# Spatial Aggregation Using Space Filling Curves

**Dimitrios S. Dimitropoulos**

**Diploma Thesis**

Supervisor: Professor Nikolaos Mamoulis

Ioannina 10/2024

**ΤΜΗΜΑ ΜΗΧ. Η/Υ & ΠΛΗΡΟΦΟΡΙΚΗΣ**
**ΠΑΝΕΠΙΣΤΗΜΙΟ ΙΩΑΝΝΙΝΩΝ**

**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING**
**UNIVERSITY OF IOANNINA**

# Dedication

To my beloved mother and father, Aspasia and Stavros, whose strength and love continue to inspire me.

# Acknowledgements

I would like to thank Thanasis Georgiadis for his unwavering support and invaluable advice throughout the implementation of my Diploma Thesis.

I would also like to thank my supervisor, Professor Nikolaos Mamoulis, for his guidance and mentorship, which have been instrumental in the completion of this work.

15/9/2024

Dimitrios S. Dimitropoulos

# Περίληψη

Δημήτριος Σ. Δημητρόπουλος, Μ.Δ.Ε στην Πληροφορική, Τμήμα Μηχανικών Η/Υ και Πληροφορικής, Πανεπιστήμιο Ιωαννίνων, Ελλάδα, Σεπτέμβριος 2024.

Χωρική συνάθροιση με χρήση καμπυλών πλήρωσης χώρου.

Σύμβουλος: Νικόλαος Μαμούλης , Καθηγητής.

Σκοπός της παρούσας διπλωματικής εργασίας είναι η ανάπτυξη ενός εργαλείου με το όνομα ΧΚ-ΠΕΕ (Χίλμπερτ Καμπύλη Προσέγγιση Ερωτήματος Εύρους)  που θα υπολογίζει προσεγγιστικά το αποτέλεσμα ενός πεδίου σε ένα ερώτημα εύρους. Πιο συγκεκριμένα, το ΧΚ-ΠΕΕ υποστηρίζει χωρικά ερωτήματα εύρους στα οποία παρέχεται ένα γεωγραφικό πλαίσιο (ορθογώνιο) και υπολογίζει προσεγγιστικά το αποτέλεσμα του ερωτήματος για ένα καθορισμένο πεδίο και λειτουργία (μέγιστη, ελάχιστη και μέση τιμή). Το ΧΚ-ΠΕΕ αποτελεί μια νέα προσέγγιση που χρησιμοποιεί διαστήματα που προέρχονται από την καμπύλη πλήρωσης χώρου Χίλμπερτ για τον  υπολογισμό των αποτελεσμάτων.  Το ΧΚ-ΠΕΕ μπορεί να είναι έως και 26.73% ταχύτερο από τους ανταγωνιστές του και παρέχει αποτελέσματα με μέση ακρίβεια που κυμαίνεται από 92.5% έως και 98.5%.

**Λέξεις Κλειδιά:** Space Filling Curves, Hilbert Space Filling Curve, Quadtree, Uniform Grid, Spatial Aggregation, Spatial Queries.

# Abstract

Dimitrios S. Dimitropoulos, M.Sc. in Computer Science, Department of Computer Science and Engineering, University of Ioannina, Greece, September 2024.

Spatial Aggregation using Space Filling Curves.

Advisor: Nikolaos Mamoulis, Professor.


The purpose of this thesis is to develop a tool called HC-ARQ (Hilbert Curve Approximate Range Queries), that will approximate the result of a field in a range query. More specifically, HC-ARQ supports spatial range queries in which a geographic frame (rectangle) is provided and computes the result of the query for the specified field and function (maximum, minimum and average). HC-ARQ is a novel approach that uses intervals derived from the Hilbert Space Filling Curve to calculate the results. It can be up to 26.73% faster than its competitors while providing results with average accuracy ranging from 92.5% until 98.5%.


**Keywords:** Space Filling Curves, Hilbert Space Filling Curve, Quadtree, Uniform Grid, Spatial Aggregation, Spatial Queries.

# Table of Contents

# Chapter 1.        Introduction

In this chapter, we explain the concepts of Spatial Aggregation and Spatial Queries, providing a clear understanding of these key terms and how we provide a novel approach, the HC-ARQ, that can be faster than its competitors while maintaining strong accuracy.

## 1.1 Spatial Aggregation

Spatial aggregation, a fundamental concept in spatial analysis, refers to the process of computing aggregate statistics (e.g. sum, average etc.) for spatial data within a defined region. This is distinct from spatial summarization which focuses on data reduction by grouping or combining objects into larger units providing a higher-level overview of the spatial data. Spatial Aggregation is widely employed in various fields, including urban planning, epidemiology, and environmental studies, to gain a broader understanding of spatial patterns and relationships [1] [2] [3] [4].

One of the primary applications of spatial aggregation is in the field of urban research. Researchers can use this approach to study the spatial structure of cities, including the behavioral patterns of residents, by examining the aggregated data. For example, the average commute time or traffic congestion levels in different regions of the same city can be aggregated and analyzed to identify areas with transportation bottlenecks. This information can help urban planners and policymakers develop targeted strategies to improve public transportation routes, optimize traffic flow, and allocate resources for road expansions or alternative transport modes.

Similarly, spatial aggregation is employed in the analysis of crime data. By aggregating crime data into spatial units, researchers can identify "hot spots" of

criminal activity, which can inform law enforcement strategies and resource allocation [5]. Additionally, the temporal dimension of crime patterns can be explored by disaggregating the data by time, such as hour of day, day of week, or season, providing a more comprehensive understanding of the dynamics of crime opportunities in urban areas [5].

Furthermore, spatial aggregation has applications in the field of epidemiology, where it can be used to study the spatial distribution of diseases and identify areas with higher disease prevalence. By aggregating health data into larger spatial units, researchers can identify clusters or patterns of disease, which can inform public health interventions and resource allocation.

In the field of environmental studies, spatial aggregation is used to analyze the distribution of environmental factors, such as air pollution, water quality, or biodiversity. By aggregating these data into larger spatial units, researchers can identify regional trends and patterns, which can inform policy decisions and guide conservation efforts.

In conclusion, spatial aggregation is a powerful technique that enables researchers and decision-makers to gain a broader understanding of spatial patterns and relationships. By aggregating spatial data into larger units, researchers can identify trends, patterns, and hot spots that would be difficult to discern at the individual data point level. Leveraging this approach can lead to more informed decision-making and the development of effective strategies in a variety of domains, from urban planning to public health and environmental conservation [5] [1] [2].

Figure 1. provides an example of aggregation based on districts. The resulting table groups schools by their district (District A or District B) and shows the population for each type of school within the district.
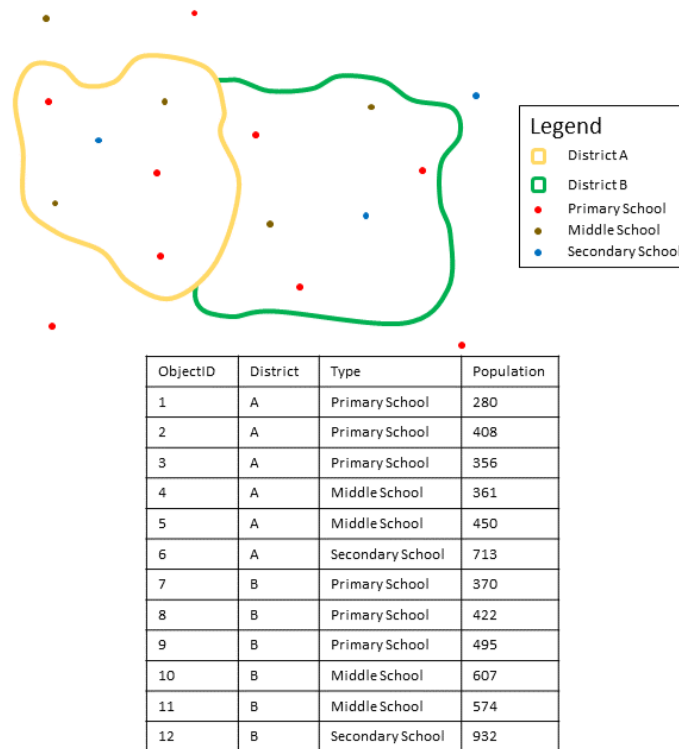
| ObjectID | District | Type | Population |
|---|---|---|---|
| 1 | A | Primary School | 280 |
| 2 | A | Primary School | 408 |
| 3 | A | Primary School | 356 |
| 4 | A | Middle School | 361 |
| 5 | A | Middle School | 450 |
| 6 | A | Secondary School | 713 |
| 7 | B | Primary School | 370 |
| 8 | B | Primary School | 422 |
| 9 | B | Primary School | 495 |
| 10 | B | Middle School | 607 |
| 11 | B | Middle School | 574 |
| 12 | B | Secondary School | 932 |

Figure 1: Spatial aggregation of school types in different districts [7]

# 1.2 Spatial Queries

Spatial queries are a fundamental component in the field of geospatial analysis and geographic information systems, as they enable users to retrieve and analyze data based on geographic or spatial relationships. These queries allow for the identification of spatial patterns, the establishment of relationships between geographical elements, and the ability to make decisions based on spatial information.

## 1.2.1 Common Types of Spatial Queries

Spatial queries can be categorized into various types based on the nature of their spatial relationships. Here are some of the most common types:

- Range Queries: Retrieve all objects within a specified spatial area (e.g., a rectangle, circle).

- k-Nearest Neighbour Queries: Retrieve the k closest objects to a given point or object based on a specified distance metric.

- Spatial Joins: Combine two spatial datasets based on their spatial relationships (e.g., intersecting, containing, or being within).

- Distance Joins: Retrieve pairs of objects that are within a certain distance from each other.

- Intersection Queries: Retrieve objects that share at least one common point with a specified object.

- Containment Queries: Retrieve objects that completely enclose another object.

- Within Queries: Retrieve objects that are entirely located within another object or region.

- Crosses Queries: Retrieve objects that intersect but do not completely contain each other.

## 1.2.2 Applications of Spatial Queries

One key application of spatial queries is in the analysis of big spatial data, which has become increasingly prevalent with the rise of technologies such as satellite imagery, GPS, and social media. In this context, range queries are commonly used to filter data within a specified region or bounding box. Additionally, spatial joins are applied to merge datasets, while k-Nearest Neighbour queries are used to find the closest features to a given location, such as identifying the nearest city to a disaster site.

To handle this large and complex data researchers have developed various approaches to handle this voluminous and complex data, including on-top, from-

scratch, and built-in methods for spatial query processing [15]. These techniques leverage spatial indexing and high-level query languages to enable efficient querying and analysis of big spatial data, which has proven valuable in domains such as urban planning, transportation, and environmental monitoring.

Another significant application of spatial queries is in the field of data visualization and visual communication. Spatial data, particularly when combined with temporal information, can provide valuable insights when represented visually. In this case, spatial aggregation queries (e.g., calculating population density in specific regions) are crucial to summarize spatial information for better visualization. Additionally, range queries and within queries help to filter data within certain geographic boundaries to make the visual representation clearer and more focused. Scholars have explored how to effectively represent and visualize the results of spatial data analysis, utilizing interactive technologies and human-computer interaction techniques to enhance the understanding of spatial patterns and relationships.

Spatial queries are increasingly being used in a wide range of real-world applications, from studying the spatial structure of cities and understanding residents' behavioral patterns to analyzing and visualizing environmental data.

As the world becomes increasingly connected and digitized, the need for efficient and effective spatial query techniques will only continue to grow.

### 1.2.3 Spatial Aggregation in Defined Areas

Beyond simple retrieval based on spatial relationships, spatial queries often involve aggregation within user-defined areas. This allows for calculating statistics and summarizing data within specific regions of interest. Common aggregation functions used in spatial queries include:

- Max: Finds the maximum value of a specific attribute within the defined area.
- Min: Finds the minimum value of a specific attribute within the defined area.

- **Average:** Calculates the average value of a specific attribute within the defined area.

- **Sum:** Calculates the sum of a specific attribute for all objects within the defined area.

- **Count:** Counts the number of objects present within the defined area.

These aggregation functions, combined with user-defined areas, provide a powerful tool for analyzing spatial patterns and trends. For example, you can use these operations to:

- **Analyze pollution levels:** Calculate the average concentration of pollutants within a city boundary.

- **Study urban density:** Determine the maximum building height within different zoning districts.

- **Assess traffic congestion:** Count the number of vehicles present within a specific road segment during peak hours.

By defining the area of interest and applying appropriate aggregation functions, spatial queries enable valuable insights from spatial data.

# 1.3 HC-ARQ

HC-ARQ maps every geographic coordinate of the dataset to the corresponding Hilbert Curve index. These indexes are grouped into continuous intervals, which are stored on disk along with statistics representing the areas they cover. When a range query is executed, the program calculates all the Hilbert Curve indexes that fall within the query area.

These indexes are then grouped into new intervals based on their continuity. The system compares the new set of intervals with the pre-stored intervals to extract all the relevant information for the query.

This innovative approach can be up to 26.73% faster than competing methods, while maintaining accuracy between 92.5% and 98.5%.

# Chapter 2.    Background

In this chapter, we aim to thoroughly discuss the primary competing methodologies to our developed tool.

Additionally, we will provide a comprehensive background on Space-Filling Curves, with a particular focus on the Hilbert Space-Filling Curve, which is central to our tool development.

## 2.1 Uniform Grid

In the realm of spatial data processing, spatial aggregation and spatial queries are fundamental operations that allow for the efficient analysis and manipulation of geospatial information. While advanced techniques such as quadtrees and space-filling curves have gained popularity, the uniform grid iteration approach offers a viable alternative that is worth considering [8] [9].

The uniform grid iteration method partitions the spatial domain into a regular grid of cells, each representing a specific geographic area. Each cell contains the spatial objects that fall withing that geographic area. This approach leverages the inherent structure of the grid to perform spatial queries and aggregations, providing a straightforward and intuitive framework for working with spatial data [8].

The main advantage of the uniform grid iteration technique is its precision – just like the Quadtree. By using a regular grid, the method can accurately represent the spatial relationships and characteristics of the underlying data [18]. This level of precision is particularly important in applications where the spatial location and extent of objects are critical, such as in geographic information systems or location-based services.

However, the uniform grid iteration approach does come with a trade-off: it can be more computationally intensive compared to more advanced techniques like Quadtrees and space-filling curves [8][15] which we describe in detail at sections 2.2 and 2.3 respectively. Quadtrees and space-filling curves are often more efficient in terms of storage and query performance, as they can effectively partition the spatial domain in a hierarchical or space-filling manner, respectively.

The time complexity and space complexity of the uniform grid iteration method can be more significant compared to quadtrees and space-filling curves [15][12][13][14].

Despite this trade-off, the uniform grid iteration approach remains a viable option, particularly in situations where the spatial data is relatively simple, and the emphasis is on maintaining high precision.

Researchers have explored various strategies to optimize the performance of the uniform grid iteration method, such as leveraging parallel processing or implementing efficient spatial indexing mechanisms.

Furthermore, the uniform grid iteration time depends on a lot in the complexity of the spatial objects and queries.

Taking everything into consideration, we need to carefully consider all the requirements of the application and the characteristics of the spatial data when deciding on the most appropriate spatial data processing technique [18] [12] [13] [14].

The Figure 2. provides an example of uniform grid iteration. More specifically, the numbers inside the cells of the figure are indicating an order (y-axis iteration) in which the uniform grid will visit the cells in an iteration.

Figure 2:  Uniform Grid Iteration Example

## 2.2 Quadtrees

One of the fundamental data structures used in spatial data management is the Quadtree, a hierarchical tree-like structure that partitions a two-dimensional space into a series of nested quadrants.

The Quadtree is a powerful tool for spatial aggregation and spatial queries, offering several benefits over alternative methods.

Spatial aggregation, the process of grouping spatial data points into larger units, is crucial for summarizing and visualizing large datasets. Quadtrees excel at this task, as they can efficiently partition the space and aggregate data at different levels of granularity.

Spatial queries, on the other hand, involve retrieving data based on spatial relationships, such as finding all points within a certain distance of a target location. Quadtrees support efficient spatial queries by providing a structured index that allows for rapid identification of relevant data [16] [17].

Furthermore, Quadtrees offer several advantages over other spatial data structures, such as the R-tree and the k-d tree. Unlike R-trees, which can suffer from overlapping bounding boxes, Quadtrees use a fixed partitioning scheme that

guarantees non-overlapping regions. This, in turn, can lead to more efficient query processing, as the search space can be more effectively pruned. In comparison to k-d trees, Quadtrees are better suited for handling data with varying spatial distributions, as they can adaptively subdivide regions with higher data density [16].

The benefits of the Quadtree approach are numerous. It offers a compact and hierarchical representation of spatial data, enabling efficient storage and retrieval. Additionally, the hierarchical structure of the Quadtree allows for multi-resolution analysis, where data can be accessed at different levels of detail depending on the query [1]. Furthermore, the Quadtree's ability to handle overlapping and irregularly shaped spatial objects makes it a versatile choice for a wide range of applications [18]. More specifically, Quadtrees have been successfully applied in fields such as computer graphics [19], geographic information systems [16], and multimedia databases [20].

However, Quadtrees are not without their drawbacks. While they excel at spatial aggregation and spatial queries, they can be less efficient for high-dimensional data or for data with highly irregular shapes. The partitioning of the space can lead to increased complexity, particularly when dealing with dynamic or irregularly shaped spatial data. Additionally, the performance of Quadtree-based spatial operations can be sensitive to the distribution of the data, with unbalanced or skewed data potentially leading to suboptimal performance [16][18].

Overall, the Quadtree is a powerful and widely used data structure for spatial data management, offering significant advantages for spatial aggregation and spatial queries.

The Figure 3. illustrates a Quadtree spatial partitioning where the 2D space is recursively divided into quadrants and the respective quadtree structure.
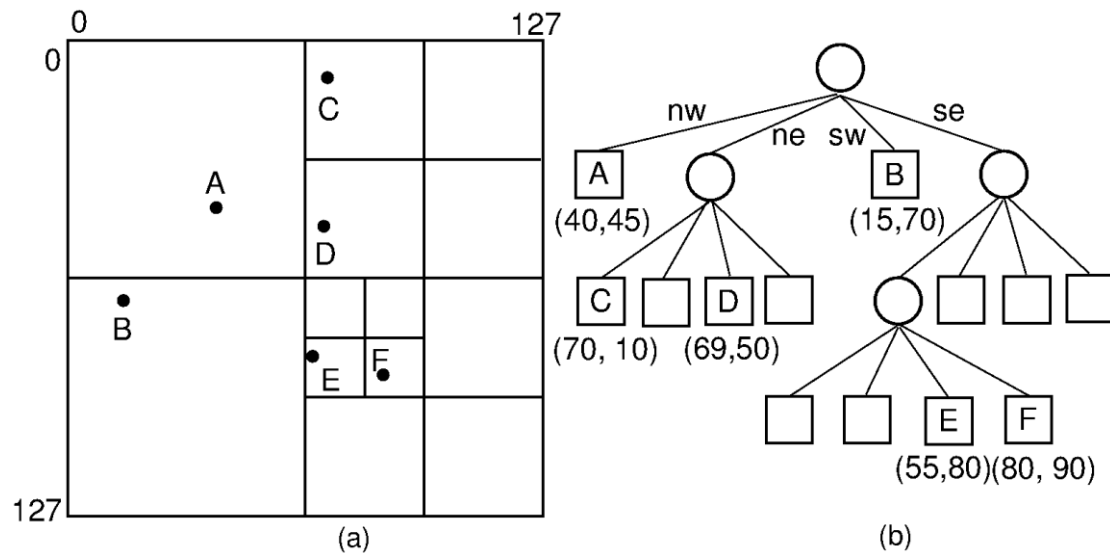
Figure 3: Quadtree Example [10]

# 2.3 Space Filling Curves

Space filling curves are an interesting mathematical concept that have found numerous applications in various fields, including computer science, physics, and geography.

These curves possess the remarkable property of mapping a one-dimensional space onto a higher-dimensional space, effectively "filling" the higher dimensional space with a continuous, self-similar curve [24].

Among the diverse family of space filling curves, the Hilbert curve has emerged as a particularly prominent and widely studied example [24]. The Hilbert curve is characterized by its ability to preserve the spatial locality of data points, meaning that neighboring points in the original space remain close to each other in the mapped, higher-dimensional space. This property makes the Hilbert curve a valuable tool for spatial aggregation and spatial queries, as it allows for efficient indexing and retrieval of data.

One of the primary applications of space filling curves, and the Hilbert curve in particular, is in the field of spatial aggregation. By mapping higher-dimensional spatial data onto a one-dimensional curve, the Hilbert curve enables efficient storage and querying of spatial information. This is particularly useful in scenarios where large datasets need to be organized and accessed, such as in geographic information systems, image processing, and sensor networks.

In the context of spatial queries, the Hilbert curve's ability to preserve spatial locality is crucial. By transforming spatial data into a one-dimensional representation, the Hilbert curve allows for the use of efficient one-dimensional indexing structures, such as B-trees or R-trees, to facilitate fast retrieval of data based on spatial proximity. Furthermore, the self-similar nature of the Hilbert curve enables the use of hierarchical indexing schemes, allowing for efficient range queries, nearest neighbor searches, and other spatial operations [25] [26].

While the Hilbert curve is a prominent example, other space filling curves, such as the Peano curve, Sierpinski curve, and Morton curve, also possess unique properties and have found applications in various domains. These curves differ in their specific characteristics, such as the degree of locality preservation, computational complexity, and visual patterns, making them suitable for different types of spatial data and applications.

The main advantages of space filling curves, including the Hilbert curve, can be summarized as follows:

- Efficient storage and indexing of spatial data: By mapping higher-dimensional data onto a one-dimensional curve, space filling curves enable the use of well-established one-dimensional indexing structures for efficient storage and retrieval.

- Preserving spatial locality: The Hilbert curve, in particular, excels at preserving the spatial relationships between data points, allowing for effective spatial aggregation and querying.

- Hierarchical organization: The self-similar nature of space filling curves enables the use of hierarchical indexing schemes (e.g. quadtrees by mapping quadtree nodes to the Hilbert Curve), further enhancing the efficiency of spatial operations.

Furthermore, the Hilbert Space Filling Curve can be effectively combined with other techniques, such as data compression (by reducing redundant spatial data through locality-preserving mapping) and parallelization (by dividing the curve into segments that can be processed independently), to further enhance its performance and scalability.

However, the usage of Space Filling Curves in spatial data management has drawbacks as well. More specifically, the transformation from multi-dimensional to one-dimensional space can lead to some loss of information (points close together in the original space may not always be close in the 1D mapping), and the optimal choice of parameters for the Hilbert Curve can be challenging in certain scenarios.

While the advantages are well-established, the choice of the appropriate spatial data management technique depends on the specific requirements of the application, such as the dimensionality of the data, the frequency and types of spatial queries, and the overall data volume.

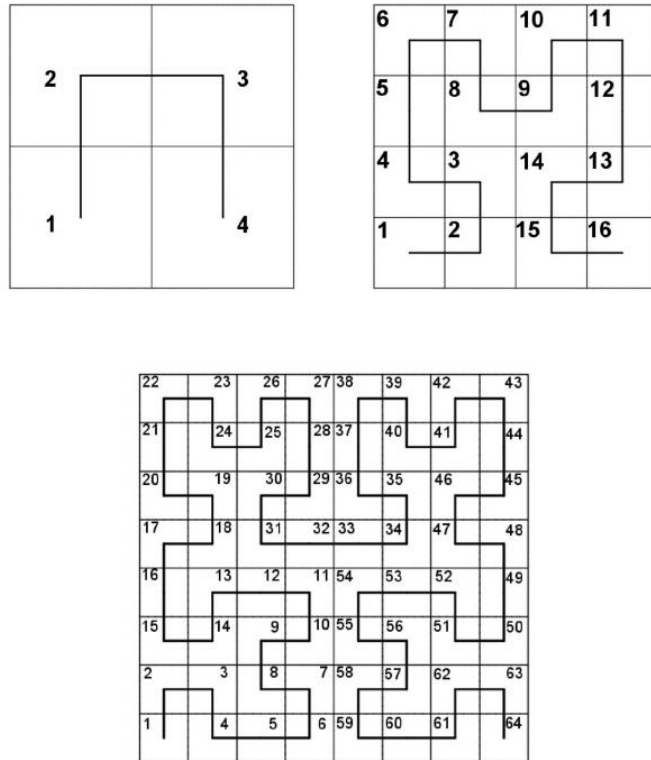Figure 4. illustrates how the Hilbert Space Filling curve traverses space at different orders.

Figure 4: Different Hilbert Curve Orders Graphs [11]

# Chapter 3. Methodology

In this section we aim to explain the methodology that we followed for the development of HC-ARQ which is a process that utilizes the Hilbert Space Filling Curve for efficient spatial aggregation and querying using intervals.

More specifically, spatial data points are aggregated into intervals based on their continuity using Hilbert curve transformations as a pre-processing step. Then spatial aggregation queries, such as calculating averages, maximums, and minimums, are performed over these intervals on specific attributes.

We selected the Hilbert Space Filling curve because it preserves spatial locality, meaning that points that are close together in 2D space will also be close together in 1D Hilbert space. This allows for evaluating distance-dependent spatial queries such as range queries, on the interval data directly.

## 3.1 Pre-processing

In the Pre-processing step we transform the longitude and latitude of the data entries into Hilbert orders and group them into intervals.

### 3.1.1 Grid Partitioning

The purpose of grid partitioning is to discretize space so that spatial objects, such as points representing crime incidents, can be grouped into cells. The space is typically partitioned by dividing the range of latitude and longitude into equal intervals.

For example, if the range of latitudes is $[lat_{min}, lat_{max})$ and longitudes $[lon_{min}, lon_{max})$ and we want to partition the space into an NxN grid, the width and height of each cell are determined via the following formulas:

$$cell\_width = \frac{lon\_max - min\_lon}{N}$$

$$cell\_height = \frac{lat\_max - lat\_min}{N}$$

Each data point is then assigned to a specific grid cell based on its latitude and longitude.

Then, we map the latitude and longitude coordinates of each non-empty cell of the grid into its corresponding Hilbert curve value.

Below in Figure 5. is the code of the method that we used for the cells x and y indexes to Hilbert Order value calculation.

```cpp
uint32_t xy2d(int n, int x, int y) {
    int rx, ry, s;
    uint32_t d = 0;
    for (s = n / 2; s > 0; s /= 2) {
        rx = (x & s) > 0;
        ry = (y & s) > 0;
        d += s * s * ((3 * rx) ^ ry);
        if (ry == 0) {
            if (rx == 1) {
                x = n - 1 - x;
                y = n - 1 - y;
            }
            std::swap(x, y);
        }
    }
    return d;
}
```

Figure 5: The Hilbert Order transformation algorithm [21].

## 3.1.2 Hilbert Order Grouping

After the calculation of the Hilbert Order of all the data entries we group the orders into intervals based on the continuity – continuous Hilbert orders are grouped together.

For example, for the Hilbert Orders [543,544,545,549,550,552] we have three intervals.

The first one is the [543,546), the second one is the [549,551) and the last one is the [552,553). Furthermore, we have the opportunity to specify the maximum number of Hilbert orders that can be grouped into intervals in order to avoid really big intervals.

After this step we have produced the intervals that we will use in spatial queries. The produced intervals store the aggregated information about the enclosed point data such as minimum, maximum and average values along with the number of points and Hilbert orders grouped into the interval. The intervals are stored on disk for persistence and are loaded only once to be used in the query step.


## 3.2 Queries Evaluation

The program reads the intervals that have been produced in the previous step (*pre-processing*) from the disk and calculates the dataspace boundaries.

Afterwards, it reads from the disk all the queries that need to be performed and then executes them one by one. The queries have the following format: *x1, y1, x2, y2, operation and attribute*, where *x1 < x2, y1 < y2, operation* can be either minimum, maximum or average and *attribute* is a data column (e.g. age).

The (x1, y1) and (x2, y2) define a rectangle which is the query window.

In order to find out the in-range Hilbert Order values of the query window we need to make an iteration through the query window boundaries.

Figure 6. show the code used for the iteration.

```
double width = (x_max - x_min) / num_cells;
double height = (y_max - y_min) / num_cells;

for (double i = y_min; i <= y_max; i+=height) {
    for (double j = x_min; j <= x_max; j+=width) {
        /* lat lon to Hilbert Order */
    }
}
```

Figure 6: Hilbert Orders in range iteration algorithm

For example in the second order of Hilbert Space Filling curve and the points
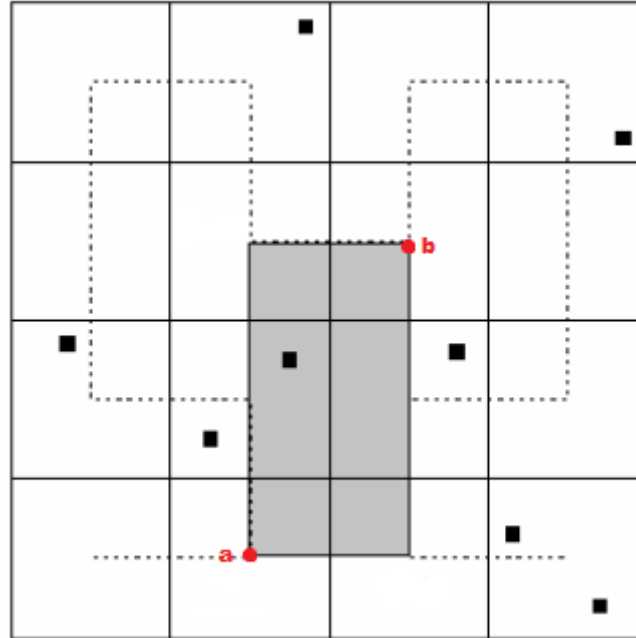a = (1,0) and b = (2,2) we have the following graph.



Figure 7: Hilbert Orders in range has part of their cell colored in gray

The iteration will take place from x = 1 until x = 2 and from y = 0 until y = 2.

The cells that fall inside the query window are the (1,0), (2,0), (1,1), (2,1), (1,2), (2,2). The graph can verify that these are indeed the in range cells.

After identifying the cells within the query area, we transform these cells into their corresponding Hilbert curve values (Hilbert orders). These Hilbert orders are then grouped into intervals based on their continuity along the Hilbert curve.

Once the intervals within the query range are determined, we retrieve the relevant intervals from the pre-processing step.

More specifically, to identify the relevant intervals, we compare the two sets of intervals (pre-processing and query intervals). If there is an overlap, we calculate the extent of the overlap and assign a weight based on the proportion of it.

# 3.3 Weight Strategies

Weighting refers to assigning different levels of importance to each data interval based on its relevance to the query. It is an important step because not all intervals contribute equally to the query (e.g. some intervals may overlap more or have higher point density).

We have developed two ways for the weighting of the intervals. The first one is a simple coverage-weight calculation based on how many common cells the range intervals have with our HC-ARQ intervals and the second one is a density-coverage calculation that takes into consideration both the point density of the overlapping intervals and the area coverage.

## 3.3.1 Coverage Weighting

The coverage-weighting method is based on the fraction of the interval that overlaps with the query region. More specifically on how much of the spatial query area it covers. This technique helps us estimate how relevant the interval for the query region is.

The coverage – weighting woks as follows:

- Each interval represents a portion of the entire spatial area being analyzed.

- The coverage of an interval is how much of that interval intersects with or lies within the spatial query.

- The more an interval or cell overlaps with the query region, the more it contributes to the result.

$$W_i = \frac{\text{Area of interval } i \cap \text{ Query Region}}{\text{Total Area of Query Region}}$$

where *Area of interval i ∩ Query Region* is the portion of the interval that overlaps with the query region and *Total Area of Query Region* is the total area being queried.

For example, if a query region covers 50% of interval 1 and 100% of interval 2, the weights for those intervals would be:

$$W_1 = \frac{0.5 \times \text{Area of Interval 1}}{\text{Total Query Area}} \qquad W_2 = \frac{1.0 \times \text{Area of Interval 2}}{\text{Total Query Area}}$$

Figure 8. displays how the query interval = [1,5) overlaps with the interval 1 = [1,3) and the interval 2 = [4,6) of the pre-processing step.



Figure 8: The overlap of the query interval and the pre-processing intervals

This method is more suitable when the spatial extent of each interval matters more than the number of points it contains. This method is ideal if the goal is to ensure that regions with greater geographic overlap with the query region have a proportionally larger influence.

### 3.3.2 Density-Coverage Weighting

Density coverage weighting builds on the idea of coverage weighting but adds an additional factor related to the density of data points within each interval.

This method adjusts the weight assigned to each interval not only by how much area it covers but also by how many data points are present in that interval.

The density-coverage – weighting woks as follows:

- Like coverage weighting, this method starts by calculating how much of an interval overlaps with the query region.

- Then it adjusts the weight based on the density of data points in the interval. Intervals with more points (higher density) are given higher weights, reflecting that these intervals contain more data and are potentially more important for the analysis.

$$W_i = \frac{\text{Area of interval } i \cap \text{ Query Region}}{\text{Total Area of Query Region}} \times \frac{\text{Number of Points in Interval } i}{\text{Total in range Points}}$$

where the first term is the coverage and the second is the point density.

For example, if a query region covers 50% of interval 1 which has 10 points and 100% of interval 2 which has 2 points, the weights for those intervals would be:

$$W_1 = \left( \frac{0.5 \times \text{Area of Interval 1}}{\text{Total Query Area}} \right) \times \frac{10}{\text{Total in range Points}}$$

$$W_2 = \left( \frac{1.0 \times \text{Area of Interval 2}}{\text{Total Query Area}} \right) \times \frac{2}{\text{Total in range Points}}$$

where Total in range Points equals to 12.

In this case even though Interval 1 covers less of the query region, it might contribute more to the final result because it has a higher point density than the Interval 2.

Figure 9. displays how the query interval = [1,5) overlaps with the interval 1 = [1,3) and the interval 2 = [4,6) of the pre-processing step along with the point density of each interval.
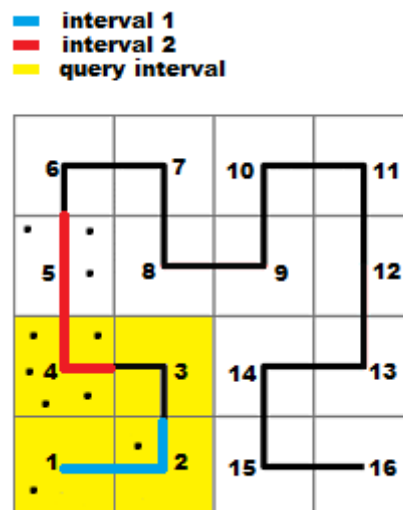


Figure 9: The overlap and the point density of the query interval and the pre-processing intervals

This method is more suitable when the density of data points within an interval is also important, not just the area it covers. This is useful in situations where intervals with more data points are considered more reliable or important for the analysis.

### 3.3.3 Weighting Methods Summary

Coverage Weighting focuses purely on the area overlap between intervals and the query region, making it suitable for spatial studies where area is a key factor.

Density Coverage Weighting incorporates data point density along with area coverage, giving higher importance to regions that are more data-rich, which is ideal for analysis where both the number of observations and the spatial extent are important factors.

The proposed weighting methods aim to improve the efficiency and accuracy of HC-ARQ by assigning importance to intervals based on coverage and data density while not adding significant computational overhead.

Taking everything into consideration, HC-ARQ is approximate due to two main factors, the use of the Hilbert Space Filling Curve, which introduces minimal distortions in spatial relationships when converting 2D space into 1D, and the grid resolution. The grid discretizes the spatial data into cells, and points within the same cell are treated as having the same location, leading to small approximations in both query results and the assignment of weights. The degree of approximation is controlled mainly by the grid resolution - finer grids provide higher accuracy, while coarser grids prioritize computational efficiency.

# Chapter 4.    Experiments

In this section, we provide detailed information about the development environment and the dataset used for conducting experiments with the tool.

Additionally, we present the results in terms of accuracy and performance time, offering a comparative analysis of our tool against the main competitors (Quadtree and Uniform Grid).

## 4.1 Datasets

The original dataset is a 52K entries dataset downloaded from the open-source data science community Kaggle [23] about homicides over the past decade in 50 of the largest American cities collected from the Washington Post.

The dataset has the following fields: id, longitude, latitude, victim age and victim race. The dataset was expanded using Python's Pandas library and the random sampling with replacement technique to 500K entries.

Random sampling with replacement is a sampling technique where each element in the population has an equal chance of being selected multiple times for a sample. In this method, after an element is selected, it is "replaced" back into the population, meaning it can be selected again in subsequent draws.

The expanded data set reflects similar distribution patterns to the original dataset as you can see at Figures 9.1-9.2 and Figures 10.1-10.2.

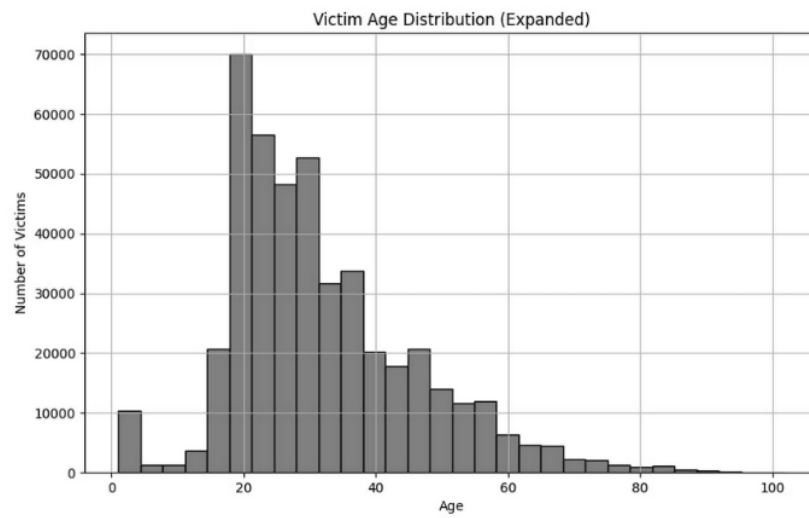Figure 9.1: The original data victim age distribution



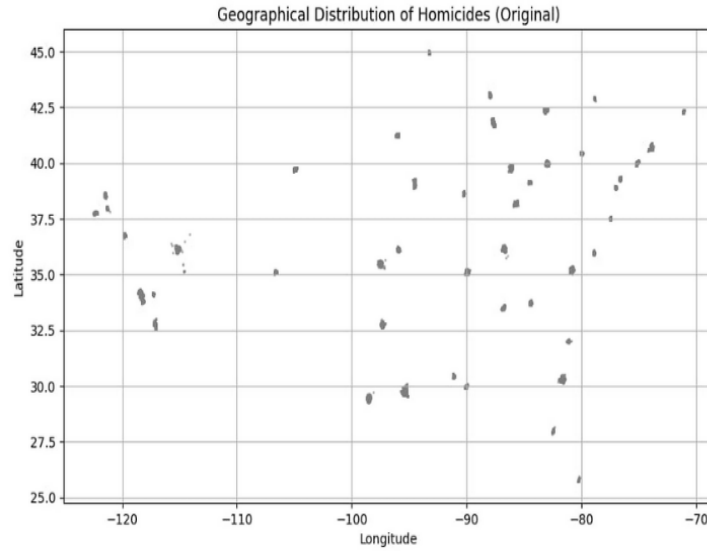Figure 9.2: The expanded data victim age distribution
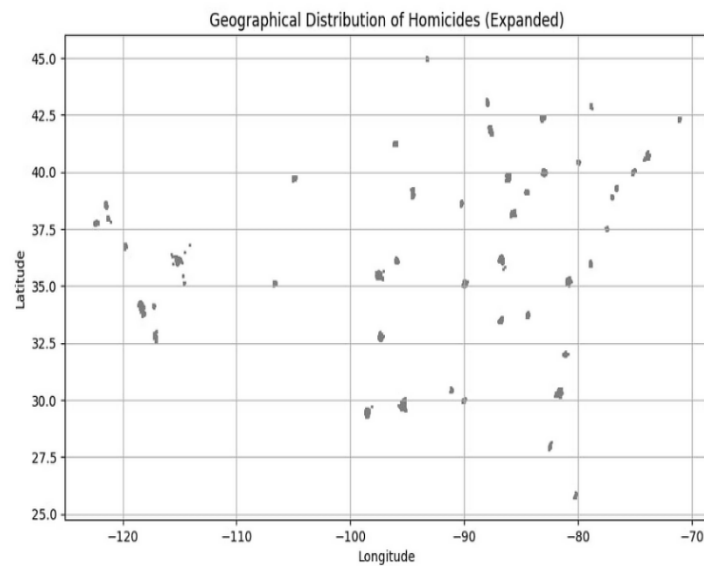
Figure 10.1: The original data point distribution



Figure 10.2: The expanded data point distribution

## 4.2 Development Environment

During development, all experiments were conducted in an HP EliteBook 8460p computer.

The operating system of the computer was the Microsoft Windows 10 Home.

The computer had 8GB of RAM and the CPU was the Intel(R) Core (TM) i5-2520M CPU at 2.50GHz with 2501 MHz, 2 Cores and 4 Logical Processors.

The code of the application was developed in the C++ programming language and GCC compiler.

# 4.3 Results

## 4.3.1 Range Queries

We executed multiple range queries for different areas and operations. For each query we selected randomly a point (x,y) from the dataspace and the query window was defined based on this point as follows : $(x+x \times i, x-x \times i)$ for the x-axis and $(y+y \times i, y-y \times i)$ for the y-axis, where $i$ takes values as 0.00001, 0.0001, 0.001 and 0.01. The parameter $i$ controls the size of the query window, with larger values of $i$ creating larger query windows.

For example, we randomly select the point x = 10 and y = 20 from the dataspace. If the value of the parameter $i$ is set at 0.001 the query window for the x-axis and y-axis would be calculated as follows:

- For the x-axis: $(10+10 \times 0.001, 10-10 \times 0.001) = (10.01, 9.99)$
- For the y-axis: $(20+20 \times 0.001, 20-20 \times 0.001) = (20.02, 19.98)$

Thus, the query window would cover the rectangular area with corners at (9.99, 19.98) and (10.01, 20.02).

### 4.3.2  Accuracy Comparison

#### 4.3.2.1  Definition of Accuracy

We compared every query entry result of the tool with the precise results of the quadtree.

The formula that was used for the accuracy was the following:

$$\text{Accuracy} = 1 - \frac{V_2 - V_1}{V_2}$$

where *V1* the *min (Interval Result, Quadtree Result)* and *V2* the *max (Interval Result, Quadtree Result)*.

For example, for *Interval Result* = 5 and *Quadtree Result* = 2 the accuracy of the tool would be 40%.

After calculating the accuracy of each entry individually, we sum up the individual values and divide them by the number of entries to find the overall accuracy of the tool in the query set.

#### 4.3.2.2  Accuracy Comparison Results

We tested the HC-ARQ in different Hilbert orders and with both weighting techniques at 100 query entries per test in all the supported operations (min, max, avg).

In all test cases, as the grid size increases from $2^9$ to $2^{13}$, the accuracy improves, with significant gains observed between the $2^9$ and $2^{11}$ grid sizes.

For larger values of *i*, the accuracy remains consistently high (around 98%) across all grid sizes, with minimal change.

However, as *i* decreases, the initial accuracy at smaller grid sizes is lower, showing more noticeable improvement as the grid size increases.

More specifically, for very small values of $i$, the accuracy starts around 92% at $2^9$ and reaches approximately 98% at $2^{13}$.

Across almost all graphs, density-coverage slightly outperforms the coverage weighting technique, indicating that the density-coverage method is more reliable regardless of the grid size or the value of $i$.

Taking everything into consideration, the grid size plays a more critical role in enhancing accuracy when $i$ is small, while for larger $i$, the accuracy remains high even for smaller grid sizes.

Figures 11.1-11.4 display the average accuracy of the developed tool across different values of $i$ for different grid resolutions.



Figure 11.1: The accuracy-grid size graph for i=0.00001
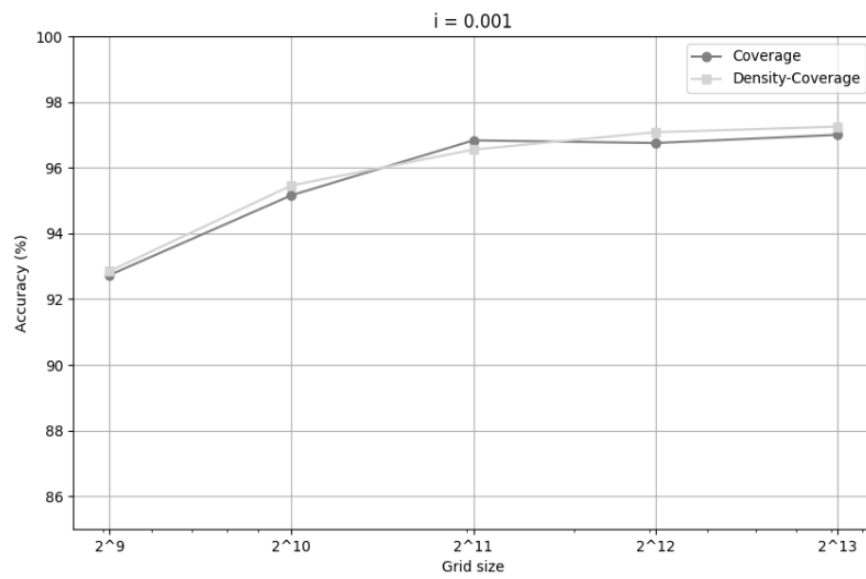
Figure 11.2: The accuracy-grid size graph for i=0.0001



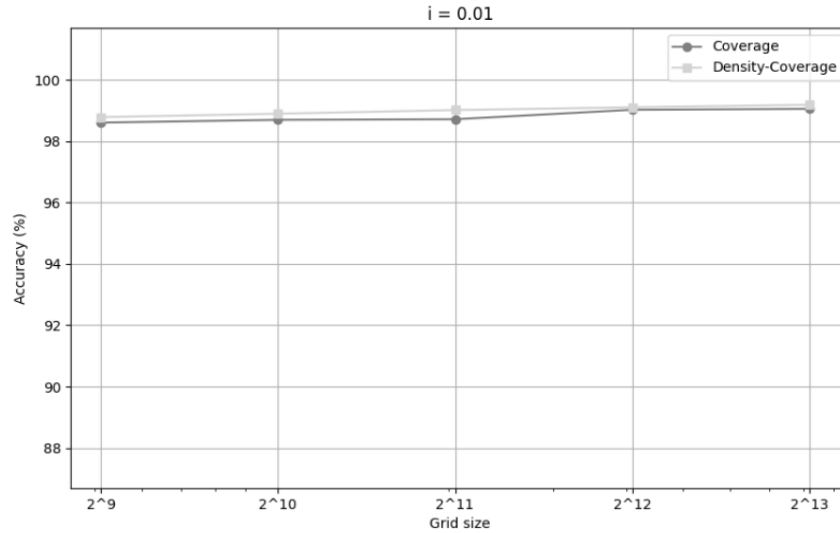Figure 11.3: The accuracy-grid size graph for i=0.001

Figure 11.4: The accuracy-grid size graph for i=0.01

## 4.3.3 Time Comparison Results

We compared the developed tool with two competitors – the Quadtree and the Uniform Grid. The comparison was divided into Pre-processing time and Query Execution time.

In all cases the data are being read from the disk entry by entry.

### 4.3.3.1 Pre-processing Time

The HC-ARQ pre-processing time increases slightly as the grid size increases.

As the grid size becomes finer, the number of Hilbert Orders and Intervals increases, leading to longer Pre-processing times.

For even finer grids the Pre-processing time of the tool surpass the Quadtree Pre-processing time.
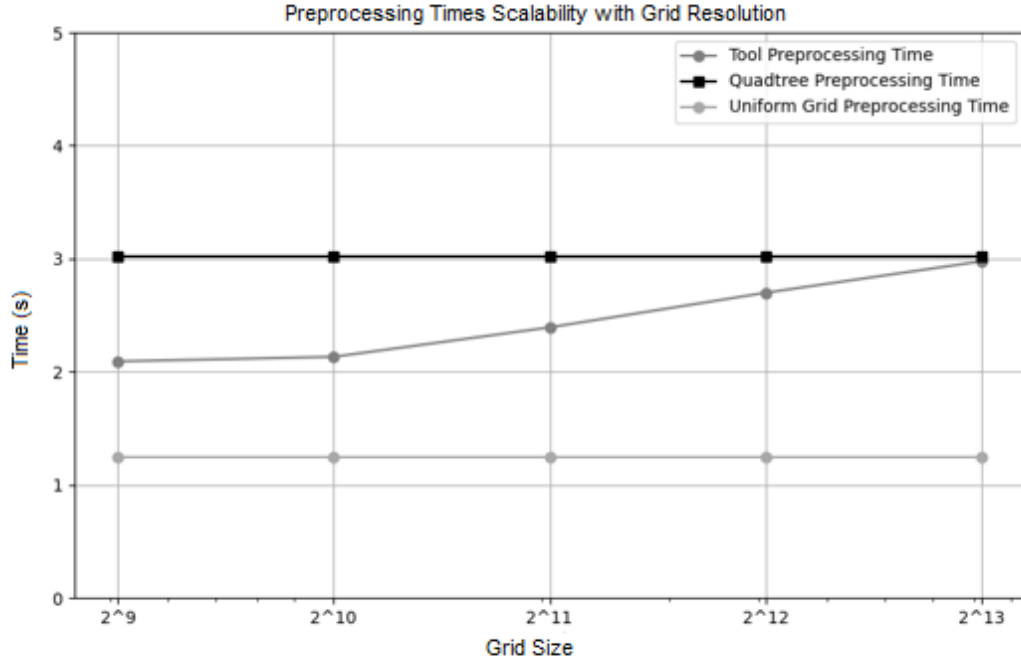
Figure 12: The Pre-processing times for the different grid sizes

### 4.3.3.2 Query Execution Time

The execution time of the tool increases significantly at finer grid sizes due to the increase of the Hilbert Orders and Intervals.

As the grid resolution increases, the computational demand rises, particularly due to the additional calculations required to transform longitude and latitude into the corresponding Hilbert Orders.

Furthermore, the larger the number of Hilbert Orders, the more iterations are needed to identify the Hilbert Orders that fall within the query range.

Figures 13.1-13.4 display the need time for the execution of the query set for the developed tool across different values of *i* for different grid resolutions compared to the main competitors of the tool, the Uniform Grid and the Quadtree.

It is clear from the figures that HC-ARQ can achieve faster query execution times in all cases for $2^9$ and $2^{10}$ grids against both competitors and better than the Uniform Grid at all cases of $2^{11}$ and $2^{12}$ grids except the case in which $i$=0.01.
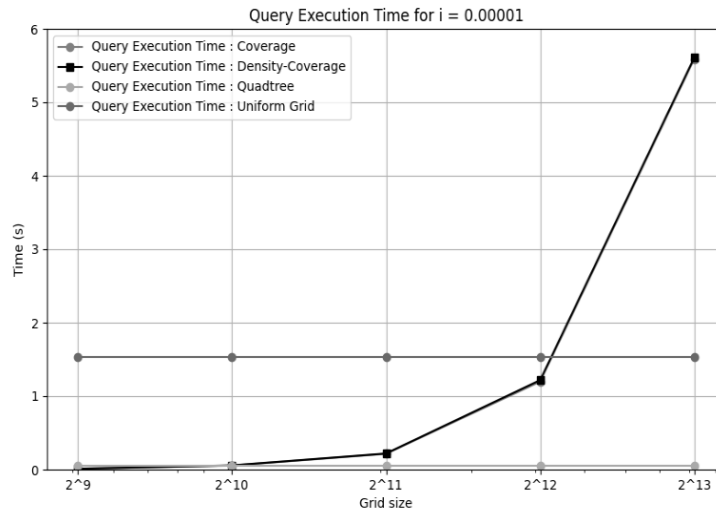


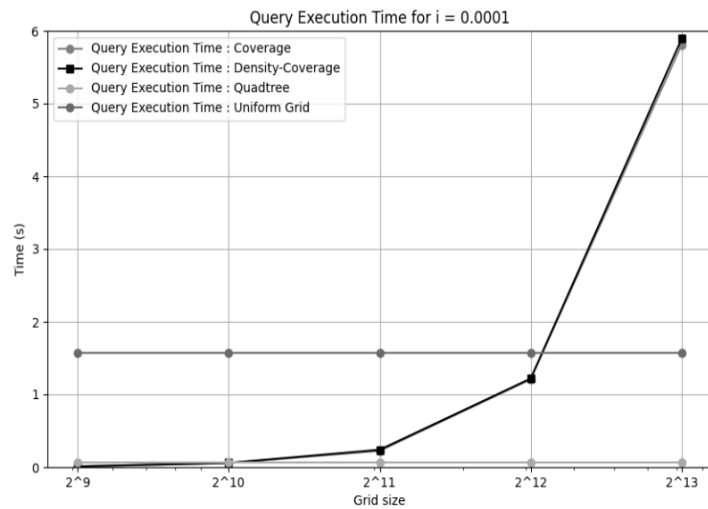Figure 13.1: The query execution times in seconds for i=0.00001



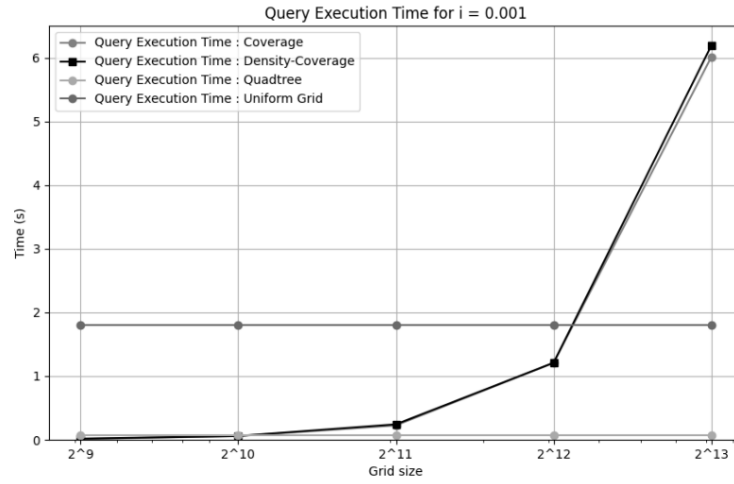Figure 13.2: The query execution times in seconds for i=0.0001

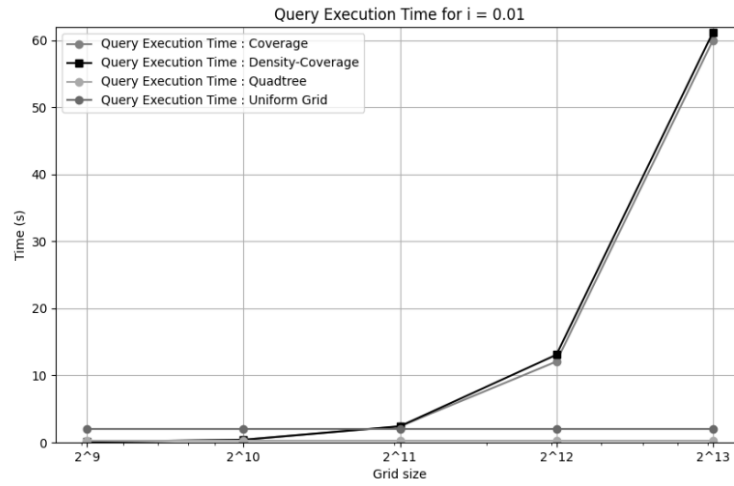Figure 13.3: The query execution time in seconds for i=0.001



Figure 13.4: The query execution time in seconds for i=0.01

### 4.3.3.3 Total (Pre-processing and Query Execution) Time

HC-ARQ experiences a significant total time increase at higher grid resolutions, especially for larger values of $i$.

This follows the trend from query execution times, where increased grid sizes lead to greater computational overhead due to the handling of Hilbert Orders.

For tasks involving smaller grid sizes or smaller values of *i* (smaller query window), the tool performs well. However, at finer grid resolutions, its time-efficiency declines due to the computational complexity associated with Hilbert Order processing. More specifically, for finer grid resolutions, the cell width and height become significantly smaller – more cells, leading to a substantial increase in the number of latitude and longitude to Hilbert Order conversions (Figure 5 and Figure 6). This introduces a significant computational overhead.

Figures 14.1-14.4 display the total execution time across different values of *i* and grid resolutions compared to the tool's main competitors - the Quadtree and the Uniform Grid.
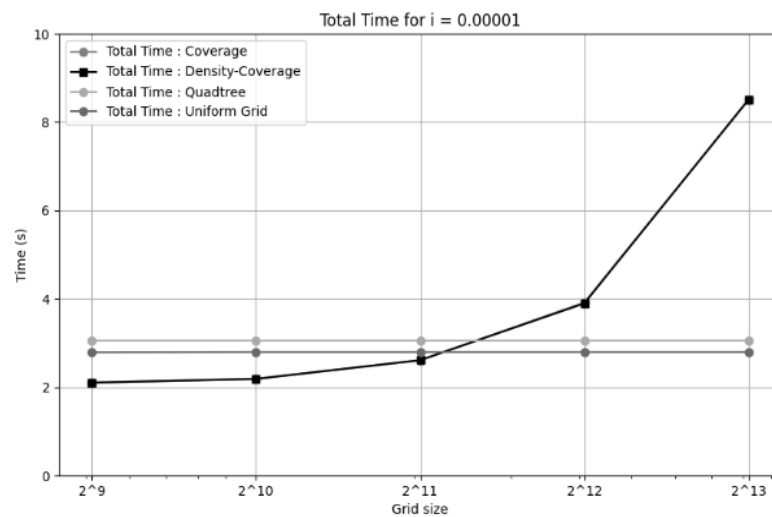


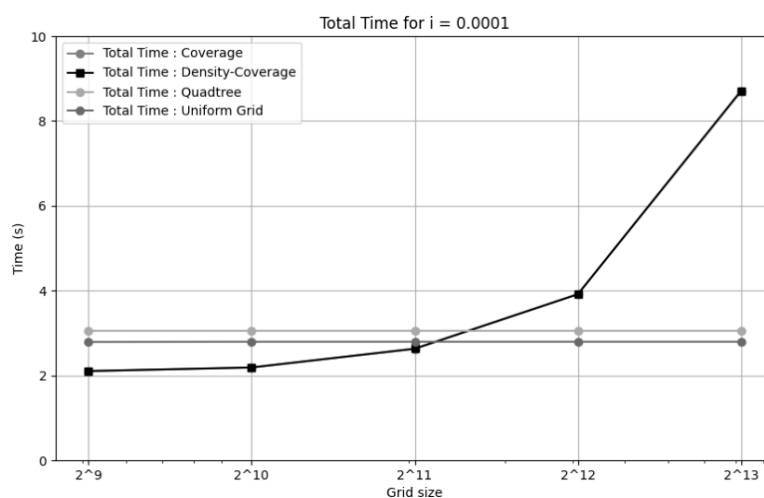Figure 14.1: The total execution time in seconds for i=0. 00001



Figure 14.2: The total execution time in seconds for i=0.0001

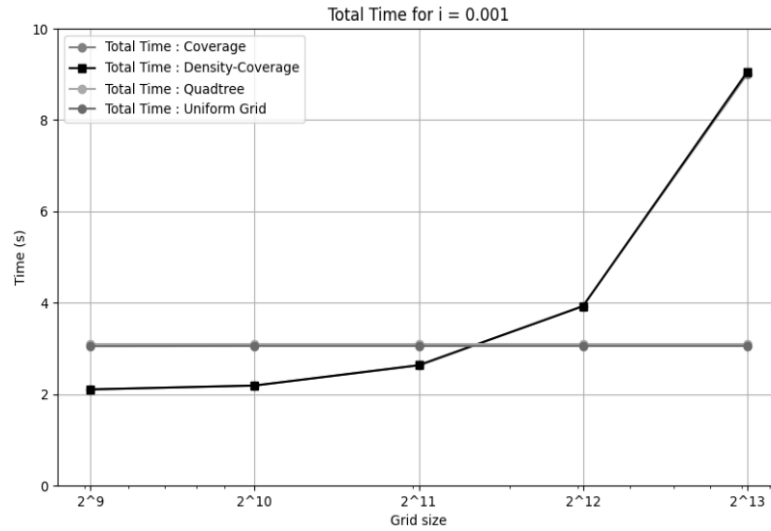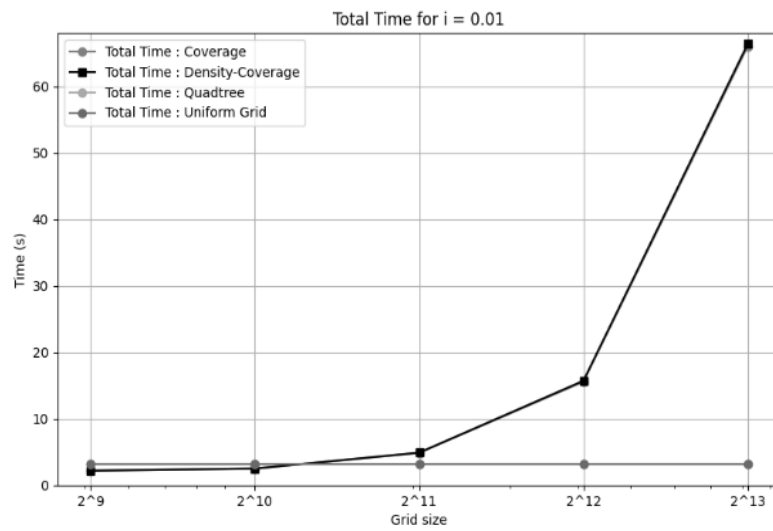Figure 14.3: The total execution time in seconds for i=0.001



Figure 14.4: The total execution time in seconds for i=0.01

It is obvious from the figures that HC-ARQ can achieve faster total times for smaller grid resolutions (in all cases from $2^9$ to $2^{10}$ and all except the case in which $i$=0.01 for $2^{11}$) than Uniform Grid and Quadtree while maintaining an average accuracy from 92.5% to 97.5%.

# Chapter 5.  Conclusion

HC-ARQ demonstrates its strengths better for smaller values of $i$, particularly in the range of i=0.00001 to i=0.001, where it maintains high accuracy across all grid sizes.  Additionally, for moderate grid sizes (e.g., $2^{10}$ and $2^{11}$), HC-ARQ tends to achieve both strong accuracy (92.5% to 97.5%) and superior total execution time, particularly when $i \leq 0.001$, outperforming both the Quadtree and Uniform Grid methods in terms of total execution time.

Taking everything into consideration, the developed tool demonstrates its strengths under specific conditions, making it a versatile solution for spatial queries under these certain conditions.

More specifically the tool usage can be summarized in the following scenarios based on the condition of each task:

- Precision-focused tasks at moderate grid sizes and smaller $i$: The tool is ideal for spatial queries that demand high accuracy, particularly when using moderately large grids (e.g. around $2^{10}$-$2^{11}$) and smaller values of $i$ (e.g., $i$=0.00001 to $i$=0.001). In these scenarios, the tool excels by maintaining high accuracy while keeping total execution time relatively low.

- Smaller grids for tasks that prioritize faster processing with reasonable accuracy: For situations where total execution time is a critical factor, smaller grid sizes (e.g., $2^9$ or $2^{10}$) allow the tool to deliver the best performance with accuracy above 92.5%, making it suitable for many practical applications.

# Chapter 6.      Future Work

Considering that spatial aggregation using space filling curves is a significantly under-researched topic, there are several possible improvements that should be considered.

**Attribute Homogeneity at Hilbert Orders Grouping** – One proposal is grouping Hilbert Orders in a way that clusters similar values based on a specified attribute (e.g. the most queried one). This means that the value of the points in this attribute within each group should not vary significantly across groups. This can be implemented by setting an attribute difference limit in the grouping of the Hilbert Orders.

For example, if we want to group the Hilbert Orders based on the average values of an attribute X then the grouping will be done based on the following check:

$$|A(i) - A(j)| \leq T$$

where A(*i*) represents the average value of the attribute X of the points within group *i* and A(*j*) represents the same value for the Hilbert Order that we check if it is qualified to add in the group *i*.

For any new Hilbert Order *j* to be added, the absolute difference in average values of attribute X between the current group (*i*) and the Hilbert Order (*j*) must not exceed a predefined threshold T.

The addition of this check will lead to results in groups that are spatially coherent due to the continuity of Hilbert orders and attribute-consistent, ensuring that no group contains data points with significantly disparate attribute values. This is particularly useful when it is necessary to avoid high variance in an attribute across different spatial regions.

However, by adding this check we may increase significantly the number of Intervals leading to a higher computational cost.

**Alternative Space Filling Curves** - We decided to use the Hilbert Space-Filling Curve because it is the most popular choice for such tasks, as it preserves spatial locality and provides a smooth mapping from 2D to 1D. However, there are other space-filling curve alternatives that may offer distinct advantages in certain applications.

Peano Space Filling Curve - One such alternative is the Peano Curve, which is known for its space-filling properties and self-similar structure [22] which is also a self-similar fractal curve that can be used for mapping higher-dimensional data to one-dimensional space. Unlike the Hilbert Curve, the Peano Curve is constructed using a set of simple rules, making it easier to implement and understand. Peano curves have the advantage of being easier to generate and compute than Hilbert curves, but they tend to have a lower degree of spatial locality, particularly in higher dimensions.

Lebesgue (Z-order) Space Filling Curve - The Lebesgue curve is a space-filling curve is particularly useful for mapping data with strong anisotropic spatial characteristics. The Lebesgue curve is relatively simple to generate and implement compared to curves like the Hilbert curve because its traversal pattern follows a straightforward, recursive structure. In lower dimensions (particularly 2D), the Lebesgue curve demonstrates good spatial locality preservation, meaning that points close together in the original space tend to remain close together in the one-dimensional representation.

However, the Lebesgue curve's traversal pattern can introduce anisotropy, meaning its locality preservation is direction dependent. Points along the diagonal of a 2D space, for instance, might be mapped farther apart than points along the horizontal or vertical axes meaning tha as the dimensionality of data increases, the Lebesgue curve's ability to maintain spatial locality can degrade compared to curves like the Hilbert curve.

In conclusion, while the Hilbert Space-Filling Curve is a widely used and effective method for mapping 2D data to 1D, there are several alternative space-filling curves that may offer unique advantages in specific applications. Researchers and practitioners should carefully evaluate the trade-offs and choose the curve that best fits the requirements of the particular problem.

**Alternative Weighting Techniques** - We selected the coverage and density-coverage weighting methods because they provide a balance weighting to the intervals with a simple and not computational costly algorithm.

However, there are other weighting techniques that should be taken into consideration.

Gaussian weighting, or kernel-based weighting, assigns weights to data points or intervals using a Gaussian (bell-shaped) function based on their distance from a central point (query center). Points near the query center receive the highest weight, with the weight decreasing smoothly as the distance increases. This method gives more influence on data points that are closer to the query center, which is important when spatial proximity is a key factor.

Moreover, Gaussian weighting is effective when you want to focus the analysis on a central point but still account for surrounding data points in a gradual way.

However, this method heavily weights data points near the query center and may underemphasize important patterns in more distant regions.

Furthermore, Gaussian weighting does not account for the density of points in an interval neglecting intervals with more points. The weight $W_i$ for each interval $i$ is based on its distance $d_i$ from the center of the query region, using a Gaussian function

$$W_i = e^{-\frac{d_i^2}{2\sigma^2}}$$

Where σ is the standard deviation that controls the spread of the weights.

One more weighting technique that could be used is Adaptive Weighting which adjusts the weights dynamically based on the characteristics of the data and the query region. Instead of using a fixed formula, this method adapts to local conditions, such as data density, variance, or other factors, to assign weights. Adaptive Weighting uses local data properties (such as variance, density, or spatial patterns) to dynamically adjust weights during analysis. This method allows for more flexibility and customization based on local conditions or data patterns. It can balance between different weighting approaches based on the characteristics of the data. However Adaptive Weighting requires sophisticated algorithms to adjust weights dynamically and may be computationally intensive.

# Chapter 7.        References

[1] W. Zhu, "A Study of Big-Data-Driven Data Visualization and Visual Communication Design Patterns". Scientific Programming, pp. 1-11, 2021.

[2] M. Oliveira and R. Menezes, "Spatial concentration and temporal regularities in crime", Understanding Crime through Science", Springer, 2019

[3] K. M. Divya, K. M. Reddy, G. S. Pujar and P. J. Rao, "Assessing the Spatial

Patterns of Geotagged MGNREGA Assets on Bhuvan Using GIS Based

Analysis". Springer Series in Geomechanics and Geoengineering pp. 227-241, 2018.

[4] T. Osaragi, "Classification and Space Cluster for Visualizing GeoInformatio" International Journal of Data Warehousing and Mining, 2019 15(1), 19-38.

[5] J. Sheikh, I. Shafique, M. Sharif, S. A. Zahra and T. Farid, "IST: Role of GIS in crime mapping and analysis". *International Conference on Communication Technologies (ComTech),* IEEE 2017, 126-131.

[6] A. Newton and M. Felson, "Editorial: crime patterns in time and space: the dynamics of crime opportunities in urban areas". Crime Science, 2015.

[7]  ArcGIS Insights. [Online]. Available:
https://doc.arcgis.com/en/insights/latest/analyze/spatial-aggregation.htm

[8] K. Clarke, P. H. Dana and J. T. Hastings, "A New World Geographic Reference System". Cartography and Geographic Information Science, pp. 355-362, 2002.

[9] C. Zhou, F. Su, F. Harvey and J. Xu, "Spatial Data Handling in Big Data Era". Advances in Geographic Information of Science, Springer Singapore, 2017.

[10] [Online] Available:
https://opendsaserver.cs.vt.edu/OpenDSA/Books/CS3/html/PRquadtree.html

[11] [Online] Available: https://www.researchgate.net/figure/First-three-stages-in-the-generation-of-Hilberts-space-filling-curve_fig1_265074953

[12] F. Heimerl, C. Chang, A. Sarıkaya and M. Gleicher, "Visual Designs for Binned Aggregation of Multi-Class Scatterplots". IEEE Transactions on Visualization and Computer Graphics, pp. 402-412, 2018.

[13] M. Calì and R. Ambu, "Advanced 3D Photogrammetric Surface Reconstruction of Extensive Objects by UAV Camera Image Acquisition". Sensors p. 2815, 2018.

[14] A. M. Ali, R. Shivapurkar and D. Söffker, "Optimal Situation-Based Power Management and Application to State Predictive Models for Multi-Source Electric Vehicles". IEEE Transactions on Vehicular Technology, pp. 11473-11482, 2019.

[15] A. Eldawy and M. F. Mokbel, "The Era of Big Spatial Data: A Survey". IEEE 32nd International Conference on Data Engineering (ICDE), pp. 1424-1427 2016.

[16] P. Delicado, R. Giraldo, C. Comas and J. Mateu, "Statistics for spatial functional data: some recent contributions". Environmetrics, pp. 224-239, 2009.

[17] S. Ying, L. Li and R. Guo, "Building 3D cadastral system based on 2D survey plans with SketchUp". Geo-spatial Information Science, pp. 129-136 2011.

[18] C. Böhm, G. Klump and H. Kriegel, "XZ-Ordering: A Space-Filling Curve for Objects with Spatial Extension". Advances in Spatial Databases: 6th International Symposium, SSD'99, pp. 75-90 1999.

[19] H. Shin, B. Moon and S. Lee, "Adaptive multi-stage distance join processing". ACM SIGMOD Record, pp. 343-354, 2000.

[20] P. Baumann, D. Mišev, V. Merticariu and B. P. Huu, "Array databases: concepts, standards, implementations". Journal of Big Data, pp. 1-8, 2021.

[21] T.Georgiadis and N.Mamoulis "Raster Intervals: An Approximation Technique for Polygon Intersection Joins". Proceedings of the ACM on Management of Data, pp. 1-18, 2023.

[22] M. Altın, İ. Ünal and F. Mofarreh, "Interpolation of surfaces with asymptotic curves in Euclidean 3-space". Applied Mathematics & Information Sciences, pp. 533-538, 2021.

[23] Kaggle Homicide Dataset. [Online].  Available:

https://www.kaggle.com/datasets/joebeachcapital/homicides

[24] M. M. Bertoldi, M. A. Yardımcı, C. M. Pistor, and S. I. Güçeri, "Domain Decomposition and Space Filling Curves in Toolpath Planning and Generation," in Proceedings of the 1998 Solid Freeform Fabrication Symposium, pp. 267–274, 1998.

[25] S. Alexander, R. L. Bishop, and R. Ghrist, "Total Curvature and Simple Pursuit on Domains of Curvature Bounded Above," Geometriae Dedicata, pp 275-290, 2010.

[26] M. Sarfraz, *Some Algorithms for Curve Design and Automatic Outline Capturing of Images*, International Journal of Image and Graphics. p.p 301-324