

Querying Probabilistic Neighborhoods in Spatial Data Sets Efficiently

Moritz von Looz and Henning Meyerhenke

{moritz.looz-corswarem, meyerhenke}@kit.edu

Institute of Theoretical Informatics, Karlsruhe Institute of Technology (KIT), Germany

Abstract. The probability that two spatial objects establish some kind of mutual connection often depends on their proximity. To formalize this concept, we define the notion of a *probabilistic neighborhood*: Let P be a set of n points in \mathbb{R}^d , $q \in \mathbb{R}^d$ a query point, dist a distance metric, and $f : \mathbb{R}^+ \rightarrow [0, 1]$ a monotonically decreasing function. Then, the probabilistic neighborhood $N(q, f)$ of q with respect to f is a random subset of P and each point $p \in P$ belongs to $N(q, f)$ with probability $f(\text{dist}(p, q))$. Possible applications include query sampling and the simulation of probabilistic spreading phenomena, as well as other scenarios where the probability of a connection between two entities decreases with their distance. We present a fast, sublinear-time query algorithm to sample probabilistic neighborhoods from planar point sets. For certain distributions of planar P , we prove that our algorithm answers a query in $O((|N(q, f)| + \sqrt{n}) \log n)$ time with high probability. In experiments this yields a speedup over pairwise distance probing of at least one order of magnitude, even for rather small data sets with $n = 10^5$ and also for other point distributions not covered by the theoretical results.

1 Introduction

In many scenarios, connections between spatial objects are not certain but probabilistic, with the probability depending on the distance between them: The probability that a customer shops at a certain physical store shrinks with increasing distance to it. In disease simulations, if the social interaction graph is unknown but locations are available, disease transmission can be modeled as a random process with infection risk decreasing with distance. Moreover, the wireless connections between units in an ad-hoc network are fragile and collapse more frequently with higher distance.

For these and similar scenarios, we define the notion of a *probabilistic neighborhood* in spatial data sets: Let a set P of n points in \mathbb{R}^d , a query point $q \in \mathbb{R}^d$, a distance metric dist , and a monotonically decreasing function $f : \mathbb{R}^+ \rightarrow [0, 1]$ be given. Then, the probabilistic neighborhood $N(q, f)$ of q with respect to f is a random subset of P and each point $p \in P$ belongs to $N(q, f)$ with probability $f(\text{dist}(p, q))$. A straightforward query algorithm for sampling a probabilistic neighborhood would iterate over each point $p \in P$ and sample for each whether it is included in $N(q, f)$. This has a running time of $\Theta(n \cdot d)$ per query point, which

is prohibitive for repeated queries in large data sets. Thus we are interested in a faster algorithm for such a *probabilistic neighborhood query* (PNQ, spoken as “pink”). We restrict ourselves to the planar case in this work, but the algorithmic principle is generalizable to higher dimensions.

While the linear-time approach has appeared before in the literature for a particular application [2] (without formulating the problem as a PNQ explicitly), we are not aware of previous work performing more efficient PNQs with an index structure. For example, the probabilistic quadtree introduced by Kraetzschmar et al. [12] is designed to store probabilistic occupancy data and gives deterministic results. Other range queries related to (yet different from) our work as well as deterministic index structures are described in Section 2.2.

Contributions. We develop, analyze, implement, and evaluate an index structure and a query algorithm that together provide fast probabilistic neighborhood queries in the Euclidean and hyperbolic plane. Our key data structure for these fast PNQs is a polar quadtree which we adapt from our previous work [19]. Preprocessing for quadtree construction requires $O(n \log n)$ time with high probability¹ (whp).

To answer PNQs, we first present a simple query algorithm (Section 3). We then improve its time complexity by treating whole subtrees as so-called virtual leaves, see Section 4. As shown by our detailed theoretical analysis, the improved algorithm yields a query time complexity of $O(|N(q, f)| + \sqrt{n} \log n)$ whp to find a probabilistic neighborhood $N(q, f)$ among n points, for n sufficiently large. This is sublinear if the returned neighborhood $N(q, f)$ is of size $o(n/\log n)$ – an assumption we consider reasonable for most applications. For our theoretical results to hold, the quadtree structure needs to be able to partition the distribution of the point positions in P , i. e. not all of the probability mass may be concentrated on a single point or line. In our case of polar quadtrees, this is achieved if the distribution is continuous, integrable, rotationally invariant with respect to the origin and non-zero only for a finite area.

Experimental results are shown in Section 5: We apply our query algorithm to generate random graphs in the hyperbolic plane [14] in subquadratic time. Graphs with millions of edges can now be generated within a few minutes sequentially. This yields an acceleration of at least one order of magnitude in practice compared to a reference implementation [2] that uses linear-time queries. Compared to our previous work on graph generation [19], our new algorithm is able to generate a more extensive model. Even if the distribution of a given point set P is unknown in practice, running times are fast: As an example of probabilistic spreading behavior, we simulate a simple disease spreading mechanism on real population density geodata. In this scenario, our fast PNQs are at least two orders of magnitude faster than linear-time queries.

¹ We say “with high probability” (whp) when referring to a probability $\geq 1 - 1/n$ for sufficiently large n .

2 Preliminaries

2.1 Notation

Let the input be given as set P of n points. The points in P are distributed in a disk \mathbb{D}_R of radius R in the hyperbolic or Euclidean plane, the distribution is given by a probability density function $j(\phi, r)$ for an angle ϕ and a radius r . Recall that, for our theoretical results to hold, we require j to be known, continuous and integrable. Furthermore, j needs to be rotationally invariant – meaning that $j(\phi_1, r) = j(\phi_2, r)$ for any radius r and any two angles ϕ_1 and ϕ_2 – and positive within \mathbb{D}_R , so that $j(r) > 0 \Leftrightarrow r < R$. Due to the rotational invariance, $j(\phi, r)$ is the same for every ϕ and we can write $j(r)$. Likewise, we define $J(r)$ as the indefinite integral of $j(r)$ and normalize it so that $J(R) = 1$ (also implying $J(0) = 0$). The value $J(r)$ then gives the fraction of probability mass inside radius r .

For the distance between two points p_1 and p_2 , we use $\text{dist}_{\mathbb{H}}(p_1, p_2)$ for the hyperbolic and $\text{dist}_{\mathbb{E}}(p_1, p_2)$ for the Euclidean case. We may omit the index if a distinction is unnecessary. As mentioned, a point p is in the probabilistic neighborhood of query point q with probability $f(\text{dist}(p, q))$. Thus, a *query pair* consists of a query point q and a function $f : \mathbb{R}^+ \rightarrow [0, 1]$ that maps distances to probabilities. The function f needs to be monotonically decreasing but may be discontinuous. Note that f can be defined differently for each query. The query result, the probabilistic neighborhood of q w. r. t. f , is denoted by the set $N(q, f) \subseteq P$.

For the algorithm analysis, we use two additional sets for each query (q, f) :

- $\text{Candidates}(q, f)$: neighbor candidates examined when executing such a query,
- $\text{Cells}(q, f)$: quadtree cells examined during execution of the query.

Note that the sets $N(q, f)$, $\text{Candidates}(q, f)$ and $\text{Cells}(q, f)$ are probabilistic, thus theoretical results about their size are usually only with high probability.

2.2 Related Work

Fast deterministic range queries. Numerous index structures for fast range queries on spatial data exist. Many such index structures are based on trees or variations thereof, see Samet’s book [17] for a comprehensive overview. I/O efficient worst case analysis is usually performed using the EM model, see e. g. [3]. In more applied settings, average-case performance is of higher importance, which popularized R-trees or newer variants thereof, e. g. [11]. Concerning (balanced) quadtrees for spatial dimension d , it is known that queries require $O(d \cdot n^{1-1/d})$ time (thus $O(\sqrt{n})$ in the planar case) [17, Ch. 1.4]. Regarding PNQs our algorithm matches this query complexity up to a logarithmic factor. Yet note that, since for general f and dist in our scenario all points in the set P could be neighbors, data structures for deterministic queries cannot solve a PNQ efficiently without adaptations.

Hu et al. [10] give a query sampling algorithm for one-dimensional data that, given a set P of n points in \mathbb{R} , an interval $q = [x, y]$ and an integer, $t \geq 1$, returns t elements uniformly sampled from $P \cap q$. They describe a structure of $O(n)$ space that answers a query in $O(\log n + t)$ time and supports updates in $O(\log n)$ time. While also offering query sampling, PNQs differ from the problem considered by Hu et al. in two aspects: We consider two dimensions instead of one and our sampling probabilities are not necessarily uniform, but can be set by the user by a distance-dependent function.

Range queries on uncertain data. During the previous decade probabilistic queries *different* from PNQs have become popular. The main scenarios can be put into two categories [16]: (i) Probabilistic databases contain entries that come with a specified confidence (e.g. sensor data whose accuracy is uncertain) and (ii) objects with an uncertain location, i. e. the location is specified by a probability distribution. Both scenarios differ under typical and reasonable assumptions from ours: Queries for uncertain data are usually formulated to return *all* points in the neighborhood whose confidence/probability exceeds a certain threshold [13], or computing points that are possibly nearest neighbors [1].

In our model, in turn, the choice of inclusion of a point p is a random choice for every different p . In particular, depending on the probability distribution, *all* nodes in the plane can have positive probability to be part of some other's neighborhood. In the related scenarios this would only be true with extremely small confidence values or extremely large query circles.

Applications in fast graph generation. One application for PNQs as introduced in Section 1 is the hyperbolic random graph model by Krioukov et al. [14]. The n graph nodes are represented by points thrown into the hyperbolic plane at random² and two nodes are connected by an edge with a probability that decreases with the distance between them. An implementation of this generative model is available [2], it performs $\Theta(n^2)$ neighborhood tests. Bringmann et al. provide an algorithm to generate hyperbolic random graphs in expected linear time [6]; to our knowledge no implementation of it exists yet.

In previous work we designed a generator [19] faster than [2] for a restricted model; it runs in $O((n^{3/2} + m) \log n)$ time whp for the whole graph with m edges. The range queries discussed there are facilitated by a quadtree which supports only deterministic queries. Consequently, the queries result in unit-disk graphs in the hyperbolic plane and can be considered as a special case of the current work (a step function f with values 0 and 1 results in a deterministic query).

Our major technical inspiration for enhancing the quadtree for probabilistic neighborhoods is the work of Batagelj and Brandes [5]. They were the first to present a random sampling method to generate Erdős-Rényi-graphs with n nodes and m edges in $O(n+m)$ time complexity. Faced with a similar problem of selecting each of n elements with a constant probability p , they designed an efficient

² The probability density in the polar model depends only on radii r and R as well as a growth parameter α and is given by $g(r) := \alpha \frac{\sinh(\alpha r)}{\cosh(\alpha R) - 1}$.

algorithm (see Algorithm 2 in Appendix A). Instead of sampling each element separately, they use random jumps of length $\delta(p)$, $\delta(p) = \ln(1 - rand) / \ln(1 - p)$, with $rand$ being a random number uniformly distributed in $[0, 1)$.

2.3 Quadtree Specifics

Our key data structure is a polar region quadtree in the Euclidean or hyperbolic plane. While they are less suited to higher dimensions as for example k-d-trees, the complexity is comparable in the plane. For the (circular) range queries we discuss, quadtrees have the significant advantage of a bounded aspect ratio: A cell in a k-d-tree might extend arbitrarily far in one direction, rendering theoretical guarantees about the area affected by the query circle difficult to impossible. In contrast, the region covered by a quadtree cell is determined by its position and level.

We mostly reuse our previous definition [19] of the quadtree: A node in the quadtree is defined as a tuple $(\min_\phi, \max_\phi, \min_r, \max_r)$ with $\min_\phi \leq \max_\phi$ and $\min_r \leq \max_r$. It is responsible for a point $p = (\phi_p, r_p)$ exactly if $(\min_\phi \leq \phi_p < \max_\phi)$ and $(\min_r \leq r_p < \max_r)$. We call the region represented by a particular quadtree node its quadtree *cell*. The quadtree is parametrized by its radius R , the \max_r of the root cell. If the probability distribution j is known (which we assume for our theoretical results), we set the radius R to $\arg \min_r J(r) = 1$, i. e. to the minimum radius that contains the full probability mass. If only the points are known, the radius is set to include all of them. While in this latter case the complexity analysis of Section 3 and 4 does not hold, fast running times in practice can still be achieved (see Section 5).

3 Baseline Query Algorithm

We begin the main technical part by describing adaptations in the quadtree construction as well as a baseline query algorithm. This latter algorithm introduces the main idea, but is asymptotically not faster than the straightforward approach. In Section 4 it is then refined to support faster queries.

3.1 Quadtree Construction

At each quadtree node v , we store the size of the subtree rooted there. We then generalize the rule for node splitting to handle point distributions j as defined in Section 2.1: As is usual for quadtrees, a leaf cell c is split into four children when it exceeds its fixed capacity. Since our quadtree is polar, this split happens once in the angular and once in the radial direction. Due to the rotational symmetry of j , splitting in the angular direction is straightforward as the angle range is halved: $\text{mid}_\phi := \frac{\max_\phi + \min_\phi}{2}$. For the radial direction, we choose the splitting radius to result in an equal division of probability mass. The total probability mass in a ring delimited by \min_r and \max_r is $J(\max_r) - J(\min_r)$. Since $j(r)$ is positive for r between R and 0, the restricted function $J|_{[0, R]}$ defined above is a

bijection. The inverse $(J|_{[0,R]})^{-1}$ thus exists and we set the splitting radius mid_r to $(J|_{[0,R]})^{-1}\left(\frac{J(\max_r)+J(\min_r)}{2}\right)$.

Figure 1 visualizes a point distribution on a hyperbolic disk with 200 points and Figure 2 its corresponding quadtree.

Two results on quadtree properties help to establish the time complexity of quadtree operations. They are generalized versions of our previous work [19, Lemmas 1 and 2] and state that each quadtree cell contains the same expected number of points and that the quadtree height is $O(\log n)$ whp (proofs in Appendix B).

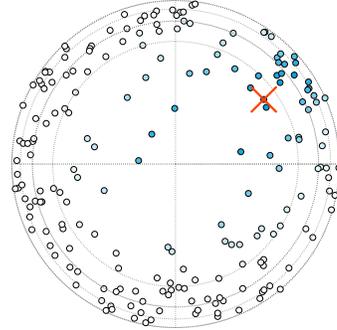


Fig. 1. Query over 200 points in a polar hyperbolic quadtree, with $f(d) := 1/(e^{(d-7.78)} + 1)$ and the query point q marked by a red cross. Points are colored according to the probability that they are included in the result. Blue represents a high probability, white a probability of zero.

Lemma 1. Let \mathbb{D}_R be a hyperbolic or Euclidean disk of radius R , j a probability distribution on \mathbb{D}_R which fulfills the properties defined in Section 2.1, p a point in \mathbb{D}_R which is sampled from j , and T be a polar quadtree on \mathbb{D}_R . Let C be a quadtree cell at depth i . Then, the probability that p is in C is 4^{-i} .

Lemma 2. Just as in Lemma 1, let \mathbb{D}_R be a hyperbolic or Euclidean disk of radius R , j a probability distribution on \mathbb{D}_R which fulfills the properties defined in Section 2.1, and T be a polar quadtree on \mathbb{D}_R . The expected number of nodes in T is then in $O(n)$.

Proposition 1. Let \mathbb{D}_R and j be as in Lemma 1. Let T be a polar quadtree on \mathbb{D}_R constructed to fit j . Then, for n sufficiently large, $\text{height}(T) \in O(\log n)$ whp.

A direct consequence from the results above and our previous work [19] is the preprocessing time for the quadtree construction. The generalized splitting rule and storing the subtree sizes only change constant factors.

Corollary 1. Since a point insertion takes $O(\log n)$ time whp, constructing a quadtree on n points distributed as in Section 2.1 takes $O(n \log n)$ time whp.

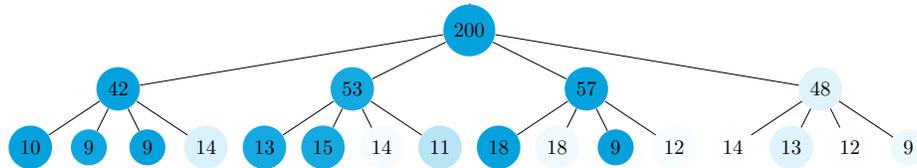


Fig. 2. Visualization of the data structure used in Figure 1. Quadtree nodes are colored according to the upper probability bound for points contained in them. The color of a quadtree node c is the darkest possible shade (dark = high probability) of any point contained in the subtree rooted at c . Each node is marked with the number of points in its subtree.

Algorithm 1: QuadNode.getProbabilisticNeighborhood

Input: query point q , prob. function f , quadtree node c
Output: probabilistic neighborhood of q

```
1 N = {};  
2  $\underline{b}$  = dist( $q, c$ );  
   /* Distance between point and cell */  
3  $\bar{b}$  =  $f(\underline{b})$ ;  
   /* Since  $f$  is monotonically decreasing, a lower bound for the  
   distance gives an upper bound  $\bar{b}$  for the probability. */  
4  $s$  = number of points in  $c$ ;  
5 if  $c$  is not leaf then  
   /* internal node: descend, add recursive result to local set */  
6   for  $child \in children(c)$  do  
7     | add getProbabilisticNeighborhood( $q, f, child$ ) to N;  
8 else  
   /* leaf case: apply idea of Batagelj and Brandes [5] */  
9   for  $i=0; i < s ; i++$  do  
10    |  $\delta = \ln(1 - rand) / \ln(1 - \bar{b})$ ;  
11    |  $i += \delta$ ;  
12    | if  $i \geq s$  then  
13    |   | break;  
14    |    $prob = f(\text{dist}(q, c.points[i])) / \bar{b}$ ;  
15    |   add  $c.points[i]$  to N with probability  $prob$   
16 return N
```

3.2 Algorithm

The baseline version of our query (Algorithm 1) has unfortunately a time complexity of $\Theta(n)$, but serves as a foundation for the fast version (Section 4). It takes as input a query point q , a function f and a quadtree cell c . Initially, it is called with the root node of the quadtree and recursively descends the tree. The algorithm returns a point set $N(q, f) \subseteq P$ with

$$\Pr[p \in N(q, f)] = f(\text{dist}(q, p)). \quad (1)$$

Algorithm 1 descends the quadtree recursively until it reaches the leaves. Once a leaf l is reached, a lower bound \underline{b} for the distance between the query point q and all the points in l is computed (Line 2). Such distance calculations are detailed in Appendix B.6. Since f is monotonically decreasing, this lower bound for the distance gives an upper bound \bar{b} for the probability that a given point in l is a member of the returned point set (Line 3). This bound is used to select *neighbor candidates* in a similar manner as Bategelj and Brandes [5]: In Line 10, a random number of vertices is skipped, so that every vertex in l is selected as a neighbor candidate with probability \bar{b} . The actual distance $\text{dist}(q, a)$ between a candidate a and the query point q is at least \underline{b} and the probability of $a \in N(q, f)$ thus at most \bar{b} . For each candidate, this actual distance $\text{dist}(q, a)$ is then calculated and

a neighbor candidate is confirmed as a neighbor with probability $f(\text{dist}(q, a))/\bar{b}$ in Line 14.

Regarding correctness and time complexity of Algorithm 1, we can state:

Proposition 2. *Let T be a quadtree as defined above, q be a query point and $f : \mathbb{R}^+ \rightarrow [0, 1]$ a monotonically decreasing function which maps distances to probabilities. The probability that a point p is returned by a PNQ (q, f) from Algorithm 1 is $f(\text{dist}(q, p))$, independently from whether other points are returned.*

Proposition 3. *Let T be a quadtree with n points. The running time of Algorithm 1 per query on T is $\Theta(n)$ in expectation.*

The proofs can be found in Appendices B.4 and B.5.

4 Queries in Sublinear Time by Subtree Aggregation

One reason for the linear time complexity of the baseline query is the fact that every quadtree node is visited. To reach a sublinear time complexity, we thus aggregate subtrees into *virtual leaf cells* whenever doing so reduces the number of examined cells and does not increase the number of candidates too much.

To this end, let S be a subtree starting at depth l of a quadtree T . During the execution of Algorithm 1, a lower bound \underline{b} for the distance between S and the query point q is calculated, yielding also an upper bound \bar{b} for the neighbor probability of each point in S . At this step, it is possible to treat S as a *virtual leaf cell*, sample jumping widths using \bar{b} as upper bound and use these widths to select candidates within S . Aggregating a subtree to a virtual leaf cell allows skipping leaf cells which do not contain candidates, but uses a weaker bound \bar{b} and thus a potentially larger candidate set. Thus, a fast algorithm requires an aggregation criterion which keeps both the number of candidates and the number of examined quadtree cells low.

As stated before, we record the number of points in each subtree during quadtree construction. This information is now used for the query algorithm: We aggregate a subtree S to a virtual leaf cell exactly if $|S|$, the number of points contained in S , is below $1/f(\text{dist}(S, q))$. This corresponds to less than one expected candidate within S . The changes required in Algorithm 1 to use the subtree aggregation are minor. Lines 5, 14 and 15 are changed to:

5 if c is inner node and $|c| \cdot \bar{b} \geq 1$ then

14 neighbor = maybeGetKthElement(q, f, i, \bar{b}, c);
15 add neighbor to N if not null

The main change consists in the use of the function maybeGetKthElement (Algorithm 5, Appendix C). Given a subtree S , an index k , q , f , and \bar{b} , this function descends S to the leaf cell containing the k th element. This element p_k is then accepted with probability $f(\text{dist}(q, p_k))/\bar{b}$.

Since the upper bound calculated at the root of the aggregated subtree is not smaller than the individual upper bounds at the original leaf cells, Proposition 2 also holds for the virtual leaf cells. This establishes the correctness.

The time complexity is given by the following theorem, whose proof can be found in Appendix D.

Theorem 1. *Let T be a quadtree with n points and (q, f) a query pair. A query (q, f) using subtree aggregation has time complexity $O(|N(q, f)| + \sqrt{n} \log n)$ whp.*

5 Application Case Studies

In order to test our algorithm for PNQs, we apply it in two application case studies, one for Euclidean, the other one for hyperbolic geometry. For the Euclidean case study we build a simple disease spread simulation as an example for a probabilistic spreading process. The probability distribution of points is in this case non-uniform and unknown. The hyperbolic application, in turn, is a generator for complex networks with a known point distribution.

5.1 Probabilistic Spreading

When both contact graph and travel patterns of a susceptible population are not known in detail, the resulting spreading behavior of an infectious disease seems probabilistic. Contagious diseases usually spread to people in the vicinity of infected persons, but an infectious person occasionally bridges larger distances by travel and spreads the disease this way. We model this effect with our probabilistic neighborhood function f , giving a higher probability for small distances and a lower but non-zero probability for larger distances. Note that this scenario is meant as an example of the probabilistic spreading simulations possible with our algorithm and not as highly realistic from an epidemiological point of view.

In the simulation, the population is given as a set P of points in the Euclidean plane. In the initial step, exactly one point (= person) from P is marked as infected. Then, in each round, a PNQ is performed for each infected person q . All points in $N(q, f)$ become infected in the next round. We use an SIR model [8], i. e. previously infected persons recover with a certain probability in each round and stay infectious otherwise. In our simulation, persons recover with a rate of 0.8 and are then immune.

5.2 Random Hyperbolic Graph Generation

Random hyperbolic graphs (RHGs, also see Section 2.2) are a generative graph model for complex networks. For graph generation one places n points (= vertices) randomly in a hyperbolic disk. The radius R of the disk can be used to control the average degree of the network. A pair of vertices is connected by an edge with

Country	5000 PDP queries	Construction QT	5000 QT queries
France	1007 seconds	1.6 seconds	1.2 seconds
Germany	1395 seconds	2.8 seconds	1.3 seconds
USA	4804 seconds	8.7 seconds	0.7 seconds

Table 1. Running time results for polar Euclidean quadtrees on population data. The query points were selected uniformly at random from P, the probabilistic neighborhood function is $f(x) := (1/x) \cdot e^7/n$.

a probability that depends on the vertices’ hyperbolic distance. This connection probability is given in [14, Eq. (41)] and parametrized by a temperature $T \geq 0$:

$$f(x) = \frac{1}{e^{(1/T) \cdot (x-R)/2} + 1} \quad (2)$$

This definition of random hyperbolic graphs is a generalized version of the one considered in our previous work, which was restricted to the special case of $T = 0$.

5.3 Experimental Settings and Results

Our implementation uses the NetworKit toolkit [18] and is written in C++ 11. It is included in NetworKit release 4.1. Running time measurements were made with g++ 4.8 -O3 on a machine with 128 GB RAM and an Intel Xeon E5-1630 v3 CPU with four cores at 3.7 GHz base frequency. Our code is sequential, as is the reference implementation for random hyperbolic graph generation [2].

Disease Spread Simulation. We experimented on three data sets taken from NASA population density raster data [7] for Germany, France and the USA. They consist of rectangles with small square cells (geographic areas) where for each cell the population from the year 2000 is given. To obtain a set of points, we randomly distribute points in each cell to fit 1/20th of the population density. Figure 4 (left) in the appendix shows an example with roughly 4 million points on the map of Germany. The data sets of France and USA have roughly 3 and 14 million points, respectively.

The number of required queries naturally depends heavily on the simulated disease. For our parameters, a number of 5000 queries is typically reached within the first dozen steps. To evaluate the algorithmic speedup, Table 1 compares running times for 5000 pairwise distance probing (PDP) queries against 5000 fast PNQs on the three country datasets. To obtain a similar total number of infections, we use a slightly different probabilistic neighborhood function for each country and divide by the population: $f(x) := (1/x) \cdot e^7/n$. This results in a slower initial progression for the US. Our algorithm achieves a speedup factor of at least two orders of magnitude, even including the quadtree construction time.

Random Hyperbolic Graph Generation. An example graph generated from hyperbolic geometry can be seen in Figure 4 (right) in the appendix. We compare our

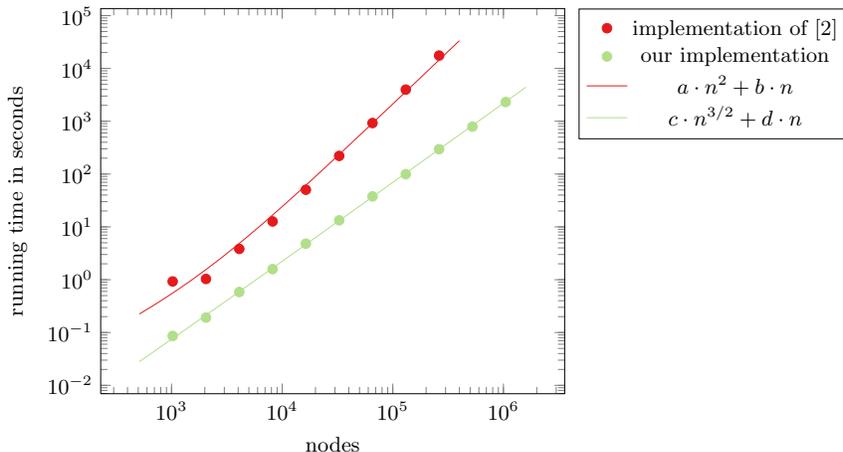


Fig. 3. Comparison of running times to generate networks with 2^{10} - 2^{20} vertices, $\alpha = 1$, $T = 0.5$ and average degree $\bar{k} = 6$. The gap between the running times widens, which in the loglog-plot implies a different exponent in the time complexities. Running times are fitted with $a = 2.089 \cdot 10^{-7}$, $b = 3.311 \cdot 10^{-4}$, $c = 2.18 \cdot 10^{-6}$ and $d = 5.6 \cdot 10^{-6}$.

generator using PNQs with the only (to our knowledge) previously existing generator for general random hyperbolic graphs [2], i. e. those not only following the threshold model. As seen in Figure 3, our implementation is faster by at least one order of magnitude and the experimental running times support our theoretical time complexity of $O((n^{3/2} + m) \log n)$. A comparison of the generated graphs with those created by the existing implementation can be found in Appendix G. The differences measured by a set of suitable network analysis metrics are within the range of random fluctuations for the sample size of 80.

6 Conclusions

After formally defining the notion of probabilistic neighborhoods, we have presented a quadtree-based query algorithm for such neighborhoods in the Euclidean and hyperbolic plane. Our analysis shows a time complexity of $O((|N(q, f)| + \sqrt{n}) \log n)$, our algorithm is to the best of our knowledge the first to solve the problem asymptotically faster than pairwise distance probing. With two example applications we have shown that our algorithm is also faster in practice by at least one order of magnitude.

Acknowledgements. This work is partially supported by German Research Foundation (DFG) grant ME 3619/3-1 within the Priority Programme 1736 *Algorithms for Big Data*. The authors thank Mark Ortman for helpful discussions.

References

1. Pankaj K Agarwal, Boris Aronov, Sariel Har-Peled, Jeff M Phillips, Ke Yi, and Wuzhou Zhang. Nearest neighbor searching under uncertainty ii. In *Proceedings of the 32nd symposium on Principles of database systems*, pages 115–126. ACM, 2013.
2. Rodrigo Aldecoa, Chiara Orsini, and Dmitri Krioukov. Hyperbolic graph generator. *Computer Physics Communications*, 2015.
3. Lars Arge and Kasper Green Larsen. I/o-efficient spatial data structures for range queries. *SIGSPATIAL Special*, 4(2):2–7, July 2012.
4. R. Arratia and L. Gordon. Tutorial on large deviations for the binomial distribution. *Bulletin of Mathematical Biology*, 51(1):125–131, 1989.
5. Vladimir Batagelj and Ulrik Brandes. Efficient generation of large random networks. *Physical Review E*, 71(3):036113, 2005.
6. Karl Bringmann, Ralph Keusch, and Johannes Lengler. Geometric inhomogeneous random graphs. *arXiv preprint arXiv:1511.00576*, 2015.
7. Center for International Earth Science Information Network CIESIN Columbia University; Centro Internacional de Agricultura Tropical CIAT. Gridded population of the world, version 3 (gpwv3): Population density grid, 2005.
8. Herbert W Hethcote. The mathematics of infectious diseases. *SIAM review*, 42(4):599–653, 2000.
9. W. Hoeffding. Probability inequalities for sums of bounded random variables. *J. American Statistical Association*, 58(301):13–30, 1963.
10. Xiaocheng Hu, Miao Qiao, and Yufei Tao. Independent range sampling. In *Proceedings of the 33rd ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 246–255. ACM, 2014.
11. Ibrahim Kamel and Christos Faloutsos. Hilbert r-tree: An improved r-tree using fractals. In *Proceedings of the 20th International Conference on Very Large Data Bases, VLDB '94*, pages 500–509, San Francisco, CA, USA, 1994. Morgan Kaufmann Publishers Inc.
12. Gerhard K Kraetzschmar, Guillem Pages Gassull, Klaus Uhl, Guillem Pags, and Gassull Klaus Uhl. Probabilistic quadtrees for variable-resolution mapping of large environments. In *Proceedings of the 5th IFAC/EURON symposium on intelligent autonomous vehicles*, 2004.
13. Hans-Peter Kriegel, Peter Kunath, and Matthias Renz. Probabilistic nearest-neighbor query on uncertain objects. In *Advances in databases: concepts, systems and applications*, pages 337–348. Springer, 2007.
14. Dmitri Krioukov, Fragkiskos Papadopoulos, Maksim Kitsak, Amin Vahdat, and Marián Boguñá. Hyperbolic geometry of complex networks. *Physical Review E*, 82(3):036106, Sep 2010.
15. Michael Mitzenmacher and Eli Upfal. *Probability and computing: Randomized algorithms and probabilistic analysis*. Cambridge University Press, 2005.
16. Jian Pei, Ming Hua, Yufei Tao, and Xuemin Lin. Query answering techniques on uncertain and probabilistic data: tutorial summary. In *Proc. 2008 ACM SIGMOD intl. conference on Management of data*, pages 1357–1364. ACM, 2008.
17. Hanan Samet. *Foundations of Multidimensional and Metric Data Structures*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2005.
18. Christian L Staudt, Aleksejs Sazonovs, and Henning Meyerhenke. NetworKit: A tool suite for large-scale complex network analysis. *arXiv:1403.3005*, 2015.
19. Moritz von Looz, Roman Prutkin, and Henning Meyerhenke. Generating random hyperbolic graphs in subquadratic time. In *ISAAC 2015: Proc. 26th Int'l Symp. on Algorithms and Computation*, 2015.

A Related Algorithmic Idea

Our approach was inspired by the following algorithm with optimal linear running time for Erdős-Rényi graph generation [5].

Algorithm 2: Efficient neighborhood generation for Erdős-Rényi graphs [5].

Input: number of vertices n , edge probability $0 < p < 1$

Output: $G = (\{0, \dots, n-1\}, E) \in \mathcal{G}(n, p)$

$E = \emptyset;$

$v = 1;$

$w = -1;$

while $v < n$ **do**

 draw $r \in [0, 1)$ uniformly at random;

$w = w + 1 + \lfloor \frac{\log(1-r)}{\log(1-p)} \rfloor;$

while $w \geq v$ **and** $v < n$ **do**

$w = w - v;$

$v = v + 1;$

if $v < n$ **then**

 add $\{u, v\}$ to E

B Proofs of Section 3

B.1 Proof of Lemma 1

Proof. Due to the similarity of Lemma 1 to [19, Lemma 1], the proof follows a similar structure. Let C be a quadtree cell at level k , delimited by \min_r , \max_r , \min_ϕ and \max_ϕ . As stated in Section 2.1, we require the point probability distribution to be rotationally invariant. The probability that a point p is in C is then given by

$$\Pr(p \in C) = \frac{\max_\phi - \min_\phi}{2\pi} \cdot (J(\max_r) - J(\min_r)). \quad (3)$$

The boundaries of the children of C are given by the splitting rules in Section 3.1.

$$\text{mid}_\phi := \frac{\max_\phi + \min_\phi}{2} \quad (4)$$

$$\text{mid}_r := (J|_{[0,R]})^{-1} \left(\frac{J(\max_r) + J(\min_r)}{2} \right) \quad (5)$$

We proceed with induction over the depth i of C . Start of induction ($i = 0$): At depth 0, only the root cell exists and covers the whole disk. Since $C = \mathbb{D}_R$, $\Pr(p \in C) = 1 = 4^{-0}$.

Inductive step ($i \rightarrow i + 1$): Let C_i be a node at depth i . C_i is delimited by the radial boundaries \min_r and \max_r , as well as the angular boundaries \min_ϕ

and \max_ϕ . It has four children at depth $i + 1$, separated by mid_r and mid_ϕ . Let SW be the south west child of C_i . With Eq. (3), the probability of $p \in SW$ is:

$$\Pr(p \in SW) = \frac{\text{mid}_\phi - \min_\phi}{2\pi} \cdot (J(\text{mid}_r) - J(\min_r)) \quad (6)$$

Using Equations (4) and (5), this results in a probability of

$$\Pr(p \in SW) = \frac{\frac{\max_\phi + \min_\phi}{2} - \min_\phi}{2\pi} \cdot \left(J \left((J|_{[0,R]})^{-1} \left(\frac{J(\max_r) + J(\min_r)}{2} \right) \right) - J(\min_r) \right) \quad (7)$$

$$\Pr(p \in SW) = \frac{\frac{\max_\phi + \min_\phi}{2} - \min_\phi}{2\pi} \cdot \left(\frac{J(\max_r) + J(\min_r)}{2} - J(\min_r) \right) \quad (8)$$

$$\Pr(p \in SW) = \frac{\frac{\max_\phi - \min_\phi}{2}}{2\pi} \cdot \left(\frac{J(\max_r) - J(\min_r)}{2} \right) \quad (9)$$

$$\Pr(p \in SW) = \frac{1}{4} \frac{\max_\phi - \min_\phi}{2\pi} \cdot (J(\max_r) - J(\min_r)) \quad (10)$$

$$(11)$$

As per the induction hypothesis, $\Pr(p \in C_i)$ is 4^{-i} and $\Pr(p \in SW)$ is thus $\frac{1}{4} \cdot 4^{-i} = 4^{-(i+1)}$. Due to symmetry when selecting mid_ϕ , the same holds for the south east child of C_i . Together, they contain half of the probability mass of C_i . Again due to symmetry, the same proof then holds for the northern children as well. \square

B.2 Proof of Lemma 2

Proof. A quadtree T containing n points can have at most n non-empty leaf cells. We can thus bound the total number of leaf cells in T by limiting the number of empty cells.

An empty leaf cell occurs when a previous leaf cell c is split. We consider two cases, depending on how many of the children of c contain points:

Case 1: All but one of the children of c are empty and all points in c are concentrated in one child. We call a split of this kind an *excess* split, since it did not result in dividing the points in c .

Case 2: At least two children of c contain points.

The number of excess splits caused by a pair of points depends on the area they are clustered in. Two sufficiently close points could cause a potentially unbounded number of excess splits. However, due to Lemma 1, each child cell contains a quarter of the probability mass of its parent cell. Given two points p, q in a cell which is split, they end up in different child cells with probability $3/4$.

The expected number of excess splits for a point p is thus at most³

$$\sum_{i=0}^{\infty} i \cdot 4^{-i} = \frac{4}{9}. \quad (12)$$

Due to the linearity of expectations, the expected number of excess splits caused by n points is then at most $4n/9$. Each excess split causes four additional quadtree nodes, three of them are empty leaf cells.

If we remove all quadtree nodes caused by excess splits and reconnect the tree by connecting the remaining leaves to their lowest unremoved ancestor, every inner node in the remaining tree T' has at least two non-empty subtrees. Since a binary tree with n leaves has $O(n)$ inner nodes [17] and the branching factor in T' is at least two, T' also contains at most $O(n)$ inner nodes.

Together with the expected $O(n)$ nodes caused by excess splits, this results in $O(n)$ nodes in T in expectation. \square

B.3 Proof of Proposition 1

Proof. We proved a similar lemma in previous work [19], for hyperbolic geometry only and a restricted family of probability distributions. The requirement for that proof was that a given point p has a probability of 4^{-i} to land in a given cell at depth i . In Lemma 1, we show that this requirement is fulfilled for the quadtrees used in this paper in both Euclidean and hyperbolic geometry. We can thus reuse the proof of [19, Lemma 2], which we include for the purpose of self-containment:

Proof of [19, Lemma 2]

Proof. In a complete quadtree, 4^i cells exist at depth i . For analysis purposes only, we construct such a complete but initially empty quadtree of height $k = 3 \cdot \lceil \log_4(n) \rceil$, which has at least n^3 leaf cells. As seen in Lemma 1, a given point has an equal chance to land in each leaf cell. Hence, we can apply [19, Lemma 6] with each leaf cell being a bin and a point being a ball. (The fact that we can have more than n^3 leaf cells only helps in reducing the average load.) From this we can conclude that, for n sufficiently large, no leaf cell of the current tree contains more than 1 point with high probability (whp). Consequently, the total quadtree height does not exceed $k = 3 \cdot \lceil \log_4(n) \rceil \in O(\log n)$ whp.

Let T' be the quadtree as constructed in the previous paragraph, starting with a complete quadtree of height k and splitting leaves when their capacity is exceeded. Let T be the quadtree created in our algorithm, starting with a root node, inserting points and also splitting leaves when necessary, growing the tree downward.

Since both trees grow downward as necessary to accommodate all points, but T does not start with a complete quadtree of height k , the set of quadtree nodes

³ Note that the real number of excess splits might be lower, since a split might separate another point from p and q .

in T is a subset of the quadtree nodes in T' . Consequently, the height of T is bounded by $O(\log n)$ whp as well. \square

B.4 Proof of Proposition 2

Proof. Note that the hyperbolic [Euclidean] distances, which are mapped to probabilities according to the function f , are calculated by Algorithm 3 [Algorithm 4], which are presented in Appendix B.6 (together with their correctness proofs). We continue the current proof with details for all three main steps.

Step 1: Between two points, the jumping width δ is given by Line 10. The probability that exactly i points are skipped between two given candidates is $(1 - \bar{b})^i \cdot \bar{b}$:

$$\Pr(i \leq \delta < i + 1) = \quad (13)$$

$$\Pr(i \leq \ln(1 - r) / \ln(1 - \bar{b}) < i + 1) = \quad (14)$$

$$\Pr(\ln(1 - r) \leq i \cdot \ln(1 - \bar{b}) \wedge \ln(1 - r) > (i + 1) \cdot \ln(1 - \bar{b})) = \quad (15)$$

$$\Pr(1 - (1 - \bar{b})^i \leq r < 1 - (1 - \bar{b})^{i+1}) = \quad (16)$$

$$1 - (1 - \bar{b})^{i+1} - 1 + (1 - \bar{b})^i = \quad (17)$$

$$(1 - \bar{b})^i (1 - (1 - \bar{b})) = \quad (18)$$

$$(1 - \bar{b})^i \cdot \bar{b} \quad (19)$$

Note that in Eq. (14) the denominator is negative, thus the direction of the inequality is reversed in the transformation. The transformation from Eq. (16) to Eq. (17) works since r is uniformly distributed.

Following from Eq. (19), the probability is \bar{b} for $i = 0$, and if a point is selected as a candidate, the subsequent point is selected with a probability of \bar{b} .

Step 2: Let p_i, p_j and p_l be points in a leaf, with $i < j < l$ and let p_i be a neighbor candidate. For now we assume that no other points in the same leaf are candidates and consider the probability that p_l is selected as a candidate depending on whether the intermediate point p_j is a candidate.

Case 2.1: If point p_j is a candidate, then point p_l is selected if $l - j$ points are skipped after selecting p_j . Due to Step 1, this probability is $(1 - \bar{b})^{l-j} \cdot \bar{b}$

Case 2.2: If point p_j is *not* a candidate, then point p_l is selected if $l - i$ points are skipped after selecting p_i . Given that p_j is not selected, at least $j - i$ points are skipped. The conditional probability is then:

$$\Pr(l - i \leq \delta < l - i + 1 | \delta > j - i) = \quad (20)$$

$$\Pr(1 - (1 - \bar{b})^{l-i} < r < (1 - (1 - \bar{b})^{l-i+1}) | \delta > j - i) = \quad (21)$$

$$(1 - \bar{b})^{l-i} \cdot \bar{b} / (1 - \bar{b})^{j-i} = \quad (22)$$

$$(1 - \bar{b})^{l-j} \cdot \bar{b} \quad (23)$$

As both cases yield the same result, the probability $\Pr(p_l \in \text{Candidates})$ is independent of whether p_j is a candidate.

Step 3: Let C be a leaf cell in which all points up to point p_i are selected as candidates. Due to Step 1, the probability that p_{i+1} is also a candidate, meaning no points are skipped, is $(1 - \bar{b})^0 \cdot \bar{b} = \bar{b}$. Due to Step 2, the probability of p_{i+1} being a candidate is independent of whether p_i is a candidate. This can be applied iteratively until the beginning of the leaf cell, yielding a probability of \bar{b} for p_i being a candidate, independent of whether other points are selected.

A neighbor candidate p_i is accepted as a neighbor with probability $f(\text{dist}(p_i, q))/\bar{b}$ in Line 14. Since \bar{b} is an upper bound for the neighborhood probability, the acceptance ratio is between 0 and 1. The probability for a point p to be in the probabilistic neighborhood computed by Algorithm 1 is thus:

$$\Pr(p \in N(q, f)) = \quad (24)$$

$$\Pr(p \in N(q, f) \wedge p \in \text{Candidates}(q, f)) = \quad (25)$$

$$\Pr(p \in N(q, f) | p \in \text{Candidates}(q, f)) \cdot \Pr(p \in \text{Candidates}(q, f)) = \quad (26)$$

$$f(\text{dist}(p, q))/\bar{b} \cdot \bar{b} = \quad (27)$$

$$f(\text{dist}(p, q)) \quad (28)$$

□

B.5 Proof of Proposition 3

Proof. The total time complexity of the query algorithm is determined by the number of recursive calls (Line 7) and the number of loop iterations (Line 9). During tree traversal, one recursive call is made for each examined quadtree node. During examination of a leaf, one loop iteration happens for every examined candidate. Let the set of neighbors ($N(q, f)$), candidates ($\text{Candidates}(q, f)$) and examined cells ($\text{Cells}(q, f)$) be as defined in Section 2.1. The time complexity of the query is then in $\Theta(|\text{Candidates}(q, f)| + |\text{Cells}(q, f)|)$.

All cells of the quadtree are examined, thus $\text{Cells}(q, f) = \text{Cells}(T)$. If the cells are split using the medians of point positions, then no leaf cell is empty and the tree contains at most n cells. If cells are split using the theoretical probability distributions, the tree contains at most $O(n)$ cells in expectation due to Lemma 2. It follows that the number of examined cells is in $\Theta(n)$ in expectation. Since the candidate set is a subset of the point set, the expected number of candidates is at most n . The query time complexity is then in $O(n) + \Theta(|\text{Cells}(T)|) = \Theta(n)$ in expectation. □

B.6 Distance between Quadtree Cell and Point

To calculate the upper bound \bar{b} used in Algorithm 1, we need a lower bound for the distance between the query point q and any point in a given quadtree cell. Since the quadtree cells are polar, the distance calculations might be unfamiliar and we show and prove them explicitly. For the hyperbolic case, the distance calculations are shown in Algorithm 3 and proven in Lemma 3. The Euclidean calculations are shown in Algorithm 4 and proven in Lemma 4.

Algorithm 3: Infimum and supremum of distance in a hyperbolic polar quadtree

Input: quadtree cell $C = (\min_r, \max_r, \min_\phi, \max_\phi)$, query point $q = (\phi_q, r_q)$
Output: infimum and supremum of hyperbolic distances q to interior of C
 /* start with corners of cell as possible extrema */
 1 cornerSet = $\{(\min_\phi, \min_r), (\min_\phi, \max_r), (\max_\phi, \min_r), (\max_\phi, \max_r)\}$;
 2 $a = \cosh(r_q)$;
 3 $b = \sinh r_q \cdot \cos(\phi_q - \min_\phi)$;
 /* Left/Right boundaries */
 4 leftExtremum = $\frac{1}{2} \ln\left(\frac{a+b}{a-b}\right)$;
 5 if $\min_r < \text{leftExtremum} < \max_r$ then
 6 | add $(\min_\phi, \text{leftExtremum})$ to cornerSet;
 7 $b = \sinh r_q \cdot \cos(\phi_q - \max_\phi)$;
 8 rightExtremum = $\frac{1}{2} \ln\left(\frac{a+b}{a-b}\right)$;
 /* Top/bottom boundaries */
 9 if $\min_r < \text{rightExtremum} < \max_r$ then
 10 | add $(\max_\phi, \text{rightExtremum})$ to cornerSet;
 11 if $\min_\phi < \phi_q \max_\phi$ then
 12 | add (ϕ_q, \min_r) and (ϕ_q, \max_r) to cornerSet;
 13 $\phi_{\text{mirrored}} = \phi_q + \pi \pmod{2\pi}$;
 14 if $\min_\phi < \phi_{\text{mirrored}} < \max_\phi$ then
 15 | add $(\phi_{\text{mirrored}}, \min_r)$ and $(\phi_{\text{mirrored}}, \max_r)$ to cornerSet;
 /* If point is in cell, distance is zero: */
 16 if $\min_\phi \leq \phi_q < \max_\phi$ AND $\min_r \leq r_q < \max_r$ then
 17 | infimum = 0;
 18 else
 19 | infimum = $\min_{e \in \text{cornerSet}} \text{dist}_{\mathbb{H}}(q, e)$;
 20 supremum = $\max_{e \in \text{cornerSet}} \text{dist}_{\mathbb{H}}(q, e)$;
 21 return infimum, supremum;

Lemma 3. Let C be a quadtree cell and q a point in hyperbolic space. The first value returned by Algorithm 3 is the distance of C to q .

Proof. When q is in C , the distance is trivially zero. Otherwise, the distance between q and C can be reduced to the distance between q and the boundary of C , δC :

$$\text{dist}_{\mathbb{H}}(C, q) = \text{dist}_{\mathbb{H}}(\delta C, q) = \inf_{p \in \delta C} \text{dist}_{\mathbb{H}}(p, q) \quad (29)$$

Since the boundary is closed, this infimum is actually a minimum:

$$\text{dist}_{\mathbb{H}}(C, q) = \inf_{p \in \delta C} \text{dist}_{\mathbb{H}}(p, q) = \min_{p \in \delta C} \text{dist}_{\mathbb{H}}(p, q) \quad (30)$$

The boundary of a quadtree cell consists of four closed curves:

- left: $\{(\min_\phi, r) | \min_r \leq r \leq \max_r\}$
- right: $\{(\max_\phi, r) | \min_r \leq r \leq \max_r\}$

Algorithm 4: Infimum and supremum of distance in a Euclidean polar quadtree

```

Input: quadtree cell  $C = (\min_r, \max_r, \min_\phi, \max_\phi)$ , query point  $q = (\phi_q, r_q)$ 
Output: infimum and supremum of Euclidean distances  $q$  to interior of  $C$ 
/* start with corners of cell as possible extrema */
1 cornerSet =  $\{(\min_\phi, \min_r), (\min_\phi, \max_r), (\max_\phi, \min_r), (\max_\phi, \max_r)\}$ ;
/* Left/Right boundaries */
2 leftExtremum =  $r_q \cdot \cos(\min_\phi - \phi_q)$ ;
3 if  $\min_r < \text{leftExtremum} < \max_r$  then
4 | add  $(\min_\phi, \text{leftExtremum})$  to cornerSet;
5 rightExtremum =  $r_q \cdot \cos(\max_\phi - \phi_q)$ ;
6 if  $\min_r < \text{rightExtremum} < \max_r$  then
7 | add  $(\max_\phi, \text{rightExtremum})$  to cornerSet;
/* Top/bottom boundaries */
8 if  $\min_\phi < \phi_q < \max_\phi$  then
9 | add  $(\phi_q, \min_r)$  and  $(\phi_q, \max_r)$  to cornerSet;
10  $\phi_{\text{mirrored}} = \phi_q + \pi \bmod 2\pi$ ;
11 if  $\min_\phi < \phi_{\text{mirrored}} < \max_\phi$  then
12 | add  $(\phi_{\text{mirrored}}, \min_r)$  and  $(\phi_{\text{mirrored}}, \max_r)$  to cornerSet;
/* If point is in cell, distance is zero: */
13 if  $\min_\phi \leq \phi_q < \max_\phi$  AND  $\min_r \leq r_q < \max_r$  then
14 | infimum = 0;
15 else
16 | infimum =  $\min_{e \in \text{cornerSet}} \text{dist}_{\mathbb{H}}(q, e)$ ;
17 supremum =  $\max_{e \in \text{cornerSet}} \text{dist}_{\mathbb{H}}(q, e)$ ;
18 return infimum, supremum;

```

- lower: $\{(\phi, \min_r) | \min_\phi \leq \phi \leq \max_\phi\}$
- upper: $\{(\phi, \max_r) | \min_\phi \leq \phi \leq \max_\phi\}$

We write the distance to the whole boundary as a minimum over the distances to its parts:

$$\text{dist}_{\mathbb{H}}(\delta C, q) = \min_{A \in \{\text{left}, \text{right}, \text{lower}, \text{upper}\}} \text{dist}_{\mathbb{H}}(A, q) \quad (31)$$

All points on an angular boundary curve A have the same angular coordinate ϕ_A . Let $d_A(r) = \text{acosh}(\cosh(r) \cosh(r_q) - \sinh(r) \sinh(r_q) \cos(\phi_q - \phi_A))$ for a fixed point q . The distance $\text{dist}_{\mathbb{H}}(A, q)$ can then be reduced to:

$$\text{dist}_{\mathbb{H}}(A, q) = \min_{\min_r \leq r \leq \max_r} d_A(r) \quad (32)$$

$$(33)$$

The minimum of d_A on A is the minimum of $d_A(\min_r)$, $d_A(\max_r)$ and the value at possible extrema. To find the extrema, we define a function $g(r) = \cosh(r) \cosh(r_q) - \sinh(r) \sinh(r_q) \cos(\phi_q - \phi_A)$. Since acosh is strictly monotone, $g(r)$ has the same extrema as $d_A(r)$.

The factors $\cosh(r_q)$ and $\sinh(r_q) \cos(\phi_q - \phi_A)$ do not depend on r , to increase readability we substitute them with the constants a and b :

$$a = \cosh(r_q) \quad (34)$$

$$b = \sinh(r_q) \cos(\phi_q - \phi_A) \quad (35)$$

$$d_A(r) = \operatorname{acosh}(\cosh(r) \cdot a - \sinh(r) \cdot b) \quad (36)$$

$$g(r) = \cosh(r) \cdot a - \sinh(r) \cdot b \quad (37)$$

The derivative of g is thus:

$$g'(r) = \sinh(r) \cdot a - \cosh(r) \cdot b = \frac{e^r - e^{-r}}{2} \cdot a - \frac{e^r + e^{-r}}{2} \cdot b \quad (38)$$

With some transformations, we get the roots of $g'(r)$:

Case $a = b$:

$$g'(r) = 0 \Leftrightarrow \quad (39)$$

$$\frac{e^r - e^{-r}}{2} \cdot a = \frac{e^r + e^{-r}}{2} \cdot a \quad (40)$$

$$e^r - e^{-r} = e^r + e^{-r} \quad (41)$$

$$-e^{-r} = e^{-r} \quad (42)$$

$$e^{-r} = 0 \quad (43)$$

$$(44)$$

For $a = b$, d_A has no extrema in \mathbb{R} .

$a \neq b$:

$$g'(r) = 0 \Leftrightarrow \quad (45)$$

$$\frac{e^r - e^{-r}}{2} \cdot a = \frac{e^r + e^{-r}}{2} \cdot b \Leftrightarrow \quad (46)$$

$$ae^r - ae^{-r} = be^r + be^{-r} \Leftrightarrow \quad (47)$$

$$(a - b)e^r - (a + b)e^{-r} = 0 \Leftrightarrow \quad (48)$$

$$(a - b)e^r = (a + b)e^{-r} \Leftrightarrow \quad (49)$$

$$e^r = \frac{a + b}{a - b} e^{-r} \Leftrightarrow \quad (50)$$

$$e^{2r} = \frac{a + b}{a - b} \Leftrightarrow \quad (51)$$

$$2r = \ln \left(\frac{a + b}{a - b} \right) \Leftrightarrow \quad (52)$$

$$r = \frac{1}{2} \ln \left(\frac{a + b}{a - b} \right) \quad (53)$$

For $a \neq b$, d_A has a single extremum at $\frac{1}{2} \ln \left(\frac{a+b}{a-b} \right)$. This extremum is calculated for both angular boundaries in Lines 4 and 8 of Algorithm 3.

If $d(r)$ has an extremum x in A , the minimum of $d_A(r)$ on A is $\min\{d_A(\min_r), d_A(\max_r), d_A(x)\}$, otherwise it is $\min\{d_A(\min_r), d_A(\max_r)\}$.

A similar approach works for the radial boundary curves. Let B be a radial boundary curve at radius r_B and angular bounds \min_ϕ and \max_ϕ . Let $d_B(\phi)$ be the distance to q restricted to radius r_B .

$$d_B : [0, 2\pi] \rightarrow \mathbb{R} \quad (54)$$

$$d_B(\phi) = \text{acosh}(\cosh(r_B) \cosh(r_q) - \sinh(r_B) \sinh(r_q) \cos(\phi_q - \phi)) \quad (55)$$

Similarly to the angular boundaries, we define some constants and a function $g(\phi)$ with the same extrema as d_B :

$$a = \cosh(r_B) \cosh(r_q) \quad (56)$$

$$b = \sinh(r_B) \sinh(r_q) \quad (57)$$

$$g(\phi) = a - b \cos(\phi_q - \phi) \quad (58)$$

Case: $b = 0$:

$$b = \sinh(r_B) \sinh(r_q) = 0 \Leftrightarrow \quad (59)$$

$$g(\phi) = a \quad (60)$$

Since g is constant, no extrema exist.

Case: $b \neq 0$: We obtain the extrema with some transformations:

$$g'(\phi) = -b \sin(\phi_q - \phi) \quad (61)$$

$$g'(\phi) = 0 \Leftrightarrow \quad (62)$$

$$\sin(\phi_q - \phi) = 0 \Leftrightarrow \quad (63)$$

$$\phi = \phi_q \pmod{\pi} \quad (64)$$

The distance function $d_B(\phi)$ thus has two extrema.

The minimum of $d_B(r)$ on B is then:

$$\min_{r \in B} d_B(r) = \min\{d_B(\min_r), d_B(\max_r)\} \cup \{d_B(\phi) \mid \min_\phi \leq \phi \leq \max_\phi \wedge \phi = \phi_q \pmod{\pi}\} \quad (65)$$

The distance $\text{dist}_{\mathbb{H}}(C, q)$ can thus be written as the minimum of four to ten point-to-point distances. Algorithm 3 collects the arguments for these distances in the variable `cornerSet` and returns the distance minimum as the first return value. \square

Lemma 4. *Let T be a polar quadtree in Euclidean space, c a quadtree cell of T and q a point in Euclidean space. The first value returned by Algorithm 4 is the distance of c to q .*

Proof. The general distance equation for polar coordinates in Euclidean space is

$$f(r_p, r_q, \phi_p, \phi_q) = \sqrt{r_p^2 + r_q^2 - 2r_p r_q \cos(\phi_p - \phi_q)} \quad (66)$$

If the query point q is within C , the distance is zero. Otherwise, the distance between q and C is equal to the distance between q and the boundary of C . We consider each boundary component separately and derive the extrema of the distance function.

Radial boundary. When considering the radial boundary, everything but one angle is fixed:

$$f(\phi_p) = \sqrt{r_p^2 + r_q^2 - 2r_p r_q \cos(\phi_p - \phi_q)} \quad (67)$$

Since the distance is positive and the square root is a monotone function, the extrema of the previous function are at the same values as the extrema of its square $g(\phi)$:

$$g(\phi_p) = r_p^2 + r_q^2 - 2r_p r_q \cos(\phi_p - \phi_q) \quad (68)$$

We set the derivative to zero to find the extrema:

$$g'(\phi_p) = 0 \Leftrightarrow \quad (69)$$

$$2r_p r_q \sin(\phi_p - \phi_q) \cdot (\phi_p - \phi_q) = 0 \quad (70)$$

$$\phi_p = \phi_q \pmod{\pi} \quad (71)$$

Angular boundary. Similar to the radial boundary, we fix everything but the radius:

$$f(r_p) = \sqrt{r_p^2 + r_q^2 - 2r_p r_q \cos(\phi_p - \phi_q)} \quad (72)$$

Again, we define a helper function with the same extrema:

$$g(r_p) = r_p^2 + r_q^2 - 2r_p r_q \cos(\phi_p - \phi_q) \quad (73)$$

We set the derivative to zero to find the extrema:

$$g'(r_p) = 0 \Leftrightarrow \quad (74)$$

$$2r_p - 2r_q \cos(\phi_p - \phi_q) = 0 \Leftrightarrow \quad (75)$$

$$r_p = r_q \cos(\phi_p - \phi_q) \Rightarrow \quad (76)$$

$$g(r_p) = r_p^2 + r_q^2 - 2r_p^2 \quad (77)$$

$$= r_q^2 - r_p^2 \quad (78)$$

$$= r_q^2(1 - \cos(\phi_p - \phi_q)) \quad (79)$$

An extremum of f on the boundary of cell c is either at one of its corners or at the points derived in Eq. (71) or Eq. (79). If $q \notin c$, the minimum over these points and the corners, as computed by Algorithm 4, is the minimal distance between q and any point in c . If q is contained in c , the distance is trivially zero. \square

C Algorithm maybeGetKthElement, used in Section 4

Algorithm 5: maybeGetKthElement

Input: query point q , function f , index k , bound \bar{b} , subtree S

Output: k th point of S or empty set

```
1 if  $S.isLeaf()$  then
2   acceptance =  $f(\text{dist}(q, S.points[k]))/\bar{b}$ ;
3   if  $1 - \text{rand}() < \text{acceptance}$  then
4     return  $S.points[k]$ ;
5   else
6     return  $\emptyset$ ;
7 else
8   /* Recursive call */
9   offset := 0;
10  for  $child \in S.children$  do
11    if  $k - \text{offset} < |child|$  then
12      /*  $|child|$  is the number of points in  $child$  */
13      return maybeGetKthElement( $q, f, k - \text{offset}, \bar{b}, child$ );
14      offset +=  $|child|$ ;
```

D Proof of Theorem 1

Proof. Similar to the baseline algorithm, the complexity of the faster query is determined by the number of recursive calls and the total number of loop iterations across the calls. The first corresponds to the number of examined quadtree cells, the second to the total number of candidates. With subtree aggregation, we obtain improved bounds: Lemma 7 limits the number of candidates to $O(|N(q, f)| + \sqrt{n})$ whp, while Lemma 8 bounds the number of examined quadtree cells to $O((|N(q, f)| + \sqrt{n}) \log n)$ whp. Together, this results in a query complexity of $O((|N(q, f)| + \sqrt{n}) \log n)$ whp. \square

For the lemmas required in the proof of Theorem 1 we need to introduce some notation: Let T be a quadtree with n points, S a subtree of T containing s points, q a query point and f a function mapping distances to probabilities. The set of neighbors ($N(q, f)$), candidates ($\text{Candidates}(q, f)$) and examined cells ($\text{Cells}(q, f)$) are defined as in Section 2.1.

For the analysis we divide the space around the query point q into infinitely many bands, based on the probabilities given by f . A point $p \in P$ is in band i exactly if the probability of it being a neighbor of q is between $2^{-(i+1)}$ and 2^{-i} :

$$p \in \text{band } i \Leftrightarrow 2^{-(i+1)} < f(\text{dist}(p, q)) \leq 2^{-i}$$

Based on these bands, we divide the previous sets into infinitely many subsets:

- $P(q, f, i) := \{v \in P \mid 2^{-(i+1)} < f(\text{dist}(v, q)) \leq 2^{-i}\}$
- $N(q, f, i) := N(q, f) \cap P(q, f, i)$
- $\text{Candidates}(q, f, i) := \text{Candidates}(q, f) \cap P(q, f, i)$
- $\text{Cells}(q, f, i) := \{c \in \text{Cells}(q, f) \mid 2^{-(i+1)} < f(\text{dist}(c, q)) \leq 2^{-i}\}$

Note that for fixed n , all but at most finitely many of these sets are empty. We call the quadtree cells in $\text{Cells}(q, f, i)$ to be *anchored* in band i . The region covered by a quadtree cell is in general not aligned with the probability bands, thus a quadtree cell anchored in band i ($c \in \text{Cells}(q, f, i)$) may contain points from higher bands (i.e. with lower probabilities).

We continue with two auxiliary results used in Lemma 7: Lemma 5 helps in bounding the number of candidates that are in the same band as their (virtual or original) quadtree cell is anchored in. Lemma 6 is used to bound the number of points in a higher band than their original quadtree cell.

Lemma 5. *Let n be a natural number and let A, B be sets with $A \subseteq B, |B| \leq n$ and the following property: $\Pr(b \in A) \geq 0.5, \forall b \in B$. Further, let the probabilities for membership in A be independent. Then, the number of points in B is in $O(|A| + \log n)$ with probability at least $1 - 1/n^3$.*

Proof. Let $X = |A|$ be a random variable denoting the size of A . Since the individual probabilities for membership in A might be different, X does not necessarily follow a binomial distribution. We define an auxiliary distribution

$Y := B(|B|, 0.5)$. Since all membership probabilities for A are at least 0.5, lower tail bounds derived for Y also hold for X .

The probability that Y is less than $0.1|B|$ is then [9]:

$$\Pr(Y < 0.1|B|) \leq \exp\left(-2 \frac{(0.5|B| - 0.1|B|)^2}{|B|}\right) \quad (80)$$

$$= \exp\left(-2 \frac{(0.4|B|)^2}{|B|}\right) \quad (81)$$

$$= \exp(-2 \cdot 0.16|B|) \quad (82)$$

$$= \exp(-0.32|B|) \quad (83)$$

$$(84)$$

Similar to the proof of Lemma 8, we conclude with a case distinction:

If $|B| > 10 \log n$: The probability $\Pr(|A| < 0.1|B|)$ is then $\Pr(|A| < 0.1|B|) \leq \Pr(Y < 0.1|B|) \leq \exp(-3.2 \log n) = n^{-3.2} < 1/n^3$. Thus $|B| \leq 10|A| \in O(|A|)$ with probability at least $1 - 1/n^3$.

If $|B| < 10 \log n$: $|B|$ is then trivially in $O(\log n)$. □

Lemma 6. *Let T be a polar hyperbolic [Euclidean] quadtree with n points and $s < n$ a natural number. Let A be a circle in the hyperbolic [Euclidean] plane and let \ominus be the disjoint set of subtrees of T that contain at most s points and are cut by A . Then, the subtrees in \ominus contain at most $24\sqrt{n \cdot s}$ points with probability at least $1 - 0.7^{\sqrt{n}}$ for n sufficiently large.*

Proof. This proof is adapted from [19, Lemma 3]. Let $k := \lfloor \log_4 n/s \rfloor$ be the minimal depth at which cells have at least s points in expectation. At most 4^k cells exist at depth k , defined by at most 2^k angular and 2^k radial divisions. When following the circumference of the query circle A , each newly cut cell requires the crossing of an angular or radial division. Each radial and angular coordinate occurs at most twice on the circle boundary, thus each division can be crossed at most twice. With two types of divisions, A crosses at most $2 \cdot 2 \cdot 2^k = 4 \cdot 2^{\lfloor \log_4 n/s \rfloor}$ cells at depth k . Since the value of $4 \cdot 2^{\lfloor \log_4 n/s \rfloor}$ is at most $4 \cdot 2^{\log_4 n/s}$, this yields $\leq 8 \cdot \sqrt{n/s}$ cut cells. We denote the set of cut cells with ζ . Since the cells in ζ cover the circumference of the circle A , a subtree S which is cut by A is either contained within one of the cells in ζ , corresponds to one of the cells or contains one. In the first two cases, all points in S are within the cells of ζ . In the second case, at least one cell of ζ is contained in S . As the subtrees are disjoint, this cell cannot be contained in any other of the considered subtrees. Thus, there are no more subtrees containing points not in ζ than there are cells in ζ , which are less than $8 \cdot \sqrt{n/s}$ many.

Due to Lemma 1, the probability that a given point is in a given cell at level k is 4^{-k} . The number of points contained in cells of ζ thus follows a binomial distribution $B(n, p)$. An upper bound for the probability p is given by $\frac{8 \cdot \sqrt{ns}}{n}$, thus

a tail bound for a slightly different distribution $B(n, \frac{8\sqrt{ns}}{n})$ also holds for $B(n, p)$. In the proof of [19, Lemma 7] a similar distribution is considered. Setting the variable c to $8\sqrt{s}$, we see that the probability of ζ containing more than $16 \cdot \sqrt{sn}$ points is smaller than $0.7\sqrt{n}$.

The subtrees in \ominus contain at most s points by definition, thus an upper bound for the number of points in these subtrees is given by $s \cdot 8 \cdot \sqrt{n/s}$ (points not in ζ) + $16 \cdot \sqrt{sn}$ (points in ζ). This results in at most $24 \cdot \sqrt{sn}$ points contained in \ominus with probability at least $1 - 0.7\sqrt{n}$. \square

The following Lemmas 7 and 8 bound the number of examined candidates and examined quadtree cells and are used in the proof of Theorem 1.

Lemma 7. *Let T be a quadtree with n points and (q, f) a query pair. The number of candidates examined by a query using subtree aggregation is in $O(|N(q, f)| + \sqrt{n})$ whp.*

Proof. For the analysis we consider each probability band i separately. As defined above, band i contains points with a neighbor probability of $2^{-(i+1)}$ to 2^{-i} . Among the cells anchored in band i , some are original leaf cells and others are virtual leaf cells created by subtree aggregation. The virtual leaf cells contain less than one expected candidate and thus less than 2^{i+1} points. The capacity of the original leaf cells is constant. All the points in cells anchored in band i have a probability between $2^{-(i+1)}$ and 2^{-i} to be a candidate. Among the points in virtual or original leaf cells, some are in the same band their cell is anchored in, others are in higher bands.

We divide the set of points within cells anchored in band i into four subsets:

1. points in band i and in original leaf cells
2. points in band i and in virtual leaf cells
3. points not in band i and in original leaf cells
4. points not in band i and in virtual leaf cells

The points in the first two sets are unproblematic. Since the probability that a point in these sets is a neighbor is at least $2^{-(i+1)}$, the probability for a given candidate to be a neighbor is at least $\frac{1}{2}$. Due to Lemma 5, the number of candidates in these sets is in $O(|N(q, f)| + \log n)$ whp, which is in $O(|N(q, f)| + \sqrt{n})$ whp.

Points in the third set are in cells cut by the boundary between band i and band $i + 1$. Since the probabilities are determined by the distance, this boundary is a circle and we can use Lemma 6 to bound the number of points to $24\sqrt{n \cdot \text{capacity}}$ with probability at least $1 - 0.7\sqrt{n}$ for n sufficiently large. The mentioned capacity is the capacity of the original leaf cells.

Likewise, points in the fourth set are in virtual leaf cells cut by the boundary between bands i and $i + 1$. A virtual leaf cell, which is an aggregated subtree, contains at most 2^{i+1} points, otherwise it would not have been aggregated. Again, using Lemma 6, we can bound the number of points in these sets to $24\sqrt{n} \cdot 2^{i+1}$ points with probability at least $1 - 0.7\sqrt{n}$.

We denote the union of the third and fourth sets with $\text{Overhang}(q, f, i)$. From the individual bounds derived in the previous paragraphs, we obtain an upper bound for the number of points in $\text{Overhang}(q, f, i)$ of $24(\sqrt{n \cdot \text{capacity}} + \sqrt{n \cdot 2^{i+1}})$ with probability at least $(1 - 0.7^{\sqrt{n}})^2$. Simplifying the bound, we get that $|\text{Overhang}(q, f, i)| \leq 24\sqrt{n} \cdot (2^{(i+1)/2} + \sqrt{\text{capacity}})$ with probability at least $1 - 2 \cdot 0.7^{\sqrt{n}}$.

Each of the points in $\text{Overhang}(q, f, i)$ is a candidate with a probability between 2^{-i} and $2^{-(i+1)}$. The candidates are sampled independently (see Step 2 of Lemma 2). While different points may have different probabilities of being a candidate and the total number of candidates does not follow a binomial distribution, we can bound the probabilities from above with 2^{-i} .

We proceed towards a Chernoff bound for the total number of candidates across all overhangs. Let X_i denote the random variable representing the candidates within $|\text{Overhang}(q, f, i)|$ and let $X = \sum_{i=0}^{\infty} X_i$ denote the total number of candidates in overhangs.

The expected value $\mathbb{E}(X)$ follows from the linearity of expectations:

$$\mathbb{E}(X) = \sum_{i=0}^{\infty} \mathbb{E}(X_i) \tag{85}$$

$$\leq \sum_{i=0}^{\infty} 24\sqrt{n} \cdot (2^{(i+1)/2} + \sqrt{\text{capacity}}) \cdot 2^{-i} \tag{86}$$

$$= 24\sqrt{n} \sum_{i=0}^{\infty} (\sqrt{2} \cdot 2^{-i/2} + 2^{-i} \sqrt{\text{capacity}}) \tag{87}$$

$$= 24\sqrt{n}((2\sqrt{2} + 2) + 2\sqrt{\text{capacity}}) \tag{88}$$

(Cells anchored in the band ∞ , which has an upper bound \bar{b} of zero for the neighborhood probability, do not have any candidates and can be omitted here.)

Since the candidates are sampled independently with a probability of at most 2^{-i} , we can treat X as a sum of independent Bernoulli random variables without losing generality. This allows us to use a multiplicative Chernoff bound [15] and we can now give an upper bound for the probability that the overhangs contain more than twice as many candidates as expected:

$$\Pr(X > 2\mathbb{E}(X)) \leq \left(\frac{e}{2^2}\right)^{\mathbb{E}(X)} \tag{89}$$

$$= \left(\frac{e}{2^2}\right)^{24\sqrt{n}((2\sqrt{2}+2)+2\sqrt{\text{capacity}})} \tag{90}$$

$$\leq \left(\frac{e}{2^2}\right)^{\sqrt{n}} \tag{91}$$

$$\leq 0.7^{\sqrt{n}} \tag{92}$$

While the random variable $X = \sum_{i=0}^{\infty} X_i$ is written as an infinite sum, all but at most n bands are empty, thus we are only applying the Chernoff bound over finitely many variables. For each of the at most n non-empty bands, we

defined two tail bounds for the number of points in the overhang. Including this last bound, we thus have a chain of $2n + 1$ tail bounds, each with a probability of at least $(1 - 0.7^{\sqrt{n}})$. The event that any of these tail bounds is violated is a union over each event that a specific tail bound is violated. With a union bound [15, Lemma 1.2], the probability that any of the individual tail bounds is violated is at most $(2n + 1)0.7^{\sqrt{n}}$. Since $\frac{1}{(2n+1)0.7^{\sqrt{n}}}$ grows faster than n for n sufficiently large, we conclude that the total number of candidates is thus bounded by $O(|N(q, f)|) + 48\sqrt{n}((2\sqrt{2}+2) + 2\sqrt{\text{capacity}})$ with probability at least $(1 - 1/n)$ for n sufficiently large. The leaf capacity is constant, thus the number of candidates evaluated during execution of a query (q, f) is in $O(|N(q, f)| + \sqrt{n})$ whp. \square

We proceed with an auxiliary result necessary for bounding the number of examined quadtree cells in a query:

Lemma 8. *Let T be a quadtree with n points and (q, f) a query pair. The number of quadtree cells examined by a query using subtree aggregation is in $O((|N(q, f)| + \sqrt{n}) \log n)$.*

To prove Lemma 8, we first introduce another auxiliary lemma:

Lemma 9. *Let \mathbb{D}_R be a hyperbolic or Euclidean disk of radius R and let T be a polar quadtree on \mathbb{D}_R containing n points distributed according to Section 2.1. Let $\mathcal{Y}(q, f)$ be the set of unaggregated quadtree cells that have only (virtual) leaf cells as children (category C2 in the proof of Lemma 8). With a query using subtree aggregation, $|\mathcal{Y}(q, f)|$ is in $O(|N(q, f)| + \sqrt{n})$ whp.*

Proof. Let $c \in \mathcal{Y}(q, f, i)$ be such an unaggregated quadtree cell anchored in band i that has only original or virtual leaf cells as children. It contains at least 2^i points and has four children, of which at least one is also anchored in band i . We denote this (virtual) leaf anchored in band i with l . Since each child of c contains the same probability mass (Lemma 1), each point of c is in l with probability $1/4$:

$$\Pr(p \in l | p \in c) = \frac{1}{4}. \quad (93)$$

A point in l is a candidate (in l) with probability $f(\text{dist}(q, l))$, which is between $2^{-(i+1)}$ and 2^{-i} since l is anchored in band i . The probability that a given point $p \in c$ is a candidate in l is then

$$\Pr(p \in l \wedge p \in \text{Candidates}(q, f, i) \mid p \in c) = \frac{1}{4} \cdot f(\text{dist}(q, l)) \geq 2^{-(i+3)} \quad (94)$$

Since the point positions and memberships in $\text{Candidates}(q, f, i)$ are independent, we can bound the number of candidates in l with a binomial distribution $B(|c|, 2^{-(i+3)})$. The probability that l contains no candidates is:

$$f(0, |c|, \frac{1}{8} \cdot 2^{-i}) = (1 - \frac{1}{8} \cdot 2^{-i})^{|c|} \quad (95)$$

$$\leq (1 - \frac{1}{8} \cdot \frac{1}{2^i})^{2^i} \quad (96)$$

Considered as a function of i , this probability is monotonically ascending. In the limit of $2^i \rightarrow \infty$, it trends to $\exp(-1/8) \approx 0.88$, a value it never exceeds. The probability that the cell c contains at least one candidate is then above $1 - \frac{1}{\sqrt[8]{e}} > 0.1$.

For each cell in \mathcal{Y} , the probability that it contains at least one candidate is > 0.1 . Let X be the random variable denoting the number of cells in \mathcal{Y} that contain at least one candidate. We define an auxiliary binomial distribution $B(|\mathcal{Y}|, 0.1)$ and use a tail bound to estimate the number of cells in \mathcal{Y} containing candidates. Let $Y \propto B(|\mathcal{Y}|, 0.1)$ be a random variable distributed according to this auxiliary distribution.

We use a tail bound from [4] to limit the probability that $Y < 0.05|\mathcal{Y}|$ to at most $\exp(-|\mathcal{Y}|/80)$. Since 0.1 was a lower bound for the probability that a cell contains a candidate, this tail bound also holds for X . The probability that the set of \mathcal{Y} contains at least $0.05|\mathcal{Y}|$ many candidates is then at least $(1 - \exp(-|\mathcal{Y}|/80))$.

We continue with a case distinction:

If $|\mathcal{Y}| \in \omega(\sqrt{n})$: The probability $(1 - \exp(-|\mathcal{Y}|/80))$ is then smaller than $(1 - \exp(-\sqrt{n}/80))$, which is $< 1/n$ for sufficiently large n . Thus the number of examined quadtree cells during a query is then linear in the number of candidates. Due to Lemma 7, this is in $O(|N(q, f)| + \sqrt{n})$.

If $|\mathcal{Y}| \in o(\sqrt{n})$: The cardinality $|\mathcal{Y}|$ is trivially in $O(\sqrt{n})$. □

The proof of Lemma 8 then follows easily:

Proof. We split the set of examined quadtree cells into three categories:

- leaf cells and root nodes of aggregated subtrees (C1)
- parents of cells in the first category (C2)
- all other (C3)

The third category (C3) then exclusively consists of inner nodes in the quadtree. When following a chain of nodes in category C3 from the root downwards, it ends with a node in category C2. The size $|C3|$ is thus at most $O(|C2| \log n)$ whp, since the number of elements in a chain cannot exceed the height of the quadtree, which is $O(\log n)$ by Proposition 1.

With a branching factor of 4, $|C1| = 4|C2|$ holds.

The number of cells in category C2 can be bounded using Lemma 9 to $O(|N| + \sqrt{n})$ with high probability. The total number of examined cells is thus in $O((|N| + \sqrt{n}) \log n)$. □

E Visualizations of Experimental Results

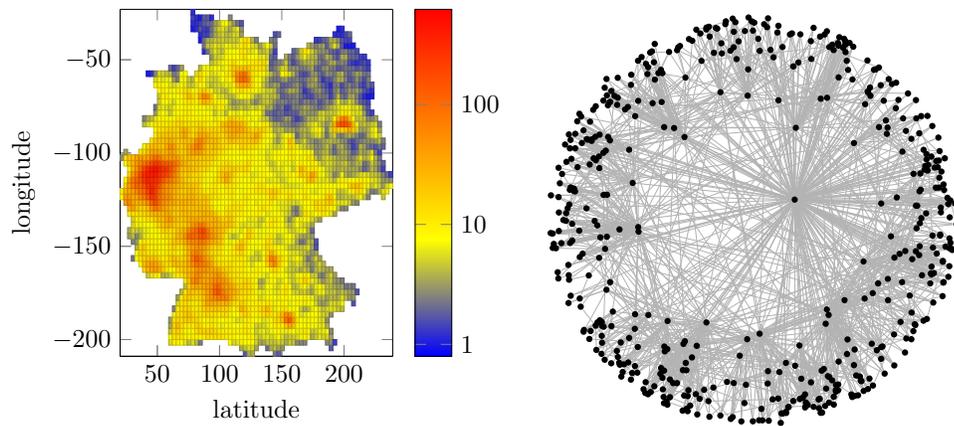


Fig. 4. Left: Twenty-third time step of a simulated disease progression through Germany. The colors indicate the number of infected persons within a cell. Right: Random hyperbolic graph with 500 nodes and average degree 12.

F Performance of Baseline Algorithm 1

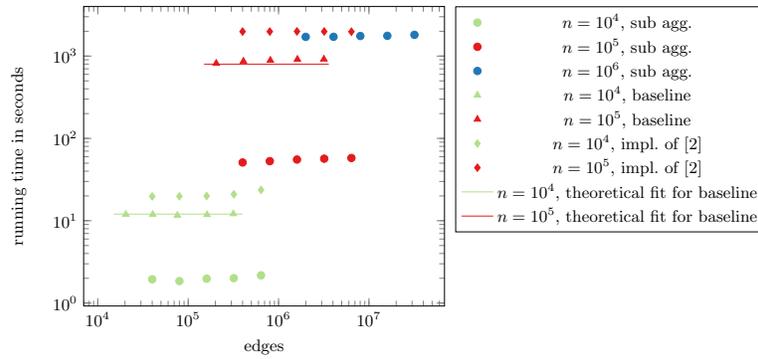


Fig. 5. Comparison of running times to generate networks with 10^4 to 10^6 vertices. Generating a graph requires n queries. Shown are running times of the baseline algorithm, queries using subtree aggregation and the implementation of [2]. The theoretical fit is given by the equation $T(n, m) = (7.94 \cdot 10^{-8}n^2 + 4.1 \cdot 10^{-4}n)$ seconds. The baseline algorithm is still faster than the previous implementation [2], but much slower than the improved query using subtree aggregation.

G Fast RHG Generator vs Reference Implementation [2]

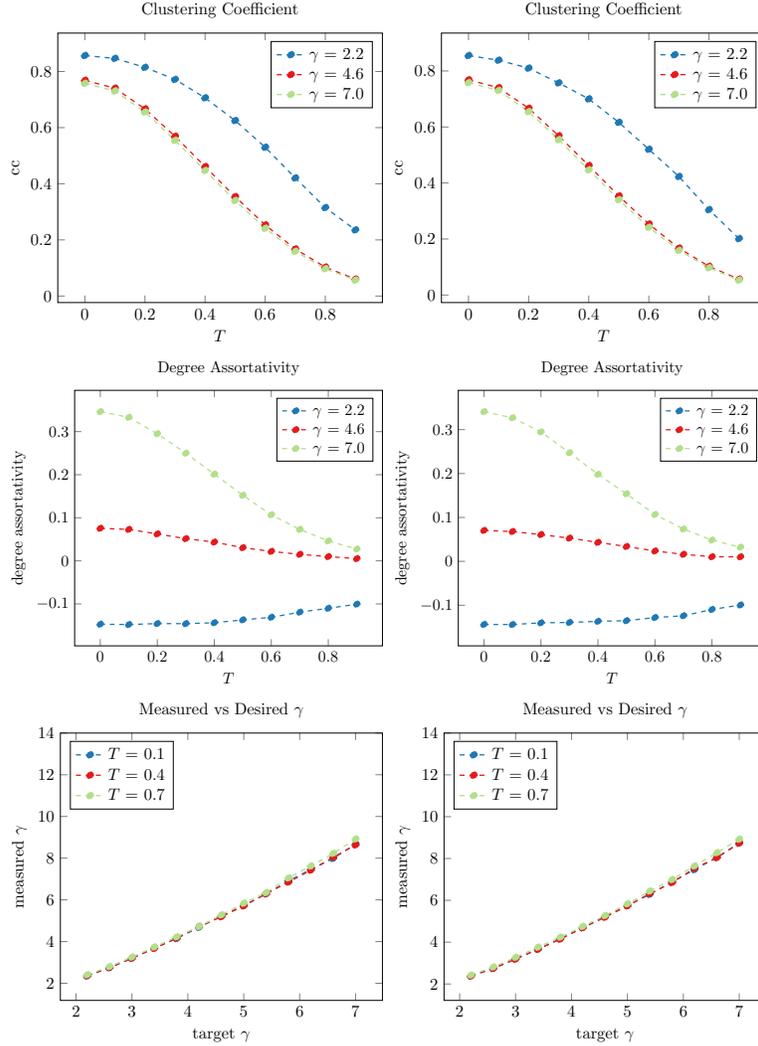


Fig. 6. Comparison of clustering coefficients, degree assortativity and measured vs desired power-law exponent γ . Shown are the implementation of [2] (left) and our implementation (right). The clustering coefficient describes the ratio of closed triangles to triads in a graph. Degree assortativity describes whether vertices have neighbors of similar degree. The degree distribution of random hyperbolic graphs follows a power law, whose exponent γ can be adjusted. In the degree distribution plot, the blue curve is almost always identical to the red curve and thus covered by it. Values are averaged over 80 runs.