

Towards Adaptive Trajectory Data Management: Modelling, Accessing, Distributing, and Query Optimization in Distributed Database

SICONG, GONG (5711932)

1ST SUPERVISOR: DR.IR. MARTIJN MEIJERS

2ND SUPERVISOR: DRS. WILKO QUAK

CO-READER: DR. KEN ARROYO OHORI

JUNE 17, 2024



Table of Contents

1 Introduction

2 Q1: Modelling

3 Q2: Accessing

4 Q3: Distributing

5 Conclusion

Msc Thesis In Geomatics

Towards Adaptive Trajectory Data Management: Modeling, Accessing, Distributing, and Query Optimization in Distributed Database

Sicong Gong
2024

Introduction

- 1 Introduction**
 - What are the trajectories?
 - What are the difficulties?
 - What are the DDBMSs?

2 Q1: Modelling

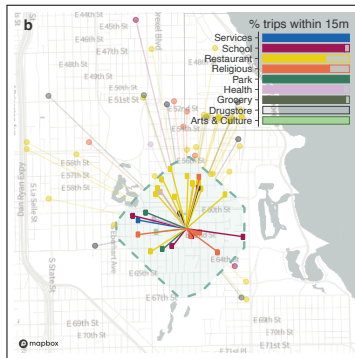
3 Q2: Accessing

4 Q3: Distributing

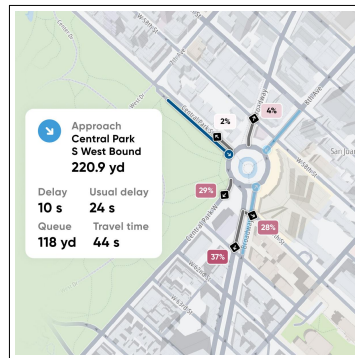
5 Conclusion

Trajectory Applications

- Urban planning → 15-minute city quantification.
- Traffic management → Traffic condition identification.



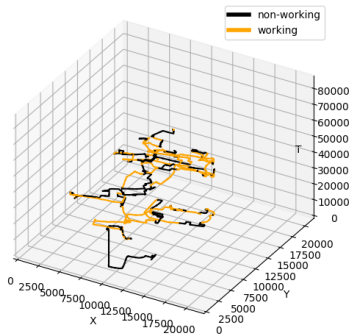
(a) 15-minute City Quantification, adopted from [Abbasov et al. \(2024\)](#)



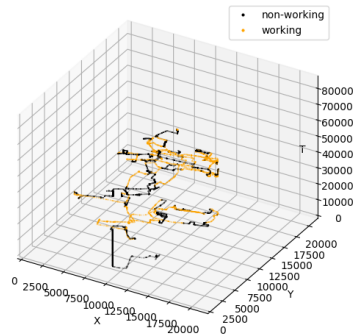
(b) Traffic Condition Identification, adopted from [TomTom](#)

Trajectory Definition

- Humans' understanding → Continuous sequences.
- Data records acquired → Discrete points.



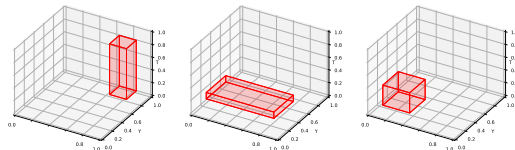
(a) Continuous Sequences



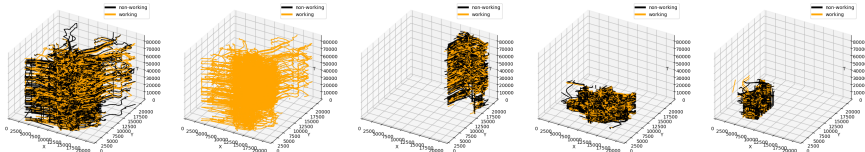
(b) Discrete Points

Trajectory Operations

■ Selection by ID + Selection by range.



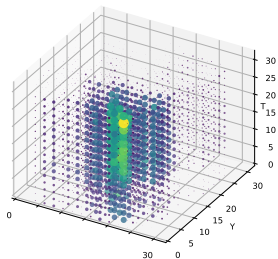
(a) t-4-c (Needle Used) (b) t-4-d (Used) (c) t-4-e (Dice Used)



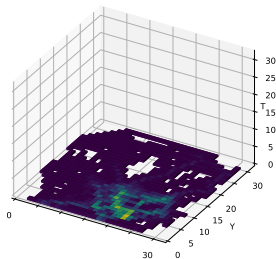
(d) t-4-a (Selection) (e) t-4-b (Selection) (f) t-4-c (Needle Selection) (g) t-4-d (Slice Selection) (h) t-4-e (Dice Selection)

Trajectory Operations (Continued)

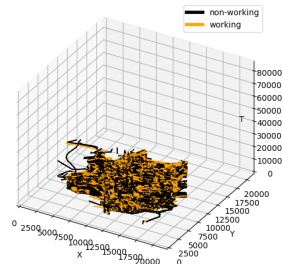
■ Aggregation + Projection + Simplification.



(a) t-6-a (Aggregation)



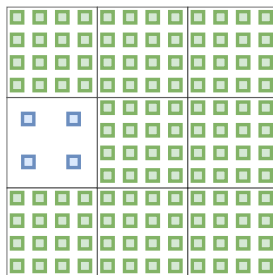
(b) t-6-b (Projection)



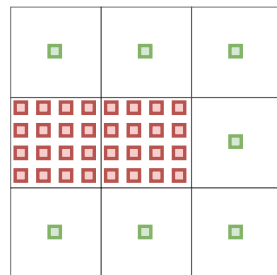
(c) t-6-c (Simplification)

Trajectory Properties

- **High frequency: Huge volume.**
- **High cardinality: Numerous entries.**
- **High dimensionality: Integration of space, time and semantics.**
- **High heterogeneity: Uneven distribution.**



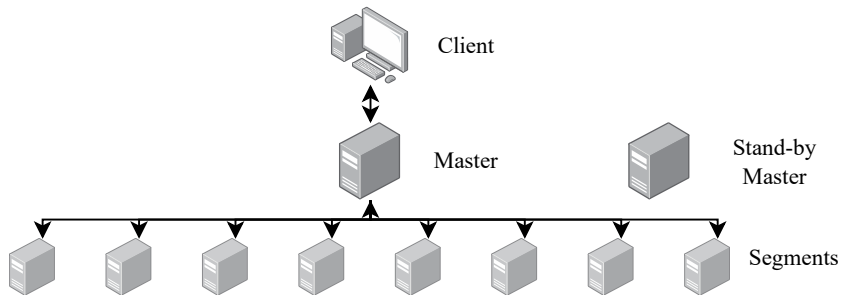
(a) Laser Scanning of Lake and Grass



(b) Traffic Density of City and Rural Area

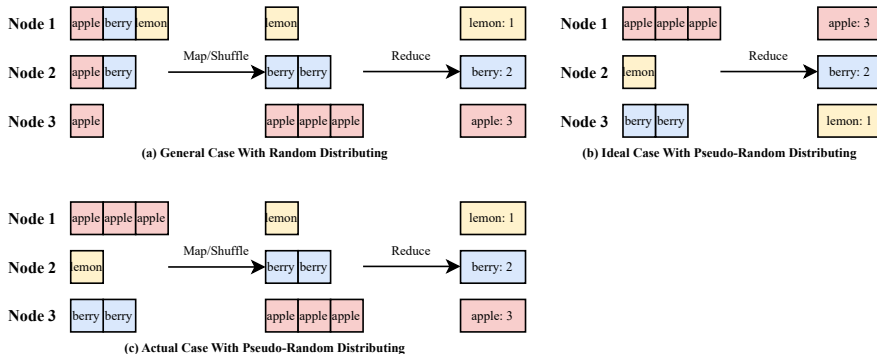
Distributed DBMS

- **Scalability (Speed-up and Scale-up) and Localization (Cluster data and localize computations)**



- **Experiment setting: 5 virtual machines (1 master, 4 nodes), each node has 2 segments.**

Map-Reduce Computation Model



- MapReduce partitions data into key-value pairs, processes them in parallel mappings and then aggregates the results through local computation in the reduce step.

Distributed Products



Feature	Traditional Database	Hadoop	Spark	MPP Database
Volume	GB-TB	PB-EB	TB-PB	TB-PB
Robustness	High	High	High	Medium
Scalability	Low	High	High	High
Latency	Medium	High	Low	Very Low
Throughput	Low	High	High	Medium
Data Type	Structured	All	All	Structured

Q1: Modelling

1 Introduction

2 Q1: Modelling

- How to model trajectory?

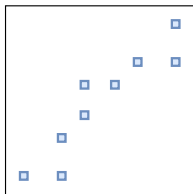
- Does the Above Methods Work? - Compression

3 Q2: Accessing

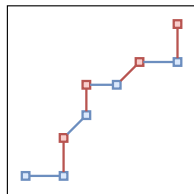
4 Q3: Distributing

5 Conclusion

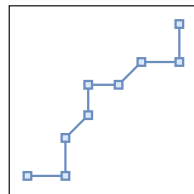
Modelling Options



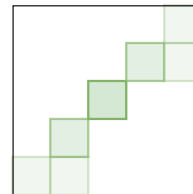
(a) Individual Point(s)



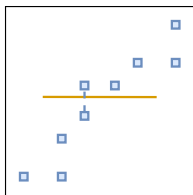
(b) Isolate Segement(s)



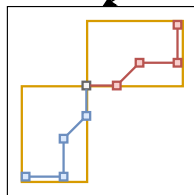
(c) Continuous Sequence(s)



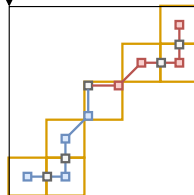
(d) Discrete Grid



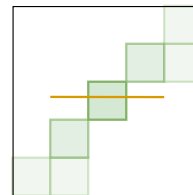
(e) No Explicit Geometry



(f) Irregular Splitting



(g) Regular Splitting

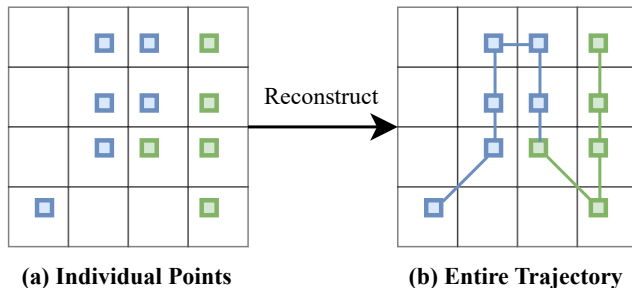


(h) No Geometry/Identifier

Reasons for Sequence-based Model

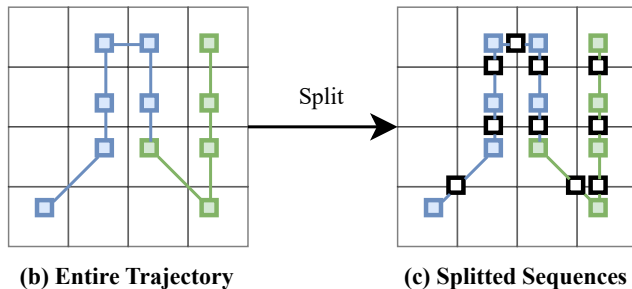
The sequence is better!

- More supported operations.
- Smaller entries cardinality.
- Higher compression potential.



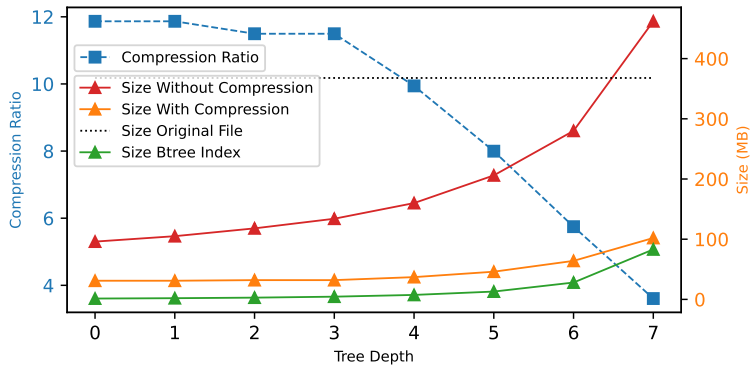
Solutions for Unstructured Nature

- Split by semantics.
- Split by spatio-temporal cube.



Does the Above Methods Work? - Compression

- Further subdivision of the space would decrease the compression ratio.



Q2: Accessing

1 Introduction

2 Q1: Modelling

3 Q2: Accessing

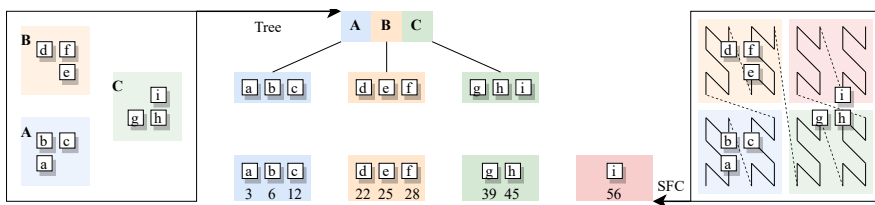
- How to access trajectory?
- Does the Above Methods Work? - Selection

4 Q3: Distributing

5 Conclusion

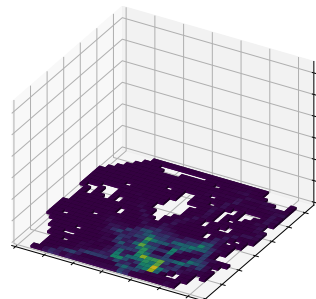
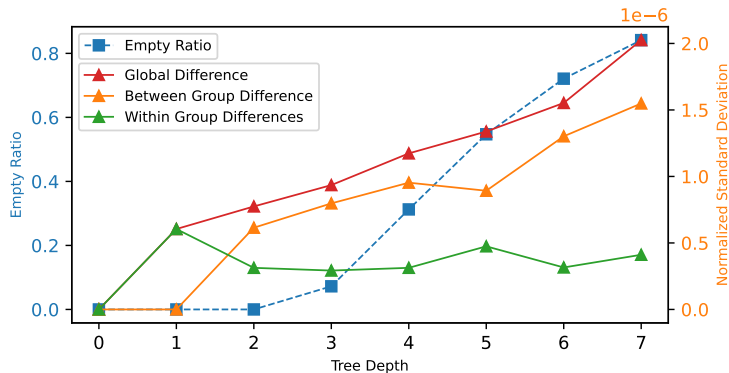
Accessing Methods

- R-tree → Adaptive but complex and costs more storage.
- Space filling curve → Rigid but simple and corresponds to the nature of modelling.



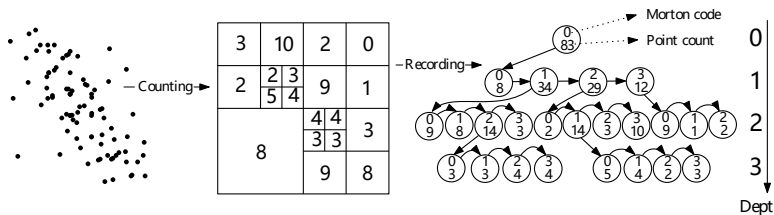
Uneven Distribution

- Further subdivision of the space would increase the empty ratio.
- Further subdivision of the space would increase the global difference (unevenness).



HistSFC Solution

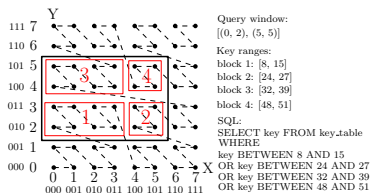
- Histogram Tree: Adaptive octree.
- Space-filling-curve (Morton): Represent a record and use b-tree for indexing.



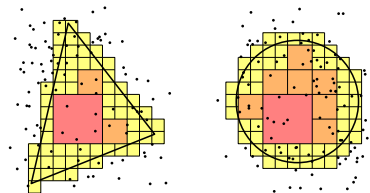
Adopted from Liu, 2022

Shape Querying

- Recursively partitioning the extent of data according to SFC regions to match different query geometries, for selecting data in the table.



(a) Executing a window query on a uniformly distributed 2D point set based on Morton encoding

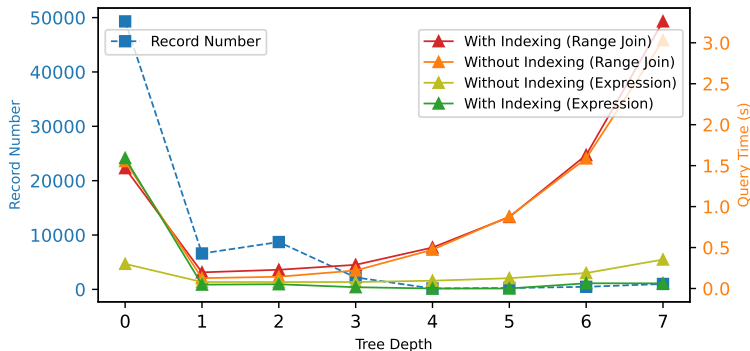


(b) Querying with a triangle and a circle: false positive points in boundary cells will be filtered out by a second filter

Adopted from Liu, 2022

Does the Above Methods Work? - Selection

- Selection performance is first increasing then decreasing with the further subdivision of the space.



Q3: Distributing

1 Introduction

2 Q1: Modelling

3 Q2: Accessing

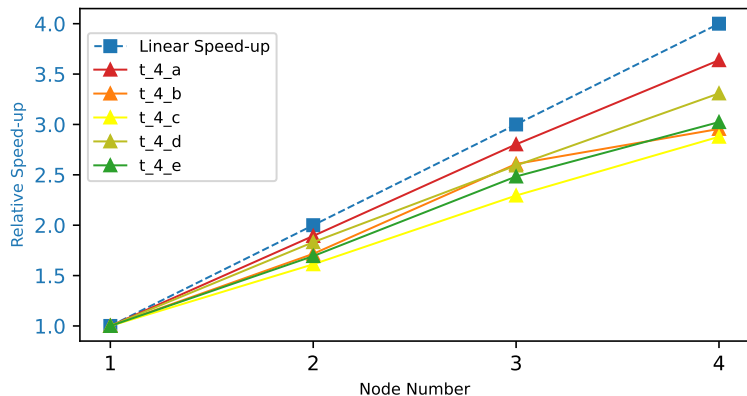
4 Q3: Distributing

- Does the Distributed DBMS Work? - Speed-up
- Does the Distributed DBMS Work? - Scale-up
- How to distribute trajectory?
- Does the Above Methods Work? - Localization

5 Conclusion

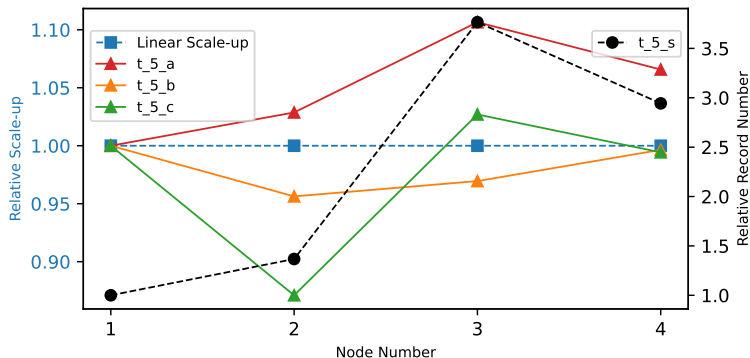
Does the Distributed DBMS Work? - Speed-up

- Five operations are designed to test the speed-up (same problem sizes with increasing resources) and the result is positive.

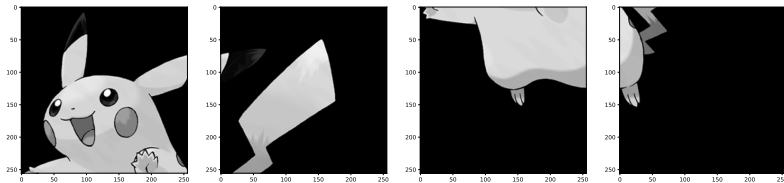


Does the Distributed DBMS Work? - Scale-up

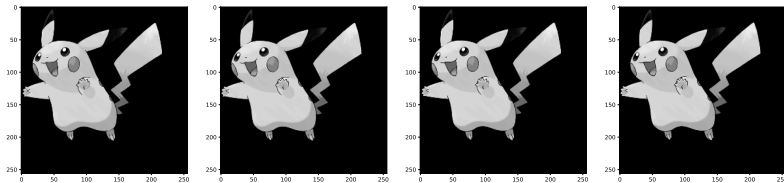
- Three operations are designed to test the scale-up (increasing problem sizes with increasing resources) and the result is positive.



Distributing Strategy



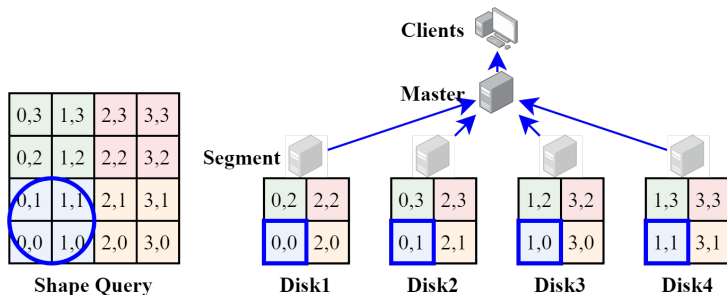
Block-based: Fold twice and split.



Sample-based: e.g. all the lower left pixels of 4 neighbours as a group.

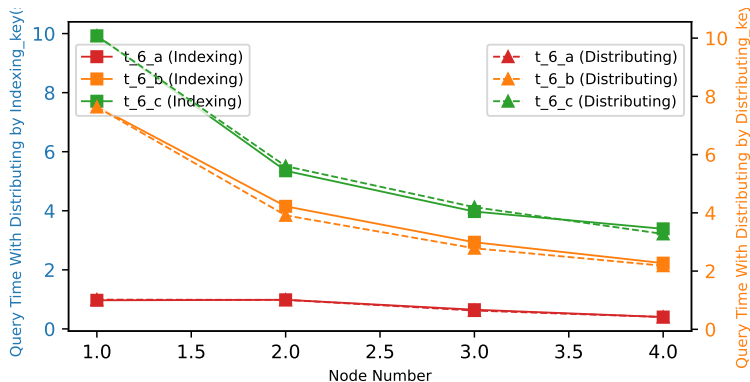
Pseudo-Sampling Distributing

- **Load-balancing:** Block-based method leads to uneven data distribution.
- **Localization:** Random sampling leads to no locality being preserved.



Does the Above Methods Work? - Localization

- Three operations (Aggregation, Projection and Simplification) are designed to test the localization but the result is negative.



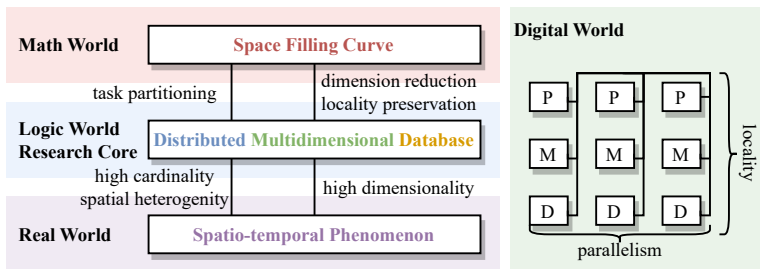
Conclusion

- 1 Introduction
- 2 Q1: Modelling
- 3 Q2: Accessing
- 4 Q3: Distributing
- 5 Conclusion**

Conclusion

In conclusion, the above methods alleviate the difficulties mentioned.

- **Modelling:** Reduce cardinality.
- **Accessing:** Reduce dimensionality and alleviate uneven distribution.
- **Distributing:** Use parallelism to speed up and scale up.



Reflection

- The main lesson learnt from this thesis is the need to adapt the data properties and platform features.

$$\text{Distribution Awareness} \left(\begin{array}{l} \text{adaptive modelling (splitting)} \\ \text{adaptive accessing (indexing)} \\ \text{adaptive distributing (partitioning)} \\ \text{adaptive querying (merging)} \end{array} \begin{array}{l} + \\ + \\ + \\ + \end{array} \right) \times \text{Distributed Architecture} \quad (1)$$

Future Work

However, there are still some limitations that should be done in future work.

- **Realistic benchmarking:** Not only in virtual machines.
- **Workflow optimization:** Adaptive splitting and range merging.
- **Mathematical proofing:** Not only by experiments.

THANK YOU!

Reference

Liu, H. (2022). nD-PointCloud Data Management: continuous levels, adaptive histograms, and diverse query geometries. A+ BE| Architecture and the Built Environment, (12), 1-206.