

Privacy and Anonymity

Dawn Song
dawnsong@cs.berkeley.edu

Slides credit: John Mitchell & Vitaly Shmatikov

Current state of the world II

- **EU directive 2006/24/EC: 3 year data retention**
 - For ALL traffic, requires EU ISPs to record:
 - » Sufficient information to identify endpoints (both legal entities and natural persons)
 - » Session duration
 - » ... but not session contents
 - Make available to law enforcement
 - » ... but penalties for transfer or other access to data
- **For info on US privacy on the net:**
 - “privacy on the line” by W. Diffie and S. Landau

Anonymous web browsing

- **Why?**
 1. Discuss health issues or financial matters anonymously
 2. Bypass Internet censorship in parts of the world
 3. Conceal interaction with gambling sites
 4. Law enforcement
- **Two goals:**
 - Hide user identity from target web site: (1), (4)
 - Hide browsing pattern from employer or ISP: (2), (3)
- **Stronger goal: mutual anonymity (e.g. remailers)**

Part 1: network-layer privacy

Goals:

- Hide user's IP address from target web site
- Hide browsing destinations from network

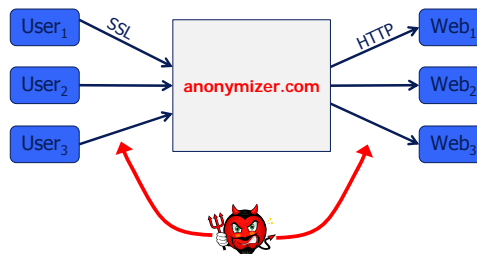
Current state of the world I

- **ISPs tracking customer browsing habits:**
 - Sell information to advertisers
 - Embed targeted ads in web pages (1.3%)
 - » Example: MetroFi (free wireless)

[Web Tripwires: Reis et al. 2008]
- **Several technologies used for tracking at ISP:**
 - NebuAd, Phorm, Front Porch
 - Bring together advertisers, publishers, and ISPs
 - » At ISP: inject targeted ads into non-SSL pages
- **Tracking technologies at enterprise networks:**
 - Vontu (symantec), Tablus (RSA), Vericept

1st attempt: anonymizing proxy

HTTPS:// anonymizer.com ? URL=target

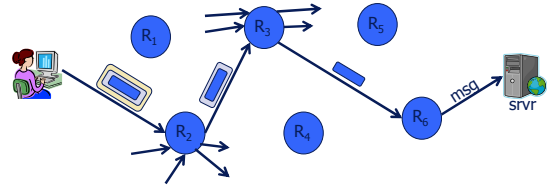


Anonymizing proxy: security

- **Monitoring ONE link:** eavesdropper gets nothing
- **Monitoring TWO links:**
 - Eavesdropper can do traffic analysis
 - More difficult if lots of traffic through proxy
- **Trust:** proxy is a single point of failure
 - Can be corrupt or subpoenaed
 - » Example: The Church of Scientology vs. anon.penet.fi
- **Protocol issues:**
 - Long-lived cookies make connections to site linkable

7

MIX nets [C'81]



- **Every router has public/private key pair**
 - Sender knows all public keys

To send packet:

- Pick random route: $R_2 \rightarrow R_3 \rightarrow R_6 \rightarrow \text{svr}$
- Prepare onion packet:

$$\text{packet} = E_{pk_2}(R_3, E_{pk_3}(R_6, E_{pk_6}(\text{svr}, \text{msg})))$$

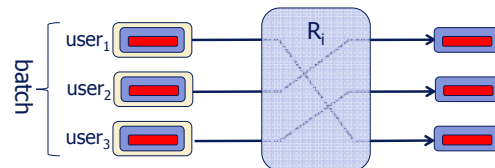
10

How proxy works

- **Proxy rewrites all links in response from web site**
 - Updated links point to anonymizer.com
 - » Ensures all subsequent clicks are anonymized
- **Proxy rewrites/removes cookies and some HTTP headers**
- **Proxy IP address:**
 - if a single address, could be blocked by site or ISP
 - anonymizer.com consists of >20,000 addresses
 - » Globally distributed, registered to multiple domains
 - » Note: chinese firewall blocks ALL anonymizer.com addresses

8

Eavesdropper's view at a single MIX



- **Eavesdropper observes incoming and outgoing traffic**
- **Crypto prevents linking input/output pairs**
 - Assuming enough packets in incoming batch
 - If variable length packets then must pad all to max len
- **Note:** router is stateless

11

2nd Attempt: MIX nets

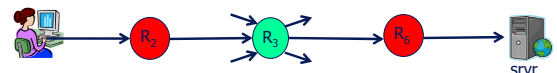
Goal: no single point of failure

9

Performance

• Main benefit:

- Privacy as long as at least one honest router on path



• Problems:

- High latency (lots of public key ops)
 - » Inappropriate for interactive sessions
 - » May be OK for email
- No forward security

• Homework puzzle: how does server respond?

- hint: user includes "response onion" in forward packet

Web-based user tracking

Browser provides many ways to track users:

1. 3rd party cookies ; Flash cookies
2. Tracking through the history file
3. Machine fingerprinting

13

Effectiveness of 3rd party blocking

- Ineffective for improving privacy
 - 3rd party can become first party and then set cookie
 - Flash cookies not controlled by browser cookie policy
- **Better proposal:**
 - Delete all browser state upon exit
 - Supported as an option in IE7

16

3rd party cookies

- **What they are:**
 - User goes to site A.com ; obtains page
 - Page contains `<iframe src="B.com">`
 - Browser goes to B.com ; obtains page
HTTP response contains cookie
 - Cookie from B.com is called a 3rd party cookie
- **Tracking:** User goes to site D.com
 - D.com contains `<iframe src="B.com">`
 - B.com obtains cookie set when visited A.com
 - ⇒ B.com knows user visited A.com and D.com

14

Tracking through the history file

- E.g., site checks hyper-link color for history
- **Applications:**
 - Context aware phishing:
 - » Phishing page tailored to victim
 - Marketing
 - Use browsing history as 2nd factor authentication

17

Can we block 3rd party cookies?

- Supported by most browsers
- **IE and Safari:** block set/write
 - Ignore the "Set-Cookie" HTTP header from 3rd parties
 - ⇒ Site sets cookie as a 1st party; will be given cookie when contacted as a 3rd party
 - Enabled by default in IE7
- **Firefox and Opera:** block send/read
 - Always implement "Set-Cookie", but never send cookies to 3rd party
 - Breaks sess. mgmt. at several sites (off by default)

15

Context-aware Phishing

- Stanford students see:



- Cal students see:



18

SafeHistory/SafeCache [JBBM'06]

- Define Same Origin Policy for all long term browser state
 - history file and web cache
 - Firefox extensions: SafeHistory and SafeCache
- **Example: history**
 - Color link as visited only when site can tell itself that user previously visited link:
 - » A same-site link, or
 - » A cross-site link previously visited from this site

19

Administrivia

- Office hour on Tue changed to 1-3pm due to schedule conflict
- Additional office hour on Thu/Fri afternoon to help students preparing the final
- In-class final: Dec 10, 306 Soda
- Final review on Wed
- Guest lecture on Mon: real-world experiences about breaking security systems

22

Machine fingerprinting

- Tracking using machine fingerprints
- User connects to site A.com
 - Site builds a fingerprint of user's machine
 - Next time user visits A.com, site knows it is the same user

20

De-anonymizing data

23

Machine fingerprints [Khono et al.'05]

- Content and order of HTTP headers
 - e.g. user-agent header:
Mozilla/5.0 (Windows; U; Windows NT 6.0; en-US; rv:1.8.1.14) Gecko/20080404 Firefox/2.0.0.14
- Javascript and JVM can interrogate machine properties:
 - Timezone, local time, local IP address
- **TCP timestamp: exploiting clock skew**
 - TCP_timestamp option: peer embeds 32-bit time in every packet header. Accurate to ≈ 100 ms
 - fingerprint = (real-time Δ between packets)
(timestamp Δ between-packets)

21

Problem statement

- An organization collects private user data
 - Wishes to make data available for research
 - Individual identities should be hidden
- Examples:
 - Search queries over a 3 month period (AOL)
 - Netflix movie rentals
 - Census data
 - Social networking data

24

Incorrect approach

- Replace “username” or “userID” by random value

Dan → a56fd863ec
John → 87649dce63

– Same value used for all appearances of userID

- **Problem:** often data can be de-anonymized by combining auxiliary information
- Examples: AOL search data
 census data

25

De-Anonymizing Netflix Records

- Average subscriber has 214 dated ratings
- How many does the attacker need to know to identify his victim’s record in the dataset?
 - **Two** is enough to reduce to 8 candidate records
 - **Four** is enough to identify uniquely (on average)
 - Works even better with relatively rare ratings
 - » “The Astro-Zombies” rather than “Star Wars”
- **Sparsity!** Negligible information leakage is sufficient for complete re-identification

Fat Tail effect helps here:
most people watch obscure crap (really!)

26

Netflix Prize Dataset

- Released in October 2006 to support research on better recommender algorithms
 - Real movie ratings of 500,000 Netflix subscribers
 - » 10% of all Netflix users as of late 2005
 - Names removed
 - Maybe perturbed
- Good target for real-world de-anonymization
 - What can you learn from someone’s movie ratings?
 - What to use as source of external information?

26

Robustness

- De-anonymization algorithm is robust to errors in attacker’s external knowledge
 - Dates and ratings may be known imprecisely
 - Some may even be completely wrong
 - Perturbation = noise in the data = doesn’t matter!
- **Why? Sparsity!!**
 - Nearest neighbor is so far, can tolerate huge amount of noise, perturbation, imprecision
- Where to find external knowledge?
 - Cross-correlating with Internet Movie DB

29

Netflix’s Take on Privacy



Even if, for example, you knew **all** your own ratings and their dates you **probably** couldn’t identify them reliably in the data because only a small sample was included (less than **one-tenth** of our complete dataset) and that data was subject to **perturbation**. Of course, since you know **all your own** ratings that really isn’t a privacy problem is it?
– Netflix Prize FAQ

27

What Did We Learn?

30

Netflix Users with Distance < 0.15

