# 802.11 User Fingerprinting

Jeffrey Pang[1]  Ben Greenstein[2]  Ramakrishna Gummadi[3]
Srinivasan Seshan[1]  David Wetherall[2,4]

[1]CMU  [2]Intel Research Seattle  [3]USC,MIT  [4]University of Washington

1

---

## Motivation: The Mobile Wireless Landscape



"Bob@Intel"  "Bob@Intel"

Why is Bob over there?

- They leave "digital fingerprints" that reveal who we are
  - And thus where we've been

2
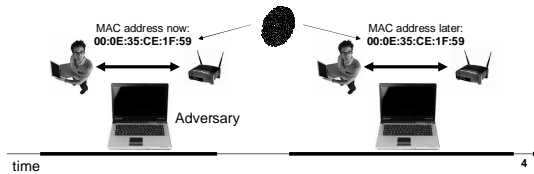
---

## Motivation: The Mobile Wireless Landscape

➔ *Location privacy* is growing concern
  - Wireless Privacy Protection Act [U.S. House Bill '05]
  - Geographic location/privacy working group [IETF]



**Man arrested for using GPS cellphone to track ex-girlfriend**

*by Brian Osborne posted on September 7, 2004 12:26 pm*

**Consumer Group Calls for Immediate Worldwide Boycott of Benetton**

**Wireless location tracking draws privacy questions**

Wireless products that can do everything from tracking your children to finding you a nearby date this weekend seem to fall outside the scope of federal privacy laws, and that may need to change, an industry group said.

3

---

## Motivation: The Mobile Wireless Landscape

- A well known technical problem
  - Devices have unique and consistent addresses
  - e.g., 802.11 devices have MAC addresses
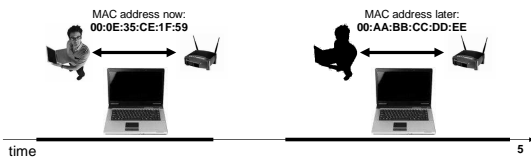  - ➔ fingerprinting them is trivial!



MAC address now:
00:0E:35:CE:1F:59

MAC address later:
00:0E:35:CE:1F:59

Adversary

time

4

---

## Motivation: The Mobile Wireless Landscape

- The widely proposed techical solution
  - **Pseudonyms**: Change addresses over time
    - 802.11: Gruteser '05, Hu '06, Jiang '07
    - Bluetooth: Stajano '05
    - RFID: Juels '04
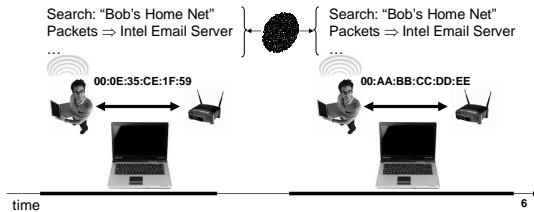    - GSM: already employed



MAC address now:
00:0E:35:CE:1F:59

MAC address later:
00:AA:BB:CC:DD:EE

time

5

---

## Motivation: The Mobile Wireless Landscape

- Our work shows: **Pseudonyms are not enough**
  - *Implicit identifiers*: identifying characteristics of traffic
  - E.G., most users identified with 90% accuracy in hotspots



Search: "Bob's Home Net"
Packets ⇒ Intel Email Server
…

Search: "Bob's Home Net"
Packets ⇒ Intel Email Server
…

00:0E:35:CE:1F:59

00:AA:BB:CC:DD:EE

time

6

1

## Contributions

- Four Novel 802.11 Implicit Identifiers

- Automated Identification Procedure

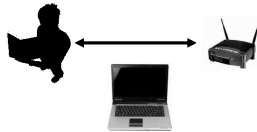- Evaluating Implicit Identifier Accuracy

7

## Contributions

- Four Novel 802.11 Implicit Identifiers

- Automated Identification Procedure

- Evaluating Implicit Identifier Accuracy

8

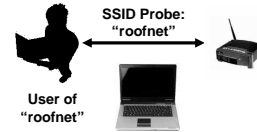## Implicit Identifiers by Example

- Consider one user at SIGCOMM 2004
  - Transferred 512MB via BitTorrent
    - Poor network etiquette?
  - Seen in a "anonymized" wireless trace
    - MAC addresses hashed, effectively a pseudonym
- Can we identify the culprit using implicit identifiers?



9

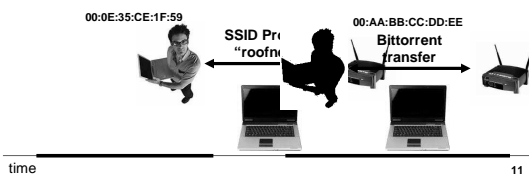## Implicit Identifiers by Example

- Implicit identifier: SSIDs in probes
  - Set of SSIDs in 802.11 probe requests
  - Many 802.11 drivers search for preferred networks
  - Usually networks you have associated with before



**SSID Probe: "roofnet"**

**User of "roofnet"**
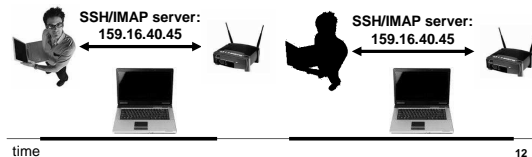
10

## Implicit Identifiers by Example

- Implicit identifier: SSIDs in probes
  - Set of SSIDs in 802.11 probe requests
  - Many 802.11 drivers search for preferred networks
  - Usually networks you have associated with before



00:0E:35:CE:1F:59

SSID Pr "roofn

00:AA:BB:CC:DD:EE
Bittorrent
transfer

time

11

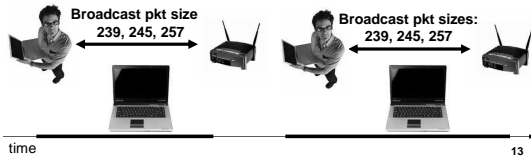## Implicit Identifiers by Example

- Implicit identifier: network destinations
  - IP <address, port> pairs in network traffic
  - At SIGCOMM, each visited by 1.15 users on average
  - Some nearly-unique destinations repeatedly visited (e.g., email server)



**SSH/IMAP server: 159.16.40.45**

**SSH/IMAP server: 159.16.40.45**

time

12

## Implicit Identifiers by Example

- Implicit identifier: broadcast packet sizes
  - Set of 802.11 broadcast packet sizes in network traffic
  - E.g., Windows machines NetBIOS naming advertisements; FileMaker and Microsoft Office advertise themselves
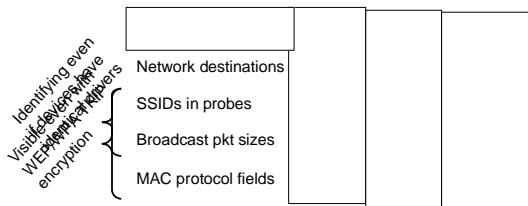
**Broadcast pkt size 239, 245, 257**

**Broadcast pkt sizes: 239, 245, 257**

time

**13**

## Implicit Identifiers by Example

- Implicit identifier: MAC protocol fields
  - Header bits (e.g., power management., order)
  - Supported rates
  - Offered authentication algorithms

**Protocol Fields: 11,4,2,1Mbps, WEP**

**Protocol Fields: 11,4,2,1Mbps, WEP**

time

**14**

## Implicit Identifier Summary

Identifying even WEP-encrypted drivers

Visible devices have identical drivers

Network destinations

SSIDs in probes

Broadcast pkt sizes

MAC protocol fields

- More implicit identifiers exist
  - ➔ Results we present establish a lower bound

**15**

## Fixing Implicit Identifiers is not Simple

- Encryption does not prevent traffic analysis
  - Cover traffic?
  - **Challenge**: Shared medium $\Rightarrow$ large performance hit

- Service discovery is done in the clear
  - Don't probe?
  - **Challenge**: Beaconing is often undesirable also

- Implementation and configuration variation
  - Standardize?
  - **Challenge**: Ambiguity of specifications

**16**

## Contributions

- Four Novel 802.11 Implicit Identifiers

- Automated Identification Procedure

- Evaluating Implicit Identifier Accuracy

**17**

## Tracking 802.11 Users

- Many potential tracking applications:
  - Was user X here today?
  - Where was user X today?
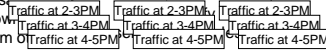  - What traffic is from user X?
  - When was user X here?
  - Etc.

**18**

## Tracking 802.11 Users

- Tracking scenario:
  - Every users changes pseudonyms every hour
  - Adversary monitors some locations
  - → One hourly traffic sample from each user in each location



Build a *profile* from training samples:
First collect some traffic known
from user X and from random of

| Traffic at 2-3PM | Traffic at 2-3PM | Traffic at 2-3PM |
| Traffic at 3-4PM | Traffic at 3-4PM | Traffic at 3-4PM |
| Traffic at 4-5PM | Traffic at 4-5PM | Traffic at 4-5PM |

19

---

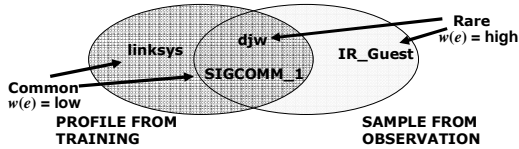## Sample Classification Algorithm

- Core question:
  - Did traffic sample $s$ come from user X?

- A simple approach: naïve Bayes classifier
  - Derive probabilistic model from training samples

  - Given $s$ with features $F$, answer "yes" if:
    **Pr[ $s$ from user X | $s$ has features $F$ ] > $T$**
    for a selected threshold $T$.

  - $F$ = feature set derived from implicit identifiers

20

---

## Sample Classification Algorithm

- Deriving features $F$ from implicit identifiers



Set similarity (Jaccard Index), weighted by frequency:

$$feature_\Box(s) = \frac{\sum_{e \in Profile_\Box \cap Set_s} w(e)}{\sum_{e \in Profile_\Box \cup Set_s} w(e)}$$

21

---

## Contributions

- Four Novel 802.11 Implicit Identifiers

- Automated Identification Procedure

- Evaluating Implicit Identifier Accuracy

22

---

## Evaluating Classification Effectiveness

- Simulate tracking scenario with wireless traces:

| | Duration | Profiled Users | Total Users |
|---|---|---|---|
| SIGCOMM conf. (2004) | 4 days | 377 | 465 |
| UCSD office building (2006) | 1 day | 153 | 615 |
| Apartment building (2006) | 14 days | 39 | 196 |

- Split each trace into training and observation phases
- Simulate pseudonym changes for each user X

23

---

## Evaluating Classification Effectiveness
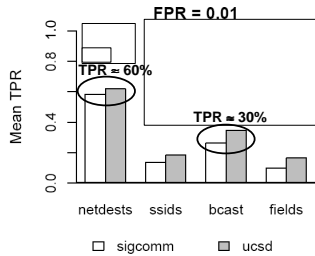
- **Question**: Is observation sample $s$ from user X?
- Evaluation metrics:
  - **True positive rate (TPR)** = ???
    Fraction of user X's samples classified correctly ↑ Measure TPR
  - **False positive rate (FPR)** = 0.01
    Fraction of other samples classified incorrectly

    Fix $T$ for FPR
    (See paper)

  **Pr[ $s$ from user X | $s$ has features $F$ ] > $T$**

24

---

4

## Results: Individual Feature Accuracy



FPR = 0.01

TPR ≈ 60%

TPR ≈ 30%

Mean TPR

netdests  ssids  bcast  fields

☐ sigcomm  ▨ ucsd

Individual implicit identifiers give evidence of identity

25

---

## Results: Multiple Feature Accuracy



Frac. users with TPR > x

Public
Home
Enterprise

True positive rate

**Users with TPR >50%:**

Public: 63%
Home: 31%
Enterprise: 27%

| | Public | Home | Enterprise |
|---|---|---|---|
| netdests | ✔ | | |
| ssids | ✔ | ✔ | ✔ |
| bcast | ✔ | ✔ | ✔ |
| fields | ✔ | ✔ | |

We can identify many users in all environments

26

---

## Results: Multiple Feature Accuracy



Frac. users with TPR > x

Public
Home
Enterprise

True positive rate

Public networks:
~20% users identified
>90% of the time

| | Public | Home | Enterprise |
|---|---|---|---|
| netdests | ✔ | | |
| ssids | ✔ | ✔ | ✔ |
| bcast | ✔ | ✔ | ✔ |
| fields | ✔ | ✔ | |

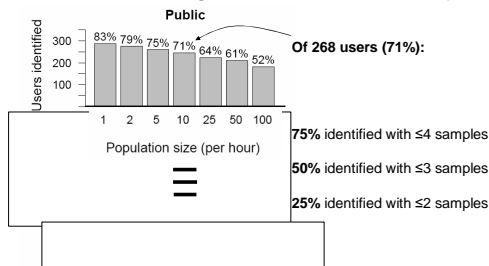Some users much more distinguishable than others

27

---

## One Application

- **Question**: Was user X here today?
- More difficult to answer:
  - Suppose $N$ users present each hour
  - Over an 8 hour day, $8N$ opportunities to misclassify
  - ➔ Decide user X is here only if *multiple* samples are classified as his

- **Revised**: Was user X here today for a few hours?

28

---

## Results: Tracking with 90% Accuracy

**Public**



Users identified

300
200
100

83%  79%  75%  71%  64%  61%  52%

1  2  5  10  25  50  100

Population size (per hour)

**Of 268 users (71%):**

**75%** identified with ≤4 samples

**50%** identified with ≤3 samples

**25%** identified with ≤2 samples

Majority of users can be identified if active long enough

29

---

## Conclusions

- Implicit identifiers can accurately identify users
  - Individual implicit identifiers give evidence of identity
  - We can identify many users in all environments
  - Some users much more distinguishable than others

- Understanding implicit identifiers is important
  - Pseudonyms are not enough
  - We establish a *lower bound* on their accuracy

- Eliminating them poses research challenges
  - Current work: Confidential service discovery
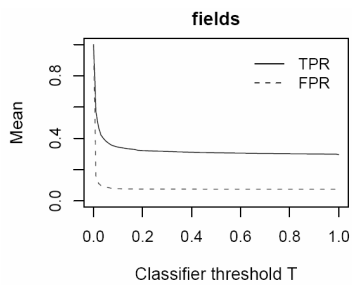  - Current work: Traffic analysis resistant MAC

30

# Extra Slides

# Related Work

- Other Implicit Identifiers
  - Device driver fingerprints [Franklin '06]
  - Clock-skew fingerprints [Kohno '05]
  - Click-prints [Padmanabhan '06]
  - RF antenna fingerprints [Hall '04]
- Our work:
  - 802.11 fingerprints for *individual users*
  - Tracking with only commodity hardware/software
  - Better coverage than some previous work
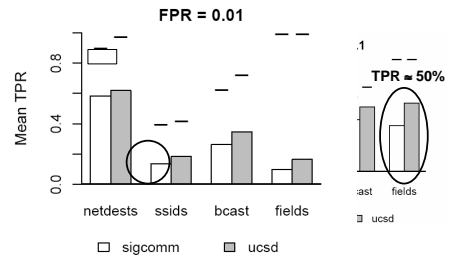  - Procedure to combine implicit identifiers
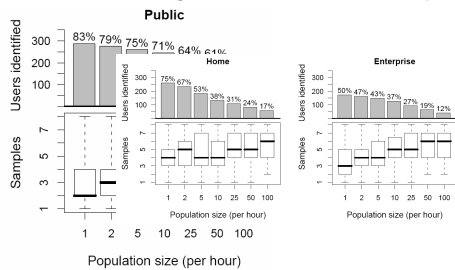
# Evaluating Classification Effectiveness

**fields**

# Results: Individual Feature Accuracy

**FPR = 0.01**



Other implicit identifiers distinguish groups of users

# Results: Tracking with 90% Accuracy



Many users can be identified in all environments

# Old Slides

## Answers to Common Questions

- Is tracking accuracy good enough to track an entire city/how does it vary with population size?
  - Our accuracy results aren't good enough to allow fine-grained tracking of most users in environments with more than 100s of users because the adversary will simply obtain too many samples and have too many opportunities for false positives. However, we note that by modifying the question a little, such as "was this person here for x hours, where x is reasonably large, such as a couple days" or "was this person here in this specific area with smaller population in the sub-100s", we can answer it fairly accurately.
- Your traces were only over a couple days, why would one expect fingerprints to remain identifying over longer periods of time?
  - You are correct and we can't say anything empirical for longer term tracking. However, we suspect that the implicit identifiers would remain identifying because they are mostly caused by automatic behavior of drivers and applications and thus wouldn't change due to human behavior or location. In the paper we show that ssids are stable over weeks at least for example.
- Do you think people really care about location privacy?
  - I think that people don't care until they realize that they have a problem and then they do (e.g., the article about a guy tracking his ex-girlfriend). One magnifying factor is that because of the low cost of 802.11 radios and the ubiquity of 802.11 in devices, almost anyone can easily monitor some locations and track users of interest.
- Don't you think people could modify their behavior in order not to be as identifying when they want to be private?
  - Sure; if they had the technical knowledge to do so, they could easily modify their behavior to "look different" and leave a different implicit fingerprint. However, I would argue that this is undesirable because it means that privacy concerns have a chilling effect on the types of act ivies we are willing to do thus limiting the wireless applications. I think the protocols themselves should protect us enough so that we don't have to worry about these things.
- What about wireless protocols other than 802.11?
  - I believe many of the implicit identifiers still exist, for the reason that the fundamental reasons behind their existence remain: you can still do traffic analysis even with encryption, all wireless protocols have some form of service discovery, and there is implementation variation.

37

---

## Outline

- Problem
  - People are worried about tracking, 802.11 is especially worrisome
  - Pseudonyms proposed, not enough
- Bittorrent Example
  - Use to explain each identifier
  - Summarize implicit identifiers
- How to train as an example
  - Points to bring up and jx:
    - 1 hour sample size
    - How to select the classifier threshold
    - How adversary could obtain training samples
  - Learning process
- Q1 results
- Q2 results
  - An attacker can use multiple features and multiple samples, so ask question…

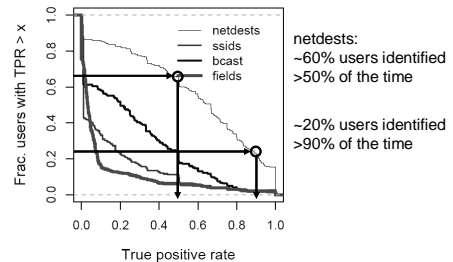38

---

## Implicit Identifiers by Example

- Implicit identifier: network destinations
  - IP <address, port> pairs in network traffic
  - At SIGCOMM, each visited by 1.15 users on average
  - Some nearly-unique destinations repeatedly visited (e.g., email server)



**WPA TKIP encryption**

39

---

## Results: Individual Feature Accuracy



netdests:
~60% users identified
>50% of the time

~20% users identified
>90% of the time

Some users more distinguishable than others

40

---

7