



INFS4205/7205 Advanced Techniques for High Dimensional Data
Managing Spatiotemporal Data

Semester 1, 2021

University of Queensland

+ Advanced Techniques for High Dimensional Data

- ❑ Course Introduction
- ❑ Introduction to Spatial Databases
- ❑ Spatial Data Organization
- ❑ Spatial Query Processing
- ❑ Managing Spatiotemporal Data
- ❑ Managing High-dimensional Data
- ❑ Other High-dimensional Data Applications
- ❑ When Spatial Temporal Data Meets AI
- ❑ Route Planning
- ❑ Trends and Course Review

+ Learning Objectives

■ What we will cover

- Spatiotemporal data and queries
- Spatiotemporal indexing and query processing
- Trajectory similarity measures
- Open issues and directions

■ Goals

- Understand the temporal dimension of spatial data
- Understand basics of spatiotemporal data management
- Learn from examples about how to deal with new data management challenges brought by new data types

+ The Temporal Dimension

- Location and time are two ubiquitous attributes of data
- RDBMS is designed to handle neither of them
- We have studied how the spatial dimension can be managed so far
- The temporal dimension of spatial data cannot be simply viewed as just another dimension
 - Data values and operations/queries are quite different

+ Spatiotemporal Data

- GPS recordings
- Sensor data, RFID data and surveillance data
- Use of smartphones and smartcards
- Use of location-aware apps
- Social media data: uploaded photos/messages and check-in data
- Much more spatiotemporal data to come with 5G
- ...and beyond the geographical domain

+ What is Trajectory Data

- Any data that record the locations of a moving object over time in a geographical space
- Simple form: $\langle id, (p_1, t_1), (p_2, t_2) \dots (p_n, t_n) \rangle$

Ordered by time: $t_1 < t_2 < \dots < t_n$

- General form:

$\langle objId, trajID, trajProperties, \\ (p_1, t_1, a_1), (p_2, t_2, a_2) \dots (p_n, t_n, a_n) \rangle$

+ Many Dimensions of Trajectory Data

■ Basic Dimensions

- Spatial dimension: locations $p_1, p_2 \dots p_n$
- Temporal dimension: time-stamps $t_1, t_2 \dots t_n$
- Attribute dimension: other data of interest $a_1, a_2 \dots a_n$

■ Other Dimensions

- Entity dimension: what type of objects?
- Environment dimension
 - Road networks, Floor plans, Water systems, Sensor networks
- Semantic dimension: what activities at a location or time?

+ How Much Trajectory Data?

- A back-of-the-envelope calculation:
 - A simple point data (x, y, t) : 24 bytes
 - A car can generate 85KB a day (10 hours a day, 10 seconds interval)
 - Beijing has 60,000 taxis, that is 5GB a day, or 1.72 TB a year
- Actual trajectory data could be much larger
 - A multiplier of X: There are much more information than just a point data: taxi ID, trip ID, job status, direction, velocity, acceleration, fuel consumption, other sensor data (OBD/M2M data)
 - Even larger once processed: original data, plus map-matched data, other derived data, other forms of representation (e.g., OpenLR - <http://www.openlr.org/>)

+ Trajectory Data in a Company

- A car navigation service provider
- Total trajectory data: 32 TB in size, 10.9 billion matched trajectories

	Current	Daily
Company X (in-car navigation provider)	17.6TB	15M trajectories
Company Y (map app provider)	14.5TB	5M trajectories
Company Z (social network)	0.68TB	18M trajectories

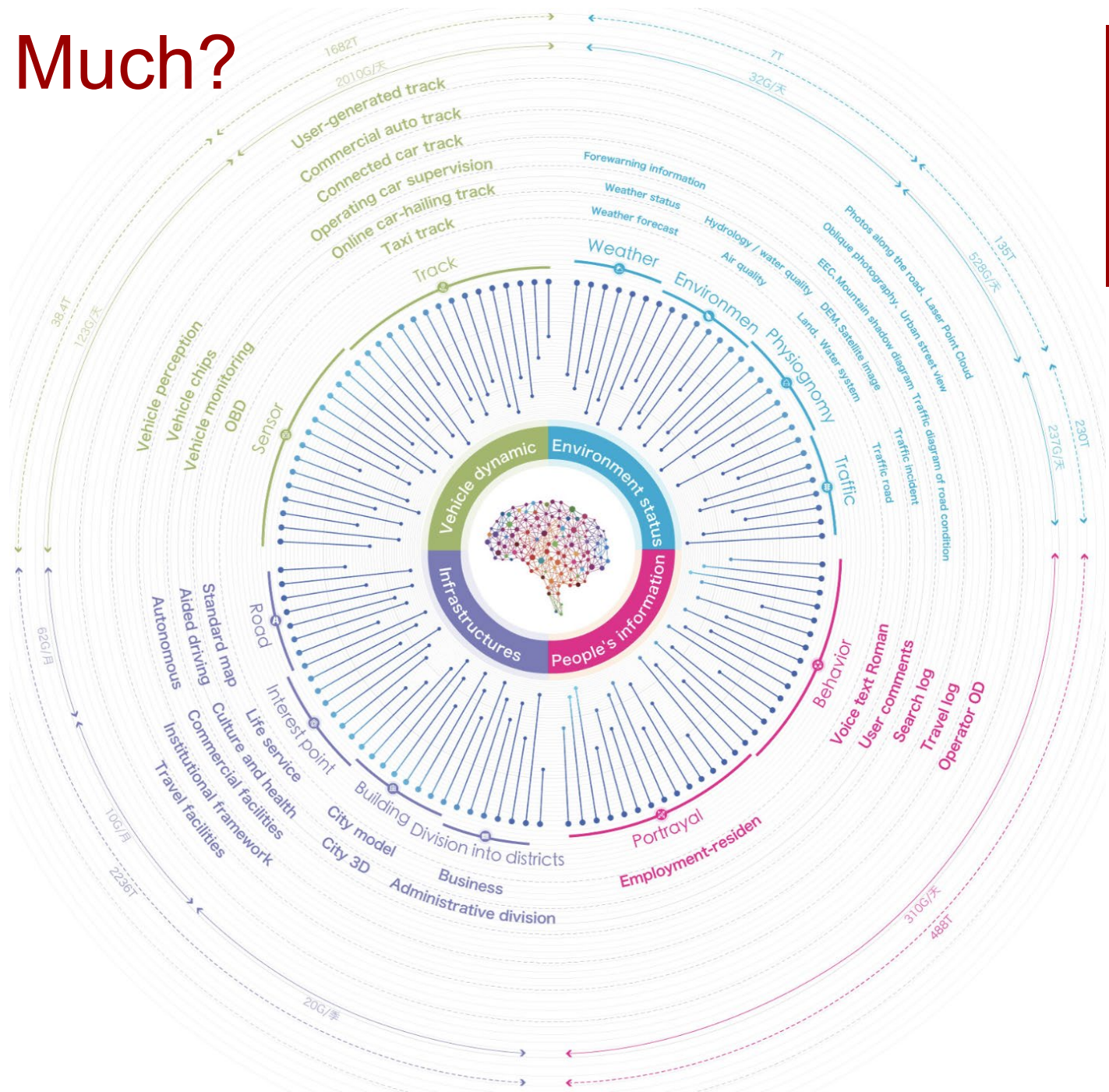
- Every day, ~40M new trajectories, ~4 billion points
- Many types of related data: maps, accident reports, roadside data, surveillance video, weather, events, social media...

+ NavInfo DataHIVE (minedata.cn)

10

Vehicle	Infrastructure	Environment	People
Trajectories:	Standard maps	Weather	Voice and text
- taxis	High res maps	Events	User comments
- uber-like	Services POIs	Air quality	Search log
- monitored	Culture POIs	Water quality	Travel log
- commercial	Commercial POIs	Land & water info	Operators' OD
- user generated	Health POIs	DEM & EEC	Workplace info
Sensor/OBD data	Travel POIs	Satellite image	
Perception data	City models	Street views	
	City 3D Models	Roadside pictures	
	Business districts	Laser point cloud	
	Admin boundaries	Road condition	
	Organization maps	Traffic condition	
		Traffic incidents	

+ How Much?



+ Applications

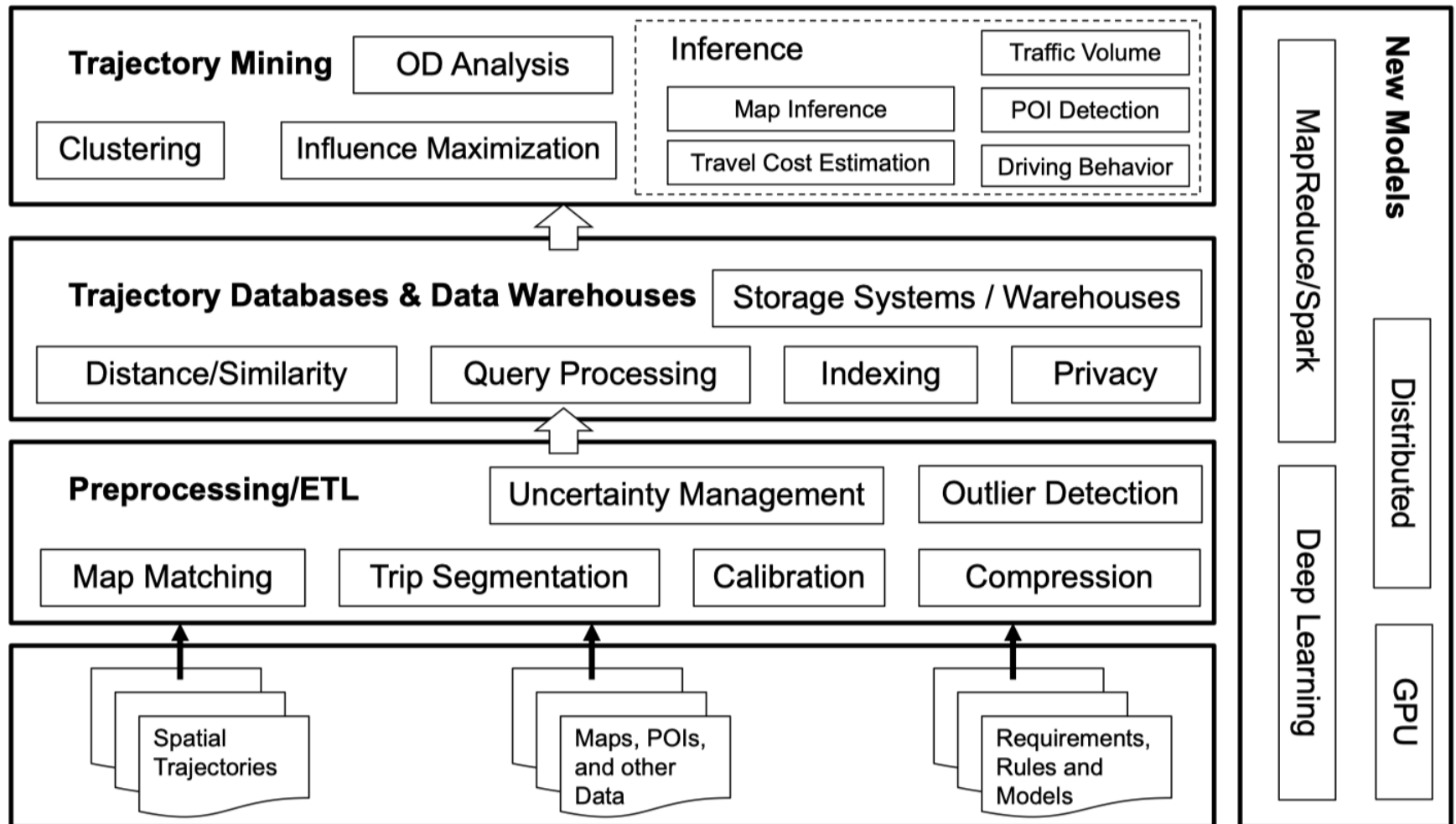
- Understanding, monitoring, predicting mobility patterns
 - ...from very large amount of movement data in real-time
- Some examples
 - Finding the nearest businesses or services
 - Location/movement-based event/service recommendations, alerts (sales, traffic jam, warnings...) and information push (e.g., advertisement)
 - Resource tracking and scheduling (e.g., for fleet management)
 - Safe drivers (e.g., for insurance industry)
 - Data-driven ITS, urban planning and smart cities
 - ...

+ Alternative Names

- **Moving objects data** concerns the current and future locations
- **Spatial trajectory data** is the movement history of moving objects
- A trajectory without the time dimension is also called a **route**

+ Spatiotemporal Data Management

14



+ Spatiotemporal Queries

- Simple point/range queries are useful
 - Spatiotemporal point/window queries: a combination of timestamp/time interval and spatial point/region
 - Examples:
 - To find where an object is at time t
 - To find all the objects at/near location p at time t
 - To find all objects inside a given region during a given period...
- Query can also take a trajectory as an input
 - To find the nearest POI for a given trajectory
 - To find top-K most similar trajectories to a given trajectory

+ More Spatiotemporal Queries

- More challenging queries
 - Monitoring queries: also called continuous queries
 - Predicative queries: at a future point of time
- Data analytics queries
 - To find trajectory clusters
 - To find where traffic jams could occur in the next 30 minutes
 - To find where people come from to a given region

...can you give some examples for each type of queries?

+ Example: Continuous NN Query

17

- Where is the nearest petrol station?
- Static NN query
 - Concerning a given location
 - Optimization goal: minimize the number of objects to be examined
- Moving NN query
 - Concerning the current location which is moving
 - Optimization goal: minimize the number of calls to the NN algorithm
 - Safe region: the results don't change within the region

how?

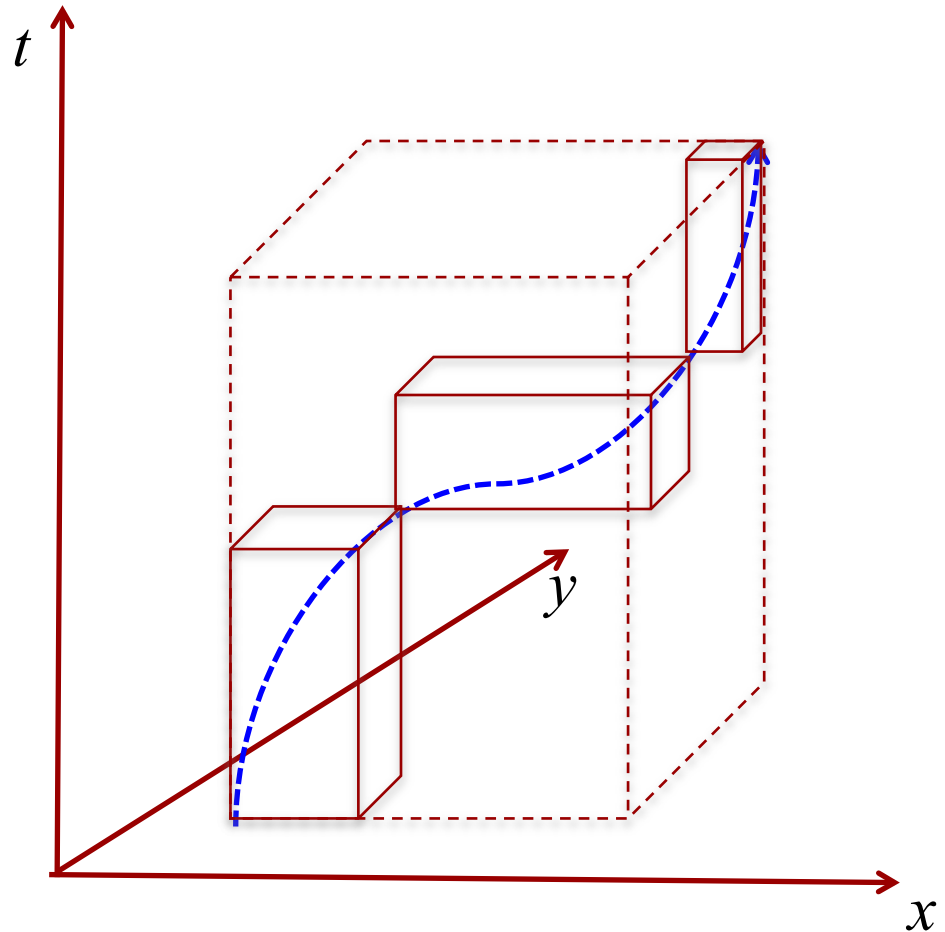
...path NN query: where is the nearest petrol station on my way from city to Surfer's Paradise?

+ Indexing Spatiotemporal Data

- R-tree has been the most efficient and widely used general purpose spatial indexing structure, so let's extend that for spatiotemporal data
- What's special now?
 - The time dimension increases monotonically and unbounded
 - An object can have a long sequence of locations (i.e., points), and those locations which are temporally close to each other are often spatially close too
 - MBR has been used as the foundation for R-trees, but now the MBR for an object or a group of objects either changes over time or occupies a huge space

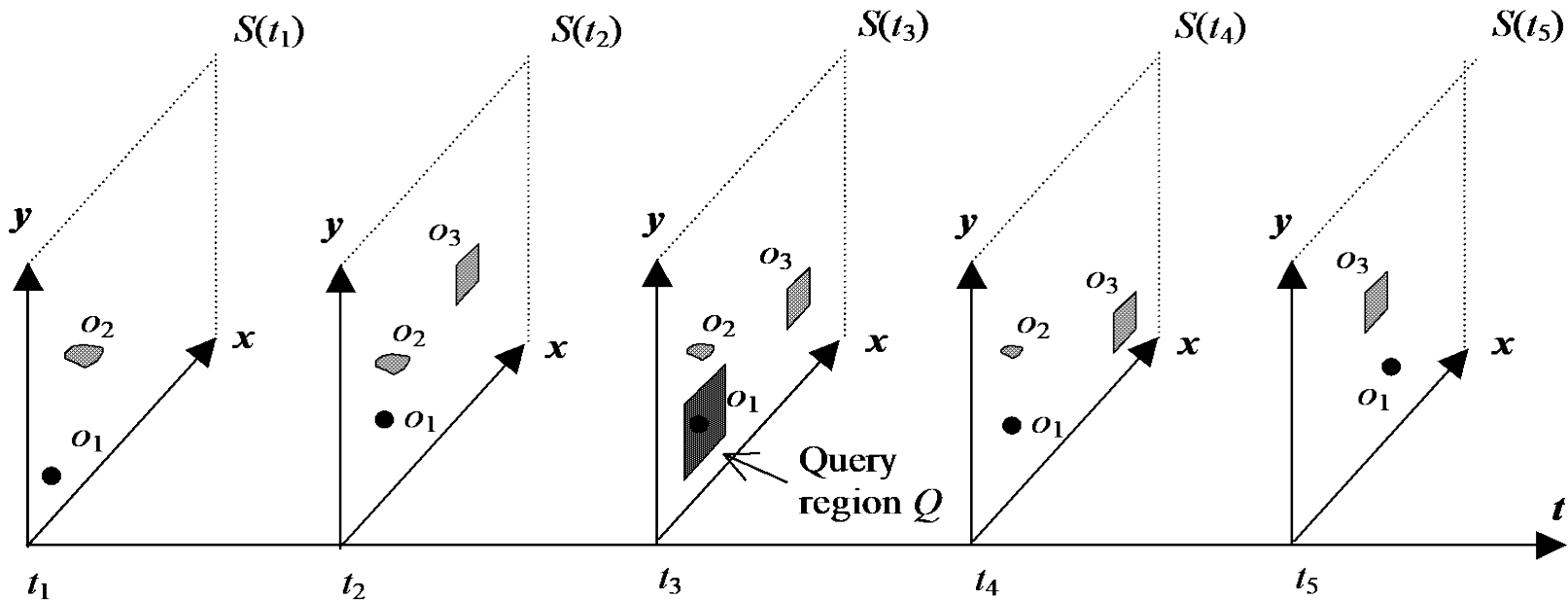
+ 3D R-tree?

- Adding time as another dimension
 - Conceptually simple: 2D for the spatial dimension + 1D for the temporal dimension
 - Problem 1: The t -dimension is unbounded
 - Problem 2: It's not effective to use boxes to approximate lines
 - Problem 3: Only efficient for coordinate-based queries (time slices and ranges), not for trajectory-based queries



+ Snapshot-Based Indexing

- One R-tree for each snapshot?

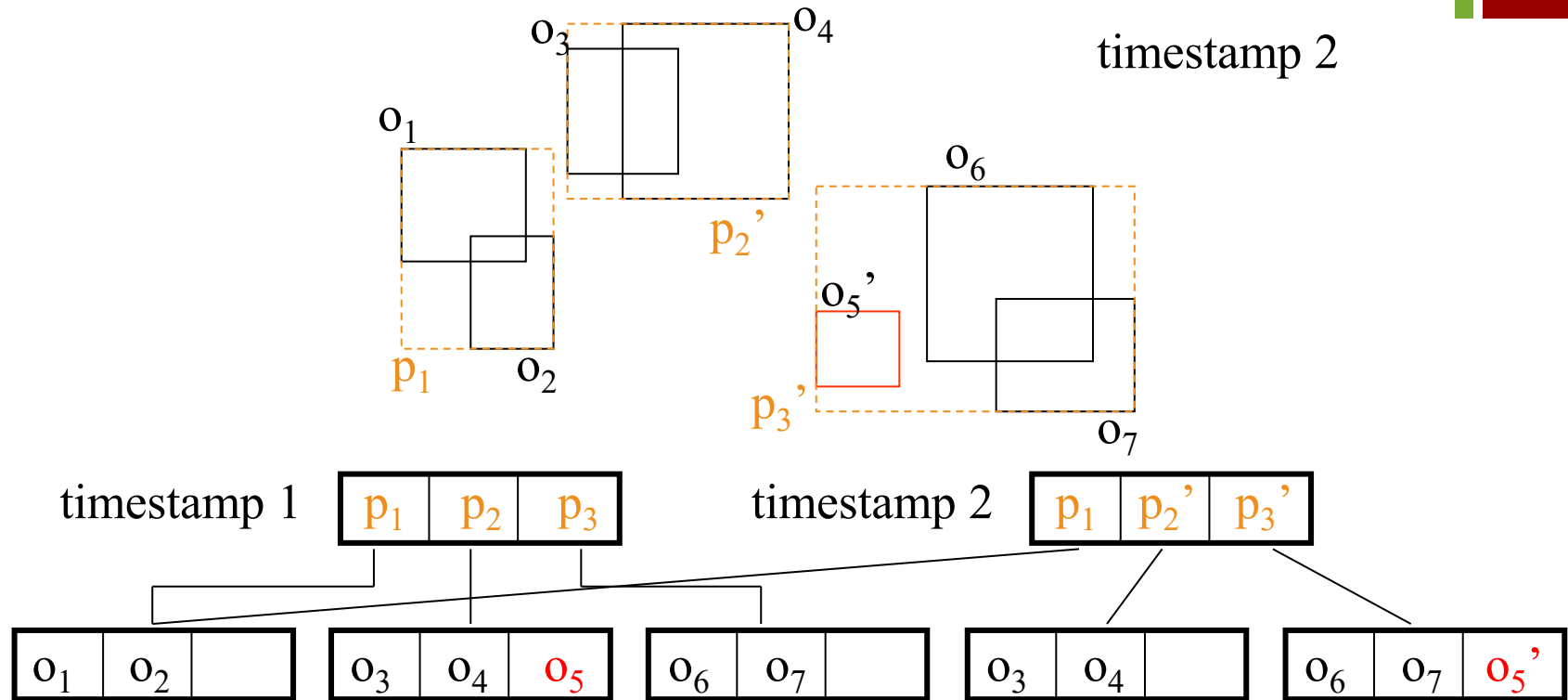


- An R-tree is maintained for each timestamp in history
- Trees at consecutive timestamps may share branches to save space



+ HR-Tree: Sharing Branches

22



...what do you think about this approach?

-

“Indexing the Position of Continuously Moving Objects”, S Saltenis, C Jensen, S T Leutenegger and M. A. Lopez, SIGMOD 2000

+ Trajectory Data

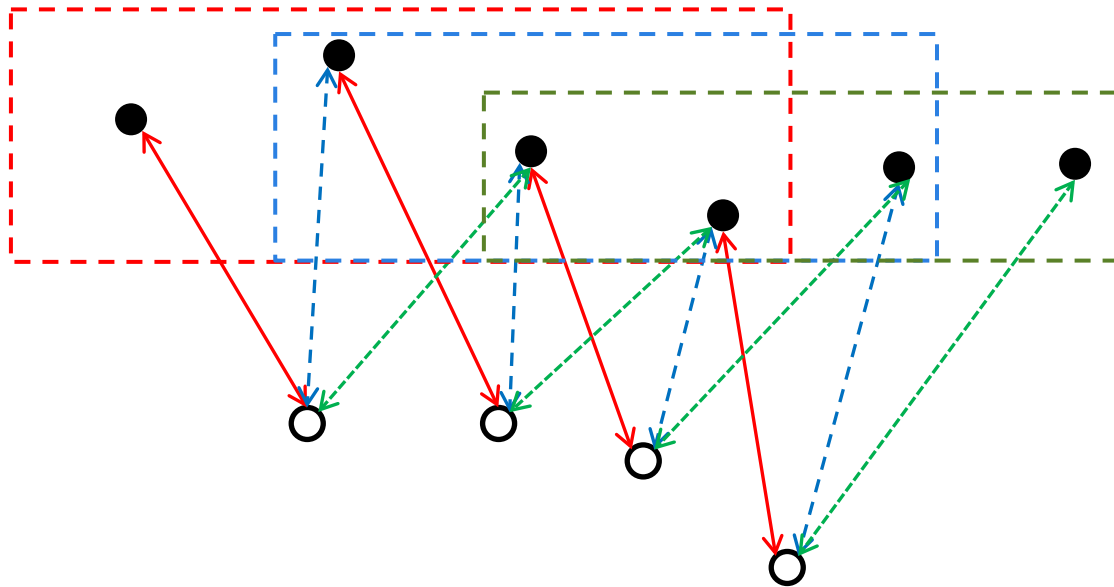
- Spatial trajectory is object movement history in a space
 - Continuous in nature, but discrete once captured and stored
- Many **location-update strategies**
 - By time, by distance, by deviation...
 - A trade-off between accuracy and other overheads
 - Variations may not always be under control
- Movement can be in a free space (e.g., birds flying), or a constrained space (e.g., cars in road networks)

+ Trajectory Similarity Measures

- The foundation to perform trajectory based queries and analytics (such as clustering)
- Many types of similarities
 - Sequence-based: passing the same sequence of points?
 - Geometry-based: similar shapes?
 - With or without time or speed considerations
- Key factors to consider
 - Alignment of sampling points (to deal with non-uniform sampling)
 - Robustness to noise (to deal with data quality issues)
- There are many trajectory similarity measures

+ Lock-Step Alignment

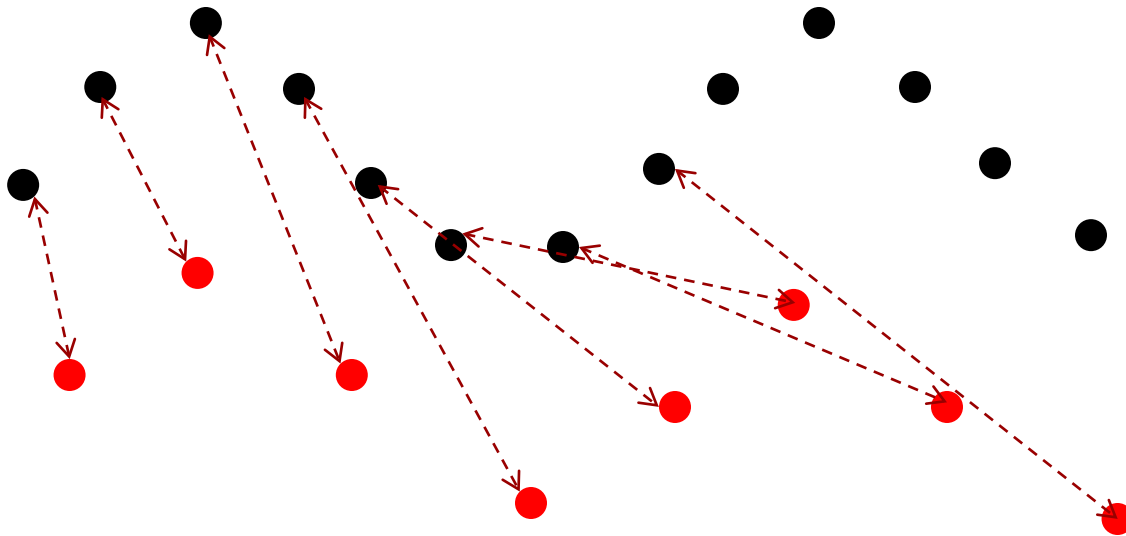
- The k^{th} point of a trajectory is aligned to the k^{th} point of the other trajectory
 - Sliding window based on the shorter trajectory
 - The distance is the best of window-based total Euclidian distance



... what is the time complexity of this approach?

+ Drawback

- Cannot find similar trajectories with different **sampling rates**, which are common in practice
- Sensitive to noise



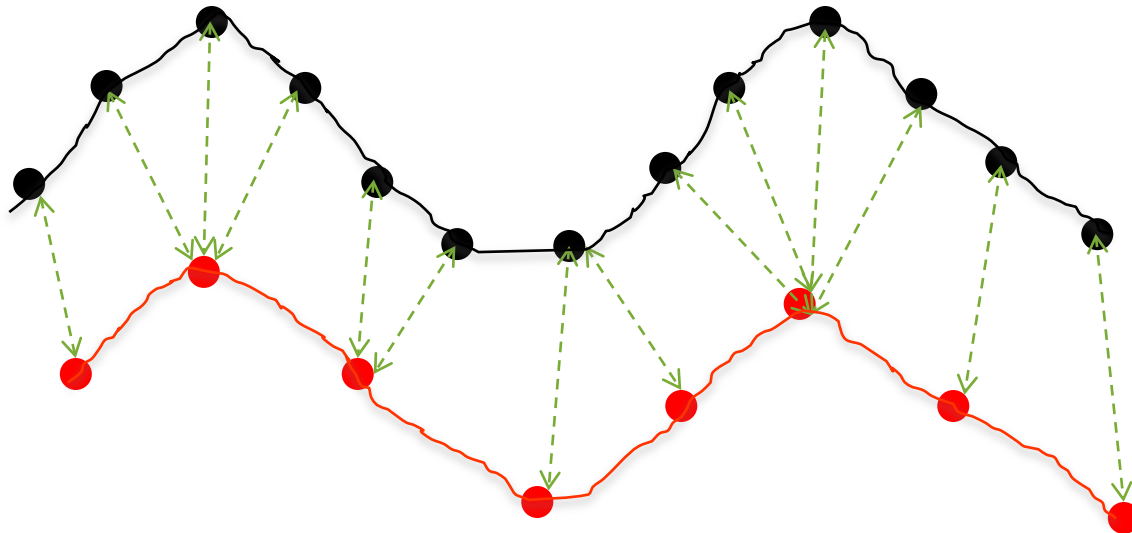
+ Adaptive Alignment

- Dynamic Time Warping (DTW) distance
 - Adaptation from time series distance measure
 - Used to handle time shift and scale in time series
- Optimal order-aware alignment between two sequences
 - Goal: minimize the aggregate distance between matched points
- 1-to-many mapping: one point in one sequence can be mapped to multiple points in another sequence

+ DTW for Trajectories

30

- Useful when detecting similar trajectories with different sampling rates



*Using Dynamic programming to compute DTW. Check the Wikipedia page
Time complexity: $O(MN)$*

+ Count-Based Similarity

- So far, geographical distance based
 - Similarity is measured by the geographical distance between matched samples
- Count based
 - Similarity is measured by the number of ‘similar’/ ‘dissimilar’ samples
 - Based on **Edit Distance**
 - The distance between two strings is the minimum number of operations (insert, edit or replace) to transform one string to another
 - Now introducing a “close” threshold so two points are considered as the same when they are close enough
 - **LCSS**: count the similar sample pairs
 - **EDR**: count the dissimilar sample pairs

Edit Distance is computed using dynamic programming. Check the Wikipedia page

■ Longest Common Sub-Sequence

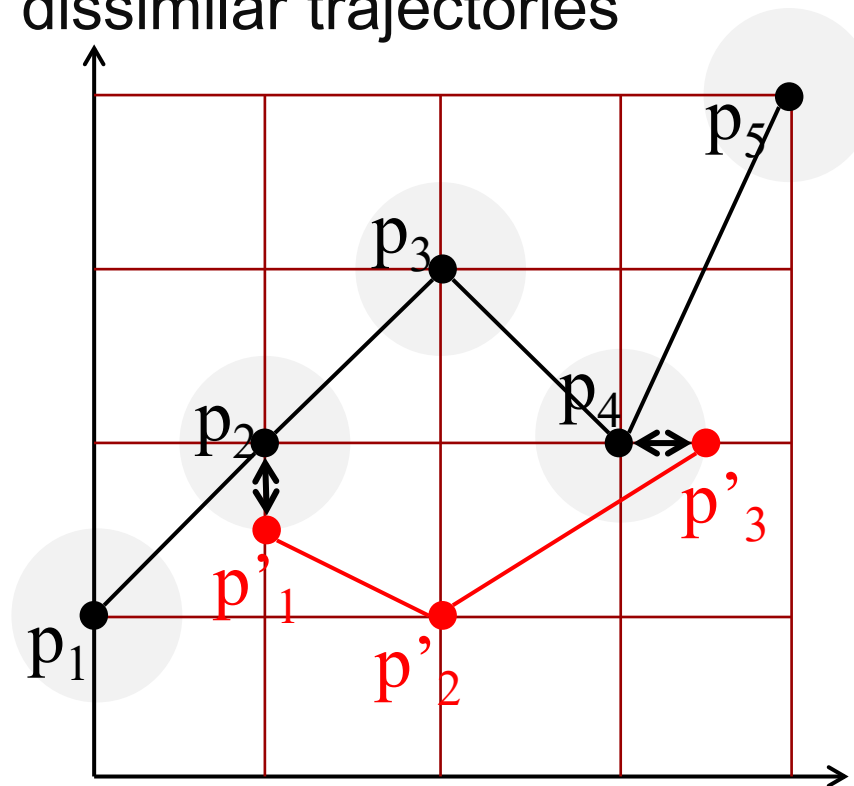
- To find the longest common subsequence (which may not be consecutive) between two strings
- $\text{LCSS}(\text{'abcde'}, \text{'bd'}) = ?$ $\text{LCSS}(\text{'abc'}, \text{'acb'}) = ?$
- 1-to-(1 or null) mapping
- This can also be computed using dynamic programming
 - Complexity $O(mn)$ where m and n are the lengths of the two strings

■ Adaptation of string similarity

- Two locations are regarded as equal if they are 'close' enough (compared to a threshold)

+ LCSS

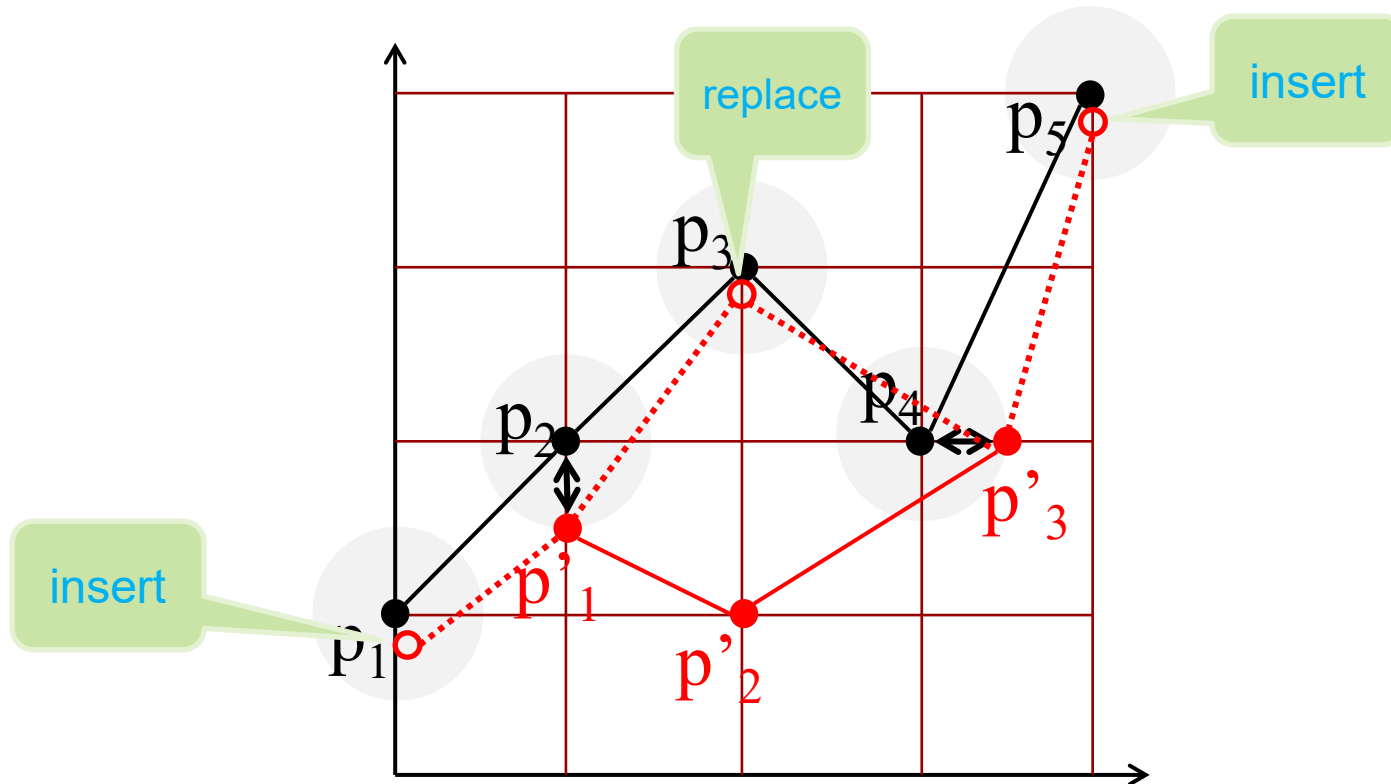
- Insensitive to noise
- Not easy to define threshold
- May return dissimilar trajectories



- Edit Distance on Real sequence
- Adaptation from Edit Distance on strings
 - Number of insert, delete, replace needed to convert one string into another
 - Two locations are regarded as equal if they're 'close' enough (compared to a threshold)

+ EDR

- Value means the number of operations, not “distance between locations”
- Insensitive to noise



+ LCSS and EDR

- They are both count-based
 - LCSS counts the number of matched pairs
 - EDR counts the cost of operations needed to fix the unmatched pairs
- Higher LCSS, lower EDR

+ Continuity

- So far, discrete measures only
 - Only consider the sample points of trajectory
 - All previous measures are in this category
- Continuous measures
 - Consider the line segments between samples
 - OWD
 - LIP

+ OWD

- One Way Distance from T1 to T2 is:
 - Integral of the distance from points of T1 to T2
 - Divided by the length of T1

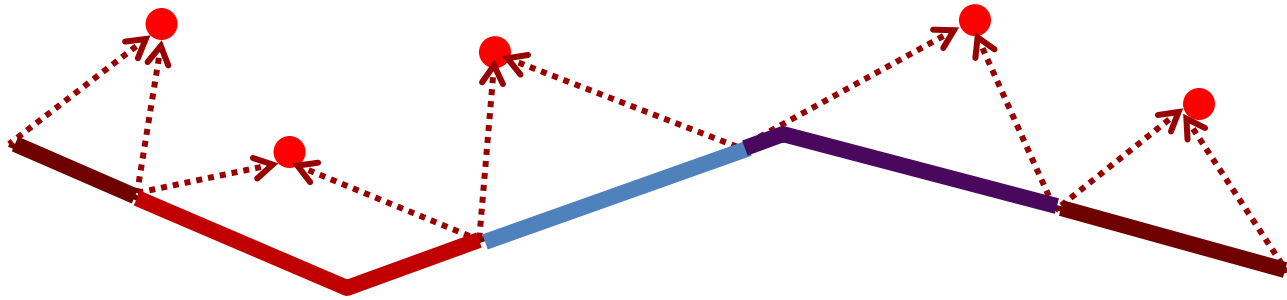
$$D_{\text{owd}}(T_1, T_2) = \frac{1}{|T_1|} \left(\int_{p \in T_1} D_{\text{point}}(p, T_2) dp \right)$$

- Make it into symmetric measure

$$D(T_1, T_2) = \frac{1}{2} (D_{\text{owd}}(T_1, T_2) + D_{\text{owd}}(T_2, T_1))$$

+ OWD example

- Consider one trajectory as piece-wise line segment, and the other as discrete samples



+ LIP distance

■ Locality In-between Polylines

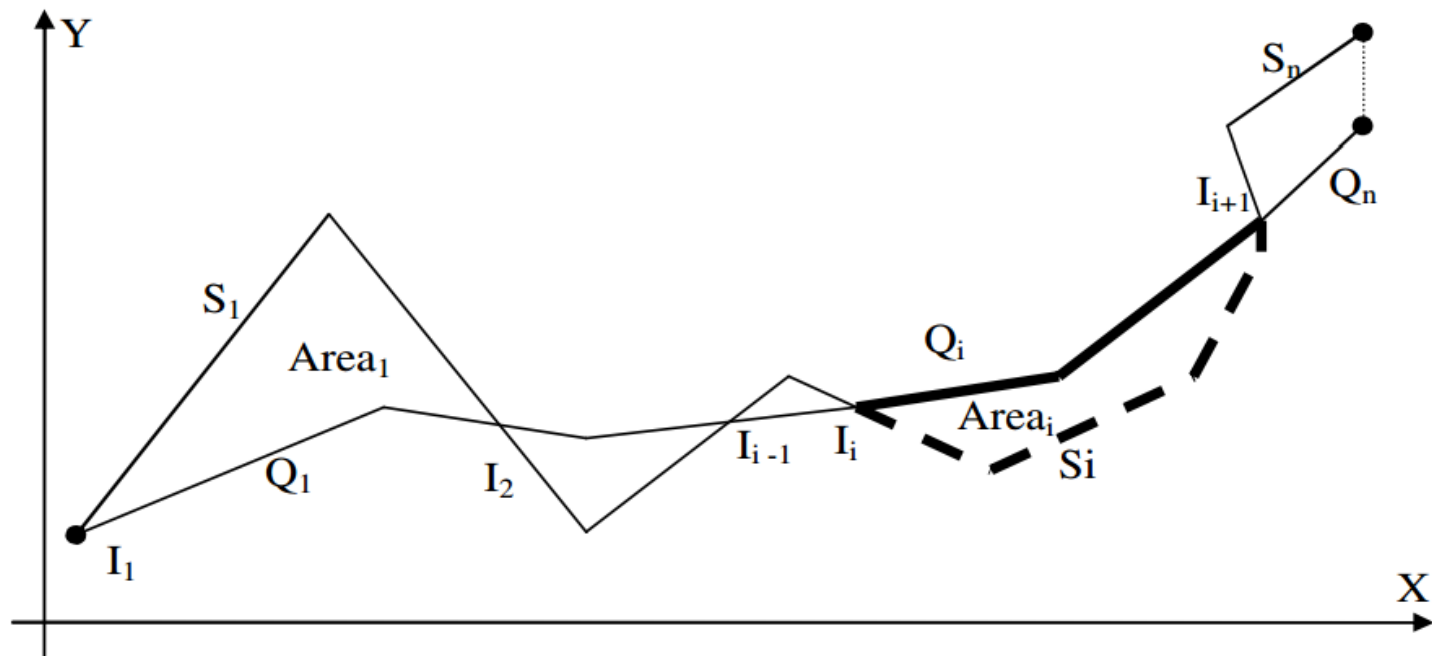
$$LIP(Q, S) = \sum_{\forall \text{ polygon}_i} Area_i \cdot w_i$$

- *Polygon* is the set of polygons formed between intersection points

- $$w_i = \frac{Length_Q(I_i, I_{i+1}) + Length_S(I_i, I_{i+1})}{Length_Q + Length_S}$$

+ LIP distance

- Only works for 2-dimensional trajectories
- Polygon \rightarrow polyhedron: non-trivial change

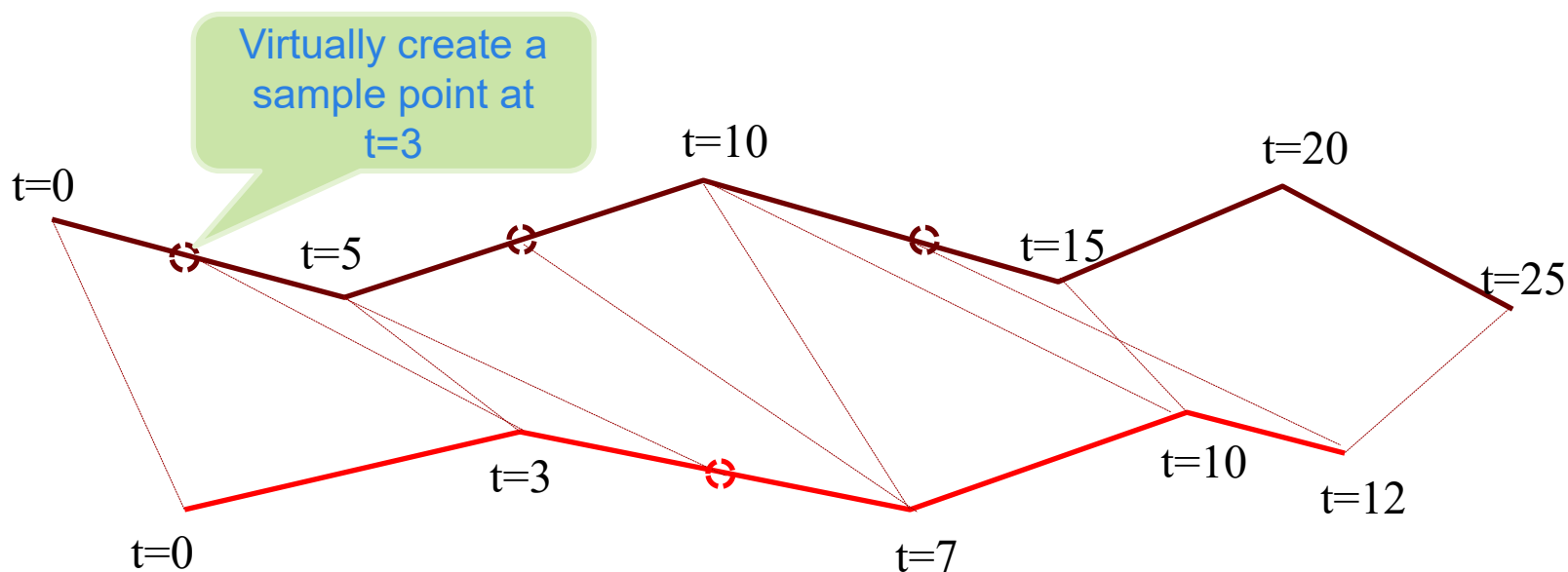


+ Spatiotemporal Distances

- So far, spatial only
 - Disregard the time information on sample points
- Spatiotemporal
 - Take the timestamp into consideration
 - Synchronous Euclidean Distance

■ Synchronous Euclidean Distance

- Euclidean distance between locations at the same time instance of two trajectories



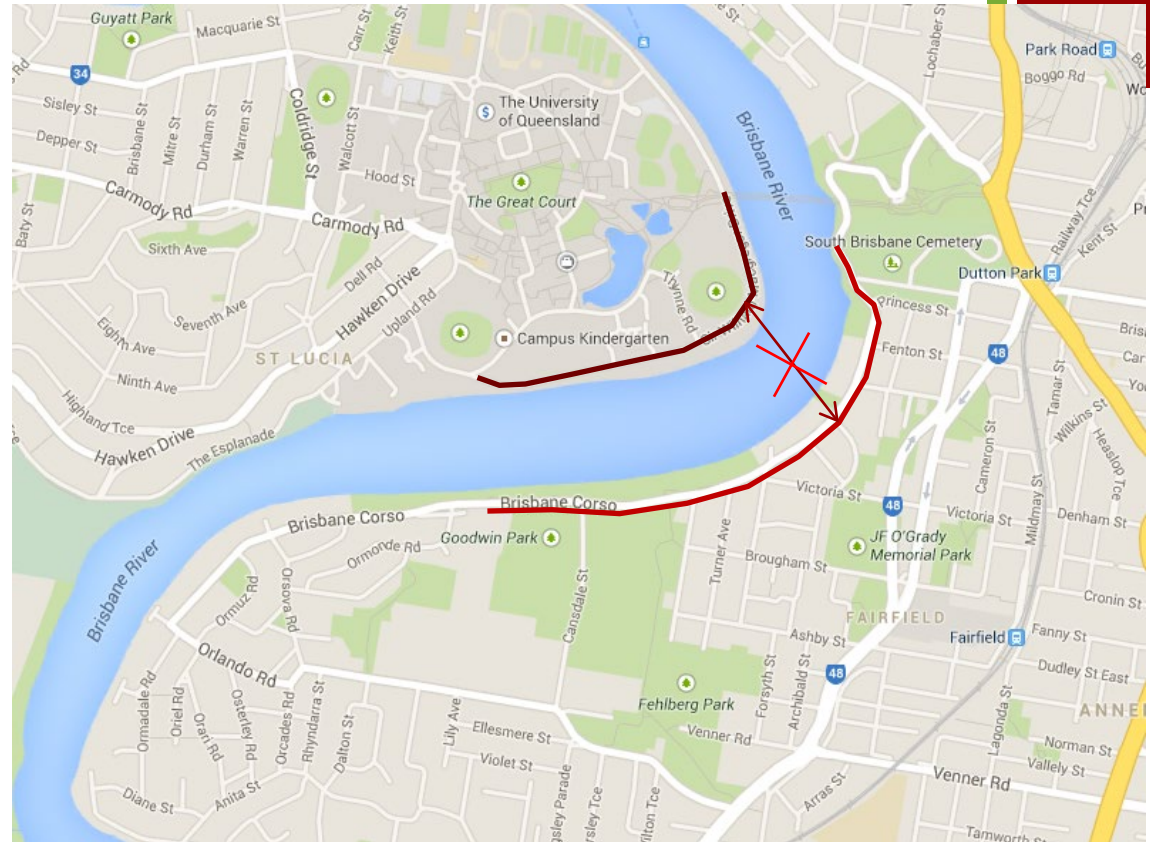
Mirco Nanni, Dino Pedreschi, Time-focused clustering of trajectories of moving objects. *Journal of Intelligent Information Systems* (2006)

POTAMIAS, M., PATROUMPAS, K., AND SELLIS, T. K. Sampling trajectory streams with spatiotemporal criteria. *SSDBM 2006*

+ Underlying Space

44

- We have considered on the Euclidean Space so far
- Road network
 - Distance is defined along shortest path



Jung-Rae Hwang, Hye-Young Kang, Ki-Joune Li, Searching for Similar Trajectories on Road Networks Using Spatio-temporal Similarity. Advances in Databases and Information Systems, pp 282 – 295, 2006

+ Trajectory Compression: The Need

45

- One record every 10 s, 24 bytes/record (min), 10 hrs/day

	Day	Month	Year
1 object	84KB	2.5MB	30MB
60,000 taxis	5GB	140GB	1.7TB

- Significant level of redundancy (i.e., the need for compression)
- Data quality can be poor (i.e., lossless compression may not be meaningful)

+ Simplification and Compression

■ Trajectory Simplification

- Removing redundant information in a trajectory

■ Trajectory Compression

- Reduce the amount of data without too much information loss

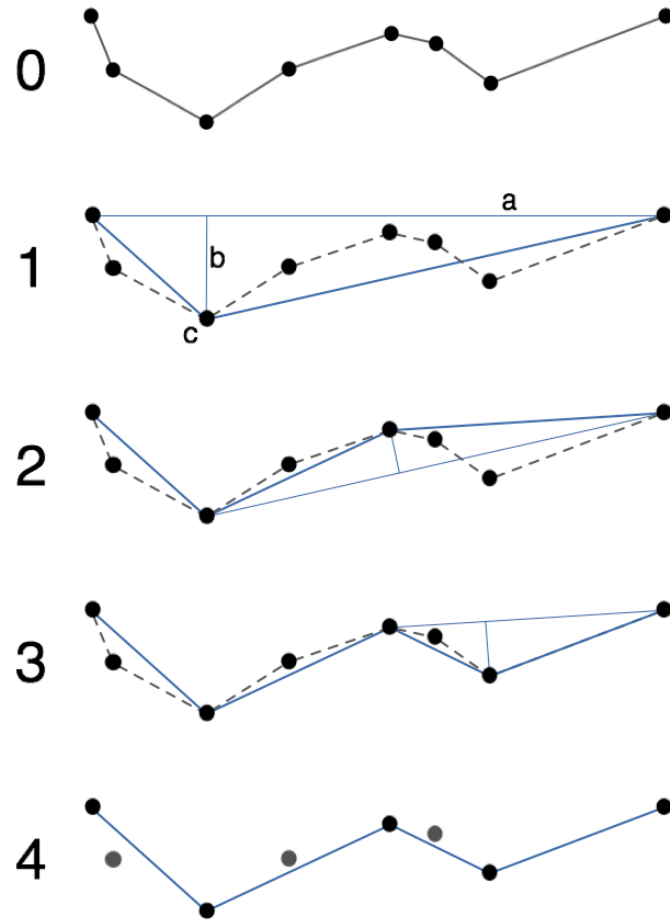
■ Questions

- Goals: size, quality, fitness for use, processing efficiency...
- Intra-, inter- or knowledge-assisted?
 - That is, within one trajectory, among a set of trajectories, and use of other information such as road network information or some kind of dictionary and patterns
- Geometric, spatiotemporal, semantic?

+ Douglas-Peucker Algorithm

52

- A classic algorithm, simple, and “most superior”
- Top-down, to check if the deviation from a straight line is acceptable
- Computationally costly: $O(n^2)$, reducible to $O(n \log n)$
- ‘Area difference’ can be used to measure too



"Algorithms for the reduction of the number of points required to represent a digitized line or its caricature", D. Douglas & T. Peucker, The Canadian Cartographer, 1973

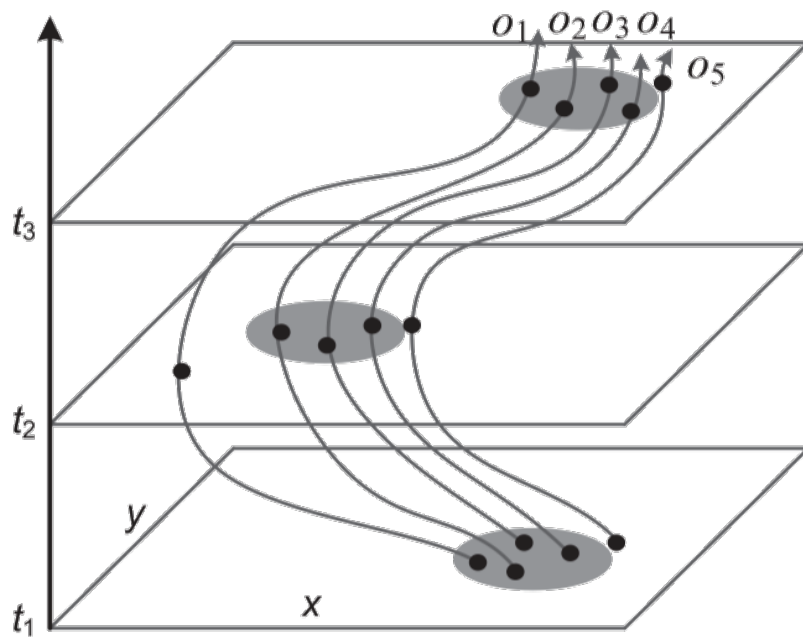
+ With Network Constraints

- Consider map-matching and trajectory compression at same time (by different orders)
 - Map-matching by Brakatsoulas, Pfooser, Salas and Wenk (VLDB 2005)
 - Trade-off between compression ratio and similarity
- New idea: using shortest path for compression
 - Between which pair of points?
 - Key technique: minimum description length (MDL)
 - $L(H) + L(D|H)$ (for compression and differences caused)
 - Need to consider all-pair combination to optimize
 - Using a greedy approach, with pre-computed requires all-pair SP

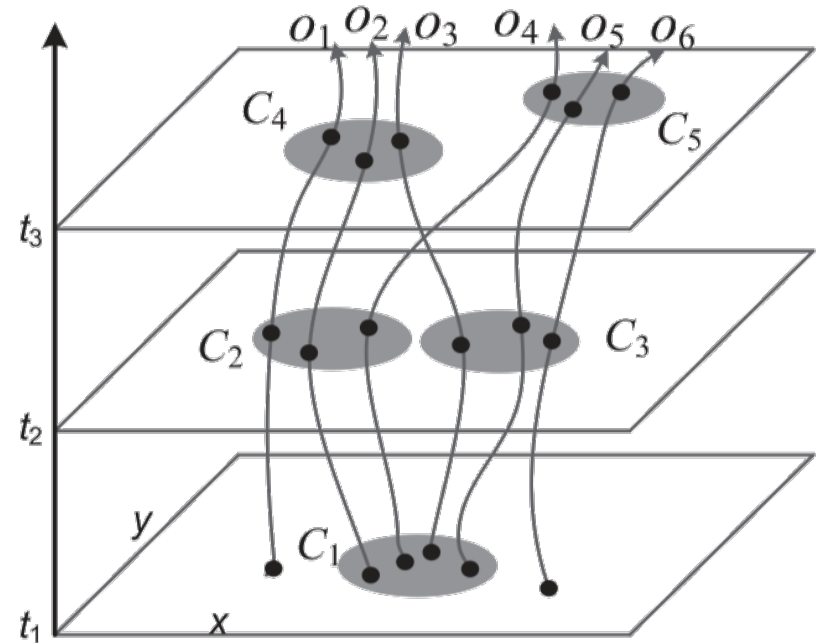
+ Trajectory Mining

- Trajectory clustering (SIGMOD 2007)
 - Partition-and-group, and partition based on MDL
- Trajectory pattern mining (ICDE 2008)
 - Long trajectories are divided based a duration T , and then they are aligned and points are density-based clustered
- Finding convoys (VLDB 2008)
 - A group of objects travelling together for long enough
 - Used line simplification algorithms for efficiency
 - Finding “swarms” (VLDB 2010)
 - Allow temporary divergence

+ Trajectory Mining



(a) Flock, convoy and swarm



(b) Gathering

+ Summary

- Spatiotemporal data is common
- Spatiotemporal data indexing and query processing are different from spatial ones
- Trajectory data is of particular importance with a wide range of applications
- Raw trajectory data cannot be compared directly
- Trajectory analytics is a new research direction