

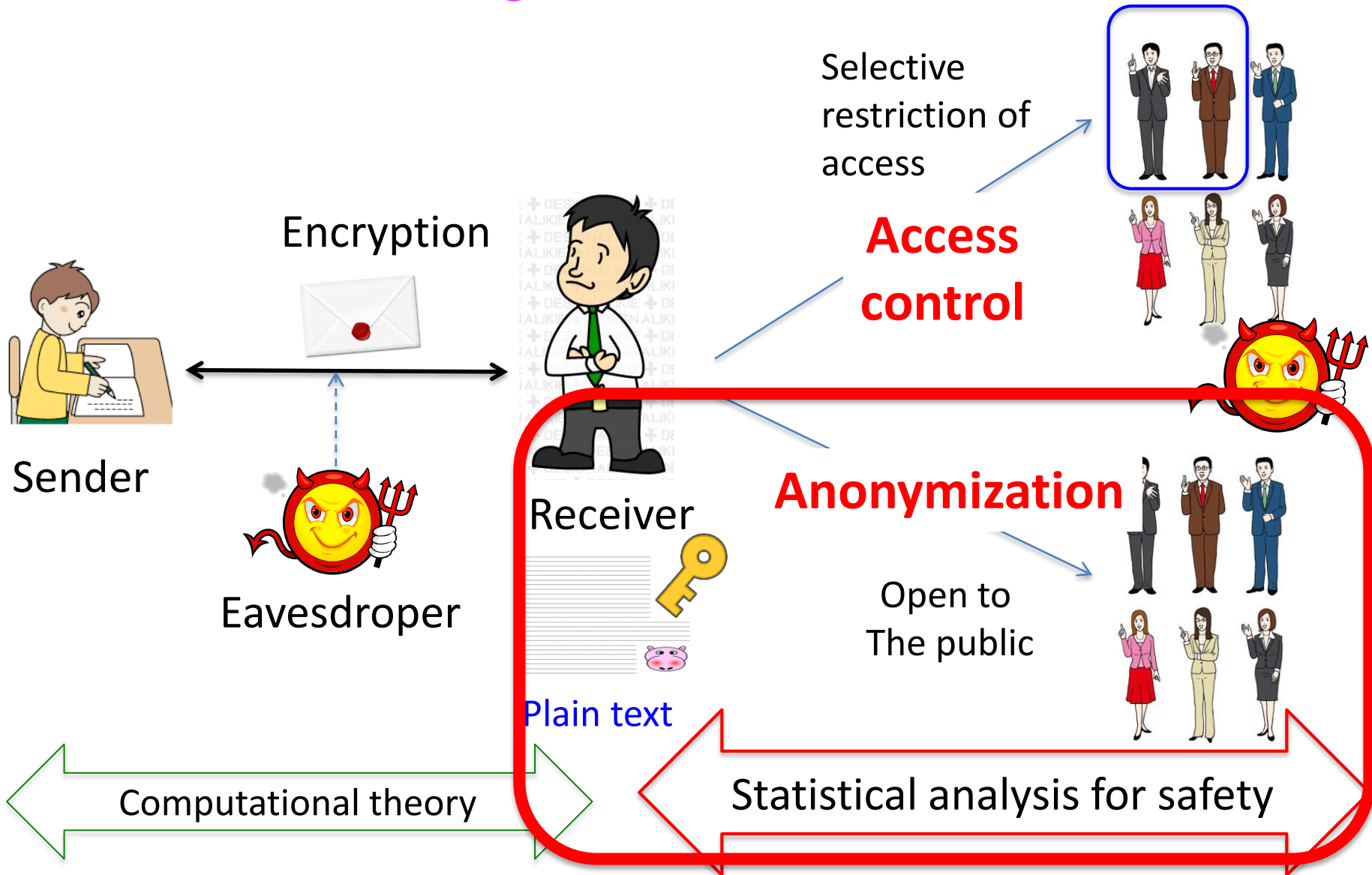
Anonymization and Location Privacy

Kazuhiro Minami

Institute of Statistical Mathematics

July 23, 2016

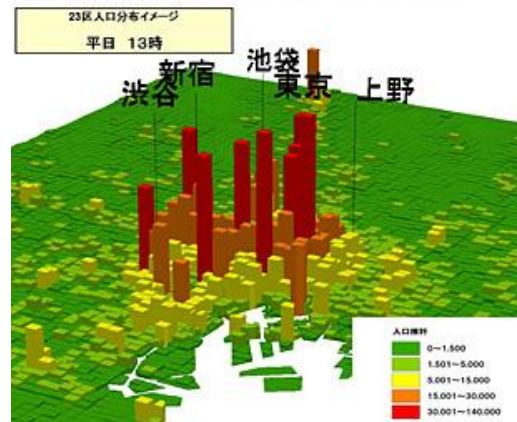
Privacy-Preserving Techniques for Sharing Useful Information



Location data is useful for many analytic purposes

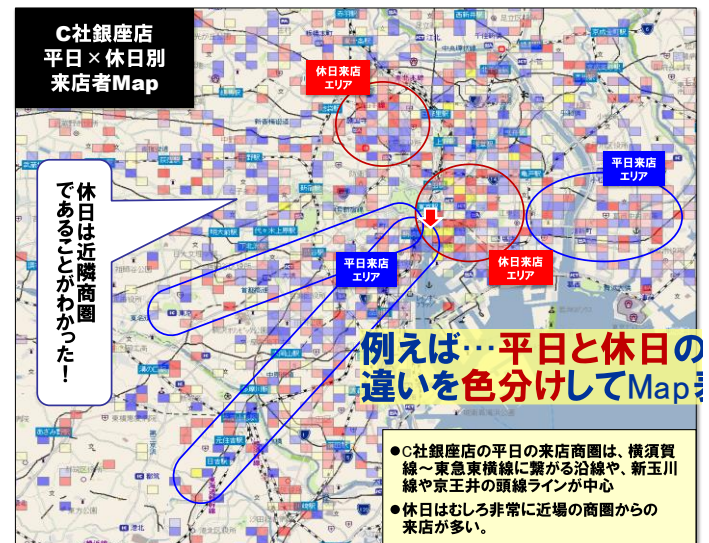
- Real-time traffic monitoring
- Dynamic population mapping
- Trade area analysis
- Disaster impact assessments

モバイル空間統計イメージ：東京23区周辺の人口分布



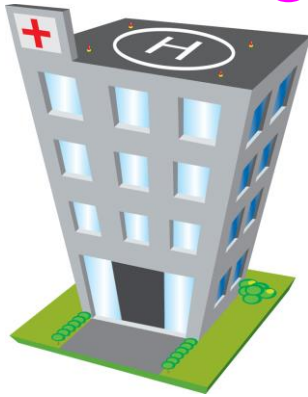
NTT Docomo

Mobile Spatial Statistics

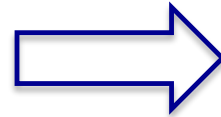


Dentsu Draffic

However, there is big concern on location privacy



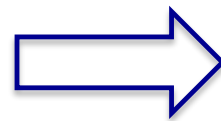
Hospital



Illness



Cafe

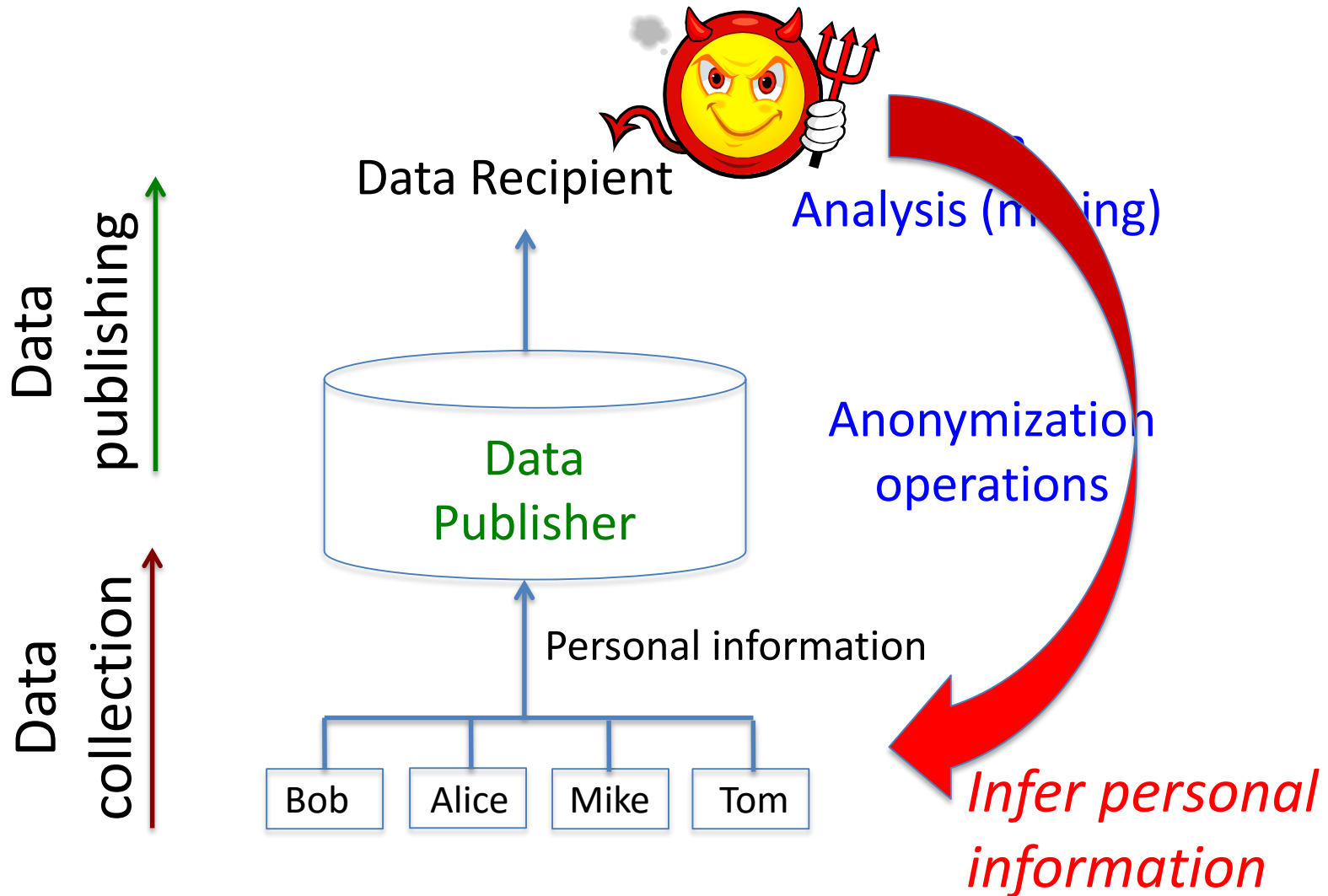


Laziness

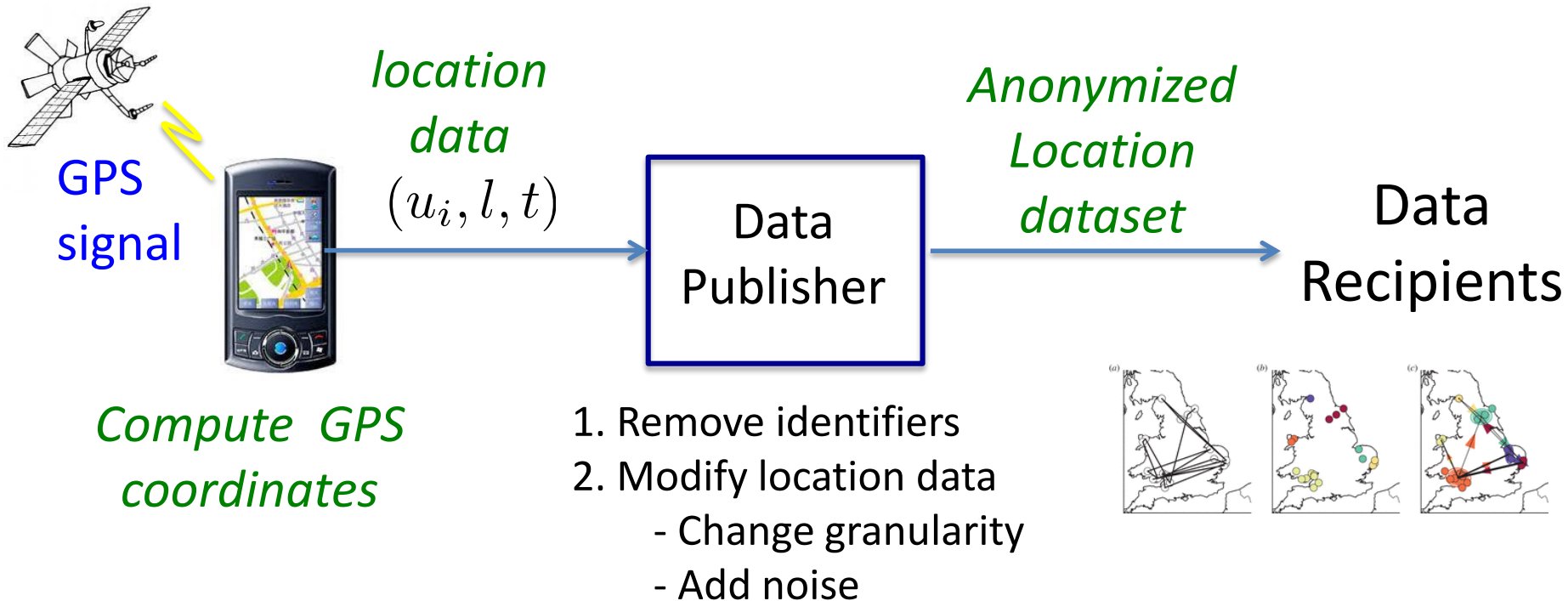


Secret meeting

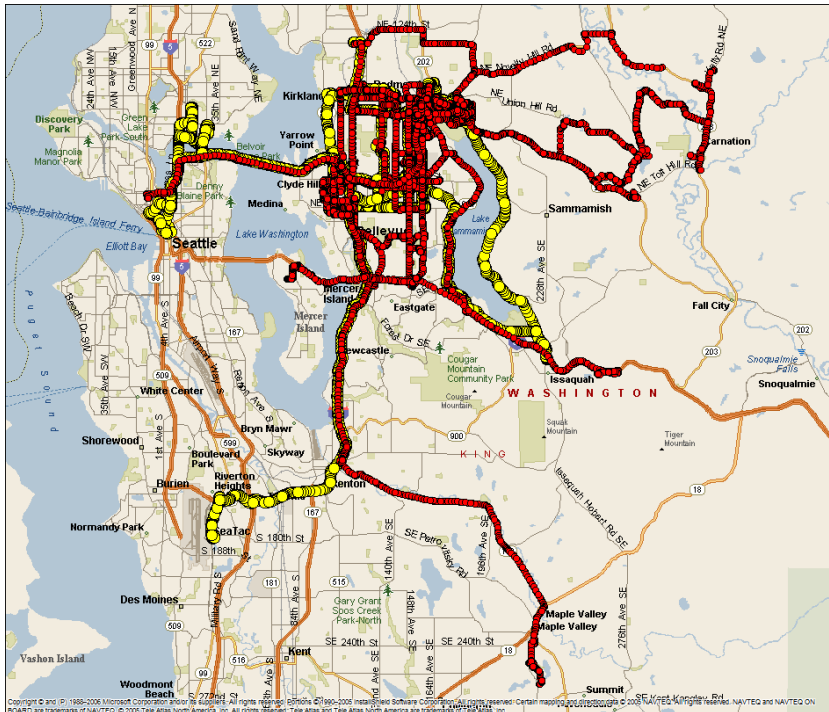
This is a special problem of Privacy-preserving data publishing (PPDP)



System model for location sharing services



Pseudonym-based approach



Pseudonymity

- Replace owner name of each point with untraceable ID
- One unique ID for each owner

Example

- “Larry Page” → “yellow”
- “Bill Gates” → “red”



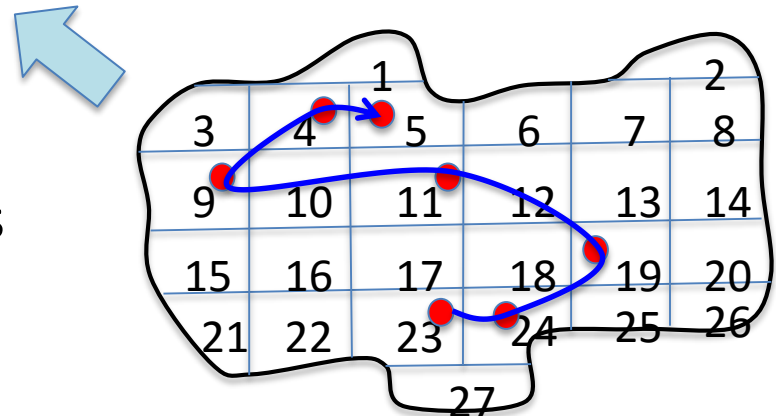
Location Traces

- Sequences of location IDs with a timestamp

July 20, 2016

氏名	8:0	8:30	9:00	9:30	10:00	10:30	11:00	
Tomoko	⁰ 1	5	4	8	12	15	9	
Gareth	10	15	24	14	21	20	19
Yoshiki	3	8	6	6	7	10	15	
Kazu	23	24	19	11	9	4	5	

Converted from
GPS coordinates

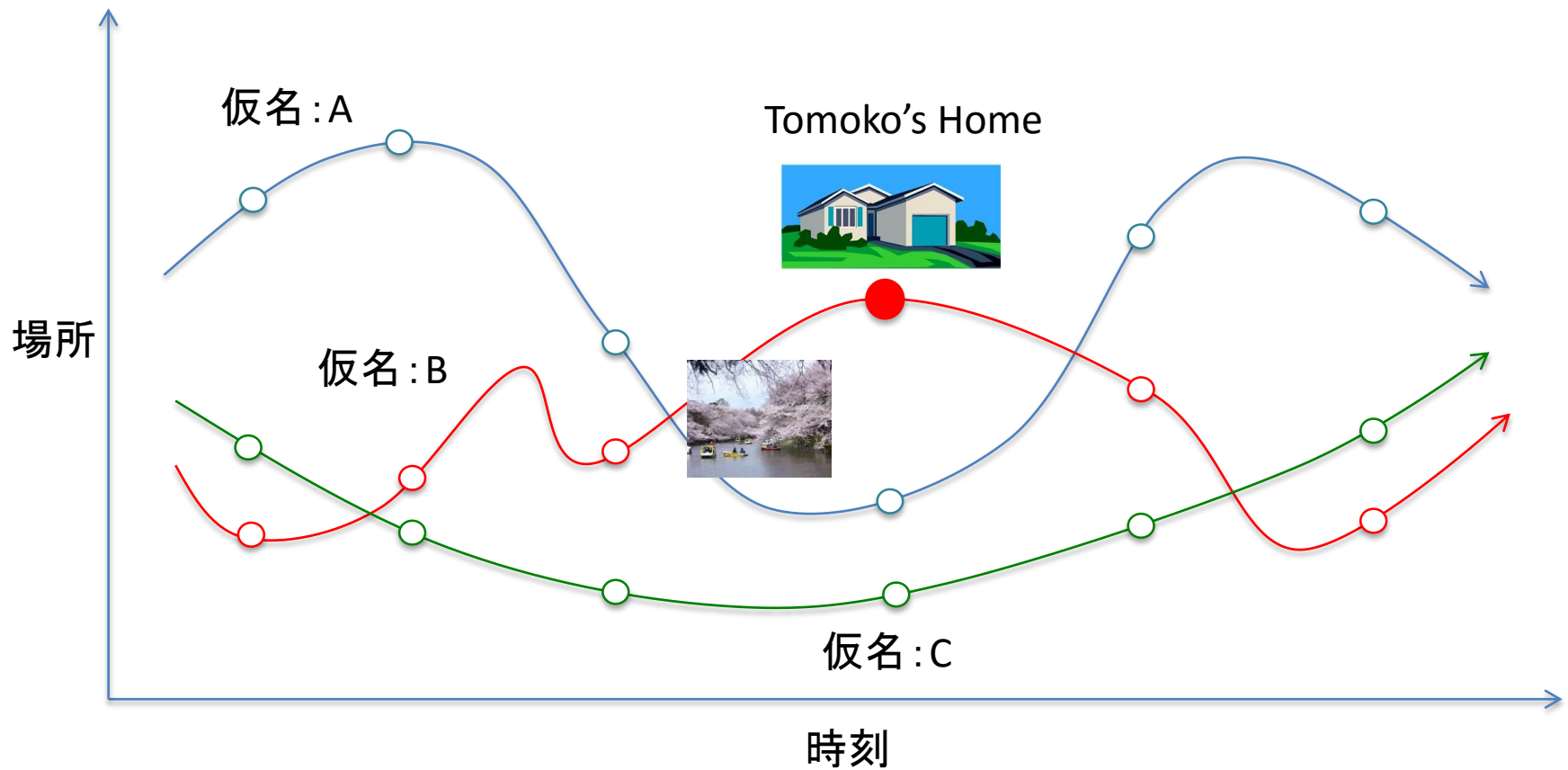


Is replacing names with pseudonyms sufficient?

July 20, 2016

Pseudonym	8:0	8:30	9:00	9:30	10:00	10:30	11:00	
A	0 1	5	4	8	12	15	9
B	10	15	24	14	21	20	19	
C	3	8	6	6	7	10	15	
D	23	24	19	11	9	4	5	

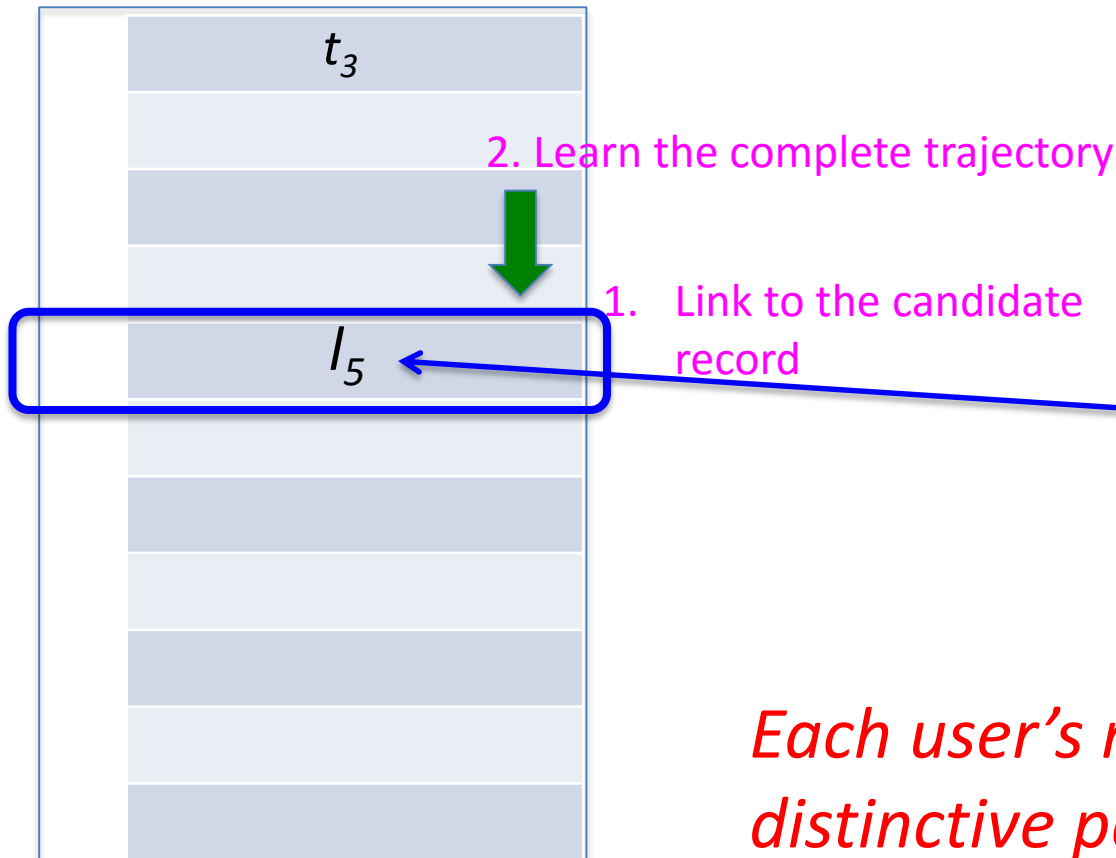
It's relatively easy to get additional information about your whereabouts



- Your home and office addresses
- Physical observations by accident or stalking

K-Anonymization for location data

Anonymized locations traces



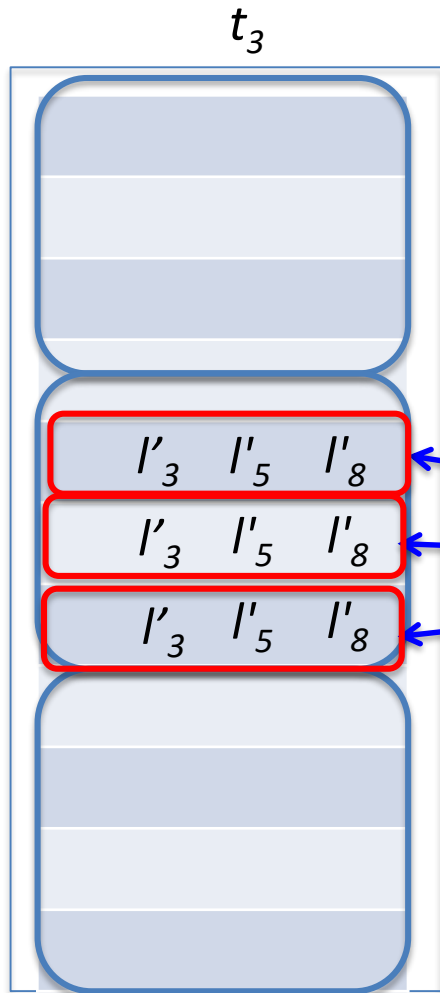
Partial information
on your location

	t_1	t_2	t_3	t_4
Bob				
Tom			l_5	
Ken	l_7			

Each user's movements have distinctive patterns

k-Anonymization of location data

Anonymized locations traces



1. Divide the table into groups of size k or more

Cannot narrow down candidate records less than k

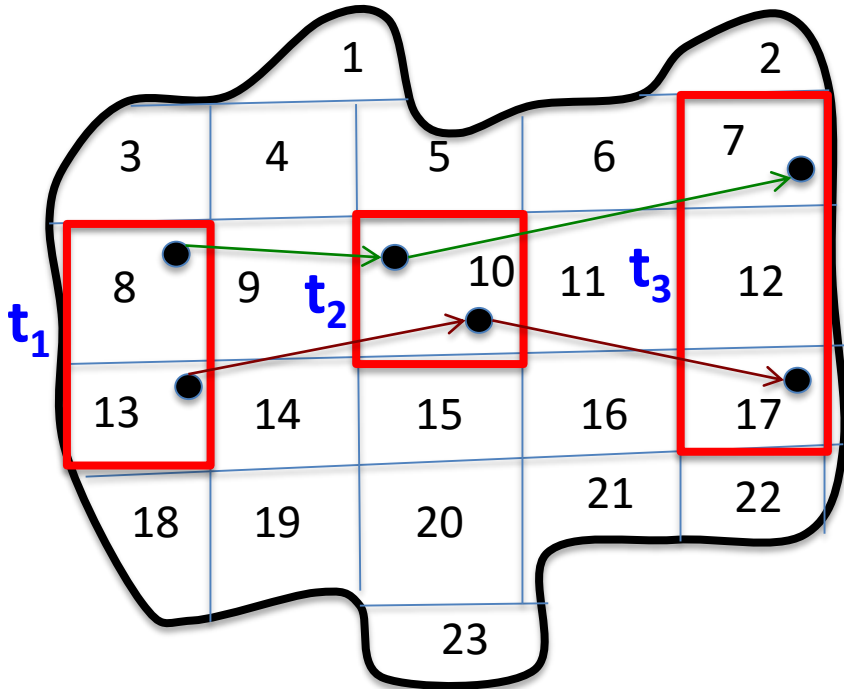
2. Generalize data to make the records identical



Partial information on your location

	t_1	t_2	t_3	t_4
Bob				
Tom			l_5	
Ken	l_7			

Example



Original table

User	t_1	t_2	t_3
Bob	{8}	{10}	{7}
Tom	{13}	{10}	{17}

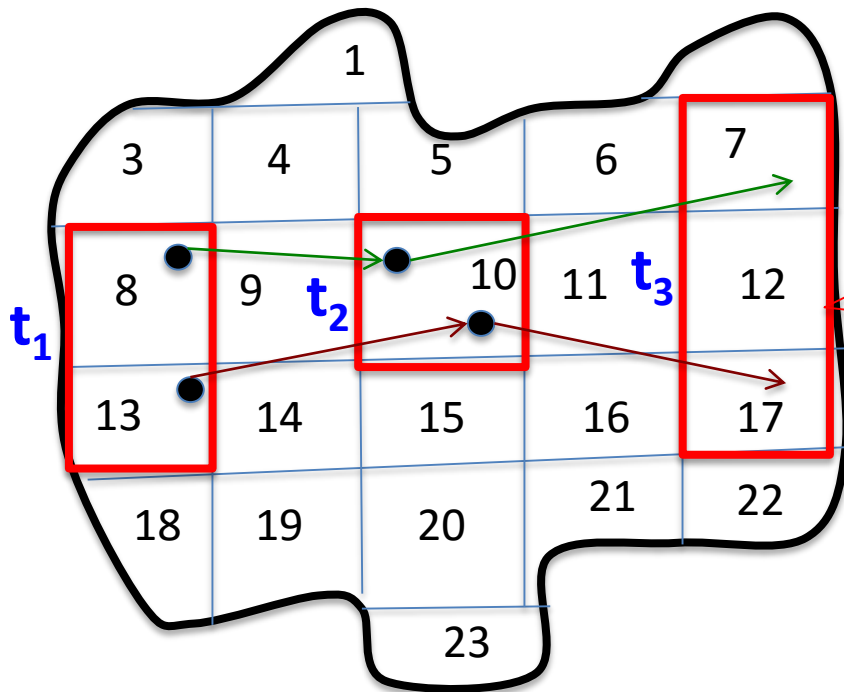
Generalization



2-Anonymous table

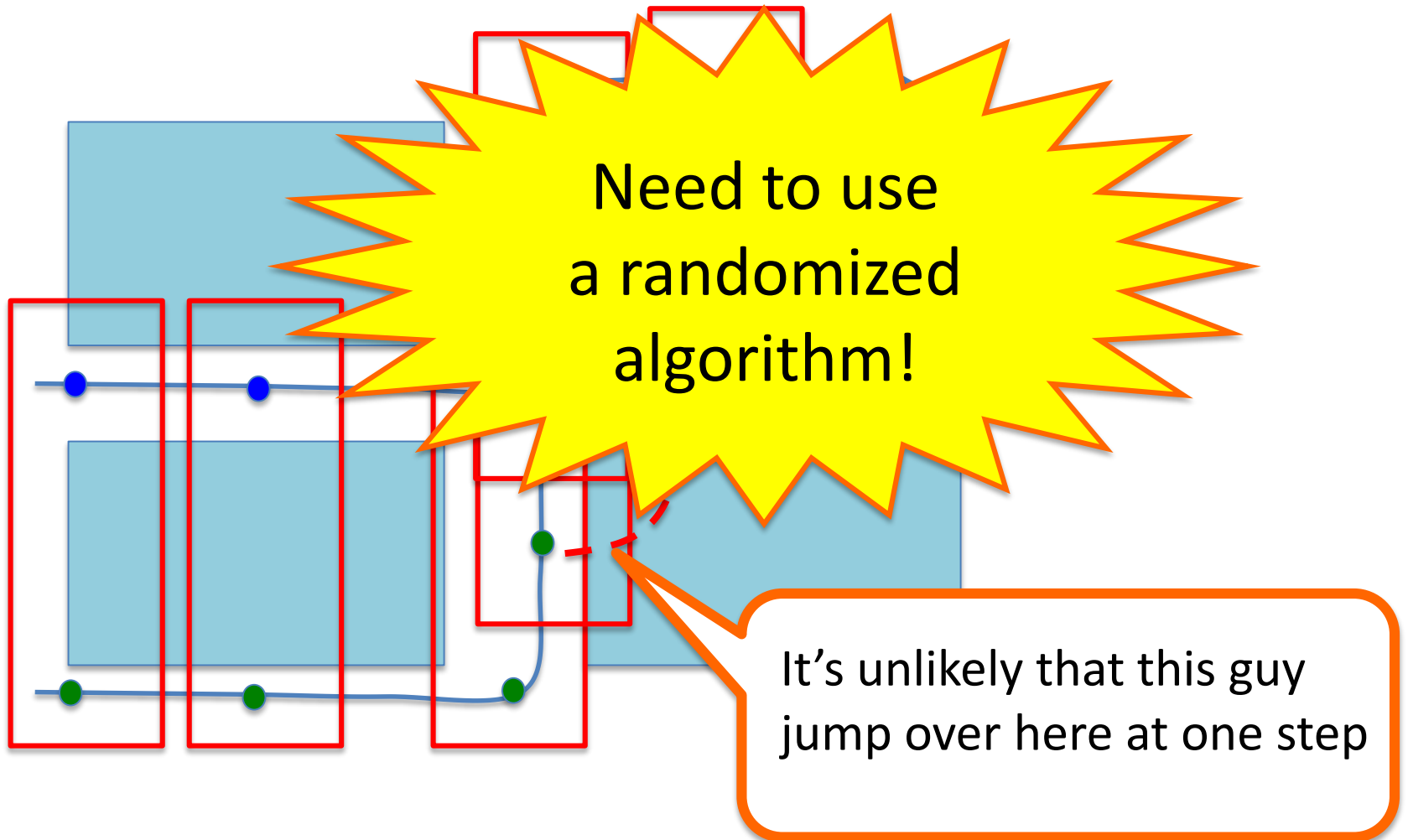
PID	t_1	t_2	t_3
A	{8, 13}	{10}	{7, 12, 17}
B	{8, 13}	{10}	{7, 12, 17}

What if an adversary knows the anonymization algorithm?



If the algorithm choose the smallest rectangle containing two users, the attacker knows one user is in grid 7 and the other in grid 17.

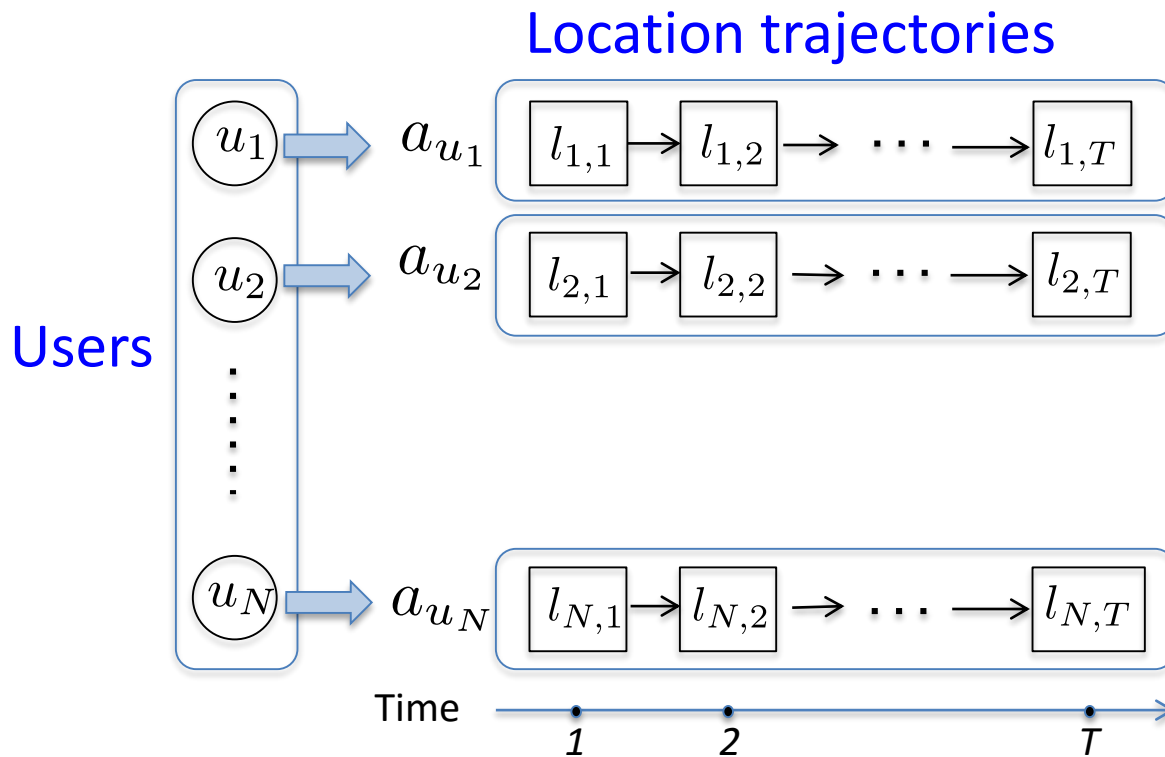
Knowing users' mobility patterns also helps



Q: How should we evaluate the safety of an anonymized location trajectories?

Problem setting [Shokri11]

- N users move around a geographical area of M regions $\mathcal{L} = \{l_1, l_2, \dots, l_M\}$ at discrete times in $\mathcal{T} = \{1, 2, \dots, T\}$



Suppose that each user's trajectory is anonymized independently, an anonymization procedure takes two steps:

1. Perturb each user's trajectory with a randomized algorithm

$$\langle a_{u_1}, a_{u_2}, \dots, a_{u_N} \rangle \implies \langle o_{u_1}, o_{u_2}, \dots, o_{u_N} \rangle$$

where for each i , $a_{u_i} = \langle a_{u_i}(1), a_{u_i}(2), \dots, a_{u_i}(T) \rangle$
 $\implies o_{u_i} = \langle o_{u_i}(1), o_{u_i}(2), \dots, o_{u_i}(T) \rangle$

2. Map user names to pseudonyms with the permutation function $\sigma: U \rightarrow U' = \{1, 2, \dots, N\}$

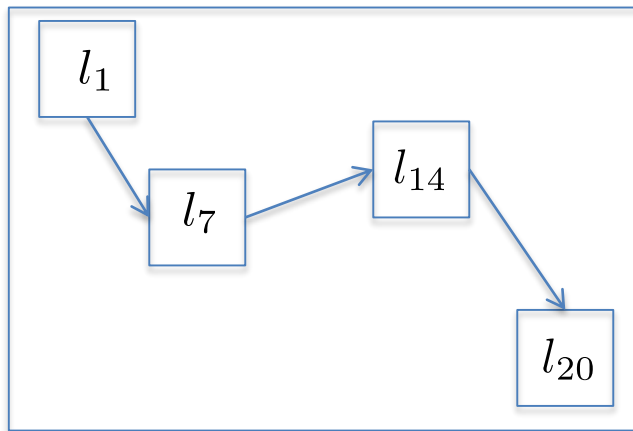
$$o_{u_i} = \langle o_{u_i}(1), o_{u_i}(2), \dots, o_{u_i}(T) \rangle$$
$$\implies \langle o_{\sigma(u_i)}(1), o_{\sigma(u_i)}(2), \dots, o_{\sigma(u_i)}(T) \rangle$$

1. Perturbation of location data

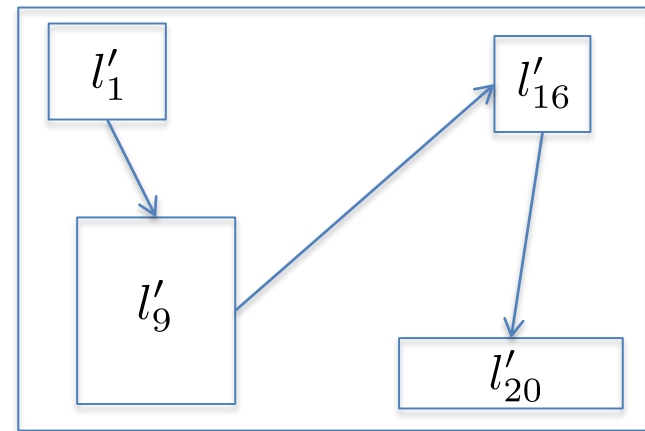
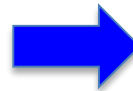
Perturb each user's trajectory with a randomized algorithm f

where $f_{a_u}(o_u) = Pr(O_u = o_u | A_u = a_u)$

- Adding noise
- Generalization (reduce precision)
- Omission (location hiding)

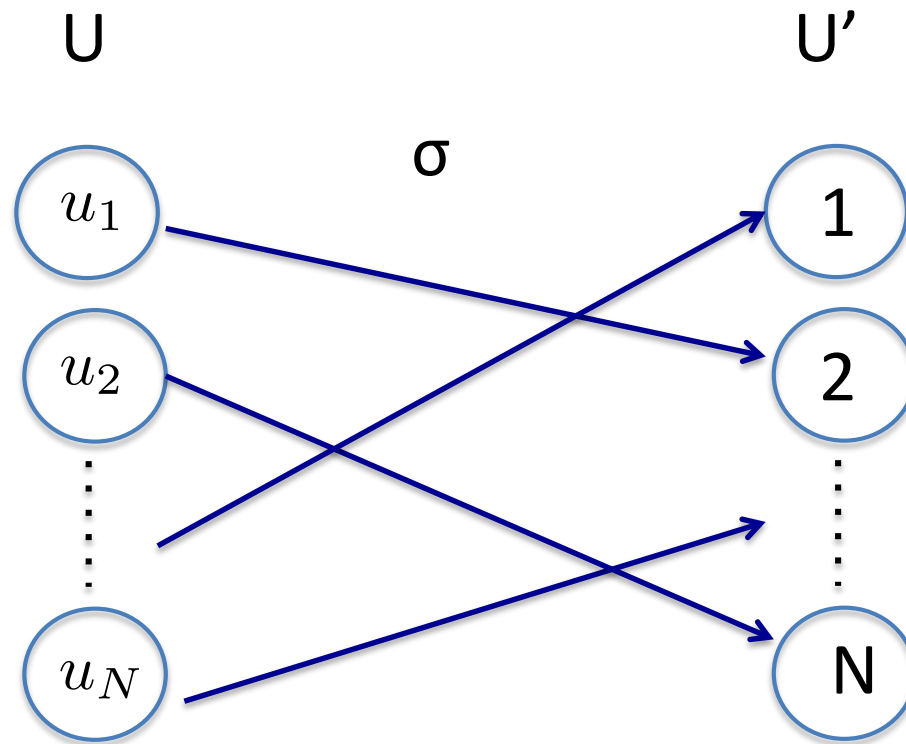


$l_i \in \mathcal{L} (= \{l_1, l_2, \dots, l_M\})$

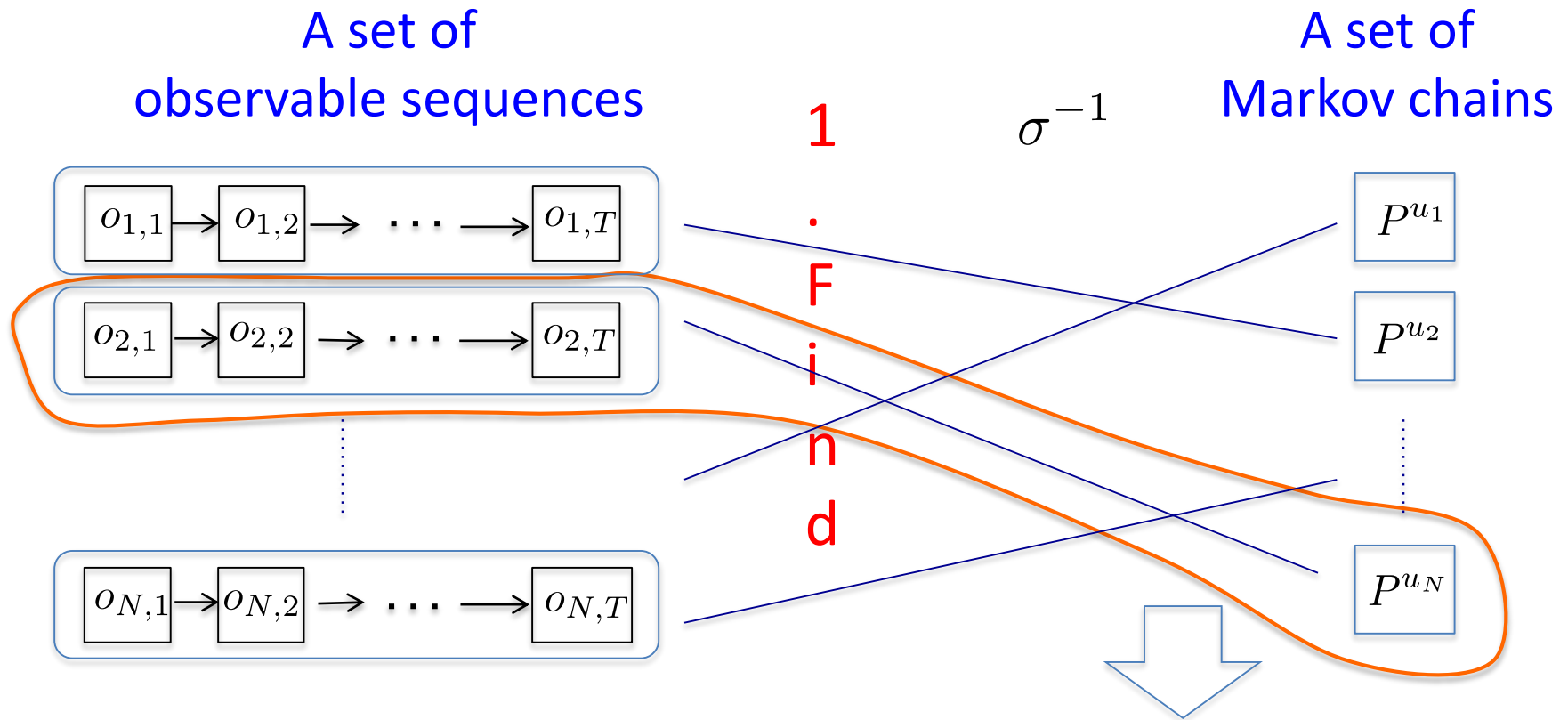


$o_i \in \mathcal{L}' (= \{l'_1, l'_2, \dots, l_{M'}\})$

2. Purmutation



So, an attacker has to do is to reverse this process with



2. For the pair of a target user, compute the prob. dist. of $Pr(a_u(t) = l | o_u, P^u)$ by HMM smoothing

Find an inverse of permutation σ

- Want $\operatorname{argmax}_{\sigma \in \Sigma} \prod_{u \in \mathcal{U}} \Pr(o_{\sigma(u)} | P^u)$ But $N!$ combinations
- Instead, for each pair (u, o_x) , compute $\Pr(o_x | P^u)$ with the forward algorithm

$$\Pr(o_x | P^u) = \sum_{l \in \mathcal{L}} \Pr(o_x(1), o_x(2), \dots, o_x(T), a_x(T) = l | P^u) = \sum_{l \in \mathcal{L}} \alpha_T(l)$$

- Consider an edge-weighted bipartite graph of traces and users and solve **the max weight assignment program** using the Hungarian algorithm

Localization attack:

Infer user u 's location at time t

- Compute $Pr(a_u(t) = l | o_u, P^u)$ with the forward-backward algorithm

$$\alpha_t(l) = Pr(o_x(1), o_x(2), \dots, o_x(T) | a_x(t) = l, P^u)$$

$$\beta_t(l) = Pr(o_x(t+1), o_x(t+2), \dots, o_x(T) | a_x(t) = l, P^u)$$

$$Pr(a_u(t) = l | o_u, P^u) = \frac{\alpha_t(l)\beta_t(l)}{Pr(o_u | P^u)}$$

Attacker's correctness as the measure of privacy risk

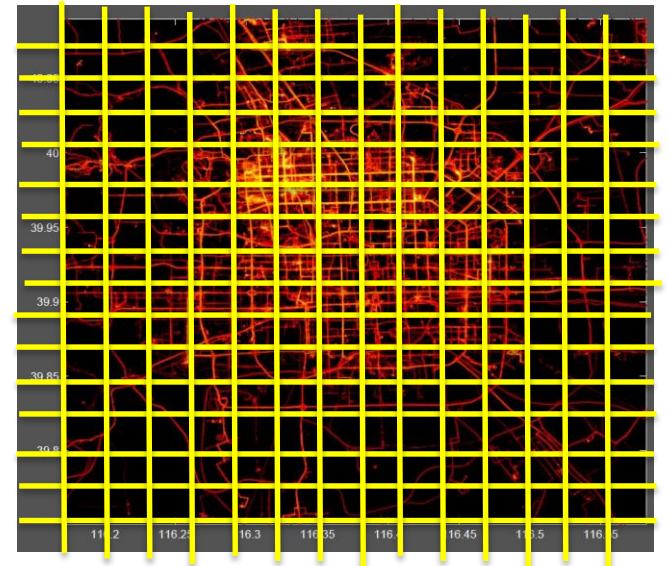
- Let \hat{l} be user u 's actual location at time t
- The probability of getting a correct answer would be a reasonable metrics

$$Pr(a_u(t) = \hat{l} \mid o_x, P^u)$$

Preliminary evaluations

Q: How many more non-sensitive locations we need to hide to protect the secrecy of private locations?

- Consider a rectangular region of 39×30 kilometers in Beijing, China
- Use top 10 users in terms of data points
- Divide the region into $140 \times 140 (=19,600)$ unit regions



- GPS dataset published by Microsoft Asia
- 178 users in the period of four years
- Logged every 1 – 5 seconds

Methods

1. Set the initial emission matrix based on users' private policies S such that

$$\text{If } l_i \in S, B_{ii} = 0, B_{i,\perp} = 1$$

$$\text{else } B_{ii} = 1, B_{i,\perp} = 0$$

2. Given a threshold δ , if the following is satisfied, exit

$$\text{For each } \hat{l} \in S : Pr(a_u(t) = \hat{l} \mid o_x, P^u) < \delta$$

Otherwise, randomly pick l_j not in S into S and set

$$B_{jj} = 0, B_{j,\perp} = 1$$

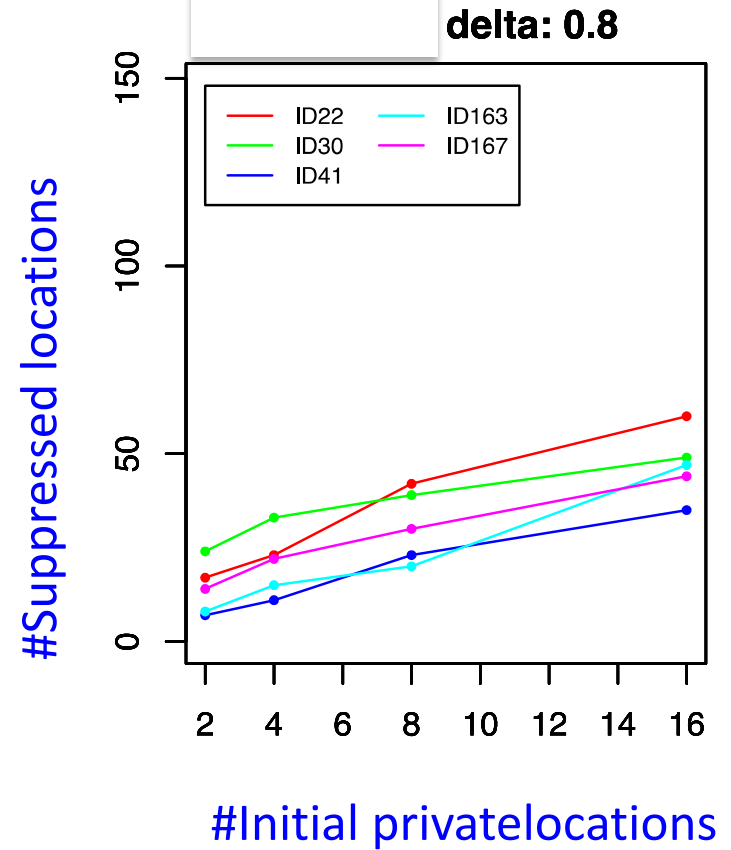
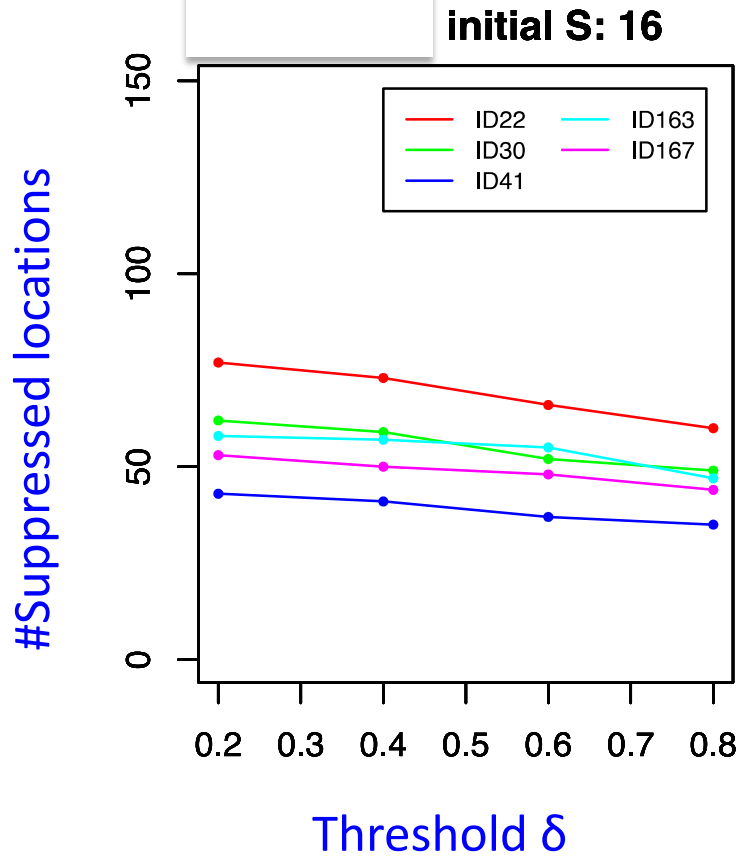
and repeat

We skipped the process of matching user IDs and pseudonyms

Initial private locations S_0

1. Pick two locations of an restaurant and a hospital, which was actually visited by users
 - China-Japan Friendship Hospital (N. latitude 39.97260, E. longitude 116.42072)
 - South Beauty Restaurant (N. latitude 39.99635, E. longitude 116.40360)
2. *Randomly* choose a given number of locations from the top most frequently visited locations

Results



However, what if the function f of perturbation depends on other records as in the case of k -Anonymization?

Summary

- Anonymizing location data has additional challenges due to spatial and temporal correlations among data points
- HMM provides a basic framework for analyzing privacy risks quantitatively
- However, further research is necessary to establish a methodology for designing a randomized function that produce observation traces