

# Deep Learning for Human Mobility Analytics

-- L3: Spatio-Temporal Data Sensing and Management

**Yuxuan Liang (梁宇轩)**

INTR & DSA Thrust

[yuxuanliang@hkust-gz.edu.cn](mailto:yuxuanliang@hkust-gz.edu.cn)



香港科技大学(广州)  
THE HONG KONG  
UNIVERSITY OF SCIENCE AND  
TECHNOLOGY (GUANGZHOU)



# Notes

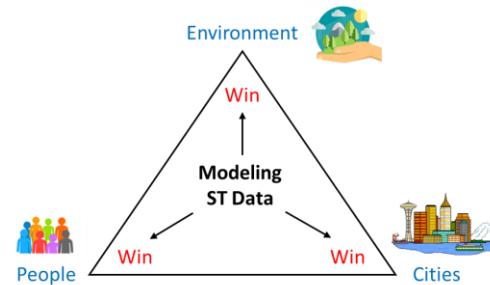
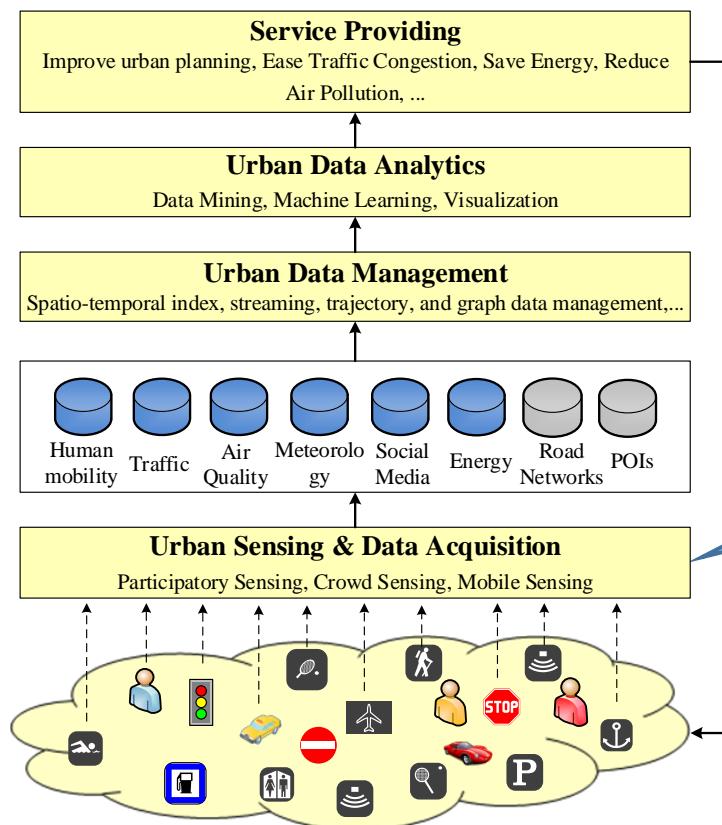
13	16 Sep	Transfer to <b>an online course for project discussion</b> (each team 20 mins) in <b>early November</b> . Details will be informed in late October		Mid-Autumn Festival
3	23 Sep	L3: Spatio-Temporal Data Sensing and Management	✓	To form team (2 people)
4	30 Sep	L4: Cross-Domain Data Fusion and Multimodal Learning	✓	



# Revisit our Last Course

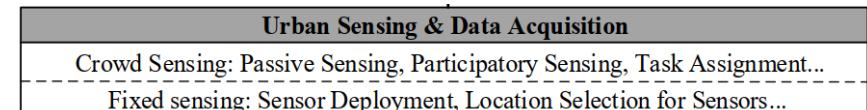


# 1<sup>st</sup> Stage: Urban Sensing & Data Acquisition



Collecting urban data through

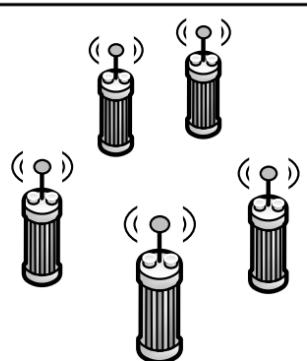
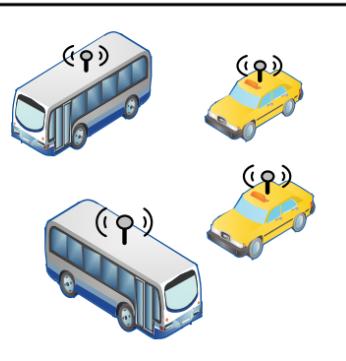
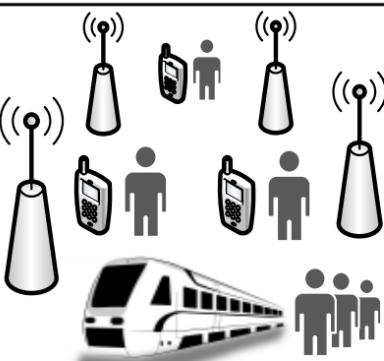
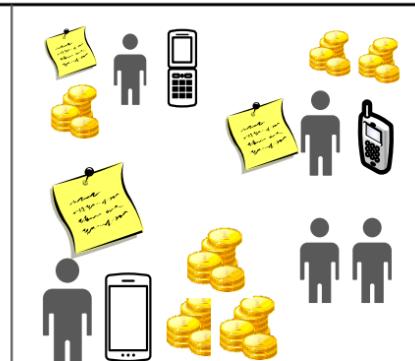
- Sensor-centric sensing
- Human-centric sensing



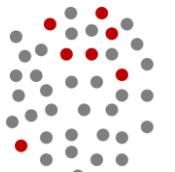


# Spatio-Temporal Sensing Modes

- Sensor-centric sensing
  - Fixed (location) sensing
  - Mobile (objects) sensing
- Human-centric sensing
  - Passive crowd sensing
  - Active crowd sensing

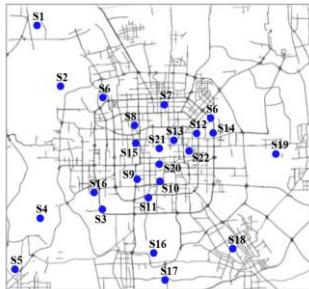
 <p>A) Fixed Sensing</p>	 <p>B) Mobile Sensing</p>	 <p>C) Passive Crowd Sensing</p>	 <p>D) Active Crowd Sensing</p>
Sensor-Centric Sensing		Human-Centric Sensing	

# Challenges of Urban Sensing

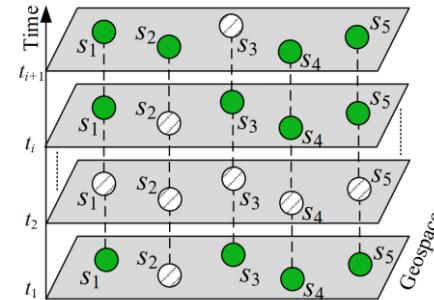


- Samples
- Other points

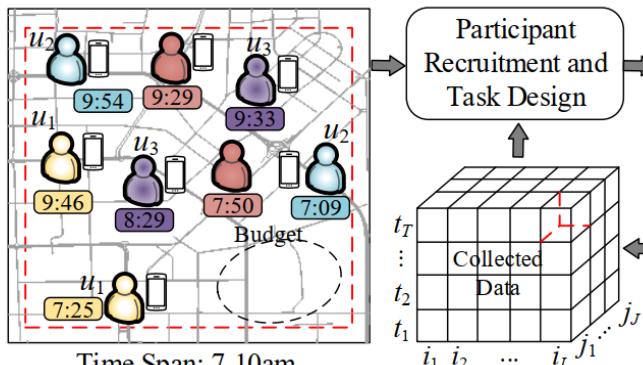
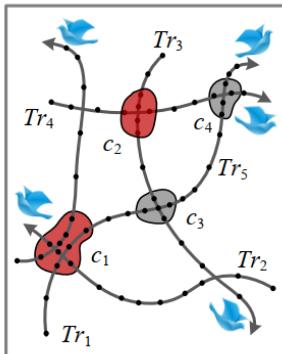
## Biased distribution



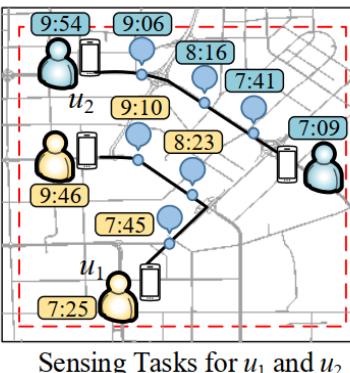
## Data sparsity



## Data missing



## Resource deployment



# Biased Distribution



- Inferring Gas Consumption and Pollution Emission of Vehicles throughout a City

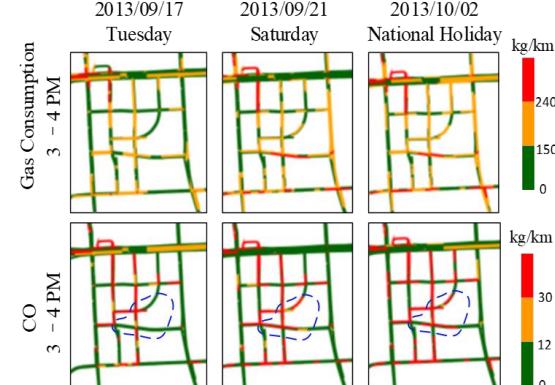
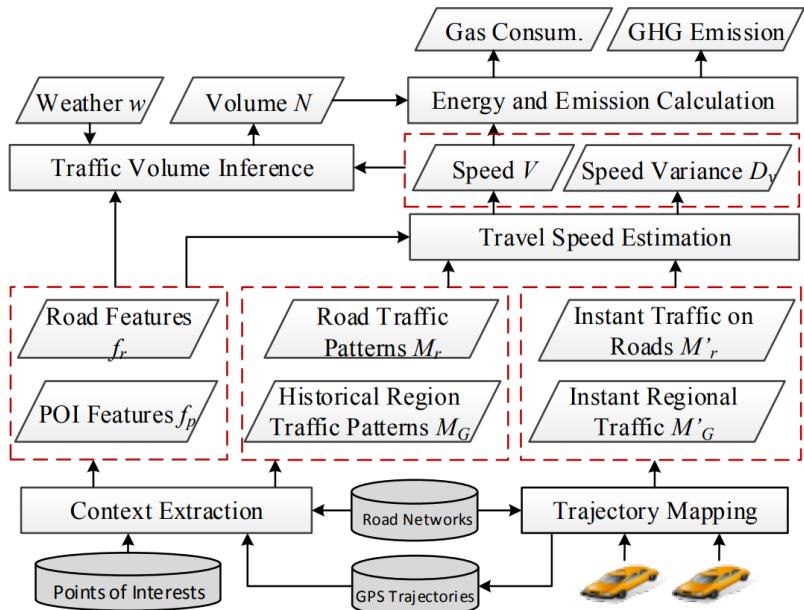


Figure 2. Framework of our method



# Data Sparsity



UrbanAir

实时·细粒度·空气质量

English | 中文

北京 ▾ PM2.5 ▾

良

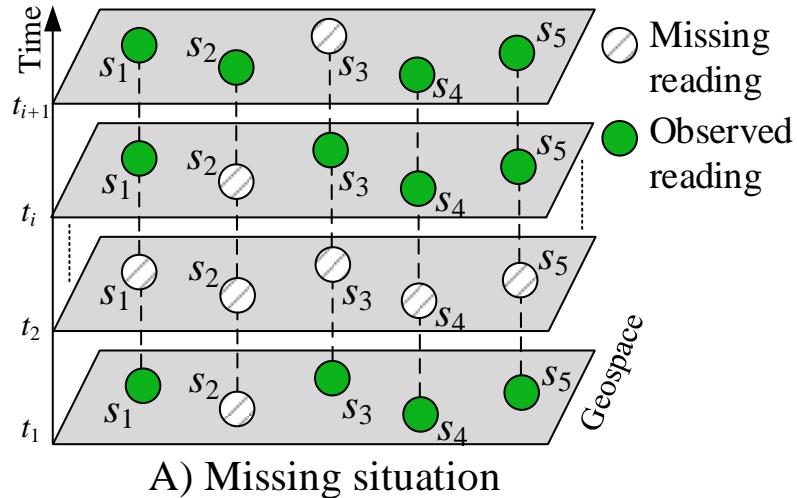


9°C

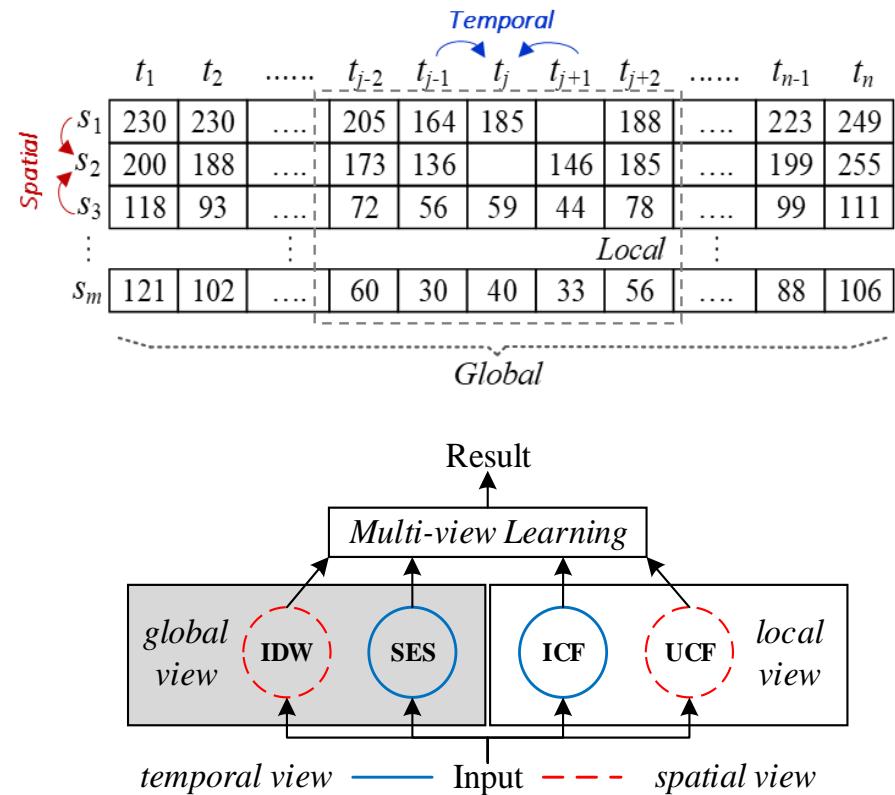




# Data Missing

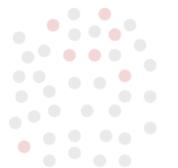


A) Missing situation



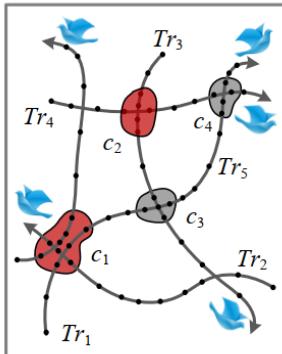


# Challenges of Urban Sensing

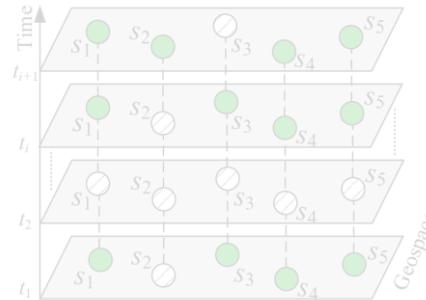


● Samples  
● Other points

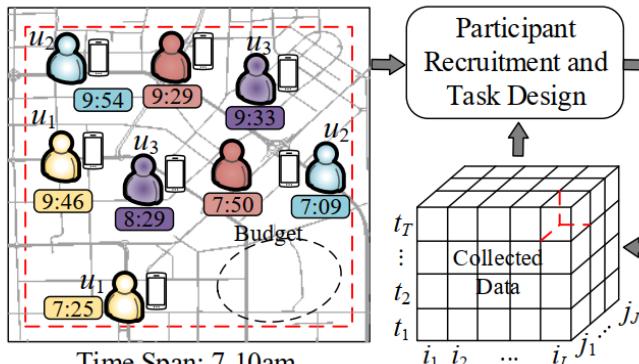
Biased distribution



Data sparsity



Data missing

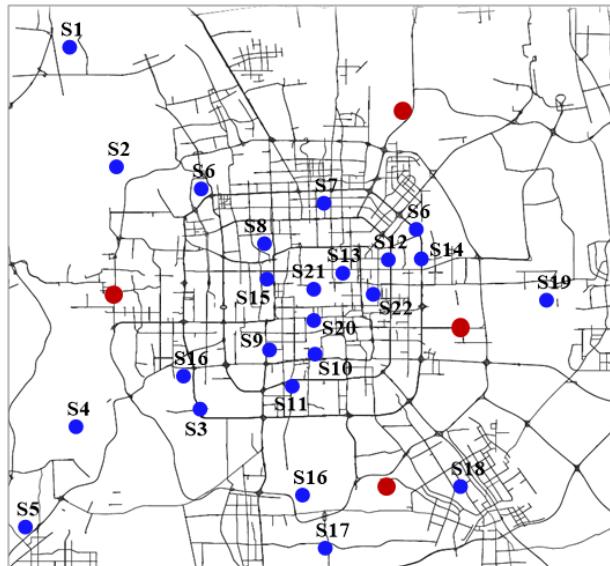


Resource deployment

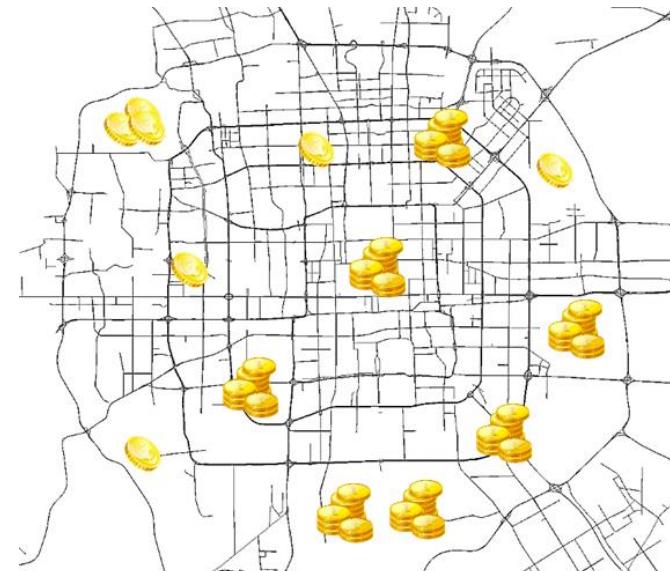


# A limited resource (budget, labors, land...)

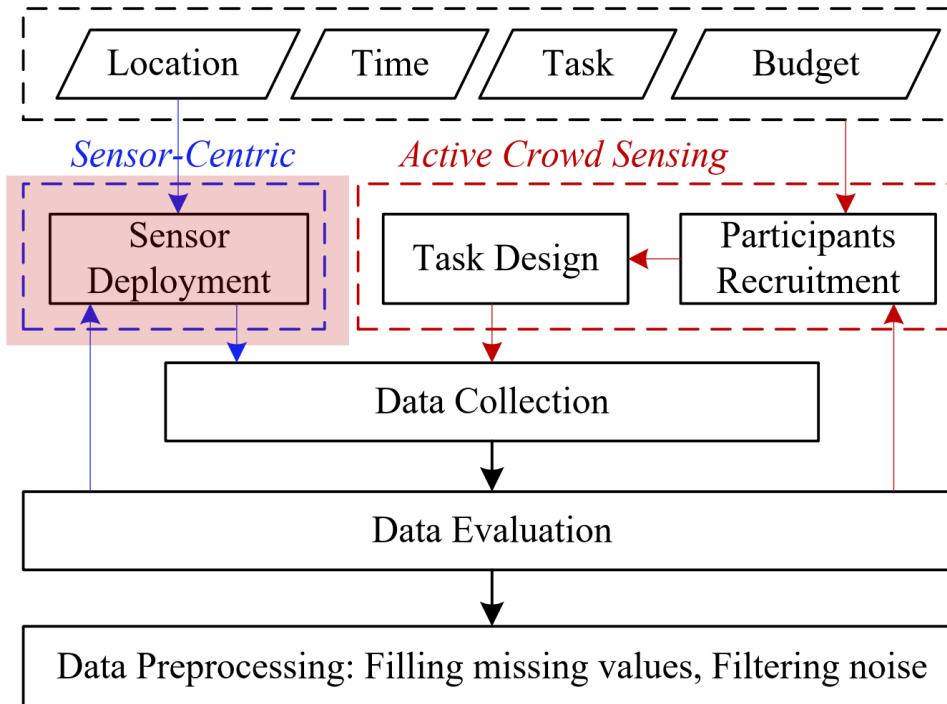
- **Static sensing:** Where to deploy sensor to maximize the gain?



- **Crowdsensing:** How to arrange the incentives dynamically?



# General Framework of Urban Sensing



# Resource Deployment and Location Selection



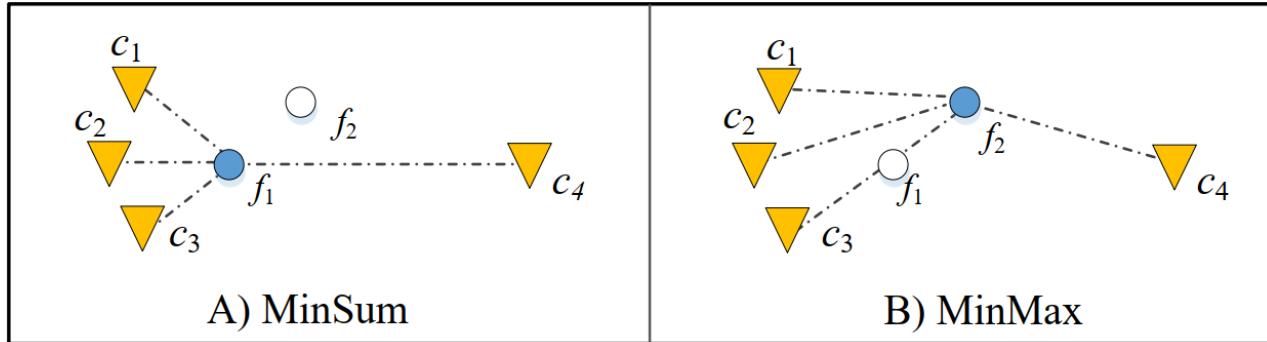
- Optimal meeting points problems
- Maximum coverage problem -- submodular maximization
- Learning-to-rank problems
- Uncertainty minimization problems (entropy)



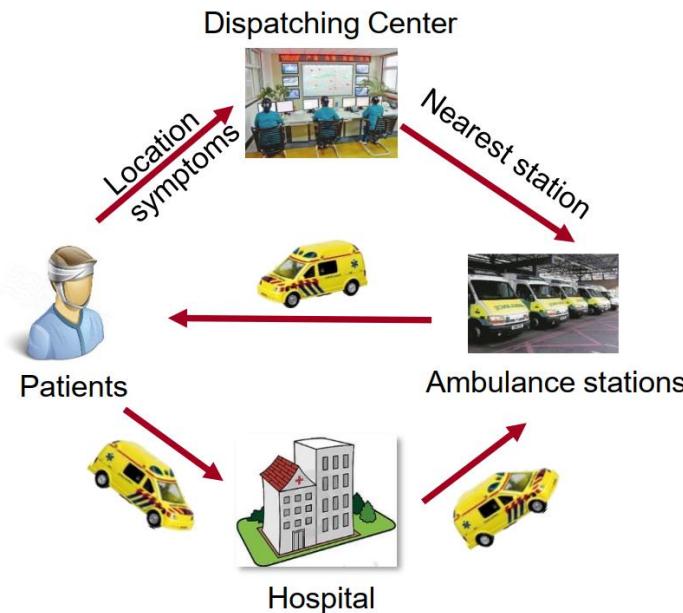
# Optimal Meeting Points

- locate  $k$  facilities such that
  - MinSum: Minimizing the average cost to reach all clients can be minimized
  - MinMax: Minimizing the maximum cost to reach those clients

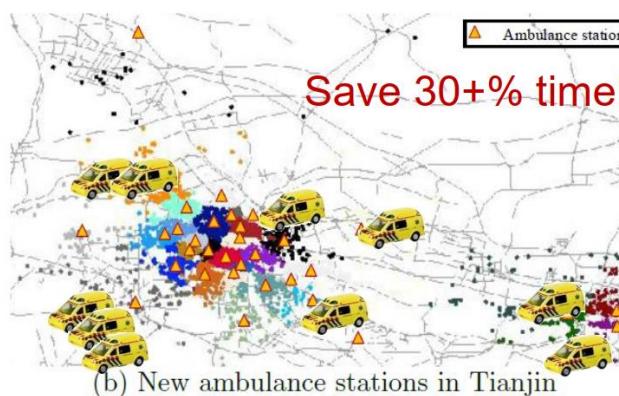
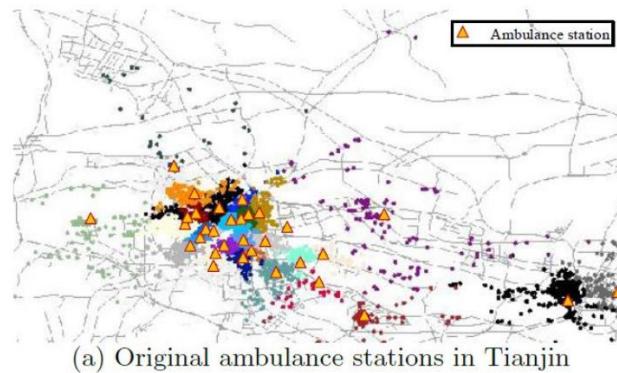
NP-hard!



# Improving Medical Emergency Services using Big Data



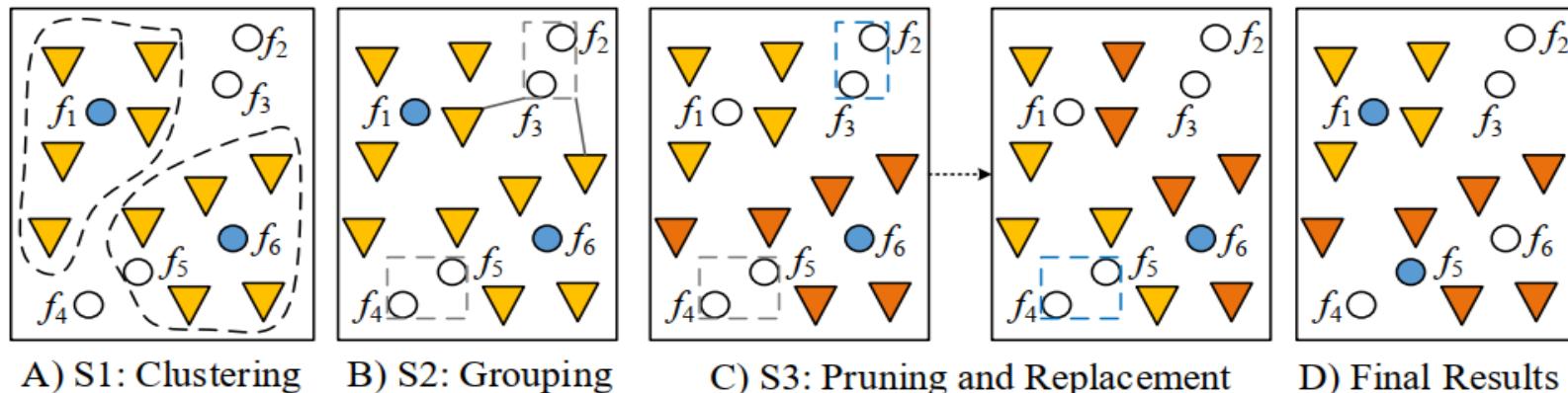
- Select locations for Ambulance Stations
- Dynamic ambulance allocation



# Deploying Ambulance Stations Based on Optimal Meeting Point Model



- Methodology
  - k-medoids clustering
  - Candidate grouping
  - Iterations: pruning and replacement

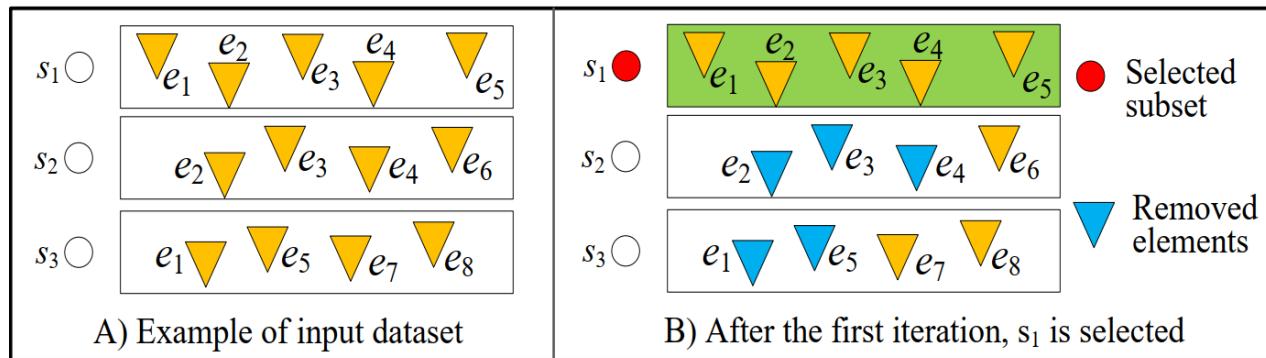




# Maximum Coverage Problem

- Select a subset of ( $k$ ) candidates to maximize
  - Nodes from a graph
  - Trajectories from a road network
  - Topic coverage from a document corpus

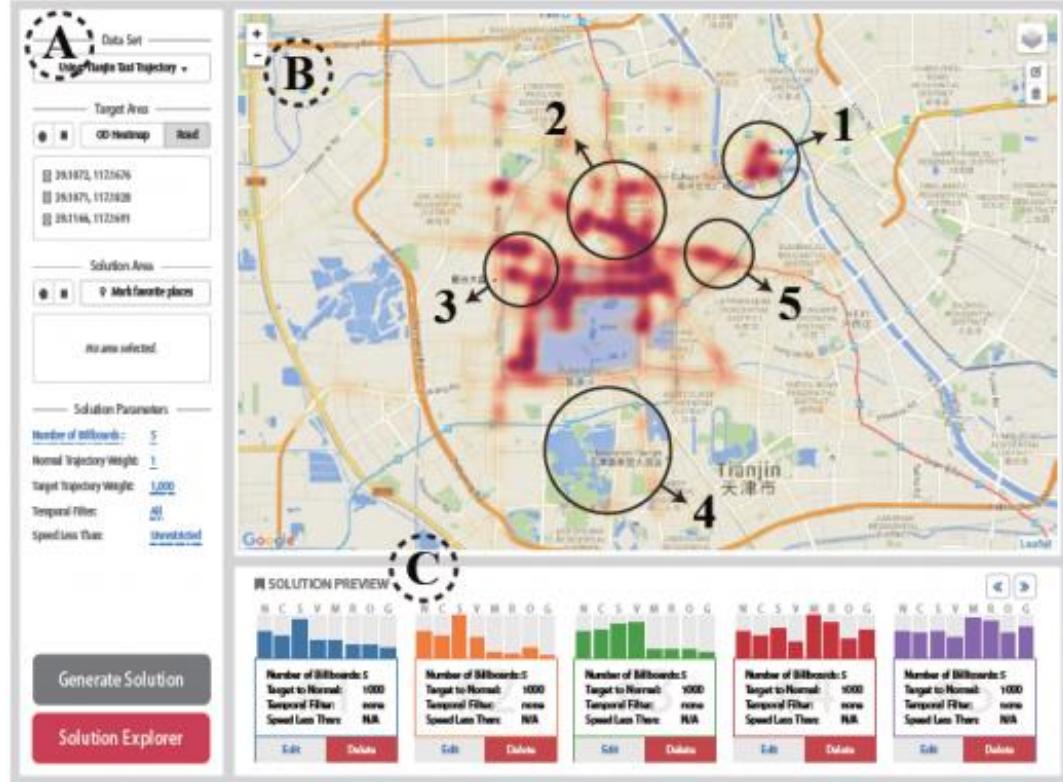
NP-hard!



$(1 - 1/e)$  approximation guarantee



# Finding Top-k Most Influential Location Set

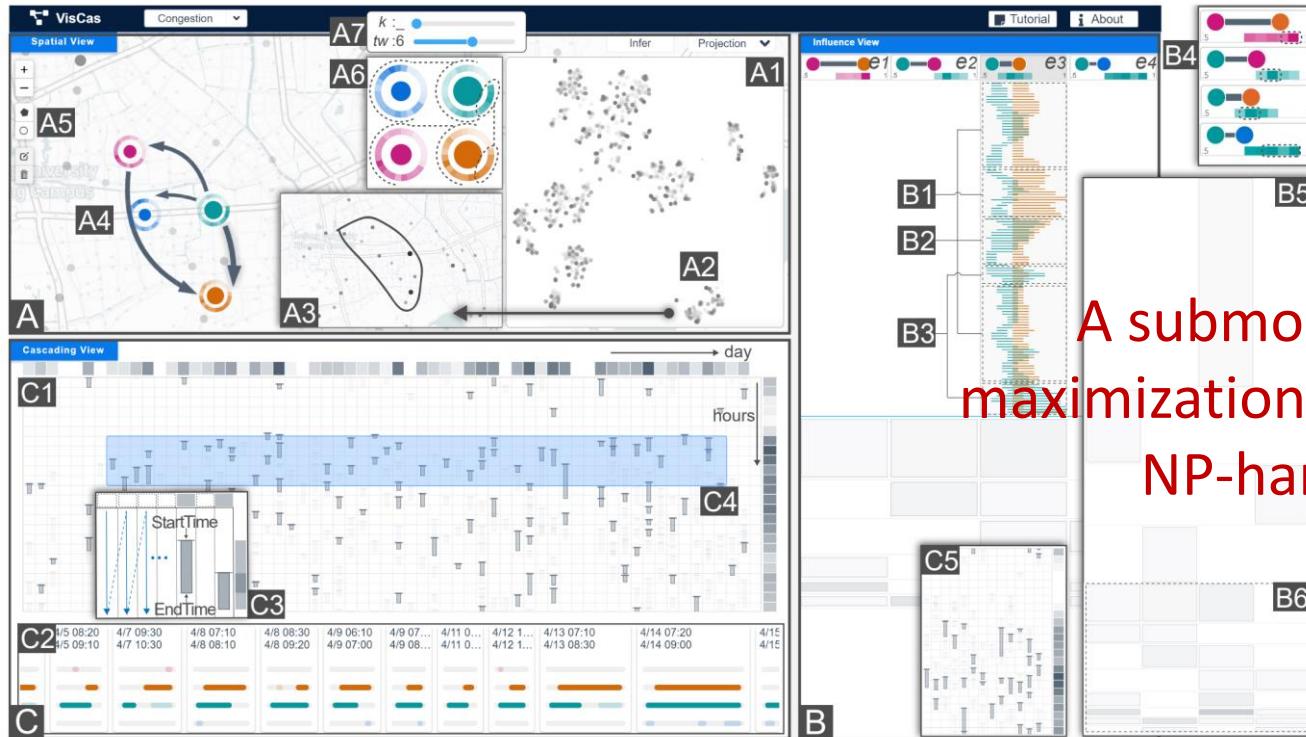


A submodular  
maximization problem,  
NP-hard!

Interactive Visual Data  
Analytics



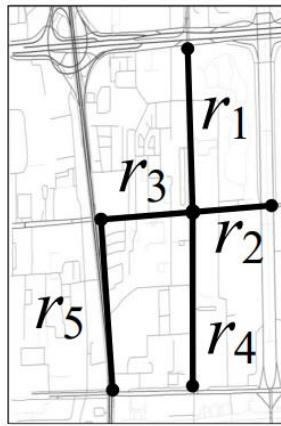
# Inferring Traffic Cascading Patterns



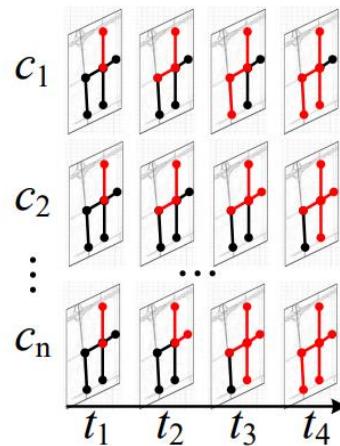


# Inferring Traffic Cascading Patterns

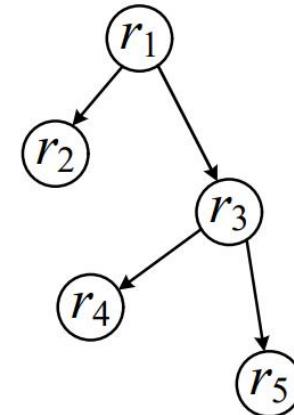
- Target: to uncover how traffic congestion propagates thru road networks



(a) Road Network



(b) Observed Cascades



(c) Cascading Pattern



# Inferring Traffic Cascading Patterns

- Also a maximum coverage problem

## Cascading Pattern Construction

### Model formulation

- Given a **propagation tree**  $T$ , the likelihood of a **cascade**  $c$ :

$$f(c|T) \xrightarrow{\text{Joint likelihood of } C} f(C|G) = \prod_{c \in C} f(c|G)$$

Conditional independence

$$f(c|G) = \sum_{T \in T_c(G)} f(c|T) f(T|G)$$

simplify

For a specific cascade, this part is equal.

$$\text{Optimization target: } \hat{G} = \arg \max_{|G| \leq k} f(C|G) \quad \text{difficult to optimize}$$

$\hat{G}$  best explains the observed cascades, and the maximization is over all possible graphs  $G$  at **most  $k$  edges**

## Alternative Optimization (CasInf)

$$\hat{G} = \arg \max_{|G| \leq k} f(C|G)$$

- matrix tree theorem
- Combine environmental intensity
- Log-likelihood

$$\hat{G} = \arg \max_{|G| \leq k} F(C|G)$$

- Alternative target
- It satisfies **submodularity**, a natural diminishing returns property.

### Optimization: greedy algorithm

#### Algorithm 1 Approximate algorithm for CasInf

##### Input:

- $k$ : the number of edges in cascading patterns we infer.
- $C$ : the set of cascades obtained in a time span.
- $D$ : a constant denoting the spatial constraint.

##### Output:

```

1:  $G \leftarrow \emptyset;$ 
2:  $P \leftarrow \text{all pairs } (j, i): \exists c \in C \text{ with } t_j < t_i \text{ and } d_{j,i} < D$ 
3: while  $|G| \leq k$  do
4:   for all  $(j, i) \in P \setminus G$  do
5:      $\delta_{j,i} = 0;$ 
6:     for all  $c: t_j < t_i$  do
7:        $e_g \leftarrow \text{environmental intensity inference}$ 
8:        $w_c(m, n) \leftarrow \text{weight of } (m, n) \text{ in } G \cup (j, i);$ 
9:       for all  $m: t_m < t_i$  and  $m \neq j$  do
10:         $\delta_{c,j,i} = \delta_{c,j,i} + w_c(m, i);$ 
11:      end for
12:       $\delta_{j,i} = \delta_{j,i} + \log(\frac{\delta_{c,j,i} + w_c(j, i)}{\delta_{c,j,i}});$ 
13:    end for
14:  end for
15:   $(j^*, i^*) \leftarrow \arg\max_{(j, i) \notin G} \delta_{j,i};$ 
16:   $G \leftarrow G \cup (j^*, i^*);$ 
17: end while

```

In each iteration, we select

$$e_i = \arg \max_{e \in G_i \setminus G_{i-1}} F(C|G_{i-1} \cup e) - F(C|G_{i-1}).$$



# Inferring Traffic Cascading Patterns

- Also a maximum coverage problem

## Cascading Pattern Construction

### Model formulation

- Given a **propagation tree**  $T$ , the likelihood of a **cascade**  $c$ :

$$f(c|T) \xrightarrow{\text{Joint likelihood of } C} f(C|G) = \prod_{c \in C} f(c|G)$$

Conditional independence

$$f(c|G) = \sum_{T \in T_c(G)} f(c|T) f(T|G)$$

simplify

$$f(c|G) \propto \sum_{T \in T_c(G)} f(c|T)$$

For a specific cascade, this part is equal.

Optimization target:  $\hat{G} = \arg \max_{|G| \leq k} f(C|G)$  difficult to optimize

$\hat{G}$  best explains the observed cascades, and the maximization is over all possible graphs  $G$  at **most  $k$  edges**

## Alternative Optimization (CasInf)

$$\hat{G} = \arg \max_{|G| \leq k} f(C|G)$$

1. matrix tree theorem

2. Combine environmental intensity

3. Log-likelihood

$$\hat{G} = \arg \max_{|G| \leq k} F(C|G)$$

◦ Alternative target

◦ It satisfies **submodularity**, a natural diminishing returns property.

Optimization: greedy algorithm

#### Algorithm 1 Approximate algorithm for CasInf

**Input:**

$k$ : the number of edges in cascading patterns we infer.  
 $C$ : the set of cascades obtained in a time span.  
 $D$ : a constant denoting the spatial constraint.

**Output:**

```

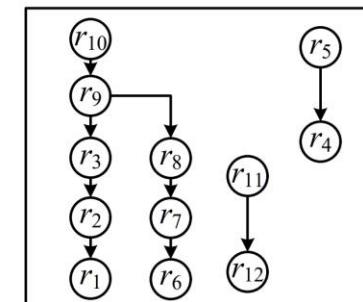
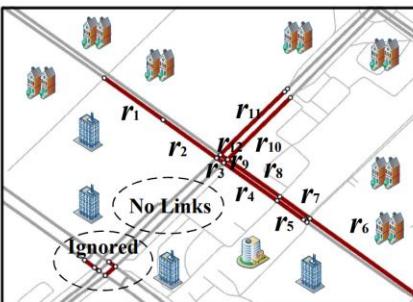
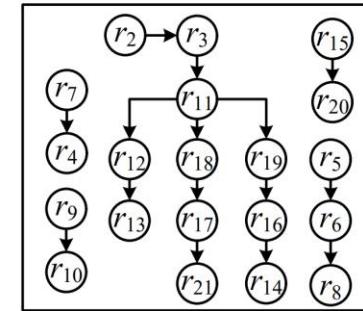
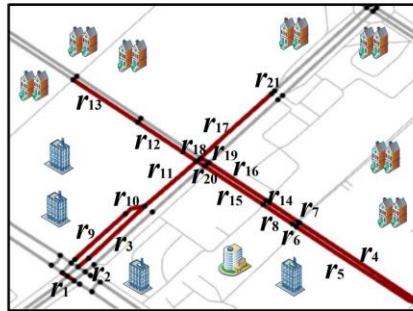
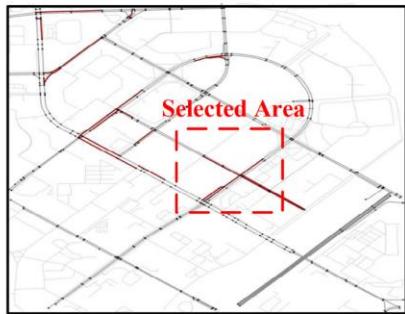
1:  $G \leftarrow \emptyset;$ 
2:  $P \leftarrow$  all pairs  $(j, i)$ :  $\exists c \in C$  with  $t_j < t_i$  and  $d_{j,i} < D$ 
3: while  $|G| \leq k$  do
4:   for all  $(j, i) \in P \setminus G$  do
5:      $\delta_{j,i} = 0;$ 
6:     for all  $c$ :  $t_j < t_i$  do
7:        $e_g \leftarrow$  environmental intensity inference
8:        $w_c(m, n) \leftarrow$  weight of  $(m, n)$  in  $G \cup (j, i)$ ;
9:       for all  $m : t_m < t_i$  and  $m \neq j$  do
10:         $\delta_{c,j,i} = \delta_{c,j,i} + w_c(m, i);$ 
11:      end for
12:       $\delta_{j,i} = \delta_{j,i} + \log(\frac{\delta_{c,j,i} + w_c(j, i)}{\delta_{c,j,i}});$ 
13:    end for
14:  end for
15:   $(j^*, i^*) \leftarrow \operatorname{argmax}_{(j,i) \notin G} \delta_{j,i};$ 
16:   $G \leftarrow G \cup (j^*, i^*);$ 
17: end while

```

In each iteration, we select

$$e_i = \arg \max_{e \in G_i \setminus G_{i-1}} F(C|G_{i-1} \cup e) - F(C|G_{i-1}).$$

# Inferring Traffic Cascading Patterns





# Learning to Rank

- Rank a set of locations so that the top-k best candidates can be selected to deploy resources or facilities

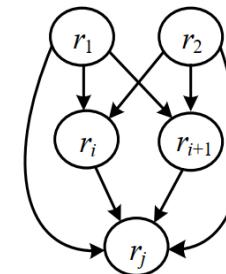
- Methods

- Pointwise
- Pairwise
- Listwise

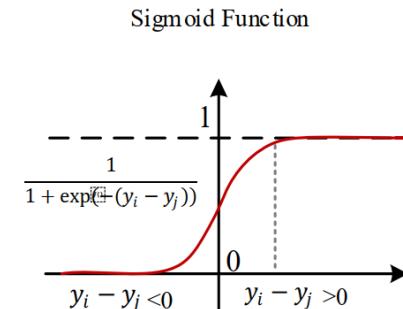
Real estates	Features	Ratio	Rank
$r_1$	$X_1$	35%	1
$r_2$	$X_2$	34%	1
$r_i$	$X_i$	25%	2
$r_{i+1}$	$X_{i+1}$	23%	2
$r_j$	$X_j$	12%	3

A) Ranking of Real Estates

A pairwise ranking representation



B) Graphical Representation of Orders

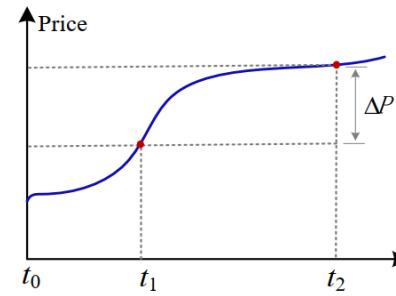


C) Keeping Pairwise Orders



# Ranking and Clustering Real Estates using Big Data

- Values (learned from big data)
  - Increase more in a rising market
  - Decrease less in a falling market



A) The price of a real estate

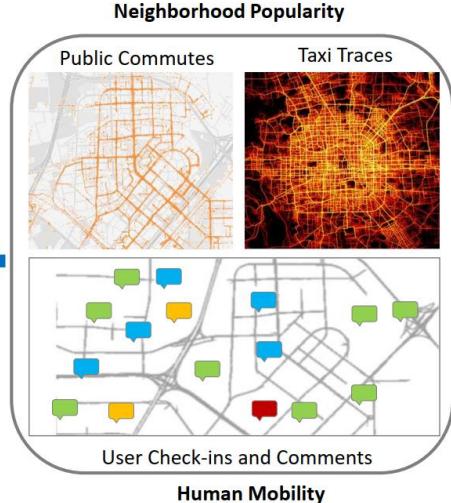
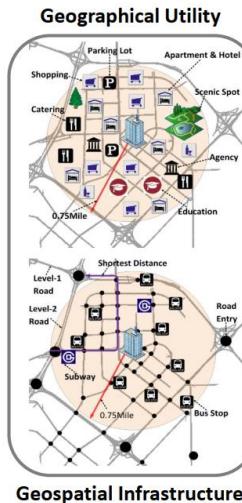
House	Increase	Rank↓
H1	35%	R1
H5	29%	R1
H4	13%	R2
...	...	...
H2	9%	R3
H3	2%	R3
H6	-1.5%	R4
H7	-6.1%	R5

B) Rank of estates by  $\Delta P$





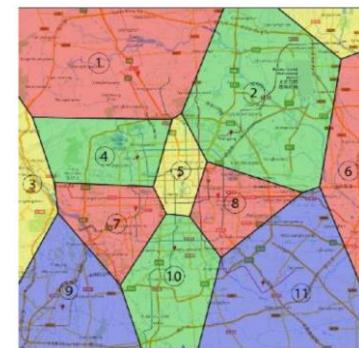
# Learning-to-Rank Problems



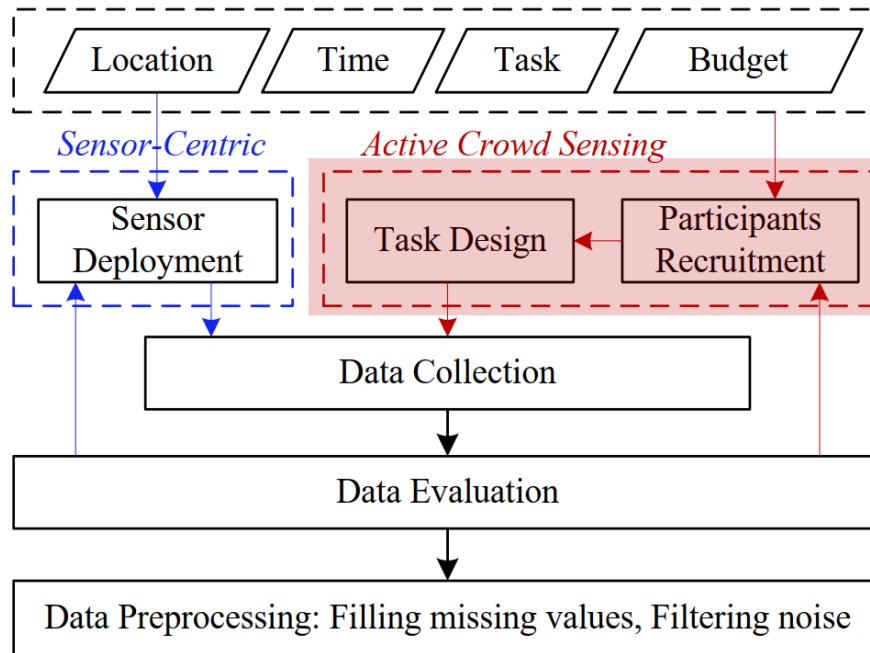
Ranking of Real Estates



Clusters of Real Estates



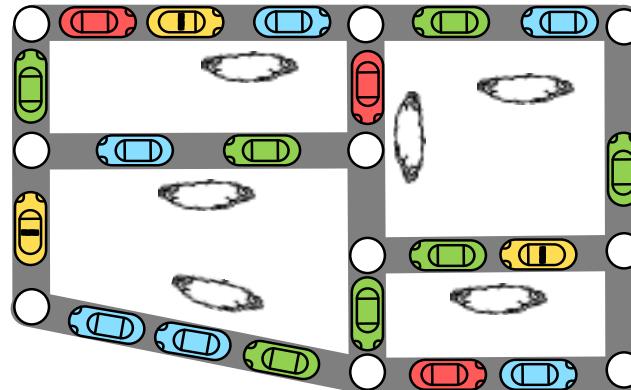
# General Framework of Urban Sensing



# Urban Sensing Based on Human Mobility

Ubicomp 2016

Shenggong Ji, et al. [Urban Sensing Based on Human Mobility](#), UBICOMP 2014

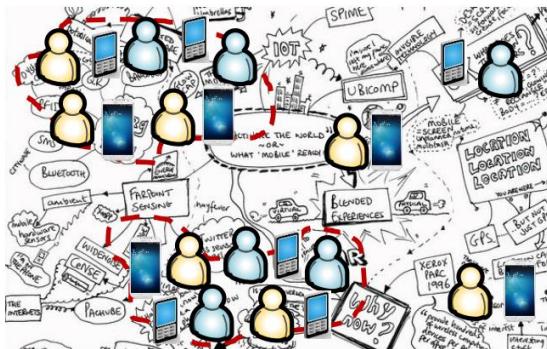


# Urban Sensing

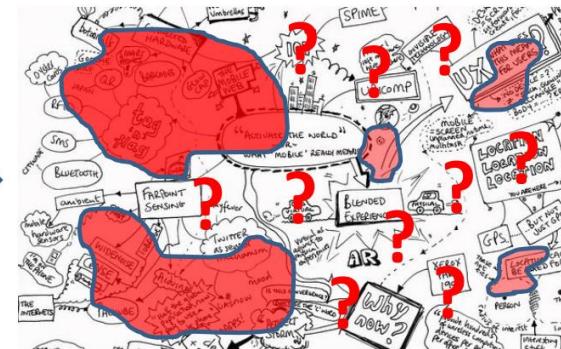


- Collecting urban data
  - Noise, temperature, air quality
  - Human as a sensor
- Brings challenges to
  - City-scale real-time monitoring
  - Further data analytics

Skewed human mobility



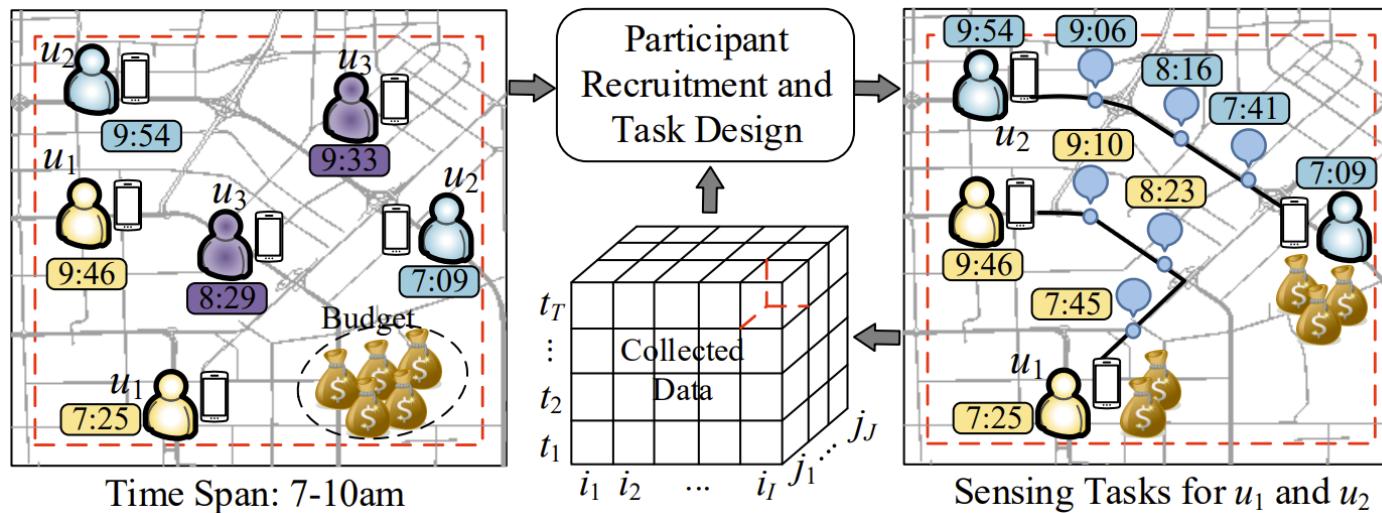
Imbalanced data coverage



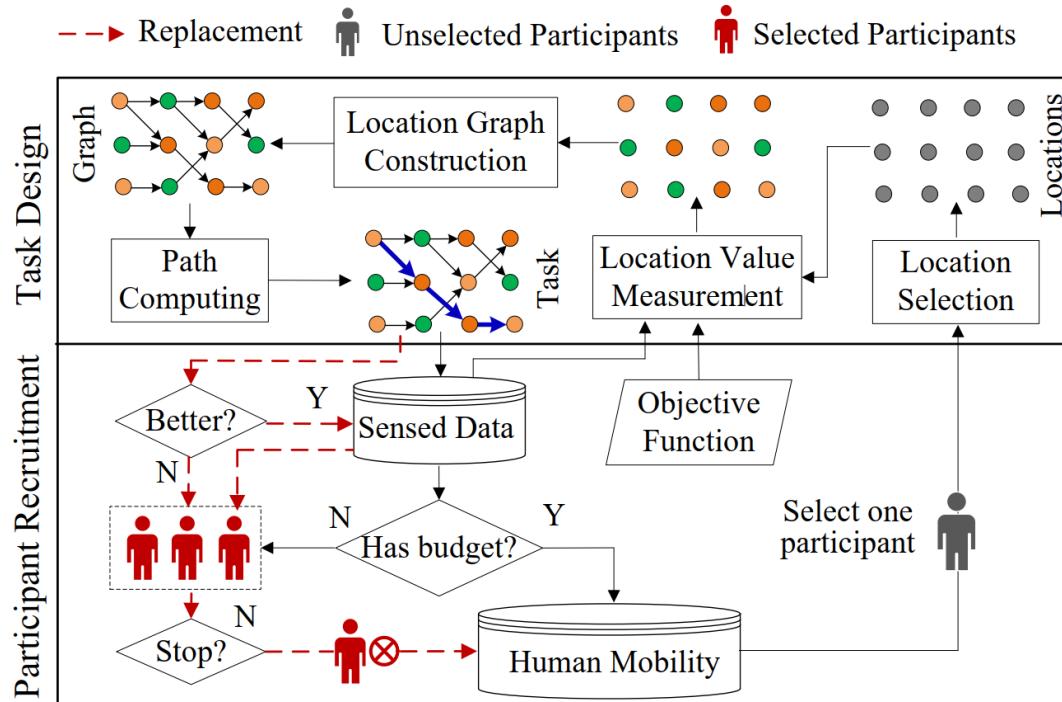


# An Urban Sensing Framework

- Consider real-world human mobility
- Maximize the amount and balance of collected data
- Given a limited budget



# Human Mobility-Based Urban Sensing Framework

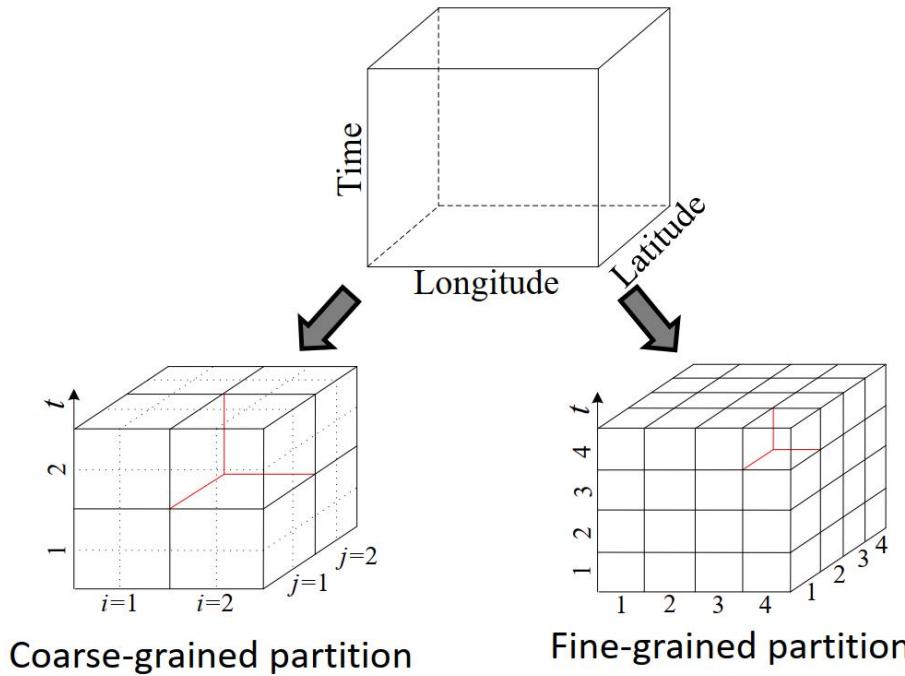


- A participant recruitment mechanism
  - random recruitment
  - replacement-based refinement

$$\max \phi = \alpha \times E + (1 - \alpha) \times \log_2 Q$$

$\alpha$ : the relative preference of **data balance** to **data amount**

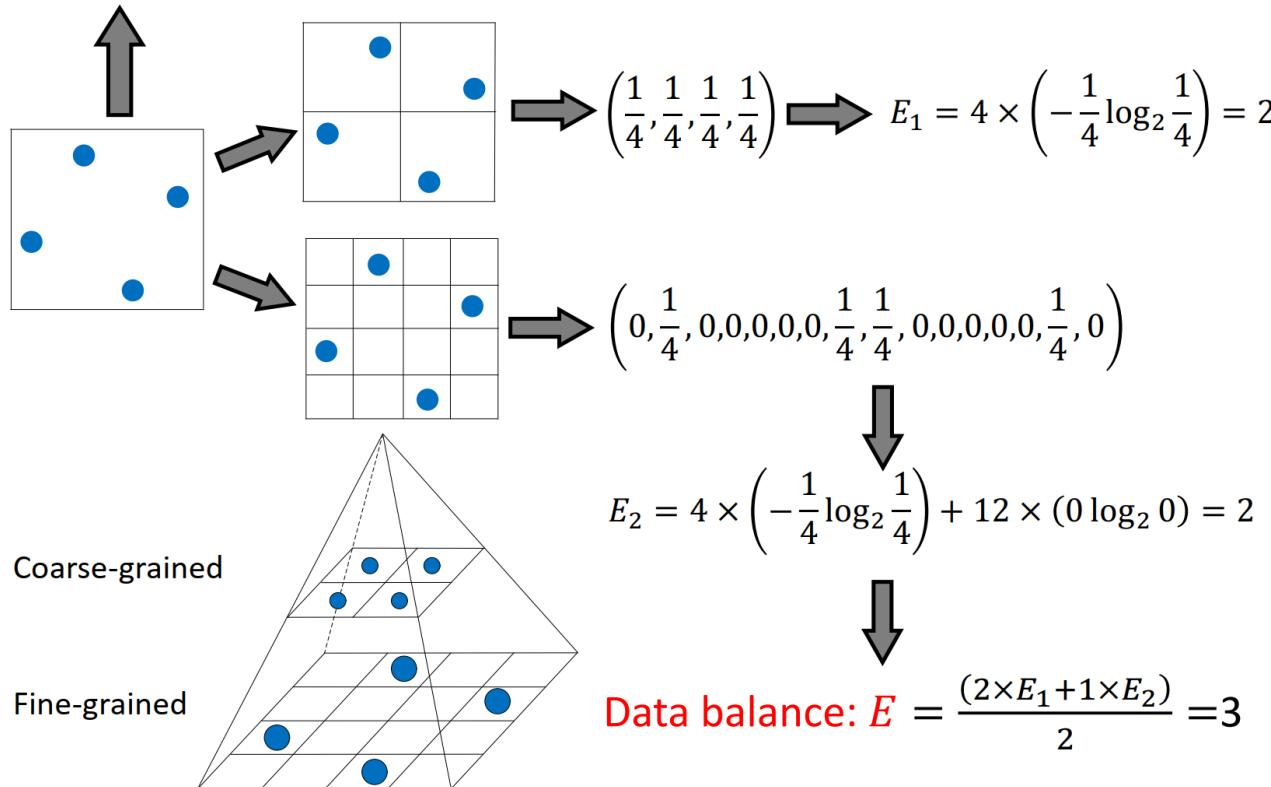
- application specific



# Hierarchical Entropy-based Objective Function



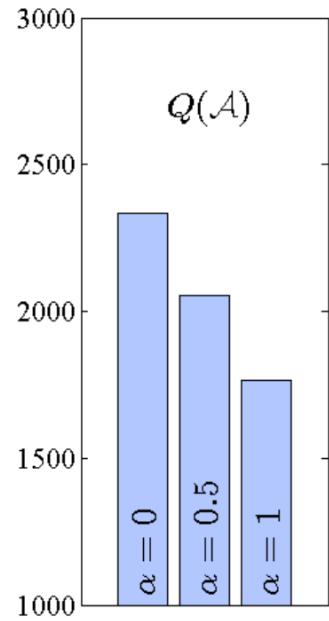
Data amount:  $Q = 4$



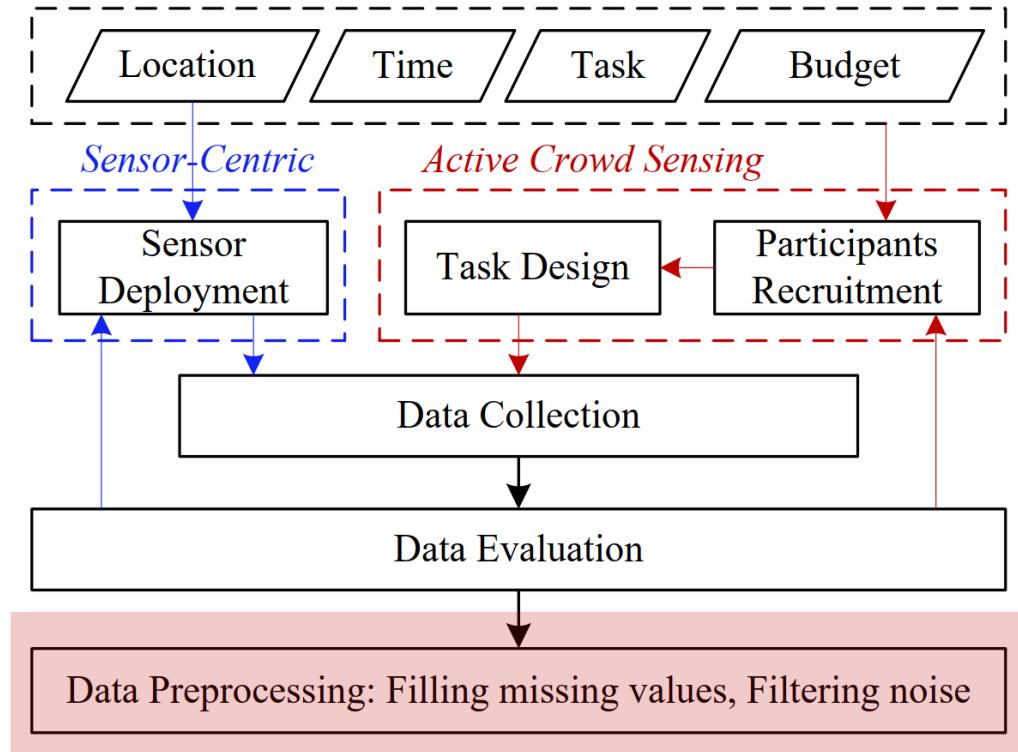
# Real-World Evaluation



- Datasets
  - Human mobility dataset from a real-world noise sensing experiment
    - Sensing region:  $6.6\text{km} \times 3.3\text{km}$
    - Sensing time interval: 6:00 am ~ 22:00 pm
  - 244 participant candidates with mobility information



# General Framework of Urban Sensing





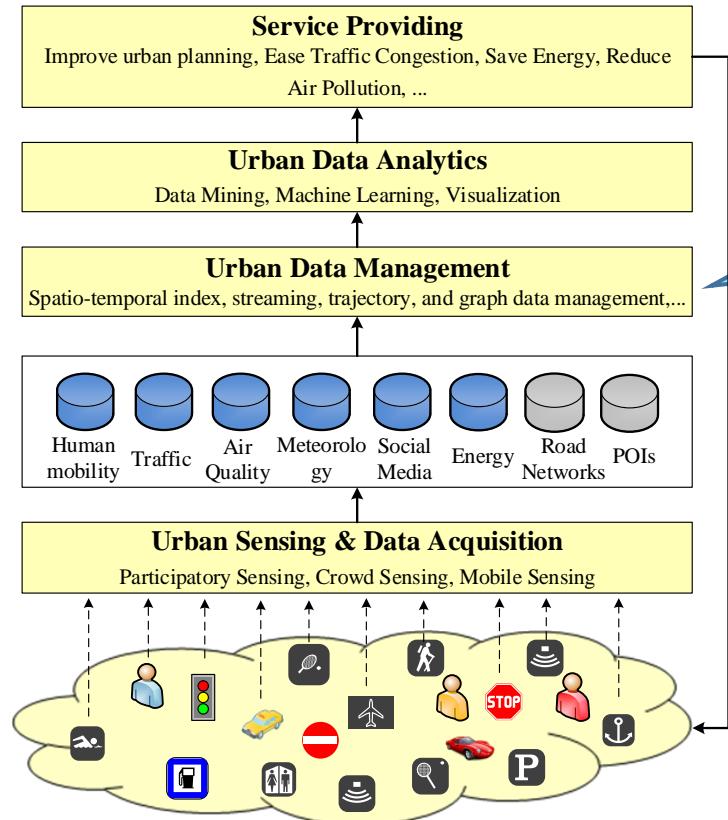
# Objectives of this Course

To introduce

- Basic data management techniques
  - Spatial databases
  - Spatio-temporal databases
- Advanced applications
  - Map matching
  - Trajectory compression
  - Large-scale dynamic ridesharing
- Managing spatio-temporal data on the cloud

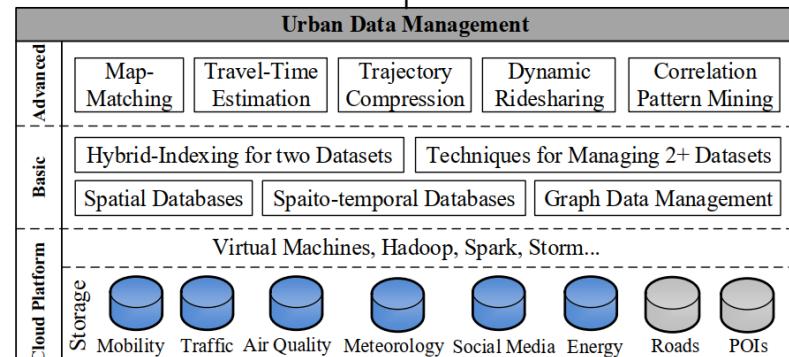


# 2<sup>nd</sup> Stage: Urban Data Management



Manage cross-domain spatio-temporal data using

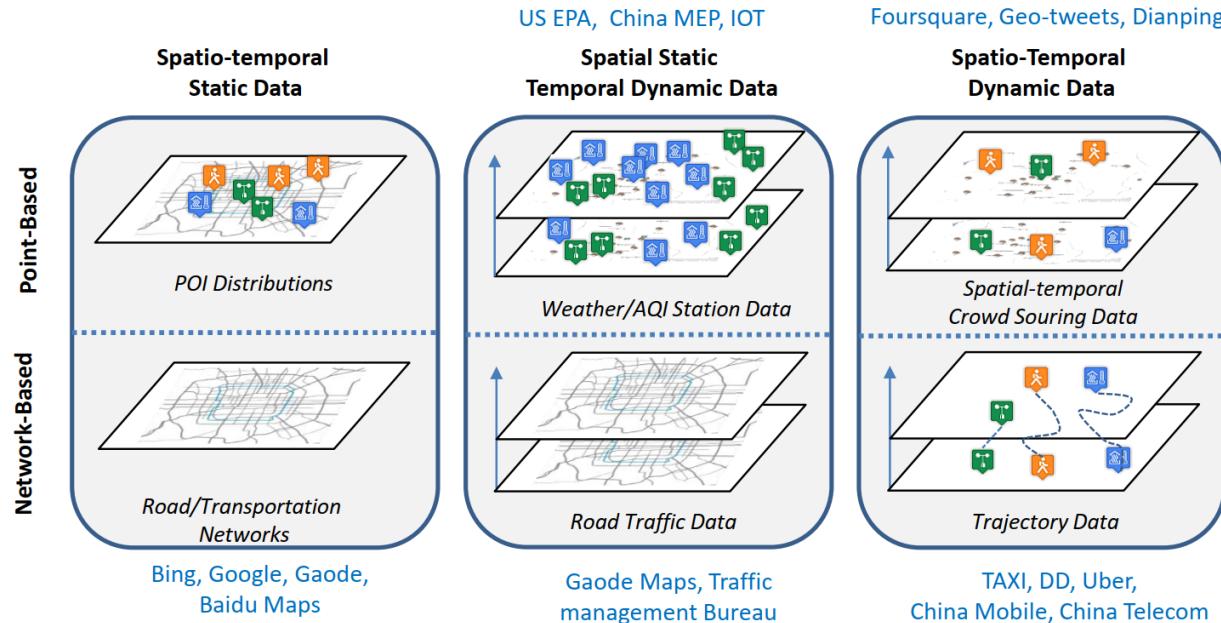
- Data: ST properties
- Indexing and retrieval algorithms
- Cloud computing platforms





# Spatio-Temporal Big Data

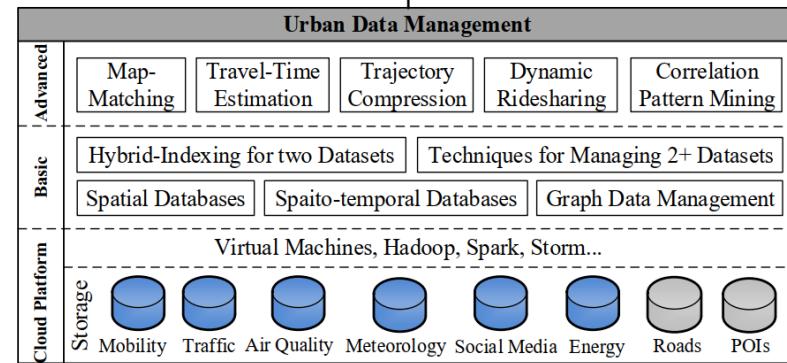
- Data Structures
- Spatio-Temporal (ST) Properties



# Spatio-Temporal Data Management



- Basic data management techniques
  - Spatial databases
  - Spatio-temporal databases
- Advanced applications
  - Map matching
  - Trajectory compression
  - Large-scale dynamic ridesharing
  - Road network recovery
- Managing spatio-temporal data on the cloud



# ST Data Management

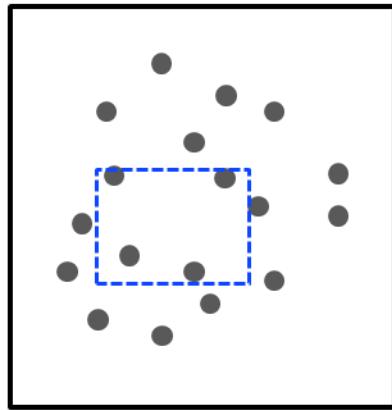


- Basic techniques
  - Part 1: Spatial databases
  - Part 2: Spatio-temporal databases

# Spatial Queries

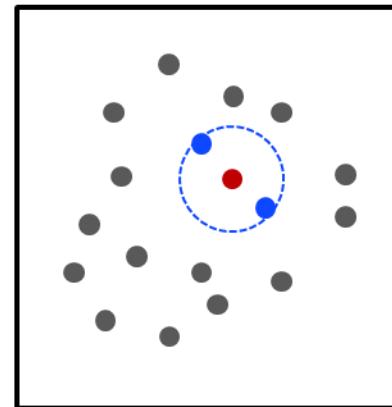


Region (Range) Query



Ask for objects that lie partially or fully inside a specified region.

Nearest Neighbour Queries



Given a point or an object, find the nearest object that satisfies given conditions

# Spatial Indexing Structures

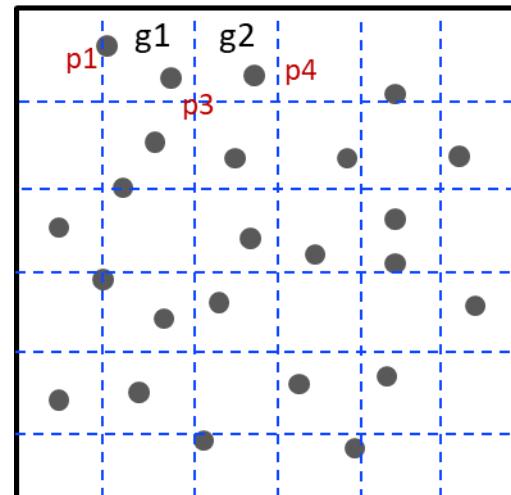
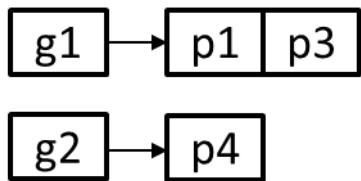


- Space Partition-Based Indexing Structures
  - Grid-based
  - Quad-tree
  - K-D tree
- Data-Driven Indexing Structures
  - R-Tree



# Grid-based Spatial Indexing

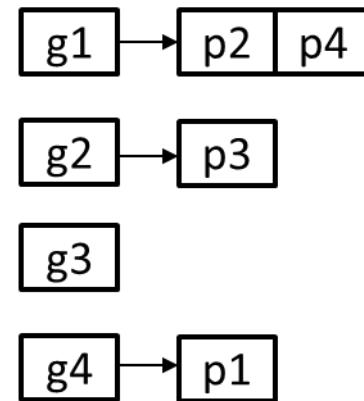
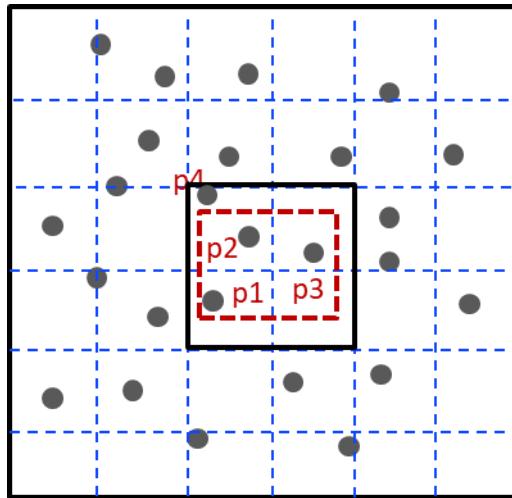
- Indexing
  - Partition the space into disjoint and uniform grids
  - Build inverted index between each grid and the points in the grid





# Grid-based Spatial Indexing

- Range Query
  - Find the grids intersecting the range query
  - Retrieve the points from the grids and identify the points in the range

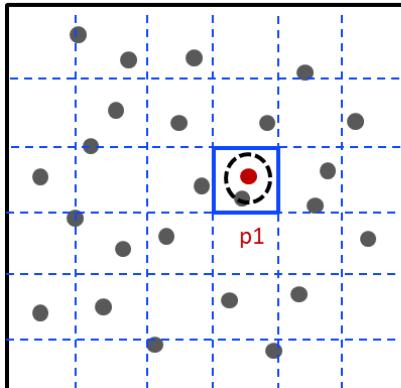




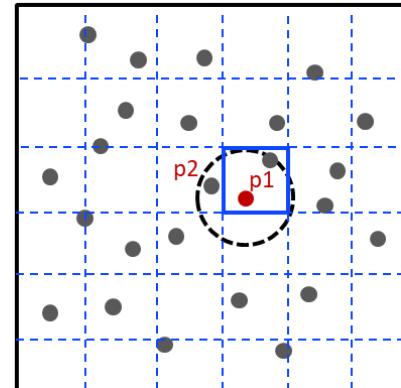
# Grid-based Spatial Indexing

- Nearest neighbor query
  - Euclidian distance
  - Road network distance is quite different

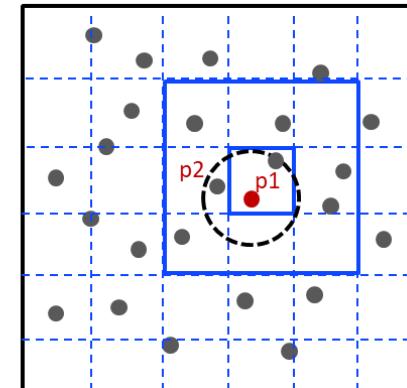
The nearest object is  
within the grid



The nearest object is  
outside the grid



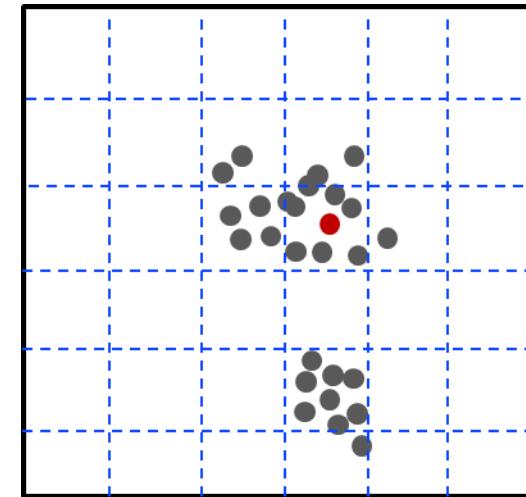
Fast approximation





# Grid-based Spatial Indexing

- Advantages
  - Easy to implement and understand
  - Very efficient for processing range and nearest queries
- Disadvantages
  - Index size could be big
  - Difficult to deal with unbalanced data

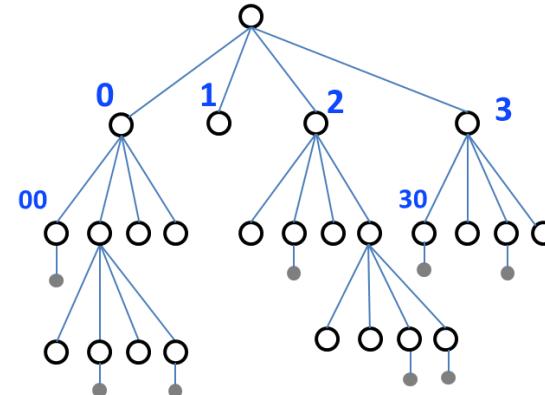
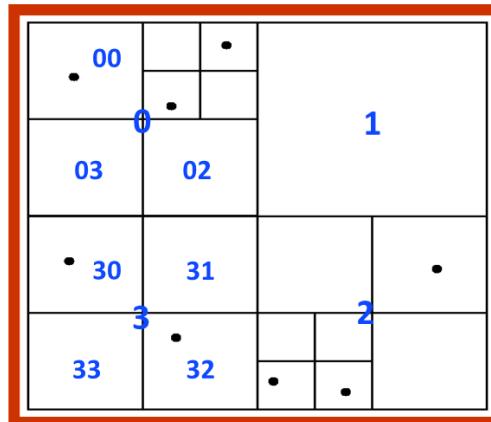




# Quad-Tree

- Indexing

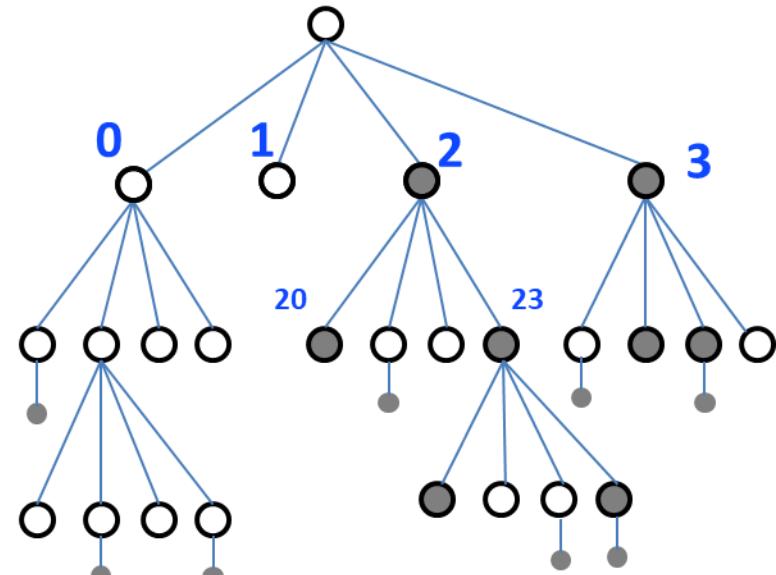
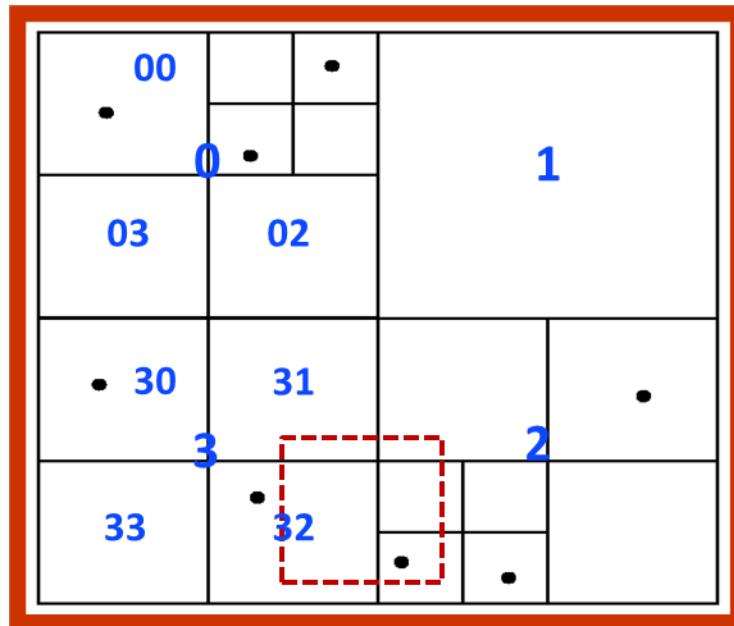
- Each node of a quad-tree is associated with a rectangular region of space; the top node is associated with the entire target space
- Each non-leaf node divides its region into four equal sized quadrants
- Leaf nodes have between zero and some fixed maximum number of points (set to 1 in example)



# Quad-Tree



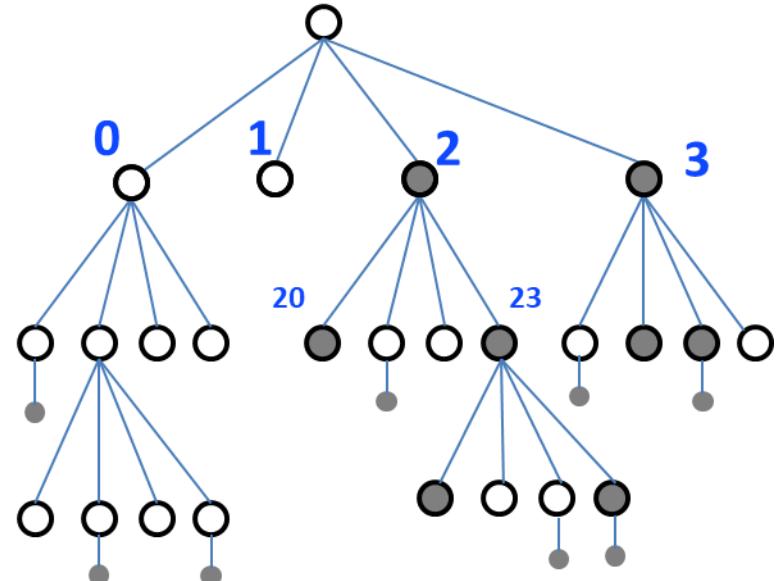
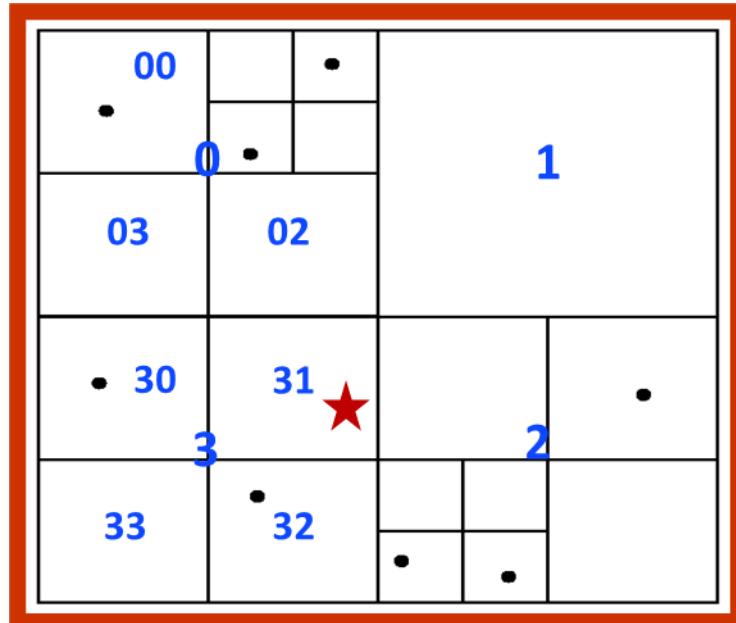
- Range query



# Quad-Tree



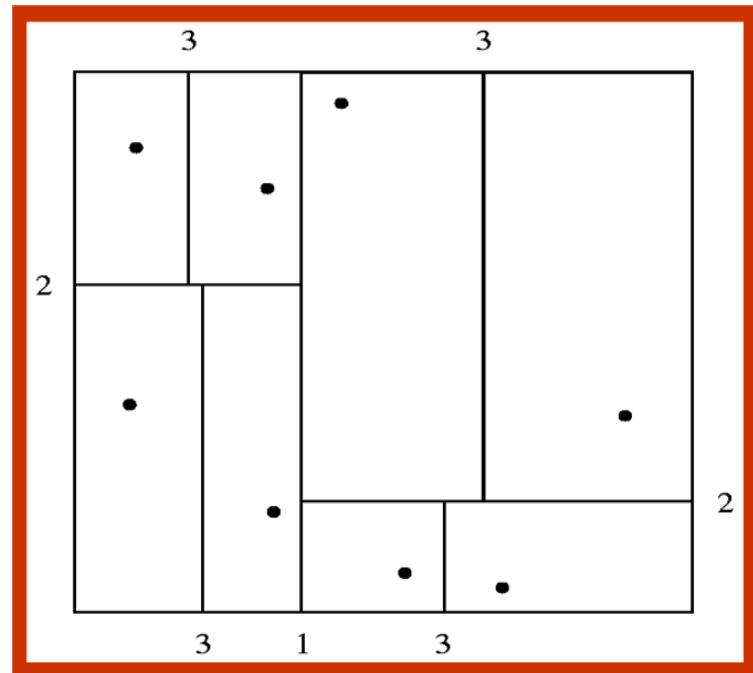
- Nearest Neighbor Query (hard)



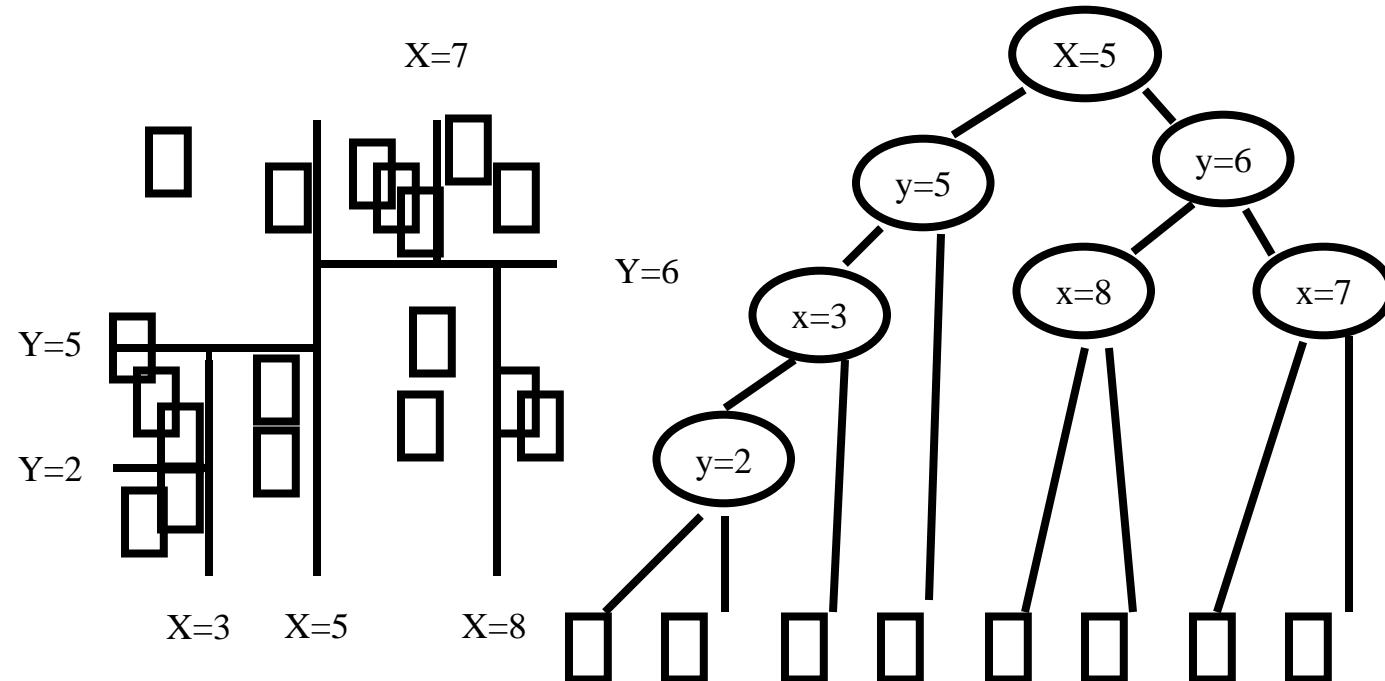


# K-D-Tree

- Each line in the figure (other than the outside box) corresponds to a node in the k-d tree
  - the maximum number of points in a leaf node has been set to 1
- The numbering of the lines in the figure indicates the level of the tree at which the corresponding node appears



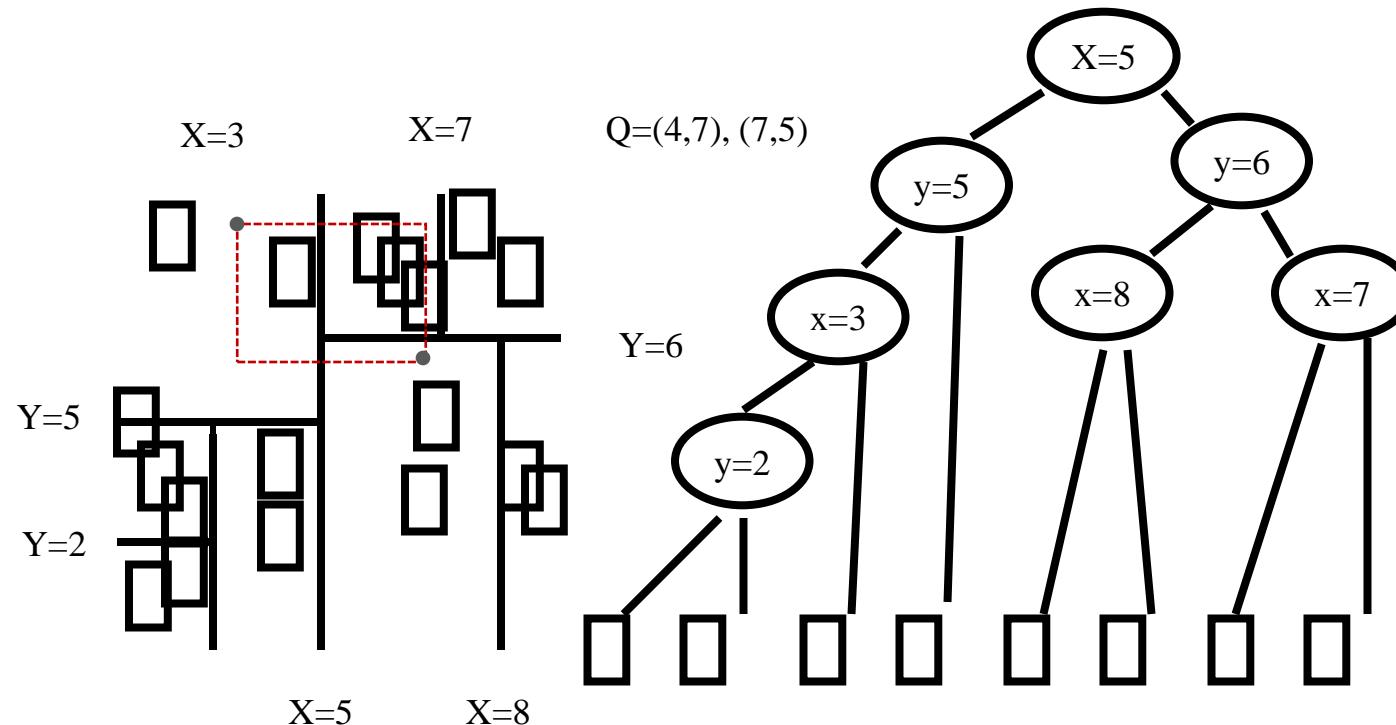
# K-D-Tree Example



# K-D-Tree Example



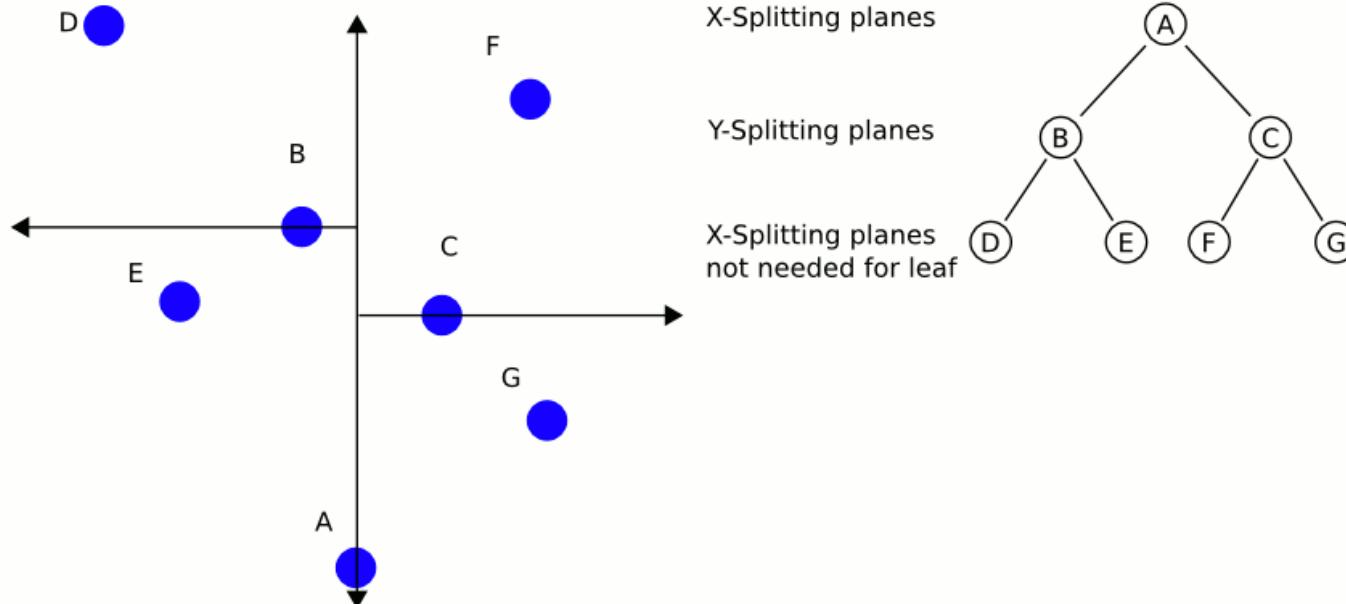
- Range query



# K-D-Tree



- Nearest neighbor query



# Spatial Indexing Structures

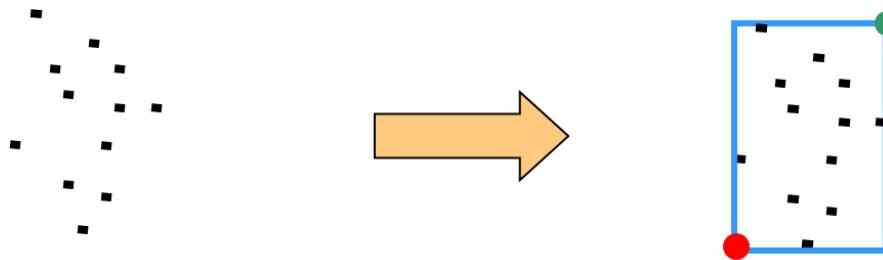


- Space Partition-Based Indexing Structures
  - Grid-based
  - Quad-tree
  - k-D tree
- Data-Driven Indexing Structures
  - R-Tree

# R-Trees



- Build a Minimum Bounding Rectangle (MBR)



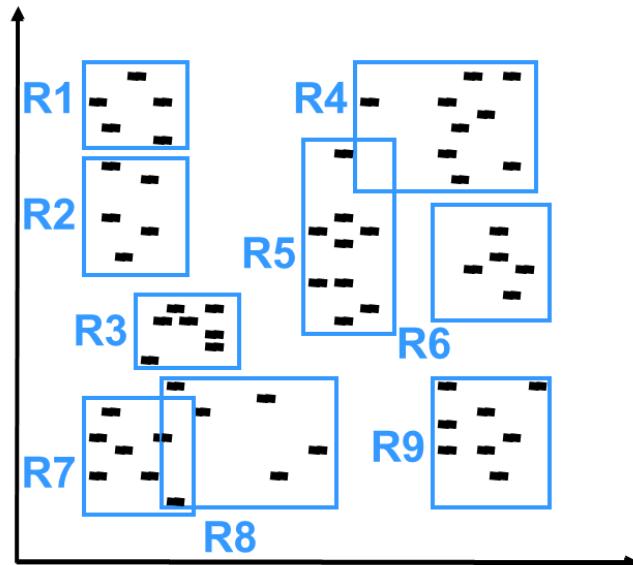
$$\text{MBR} = \{(L.x, L.y), (U.x, U.y)\}$$

Note that we only need two points to describe an MBR, we typically use lower left, and upper right.

# R-Trees



- We can group clusters of data points into MBRs
  - Can also handle line-segments, rectangles, polygons, in addition to points

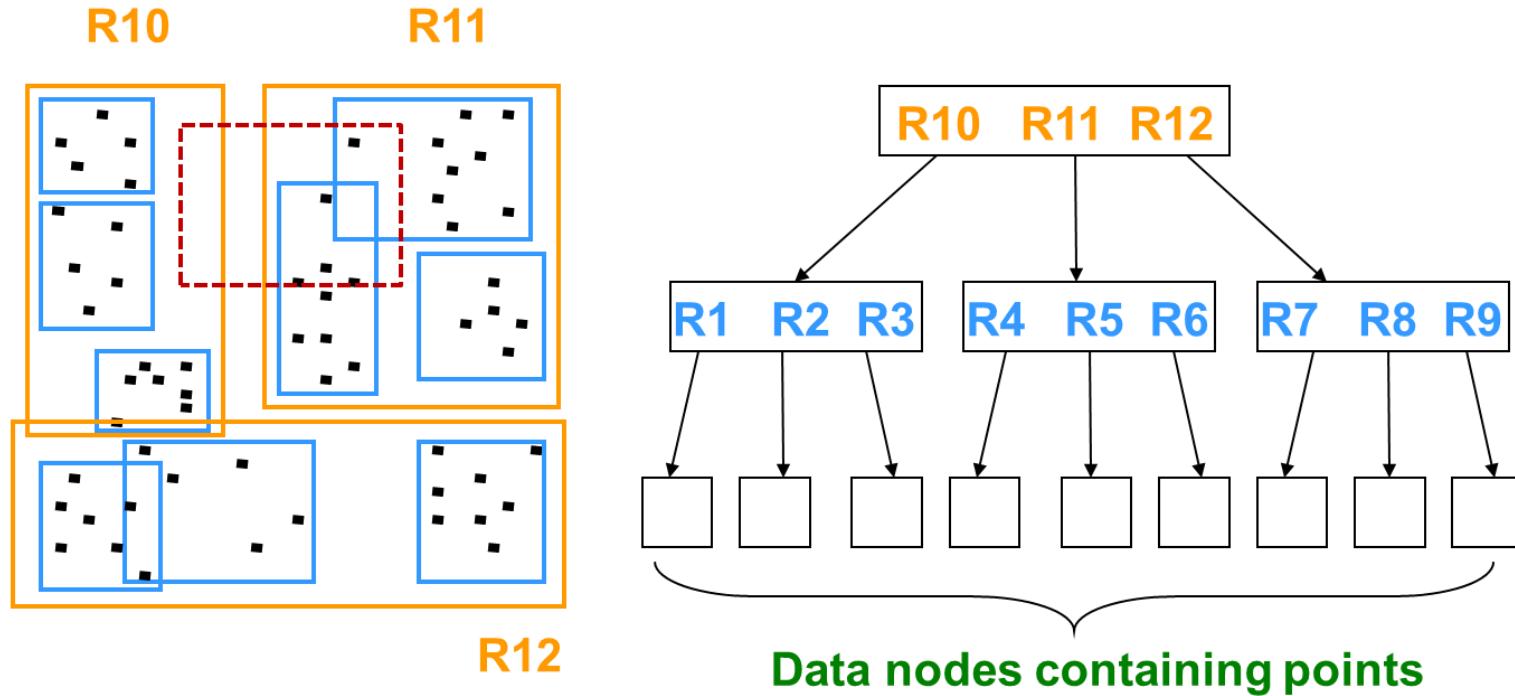


We can further recursively group MBRs into larger MBRs....

# R-Tree Structure



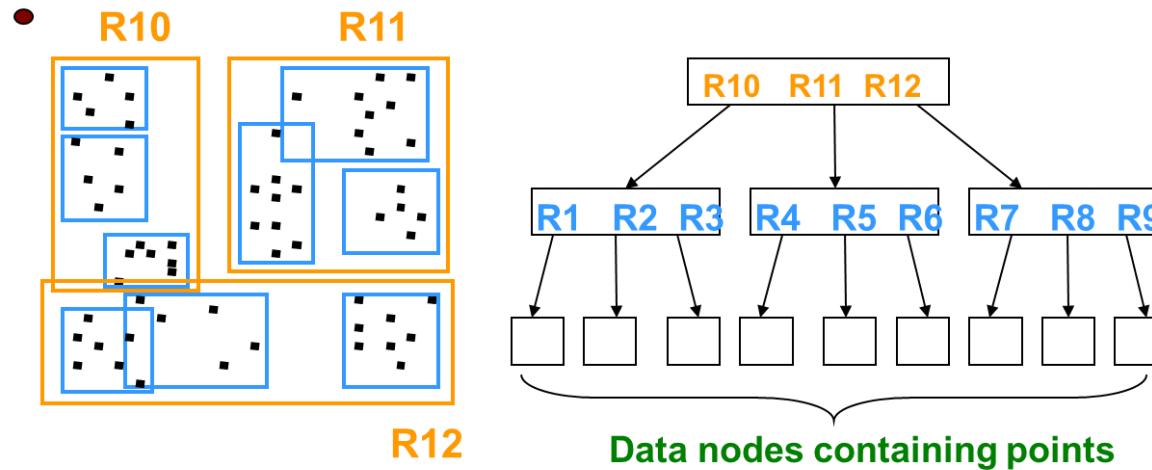
- Nested MBRs are organized as a tree





# Nearest Neighbour Search

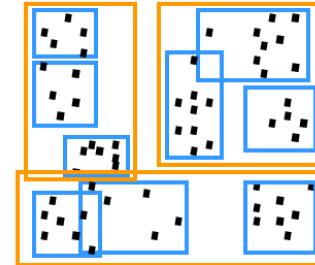
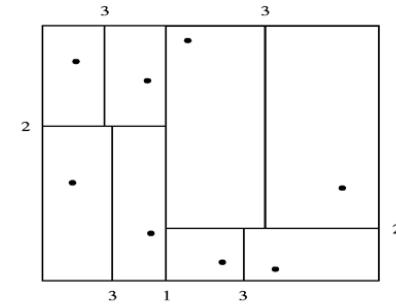
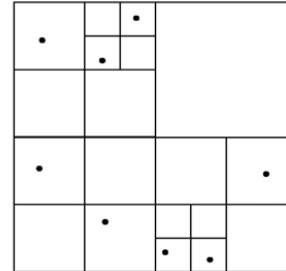
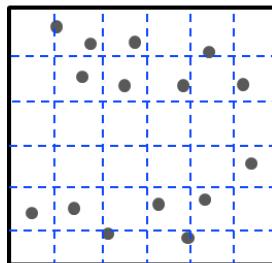
- Given an MBR, we can compute lower bounds on nearest object
- Once we know there is an item within some distance  $d$ , we can prune away all items/MBRs at distance  $> d$ 
  - Even if we haven't actually found the nearest item yet
  - Similar technique possible for k-d trees and quad-trees as well





# Comparison among Spatial Indices

	Unbalanced data	Range query	Nearest neighbor	Construction	Balanced structure	Storage
Grid-based	Poor	Good	Normal	Easy	Yes	Big
Quad-Tree	Good	Best	Poor	Easy	No	Median
KD-Tree	Good	Normal	Good	Easy	Almost	Median
R-Tree	Good	Normal	Best	Difficult	Yes	Small



# ST Data Management



- Basic techniques
  - Part 1: Spatial databases
  - Part 2: Spatio-temporal databases



# Part 2: Spatio-Temporal Databases

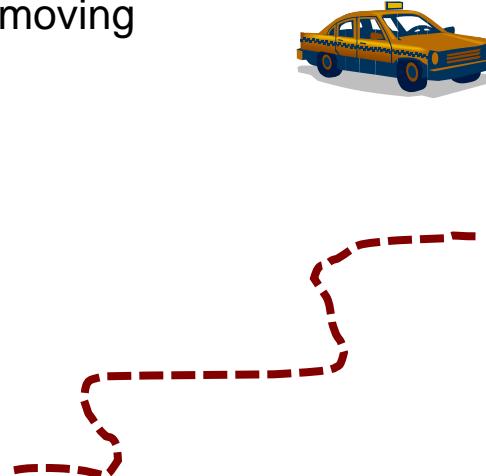
- Many moving objects
- Current status:  $(x, y, t)$
- Past: a trajectory  $(x_1, y_1, t_1), (x_2, y_2, t_2), \dots,$



# Spatio-Temporal Databases



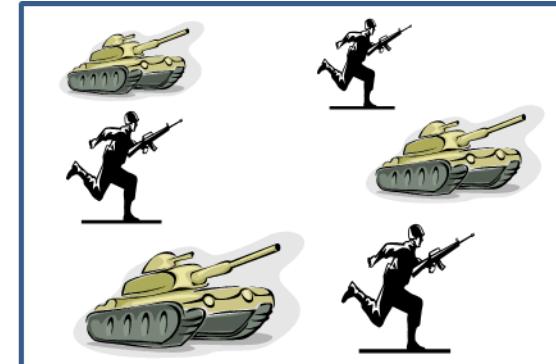
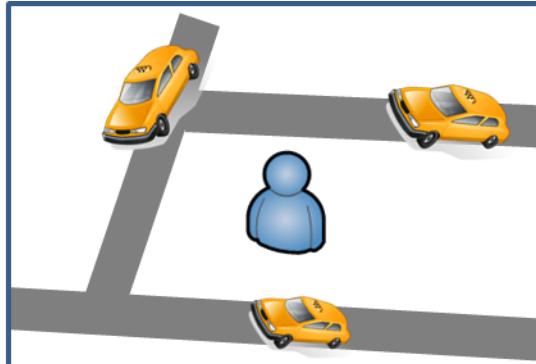
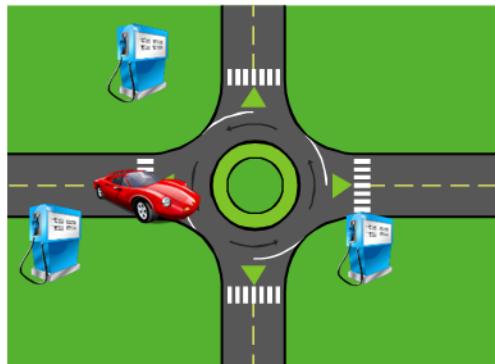
- Querying current status of a moving object
  - Scenarios:
    - Query point is moving but objects fixed
    - Query point is stationary but query objects are moving
    - Both query points and objects are moving
  - Type of queries: KNN and Range queries
- Query the past of a of a moving object
  - Trajectory management
  - KNN and Range Queries





# Querying current status of a moving object

- Examples:
  - Show me the closest gas station when I am driving
  - Tell me the top-2 closest vacant taxis that would pass me
  - How many moving soldiers are around a moving tank?



# Querying Current Status of a Moving Object

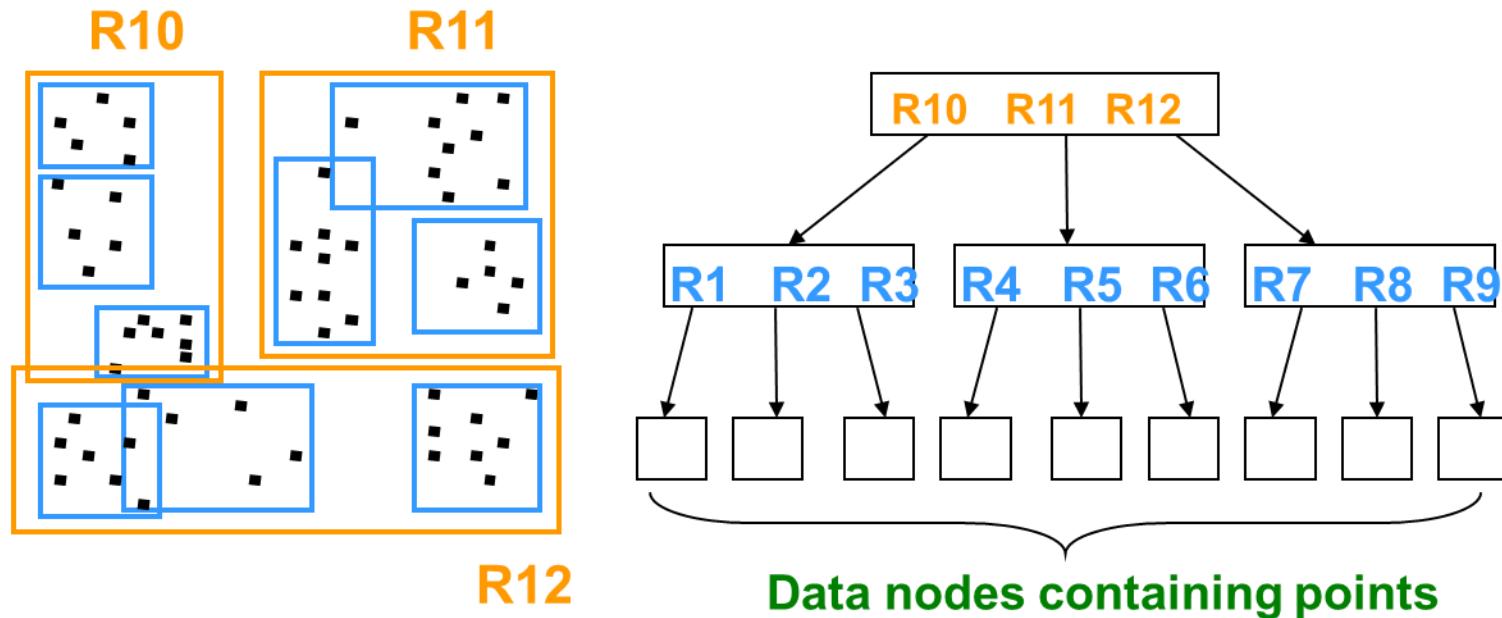


- KNN and Range queries
- Instantaneous and continuous
- Indexing structures
  - View temporal as an additional dimension
    - 3D R-Tree
    - ST R-Tree
    - TB-Tree
  - Multiple version-based
    - HR-tree
    - MR-tree
    - HR+-tree
    - MV3R-tree



# Index for the Query of Moving Objects

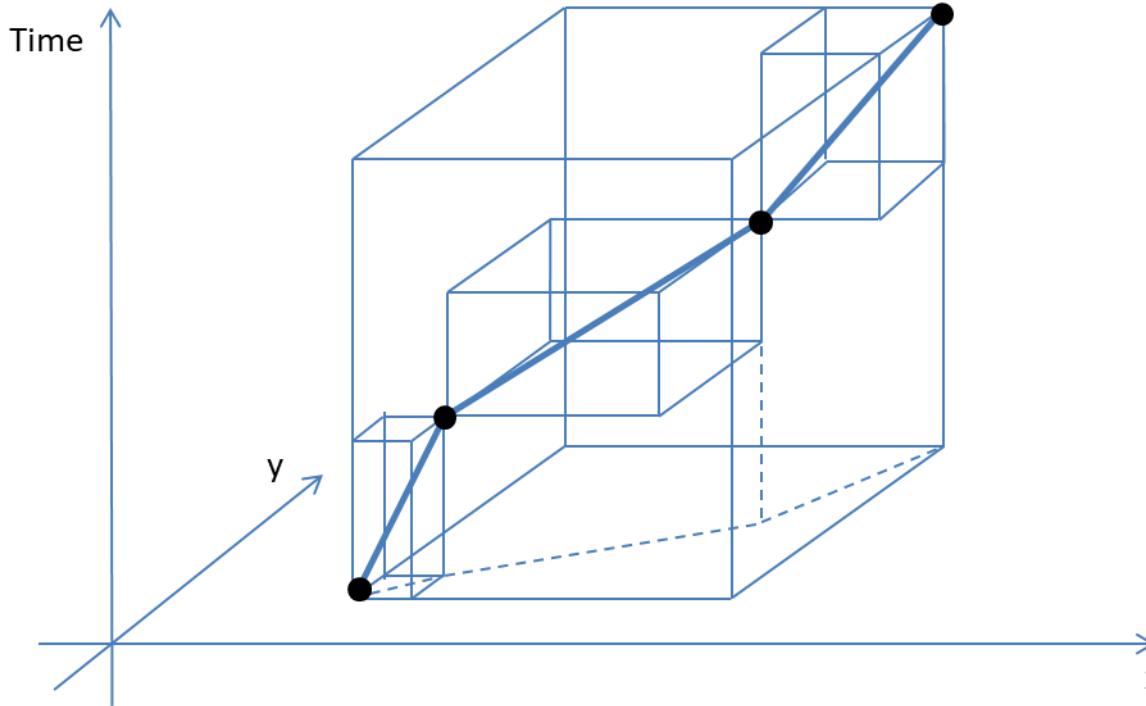
- R-Tree



# Index for the Query of Moving Objects



- 3D R-tree

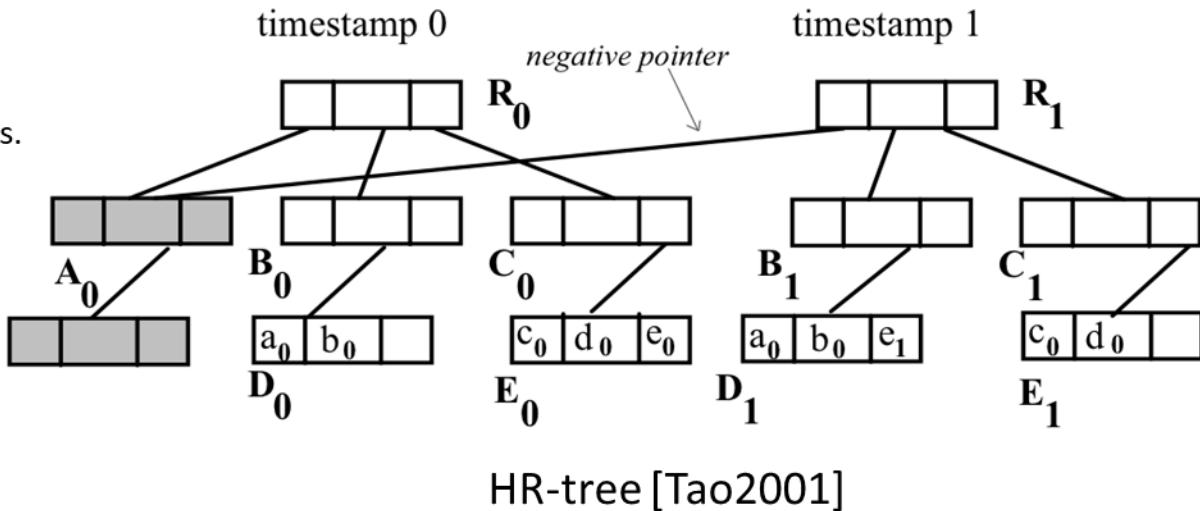


# Index for the Query of Moving Objects



- Multi-version R-tree (HR-tree [Tao2001a], HR+-tree[Tao2001b], MR-tree[Xu2005])

For each timestamp, an R-tree is created. So, there are many R-trees. These R-trees are indexed.



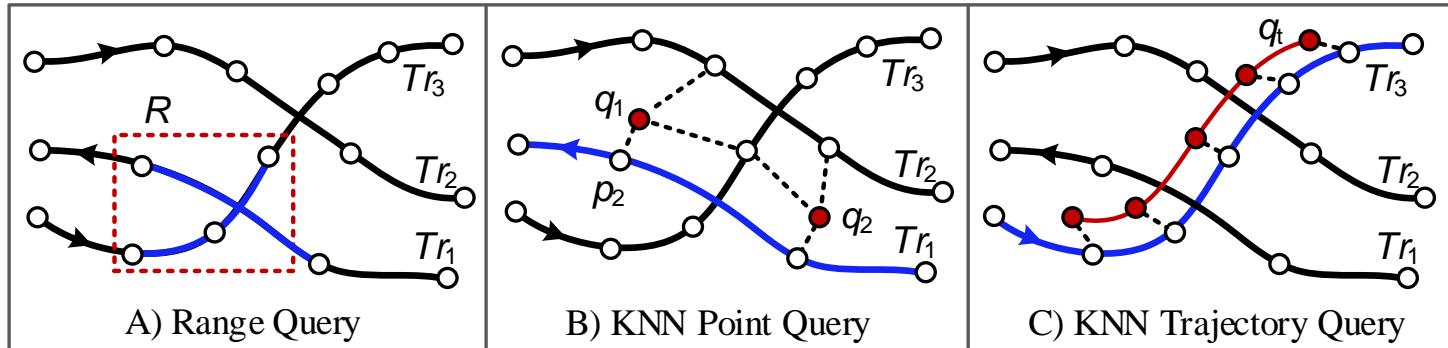
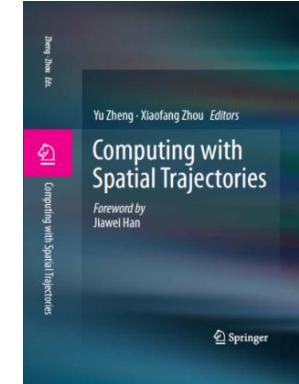
Query for trajectories in a given region and in a given time interval:

1. The R-tree at the timestamp is found first
  2. The trajectories in the specified region are retrieved from the R-tree.



# Trajectory Data Management

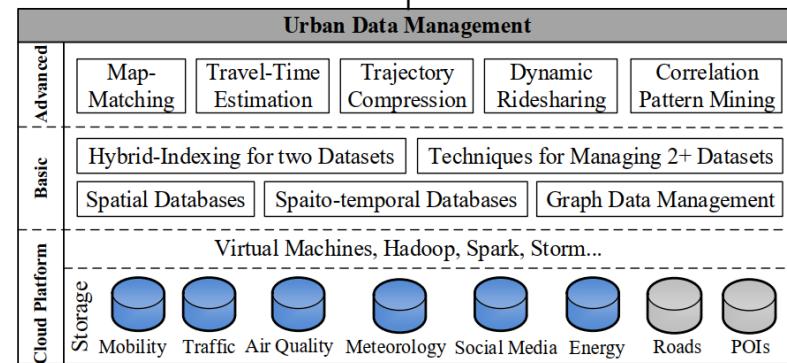
- (ST) Range queries
- KNN queries
  - Point KNN queries
  - Trajectory KDD queries



# Urban Data Management



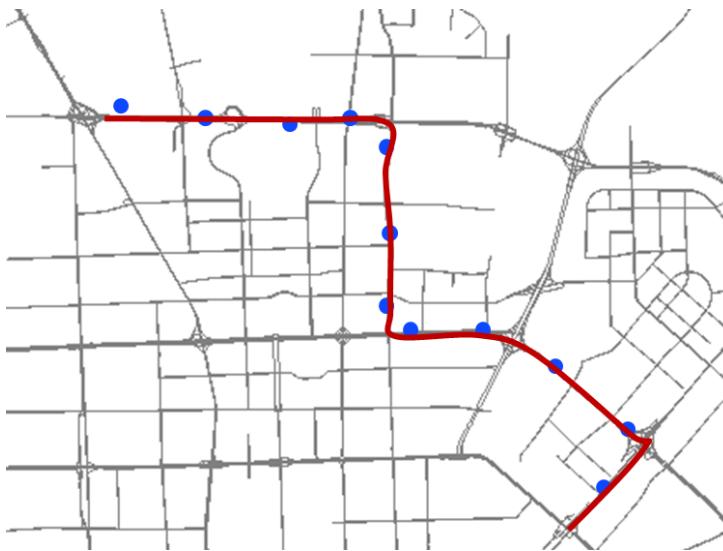
- Basic Data management Techniques
  - Spatial Databases
  - Spatio-Temporal Databases
- Advanced applications
  - Map Matching
  - Trajectory compression
  - Large-Scale Dynamic Ridesharing
- Managing Spatio-Temporal Data on the Cloud



# Map-matching



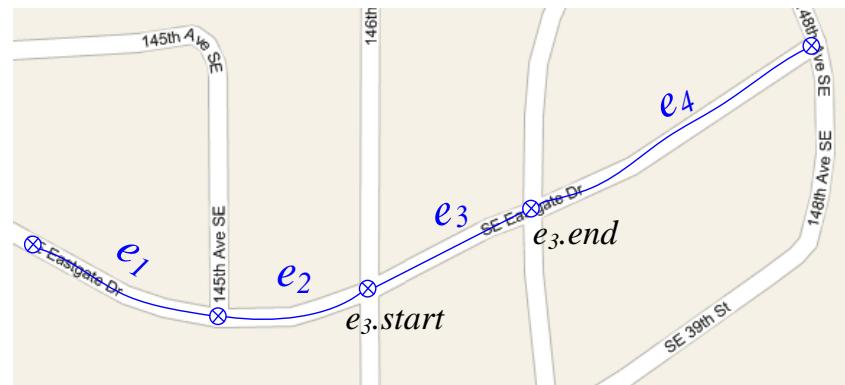
- Problem
  - Map a GPS trajectory onto a road network
  - a sequence of GPS points → a sequence of road segments





# Spatial Data

- Road network:  $G=(V, E)$ 
  - $V$  is a set of nodes
  - $E$  is a set of road segments
  - $e \in E$ , consists of two terminal nodes and a sequence of intermediate points describing the segment with a polyline
  - Properties:  $e.\text{len}$ ,  $e.\text{dir}$ ,  $e.\text{lanes}$



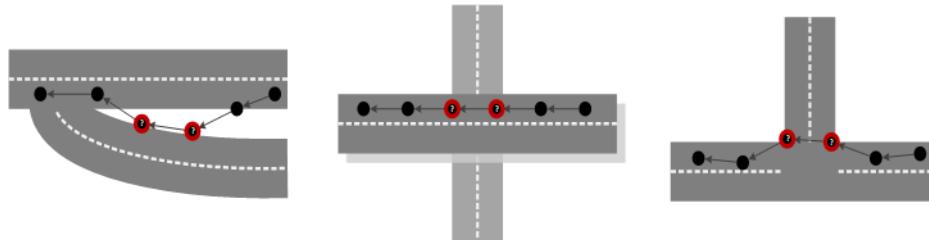
# Map-Matching



- Why it is important
  - A fundamental step in many transportation applications
    - Navigation and driving
    - Traffic analysis
    - Taxi dispatching and recommendations
  - Examples:
    - Find the vehicles passing No.1 Duxue Road
    - Calculate the average travel time from HKUST(GZ) to Dongchong
    - When will the NO. 33 bus arrive at HKUST(GZ) stop?
    - ....



# Why difficult



(a) Parallel roads

b) Overpass

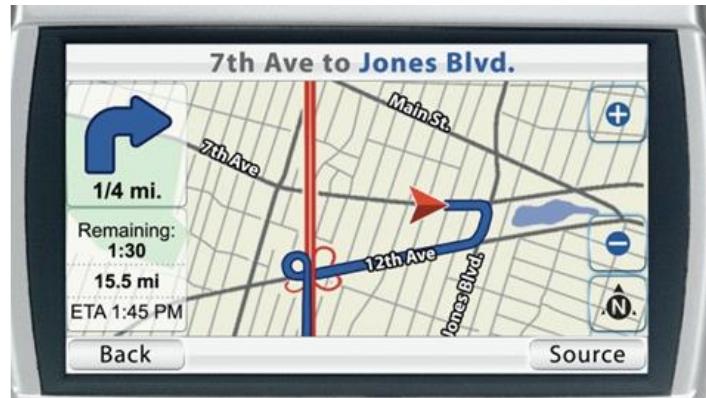
c) Spur





# Map-Matching

- Simple solution for high-sampling-rate data
  - Weighted distance



# Map-Matching

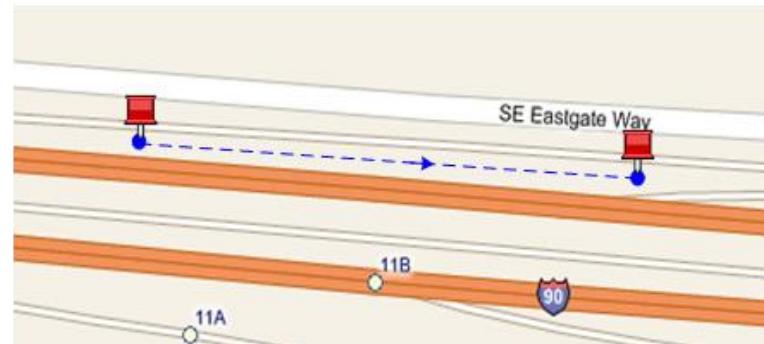
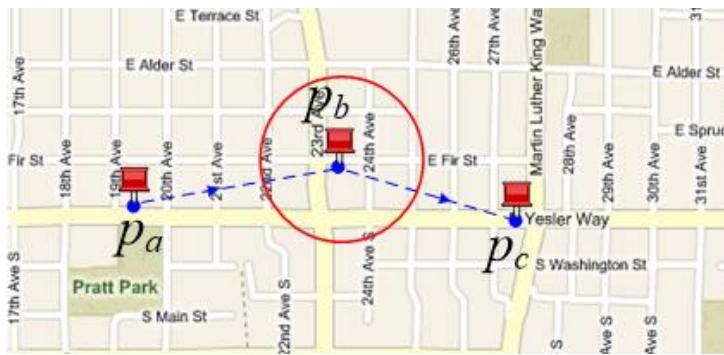


- According to the additional information used
  - Geometric
  - Topological
  - Probabilistic
  - Advanced techniques
- According to the range of sampling points
  - Local/incremental
  - Global
  - Advanced



# Map-matching

- Insights
  - Consider both local and global information
  - Incorporating both spatial and temporal features

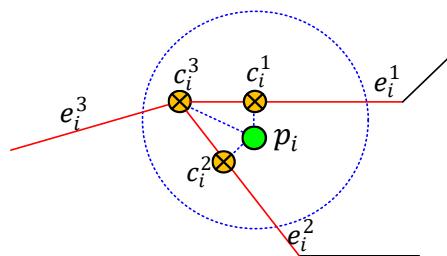




# Map-matching

- Solution (incorporating spatial information)

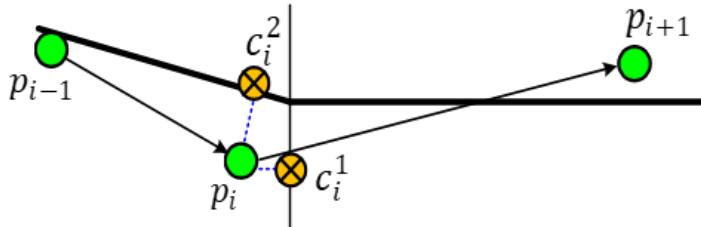
- Model local possibility



$$N(c_i^j) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i^j - \mu)^2}{2\sigma^2}}$$

where  $x_i^j = \text{dist}(c_i^j, p_i)$  is the distance between  $p_i$  and  $c_i^j$ .

- Considering context (global)



$$V(c_{i-1}^t \rightarrow c_i^s) = \frac{d_{i-1 \rightarrow i}}{w_{(i-1,t) \rightarrow (i,s)}}$$

where  $d_{i-1 \rightarrow i} = \text{dist}(p_i, p_{i-1})$  is the Euclidean distance between  $p_i$  and  $p_{i-1}$ , and  $w_{(i-1,t) \rightarrow (i,s)}$  is the length of shortest path from  $c_{i-1}^t$  to  $c_i^s$ .

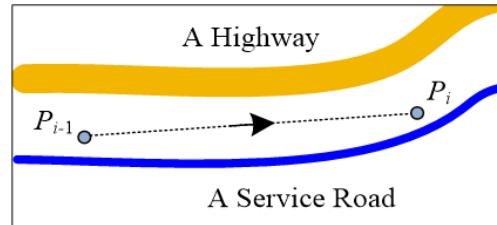
$$F_s(c_{i-1}^t \rightarrow c_i^s) = N(c_i^s) * V(c_{i-1}^t \rightarrow c_i^s), \quad 2 \leq i \leq n$$



# Map-matching

- Solution
  - Considering temporal information

$$F_t(c_{i-1}^t \rightarrow c_i^s) = \frac{\sum_{u=1}^k (e'_u \cdot v \times \bar{v}_{(i-1,t) \rightarrow (i,s)})}{\sqrt{\sum_{u=1}^k (e'_u \cdot v)^2} \times \sqrt{\sum_{u=1}^k \bar{v}_{(i-1,t) \rightarrow (i,s)}^2}}$$



$$\bar{v}_{(i-1,t) \rightarrow (i,s)} = \frac{\sum_{u=1}^k l_u}{\Delta t_{i-1 \rightarrow i}}$$



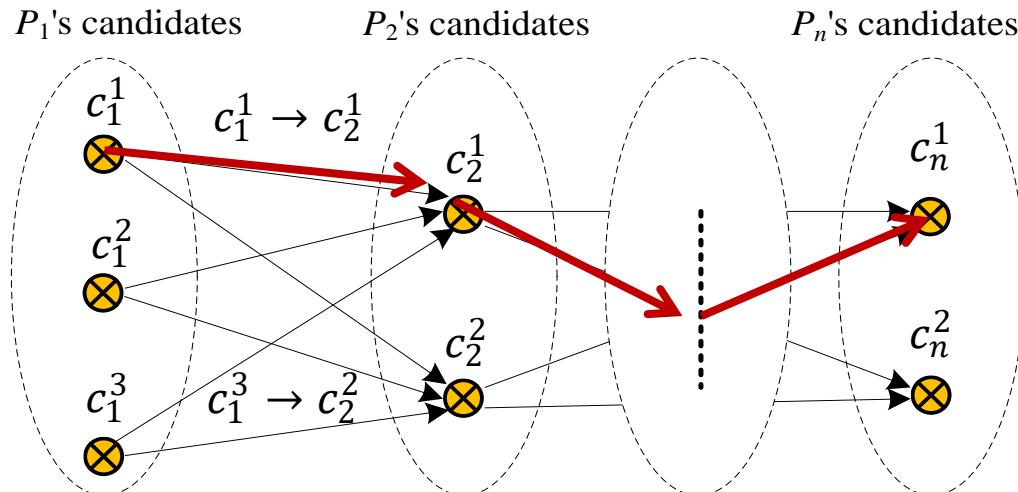
# Map-matching

## ➤ Aggregating

- Spatial and temporal information
- Local and global information

$$F(c_{i-1}^t \rightarrow c_i^s) = F_s(c_{i-1}^t \rightarrow c_i^s) * V(c_{i-1}^t \rightarrow c_i^s)$$

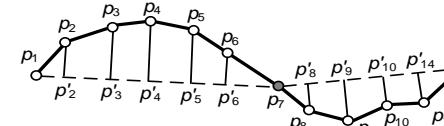
## ➤ Dynamic programming



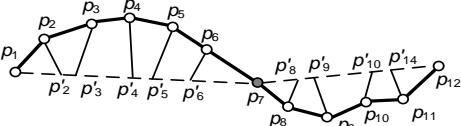


# Trajectory Compression

- Two categories of compression
  - Offline compression (a.k.a. batch mode)
    - Reduces the size of trajectory after the trajectory has been fully generated
    - Applications: content sharing sites, such as Everytrail and Bikely
  - Online compression
    - Compressing a trajectory instantly as an object travels
    - Applications: traffic monitoring and fleet management
    - Algorithms: Sliding Window, Open Window, and safe area-based methods
- Two distance metrics
  - Perpendicular Euclidean Distance
  - Time Synchronized Euclidean Distance



A) Perpendicular Euclidean Distance

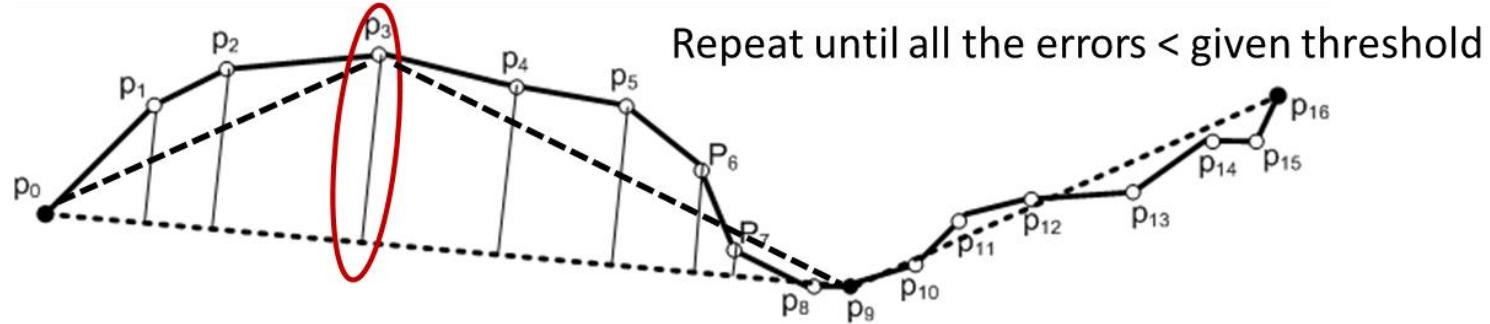
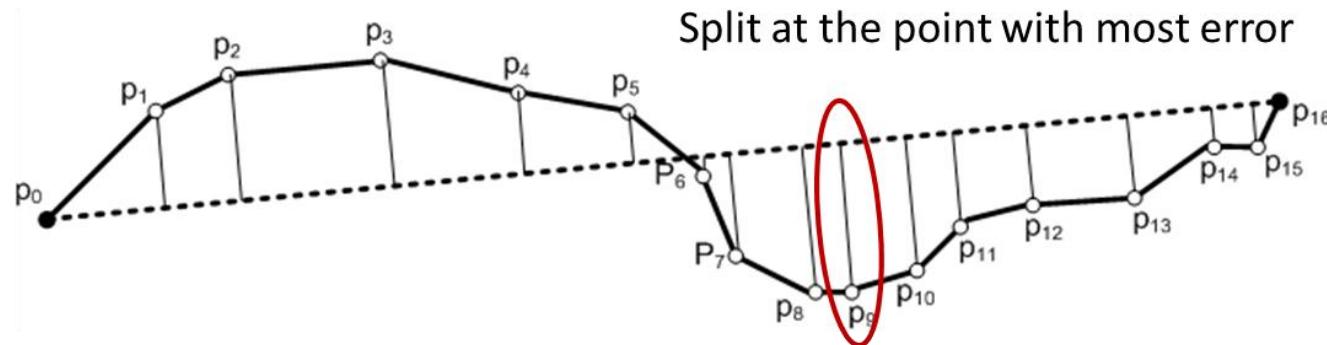


B) Time Synchronized Euclidean Distance



# Offline Trajectory Compression

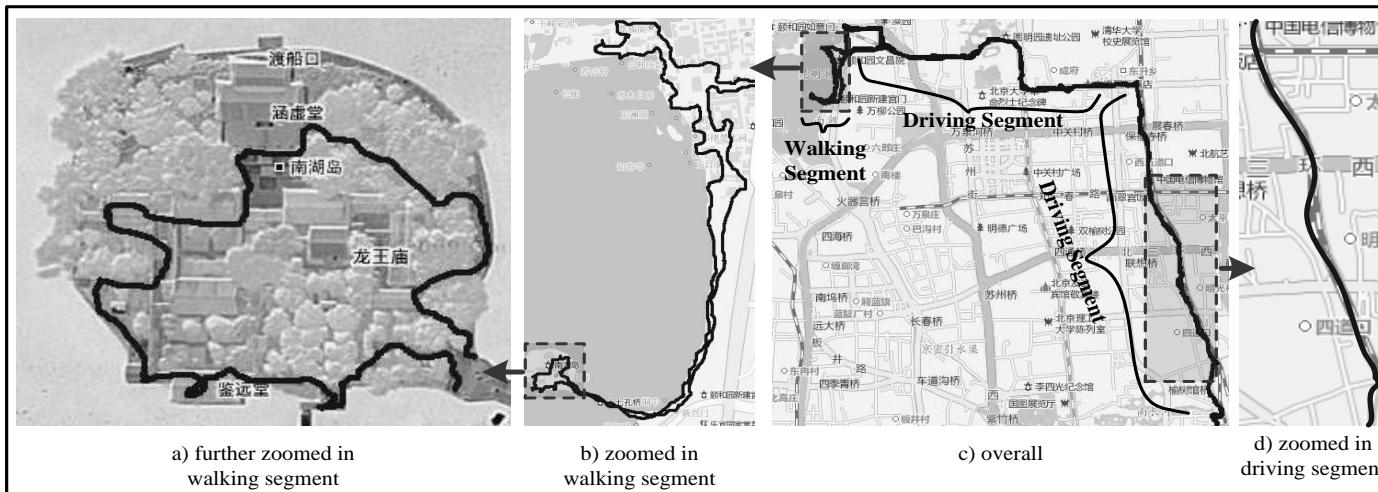
- Example



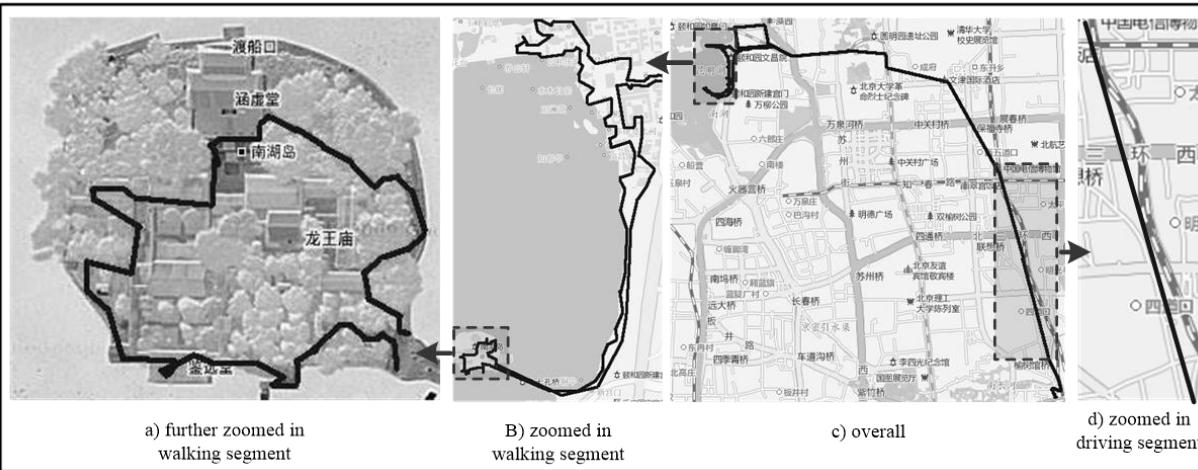
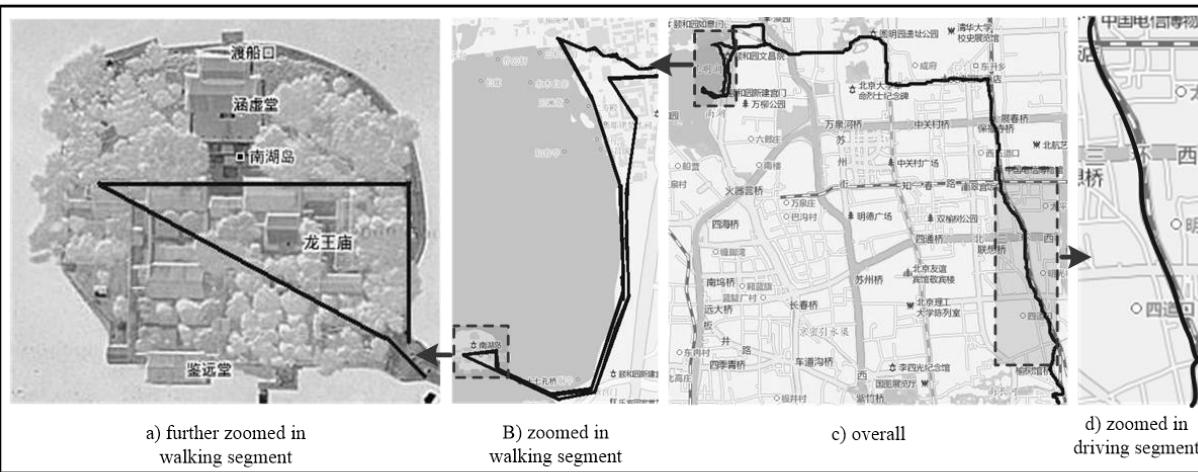
# Semantic Meaning-Based Methods (TS)



- Keep the semantic meanings of a trajectory
  - Some special points would be more significant
    - a user stayed, took photos, or changed direction greatly



# Douglas-Peucker algorithm



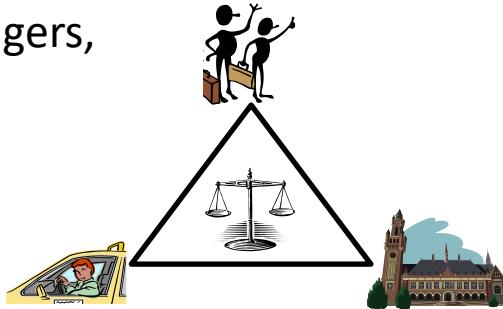
# T-Share: A Large-Scale Dynamic Taxi Ridesharing Service

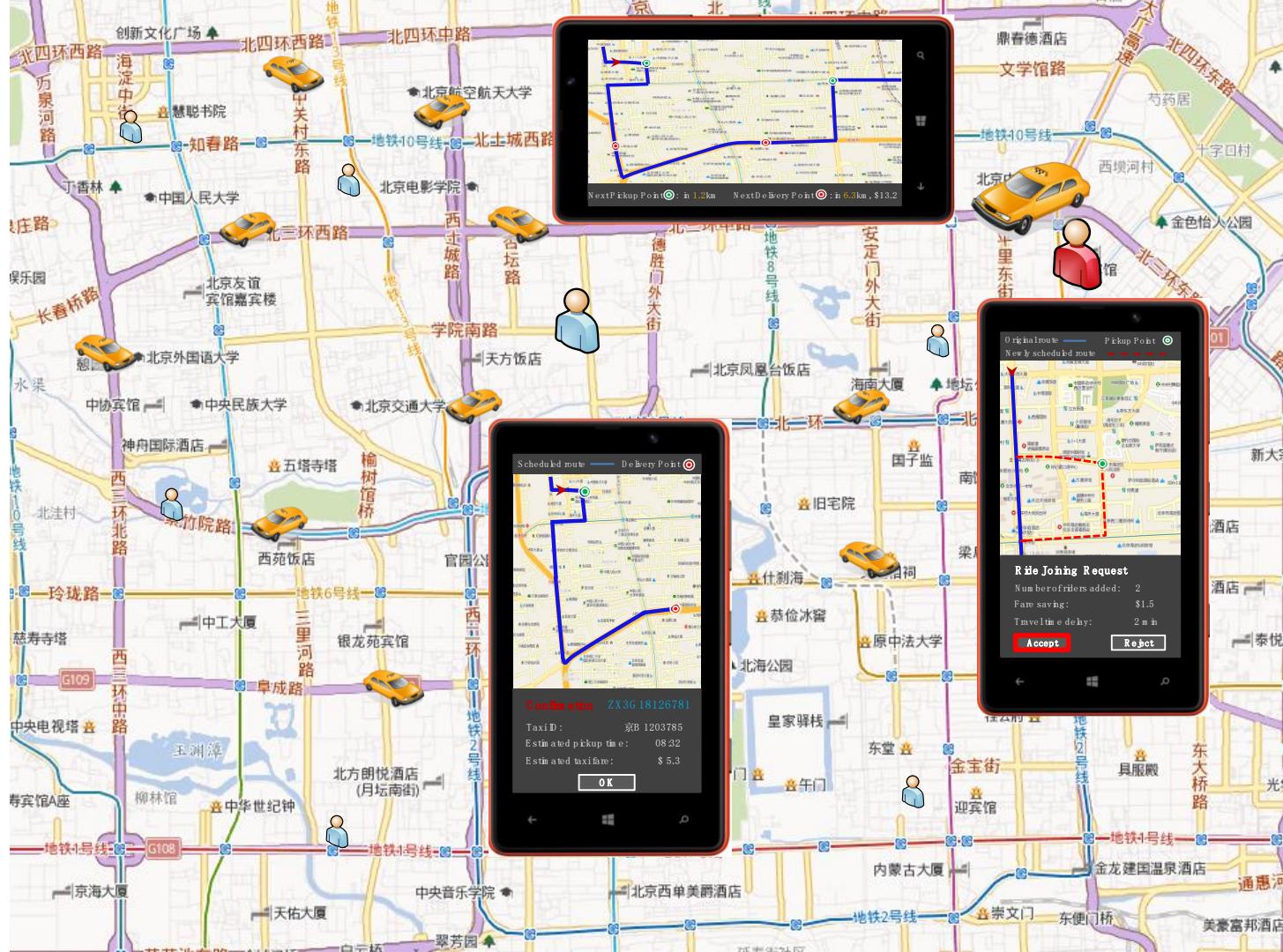
ICDE 2013 Best Paper Runner-Up Award

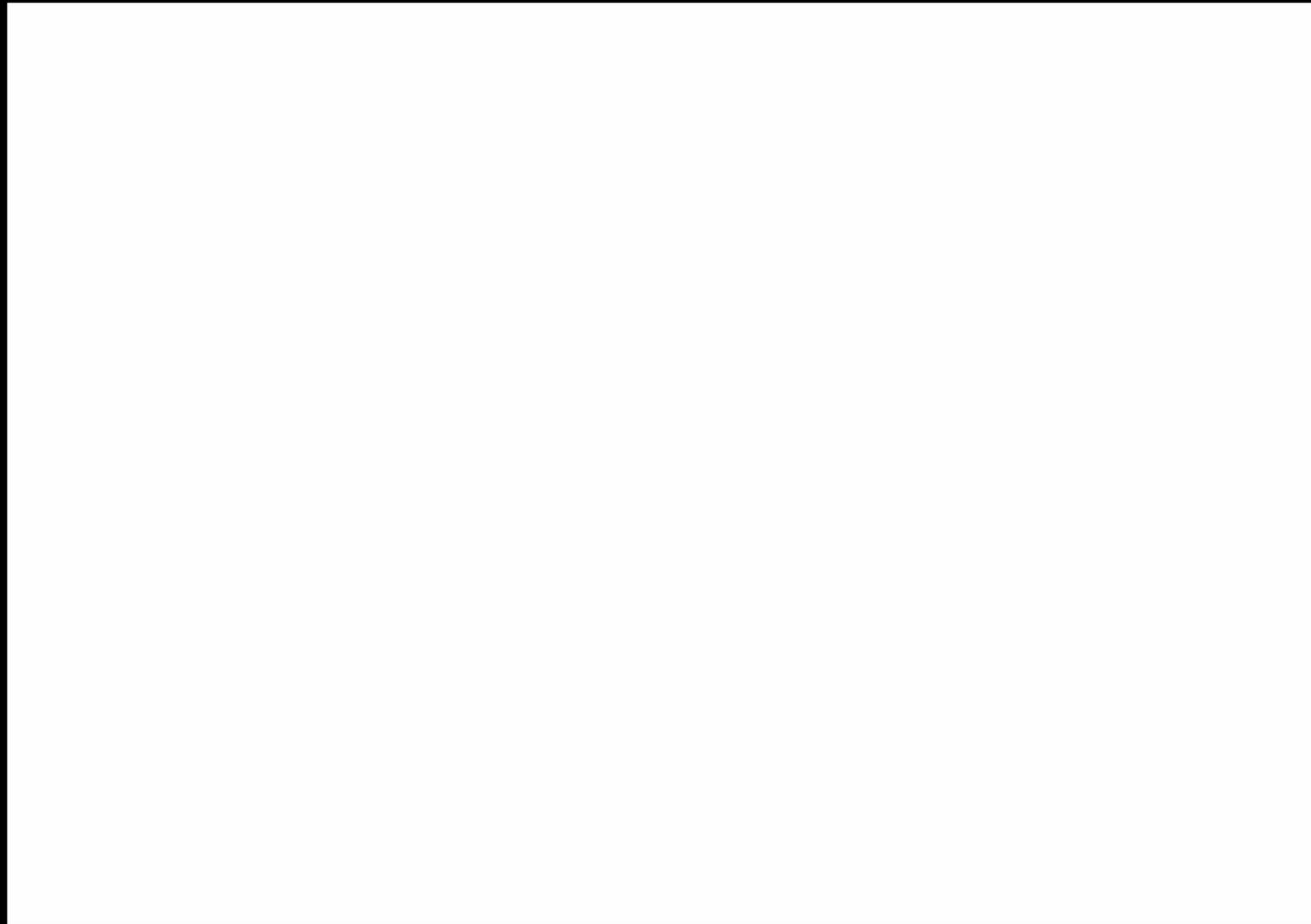


# Difficult to take a taxi!

- A social problem among passengers, taxi drivers and government
- Possible Solutions
  - Increasing taxis?
  - Taxi dispatching? ?

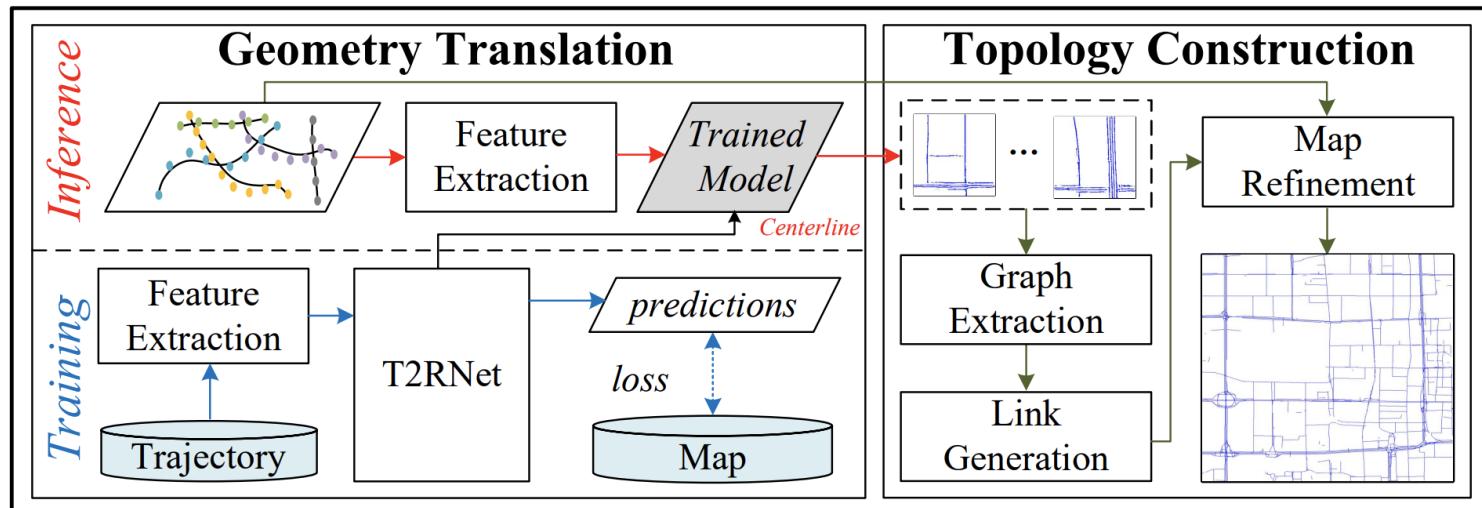




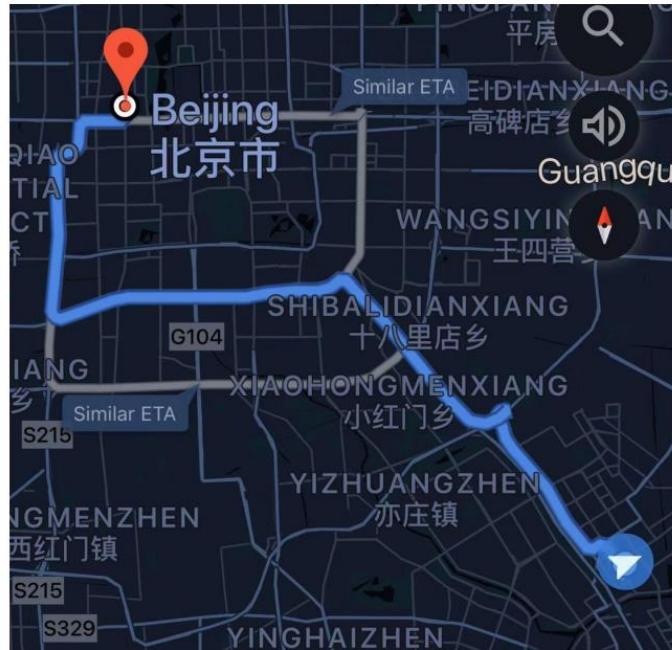


# Learning to Generate Maps from Trajectories

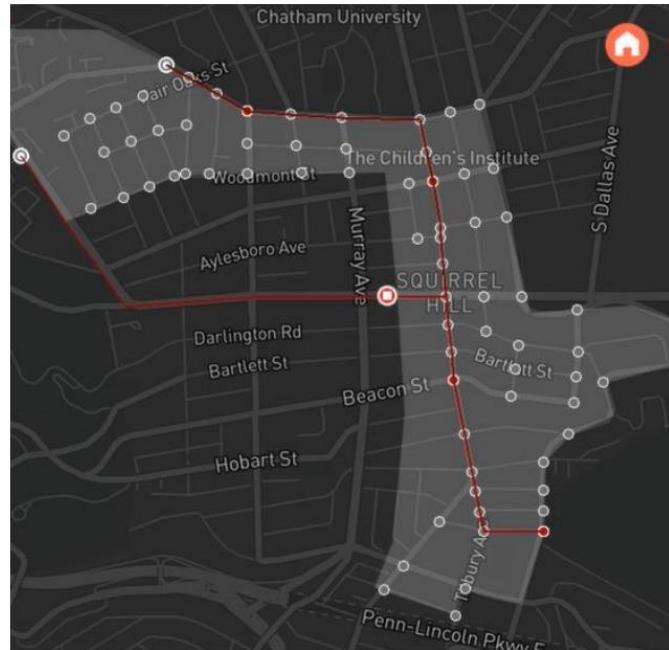
AAAI 2020



# Increasing Demands of Accurate & Updated Maps



Navigation



Route Planning & Real-time Scheduling



# Traditional Map Collection Methods



On-field Study



Low Spatial Coverage



Dynamic Traffic Status





# Crowdsourcing GPS Trajectories



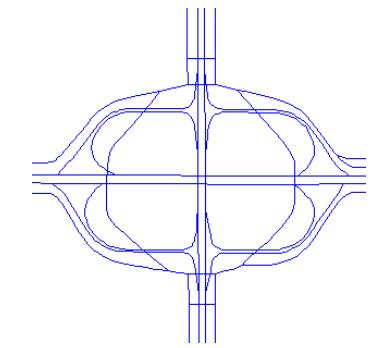
Massive Moving Objects



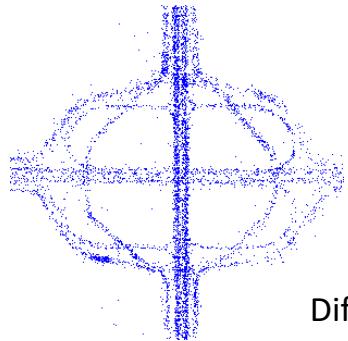
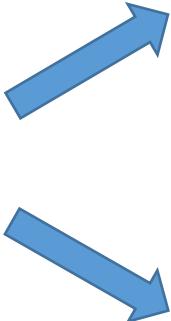
GPS Devices

Real-time City-wide Map Information

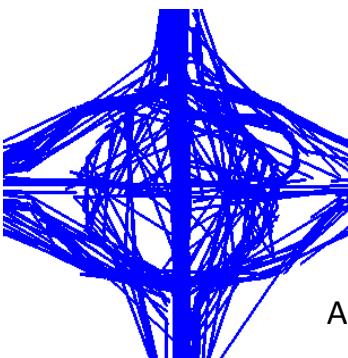
# Challenges



Ground Truth Map



Positioning Error (~15 meters)



Low Sampling Rate (30s/pt)

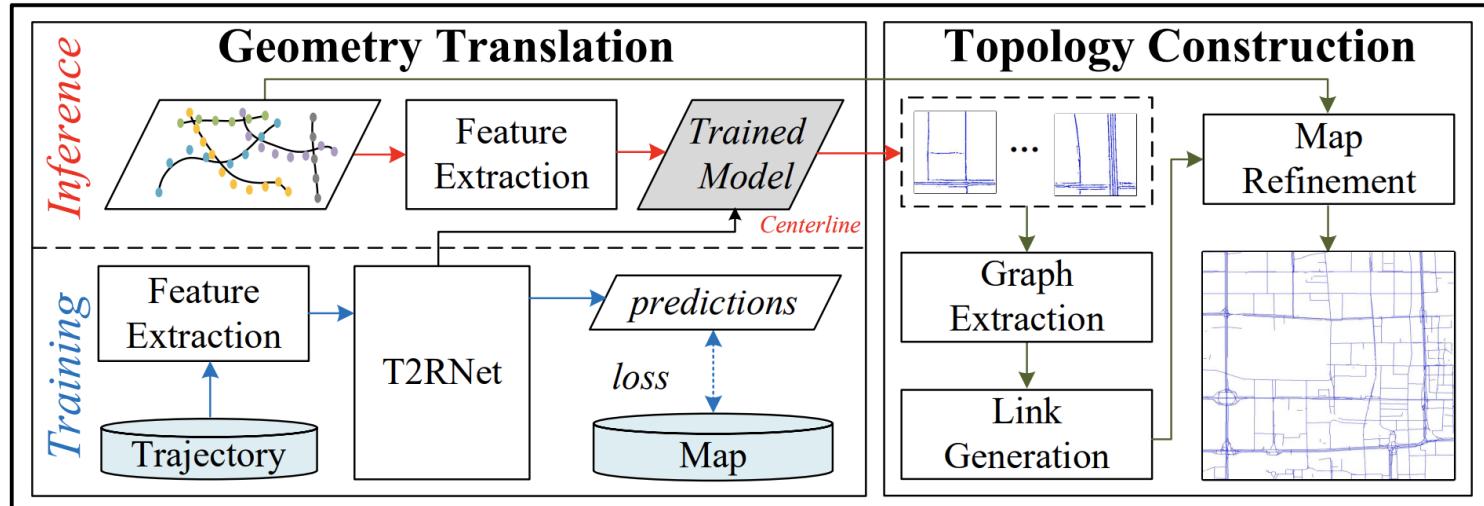
Difficult to Distinguish Spatially Close Roads

Ambiguous Underlying Traversing Routes



# Framework

- Two-stage solution (AI + DB)

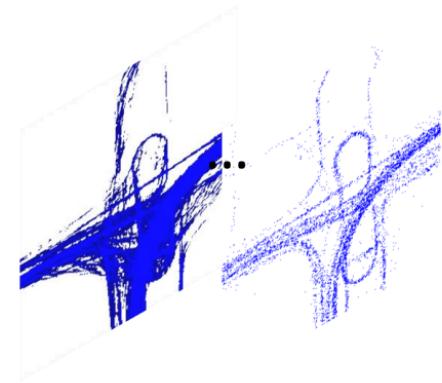




# Geometry Translation

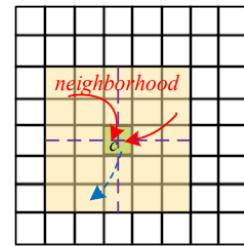
- Feature Extraction
  - For each  $I \times J$  Region Tile

Spatial View ( $\mathbb{R}^{11 \times I \times J}$ )



Point, Line, Speed, Direction (8 channels)

Transition View ( $\mathbb{Z}_2^{2 \times T \times T \times I \times J}$ )



(a) Transitions S/E at c. (b) Local Incoming Matrix. (c) Local Outgoing Matrix.

1	1
0	0

0	0
1	0

Each grid cell has an incoming and an outgoing matrix



# Geometry Translation

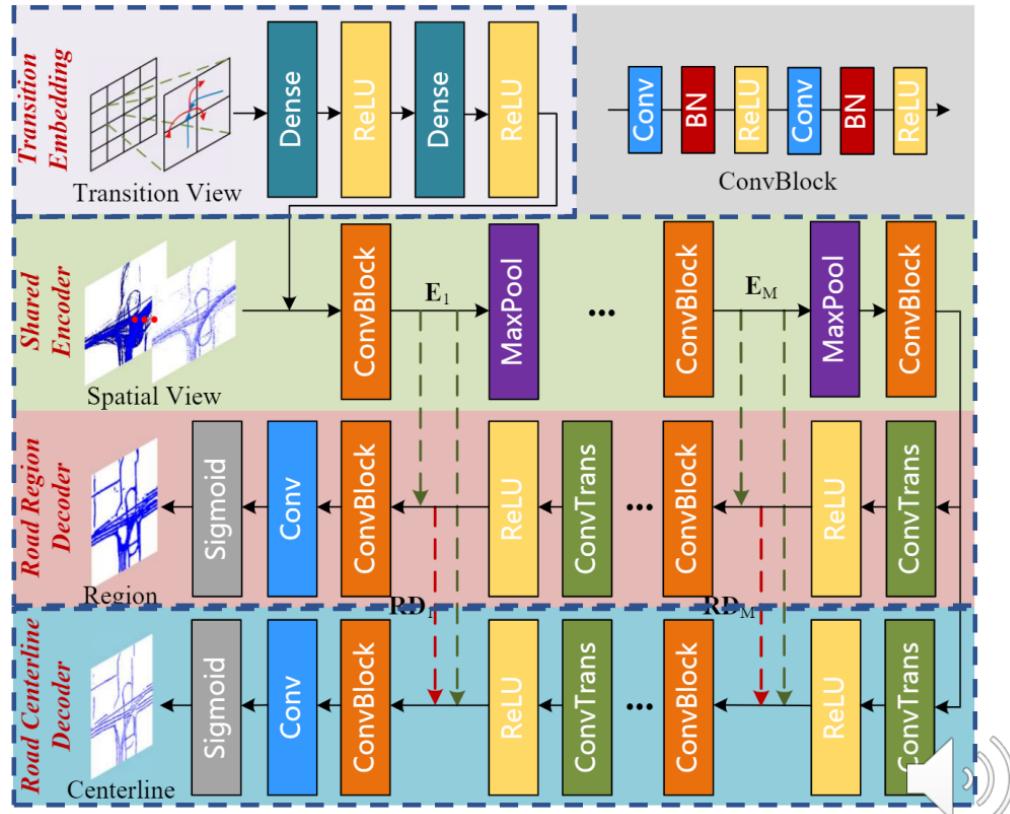


- T2RNet
  - Transition Embedding
  - Shared Encoder
  - Road Region Decoder
  - Road Centerline Decoder
- Optimization
  - Dice Loss [Milletari F, et al. 2016]

$$\mathcal{L}_{Dice}(\hat{\mathbf{Y}}, \mathbf{Y}) = 1 - \frac{2 \sum_i^I \sum_j^J \hat{Y}_{ij} Y_{ij} + \epsilon}{\sum_i^I \sum_j^J \hat{Y}_{ij} + \sum_i^I \sum_j^J Y_{ij} + \epsilon}$$

- Multi-task Loss

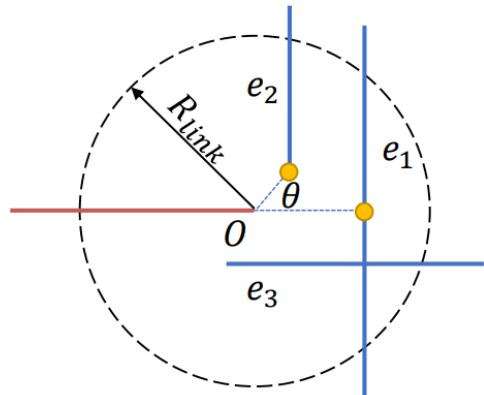
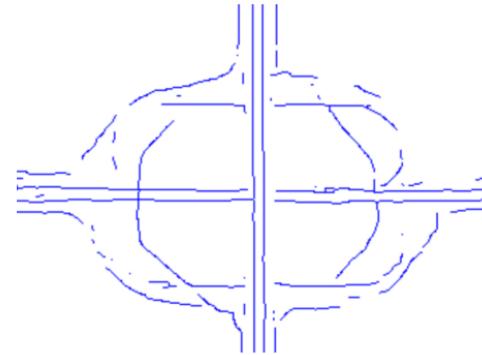
$$\mathcal{L}(\theta) = (1 - \lambda)\mathcal{L}_{Dice}(\hat{\mathbf{Y}}_c, \mathbf{Y}_c) + \lambda\mathcal{L}_{Dice}(\hat{\mathbf{Y}}_r, \mathbf{Y}_r)$$





# Topology Construction

- Graph Extraction
  - Merge predicted tiles
  - Extract road segments
- Link Generation
  - For each dead end



- Case 1: intersects another edge on the extension
- Case 2: has smooth transition to the closest dead end of another edge

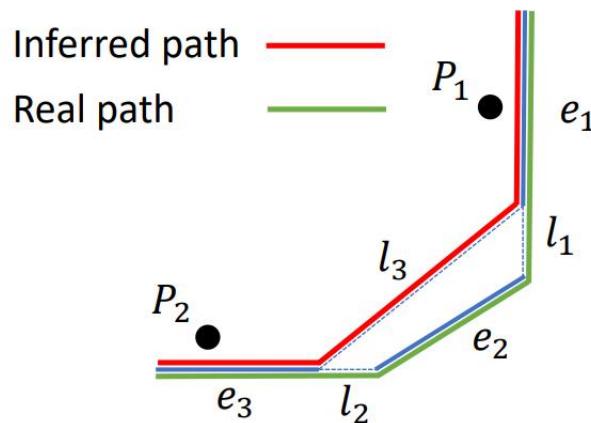


# Topology Construction

- Map Refinement

- Perform trajectory map matching on the linked map [Yuan J, et al. 2010]
- Remove edges and links with low support

If the map matching is directly applied...



Proposed solution

$$dist(P) = \sum_{e_i \in P} \omega(e_i) \cdot len(e_i)$$

$$\omega(e_i) = \begin{cases} \alpha, & e_i \text{ is a generated link} \\ 1, & e_i \text{ is a predicted edge} \end{cases} \quad \alpha > 1$$



# Evaluation

- Datasets
  - Trajectory
    - <oid, timestamp, latitude, longitude>
  - Map
    - Node: <latitude, longitude>
    - Edge: <start\_node, end\_node>

Dataset	TaxiBJ	TaxiJN
#Days	30	30
#Vehicles	500	70
Sampling Rate	~30s	~3s
Size (km <sup>2</sup> )	16×16 (5×5)	16×26 (5×5)
#Points	3.1M (304K)	5.7M (322K)
#Trajectories	66,124 (13,462)	29,556 (3,954)
Roads (km)	2,772 (284)	2,048 (123)

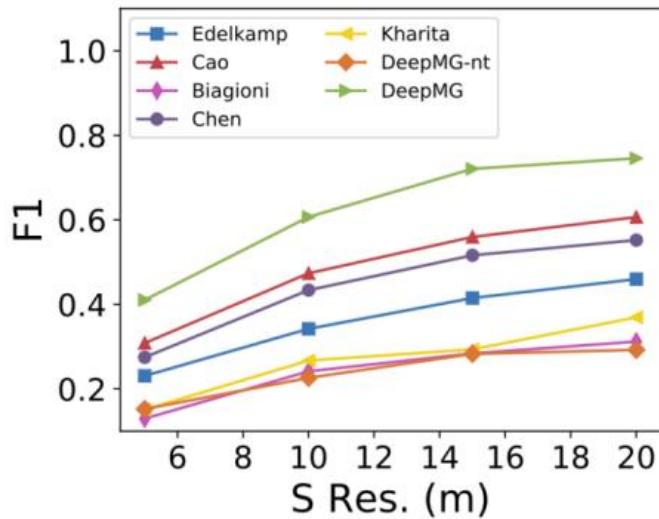
## • Evaluation Metrics

- Topological F1 [Biagioni, et al. 2012]
  - Repeat  $N$  times
    - a. Select a random starting location
    - b. Find reachable area within a maximum radius
    - c. Compare generated map with GT using F1
  - Report the average F1 score

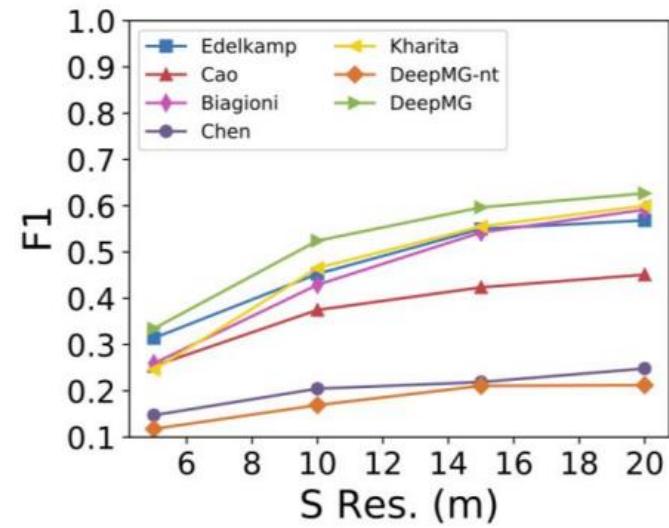




# Quantitative Comparison



TaxiBJ (32.3% improvement)



TaxiJN (6.5% improvement)

# Case Study

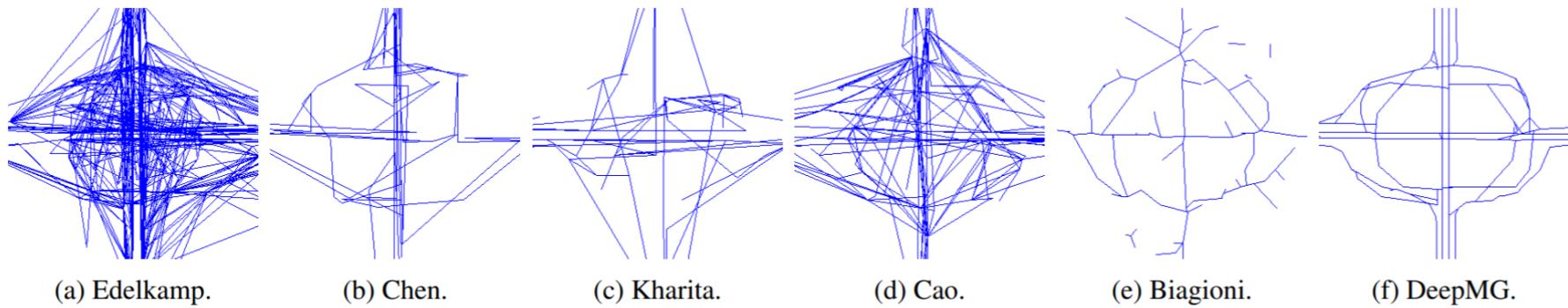


Figure 5: Seven algorithm at a roundabout near Gongzhufen area in TaxiBJ.

## Practice in JD: Community-level Map Generation from JD Couriers' Trajectories

物流轨迹地图修复

地图修复 场景模拟

马驹桥



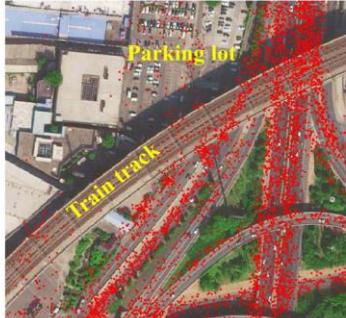
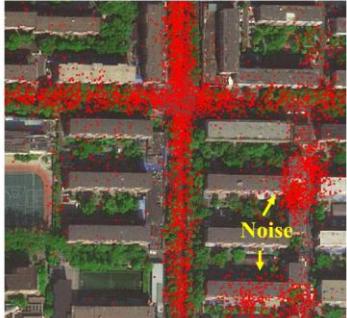
# Conclusion



- A valuable but challenging task
  - Position errors
  - Low sampling rate
- Our method
  - Geometry translation (T2RNet)
    - Convolutional network: learns the structure of the road network
    - Auxiliary task: helps the centerline inference
  - Topology construction (Link + Prune)
    - Trajectories as transition evidences
- Results
  - Superior than traditional methods
  - More effective on low-sampling rate datasets

# Related Work

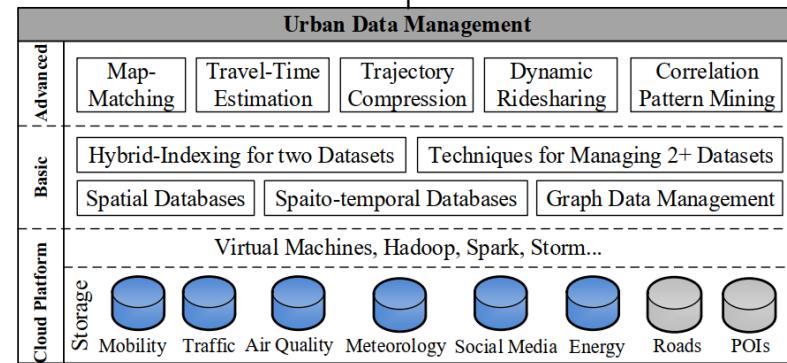
- GPS trajectories + Satellite Images



# Urban Data Management

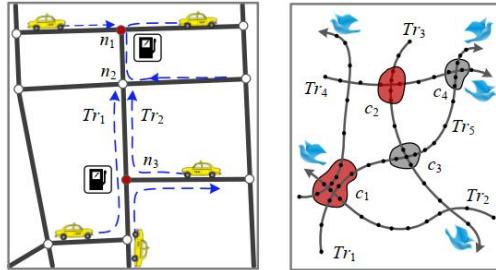
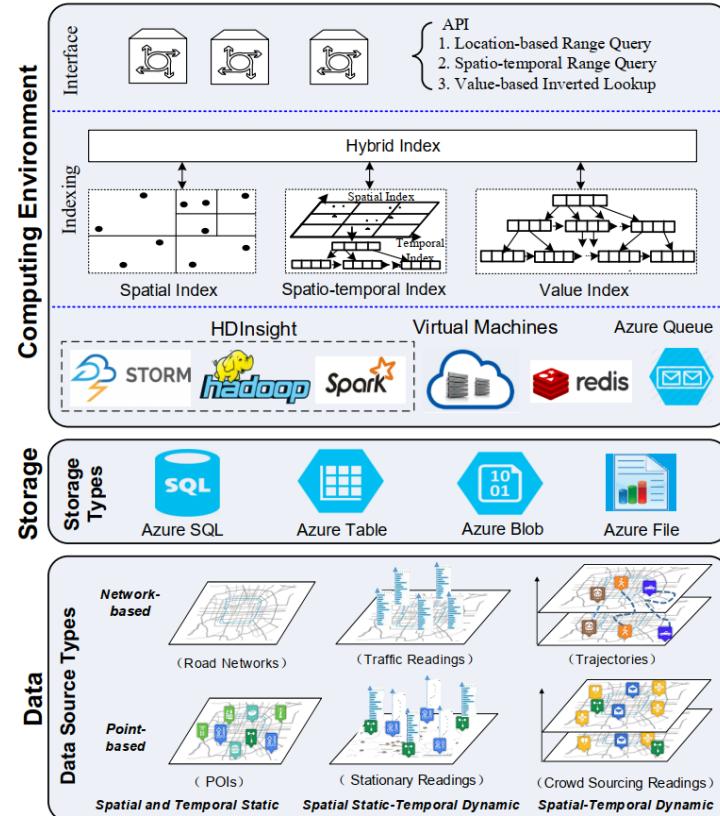


- Basic Data management Techniques
  - Spatial Databases
  - Spatio-Temporal Databases
- Advanced applications
  - Map Matching
  - Trajectory compression
  - Large-Scale Dynamic Ridesharing
- Managing Spatio-Temporal Data on the Cloud

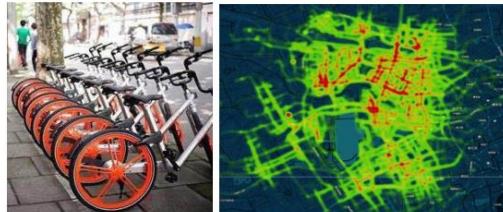




# Cloud Computing for ST Data



Y. Li, J. Bao, Y. Li, Z. Gong, **Y. Zheng**.  
[Mining the Most Influential k-Location Set From Massive Trajectories](#).  
 IEEE Transactions on Big Data. 2017



J. Bao, T. He, S. Ruan, Y. Li, and Y. Zheng et al.  
[Planning bike lanes based on Sharing-bike's trajectories](#),  
 KDD 2017

# Presented Papers



Spatio-Temporal Data Sesning and Management	1	Transformer-based Map Matching Model with Limited Ground-Truth Data using Transfer-L	<a href="#">Link</a>	TRC	2022	McGill University	Jiahui LIANG
	2	TraSS: Efficient Trajectory Similarity Search Based on Key-Value DataStores	<a href="#">Link</a>	ICDE	2022	JD	Gangyong ZHU
	3	Learning to Generate Maps from Trajectories	<a href="http://cdn.aaai.org/ojs/5435/5">cdn.aaai.org/ojs/5435/5</a>	AAAI	2020	Beijing Institute of Technology	Weilin RUAN



# Thanks!

CityMind Lab



Tencent



CAL  
NIAO 菜鸟