**University of Stuttgart**

Institute of Parallel and
Distributed Systems (IPVS)

Universitätsstraße 38
D-70569 Stuttgart

# On the Privacy of Frequently Visited User Locations

**Zohaib Riaz**, Frank Dürr, Kurt Rothermel

*International Conference on Mobile Data Management 2016 (MDM'16)*

15th June 2016

# Motivation

- Today's Location-sharing apps exploit *Location Semantics*

  - **App Examples:** Family locators, Friend finders, Geo-social networks

  - **Meaningful sharing:** (48.778786, 9.177867) → *Starbucks (Coffee Shop)*
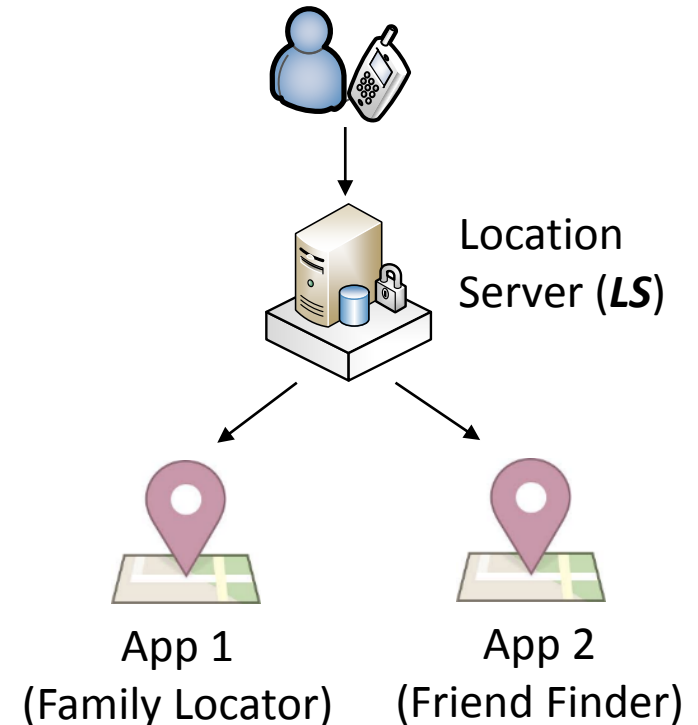  - Trivial to label any unlabeled trips (via Foursquare, Yelp etc.)

Church    Bar    Gym    Supermarket

# Motivation

- Today's Location-sharing apps exploit *Location Semantics*

  - **App Examples:** Family locators, Friend finders, Geo-social networks

  - **Meaningful sharing:** (48.778786, 9.177867) → *Starbucks (Coffee Shop)*
  - Trivial to label any unlabeled trips (via Foursquare, Yelp etc.)



- *Privacy threat:* Prolonged location sharing reveals <u>visit-frequency profile</u>!

  - *"A person who knows **all of another's travels** can deduce whether he is **a weekly church goer, a heavy drinker, a regular at the gym**,…— <u>and not just one such fact about a person, but all such facts</u>."* [United States v. Jones]
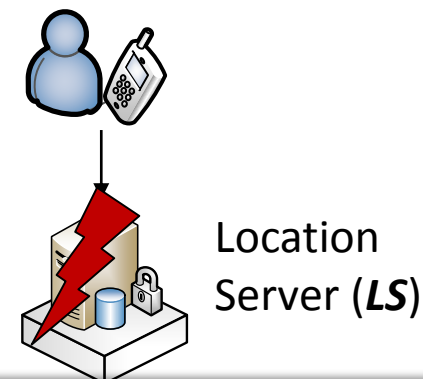
# Motivation

- Typical Location-sharing apps rely on backend *Location Server Infrastructure*

  ○ LSs store and manage user positions

  ○ Applications query user positions from LSs

- **Can the LSs provider be trusted for the security of users' location data?**

Location Server (*LS*)

App 1
(Family Locator)

App 2
(Friend Finder)

# Motivation

- Typical Location-sharing apps rely on backend *Location Server Infrastructure*

  ○ LSs store and manage user positions

  ○ Applications query user positions from LSs

- **Can the LSs provider be trusted for the security of users' location data?**

→ **No service provider can guarantee that personal information is safe!**

  ○ LS become single-point-of-failure w.r.t. privacy

*"The alarming part is that the information is so concentrated,"*



Location Server (**LS**)

**The Washington Post**

*eBay* *asks 145 million users to change passwords after data breach (2014)*

------------------------------------------

**REUTERS**

*Database of 191 million U.S. voters exposed on Internet: researcher (2015)*

------------------------------------------

**THE WALL STREET JOURNAL.**

*Twitter: Passwords Leaked for **Millions of Accounts** (6 days ago!)*

# Contributions

- *A study of real-world check-ins dataset* to show that ***frequent locations*** pose a serious privacy threat (***next 3 slides …***)

- *An approach to protect frequent locations* while avoiding a single-point-of-failure in the LS infrastructure

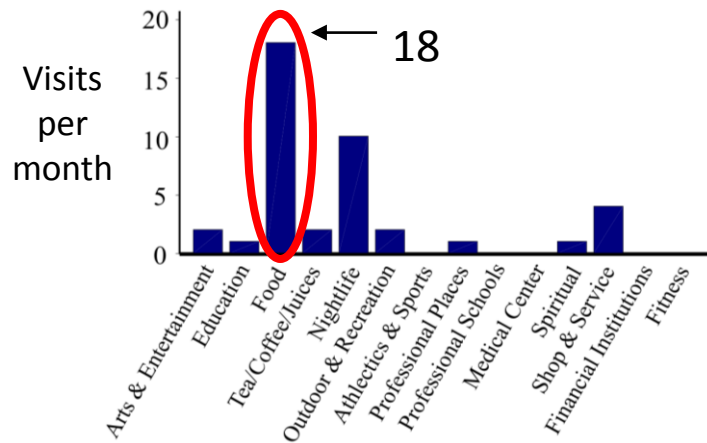- *Evaluation of the approach* for achieved Privacy and Quality-of-Service (QoS) for location-sharing apps.

# Study of Check-in Dataset: *Preprocessing*

- **Goal:** Show that visit-frequency information poses a privacy threat

- **Dataset:** *22,506,721* **Geo-tagged tweets** provided by Cheng et al. 2011

- **Selected user *Population*:** criteria
  - \>= 1 location check-in per day
  - \>= 30 days of reported location data
  - **10,306 users selected**

- Venue information, e.g., category, retrieved using Foursquare's free API

| No. | Category |
|-----|----------|
| 1 | Arts & Entertainment |
| 2 | Education |
| 3 | Food |
| 4 | Tea/Coffee/Juices |
| 5 | Nightlife |
| 6 | Outdoor & Recreation |
| 7 | Athletics & Sports |
| 8 | Professional Places |
| 9 | Professional Schools |
| 10 | Medical Center |
| 11 | Spiritual |
| 12 | Shop & Service |
| 13 | Financial Institutions |
| 14 | Fitness |

# Study of Check-in Dataset: *Analysis*
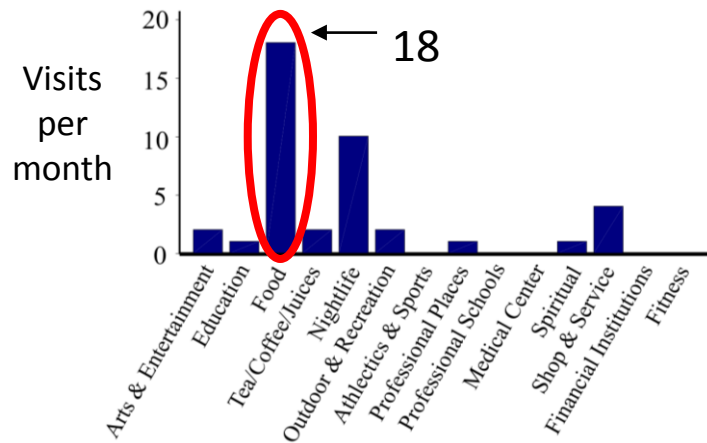


**Step 1:** Get frequency-profile

Population's Frequency Distribution **(Food)**

**Percentile Rank:**

$$PR(18) = 80th$$

**Step 2:** Determine Significance of visit-frequency

**Output:** Frequency rank-profile

# Study of Check-in Dataset: *Analysis*



18

Visits per month

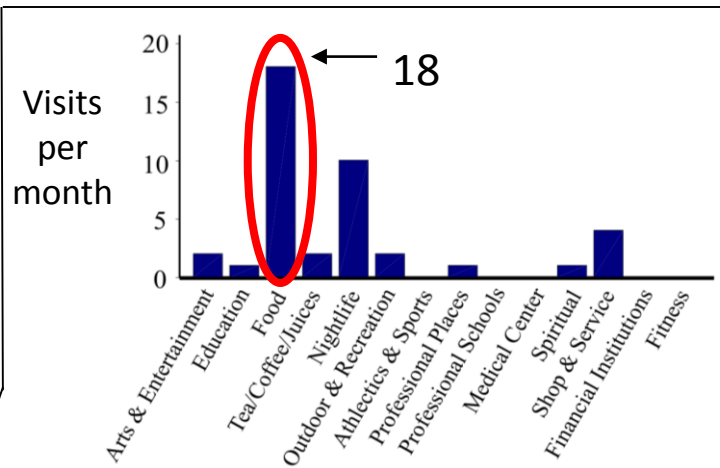**Step 1:** Get frequency-profile

Population's Frequency Distribution **(Food)**

No. of users

**Percentile Rank:**
$$PR(18) = 80th$$

Visits/Month

**Step 2:** Determine Significance of visit-frequency

$$th_{critiality} = 70$$

Percentile Ranks

**Critical locations:** *visited with <u>critically high frequencies</u>*
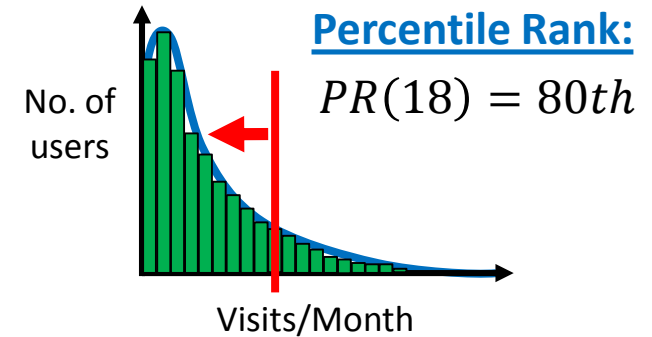
**Output:** Frequency rank-profile
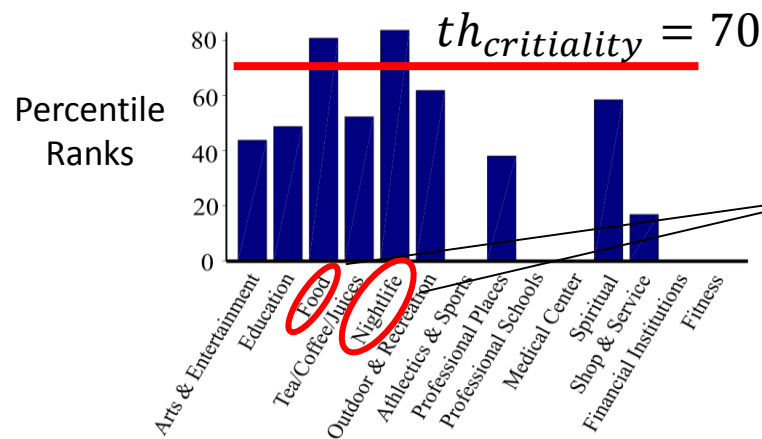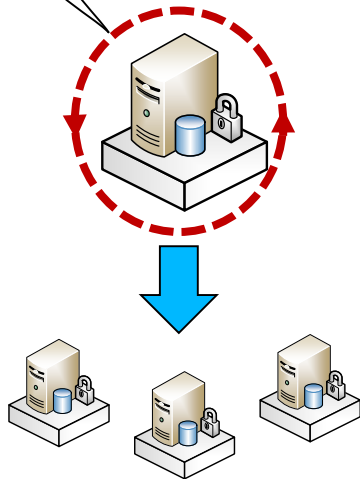
6

# Study of Check-in Dataset: *Analysis*



**Step 1:** Get frequency-profile

Population's Frequency Distribution **(Food)**

**Percentile Rank:**

$$PR(18) = 80th$$

No. of users

Visits/Month

**Step 2:** Determine Significance of visit-frequency

$$th_{critiality} = 70$$

Percentile Ranks

**Critical locations:** *visited with <u>critically high frequencies</u>*

**Output:** Frequency rank-profile

# Study of Check-in Dataset: *Analysis*



18

**Step 1:** Get frequency-profile

Population's Frequency Distribution **(Food)**

**Percentile Rank:**

$$PR(18) = 80th$$

No. of users

Visits/Month

**Step 2:** Determine Significance of visit-frequency
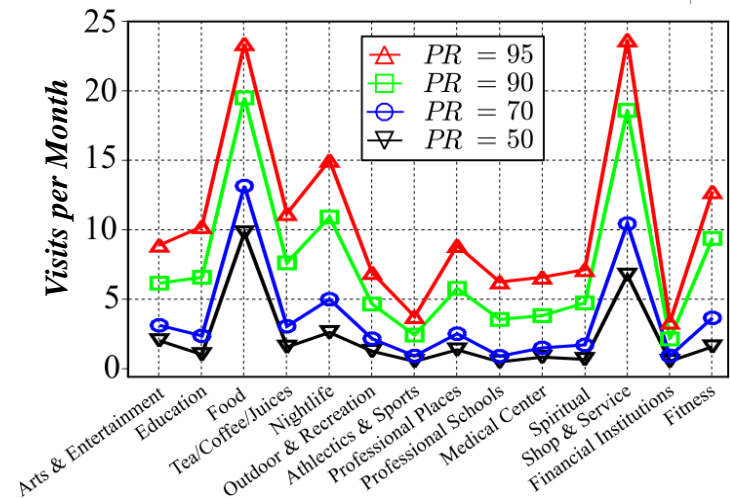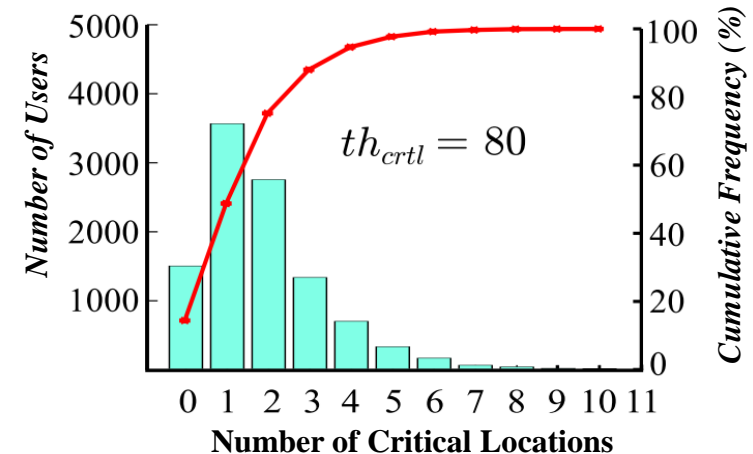
$$th_{critiality} = 70$$

Percentile Ranks

**Critical locations:**
*visited with critically high frequencies*

**Output:** Frequency rank-profile

# Study of Check-in Dataset: *Evidence of Privacy Threat!*
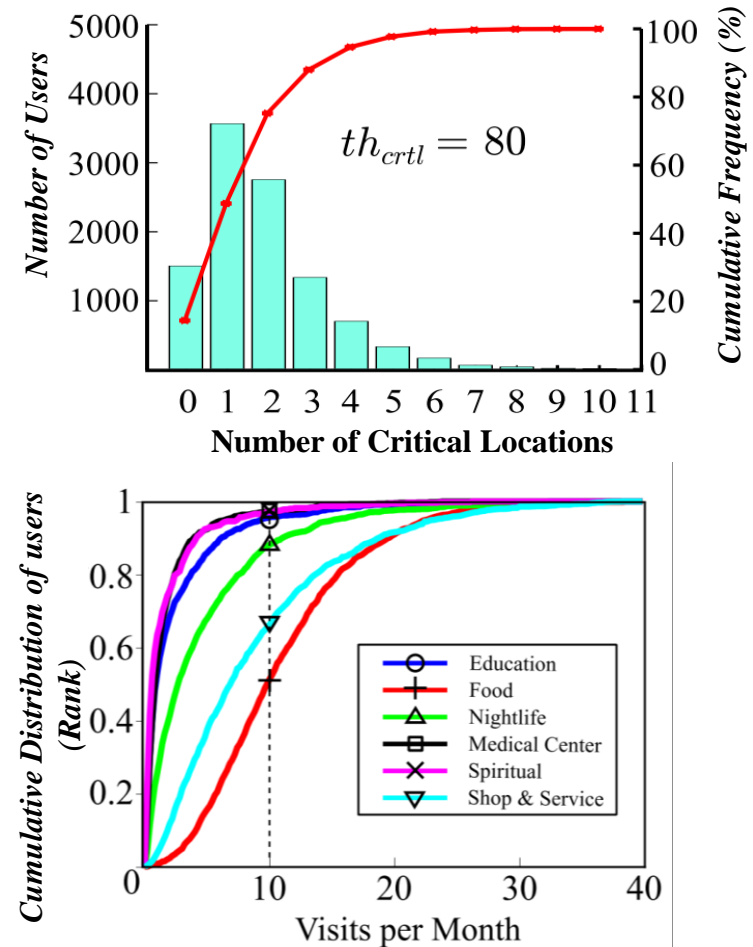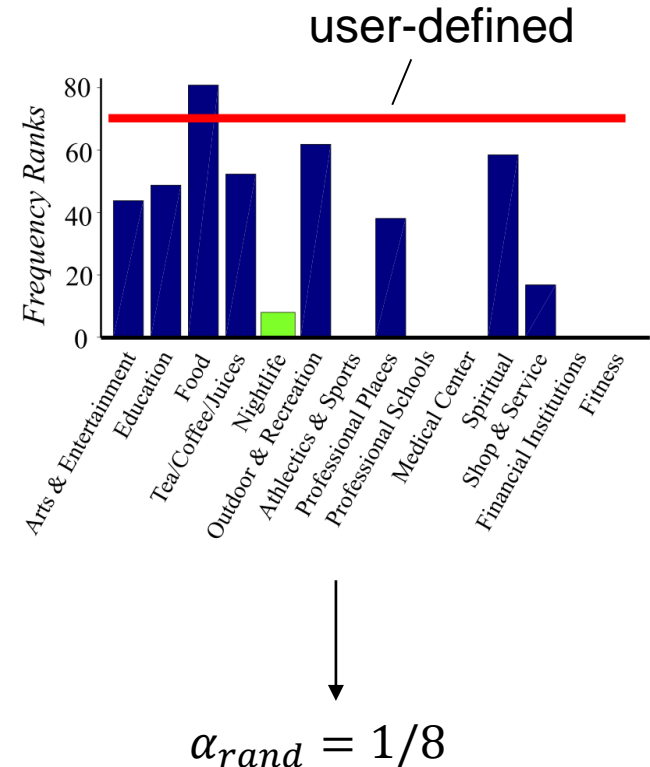
- **Critical locations are prevalent!**

  - ~85% users have <u>at least 1</u> critical location

  - ~50% users have <u>2 or more</u> critical locations

- **Visiting characteristics of locations**
  - Same percentile rank → different frequencies for diff. categories
  - High percentile-rank → a reasonable measure of user interest



$$th_{crtl} = 80$$

# Study of Check-in Dataset: *Evidence of Privacy Threat!*

- **Critical locations are prevalent!**
  - ~85% users have <u>at least 1</u> critical location
  - ~50% users have <u>2 or more</u> critical locations

- **Visiting characteristics of locations**
  - Same percentile rank → different frequencies for diff. categories
  - High percentile-rank → a reasonable measure of user interest

→ **We assume that "Frequency ⇔ Rank" relationship is publicly known**



$th_{crtl} = 80$

Number of Users / Cumulative Frequency (%) vs Number of Critical Locations



Cumulative Distribution of users (Rank) vs Visits per Month

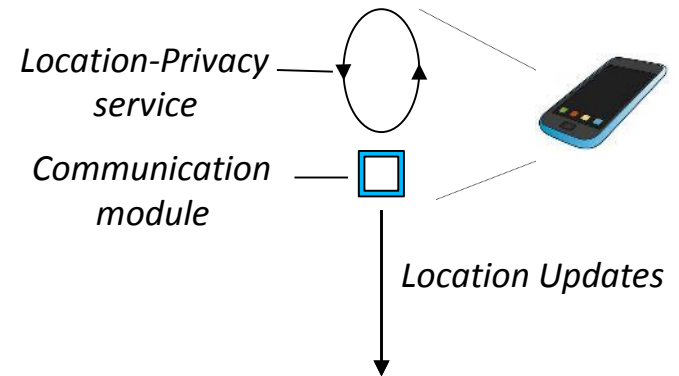Legend: Education, Food, Nightlife, Medical Center, Spiritual, Shop & Service

# Problem Statement

- **Privacy-preferences:** *(Persona, App)* pairs

- **User-Personas:** Define location categories whose criticality may be revealed e.g.:

  - *"Friends"* → {all} \ {Medical}

  - *"Colleagues"* → {all} \ {Spiritual, Medical, Night-life}

- **Privacy Requirement:** Implement privacy-preferences in location-sharing

  - Reveal unshared critical locations such that they *appear non-critical*

  - Avoid attacks to reveal critical locations

    - Attacker know our algorithm + additional knowledge

    - Baseline attack: ***random guess!***

user-defined



$$\alpha_{rand} = 1/8$$

# Architecture (1)

- **User's device:** a smart-phone

  ○ Runs **Location-Privacy service:**

    ▪ Executes our privacy algorithm

    ▪ Performs location updates to LSs

  ○ **Assumption:** Encrypted communication channel

*Location-Privacy service*

*Communication module*

*Location Updates*

# Architecture (2)

- **A set of Location Servers (LSs)**

  - from different *third-party providers*

    - Example: *Backendless, App42, Heroku etc.*

  - manage location updates

  - implement Access-control mechanism

- **Location Based Applications (Apps)**

  - Get access authorization to LSs from users

  - Access user location from LSs or subscribe for update notifications

  - May aggregate frequency-profile of user

  - User-profile' precision $\propto$ no. of accessible LSs



Location-Privacy service

Communication module

Access Rights

App 1

App 2

# Basic Privacy Algorithm (1)

**Actual *Rank* Profile**

1. **On-device determination of critical locations:**

   ◦ $S = \{s_1, \ldots, s_{14}\}$, set of location categories

   ◦ $f_u = \{f_{s_1}, f_{s_2}, \ldots, f_{s_{14}})$ and $r_u = \{r_{s_1}, r_{s_2}, \ldots, r_{s_{14}}\}$

   ◦ Critical locations: $C_u = s_i | r_{s_i} > th_{crtl}$

2. **Determine desired ranks**

   ◦ Deterministic new ranks → reversible by attacker

   ◦ **Randomized selection** of desired ranks

     ▪ Avoids advanced attacks!

**Desired *Rank* Profile**

**Research Group**

**"Distributed Systems"**

**University of Stuttgart**

**IPVS**

# Basic Privacy Algorithm (2)

## 3. Enforce desired ranks for all $s_i \in C_u$:

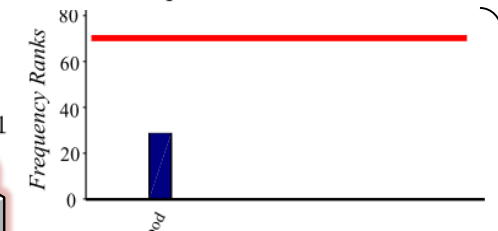- Divide trips for $s_i$ among LSs
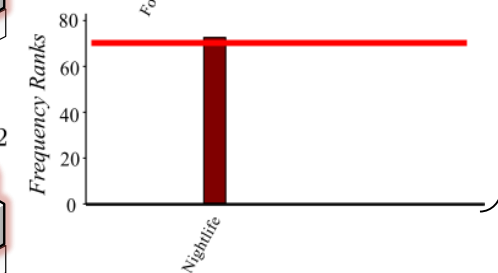- $LS_0$ hosts a safe-profile of user



Actual Rank Profile

Location-Privacy service

Protected Rank profiles at LSs

$LS_0$ — Shared with all Apps
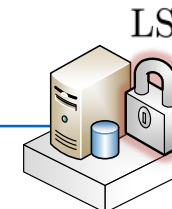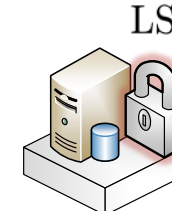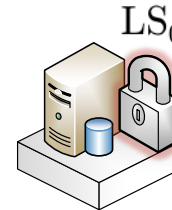
$LS_1$

$LS_2$ — Shared as per Personas of Apps

# Basic Privacy Algorithm (2)
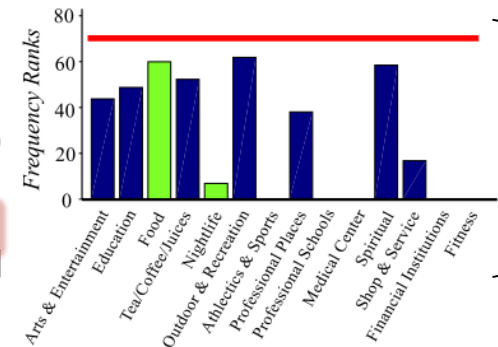
## 3. Enforce desired ranks for all $s_i \in C_u$:

- ○ Divide trips for $s_i$ among LSs
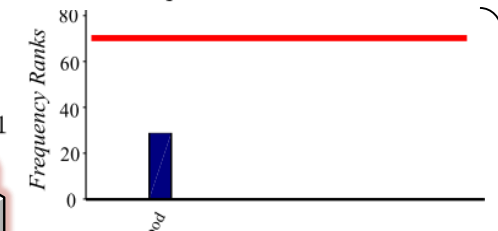
- ○ $LS_0$ hosts a <u>safe-profile</u> of user

Protected *Rank* profiles at LSs



Location-Privacy service

*User-aware Attacker: can monitor <u>Network-traffic statistics</u> → knows update timings!*

# Basic Privacy Algorithm (2)

**3.** **Enforce desired ranks for all $s_i \in C_u$:**

- Divide trips for $s_i$ among LSs

- $LS_0$ hosts a <u>safe-profile</u> of user

<u>***Population-aware Attacker:*** *Also possesses location information from the other users*</u>
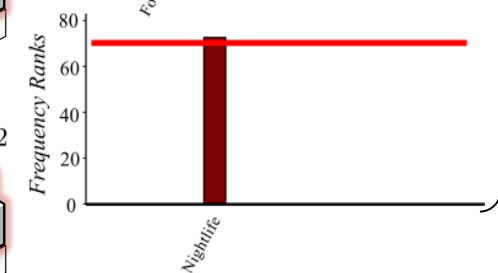
*Location-Privacy service*

Protected *Rank* profiles at LSs



Shared with all Apps

Shared as per *Personas* of Apps

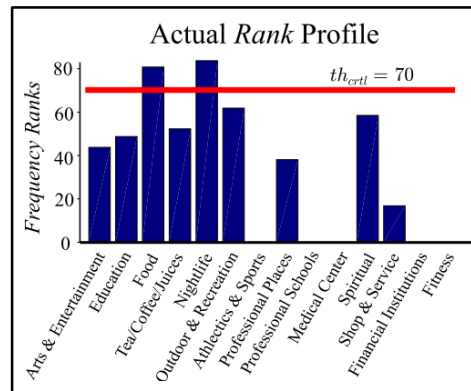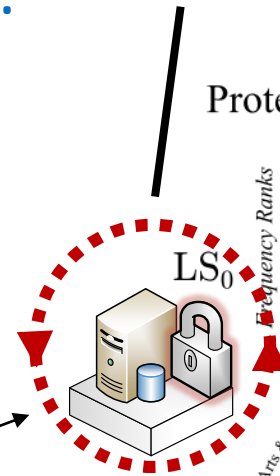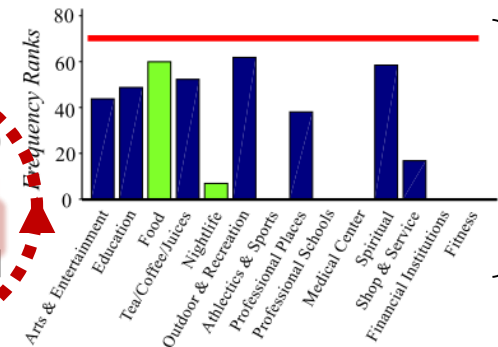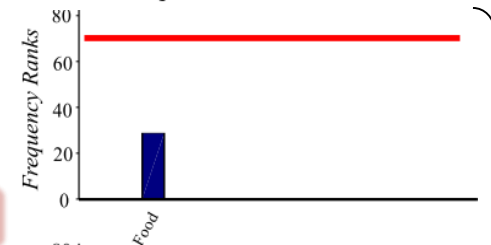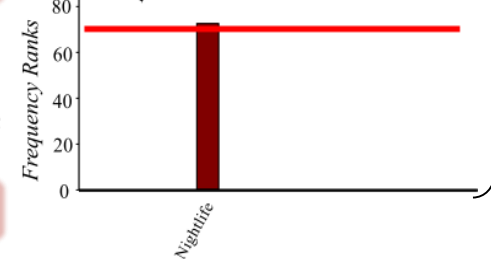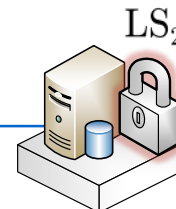<u>***User-aware Attacker:*** *can monitor <u>Network-traffic statistics</u>* → *knows update timings!*</u>
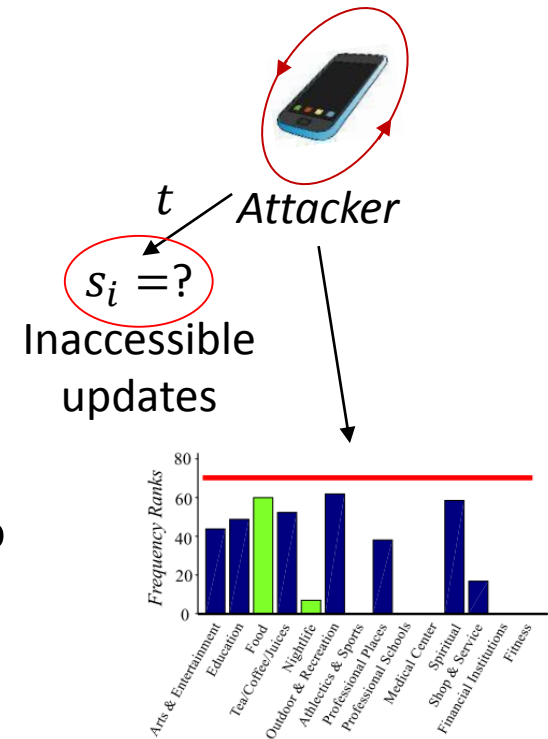
# Advanced Privacy Algorithm (1): Against _user-aware_ Attacker

- _Attacker_: has access to a few LSs
  - Knows timings of inaccessible updates

- Trail of location updates → **_Mobility Model_ $\Omega$**

| $t_1$ | $t_2$ | $t_3$ | $t_4$ | ... | $t_{150}$ | $t_{151}$ | $t_{152}$ | $t_{153}$ | ... |
|-------|-------|-------|-------|-----|-----------|-----------|-----------|-----------|-----|
| A | F | E | N | | ? | E | A | ? | |

$t$ /Attacker

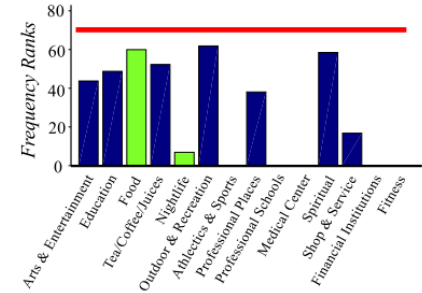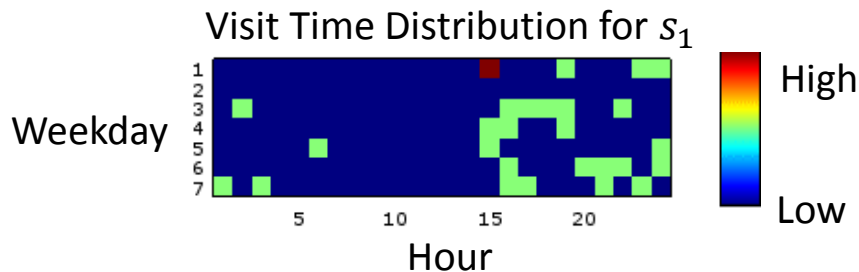$s_i = ?$

Inaccessible updates

- **Attack inaccessible updates:** Maximize $P(s_i|t)$ for $s_i$ to predict visited location using $\Omega$

- Bayes theorem: $P(s_i|t) = \dfrac{P(t|s_i)\, P(s_i)}{P(t)}$

**Prior**: Changed by our algorithm (<u>unreliable</u>)

**Normalizer**: constant for all $s_i$ (<u>unimportant</u>)

**Likelihood** of visiting time $s_i$ at time $t$ over all possible times $T$

# Advanced Privacy Algorithm (2): Defense

- **Defense:** Generate *fake events* for each location as if it were critical!

    ○ Fake events → *garbage data* → *discarded by LSs!*

    ○ **Desired effect:** Rank of all locations should "appear" equal

- **Algorithmic steps:**

    1. Keep track of temporal likelihood of each category

    2. Accordingly schedule enough fake events to meet maximum rank in the rank-profile



Visit Time Distribution for $s_1$

+

**fake**/inaccessible updates

**"Attacker's View"**

# Evaluation: *Population-aware* Attacker Model

- **Attacker:** Aims to find all critical locations

  1. knows '$k$' out of '$n$' critical locations from **authorized or compromised LSs**

  2. Knows **correlations among visit-frequencies of different location categories** *(Acquired from the population)*

- **Frequency-correlation attack**

  - Learn correlations using **Machine Learning techniques**

  - **Data:** Frequency-profiles of 10,036 users



Classification → Nightlife

Regression

# Evaluation: Privacy results for <u>Classification Attacks</u>

- **Classifiers:** *Random Forest* (RF) & *Support-Vector Machine* (SVM)

- **Training:**
  - On frequency-profiles with *one critical location*
  - 10-fold cross-validation

- **Results:**
  - Low classification accuracy: ***25%***

- **Repeated experiment:**
  - Added frequency-profiles with no critical locations!
  - Again, low accuracy for critical locations: ***22%***
  - ***High accuracy for non-critical: 87%***

Protected profile (altered)



non-critical profile (unaltered)

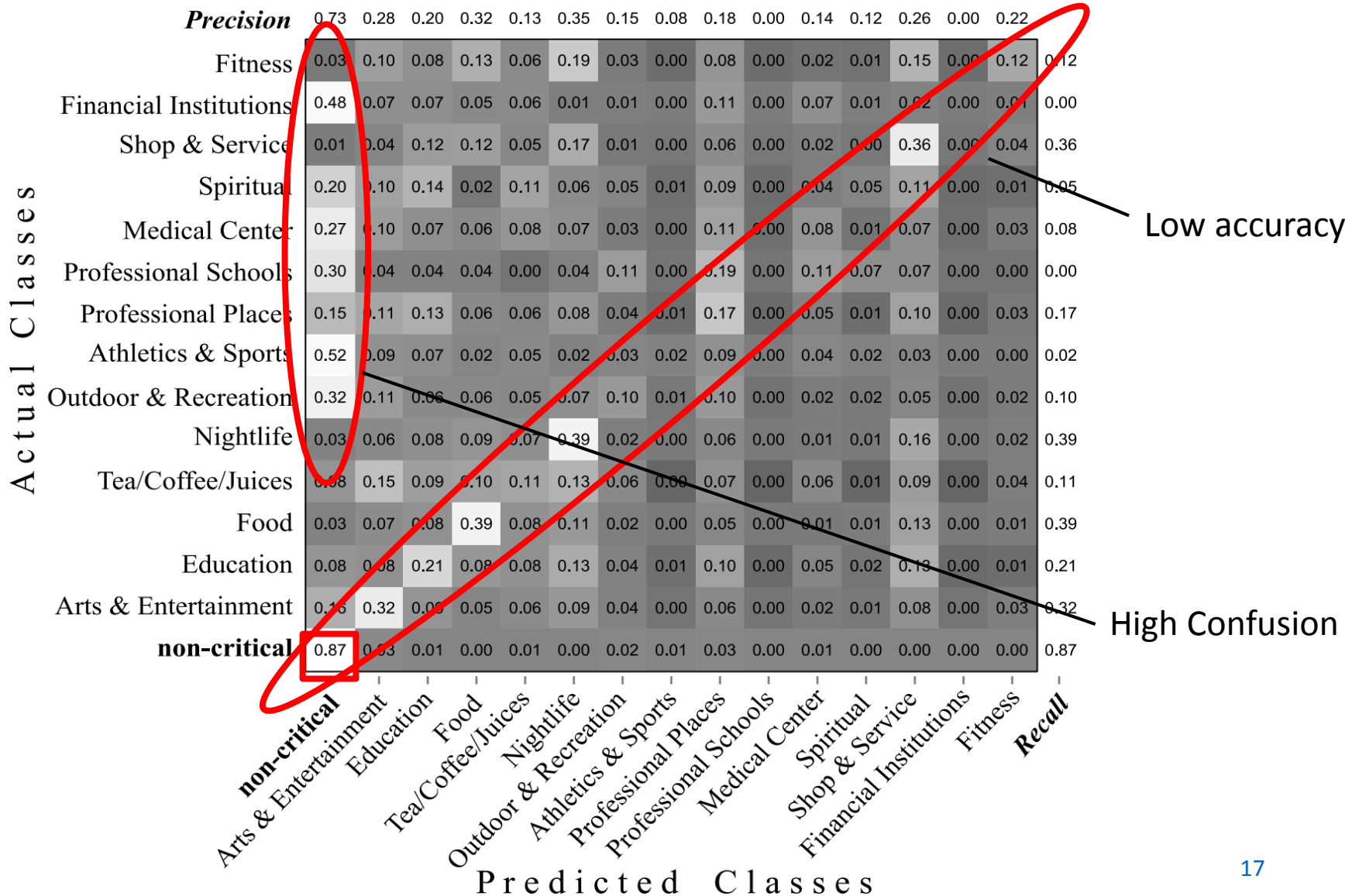# Evaluation: Privacy results for Classification Attacks

# Evaluation: Privacy results for <u>Regression Attacks</u>

- **Regression Models:** *RF, SVM and Gaussian Mixture Regression* (GMR)

    ○ Percentage prediction error: < 5% for each semantic location

    ○ Attack performance on <u>protected frequency-profiles</u>:



$P_{attack}(k)$ - probability of correct detection of a critical location when $k$ out of $n$ critical locations are already known
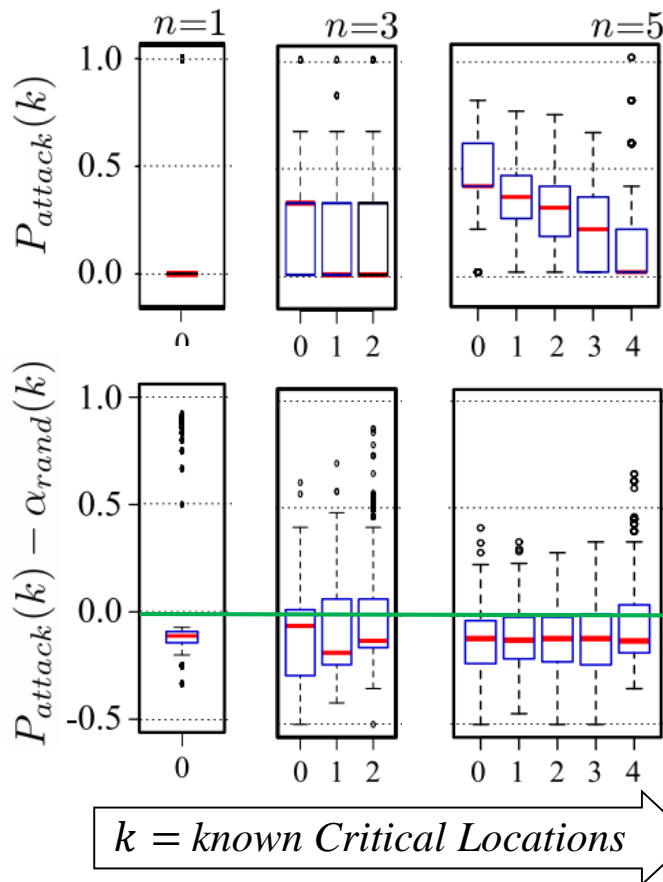
# Evaluation: Privacy results for Regression Attacks

- **Regression Models:** *RF, SVM and Gaussian Mixture Regression* (GMR)

  ○ Percentage prediction error: < 5% for each semantic location

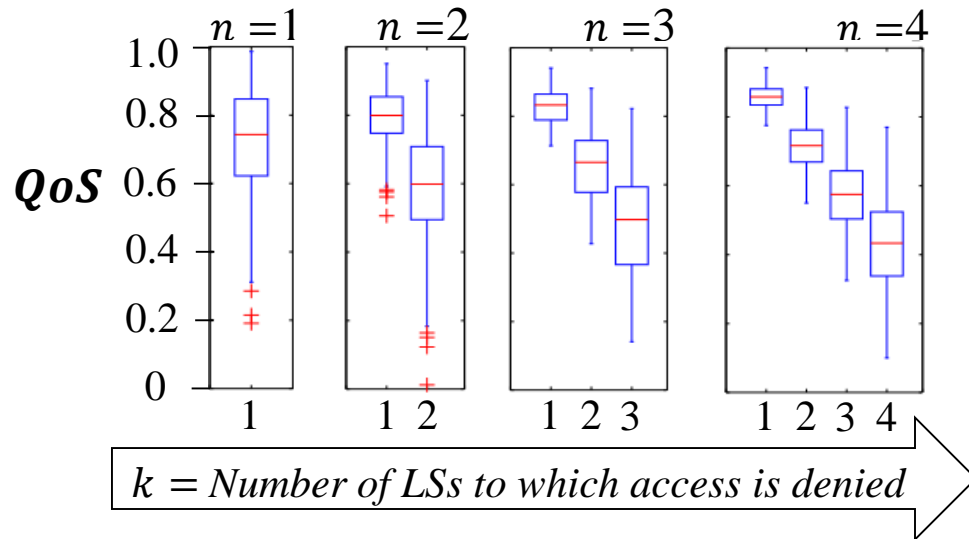  ○ Attack performance on protected frequency-profiles:



$P_{attack}(k)$ - probability of correct detection of a critical location when $k$ out of $n$ critical locations are already known

$\alpha_{rand}(k)$ - probability of randomly selecting a critical location

← Zero knowledge gain

$k = known\ Critical\ Locations$

# Evaluation: QoS and Communication Overhead

- **QoS** = proportion of available location updates



$k = 1, QoS \sim 80\%$
$k = 2, QoS \sim 70\%$
$k = 3, QoS \sim 60\%$

$k = Number\ of\ LSs\ to\ which\ access\ is\ denied$

$th_{crtl} = 80$

- ○ QoS is reasonably high given 60% population has 1 or 2 critical locations

- **Communication Cost** = no. of fake message per day

  - ○ 1-2 messages a day for most users!

# Related Work

- **Semantic location obfuscation (PROBE framework by Damiani et al. 2010)**

  + Cloak individual sensitive visits with neighboring non-sensitive venues

  – Sensitivity of location categories is not related to an individual's visit-frequency

- **Venue Recommendation techniques (Riboni et al. 2014, Zhang et al. 2014)**

  + Offline publishing of check-in history statistics in a differentially private manner

  – Require Trusted parties for implementing the privacy algorithm

  – Cannot be used for online location sharing

- **Distributed Location Management (Duerr et al. at Percom 2011)**

  + No *single-point-of-failure*

  – For single locations without considering location semantics

# Conclusion & Future Work

- *Frequent locations* naturally pose a privacy threat by revealing user interests

- *Distributing location information* in LS infrastructure → *promising privacy solution*

- Proposed an algorithm for controlled sharing of frequent locations
  - Hides frequent locations from:
    - *User-aware attackers*
    - *Population-aware attackers*

- **Future Work**
  - Integrate existing *single-location* semantic obfuscation approaches for forming a comprehensive privacy mechanism

# Contact and Discussion



*www.priloc.de*

**Zohaib Riaz**

Institute for Parallel and Distributed Systems,

University of Stuttgart, Germany

*zohaib.riaz@ipvs.uni-stuttgart.de*