

# DIMACS Working Group on Measuring Anonymity

## Notes from Session 3: Information Theoretic and Language-based Approaches

Scribe: Matthew Wright

In this session, we had three 15-minute talks based on submitted abstracts and about 45 minutes of "panel" discussion with the three speakers as panelists. The focus of the session was on information theoretic approaches to measuring anonymity, plus a novel language-based approach to using anonymity.

Talk #1: **Parv Venkatasubramaniam**. [\*Anonymity in Tor-like systems under Timing Analysis: An Information Theoretic Perspective\*](#). (15 min.)

### Basic Idea/Issue

- End-to-end timing attacks (confirmation)
  - Adversary monitors all the links?
  - What is the likelihood of this?
- Need an objective to evaluate the building blocks = the routing/mixing
- Adversary assumptions
  - Prior knowledge: traffic statistics, likelihood of a link
  - System observation: timing on all links
  - Network strategy: (observation points? missed this one)
  - Observation includes all timing (past, present, future)
  - Adversary is aware of the strategy
- Information theoretic approach: JOINT distribution
  - $\Pr(\text{these are the specific sources of the packets} \mid \text{obs., etc.})$
- Entropic measure is useful
  - Fano's inequality: Shannon entropy provides lower bound on probability of error
  - Can integrate prior information using Bayes' Theorem.

### Our contributions

- We use a Poisson processes for arrivals and create optimal reordering strategies for a latency- and buffer-constrained system.
  - Fundamental trade-off between anonymity and QoS.
- Anonymity of a mix-net is a linear fn. of anon. of the individual mixes.
- Packets vs. Streams:
  - For long streams, need dummies. Lots!, e.g. for DLP.
  - What if the stream is short lived? Maybe less padding?
  - *Admissible length*: how long can the stream live w/ perfect anon?
- Poisson assumption? Maybe this is not so realistic -- neither users nor websites are memoryless. Would be interesting to see this approach using other assumptions.
- Reordering: is this really possible?

- What is important is that order is hidden.
- [Matt's note: burstiness in traffic makes this hard, too]

Talk #2: **Kostas Chatzikokolakis**. *Information theory and decision theory to measure information leakage [using gain functions]*. (15 min.)

Information theoretic definitions must say something meaningful about your application to be viable.

Quantitative information flow: measure how much information is leaked by the system. [think covert channels or data query privacy]

Model

- Channel  $C$ , input  $x$ , output  $y$  --  $C[x, y]$ : prob. of  $y$  given  $x$
- Inputs governed by a prior probability. apply Bayes.

Leakage

- Attacker tries to guess the secret ( $x$ ) in one try. Chance of success?
- Prior vuln: just given the prior prob.
- Post: add the channel leakage
- Leakage: difference from Post and Prior (*min-entropy leakage*)

Limitations

- What about partial guessing, a property/part of a secret, multiple tries, or other aspects?
- Example: a toy channel that reveals exactly one sender (no more or less)
  - If attacker needs to guess the whole list: lost 10 of 10240 bits
  - Guess the receiver of a particular sender: lost 1.016 out of 10 bits
  - Guess just the receiver of *any* sender: lost all bits!

Gain functions

- Attacker makes one guess about the secret
- The benefit is a gain function
- Success measure: The expected gain of a best guess.
- Benefit: Model a variety of attackers and operational scenarios.
  - e.g.: can do approximate guessing, property of a secret (gender/country), part of a secret (part of a location, IP), multiple tries.
- Theorem:  $g\text{-capacity} \leq \text{min-capacity}$  for all gain functions  $g$
- Min-capacity is an upper bound on Shannon capacity

This essentially brings the existing work in decision theory to anonymity.

Talk #3: **Aslan Askarov** and Stephen Chong. [Towards Language-Based Network Anonymity](#). (15 min.)

Applications (browser, ssh, etc.)

- What if the application knows about anonymity underneath?
  - App may realize that anonymity isn't needed
  - App may realize that direct communication is required
- Why anonymity?
  - Want to hide communications from network (e.g. ISP)
  - Want to remain anonymous from the receiver

Programming languages techniques

- ISP case: soundly infer such messages
- Receiver case: ... [what was this?]

Model

- Anonymous communications as a primitive
- Ensure that they are used securely
  - Ex: online auction. participation is public, winner is secret
- Can infer using information flow that the winner declaration must be anonymous.
  - Ex2: EasyChair. Author button should be anonymous connection.
- App figures out when you need privacy and uses an anon connection then (and only then) to improve performance.

Measuring anonymity

- Given a network anonymity metric X
- Can we be sure that the anonymity doesn't go down?

**Panel Discussion** (45 min.)

Note: Discussion participants are labeled 'A' to 'Z' for each question.

To Aslan: Suppose I don't want NYT advertisers to know who I am, but NYT is OK. Can I do this without a leak given possible side channels due to simultaneous loading?

- A: [yes] The inputs are governed by a prior probability. Apply Bayes.

To Aslan: If you weren't able to split anon and non-anon, this wouldn't help much, right?

- A: yes. Reframed to be positive -- it's an anonymity-preserving optimization of your application.
- B: Does the programmer need to know something about the anonymity service? A: Yes.
- C: there can be different specialized applications requiring different anonymity levels, and this could be taken to the transport level too.

[To Aslan] If all the traffic is especially sensitive, did we make it a bigger target to attack? e.g. in the auction protocol, if you look at the timing, you can see the auction winner is the only person who's going to use anonymous channels.

- A: To some degree, people already do this. You don't just turn on Tor 100% of the time. May need to move away from this idealistic goal.
- B: An example of this is Tor set up w/ DNS not Tor-ized. This obviously leaks your connections to the DNS (e.g. your ISP).
- C: If you start at 100% Tor, and you run our idea on your application, you can back off without losing some anonymity.
- D: Isn't the point of something like Tails that you can be sure to run everything over Tor without screwing up? This approach seems backwards relative to that.
- E: When crypto was weak, you had to use message discipline and be cautious about what you sent. Same here -- the more you use the system, the bigger/clearer the fingerprint you leave behind.
- F: That depends on the attacker model.
- G: e.g. if you have two gmail accounts, one is public and one is sensitive, it's better to access the public one without using Tor.
- A: You have an anonymity budget. The more you use, the more you lose, and then you need a refresh.
- H: trying to do this [what?] in the Dissent project.

To Kostas: Is there a reason to use expectation in the gain functions (vs. min)?

- A: This depends on the application.
- Very unlikely, but revealing, events are OK (or we can live with them).
  - Differential privacy will not say that. It will say the unlikely event is very bad.
- This is an extension of information theoretic definitions.
- B: (to Kostas) you should not use putting in your password as an example

Should we really be developing an anonymity-usability metric?

- A: I agree with this idea. There is no free lunch. Where on the spectrum do you land [b/w high anonymity and high usability]?
- B: Anonymity depends on other users' activity. If you have a system with very sensitive users, the anonymity set is small. You need to make the cost for users as low as possible to get users who care less.
- C: Sure, if you divide people up, it's a problem.
- B: The system can never get started.
- A: What about a model in which the users also provide the resources (P2P anonymity)?
  - D: Oh, that's a world of pain.
    - Ref. to the many broken designs in P2P anonymity.
  - E: And not so friendly to many users.

None of the metrics we've talked about today are any good at dealing w/ active adversaries. This is a big problem.

- A: that's not entirely true. Strategies have some tolerance to some level of active attacks. That tolerance is a measure.\*
- B: What about DoSing, degrading, etc.? These don't account for that.
- C: If it is assuming a distribution in advance, then yes. But if the metric can be done based on what is observed, then it can be OK.

\* Diaz+ have a HotPETs paper that examines this:

<http://www.cosic.esat.kuleuven.be/publications/article-2320.pdf>