

## Spatial Trajectory Analytics: Past, Present and Future

Xiaofang Zhou



## Trajectory Data

...data about moving objects

### + 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:

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

### + Massive Amount of GPS Data



### + Other Types of Trajectory Data



### + More about Trajectory Data

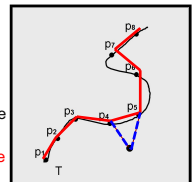
- A trajectory is obtained from sampling the movement of an object

- Some **sampling strategies** are used  $\rightarrow$  not only data, but also models to generate data

- Objects movement with **constraints** (e.g., by map)  $\rightarrow$  not only data, but also environment data

- There are many other factors which cannot be controlled  $\rightarrow$  **data quality** issues

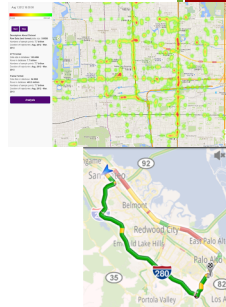
- Data can be both **redundant** as well as **sparse**  $\rightarrow$  compression, alignment and prediction



- It is non-trivial even to restore the original trace from a trajectory  $\rightarrow$  harder to compare

## + Are Trajectory Data Useful?

- Route planning
- POI recommendation
- LBS and mobile advertisement
- Resource tracking and scheduling
- Fleet management
- Road safety
- Emergency responses
- Environment monitoring...
- Urban planning and smart cities

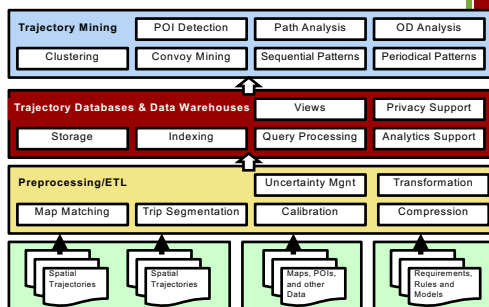


Trajectory analytics now becomes a new frontier for business intelligence, especially in real-time and in combination with other types of data

## Processing Trajectory Data

...monitoring, managing and processing

## + Trajectory Processing Framework



## The Past

...driven by curiosity

## + Moving Objects/Trajectory Work

- Initially on foundations
  - Data representation, query languages and basic operations, indexing methods etc.
- Curiosity-driven
  - Imagine a special "novel" type of query, find a "novel" indexing method and then use "standard" methods to improve efficiency
- Not directly useful
  - Strong assumptions (not useful in practice)
  - Highly specialized indexes (cannot be implemented)
- Also, data mining, social networks, recommender...

## + An Introduction Book

### ■ *Computing with Spatial Trajectories*

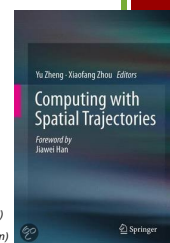
■ Yu Zheng and Xiaofang Zhou, 2011

### ■ Part I Foundations

- Trajectory Preprocessing (W.-C. Lee, J. Krumm)
- Trajectory Indexing and Retrieval (X. Zhou et al)

### ■ Part II Advanced Topics

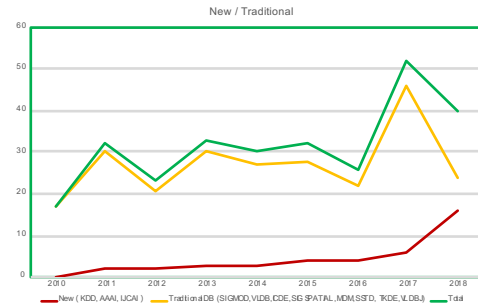
- Uncertainty in Spatial Trajectories (G. Trajcevski)
- Privacy of Spatial Trajectories (C.-Y. Chow, M. Mokbel)
- Trajectory Pattern Mining (H. Jeung, K. L. Yiu, C. Jensen)
- Activity Recognition from Trajectory Data (Y. Zhu, V. Zheng, Q. Yang)
- Trajectory Analysis for Driving (J. Krumm)
- Location-Based Social Networks: Users (Y. Zheng)
- Location-Based Social Networks: Locations (Y. Zheng and X. Xie)



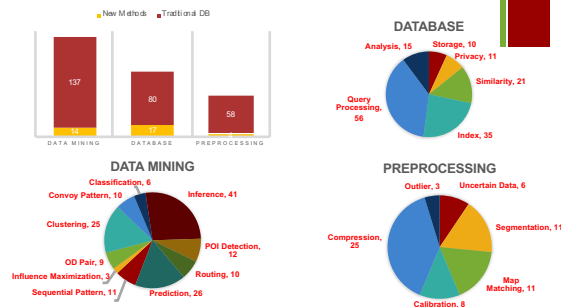
## The Present

...continues the past and also driven by new trends

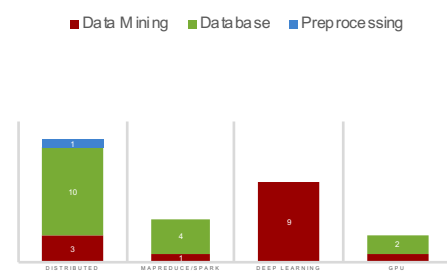
## + Trajectory Data Research: Trends



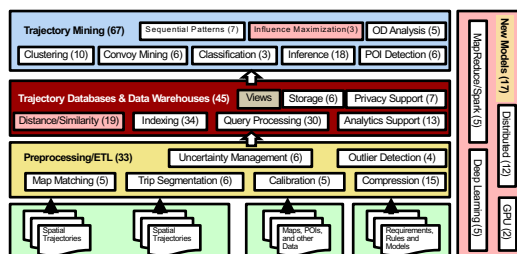
## + Trajectory Research: Topics



## + Trajectory Research: What's New



## + A Systematic View



## + Some of Our Work

Prediction of movement [ICDE08] and paths [VLDBJ10], trajectory simplification with error bound [VLDB08], path nearest neighbor query [SIGMOD09], searching trajectory by locations [SIGMOD10], most popular routes [ICDE11], probabilistic range query [EDBT11, ICDE12], materialized shortest paths [TODS12], spatial keyword search for trajectories [ICDE13,15,16], clue-based queries [VLDB17], **minimum on-road time routing** [VLDB17, VLDBJ18], **trajectory calibration** [SIGMOD13, VLDBJ15], route and location recommendation [ICDE14, SIGKDD15, ICDE16, TOIS16, TIST18], trajectory exploration and summarization [ICDE15], in-memory trajectory databases [CIKM14, SIGMOD15], privacy-preserving trajectory search [ICDE15], **data sparsity** [MDM18], ML for **speed prediction** [IJCAI18]

# Minimum On-Road Time Route Planning

...yet another “new” type of queries

## + Time Dependent Graph

20

- Time-Dependent Graph is defined as  $G_T(V, E, W)$ 
  - $V = \{v_i\}$  vertex set
  - $E \subseteq V \times V$  directed edge set
  - $W$  a set of cost functions
    - $\forall (v_i, v_j) \in E, w(v_i, v_j, t) \in W$
    - Tells how much time it costs to travel from  $v_i$  to  $v_j$  at time  $t$
    - Also known as **Speed Profile**
- Can Dijkstra algorithm still work?
  - FIFO Time-Dependent Graph
    - Given a time dependent graph  $G_T(V, E, W)$ . It is a FIFO Graph iff  $\forall (v_i, v_j) \in E, w(v_i, v_j, t) \leq \Delta + w(v_i, v_j, t + \Delta)$ , where  $\Delta > 0$ .



## + Time Dependent Path

21

- A Path  $p = \langle v_1, \dots, v_k \rangle$  from  $v_1$  to  $v_k$ 
  - $\alpha(v_i)$ : arrival time
  - $\beta(v_i)$ : departure time
  - $\forall v_i \in V: \beta(v_i) - \alpha(v_i) \geq 0$
- Total Travel Time
  - $w_{TOT} = \alpha(v_k) - \beta(v_1)$
- Fastest Path
  - $\min(\alpha(v_k) - \beta(v_1))$
  - Because of FIFO, only when  $\beta(v_i) - \alpha(v_i) = 0$  can lead to optimal solution

## + Fastest Path

22

- Single Starting-Time Fastest Path
  - $Q(v_s, v_d, t_s)$
  - Starting time  $t_s$  on  $v_s$  is fixed
  - Can use Dijkstra directly [1]
  - Time Complexity:  $O(|V| \log |V| + |E|)$
- Variations
  - Earliest Arrival Path
    - $\min(\alpha(v_d))$
    - Same as single starting-time fastest path
  - Latest Departure Path
    - $\max(\beta(v_s))$
    - Reverse

[1] Ding B, Yu J X, Qin L. Finding time-dependent shortest paths over large graphs. EDBT 2008

## + Fastest Path

23

- Interval Starting-Time Fastest Path
  - $Q(v_s, t_{s1}, t_{s2}, v_d, t_d)$ 
    - $[t_{s1}, t_{s2}]$ : starting time interval
    - $t_d$  latest arrival time
    - $\min(\alpha(v_d) - \beta(v_s))$
  - Find the optimal starting time within  $[t_{s1}, t_{s2}]$ 
    - Solution in [1]
    - Only waiting on  $v_s$  is considered.
- Problem Complexity
  - $\Omega(T(|V| \log |V| + |E|))$  [2]
  - If  $w$  is linear piecewise functions,  $T$  is the turning points in the final result

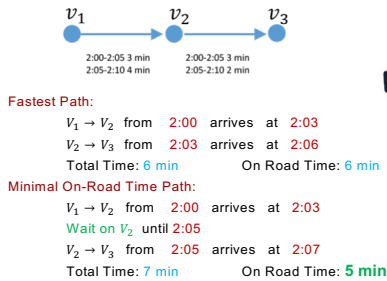
[1] B Ding, J X Yu, L Qin. Finding time-dependent shortest paths over large graphs. EDBT 2008  
 [2] L Foschini, J Hershberger, S Suri, On the complexity of time-dependent shortest paths, Algorithmica 2014

## + But...

24

- Is the Fastest Path always best?
  - For Logistics company
    - 
    - Major operation cost  $\leftarrow$  Fuel consumption  $\leftarrow$  Time spent on road
  - For tourists
    - Reduce their time spent on road so that they can spend more time on the tourist attractions.
- Spending less time on road rather than arriving final destination earlier

## + Example



25

## + Time Dependent Graph with Parking Nodes

Time-Dependent Graph is defined as  $G_T(V, E, W)$

- $V = \{v_i\}$  vertex set
- $V' \subseteq V$ , parking vertex set
- $E \subseteq V \times V$  directed edge set
- $W$  a set of cost functions

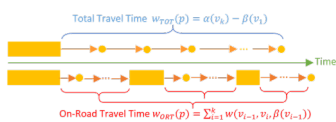
26

## + On-Road Time

- A Path  $p = \langle v_1, \dots, v_k \rangle$  from  $v_1$  to  $v_k$
- $\alpha(v_i)$ : arrival time
- $\beta(v_i)$ : departure time
- $\forall v_i \in V - V'$ :  $\beta(v_i) = \alpha(v_i)$
- $\forall v_i \in V'$ :  $\beta(v_i) \geq \alpha(v_i)$

## ■ On-road Travel Time

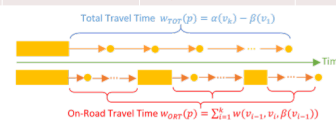
$$w_{ORT} = \sum_{i=1}^k (w(v_{i-1}, v_i, \beta(v_{i-1})))$$



27

## + Time-Dependent Path Problems

Graph Type	Path Problem		Objective	Waiting	
Static Graph	Shortest Path		$\sum_{i=2}^d Weight(p_{i,i+1})$	Total Static Weight	No
Time Dependent Graph	Single start time fastest path	General	$\sum_{i=2}^d Time(p_{i,i+1})$	Total Temporal Weight	No
		Earliest Arrival	$Min(Arrival(v_d))$		
		Latest Departure	$Max(Depart(v_2))$		
	Interval start time fastest path	$Min(Arrival(v_d) - Depart(v_2))$		Minimum Total Time	Source Vertex $V' = \{v_s\}$
	MORT	$Min(\sum_{i=1}^{k-1} w(v_{i-1}, v_i, \beta(v_{i-1})))$	Minimum On-Road Time	A set of parking vertices $V' \subseteq V$	



28

## + Naïve Approach

- Find the fastest path from  $v_s$  to  $v_d$   
 $\langle v_s, v_{s+1}, v_{s+2}, \dots, v_d \rangle$  with the optimal starting time  $t_s$  on  $v_s$
- Find the fastest path from  $v_{s+1}$  to  $v_d$   $\langle v_{s+1}, v'_{s+2}, \dots, v_d \rangle$  with optimal starting time  $t_{s+1}$  on  $v_{s+1}$
- Run repeatedly at every node

Any problems?

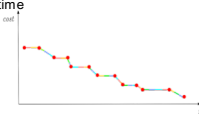
29

## + Algorithm Overview

Problem: Find a path with  $\text{Min}(\sum_{i=1}^{k-1} w(v_{i-1}, v_i, \beta(v_{i-1})))$

## ■ MORT Algorithm

- Active Time Interval
- Path Expansion
  - $C_i(t)$  The minimum on-road travel time from  $v_s$  to  $v_i$  that arrives  $v_i$  on  $t$ .
  - Update  $C_i(t)$  for all  $v_i$  until the optimal is reached
  - Non-Increasing property of parking vertices
    - Waiting will increase the total travel time
    - Waiting will not increase the on-road travel time
- Basic Algorithm:  $O(T|V| \log|V| + T^2|E|)$
- Incremental Algorithm:  $O(L(|V| \log|V| + |E|))$
- Approximation
- Route Retrieval



L. Li, W. Hua, X. Du and X. Zhou, "Minimal On Road Time Route Scheduling on Time-Dependent Graph", VLDB 2017.

30

## Road Speed Profile Generation

...another example of large-scale space problem

### + Speed Profile Generation

1. Map Matching
  - Attach the GPS points on the roads
  - Use the length and temporal information to get speed of each road at different time
2. Speed Data Collection
  - 5-Minute Histogram
3. Missing Value Estimation
  - Cosine Similarity, Matrix Factorization based Collaborative Filtering, HMM....
4. Compression
  - PLA (Piecewise Linear Approximation)

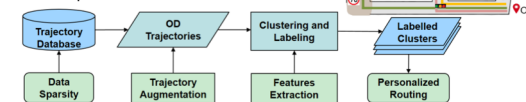
L. Li, K. Zheng, X. Zhou and S. Wang, "Go Slow to Go Fast: Minimal On-Road Time Route Scheduling with Parking Facilities Using Historical Trajectory", in *VLDB Journal* (accepted).

## Data Sparsity

...no matter how much data you have, you don't have enough

### + Trajectory Clustering and Labeling

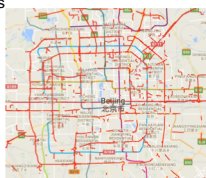
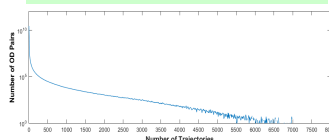
- Applications
  - Moving Behaviors Analysis
  - Personalized Routing
- Clustering
  - Data Sparsity: Origin-Destination Trajectory
  - Clusters Identification: Efficient Algorithms
- Labeling
  - Features: Fastest, Shortest, Most popular, Time-dependent



### + Trajectory Augmentation

- Data augmentation approach
  - Factorization-based [1]: extra data sources (geo-spatial, temporal, and historical correlation)
  - Concatenation-based [2]: sub-trajectories

No Robustness Validation Guarantee



[1] Wang Yilun, Yu Zheng, and Yexiang Xue, "Travel time estimation of a path using sparse trajectories," *SIGKDD*, 2014.  
[2] Dai Jian, Bin Yang, Chenjuan Guo, and Zhiming Ding, "Personalized route recommendation using big trajectory data," *ICDE*, 2015.

D. He, B. Ruan, B. Zheng and X. Zhou, Origin-Destination Trajectory Diversity Analysis: Efficient Top-k Diversified Search, *MDM* 2018

## Deep-Learning for Speed Predication

...many things we do is about using spatial data for prediction

## Deep Learning + ITS = ?



To Deriving Knowledge from Data

+ CV =



+ GO =



+ ITS = ?

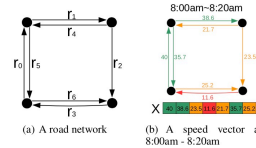
## + The Problem

Given:

- A road map (as a directed graph)
- A sequence of **speed vectors**, each vector is the speed at each road segment during a time interval

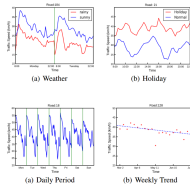
$$X_t = [x_t^{r_0}, x_t^{r_1}, \dots, x_t^{r_{|E|-1}}]$$

**Problem:** Given the historical observations  $\{X_i | i = 1, \dots, t\}$ , this paper aims to predict  $Y_t = \{X_j | j = t+1, \dots, t+z\}$ , where  $z$  is the number of time intervals to be predicted.

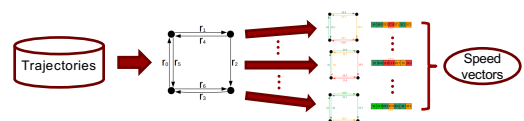


## + Challenges

- Spatiotemporal characteristics**
  - Speed on a road tends to slow down if congestions occur in its surrounding area
  - Speed in the downtown area is relatively low during rush hours
- Restricted by road network**
  - Traffic in a certain road just affects its adjacent roads
  - Traditional CNN cannot capture topology of road network
- Other factors**
  - Context Information
    - e.g. weather, holiday
  - Periodic Law
    - e.g. daily period, weekly trend

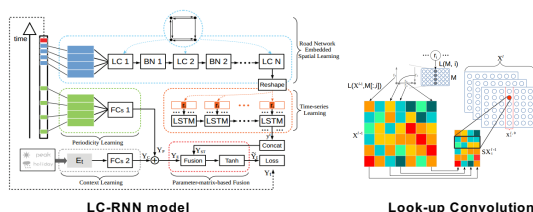


## + Generating Speed Vectors



## + LC-RNN Model

- Look-up Convolution (LC):** learn the latent features of surrounding area
- LSTM components:** learn the time-series pattern that is aware of surrounding area dynamics

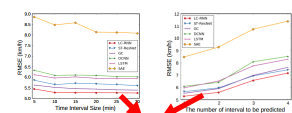


## + Results

Beijing Dataset

- A big road network about 10 thousand main roads
- Trajectory data comes from 1<sup>st</sup> Mar to 31<sup>st</sup> Jul in 2016
- The first 4 months as train and the last 1 month as test

Methods	RMSE
SVR	10.245
HAIRMA	11.867
SAE	8.471
LSTM	5.958
DCNN	6.085
GC	5.514
ST-ResNet	5.749
LC-3	5.437
LC-3-RNN	5.328
LC-3-RNN-E	5.296
LC-3-RNN-E-BN	5.392
LC-2-RNN-E-BN	5.319
LC-4-RNN-E-BN	5.285



Shanghai Dataset

- A small road network about 1.5 thousand main roads
- Trajectory data comes from 1<sup>st</sup> Mar to 31<sup>st</sup> Apr in 2015
- The first 45 days as train and the last 15 days as test

Methods	RMSE
SVR	9.083
HAIRMA	9.317
SAE	7.421
LSTM	5.274
DCNN	5.269
GC	4.921
ST-ResNet	5.087
LC-3-RNN-E-BN (ours)	4.686

Z. Lv, J. Xu, K. Zheng, P. Zhao, H. Yin and X. Zhou., "LC-RNN: A Deep Learning Model for Traffic Speed Prediction", *IJCAI* 2018.





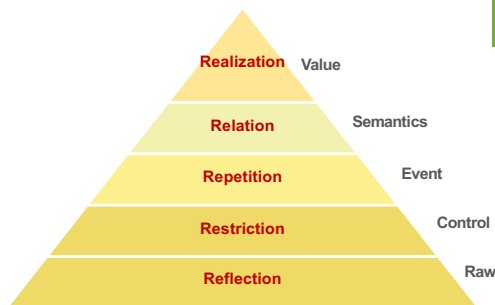
# Trajectory Data Management System

...a common platform and API

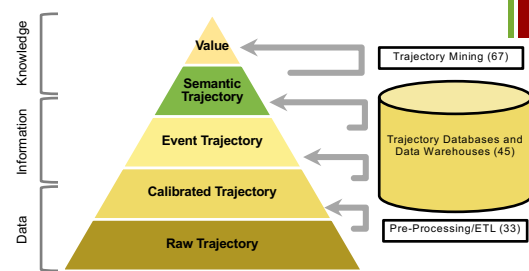
## + Why Common Platform?

- Universal
  - GPS, telecom tokens, social apps...
- Shared enterprise data
  - For monitoring, predication, business insights...
- Separation of conceptual, logical and physical design
  - Especially many computing platforms to consider today
- Other benefits we took for granted
  - Optimization for data storage and query processing, scheduling, concurrency control...

## + The 5R Approach

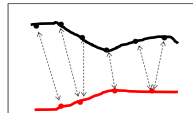


## + A Spatiotemporal Pyramid



## + Trajectory Calibration

- Popular trajectory distance measures
  - Euclidean distance
  - LCSS (longest common sequence)
  - DTW (Dynamic time warping)
  - EDR (Edit distance on Real sequences)
- How distance measures work?
  - Sample points alignment
  - Aggregating differences of aligned pairs
- Experiments
  - Ground Truth: 11,000 high-sampling-rate real trajectories
  - Derived Trajectory Datasets: re-sampling, shifting, jumping
  - Results?



H. Su, K. Zheng, H. Wang and X. Zhou, Calibrating Trajectory Data for Similarity-based Analysis, **SIGMOD** 2013

## + SparkDB

- A time-centric storage and processing system for trajectories
- Designed for in-memory computers
- A more ambitious system called **Traminer** is under development, following the proposed processing framework
- Now supported by several companies

H. Wang, K. Zheng, X. Zhou and S. Sadiq, "SharkDB: An In-memory Column-oriented Trajectory Storage", **CIKM** 2014

## Batch Fastest Path Queries

...can we do better if we can queries in batch?

### + A Real Problem

- At any time, 100K-1M OD pairs are given for route planning
- Options:
  - Processing them in parallel
  - Materializing all-pair shortest path information
  - Batch processing?
- Additional dimensions:
  - Ridesharing
  - Streaming requests (requests come continuously and cars are moving)

### + Batch Shortest Path Queries

- Case of 1:N
  - Dijkstra can be used straightaway
  - A\* can be generalized
  - A good partition of N can improve efficiency
- Case of M:1
  - This can be done by reversing the above case
- Case of M:N
  - One-directional
  - Bidirectional
  - Partition-based

## Trajectory-based Entity Linking

...everyone's movement pattern is unique

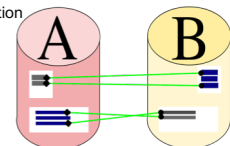
### + People Have many Trajectory Datasets

- Broadly available human mobility data:
  - Check-ins
  - Credit card transactions
  - Phone call records
  - Go-card records
  - Vehicle trajectories
  - Many social networks...



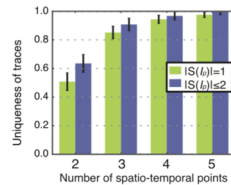
### + Entity Resolution using Trajectory Data?

- Understanding the extent to which spatiotemporal data are distinctive is crucial to:
  - Location privacy protection
  - Entity resolution
  - Data integration
  - Spatial content retrieval
  - Personalized location recommendation
  - Driver performance assessment
  - ...



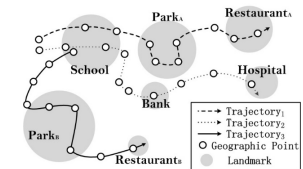
### + Only 4 Points, Really?

- "Unique in the Crowd: The privacy bounds of human mobility", Yves-Alexandre de Montjoye, César A. Hidalgo, Michel Verleysen & Vincent D. Blondel, *Scientific Reports* (2013)
- "four spatiotemporal points are enough to uniquely identify 95% of the individuals"



### + Everyone Has Mobility Signature?

- Spatial signature?
  - Regular: you visit frequently, such as your office building
  - Unique: you can be distinguished from others, like personal home address



### + Yes!

Top K	Precision
1	0.95987964
2	0.97241725
3	0.9779338
4	0.98194584
5	0.98445336
6	0.98545637
7	0.98696088
8	0.98746239
9	0.9884654
10	0.9889669

(1000 taxis)

- Signature: TFIDF-weighted point set
  - Many possible reduction methods (simple frequency-based truncating, PCA, LSH)
  - Several alternatives (may need different similarity measures)
- Cosine similarity
- Tested with 11K taxi data over 3 months in 3 major cities
- Accuracy: now extremely high
- Performance is bad: we are improving it

## Some Other New Research Problems

...from discussions with companies doing trajectory analytics as their bread-and-butter business

### + A List of Problems

- ETA: O-side and D-side
- Map matching + map inferencing: an integrated approach
- Data fusion: among trajectory datasets and with others
- Similarity based search
- Traffic prediction: for prevention and intervention
- Transport mode detection: both large/small scale
- Personalized/constrained routing algorithms
- Privacy: can you really protect trajectory privacy?
- Smart city – a holistic traffic solution
- ...

### + Conclusions

- New problems
  - More data, more queries, more applications, more tools
  - From SDBMS, spatial data mining to spatial learning
- Some current research problems
  - Large-scale space problems
  - Dynamic road networks
  - Massive batch queries
  - Personalization and privacy issues
- We need a DBMS approach!

