

# CSE 484 / CSE M 584: Web Security

Winter 2023

Tadayoshi (Yoshi) Kohno  
yoshi@cs.Washington.edu

UW Instruction Team: David Kohlbrenner, Yoshi Kohno, Franziska Roesner. Thanks to Dan Boneh, Dieter Gollmann, Dan Halperin, John Manferdelli, John Mitchell, Vitaly Shmatikov, Bennet Yee, and many others for sample slides and materials ...

# Announcements

- Wednesday, 3/1: My office hours 12:30-1 (changed); TA office hours still available, e.g., for HW2
- Friday, 3/3: Physical Security
  - Details discussed on 2/27/2023 lecture
  - Amazon.com (and other places) sell equipment; note local laws
  - (We will bring supplies, not necessary to buy any)
- Friday, 3/10: Charlie Reis (Google, Chrome, Site Isolation)

# Example Set (Amazon.com)



# Announcements

- HW3: There will be one, but only for extra credit
- HW1: Impressive!! I am still trying to slowly go through all of them and may adjust some grades + add more comments
- Final Project Checkpoint #1: Great ideas and directions!
- Final Project Checkpoint #2: March 8 – outline and list of references
- Final Project Final: Asking CSE Advising (may update with more information)



The New Yorker,  
1993

*"On the Internet, nobody knows you're a dog."*

The Joy of Tech™



© 2013 Geek Culture

by Nitrozac & Snaggy



joyoftech.com

<http://www.geekculture.com/joyoftech/joyarchives/1862.html>

# Privacy on Public Networks

- Internet is designed as a public network
  - Machines on your LAN may see your traffic, network routers see all traffic that passes through them
- Routing information is public
  - IP packet headers identify source and destination
  - Even a passive observer can figure out **who is talking to whom**
- Encryption does not hide identities
  - **Encryption hides payload, but not routing information**
  - Even IP-level encryption (tunnel-mode IPSec/ESP) reveals IP addresses of IPSec gateways
- **Modern web: Accounts, web tracking, etc. ...**

# Questions

**Q1:** What is anonymity?

**Q2:** Why might people **want** anonymity on the Internet? (E.g., want anonymity for themselves or others.)

**Q3:** Why might people **not want** anonymity on the Internet? (E.g., not want anonymity for themselves or others.)



# What is Anonymity?

- Anonymity is the state of being not identifiable within a **set of subjects**
  - You cannot be anonymous by yourself!
    - Big difference between anonymity and confidentiality
  - Hide your activities among others' similar activities
- Unlinkability of action and identity
  - For example, sender and email they send are no more related after observing communication than before
- Unobservability (hard to achieve)
  - Observer cannot even tell whether a certain action took place or not

# Applications of Anonymity (I)

- Privacy
  - Hide online transactions, Web browsing, etc. from intrusive governments, marketers and archivists
- Untraceable electronic mail
  - Corporate whistle-blowers
  - Political dissidents
  - Socially sensitive communications
  - Confidential business negotiations
- Law enforcement and intelligence
  - Sting operations and honeypots
  - Secret communications on a public network

# Applications of Anonymity (II)

- Digital cash
  - Electronic currency with properties of paper money (online purchases unlinkable to buyer's identity)
  - (Anonymous) digital cash long predates crypto currencies (e.g., Chaum's digital cash from 1983)
- Anonymous electronic voting
- Censorship-resistant publishing

# Part 1: Anonymity in Datasets

# How to release an anonymous dataset?

## A Face Is Exposed for AOL Searcher No. 4417749

By MICHAEL BARBARO and TOM ZELLER Jr.; Saul Hansell contributed reporting for this article.  
Published: August 9, 2006

Buried in a list of 20 million Web search queries collected by AOL and recently released on the Internet is user No. 4417749. The number was assigned by the company to protect the searcher's anonymity, but it was not much of a shield.


No. 4417749 conducted hundreds of searches over a three-month period on topics ranging from "numb fingers" to "60 single men" to "dog that urinates on everything."


And search by search, click by click, the identity of AOL user No. 4417749 became easier to discern. There are queries for "landscapers in Lilburn, Ga," several people with the last name Arnold and "homes sold in shadow lake subdivision gwinnett county georgia."


It did not take much investigating to follow that data trail to Thelma Arnold, a 62-year-old widow who lives in Lilburn, Ga., frequently researches her friends' medical ailments and loves her three dogs. "Those are my searches," she said, after a reporter read part of the list to her.

 FACEBOOK

 TWITTER

 GOOGLE+

 EMAIL

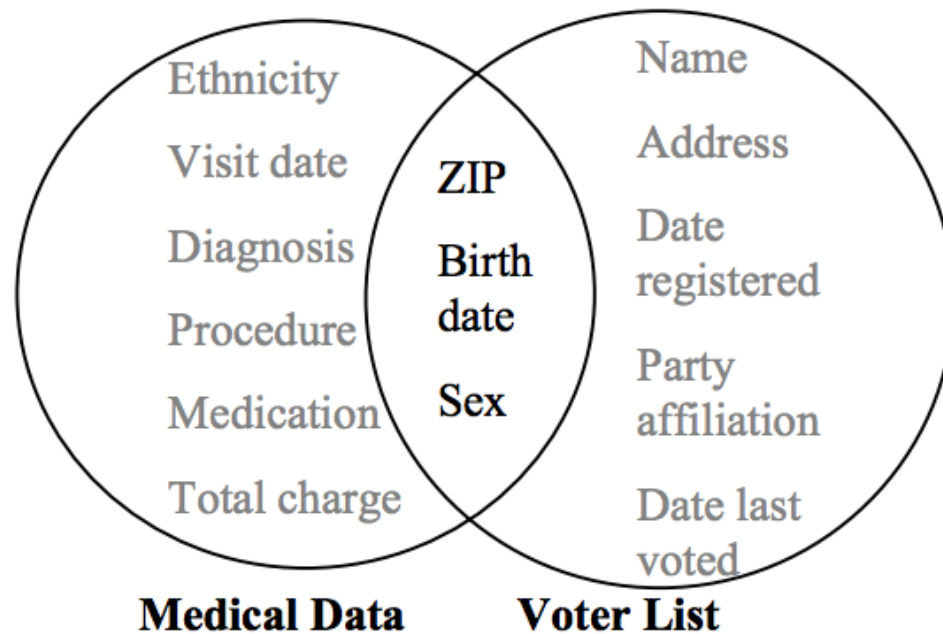
 SHARE

 PRINT

 REPRINTS

# How to release an anonymous dataset?

- Possible approach: remove identifying information from datasets?



Massachusetts  
medical+voter data  
[Sweeney 1997]

**Figure 1 Linking to re-identify data**

[Sweeney 2002]

# k-Anonymity

- Each person contained in the dataset cannot be distinguished from at least k-1 others in the data.

Name	Age	Gender	State of domicile	Religion	Disease
*	20 < Age ≤ 30	Female	Tamil Nadu	*	Cancer
*	20 < Age ≤ 30	Female	Kerala	*	Viral infection
*	20 < Age ≤ 30	Female	Tamil Nadu	*	TB
*	20 < Age ≤ 30	Male	Karnataka	*	No illness
*	20 < Age ≤ 30	Female	Kerala	*	Heart-related
*	20 < Age ≤ 30	Male			
*	Age ≤ 20	Male			
*	20 < Age ≤ 30	Male			
*	Age ≤ 20	Male			
*	Age ≤ 20	Male	Kerala	*	Viral infection

Doesn't work for high-dimensional datasets (which tend to be **sparse**)

## Robust De-anonymization of Large Sparse Datasets

Arvind Narayanan and Vitaly Shmatikov

The University of Texas at Austin

[Narayanan and Shmatikov 2008]

# Netflix Challenge

- Netflix released a (non-uniform) random sample of user's movie ratings
- Challenge was to build a better recommendation system
- Data was “anonymous”
  - ID # only
  - Random selection of a given user's ratings
  - “noise” added (actually, appears that there was no noise)



[Narayanan and Shmatikov 2008]

## Result: No real anonymity

- Deanonymization technique: Cross-correlate with IMBD ratings
- A handful (6 or fewer) ratings of non-top 500 movies is enough!

# Differential Privacy

- **Setting:** Trusted party has a database
- **Goal:** allow queries on the database that are useful but preserve the privacy of individual records
- **Differential privacy intuition:** add noise so that an output is produced with similar probability whether any single input is included or not
- Privacy of the computation, not of the dataset

# Part 2: Anonymity in Communication

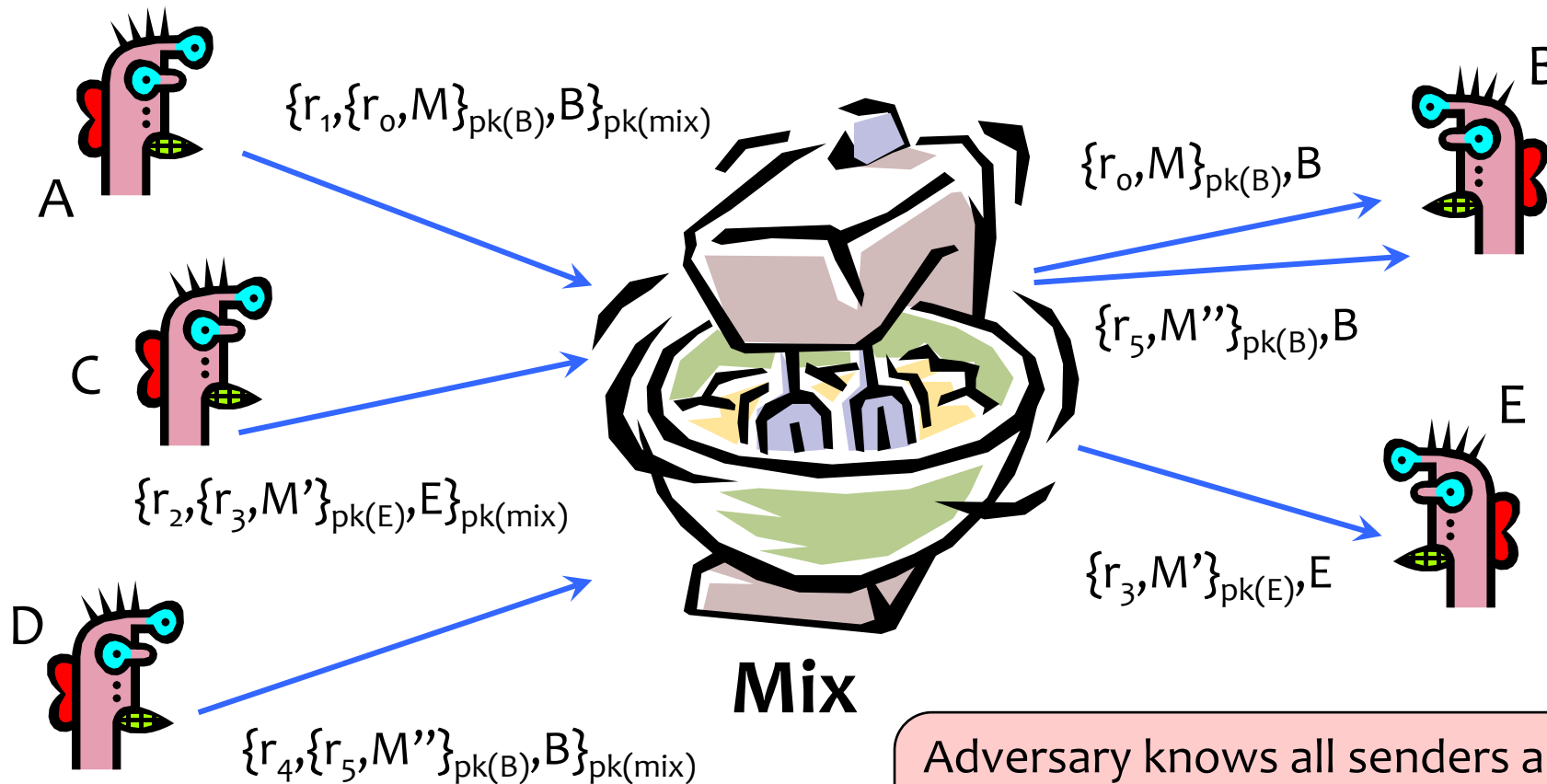
# Chaum's Mix

- Early proposal for anonymous email
  - David Chaum. “Untraceable electronic mail, return addresses, and digital pseudonyms”. Communications of the ACM, February 1981.

Before spam, people thought anonymous email was a good idea 😊

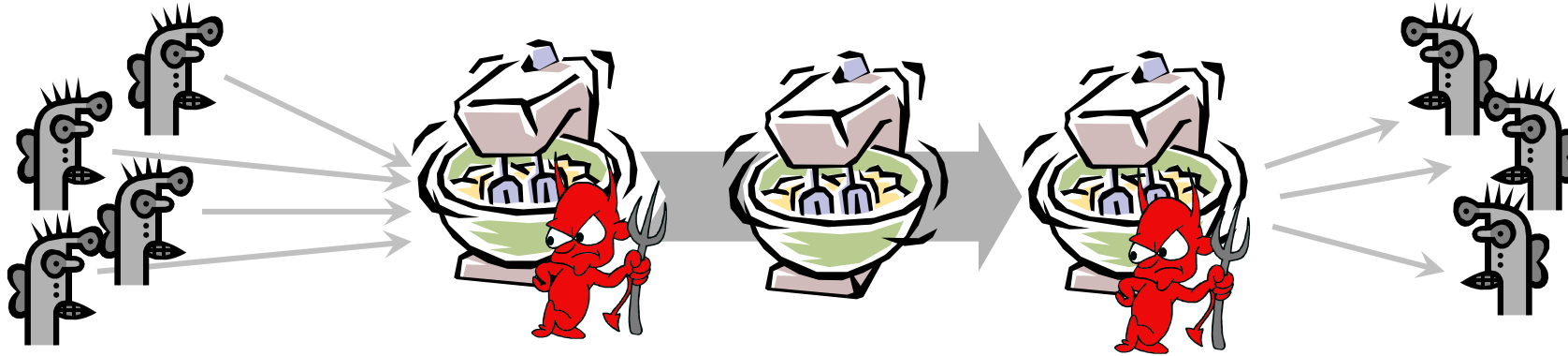
- Modern anonymity systems use Mix as the basic building block

# Basic Mix Design



Adversary knows all senders and all receivers, but cannot link a sent message with a received message

# Mix Cascades and Mixnets



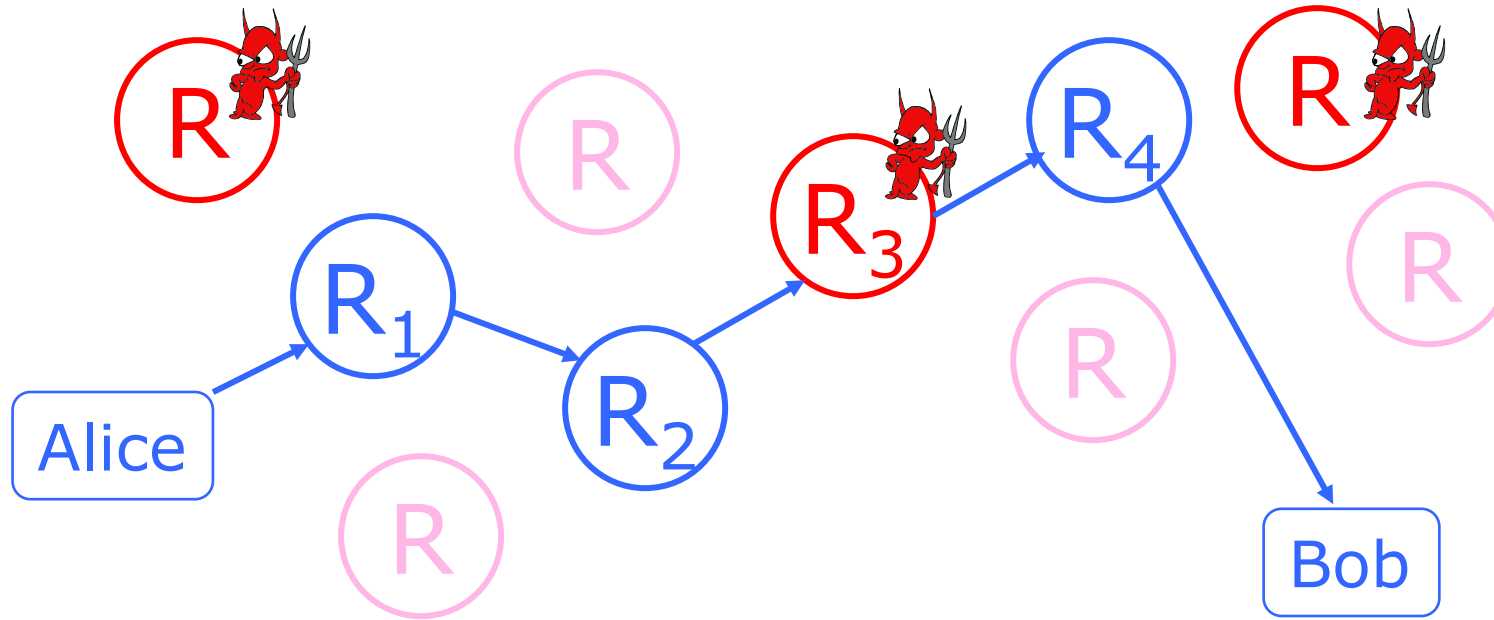
- Messages are sent through a **sequence of mixes**
  - Can also form an arbitrary network of mixes (“mixnet”)
- Some of the mixes may be controlled by attacker, but even a single good mix ensures anonymity
- Pad and buffer traffic to foil **correlation attacks**

# Disadvantages of Basic Mixnets

- Public-key encryption and decryption at each mix are **computationally expensive**
- Basic mixnets have **high latency**
  - OK for email, not OK for anonymous Web browsing
- Challenge: **low-latency anonymity network**

# Another Idea: Randomized Routing

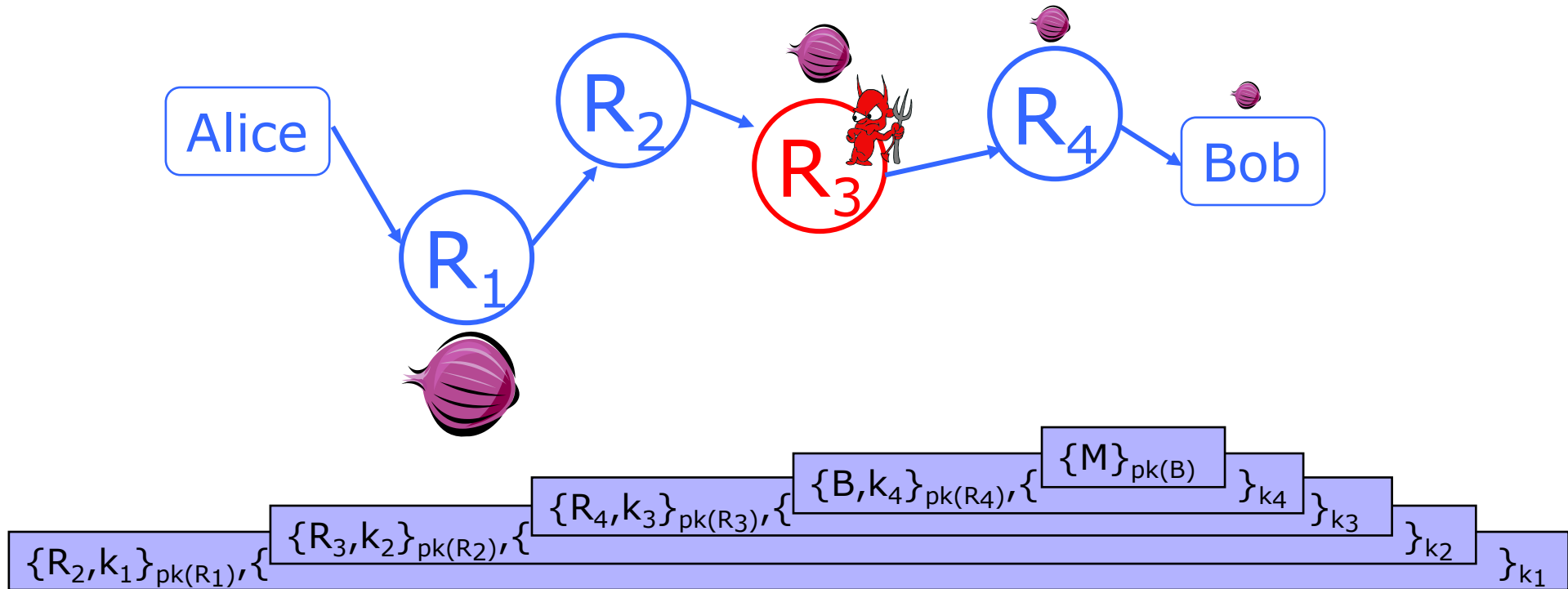
e.g., Onion Routing



- Sender chooses a random sequence of routers
  - Some routers are honest, some controlled by attacker
  - Sender controls the length of the path



# Onion Routing



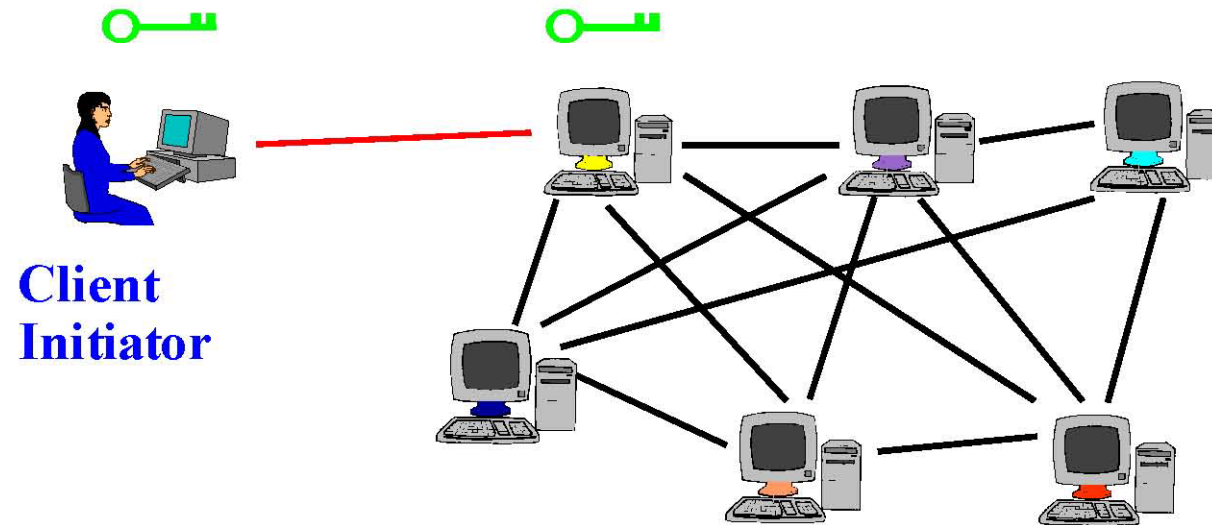
- Routing info for each link encrypted with router's public key
- Each router learns only the identity of the next router

# Tor

- Second-generation onion routing network
  - <http://tor.eff.org>
  - Developed by Roger Dingledine, Nick Mathewson and Paul Syverson
  - Specifically designed for **low-latency** anonymous Internet communications
- Running since October 2003
- “Easy-to-use” client proxy
  - Freely available, can use it for anonymous browsing

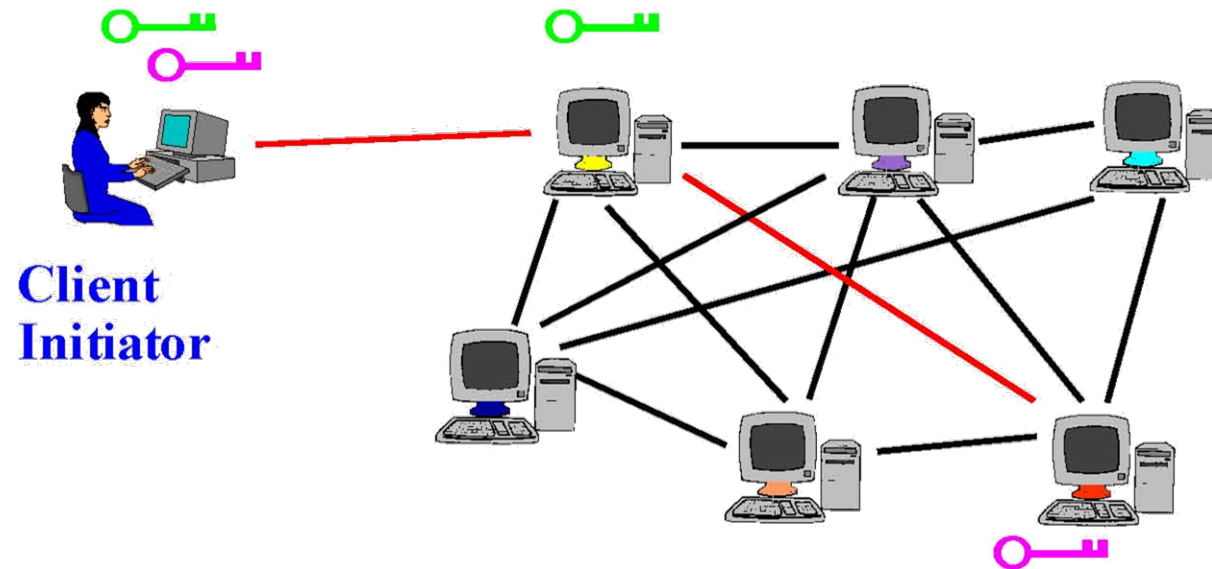
# Tor Circuit Setup (1)

- Client proxy establishes a symmetric session key and circuit with Onion Router #1



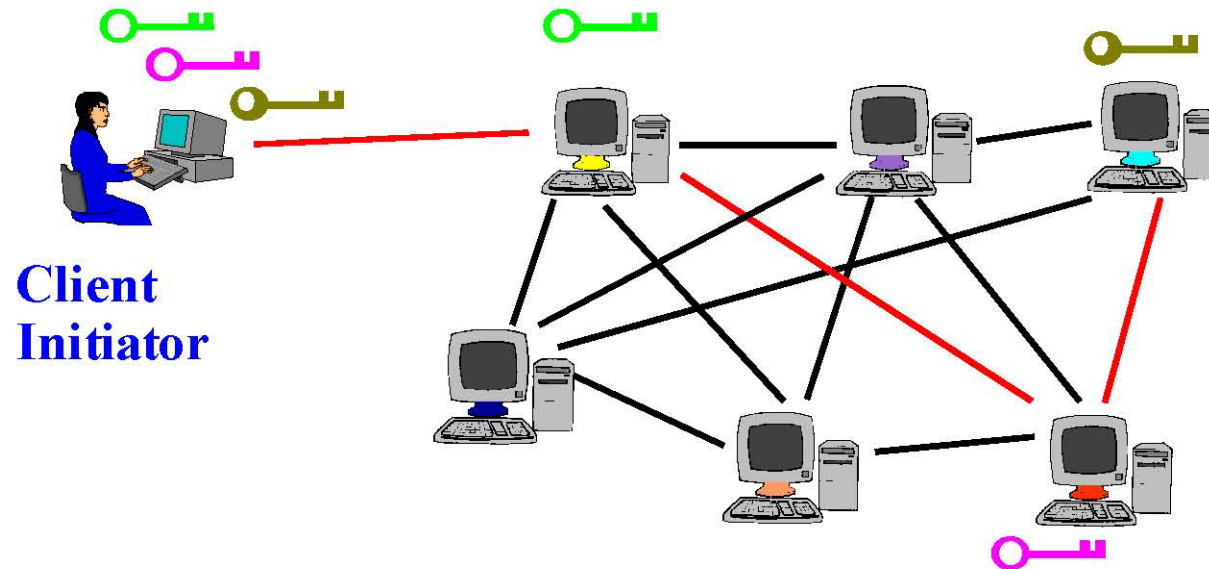
# Tor Circuit Setup (2)

- Client proxy extends the circuit by establishing a symmetric session key with Onion Router #2
  - Tunnel through Onion Router #1



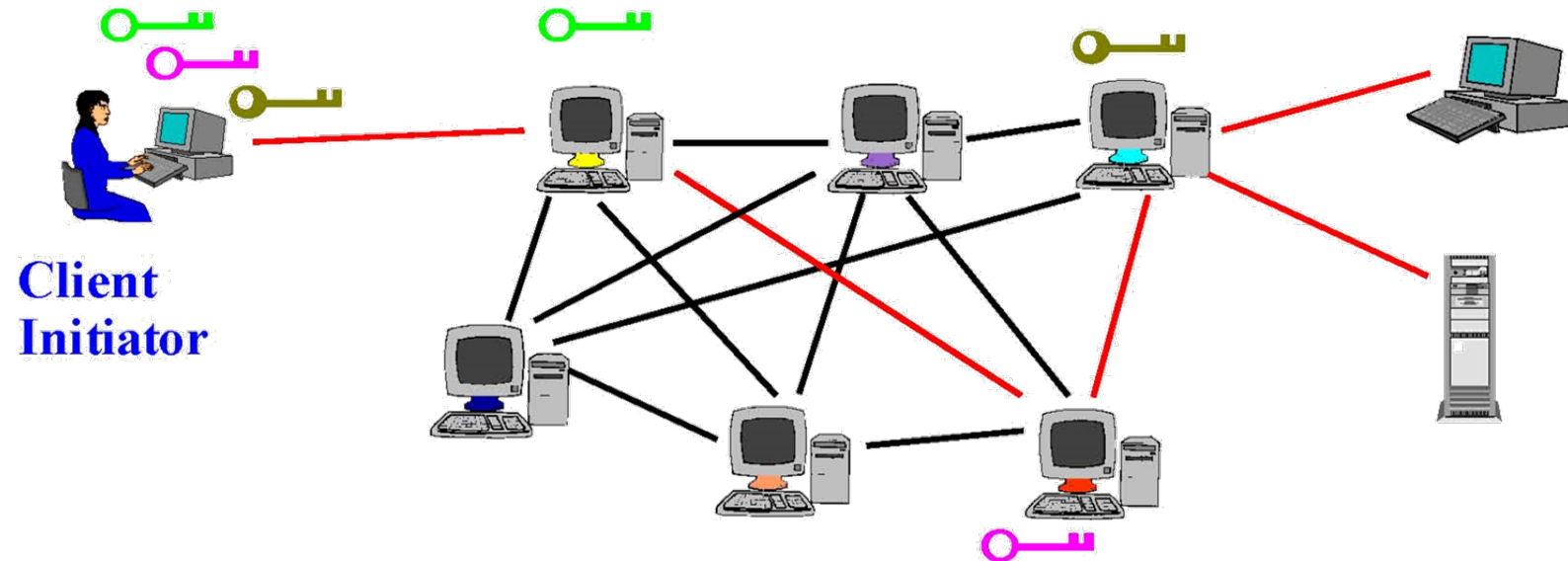
# Tor Circuit Setup (3)

- Client proxy extends the circuit by establishing a symmetric session key with Onion Router #3
  - Tunnel through Onion Routers #1 and #2



# Using a Tor Circuit

- Client applications connect and communicate over the established Tor circuit.



# How do you know who to talk to?

- Directory servers
  - Maintain lists of active onion routers, their locations, current public keys, etc.
  - Control how new routers join the network
    - “Sybil attack”: attacker creates a large number of routers
  - Directory servers’ keys ship with Tor code

# Tor Management Issues

- Many applications can share one circuit
  - Multiple TCP streams over one anonymous connection
- Tor router doesn't need root privileges
  - Encourages people to set up their own routers
  - More participants = better anonymity for everyone
- Directory servers
  - Maintain lists of active onion routers, their locations, current public keys, etc.
  - Control how new routers join the network
    - “Sybil attack”: attacker creates a large number of routers
  - Directory servers' keys ship with Tor code

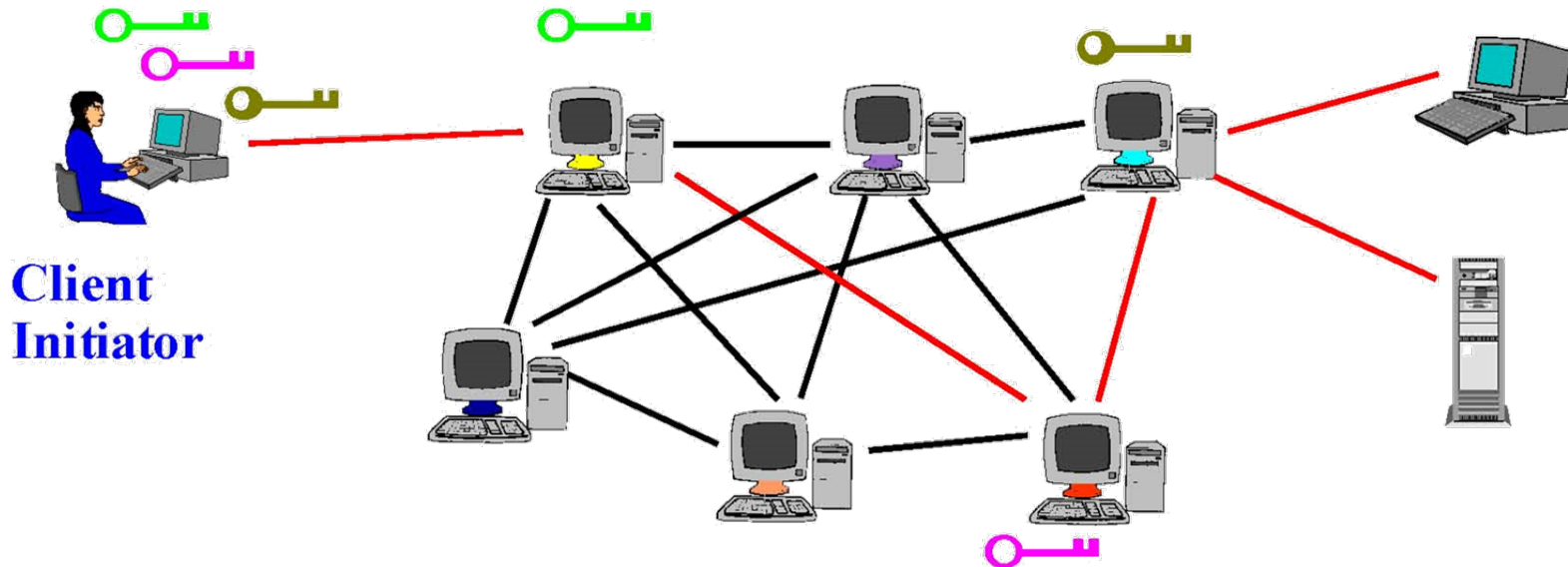


# Issues and Notes of Caution

- Passive traffic analysis
  - Infer from network traffic who is talking to whom
  - To hide your traffic, must carry other people's traffic!
- Active traffic analysis
  - Inject packets or put a timing signature on packet flow
- Compromise of network nodes
  - Attacker may compromise some routers
    - Powerful adversaries may compromise “too many”
  - It is not obvious which nodes have been compromised
    - Attacker may be passively logging traffic
  - Better not to trust any individual router
    - Assume that some fraction of routers is good, don't know which

# Issues and Notes of Caution

- Tor isn't completely effective by itself
  - Tracking cookies, fingerprinting, etc.
  - Exit nodes can see everything!



# Issues and Notes of Caution

- The simple act of using Tor could make one a **target for additional surveillance**
- Hosting an exit node could result in **illegal activity coming from your machine**
- Tor not designed to protect against **adversaries with the capabilities of a nation state** (public statement by designers, at least in the past, and government(s) have compromised Tor in the past)