

MAPAP

Measurement and Analysis of P2P Activity Against Paedophile Content

Raphaëlle Nollez-Goldbach

Project Administrator - CNRS

Safer Internet Plus

<http://antipaedo.lip6.fr>

October 2007 – October 2009

Project Presentation

Project Presentation

Partners:

- ▶ CNRS (France), Computer Science
- ▶ University College Cork (Ireland), Psychology
- ▶ University of Ljubljana (Slovenia), Social Sciences
- ▶ Nobody's Children Foundation (Poland), NGO

Project Presentation

Partners:

- ▶ CNRS (France), Computer Science
- ▶ University College Cork (Ireland), Psychology
- ▶ University of Ljubljana (Slovenia), Social Sciences
- ▶ Nobody's Children Foundation (Poland), NGO

Main goals:

- ▶ Methods and tools to protect peer-to-peer users
- ▶ Accurate information on paedophile exchanges
- ▶ Helping law enforcement institutions and NGOs

Project Presentation

Partners:

- ▶ CNRS (France), Computer Science
- ▶ University College Cork (Ireland), Psychology
- ▶ University of Ljubljana (Slovenia), Social Sciences
- ▶ Nobody's Children Foundation (Poland), NGO

Main goals:

- ▶ Methods and tools to protect peer-to-peer users
- ▶ Accurate information on paedophile exchanges
- ▶ Helping law enforcement institutions and NGOs

Actions:

- ▶ Large-scale measurements
- ▶ Analysis of P2P exchanges

Outline

1. P2P System: eDonkey
2. Measurement
3. Data Processing

1. P2P System: eDonkey

Peer → *Keywords* → *Server*
Server → *File* → *Peer*
Peer → *File_{id}* → *Server*
Server → *List of peers* → *Peer*

1. P2P System: eDonkey

Peer → *Keywords* → *Server*

1. P2P System: eDonkey

Peer → *Keywords* → *Server*
Server → *File* → *Peer*

1. P2P System: eDonkey

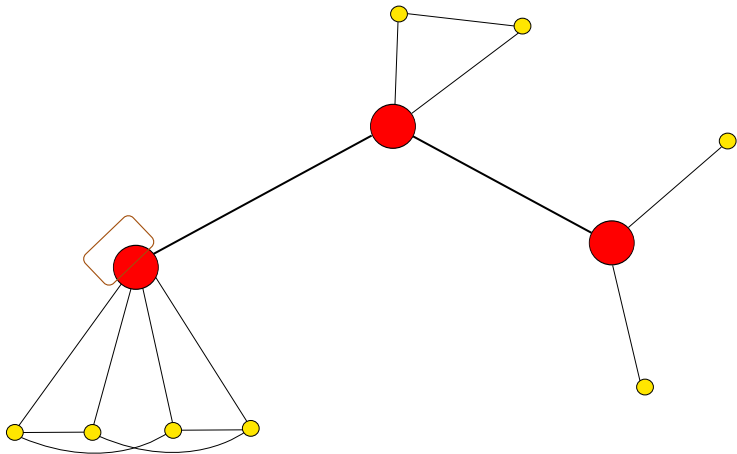
Peer → *Keywords* → *Server*
Server → *File* → *Peer*
Peer → *File_{id}* → *Server*

1. P2P System: eDonkey

Peer \longrightarrow *Keywords* \longrightarrow *Server*
Server \longrightarrow *File* \longrightarrow *Peer*
Peer \longrightarrow *File_{id}* \longrightarrow *Server*
Server \longrightarrow *List of peers* \longrightarrow *Peer*

2. Measurement

Method 1: measure on server



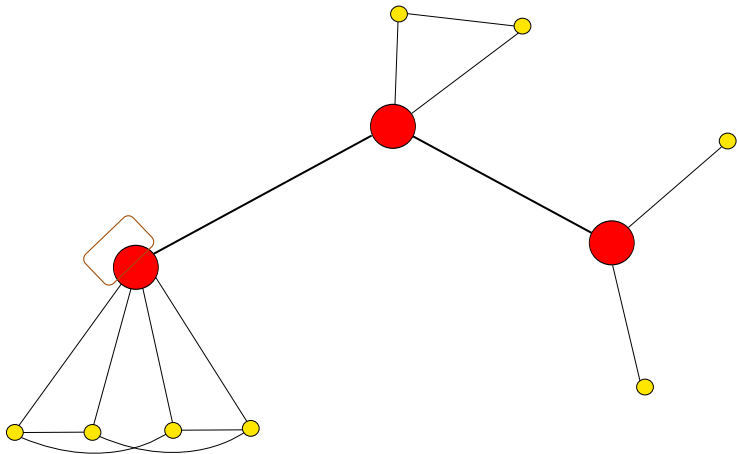
2. Measurement

Method 1: measure on server

- ▶ 10 weeks of measurement
- ▶ 90 millions of peers
- ▶ 275 millions of file-ids
- ▶ Largest measurement of P2P

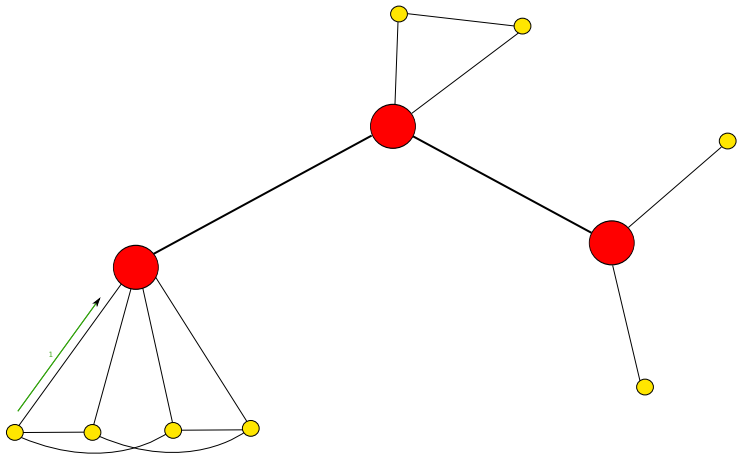
2. Measurement

Method 1: measure on server



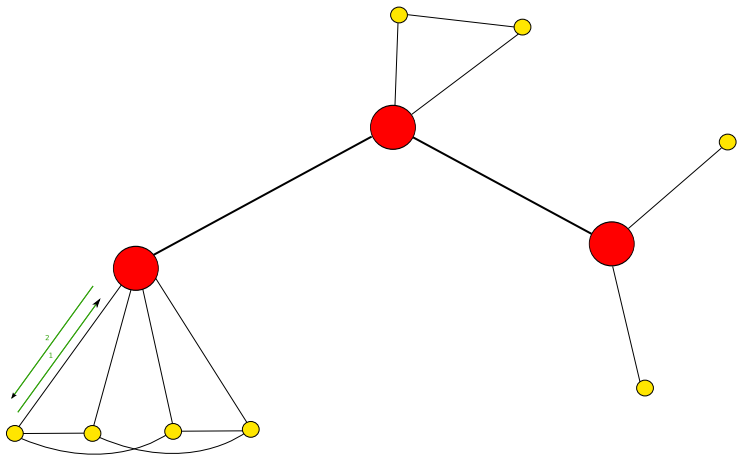
2. Measurement

Method 1: measure on server



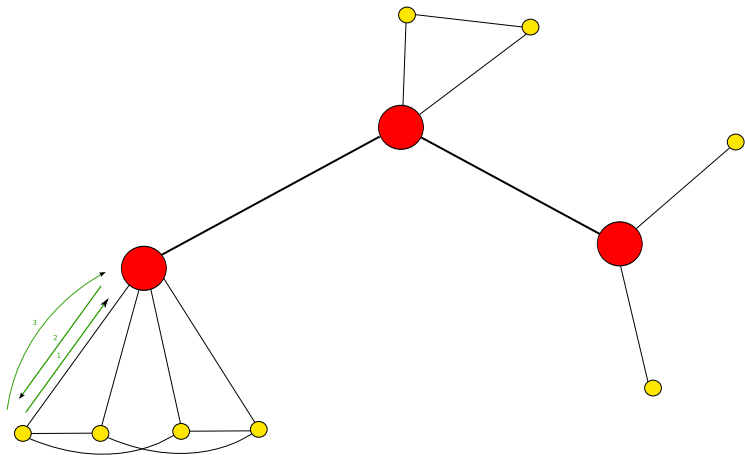
2. Measurement

Method 1: measure on server



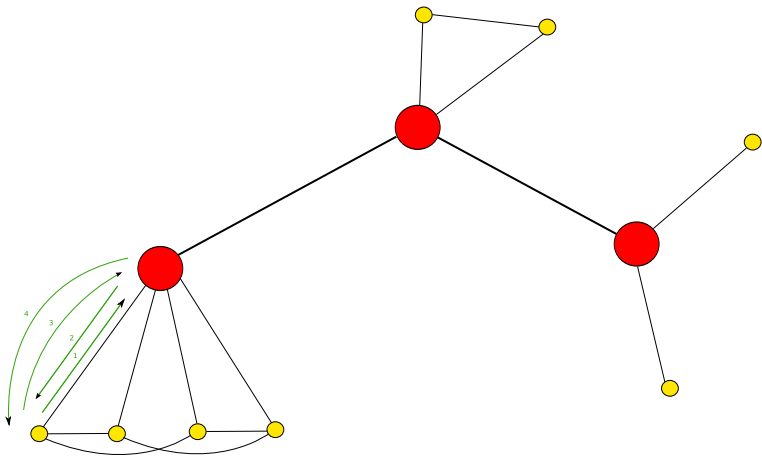
2. Measurement

Method 1: measure on server



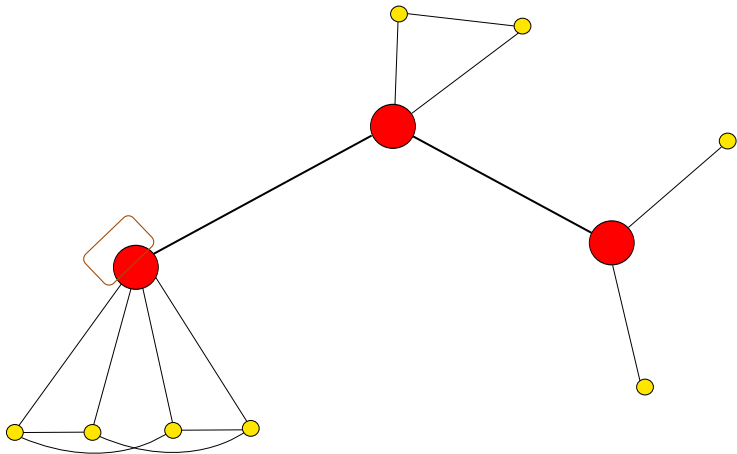
2. Measurement

Method 1: measure on server



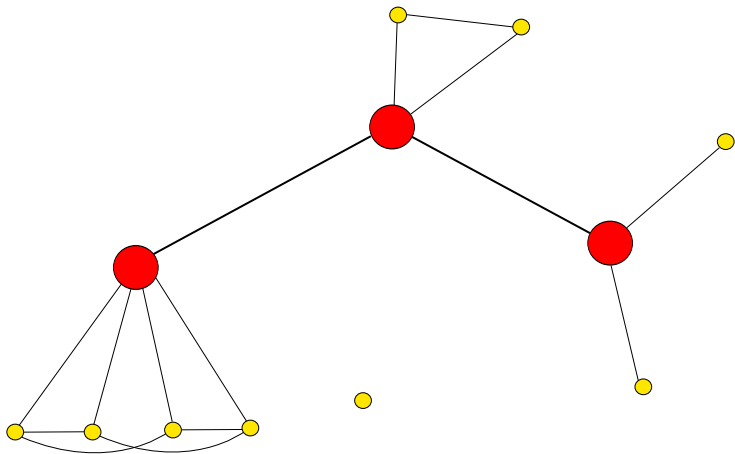
2. Measurement

Method 1: measure on server



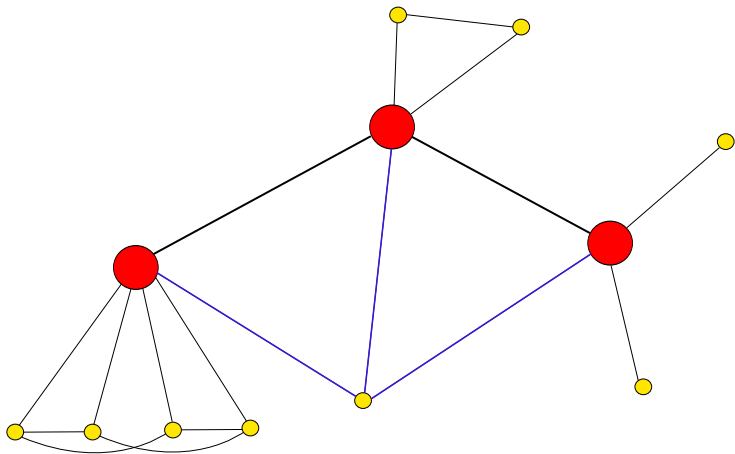
2. Measurement

Method 2: measure by clients sending queries



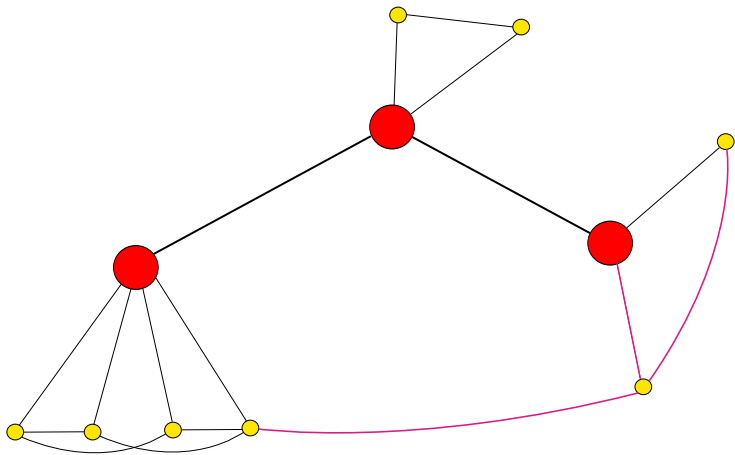
2. Measurement

Method 2: measure by clients sending queries

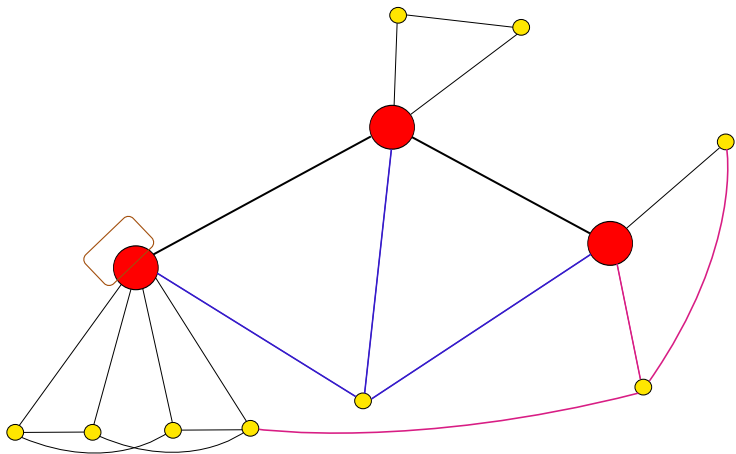


2. Measurement

Method 3: honey pot



2. Measurement



3. Data Processing

- ▶ Raw data:

$$T_1 \quad P \quad Q$$

$$T_2 \quad F_1 \quad F_2 \quad F_3 \longrightarrow P$$

$$T_3 \quad P \longrightarrow F_1 \quad F_3$$

$$T_4 \quad F_1 : P_1 \quad P_{15} \longrightarrow P$$

$$F_3 : P_1 \quad P_3 \quad P_7 \longrightarrow P$$

3. Data Processing

- ▶ Raw data:

T_1 P Q

3. Data Processing

- ▶ Raw data:

$$\begin{array}{l} T_1 \quad P \quad Q \\ T_2 \quad F_1 \quad F_2 \quad F_3 \longrightarrow P \end{array}$$

3. Data Processing

- ▶ Raw data:

$$\begin{array}{l} T_1 \quad P \quad Q \\ T_2 \quad F_1 \quad F_2 \quad F_3 \longrightarrow P \\ T_3 \quad P \longrightarrow F_1 \quad F_3 \end{array}$$

3. Data Processing

- ▶ Raw data:

$$T_1 \quad P \quad Q$$

$$T_2 \quad F_1 \quad F_2 \quad F_3 \longrightarrow P$$

$$T_3 \quad P \longrightarrow F_1 \quad F_3$$

$$T_4 \quad F_1 : P_1 \quad P_{15} \longrightarrow P$$

3. Data Processing

- ▶ Raw data:

$$T_1 \quad P \quad Q$$

$$T_2 \quad F_1 \quad F_2 \quad F_3 \longrightarrow P$$

$$T_3 \quad P \longrightarrow F_1 \quad F_3$$

$$T_4 \quad F_1 : P_1 \quad P_{15} \longrightarrow P$$

$$F_3 : P_1 \quad P_3 \quad P_7 \longrightarrow P$$

3. Data Processing

- ▶ Web interface

3. Data Processing

▶ Web interface

For each file id:

- file names
- number of providers
- first provider (date and peer id)
- main provider (list of peers id)
- content rating
- fake

3. Data Processing

▶ Web interface

For each peer:

- date of arrival
- number of files shared
- sent requests
- provided files

3. Data Processing

- ▶ Anonymisation

3. Data Processing

- ▶ Anonymisation

Strings $\xrightarrow{\text{normalisation}}$ Normalised strings $\xrightarrow{\text{cutting}}$ Words $\xrightarrow{\text{computation}}$

Word frequency:

≤ 100 anonymisation

≥ 100 stored data

3. Data Processing

- ▶ ≥ 100 stored data: examples

3. Data Processing

- ▶ ≥ 100 stored data: examples

Top “interesting” words (in file names and queries):

you, love, sex, xxx

3. Data Processing

- ▶ ≥ 100 stored data: examples

Top “interesting” words (in file names and queries):

you, love, sex, xxx

Top 15 words inquiries (not in file names):

girl, girls, boy, boys, child, children, playboy, pedo, attack, fille, enfants, boyz, enfant, incest, preteen

4. Future

- ▶ More measurements (other methods)
- ▶ Content Rating and Fake Detection System
- ▶ Keywords analysis
- ▶ Better knowledge of paedophile activities