 1-entries. We now apply this technique to the non λ -neighbour problem. We wish to derive asymptotic tradeoffs for “large” dimensions. Thus we consider an infinite sequence of problems for all sufficiently large values of n , and dimension $d = d(n)$. Here we assume $\log n \ll d$.

Theorem 2.2. *For every $0 < \varepsilon < 1$ there exist δ and ν , such that for every $n > \nu$ there exists λ for which in every protocol for non λ -neighbor problem in C_d , either the query side sends at least $(1 - \varepsilon)d$ bits, or the database side sends at least δnd bits.*

Let $B_d(\lambda)$ denote the Hamming ball of radius λ around an arbitrary vector in C_d and let $B_d(A, \lambda)$ denote the set of cube points at distance at most λ from a point in $A \subseteq C_d$. In the proof of this theorem, we use the following standard inequalities (see, for example, [5]) :

Fact 2.3 (Entropy bound). *For $0 < p < \frac{1}{2}$, $|B_d(pd)| \leq 2^{(H(p)+o(1))d}$, where $H(p) = -p \log p - (1 - p) \log(1 - p)$ is the entropy function.*

Fact 2.4 (Harper’s isoperimetric inequality [33]). *For $A \subseteq C_d$, let $r > 0$ be such that $|A| \geq |B_d(r)|$. Then, for every $\lambda > 0$, $|B_d(A, \lambda)| \geq |B_d(r + \lambda)|$.*

Fact 2.5 (Chernoff’s bound). *For every $a > 0$, $|B_d(d/2 - a)| \leq e^{-a^2/d} \cdot 2^d$.*

Proof of Theorem 2.2. We apply the richness technique. Take $\lambda = \frac{d}{2} - \sqrt{d \ln(2n)}$. (The hardest case seems to be to distinguish between a distance of at most $\frac{d}{2}$ and a distance of at least $\frac{d}{2} + 1$.) Using Claim 2.5, $n|B_d(\lambda)| \leq 2^{d-1}$. Thus, for every database at most half the queries are within a distance of λ of one of the points in the database. Therefore, for our choice of λ , non λ -neighbor is $[2^{d-1}, 2^{nd}]$ -rich. We will show that for every $\varepsilon > 0$, there exists δ such that every $2^{\varepsilon d} \times 2^{(1-\delta)nd}$ submatrix of the communication matrix of the non λ -neighbor problem contains a zero entry.

Consider a set Q of queries of cardinality $2^{\varepsilon d}$. Let λ_Q be the largest integer such that $|B_d(\lambda_Q)| \leq |Q|$. By Fact 2.3, there exists a constant $p = p(\varepsilon)$ such that $\lambda_Q \geq pd$. Let ξ be a constant such that $0 < \xi < p - \sqrt{\ln(2n)/d}$. By our assumption that $\log n \ll d$, we can choose ξ arbitrarily close to p when n is sufficiently large. Then,

$$\begin{aligned}
|B_d(Q, \lambda)| &\geq |B_d(pd + \lambda)| \\
&= \left| B_d \left(\left(\frac{1}{2} + p \right) d - \sqrt{d \ln(2n)} \right) \right| \\
&\geq \left| B_d \left(\left(\frac{1}{2} + \xi \right) d \right) \right| \\
&> 2^d - \left| B_d \left(\left(\frac{1}{2} - \xi \right) d \right) \right| \\
&\geq 2^d - e^{-\xi^2 d} 2^d \\
&= 2^d - 2^{(1-\xi^2)d},
\end{aligned}$$

where the first inequality follows from Fact 2.4, and the fourth inequality follows from Fact 2.5. Consider a probability distribution over C_d , with all points equally likely. From the final inequality above, the probability a random point from this distribution does not fall in $B_d(Q, \lambda)$ is less than $\frac{2^{(1-\xi^2)d}}{2^d} = 2^{-\xi^2 d}$. Consider now a distribution over C_d^n , with all databases equally likely. This distribution is equivalent to taking n independent samples from the uniform distribution over C_d . Thus, the probability that none of the points of a random n point database fall in $B_d(Q, \lambda)$ is less than $2^{-\xi^2 nd}$. Put $\delta = \xi^2$, and the theorem follows. ■

3 A Communication Complexity Lower Bound for the partial match problem

In the exact partial match problem, the database consists of n vectors v^1, v^2, \dots, v^n in C_d , and a query is a vector in $\tilde{C}_d = \{0, 1, *\}^d$. A query q matches a vector $v \in C_d$ iff for all $j \in \{1, 2, \dots, d\}$, either $q_j = *$ or $q_j = v_j$. A query q matches the database iff there exists $i \in \{1, 2, \dots, n\}$ such that q matches v^i . We say that the coordinates j for which $q_j \neq *$ are *exposed* in q . We also refer to any j such that $q_j = *$ as a *don't care coordinate*, or simply as a *don't care*. For simplicity, we assume that n is a power of 2 (thus $\log n$ is an integer). We shall see that lower bounds for the partial match problem imply corresponding lower bounds for the exact NNS problem. (As far as we know, lower bounds for exact NNS cannot be used to imply lower bounds for the partial match problem.)

We wish to derive asymptotic tradeoffs for the partial match problem for “large” dimensions. Thus we again consider an infinite sequence of problems for all sufficiently large values of n , each of dimension $d = d(n)$. We assume that $\log n \ll d \leq n^\kappa$ for a positive $\kappa < 1$.¹⁰

The set of possible databases again includes all 2^{nd} choices in C_d^n . The set of possible queries is restricted to

$$Q_{n,d} = \{q; |\{i; q_i \neq *\}| = \log n + 1\}.$$

In words, the queries are restricted to have exactly $\log n + 1$ exposed coordinates. Notice that the number of possible queries is exactly

$$2n \binom{d}{\log n + 1} = 2^{O(\log n \log d)}.$$

Our main result is the following theorem:

Theorem 3.1. *Let $0 < \varepsilon < 1 - \kappa$ be fixed. Suppose there is an $[a, b]$ (deterministic or randomized, two-sided error) communication protocol for NPM. Then, either $a = \Omega(\log n \log d)$, or $b = \Omega(n^{1-\varepsilon})$.*

We first present the proof for one-sided error protocols. We then extend the proof to handle two-sided error protocols. The latter is somewhat more involved, yet it builds on the ideas for the one-sided error case. Throughout this section ε refers to the ε as stated in the theorem.

Lemma 3.2. *NPM is $[\frac{R}{5}, \frac{C}{6}]$ -rich where $R = 2n \binom{d}{\log n + 1}$ is the number of rows in M_{NPM} and $C = 2^{nd}$ is the number of columns in M_{NPM} .*

Proof. As in the previous section, consider a probability distribution over C_d , with all points equally likely. For $q \in Q_{n,d}$, the probability that q matches a random vector from this distribution is $\frac{1}{2n}$. As before, consider now a distribution over C_d^n , with all databases equally likely. The probability that q does not match a random database from this distribution is

$$\left(1 - \frac{1}{2n}\right)^n \geq e^{-1}.$$

If, however, less than $\frac{1}{6}$ of the columns (databases) contain at least a fraction of $\frac{1}{5}$ ones entries, then the fraction of ones entries in the communication matrix does not exceed $\frac{1}{3}$, a contradiction. ■

We say that two queries $q, q' \in Q_{n,d}$ are *consistent* if they agree on all coordinates which are exposed in both vectors. We say that q and q' are ε -*neighbors* iff they are consistent and the number of coordinates which are exposed in both of them is at least $\varepsilon \log n$. The notion of ε -neighbors is useful because of the following lemma.

Lemma 3.3. *If $q, q' \in Q_{n,d}$ are not ε -neighbors, then the fraction of vectors $v \in C_d$ such that both q and q' match v is at most $\frac{1}{4n^{2-\varepsilon}}$.*

¹⁰The assumption that the dimension $d \leq n^\kappa$ is only used for the case of 2-sided errors.

Proof. If q and q' are not consistent, then by definition there is no vector that they both match. Otherwise, all of the coordinates that are exposed in both queries must have the same value in both. Furthermore, if both q and q' match a vector v , then for every j which is exposed in either q or q' , v_j must equal the exposed bit. The total number of different coordinates exposed in either q or q' is at least $(2 - \varepsilon) \log n + 2$. The lemma follows from computing the fraction of vectors with these bits fixed. ■

Lemma 3.4. *For every $\delta > 0$ there exists ν such that for all $n > \nu$, for $d = d(n)$ the following holds. For every $q \in Q_{n,d}$, the number of its ε -neighbors is less than*

$$\binom{d}{\log n + 1}^{1-\varepsilon+\delta}.$$

Proof. The number N of neighbors of q is given by

$$\begin{aligned} \sum_{j=0}^{(1-\varepsilon)\log n+1} \binom{\log n + 1}{j} \binom{d - \log n - 1}{j} 2^j &\leq 2n^{1-\varepsilon} \binom{d - \log n - 1}{(1-\varepsilon)\log n + 1} \sum_{j=0}^{\log n+1} \binom{\log n + 1}{j} \\ &= 4n^{2-\varepsilon} \binom{d - \log n - 1}{(1-\varepsilon)\log n + 1} \\ &< 4n^{2-\varepsilon} \binom{d}{(1-\varepsilon)\log n + 1}. \end{aligned}$$

Now,

$$\begin{aligned} \frac{\binom{d}{\log n+1}}{\binom{d}{(1-\varepsilon)\log n+1}} &= \frac{((1-\varepsilon)\log n + 1)! (d - (1-\varepsilon)\log n - 1)!}{(\log n + 1)! (d - \log n - 1)!} \\ &\geq \left(\frac{d - \log n}{\log n + 1} \right)^{\varepsilon \log n} \\ &= \left(1 - \frac{\log n}{d} \right)^{\varepsilon \log n} \left(\frac{d}{\log n + 1} \right)^{\varepsilon \log n} \\ &\geq n^{-\varepsilon} \left(\frac{d}{\log n + 1} \right)^{\varepsilon \log n} \\ &\geq n^{-\varepsilon} \left(\frac{ed}{\log n + 1} \right)^{-\varepsilon} n^{-\varepsilon \log e} \left(\frac{d}{\log n + 1} \right)^{\varepsilon}. \end{aligned}$$

Therefore,

$$N < 4n^{2+\varepsilon \log e} \left(\frac{ed}{\log n + 1} \right)^{\varepsilon} \left(\frac{d}{\log n + 1} \right)^{1-\varepsilon}.$$

Let $\delta > 0$. As $d \gg \log n$, for n sufficiently large,

$$4n^{2+\varepsilon \log e} \left(\frac{ed}{\log n + 1} \right)^{\varepsilon} \leq \left(\frac{d}{\log n + 1} \right)^{\delta}. \quad \blacksquare$$

We call a set $I \subseteq Q_{n,d}$ ε -independent iff for every $q, q' \in I$ such that $q \neq q'$, q and q' are not ε -neighbors. We conclude from the above lemma

Lemma 3.5. *For every δ such that $0 < \delta < \varepsilon/2$, there is ν such that for every $n > \nu$, every $R \subseteq Q_{n,d}$ of cardinality at least $\binom{d}{\log n+1}^{1-\delta}$ contains an ε -independent subset $\text{ind}(R)$ of cardinality $n^{1-\varepsilon}$.*

Proof. Let n be sufficiently large so that Lemma 3.4 holds for δ , and $\binom{d}{\log n+1}^{\varepsilon-2\delta} \geq n^{1-\varepsilon}$. Consider the graph whose nodes are the elements of R , and a pair of nodes is an edge iff its endpoints are ε -neighbors. By Lemma 3.4, the maximum degree in this graph is less than $\binom{d}{\log n+1}^{1-\varepsilon+\delta}$. Therefore, it has an ε -independent set of size at least $\binom{d}{\log n+1}^{\varepsilon-2\delta} \geq n^{1-\varepsilon}$. ■

For $q \in Q_{n,d}$ we denote $\mathcal{D}(q) = \{D \in C_d^n; q \text{ does not match } D\}$. For $R \subseteq Q_{n,d}$ we abuse notation and denote $\mathcal{D}(R) = \bigcap_{q \in R} \mathcal{D}(q)$.

Lemma 3.6. *For every δ such that $0 < \delta < \varepsilon/2$, there is ν such that for every $n > \nu$, if $R \subseteq Q_{n,d}$ has cardinality at least $\binom{d}{\log n+1}^{1-\delta}$, then $\mathcal{D}(R)$ has cardinality less than $2^{nd-n^{1-\varepsilon}/4}$.*

Proof. We examine the subset $\text{ind}(R)$ of cardinality $n^{1-\varepsilon}$ from Lemma 3.5. Consider a distribution over C_d , with all points equally likely. The probability that any $q \in \text{ind}(R)$ matches a random point from this distribution is exactly $\frac{1}{2n}$. Let $q, q' \in \text{ind}(R)$, $q \neq q'$. By Lemma 3.3, the probability that both q and q' match a random point from this distribution is at most $\frac{1}{4n^{2-\varepsilon}}$. Therefore, by the inclusion-exclusion principle, the probability that a random point from the distribution is matched by at least one point in $\text{ind}(R)$ is at least

$$n^{1-\varepsilon} \cdot \frac{1}{2n} - \binom{n^{1-\varepsilon}}{2} \cdot \frac{1}{4n^{2-\varepsilon}} \geq \frac{1}{2n^\varepsilon} - \frac{1}{8n^\varepsilon} = \frac{3}{8n^\varepsilon}.$$

Now, consider a distribution over C_d^n with all databases equally likely. The probability that none of the points in $\text{ind}(R)$ match a random database from this distribution is at most $(1 - 3/8n^\varepsilon)^n \leq e^{-n^{1-\varepsilon}/4}$. ■

Lemma 3.6 implies the following

Corollary 3.7. *For every δ such that $0 < \delta < \varepsilon/2$, there is ν such that for every $n > \nu$, the communication matrix of NPM does not contain a $\binom{d}{\log n+1}^{1-\delta} \times 2^{nd-n^{1-\varepsilon}/4}$ 1-monochromatic rectangle.*

Proof. Otherwise, we have a set $R \subseteq Q_{n,d}$ with $|R| \geq \binom{d}{\log n+1}^{1-\delta}$, and $|\mathcal{D}(R)| \geq 2^{nd-n^{1-\varepsilon}/4}$, in contradiction with Lemma 3.6. ■

We now can conclude the proof of Theorem 3.1 for one-sided errors by applying the richness Lemma 2.1.

It remains to show how to extend this proof to the case of two-sided error. Miltersen *et al.* [46] prove a second form of the richness lemma which makes it possible to prove lower bounds for randomized algorithms having two-sided error. Essentially, instead of showing that every sufficiently large submatrix is not 1-monochromatic, we now need to show that every sufficiently large submatrix has a constant fraction of zeros. We now indicate how to apply this form of the richness lemma to establish Theorem 3.1 for two-sided error protocols.

Lemma 3.8. *For every δ such that $0 < \delta < \varepsilon/2$, there is ν such that for every $n > \nu$, every $R \subseteq Q_{n,d}$ of cardinality at least $\binom{d}{\log n+1}^{1-\delta}$ can be partitioned into sets $I_0, I_1, I_2, \dots, I_f$ such that the following hold:*

1. I_0 contains at most half of R ; and,
2. for $j = 1, 2, \dots, f$, I_j is an ε -independent set with $|I_j| = 2n^{1-\varepsilon}$, and
3. $f \ll 2^{n^\kappa}$. (Recall $\kappa < 1 - \varepsilon$ and $d \leq n^\kappa$.)

Proof. As long as at least half of R remains, repeatedly apply Lemma 3.5 to pick a set I_j with the desired properties, then remove it from R . (To be more precise, we have to slightly modify Lemma 3.5 to make each I_j have size $2n^{1-\varepsilon}$ assuming $R \subseteq Q_{n,d}$ is of cardinality at least $\frac{1}{2} \binom{d}{\log n+1}^{1-\delta}$.) Now f is trivially smaller than the total number of queries. Thus, $f \leq 2n \binom{d}{\log n+1} \ll 2^{n^\kappa}$, for sufficiently large n (as we assume that $d \leq n^\kappa$). ■

We want most databases to match many points in the sets I_j , $j > 0$. For $I \subseteq Q_{n,d}$, we denote by $\mathcal{D}(I)$ the set of databases that match less than $\frac{1}{100}$ of the queries in I .¹¹ We have

¹¹ $\frac{1}{100}$ is a somewhat arbitrary constant.

Lemma 3.9. *For every $\varepsilon > 0$, there is ν such that for every $n > \nu$ the following holds: For every ε -independent $I \subseteq Q_{n,d}$ such that $|I| = 2n^{1-\varepsilon}$, we have $\mathcal{D}(I) \leq 2^{nd-n^{1-\varepsilon}/50}$.*

Proof. Consider a database chosen at random, all databases equally likely. We think of the database as being chosen in sequence, one point at a time, each point chosen independently of the others from a uniform distribution over C_d . Let x_1, x_2, \dots, x_n be the database points. Let $M_0, M_1, M_2, \dots, M_{n^{1-\varepsilon}}$ be the following subsets of I : $M_0 = \emptyset$. M_i includes all of M_{i-1} , and if any query q in $I \setminus M_{i-1}$ matches one of the database points $x_{(i-1)n^\varepsilon+1}, \dots, x_{in^\varepsilon}$, then M_i also contains one (arbitrarily chosen) such matching query q ; thus, $|M_{i-1}| \leq |M_i| \leq |M_{i-1}| + 1$. Let $X_i = |M_i|$.

We show that for all i , $0 \leq i < n^{1-\varepsilon}$, $\Pr[X_{i+1} > X_i] \geq 1 - e^{-1/4}$. For all i , $|I \setminus M_i| \geq n^{1-\varepsilon}$. Therefore, by the proof of Lemma 3.6, the probability that a random database point is matched by one of the queries in $I \setminus M_i$ is at least $3/8n^\varepsilon$. For n^ε random points we get that the probability that none of these points are matched by a query in $I \setminus M_i$ is $(1 - 3/8n^\varepsilon)^{n^\varepsilon} \leq e^{-1/4}$, for sufficiently large n .

Now, define $Y_0, Y_1, \dots, Y_{n^{1-\varepsilon}}$: Y_0 is the initial expected size of $M_{n^{1-\varepsilon}}$. Y_i is the same expectation after choosing $x_1, \dots, x_{in^\varepsilon}$. Notice that all these expectations are random variables (depending on the choice of $x_1, \dots, x_{in^\varepsilon}$ with Y_0 being a constant). Further notice that $Y_{n^{1-\varepsilon}} = |M_{n^{1-\varepsilon}}|$. By definition, $Y_i = E[Y_{i+1} | Y_i]$. Also, $|Y_i - Y_{i+1}| \leq 1$. Therefore, the sequence $Y_0, Y_1, \dots, Y_{n^{1-\varepsilon}}$ is a martingale. By Azuma's inequality, $\Pr[Y_{n^{1-\varepsilon}} < Y_0 - \lambda n^{(1-\varepsilon)/2}] < e^{-\lambda^2/2}$. Now, by the linearity of expectation, $Y_0 \geq (1 - e^{-1/4}) n^{1-\varepsilon}$. Set $\lambda = \frac{1}{5} n^{(1-\varepsilon)/2}$. We get, $\Pr[Y_{n^{1-\varepsilon}} < (1 - e^{-1/4} - 1/5) n^{1-\varepsilon}] < e^{-n^{1-\varepsilon}/50}$. Finally, notice that $1 - e^{-1/4} - \frac{1}{5} > \frac{1}{50}$ so that

$$\left(1 - e^{-1/4} - 1/5\right) n^{1-\varepsilon} > \frac{1}{100} |I|. \quad \blacksquare$$

Theorem 3.10. *For every δ such that $0 < \delta < \varepsilon/2$, there exists ν such that for all $n > \nu$ the following holds: In the communication matrix of NPM, in every rectangle $R \times D$ with $|R| \geq \binom{d}{\log n+1}^{1-\delta}$ and $|D| \geq 2^{nd-n^{1-\varepsilon}/100}$, at least a fraction of $\frac{1}{400}$ of the entries are zeros.*

Proof. Let n be sufficiently large, and let $R \times D$ be a rectangle satisfying the conditions of the theorem. By Lemma 3.8, at least half the queries in R can be partitioned into disjoint independent subsets I_1, I_2, \dots, I_f , $f \ll 2^{n^\kappa}$, for all $\kappa > 0$ and for n sufficiently large. For any j , $1 \leq j \leq f$, the number of databases that match less than $\frac{1}{100}$ of the queries in I_j is at most $2^{nd-n^{1-\varepsilon}/50}$, by Lemma 3.9. Therefore, the number of databases that match less than $\frac{1}{100}$ of the queries in any of the sets I_j is less than $2^{nd-n^{1-\varepsilon}/50+n^\kappa} \leq 2^{nd-(n^{1-\varepsilon}/100)-1}$, for n sufficiently large and $0 < \kappa < 1 - \varepsilon$. Thus, if we take all $2^{nd-n^{1-\varepsilon}/100}$ databases in D , at least half of them match at least $\frac{1}{100}$ of the queries in every set I_j . The theorem follows because the number of queries in these sets is at least half the total number of queries in R \blacksquare

4 Consequences

Miltersen [44] shows that asymmetric communication complexity lower bounds can be used to derive lower bounds for the cell probe model. Specifically, if there is a deterministic (respectively, randomized) cell probe model solution to a ‘‘data structure’’ problem with parameters m (the number of cells), b (the maximum cell size) and t (the number of probes of the data structure), then there is a deterministic (respectively, randomized) asymmetric communication protocol for this problem with $2t$ rounds of communication¹² in which Alice sends $\log m$ bits in *each* of her messages and Bob sends b bits in *each* of his messages. That is, Alice (respectively, Bob) sends a total of at most $t \log m$ bits (respectively, tb bits). Using Theorem 2.2 and the connection to the cell probe model, we get the following lower bound for the λ -neighbor problem in the cell probe model.

¹²In the asymmetric communication model, a message passed by either Alice or Bob is considered a round.

Theorem 4.1. *Any one-sided randomized algorithm for the non λ -neighbor problem (and hence deterministic algorithm for the λ -neighbor problem or its complement) that makes t probes, either uses $2^{\Omega(d/t)}$ cells, or uses cells of size $\Omega(nd/t)$.*

We point out two extremes of Theorem 4.1:

- If the cell size is $d^{O(1)}$ and the number of cells is polynomial (in n) then the algorithm must make $\Omega(d/\log n)$ probes.
- If the algorithm answers a query in a constant number of probes, then either it uses $2^{\Omega(d)}$ cells, or requires the processing of a cell containing $\Omega(nd)$ bits.

In the same way, using Theorem 3.1 we obtain the following lower bound for the partial match problem.

Theorem 4.2. *Any randomized (two-sided error) cell probe algorithm for the exact partial match problem that makes t probes, either uses $2^{\Omega(\log n \log d/t)}$ cells, or uses cells of size $\Omega(n^{1-\varepsilon}/t)$.*

And now we get the following two extremes of Theorem 4.2:

- If the cell size is $d^{O(1)}$ and the number of cells is polynomial (in n) then the algorithm must make $\Omega(\log d)$ probes.
- If the algorithm answers a query in a constant number of probes, then either it uses $2^{\Omega(\log n \log d)}$ cells, or requires the processing of a cell containing $\Omega(n^{1-\varepsilon})$ bits.

Lower bounds corresponding to Theorems 3.1 and 4.2 can be obtained for the λ -neighbor problem by applying the following reduction:

Theorem 4.3. *Let $Q(d, \lambda)$ denote the set of points in \tilde{C}_d with exactly λ don't cares. Then, there exist functions $\varphi_A : Q(d, \lambda) \rightarrow C_{2d}$ and $\varphi_B : C_d^n \rightarrow C_{2d}^n$ with the following property: If $(q, D) \in Q(d, \lambda) \times C_d^n$, then q matches D iff for $(q', D') = (\varphi_A(q), \varphi_B(D))$, q' is a λ -neighbor of D' . Furthermore, both φ_A and φ_B can be computed efficiently (in linear time).*

Proof. For $x \in \tilde{C}_d$, define $y = \varphi_A(x) \in C_{2d}$ as follows. For $i \in \{1, 3, 5, \dots, 2d-1\}$,

$$y_i y_{i+1} = \begin{cases} x_{\lceil i/2 \rceil} x_{\lfloor i/2 \rfloor} & x_{\lceil i/2 \rceil} \neq * \\ 01 & x_{\lceil i/2 \rceil} = * \end{cases}$$

Now, define φ_B by applying the transformation φ_A to each of the n points in D . Consider any point $x \in D$, and its image $x' \in D'$. Each don't care in q produces one mismatch between q' and x' , regardless of the value of the corresponding coordinate in x . If q matches x , no additional mismatches are produced between q' and x' . Otherwise, there are at least two additional mismatches. ■

Now, consider the r -neighbor problem in ℓ_p^d , $1 \leq p < \infty$, for $0 < r \in \mathbb{R}$. Analogous to its definition for the cube, this problem requires deciding whether or not the minimum distance between a query point and a database of n points is at most r . We have

Theorem 4.4. *For every $p \in \mathbb{R}$, $1 \leq p < \infty$, for every $\lambda \in \{0, 1, 2, \dots, d\}$, there exists $r \in \mathbb{R}$, $r > 0$, such that exact partial match with queries in $Q(d, \lambda)$ reduces to the r -neighbor problem in ℓ_p^{2d} .*

Proof. For $p = 1$, the theorem follows from Theorem 4.3, as the points of C_{2d} are a subset of \mathbb{R}^{2d} , and the Hamming distance is equivalent to the L_1 distance for these points. For $p > 1$, the theorem follows from a monotonicity property of the L_p norm on C_{2d} (viewed as a subset of \mathbb{R}^{2d}): If $w, x, y, z \in C_{2d}$, then $\|w - x\|_1 < \|y - z\|_1$ iff $\|w - x\|_p < \|y - z\|_p$, where $\|\cdot\|_p$ denotes the L_p norm. (Notice that this monotonicity property does not hold for the L_∞ norm.) ■

Finally, we mention some implications to other geometric search problems: First, exact NNS in Euclidean space is a special case of point location in an arrangement of hyperplanes (with $\binom{n}{2}$ hyperplanes, defining the Voronoi diagram). Therefore, our results imply lower bounds for point location.

Next, consider the cube C_d . Notice that the reduction in Theorem 4.3 proves a somewhat stronger claim than mentioned. It shows that exact partial match reduces to the problem of determining whether or not there is a database point at distance *precisely* λ from the query. So, we get a lower bound for this problem as well (for the Hamming cube). Now, consider the cube as a subset of \mathbb{R}^d (for simplicity, we'll use the vectors $\{\pm 1\}^d$ here). The set of cube points at Hamming distance exactly λ from a cube point v lies on the hyperplane $x \cdot v = d - 2\lambda$. Therefore, we get a lower bound for the problem of determining whether or not a query point v lies on one of a collection of n hyperplanes (one hyperplane for each of the n database points).¹³ As Chazelle points out [17], this problem can be viewed as a multi-dimensional generalization of the dictionary problem. The dictionary problem can be stated as follows: In the one-dimensional real line, we have a database of hyperplanes (zero-dimensional flats, i.e. points, in this case), queries are points, and the answer to a query is whether or not it is contained in the database. The problem can be solved in $O(1)$ probes per query via hashing. Our lower bounds show that a similar result for the multi-dimensional generalization is impossible. Erickson [24] proves lower bounds on the related *Hopcroft's problem* of deciding for a set of points and a set of hyperplanes whether or not there is a point that lies on one of the hyperplanes. This is not considered as a data structure problem, and the bounds are on the computation time as a function of the number n of points and the number m of hyperplanes. In a subsequent paper, Erickson [25] proves strong time space lower bounds in a structured model (called *partition graphs*) for an online version of Hopcroft's problem where the data base consists of points and the queries are hyperplanes.

Last, consider a geometric interpretation of exact partial match. The database points are vectors in C_d (viewed either as vectors in Z_2^d , or as vectors in \mathbb{R}^d). The query is an affine subspace of C_d (in either view), defined by the linear equations $x_i = q_i$ for all i such that $q_i \neq *$. For definiteness, assume that this subspace is given by an orthogonal basis plus a shift — this representation can be computed easily from a partial match query. Thus we have lower bounds for the problem of determining whether or not a query affine subspace contains at least one database point. (Miltersen et al. give rather strong lower bounds for the SPAN problem of determining whether or not a given linear subspace contains a given point.)

5 Limitations of the Method

In this section we explain why a direct application¹⁴ of the richness technique does not appear to provide stronger lower bounds.

First, we consider why we cannot apply the richness technique directly to the the exact NNS problem, or, more precisely, the λ -neighbor decision problem.¹⁵ (As noted before, the hardest case seems to be to distinguish between a distance of at most $\frac{d}{2}$ and a distance of at least $\frac{d}{2} + 1$.) We again let $B_d(\lambda)$ denote the Hamming ball of radius λ (say around the all-zeros vector $0^d \in C_d$).

Claim 5.1. *For every $\lambda \in \{0, 1, \dots, d\}$, the communication matrix for the λ -neighbor problem contains a 1-monochromatic rectangle of size $|B_d(\lambda)| \times 2^{n-d}$.*

Before we prove this claim, we point out its consequence. In the λ -neighbor problem, each database is close enough to at most $n |B_d(\lambda)|$ queries. Thus, the problem cannot be richer than $[n |B_d(\lambda)|, 2^{nd}]$. Therefore, the best lower bound we can hope to prove this way is the very weak conclusion: Either the query side sends $\Omega(\log n)$ bits or the database side sends $\Omega(d)$ bits.

Proof of Claim 5.1. Take all the queries in $B_d(\lambda)$. If 0^d is contained in the database then it produces a value of 1 with all these queries. If we pick a database at random, all databases equally likely, the database contains 0^d with probability more than 2^{-d} . ■

We now consider the natural idea of restricting exact partial match to instances with fewer don't cares, in an attempt to prove better lower bounds. The hardest case seems to be when queries have

¹³By duality, we can interchange the role played by points (which become the data base) and hyperplanes (which become the queries)

¹⁴By "direct application" we mean that we do not restrict the set of queries.

¹⁵By "directly" we mean that we do not restrict the set of queries.

exactly $\frac{d}{2}$ don't cares. In this case NPM is extremely rich. Almost all entries in the communication matrix are one. However, perhaps not surprisingly, we have the following:¹⁶

Claim 5.2. *The communication matrix for NPM restricted to queries with exactly $\frac{d}{2}$ don't cares contains a 1-monochromatic rectangle of size $n^{-2} \binom{d}{d/2} 2^{d/2} \times e^{-1} 2^{nd}$.*

The consequence here is obvious. The total number of possible queries is $\binom{d}{d/2} 2^{d/2}$. Thus, the best lower bound we can prove by the richness technique is the rather pathetic “either the query sends $\Omega(\log n)$ bits or the database sends $\Omega(1)$ bits”.

Proof of Claim 5.2. Take the set of queries to be all possible queries with $\frac{d}{2}$ don't cares, and the first k bits fixed as zeros (k to be determined shortly). The number of such queries is

$$\binom{d-k}{d/2} 2^{d/2-k} \geq \binom{d}{d/2} 2^{d/2-2k}.$$

The number of cube points matched by at least one query is exactly 2^{d-k} . Therefore, the number of databases that are not matched by any query is

$$(2^d - 2^{d-k})^n = (1 - 2^{-k})^n 2^{nd} \geq e^{-n/2^k} 2^{nd}.$$

Now, take $k = \log n$. ■

Returning to the case of $\log n + 1$ exposed bits, is it possible to improve upon the proven bounds? If all possible queries are enumerated in some predefined order, the database can store the answer to all possible queries and the query player can then simply send the index of the query using $O(\log n \log d)$ bits (and the database player responds with the correct answer using one bit). Hence the bound on the query player is optimal. Finally, we ask if we can improve our lower bound on the database side to $\Omega(n)$? The following claim shows that our analysis cannot be improved significantly.

Claim 5.3. *For every integer c , there is $\nu > 0$ such that for every $n \geq \nu$, the communication matrix of NPM restricted to queries from $Q_{n,d}$ contains a 1-monochromatic rectangle of size $2^{-c(\log d+1)} |Q_{n,d}| \times 2^{nd-n \log e/2^c}$.*

Proof. We may assume that n is sufficiently large so that $\log n \gg c$. Take all queries in $Q_{n,d}$ with the first c bits fixed as zeros. The number of such queries is

$$\frac{n}{2^{c-1}} \binom{d-c}{\log n + 1 - c} \geq \frac{2n}{2^{c(\log d+1)}} \binom{d}{\log n + 1}.$$

The number of cube points matched by at least one of these queries is 2^{d-c} . Therefore, the number of databases not matched by any of these queries is

$$(1 - 2^{-c})^n 2^{nd} \geq e^{-n/2^c} 2^{nd}. \quad \blacksquare$$

References

- [1] P.K. Agarwal and J. Matoušek. Ray shooting and parametric search. *SICOMP*, 22:794-806, 1993
- [2] M. Ajtai. A lower bound for finding predecessors in Yao's cell probe model. *Combinatorica*, 8:235-247, 1988.
- [3] M. Ajtai. Determinism versus non-determinism for linear-time RAMs. In *Proc. of 31st STOC*, pp. 632-641, 1999.

¹⁶Using similar arguments, one can show that considering the complement function does not help in this case.

- [4] M. Ajtai. A non-linear time lower bound for Boolean branching programs. In *Proc. of 40th FOCS*, pp. 60-70, 1999.
- [5] N. Alon and J.H. Spencer. *The Probabilistic Method*. Wiley, 1992.
- [6] S. Arya and D. Mount. Approximate nearest neighbor queries in fixed dimensions. In *Proc. of 4th SODA*, pp. 271–280, 1993.
- [7] S. Arya, D. Mount, N. Netanyahu, R. Silverman, and A. Wu. An optimal algorithm for approximate nearest neighbor searching in fixed dimensions. *JACM*, 45(6):891-923, 1998.
- [8] O. Barkol and Y. Rabani. Tighter bounds for nearest neighbor search and related problems in the cell probe model. In *Proc. of the 32nd STOC*, pp. 388-396, 2000.
- [9] P. Beame, M. Saks, and J.S. Thathachar. Time-space tradeoffs for branching programs. In *Proc. of 39th FOCS*, pp. 254–263, 1998.
- [10] J.S. Beis and D.G. Lowe. Shape indexing using approximate nearest-neighbor search in high-dimensional spaces. In *Proc. IEEE Conf. Comp. Vision Patt. Recog.*, pages 1000–1006, 1997.
- [11] J.L. Bentley and R. Sedgwick. Fast algorithms for sorting and searching strings. In *Proc. of 8th SODA*, pp. 360–369, 1997.
- [12] M. Ben-Or. Lower bounds for algebraic computation trees. In *Proc. of 15th STOC*, pp. 80-86, 1983.
- [13] M. de Berg, M. van Kreveld, M. Overmars, O. Schwarzkopf. *Computational Geometry, Algorithms and Applications*. Springer, 1997.
- [14] M. Bern. Approximate closest-point queries in high dimensions. *Information Processing Lett.*, 45:95–99, 1993.
- [15] A. Chakrabarti, B. Chazelle, B. Gum, A. Lvov. A good neighbor is hard to find. In *Proc. of 31st STOC*, pp 305-311, 1999.
- [16] T.M. Chan. Approximate nearest neighbor queries revisited. *Discrete and Computational Geometry*, 20:359-373, 1998.
- [17] B. Chazelle. Private communication.
- [18] K. Clarkson. A randomized algorithm for closest-point queries. *SIAM J. Computing*, 17:830-847, 1988.
- [19] K. Clarkson. An algorithm for approximate closest-point queries. In *Proc. of 10th SCG*, pp. 160–164, 1994.
- [20] S. Deerwester, S.T. Dumais, G.W. Furnas, T.K. Landauer, and R. Harshman. Indexing by latent semantic analysis. *J. Amer. Soc. Info. Sci.*, 41(6):391–407, 1990.
- [21] D. Dobkin and R. Lipton. Multidimensional search problems. *SIAM J. Computing*, 5:181-186, 1976.
- [22] D. Dolev, Y. Harari, and M. Parnas. Finding the neighborhood of a query in a dictionary. In *Proc. of 2nd ISTCS*, 1993.
- [23] D. Dolev, Y. Harari, N. Linial, N. Nisan, and M. Parnas. Neighborhood preserving hashing and approximate queries. In *Proc. of 5th SODA*, pp. 251–259, 1994.
- [24] J. Erickson. New lower bounds for Hopcroft’s problem. *Discrete Comput. Geom.*, 16:389–418, 1996.

- [25] J. Erickson. Space-time tradeoffs for emptiness queries. *SIAM J. Computing*, 29(6):1968-1996, 2000.
- [26] R. Fagin. Fuzzy queries in multimedia database systems. In *Proc. of PODS*, pp 1-10, 1998.
- [27] M. Flickner, H. Sawhney, W. Niblack, J. Ashley, Q. Huang, B. Dom, M. Gorkani, J. Hafner, D. Lee, D. Petkovic, D. Steele, and P. Yankner. Query by image and video content: the QBIC system. *IEEE Computer*, 28:23–32, 1995.
- [28] M.L. Fredman. A lower bound on the complexity of orthogonal range queries. *Journal of the ACM*, 28:696-705, 1981
- [29] M.L. Fredman. Lower bounds on the complexity of some optimal data structures. *SIAM J. Computing*, 10:1-10, 1981
- [30] M.L. Fredman and D.J. Volper. The complexity of partial match retrieval in a dynamic setting. *Journal of Algorithms*, 3:68-78, 1982.
- [31] D. Grigoriev. Randomized complexity lower bounds for arrangements and polyhedra. *Discrete and Computational Geometry*, 21:329-344, 1999.
- [32] D. Grigoriev and M. Karpinski. Randomized $\Omega(n^2)$ lower bound for knapsack. *Proc. of 29th STOC*, pp. 76-85, 1997.
- [33] L. Harper. Optimal numberings and isoperimetric problems on graphs. *Journal of Combinatorial Theory* 1:385–394, 1966.
- [34] T. Hastie and R. Tibshirani. Discriminant adaptive nearest neighbor classification and regression. *Advances in Neural Information Processing Systems*, 8:409-415, 1996.
- [35] J. Hellerstein, E. Koutsoupias, and C.H. Papadimitriou. On the analysis of indexing schemes. In *Proc. of PODS*, 1997.
- [36] P. Indyk. On approximate nearest neighbors in non-Euclidean spaces. In *Proc. of 39th FOCS*, 1998.
- [37] P. Indyk and R. Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In *Proc. of 30th STOC*, 1998.
- [38] J. Kleinberg. Two algorithms for nearest-neighbor search in high dimensions. In *Proc. of 29th STOC*, pp. 599–608, 1997.
- [39] E. Kushilevitz and N. Nisan. *Communication Complexity*. Cambridge University Press, 1997.
- [40] E. Kushilevitz, R. Ostrovsky, and Y. Rabani. Efficient search for approximate nearest neighbor in high dimensional spaces. *SICOMP*, 30(2): 457-474, 2000.
- [41] J. Matoušek. Reporting points in halfspaces. *Computational Geometry and Applications*, 2:169-186, 1992.
- [42] S. Meiser. Point location in arrangements of hyperplanes. *Information and Computation*, 106(2):286–303, 1993.
- [43] P.B. Miltersen. The bit probe complexity measure revisited. In *Proc. of 10th STACS*, pp. 662–671, 1993.
- [44] P.B. Miltersen. Lower bounds for union-split-find related problems on random access machines. In *Proc. of 26th STOC*, pp. 625–634, 1994.
- [45] P.B. Miltersen. On the cell probe complexity of polynomial evaluation. *Theoretical Computer Science*, 143:167–174, 1995.

- [46] P.B. Miltersen, N. Nisan, S. Safra, and A. Wigderson. On data structures and asymmetric communication complexity. *JCSS*, 57(5): 37-49, 1998.
- [47] A. Pentland, R.W. Picard, and S. Sclaroff. Photobook: tools for content-based manipulation of image databases. In *Proc. SPIE Conf. on Storage and Retrieval of Image and Video Databases II*, 1994.
- [48] R. Rivest. Partial-match retrieval algorithms. *SIAM J. Computing*, 5:19–50, 1976.
- [49] G. Salton. *Automatic Text Processing*. Addison-Wesley, 1989.
- [50] A.W.M. Smeulders and R. Jain (eds). *Proc. 1st Workshop on Image Databases and Multi-Media Search*, 1996.
- [51] B. Xiao. New Bounds in Cell Probe Model. Ph.D. thesis, UC San Diego, 1992.
- [52] A.C. Yao. Should tables be sorted? *J. Assoc. Comput. Mach.*, 28:615–628, 1981.
- [53] A.C. Yao and F.F. Yao. A general approach to d -dimension geometric queries. In *Proc. of 17th STOC*, pp. 163–168, 1985.

Acknowledgments

Allan Borodin is at the Department of Computer Science, University of Toronto, Toronto M5S 3G4, Canada. borodin@cs.toronto.edu. Rafail Ostrovsky is at the Math Sciences Research Center, Telcordia Technologies, 445 South Street, Morristown, NJ 07960-6438, USA. rafail@research.telcordia.com. Yuval Rabani is at the Computer Science Department, Technion — Israel Institute of Technology, Haifa 32000, Israel. rabani@cs.technion.ac.il.

Part of this work was done while the first and third authors were visiting Telcordia Technologies and while the first and second authors were visiting the Technion. Work at the Technion supported by BSF grant 96-00402, by Ministry of Science contract number 9480198, and by a grant from the Fund for the Promotion of Research at the Technion.

The authors would like to thank Bernard Chazelle for many helpful comments. We are also indebted to the referees for their very constructive suggestions. In particular, we thank one of the referees for the clarification of Meiser's [42] space bound and for the reference to Erickson [25].