

# Anonymous Communication and Traffic Analysis II

Privacy Technology in Context David Sidi (dsidi@email.arizona.edu)



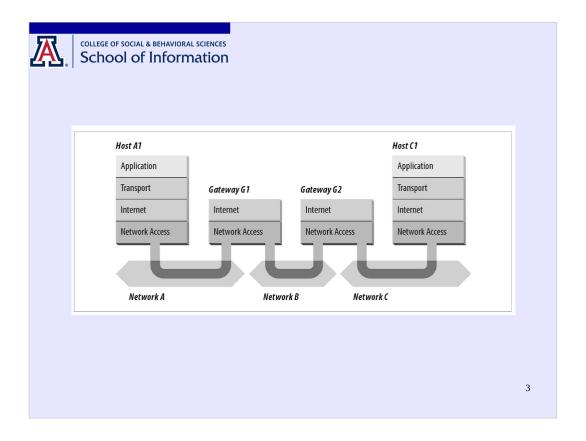
we're going to finish up some background before moving to anonymity, remaining fairly high level. Next time we will discuss further details.



# Small mention of interesting things

- Grading for assignment 4
- MIT student work on freehaven led to Tor

2

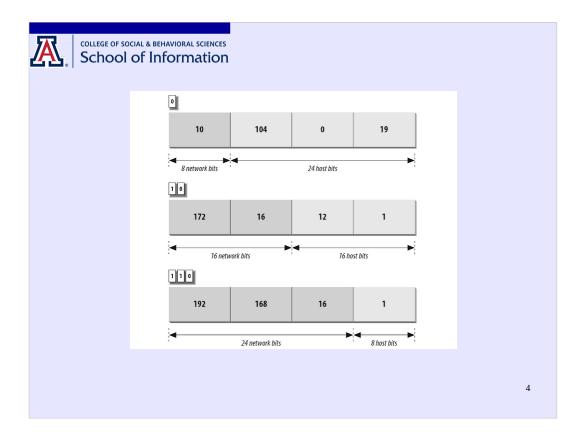


A bunch of physically heterogeneous networks that use TCP/IP form an internet. The Internet is, basically but not exactly, all the TCP/IP internets.

IP is central to the internet.

Hosts go through the whole protocol stack; gateways just go up to IP

Thus packets providing routing info get data from one network to another, then to a host on that network

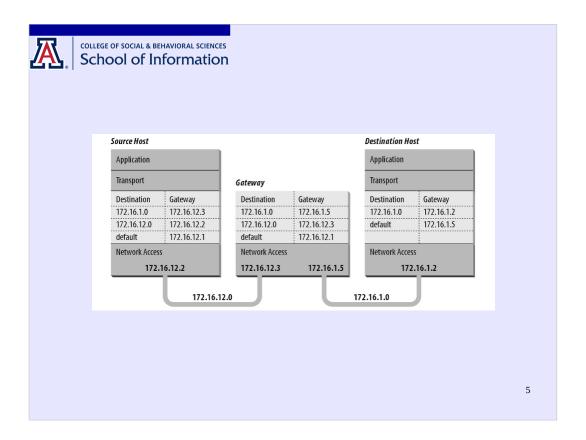


IP defines an addressing scheme, and uses it in IP packets (called datagrams) to move data hop-by-hop to the network that has the destination host on it.

This diagram is a little like what I drew on the board last time, showing the IP addressing scheme.

There is a network part and a host part to the address. The network part takes up less of the address for larger networks, and less of the address for smaller networks. The size of the network part is indicated using subnet masks (see: CIDR).

(In practice, big blocks of contiguous addresses are given to service providers who are best placed to provide routing information. Example: UA has 150.135.0.0; these are then delegated to groups in the university (the iSchool has 150.135.15.0, I think). The same is true of your ISP--- everyone is assigned an IP by their ISP in order to connect to the internet.

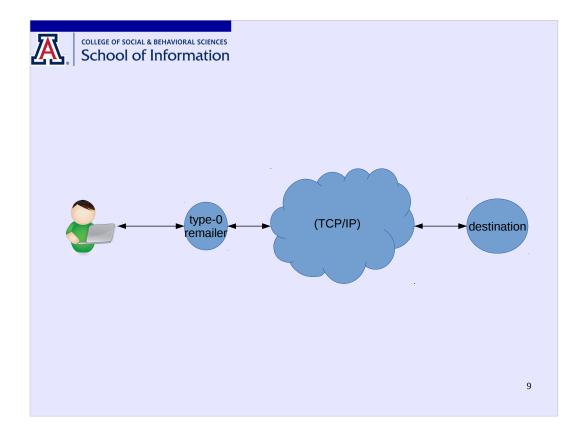


How does the IP addressing scheme support routing? Essentially, with (deferred) table lookups

here 172.16.12.2 sends to 172.16.1.2 (walk through).

notice the source and gateway determine the next hop with a table lookup. They also have a default route for addresses they can't look up.

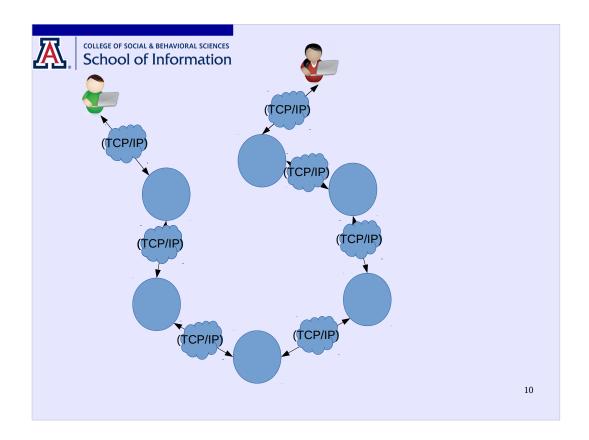
routing information isn't hidden in IP: its only used bit-by-bit, but its all available in principle. But censorship systems often work by ensuring they are on the route, and filtering. So what to do?



One simple answer is to proxy: encrypted traffic is sent to an intermediate node outside of which the censor is not on the route. The censor just sees your connection to some box (ideally, it doesn't know its a proxy), and the proxy handles forwarding your traffic for you.

A variant of this idea is this simplest (degenerate) form of mixnet, which forwards mail.

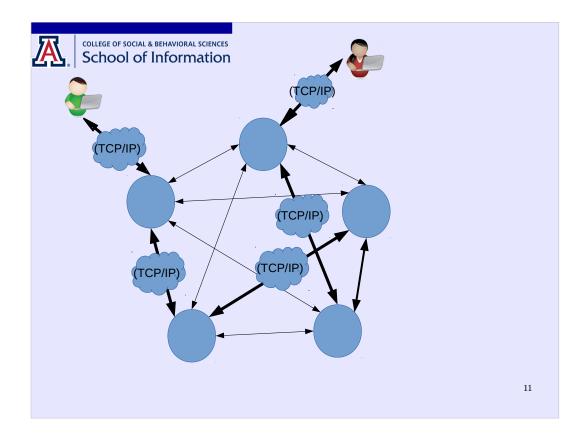
Problems here?



A more sophisticated approach combines several intermediate nodes to ameliorate the problem of centralized trust in one intermediate. Even if some of these intermediates are malicious, the traffic can continue.

Here's the classic topology for mixnets: a cascade. These have a fixed route through all the mixes

Problems here? (Suddenly lots of traffic going to an intermediate node---might it attract attention of censors? Are you hidden as a user of such a system?)



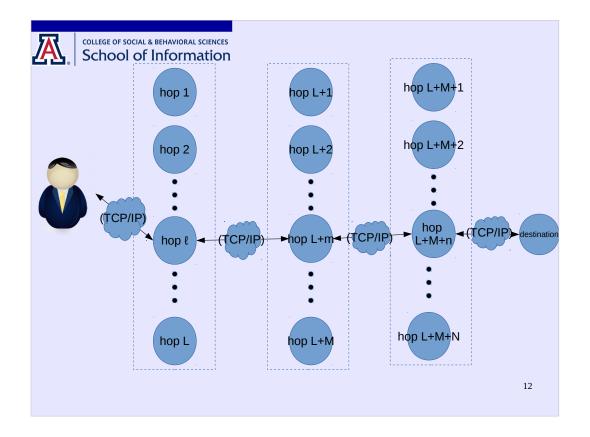
#### free route topology

it is well known that this topology is the most heavy metal (see pentagram). (If you are reading this and are confused: it's a joke).

Old-school onion routing is free-route like this.

more scalable, so more anonymous in practice in some ways; but also less anonymous in some ways than mixes.

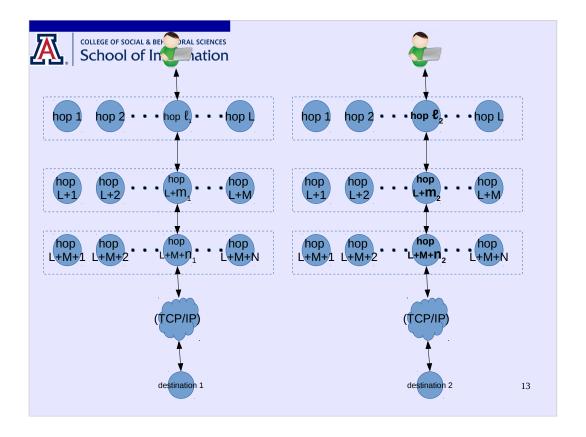
Problems? are the entry points to the network known? scalability and anonymity are tied together



stratified topology. The choice at each layer is random (see next slide)

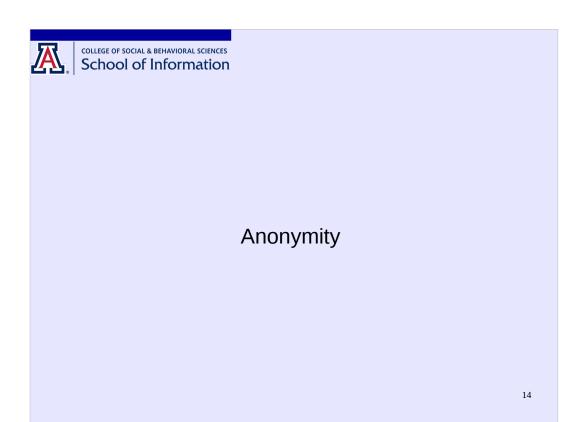
balances scalability with anonymity. Tor does this with helper nodes ("guards"), intermediate nodes ("relays"), and exit nodes.

To address public knowability of entry nodes, Tor uses bridges. An elaboration of this to avoid recognition of tor traffic by traffic analysis is pluggable transports.

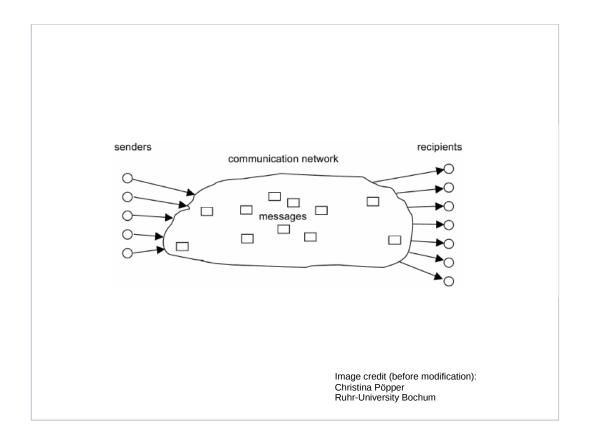


random choice at each layer.

we have talked



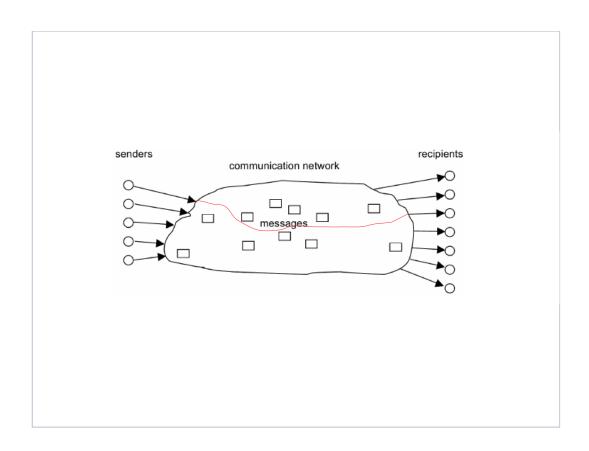
so far I've relied on an intuitive understanding of anonymity, but when we're building technologies for anonymity, it's better to be precise about what is protected and what is not.

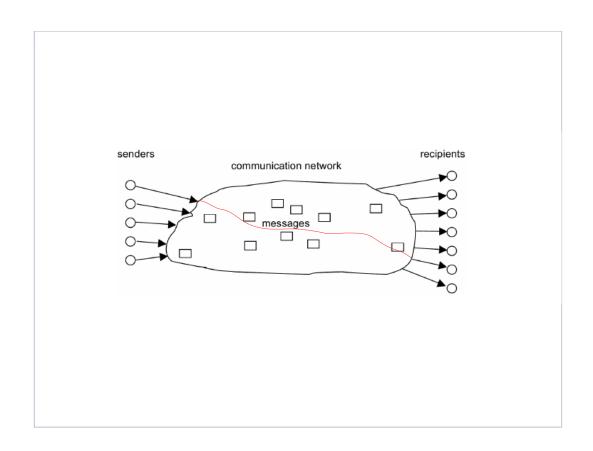


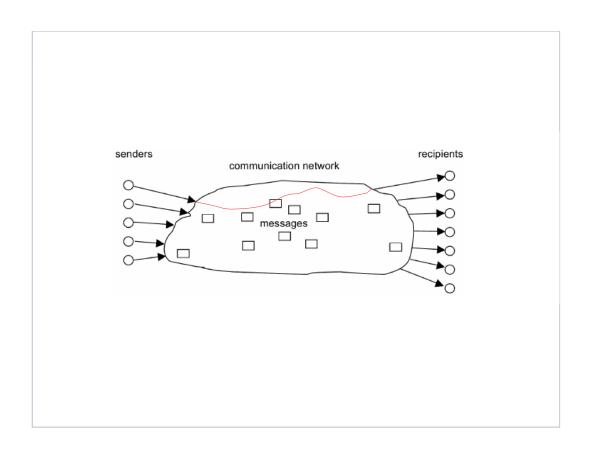
<u>senders</u> communicating <u>messages</u> over a <u>channel</u> with <u>recipient</u>.

An <u>adversary</u> (or <u>attacker</u>) tries to <u>reduce the</u> <u>anonymity</u> of some or all of the parties to the communication

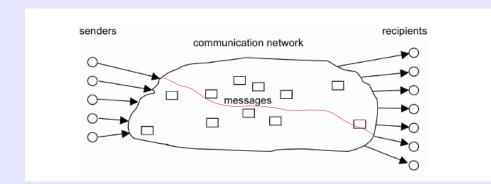
Anonymity is a property of a channel.













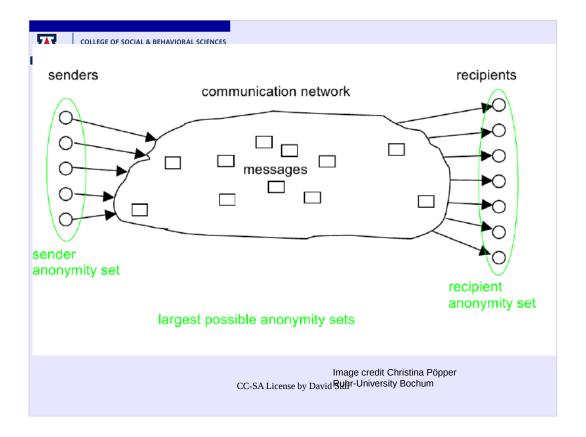
### Anonymity set

 Can you clearly describe the limiting cases for the anonymity set?

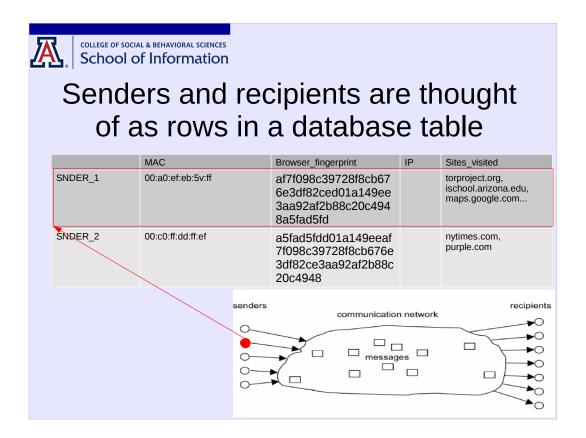
CC-SA License by David Sidi

- 1. singleton. The subject has no anonymity; he is perfectly identifiable relative to the set of subjects.
- 2. universe. The subject has perfect anonymity in the set of users of the system.

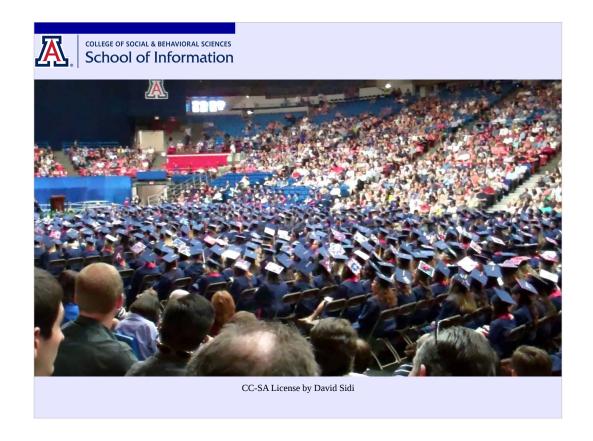
Notice what this means (Berthold's metric)



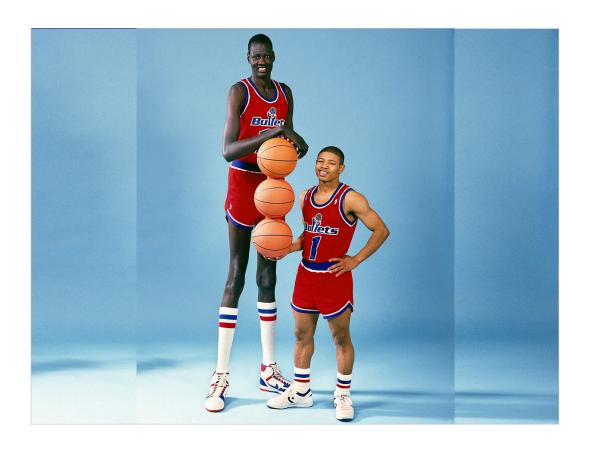
- What is a sender? i.e., how do we get the set of all senders? (Think about the definitions)
  - something that sends messages over the network to recipients (implements protocols, etc.).
     People, personal computers, cameras, phones, etc.?
- if that were all, all senders would be the same! But the anonymity set is intended to be useful, not trivial, in its separation of senders that cannot be distinguished from those that can be
  - senders should have attributes to distinguish them



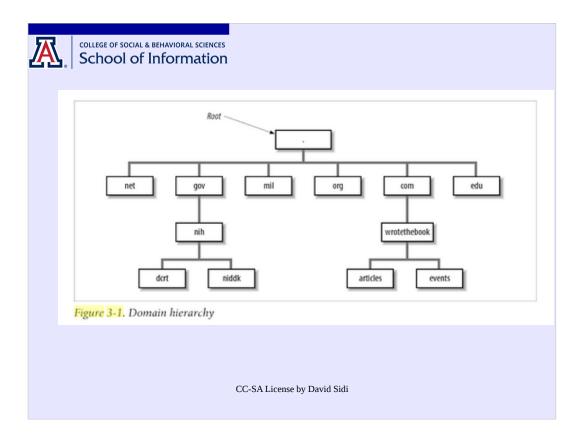
 Behind the anonymity set in Torra's discussion (which is typical) is the database: things that can be identified ("subjects," "persons," etc...) are records in a database---descriptions of people via a set of attribute values



 Suppose I include in a record of UA students a person's weight and height as 150 lbs, 5'3". Is the person anonymous? (think: anonymity set)



- Now suppose further that I do so for a database of male UA basketball players. Obviously, the player is not (as) anonymous.
- Where might you find combinations of attributes that identify people using computer networks?
  - IP, Ethernet, Tor user (Harvard bomb threat example), DNS, third-party tracking, browser fingeprinting
  - Lets look at DNS for a second



- DNS puts a human-friendly layer on top of destination IP addresses, by associating structured strings of words (and abbreviations) with IPs.
- In the beginning was the host table. And it was good (for some things, but not for scalability)
- Distributed, hierarchical
  - again with the layers: root servers have information about TLD servers beneath them
  - Registering a domain name involves telling the TLD servers about your server
  - you can then do subdomains at will
- Forwarding DNS server
- Recursive DNS server (or resolver)
- (Root nameserver)
- (Top Level Domain nameserver)
- Authoritative nameserver
- quick task: get the DNS for arizona.edu using the



#### **DNS**

#### \$dig arizona.edu

- ; <<>> DiG 9.9.5-9+deb8u14-Debian <<>> arizona.edu
- ;; global options: +cmd
- ;; Got answer:
- ;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 14058
- ;; flags: qr rd ra; QUERY: 1, ANSWER: 1, AUTHORITY: 0, ADDITIONAL: 1
- ;; OPT PSEUDOSECTION: ; EDNS: version: 0, flags:; udp: 4096
- ;; QUESTION SECTION:
- ;arizona.edu. IN A
- ;; ANSWER SECTION:
- arizona.edu. 6813IN A 128.196.128.233
- ;; Query time: 32 msec ;; SERVER: 208.67.222.222#53(208.67.222.222)
- ;; WHEN: Mon Oct 23 12:43:03 MST 2017 ;; MSG SIZE rcvd: 56 CC-SA License by David Sidi



## **DNS**

 Suppose I use a VPN to tunnel my traffic to a server I control, but you serve my DNS requests. What can you learn about me?



# Browser fingerprinting

- UserAgent
- Language
- · Color Depth
- Screen Resolution
- Timezone
- · Has session storage or not
- · Has local storage or not
- · Has indexed DE
- Has IE specific 'AddBehavior'
- · Has open DB
- · CPU class
- Platform
- DoNotTrack or not
- Full list of installed fonts (maintaining their order, which increases the entropy), implemented with Flash.

- A list of installed fonts, detected with JS/CSS (sidechannel technique) - can detect up to 500 installed fonts without flash
- · Canvas fingerprinting
- · WebGL fingerprintingPlugins (IE included)
- · Is AdBlock installed or not
- Has the user tampered with its languages 1
- · Has the user tampered with its screen resolution 1
- · Has the user tampered with its OS 1
- Has the user tampered with its browser 1
- · Touch screen detection and capabilities
- Pixel Ratio
- System's total number of logical processors available to the user agent.



# Browser fingerprinting

- · Multi-monitor detection,
- · Internal HashTable implementation detection
- · WebRTC fingerprinting
- Math constants
- Accessibility fingerprinting
- · Camera information
- DRM support
- Accelerometer support
- Virtual keyboards
- List of supported gestures (for touch-enabled devices)
- · Pixel density
- · Video and audio codecs availability
- · Audio stack fingerprinting

CC-SA License by David Sidi

note: this has some uncertainty associated with it. comparison on panopticlick.



# Discussion: "Identifiability"

- Why might it be too simple to say that for a sender S, every other potentially-different sender is either completely indistinguishable from S or not?
- 2 minutes alone, 2 minutes with a partner, then we'll talk as a class

#### Question

 What is problematic about this definition of anonymity? "Anonymity is thus defined as the state of being not identifiable within a set of subjects, the anonymity set." (Danezis and Diaz 3)

 Later, they say "A subject carries on the transaction anonymously if he cannot be distinguished (by an adversary) from other subjects. This definition of anonymity captures the probabilistic information often obtained by adversaries trying to identify anonymous subjects.

#### Question

- What is problematic about this definition of anonymity? "Anonymity is thus defined as the state of being not identifiable within a set of subjects, the anonymity set." (Danezis and Diaz 3)
- We need to say something about an adversary and an attack model

 Later, they say "A subject carries on the transaction anonymously if he cannot be distinguished (by an adversary) from other subjects. This definition of anonymity captures the probabilistic information often obtained by adversaries trying to identify anonymous subjects.



 Definition 1.2 From an adversary's perspective, anonymity of a subject s means that the adversary cannot achieve a certain level of identification for the subject s within the anonymity set. (Torra)



- Torra's terminology is confusing. His own idea of `n-Confusion' is in the background here :-).
- Simplifying, the point is: for each row in the database and some auxiliary information, we have a distribution over the subjects.
- The closer to uniform this distribution is, the "less identifiable" an entity is within the anonymity set, and the better the anonymization



# **Anonymity metrics**

- Reviewing the discussion from a recent breakout session in a workshop on Privacy as Engineering Practice, Deirdre Mulligan said that there is a need for more formal measures of privacy (including anonymity)
- privacy loss in terms of information flow analysis, measures that take into account inference and not only disclosure, etc.



# **Anonymity metrics**

- Degree of anonymity is a distribution 1 p, where p is the probability assigned to the senders in the anonymity set (Reiter and Rubin 1998)
  - I talk about senders just for convenience, this applies more generally
- Think "worst case" -- who's got the highest probability of being identified, and how high is that probability?
- This doesn't account for how evenly distributed the probability is over the anonymity set, in the sense that it just depends on the greatest probability



# **Anonymity metrics**

- Suppose a user u has 0.1 degree of anonymity.
   Consider two scenarios s1 and s2
  - s1: 2 users u and v, with v also 0.1 degree anonymity
  - s2: 1000 users, all users distinct from u with the same degree anonymity (which is less than 0.001)
- With degree of anonymity as measure, both have equal anonymity



# Degree of anonymity

- One way to account for evenness of distribution is with Information Theory
- "Effective size": Roughly, how many bits does the attacker need to identify a member of the anonymity set?

# Degree of anonymity

- "Effective size": Roughly, how many bits does the attacker need to identify a member of the anonymity set?
- Less roughly, use the Shannon entropy. For  $\Psi$  the set of users,  $\mathcal U$  the posterior of a user being the sender given a message,

**Definition 2.** We define the effective size S of an r anonymity probability distribution U to be equal to the entropy of the distribution. In other words

$$S = -\sum_{u \in \Psi} p_u \log_2(p_u)$$

where  $p_u = \mathcal{U}(u, r)$ .



### Degree of anonymity

- Less roughly, use the Shannon entropy.
  - This gives an expected value
- Danezis and Diaz also mention defining degree of anonymity as log<sub>2</sub>(N), for N the number of users
- Can also use min entropy for worst case
- They don't mention it, but all these are related



## Degree of anonymity

- $H_0 = \log |X|$
- $\bullet \quad H_1 = -\sum_{i=1}^n p_i \log(p_i)$
- $H_{\infty} = -\log(\max_{i} p_{i})$

### {Sender, Receiver, Relationship} Anonymity

- Sender/receiver anonymity: For a given message m, what is the probability that m came from sender/receiver A?
  - Note this depends on who the adversary is: the recipient? A global passive adversary? An active adversary? ...
- Relationship anonymity: For a given message m, what is the probability that m came from sender A and went to destination B?
- Relationship\* anonymity: What is the probability that sender A is communicating with destination B?
  - a persistent relationship, not a single message or request exchange
- Question: how can you have sender anonymity but not relationship\* anonymity? Hint: a universal generalization implies all instantiations.

"Suppose the network provides perfect sender anonymity, i.e., any message exiting the network is equally likely to have originated from any active sender. By observing these messages, however, the attacker can easily infer that all of them have the same destination. For every active sender, the attacker can thus determine with 100% certainty that this sender is communicating with the website, completely breaking relationship anonymity."

Shmatikov and Wang, 'Measuring Relationship Anonymity in Mix Networks'



### Anonymity networks



# Anonymity networks address the traffic analysis problem

- Chaum: "Keeping confidential who converses with whom, and when they converse"
- Contrast with secrecy of message content



# Anonymity networks can involve trusted or semi-trusted relays

- Trusted parties are not adversaries: they can break anonymity
- Semi-trusted parties don't all collude



### Trusted relays

- Example: Nym servers
  - a server keeps a dictionary between real and pseudonymous emails
  - request comes to the remailer, which forwards it, gets the response, and returns it to the user
  - Example: anon.penet.fi
- Other Examples: Anonymous proxies (startpage.com), VPNs



### Trusted relays

- Problem: messages are all linked
  - Stylometric attacks: the frequency of function words in the English language can be used in the long term to identify users (Rao & Rohatgi (2000), "Can Pseudonymity Really Guarantee Privacy?")
  - Correspondent sets of each nym
- Anonymity is compromised if one node is compromised. ("Single point of failure.")
  - lots of incentive to coerce
  - or if the node is not honest
- Fails bitwise indistinguishability: sometimes traffic analysis can deanonymize
  - http proxy example
  - timing correlation



### Semi-trusted relays

#### Strengths

- Compromise of more than one is needed, so more coercion resistant than trusted-relay approaches
- "any single mix is able to provide the secrecy of the correspondence between the input and the outputs of the entire cascade" (Chaum)

#### Weaknesses

- Tagging attacks violate unlinkability (blind signing attack)
- replay attacks
- slow (public-key cryptography)

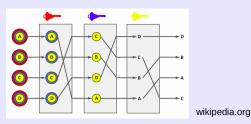


### Semi-trusted relays

 What are the problems with a mixnet with only one node? (Chaum)



- Routing protocol with a cascade of cryptographic relays called 'mixes'
- Mixes only know their neighbors
- User-specifiable routing (Chaum's "new kind of mix")





- Suppose we are at a mix A1, which receives message m.
- m is split into a fixed number blocks, ℓ

```
A_1: [K_{A_1}(R_{A_1}, A_2)], [R_{A_1}^{-1}(K_{A_2}(R_{A_3}, A_3))], \dots, 
[R_{A_1}^{-1}(R_{A_2}^{-1} \cdots R_{A_{n-1}}^{-1}(K_{A_n}(R_{A_n}, A)) \cdots)], 
[R_{A_1}^{-1}(R_{A_2}^{-1} \cdots R_{A_n}^{-1}(M_1) \cdots], \dots, 
[R_{A_1}^{-1}(R_{A_2}^{-1} \cdots R_{A_n}^{-1}(M_{l-n}) \cdots)] \rightarrow.
```

```
A_{2}: [K_{A_{2}}(R_{A_{3}}, A_{3})], [R_{A_{2}}^{-1}(K_{A_{3}}(R_{A_{3}}, A_{4}))], \dots, 
[R_{A_{2}}^{-1}(R_{A_{3}}^{-1} \cdots R_{A_{n-1}}^{-1}(K_{A_{n}}(R_{A_{n}}, A)) \cdots)], 
[R_{A_{2}}^{-1}(R_{A_{3}}^{-1} \cdots R_{A_{n}}^{-1}(M_{1}) \cdots)], \dots, 
[R_{A_{2}}^{-1}(R_{A_{3}}^{-1} \cdots R_{A_{n}}^{-1}(M_{l-n}) \cdots)], [R_{A_{1}}(J_{A_{1}})] \rightarrow,
```

```
A: [M_1], [M_2], \ldots, [M_{l-n}],

[R_{A_n}(R_{A_{n-1}} \cdots R_{A_1}(J_{A_1}) \cdots)], \ldots, [R_{A_n}(J_{A_n})].
```



 The first block is like a header: it contains the key R<sub>A1</sub> and address A2 for the next hop. This is stripped off of the message, and a padding ("junk") block is added to the end

```
A_{1}: [K_{A_{1}}(R_{A_{1}}, A_{2})], [R_{A_{1}}^{-1}(K_{A_{2}}(R_{A_{2}}, A_{3}))], \dots, 
[R_{A_{1}}^{-1}(R_{A_{2}}^{-1} \cdots R_{A_{n-1}}^{-1}(K_{A_{n}}(R_{A_{n}}, A)) \cdots)], 
[R_{A_{1}}^{-1}(R_{A_{2}}^{-1} \cdots R_{A_{n}}^{-1}(M_{1}) \cdots], \dots, 
[R_{A_{1}}^{-1}(R_{A_{2}}^{-1} \cdots R_{A_{n}}^{-1}(M_{l-n}) \cdots)] \rightarrow.
```

```
A_{2}: [K_{A_{2}}(R_{A_{2}}, A_{3})], [R_{A_{2}}^{-1}(K_{A_{3}}(R_{A_{3}}, A_{4}))], \dots, 
[R_{A_{2}}^{-1}(R_{A_{3}}^{-1} \cdots R_{A_{n-1}}^{-1}(K_{A_{n}}(R_{A_{n}}, A)) \cdots)], 
[R_{A_{2}}^{-1}(R_{A_{3}}^{-1} \cdots R_{A_{n}}^{-1}(M_{1}) \cdots)], \dots, 
[R_{A_{2}}^{-1}(R_{A_{3}}^{-1} \cdots R_{A_{n}}^{-1}(M_{1-n}) \cdots)], [R_{A_{1}}(J_{A_{1}})] \rightarrow,
```

```
A: [M_1], [M_2], ..., [M_{i-n}], [R_{A_n}(R_{A_{n-1}} \cdots R_{A_1}(J_{A_1}) \cdots)], ..., [R_{A_n}(J_{A_n})].
```

CC-SA License by David Sidi

•



 The rest of the blocks are, first, the header blocks for all remaining routers in the cascade, and next, the message. All of these are encoded using .

```
A_{1}: [K_{A_{1}}(R_{A_{1}}, A_{2})], [R_{A_{1}}^{-1}(K_{A_{2}}(R_{A_{2}}, A_{3}))], \dots, \\ [R_{A_{1}}^{-1}(R_{A_{2}}^{-1} \cdots R_{A_{n-1}}^{-1}(K_{A_{n}}(R_{A_{n}}, A)) \cdots)], \\ [R_{A_{1}}^{-1}(R_{A_{2}}^{-1} \cdots R_{A_{n}}^{-1}(M_{1}) \cdots], \dots, \\ [R_{A_{1}}^{-1}(R_{A_{2}}^{-1} \cdots R_{A_{n}}^{-1}(M_{l-n}) \cdots)] \rightarrow. 
A_{2}: [K_{A_{2}}(R_{A_{2}}, A_{3})], [R_{A_{2}}^{-1}(K_{A_{3}}(R_{A_{3}}, A_{4}))], \dots, \\ [R_{A_{2}}^{-1}(R_{A_{3}}^{-1} \cdots R_{A_{n-1}}^{-1}(K_{A_{n}}(R_{A_{n}}, A)) \cdots)], \\ [R_{A_{2}}^{-1}(R_{A_{3}}^{-1} \cdots R_{A_{n}}^{-1}(M_{1}) \cdots)], \dots, \\ [R_{A_{2}}^{-1}(R_{A_{3}}^{-1} \cdots R_{A_{n}}^{-1}(M_{l-n}) \cdots)], [R_{A_{1}}(J_{A_{1}})] \rightarrow, 
A: [M_{1}], [M_{2}], \dots, [M_{l-n}], \\ [R_{A_{n}}(R_{A_{n-1}} \cdots R_{A_{1}}(J_{A_{1}}) \cdots)], \dots, [R_{A_{n}}(J_{A_{n}})].
```



- A1 uses the R<sub>A1</sub> it now has to decode the (ℓ-1) blocks after the header in the original message: these are the first part of the message sent out from A1, they contain the headers for A2, the encoded headers for A3,...An, and then the encoded message
- The blocks are passed to the next node, which could be another mix

```
A_1: [K_{A_1}(R_{A_1}, A_2)], [R_{A_1}^{-1}(K_{A_2}(R_{A_2}, A_3))], \dots, \\ [R_{A_1}^{-1}(R_{A_2}^{-1} \cdots R_{A_{n-1}}^{-1}(K_{A_n}(R_{A_n}, A)) \cdots)], \\ [R_{A_1}^{-1}(R_{A_2}^{-1} \cdots R_{A_n}^{-1}(M_1) \cdots], \dots, \\ [R_{A_1}^{-1}(R_{A_2}^{-1} \cdots R_{A_n}^{-1}(M_{l-n}) \cdots)] \rightarrow.
```

```
A: [M_1], [M_2], ..., [M_{l-n}], [R_{A_n}(R_{A_{n-1}} \cdots R_{A_1}(J_{A_1}) \cdots)], ..., [R_{A_n}(J_{A_n})].
```



- Mixes only know their neighbors. (Question: Why?)
- All nodes have a public key
- Weaknesses
  - active attacks: tagging attacks (blind signing attack), replay attacks
  - slow (public-key cryptography, latency in anonymous remailers).

```
A_{1}: [K_{A_{1}}(R_{A_{1}}, A_{2})], [R_{A_{1}}^{-1}(K_{A_{2}}(R_{A_{2}}, A_{3}))], \dots, 
[R_{A_{1}}^{-1}(R_{A_{2}}^{-1} \cdots R_{A_{n-1}}^{-1}(K_{A_{n}}(R_{A_{n}}, A)) \cdots)], 
[R_{A_{1}}^{-1}(R_{A_{2}}^{-1} \cdots R_{A_{n}}^{-1}(M_{1}) \cdots], \dots, 
[R_{A_{1}}^{-1}(R_{A_{2}}^{-1} \cdots R_{A_{n}}^{-1}(M_{l-n}) \cdots)] \rightarrow.
```

```
A_{2}: [K_{A_{2}}(R_{A_{2}}, A_{3})], [R_{A_{2}}^{-1}(K_{A_{3}}(R_{A_{3}}, A_{4}))], \dots, 
[R_{A_{2}}^{-1}(R_{A_{3}}^{-1} \cdots R_{A_{n-1}}^{-1}(K_{A_{n}}(R_{A_{n}}, A)) \cdots)], 
[R_{A_{2}}^{-1}(R_{A_{3}}^{-1} \cdots R_{A_{n}}^{-1}(M_{1}) \cdots)], \dots, 
[R_{A_{2}}^{-1}(R_{A_{3}}^{-1} \cdots R_{A_{n}}^{-1}(M_{l-n}) \cdots)], [R_{A_{1}}(J_{A_{1}})] \rightarrow,
```

```
A: [M_1], [M_2], \ldots, [M_{i-n}],

[R_{A_n}(R_{A_{n-1}} \cdots R_{A_1}(J_{A_1}) \cdots)], \ldots, [R_{A_n}(J_{A_n})].
```



## Mixing techniques for Mixnets

- Cascading: All nodes are always used, in the same order
- Scalability is a problem, requires setting up a fixed route with all nodes
- Only requires one honest node to preserve anonymity



### Mixing techniques for Mixnets

- User specified: user arbitrarily picks its route through the network
- Scalable, does not require initial configuration of a route
- Not anonymous if only one node is honest (nodes can figure out their positions)

CC-SA License by David Sidi

relays can determine their position in the chain, which can be used to deanonymize with enough collusion