

SkipTree: A Scalable Range-Queryable Distributed Data Structure for Multidimensional Data

Saeed Alaei, Mohammad Toossi

Sharif University of Technology

alaei@ce.sharif.edu, mohammad@bamdad.org

Abstract

This paper presents the SkipTree, a new balanced, distributed data structure for storing data with multidimensional keys in a peer-to-peer network. The SkipTree supports range queries as well as single point queries which are routed in $O(\log n)$ hops. SkipTree is fully decentralized with each node being connected to $O(\log n)$ other nodes. The memory usage for maintaining the links at each node is $O(\log n \log \log n)$ on average and $O(\log^2 n)$ in the worst case. Load balance is also guaranteed to be within a constant factor.

1 Introduction and Related Work

Over the past few years, there has been a trend to move from centralized server based network architectures toward decentralized and distributed architectures and peer to peer networks. The term *Scalable Distributed Data Structure* (SDDS) first introduced by Litwin et al. in LH* [15] refers to this class of data structures which hold the following properties:

- There is no central directory.
- Client images (i.e. client information on where data is located) may be outdated, and is only adjusted in response to read queries.
- A client may send a request to the incorrect server, which will be forwarded to the correct server and the client image will be updated.

Litwin et al. modified the original hash-based LH* [15] structure to support range queries in RP* [14, 15]. Based on the previous work of distributed data structures like LH* [15], RP* [14] and Distributed Random Tree (DRT) [11], new data structures based on either hashing or key comparison have been proposed like Chord [20], Viceroy [16], Koorde [9], tapestry [22], Pastry [19], PeerDB [18] and P-Grid [3]. Most existing peer-to-peer overlays require $\Theta(\log n)$ links per node in order to achieve $O(\log n)$ hops for routing. Viceroy [16] and Koorde [9] which are based on DHTs are the remarkable exceptions in that they achieve $O(\log n)$ hops with only $O(1)$ links per node at the cost of restricted or no load balancing. Family Tree [21] is the first overlay network which does not use hashing but supports routing in $O(\log n)$ hops with only $O(1)$ links per node.

Typically, those systems which are based on DHTs and hashing lack range-query, locality properties and control over distribution of keys due to hashing. In contrast, those which are based

on key comparison, although requiring more complicated load balancing techniques, do better in those respects. P-Grid [3] by Aberer et al. is one of the systems based on key comparison which uses a distributed binary tree to partition a single dimensional space with network nodes representing the leaves of the tree and each node having a link to some node in every sibling subtree along the path from the root to that node. Gridella [4] a P2P system based on P-Grid working on Gnutella has also been developed. Other systems like P-Tree [5] have been proposed that provide range queries in single dimensional space. Besides, Some data structures like dE-Trees [8] based on B-Trees have been developed for distributed environments.

SkipNet [7] on which our new system relies heavily, is another system for single dimensional spaces based on an extension to skip lists.

G-Grid [6] is a solution proposed for the multidimensional case which is also based on partitioning the space into regions. However, regions in G-Grid are restricted in that they can only be split to two regions of equal size. So, their boundaries cannot take arbitrary values and are restricted to multiples of their size. Their size are also restricted to negative powers of 2.

RAQ [17] is also another solution for the multidimensional case which incorporates a distributed partition tree structure to partition the space. Its network model is similar to that of the P-Grid [3], Therefore it requires $O(h)$ links at each node and routes in $O(h)$ hops where h is the height of the partition tree which can be of $O(n)$ for an unbalanced tree. Although it has been shown [2, 1] that even for such unbalanced trees the number of messages required to resolve a query still remains of $O(\log n)$ on average if the links are chosen randomly, the number of links that a node should maintain and the memory requirement at each node for storing information about the path from that node to the root still remain of $O(h)$ which is as bad as $O(n)$ for unbalanced trees.

In this paper we propose a new efficient scalable distributed data structure called the *SkipTree* for storage of keys in multidimensional spaces. Our system uses a distributed partition tree to partition the space into smaller regions with each network node being a leaf node of that tree and responsible for one of the regions. In contrast to similar tree-based solutions the partition tree here is used only to define an ordering between the regions. The routing mechanism and link maintenance is similar to that of SkipNet and independent of the shape of the partition tree, so in general our system does not need to balance the partition

tree (in fact, it has been shown [12] that such a tree cannot be balanced efficiently by means of rotation). Our system, maintains a SkipNet by the leaves of the tree in which the sequence of nodes in the SkipNet is the same sequence defined by the leaves of the partition tree from left to right. Handling a single key query is almost similar to that of an ordinary SkipNet while range queries are quite different due to the multidimensional nature of the SkipTree. From another point of view, our system can be seen as an extension to the SkipNet for the multidimensional spaces.

In section 2 we explain the basic structure of the SkipTree including the structure of the partition tree, its associated SkipNet and the additional information that needs to be stored in each node. In section 3, single and range queries are explained. In section 4, the procedure for joining and leaving the network is described. In section 5, some techniques for load balancing in SkipTrees are discussed. In section 6 we modify the SkipTree structure to reduce the amount of information that needs to be stored in each node about the partition tree and finally section 7 concludes the paper.

2 Basic SkipTree Structure

The distributed data structure used in the SkipTree consists of two parts. First, a *Partition Tree* is used to divide the search space among the nodes. This is described in subsection 2.1. Then, as is shown in subsection 2.2, nodes are linked together using a technique similar to SkipNet.

2.1 Space Partitioning

We assume that each data element has a key which is a point in our k -dimensional search space. This space is split into n regions corresponding to the n network nodes. Let $S(v)$ denote the region assigned to node v . v is the node responsible for every data element whose key is in $S(v)$.

We use *Partition Tree*, a binary tree, to perform this assignment. Although only leaves in the partition tree represent actual nodes in our overlay network, each node in this tree has a corresponding section in the search space. Thus, we extend the definition of $S(v)$ to also denote the region assigned to an internal node v in this tree.

Assuming r is the root of our Partition Tree, $S(r)$ is always the whole search space. Each internal node then recursively splits its region into two smaller regions using a hyperplane equation. That is, if an internal node v has two children, l and r , which are its left and right children respectively, $S(l)$ will be the portion of $S(v)$ left on one side of the hyperplane specified by v and $S(r)$ will be the space to the other side. A sample partition tree and its corresponding space partitioning are depicted in Figure 1.

For network node u , which corresponds to a leaf in the partition tree, we call the path connecting the root of the tree to u the *Principal Path* of node u . We refer to the hyperplane equations assigned to the nodes of a principal path of a node u (including information about on which side of those hyperplanes u resides) as the *Characteristic Plane Equations* of u or *CPE* of u for short. Every leaf node in the SkipTree stores its own CPE as well as the CPE of its links. Using these CPE information, every node like

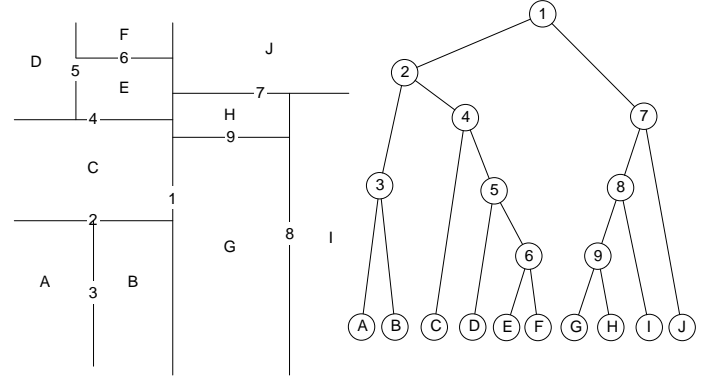


Figure 1: A sample two dimensional partition tree and its corresponding space partitioning. Each internal node in the partition tree, labelled with a number, divides a region using the line labelled with the same number. Each leaf of the partition tree is a network node responsible for the region labelled with the same letter.

u can locally identify if a given point belongs to a node to the left or the right of u or to the left or right of any of its links in the partition tree. The latter is useful in routing queries as explained in section 3. Hereafter, whenever we refer to a plane, we actually refer to a hyperplane of $k - 1$ dimensions in a k -dimensional space.

Storing the above CPEs, however, requires $O(\log nh)$ memory at every leaf node, where h is the height of the tree. While this is of $O(\log^2 n)$ for a balanced tree, it may require as much as $O(n)$ memory in the worst case. We will provide a method for reducing the memory requirement in section 6.

2.2 Network Links

We link the network nodes in the SkipTree together by forming a SkipNet among the leaves of the partition tree described in the previous subsection. However, using the SkipNet requires a total ordering to be defined among the nodes. We define this ordering to be the order in which the nodes appear as the leaves of the partition tree from left to right. We also make this sequence circular by considering the rightmost leaf of the tree as the node to the left of the leftmost leaf and vice versa.

In an SkipTree in its ideal form, a node v keeps $2 \log_2 n - 1$ links to other nodes. These are the 2^i th nodes to the left and right of v for every i from 0 to $\log_2 n$ as shown in Figure 2. Unfortunately, maintaining this form is very inefficient when handling node arrivals and departures. As a result, only an approximation to these ideal links is maintained in SkipNet.

For a given i , if we start from any node and follow the links that jump 2^i nodes in a specific direction in the ideal form, we will find a loop of length $n/2^i$. Let's call this loop a level i ring. There are 2^i level i rings. For example, there is only one level 0 ring, a circular doubly-linked list that connects every node in the aforementioned order. On the other hand, there are n last level rings consisting only of individual nodes. In general, there are 2^i level i rings. As illustrated in Figure 3, the nodes in each

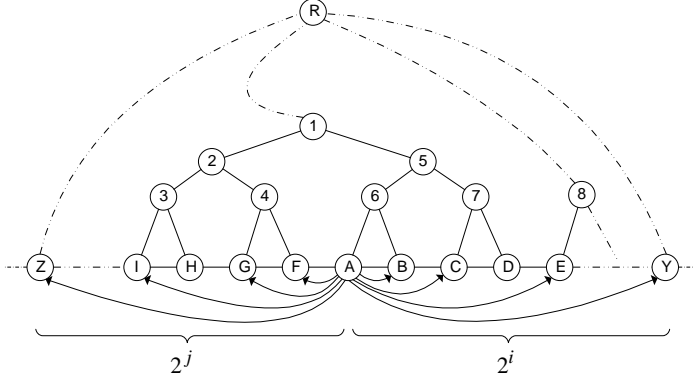


Figure 2: The links maintained by node A in the ideal SkipTree. The target nodes are independent of the tree structure. The tree only helps us to put an ordering on the nodes. The i^{th} link in each direction skips over $2^{i-1} - 1$ nodes in that direction.

level i ring form exactly two disjoint level $i + 1$ rings. This is the property which will be conserved when nodes are inserted into or deleted from the SkipTree in section 4.

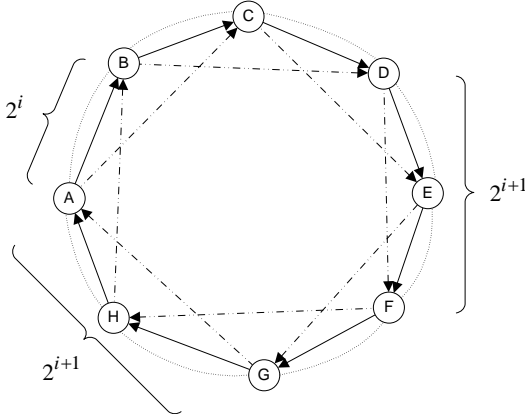


Figure 3: The nodes in each level i ring are split between two level $i + 1$ rings. Solid arrows, representing level i links, form the level i ring. Dashed arrows are the next level links that form two disjoint level $i + 1$ rings.

Finally, we note that a real number p_v is assigned to each node v . p_v is randomly generated when v joins the SkipTree so that $p_a < p_v < p_b$ where a and b are v 's predecessor and successor in the total ordering. This number is used in subsection 3.2 to handle range queries more efficiently.

3 Handling Queries

Queries in a SkipTree can take two forms, either a single point query or a range query. We will discuss them separately on the following subsections.

3.1 Single Point Query

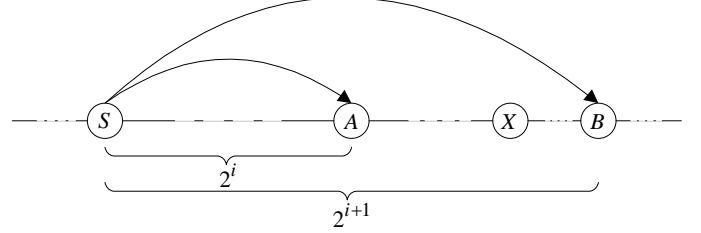


Figure 4: A point query is routed through the farthest link which does not point past the destination node. Here, S receives a query targeting node X , so it routes the query to A . The distance to the destination node is at least halved at each hop.

Whenever a node in the network receives a single point query, it must route the query to the node which is responsible for the region containing the point. The routing algorithm is essentially the same algorithm used in the SkipNet, that is every node receiving the query along the path, sends it through its farthest link which does not point past the destination node. This is shown in Figure 4 where node S is about to route a single point query to some unknown node X which lies somewhere between node A and node B . Here A and B are two consecutive nodes in the list of nodes to which A has a direct link and are at distance 2^i and 2^{i+1} from A for some i respectively. Node S routes the query to A and then A routes the query again in the same way until the query reaches its destination node. Note that the distance from A to the destination node is less than 2^i , so the next hop is at most 2^{i-1} nodes away from A . In fact, the distance to the destination node is at least halved at each hop. This implies that the query reaches the destination after at most $\log_2 n$ hops. However, because SkipNet uses a probabilistic method for selecting and maintaining links in the network, it guarantees routing in $O(\log n)$ hops w.h.p.¹. A formal proof of this can be found in [7].

For the above procedure to be effective, given a point, every node must be able to identify whether the node which is responsible for that point lies to its left or to its right side. In other words, we must be able to compare points against nodes to identify whether the node containing a given point lies before or after another node in the sequence. This is where the tree structure helps us. To do so, a node compares the point against the planes in its own CPE in the order they appear in its principal path starting from the root until it finds the first plane where the current node and the point lie on different sides of the plane. This is where the point is contained in a region belonging to a sibling

¹An event is said to be occurring with high probability (w.h.p) if for any constant value of α the event occurs with a probability of at least $1 - O(\frac{1}{n^\alpha})$.

subtree. If that subtree is a left (right) subtree, all of its nodes as well as the node containing the point must also be to the left (right) of the current node. That is why every node in the network must also store the CPE of its link nodes in addition to its own CPE to be able to compare queries against its links too. The above procedure leads to $O(\min(h \log n, n))$ memory usage at each node for storing the CPE, where h is the height of the tree. This may be as bad as $O(n)$ memory for an unbalanced tree. We will modify the tree structure in [section 6](#) to overcome this problem and guarantee $O(\log h \log n)$ memory usage at each node for the storage of CPE, which means $O(\log n \log \log n)$ on average and $O(\log^2 n)$ memory usage in the worst case for an unbalanced tree.

3.2 Range Query

A range query in the SkipTree is a 3-tuple of the form (R, fs, ls) where R is the query range and fs and ls are two real numbers which define the range of nodes in the sequence of nodes to be searched. That is, only the network nodes whose sequence numbers reside in the interval $[fs, ls]$ are searched. Using this form of queries one can perform a complete range query for a region R using the 3-tuple $(R, -\infty, +\infty)$, so that all of the nodes are included in the search regardless of their sequence number. Note that the region defined by R can be of any shape as long as every node can locally identify whether R intersects with a given hypercube.

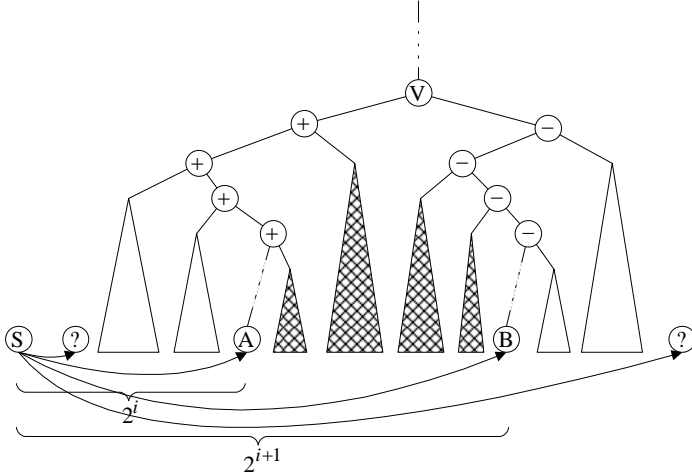


Figure 5: A range query is propagated through each of the links maintained by S whenever there is node which intersects R between that link and the next link. Here, S has consecutive links to A and B . A copy of the query is propagated to A if any of the nodes between A and B intersects with R .

Handling a range query is very similar to that of a single point query with some minor differences. Suppose that a node S receives a range query (R, fs, ls) . To handle this query then S breaks the range query to several queries (at most $O(\log n)$ new queries) with each targeting fewer number of nodes relative to the original query. A range query is propagated to each of the links maintained by S whenever there is node which intersects with the region R between that link and the next link. Assume that A and B shown in [Figure 5](#) are two nodes corresponding to

some two consecutive links maintained by S . S sends a copy of the query to A if there is any node between A and B which intersects R . Every such node, if any, must reside in one of the crosshatched subtrees illustrated in the figure. In fact, such a node must be to the right of the nodes marked with $+$ and to the left of the node marked with $-$ and because S has all of CPEs corresponding to its links, it also has the plane equations corresponding to the internal nodes marked with a $+$ or $-$ sign. So, it can easily identify from those equations the regions in the multi-dimensional space associated with each of the subtrees between A and B and from that it can determine whether there is any subtree between A and B whose region intersects with R and if there is such a subtree, it must also contain a node whose region intersects with R . In this way a query is broken by S to several queries and is propagated until it reaches its targets. Note that the fs and ls fields of the query are modified appropriately before a copy of the query is sent through a link. The reason is to restrict the sequence of nodes to be searched to prevent duplicate queries. For example in [Figure 5](#), suppose that a copy of the form (R, fs, ls) is to be sent from S to A . Also assume that $A.seq$ and $B.seq$ are the sequence number of A and B respectively. Then S computes the interval $[fs', ls']$ as the intersection of $[fs, ls]$ and $[A.seq, B.seq]$ and it sends the query (R, fs', ls') to A . This will ensure that no nodes in the network receives the query more than once.

For the above procedure note that the length of the path that a query travels through is of $O(\log n)$ regardless of the width of propagation at each hop. The proof is basically the same as the case with the single point query.

It is worth mentioning that this is the only place where sequence numbers are actually used. Sequence numbers make it possible to determine the relative ordering of two nodes in the network without knowing their corresponding plane equations or the regions they represent.

4 Node Join and Departure

Join and Departure operations are described in the following subsections. For each operation the node has to perform two relatively independent actions. Update the Partition Tree and update the network connections.

4.1 Joins

To join the SkipTree, a new node v has to be able to contact an existing node u in the SkipTree. v then splits the space assigned to u using a new plane. This allows u to transfer the control of one of the new regions along with its stored data items to v .

The algorithm is shown in function $v.join(u)$ of [Figure 6](#). In lines 1 and 2, the new node copies the CPE of the existing node. Then, a new plane equation is generated to split the region formerly assigned to u . This plane can be arbitrarily chosen as our load balancing protocol will gradually change the partitioning to a more balanced configuration. Each of the two nodes then selects one of the two newly created regions. This is done in lines 3-5 by extending the principal paths using the `add_to_cpe()` function.

The provided algorithm inserts v immediately after u and chooses the side of the plane containing the origin

(ORIGIN_SIDE) for u and the other side (OTHER_SIDE) for v . This will help us in performing memory optimization in [section 6](#).

```

// Join new node v to the SkipTree.
// u is an arbitrary node in SkipTree.
v.join(u){
1   v.cpe ← u.cpe;
2   v.height ← u.height;
3   new_plane ← choose_plane(v.cpe, v.height);
4   u.add_to_cpe(new_plane, ORIGIN_SIDE);
5   v.add_to_cpe(new_plane, OTHER_SIDE);
6   v.seq ← random number in
   (u.sequence, u.successor.sequence);
7   v.join_SkipNet(u);
8   u.transfer_data(v);
}

// Append a plane to the CPE.
v.add_to_cpe(plane, side){
   v.cpe[v.height].plane ← plane;
   v.cpe[v.height].side ← side;
   v.height ← v.height + 1;
}

```

Figure 6: Joining the SkipTree

After updating the Partition Tree, v has to establish its network connections. This is done by inserting v into approximately $\log n$ rings mentioned in [subsection 2.2](#). Starting with the level 0 ring, v randomly selects one of the two level $i + 1$ rings derived from the selected level i ring until it reaches a ring in which there is no node except the new node itself. v then moves backward, inserting itself in each of these rings with regard to the total ordering defined in [subsection 2.2](#) which is the same as using p_v . The exact algorithm is described in [7] and involves only $O(\log n)$ steps w.h.p.

Finally, u transfers the data items which are no longer in its assigned region to v in line 8.

4.2 Departures

When node v is leaving the SkipTree, it has to follow three steps. First, update the Partition Tree; second, transfer its data items to the appropriate nodes; and third, leave the SkipNet.

Suppose that v is responsible for region R and that the nodes in its sibling subtree are collectively responsible for the region S . In other words, the last plane in the node v 's CPE, called P , splits its parent's region into regions R and S . To update the Partition Tree, node v sends a special range query to the nodes in region S and instructs them to remove the plane P from their CPE. This will effectively remove v from the partition tree and shift every node in S one level closer to the root by removing their common parent.

To transfer the data items, v can simply find the node responsible for each item using a single point query and transfer the item accordingly. However, a more efficient method is to create a collection of the regions associated with every possible target

node for v 's data items and perform the single point queries for these items locally. This collection can be created by asking every node (as part of the previous special range query) in S to send its newly associated region to v if this region intersects with R .

In the last step, v has to close its $O(\log n)$ connections. As [7] points out, all pointers except the ones forming the level 0 ring can be regarded as redundant routing optimization hints and can be updated using a background repair process similar to Chord and Pastry. Therefore, v only needs to cleanly remove itself from the level 0 ring before leaving the SkipTree.

5 Load Balancing

Many distributed lookup protocols use hashing to distribute keys uniformly in the search space and achieve some degree of load balance. Hashing cannot be used in the SkipTree as it makes range queries impossible. As a result, a load balancing mechanism is necessary to deal with the nonuniform key distribution.

Our load balancing protocol is derived from the *Item Balancing* technique in [10]. Load balancing is achieved using a randomized algorithm that requires a node to be able to contact random nodes in the network. This can be implemented either using the existing network connections in SkipNet or using the underlying peer-to-peer routing framework. The second approach is preferred because of its higher speed and lower network traffic.

Let l_i , the load on node i , be the number of data items stored on i and α be a constant number so that $\alpha > 1$. We will prove that the SkipTree's load will be balanced w.h.p. if each node performs a minimum number of *load balancing tests* as per system *half-life*².

Load Balancing Test In a load balancing test, node i asks a randomly chosen node j for l_j . If $l_j \geq \alpha l_i$ or $l_i \geq \alpha l_j$, i performs a *load balancing operation*.

Load Balancing Operation Assume w.l.o.g that $l_i < l_j$. First, node i normally leaves the SkipTree using the algorithm given in [subsection 4.2](#). Then, i joins the network once again at node j and selects a hyperplane for the newly created internal node in the partition tree in a way that the number of data elements is halved at both sides of the hyperplane. This makes both l_i and l_j to become equal to half the old value of l_j .

Theorem 1 *If each node performs $\Omega(\log n)$ load balancing operations per half-life as well as whenever its own load doubles, then the above protocol has the following properties where N is the total number of stored data items.*

- With high probability, the load of all nodes is between $\frac{N}{8\alpha n}$ and $\frac{16\alpha N}{n}$.
- The amortized number of items moved due to load balancing is $O(1)$ per insertion or deletion, and $O(N/n)$ per node insertion or deletion.

The proof of this theorem using potential functions can be found in [10].

²A half-life is the time it takes for half the nodes or half the items in the system to arrive or depart. [13]

6 Memory Optimization

Throughout the previous sections we assumed that every node in the network must store the CPE of all of the nodes to which it maintains a link as well as its own CPE. As we mentioned earlier, in a SkipTree of height h , this requires $O(h \log n)$ memory for each node to store its own CPE as well as CPEs of its link nodes. So, a node may require $O(n)$ memory with an unbalanced SkipTree in the worst case. In this section we enforce some constraints on the plane equations that a node may choose when joining the network and splitting another node, so that for a SkipTree of height h only $O(\log h)$ of the plane equations of any CPE will be needed for the correct operation of the SkipTree. The constraints that we enforce are the following:

- The planes must be perpendicular to a principal axis. So, in a k -dimensional space of (x_1, x_2, \dots, x_k) it must take the form of $x_i = c$ for some $1 \leq i \leq k$ and some value of c . This effectively means that every such plane partitions the keys in the space based on the value of x_i for some i . We put further constraints on how such a plane like $x_i = c$ partitions a region into two smaller region by requiring that every such plane partitions the region associated with an internal node of the SkipTree like U in such a way that the region containing the points with $x_i \leq c$ is assigned to the left subtree of U and the region containing the points with $c < x_i$ is assigned to the right subtree of U .
- If the search space is k -dimensional, we precisely define the form of the plane equation that may be assigned to an internal node depending on the depth of that node. We first introduce the following notation:

d_A : for a node A in the SkipTree, the depth of A is represented by d_A and is defined to be the length of the principal path corresponding to A plus one. For an example see the SkipTree of Figure 7.

l_A : for every node A in the SkipTree, the level of A is indicated by l_A where $l_A = \lceil \log_2 (\frac{d_A}{k} + 1) \rceil$. This means that all of the nodes whose depth are in the interval $[k(2^i - 1) + 1, k(2^{i+1} - 1)]$ belong to level $i + 1$. This implies that on any principal path, the first k nodes are in level 1, the next $2k$ nodes are in level 2, the next $4k$ are on the next level and so on. For an example see the SkipTree of Figure 7.

d'_A : for a node A in the SkipTree, the relative depth of A is represented by d'_A and is defined so that $d'_A = d_A - d_B + 1$ where B is the highest node which has the same level as A , or alternatively we can define $d'_A = d_A - k(2^{l_A-1} - 1)$. For an example see the SkipTree of Figure 7.

s_A : for a node A in the SkipTree, the section number of A is represented by s_A where $s_A = \lceil \frac{d'_A}{k} \rceil$. In fact, nodes at every level are partitioned to k sections. This implies that on the i -th level of any principal path, the first 2^{i-1} nodes have section number 1, the next 2^{i-1} nodes are in section 2 and so on.

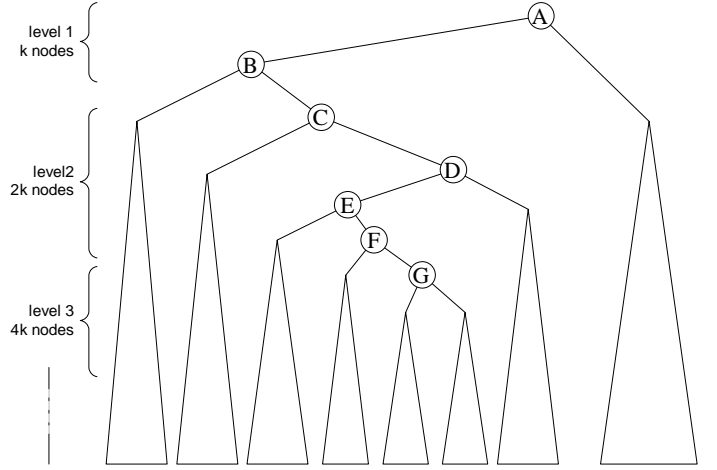


Figure 7: A sample SkipTree for a two dimensional space. Nodes A to G have depths 1 to 7 respectively. A and B are on level 1; C , D , E and F are on level 2 and G is on level 3. The relative depth are: $d'_A = 1$, $d'_B = 2$, $d'_C = 1$, $d'_D = 2$, $d'_E = 3$, $d'_F = 4$, $d'_G = 1$.

We are now ready to state the last constraint:

If A is an internal node, the plane equation assigned to A must be of the form $x_{s_A} = c$ for an arbitrary value of c , that is for any given i , all of the nodes whose section numbers are i are assigned plane equations of the form $x_i = c$. This implies that whenever a new node joins the SkipTree and splits the region of another node which leads to a new internal node, the plane equation of that internal node must obey the above schema. So the only parameter that the new node can define to balance the load when it splits a region is the value of the constant c which should be enough for that purpose. A typical 2-dimensional space partitioned under the above constraints and its associated tree are shown in Figure 8.

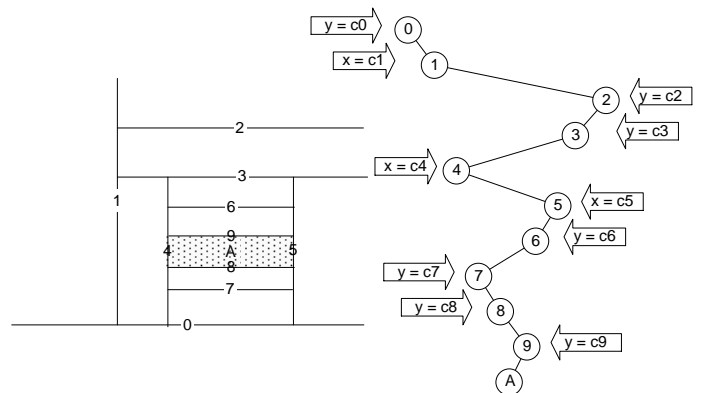


Figure 8: The left is a sample partitioning of a 2-dimensional space under the memory optimization constraints, from the view point of node A and the right is the principal path of node A . The plane equations assigned to the internal nodes are shown in the arrows.

Lemma 1 *In any principal path of length h nodes are partitioned to at most $k \lceil \log_2(\frac{h}{k} + 1) \rceil$ different sections.*

Proof: Since we defined the level of a node at depth d to be $\lceil \log_2(\frac{d}{k} + 1) \rceil$, nodes in any principal cannot be partitioned to more than $\lceil \log_2(\frac{h}{k} + 1) \rceil$ levels. Nodes at each level are further partitioned to k sections so there can be at most $k \lceil \log_2(\frac{h}{k} + 1) \rceil$ sections in any principal path.

Lemma 2 *For any leaf node A in a SkipTree, A needs to store only two plane equations for each section of its principal path. we call the sequence of these pairs of plane equations that node A stores, the Reduced-Characteristic Plane Equations of node A or for short the RCPE of node A .*

Proof: All of the planes on the same section partition the space based on the value of the same field x_i . For example in Figure 8, in the section 1 in the 3rd level of principal path of A , all of the internal nodes are assigned a plane equation of the form $y = c$ for different values of c and the region associated with node A or any other leaf for that matter is between at most two of the planes of that section. That is, for every section in the principal path for any node A , there are at most two planes which best represent the region in which A 's region lies. So for each of the sections, A needs to store an inequality of the form $a \leq x_i < b$. Therefore an RCPE can be stored as an ordered sequence of inequalities of the form $a \leq x_i < b$, one for each section in the principal path. When a node like A receives a point query it finds the first inequality in the RCPE sequence that does not hold for the queried point. Then the first constraint we introduced on the beginning of this section ensures that the destination node which is responsible for the queried point will be to the left of the current node if the point is to the left of the interval represented by the first unsatisfied inequality and the destination node will be to the right of the current node otherwise. The situation with range queries is quite similar. The sequence of inequalities in the RCPE for the node A in Figure 8 is shown below:

- level=1, section=1 : $c_0 \leq y < +\infty$
- level=1, section=2 : $c_1 \leq x < +\infty$
- level=2, section=1 : $c_0 \leq y < c_3$
- level=2, section=2 : $c_4 \leq x < c_5$
- level=3, section=1 : $c_8 \leq y < c_9$

6.1 Node Join and Departure

Joining mechanism is the same as before except that a new node must obey the constraints mentioned earlier. However leaving is a bit tricky since when some node A is about to leave, it must then remove the internal node which is its direct parent which causes the plane associated with that internal node to be removed from the RCPE of the nodes in its sibling subtree as well. This becomes a problem only when that plane belongs to a level which is not the last level in the principal path of some node B in the sibling subtree of A because removing it will contradict the second constraint of memory optimization since it is

removed from a level which is not the last level for some principal path. To overcome this problem A sends a special form of range query containing the region associated with its sibling subtree and by this special query it tries to find some node X in its sibling subtree such that the sibling subtree of X does not contain a node whose level is greater than that of X , then node A and X swap their roles, that is they swap their associated regions as well as the keys that they store and their position in the SkipTree. After this swapping, A will be in place of the X in the SkipTree so it can then leave the network using the procedure described in the previous sections and this time the above problem will not occur for X because of the way the node X was chosen. It is easy to see that such a node X must always exist. For example the lowest node in the sibling subtree of A always has the desired property although it is not the only node with such property.

6.2 Complexity

The memory requirement of any node A for storing its RCPE as well as the RCPE of its links as described earlier is of $O(\log h \log n)$ where h is the height of the tree which is a major improvement over the $O(\min(h \log n, n))$ memory requirement in the default case.

In addition to the memory requirement guarantee, the constraints that we enforced in section 6, guarantee the following strong invariant on the distribution of planes in each direction for every principal path:

Theorem 2 *For every principal path in a SkipTree if m_i is the number of plane equations of the form $x_i = c$ and m_j is the number of plane equations of the form $x_j = c$ for possibly different values of c , then the inequality $x_i \leq 2x_j + 1$ must always hold.*

Proof: For every principal path in a SkipTree, there are equal number of plane in each direction at each level except possibly for the last level (the level with highest number, that is the level of the lowest part of the path), because every level except the last level consists of exactly k different sections of equal size with all of the plane of each section being in the same direction. Besides, the number of planes in a single section at the i^{th} level is 2^{i-1} , so if a principal path consists of r levels, for each direction, the total number of planes in that direction at all levels except the last level is $2^0 + 2^1 + 2^2 + \dots + 2^{l-2}$ which is $2^{l-1} - 1$. Also, for each direction, the number of planes in that direction at the last level is between 0 to 2^{l-1} . So, in any principal path, for each direction, the total number of planes in that direction at all levels is between $2^{l-1} - 1$ to $2^l - 1$ and the above inequality obviously results.

From Theorem 2, it is implied that the planes are distributed almost uniformly in each direction which is an advantage over the default case where the plane equations are chosen randomly.

7 Conclusion and Future Work

In this paper we introduced the *SkipTree* which is designed to handle point and range queries over a multidimensional space in a distributed environment. Our data structure maintains $O(\log n)$ links at each node and guarantees an upper bound of

$O(\log n)$ messages w.h.p for point queries and also guarantees range queries with depth of $O(\log n)$ message w.h.p. It improves the previously proposed data structures for multidimensional space which were based on binary trees in the following aspects:

- **Links:** every node in a SkipTree needs to keep track of $O(\log n)$ links regardless of the shape of the tree in contrast to other tree based structures where each node should keep track of $O(h)$ links, where h (the height of the tree) can be of $O(n)$ for an unbalanced tree.
- **Query depth:** the maximum depth of a point and range query in a SkipTree is of $O(\log n)$ regardless of the shape of the tree, in contrast to other tree based structures where a query may travel $O(h)$ hops in the worst case where h can be of $O(n)$ for an unbalanced tree.
- **Memory requirement:** using the memory optimization of [section 6](#), each node needs only to store the RCPE of itself and its links that requires $O(\log h \log n)$ which is quite an improvement over memory requirement of similar tree based structures where each node maintains information for every node along its principal path which requires $O(h)$ memory that can be as bad as $O(n)$ for unbalanced trees.

In addition to the above improvements we also adapted some load balancing techniques to improve our data structure. However it seems that the load balancing procedure and the memory optimization technique may be conflicting. In fact in some situation the node swapping method described in [subsection 6.1](#) may cancel out the effect of the load balancing method. This is one important area which needs further investigation. Another important area which needs further improvement is about the fault tolerance of the structure in presence of node failures. Also, as mentioned above load balancing and memory optimization need more improvements.

References

- [1] K. Aberer. Efficient search in unbalanced, randomized peer-to-peer search trees. Technical Report IIC/2002/79, EPFL, 2002.
- [2] K. Aberer. Scalable data access in p2p systems using unbalanced search trees. In *Proceedings of Workshop on Distributed Data and Structures (WDAS-2002)*, 2002.
- [3] K. Aberer, P. Cudr-Mauroux, A. Datta, Z. Despotovic, M. Hauswirth, M. Puceva, and R. Schmidt. P-grid: a self-organizing structured p2p system. *SIGMOD Rec.*, 32(3):29–33, 2003.
- [4] K. Aberer, M. Puceva, M. Hauswirth, and R. Schmidt. Improving data access in p2p systems. *IEEE Internet Computing*, 6(1):58–67, 2002.
- [5] A. Crainiceanu, P. Linga, J. Gehrke, and J. Shanmugasundaram. Querying peer-to-peer networks using p-trees. In *WebDB '04: Proceedings of the 7th International Workshop on the Web and Databases*, pages 25–30. ACM Press, 2004.
- [6] D. S. For. G-grid: A class of scalable and self-organizing data structures for multi-dimensional querying and content routing in p2p networks.
- [7] N. HARVEY, M. JONES, S. SAROIU, M. THEIMER, and A. WOLMAN. Skipnet: A scalable overlay network with practical locality properties, 2003.
- [8] T. Johnson and A. Colbrook. A distributed data-balanced dictionary based on the b-link tree. In *Proceedings of the 6th International Parallel Processing Symposium*, pages 319–324. IEEE Computer Society, 1992.
- [9] M. Kaashoek and D. Karger. Koorde: A simple degree-optimal distributed hash table, 2003.
- [10] D. R. Karger and M. Ruhl. Simple efficient load balancing algorithms for peer-to-peer systems. In *SPAA '04: Proceedings of the sixteenth annual ACM symposium on Parallelism in algorithms and architectures*, pages 36–43. ACM Press, 2004.
- [11] B. Kroll and P. Widmayer. Distributing a search tree among a growing number of processors. In *SIGMOD '94: Proceedings of the 1994 ACM SIGMOD international conference on Management of data*, pages 265–276. ACM Press, 1994.
- [12] B. Kroll and P. Widmayer. Balanced distributed search trees do not exist. In *WADS '95: Proceedings of the 4th International Workshop on Algorithms and Data Structures*, pages 50–61. Springer-Verlag, 1995.
- [13] D. Liben-Nowell, H. Balakrishnan, and D. Karger. Analysis of the evolution of peer-to-peer systems. In *PODC '02: Proceedings of the twenty-first annual symposium on Principles of distributed computing*, pages 233–242. ACM Press, 2002.
- [14] W. Litwin, M.-A. Neimat, and D. A. Schneider. Rp*: A family of order preserving scalable distributed data structures. In J. B. Bocca, M. Jarke, and C. Zaniolo, editors, *VLDB'94, Proceedings of 20th International Conference on Very Large Data Bases*, sep 1994.
- [15] W. Litwin, M.-A. Neimat, and D. A. Schneider. Lh* – a scalable, distributed data structure. *ACM Trans. Database Syst.*, 21(4):480–525, 1996.
- [16] D. Malkhi, M. Naor, and D. Ratajczak. Viceroy: a scalable and dynamic emulation of the butterfly. In *PODC '02: Proceedings of the twenty-first annual symposium on Principles of distributed computing*, pages 183–192. ACM Press, 2002.
- [17] H. Nazerzadeh. RAQ: A range-queriable distributed data structure (extended version). In *Proceeding of Sofsem 2005, 31st Annual Conference on Current Trends in Theory and Practice of Informatics, LNCS 3381*, pp. 264–272, February 2005.
- [18] W. S. Ng. Peerdb: A p2p-based system for distributed data sharing, 2003.
- [19] A. I. T. Rowstron and P. Druschel. Pastry: Scalable, decentralized object location, and routing for large-scale peer-to-peer systems. In *Middleware 2001: Proceedings of the IFIP/ACM International Conference on Distributed Systems Platforms Heidelberg*, pages 329–350. Springer-Verlag, 2001.
- [20] I. Stoica, R. Morris, D. Karger, F. Kaashoek, and H. Balakrishnan. Chord: A scalable Peer-To-Peer lookup service for internet applications. pages 149–160.
- [21] K. C. Zatloukal and N. J. A. Harvey. Family trees: an ordered dictionary with optimal congestion, locality, degree, and search time. In *SODA '04: Proceedings of the fifteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 308–317. Society for Industrial and Applied Mathematics, 2004.
- [22] B. Y. Zhao, L. Huang, J. Stribling, S. C. Rhea, A. D. Joseph, and J. D. Kubiatowicz. Tapestry: A global-scale overlay for rapid service deployment. *IEEE Journal on Selected Areas in Communications*, 2003. Special Issue on Service Overlay Networks, to appear.