

Organizing and Indexing Non-Convex Regions

Eric Perlman, Randal Burns, Michael Kazhdan
Department of Computer Science
Johns Hopkins University
Baltimore, MD 21218
{eric,randal,misha}@cs.jhu.edu

1. OVERVIEW

We demonstrate data indexing and query processing techniques that improve the efficiency of comparing, correlating, and joining data contained in non-convex regions. We use computational geometry techniques to automatically characterize the region of space from which data are drawn, partition the region based on that characterization, and create an index from the partitions. Our motivating application performs distributed data analysis queries among federated database sites that store scientific data sets from the Chesapeake Bay. Our preliminary findings indicate that these techniques often reduce the number of I/Os needed to serve a query by a factor of five—depending on the geometry of the query region.

Our approach automatically extracts the structure of the region to facilitate data organization and query processing. We use an approximate medial axis transform to discover and describe the region. The medial axis is generated using a constrained Delaunay triangulation [3]. We represent the medial axis as a tree and enumerate the nodes of the tree to generate the labels for our spatial index. The index defines a linear order on the Delaunay triangles. Data from each triangle are then placed on disk in index order, e.g. using a B+-tree. When using this index, structures (in our case, estuaries, rivers, and bays) occupy contiguous regions of the index, which makes the data from these regions contiguous on disk.

We will demonstrate an interactive application in which the user can select from a variety of shape files, including the Chesapeake Bay, other estuaries, and artificial shapes. The demo will show the generation of an index, including the Delaunay triangulation and medial axis transform. The user will then select query regions in the space and the demo will show both the portions of the index associated with these regions and, for the Chesapeake Bay, the data layout on disk for multiple data sets. The demo will compare the index regions and data layout among our medial axis index and regular spatial decompositions, such as region trees, tessellations, and space filling curves.

Permission to make digital or hard copies of portions of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyright for components of this work owned by others than VLDB Endowment must be honored.

Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists requires prior specific permission and/or a fee. Request permission to republish from: Publications Dept., ACM, Inc. Fax +1 (212)869-0481 or permissions@acm.org.

PVLDB '08, August 23-28, 2008, Auckland, New Zealand
Copyright 2008 VLDB Endowment, ACM 978-1-60558-306-8/08/08

1.1 Design Goals

Our original goal was to find a **spatially-derived data-independent decomposition** suitable for the Chesapeake bay. We need a data-independent decomposition of the region so that the organization of the region is uniform across all data sets. With a data-independent decomposition, data may be indexed at different sites, at different times, and incrementally, and the indexes may be used to join data, find nearest neighbors, etc. In contrast, data-derived indexes, e.g. k-d trees, R-trees, Voronoi diagrams, give irregular decompositions so that separate indexes of data drawn from the same region are incomparable. Data-independent decompositions are particularly important in federated data systems, in which different data sets are stored and managed independently, or when data are too large to build a single index for all data, e.g. in parallel database systems.

The need for a more complex spatial decomposition becomes clear with estuarine data. Available data-independent decompositions tend to be regular, e.g. space-filling curves, and region trees [6], which either divide a region into a series of regular polygons or recursively divide the region into squares (in 2-d). Estuaries are long and skinny and exhibit a large number of winding, tendril-like tributaries, which are not likely to be wholly encapsulated by the cells of a regular decomposition. The distribution of tributary fragments will not conform to the regular indexing of the spatial partition and the data will be spread across non-contiguous regions of the disk and database indexes. For example, Figure 1 shows how a Hilbert curve, one of the most-frequently used space-filling curves, performs very poorly in clustering the region of space associated with the Potomac River (shaded). The Hilbert curve breaks up data from the river into many disjoint regions in the index: the linear ordering following the curve. In contrast, we would like to maintain the data from contiguous spatial regions contiguously on disk. We need to preserve this property at multiple scales, because estuary data are self-similar in addition to being non-convex.

Our techniques are applicable beyond estuaries. Any data set drawn from a non-convex region will realize performance benefits. Examples include medical applications (the circulatory system), manufactured systems (road networks and indoor air systems), and other natural systems (turbulent structure and flow in porous media). At present, we support two-dimensional spatial data and are pursuing three-dimensional data. However, the medial axis becomes more difficult to work with in that the 3-d medial axis consists of complex unions of surfaces and curves.

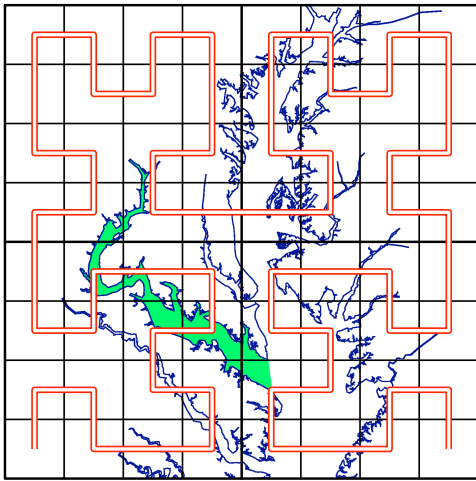


Figure 1: Example of a Hilbert curve ordering applied to the Chesapeake Bay.

2. IMPLEMENTATION

We implemented index generation and query processing in the Chesapeake Bay Environmental Observatory (CBEO): a platform for the data exploration of more than 15 Chesapeake Bay data sets. The CBEO makes the output of hydrodynamic and water-quality models available alongside observational data from buoys, cruises, shallow water monitoring, satellites, and overflight. Exploring correlations among these disparate data sets enables new types of discovery in hydrology and environmental engineering [1].

We programmed the indexing routines and query evaluation in C++, using the CGAL [2] computational geometry algorithm library. These routines have also been integrated into Microsoft SQL Server—the platform for the CBEO testbed—as user-defined functions. The maps are pre-computed from ArcGIS shape files and loaded into SQL Server as dynamic libraries. The demonstration software has been implemented in OBJECTIVE-C that calls the C++ indexing and query routines directly.

2.1 Index Generation

We use computational geometry techniques to generate our spatial index from a region. We begin by looking at the outline of our region, a simple, closed polygon (Figure 2(a)). (The demo does not support polygons with holes. It requires the medial axis to be a tree. More complex medial axes could be simplified using the minimal spanning tree.) Next, we generate a triangulation of the region (Figure 2(b)). We use a constrained Delaunay triangulation (CDT) [3] to produce a triangulation that partitions the entire region exactly, i.e. without including portions of space outside the region’s boundary. The triangulation is constrained in that the edges that form the perimeter of the shape must be included. The constraints results in non-Delaunay triangles in the output. At this point, we compute the medial axis by generating a tree connecting the centers of the circumcircles of all triangles (Figures 2(c) and Figure 2(d)).

We assign unique, incremental, index addresses to each of the triangles through a traversal of the medial axis. Figures 3(a) and 3(b) show the resulting map. The medial axis forms a tree for two-dimensional objects with no holes.

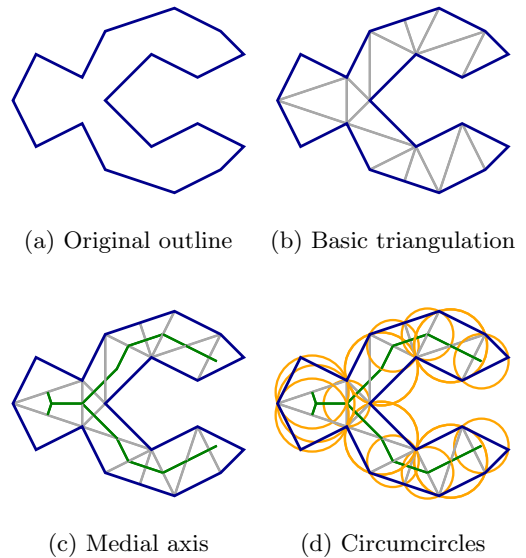


Figure 2: Delaunay triangulation of a polygon constrained by outline segments.

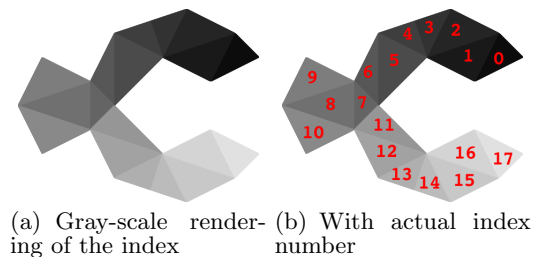


Figure 3: Index for our example polygon.

The order in which we traverse the medial-axis can result in quite different indexes. Figure 4(a) shows the medial-axis computed on the Chesapeake Bay. Figures 4(b) and 4(c) shows two different orderings of index addresses. We start the traversal of the medial axis at the Nansemond River, the southern most tributary of the bay. Any randomly chosen extrema will work well; starting at a central node partitions a central structure into two parts. These two renderings of the index-space are generated by traversing the medial axis with respect to the weight of the subtree, i.e. sum of the area of the descendant triangles. Figure 4(b) traverses the smallest area subtrees first, which gives a smoother indexing of the shape, including all of the micro-structure on the sides of each tributary. Figure 4(c) traverses the greatest area subtrees first, which gives us a consecutive “main stem” of the bay and each of the tributaries, but the small structures are entirely disjoint in the index space. The preferred ordering depends on the application’s data distribution and access patterns. For example, scientists studying hydrodynamics are mostly concerned with deeper water and prefer to keep the main stem of the bay as contiguous as possible. Scientists looking at nitrogen are concerned with how shallow and deep water interact and would want small structures to be contiguous to nearby larger structures.

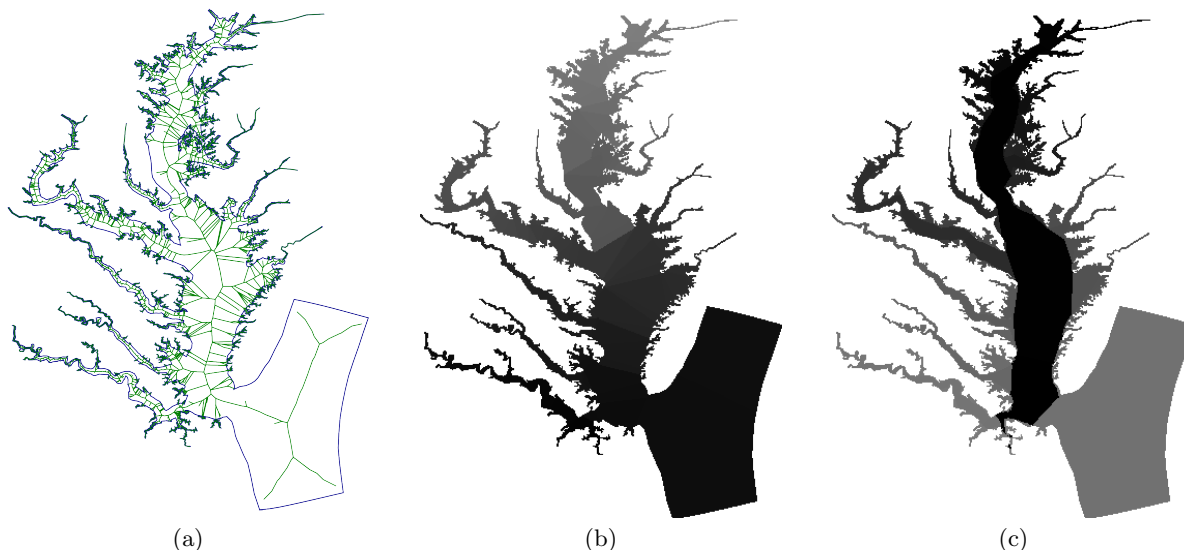


Figure 4: Renderings of the Chesapeake Bay: (a) shows the medial-axis, (b) and (c) are gray-scale images of the index generated by traversing the medial-axis by smallest and largest area subtrees first.

The centerlines of river networks have long been used in the field of hydrology for feature extraction [5]. This idea was extended to the medial-axis transform by McAllister and Snoeyink [4] to further characterize rivers, e.g. match shores, calculate width, and estimate volume. However, they do not apply the medial axis to the physical organization of data or indexing.

2.2 Query Evaluation

The medial axis spatial index cannot be used directly for query evaluation. The index encodes spatial locality in a linear ordering of the Delaunay triangulation and defines an organization on disk based on that linear ordering. However, it does naturally associate points in space with that index.

We use an additional lookup table to determine the triangle(s) associated with a point or region of space. We generate the lookup table by rasterizing the region onto a low-resolution pixel structure. For each pixel, we store a list of the triangles that intersect the square region covered by that pixel. The table is generated once for each shape; it too is data independent.

The size of the lookup table needs to be chosen with consideration and depends on the geometric complexity of the region. Choosing too fine a pixelation produces an overly large lookup table, which consumes cache space and, in the extreme, cannot fit in memory. Choosing too coarse a pixelation results in many triangles intersecting each pixel, which slows query evaluation. For the Chesapeake Bay, we choose a resolution of 350 by 500 so that there are no more than 38 triangles intersecting any pixel and an average of 2.3 triangles per pixel.

Index Lookup: The basic operation finds the index value for a coordinate location. This is used when ingesting data into a database to generate the index labels or primary keys for records. It is also the first step in nearest-neighbor queries, because it identifies where in the index to start the nearest neighbor search.

To identify a coordinate's index, we map that coordinate to a pixel in our lookup table and retrieve the list of triangles at that pixel. Then, we iterate over the list of triangles to determine which triangle contains the point.

Alternatively, this query may be implemented using other spatial partitioning techniques, e.g B-trees, BSP-trees, etc. We chose to use a lookup table for ease of implementation and acceptable performance on our data.

Range Query: Euclidean distance range queries are processed by identifying candidate triangles that might contribute data because they intersect the query range. To do so, we consult the lookup table to find pixels that overlap the range. For pixels entirely inside the query range, their triangles intersect the query range and must be included. For pixels that intersect the range boundary, their triangles may or may not intersect the range. We evaluate these triangles individually for intersection. Then, we obtain the requested data by joining the candidate indexes with the data tables and applying a filter to eliminate data outside the query range.

3. DEMONSTRATION

Our demonstration allows the user to both interact with the construction of medial axis spatial indexes and to visualize the benefits of the technique through a graphical-user interface. The user first selects a shape file from among the Chesapeake Bay, other estuaries and rivers (e.g., the Corpus Christi bay, the Charleston river, the Camargue), and geometric shapes. The demo generates the corresponding spatial index. The user can visualize this process by adding and removing layers, including the triangulation, map, medial axis, and circumcircles (Figure 5(a)).

The user then explores and visualizes how the region and range queries relate to the index. The user selects a query region with the mouse and the demo displays the portions of the index that intersect the selection (Figure 5(b)).

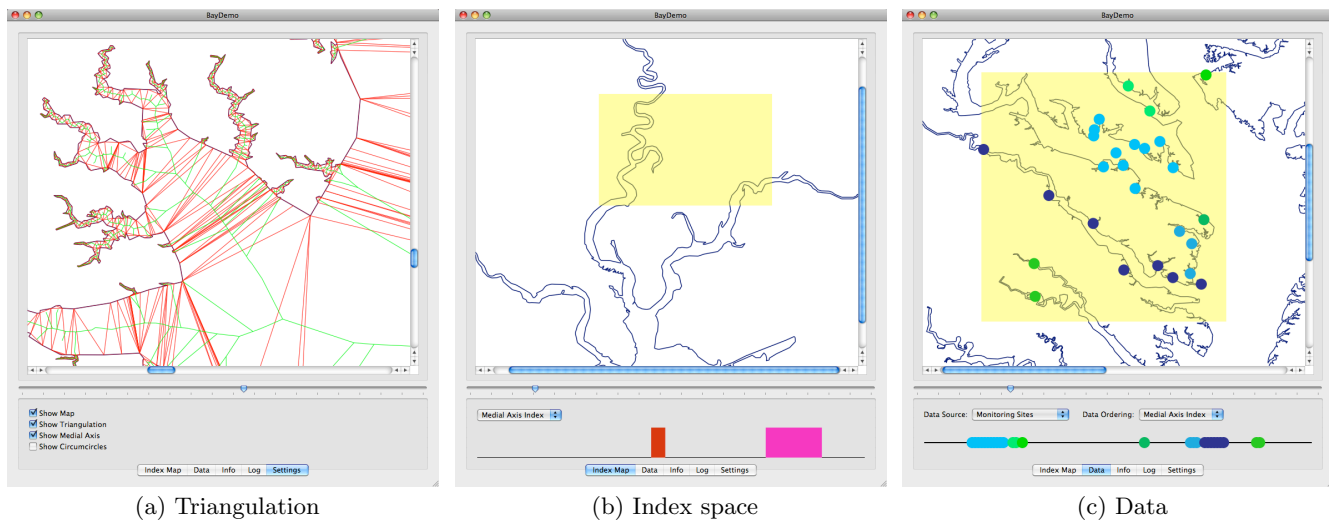


Figure 5: Screenshots of the demo showing (a) the triangulation used to generate the medial axis index, (b) the regions of index space needed to resolve a range query and (c) the data covered by a range query.

For the Chesapeake Bay, the user can perform queries on sample data sets. The demo will show the sites on the map inside the query region that hold data and display the location of that data (Figure 5(c)). The interface uses color to indicate which data are stored contiguously and can be read in a single I/O. For the same query region, the data layout often differs substantially from the index, because some index regions have little data whereas other index regions have lots. Disjoint sub-regions in the index reflect the local geometric complexity of the query region in the medial axis, whereas the data layout is data dependent and reflects the amount of information in each disjoint region. Large simple regions have simple medial axes and may contain lots of data, whereas small tributaries with complex medial axes may have little data. We have several data sources available, including subsets of measured and modeled environmental data from the Chesapeake Bay.

Finally, the demo allows users to compare the effectiveness of the medial axis spatial index with regular decompositions by selecting an indexing strategy in a drop down listbox. For a selected query region, the user can toggle between, for example, the medial axis index and a Morton-order space-filling curve and witness the relative contiguity in the index or of the data on the disk.

4. ACKNOWLEDGMENTS

We thank Rebecca Murphy, Bill Ball, and the other members of CBEO team for their expertise on estuarine science and Chesapeake Bay data. This work was supported by the National Science Foundation award ATM-0618986.

5. REFERENCES

- [1] W. P. Ball et al. A prototype system for multi-disciplinary shared cyberinfrastructure—Chesapeake Bay Environmental Observatory (CBEO). To appear in *Journal of Hydrological Engineering*, Accepted 2007.
- [2] CGAL, Computational Geometry Algorithms Library. <http://www.cgal.org>.
- [3] L. P. Chew. Constrained Delaunay triangulations. *Algorithmica*, 4(1):97–108, 1989.
- [4] M. McAllister and J. Snoeyink. Medial axis generalization of river networks. *CaGIS*, 27(2):129–138, 2000.
- [5] B. G. Nickerson. Automated cartographic generalization for linear features. *Cartographica*, 25(3):15–66, 1988.
- [6] H. Samet. *Foundations of Multidimensional and Metric Data Structures*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2006.