
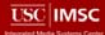


Blind Evaluation of Nearest Neighbor Queries Using Space Transformation to Preserve Location Privacy

Ali Khoshgozaran and Cyrus Shahabi
 University of Southern California
 Los Angeles, CA 90089-0781
 [ajfkhosh, shahabi]@usc.edu
 http://infolab.usc.edu


NEWS Privacy worry over location data

The Washington Post Online Data-Coin Personal: Call Phone Records for Sale

Obama's BlackBerry brings personal safety risks

POI Location Server (LS)

Which is nearby?

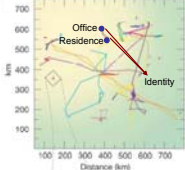
They Know

Isn't Confidentiality Enough?

Sensitive information obtained by anonymous location data


- Baraba'si et al., Nature'08

Human Mobility \leftrightarrow Spatial Probability Distribution



- Anonymous queries leak information

Location Queries \leftrightarrow Affiliations (political, religious, etc.)



3/42

Problem Definition

Objects $S = \{o_1, o_2, \dots, o_n\}$

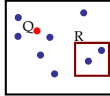
*k*NN Query

~~kNN with respect to query point Q~~
 Prevent Spatial Cores in Location-Based Services

$S' \subseteq S$ of *k* objects where for any object $o' \in S'$ and $o \in S - S'$
 $D(o', Q) \leq D(o, Q)$

Range (Window) Query

Range with respect to query window R:
 $S' \subseteq S$ of objects where for any object $o' \in S'$ o' is within R



What is required to make these queries "Privacy Aware"?


Blind Evaluation Criteria

When querying LS, the location of the querying user should not be revealed to untrusted entities through R, Q, or the query result set.

4/42

Trust and Threat Model

- Users subscribe to LS's services
 - LS owns a publicly available POI database DB
- LS is *honest* but *curious*
 - Database software is trusted
 - LS might passively exploit sensitive information
- Any entity in the system can be adversarial
 - The LS and other clients
 - Slightly different for querying other users
- Secure client/server communication channel
 - Any privacy violation should have included LS
 - We focus on LS as the most powerful adversary



5/42

Privacy/Efficiency Dilemma

- Privacy:** Hiding knowledge of object & query *locations* from LS
- Efficiency:** LS requiring this knowledge for efficient query processing

Privacy \longleftrightarrow Our Contribution \longleftrightarrow Efficiency

Information-theoretic secrecy <ul style="list-style-type: none"> Privacy against an adversary with unbounded computational resources and infinite time Lower communication & computation bound: linear w.r.t. database size 	Server knowing all information about queries and object locations <ul style="list-style-type: none"> No user location privacy possible
---	---

6/42

Related Work

Privacy

Cryptographic Techniques

- S. Zhong et al., TR'04
- Indyk et al. TCC'06
- G. Zhong et al., PET'07

No spatial query processing (MPC schemes)
 $O(n)$ computation and/or communication

Our Goal: Avoiding a linear scan of the entire DB

Efficiency

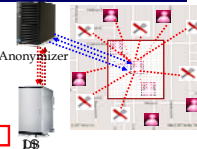
K-anonymity/Cloaking Approaches

- Gruteser et al. MobiSys'03
- Gedik et al. ICDCS'05
- Mokbel et al. VLDB'06
- Kido et al. ICPS'05
- Chow et al. GIS'06
- Ghinita et al. WWW'07 & SSTD'07

Trusting an Anonymizer
 Single point of failure/attack
 Sensitive to number of subscribed users
 No query processing

Assuming all users are trustworthy
 Dependence on other user locations

Our Goal: Complete cloaking and anonymity



7/42

Space Encoding

Offline Process

Data Owner: Original Space → Points of Interest → Space Encoder/Decoder → Encoded Locations → Transformed Space

Query Time

Client: Original Space → Query Encoder/Decoder → Encoded Query → Transformed Space

Transformation Key

Transformation Properties:

- ✓ Efficiency (locality preserving)
- ✓ Privacy (irreversible)

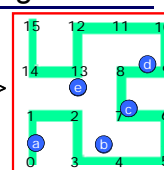
Space-Filling Curves

8/42

Background: Space Filling Curves


- Passing through (indexing) all points without crossing itself
- Example: $\langle a, b, c, d, e \rangle \rightarrow \langle 0, 4, 7, 9, 13 \rangle$

H-values



N=2 → H: 0-15
Hilbert Curves

- Proximity & distance preserving



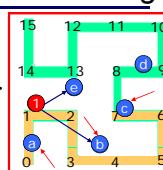
$$H = \ell(P): [0, 2^N - 1]^d \rightarrow [0, 2^{Nd} - 1]$$

$$d=2: H = \ell(P): [0, 2^N - 1]^2 \rightarrow [0, 2^{2N} - 1]$$

9/42

Hilbert Curves: Proximity Preserving

- Proximity in Hilbert space
- Example: $\langle a, b, c, d, e \rangle \rightarrow \langle 0, 4, 7, 9, 13 \rangle$
- $2NN(Q) = e$ because $D(Q, e) < D(Q, b)$



N=2 → H: 0-15

Approximate distance preservation

Complexity: Constant computation and communication

- Each node visited contains at least one object

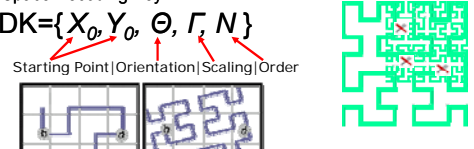
10/42

Hilbert Curves: One-wayness

- Five parameters decide how points are traversed (indexed)
- Possible when curve parameters are unknown
 - Space Decoding Key

$$SDK = \{ X_0, Y_0, \Theta, \Gamma, N \}$$

Starting Point | Orientation | Scaling | Order

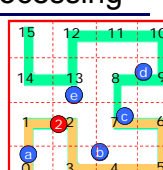


- Linear increase in N results in exponential increase in H-values
 - $3 \times 2^2 = 12$ increase in H positions for 1 H-values
- Exponential complexity for LS to reverse the transformation

11/42

2-Phase kNN Query Processing


- **Offline Space Encoding**
 - Encoding points of original space
 - Trent chooses SDK
 - Trent constructs a lookup table DB
 - $DB = \{ \langle a, 0 \rangle, \langle b, 4 \rangle, \langle c, 7 \rangle, \langle d, 9 \rangle, \langle e, 13 \rangle \}$
 - Trent encrypts objects identifiers
 - Trent uploads DB to LS
- **Online Query Processing**
 - Alice encodes her query point Q: $H = \ell(X_Q, Y_Q)$
 - Knowing H and k, LS computes the result set
 - $H=2, k=3 \rightarrow RS^* = \{0, 4, 7\} = \{ \ell(X_a, Y_a), \ell(X_b, Y_b), \ell(X_c, Y_c) \}$
 - Knowing SDK, Alice gets $RS = \{ (X_a, Y_a), (X_b, Y_b), (X_c, Y_c) \}$



12/42

Curve Rotation & kNN Search

- Issue: Approximation due to dimension reduction
 - Hilbert curves widely used for dimension reduction
- Indexing data with a **rotated** dual Hilbert curve
- Drawbacks of using a single curve:
 - 1. $N \uparrow$ (linear) \rightarrow Missed Sides \uparrow (exponential)



- 2. Reducing number of neighbors from 4 to 2

13/42

Dual Curve Query Resolution (DCQR)

- Trent indexes objects using both curves (SDK/SDK')
- Queries are evaluated on both curves
- kNN Search:
 - Alice computes $H = \ell(X_Q, Y_Q)$ & $H' = \ell(X_Q, Y_Q)$ for Q
 - LS runs two separate queries and returns 2k points to Alice
 - Alice sorts the result sets and pick the top k
- Query complexity is not affected by DCQR

14/42


Dual Curve Indexing

- We use a dual curve which is a replication of the original curve **rotated and shifted**
 - Rotation improves kNN search precision with no effect on range search
 - Translation reduces server throughput in processing range queries with positive effect on kNN search

15/42

Performance Evaluation

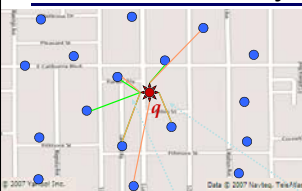
- Methodology:** issuing 1000 kNN queries with random origin
- Datasets (10000 data points):**
 - Uniform Distribution
 - Real-world**
 - Restaurants from NAVTEQ in a 26 by 26 mile area in Los Angeles
 - Skewed
 - Four clusters of points: 99% Gaussian with ($\sigma=0.05$ and Random μ) and 1% uniform
- Evaluations:**
 - Q_i
 - A_i
- Parameters:**
 - C_i
 - D_i
- Assumptions:**



(a) Uniform (b) Real-world (c) Skewed

16/42

Accuracy Metrics



- Actual Query Results
 $R = \{o_1, o_2, \dots, o_k\}$
- Approximated Query Results
 $R' = \{o'_1, o'_2, \dots, o'_k\}$
- $R \cap R'$ = Common Results

Metric 1: The Resemblance: $\rho = \frac{|R \cap R'|}{|R'|}$

Metric 2: The Displacement: $\rho = \frac{1}{K} \left(\sum_{i=1}^K \|Q - o'_i\| - \sum_{i=1}^K \|Q - o_i\| \right)$

Accuracy vs. Displacement

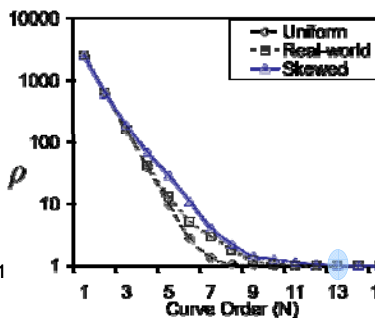
17/42

Effect of the Curve Order (N)

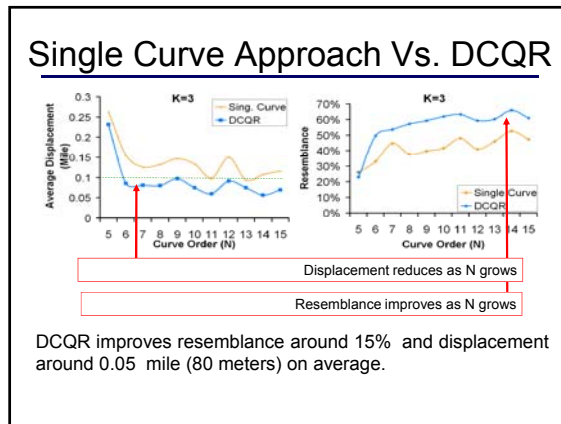
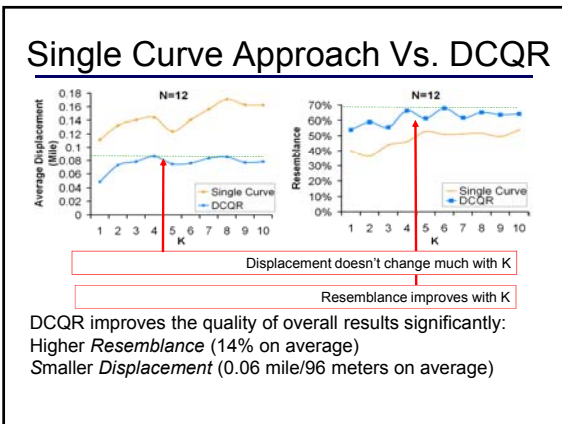
$$\rho = \frac{\sum_{i=1}^N o_i}{2^{2N}}$$

Ideally $\rho \leq 1$

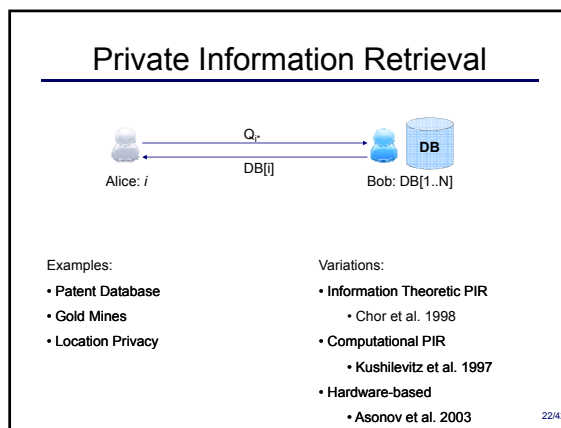
Uniform (skewed):
First (last) to hit $\rho=1$



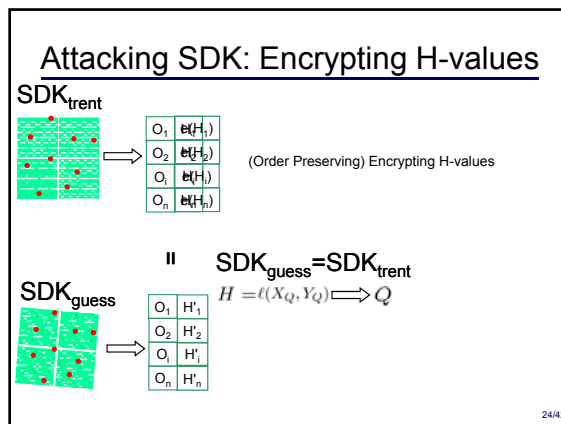
18/42



- ### Location Privacy through Information Hiding
- Achieving Location Privacy by
 - Hiding user **identity**
 - Who's accessing? (orthogonal to our work)
 - What is being accessed?
 - Developing a secure and **privacy aware spatial index**
 - Developing such privacy index reduces to
 - 1. secure index navigation
 - 2. private object retrieval

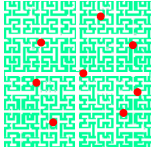


- ### Discussion
- Strengths?
 - Computation/Communication Cost
 - Lightweight client overhead
 - Weaknesses?
 - Approximate
 - Prior Knowledge
 - Object distributions
 - Correlation queries
- Privacy ← Efficiency



Attacking SDK: Random Translation

- Before indexing, points are first translated using a random vector $\langle \epsilon, \hat{\epsilon} \rangle$
 - Analogous to the notion of salt in cryptography



25/42

Approximating SDK

- Assume LS knows precise values of N, Θ, Γ and X_0 and wants to guess Y_0 by Y'_0
- LS indexes objects with SDK_{guess} and compares DB_{guess} with DB

	N				
$ Y_0 - Y'_0 $	13	14	15	20	30
0.00	0.05%	0.05%	0.05%	0.05%	0.05%
0.0001	0.05%	0.025%	0.05%	0.05%	0.05%
0.0001	0.025%	0.05%	0.05%	0.025%	0.025%
0.00001	3.26%	0.29%	0.22%	0.05%	0.05%

	N				
Γ/Γ'	13	14	15	20	30
0.00	0.05%	0.05%	0.05%	0.05%	0.05%
0.0001	0.05%	0.025%	0.05%	0.05%	0.05%
0.0001	0.025%	0.05%	0.05%	0.025%	0.025%
0.00001	0.48%	0.12%	0.03%	0.05%	0.05%

10^{-5} mile \sim 1.6cm

26/42

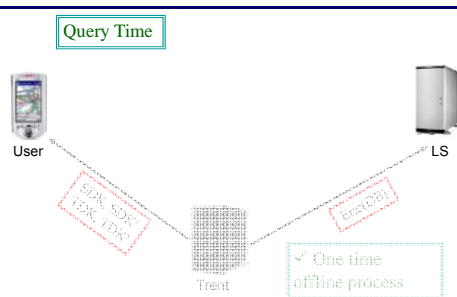
LS & External Adversary Collusion

We assume unmolested program execution on users' client devices that prevents adversaries from breaching into a client device

- Running code securely on an untrusted client is an open problem
- 100% utilization of server
 - Hard to map an H-value request to an external adversary's location
- Using SALT, makes it impossible for the attacker and LS to find the entire mapping

27/42

End to End Architecture

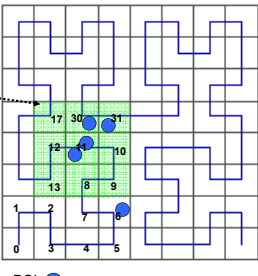


28/42

Range Queries

A window (range) query

Querying objects O such that $\ell(O)$ belongs to the set $RS = \{8, 9, 10, 11, 12, 13, 17, 30, 31\}$



POI ●

29/42

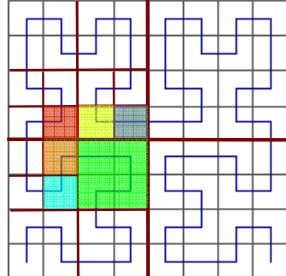
Steps to Answer a Range Query

The Hilbert space is recursively decomposed into each piece is fully contained in the range.

Result: maximal quad-tree blocks

Property: H-values inside a maximal block form a continuously increasing sequence.

Tsai et al. A strip-splitting-based optimal algorithm for decomposing a query window into maximal quadtree blocks, ICDE'04



30/42

Steps to Answer a Range Query

8-11, 17, 13, 12, 30, 31

Sort

8-11, 12, 13, 17, 30, 31

Merge

8-11, 12-13, 17, 30-31

Each sequence is called a **run**

Chung et al. Space-filling approach for fast window query on compressed images, Transactions on Image Processing '00

31/42

Example

A range query is decomposed into its maximal quad tree blocks (each of the colored squares is a maximal quad tree block)

The final runs:
Each colored part is a single run (7 runs total)

Range Query Result Set Is Exact but May Contain Excessive Objects

32/42

Privacy Aware Range Query Search

- Packing all runs
 - The server can
 - The cardinality
 - Number of runs
 - Server learns a [RS] points
- Query runs are
 - If each run is q_i
 - r_{xy} (run length)

33/42

Algorithm Complexity

- Range algorithm takes $O(n_i \log T)$ time where $n_i = \max(n_1, n_2)$ for a query of size $n_1 * n_2$ and $T = 2^N$ (N is the curve order).
 - $O(n_i)$ for decomposition
 - $O(n_i * N)$ for finding α and β
 - $O(n_i * \log n_i)$ for sorting sub runs
 - $O(n_i)$ for merging runs
- Search:
 - Alice performs quadtree decomposition on both curves and chooses the one with fewer runs and sends runs to LS
 - LS returns the encoded result set to Alice

34/42

Curves Translation & Range Search

- A range query maps into many runs
 - It is desirable to minimize the number of runs (quadtree blocks)
- Indexing the data with a second shifted curve can achieve this

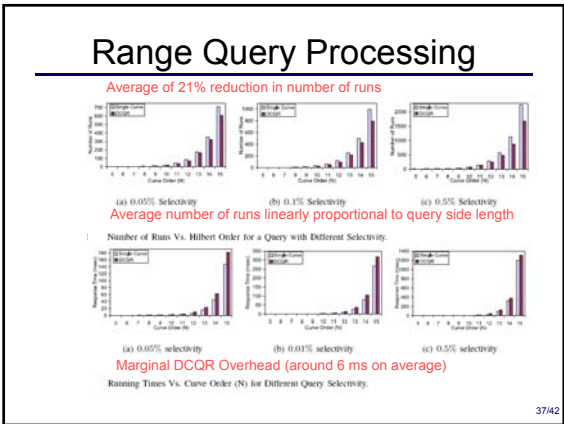
35/42

Range Query Processing

- Range queries are exact
- Include excessive objects
- Measuring precision $|\text{relevant}|/|\text{returned}|$
 - Higher precision for larger selectivity

Precision reaches 100% for $N \geq 13$ for real-world data

36/42



- ### Attacks on Cloaking and Anonymity
- Center of the cloaked region
 - Single point of failure and attack
 - Cloaking failure under certain distributions
 - Availability of all user locations to LS in anonymity approaches
 - Huge performance penalty for privacy-paranoid users.
- 38/42

