

Approximate k -flat Nearest Neighbor Search*

Wolfgang Mulzer[†] Huy L. Nguyễn[‡] Paul Seiferth[†] Yannik Stein[†]

Abstract

Let k be a nonnegative integer. In the *approximate k -flat nearest neighbor* (k -ANN) problem, we are given a set $P \subset \mathbb{R}^d$ of n points in d -dimensional space and a fixed approximation factor $c > 1$. Our goal is to preprocess P so that we can efficiently answer *approximate k -flat nearest neighbor queries*: given a k -flat F , find a point in P whose distance to F is within a factor c of the distance between F and the closest point in P . The case $k = 0$ corresponds to the well-studied approximate nearest neighbor problem, for which a plethora of results are known, both in low and high dimensions. The case $k = 1$ is called *approximate line nearest neighbor*. In this case, we are aware of only one provably efficient data structure, due to Andoni, Indyk, Krauthgamer, and Nguyễn (AIKN) [2]. For $k \geq 2$, we know of no previous results.

We present the first efficient data structure that can handle approximate nearest neighbor queries for arbitrary k . We use a data structure for 0-ANN-queries as a black box, and the performance depends on the parameters of the 0-ANN solution: suppose we have an 0-ANN structure with query time $O(n^\rho)$ and space requirement $O(n^{1+\sigma})$, for $\rho, \sigma > 0$. Then we can answer k -ANN queries in time $O(n^{k/(k+1-\rho)+t})$ and space $O(n^{1+\sigma k/(k+1-\rho)} + n \log^{O(1/t)} n)$. Here, $t > 0$ is an arbitrary constant and the O -notation hides exponential factors in k , $1/t$, and c and polynomials in d .

Our approach generalizes the techniques of AIKN for 1-ANN: we partition P into *clusters* of increasing radius, and we build a low-dimensional data structure for a random projection of P . Given a query flat F , the query can be answered directly in clusters whose radius is “small” compared to $d(F, P)$ using a grid. For the remaining points, the low dimensional approximation turns out to be precise enough. Our new data structures also give an improvement in the space requirement over the previous result for 1-ANN: we can achieve near-linear space and sublinear query time, a further step towards practical applications where space constitutes the bottleneck.

*WM and PS were supported in part by DFG Grants MU 3501/1 and MU 3501/2. YS was supported by the Deutsche Forschungsgemeinschaft within the research training group ‘Methods for Discrete Structures’ (GRK 1408).

[†]Institut für Informatik, Freie Universität Berlin, {mulzer,pseiferth,yannikstein}@inf.fu-berlin.de.

[‡]Simons Institute, UC Berkeley hlnguyen@cs.princeton.edu.

1 Introduction

Nearest neighbor search is a fundamental problem in computational geometry, with countless applications in databases, information retrieval, computer vision, machine learning, signal processing, etc. [10]. Given a set $P \subset \mathbb{R}^d$ of n points in d -dimensional space, we would like to preprocess P so that for any query point $q \in \mathbb{R}^d$, we can quickly find the point in P that is closest to q .

There are efficient algorithms if the dimension d is “small” [7, 18]. However, as d increases, these algorithms quickly become inefficient: either the query time approaches linear or the space grows exponentially with d . This phenomenon is usually called the “curse of dimensionality”. Nonetheless, if one is satisfied with just an *approximate* nearest neighbor whose distance to the query point q lies within some factor $c = 1 + \varepsilon$, $\varepsilon > 0$, of the distance between q and the actual nearest neighbor, there are efficient solutions even for high dimensions. Several methods are known, offering trade-offs between the approximation factor, the space requirement, and the query time (see, e.g., [1, 3] and the references therein).

From a practical perspective, it is important to keep both the query time and the space small. Ideally, we would like algorithms with almost linear (or at least sub-quadratic) space requirement and sub-linear query time. Fortunately, there are solutions with these guarantees. These methods include *locality sensitive hashing* (LSH) [11, 12] and a more recent approach that improves upon LSH [3]. Specifically, the latter algorithm achieves query time $n^{7/(8c^2)+O(1/c^3)}$ and space $n^{1+7/(8c^2)+O(1/c^3)}$, where c is the approximation factor.

Often, however, the query object is more complex than a single point. Here, the complexity of the problem is much less understood. Perhaps the simplest such scenario occurs when the query object is a k -dimensional flat, for some small constant k . This is called the *approximate k -flat nearest neighbor* problem [2]. It constitutes a natural generalization of approximate nearest neighbors, which corresponds to $k = 0$. In practice, low-dimensional flats are used to model data subject to linear variations. For example, one could capture the appearance of a physical object under different lighting conditions or under different viewpoints (see [4] and the references therein).

So far, the only known algorithm with worst-case guarantees is for $k = 1$, the *approximate line nearest neighbor* problem. For this case, Andoni, Indyk, Krauthgamer, and Nguyễn (AIKN) achieve sub-linear query time $d^{O(1)}n^{1/2+t}$ and space $d^{O(1)}n^{O(1/\varepsilon^2+1/t^2)}$, for arbitrarily small $t > 0$. For the “dual” version of the problem, where the query is a point but the data set consists of k -flats, three results are known [4, 14, 15]. The first algorithm is essentially a heuristic with some control of the quality of approximation [4]. The second algorithm provides provable guarantees and a very fast query time of $(d + \log n + 1/\varepsilon)^{O(1)}$ [14]. The third result, due to Mahabadi, is very recent and improves the space requirement of Magen’s result [15]. Unfortunately, these algorithms all suffer from very high space requirements, thus limiting their applicability in practice. In fact, even the basic LSH approach for $k = 0$ is already too expensive for large datasets and additional theoretical work and heuristics are required to reduce the memory usage and make LSH suitable for this setting [13, 19]. For $k \geq 2$, we know of no results in the theory literature.

Our results. We present the first efficient data structure for general approximate k -flat nearest neighbor search. Suppose we have a data structure for approximate *point* nearest neighbor search with query time $O(n^\rho + d \log n)$ and space $O(n^{1+\sigma} + d \log n)$, for some constants $\rho, \sigma > 0$. Then our algorithm achieves query time $O(d^{O(1)}n^{k/(k+1-\rho)+t})$ and space $O(d^{O(1)}n^{1+\sigma k/(k+1-\rho)} + n \log^{O(1/t)} n)$, where $t > 0$ can be made arbitrarily small. The constant factors for the query time depend on k , c , and $1/t$. Our main result is as follows.

Theorem 1.1. *Fix an integer $k \geq 1$ and an approximation factor $c > 1$. Suppose we have a data structure for approximate point nearest neighbor search with query time $O(n^\rho + d \log n)$ and space $O(n^{1+\sigma} + d \log n)$, for some constants $\rho, \sigma > 0$. Let $P \subset \mathbb{R}^d$ be a d -dimensional n -point set. For*

any parameter $t > 0$, we can construct a data structure with $O(d^{O(1)}n^{1+k\sigma/(k+1-\rho)} + n \log^{O(1/t)} n)$ space that can answer the following queries in expected time $O(d^{O(1)}n^{k/(k+1-\rho)+t})$: given a k -flat $F \subset \mathbb{R}^d$, find a point $p \in P$ with $d(p, F) \leq cd(P, F)$.

Algorithm	ρ	σ
AINR [3]	$7/8c^2 + O(1/c^3)$	$7/8c^2 + O(1/c^3)$
LSH1 [1, Theorem 3.2.1]	$1/c^2$	$1/c^2$
LSH2 [1, Theorem 3.4.1]	$O(1/c^2)$	0

The table above gives an overview of some approximate point nearest neighbor structures that can be used in Theorem 1.1. The result by AINR gives the current best query performance for large enough values of c . For smaller c , an approach using locality sensitive hashing (LSH1) may be preferable. With another variant of locality sensitive hashing (LSH2), the space can be made almost linear, at the expense of a slightly higher query time. The last result (and related practical results, e.g., [13]) is of particular interest in applications as the memory consumption is a major bottleneck in practice. It also improves over the previous algorithm by AIKN for line queries.

Along the way towards Theorem 1.1, we present a novel data structure for k -flat *near* neighbor reporting when the dimension d is constant. The space requirement in this case is $O_d(n \log^{O(d)} n)$ and the query time is $O_d(n^{k/(k+1)} \log^{d-k-1} n + |R|)$, where R is the answer set. We believe that this data structure may be of independent interest and may lead to further applications. Our results provide a vast generalization of the result in AIKN and shows for the first time that it is possible to achieve provably efficient nearest neighbor search for higher-dimensional query objects.

Our techniques. Our general strategy is similar to the approach by AIKN. The data structure consists of two main structures: the *projection structure* and the *clusters*. The projection structure works by projecting the point set to a space of constant dimension and by answering the nearest neighbor query in that space. As we will see, this suffices to obtain a rough estimate for the distance, and it can be used to obtain an exact answer if the point set is “spread out”.

Unfortunately, this does not need to be the case. Therefore, we partition the point set into a sequence of *clusters*. A cluster consists of m points and a k -flat K such that all points in the cluster are “close” to K , where m is a parameter to be optimized. Using a rough estimate from the projection structure, we can classify the clusters as *small* and *large*. The points in the large clusters are spread out and can be handled through projection. The points in the small clusters are well behaved and can be handled directly in high dimensions using grids and discretization.

Organization. In order to provide the curious reader with quick gratification, we will give the main data structure together with the properties of the cluster and the projection structure in Section 2. Considering these structures as black boxes, this already proves Theorem 1.1.

In the remainder of the paper, we describe the details of the helper structures. The necessary tools are introduced in Section 3. Section 4 gives the approximate nearest neighbor algorithm for small clusters. In Section 5, we consider approximate near neighbor reporting for k -flats in constant dimension. This data structure is then used for the projection structures in Section 6.

2 Main Data Structure and Algorithm Overview

We describe our main data structure for approximate k -flat nearest neighbor search. It relies on various substructures that will be described in the following sections. Throughout, P denotes a d -dimensional n -point set, and $c > 1$ is the desired approximation factor.

Let K be a k -flat in d dimensions. The *flat-cluster* C (or cluster for short) of K with radius α is the set of all points with distance at most α to K , i.e., $C = \{p \in \mathbb{R}^d \mid d(p, K) \leq \alpha\}$. A cluster is *full* if it contains at least m points from P , where m is a parameter to be determined.

We call P α -cluster-free if there is no full cluster with radius α . Let $t > 0$ be an arbitrarily small parameter. Our data structure requires the following three subqueries.

Q1: Given a query flat F , find a point $p \in P$ with $d(p, F) \leq n^t d(P, F)$.

Q2: Assume P is contained in a flat-cluster with radius α . Given a query flat F with $d(P, F) \geq \alpha/n^{2t}$, return a point $p \in P$ with $d(p, F) \leq cd(P, F)$.

Q3: Assume P is $\alpha n^{2t}/(2k+1)$ -cluster free. Given a query flat F with $d(P, F) \leq \alpha$, find the nearest neighbor $p^* \in P$ to F .

Briefly, our strategy is as follows: during the preprocessing phase, we partition the point set into a set of full clusters of increasing radii. To answer a query F , we first perform a query of type **Q1** to obtain an n^t -approximate estimate \tilde{r} for $d(P, F)$. Using \tilde{r} , we identify the “small” clusters. These clusters can be processed using a query of type **Q2**. The remaining point set contains no “small” full cluster, so we can process it with a query of type **Q3**.

We will now describe the properties of the subqueries and the organization of the data structure in more detail. The data structure for **Q2**-queries is called the *cluster structure*. It is described in Section 4, and it has the following properties.

Theorem 2.1. *Let Q be a d -dimensional m -point set that is contained in a flat-cluster of radius α . Let $c > 1$ be an approximation factor. Using space $O_c(m^{1+\sigma} + d \log^2 m)$, we can build a data structure with the following property. Given a query k -flat F with $d(P, F) \geq \alpha/n^{2t}$ and an estimate \tilde{r} with $d(P, F) \in [\tilde{r}/n^t, \tilde{r}]$, we can find a c -approximate nearest neighbor for F in Q in total time $O_c((n^{2t}k^2)^{k+1}(m^{1-1/k+\rho/k} + (d/k) \log m))$.*

The data structures for **Q1** and **Q3** are very similar, and we cover them in Section 6. They are called *projection structures*, since they are based on projecting P into a low dimensional subspace. In the projected space, we use a data structure for approximate k -flat near neighbor search to be described in Section 5. The projection structures have the following properties.

Theorem 2.2. *Let P be a d -dimensional n -point set, and let $t > 0$ be a small enough constant. Using space and time $O(n \log^{O(1/t)} n)$, we can obtain a data structure for the following query: given a k -flat F , find a point $p \in P$ with $d(p, F) \leq n^t d(P, F)$. A query needs $O(n^{k/(k+1)} \log^{O(1/t)} n)$ time, and the answer is correct with high probability.*

Theorem 2.3. *Let P be a d -dimensional n -point set, and let $t > 0$ be a small enough constant. Using space and time $O(n \log^{O(1/t)} n)$, we can obtain a data structure for the following query: given a k -flat F and $\alpha > 0$ such that $d(F, P) \leq \alpha$ and such that P is $\alpha n^t/(2k+1)$ -cluster-free, find an exact nearest neighbor for F in P . A query needs $O(n^{k/(k+1)} \log^{O(1/t)} n + m)$ time, and the answer is correct with high probability. Here, m denotes the size of a full cluster.*

2.1 Constructing the Data Structure

First, we build a projection structure for **Q1** queries on P . This needs $O(n \log^{O(1/t)} n)$ space, by Theorem 2.2. Then, we repeatedly find the full flat-cluster C with smallest radius. The m points in C are removed from P , and we build a cluster structure for **Q2** queries on this set. By Theorem 2.1, this needs $O_c(m^{1+\sigma} + d \log^2 m)$ space. To find C , we check all flats K spanned by $k+1$ distinct points of P . In Lemma 3.2 below, we prove that this provides a good enough approximation. In the end, we have n/m point sets $Q_1, \dots, Q_{n/m}$ ordered by decreasing radius, i.e., the cluster for Q_1 has the largest radius. The total space occupied by the cluster structures is $O(nm^\sigma + (n/m)d \log^2 n)$.

Finally, we build a perfect binary tree T with n/m leaves labeled $Q_1, \dots, Q_{n/m}$, from left to right. For a node $v \in T$ let Q_v be the union of all Q_i assigned to leaves below v . For each $v \in T$ we build a data structure for Q_v to answer **Q3** queries. Since each point is contained in $O(\log n)$ data structures, the total size is $O(n \log^{O(1/t)} n)$, by Theorem 2.3. For pseudocode, see Algorithm 1.

Input: point set $P \subset \mathbb{R}^d$, approximation factor c , parameter $t > 0$

- 1 $Q \leftarrow P$
- 2 **for** $i \leftarrow n/m$ **downto** 1 **do**
- 3 For each $V \in \binom{Q}{k+1}$, consider the k -flat K_V defined by V . Let α_V be the radius of the smallest flat-cluster of K_V with exactly m points of Q .
- 4 Choose the flat $K = K_V$ that minimizes α_V and set $\alpha_i = \alpha_V$.
- 5 Remove from Q the set Q_i of m points in Q within distance α_i from K .
- 6 Construct a cluster structure C_i for the cluster (K, Q_i) .
- 7 Build a perfect binary tree T with n/m leaves, labeled $Q_1, \dots, Q_{n/m}$ from left to right.
- 8 **foreach** node $v \in T$ **do**
- 9 Build data structure for **Q3** queries as in Theorem 2.3 for the set Q_v corresponding to the leaves below v .

Algorithm 1: Preprocessing algorithm. Compared with AIKN [2], we organize the projection structure in a tree to save space.

2.2 Performing a Query

Suppose we are given a k -flat F . To find an approximate nearest neighbor for F we proceed similarly as AIKN [2]. We use **Q2** queries on “small” clusters and **Q3** queries on the remaining points; for pseudocode, see Algorithm 2.

Input : query flat F

Output: a c -approximate nearest neighbor for F in P

- 1 Query the root of T for a n^t -approximate nearest neighbor p_1 to F . /* type Q1 */
- 2 $\tilde{r} \leftarrow d(p_1, F)$
- 3 $i^* \leftarrow$ maximum $i \in \{1, \dots, n/m\}$ with $\alpha_i > \tilde{r}n^t$, or 0 if no such i exists
- 4 **for** $i \leftarrow i^* + 1$ **to** n/m **do**
- 5 /* type Q2; we have $d(Q_i, F) \geq \tilde{r}/n^t \geq \alpha_i/n^{2t}$ */
- 5 Query cluster structure C_i with estimate \tilde{r} . */
- /* type Q3 */
- 6 Query projection structure for a \tilde{r} -thresholded nearest neighbor of F in $Q = \bigcup_{i=1}^{i^*} U_i$.
- return** closest point to F among query results.

Algorithm 2: Algorithm for finding approximate nearest neighbor in high dimensions.

First, we perform a query of type **Q1** to get a n^t -approximate nearest neighbor p_1 for F in time $O(n^{k/(k+1)} \log^{O(1/t)} n)$. Let $\tilde{r} = d(p_1, F)$. We use \tilde{r} as an estimate to distinguish between “small” and “large” clusters. Let $i^* \in \{1, \dots, n/m\}$ be the largest integer such that the cluster assigned with Q_{i^*} has radius $\alpha_{i^*} > \tilde{r}n^t$. For $i = i^* + 1, \dots, n/m$, we use \tilde{r} as an estimate for a **Q2** query on Q_i . Since $|Q_i| = m$ and by Theorem 2.1, this needs total time $O(n^{2t(k+1)+1} m^{-1/k+\rho/k} + (n/m)d \log^2 m)$.

It remains to deal with points in “large” clusters. The goal is to perform a type **Q3** query on $\bigcup_{1 \leq i \leq i^*} Q_i$. For this, we start at the leaf of T labeled Q_{i^*} and walk up to the root. Each time we encounter a new node v from its right child, we perform a **Q3** query on Q_u , where u denotes the left child of v . Let L be all the left children we find in this way. Then clearly we have $|L| = O(\log n)$ and $\bigcup_{u \in L} Q_u = \bigcup_{1 \leq i \leq i^*} Q_i$. Moreover, by construction, there is no full cluster with radius less than $\tilde{r}n^t$ defined by $k+1$ vertices of Q_u for any $u \in L$. We will see that this implies every Q_u to be $\tilde{r}n^t/(2k+1)$ -cluster-free, so Theorem 2.3 guarantees a total query time of $O(n^{k/(k+1)} \log^{O(1/t)} n + m)$ for this step. Among all the points we obtained during the queries, we return the one that is closest to F . A good trade-off point is achieved for $m = nm^{-1/k+\rho/k}$, i.e., for $m = n^{k/(k+1-\rho)}$. This gives the bounds claimed in Theorem 1.1.

Correctness. Let p^* be a point with $d(p^*, F) = d(P, F)$. First, suppose that $p^* \in Q_i$, for some

$i > i^*$. Then, we have $d(p^*, F) \geq \tilde{r}/n^t \geq \alpha_i/n^{2t}$, where α_i is the radius of the cluster assigned to Q_i . Since \tilde{r} is a valid n^t -approximate estimate for $d(F, Q_i)$, a query of type **Q2** on Q_i gives a c -approximate nearest neighbor, by Theorem 2.1. Now, suppose that $p^* \in Q_i$ for $1 \leq i \leq i^*$. Let u be the node of L with $p^* \in Q_u$. Then Theorem 2.3 guarantees that we will find p^* when doing a **Q3** query on Q_u .

3 Preliminaries

Partition Trees. Our algorithm relies on *partition trees* [6, 16]. We use the optimal version due to Chan [5], as summarized in the following theorem.

Theorem 3.1 (Optimal Partition Trees [5]). *For any d -dimensional n -point set $P \subset \mathbb{R}^d$, and for any large enough constant r , there is a partition tree \mathcal{T} with the following properties: (i) the tree \mathcal{T} has degree $O(r)$ and depth $\log_r n$; (ii) each node is of the form (Q, Δ) , where Q is a subset of P and Δ a relatively open simplex that contains Q ; (iii) for each node (Q, Δ) , the simplices of the children of Q are contained in Δ and are pairwise disjoint; (iv) the point set associated with a node of depth ℓ has size at most n/r^ℓ ; (v) for any hyperplane h in \mathbb{R}^d , the number m_ℓ of simplices in \mathcal{T} that h intersects at level ℓ obeys the recurrence*

$$m_\ell = O(r^{\ell(d-1)/d} + r^{\ell(d-2)/(d-1)}m_{\ell-1} + r\ell \log r \log n).$$

Thus, h intersects $O(n^{1-1/d})$ simplices in total. The tree \mathcal{T} can be build in expected time $O(n \log n)$.

k -flat Discretization. For our cluster structure we must find k -flats that are close to many points. The following lemma shows that it suffices to check “few” k -flats for this.

Lemma 3.2. *Let $P \subset \mathbb{R}^d$ be a finite point set with $|P| \geq k+1$, and let $F \subset \mathbb{R}^d$ be a k -flat. There is a k -flat F' such that F' is the affine hull of $k+1$ points in P and $\delta_{F'}(P) \leq (2k+1)\delta_F(P)$, where $\delta_{F'}(P) = \max_{p \in P} d(p, F')$ and $\delta_F(P) = \max_{p \in P} d(p, F)$.*

Proof. This proof generalizes the proof of Lemma 2.3 by AIKN [2].

Let Q be the orthogonal projection of P onto F . We may assume that F is the affine hull of Q , since otherwise we could replace F by $\text{aff}(Q)$ without affecting $\delta_F(P)$. We choose an orthonormal basis for \mathbb{R}^d such that F is the linear subspace spanned by the first k coordinates. An affine basis for F' is constructed as follows: first, take a point $p_0 \in P$ whose x_1 -coordinate is minimum. Let q_0 be the projection of p_0 onto F , and translate the coordinate system such that q_0 is the origin. Next, choose k additional points $p_1, \dots, p_k \in P$ such that $|\det(q_1, \dots, q_k)|$ is maximum, where q_i is the projection of p_i onto F , for $i = 1, \dots, k$. That is, we choose k additional points such that the volume of the k -dimensional parallelogram spanned by their projections onto F is maximized. The set $\{q_1, \dots, q_k\}$ is a basis for F , since the maximum determinant cannot be 0 by our assumption that F is spanned by Q .

Now fix some point $p \in P$ and let q be its projection onto F . We write $q = \sum_{i=1}^k \mu_i q_i$. Then, the point $r = \sum_{i=1}^k \mu_i p_i + (1 - \sum_{i=1}^k \mu_i) p_0$ lies in F' . By the triangle inequality, we have

$$d(p, r) \leq d(p, q) + d(q, r) \leq \delta_F(P) + d(q, r). \quad (1)$$

To upper-bound $d(q, r)$ we first show that all coefficients μ_i lie in $[-1, 1]$.

Claim 3.3. *Take $p \in P$, $q \in Q$ and q_1, \dots, q_k as above. Write $q = \sum_{i=1}^k \mu_i q_i$. Then for $i = 1, \dots, k$, we have $\mu_i \in [-1, 1]$, and $\mu_j \geq 0$ for at least one $j \in \{1, \dots, k\}$.*

Proof. We first prove that all coefficients μ_i lie in the interval $[-1, 1]$. Suppose that $|\mu_i| > 1$ for some $i \in \{1, \dots, k\}$. We may assume that $i = 1$. Using the linearity of the determinant,

$$|\det(q, q_2, \dots, q_k)| = |\det(\mu_1 q_1, q_2, \dots, q_k)| = |\mu_1| \cdot |\det(q_1, q_2, \dots, q_k)| > |\det(q_1, q_2, \dots, q_k)|,$$

contradicting the choice of q_1, \dots, q_k .

Furthermore, by our choice of the origin, all points in Q have a non-negative x_1 -coordinate. Thus, at least one coefficient μ_j , $j \in \{1, \dots, k\}$, has to be non-negative. \square

Using Claim 3.3, we can now bound $d(q, r)$. For $i = 1, \dots, k$, we write $p_i = q_i + q_i^\perp$, where q_i^\perp is orthogonal to F . Then,

$$\begin{aligned} d(q, r) &= \left\| \sum_{i=1}^k \mu_i q_i - \sum_{i=1}^k \mu_i (q_i + q_i^\perp) - \left(1 - \sum_{i=1}^k \mu_i\right) p_0 \right\| \\ &= \left\| \sum_{i=1}^k \mu_i q_i^\perp + \left(1 - \sum_{i=1}^k \mu_i\right) p_0 \right\| \leq \left(\sum_{i=1}^k |\mu_i| + \left|1 - \sum_{i=1}^k \mu_i\right| \right) \delta_F(P) \leq 2k\delta_F(P), \quad (2) \end{aligned}$$

since $\|q_1^\perp\|, \dots, \|q_k^\perp\|, \|p_0\| \leq \delta_F(P)$, and since $\left|1 - \sum_{i=1}^k \mu_i\right| \leq k$ follows from fact that at least one μ_i is non-negative. By (1) and (2), we get $d(p, F') \leq (2k + 1)\delta_F(P)$. \square

Remark 3.4. For $k = 1$, the proof of Lemma 3.2 coincides with the proof of Lemma 2.3 by AIKN [2]. In this case, one can obtain a better bound on $d(q, r)$ since q is a convex combination of q_0 and q_1 . This gives $\delta_{F'}(P) \leq 2\delta_F(P)$.

4 Cluster Structure

A k -flat cluster structure consists of a k -flat K and a set Q of m points with $d(q, K) \leq \alpha$, for all $q \in Q$. Let $K : u \mapsto A'u + a$ be a parametrization of K , with $A' \in \mathbb{R}^{d \times k}$ and $a \in \mathbb{R}^d$ such that the columns of A' constitute an orthonormal basis for K and such that a is orthogonal to K . We are also given an approximation parameter $c > 1$. The cluster structure uses a data structure for approximate point nearest neighbor search as a black box. We assume that we have such a structure available that can answer c -approximate point nearest neighbor queries in d dimensions with query time $O_c(n^\rho + d \log n)$ and space requirement $O_c(n^{1+\sigma} + d \log n)$ for some constants $\rho, \sigma > 0$. As mentioned in the introduction, the literature offers several data structures for us to choose from.

The cluster structure distinguishes two cases: if the query flat F is close to K , we can approximate F by few “patches” that are parallel to K , such that a good nearest neighbor for the patches is also good for K . Since the patches are parallel to K , they can be handled through 0-ANN queries in the orthogonal space K^\perp and low-dimensional queries inside K . If the query flat is far from K , we can approximate Q by its projection onto K and handle the query with a low-dimensional data structure.

4.1 Preprocessing

Let K^\perp be the linear subspace of \mathbb{R}^d that is orthogonal to K . Let Q_a be the projection of Q onto K , and let Q_b be the projection of Q onto K^\perp . We compute a k -dimensional partition tree \mathcal{T} for Q_a . As stated in Theorem 3.1, the tree \mathcal{T} has $O(m)$ nodes, and it can be computed in time $O(m \log m)$.

For each node (S_a, Δ) of \mathcal{T} , we do the following: we determine the set $S \subseteq Q$ whose projection onto K gives S_a , and we take the projection S_b of S onto K^\perp . Then, we build a $d - k$ dimensional c' -ANN data structure for S_b , as given by the assumption, where $c' = (1 - 1/\log n)c$. See Algorithm 3 for pseudocode.

Lemma 4.1. *The cluster structure can be constructed in total time $O_c(m^{2+\rho} + md \log^2 m)$, and it requires $O_c(m^{1+\sigma} + d \log^2 m)$ space.*

Input: k -flat $K \subseteq \mathbb{R}^d$, point set $Q \subset \mathbb{R}^d$ with $d(q, K) \leq \alpha$ for all $q \in Q$, approximation parameter $c > 1$

- 1 $Q_a \leftarrow$ projection of Q onto K
- 2 $Q_b \leftarrow$ projection of Q onto K^\perp
- 3 Build a k -dimensional partition tree \mathcal{T} for Q_a as in Theorem 3.1.
- 4 $c' \leftarrow (1 - 1/\log n)c$
- 5 **foreach** *node* $(S_a, \Delta) \in \mathcal{T}$ **do**
- 6 $S_b \leftarrow$ projection of the points in Q corresponding to S_a onto K^\perp
- 7 Build a $(d - k)$ -dimensional c' -ANN structure for S_b as given by the assumption.

Algorithm 3: CreateClusterStructure

Proof. By Theorem 3.1, the partition tree can be built in $O(m \log m)$ time. Thus, the preprocessing time is dominated by the time to construct the c' -ANN data structures at the nodes of the partition tree \mathcal{T} . Since the sets on each level of \mathcal{T} constitute a partition of Q , and since the sizes of the sets decrease geometrically, the bounds on the preprocessing time and space requirement follow directly from our assumption. Note that by our choice of $c' = (1 - 1/\log n)c$, the space requirement and query time for the ANN data structure change only by a constant factor. \square

4.2 Processing a Query

We set $\varepsilon = 1/100 \log n$. Let F be the query k -flat, given as $F : v \mapsto B'v + b$, with $B' \in \mathbb{R}^{d \times k}$ and $b \in \mathbb{R}^d$ such that the columns of B' are an orthonormal basis for F and b is orthogonal to F . Our first task is to find bases for the flats K and F that provide us with information about the relative position of K and F . For this, we take the matrix $M = A'^T B' \in \mathbb{R}^{k \times k}$, and we compute a *singular value decomposition* $M = U \Sigma V^T$ of M [9, Chapter 7.3]. Recall that U and V are orthogonal $k \times k$ matrices and that $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_k)$ is a $k \times k$ diagonal matrix with $\sigma_1 \geq \dots \geq \sigma_k \geq 0$. We call $\sigma_1, \dots, \sigma_k$ the *singular values* of M . The following lemma summarizes the properties of the SVD that are relevant to us.

Lemma 4.2. *Let $M = A'^T B'$, and let $M = U \Sigma V^T$ be a singular value decomposition for M . Let u_1, \dots, u_k be the columns of U and v_1, \dots, v_k be the columns of V . Then, (i) u_1, \dots, u_k is an orthonormal basis for K (in the coordinate system induced by A'); (ii) v_1, \dots, v_k is an orthonormal basis for F (in the coordinate system induced by B'); and (iii) for $i = 1, \dots, k$, the projection of v_i onto K is $\sigma_i u_i$ and the projection of u_i onto F is $\sigma_i v_i$ (again in the coordinate systems induced by A' and B'). In particular, we have $\sigma_1 \leq 1$.*

Proof. Properties (i) and (ii) follow since U and V are orthogonal matrices. Property (iii) holds because $M = A'^T B'$ describes the projection from F onto K (in the coordinate systems induced by A' and B') and because $M^T = B'^T A' = V \Sigma U^T$ describes the projection from K onto F . \square

We reparametrize K according to U and F according to V . More precisely, we set $A = A'U$ and $B = B'V$, and we write $K : u \mapsto Au + a$ and $F : v \mapsto Bv + b$. The new coordinate system provides a simple representation for the distances between F and K . We begin with a technical lemma that is a simple corollary of Lemma 4.2.

Lemma 4.3. *Let a_1, \dots, a_k be the columns of the matrix A ; let $a_1^\parallel, \dots, a_k^\parallel$ be the columns of the matrix $BB^T A$, and $a_1^\perp, \dots, a_k^\perp$ the columns of the matrix $A - BB^T A$. Then, (i) for $i = 1, \dots, k$, the vector a_i^\parallel is the projection of a_i onto F and the vector a_i^\perp is the projection of a_i onto F^\perp ; (ii) for $i = 1, \dots, k$, we have $\|a_i^\parallel\| = \sigma_i$ and $\|a_i^\perp\| = \sqrt{1 - \sigma_i^2}$; and (iii) the vectors $a_1^\parallel, \dots, a_k^\parallel, a_1^\perp, \dots, a_k^\perp$ are pairwise orthogonal. An analogous statement holds for the matrices $B, AA^T B$, and $B - AA^T B$.*

Proof. Properties (i) and (ii) are an immediate consequence of the definition of A and B and of Lemma 4.2. The set $a_1^\parallel, \dots, a_k^\parallel$ is orthogonal by Lemma 4.2(ii). Furthermore, since for any $i, j \in \{1, \dots, k\}$, the vector a_i^\parallel lies in F and the vector a_j^\perp lies in F^\perp , a_i^\parallel and a_j^\perp are orthogonal. Finally, let $1 \leq i < j \leq k$. Then,

$$\langle a_i^\perp, a_j^\perp \rangle = \langle a_i^\perp, a_j^\perp \rangle + \langle a_i^\perp, a_j^\parallel \rangle + \langle a_i^\parallel, a_j^\perp \rangle + \langle a_i^\parallel, a_j^\parallel \rangle = \langle a_i, a_j \rangle = 0,$$

since we already saw that $\langle a_i^\perp, a_j^\parallel \rangle = \langle a_i^\parallel, a_j^\perp \rangle = \langle a_i^\parallel, a_j^\parallel \rangle = \langle a_i, a_j \rangle = 0$. The argument for the other matrices is completely analogous. \square

The next lemma shows how our choice of bases gives a convenient representation of the distances between F and K .

Lemma 4.4. *Take two points $x_F \in K$ and $y_K \in F$ such that $d(F, K) = d(y_K, x_F)$. Write $x_F = Au_F + a$ and $y_K = Bv_K + b$. Then, for any point $x \in K$ with $x = Au + a$, we have*

$$d(F, x)^2 = \sum_{i=1}^k (1 - \sigma_i^2) (u - u_F)_i^2 + d(F, K)^2,$$

and for any point $y \in F$ with $y = Bv + b$, we have

$$d(y, K)^2 = \sum_{i=1}^k (1 - \sigma_i^2) (v - v_K)_i^2 + d(F, K)^2.$$

Proof. We show the calculation for $d(F, x)$. The calculation for $d(y, K)$ is symmetric. Let $x \in K$ with $x = Au + a$ be given. Let $y_x \in F$ be the projection of x onto F . Then,

$$d(F, x)^2 = \|x - y_x\|^2 = \|(x - x_F) + (x_F - y_K) + (y_K - y_x)\|^2 = \|(x - x_F) - (y_x - y_K)\|^2 + \|x_F - y_K\|^2,$$

where the last equality is due to Pythagoras, since $x - x_F$ lies in K , $y_x - y_K$ lies in F , and $x_F - y_K$ is orthogonal to both K and F . Now, we have $y_x = BB^T x + b$. Similarly, since y_K is the projection of x_F onto F , we have $y_K = BB^T x_F + b$. Thus,

$$d(F, x)^2 = \|(x - x_F) - BB^T(x - x_F)\|^2 + d(F, K)^2 = \|(A - BB^T A)(u - u_F)\|^2 + d(F, K)^2,$$

using the definition of x and x_F . By Lemma 4.3, the columns $a_1^\perp, \dots, a_k^\perp$ of the matrix $A - BB^T A$ are pairwise orthogonal and for $i = 1, \dots, k$, we have $\|a_i^\perp\|^2 = 1 - \sigma_i^2$. Pythagoras gives

$$d(F, x)^2 = \sum_{i=1}^k \|a_i^\perp\|^2 (u - u_F)_i^2 + d(F, K)^2 = \sum_{i=1}^k (1 - \sigma_i^2) (u - u_F)_i^2 + d(F, K)^2.$$

\square

We now give a brief overview of the query algorithm, refer to Algorithm 4 for pseudocode. First, we check for the special case that F and K are parallel, i.e., that $\sigma_1 = \dots = \sigma_k = 1$. In this case, we need to perform only a single c' -ANN query in Q_b to obtain the desired result. If F and K are not parallel, we distinguish two scenarios: if F is far from Q , we can approximate Q by its projection Q_a onto K . Thus, we take the closest point x_F in K to F , and we return an approximate nearest neighbor for x_F in Q_a according to an appropriate metric derived from Lemma 4.4. Details can be found in Section 4.2.2. If F is close to Q , we use Lemma 4.4 to discretize the relevant part of F into *patches*, such that each patch is parallel to K and such that the best nearest neighbor in Q for the patches provides an approximate nearest neighbor for F . Each patch can then be handled essentially by an appropriate nearest neighbor query in K^\perp . Details follow in Section 4.2.1. We say F and Q are *close* if $d(F, Q) \leq \alpha/\varepsilon$, and *far* if $d(F, Q) > \alpha/\varepsilon$. Recall that we chose $\varepsilon = 1/100 \log n$.

<p>Input: query k-flat $F \subseteq \mathbb{R}^d$; an estimate \tilde{r} with $d(F, Q) \in [\tilde{r}/n^t, \tilde{r}]$.</p> <ol style="list-style-type: none"> 1 $M \leftarrow A'^T B'$. 2 Compute an SVD $M = U\Sigma V^T$ of M with singular values $1 \geq \sigma_1 \geq \dots \geq \sigma_k \geq 0$. 3 if $\sigma_k = 1$ then 4 $f \leftarrow$ projection of F onto K^\perp; /* F and K are parallel; f is a point */ 5 $r \leftarrow c'$-ANN for f in Q_b 6 return r 7 Reparametrize K according to U and F according to V. /* Near case */ 8 $\mathcal{G} \leftarrow$ set of approximate patches obtained by combining Lemma 4.6 and 4.7 9 $R \leftarrow \emptyset$ 10 foreach $G \leftarrow \mathcal{G}$ do 11 $R \leftarrow R \cup$ result of approximate nearest-neighbor query for G as in Lemma 4.8 /* Far case */ 12 $R \leftarrow R \cup$ result of approximate nearest-neighbor for G as in Lemma 4.11 13 return point in R that minimizes the distance to F
--

Algorithm 4: QueryClusterStructure

4.2.1 Near: $d(F, Q) \leq \alpha/\varepsilon$

We use our reparametrization of F and K to split the coordinates as follows: recall that $1 \geq \sigma_1 \geq \dots \geq \sigma_k \geq 0$ are the singular values of $M = A'^T B'$. Pick $l \in \{0, \dots, k\}$ such that $1 \geq \sigma_i \geq \sqrt{1 - \varepsilon}$, for $i = 1, \dots, l$, and $\sqrt{1 - \varepsilon} > \sigma_i \geq 0$, for $i = l + 1, \dots, k$. For a $d \times k$ matrix X , let $X_{[i]}$ denote the $d \times i$ matrix with the first i columns of X , and $X_{-[i]}$ the $d \times (k - i)$ matrix with the remaining $k - i$ columns of X . Similarly, for a vector $v \in \mathbb{R}^k$, let $v_{[i]}$ be the vector in \mathbb{R}^i with the first i coordinates of v , and $v_{-[i]}$ the vector in \mathbb{R}^{k-i} with the remaining $k - i$ coordinates of v .

The following lemma is an immediate consequence of Lemma 4.4. It tells us that we can partition the directions in F into those that are almost parallel to K and those that are almost orthogonal to K . Along the orthogonal directions, we discretize F into few lower-dimensional flats that are almost parallel to K . After that, we approximate these flats by few patches that are actually parallel to K . These patches are then used to perform the query.

Lemma 4.5. *Let $y \in F$ be a point and $y_K \in F$ with $d(F, K) = d(y_K, K)$. Write $y_K = Bv_K + b$ and $y = Bv + b$. Then, $\|(v - v_K)_{-[l]}\| \leq d(y, K)/\sqrt{\varepsilon}$.*

Proof. By Lemma 4.4 and the choice of l ,

$$d(y, K)^2 = \sum_{i=1}^k (1 - \sigma_i^2) (v - v_K)_i^2 + d(F, K)^2 \geq \sum_{i=l+1}^k (1 - \sigma_i^2) (v - v_K)_i^2 \geq \varepsilon \|(v - v_K)_{-[l]}\|^2.$$

□

Using Lemma 4.5, we can discretize the query F into a set of l -flats that are almost parallel to the cluster flat K .

Lemma 4.6. *There is a set \mathcal{L} of $O((n^{2t} k^{1/2} \varepsilon^{-5/2})^{k-l})$ l -flats such that the following holds: (i) for every $L \in \mathcal{L}$, we have $L \subseteq F$; (ii) for every $L \in \mathcal{L}$ and for every unit vector $u \in L$, the projection of u onto K has length at least $\sqrt{1 - \varepsilon}$; and (iii) if $d(F, Q) \in [\alpha/n^{2t}, \alpha/\varepsilon]$, then there is an l -flat $L \in \mathcal{L}$ with $d(L, Q) \leq (1 + \varepsilon)d(F, Q)$.*

Proof. Let $y_K = Bv_K + b \in F$ be a point in F with $d(F, K) = d(y_K, K)$. Furthermore, let

$$\tau = \frac{\alpha\varepsilon}{n^{2t}\sqrt{k}} \quad \text{and} \quad o_\tau = \left\lceil \frac{n^{2t}\sqrt{k}}{\varepsilon^{5/2}} \right\rceil.$$

Using τ and o_τ , we define a set I of *index vectors* with $I = \{-o_\tau\tau, (-o_\tau + 1)\tau, \dots, o_\tau\tau\}^{k-l}$ and $|I| = O(o_\tau^{k-l}) = O((n^{2t}k^{1/2}\varepsilon^{-5/2})^{k-l})$. For each $i \in I$, we define the l -flat L_i as

$$L_i : w \mapsto B_{[l]}w + B_{-[l]}((v_K)_{-[l]} + i) + b.$$

Our desired set of approximate query l -flats is now $\mathcal{L} = \{L_i \mid i \in I\}$.

The set \mathcal{L} meets properties (i) and (ii) by construction, so it remains to verify (iii). For this, we take a point $y_Q \in F$ with $d(F, Q) = d(y_Q, Q)$. We write $y_Q = Bv_Q + b$, and we define $s = (v_Q - v_K)_{-[l]}$. We assumed that $d(y_Q, K) \leq \alpha/\varepsilon$, so Lemma 4.5 gives $\|s\| \leq \alpha/\varepsilon^{3/2}$. It follows that by rounding each coordinate of s to the nearest multiple of τ , we obtain an index vector $i_Q \in I$ with $\|i_Q - s\| \leq \tau\sqrt{k} = \varepsilon\alpha/n^{2t}$. Hence, considering the point in L_{i_Q} with $w = (v_Q)_{[l]}$, we get

$$\begin{aligned} d(L_{i_Q}, Q) &\leq d(L_{i_Q}, y_Q) + d(y_Q, Q) \\ &\leq \|B_{[l]}(v_Q)_{[l]} + B_{-[l]}((v_K)_{-[l]} + i_Q) + b - Bv_Q - b\| + d(F, Q) \\ &= \|B_{-[l]}((v_K)_{-[l]} + i_Q - (v_Q)_{-[l]})\| + d(F, Q) \\ &= \|(v_K - v_Q)_{-[l]} + i_Q\| + d(F, Q) \tag{*} \\ &= \|i_Q - s\| + d(F, Q) \\ &\leq \varepsilon\alpha/n^{2t} + d(F, Q) \\ &\leq (1 + \varepsilon)d(F, Q), \tag{**} \end{aligned}$$

where in (*) we used that the columns of $B_{-[l]}$ are orthonormal and in (**) we used the assumption $d(F, Q) \geq \alpha/n^{2t}$. \square

From now on, we focus on an approximate query l -flat L . Our next goal is to approximate L by a set of patches such that each is parallel to K .

Lemma 4.7. *There is a set \mathcal{G} of $O((n^{2t}k^{1/2}\varepsilon^{-2})^l)$ patches such that the following holds: (i) every $G \in \mathcal{G}$ is an l -dimensional polytope, given by $O(l)$ inequalities; (ii) for every $G \in \mathcal{G}$, the affine hull of G is parallel to K ; (iii) if $d(L, Q) \in [\alpha/n^{2t}, 2\alpha/\varepsilon]$, then there exists $G \in \mathcal{G}$ with $d(G, Q) \leq (1 + \varepsilon)d(L, Q)$; (iv) for all $G \in \mathcal{G}$ and for all $q \in Q$, we have $d(G, q) \geq (1 - \varepsilon)d(L, q)$.*

Proof. Let $C = AA^T B_1$ be the $d \times l$ matrix whose columns $b_1^\parallel, \dots, b_l^\parallel$ constitute the projections of the columns of B onto K . By Lemma 4.3, the vectors b_i^\parallel are orthogonal with $\|b_i^\parallel\| = \sigma_i$, for $i = 1, \dots, l$, and the columns $b_1^\perp, \dots, b_l^\perp$ of the matrix $B_1 - C$ also constitute an orthogonal set, with $\|b_i^\perp\|^2 = 1 - \sigma_i^2$, for $i = 1, \dots, l$. Let z_K be a point in L that minimizes the distance to K , and write $z_K = B_1 w_K + b_1$. Furthermore, let

$$\tau_i = \frac{\alpha\varepsilon}{n^{2t}\sqrt{l(1 - \sigma_i^2)}}, \quad \text{for } i = 1, \dots, l, \quad \text{and} \quad o_\tau = \left\lceil \frac{2n^{2t}\sqrt{l}}{\varepsilon^2} \right\rceil.$$

We use the τ_i and o_τ to define a set I of *index vectors* as $I = \prod_{i=1}^l \{-o_\tau\tau_i, (-o_\tau + 1)\tau_i, \dots, o_\tau\tau_i\}$. We have $|I| = O(o_\tau^l) = O((n^{2t}k^{1/2}\varepsilon^{-2})^l)$. For each index vector $i \in I$, we define the patch G_i as

$$G_i : w \mapsto Cw + B_1(w_K + i) + b_1, \quad \text{subject to } w \in \prod_{i=1}^l [0, \tau_i].$$

Our desired set of approximate query patches is now $\mathcal{G} = \{G_i \mid i \in I\}$. The set \mathcal{G} fulfills properties (i) and (ii) by construction, so it remains to check (iii). Fix a point $z \in L$. Since $L \subseteq F$, we can write $z = B_1 w + b_1 = Bv + b$, where the vector w represents the coordinates of z in L and the vector v represents the coordinates of z in F . By Lemma 4.4,

$$d(z, K)^2 = \sum_{i=1}^k (1 - \sigma_i^2)(v - v_K)_i^2 + d(F, K)^2,$$

where the vector v_K represents the coordinates of a point in F that is closest to K . By definition of L , the last $k - l$ coordinates $v_{-[l]}$ in F are the same for all points $z \in L$, so we can conclude that the coordinates for a closest point to K in L are given by $w_K = (v_K)_{[l]}$ and that

$$d(z, K)^2 = \sum_{i=1}^l (1 - \sigma_i^2)(w - w_K)_i^2 + d(L, K)^2. \quad (3)$$

Now take a point z_Q in L with $d(z_Q, Q) = d(L, Q)$ and write $z_Q = B_1 w_Q + b_1$. Since we assumed $d(L, Q) \leq 2\alpha/\varepsilon$, (3) implies that for $i = 1, \dots, l$, we have $|(w_Q - w_K)_i| \leq 2\alpha/(\varepsilon\sqrt{1 + \sigma_i^2})$. Thus, if for $i = 1, \dots, l$, we round $(w_Q - w_K)_i$ down to the next multiple of τ_i , we obtain an index vector $i_Q \in I$ with $(w_Q - w_K) - i_Q \in \prod_{i=1}^l [0, \tau_i]$. We set $s_Q = (w_Q - w_K) - i_Q$. Considering the point $Cs_Q + B_1(u_K + i_Q) + b_1$ in G_{i_Q} , we see that

$$\begin{aligned} d(G_{i_Q}, z_Q)^2 &\leq \|Cs_Q + B_1(u_K + i_Q) + b_1 - B_1 w_Q - b_1\|^2 = \|Cs_Q - B_1((w_Q - w_K) - i_Q)\|^2 \\ &= \|(C - B_1)s_Q\|^2 = \sum_{i=1}^l (1 - \sigma_i^2)(s_Q)_i^2 \leq \sum_{i=1}^l (1 - \sigma_i^2)\tau_i^2 = \varepsilon^2 \alpha^2 / n^{4t}, \end{aligned}$$

using the properties of the matrix $B_1 - C$ stated above. It follows that

$$d(G_{i_Q}, Q) \leq d(G_{i_Q}, z_Q) + d(z_Q, Q) \leq \varepsilon\alpha/n^{2t} + d(L, Q) \leq (1 + \varepsilon)d(L, Q),$$

since we assumed $d(L, Q) \geq \alpha/n^{2t}$. This proves (iii). Property (iv) is obtained similarly. Let $G_i \in \mathcal{G}$, $q \in Q$ and let z be a point in G_i . Write $z = Cw + B_1(w_K + i) + b_1$, where $w \in \prod_{i=1}^t [0, \sigma_i]$. Considering the point $z_x = B_1(w + w_K + i) + b_1$ in L , we see that

$$d(G_i, z_x)^2 \leq \|z - z_x\|^2 = \|(C - B_1)w\|^2 \leq \varepsilon^2 \alpha^2 / n^{4t}.$$

Thus,

$$d(G_i, q) \geq d(z_x, q) - d(G_i, z_x) \geq d(L, q) - \varepsilon\alpha/n^{2t} \geq (1 - \varepsilon)d(L, q).$$

□

Finally, we have a patch G , and we are looking for an approximate nearest neighbor for G in Q . The next lemma states how this can be done.

Lemma 4.8. *Suppose that $d(G, Q) \in [\alpha/2n^{2t}, 3\alpha/\varepsilon]$. We can find a point $\tilde{q} \in Q$ with $d(G, \tilde{q}) \leq (1 - 1/2 \log n)cd(G, Q)$ in total time $O_c((k^2 n^{2t}/\varepsilon^2)(m^{1-1/k+\rho/k} + (d/k) \log m))$.*

Proof. Let G_a be the projection of G onto K , and let g be the projection of G onto K^\perp . Since G and K are parallel, g is a point, and G_a is of the form $G_a : w \mapsto Cw + a_2$, with $a_2 \in K$ and $w \in \prod_{i=1}^t [0, \tau_i]$. Let $G_a^+ = \{x \in K \mid d_\infty(x, G_a) \leq 3\alpha\sqrt{k}/\varepsilon\}$, where $d_\infty(\cdot, \cdot)$ denotes the ℓ_∞ -distance with respect to the coordinate system induced by A . We subdivide the set $G_a^+ \setminus G_a$, into a collection \mathcal{C} of axis-parallel cubes, each with diameter $\varepsilon\alpha/2n^{2t}$. The cubes in \mathcal{C} have side length $\varepsilon\alpha/2n^{2t}\sqrt{k}$, the total number of cubes is $O((kn^{2t}/\varepsilon^2)^k)$, and the boundaries of the cubes lie on $O(k^2 n^{2t}/\varepsilon^2)$ hyperplanes.

We now search the partition tree \mathcal{T} to find the highest nodes (Δ, Q) in \mathcal{T} whose simplices Δ are completely contained in a single cube of \mathcal{C} . This is done as follows: we begin at the root of \mathcal{T} , and we check for all children (Δ, Q) and for all boundary hyperplanes h of \mathcal{C} whether the simplex Δ crosses the boundary h . If a child (Δ, Q) crosses no hyperplane, we label it with the corresponding cube in \mathcal{C} (or with G_a). Otherwise, we recurse on (Δ, Q) with all the boundary hyperplanes that it crosses.

In the end, we have obtained a set \mathcal{D} of simplices such that each simplex in \mathcal{D} is completely contained in a cube of \mathcal{C} . The total number of simplices in \mathcal{D} is $s = O((k^2 n^{2t}/\varepsilon^2)m^{1-1/k})$, by Theorem 3.1. For each simplex in \mathcal{D} , we query the corresponding c' -ANN structure. Let $R \subseteq Q_b$ be the set of the query results. For each point $q_b \in R$, we take the corresponding point $q \in Q$, and we compute the distance $d(q, G)$. We return a point \tilde{q} that minimizes $d(q, G)$. The query time is dominated by the time for the ANN queries. For each $\Delta \in \mathcal{D}$, let m_Δ be the number of points in the corresponding ANN structure. By assumption, an ANN-query takes time $O_c(m_\Delta^\rho + d \log m_\Delta)$, so the total query time is proportional to

$$\begin{aligned} \sum_{\Delta \in \mathcal{D}} m_\Delta^\rho + d \log m_\Delta &\leq s \left(\sum_{\Delta \in \mathcal{D}} m_\Delta / s \right)^\rho + s d \log \left(\sum_{\Delta \in \mathcal{D}} m_\Delta / s \right) \\ &\leq O_c \left((k^2 n^{2t} / \varepsilon^2) (m^{1-1/k+\rho/k} + (d/k) \log m) \right), \end{aligned}$$

using the fact that $m \mapsto m^\rho + d \log m$ is concave and that $\sum_{\Delta \in \mathcal{D}} m_\Delta \leq m$.

It remains to prove that approximation bound. Take a point q^* in Q with $d(q^*, Q) = d(Q, G)$. Since we assumed that $d(Q, G) \leq 3\alpha/\varepsilon$, the projection q_a^* of q^* onto K lies in G_a^+ . Let Δ^* be the simplex in \mathcal{D} with $q_a^* \in \Delta^*$. Suppose that the ANN-query for Δ^* returns a point $\hat{q} \in Q$. Thus, in K^\perp , we have $d(\hat{q}_b, g) \leq c' d(Q_{b\Delta^*}, g) \leq c' d(q_b^*, g)$, where \hat{q}_b and q_b^* are the projections of \hat{q} and q^* onto K^\perp and $Q_{b\Delta^*}$ is the point set stored in the ANN-structure of Δ^* . By the definition of \mathcal{C} , in K , we have $d(\hat{q}_a, G_a) \leq d(q_a^*, G_a) + \varepsilon\alpha/2n^{2t} \leq d(q_a^*, G_a) + \varepsilon d(q^*, G)$, where \hat{q}_a is the projection of \hat{q} onto K . By Pythagoras,

$$\begin{aligned} d(\hat{q}, G)^2 &= d(\hat{q}_b, g)^2 + d(\hat{q}_a, G_a)^2 \\ &\leq c'^2 d(q_b^*, g)^2 + (d(q_a^*, G_a) + \varepsilon d(q^*, G))^2 \\ &\leq c'^2 d(q_b^*, g)^2 + d(q_a^*, G_a)^2 + (2\varepsilon + \varepsilon^2) d(q^*, G)^2 \\ &\leq (c'^2 + 3\varepsilon) (q^*, G)^2 \\ &\leq ((1 - 1/\log n)^2 c^2 + 3/100 \log n) (q^*, G)^2 \\ &\leq (1 - 1/2 \log n)^2 c^2 (q^*, G)^2, \end{aligned}$$

recalling that $c' = (1 - 1/\log n)c$ and $\varepsilon = 1/100 \log n$. Since $d(\tilde{q}, G) \leq d(\hat{q}, G)$, the result follows. \square

Of all the candidate points obtained through querying patches, we return the one closest to F . The following lemma summarizes the properties of the query algorithm.

Lemma 4.9. *Suppose that $d(F, Q) \in [\alpha/n^{2t}, \alpha/\varepsilon]$. Then the query procedure returns a point $\tilde{q} \in Q$ with $d(F, \tilde{q}) \leq cd(F, Q)$ in total time $O_c((k^2 n^{2t} \varepsilon^{-5/2})^{k+1} (m^{1-1/k+\rho/k} + (d/k) \log m))$.*

Proof. By Lemmas 4.6 and 4.7, there exists a patch G with $d(G, Q) \leq (1 + \varepsilon)^2 d(F, Q)$. For this patch, the algorithm from Lemma 4.8 returns a point \hat{q} with $d(\hat{q}, G) \leq (1 + 1/2 \log n)cd(G, Q)$. Thus, using Lemma 4.7(iv), we have

$$(1 - \varepsilon)d(\hat{q}, L) \leq d(\hat{q}, G) \leq (1 - 1/2 \log n)c(1 + \varepsilon)^2 d(F, Q)$$

and by our choice of $\varepsilon = 1/100 \log n$, we get

$$\begin{aligned} (1 - 1/2 \log n)(1 + \varepsilon)^2 / (1 - \varepsilon) &\leq (1 - 1/2 \log n)(1 + 3\varepsilon)(1 + 2\varepsilon) \\ &\leq (1 - 1/2 \log n)(1 + 6/100 \log n) \leq 1. \end{aligned}$$

□

4.2.2 Far: $d(F, Q) \geq \alpha/\varepsilon$

If $d(F, Q) \geq \alpha/\varepsilon$, we can approximate Q by its projection Q_a onto K without losing too much. Thus, we can perform the whole algorithm in K . This is done by a procedure similar to Lemma 4.8.

Lemma 4.10. *Suppose we are given an estimate \tilde{r} with $d(F, Q_a) \in [\tilde{r}/2n^t, 2\tilde{r}]$. Then, we can find a point $\tilde{q} \in Q_a$ with $d(F, \tilde{q}) \leq (1 + \varepsilon)d(F, Q_a)$ in time $O((k^{3/2}n^t/\varepsilon)m^{1-1/k})$.*

Proof. Let x_F be a point in K with $d(F, K) = d(F, x_F)$. Write $x_F = Au_F + a$. Define

$$C = \prod_{i=1}^k \left((u_F)_i + \left[0, 2\tilde{r}/\sqrt{1 - \sigma_i^2} \right] \right)$$

If we take a point $x \in K$ with $d(x, F) \in [\tilde{r}/2n^t, 2\tilde{r}]$ and write $x = Au + a$, then Lemma 4.4 gives

$$d(F, x)^2 = \sum_{i=1}^k (1 - \sigma_i^2)(u - u_F)_i^2 + d(F, K)^2,$$

so $u \in C$. We subdivide C into copies of the hyperrectangle $\prod_{i=1}^k [0, \varepsilon\tilde{r}/2n^t \sqrt{k(1 - \sigma_i^2)}]$. Let \mathcal{C} be the resulting set of hyperrectangles. The boundaries of the hyperrectangles in \mathcal{C} lie on $O(k^{3/2}n^t/\varepsilon)$ hyperplanes. We now search the partition tree \mathcal{T} in order to find the highest nodes (Δ, Q) in \mathcal{T} whose simplices Δ are completely contained in a single hyperrectangle of \mathcal{C} . This is done in the same way as in Lemma 4.8.

This gives a set \mathcal{D} of simplices such that each simplex in \mathcal{D} is completely contained in a hyperrectangle of \mathcal{C} . The total number of simplices in \mathcal{D} is $O((k^{3/2}n^t/\varepsilon)m^{1-1/k})$, by Theorem 3.1. For each simplex $\Delta \in \mathcal{D}$, we pick an arbitrary point $q \in Q_a$ that lies in Δ , and we compute $d(F, q)$. We return the point $\tilde{q} \in Q_a$ that minimizes the distance to F . The total query time is $O((k^{3/2}n^t/\varepsilon)m^{1-1/k})$.

Now let q^* be a point in Q_a with $d(F, Q_a) = d(F, q^*)$, and let Δ^* be the simplex \mathcal{D} that contains q^* . Furthermore, let $\hat{q} \in Q_a$ be the point that the algorithm examines in Δ^* . Write $q^* = Au^* + a$ and $\hat{q} = A\hat{u} + a$. Since q^* and \hat{q} lie in the same hyperrectangle and by Lemma 4.4,

$$\begin{aligned} d(F, \hat{q})^2 &= \sum_{i=1}^k (1 - \sigma_i^2)(\hat{u} - u_F)_i^2 + d(F, K)^2 \leq \\ &\sum_{i=1}^k (1 - \sigma_i^2)(u^* - u_F)_i^2 + \varepsilon^2 \tilde{r}^2 / 4n^{2t} + d(F, K)^2 \leq (1 + \varepsilon)^2 d(F, q^*)^2. \end{aligned}$$

Since $d(F, \tilde{q}) \leq d(F, \hat{q})$, the result follows. □

Lemma 4.11. *Suppose we are given an estimate \tilde{r} with $d(F, Q) \in [\tilde{r}/n^t, \tilde{r}]$. Suppose further that $d(F, Q) \geq \alpha/\varepsilon$. Then we can find a $\tilde{q} \in Q$ with $d(F, \tilde{q}) \leq cd(F, Q)$ in time $O((k^{3/2}n^{2t}/\varepsilon)m^{1-1/k})$.*

Proof. For any point $q \in Q$, let $q_a \in Q$ be its projection onto K . Then, $d(q_a, q) \leq \alpha \leq \varepsilon d(F, Q)$. Thus, $d(F, Q_a) \in [(1 - \varepsilon)d(F, Q), (1 + \varepsilon)d(F, Q)]$, and we can apply Lemma 4.10. Let $\tilde{q}_a \in Q_a$ be the result of this query, and let \tilde{q} be the corresponding point in Q . We have

$$\begin{aligned} d(F, \tilde{q}) &\leq d(\tilde{q}, \tilde{q}_a) + d(F, \tilde{q}_a) \leq \varepsilon d(F, Q) + (1 + \varepsilon)d(F, Q_a) \\ &\leq \varepsilon d(F, Q) + (1 + \varepsilon)^2 d(F, Q) \leq (1 + 4\varepsilon)d(F, Q) \leq cd(F, Q), \end{aligned}$$

by our choice of ε . □

By combining Lemmas 4.1, 4.9, and 4.11, we obtain Theorem 2.1.

5 Approximate k -flat Range Reporting in Low Dimensions

In this section, we present a data structure for low dimensional k -flat approximate near neighbor reporting. In Section 6, we will use it as a foundation for our projection structures. The details of the structure are summarized in Theorem 5.1. Throughout this section, we will think of d as a constant, and we will suppress factors depending on d in the O -notation.

Theorem 5.1. *Let $P \subset \mathbb{R}^d$ be an n -point set. We can preprocess P into an $O(n \log^{d-k-1} n)$ space data structure for approximate k -flat near neighbor queries: given a k -flat F and a parameter α , find a set $R \subseteq P$ that contains all $p \in P$ with $d(p, F) \leq \alpha$ and no $p \in P$ with $d(p, F) > ((4k + 3)(d - k - 1) + \sqrt{k + 1})\alpha$. The query time is $O(n^{k/(k+1)} \log^{d-k-1} n + |R|)$.*

5.1 Preprocessing

Let $E \subset \mathbb{R}^d$ be the $(k + 1)$ -dimensional subspace of \mathbb{R}^d spanned by the first $k + 1$ coordinates, and let Q be the projection of P onto E .¹ We build a $(k + 1)$ -dimensional partition tree \mathcal{T} for Q , as in Theorem 3.1. If $d > k + 1$, we also build a *slab structure* for each node of \mathcal{T} . Let v be such a node, and let Ξ be the simplicial partition for the children of v . Let $w > 0$. A w -slab S is a closed region in E that is bounded by two parallel hyperplanes of distance w . The *median hyperplane* \hat{h} of S is the hyperplane inside S that is parallel to the two boundary hyperplanes and has distance $w/2$ from both. A w -slab S is *full* if there are at least $r^{2/3}$ simplices Δ in Ξ with $\Delta \subset S$.

```

Input: point set  $P \subset \mathbb{R}^d$ 
1 if  $|P| = O(1)$  then
2   | Store  $P$  in a list and return.
3  $Q \leftarrow$  projection of  $P$  onto the subspace  $E$  spanned by the first  $k + 1$  coordinates.
4  $\mathcal{T} \leftarrow$   $(k + 1)$ -dimensional partition tree for  $Q$  as in Theorem 3.1.
5 if  $d > k + 1$  then
6   | foreach node  $v \in \mathcal{T}$  do
7     |  $\Xi_1 \leftarrow$  simplicial partition for the children of  $v$ 
8     | for  $j \leftarrow 1$  to  $\lfloor r^{1/3} \rfloor$  do
9       |    $D_j \leftarrow$  CreateSlabStructure $(\Xi_j)$ 
10      |    $\Xi_{j+1} \leftarrow \Xi_j$  without all simplices inside the slab for  $D_j$ 

```

Algorithm 5: CreateSearchStructure

The slab structure for v is constructed in several iterations. In iteration j , we have a current subset $\Xi_j \subseteq \Xi$ of pairs in the simplicial partition. For each $(k + 1)$ -set v_0, \dots, v_k of vertices of simplices in Ξ_j , we determine the smallest width of a full slab whose median hyperplane is spanned by v_0, \dots, v_k . Let S_j be the smallest among those slabs, and let \hat{h}_j be its median

¹ We assume general position: any two distinct points in P have distinct projections in Q .

Input: $\Xi_j = (Q_1, \Delta_1), \dots, (Q_{r'}, \Delta_{r'})$

- 1 $V_j \leftarrow$ vertices of the simplices in Ξ_j
- 2 For each $(k+1)$ -subset $V \subset V_j$, find the smallest $w_V > 0$ such that the w_V -slab with median hyperplane $\text{aff}(V)$ is full.
- 3 Let w_j be the smallest w_V ; let S_j be the corresponding full w_j -slab and $\hat{h}_j = \text{aff}(V)$ its median hyperplane.
- 4 Find the set \mathcal{D}_j of $r^{2/3}$ simplices in S_j ; let $\mathcal{Q}_j \leftarrow \bigcup_{\Delta_i \in \mathcal{D}_j} Q_i$ and let \mathcal{P}_j be the d -dimensional point set corresponding to \mathcal{Q}_j .
- 5 $h_j \leftarrow$ the hyperplane orthogonal to E through \hat{h}_j
- 6 $P' \leftarrow$ projection of \mathcal{P}_j onto h_j /* P' is $(d-1)$ -dimensional */
- 7 **CreateSearchStructure**(P')

Algorithm 6: CreateSlabStructure

hyperplane. Let \mathcal{D}_j be the $r^{2/3}$ simplices that lie completely in S_j . We remove \mathcal{D}_j and the corresponding point set $\mathcal{Q}_j = \bigcup_{\Delta_i \in \mathcal{D}_j} Q_i$ from Ξ_j to obtain Ξ_{j+1} . Let $\mathcal{P}_j \subseteq P$ be the d -dimensional point set corresponding to \mathcal{Q}_j . We project \mathcal{P}_j onto the d -dimensional hyperplane h_j that is orthogonal to E and goes through \hat{h}_j . We recursively build a search structure for the $(d-1)$ -dimensional projected point set. The j th slab structure D_j at v consists of this search structure, the hyperplane h_j , and the width w_j . This process is repeated until less than $r^{2/3}$ simplices remain; see Algorithms 5 and 6 for details.

Denote by $S(n, d)$ the space for a d -dimensional search structure with n points. The partition tree \mathcal{T} has $O(n)$ nodes, so the overhead for storing the slabs and partitions is linear. Thus,

$$S(n, d) = O(n) + \sum_D S(n_D, d-1),$$

where the sum is over all slab structures D and where n_D is the number of points in the slab structure D . Since every point appears in $O(\log n)$ slab structures, and since the recursion stops for $d = k+1$, we get

Lemma 5.2. *The search structure for n points in d dimensions needs space $O(n \log^{d-k-1} n)$. \square*

5.2 Processing a Query

For a query, we are given a distance threshold $\alpha > 0$ and a k -flat F . For the recursion, we will need to query the search structure with a k -dimensional polytope. We obtain the initial query polytope by intersecting the flat F with the bounding box of P extended by α in each direction. With slight abuse of notation, we still call this polytope F .

A query for F and α is processed by using the slab structures for small enough slabs and by recursing in the partition tree for the remaining points. Details follow.

Suppose we are at some node v of the partition tree, and let j^* be the largest integer with $w_{j^*} \leq (4k+2)\alpha$. For $j = 1, \dots, j^*$, we recursively query each slab structure D_j as follows: let $\tilde{F} \subseteq F$ be the polytope containing the points in F with distance at most $\alpha + w_j/2$ from h_j , and let F_h be the projection of \tilde{F} onto h_j . We query the search structure in D_j with F_h and α . Next, we project F onto the subspace E spanned by the first $k+1$ coordinates. Let \mathcal{D} be the simplices in Ξ_{j^*+1} with distance at most α from the projection. For each simplex in \mathcal{D} , we recursively query the corresponding child in the partition tree. Upon reaching the bottom of the recursion (i.e., $|P| = O(1)$), we collect all points within distance α from F in the set R .

If $d = k+1$, we approximate the region of interest by the polytope $F_\diamond = \{x \in \mathbb{R}^d \mid d_1(x, F) \leq \alpha\}$, where $d_1(\cdot, \cdot)$ denotes the ℓ_1 -metric in \mathbb{R}^d . Then, we query the partition tree \mathcal{T} to find all points of P that lie inside F_\diamond . We prove in Lemma 5.4 that F_\diamond is a polytope with $O(d^{O(k^2)})$

```

Input : polytope  $F$ , distance threshold  $\alpha > 0$ 
Output: point set  $R \subseteq P$ 
1  $R \leftarrow \emptyset$ 
2 if  $|P| = O(1)$  then
3    $R \leftarrow \{p \in P \mid d(p, F) \leq \alpha\}$ 
4 else if  $d = k + 1$  then
5   Compute polytope  $F_\diamond$  as described.
6    $R \leftarrow R \cup$  all points of  $P$  inside  $F_\diamond$ 
7 else
8    $j^* \leftarrow$  the largest integer with  $w_{j^*} \leq (4k + 2)\alpha$ 
9   for  $j \leftarrow 1$  to  $j^*$  do
10     $F_h \leftarrow$  projection of  $\tilde{F}$  onto  $h_j$  as described
11     $R \leftarrow R \cup D_j.\text{query}(F_h, \alpha)$ 
12   $\hat{F} \leftarrow$  projection of  $F$  onto the subspace  $E$  spanned by the first  $k + 1$  coordinates
13   $\mathcal{D} \leftarrow$  simplices in  $\Xi_{j^*+1}$ 
14   $\mathcal{D}' \leftarrow \{\Delta \in \mathcal{D} \mid d(\Delta, \hat{F}) \leq \alpha\}$ 
15  foreach  $\Delta \in \mathcal{D}'$  do
16     $R \leftarrow R \cup$  result of recursive query to partition tree node for  $\Delta$ .
17 return  $R$ 

```

Algorithm 7: Find a superset R of all points in P with distance less than α from a query polytope F .

facets; see Algorithm 7 for details. The following two lemmas analyze the correctness and query time of the algorithm.

Lemma 5.3. *The set R contains all $p \in P$ with $d(p, F) \leq \alpha$ and no $p \in P$ with $d(p, F) > \kappa\alpha$, where $\kappa = (4k + 3)(d - k - 1) + \sqrt{k + 1}$.*

Proof. The proof is by induction on the size n of P and on the dimension d . If $n = O(1)$, we return all points with distance at most α to F . If $d = k + 1$, we report the points inside the polytope F_\diamond (lines 4–6) using \mathcal{T} . Since $\|x\|_2 \leq \|x\|_1 \leq \sqrt{k + 1}\|x\|_2$ holds for all $x \in \mathbb{R}^{k+1}$, the polytope F_\diamond contains all points with distance at most α from F and no point with distance more than $\alpha\sqrt{k + 1}$ from F . Thus, correctness also follows in this case.

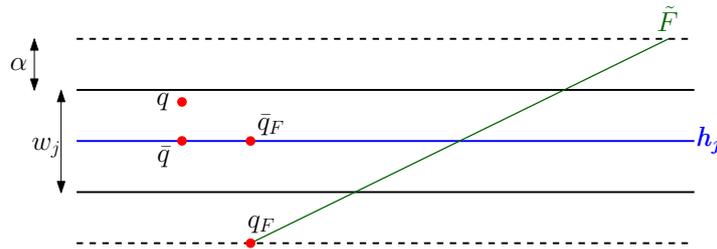


Fig. 1: The distance error due to the reduction of the dimension in the slab structure D_j .

In the general case ($d > k + 1$ and n not constant), we prove the lemma individually for the slab structures and for the partition tree. Let D_j be a slab structure and P_j the corresponding d -dimensional point set. Fix some point $p \in P_j$ with $d(p, F) \leq \alpha$. To query D_j , we take the subpolytope $\tilde{F} \subseteq F$ with distance at most $\alpha + w_j/2$ from the median hyperplane h_j , and we project it onto h_j . Let F_h be this projection. Since orthogonal projections can only decrease distances, we have $p \in D_j.\text{query}(F_h, \alpha)$ by induction. Now fix a point $q \in D_j.\text{query}(F_h, \alpha)$. We must argue that $d(q, F) \leq (4k + 3)(d - k - 1) + \sqrt{k + 1}$. Let $\bar{q} \in \mathbb{R}^{d-1}$ be the projection of q

onto h_j and $\bar{q}_F \in F_h$ the closest point to \bar{q} in F_h . Let $q_F \in \tilde{F}$ be the corresponding d -dimensional point (see Fig. 1). By triangle inequality and the induction hypothesis,

$$\begin{aligned} d(q, F) &= d(q, \tilde{F}) \leq d(q, \bar{q}) + d(\bar{q}, \bar{q}_F) + d(\bar{q}_F, q_F) \\ &\leq w_j/2 + ((4k+3)(d-k-2) + \sqrt{k+1})\alpha + (\alpha + w_j/2). \end{aligned}$$

By construction, we have $w_j \leq (4k+2)\alpha$, so $d(q, F) \leq ((4k+3)(d-k-1) + \sqrt{k+1})\alpha$, as claimed.

Consider now a child in the partition tree queried in line 16, and let P_j be the corresponding d -dimensional point set. Since $|P_j| < |P|$, the claim follows by induction. \square

Lemma 5.4. *The query time is $O(n^{k/(k+1)} \log^{d-k-1} n + |R|)$.*

Proof. Let $Q(n, d)$ be the total query time.

First, let $d > k + 1$. We bound the time to query the partition tree \mathcal{T} . Let \hat{F} be the projection of F onto E . Furthermore, let V be the set of nodes in \mathcal{T} that are visited during a query, and let \mathcal{D} be the corresponding simplices. By construction, all simplices in \mathcal{D} have distance at most α from \hat{F} . Consider the 2α -slab S whose median hyperplane contains \hat{F} . We partition V into two sets: the nodes V_B whose simplices intersect ∂S , and the nodes V_C whose simplices lie completely in S . First, since the simplex for each node in \mathcal{T} is contained in the simplex for its parent node, we observe that V_B constitutes a connected subtree of \mathcal{T} , starting at the root. The nodes of V_C form several connected subtrees, each hanging off a node in V_B . Furthermore, by construction, each node from V has at most $r^{1/3}$ children from V_B . Let V_ℓ be the set of nodes in V with level ℓ , for $\ell = 0, \dots, \log_r n$, and let $m_\ell = |V_\ell|$. By Theorem 3.1, we have $|V_\ell \cap V_B| \leq O(r^{\ell k/(k+1)} + r^{(k-1)/k} m_{\ell-1} + r\ell \log r \log n)$. Since $|V_\ell \cap V_C| \leq r^{1/3} m_{\ell-1}$, we get

$$m_\ell = |V_\ell| = O\left(r^{\ell k/(k+1)} + (r^{(k-1)/k} + r^{1/3})m_{\ell-1} + r\ell \log r \log n\right).$$

For any $\delta > \max(0, 1/3 - (k-1)/k)$, if we choose r large enough depending on δ , this solves to $m_\ell = O(r^{\ell k/(k+1)} + r^{\ell((k-1)/k+\delta)} \log n)$. Thus, we get

$$Q(n, d) = \sum_{\ell=0}^{\log_r n} O(r^{\ell k/(k+1)} + r^{\ell((k-1)/k+\delta)} \log n)(O(r) + Q(n/r^\ell, d-1)). \quad (4)$$

For $d = k + 1$, we use \mathcal{T} directly. Thus, by Theorem 3.1, the query time $Q(n, k+1)$ is $O(f_{k+1}n^{1-1/k} + |R_F|)$, where f_{k+1} is the number of facets of F_\diamond and R_F is the answer set. We claim that f_{k+1} is bounded by $((2k+2)(5d_0 - 2k)^{k/2})^{(k+1)/2}$. Recall that F_\diamond is the Minkowski sum of F and the ℓ_1 -ball with radius α as in Algorithm 7 line 5. Initially F is the intersection of the query k -flat with the extended bounding box of P . This intersection can be described by at most $d_0 - k + 2d_0 = 3d_0 - k$ oriented half-spaces, where d_0 denotes the initial dimension. In each recursive step, we intersect F with the 2 bounding hyperplanes of a slab. Therefore, the descriptive complexity of F in the base case is at most $5d_0 - 2k$. By duality and the Upper Bound theorem [17], the \mathcal{V} -description of F consists of at most $(5d_0 - 2k)^{k/2}$ vertices. Using that the Minkowski sum of two polytopes with v_1 and v_2 vertices has at most $v_1 v_2$ vertices, we deduce that F_\diamond has at most $(2k+2)(5d_0 - 2k)^{k/2}$ vertices. Applying the upper bound theorem again, it follows that $f_{k+1} = ((2k+2)(5d_0 - 2k)^{k/2})^{(k+1)/2}$, as claimed.

Thus, plugging the base case into (4), we get that the overall query time $Q(n, d)$ is bounded by $O(n^{k/(k+1)} \log^{d-k-1} n + |R|)$. \square

Theorem 5.1 follows immediately from Lemmas 5.2, 5.3, and 5.4.

5.3 Approximate k -Flat Nearest Neighbor Queries

We now show how to extend our data structure from Section 5.1 for approximate k -flat nearest neighbor queries with multiplicative error $(4k + 3)(d - k - 1) + \sqrt{k + 1}$. That is, given an n -point set $P \subset \mathbb{R}^d$, we want to find for any given query flat $F \subset \mathbb{R}^d$ a point $p \in P$ with $d(p, F) \leq ((4k + 3)(d - k - 1) + \sqrt{k + 1})d(P, F)$. We reduce this problem to a near neighbor query by choosing an appropriate threshold α that ensures $|R| = O(\sqrt{n})$, using random sampling. For preprocessing we build the data structure D from Theorem 5.1 for P .

Let a query flat F be given. The F -rank of a point $p \in P$ is the number of points in P that are closer to F than p . Let $X \subseteq P$ be a random sample obtained by taking each point in P independently with probability $1/\sqrt{n}$. The expected size of X is \sqrt{n} , and if $x \in X$ is the closest point to F in X , then the expected F -rank of x is \sqrt{n} . Set $\alpha = d(x, F)/((4k + 3)(d - k - 1) + \sqrt{k + 1})$. We query D with F and α to obtain a set R . If $d(P, F) \leq \alpha$, then R contains the nearest neighbor. Otherwise, x is a $((4k + 3)(d - k - 1) + \sqrt{k + 1})$ -approximate nearest neighbor for F . Thus, it suffices to return the nearest neighbor in $R \cup \{x\}$. Since with high probability all points in R have F -rank at most $O(\sqrt{n} \log n)$, we have $|R| = O(\sqrt{n} \log n)$, and the query time is $O(n^{k/(k+1)} \log^{d-k-1} n)$. This establishes the following corollary of Theorem 5.1.

Corollary 5.5. *Let $P \subset \mathbb{R}^d$ be an n -point set. We can preprocess P into an $O(n \log^{d-k-1} n)$ space data structure for approximate k -flat nearest neighbor queries: given a flat F , find a point $p \in P$ with $d(p, F) \leq ((4k + 3)(d - k - 1) + \sqrt{k + 1})d(P, F)$. The expected query time is $O(n^{k/(k+1)} \log^{d-k-1} n)$.*

6 Projection Structures

We now describe how to answer queries of type **Q1** and **Q3** efficiently. Our approach is to project the points into random subspace of constant dimension and to solve the problem there using our data structures from Theorem 5.1 and Corollary 5.5. For this, we need a Johnson-Lindenstrauss-type lemma that bounds the distortion, see Section 6.1.

Let $0 < t \leq 2/(2 + 40k)$ be a parameter and let $P \subset \mathbb{R}^d$ be a high dimensional n -point set. Set $d' = 2/t + 2$ and let $M \in \mathbb{R}^{d' \times d}$ be a random projection from \mathbb{R}^d to $\mathbb{R}^{d'}$, scaled by $\sqrt{d/4d'}$. We obtain $\bar{P} \subset \mathbb{R}^{d'}$ by projecting P using M . We build for \bar{P} the data structure D_1 from Corollary 5.5 to answer **Q1** queries and D_2 from Theorem 5.1 to answer **Q3** queries. This needs $O(n \log^{O(d')} n) = O(n \log^{O(1/t)} n)$ space. For each $p \in P$ we write \bar{p} for the d' -dimensional point Mp and \bar{F} for the projected flat MF .

6.1 Dimension Reduction

We use the following variant of the Johnson-Lindenstrauss-Lemma, as proved by Dasgupta and Gupta [8, Lemma 2.2].

Lemma 6.1 (JL-Lemma). *Let $d' < d$, and let $M \in \mathbb{R}^{d' \times d}$ be the projection matrix onto a random d' -dimensional subspace, scaled by a factor of $\sqrt{d/d'}$. Then, for every vector $x \in \mathbb{R}^d$ of unit length and every $\beta > 1$, we have*

1. $\Pr[\|Mx\|^2 \geq \beta] \leq \exp(-\frac{1}{2}d'(1 - \beta + \ln \beta))$, and
2. $\Pr[\|Mx\|^2 \leq 1/\beta] \leq \exp(-\frac{1}{2}d'(1 - 1/\beta - \ln \beta)) \leq \exp(-\frac{1}{2}d'(1 - \ln \beta))$. □

Lemma 6.2. *Let $p \in \mathbb{R}^d$ be a point and let $F \subset \mathbb{R}^d$ be a k -flat. For $d' \in \{40k, \dots, d - 1\}$, let $M \in \mathbb{R}^{d' \times d}$ be the projection matrix into a random d' -dimensional subspace, scaled by $\sqrt{d/4d'}$. Let $\bar{p} = Mp$ and $\bar{F} = MF$ be the projections of p and of F , respectively. Then, for any $\beta \geq 40k$, (i) $\Pr[d(\bar{p}, \bar{F}) \leq d(p, F)] \geq 1 - e^{-d'/2}$; and (ii) $\Pr[d(\bar{p}, \bar{F}) \geq d(p, F)/\beta] \geq 1 - \beta^{-d'/2}$.*

Proof. Let $N = 2M$, and set $q = Np$ and $K = NF$. Defining $\Delta_p = d(p, F)$ and $\Delta_q = d(q, K)$, we must bound the probabilities $\Pr[\Delta_q \leq 2\Delta_p]$ for (i) and $\Pr[\Delta_q \geq 2\Delta_p/\beta]$ for (ii).

We begin with (i). Let p^\parallel be the orthogonal projection of p onto F , and let $p^\perp = p - p^\parallel$. Let $q^\perp = Np^\perp$. Then, $\Delta_p = \|p^\perp\|$ and $\Delta_q \leq \|q^\perp\|$. By Lemma 6.1(1),

$$\begin{aligned} \Pr \left[\|q^\perp\| \geq 2\Delta_p \right] &= \Pr \left[\|Np^\perp\|/\|p^\perp\| \geq 2 \right] = \Pr \left[\|N(p^\perp/\|p^\perp\|)\|^2 \geq 4 \right] \\ &\leq \exp \left(\frac{1}{2}d'(1 - 4 + \ln 4) \right) \leq \exp(-d'/2). \end{aligned}$$

Thus, $\Pr[\Delta_q \leq 2\Delta_p] \geq \Pr[\|q^\perp\| \leq 2\Delta_p] \geq 1 - \exp(-d'/2)$, as desired.

For (ii), choose k orthonormal vectors e_1, \dots, e_k such that $F = \left\{ p^\parallel + \sum_{i=1}^k \lambda_i e_i \mid \lambda_i \in \mathbb{R} \right\}$. Set $u_i = e_i \|p^\perp\|/\beta^2$, and consider the lattice $L = \left\{ p^\parallel + \sum_{i=1}^k \mu_i u_i \mid \mu_i \in \mathbb{Z} \right\} \subset F$. Let $\bar{L} = NL$ be the projected lattice. We next argue that with high probability (i) all points in \bar{L} have distance at least $3\|p^\perp\|/\beta$ from q ; and (ii) for $i = 1, \dots, k$, we have $\|Nu_i\| < \|p^\perp\|/\beta\sqrt{k}$.

To show (i), we partition L into *layers*: for $j \in \{0, 1, \dots\}$, let

$$L_j = \left\{ p^\parallel + \sum_{i=1}^k \mu_i u_i \mid \mu_i \in \{-j, \dots, j\}, \max_i |\mu_i| = j \right\} \subset L.$$

Now for any $j \in \mathbb{N}$ and $r = p^\parallel + \sum_{i=1}^k \mu_i u_i \in L_j$, Pythagoras gives

$$\|p - r\| = \sqrt{\|p^\perp\|^2 + \left\| \sum_{i=1}^k \mu_i u_i \right\|^2} = \|p^\perp\| \sqrt{1 + \sum_{i=1}^k |\mu_i|^2/\beta^4} \geq \|p^\perp\| \sqrt{1 + j^2/\beta^4}.$$

Thus, using Lemma 6.1(2),

$$\begin{aligned} \Pr \left[\|N(p - r)\| \leq 3\|p^\perp\|/\beta \right] &= \Pr \left[\|N(p - r)\|/\|p - r\| \leq 3\|p^\perp\|/\beta\|p - r\| \right] \\ &\leq \Pr \left[\|N(p - r)\|/\|p - r\|^2 \leq 9/\beta^2(1 + j^2/\beta^4) \right] \\ &\leq \exp \left(\frac{1}{2}d'(1 + \ln(9/\beta^2(1 + j^2/\beta^4))) \right) \\ &\leq (5/\beta)^{d'} (1 + j^2/\beta^4)^{-d'/2}, \end{aligned}$$

as $\sqrt{9e} \leq 5$. Now we use a union bound to obtain

$$\Pr \left[\exists r \in L : \|N(p - r)\| \leq 3\|p^\perp\|/\beta \right] = \sum_{j=0}^{\infty} \Pr \left[\exists r \in L_j : \|N(p - r)\| \leq 3\|p^\perp\|/\beta \right].$$

Grouping the summands into groups of β^2 consecutive terms, this is

$$\begin{aligned} &= \sum_{l=0}^{\infty} \sum_{j=l\beta^2}^{l\beta^2+\beta^2-1} \Pr \left[\exists r \in L_j : \|N(p - r)\| \leq 3\|p^\perp\|/\beta \right] \\ &\leq (5/\beta)^{d'} \sum_{l=0}^{\infty} |L_{\leq (l+1)\beta^2}| (1 + (l\beta^2)^2/\beta^4)^{-d'/2}, \end{aligned}$$

where $L_{\leq a} = \bigcup_{i=0}^a L_i$. Using the rough bound $|L_{\leq a}| \leq (3a)^k$, this is

$$\begin{aligned} &\leq (5/\beta)^{d'} \sum_{l=0}^{\infty} (3(l+1)\beta^2)^k (1 + l^2)^{-d'/2} \\ &= 5^{d'} 3^k \beta^{2k-d'} \sum_{l=0}^{\infty} (l+1)^k (1 + l^2)^{-d'/2}. \end{aligned}$$

For $d' \geq 4k$, we have $(l+1)^k(1+l^2)^{-d'/2} \leq (l+1)^k(1+l^2)^{-2k} \leq 1/(1+l^2)$, so we can bound the sum by $\sum_{l=1}^{\infty} 1/(1+l^2) \leq \pi^2/6$. Thus, we have derived

$$\Pr \left[\exists r \in L : \|N(p-r)\| \leq 3\|p^\perp\|/\beta \right] \leq 5^{d'} 3^k \beta^{2k-d'} (\pi^2/6) \leq 5^{d'+k} \beta^{2k-d'}, \quad (5)$$

since $\pi^2/6 \leq 5/3$.

To show (ii), we use a union bound with Lemma 6.1(1). Recalling $\|u_1\| = \|p^\perp\|/\beta^2$,

$$\begin{aligned} \Pr \left[\exists i = 1, \dots, k : \|Nu_i\| > \|p^\perp\|/\sqrt{k}\beta \right] &\leq k \Pr \left[\|Nu_1\| > \|p^\perp\|/\sqrt{k}\beta \right] \\ &= k \Pr \left[\|N(u_1/\|u_1\|)\|^2 > \beta^2/k \right] \\ &\leq k \exp \left(\frac{1}{2} d' (1 - \beta^2/k + \ln(\beta^2/k)) \right) \\ &\leq k \exp \left(-\beta^2 d'/4k \right), \end{aligned} \quad (6)$$

since $c^2/k \geq 2(1 + \ln(\beta^2/k))$ for $\beta^2/k \geq 6$. By (5) and (6), and recalling $\beta, d' \geq 40k$, the probability that events (i) and (ii) do not both happen is at most

$$\begin{aligned} 5^{d'+k} \beta^{2k-d'} + k e^{-\beta^2 d'/4k} &\leq 5^{(1+1/40)d'} \beta^{(1/20-1)d'} + (\beta/40) e^{-10cd'} \\ &\leq \left(\frac{5^{41/40}}{40^{9/20}} \right)^{d'} \beta^{-d'/2} + \frac{1}{10} \beta^{-d'/2} \leq \beta^{-d'/2}. \end{aligned}$$

Suppose (i) and (ii) happen. Fix a point $w \in K$, and let $\bar{r} \in \bar{L}$ be the point in the projected lattice that is closest to w . By (i), $d(q, \bar{r}) > 3\|p^\perp\|/\beta$. By (ii) and the choice of \bar{r} , the k -dimensional cube with center \bar{r} and side length $\|p^\perp\|/\beta\sqrt{k}$ contains w . This cube has diameter $\|p^\perp\|/\beta$. By triangle inequality, $d(q, w) > d(q, \bar{r}) - d(\bar{r}, w) \geq (3/\beta - 1/\beta)\|p^\perp\| = 2\|p^\perp\|/\beta$. \square

6.2 Queries of Type Q1

Let a query flat F be given. To answer **Q1** queries, we compute \bar{F} and query D_1 with \bar{F} to obtain a $((4k+3)(d'-k-1) + \sqrt{k+1})$ -nearest neighbor \bar{p} . We return the original point p . To obtain Theorem 2.2, we argue that if \bar{p} is a $((4k+3)(d'-k-1) + \sqrt{k+1})$ -nearest neighbor for \bar{F} , then p is a n^t -nearest neighbor for F with high probability.

Let $p^* \in P$ be a point with $d(p^*, F) = d(P, F)$. Set $\delta_{p^*} = d(p^*, F)$ and $\bar{\delta}_{p^*} = d(\bar{p}^*, \bar{F})$. Denote by A_1 the event that $\bar{\delta}_{p^*} \leq \delta_{p^*}$. By Lemma 6.2, $\Pr[A_1] \geq 1 - e^{-d'/2} = 1 - e^{-1/t-1}$. Let A_2 be the event that for all points $p \in P$ with $\delta_p = d(p, F) > n^t \delta_{p^*}$ we have $\bar{\delta}_p = d(\bar{p}, \bar{F}) > ((4k+3)(d'-k-1) + \sqrt{k+1})\delta_{p^*}$. For a fixed $p \in P$, by setting $\beta = n^t/((4k+3)(d'-k-1) + \sqrt{k+1})$ in Lemma 6.2, this probability is seen to be at least $1 - n^{-1-t/2}$, for n large enough. By the union bound, we get $\Pr[A_2] \geq 1 - n^{-t/2}$, so the event $A_1 \cap A_2$ occurs with constant probability. Then, p is a n^t -approximate nearest neighbor for F , as desired.

6.3 Queries of Type Q3

To answer a query of type **Q3**, we compute the projection \bar{F} and query D_2 with parameter α . We obtain a set $\bar{R} \subset \bar{P}$ in time $O(n^{k/(k+1)} \log^{O(1/t)} n + |\bar{R}|)$. Let $R \subset P$ be the corresponding d -dimensional set. We return a point $p \in R$ that minimizes $d(p, F)$. If $\delta_{p^*} \leq \alpha$, the event A_1 from above implies that $\bar{p}^* \in \bar{R}$, and we correctly return p^* .

To bound the size of $|\bar{R}|$, and thus the running time, we use that P is $\alpha n^t/(2k+1)$ -cluster-free. Let A_3 be the event that for all $p \in P$ with $d(p, F) > \alpha n^t/(2k+1)$, we have $d(\bar{p}, \bar{F}) > ((4k+$

$3)(d' - k - 1) + \sqrt{k + 1})\alpha$. By the definition of cluster-freeness and the guarantee of Theorem 5.1, we have $|\bar{R}| = m$ in the case of A_3 . Using $\beta = n^t / ((2k + 1)((4k + 3)(d - k - 1) + \sqrt{k + 1}))$ in Lemma 6.2 and doing a similar calculation as above yields again $\Pr[A_3] \geq 1 - n^{-t/2}$. Thus, we can answer queries of type **Q3** successfully in time $O(n^{k/(k+1)} \log^{O(1/t)} n + m)$ with constant probability, as claimed in Theorem 2.3.

7 Conclusion

We have described the first provably efficient data structure for general k -ANN. Our main technical contribution consists of two new data structures: the cluster data structure for high-dimensional k -ANN queries, and the projection data structure for k -ANN queries in constant dimension. We have only presented the latter structure for a constant approximation factor $(4k + 3)(d - k - 1) + \sqrt{k + 1}$, but we believe that it is possible to extend it to any fixed approximation factor $c > 1$. For this, one would need to subdivide the slab structures by a sufficiently fine sequence of parallel hyperplanes.

Naturally, the most pressing open question is to improve the query time of our data structure. Also, a further generalization to more general query or data objects would be of interest.

Acknowledgments

This work was initiated while WM, PS, and YS were visiting the Courant Institute of Mathematical Sciences. We would like to thank our host Esther Ezra for her hospitality and many enlightening discussions.

References

- [1] A. Andoni. *Nearest Neighbor Search: the Old, the New, and the Impossible*. PhD thesis, Massachusetts Institute of Technology, 2009.
- [2] A. Andoni, P. Indyk, R. Krauthgamer, and H. L. Nguyen. Approximate line nearest neighbor in high dimensions. In *Proc. 20th Annu. ACM-SIAM Sympos. Discrete Algorithms (SODA)*, pages 293–301, 2009.
- [3] A. Andoni, P. Indyk, H. L. Nguyen, and I. Razenshteyn. Beyond locality-sensitive hashing. In *Proc. 24th Annu. ACM-SIAM Sympos. Discrete Algorithms (SODA)*, pages 1018–1028, 2014.
- [4] R. Basri, T. Hassner, and L. Zelnik-Manor. Approximate nearest subspace search with applications to pattern recognition. In *Proc. Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–8, 2007.
- [5] T. M. Chan. Optimal partition trees. *Discrete Comput. Geom.*, 47(4):661–690, 2012.
- [6] B. Chazelle. *The discrepancy method. Randomness and complexity*. Cambridge University Press, 2000.
- [7] K. L. Clarkson. A randomized algorithm for closest-point queries. *SIAM J. Comput.*, 17(4):830–847, 1988.
- [8] S. Dasgupta and A. Gupta. An elementary proof of a theorem of Johnson and Lindenstrauss. *Random Structures Algorithms*, 22(1):60–65, 2003.
- [9] R. A. Horn and C. R. Johnson. *Matrix analysis*. Cambridge University Press, second edition, 2013.

-
- [10] P. Indyk. Nearest neighbors in high-dimensional spaces. In J. E. Goodman and J. O'Rourke, editors, *Handbook of Discrete and Computational Geometry*, chapter 39. CRC Press, 2nd edition, 2004.
- [11] P. Indyk and R. Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proc. 30th Annu. ACM Sympos. Theory Comput. (STOC)*, pages 604–613, 1998.
- [12] E. Kushilevitz, R. Ostrovsky, and Y. Rabani. Efficient search for approximate nearest neighbor in high dimensional spaces. In *Proc. 30th Annu. ACM Sympos. Theory Comput. (STOC)*, pages 614–623, 1998.
- [13] Q. Lv, W. Josephson, Z. Wang, M. Charikar, and K. Li. Multi-probe LSH: Efficient indexing for high-dimensional similarity search. In *Proc. 33rd International Conference on Very Large Data Bases (VLDB)*, pages 950–961, 2007.
- [14] A. Magen. Dimensionality reductions in ℓ_2 that preserve volumes and distance to affine spaces. *Discrete Comput. Geom.*, 38(1):139–153, 2007.
- [15] S. Mahabadi. Approximate nearest line search in high dimensions. In *Proc. 26th Annu. ACM-SIAM Sympos. Discrete Algorithms (SODA)*, page to appear, 2015.
- [16] J. Matoušek. Efficient partition trees. *Discrete Comput. Geom.*, 8(3):315–334, 1992.
- [17] J. Matoušek. *Lectures on discrete geometry*. Springer-Verlag, 2002.
- [18] S. Meiser. Point location in arrangements of hyperplanes. *Inform. and Comput.*, 106(2):286–303, 1993.
- [19] R. Panigrahy. Entropy based nearest neighbor search in high dimensions. In *Proc. 17th Annu. ACM-SIAM Sympos. Discrete Algorithms (SODA)*, pages 1186–1195, 2006.