

Structured Overlay for Multi-dimensional Range Queries^{*}

Thorsten Schütt, Florian Schintke, and Alexander Reinefeld

Zuse Institute Berlin

Abstract. We introduce SONAR, a structured overlay to store and retrieve objects addressed by multi-dimensional names (*keys*). The overlay has the shape of a multi-dimensional torus, where each node is responsible for a contiguous part of the data space. A uniform distribution of nodes on the data space is not necessary, because denser areas get assigned more nodes. To nevertheless support logarithmic routing, SONAR maintains, per dimension, fingers to other nodes, that span an exponentially increasing number of *nodes*. Most other overlays maintain such fingers in the *key-space* instead and therefore require a uniform data distribution. SONAR, in contrast, avoids hashing and is therefore able to perform range queries of arbitrary shape in a logarithmic number of routing hops—independent of the number of system- and query-dimensions. SONAR needs just one hop for updating an entry in its routing table: a longer finger is calculated by querying the node referred to by the next shorter finger for its shorter finger. This doubles the number of spanned nodes and leads to exponentially spaced fingers.

Introduction

Efficient handling of multi-dimensional range queries in Internet-scale distributed systems is still an open issue. Several approaches exist, but their lookup costs are either expensive (space-filling curves) [2] or use probabilistic space-like consistent hashing [10] to build the overlay.

Imagine a system for storing and retrieving objects with d -dimensional keys in a peer-to-peer network. SONAR (Structured Overlay Network with Multi-dimensional Range-queries) directly maps the multi-dimensional data space to a multi-dimensional torus. It supports range queries of arbitrary shape, which are useful, for example, in geo-information systems where objects in a given distance of a point are sought. SONAR can also be employed in Internet games with millions of online-players who concurrently interact in a virtual space and need quick access to the local surroundings of their avatars. In a broader context, SONAR can be employed as a hierarchical publish/subscribe system, where published events are categorized by several independent attributes. The category of published events addresses a data point in the d -dimensional space and consumers belonging to subareas will receive all events published in their subarea.

. Schütt, F. Schintke, and A. Reinefeld

paper is organized as follows: First, we discuss related work. Then, in Section 3, we introduce SONAR. In Section 4, we present empirical results and in Section 5 we conclude the paper with a brief summary.

Related Work

Systems [1] have been proposed that support complex queries with multi-attribute keys and ranges. They can be split into two groups.

Filling Curves. These systems [2,9,16] use locality preserving space-filling curves to map multi-dimensional to one-dimensional keys. They provide efficient range queries than the space partitioning schemes described below, as a single range query may cover several parts of the curve, which have to be queried separately (Fig. 5a). Chawathe et al. [7] present performance results for a real-world application using Z-curves on top of OpenDHT. The query latency (≈ 2 sec. for ≤ 30 nodes) is rather low due to the layered approach.

Partitioning. The schemes using space partitioning split the key-space into several nodes. SONAR belongs to this group of systems. The proposed systems only differ by their routing strategies.

OpenDHT [14] was one of the very first DHTs. It hashes the key-space onto a multi-dimensional torus. While the topology resembles that of SONAR and MURK (see below), CAN uses just the neighbors for routing and it does not support multi-attribute queries.

OpenDHT [4] employs a Voronoi-based space partitioning scheme and uses a multi-dimensional graph overlay with routing tables of size $O(1)$. The overlay is not based on some regular partitioning scheme (e.g. kd-tree [5]) but uses a sample of nodes to place the fingers.

Multi-attribute range queries were also addressed by Mercury [6] which needs a large number of replicas per item to achieve logarithmic routing performance. OpenDHT [12] uses super-peers and query-caching to allow multi-attribute range queries on top of the Bamboo-DHT [15].

OpenDHT et al. [9] proposed two systems for multi-dimensional range queries in peer-to-peer systems: SCRAP and MURK. SCRAP uses the traditional approach of mapping multi-dimensional to one-dimensional data with space-filling curves which destroys the data locality. Consequently, each single multi-dimensional range-query is mapped to several one-dimensional queries. MURK is more similar to our approach, as it divides the data space into hypercuboids with each cuboid assigned to one node. In contrast to SONAR, MURK uses a heuristic for routing based on skip graphs [3] to set routing fingers.

System Design

In this section we present the overall topology of SONAR and then discuss its routing and

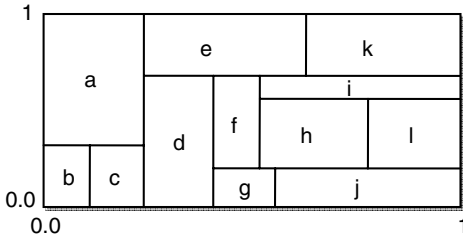


Fig. 1. Example two-dimensional overlay with attribute domains $[0, 1]$

R is used to store and retrieve *objects*. It works on a d -dimensional torus, *space*. Objects have a name, the *key*, which is a vector of d components, *coordinates* of the key. Each *dimension* of the torus is responsible for one *attribute domain*. Figure 1 illustrates a two-dimensional key-space $([0, 1]^2)$. Allocated computers, the *nodes*, are each responsible for a dedicated area (rectangle) in the key-space of the overlay (rectangles in Fig. 1). The *node-space* has the same extent as the key-space, but is completely filled with nodes. Nodes in the node-space are *adjacent* (or *neighbors*) when their key-space rectangles are adjacent. The direct mapping between key-space and node-space guarantees that keys to be stored on the same or adjacent nodes, which enables efficient range queries across node boundaries by local query propagation.

Nodes are dynamically assigned to the key-space such that each node serves approximately the same number of objects. Load-balancing is done by changing the number of nodes instead of moving around objects in the key-space. That is necessary when the number of objects or nodes in the system changes (see Section 4).

Overlay Topology

As illustrated in Figure 1, the two-dimensional key-space is covered by rectangles, each containing about the same number of objects. Because the keys are not uniformly distributed, the rectangles have different sizes and shapes. A node may have more than one neighbor per direction. The neighbors are stored in *neighbor lists*, one per dimension.

The overlay described so far resembles that of CAN [14] except for the hashing which prevents efficient range queries. Consequently, SONAR would also require \sqrt{N} network hops if it would just use the neighbors for routing. In the next section, we introduce routing tables to achieve logarithmic routing performance.

Routing

In SONAR, we use separate *routing tables*, one per dimension. Each routing table contains a mapping from key-space coordinates to the ID of a node.

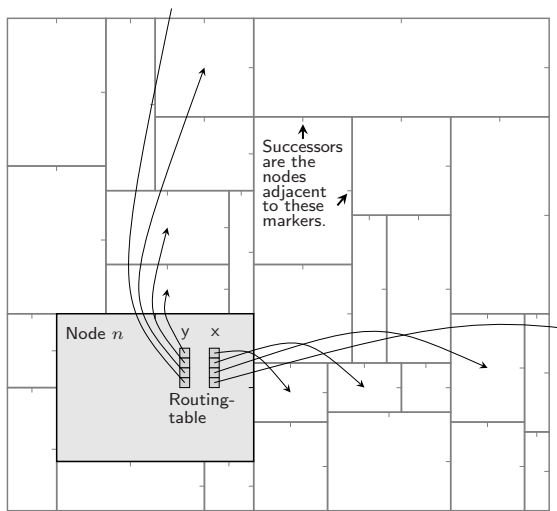


Fig. 2. Routing table for the two-dimensional case

calculate its i^{th} finger in the routing table, a node looks at its $(i - 1)^{\text{th}}$ finger and asks the remote node listed there for the $(i - 1)^{\text{th}}$ finger. At the lowest fingers point to the successor.

$$finger_i = \begin{cases} successor & : i = 0 \\ finger_{i-1}.getFinger(i - 1) & : i \neq 0 \end{cases}$$

update process works in a running system, but also during startup. All fingers are set to *unknown* except for the finger to the successor. The second entry will always succeed, because the successor knows its own successor. Filling further entries may fail (result *unknown*), because the remote node may not have determined the corresponding entry yet. But with subsequent updates, eventually all nodes will get their entries filled. The resulting routing table is similar to skip lists [13], but the behavior is more deterministic.

Used for Routing. A node may have more than one neighbor per side. We define the node adjacent to the middle of the respective side to be the *successor*. Successors are marked by small ticks in Figure 2.

The different box sizes and the calculation of longer fingers from shorter fingers are not necessarily straight in one direction. Slight deviations in y -direction might occur when following the fingers of the x -direction (and vice versa) as shown in Figure 2. Our empirical results indicate, however, that this does not affect the logarithmic routing performance (Sect. 4).

Table Size. Each node holds approximately $\log N$ fingers in its routing table.

```

calculates the entries of a routing table
d updateRoutingTable(int dim) {
    int i = 1;
    bool done = false;

    rt[dim][0] = this.Successor[dim];

    while (!done) {
        Node candidate = rt[dim][i - 1].getFinger(dim, i - 1);
        if (IsBetween(dim, rt[dim][i - 1].Key, candidate.Key, this.Key)){
            rt[dim][i] = candidate;
            i++;
        } else
            done = true;
    }

    checks whether the resp coordinate of pos lies between start and end
    IsBetween(Dim dim, Key start, Key pos, Key end);

```

Fig. 3. Finger calculation for dimension **dim**

[6] predicts the system size N by estimating the key density. SONAR simpler, deterministic solution with less overhead.

Each dimension dim , SONAR’s finger update algorithm (Fig. 3) inserts a new finger $finger_i$ as long as its position is between that of the last routing table entry $finger_{i-1}$ and that of the node itself. Otherwise the new finger circles the ring and is not inserted.

Results in Section 4 confirm that each node holds indeed $\log N$ fingers. The construction process guarantees—in contrast to Chord [18]—that no two fingers point to the same node. Since the fingers in the routing tables span an exponentially increasing number of nodes, the routing table of each dimension contains $\lceil \log D \rceil$ entries on the average, where D is the number of nodes in the system on the torus.

Finger Update. Our periodically running finger update algorithm needs only one network hop to determine an entry in the routing table. Chord in contrast needs $O(\log n)$ for the same operation, because it performs a DHT lookup to locate a finger.

Lookup and Range Queries

Like other DHTs, SONAR uses greedy routing. In each node the finger that is closest to the target reduces the Euclidean distance to the target in the key-space is found independently of the dimension (see Fig. 4).

SONAR supports range queries with multiple attributes. In its most basic form a range query is defined by d intervals for the d attribute domains. The query finds all keys whose attributes match the respective intervals and returns the corresponding objects. Because of their shape, such range queries are not supported by d -dimensional systems like range queries.

. Schütt, F. Schintke, and A. Reinefeld

```
// find the responsible node for a given key
Node find(Point target) {
    Node nextHop = findNextHop(target);
    if (nextHop == this)
        return this;
    else
        return nextHop.Find(target);
}

double getDistance(Node a, Point b);

Node findNextHop(Point target) {
    Node candidate = this;
    double distance = getDistance(this, target);

    if (distance == 0.0)
        // found target
        return this;

    for (int d = 0; d < dimensions; d++) {
        for (int i = 0; i < rt[d].Size; i++) {
            double dist = getDistance(rt[d][i], target);
            if (dist < distance) {
                // new candidate
                candidate = rt[d][i];
                distance = dist;
            }
        }
    }
    // will never happen:
    Assert(candidate != this);
    return candidate;
}
```

Fig. 4. Lookup for a target

er and a radius. Here, we assume a person located in the governmental of Berlin searching for a hotel in ‘walking distance’ (circle around the center). The query is first routed to the node responsible for the center of the circle, then forwarded to all neighbors that partially cover the circle (Fig. 5b). The query is checked against the local data and the results are returned to the originating node. Figure 6 shows the pseudocode of this algorithm. Note that duplicate messages are eliminated. `op` is an additional check for objects in the area—in this case for type `hotel`.

RR performs a range query with a single lookup. When the target node does not hold the complete key range, the query is locally forwarded. Systems using space-filling curves, in contrast, usually require more than one lookup for a range query because they map connected areas to multiple independent segments, see Figure 5a.

Topology Maintenance to Handle Churn

5.1. *Introduction*. When a node joins the system, the key-space of a participating node has to be split and the data responsibility redistributed. To achieve this, the

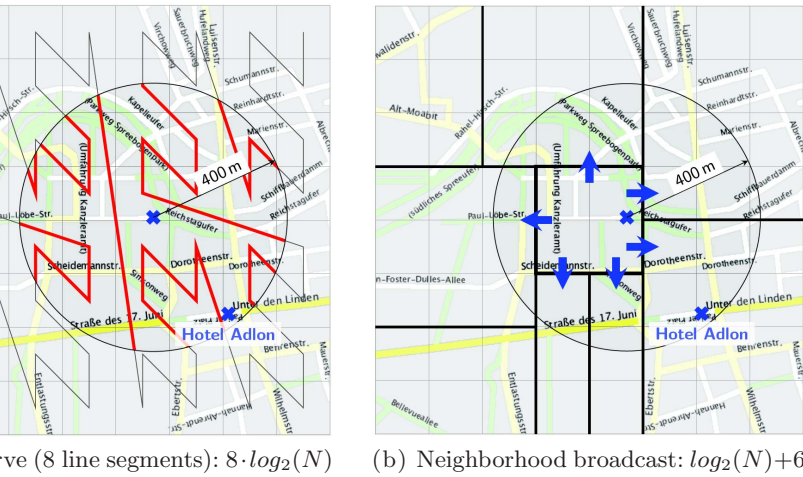


Fig. 5. Circular range query

```

// perform a range query
void queryRange(Range r, Operation op) {
    Node center = Find(r.Center);
    center.doRangeQuery(r, op, newId());
}

void doRangeQuery(Range r, Operation op, Id id) {
    // avoid redundant executions
    if (pastQueries.Contains(id))
        return;
    pastQueries.add(id);

    foreach (Node neighbor in this.Neighbors)
        if (r ∩ neighbor.Range != ∅)
            neighbor.doRangeQuery(r \ this.Range, op, id);

    // execute operation locally
    op(this, r);
}

```

Fig. 6. Range query algorithm

at a random target node: A random position in the key-space is routed and a random walk is started from there. The final target node of this candidate to be split. The random walk ensures that nodes responsible for larger areas of the key-space are not preferred over smaller ones.

the key-space and transfer one part to the new node: Splits are parallel to the coordinate system axes. The selection of the axis to be split does not strictly favor one dimension over the others because the number of nodes to be contacted for a range query could become disproportionately

ve. Handling a leaving node is more difficult, because it is not always which node can fill the area of the leaving node. For example, in Figure 1, of node *f* cannot be merged with any of its neighbors, because this would a non-rectangular node-space.

ore the node-space is constructed in such a way that the splitting plane *d-tree* [5]. KD-trees are used only for topology maintenance, similar as [9], but not as index structures like in database systems. The space of node can be taken over by a neighboring node which is also a sibling *-tree*. By keeping the tree balanced the probability of having a sibling labor increases. Each node must remember its position in the kd-tree, a describing the path from the root of the tree to the node itself.

neighboring nodes are siblings in the kd-tree, another node must be fill the gap. Either a neighboring node additionally takes over the re-ty of the separate area until a free node can be found, or two completely ent nodes that are siblings in the kd-tree have to be found to merge d thus free a node that takes over the free area. The former concept, *tual nodes*, is also used for load-balancing in other systems.

ancing. Load-balancing can be implemented by either adjusting the es of the responsibilities locally or freeing nodes in underloaded areas ng them to overloaded areas. The former has similar issues as a node e boundaries are interlocked with limited room for adjustments. The s shown to be converging [11] with predictable performance. alancing can be based on different metrics for load, like object or query combination of both. To avoid thrashing effects a threshold for per- load-balancing round must be introduced.

Empirical Results

ng the performance of SONAR we used a traveling salesman data set 4,711 cities¹. The cities' geographical locations follow a Zipf distribu- which is also common in other scenarios.

igned the responsibility of nodes by recursively splitting the key-space ger side, so that each part gets half of the cities until enough rectangles ed. Figure 7 shows a sample splitting for 256 nodes.

ordinates were mapped onto a doughnut-shaped torus rather than a cause in a globe all vertical rings meet at the poles. This would not only outing bottleneck at the poles but would also result in different ring s for the western and eastern hemisphere (southwards vs. northwards). 8 shows the results for various all-to-all searches in networks of differ- The routing performance, depicted by the '+' ticks, almost perfectly he expected $0.5 \log_2 N$ hops. Only in the larger networks the expected

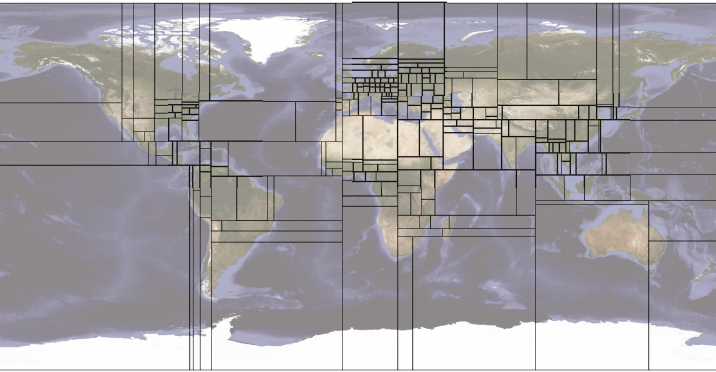


Fig. 7. 1,904,711 cities split evenly into 256 rectangular nodes

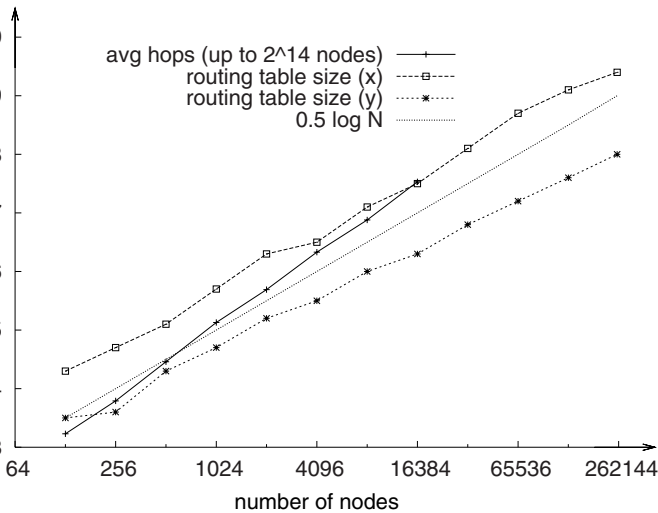


Fig. 8. SONAR results for increasing system sizes (2-dimensional)

to checked whether the number of fingers in the routing tables, which are built without global information (Fig. 4), meets our expectations. The ‘□’ ticks represent the routing table sizes in horizontal direction, and the ‘*’ ticks represent the routing table sizes in vertical direction. As expected, both graphs have the same slope of 0.5: One lies consistently above, the other below. This is attributed to the domain sizes of the coordinate system (360 versus 180 degrees) and to the number of splitting planes.

Figure 9 gives further insight into the characteristics of SONAR’s routing tables. When comparing various network sizes, the deviation of the table sizes from the theoretical values is small. The deviation of the table

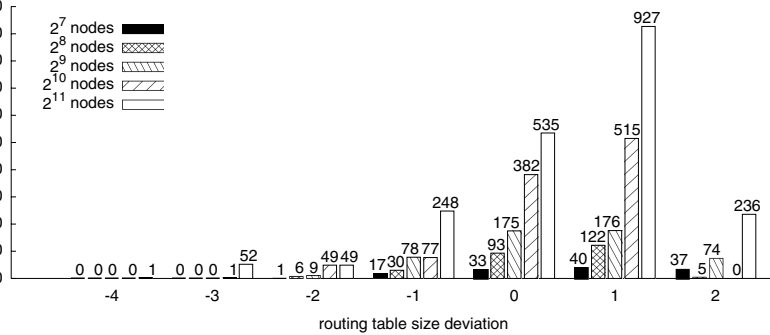


Fig. 9. Routing table size deviation from the expected value

out 25% meet the expected size of $\log_2 N$, while there is a decreasing of tables with fewer entries. These deviations are caused by the uneven distribution and by SONAR’s finger update algorithm which has a tendency in some cases an extra finger that is more than halfway around the ring (‘left’ of the own node).

Conclusion

efficiently supports range-queries on multi-dimensional data in structured overlay networks. It needs $O(\log N)$ routing steps for processing range-queries of arbitrary shapes and an arbitrary number of attribute domains. The update operation needs just one hop for updating an entry in the routing table. The presented empirical results from a Zipf distributed data set with approximately two million keys. The results confirm that SONAR does its routing with a constant number of hops—even in skewed data distributions. Additional tests for practical and uniform distributions (not shown here) gave the same routing performance. Furthermore, we observed that the sizes of the routing tables are always $O(\log N)$ although they are autonomously updated by the nodes with local information only.

Acknowledgements

Thanks to the anonymous reviewers for their valuable comments. The topographic images were taken from the ‘Blue Marble next generation’ project of NASA Earth Observatory. Thanks to Slaven Rezić for the street map of Berlin.

- ejak, A., Xu, Z.: Scalable, efficient range queries for Grid information ser-
In: P2P 2002 (2002)
- s, J., Shah, G.: Skip graphs. In: SODA (January 2003)
- i-Kashani, F., Shahabi, C.: SWAM: A family of access methods for similarity-
in peer-to-peer data networks. In: CIKM (November 2004)
- ey, J.: Multidimensional binary search trees used for associative searching.
Communications of the ACM 18(9) (1975)
- mbe, A., Agrawal, M., Seshan, S.: Mercury: Supporting scalable multi-
te range queries. In: ACM SIGCOMM 2004 (August 2004)
- athe, Y., Ramabhadran, S., Ratnasamy, S., LaMarca, A., Shenker, S., Heller-
J.: A Case Study in building layered DHT applications. In: SIGCOMM'05
st 2005)
- e, V., Günther, O.: Multidimensional access methods. ACM Computing Sur-
30(2) (1998)
- an, P., Yang, B., Garcia-Molina, H.: One torus to rule them all: Multi-
sional queries in P2P systems. In: WebDB2004 (2004)
- r, D., Lehman, E., Leighton, T., Panigrah, R., Levine, M., Lewin, D.: Con-
t hashing and random trees: Distributed caching protocols for relieving hot
on the World Wide Web. In: ACM Sympos. Theory of Comp. (May 1997)
- r, D., Ruhl, M.: Simple efficient load balancing algorithms for peer-to-peer
s. In: Voelker, G.M., Shenker, S. (eds.) IPTPS 2004. LNCS, vol. 3279,
ger, Heidelberg (2005)
- heimer, D., Albrecht, J., Patterson, D., Vahdat, A.: Design and implementa-
adeoffs for wide-area resource discovery. In: 14th IEEE Symposium on High
mance Distributed Computing (HPDC-14) (July 2005)
- W.: Skip lists: A probabilistic alternative to balanced trees. Communications
ACM (June 1990)
- samy, S., Francis, P., Handley, M., Karp, R., Shenker, S.: A scalable content-
ssable network. In: ACM SIGCOMM 2001 (August 2001)
- S., Geels, D., Roscoe, T., Kubiawicz, J.: Handling churn in a DHT. In:
edings of the USENIX Annual Technical Conference (June 2004)
- dt, C., Parashar, M.: Enabling flexible queries with guarantees in P2P sys-
IEEE Internet Computing, 19–26 (May/June 2004)
- t, T., Schintke, F., Reinefeld, A.: Structured overlay without consistent hash-
mpirical results. In: GP2PC'06 (May 2006)
- , I., Morris, R., Kaashoek, M.F., Karger, D., Balakrishnan, H.: Chord: A
le peer-to-peer lookup service for Internet application. In: ACM SIGCOMM
August 2001)
- G.: Relative frequency as a determinant of phonetic change. Harvard Studies
ssical Philology (1929)