

Computing Spatial Distance Histograms Efficiently in Scientific Databases

Yi-Cheng Tu¹, Shaoping Chen², and Sagar Pandit³

¹*Department of Computer Science and Engineering, University of South Florida
4202 E. Fowler Ave., ENB118, Tampa, FL 33620, U.S.A.
ytuc@cse.usf.edu*

²*Department of Mathematics, Wuhan University of Technology
122 Luosi Road, Wuhan, Hubei, 430070, P. R. China
chensp@whut.edu.cn*

³*Department of Physics, University of South Florida
4202 E. Fowler Ave., PHY114, Tampa, FL 33620, U.S.A.
pandit@cas.usf.edu*

Abstract—Particle simulation has become an important research tool in many scientific and engineering fields. Data generated by such simulations impose great challenges to database storage and query processing. One of the queries against particle simulation data, the spatial distance histogram (SDH) query, is the building block of many high-level analytics, and requires quadratic time to compute using a straightforward algorithm. In this paper, we propose a novel algorithm to compute SDH based on a data structure called density map, which can be easily implemented by augmenting a Quad-tree index. We also show the results of rigorous mathematical analysis of the time complexity of the proposed algorithm: our algorithm runs on $\Theta(N^{\frac{3}{2}})$ for two-dimensional data and $\Theta(N^{\frac{5}{3}})$ for three-dimensional data, respectively. We also propose an approximate SDH processing algorithm whose running time is unrelated to the input size N . Experimental results confirm our analysis and show that the approximate SDH algorithm achieves very high accuracy.

I. INTRODUCTION

Many scientific fields have undergone a transition to data/computation intensive science, as the result of automated experimental equipments and computer simulations. In recent years we have witnessed many efforts in building data management tools suitable for processing scientific data [1], [2], [3], [4], [5]. Scientific data imposes great challenges to the design of database management systems that are traditionally optimized toward handling business applications. First, scientific data often come in large volumes, this requires us to rethink the storage, retrieval, and replication techniques in current DBMSs. Second, user accesses to scientific databases are focused on complex high-level analytics and reasoning that go beyond simple aggregate queries. While many types of domain-specific analytical queries are seen in scientific databases, the DBMS should be able to support those that are frequently used as building blocks for more complex analysis.

S. Chen is currently a visiting professor in the Department of Computer Science and Engineering at the University of South Florida (USF). His email at USF is: schen11@cse.usf.edu

However, many of such basic analytical queries need super-linear processing time if handled in a straightforward way, as they are handled in current scientific databases. In this paper, we report our efforts to design efficient algorithms for a type of query that are extremely important in the analysis of **particle simulation data**.

Particle simulations are computer simulations in which the basic components of large systems (e.g., atoms, molecules, stars, galaxies ...) are treated as classical entities that interact for certain duration under postulated empirical forces. For example, molecular simulations (MS) explore relationship between molecular structure, movement and function. These techniques are primarily applicable in modeling of complex chemical and biological systems that are beyond the scope of theoretical models. MS are most frequently used in material sciences, biomedical sciences, and biophysics, motivated by a wide range of applications. In astrophysics, the N-body simulations are predominantly used to describe large scale celestial structure formation [6], [7], [8], [9]. Similar to MS in applicability and simulation techniques, the N-body simulation comes with even larger scales in terms of total number of particles simulated.

Results of particle simulations form large datasets of particle configurations. Typically, these configurations store information about the particle types, their coordinates and velocities - the same type of data we have seen in spatial-temporal databases [10]. While snapshots of configurations are interesting, quantitative structural analysis of inter-atomic structures are the mainstream tasks in data analysis. This requires the calculation of statistical properties or functions of particle coordinates [6]. Of special interest to scientists are those quantities that require coordinates of two particles simultaneously. In their brute force form these quantities require $O(N^2)$ computations for N particles [7]. In this paper, we focus on one such analytical query: the **Spatial Distance Histogram (SDH)** query, which asks for a histogram of the distances of all pairs of particles in the simulated system.

A. Motivation

The SDH is a fundamental tool in the validation and analysis of particle simulation data. It serves as the main building block of a series of critical quantities to describe a physical system. Specifically, SDH is a direct estimation of a continuous statistical distribution function called *radial distribution functions* (RDF) [6], [11], [12]. The RDF is defined as

$$g(r) = \frac{N(r)}{4\pi r^2 \delta r \rho} \quad (1)$$

where $N(r)$ is the number of atoms in the shell between r and $r + \delta r$ around any particle, ρ is the average density of particles in the whole system, and $4\pi r^2 \delta r$ is the volume of the shell. The RDF can be viewed as a normalized SDH.

The RDF is of great importance in computation of thermodynamic quantities about the system. Some of the important quantities like total pressure,

$$p = \rho kT - \frac{2\pi}{3} \rho^2 \int dr r^3 u'(r) g(r, \rho, T)$$

and energy

$$\frac{E}{NkT} = \frac{3}{2} + \frac{\rho}{2kT} \int dr 4\pi r^2 u(r) g(r, \rho, T)$$

cannot be calculated without $g(r)$. For mono-atomic systems, the RDF can also be directly related to the structure factor of the system [13], via

$$S(k) = 1 + \frac{4\pi\rho}{k} \int_0^\infty (g(r) - 1) r \sin(kr) dr.$$

We skip the definitions of all notations in the above formulae, as the purpose is to show the importance of SDH in particle simulations. In current solutions, we have to calculate distances between all pairs of particles and put the distances into bins with a user-specified width, as done in state-of-the-art simulation data analysis software packages [?], [12]. MS or N-body techniques generally consist of large number of particles. For example, the Virgo consortium has accomplished a simulation containing 10 billion particles to study the formation of galaxies and quasars [14]. MS systems also hold up to millions of atoms. This kind of scale prohibits the analysis of large datasets following the brute-force approach. From a database viewpoint, it would be desirable to make SDH a basic query type with the support of scalable algorithms.

B. Contributions and roadmap

We claim the following contributions via this work:

1. We propose an innovative algorithm to solve the SDH problem based on a Quadtree-like data structure we call *density map*;
2. We accomplish rigorous performance analysis of our algorithm and prove its time complexity to be $\Theta(N^{\frac{3}{2}})$ and $\Theta(N^{\frac{5}{3}})$ for 2D and 3D data, respectively;
3. Our analytical results on the algorithm gives rise to an approximate SDH solution whose time complexity is independent to the size of the dataset. In practice, this algorithm computes SDH with very low error rates.

We continue this paper by formally defining the SDH problem and listing important notations in Section II; we introduce our SDH processing algorithm in Section III; performance analysis of our algorithm is sketched in Section IV; we discuss an approximate SDH solution in Section V; Section VI is dedicated to the evaluation of our algorithms; we briefly survey related work in Section VII, and conclude this paper by Section VIII.

II. PROBLEM STATEMENT AND LIST OF NOTATIONS

The SDH problem can be defined as follows: given the coordinates of N points in space, we are to compute the counts of point-to-point distances that fall into a series of l ranges in the \mathbb{R} domain: $[r_0, r_1), [r_1, r_2), [r_2, r_3), \dots, [r_{l-1}, r_l]$. A range $[r_i, r_{i+1})$ in such series is called a *bucket*, and the span of the range $r_{i+1} - r_i$ is called the *width* of the bucket. In this paper, we focus our discussions on the case of *standard SDH query* where all buckets have the same width p and $r_0 = 0$, which gives the following series of buckets: $[0, p), [p, 2p), \dots, [(l-1)p, lp]$. Generally, the boundary of the last bucket lp is set to be the maximum distance of any pair of points. Although almost all scientific data analysis only require the computation of standard SDH queries, our solutions can be easily extended to handle histograms with non-uniform bucket width and/or arbitrary values of r_0 and r_l .¹ The SDH is basically a series of non-negative integers $\mathbf{h} = (h_1, h_2, \dots, h_l)$ where h_i ($0 < i \leq l$) is the number of pairs of points whose distances are within the bucket $[(i-1)p, ip)$.

In Table I, we list the notations that are used throughout this paper. Note that symbols defined and referenced in a local context are not listed here.

TABLE I
SYMBOLS AND NOTATIONS.

Symbol	Definition
p	width of histogram buckets
l	total number of histogram buckets
\mathbf{h}	the histogram with elements h_i ($0 < i \leq l$)
N	total number of particles in data
i	an index symbol for any series
DM_i	the i -th level density map
d	number of dimensions of data
δ	side length of a cell
S	area of a region in 2D space
ϵ	error bound for the approximate algorithm
H	total level of density maps, i.e., tree height

III. OUR APPROACH

A. Overview

In processing SDH using the naive approach, the difficulty comes from the fact that the distance of any pair of points

¹The only complication of non-uniform bucket width is that, given a distance value, we need $O(\log l)$ time to locate the bucket instead of constant time for equal bucket width.

is calculated to determine which bucket a pair belongs to. An important observation here is: a histogram bucket always has a non-zero width. Given a pair of points, their bucket membership could be determined if we only know a range that the distance belongs to and this range is contained in a histogram bucket. With the bucket width p increases (i.e., user sends in a coarser SDH query), the chance that any such range with a fixed span will fall into a bucket also increases. In other words, we need to save time in our algorithm by calculating point-to-point distances approximately.

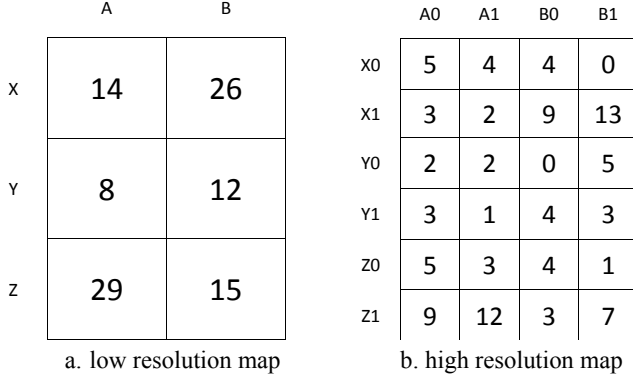


Fig. 1. Two density maps of different resolutions.

The central idea of our approach is a conceptual data structure called *density map*. For a 3D space, a density map is essentially a 3D grid that divides the simulated space into cubes of equal volumes. For a 2D space, it consists of squares of equal size. From now on, we use 2D data and grids to elaborate and illustrate our ideas unless specified otherwise. Note that extending our discussions to 3D data/space would be straightforward. In every cell of the grid, we record the number of particles that are located in the space represented by that cell as well as the four coordinates that determine the exact boundary of the cell in space. The reciprocal of the cell size in a density map is called the *resolution* of the density map. In order to process SDH, we build a series of density maps with different resolutions. We organize the array of density maps in a way such that the resolution of a density map is always doubled as compared to the previous one in the series. Consequently, any cell in a density map is divided into exactly four (eight for a 3D space) disjoint cells in the next density map. In Figure 1, we illustrate two density maps with different resolutions built for the same dataset. For example, the simulated space is divided into six cells in Fig. 1a, each with side length 2, and cell XA has 14 particles in it. The next density map is shown in Fig. 1b: cells are of side length 1 and XA is divided into 4 cells on this map: X0A0, X0A1, X1A0, and X1A1. A natural way to organize the density maps is to connect all cells in a quad-tree. We will elaborate more on the implementation of the density maps in Section III-C.

Algorithm DM-SDH

Inputs: all data points, density maps built beforehand, and bucket width p

Output: an array of counts \mathbf{h}

```

1 initialize all elements in  $\mathbf{h}$  to 0
2 find the first density map  $DM_i$  whose cells have
   diagonal length  $k \leq p$ 
3 for all cells in  $DM_i$ 
4   do  $n \leftarrow$  number of particles in the cell
5      $h_1 \leftarrow h_1 + \frac{1}{2}n(n-1)$ 
6 for any two cells  $m_j$  and  $m_k$  in  $DM_i$ 
7   do RESOLVETWOCELLS ( $m_j, m_k$ )
8 return  $\mathbf{h}$ 

```

Procedure RESOLVETWOCELLS (m_1, m_2)

```

0 check if  $m_1$  and  $m_2$  are resolvable
1 if  $m_1$  and  $m_2$  are resolvable
2   then  $i \leftarrow$  index of the bucket  $m_1$  and  $m_2$  resolve into
3      $n_1 \leftarrow$  number of particles in  $m_1$ 
4      $n_2 \leftarrow$  number of particles in  $m_2$ 
5      $h_i \leftarrow h_i + n_1 n_2$ 
6 else if  $m_1$  and  $m_2$  are on the last density
   map (i.e., the one with highest resolution)
7   for each particle A in  $m_1$ 
8     for each particle B in  $m_2$ 
9       do  $f \leftarrow$  distance between A and B
10       $i \leftarrow$  the bucket  $f$  falls into
11       $h_i \leftarrow h_i + 1$ 
12 else
13    $DM' \leftarrow$  next density map with higher resolution
14   for each partition  $m'_1$  of  $m_1$  on  $DM'$ 
15     for each partition  $m'_2$  of  $m_2$  on  $DM'$ 
16       do RESOLVETWOCELLS ( $m'_1, m'_2$ )

```

Fig. 2. The density-map-based SDH algorithm.

B. The Density-Map-based SDH (DM-SDH) algorithm

In this section, we describe how to use the density maps to process the SDH query. The details of the algorithm are shown in Fig. 2. The core of the algorithm is a procedure named RESOLVETWOCELLS, which is given as inputs a pair of cells m_1 and m_2 on the same density map.

In RESOLVETWOCELLS, we first compute the minimum and maximum distances between any particle from m_1 and any one from m_2 (line 1). This can be accomplished in constant time given the corner coordinates of two cells stored in the density map (only three cases are possible, as shown in Fig. 3). When the minimum and maximum distances between m_1 and m_2 fall into the same histogram bucket i , we say these two cells are *resolvable* on this density map, and they *resolve* into bucket i . If this happens, the histogram is updated (lines 2 - 5) by incrementing the count of the specific bucket i by $n_1 n_2$ where n_1, n_2 are the particle counts in cells m_1 and

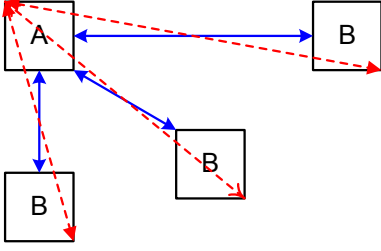


Fig. 3. Three scenarios to consider when computing the minimum and maximum distance between two cells A and B, with solid (dotted) line representing minimum (maximum) distance in each case.

m_2 , respectively. If the two cells do not resolve on the current density map, we move to a density map with higher (doubled) resolution and repeat the previous step. However, on this new density map, we try resolving all four partitions of m_1 with all those of m_2 (lines 12 - 16). In other words, there are $4 \times 4 = 16$ recursive calls to RESOLVETWOCELLS if m_1 and m_2 are not resolvable on the current density map. In another scenario where m_1 and m_2 are not resolvable yet no more density maps are available, we have to calculate the distances of all particles in the non-resolvable cells (lines 6 - 11). The DM-SDH algorithm starts (line 2) at the first density map DM_i whose cell diagonal length is smaller than the bucket width p (i.e., cell side length $\delta \leq \frac{p}{\sqrt{2}}$). It is easy to see that no pairs of cells are resolvable in density maps with resolution lower than that of DM_i . Within each cell on M_i , we are sure that any intra-cell point-to-point distance is smaller than p thus all such distances are counted into the first bucket with range $[0, p)$ (lines 3 - 5). The algorithm proceeds by resolving inter-cell distances (i.e., calling RESOLVETWOCELLS) for all pairs of cells in M (lines 6 - 7).

Clearly, the main idea behind our algorithm is to avoid computing any point-to-point distances. By only considering atom counts in the density map cells, we are able to process multiple point-to-point distances between two cells in one shot. This translates into significant improvements over the brute-force histogram construction approach.

A case study. Let us study an example by revisiting Fig. 1. Suppose the query asks for SDH with a bucket width of 3 (i.e., histogram buckets are $[0, 3)$, $[3, 6)$, $[6, 9)$, \dots) and we start on the low-resolution map in Fig. 1a. First, since all particles in XA are within a distance $2\sqrt{2} < 3$, we can safely increase the count of the first bucket (with range 0-3) by $14 \times (14 - 1)/2 = 91$, and we do this for all other cells in Fig. 1a. Then we try to resolve XA with each and every other cell in the same density map, e.g., cell ZB. However, we cannot draw any conclusions as the distances between a particle in XA and one in ZB are within the range of $[2, \sqrt{52}] \approx [2, 7.2111]$, which overlaps with the first and second buckets. In this case, we turn to the next density map in Fig. 1b, in which a cell in Fig. 1a is divided into four smaller cells. We start comparing counts of all XA cells (i.e., X0A0, X0A1, X1A0, and X1A1) with all ZB cells (i.e., Z0B0, Z0B1, Z1B0, and Z1B1). Out

TABLE II
INTER-CELL DISTANCE RANGES ON DENSITY MAP SHOWN IN FIG. 1B.
RANGES MARKED WITH * ARE RESOLVABLE INTO BUCKETS OF WIDTH 3.

ZB cells	XA cells			
	Z0B0	Z0B1	Z1B0	Z1B1
X0A0	$[\sqrt{10}, \sqrt{34}]^*$	$[\sqrt{13}, \sqrt{41}]$	$[\sqrt{4}, \sqrt{45}]$	$[\sqrt{20}, \sqrt{52}]$
X0A1	$[3, \sqrt{29}]^*$	$[\sqrt{10}, \sqrt{34}]^*$	$[\sqrt{4}, \sqrt{40}]$	$[\sqrt{17}, \sqrt{45}]$
X1A0	$[\sqrt{5}, \sqrt{25}]$	$[\sqrt{8}, \sqrt{32}]$	$[\sqrt{10}, \sqrt{34}]^*$	$[\sqrt{13}, \sqrt{41}]$
X1A1	$[2, \sqrt{20}]$	$[\sqrt{5}, \sqrt{24}]$	$[3, \sqrt{29}]^*$	$[\sqrt{10}, \sqrt{34}]^*$

of the 16 pairs of cells, six can be resolved (Table II). For example, since the distances between any particle in X0A0 and any one in Z0B0 are within $[\sqrt{10}, \sqrt{34}] \approx [3.162, 5.831]$, we increment the count of the second bucket (with range $[3, 6)$) in the histogram by $5 \times 4 = 20$. For those that are not resolvable, we need to visit a density map with an even higher resolution, or, calculate all the inter-cell point-to-point distances when no such density maps exist. Note that those cells with a zero particle count (e.g., cell Y0B0) can be ignored in this process.

C. Implementation of density maps

In DM-SDH, we assume that there are a series of density maps built beforehand for the dataset. In this section, we describe important details on the implementation and maintenance of the density maps.

1) *Tree structure:* As mentioned earlier, we organize the cells on different density maps into a tree structure, much like the point region (PR) Quad-tree presented in [15]. The nodes in the tree hold the following information:

(p-count, x1, x2, y1, y2, child, p-list, next)

where p-count is the number of particles in the cell, x1 to y2 are the four coordinates that define the region of the cell (for 3D data, we need two more coordinates for the 3rd dimension), child is a pointer to the first child on the next level.² The p-list element is the head of a list of data structures that store the real particle data. Obviously, p-list is meaningful only for leaf nodes of the tree. Unlike a regular Quad-tree, we add a next pointer to chain the sibling nodes together (the order of the four siblings in the list can be arbitrarily determined). Furthermore, for the last of the four siblings, its next pointer is used to point to its cousin. By this, all nodes on the same level are connected - such a connected list essentially forms a density map with a specific resolution. The head of all listed can be stored in an array for the ease of locating the appropriate density map to start the DM-SDH algorithm (line 2, Fig. 2). From now on, we use the phrases “density map” and “tree level”, “cell” and “tree

²A parent pointer could be added to each node to achieve logarithmic data insertion/deletion time. However, we assume the scientific dataset is static (no sporadic insertions, no deletions) therefore a child pointer is sufficient for efficient bulk loading.

node” interchangeably. In building the tree, we use the most straightforward space decomposition approach: 1) the space represented by each node is strictly set to be squares (cubes for 3D space), i.e., we have $|x_1 - x_2| = |y_1 - y_2|$; and 2) we always partition by dividing each dimension into exactly TWO equal segments. In other words, each partitioning process will generate four (eight for 3D space) partitions in the next tree level. Space partition using shapes other than squares (e.g., rectangles, triangles) and/or partitioning into more than four children (e.g., any n^2 ($n \in \mathbb{Z}^+$ and $n > 2$) partitions) would make interesting topics for future research and are beyond the scope of this paper.

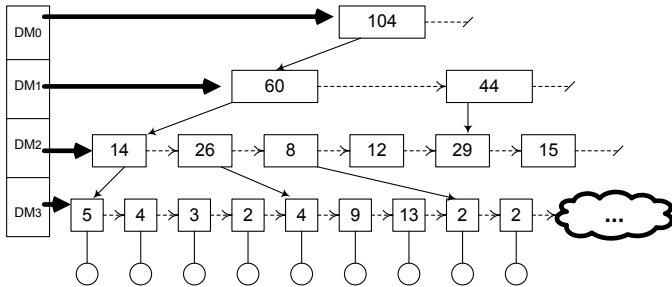


Fig. 4. Tree structure to organize the density maps in Fig. 1. Here we show the p -count (number in each node), next (dotted lines), child (thin solid lines), and p -list (lines connecting to a ball).

The density maps in Fig. 1 can be put into a tree structure as shown in Fig. 4. Due to space limitations, we only show part of the tree. Each node is shown with its p -count field. The root node represent a square with side length 8. The first node on each level is stored in an array from DM_0 to DM_3 , in which DM_2 corresponds to the density map in Fig. 1a, and DM_3 the one in Fig. 1b. In this case, DM_3 has the highest resolution so all DM_3 nodes connect to the data of the particles they contain via p -list.

2) *Tree height.*: To be able to answer SDH queries with different parameters (e.g., bucket width p , subregion of the simulated space), we need to build a series of density maps from the most coarse resolution to the finest. On the coarsest end, we can build a single node map that covers the whole simulated space. The question is from the other end: what should be the highest resolution in the maps? This is a subtle issue: first, given any bucket width p , the percentage of resolvable cells increases with the level of the tree. However, the number of pairs of cells also increases dramatically (i.e., by a factor of 2^d).

Recall that DM-SDH saves our time of processing SDH by resolving cells such that we need not calculate the point-to-point distances one by one. However, when the p -count of a cell decreases, the time we save by resolving that cell also decreases. Imagine a cell with a p -count of 4 or smaller (8 or smaller for 3D data/space), it does not give us any benefit in processing SDH to further partition this cell on the next level: the cost of resolving the partitions could be higher than directly retrieving the particles and calculating distances (lines

7 - 11 in RESOLVETWOCCELLS). Based on this observation, the total level of density maps H is set to be

$$H = \left\lceil \log_{2^d} \frac{N}{\beta} \right\rceil + 1 \quad (2)$$

where d is the number of dimensions and 2^d is essentially the degree of tree nodes (4/8 for 2D/3D data), β is the average number of particles we desire in each leaf node. In practice, we set β to be slightly greater than 4 in 2D (8 for 3D data) since the CPU cost of resolving two cells is higher than computing the distance between two points.

3) *Other issues.*: In addition to the bucket width p , which is the most important parameter, user can attach other conditions to a SDH query. Two common varieties of the regular SDH query are:

1. Compute the SDH of a specific region of the whole simulated space;
2. Compute the SDH of all particles of a specific type (e.g., carbon atoms) in the dataset.

The first variety requires modifications to our algorithm: in RESOLVETWOCCELLS, we need to check if both cells are contained by the query region and add one more case for the recursive call, that is, if the cells are resolvable but at least one of the cells overlaps with, or locates out of, the query region, we still need to go to the next density map. If both cells are out of the query region, nothing needs to be done. In calculating the distances of particles (lines 7 - 11), again, we only consider those particles that are within the query region. The second variety requires more information be stored in the density map cells (i.e., tree nodes): in addition to the p -count field, we keep a list of other counts, one for each possible type of particles in the data. Fortunately, the number of particle types is not very large in the sciences of interest (e.g., about 10 for molecular simulation).

Another piece of information we can store in the tree nodes is the minimum bounding rectangle (MBR) formed by of all the particles in a node. In RESOLVETWOCCELLS, we can use the MBR of the two cells to compute the minimum and maximum point-to-point distances. As compared to the theoretical bounds of the space occupied by a tree node, the MBR will cover a smaller area. Intuitively, the chance of a cell’s being resolvable under a given p increases as the cell shrinks. The use of MBR can thus shorten the running time by making more cells resolvable at a higher level on the tree. The MBR can be easily computed when data points are loaded to the tree, with the storage overhead of four extra coordinates in each node.

In general, we build balanced trees that are exactly H levels in height. One possible optimization is to use Equation (2) as a guideline instead of a strict rule: we can further partition a node if it has a large number of particles in it, giving rise to a unbalanced tree. This can benefit those datasets with highly skewed distribution of particles. With an unbalanced tree, we need to modify the last three lines of code in RESOLVETWOCCELLS to take resolving two cells on different levels of the tree into account.

IV. ANALYSIS OF THE ALGORITHM

The running time of DM-SDH consists of two main parts:

1. the time spent to check if two cells are resolvable (line 0 in RESOLVETWOCELLS, constant time needed for each operation); and
2. distance calculation for data in cells non-resolvable even on the finest density map (lines 7 - 11 in RESOLVETWOCELLS, constant time needed for each distance).

As compared to the brute-force algorithm, we save time by performing operation 1 in hope of handling multiple distances in one shot. However, it is not clear how much overhead this bears. Consider a tree illustrated in Fig. 5 where each level represent a density map but each node represents a pair of cells in the original density map. Given a histogram bucket width p_1 , we start from a density map DM_i with c_i cells. Thus, there are $O(c_i^2)$ entries on the corresponding level of the tree shown in Fig. 5. On the next map DM_{i+1} , there are $4 \times 4 = 16$ times of cell pairs to resolve. However, some of the cells in DM_{i+1} do not need to be considered as their parents are resolved on DM_i . Consider a resolvable entry a in DM_i , the whole subtree rooted at a needs no further consideration, leaving a “hole” on the leaf level. Similarly, if b is a resolvable entry on a lower level DM_j , the time needed for resolving everything in the subtree of b is also saved. However, this subtree is smaller than that of a , leading to less savings of time. From this we can easily see that the running time depends on the bucket width: if we are given another query with bucket width $p_2 < p_1$ such that we will start DM-SDH from level DM_j of the tree, more cell comparisons have to be done, giving rise to longer running time. In analyzing the time complexity of DM-SDH, we are interested in how the running time increases as the total number of particle N increases. Qualitatively, as N increases, the height of the trees also increases (as we fix β in Equation (2)), thus a higher percentage of particle pairs can be resolved in the cells. However, the total number of entries on the leaf level in Fig. 5 also increases (quadratically). Therefore, a quantitatively study on the percentage of resolvable cells on a given level is essential in such analysis.

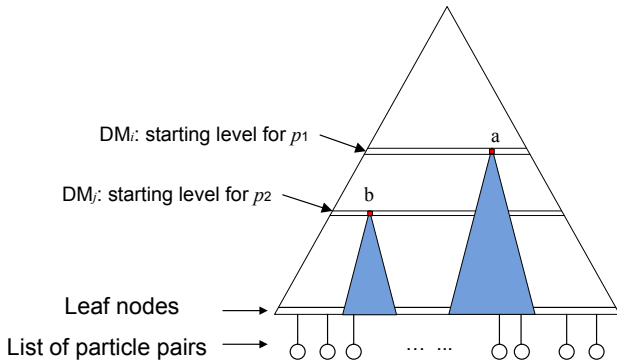


Fig. 5. A conceptual tree structure with each node representing a pair of cells in a density map. Similarly, the data nodes hold pairs of particles.

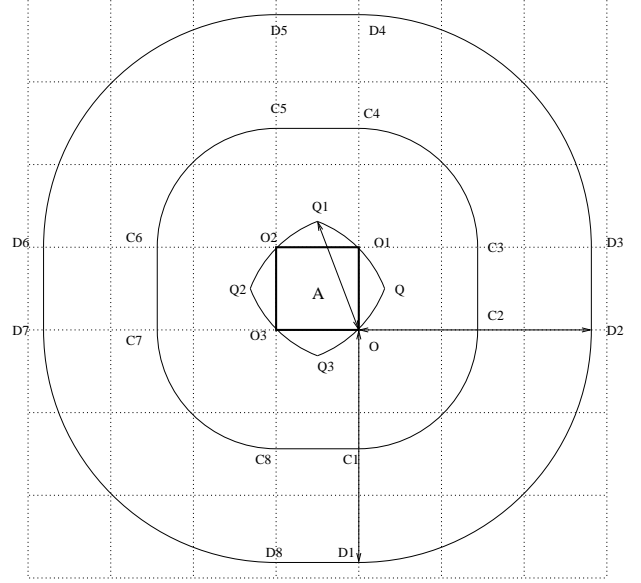


Fig. 6. Boundaries of bucket 1 and bucket 2 regions of cell A, with the bucket width p being exactly $\sqrt{2}\delta$. Here we show Q_1Q_2 , C_1C_2 , and D_1D_2 are arcs all centered at point O.

We have accomplished a quantitative analysis on the performance of our algorithm, which involves non-trivial geometric modeling and algebraic manipulation of the models.

A. Basics of our model

Essentially, our analysis needs to answer the following question: given a cell A on the first density map DM_i , how many particles are contained by those resolvable cells related to A as we visit more and more levels of density maps? Although this has something to do with the spatial distribution of the particles, we start by analyzing how much area are covered by the resolvable cells to simplify the process. To achieve this, we first need to define a theoretical region in which a particle can have distance (to a point in A) that falls into a specific bucket i . We call this region the *bucket i region* of cell A .

In Fig. 6, a cell A is drawn with four corner points O, O_1, O_2 , and O_3 . The side length of A is exactly $\delta = \frac{p}{\sqrt{2}}$. The bucket 1 region of A is bounded by a curve connected by points C_1 to C_8 . This region is drawn as follows: C_1C_2, C_3C_4, C_5C_6 , and C_7C_8 are all arcs of 90 degrees centered at the four corners of cell A and their radii are p ; C_2C_3, C_4C_5, C_6C_7 , and C_8C_1 are line segments. Note that this is a theoretical “maximum” region where a point can resolve with any point in A . It is easy to see that the area of this region is $\pi p^2 + 4p\delta + \delta^2$. Let us continue to consider distances that fall into the second bucket (i.e., $[p, 2p)$). Again, the bucket 2 region of A is of similar shape to the bucket 1 region except the radii of the arcs are $2p$, as drawn in Fig. 6 with a curve connected by points D_1 to D_8 . However, if a point is too close to the center, it may never resolve into bucket 2. These points are contained in a region as follows:

$$g(i) = \begin{cases} (2\pi + 4\sqrt{2} + 1)\delta^2 & i = 1 \\ \left[2\pi i^2 + 4\sqrt{2}i - (i-1)^2(8 \arctan \sqrt{8(i-1)^2 - 1} - 2\pi) + \sqrt{8(i-1)^2 - 1}\right] \delta^2 & i > 1 \end{cases} \quad (3)$$

on each corner point of \mathbf{A} , we draw an arc with radius p on the opposite corner (i.e., arcs QQ_1, Q_1Q_2, Q_2Q_3 , and Q_3Q_4). For any point in this region, its distance to any point in \mathbf{A} is always smaller than p . Therefore, the bucket 2 region should not include this inner region (denoted as region \mathbf{B} hereafter). A more detailed illustration of region \mathbf{B} is shown in Fig. 7.

The area of the bucket 2 region is $\pi(2p)^2 + 8p\delta$ less the area of region \mathbf{B} , which consists of eight identical smaller regions such as $\widehat{Q_1O_2D}$ shown in Fig. 7. To get the area of $\widehat{Q_1O_2D}$, we first need to know the magnitude of the angle $\angle Q_1OO_2$, which can be determined by

$$\begin{aligned} \angle Q_1OO_2 &= \angle Q_1OE - \angle COE \\ &= \arctan \frac{Q_1E}{EO} - \frac{\pi}{4} \\ &= \arctan \frac{\sqrt{p^2 - \left(\frac{\delta}{2}\right)^2}}{\frac{\delta}{2}} - \frac{\pi}{4} \end{aligned}$$

Thus, the area of sector $\widehat{Q_1O_2O}$ is $\frac{1}{2}p^2\angle Q_1OO_2$. The area of region $\widehat{Q_1O_2D}$ can be obtained by the area of this sector less the area of triangles O_2DC and Q_1CO . By this, we get

$$\begin{aligned} S_{\widehat{Q_1O_2D}} &= S_{\widehat{Q_1O_2O}} - S_{\Delta O_2DC} - S_{\Delta Q_1CO} \\ &= \frac{1}{2}p^2 \left[\arctan \frac{\sqrt{p^2 - \left(\frac{\delta}{2}\right)^2}}{\frac{\delta}{2}} - \frac{\pi}{4} \right] - \frac{1}{2} \left(\frac{\delta}{2}\right)^2 \\ &\quad - \frac{1}{2} \left[\sqrt{p^2 - \left(\frac{\delta}{2}\right)^2} - \frac{\delta}{2} \right] \frac{\delta}{2} \\ &= \frac{1}{2}p^2 \left[\arctan \frac{\sqrt{p^2 - \left(\frac{\delta}{2}\right)^2}}{\frac{\delta}{2}} - \frac{\pi}{4} \right] \\ &\quad - \frac{\delta}{4} \sqrt{p^2 - \left(\frac{\delta}{2}\right)^2} \end{aligned}$$

and we have $\pi(2p)^2 + 8p\delta - 8S_{\widehat{Q_1O_2D}} - S_{\mathbf{A}}$ as the area of the bucket 2 region.

The approach to obtain the area of bucket i ($i > 2$) regions is the same as above. For the area of the region formed by the outer boundary, we only need to consider that the arcs in Fig. 7 are of radii ip . The development of a general formula for the area of region \mathbf{B} is trickier. Our efforts lead to the following formula for the bucket i region:

$$g(i) = \begin{cases} \pi p^2 + 4p\delta + \delta^2 & i = 1 \\ \left(\pi ip^2 + 4ip\delta - [8A(i) + B(i)\delta^2]\right) & i \geq 2 \end{cases}$$

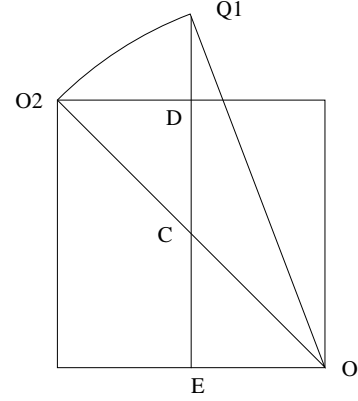


Fig. 7. An illustration on how to compute the area of region $QQ_1Q_2Q_3$ formed by four arcs in Fig. 6. Here we only show half of one of the arcs.

where $B(i) = [2(i-1) - 1]^2 - 1$ and

$$\begin{aligned} A(i) &= \frac{1}{2}[(i-1)p]^2 \left[\arctan \frac{\sqrt{[(i-1)p]^2 - \left(\frac{\delta}{2}\right)^2}}{\frac{\delta}{2}} - \frac{\pi}{4} \right] \\ &\quad - \frac{1}{2} \frac{\delta}{2} \left[\sqrt{[(i-1)p]^2 - \left(\frac{\delta}{2}\right)^2} - \frac{\delta}{2} \right] \\ &\quad - \frac{1}{2} \left[(i-2)\delta + \frac{\delta}{2} \right]^2 \end{aligned}$$

Since we have $p = \sqrt{2}\delta$, the above equation becomes Eq. (3) shown on top of this page.

B. Coverable regions

Eq. (3) gives the area of a theoretical region that contains all particles that could have distance within a given bucket to a given cell \mathbf{A} . Now let us study how much of this region can be resolved in our algorithm under different levels of density maps. We call the region that consists of all resolvable cells the *coverable region*.

1) *Case 1: the first bucket:* Let us start our discussions on the situation of bucket 1. In Fig. 8, we show the coverable regions of three different density map levels: $m = 1$, $m = 2$, and $m = 3$, as represented by blue-colored lines and denoted as \mathbf{A}' in all subgraphs. For $m = 1$, the resolvable cells are only those surrounding \mathbf{A} . All other cells, even those entirely contained by the bucket 1 region, do not resolve with any level 1 subcell of \mathbf{A} . As we increase m , the region \mathbf{A}' grows in area, with its boundary approaching that of the bucket 1 region. To represent the area of \mathbf{A}' , we need to develop a continuous line to approximate its boundary. One critical observation here is: the furthest cells in \mathbf{A}' are those that can resolve with cells on the outer rim of \mathbf{A} . For example, the cell cornered at point

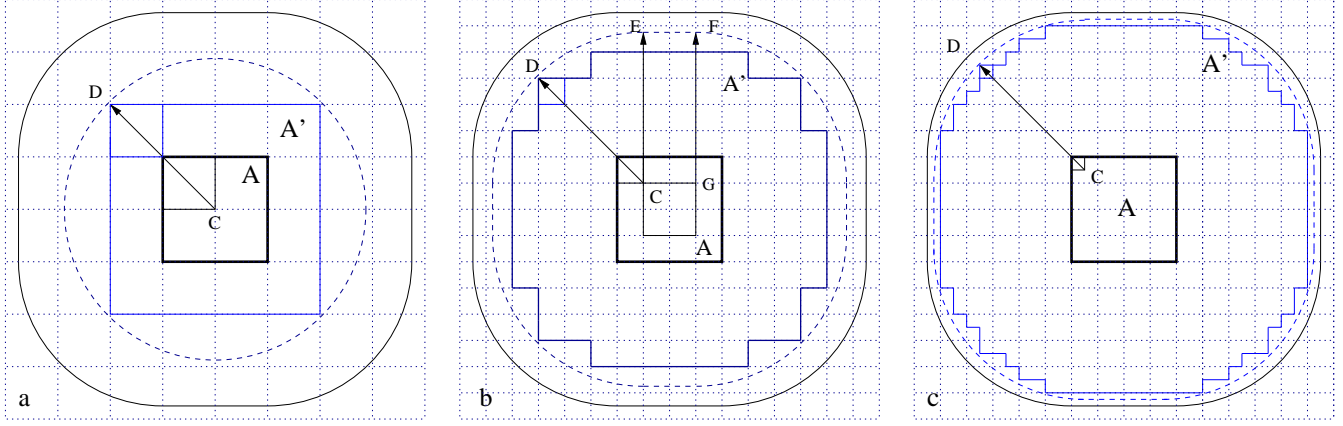


Fig. 8. Actual (solid blue line) and approximated (dotted blue line) coverable regions for bucket 1 under: a. $m = 1$; b. $m = 2$; and c. $m = 3$. Outer solid black lines represent the theoretical bucket 1 region. All arrowed line segments are drawn from the centers to the corresponding arcs with radius p .

D resolves with the cell cornered at point C in A . If we draw a 90-degree arc centered at C , the arc goes through D and all cells on the northwestern corner of A' are bounded by this arc. To approximate the boundary of A' , we can draw such an arc at all four corners of the graph and connect them with line segments (e.g., EF connecting the northwestern and northeastern arcs centered at point G in Fig. 8b), as shown by the blue dotted line. Obviously, this line approaches the theoretical boundary as m increases because the center of the arcs (e.g., point C) move further to the corner points of A as the cells become smaller. Note this line gives rise to an optimistic approximation of A' . In a moment, we will show that this overestimation will not harm our analysis on the running time of DM-SDH. The area of coverable region for bucket 1 at level m can be expressed by the following:

$$S_{A'} = \pi p^2 + 4p \left(\delta - \frac{2\delta}{2^m} \right) + \left(\delta - \frac{2\delta}{2^m} \right)^2 \quad (4)$$

where the first item πp^2 is the area of the four 90-degree sectors centered at point C , the second item is the area of the four rectangles (e.g., $EFGC$ in Fig. 8b) connecting the four sectors. We also need to add the area of the small square (with side CG in Fig. 8b) within cell A , which is given by the last item.

2) *Case 2: the second bucket and beyond:* The cases of buckets beyond the first one are more complicated. First of all, the outer boundary of the bucket i ($i \geq 2$) regions can be approximated using the same techniques we introduced for bucket 1 (Section IV-B.1). Therefore, we can use the following generalized form of Eq. (4) to quantify the region formed by the outer boundaries only.

$$S_{out}(i) = \pi(ip)^2 + 4ip \left(\delta - \frac{2\delta}{2^m} \right) + \left(\delta - \frac{2\delta}{2^m} \right)^2 \quad (5)$$

However, we also need to disregard the cells that lie in the inner boundary (e.g., those within or near region B). To quantify the area of the region contained by the inner

boundary, we need to consider the cases of $m = 1$ and $m > 1$ separately.

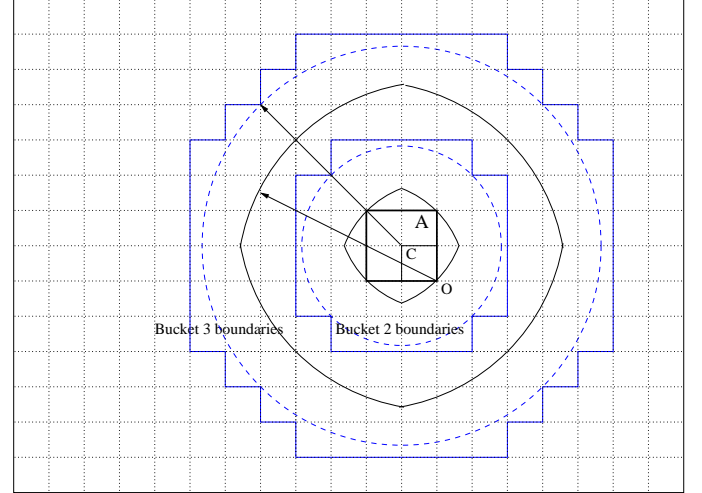


Fig. 9. Inner boundaries of the coverable regions of buckets 2 and 3 under $m = 1$. All arrowed line segments are of length $2p$.

Let us first study the case of $m = 1$. Fig. 9 shows examples with $m = 1$ with respect to the second and the third buckets. It is easy to see that any cell that contains a segment of the theoretical region B boundary will not resolve into bucket i because they can only resolve into bucket $i - 1$. Furthermore, there are more cells that resolve into neither bucket $i - 1$ nor bucket i . Here our task is to find a boundary to separate those $m = 1$ cells that can resolve into bucket i with any subcell in A and those that cannot. Such boundaries for buckets 2 and 3 are shown in Fig. 9 as solid blue lines. The boundary can be generated as follows: on each quadrant (e.g., northwest) of cell A , we draw an arc (dotted blue line) centered at the corner point C of the furthest (e.g., southeast) subcell of A with radius $(i - 1)p$. Any cell that contains a segment of this arc cannot resolve into bucket i (because they are too close to A) but the cells beyond this line can. Therefore, we can also

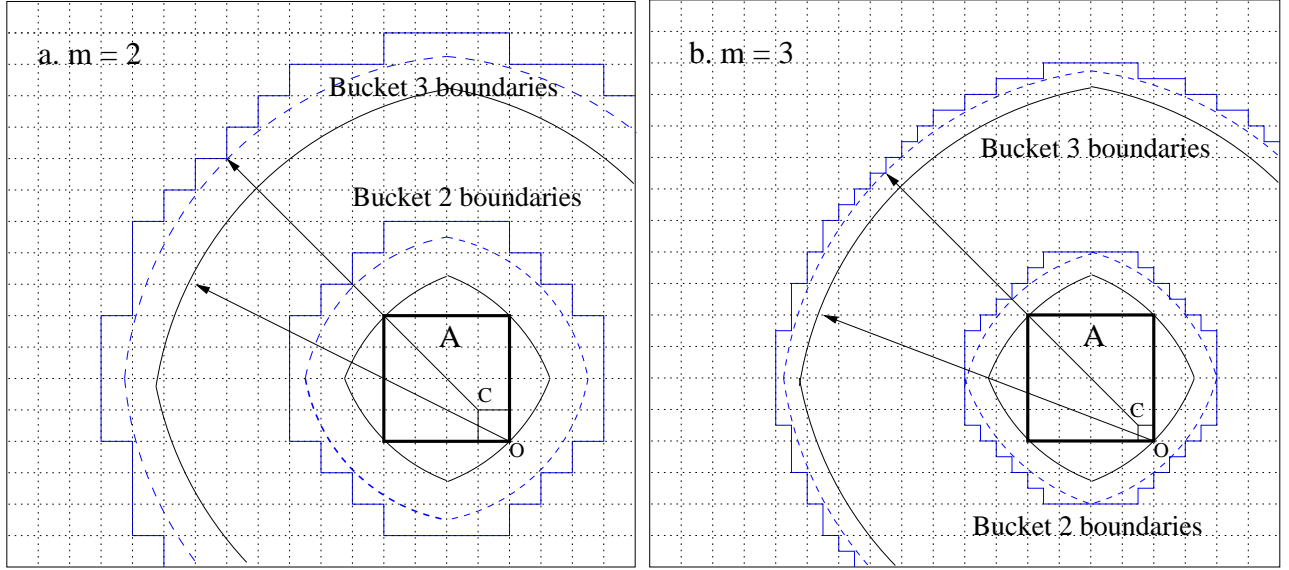


Fig. 10. Inner boundaries of the coverable regions of buckets 2 and 3 under $m = 2$ and $m = 3$. All arrowed line segments are of length $2p$.

use these arcs to approximate the zigzagged real boundaries. Let us denote the region bounded by this approximate curve as region B' . For $m = 1$, the arcs on all four quadrants share the same center C therefore they form a circle as region B' . The radii of the circles are exactly $(i - 1)p$ for bucket i . Note that this, again, could give rise to an optimistic approximation of the area of coverable regions. Therefore, the area of the coverable region for $m = 1$ and $n \geq 2$ is:

$$S_{A'} = \pi(ip)^2 - \pi[(i - 1)p]^2 \quad (6)$$

where the first item is the area of the region formed by the approximated outer boundary, which is given as a special case of Eq. (5) for $m = 1$ and happens to be a circle; and the second item is that of the region formed by the approximated inner boundary (i.e., region B').

For the case of $m > 1$, we can use the same technique described for the case of $m = 1$ to generate the curves to form region B' . However, these curves are no longer a series of circles. In Fig .10, we can find such curves for buckets 2 and 3 under m values of 2 and 3. As the four arcs on different quadrants no longer share the same center, the region B' boundaries (dotted blue lines) are of similar shapes to the theoretical region B boundaries (solid black lines). From the graphs, it is easy to see that the approximated curve fits the actual boundary better as m increases. Here we skip the formal proof as it is straightforward. Furthermore, it also converges to the region B boundary when m gets bigger. This is because the centers of the two arcs (with the same radii), points C and O , become closer and closer when the cell size decreases (as a result of the increase of m).

The area of region B' can be computed in the same way as that of region B , with the help of an illustration in Fig. 11.

First, we get the magnitude of angle BCD by

$$\begin{aligned} \angle BCD &= \angle DCE - \angle FCE \\ &= \arctan \frac{DE}{EC} - \frac{\pi}{4} \\ &= \arctan \frac{\sqrt{[(i - 1)p]^2 - \left(\frac{\delta}{2} - \frac{\delta}{2^m}\right)^2}}{\frac{\delta}{2} - \frac{\delta}{2^m}} - \frac{\pi}{4} \end{aligned}$$

The area of the sector \widehat{BDC} is $\frac{1}{2}[(i - 1)p]^2 \angle BCD$, and the area of the region \widehat{BDGF} is

$$\begin{aligned} S_{\widehat{BDGF}} &= S_{\widehat{BDC}} - S_{\triangle DHC} - S_{\triangle FGH} \\ &= \frac{1}{2}[(i - 1)p]^2 \angle BCD - \frac{1}{2}EC(DE - HE) - \frac{\delta^2}{8} \\ &= \frac{1}{2}[(i - 1)p]^2 \left[\arctan \frac{\sqrt{[(i - 1)p]^2 - \delta^2 \theta_m^2}}{\delta \theta_m} - \frac{\pi}{4} \right] \\ &\quad - \frac{\delta}{2} \theta_m \left[\sqrt{[(i - 1)p]^2 - (\delta \theta_m)^2} - \delta \theta_m \right] - \frac{\delta^2}{8} \end{aligned}$$

where we have $\theta_m = \frac{1}{2} - \frac{1}{2^m}$ for convenience.

Finally, we get the area of the coverable region for $i \geq 2, m > 1$ as

$$\begin{aligned} S_{A'} &= S_{out}(i) - 8S_{\widehat{BDGF}} - S_A \\ &= \pi(ip)^2 + 4ip \left(\delta - \frac{2\delta}{2^m} \right) + \left(\delta - \frac{2\delta}{2^m} \right)^2 \\ &\quad - 4[(i - 1)p]^2 \left[\arctan \frac{\sqrt{[(i - 1)p]^2 - \delta^2 \theta_m^2}}{\delta \theta_m} - \frac{\pi}{4} \right] \\ &\quad + 4\delta \theta_m \left[\sqrt{[(i - 1)p]^2 - (\delta \theta_m)^2} - \delta \theta_m \right] \quad (7) \end{aligned}$$

$$f(i, m) = \begin{cases} \left[2\pi + 4\sqrt{2} + 1 - (8\sqrt{2} + 4)\frac{1}{2^m} + \frac{4}{2^{2m}} \right] \delta^2 & i = 1, m \geq 1 \\ [2\pi(2i - 1)] \delta^2 & i \geq 2, m = 1 \\ \left\{ 2\pi i^2 + 4\sqrt{2}i - (8\sqrt{2}i + 4)\frac{1}{2^m} + \frac{4}{2^{2m}} - 8 \left[(i - 1)^2 \left(\arctan \frac{\gamma_m}{\theta_m} - \frac{\pi}{4} \right) - \frac{1}{2}\theta_m(\gamma_m - \theta_m) \right] + 1 \right\} \delta^2 & i \geq 2, m > 1 \end{cases} \quad (8)$$

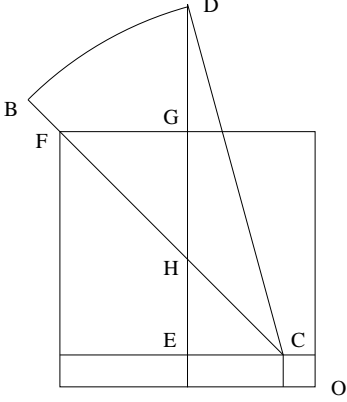


Fig. 11. An illustration on how to compute the area of region formed by four arcs in Fig. 10. Here we only show half of one of the arcs.

In summary, let us denote the area of the coverable region A' under different i and m values as $f(i, m)$. By combining and simplifying Equations 4, 6, and 7 (considering $p = \sqrt{2}\delta$), we get Equation (8), in which $\gamma_m = \sqrt{2(i-1)^2 - \theta_m^2}$.

C. Covering factor

In this section, we give a quantitative analysis on the relationship between $f(i, m)$ and the area of the theoretical region $g(i)$. For that purpose, given any density map level m , we define the *covering factor* $c(m)$ as the ratio of the total area of the coverable regions to that of the theoretical bucket i regions for all i . However, the quantity that is more related to our analysis is the *noncovering factor* $\alpha(m) = 1 - c(m)$. Specifically, we have

$$\alpha(m) = \frac{\sum_{i=1}^l [g(i) - f(i, m)]}{\sum_{i=1}^l g(i)} \quad (9)$$

The quantity $\alpha(m)$ is important in that it can directly tell how many cell pairs are resolvable on a given density map level (as the total number of cell pairs is always known for each level). Before investigating the features of $\alpha(m)$, let us define two relevant quantities, the total area of bucket regions for all buckets G , and that of all coverable regions F . Being summations over all buckets of $g(i)$ and $f(i, m)$, they can be expressed as functions of the total bucket number l . We also remove the common factor δ^2 from both Eq. (3) and Eq. (8)

for the convenience of displaying equations. First, we have

$$\begin{aligned} G(l) &= \frac{\sum_{i=1}^l g(i)}{\delta^2} \\ &= 1 + \sum_{i=1}^l (2\pi i^2 + 4\sqrt{2}) \\ &\quad - \sum_{i=2}^l [(i-1)^2 (8 \arctan \sigma_i - 2\pi) - \sigma_i] \\ &= 1 + \frac{2}{3}l (3\sqrt{2} + 3\sqrt{2}l + \pi + 2l^2\pi) \\ &\quad - \sum_{i=2}^l [(i-1)^2 (8 \arctan \sigma_i - 2\pi) - \sigma_i] \end{aligned} \quad (10)$$

where $\sigma_i = \sqrt{8(i-1)^2 - 1}$. The area of total coverable regions is considered in two cases. For $m = 1$, we get

$$\begin{aligned} F(l, 1) &= \frac{\sum_{i=1}^l f(i, 1)}{\delta^2} \\ &= 2\pi + 2\pi \sum_{i=2}^l (2i - 1) = 2\pi l^2 \end{aligned} \quad (11)$$

and for $m > 1$, we have the following formula:

$$\begin{aligned} F(l, m) &= \frac{\sum_{i=1}^l f(i, m)}{\delta^2} \\ &= 2^{2-2m} - 2^{2-m} + 1 + 2\sqrt{2}l - 2^{\frac{5}{2}-m}l \\ &\quad + 2\sqrt{2}l^2 - 2^{\frac{5}{2}-m}l^2 + \frac{3}{2}l\pi + \frac{4}{3}l^3\pi \\ &\quad - 8 \sum_{i=2}^l (i-1)^2 \arctan \frac{\sqrt{2(i-1)^2 - \theta_m^2}}{\theta_m} \\ &\quad + 4 \sum_{i=2}^l \theta_m \sqrt{2(i-1)^2 - \theta_m^2} \end{aligned} \quad (12)$$

With the above definitions, we develop the most important result in our analysis in the following lemma.

Lemma 1: For any given standard SDH query with bucket width p , let DM_i be the first density map our DM-SDH algorithm starts running, and $\alpha(m)$ be the noncovering factor of a density map that lies m levels below DM_i (i.e., map DM_{i+m}). We have

$$\lim_{p \rightarrow 0} \frac{\alpha(m+1)}{\alpha(m)} = \frac{1}{2}.$$

Proof: From Eq. (9), we easily get

$$\frac{\alpha(m+1)}{\alpha(m)} = \frac{G(l) - F(l, m+1)}{G(l) - F(l, m)}$$

Plugging Eq. (10), Eq. (11), and Eq. (12) into the above formula, we get $\frac{\alpha(m+1)}{\alpha(m)} = \frac{A(m)}{B(m)}$ where

$$\begin{aligned} A(m) &= \frac{2}{2^m} - \frac{1}{4^m} + \frac{2^{\frac{3}{2}}}{2^m}(l+l^2) + \sum_{i=2}^l \sqrt{8(i-1)^2 - 1} \\ &\quad - 4 \sum_{i=2}^l \theta_{m+1} \sqrt{2(i-1)^2 - \theta_{m+1}^2} \\ &\quad + 8 \sum_{i=2}^l (i-1)^2 \arctan \frac{\sqrt{8(i-1)^2 - \theta_{m+1}^2}}{\theta_{m+1}} \\ &\quad - 8 \sum_{i=2}^l (i-1)^2 \arctan \sqrt{8(i-1)^2 - 1} \end{aligned} \quad (13)$$

and

$$\begin{aligned} B(m) &= \frac{4}{2^m} - \frac{4}{4^m} + \frac{2^{\frac{5}{2}}}{2^m}(l+l^2) + \sum_{i=2}^l \sqrt{8(i-1)^2 - 1} \\ &\quad - 4 \sum_{i=2}^l \theta_m \sqrt{2(i-1)^2 - \theta_m^2} \\ &\quad + 8 \sum_{i=2}^l (i-1)^2 \arctan \frac{\sqrt{8(i-1)^2 - \theta_m^2}}{\theta_m} \\ &\quad - 8 \sum_{i=2}^l (i-1)^2 \arctan \sqrt{8(i-1)^2 - 1}, \end{aligned} \quad (14)$$

in which $\theta_m = \frac{1}{2} - \frac{1}{2^m}$ and $\theta_{m+1} = \frac{1}{2} - \frac{1}{2^{m+1}}$.

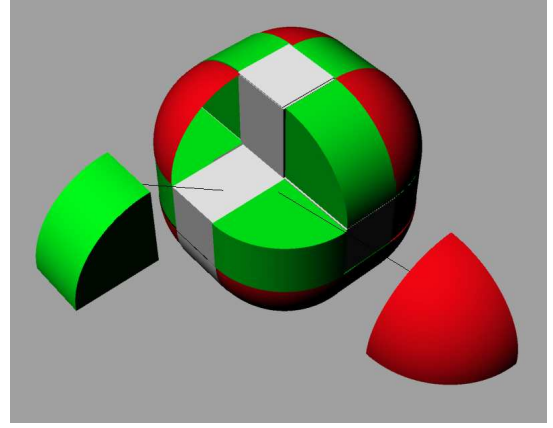
The case of $p \rightarrow 0$ is equivalent to $l \rightarrow \infty$. Despite their formidable length and complexity, $A(m)$ and $B(m)$ have the following feature

$$\lim_{l \rightarrow \infty} \frac{A(m)}{B(m)} = \frac{1}{2} \quad (15)$$

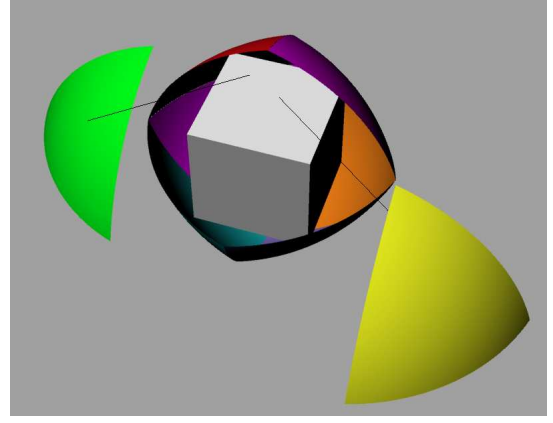
and this concludes the proof. More details on derivation of Eq. (15) can be found in Appendix I. ■

What Lemma 1 tells us is: the chance that any pair of cells is not resolvable decreases by half with the density map level increases by one. In other words, for a pair of non-resolvable cells on DM_i (or any level lower than DM_i), among the 16 pairs of subcells on the next level, we expect $16 \times 0.5 = 8$ pairs to be resolvable. One thing to point out is: Lemma 1 not only works well for large l (i.e., smaller p), it **quickly converges even when l is reasonably small**. This can be verified by our numerical results shown in Table III. Furthermore, the above result is also true for **3D data**, although we can only give numerical results due to complicated formulae developed in the 3D analysis (Section IV-D).

We have mentioned that our analysis is done based on an overestimation of the coverable regions on each density map,



(a) Outer boundary of the bucket 1 region.



(b) Inner boundary of the bucket 2 region.

Fig. 12. Geometric structure of the bucket 1/2 regions for 3D data.

and the estimation error decreases as m increases. Relate this to Lemma 1, we have an underestimated non-covering factor α on each level. Since the estimation is more accurate on larger m , the real ratio of $\alpha(m+1)$ to $\alpha(m)$ can only be smaller than the one given by Lemma 1. And the fact that $\frac{1}{2}$ being the upper bound has positive effects on the evaluation of the time complexity of DM-SDH: the result shown in Section IV-E also becomes an upper bound.

D. 3D analysis

The strategies used to achieve the above analysis can be extended to 3D. The outer and inner boundaries of bucket i regions are illustrated in Fig. 12. The analysis should be based on the volume of relevant regions surrounding a cube \mathbf{A} with side length δ . The bucket 1 region (Fig.12(a)) of \mathbf{A} consists of the following components: 1) quarter cylinders (green) with length δ and radius $p = \sqrt{3}\delta$; 2) one-eighth of a sphere (red) with radius p ; 3) cuboids (white) with dimensions δ , δ , and p ; and 4) cube \mathbf{A} itself (not shown). There are eight pieces of each of the first three items. The inner boundary (region \mathbf{B}) of the bucket 2 region (Fig. 12(b)) consists of eight identical portions of a spherical surface centered at the opposite corner of \mathbf{A} with radius p . Note that the projection of these regions on 2D are exactly those found in Fig. 6. Again, the shape of

TABLE III

VALUES OF $\alpha(m+1)/\alpha(m)$ OF 2D DATA UNDER DIFFERENT VALUES OF m AND l . COMPUTED WITH MATHEMATICA 6.0. PRECISION UP TO THE 6TH DIGIT AFTER DECIMAL POINT.

Map levels	Total Number of Buckets (l)							
	2	4	8	16	32	64	128	256
m=1	0.508709	0.501837	0.50037	0.50007	0.500012	0.500002	0.5	0.5
m=2	0.503786	0.500685	0.500103	0.500009	0.499998	0.499999	0.499999	0.5
m=3	0.501749	0.500282	0.500031	0.499998	0.499997	0.499999	0.5	0.5
m=4	0.500838	0.500126	0.50001	0.499997	0.499998	0.499999	0.5	0.5
m=5	0.50041	0.500059	0.500004	0.499998	0.499999	0.5	0.5	0.5
m=6	0.500203	0.500029	0.500002	0.499999	0.499999	0.5	0.5	0.5
m=7	0.500101	0.500014	0.500001	0.499999	0.5	0.5	0.5	0.5
m=8	0.50005	0.500007	0.5	0.5	0.5	0.5	0.5	0.5
m=9	0.500012	0.500003	0.5	0.5	0.5	0.5	0.5	0.5
m=10	0.500025	0.500002	0.5	0.5	0.5	0.5	0.5	0.5

TABLE IV

VALUES OF $\alpha(m+1)/\alpha(m)$ OF 3D DATA UNDER DIFFERENT VALUES OF m AND l . COMPUTED WITH MATHEMATICA 6.0. PRECISION UP TO THE 6TH DIGIT AFTER DECIMAL POINT.

Map levels	Total Number of Buckets (l)							
	2	4	8	16	32	64	128	256
m=1	0.531078	0.509177	0.502381	0.500598	0.50015	0.500038	0.50001	0.500002
m=2	0.514551	0.504128	0.50102	0.500247	0.50006	0.500013	0.500004	0.5
m=3	0.505114	0.500774	0.500051	0.499987	0.499991	0.501551	0.499996	0.500004
m=4	0.498119	0.497695	0.499076	0.499717	0.499931	0.498428	0.5	0.5
m=5	0.490039	0.49337	0.496703	0.499313	0.499811	0.499966	0.5	0.499983
m=6	0.47651	0.485541	0.49586	0.498521	0.499586	0.499897	0.499931	0.499897
m=7	0.448987	0.469814	0.48972	0.497032	0.499241	0.499793	0.499931	0.500138
m=8	0.38559	0.435172	0.478726	0.494029	0.49848	0.499448	0.499862	0.5

the region does not change with respect to bucket number i - we only need to change p to ip .

The volume of the bucket i region can thus be expressed as

$$g(i) = \begin{cases} \frac{4}{3}\pi p^3 + 6p\delta^2 + 3\pi p^2\delta + \delta^3, & i = 1 \\ \frac{3}{4}\pi(ip)^3 + 6ip\delta^2 + 3\pi(ip)^2\delta + \delta^3 - v(i, p, \delta), & i > 1 \end{cases}$$

where the first four items in both cases represent the volumes of the four components listed above and $v(i, p, \delta)$ is that for the region formed by half of a spherical surface in Fig. 12(b). With $p = \sqrt{3}\delta$, the above equation becomes

$$g(i) = \begin{cases} (4\sqrt{3}\pi + 6\sqrt{3} + 9\pi + 1)\delta^3 & i = 1 \\ [4\sqrt{3}\pi i^3 + 6\sqrt{3}i + 9\pi i^2 + 1 - v(i, p, \delta)]\delta^3 & i > 1 \end{cases}$$

where $v(i, p, \delta) = 16V_{\mathbf{B}}$ and

$$\begin{aligned} V_{\mathbf{B}} &= \iint_{\mathbf{B}} dx dy \int_{\delta/2}^{\sqrt{p^2 - x^2 - y^2}} dz \\ &= \iint_{\mathbf{B}} \left(\sqrt{p^2 - x^2 - y^2} - \frac{\delta}{2} \right) dx dy \\ &= \int_a^{\frac{\pi}{4}} d\theta \int_b^c \left(\sqrt{p^2 - r^2} - \frac{\delta}{2} \right) r dr \\ &= \int_a^{\frac{\pi}{4}} \left[-\frac{1}{3}(p^2 - r^2)^{\frac{3}{2}} - \frac{\delta}{4}r^2 \right] \Big|_b^c d\theta \\ &= \int_a^{\frac{\pi}{4}} \left[-\frac{\delta^3}{24} + \frac{1}{3}(p^2 - b^2)^{\frac{3}{2}} - \frac{\delta}{4}c^2 + \frac{1}{16} \frac{\delta^3}{(\sin \theta)^2} \right] d\theta \end{aligned}$$

where $a = \arctan \frac{\frac{\delta}{2}}{\sqrt{p^2 - 2\left(\frac{\delta}{2}\right)^2}}$, $c = \sqrt{p^2 - \left(\frac{\delta}{2}\right)^2}$, and $b = \frac{\delta}{2 \sin \theta}$.

We continue to develop formulae for the coverable regions $f(i, m)$ and non-covering factor $\alpha(m)$ as we do in Section IV-B and Section IV-C. These formulae can be found in Appendix II. The complexity of such formulae³ hinders an analytical conclusion on the convergence of $\alpha(m+1)/\alpha(m)$ towards $\frac{1}{2}$. Fortunately, we are able to compute the numerical values of $\alpha(m+1)/\alpha(m)$ under a wide range of inputs. These results (listed in Table IV) clearly show that it does converge to $\frac{1}{2}$.

E. Time complexity of DM-SDH

With Lemma 1, we achieve the following analysis of the time complexity of DM-SDH.

Theorem 1: In DM-SDH, the time spent on operation 1 (i.e., resolving two cells) is $\Theta(N^{\frac{2d-1}{d}})$ where $d \in \{2, 3\}$ is the number of dimensions of the data.

Proof: Given a SDH query, the starting level DM_i is fixed in DM-SDH. Assume there are I pairs of cells to be resolved on DM_i . On the next level DM_{i+1} , total number of cell pairs becomes $I2^{2d}$. According to Lemma 1, half of them will be resolved, leaving only $I2^{2d-1}$ pairs to resolve. On level DM_{i+2} , this number becomes $I2^{2d-1} \frac{1}{2} 2^{2d} = I2^{2(2d-1)}$. Therefore, the number of calls to resolve cells on the different density maps form a geometric progression

$$I, I2^{2d-1}, I2^{2(2d-1)}, \dots, I2^{n(2d-1)}$$

where n is the total number of density maps visited. The time spent on all cell-resolving operations T_c is basically the sum of all items in this progression :

$$T_c(N) = \frac{I[2^{(2d-1)(n+1)} - 1]}{2^{2d-1} - 1}. \quad (16)$$

We use $T_c(N)$ to denote the time under a given size N of the dataset. According to Equation (2), one more level of density map will be built when N increases to $2^d N$. Revisiting Equation (16), we have the following recurrence:

$$T_c(2^d N) = \frac{I[2^{(2d-1)(n+2)} - 1]}{2^{2d-1} - 1} = 2^{2d-1} T_c(N) - o(1) \quad (17)$$

Based on the master theorem, the above recurrence gives

$$T_c(N) = \Theta(N^{\log_{2^d} 2^{2d-1}}) = \Theta(N^{\frac{2d-1}{d}}).$$

Now let us investigate the time complexity for performing operation 2, i.e., calculating distance between particles. We have similar results as in Theorem 1.

Theorem 2: In DM-SDH, the time spent on operation 2 (i.e., distance calculation) is also $\Theta(N^{\frac{2d-1}{d}})$.

³We use Mathematica to solve the integration in Eq. (25) and it ended up an equation that occupies 120 pages!



Fig. 13. Two cells that are non-resolvable are divided into four subcells.

Proof: Like in the derivation of Equation (17), we consider the situation of increasing the dataset size from N to $2^d N$. For any pair of cells on the last density map when dataset is of size N , we have the chance to divide each cell into 2^d smaller cells when another density map is built (as a result of the increase of N). Altogether we have on the new density map $2^d 2^d = 2^{2d}$ pairs to resolve, among which half are expected to be resolved (Lemma 1). This leaves half of the distances in the unresolved cells and they need to be calculated.* By changing the size of the dataset from N to $2^d N$, the total number of distances between any pair of cells increases by 2^{2d} times. As half of these distances need to be calculated, the total number of distance calculations T_d increases by 2^{2d-1} . Therefore, we have the following recurrence:

$$T_d(2^d N) = 2^{2d-1} T_d(N),$$

which is essentially the same as Equation (17), and this concludes the proof. ■

Theorem 3: The time complexity of the DM-SDH algorithm is $\Theta(N^{\frac{2d-1}{d}})$.

Proof: Proof is concluded by combining Theorem 1 and Theorem 2. ■

F. Effects of particle spatial distribution

Theorem 1 is not affected by the locations of individual particles since the results are based the geometric locations of cells. However, the number of distance calculations are related to the distribution of distances, which, in turn, is determined by the spatial distribution of particles. In the proof of Theorem 2 (where we marked a ‘*’), we extend Lemma 1 from percentage of cell pairs to that the ratio of resolvable distances. This extension is obviously true for uniformly distributed particles for which the expected number of distances any cell involves is proportional to the cell size. In this section, we show that the uniform (spatial) distribution of particles is not a necessary condition for Theorem 2 to be true.

Let us consider any pair of 2D cells **A** and **B** that are non-resolvable on density map level k , which is the lowest level for a dataset with N particles. Let the expected number of particles in **A** and **B** be a and b , respectively. Note that $a \neq b$ in general (due the skewed data distribution), and we expect to have ab distances to calculate for this two cells. When the number of particles increases from N to $4N$, we can build another level of density map (level $k+1$). On this level, **A** and **B** are both divided into four cells. Let us denote the expected number of particles before the increase of N in

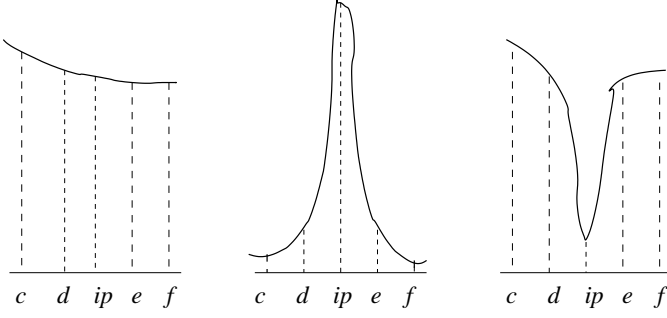


Fig. 14. Three cases of distribution of distances around the edge of buckets i and $i+1$, with the solid curves representing portions of the density function of the distances; $[c, d]$ and $[e, f]$ are examples of distance ranges of resolvable subcells. Those of the non-resolvable subcells are not shown.

the subcells as a_i ($i \in \{1, 2, 3, 4\}$) and b_j ($j \in \{1, 2, 3, 4\}$). We immediately have $a = \sum_{i=1}^4 a_i$ and $b = \sum_{j=1}^4 b_j$. When N becomes $4N$, a_i and b_j all get a four-fold increase and the number of expected distances to calculate becomes

$$T_{k+1} = \sum_{i,j} P_{i,j} 4a_i 4b_j \quad (18)$$

where $P_{i,j}$ is a binary variable that tells whether cell pair a_i and b_j is non-resolvable on the new density map $k+1$. Without any assumptions, we only know that the average of $P_{i,j}$ over all combinations of i and j is 0.5 (i.e., Lemma 1). For Theorem 2 to be true, we need to show that $T_{k+1} = \frac{16}{2} T_k = 8ab$.

We first see that, if the distribution of particles is cell-wise uniform on density map k , we can achieve the above condition. Being cell-wise uniform means that the data are uniformly distributed within the cells, i.e., we should have $a_1 = a_2 = a_3 = a_4 = \frac{a}{4}$ and $b_1 = b_2 = b_3 = b_4 = \frac{b}{4}$, which easily leads to $T_{k+1} = \frac{\sum P_{i,j}}{16} 16ab = 8ab$. The cell-wise uniform distribution is a weaker assumption than the dataset-wise uniform distribution (which requires $a = b$). In simulation data, this can be a safe assumption as the particles will not be indefinitely close to each other due to the existence of bonds and inter-molecular forces. Note that we only need to make this assumption for the smallest cells (i.e., those on the leaf nodes of the tree). Cell-wise uniform is also a popular assumption in current spatial-temporal database studies [16].

A more general discussion on the necessary condition of Theorem 2 would be helpful in identifying its limitations. Revisiting Eq. (18), we see that T_{k+1} is basically a sum of $16a_i b_j$ weighted by $P_{i,j}$, which has an average of 0.5. Therefore, we conclude that, for $T_{k+1} \leq 8ab$ to hold true, **the spatial distribution of particles should NOT be strongly (positively) correlated to the cells that are non-resolvable**. In other words, we cannot have the situation where the case of $P_{i,j} = 1$ are always associated with large $a_i b_j$ values. If we look at the distribution of the distances, this also means that we cannot have high density of distances centering around the bucket boundaries, as shown in the middle graph of Fig. 14. Suppose two cells (e.g., **A** and **B** in Fig. 13) has a distance range $[c, f]$, which overlaps with buckets i and $i+1$, as shown

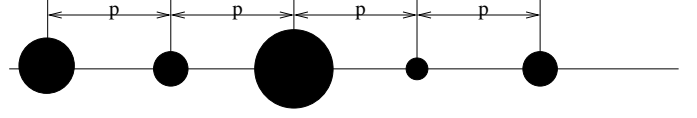


Fig. 15. A spatial distribution of particles that leads to large number of non-resolvable distances. Each ball represents a cluster of particles.

in Fig. 14. With one more density map, their subcells could generate resolvable distance ranges such as $[c, d]$ and $[e, f]$ on the two sides of ip - the boundary of the two buckets. It also generates non-resolvable distance ranges centering around ip . If the distribution of distances has heavy density around ip , most of the area under the density curve will fall into the non-resolvable ranges. On the contrary, if the density curve around ip is not a sharp peak (left hand side in Fig. 14), we could have an equal amount of area under the resolvable and non-resolvable ranges. Or, in another extreme case (right hand side of Fig. 14) where the density is very low around ip , most of the distances will be in the resolvable range.

Having a distance distribution like the one in the middle graph of Fig. 14 means large number of particles are in high-density clusters that are ip in distance. For this to be true in more than one i , the particles must be organized in a linear pattern as shown in Fig. 15. Fortunately, real simulation data will not likely generate such distance distributions because the particles in nature tend to spread out in space (instead of aligning in a line). Even if such distributions are encountered, there is an easy remedy: we can compute another histogram by moving all buckets to the left or right by $\frac{p}{2}$. By this, we can generate a histogram that shows all the trends in the distance distribution (exactly what we need in a histogram) yet most of the distance calculations are avoided.

G. General tiling approach in space partition

We use a regular tiling approach to partition the space in building the trees, i.e., the subcells are of the same shape (square/cube) as the parent cell. In the previous analysis, for each node, we evenly cut each dimension by half, leading to 2^d partitions (child nodes) on the next level. However, in general, we could cut each dimension into $s > 2$ equal segments, giving rise to s^d equal-sized squares or cubes. Interestingly, the value of s does not affect time complexity of the DM-SDH algorithm.

First, the theoretical bucket regions given by Eq. (3) are not affected. For the coverable regions, we incorporate the tiling factor s into the same reasoning as what we utilize to obtain Eq. (8). One exception here is the case of $m = 1, i \geq 2$: the approximate coverable region does not form a series of circles when $s > 2$, therefore Eq. (6) does not hold and this case should be handled in the same way as the case of $m > 1, i \geq 2$. Skipping the details, we get an improved version of Eq. (8) for $s > 2$ as Eq. (19), where $\theta'_m = \frac{1}{2} - \frac{1}{s^m}$ and $\gamma'_m = \sqrt{2(i-1)^2 - \theta'_m{}^2}$. With Eq. (19) to describe the coverable regions, we can easily generate new equations for

TABLE V
 EXPECTED PERCENTAGE OF PAIRS OF CELLS THAT CAN BE RESOLVED UNDER DIFFERENT LEVELS OF DENSITY MAPS AND TOTAL NUMBER OF HISTOGRAM BUCKETS. COMPUTED WITH MATHEMATICA 6.0.

Map levels	Total Number of Buckets (l)							
	2	4	8	16	32	64	128	256
m=1	50.6565	52.1591	52.5131	52.5969	52.6167	52.6214	52.6225	52.6227
m=2	74.8985	75.9917	76.2390	76.2951	76.3078	76.3106	76.3112	76.3114
m=3	87.3542	87.9794	88.1171	88.1473	88.1539	88.1553	88.1556	88.1557
m=4	93.6550	93.9863	94.0582	94.0737	94.0770	94.0777	94.0778	94.0778
m=5	96.8222	96.9924	97.0290	97.0369	97.0385	97.0388	97.0389	97.0389
m=6	98.4098	98.4960	98.5145	98.5184	98.5193	98.5194	98.5195	98.5195
m=7	99.2046	99.2480	99.2572	99.2592	99.2596	99.2597	99.2597	99.2597
m=8	99.6022	99.6240	99.6286	99.6296	99.6298	99.6299	99.6299	99.6299
m=9	99.8011	99.8120	99.8143	99.8148	99.8149	99.8149	99.8149	99.8149
m=10	99.9005	99.9060	99.9072	99.9074	99.9075	99.9075	99.9075	99.9075

$$f(i, m, s) = \begin{cases} \left[2\pi + 4\sqrt{2} + 1 - (8\sqrt{2} + 4)\frac{1}{s^m} + \frac{4}{s^{2m}} \right] \delta^2 & i = 1, m \geq 1 \\ \left\{ 2\pi i^2 + 4\sqrt{2}i - (8\sqrt{2}i + 4)\frac{1}{s^m} + \frac{4}{s^{2m}} - 8 \left[(i-1)^2 \left(\arctan \frac{\gamma'_m}{\theta'_m} - \frac{\pi}{4} \right) - \frac{1}{2}\theta'_m (\gamma'_m - \theta'_m) \right] + 1 \right\} \delta^2 & i > 1, m > 1 \end{cases} \quad (19)$$

the covering factor as a function of m and s . By studying these functions, we get the following lemma.

Lemma 2: With a tiling factor s ($s \in \mathbb{Z}^+$), the non-covering factors have the following property

$$\lim_{l \rightarrow \infty} \frac{\alpha(m+1, s)}{\alpha(m, s)} = \frac{1}{s}.$$

Proof: See Appendix III for details. ■

Lemma 2 is obviously a nicely-formatted extension of Lemma 1. As Lemma 1, it is well supported by numerical results even under smaller values of l (details not shown in this paper). In Section ??, we will discuss the effects of s on the time complexity of DM-SDH.

H. Other costs

I/O costs. In the previous analysis, we focus on the CPU time of the algorithm. Depending on the blocking strategy we use for retrieving data from disk, the exact I/O cost of DM-SDH varies. The bottomline, however, is that the I/O complexity will be asymptotically lower than the quadratic I/O cost needed for calculating all distances.⁴ A straightforward implementation of DM-SDH will give us an I/O complexity $O\left(\left(\frac{N}{b}\right)^{\frac{2d-1}{d}}\right)$ where b is the number of records in each page. To be specific:

⁴To be specific, it is $O\left(\left(\frac{N}{b}\right)^2 \frac{1}{B}\right)$ where b is the page factor and B is the blocking factor if we use a strategy like in block-based nested-loop join.

1. the distance calculations will happen between data points organized in data pages of associated density map cells (i.e., no random reading is needed). On average, one data page only needs to be paired with $O(\sqrt{N})$ other data pages for distance calculation (Theorem 2) in 2D space;
2. I/O complexity for reading density map cells will be the same as in 1. In practice, it will be much smaller, as the size of the nodes is small.

It would be interesting to investigate how we can improve our algorithm to take advantage of blocking and prefetching.

Storage overhead. The storage cost of our algorithm is bound by the size of the density map of the highest resolution we store (leaf nodes in the Quad-tree), as map size decreases exponentially with the decrease of resolution. Obviously, the space complexity is $O(N)$. The total storage overhead will be really small if we have a Quadtree index built for the dataset, which is a popular practice in scientific databases [18]. In this case, we only need to add a p-count field and next pointer to each index node.

V. APPROXIMATE SDH QUERY PROCESSING

While the DM-SDH algorithm is more efficient than current SDH processing methods, its running time for large datasets is still undesirably long. Actually, there are cases where even a coarse SDH will greatly help the fine-tuning of simulation programs [6]. On the other hand, the main motivation to process SDHs is to study the statistical distribution of point-to-point distances in the simulated system [6]. Since a histogram

by itself is an approximation of the underlying distribution $g(r)$ (Equation 1), an inaccurate histogram generated from a given dataset will still be useful in a statistical sense. In this section, we introduce a modified SDH algorithm to give such approximate results to gain better performance in return. Two must-have features for a decent approximate algorithm are :1) provable and controllable error bounds such that the users can have an idea on how close the results are to the fact; and 2) analysis of costs to reach (below) a given error bound, which enables desired performance/correctness tradeoffs. Fortunately, our analytical results shown in Section IV makes the derivation of such error bounds and cost model an easy task.

In the DM-SDH algorithm, we have to : 1) keep resolving cells till we reach the lowest level of the tree; 2) calculate point-to-point distances when we cannot resolve two cells on the leaf level of the tree. Our idea for approximate SDH processing is: stop at a certain tree level and totally skip all distance calculations if we are sure that the number of distances in the unvisited cell pairs fall below some error tolerance threshold.

Recall that, for any given density map DM_{i+m} and total number of buckets l , our analytical model gives the percentage of non-resolvable cell pairs $\alpha(m)$ (Equation (9)). Due to the existence of a closed-form formula, $\alpha(m)$ can be efficiently computed. We list some values of $1 - \alpha(m)$, the percentage of *resolvable* cell pairs, in Table V. Given a user-specified error bound ϵ , we can find the appropriate levels of density maps to visit such that the unvisited cell pairs only contain less than $\epsilon \frac{N(N-1)}{2}$ distances. For example, for a SDH query with 128 buckets and error bound of $\epsilon = 3\%$, we get $m = 5$ by consulting the table. This means, to ensure the 3% error bound, we only need to visit five levels of the tree (excluding the starting level DM_i), and no distance calculation is needed. Table V serves as an excellent validation of Lemma 1: $\alpha(m)$ almost exactly halves itself when m increases by 1, even when l is as small as 2. Since the numbers on the first row (i.e., values for $1 - \alpha(1)$) are also close to 0.5, a rule-of-thumb for choosing m is

$$m = \lg \frac{1}{\epsilon}.$$

The cost of the approximate algorithm only involves resolving cells on the $m + 1$ levels of density maps. Borrowing Equation (17), we obtain the time complexity of the new algorithm

$$T_c(N) \approx I2^{(2d-1)m} = I2^{(2d-1)\lg \frac{1}{\epsilon}} = I\left(\frac{1}{\epsilon}\right)^{2d-1} \quad (20)$$

where I is the number of cell pairs on the starting density map DM_i , and it is solely determined by the query parameter p . Apparently, the running time of this algorithm is not related to the input size N .

Now let us discuss how to deal with those non-resolvable cells after visiting $m + 1$ levels on the tree. In giving the error bounds in our approximate algorithm, we are conservative in assuming the distances in all the unresolved cells will be placed into the wrong bucket. In fact, this almost will never happen because we can distribute the distance counts in the

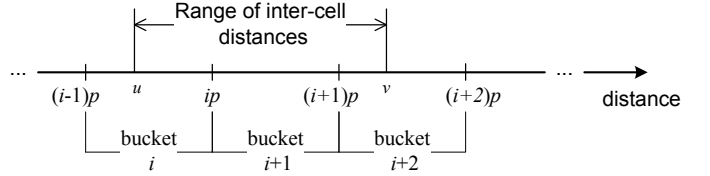


Fig. 16. Distance range of two resolvable cells overlap with three buckets.

unvisited cells to the histogram buckets heuristically and some of them will be done correctly. Consider two non-resolvable cells in a density map with particle counts n_1 and n_2 (total number of n_1n_2 distances between them), respectively. We know their minimum and maximum distances u and v (these are calculated anyway in our attempt to resolve them) fall into multiple buckets. Fig. 16 shows an example that spans three buckets. Using this example, we describe the following heuristics to distributed the n_1n_2 total distance counts into the relevant buckets. These heuristics are ordered in their expected correctness.

1. Put all n_1n_2 distance counts into one bucket;
2. Evenly distribute the distance counts into the three buckets involved, i.e., each bucket gets $\frac{1}{3}n_1n_2$;
3. Distribute the distance counts based on the overlaps between range $[u, v]$ and the buckets. In Fig. 16, the distances put into buckets i , $i + 1$, and $i + 2$ are $n_1n_2 \frac{ip - u}{v - u}$, $n_1n_2 \frac{p}{v - u}$, and $n_1n_2 \frac{v - (i + 1)p}{v - u}$, respectively. Apparently, by adapting this approach, we assume the (statistical) distribution of the point-to-point distances between the two cells is uniform;
4. Assuming a spatial distribution model (e.g., uniform) of particles within individual cells, we can generate the statistical distribution of the distances either analytically or via simulations, and put the n_1n_2 distances to involved buckets based on this distribution.

Note that all four methods need constant time to compute a solution for two cells (In the fourth one, the distribution of the distances can be derived offline). According to our experiments (Section VI-B), they generate much less error than we expect from the theoretical bounds shown in Table V.

VI. EXPERIMENTAL RESULTS

We have implemented the algorithms using the C programming language and tested it with various synthetic/real datasets. The experiments are run at an Apple Mac Pro workstation with two dual-core 2.66GHz Intel Xeon CPUs, and 8GB of physical memory. The operating system is OS X 10.5 Leopard.

A. Exact PDH processing using DM-SDH

The main purpose of this experiment is to verify the time complexity of DM-PDH. In Fig. 17, the running time of our algorithm are plotted against the size of 2D experimental datasets. Fig. 17a shows the results of using synthetic datasets where the locations of individual atoms are distributed

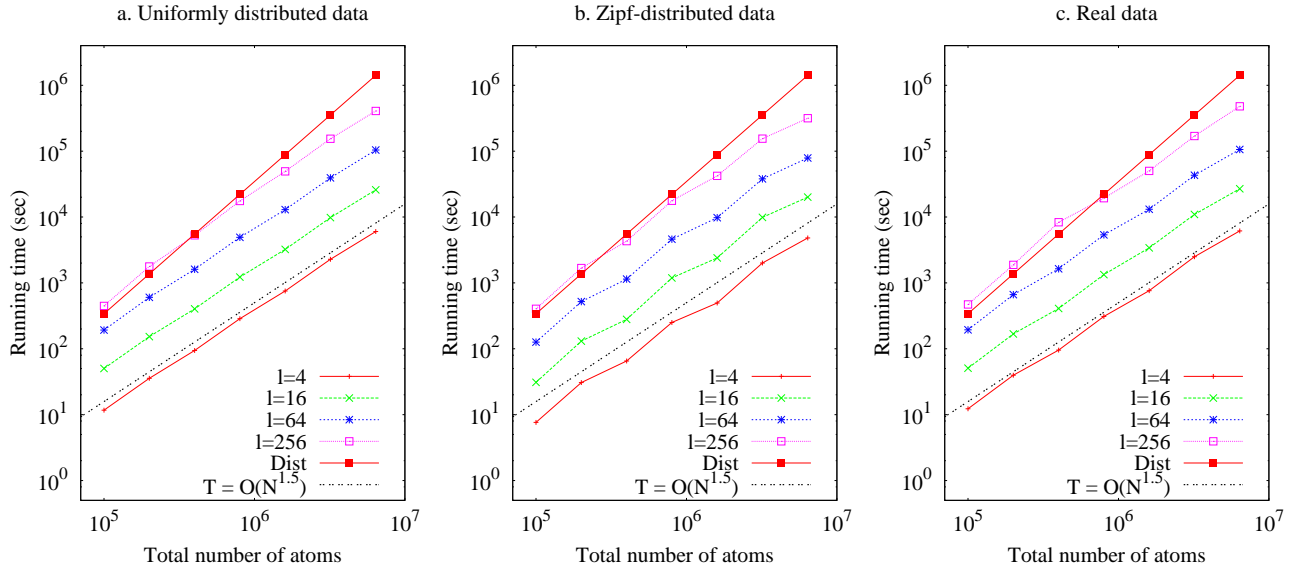


Fig. 17. Running time of the SDH processing algorithms with 2D data.

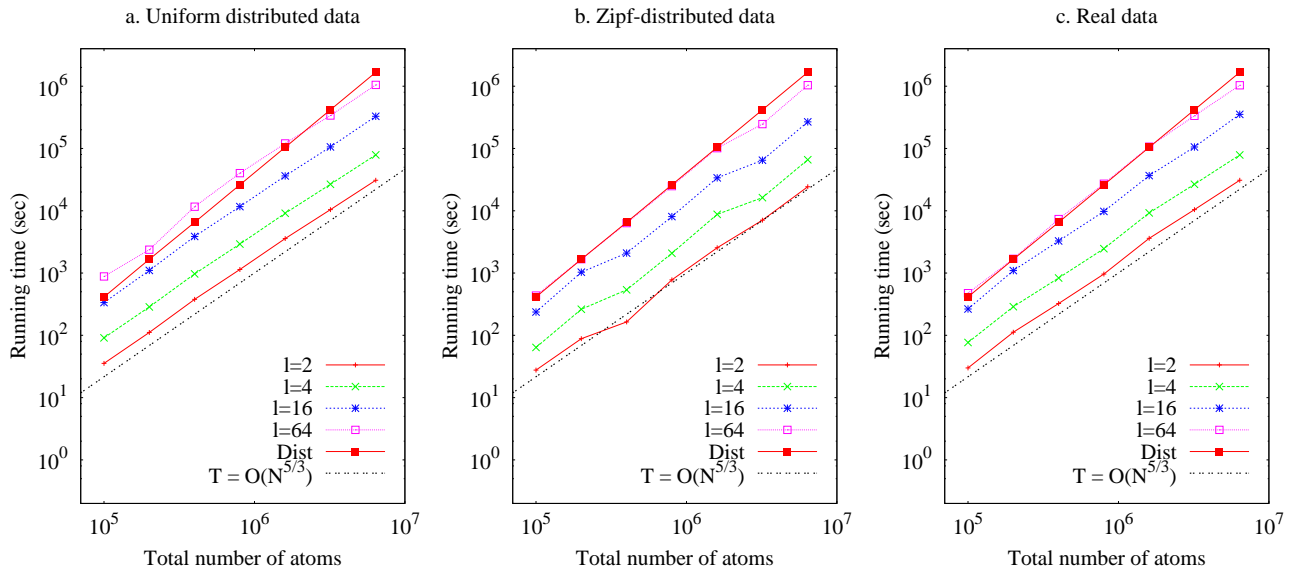


Fig. 18. Running time of the SDH processing algorithms with 3D data.

uniformly in the simulated space and Fig. 17b for those following a Zipf distribution with order one, and Fig. 17c for data generated from Pandit's previous work to simulate a bilayer membrane lipid system in NaCl and KCl solutions, as illustrated in Fig. 19. This dataset has 286000 atoms in it, we randomly choose and duplicate atoms in this dataset to reach different total number of atoms to make the experiments comparable to those in Fig. 17a and 17b. Note that both the running time and data size are plotted on logarithmic scales therefore the gradient of the lines reflects the time complexity of the algorithms. In all graphs, we show the results of a series

of doubling N values ranging from 100,000 to 6,400,000.⁵

For all three datasets, the brute-force approach ('Dist') shows an exact quadratic running time (i.e., the gradient of the line is 2). The other lines (with spots) represent experiments using our algorithm under different bucket numbers. Clearly, the running time of our algorithm grows less dramatically - they all have a gradient of about 1.5. For comparisons, we draw a dotted line in each graph with a slope of exactly 1.5. When bucket size decreases, it takes more time to run our algorithm, although the time complexity is still $\Theta(N^{1.5})$. The cases of

⁵Each point in the graphs shows the result of one single run of DM-PDH as the long running time under large N prohibits having multiple runs. However, we did run multiple experiments with different random seeds for the cases of smaller N and observed very little variances in running time.

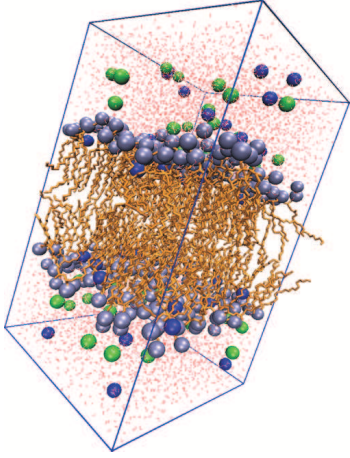


Fig. 19. The simulated hydrated dipalmitoylphosphatidylcholine bilayer system. We can see two layers of hydrophilic head groups (with higher atom density) connected to hydrophobic tails (lower atom density) are surrounded by water molecules (red dots) that are almost uniformly distributed in space.

large bucket numbers ($l = 256$) are worth some attention: its running time is similar to that of the brute-force approach when N is small. However, as N increases, the gradient of the line changes to around 1.5. The reason for this is: when N is small, we have a tree with very few levels; when the query comes with a very small bucket size p , we ended up starting DM-SDH from the leaf level of the tree (recall Fig. 5) and have to essentially calculate most or all distances. However, the same query will get the chance to resolve more cells when the tree becomes taller, as a result of larger N . Another interesting discovery is that the actual running time for the skewed data (Zipf) is always lower than the uniform data. This can be seen by the relative positions of colored lines to the ‘ $T = O(N^{1.5})$ ’ line. This gain of performance comes from the larger number of empty cells on each density map when the particles are clustered. This also confirms our argument that skewed distribution does not affect the correctness of Theorem 2. The results of the real dataset are almost the same as those for the uniform data.

We have similar results for 3D data (Fig. 18): the corresponding lines for DM-SDH have slopes that are very close to $\frac{5}{3}$, confirming our asymptotic analysis. Again, the cases for large l values are worth more discussions. For $l = 64$, we started to see the scenarios of $l = 256$ in the 2D case: the running time grows quadratically till N becomes fairly large (1,600,000) and then the line changes its slope to $\frac{5}{3}$. One thing to notice is the slope of the last segment of $l = 64$ in Fig. 18b is almost 2. This does not mean the time complexity is going back to quadratic. In fact, it has something to do with the zigzag pattern of running time change in the Zipf data: for three consecutive doubling N values (8-fold increase), the running time increases by 2, 4, and 4 times, which still gives a $2 \times 4 \times 4 = 32$ fold increase in total running time (vs. a 64-fold increase in algorithms with quadratic time).

B. Approximate histogram processing

Figure 20 shows the results of running the approximate algorithm. In these experiments, we set the programs to stop after visiting m levels of density maps and distribute the distances using the first three heuristics (Section V). We then compare the approximate histogram with those generated by regular DM-PDH. The error rate is calculated as $\frac{\sum_i |h_i - h'_i|}{\sum_i h_i}$ where for any bucket i , h_i is the accurate count and h'_i the count given by the approximate algorithm.

According to Fig. 20a, the running time does not change with the increase of dataset size for $m = 1, 2, 3$. When m is 4 or 5, the running time increases when N is small and then stays as a constant afterwards. Note that the ‘unlimited’ case shows the results of the basic SDH algorithm. Again, this is because the algorithm has less than 4 or 5 levels to visit in a short tree resulted from small N values. In these cases, our algorithm only saves the time to calculate the unresolved distances. When N is large enough, running time no longer changes with the increase of N .

We observe surprising results on the error rates (Fig. 20 c-d): all experiments have error rates under 3%, even for the cases of $m = 1$! These are much lower than the error bounds we get from Table V. The correctness achieved by heuristic 1 is significantly lower than those by heuristic 2 and 3, as expected. The performance of the latter two are very similar except in the case of $m = 1$ where heuristic 2 had error rates around 0.5%. When $m > 2$, the error rate approaches zero with the dataset becomes larger. Heuristic 3 achieves very low error rates even in scenarios with small m values.

C. Discussions

At this point, we can conclude with confidence that the DM-SDH algorithm has running time in conformity with our analytical results. On the other hand, we also see that, in processing exact SDHs, it shows advantages over the brute-force approach only when N is large (especially when l is also big). However, we would not call this a major limitation of the algorithm as the SDH problem is less meaningful when N is small (especially for large l). Its limitation, in our opinion, is that the time complexity, although superior to quadratic, is still too high to compute SDH for reasonably large simulation dataset. Fortunately, our approximate algorithm provides an elegant practical solution to the problem. According to our experiments, extremely low error rates can be obtained even we only visit as few as three levels of density maps. Note that for large N , the trees are very likely to have more than three levels to explore even when l is large.

The potential of the approximate algorithm shown by current experiments is very exciting. Our explanation for the surprisingly low error rate is: in an individual operation to distribute the distance counts heuristically, we could have made a big mistake by putting too many counts into a bucket (e.g., bucket i in Fig. 16) than needed. But the effects of this mistake could be canceled out by a subsequent mistake where too few counts are put into bucket i . The error rate is measured after the binning is done, thus reflecting the net effects of all

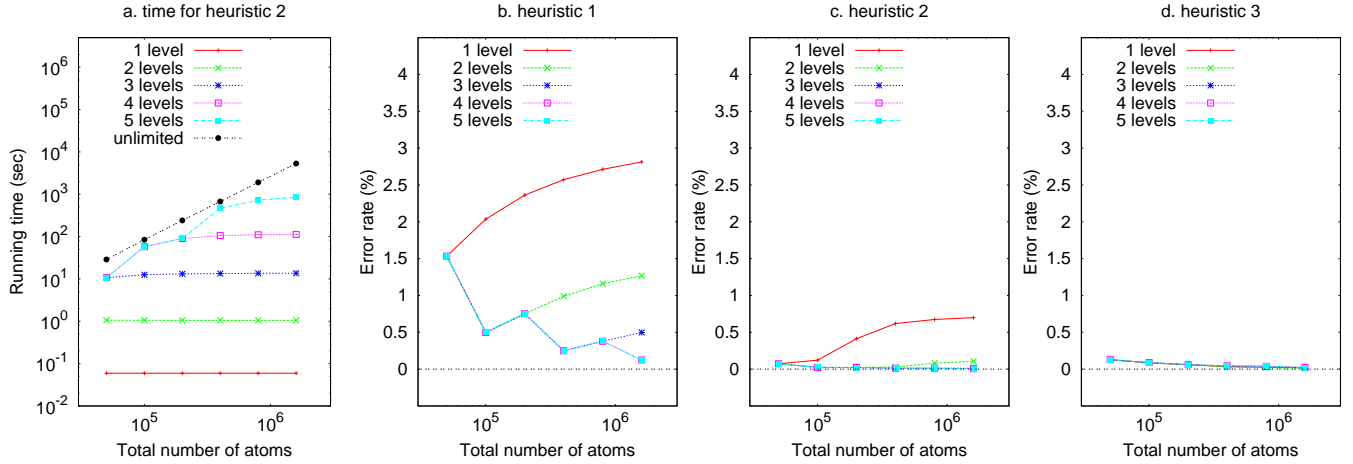


Fig. 20. Running time and correctness of the approximate SDH processing algorithm.

positive and negative mistakes. While more experiments under different scenarios are obviously needed, investigations from an analytical angle are more urgent. We understand that the bound given by Table V is a loose bound. The real error bound should be described as $\epsilon = \epsilon_1 \epsilon_2$ where ϵ_1 is the percentage given by Table V and ϵ_2 is the error rate created by the heuristic binning.

VII. RELATED WORK

Conventional (relational) database systems are designed and optimized towards data and applications from the business world. In recent years, the database community has invested much efforts into constructing database systems that are suitable for handling scientific data. The following are well-known examples of such projects: the GenBank (<http://www.ncbi.nlm.nih.gov/Genbank>) database for biological sequences; the Sloan Digital Sky Survey [2] to explore over 200 million astronomical objects in the sky; the QBISM project [19] for querying and visualizing 3D medical images; the BDBMS project [3] for handling annotation and provenance of biological sequence data; and the PeriScope [5] project for declarative queries against biological sequences. The main challenges and possible solutions of scientific data management are discussed in [1]. Traditionally, molecular simulation data are stored in large files and queries are implemented in standalone programs, as represented by popular simulation/analytics packages [?], [20]. The scientific community has gradually moved towards using database systems for the storage, retrieval, and analysis of large-scale simulation data, as represented by the BioSimGrid [4] and SimDB [21] projects developed for molecular simulations. However, such systems are still in short of efficient query processing strategies. To the best of our knowledge, the computation of SDH in such software packages is done in a brute-force way.

Although the SDH problem has not been studied in the database community, our work is deeply rooted in the philosophy of using tree-based indexing for pruning the information that is not needed. Quadtree has been a well-studied

data structure that has applications spanning databases, image processing, GIS, and computer graphics [22] and has been a topic for research till now: a recent work [23] showed some performance advantages of Quadtree over R-tree in query processing. Various metric trees [24], [25] are closely related to our density map construction: in these structures, space is partitioned into two subsets on each level of the tree. However, the main difference between our work and these tree-based strategies is: we consider the relative distance of cells to group the distance between the data in the cells while they focus more on searching based on a similarity measure.

In particle simulations, the computation of (gravitational/electrostatic) force is of similar flavor to the SDH query. Specifically, the force (or potential) is the sum of all pairwise interactions in the system, thus requires $O(N^2)$ steps to compute. The simulation community has adopted approximate solutions represented by the Barnes-Hut algorithm that runs on $O(N \log N)$ time [26] and the Multi-pole algorithm [27] with linear running time. Although both algorithms use a Quad-tree-like data structure to hold the data, they provide little insights on how to solve the SDH problem. The main reason is that these strategies take advantage of two features of force: 1). for any pairwise interaction, its contribution to the force decreases dramatically when particle distance increases; 2). the effects of symmetric interactions cancel out. However, neither features are applicable to SDH computation, in which every pairwise interaction counts and all are equally important.

VIII. CONCLUSIONS AND FUTURE WORK

In this paper, we argue that the SDH query is critical in analyzing particle simulation data. To improve the efficiency of processing this query, we take advantage of the fact that distance calculation can be processed in a batch instead of individually. We build a data structure based on a point region Quadtree to systematically solve this problem. Our analysis shows that the time complexity of our basic algorithm beats current solutions: it runs at $\Theta(N^{\frac{2d-1}{d}})$ with d being the data dimension number. An approximate algorithm derived from

our approach runs at constant time while giving surprisingly low error rates. We believe our work has provided valuable theoretical and practical insights on the problem such that computing PDH in large simulations has become a reality.

Our work on this topic can be extended in multiple directions. First, the approximate algorithm has shown great potential. Based on the experimental results shown in this paper, we strongly believe there is a tighter bound on the level of errors. Sophisticated (statistical) models should be generated to study this error bound. Second, we should explore more space partitioning plans in building the Quadtree in hope to find one with the “optimal” (or just better) cell resolving percentage. Another topic that we did not pay much attention to is the optimization targeting at the I/O costs. The main issue is how to pack tree nodes into pages and what prefetching strategy we can adopt to improve I/O performance. Finally, our discussions totally ignored another dimension in the data - time. Simulation data are essentially continuous snapshots (called *frames* in simulation terminology) of the simulated system. With large number of frames, processing SDH separately for each frame will take intolerably long time for any meaningful simulation dataset. Incremental solutions need to be developed, taking advantage of the similarity between neighboring frames.

ACKNOWLEDGMENT

The authors would like to thank Prof. Dan Lin of the Department of Computer Science at Purdue University for sharing her insights on this project.

REFERENCES

- [1] J. Gray, D. Liu, M. Nieto-Santesteban, A. Szalay, D. DeWitt, and G. Heber, “Scientific Data Management in the Coming Decade,” *SIGMOD Record*, vol. 34, no. 4, pp. 34–41, December 2005.
- [2] A. S. Szalay, J. Gray, A. Thakar, P. Z. Kunszt, T. Malik, J. Raddick, C. Stoughton, and J. vandenBerg, “The SDSS Skyserver: Public Access to the Sloan Digital Sky Server Data,” in *Proceedings of International Conference on Management of Data (SIGMOD)*, 2002, pp. 570–581.
- [3] M. Y. Eltabakh, M. Ouzzani, and W. G. Aref, “BDBMS - A Database Management System for Biological Data,” in *Proceedings of the 3rd Biennial Conference on Innovative Data Systems Research (CIDR)*, 2007, pp. 196–206.
- [4] M. H. Ng, S. Johnston, B. Wu, S. E. Murdock, K. Tai, H. Fangohr, S. J. Cox, J. W. Essex, M. S. P. Sansom, and P. Jeffreys, “BioSimGrid: Grid-enabled Biomolecular Simulation Data Storage and Analysis,” *Future Generation Computer Systems*, vol. 22, no. 6, pp. 657–664, June 2006.
- [5] J. M. Patel, “The Role of Declarative Querying in Bioinformatics,” *OMICS: A Journal of Integrative Biology*, vol. 7, no. 1, pp. 89–91, 2003.
- [6] D. Frenkel and B. Smit, *Understanding Molecular Simulation From Algorithm to Applications*, ser. Computational Science Series. Academic Press, 2002, vol. 1.
- [7] M. P. Allen and D. J. Tildesley, *Computer Simulations of Liquids*. Clarendon Press, Oxford, 1987.
- [8] J. M. Haile, *Molecular Dynamics Simulation: Elementary Methods*. Wiley, New York, 1992.
- [9] D. P. Landau and K. Binder, *A Guide to Monte Carlo Simulation in Statistical Physics*. Cambridge University Press, Cambridge, 2000.
- [10] P. K. Agarwal, L. Arge, and J. Erikson, “Indexing Moving Objects,” in *Proceedings of International Conference on Principles of Database Systems (PODS)*, 2000, pp. 175–186.
- [11] M. Bamdad, S. Alavi, B. Najafi, and E. Keshavarzi, “A new expression for radial distribution function and infinite shear modulus of lennard-jones fluids,” *Chem. Phys.*, vol. 325, pp. 554–562, 2006.
- [12] J. L. Stark and F. Murtagh, *Astronomical Image and Data Analysis*. Springer, 2002.
- [13] A. Filippini, “The radial distribution function probed by X-ray absorption spectroscopy,” *J. Phys.: Condens. Matter*, vol. 6, pp. 8415–8427, 1994.
- [14] V. Springel, S. D. M. White, A. Jenkins, C. S. Frenk, N. Yoshida, L. Gao, J. Navarro, R. Thacker, D. Croton, J. Helly, J. A. Peacock, S. Cole, P. Thomas, H. Couchman, A. Evrard, J. Colberg, and F. Pearce, “Simulations of the Formation, Evolution and Clustering of Galaxies and Quasars,” *Nature*, vol. 435, pp. 629–636, June 2005.
- [15] J. A. Orenstein, “Multidimensional Tries used for Associative Searching,” *Information Processing Letters*, vol. 14, no. 4, pp. 150–157, 1982.
- [16] Y. Tao, J. Sun, and D. Papadias, “Analysis of predictive spatio-temporal queries,” *ACM Trans. Database Syst.*, vol. 28, no. 4, pp. 295–336, 2003.
- [17] I. Csabai, M. Trencseni, L. Dobos, P. Jozsa, G. Herczegh, N. Purger, T. Budavari, and A. S. Szalay, “Spatial Indexing of Large Multidimensional Databases,” in *Proceedings of the 3rd Biennial Conference on Innovative Data Systems Research (CIDR)*, 2007, pp. 207–218.
- [18] M. Arya, W. F. Cody, C. Faloutsos, J. Richardson, and A. Toya, “QBISM: Extending a DBMS to Support 3D Medical Images,” in *ICDE*, 1994, pp. 314–325.
- [19] B. Hess, C. Kutzner, D. van der Spoel, and E. Lindahl, “GROMACS 4: Algorithms for Highly Efficient, Load-Balanced, and Scalable Molecular Simulation,” *Journal of Chemical Theory and Computation*, vol. 4, no. 3, pp. 435–447, March 2008.
- [20] M. Feig, M. Abdullah, L. Johnsson, and B. M. Pettitt, “Large Scale Distributed Data Repository: Design of a Molecular Dynamics Trajectory Database,” *Future Generation Computer Systems*, vol. 16, no. 1, pp. 101–110, January 1999.
- [21] H. Samet, *Foundations of Multidimensional and Metric Data Structure*. Morgan Kaufman, 2006.
- [22] Y. J. Kim and J. M. Patel, “Rethinking Choices for Multi-dimensional Point Indexing: Making the Case for the Often Ignored Quadtree,” in *Proceedings of the 3rd Biennial Conference on Innovative Data Systems Research (CIDR)*, 2007, pp. 281–291.
- [23] J. K. Uhlman, “Metric Trees,” *Applied Mathematics Letters*, vol. 4, no. 5, pp. 61–62, 1991.
- [24] P. N. Yianilos, “Data Structures and Algorithms for Nearest Neighbor Search in Metric Spaces,” in *Proceedings of ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 1993, pp. 311–321.
- [25] J. Barnes and P. Hut, “A Hierarchical O(N log N) Force-Calculation Algorithm,” *Nature*, vol. 324, no. 4, pp. 446–449, 1986.
- [26] L. Greengard and V. Rokhlin, “A Fast Algorithm for Particle Simulations,” *Journal of Computational Physics*, vol. 135, no. 12, pp. 280–292, 1987.

APPENDIX I
THE DERIVATION OF EQ. (15)

We accomplish this proof by studying the difference between $\frac{A(m)}{B(m)}$ and $\frac{1}{2}$. First, we see

$$\begin{aligned}
A(m) - \frac{B(m)}{2} &= \sum_{i=2}^l \sqrt{2(i-1)^2 - \frac{1}{4}} - 4 \sum_{i=2}^l \theta_{m+1} \sqrt{2(i-1)^2 - \theta_{m+1}^2} \\
&+ 2 \sum_{i=2}^l \theta_m \sqrt{2(i-1)^2 - \theta_m^2} + 8 \sum_{i=2}^l (i-1)^2 \arctan \frac{\sqrt{8(i-1)^2 - \theta_{m+1}^2}}{\theta_{m+1}} \\
&- 4 \sum_{i=2}^l (i-1)^2 \arctan \frac{\sqrt{8(i-1)^2 - \theta_m^2}}{\theta_m} - 4 \sum_{i=2}^l (i-1)^2 \arctan \sqrt{8(i-1)^2 - 1} \quad (21)
\end{aligned}$$

When $l \rightarrow \infty$, we have the results shown in (22).

$$\begin{aligned}
\sum_{i=2}^l \sqrt{2(i-1)^2 - \frac{1}{4}} &\longrightarrow \sum_{i=2}^l \sqrt{2}(i-1) \\
\sum_{i=2}^l \theta_{m+1} \sqrt{2(i-1)^2 - \theta_{m+1}^2} &\longrightarrow \sum_{i=2}^l \theta_{m+1} \sqrt{2}(i-1) \\
\sum_{i=2}^l \theta_m \sqrt{2(i-1)^2 - \theta_m^2} &\longrightarrow \sum_{i=2}^l \theta_m \sqrt{2}(i-1) \\
\sum_{i=2}^l (i-1)^2 \arctan \frac{\sqrt{8(i-1)^2 - \theta_{m+1}^2}}{\theta_{m+1}} &\longrightarrow \sum_{i=2}^l (i-1)^2 \arctan 2\sqrt{2}(i-1) \\
\sum_{i=2}^l (i-1)^2 \arctan \frac{\sqrt{8(i-1)^2 - \theta_m^2}}{\theta_m} &\longrightarrow \sum_{i=2}^l (i-1)^2 \arctan 2\sqrt{2}(i-1) \\
\sum_{i=2}^l (i-1)^2 \arctan \sqrt{8(i-1)^2 - 1} &\longrightarrow \sum_{i=2}^l (i-1)^2 \arctan 2\sqrt{2}(i-1) \quad (22)
\end{aligned}$$

Plugging the left-hand side of six formulae in (22) into Eq. (21), we get $A(m) - \frac{B(m)}{2} \longrightarrow 0$ and thus $A(m) \longrightarrow \frac{B(m)}{2}$.

APPENDIX II
RELEVANT QUANTITIES IN 3D ANALYSIS

These formulae are listed on the last page of this paper as Eq. (23) to Eq. (25).

$$\begin{aligned}
V_{\mathbf{B}'}(m) &= \iint_{\mathbf{B}'} dx dy \int_{\frac{\delta}{2}}^{\sqrt{p^2 - (x - \frac{\delta}{2m})^2 - (y - \frac{\delta}{2m})^2 + \frac{\delta}{2m}}} dz \\
&= \iint_{\mathbf{B}'} \left[\sqrt{p^2 - \left(x - \frac{\delta}{2m}\right)^2 - \left(y - \frac{\delta}{2m}\right)^2} - \delta\theta_m \right] dx dy \\
&= \int_a^{\frac{\pi}{4}} d\phi \int_b^c \left(\sqrt{p^2 - r^2} - \delta\theta_m \right) r dr \\
&= \int_a^{\frac{\pi}{4}} \left[-\frac{1}{3}(p^2 - r^2)^{\frac{3}{2}} - \frac{\delta\theta_m}{2} r^2 \right] \Big|_b^c d\phi \\
&= \int_a^{\frac{\pi}{4}} \left[-\frac{(\delta\theta_m)^3}{3} + \frac{1}{3} (p^2 - b^2)^{\frac{3}{2}} - \frac{\delta\theta_m}{2} [p^2 - (\delta\theta_m)^2] + \frac{\delta\theta_m}{2} b^2 \right] d\phi \quad (23)
\end{aligned}$$

Here we have $a = \arctan \frac{\delta\theta_m}{\sqrt{p^2 - 2(\delta\theta_m)^2}}$, $b = \frac{\delta\theta_m}{\sin \phi}$, and $c = \sqrt{p^2 - (\delta\theta_m)^2}$.

The coverable region is

$$f(i, m) = \begin{cases} \frac{4}{3}\pi p^3 + 6p \left(\delta - \frac{2\delta}{2m} \right)^2 + 3\pi p^2 \delta \left(\delta - \frac{2\delta}{2m} \right) + \left(\delta - \frac{2\delta}{2m} \right)^3 & n = 1, m \geq 1 \\ \frac{4}{3}\pi(ip)^3 - \frac{4}{3}\pi[(i-1)p]^3 & n \geq 2, m = 1 \\ \frac{4}{3}\pi(ip)^3 + 6ip \left(\delta - \frac{2\delta}{2m} \right)^2 + 3\pi(ip)^2 \delta \left(\delta - \frac{2\delta}{2m} \right) + \left(\delta - \frac{2\delta}{2m} \right)^3 - v(i, m, p, \delta) & n \geq 2, m > 1 \end{cases}$$

Simplifying the above with $p = \sqrt{3}\delta$, we get

$$f(i, m) = \begin{cases} \left[4\sqrt{3}\pi + 6\sqrt{3} \left(1 - \frac{2\delta}{2m} \right)^2 + 9\pi \left(1 - \frac{2}{2m} \right) + \left(1 - \frac{2}{2m} \right)^3 \right] \delta^3 & n = 1, m \geq 1 \\ [4\sqrt{3}\pi i^3 - 4\sqrt{3}\pi(i-1)^3] \delta^3 & n \geq 2, m = 1 \\ \left[4\sqrt{3}\pi i^3 + 6\sqrt{3}i \left(1 - \frac{2\delta}{2m} \right)^2 + 9\pi i^2 \left(1 - \frac{2}{2m} \right) + \left(1 - \frac{2}{2m} \right)^3 - v(i, m, p, \delta) \right] \delta^3 & n \geq 2, m > 1 \end{cases}$$

where $v(i, m, p, \delta) = 16V_{\mathbf{B}'}(m)$.

Continue with the same reasoning as in Section IV-C, we have

$$G(l) = \frac{\sum_{i=1}^l g(i)}{\delta^3} = \sum_{i=1}^l \left(4\sqrt{3}\pi i^3 + 6\sqrt{3}i + 9\pi i^2 + 1 \right) - 16 \sum_{i=2}^l \int_q^{\frac{\pi}{4}} \left[-\frac{1}{24} + \frac{1}{3} \left(3(i-1)^2 - \left(\frac{1}{2\sin\phi} \right)^2 \right)^{\frac{3}{2}} - \frac{1}{4} \left(3(i-1)^2 - \left(\frac{1}{2} \right)^2 \right) + \frac{1}{16} \frac{1}{(\sin\phi)^2} \right] d\phi \quad (24)$$

where $q = \arctan \frac{\frac{1}{2}}{\sqrt{3(i-1)^2 - 2\left(\frac{1}{2}\right)^2}}$, and the following formulae for the accumulated volume for all coverable regions F .

$$\frac{\sum_{i=1}^l f(i, m)}{\delta^3} = \begin{cases} 4\sqrt{3}\pi + \sum_{i=2}^l [4\sqrt{3}i^3 - 4\sqrt{3}\pi(i-1)^3] = 4\sqrt{3}\pi l^3, & m = 1 \\ \sum_{i=1}^l \left[4\sqrt{3}\pi i^3 + 6\sqrt{3}i \left(1 - \frac{2}{2m} \right)^2 + 9\pi i^2 \left(1 - \frac{2}{2m} \right) + \left(1 - \frac{2}{2m} \right)^3 \right] - 16 \sum_{i=2}^l \int_s^{\frac{\pi}{4}} \left[-\frac{\theta_m^3}{3} + \frac{1}{3} \left(3 - \left(\frac{\theta_m}{\sin\phi} \right)^2 \right)^{\frac{3}{2}} - \frac{\theta_m}{2} (3 - \theta_m^2) + \frac{\theta_m}{2} \left(\frac{\theta_m}{\sin\phi} \right)^2 \right] d\phi, & m > 1 \end{cases} \quad (25)$$

in which $s = \arctan \frac{\theta_m}{\sqrt{3(i-1)^2 - 2\theta_m^2}}$.

APPENDIX III
PROOF OF LEMMA 2

Proof: Proof is accomplished in a similar way to that of Lemma 1. While the total area of all bucket regions Eq. (10) is still the same, Eq. (11) and Eq. (12) become the following equation for all $m \geq 1$:

$$\begin{aligned}
 F(l, m, s) &= \frac{\sum_{i=1}^l f(i, m, s)}{\delta^2} \\
 &= \sum_{i=1}^l \left[\pi(ip)^2 + 4ip \left(\delta - \frac{2\delta}{s^m} \right) + \left(\delta - \frac{2\delta}{s^m} \right)^2 \right] \\
 &\quad - \sum_{i=2}^l (i-1)^2 \left[8 \arctan \frac{\sqrt{2(i-1)^2 - \theta'_m{}^2}}{\theta'_m{}^2} - 2\pi \right] + 4 \sum_{i=2}^l \left[\theta'_m \sqrt{2(i-1)^2 - \theta'_m{}^2} - \theta'_m{}^2 \right],
 \end{aligned} \tag{26}$$

which gives $\frac{\alpha(m+1, s)}{\alpha(m, s)} = \frac{A(m, s)}{B(m, s)}$ where

$$\begin{aligned}
 A(m, s) &= 1 + \frac{4\sqrt{2}(l+l^2)}{s^{1+m}} - l \left(1 - \frac{2}{s^{1+m}} \right)^2 + 4(l-1) \left(\frac{1}{2} - \frac{1}{s^{1+m}} \right)^2 \\
 &\quad - 4 \sum_{i=2}^l \theta'_{m+1} \sqrt{2(i-1)^2 - \theta'_{m+1}{}^2} + 8 \sum_{i=2}^l (i-1)^2 \arctan \frac{\sqrt{2(i-1)^2 - \theta'_{m+1}{}^2}}{\theta'_{m+1}} \\
 &\quad + \sum_{i=2}^l \sqrt{8(i-1)^2 - 1} - 8 \sum_{i=2}^l (i-1)^2 \arctan \sqrt{8(i-1)^2 - 1}
 \end{aligned} \tag{27}$$

and

$$\begin{aligned}
 B(m, s) &= 1 + \frac{4\sqrt{2}(l+l^2)}{s^m} - l \left(1 - \frac{2}{s^m} \right)^2 + 4(l-1) \left(\frac{1}{2} - \frac{1}{s^m} \right)^2 \\
 &\quad - 4 \sum_{i=2}^l \theta'_m \sqrt{2(i-1)^2 - \theta'_m{}^2} + 8 \sum_{i=2}^l (i-1)^2 \arctan \frac{\sqrt{2(i-1)^2 - \theta'_m{}^2}}{\theta'_m} \\
 &\quad + \sum_{i=2}^l \sqrt{8(i-1)^2 - 1} - 8 \sum_{i=2}^l (i-1)^2 \arctan \sqrt{8(i-1)^2 - 1}
 \end{aligned} \tag{28}$$

Following the reasoning in Appendix I, we compare the value of $\frac{A(m, s)}{B(m, s)}$ to $\frac{1}{s}$. And we have

$$\begin{aligned}
 A(m, s)s - B(m, s) &= (s-1) \sum_{i=2}^l \sqrt{8(i-1)^2 - 1} - 8(1-s) \sum_{i=2}^l (i-1)^2 \arctan \sqrt{8(i-1)^2 - 1} \\
 &\quad - 4(1-s) \sum_{i=2}^l \theta'_{m+1} \sqrt{2(i-1)^2 - \theta'_{m+1}{}^2} \\
 &\quad + 8(s-1) \sum_{i=2}^l (i-1)^2 \arctan \frac{\sqrt{2(i-1)^2 - \theta'_{m+1}{}^2}}{\theta'_{m+1}}
 \end{aligned} \tag{29}$$

When $l \rightarrow \infty$, we have the results shown in (30).

$$\begin{aligned}
 \sum_{i=2}^l \sqrt{2(i-1)^2 - \theta'_{m+1}{}^2} &\longrightarrow \frac{1}{2} \sum_{i=2}^l \sqrt{8(i-1)^2 - 1} \\
 \sum_{i=2}^l (i-1)^2 \arctan \frac{\sqrt{2(i-1)^2 - \theta'_{m+1}{}^2}}{\theta'_{m+1}} &\longrightarrow \sum_{i=2}^l (i-1)^2 \arctan \sqrt{8(i-1)^2 - 1}
 \end{aligned} \tag{30}$$

Plugging the left-hand side of the above two formulae in (30) into Eq. (29), we get $sA(m, s) - B(m, s) \rightarrow 0$ and thus $A(m, s) \rightarrow \frac{B(m, s)}{s}$. ■