# Content Availability, Pollution, and Poisoning in File Sharing Peer-to-Peer Networks

Nicolas Christin

SIMS, UC Berkeley

christin@sims.berkeley.edu

Andreas S. Weigend

Weigend Associates LLC

andreas@weigend.com

John Chuang

SIMS, UC Berkeley

chuang@sims.berkeley.edu

100 x 100

# Background

- Several petabytes of content present at any time in file sharing networks, but…

- Vast amounts of useless files (Liang *et al.*, 2005)
  - Poorly encoded or corrupted
  - Incorrect or misleading metadata
  - …

- Signal-to-noise ratio can be extremely low…

*Can we rely on injecting useless content to impact usage of file sharing networks?*

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*2*

100 X 100

# Motivation

- **Possible defense mechanism against copyright infringement in P2P networks**
  - Some companies specialize in injection of noise
    - Overpeer, Retspan, Macrovision…

- **Viable technological alternative to legal recourse?**
  - Difficult to prosecute individual users

- **Injection of useless content does *not* require monitoring, or intrusion**
  - Probably much more acceptable than most other interdiction methods in the eye of the general public
  - Does not require new "safe harbor" laws (H.R. 5211)

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*3*

# Related work

- **Bird's eye view of network measurements**
  - Effect on backbone (Sen and Wang, 2002)
  - Prevalence of P2P traffic (Saroiu *et al.*, 2002)
  - Traffic not decreasing (Karagiannis *et al.*, 2004)

- **Topological properties of P2P file sharing networks**
  - Gnutella (Loo *et al.*, 2003)
  - KaZaA (Liang *et al.*, 2004)
  - eDonkey (Tutschku, 2004, Le Fessant *et al.*, 2004)
  - …

- **Works on pollution/poisoning still rare**
  - Quantification of the phenomenon (Liang *et al.*, 2005)
  - Theoretical studies of potential attacks on P2P networks (Castro *et al.*, 2002, Dumitriu *et al.*, 2005)

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*4*

100 X 100

# Pollution vs. Poisoning

- ## Network pollution
  - *Accidental* injection of unusable or low quality files
    - Happens with most (all?) content
    - Truncated, poorly encoded, …
    - Difficulties in properly "ripping" content

- ## Item poisoning
  - *Deliberate* injection of decoys to render usable files hard to find
    - Targets specific content
    - e.g., "American Life" by Madonna
  - Currently most popular interdiction technique

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*5*

# Research questions

- **Above which level does pollution pose serious problems?**

- **Which (if any) poisoning techniques are effective?**
  - Flooding?
  - More elaborate techniques?

- **We'll look at the most popular P2P networks**
  - FastTrack (KaZaA), eDonkey, Overnet, Gnutella
  - not BitTorrent – does not have built-in search mechanism (yet)

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*6*

100 X 100

# Availability vs. perceived availability



Content replication = Number of peers that share a given file on the network

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

7

# Availability vs. perceived availability



Perceived content replication = Number of peers *that I see* sharing a given file on the network

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

8

# Availability vs. perceived availability



*What matters is not what **is** in the network, but what users **see** from the network*

# Differing perceptions of content

- Ideally all P2P nodes should have same view of content available on the network

- In practice, different nodes have very different perceptions of content availability
  - Peers coming and going $\Longrightarrow$ Content volatility
  - Size of the network/decentralized nature imposes fish-eye view

- User view of the network conditioned by query returns

- Query returns highly dependent on P2P network topology

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*10*

# P2P topologies

- Most modern P2P networks use 2-level hierarchical structure
  - Leaf nodes
  - Hubs (a.k.a. supernodes, ultrapeers, servers)
    - Higher processing power, link capacity, longer uptime…
    - Act as a centralized index for a number of leaf nodes

- Exception: Overnet
  - Distributed Hash Table (all peers are equal)
  - However, Overnet clients are also part of the eDonkey network

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*11*

# Differences in topological structures

| | eDonkey | FastTrack | Gnutella |
|---|---|---|---|
| # of nodes | ≈ 2,800,000 | ≈ 2,500,000 | ≈ 1,000,000 |
| # of hubs | 40—90 | 25,000—40,000 | 10,000—100,000 |
| Fraction of hubs | ≈ 0.00002 | ≈ 0.015 | ≈ 0.05 |
| Avg. leaf-hub connection lifetime | ≈ 24 hours | ≈ 30 minutes | ≈ 90 minutes |
| Leaf promotion | Voluntary | Election | Election |

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*12*

# Differences in topological structures

| | eDonkey | FastTrack | Gnutella |
|---|---|---|---|
| # of nodes | ≈ 2,800,000 | ≈ Semi-centralized network | 000,000 |
| # of hubs | 40—90 | 25,000—40,000 | 10,000—100,000 |
| Fraction of hubs | ≈ 0.00002 | ≈ 0.015 | ≈ 0.05 |
| Avg. leaf-hub connection lifetime | ≈ 24 hours | ≈ 30 minutes | ≈ 90 minutes |
| Leaf promotion | Voluntary | Election | Election |

**Semi- centralized network**

**Hubs are much more stable**

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005
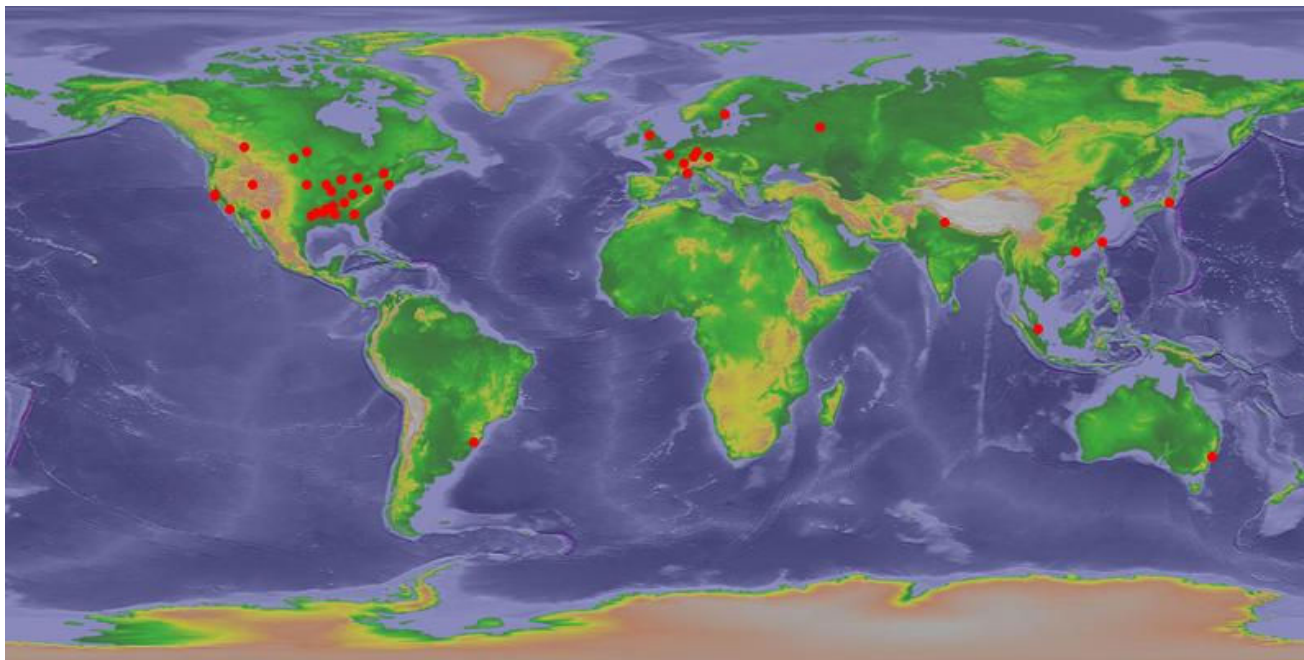
*13*

100 x 100

# Methodology

- **Perception of availability depends on time and origin of a query**
    - Need to measure from different vantage points and at different times

1. Measure content availability *in absence* of poisoning
2. Evaluate effect of pollution and poisoning on measured data by numeric simulation

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*14*

**100 X 100**

# Measurement infrastructure

- giFT-FastTrack and MLDonkey clients
  - Linux console (text-based) applications
  - Allows for scripting
- Easy to run large scale experiments
  - 50 host machines over 18 different countries (PlanetLab)

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*15*

# Active measurements

- Present network with input (queries)
  - 6 movies, 6 songs, 3 software titles
  - Specialized queries (e.g., "filetype = MP3") whenever possible
  - Content not subject to any (noticeable) ongoing poisoning attack
  - Each query is issued every half-hour for 36 hours
  - For each of the four P2P networks considered, each query is sent from at least six machines
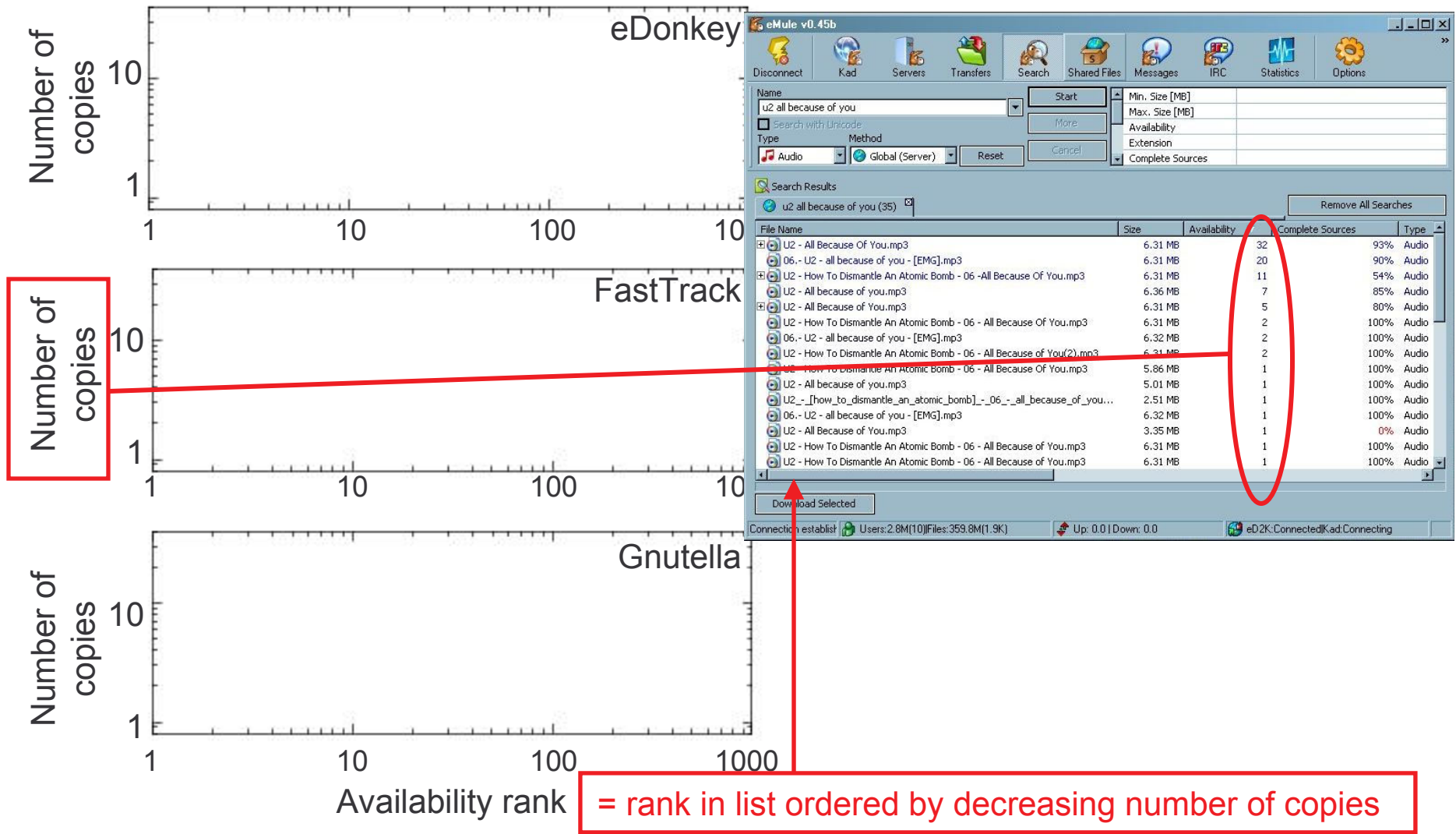
*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*16*

# Summary of measurements w/o poisoning

- ## Semi-centralized topologies (eDonkey)
  - Content remains present in the network for a while
  - Faster responses to queries
- ## FastTrack and Gnutella
  - Relatively low content stability
    - content comes and goes frequently
  - Apparently high levels of pollution
    - even when no poisoning
  - Manage to only download a few files
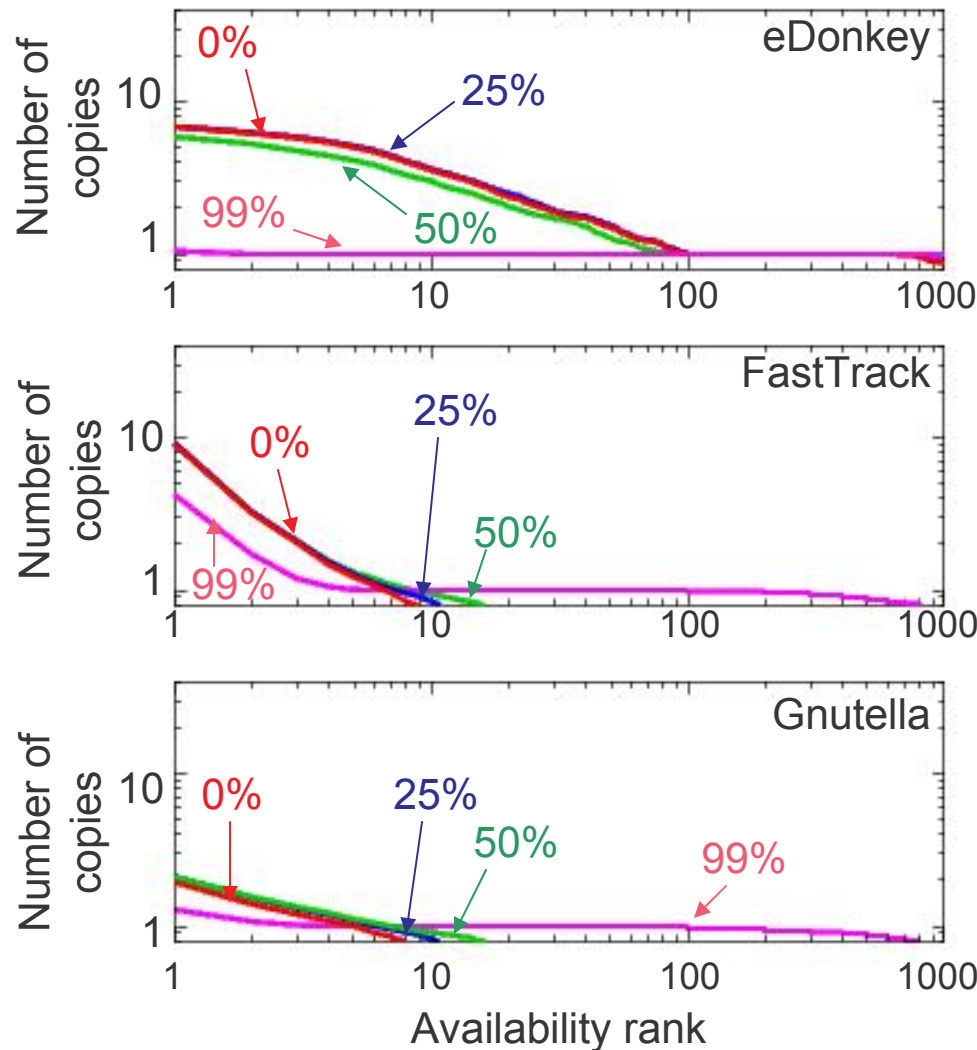  - Confirms findings of (Liang *et al.*, 2005)

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*17*

100 X 100

# Effects of pollution

- **Pollution modeled as injection of random noise in the system**
  - Make $x\%$ of the query returns (uniformly) random for each measurement sample
  - Neglects propagation effects of polluted content

- **Simplest poisoning technique (flooding) is nothing more than pollution at high levels**
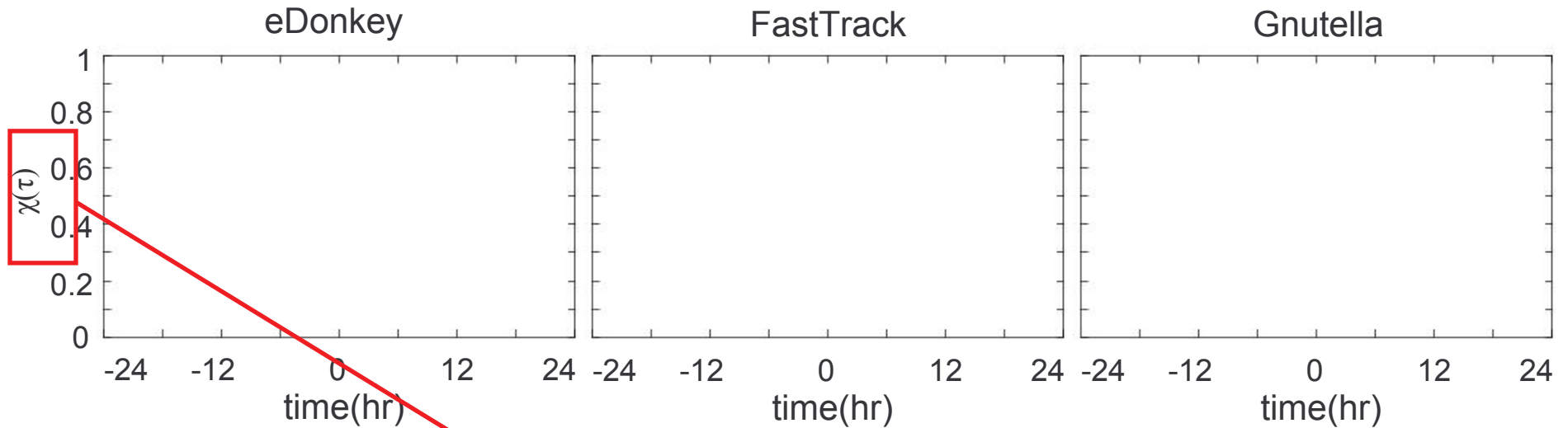  - Should not, *in theory*, reduce availability of useful files

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*18*

100 X 100

# Pollution and perceived availability



Number of copies — eDonkey

Number of copies — FastTrack

Number of copies — Gnutella

Availability rank = rank in list ordered by decreasing number of copies

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

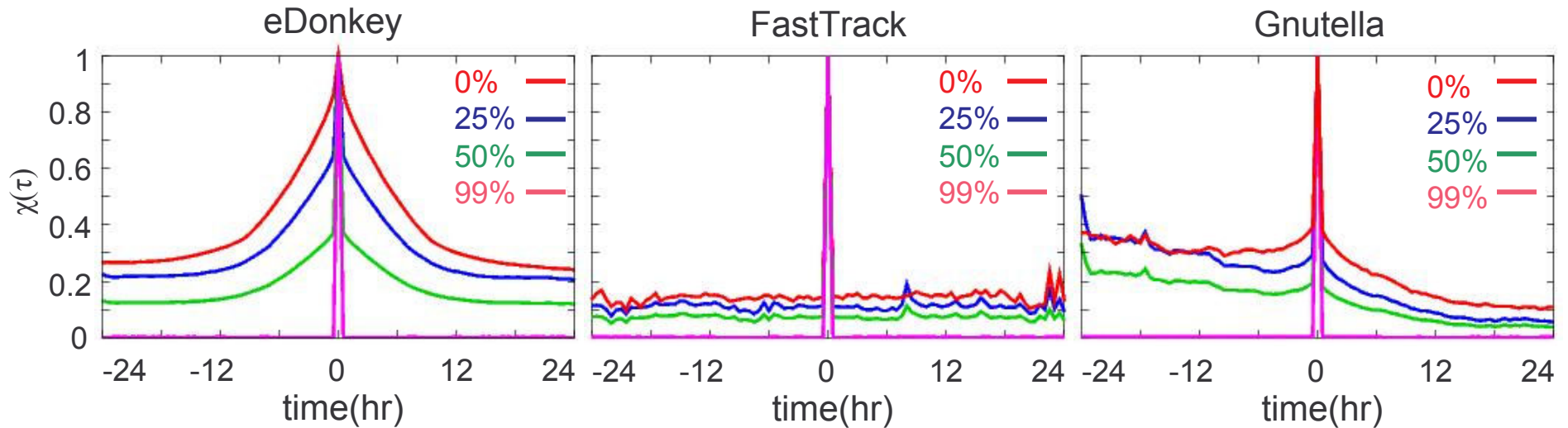*19*

# Pollution and perceived availability



- Pollution only harmful at (very) high levels
- Decoys *may* drive usable files out of the query returns
  - Number of query returns is limited
    - FastTrack example:
      - At most 200 returns for a given query
      - No more than 5 queries in a row
- Poisoning by flooding not particularly efficient
  - e.g., need to insert 99 times as many decoys as existing files
  - … at each hub

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*20*

# Flooding signature

## eDonkey

## FastTrack

## Gnutella

$\chi(\tau)$

time(hr)

time(hr)

time(hr)

$\chi(\tau)$: average probability (over all times, all clients) that an item (specific file) returned at a given time $T$ is also returned at time $T+\tau$

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*21*

100 x 100

# Flooding signature



eDonkey     FastTrack     Gnutella

- **High-levels of pollution (or poisoning by flooding) completely destroys temporal stability**

- **Flooding attack easy to thwart by giving precedence to items that have been seen in the network for some time**

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*22*

# Alternatives to flooding

- **More advanced poisoning techniques can be much less expensive and more efficient than flooding**
  - A (rather detailed) list of attacks is available in a patent application from Macrovision
    - Discussed at http://mvsn-patent-app.notlong.com

  - Chunk corruption
  - Malicious routing
  - Skewing perceived availability to bias users towards downloading useless content
  - …

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*23*

# Targeting perceived availability



*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*24*

# Targeting perceived availability

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

25

# Targeting perceived availability

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*26*

# Targeting perceived availability

- Inject a few highly replicated decoys rather than random files

- Can in addition make replicated decoys harder to detect by frequently changing them (transient decoys)

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*27*

# Replicated decoy injection



- **Insert 30 decoys with the same number of copies as most replicated file**

- **Drives useful files out of the picture**

- **Here only requires about 300 decoys**
  - as opposed to ~9900 for flooding

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*28*

# Temporal signatures



- ## Using permanent replicated decoys leaves a rather obvious signature on the temporal stability

- ## Can be solved by frequently changing the (replicated) decoys

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*29*

# Poisoning antidotes

- **Ranking by availability**
  - Simplest technique
  - Efficient against random noise (if no propagation)

- **Static reputation system**
  - "File X is useless," "IP address Y injects useless content"
  - Needs manual input, far from comprehensive
  - http://www.jugle.net, http://bitzi.com
- **Dynamic ((semi-)automated) reputation system**
  - Weighs reputation of a file as a number of factors
    - Manual input
    - Time present in the system
  - Semi-automate ban of poisoning sources
  - Unlikely such systems are *currently* deployed

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*30*

100 X 100

# Antidotes and their effectiveness

| | Pollution | Flooding | Replicated decoys | Replicated, transient decoys |
|---|---|---|---|---|
| Ranking by number of replicas found | Yes | Somewhat | No | No |
| Static reputation | Somewhat | No | Yes | No |
| Dynamic reputation | Somewhat | Somewhat | Yes | Somewhat |

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*31*

# The poisoning arms race

## P2P designers

- Need to use several antidotes in conjunction
  - e.g., ranking by number of replicas with reputation
- Efficiency of reputation systems improved by looking at statistical characteristics
  - Temporal stability signatures

## Copyright holders

- Brute force never a bad choice
  - Can be devastating if used with proper (combination of) strategies
- Clever techniques can use the reputation system to catalyze poisoning
  - False positives
  - False negatives

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*32*

# Summary

- **Network topology plays a crucial role in how users perceive content**
  - (Semi-)centralized topologies provide more stable content
- **Easy to combat (involuntary) pollution**
  - E.g., ranking results by number of replica found
- **More advanced poisoning strategies harder to thwart**
  - Arms race between poisoning techniques and reputation systems

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*33*

# Conclusion

*Can we rely on injecting useless content to impact usage of file sharing networks?*

It is far from impossible…

... and it avoids putting anyone in jail!

*Content Availability, Pollution and Poisoning in File Sharing Peer-to-Peer Networks*
Sixth ACM Conference on Electronic Commerce (EC'05) - Vancouver, BC, Canada, June 6, 2005

*34*