

Protecting privacy in social media via controlling exposure

Saptarshi Ghosh
and Mainack Mondal

CS 60017
Autumn 2019



Last time we talked about...

- Anonymity: one aspect of privacy
 - What happens if you give users anonymity?
 - They become more disinhibited
 - Case study on Whisper / 4chan /b/
 - Important for internet culture

Now a few systems on how to better control privacy in social media

- Recall Exposure control
 - Extension of access control
 - Users care about the exposure of the content, i.e., the set of users who will eventually see a post
 - Users have an expectation about exposure
 - Possible privacy violation = If the actual exposure did not meet expected exposure

Now a few systems on how to better control privacy in social media

- Recall Exposure control
 - Extension of access control
 - Users care about the exposure of the content, i.e., the set of users who will eventually see a post
 - Users have an expectation about exposure
 - Possible privacy violation = If the actual exposure did not meet expected exposure
- Two systems to control exposure
 - Strengthening access control in social media
 - Limiting data aggregators in social media

Strengthening access control in social media

- Key idea
 - We can control exposure if the users can **easily** and **accurately** specify the access control lists
 - Then the users who has access to these lists are always in user's expected exposure set

Understanding and Specifying Social Access control lists

(based on Mondal et al.'s SOUPS'14 paper)

Privacy in Online Social Networks (OSNs)

Prior to OSNs (**also called Online Social media sites or OSMS**)

Users were largely **content consumers**

In OSNs

Users expected to be **content creators and managers**

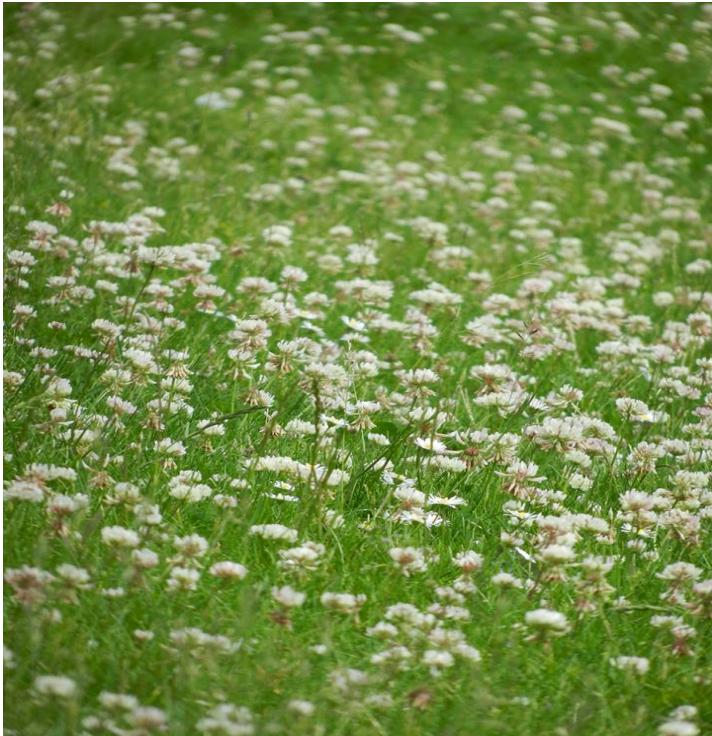
In sites like Facebook

Must enumerate who is able to access every uploaded content

Avg. 130 friends, 90 pieces of content every month per user

Really difficult to get the privacy right for “cognitive burden”

Privacy sensitive content in OSMs



Non privacy sensitive content:
all friends should be able to
access

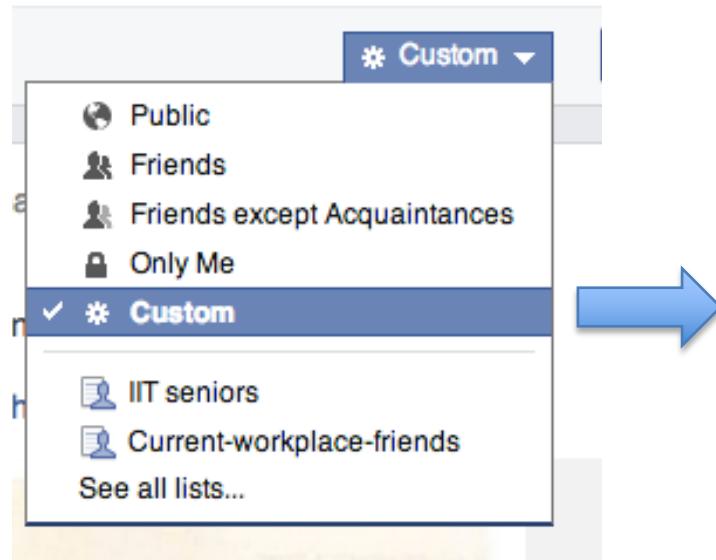


Privacy sensitive Content:
only select friends should
be able to access

How do users manage access currently?

Users specify **Social Access Control Lists** (SACLs)

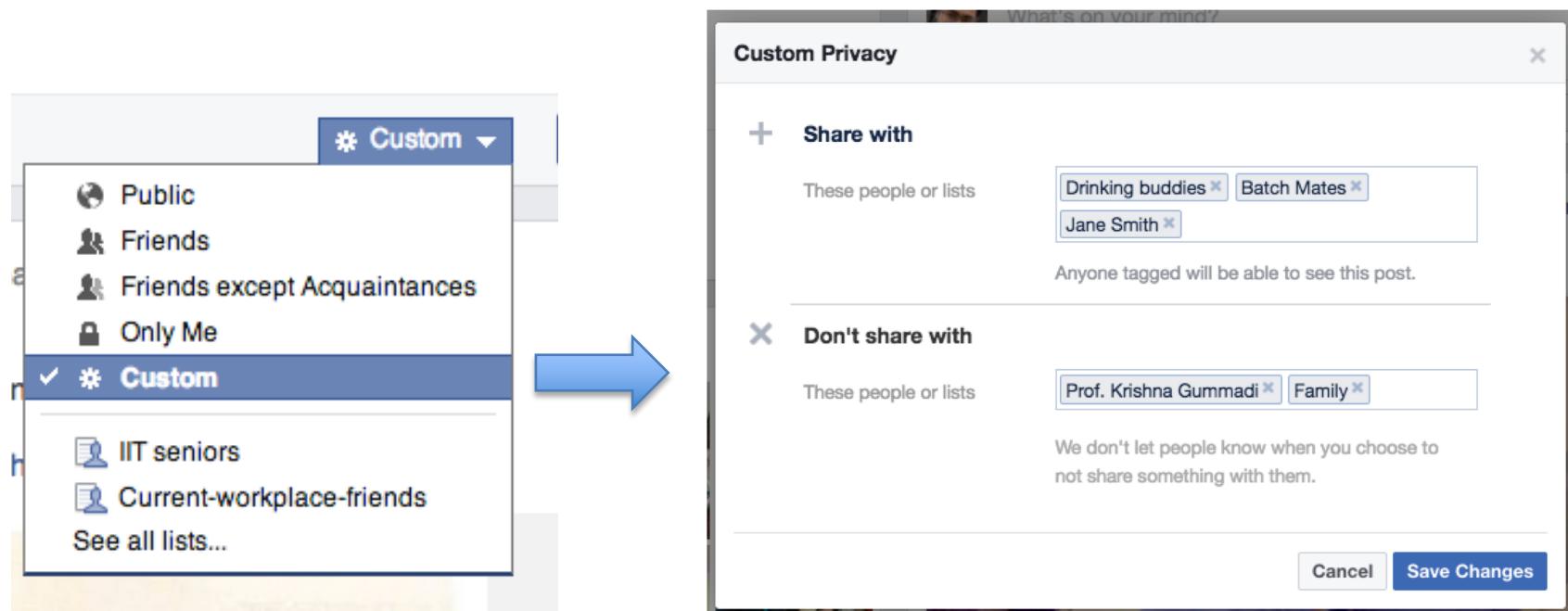
SACLs: Share with a subset of friends



How do users manage access currently?

Users specify **Social Access Control Lists** (SACLs)

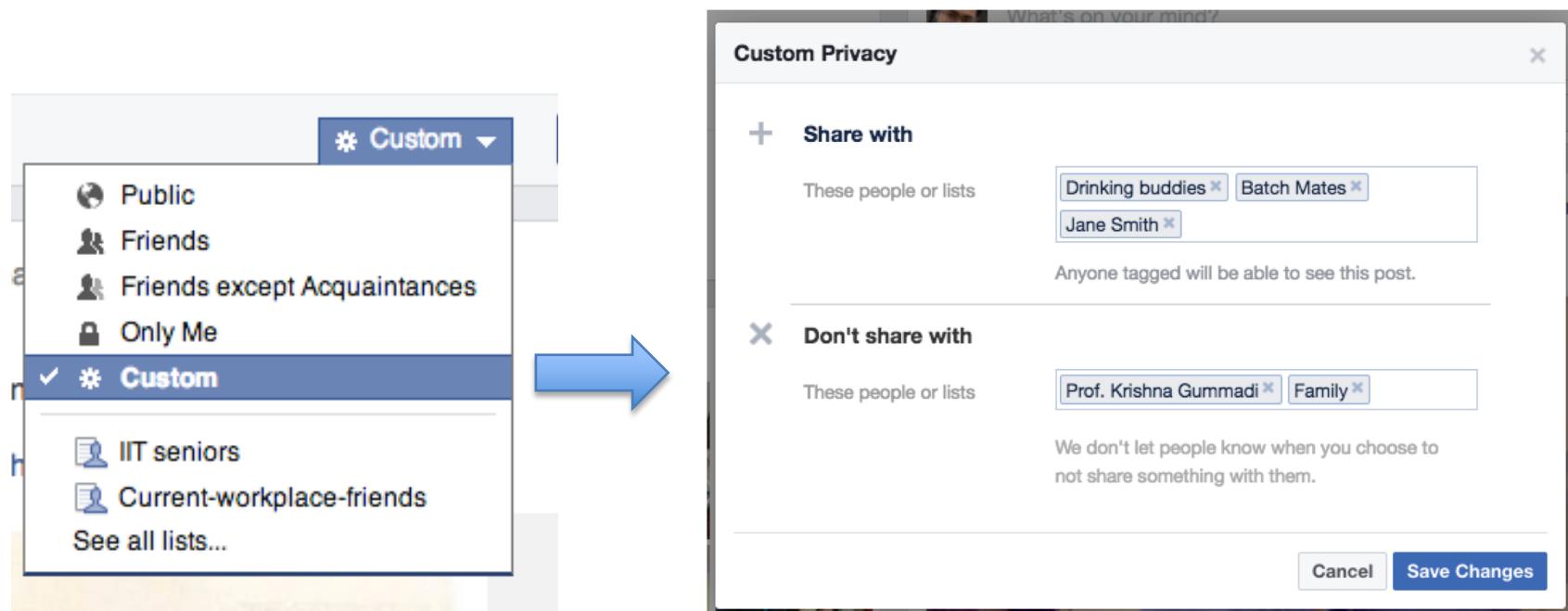
SACLs: Share with a subset of friends



How do users manage access currently?

Users specify **Social Access Control Lists** (SACLs)

SACLs: Share with a subset of friends



Each time users have to specify these SACLs manually!

State of the art for helping SACL specification

Provide users **automatically detected groups**

Network community based group detection

User profile attribute based group detection

User activity based group detection

Assumption by existing work:

Automatically detected groups are similar to SACLs

They evaluated their proposals based on small scale user interviews

State of the art for helping SACL specification

Provide users **automatically detected groups**

[WWW 2010]

Network community based group detection

[DBSocial 2012]

User profile attribute based group detection

[SOUPS 2012]

User activity based group detection

Assumption by existing work:

Automatically detected groups are similar to SACLs

They evaluated their proposals based on small scale user interviews

State of the art for helping SACL specification

Provide users **automatically detected groups**

Network community based group detection

User profile attribute based group detection

User activity based group detection

Assumption by existing work:

Automatically detected groups are similar to SACLs

They evaluated their proposals based on small scale user interviews

No validation before us using large scale real world SACLs in-use

State of the art for helping SACL specification

Provide users **automatically detected groups**

Network community based group detection

User profile attribute based group detection

User activity based group detection

Our goal is to better understand **real world SACL usage & simplify SACL specification**

They evaluated their proposals based on small scale user interviews

No validation before us using large scale real world SACLs in-use

Friendlist Manager: app for data collection

We **built and deployed** Friendlist Manager (FLM) Facebook app
Divides user's friends into groups by **network community detection**

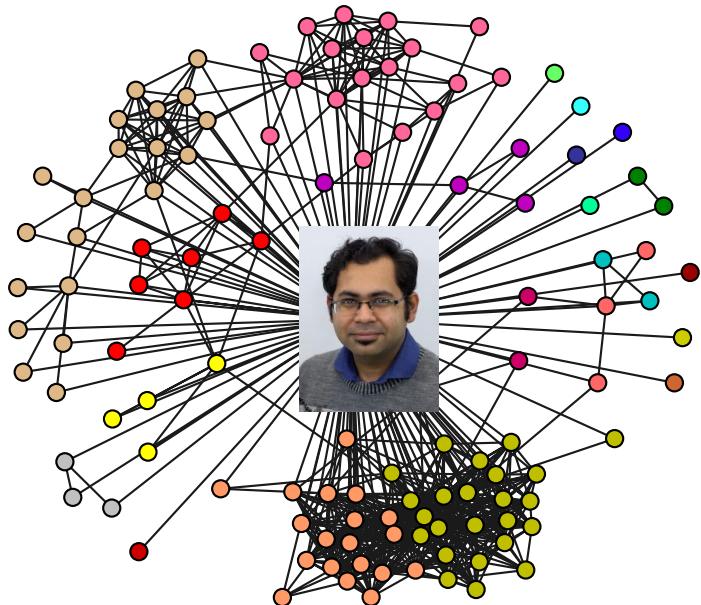
Friendlist Manager: app for data collection

We **built and deployed** Friendlist Manager (FLM) Facebook app
Divides user's friends into groups by **network community detection**



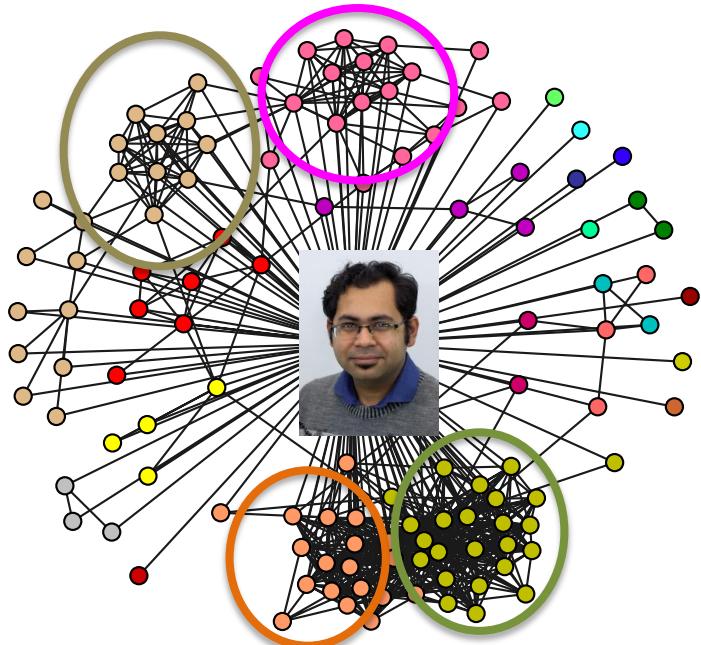
Friendlist Manager: app for data collection

We **built and deployed** Friendlist Manager (FLM) Facebook app
Divides user's friends into groups by **network community detection**



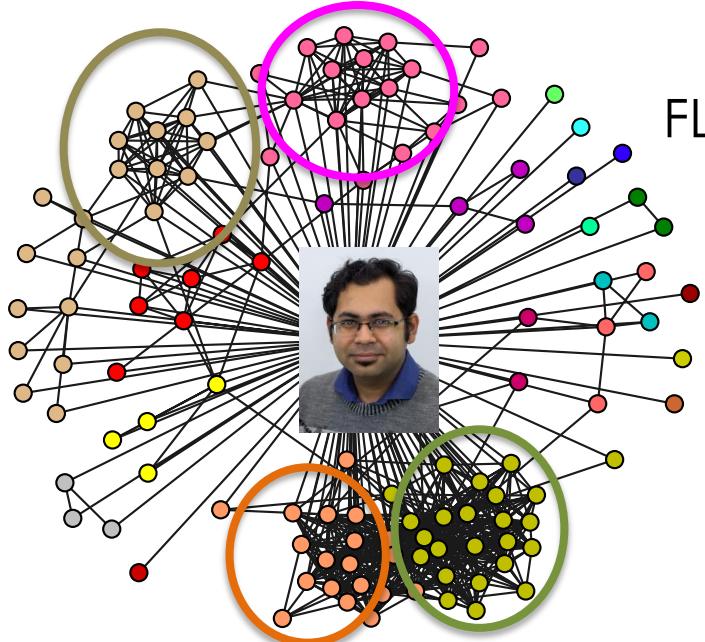
Friendlist Manager: app for data collection

We **built and deployed** Friendlist Manager (FLM) Facebook app
Divides user's friends into groups by **network community detection**



Friendlist Manager: app for data collection

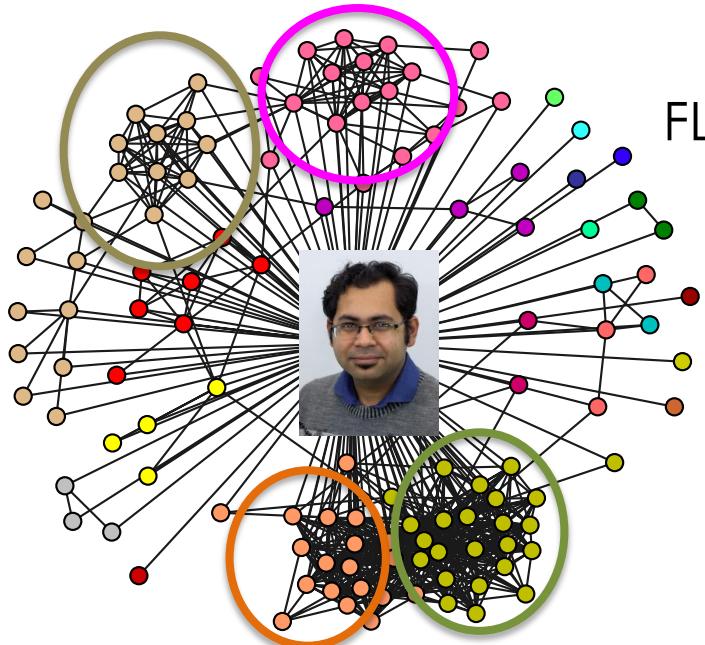
We **built and deployed** Friendlist Manager (FLM) Facebook app
Divides user's friends into groups by **network community detection**



FLM used unsupervised learning
Collection and community detection of 1-hop social subgraph in **real time**

Friendlist Manager: app for data collection

We **built and deployed** Friendlist Manager (FLM) Facebook app
Divides user's friends into groups by **network community detection**

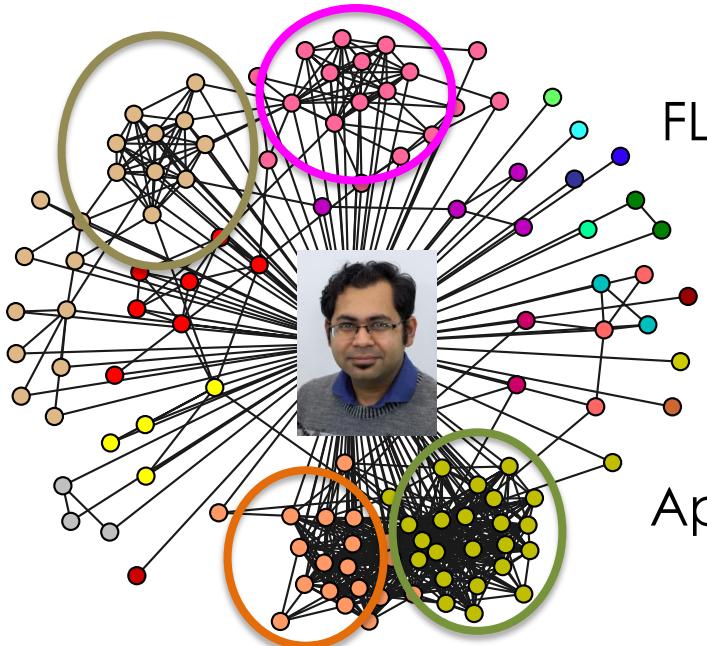


FLM used unsupervised learning
Collection and community detection of 1-hop social subgraph in **real time**

Users can create these groups as **friendlists** on Facebook

Friendlist Manager: app for data collection

We **built and deployed** Friendlist Manager (FLM) Facebook app
Divides user's friends into groups by **network community detection**



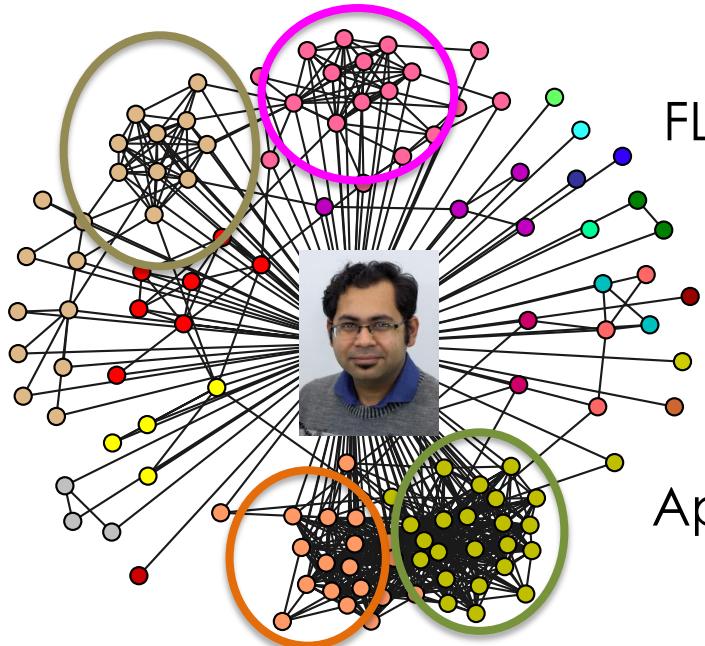
FLM used unsupervised learning
Collection and community detection of 1-hop social subgraph in **real time**

Application functionality
<https://friendlist-manager.mpi-sws.org/>

Users can create these groups as **friendlists** on Facebook

Friendlist Manager: app for data collection

We **built and deployed** Friendlist Manager (FLM) Facebook app
Divides user's friends into groups by **network community detection**



FLM used unsupervised learning
Collection and community detection of 1-hop social subgraph in **real time**

Application functionality
<https://friendlist-manager.mpi-sws.org/>

Users can create these groups as **friendlists** on Facebook

1,200+ users have installed FLM in 2013 - 2014!

How to collect data from FLM users?

We asked consent from users to access their data

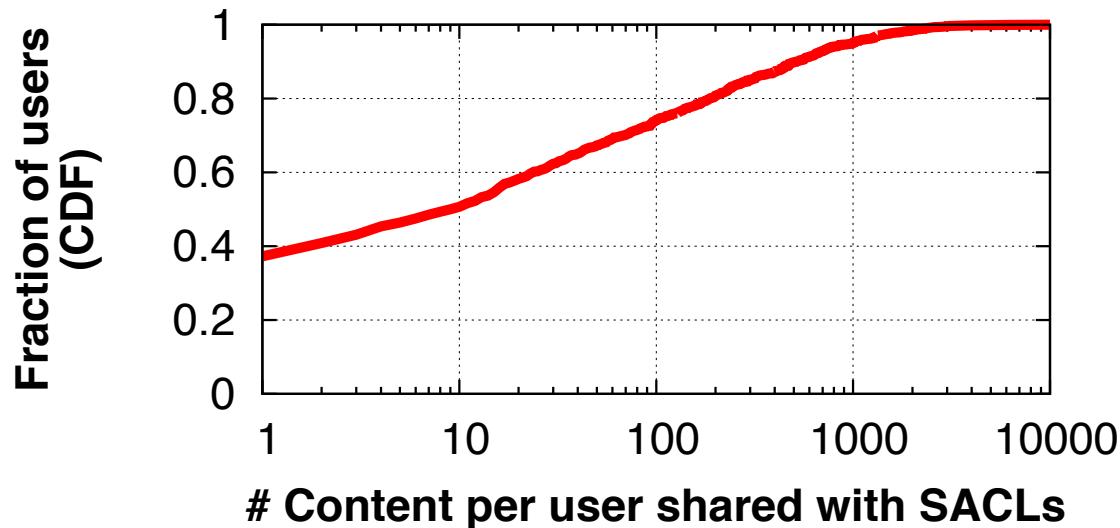
1,100+ users gave consent

Collected a snapshot of all their profile and SACL data

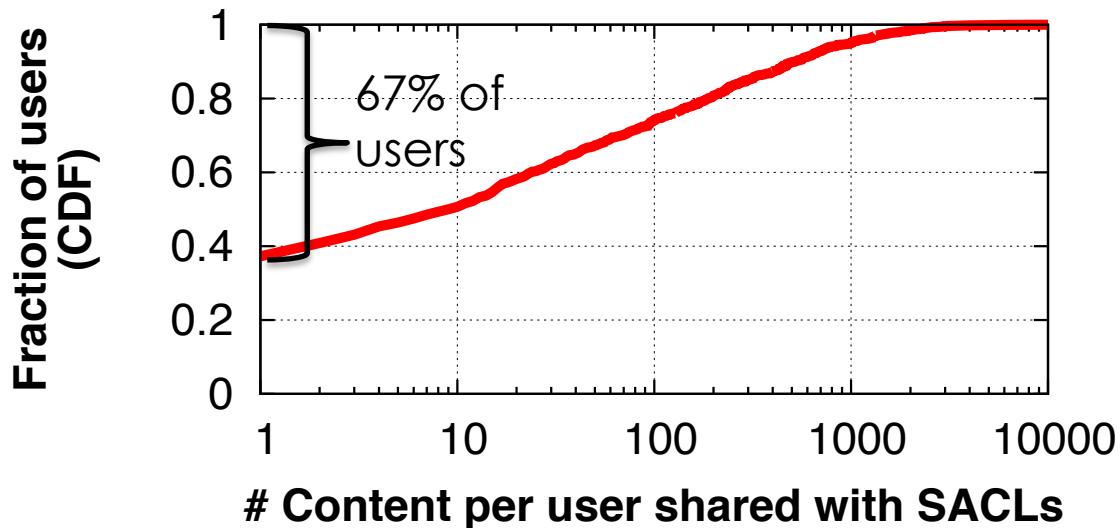
First large scale dataset of real in-use SACLs



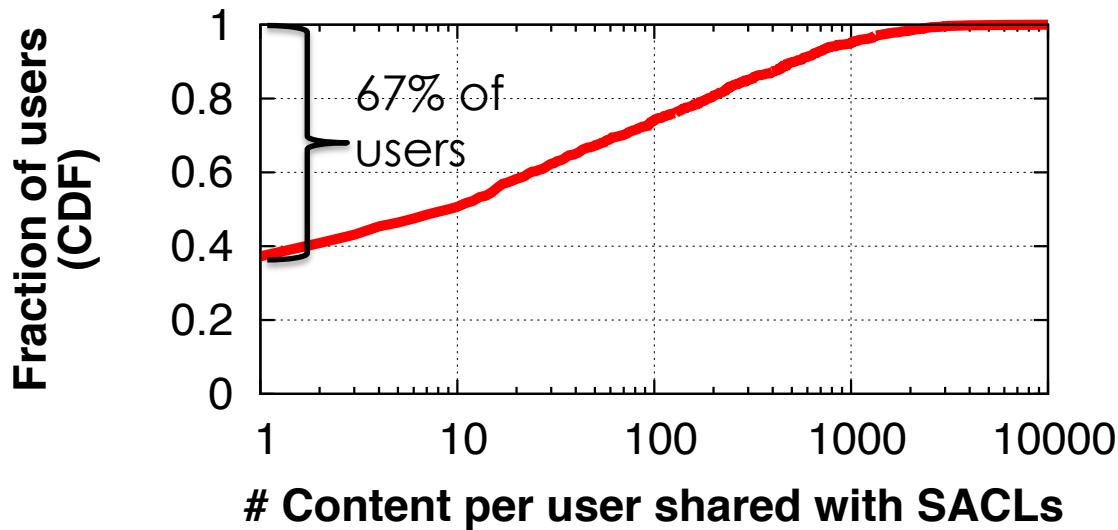
Do users leverage SACLs to share content?



Do users leverage SACLs to share content?

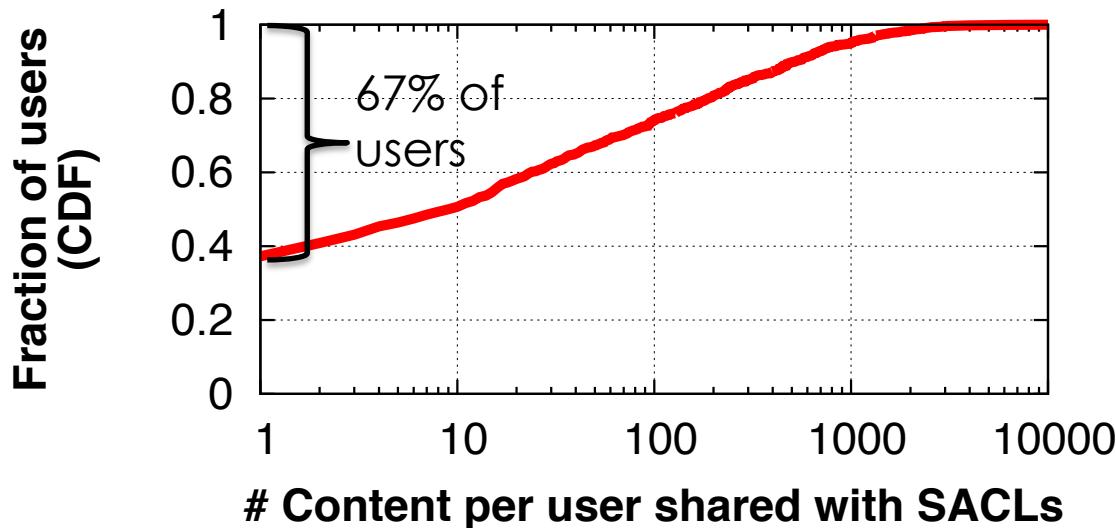


Do users leverage SACLs to share content?



Total 200K content is shared with 7.6k unique SACLs!

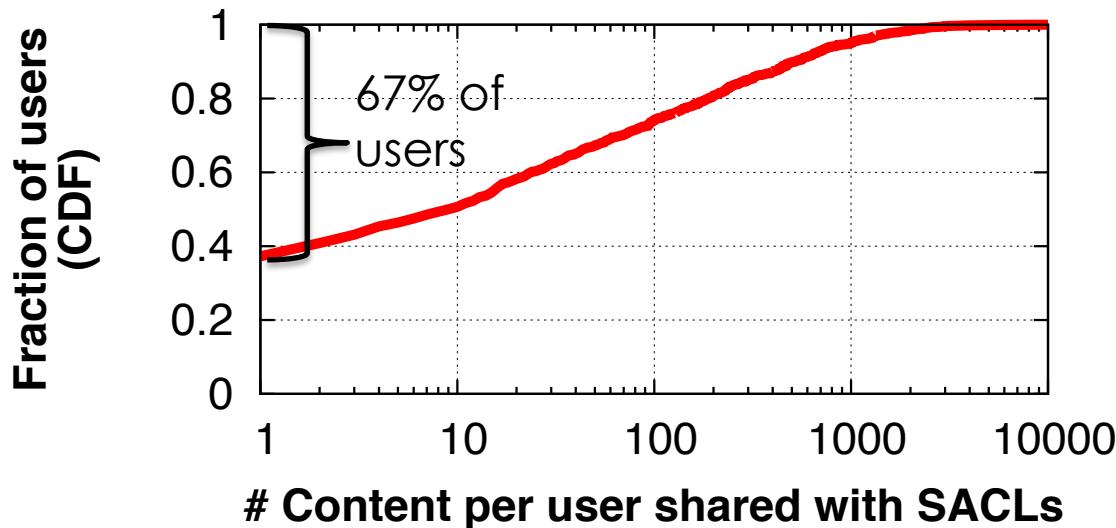
Do users leverage SACLs to share content?



Total 200K content is shared with 7.6k unique SACLs!

Majority of users used SACLs for at least one of their content

Do users leverage SACLs to share content?

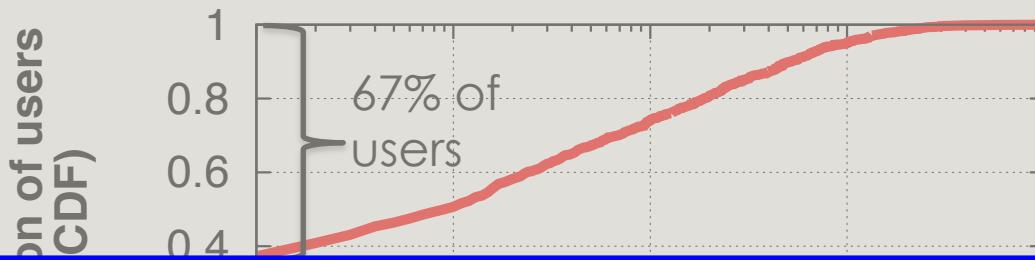


Total 200K content is shared with 7.6k unique SACLs!

Majority of users used SACLs for at least one of their content

It is important to look for ways to simplify SACL specification

Do users leverage SACLs to share content?



How can we design systems for reducing user overhead of specifying SACLs?

Total 200K content is shared with 7.6k unique SACLs!

Majority of users used SACLs for at least one of their content

It is important to look for ways to simplify SACL specification

How to quantify user overhead of SACL specification?

Custom Privacy

✓ Share this with _____

These people or lists

Note: Anyone tagged can also see this post.

✗ Don't share this with _____

These people or lists

Facebook never reveals when you choose not to share a post with somebody.

How to quantify user overhead of SACL specification?



Total terms: 5

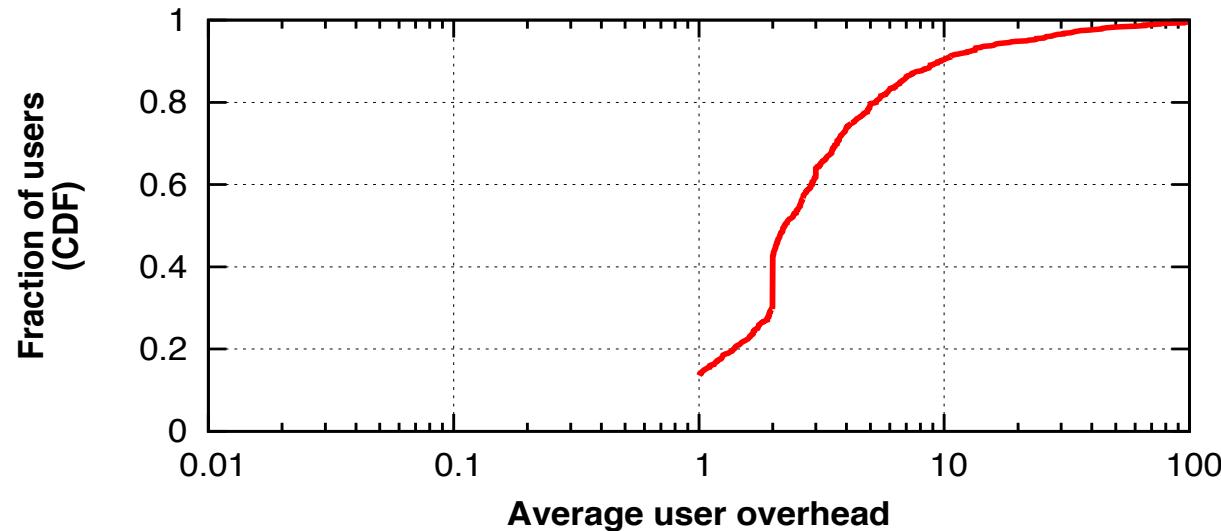
How to quantify user overhead of SACL specification?



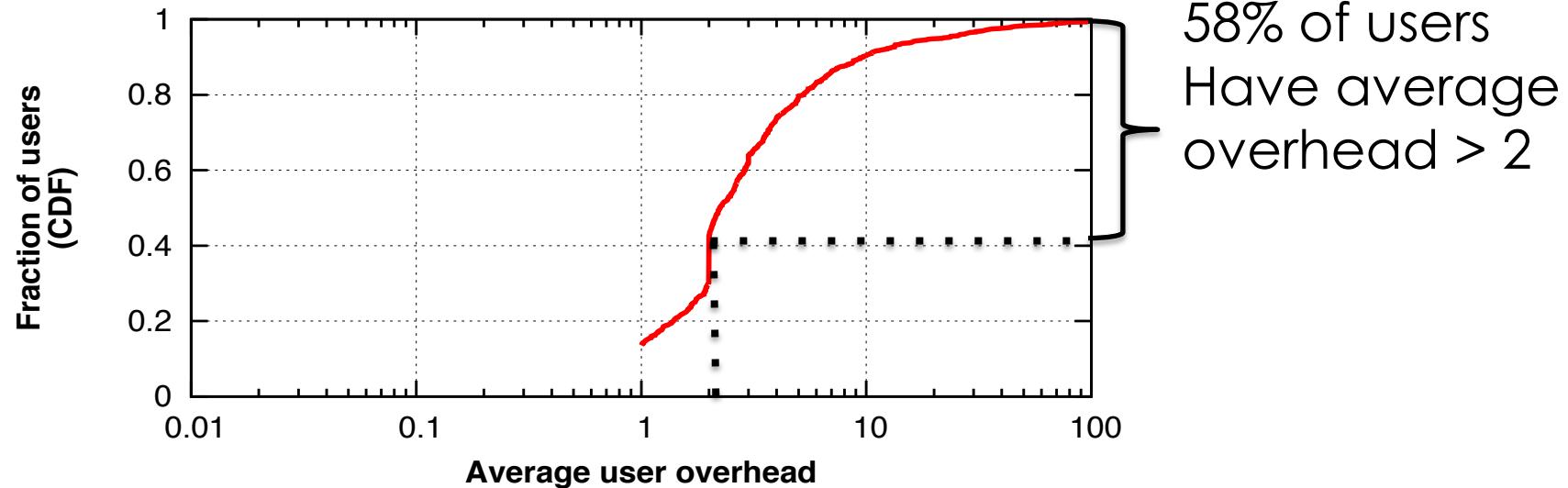
Total terms: 5

Average user overhead = average #terms per content

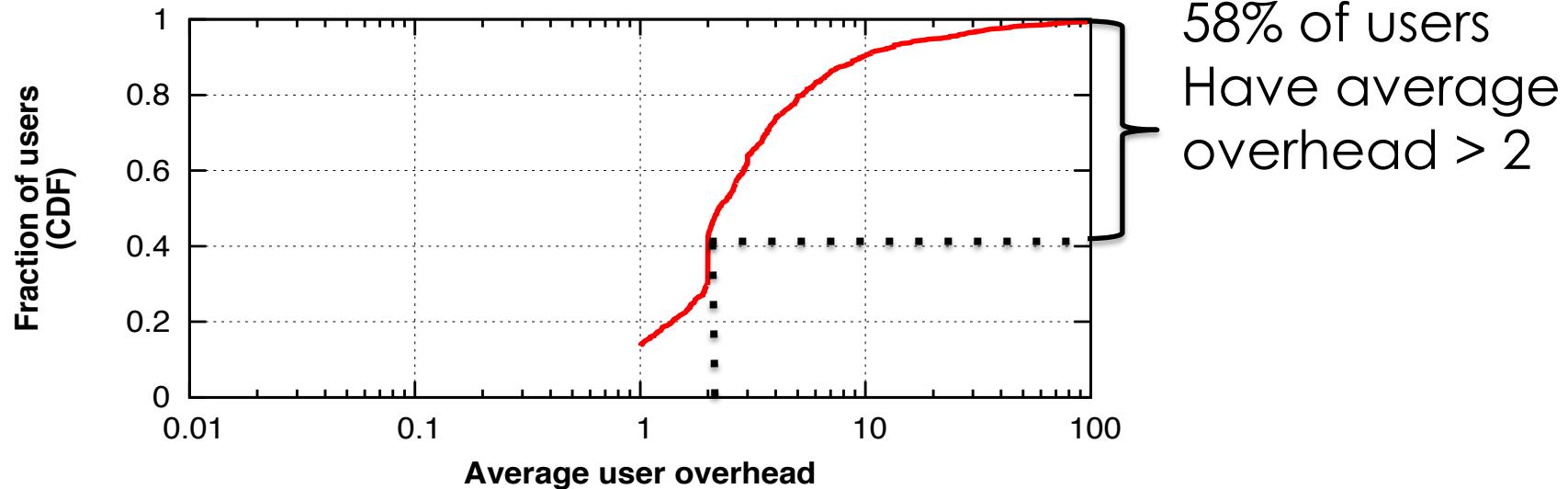
What is the current overhead for users?



What is the current overhead for users?

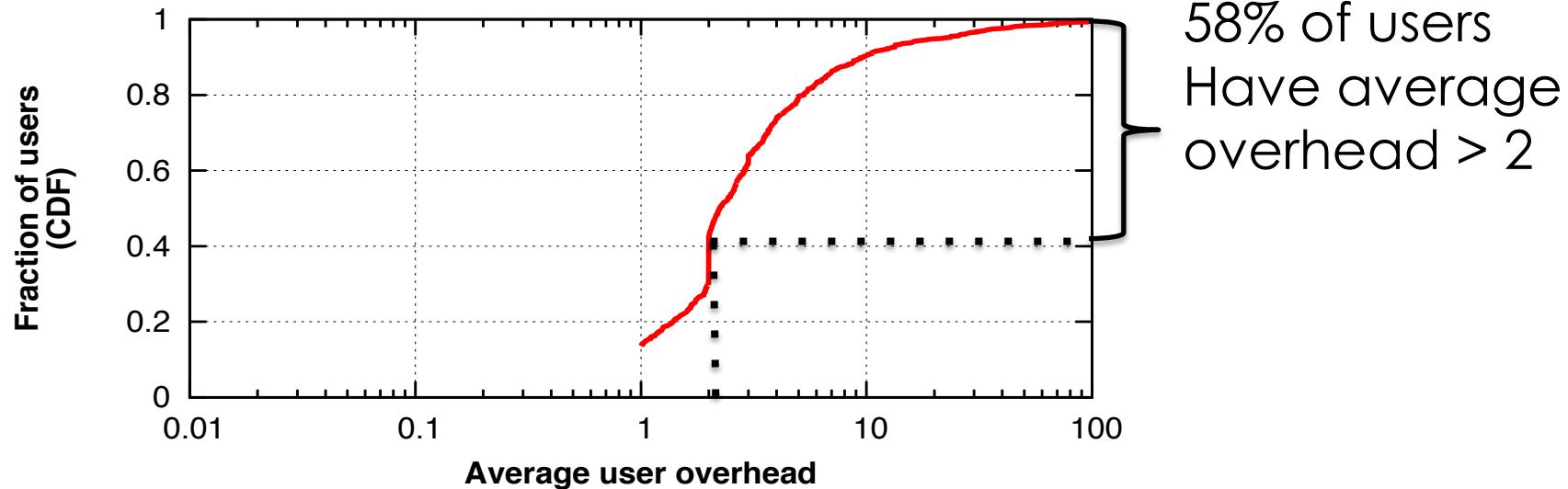


What is the current overhead for users?



More than 150 users have overhead more than 5

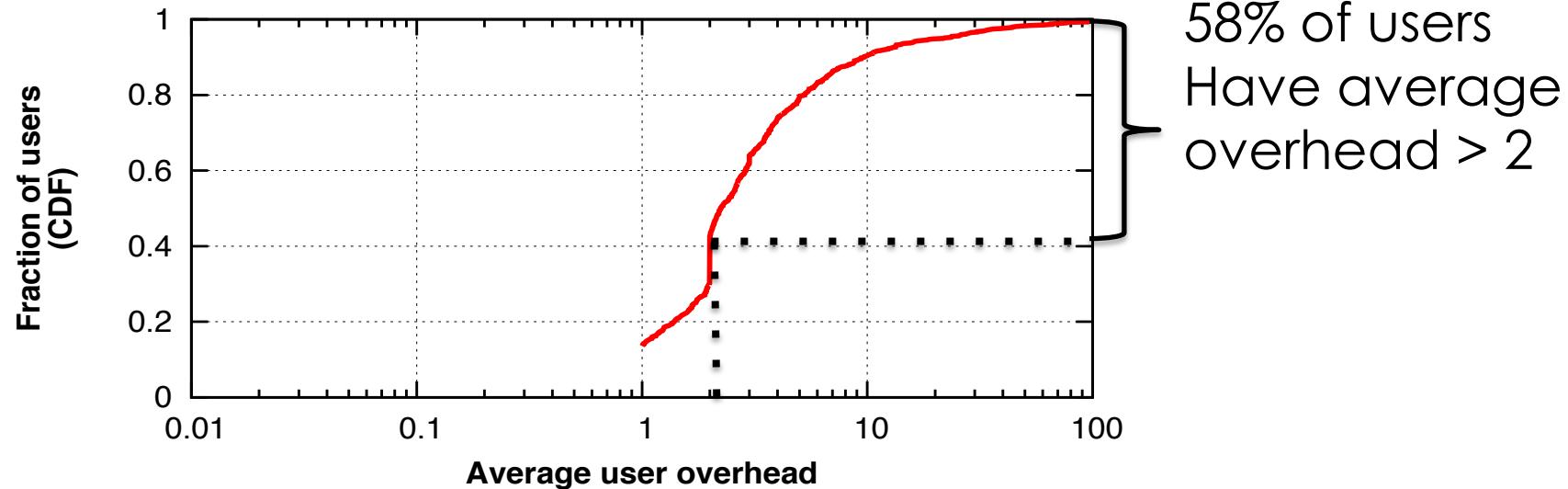
What is the current overhead for users?



More than 150 users have overhead more than 5

User overhead is high for specifying SACLs!

What is the current overhead for users?



More than 150 users have overhead more than 5

User overhead is high for specifying SACLs!

Can we use automated groups to reduce this overhead?

Can we reduce overhead using automatically detected groups?

In the existing work there are three types of group detection for SACLs:

Based on network structure, user profile attributes, user activity

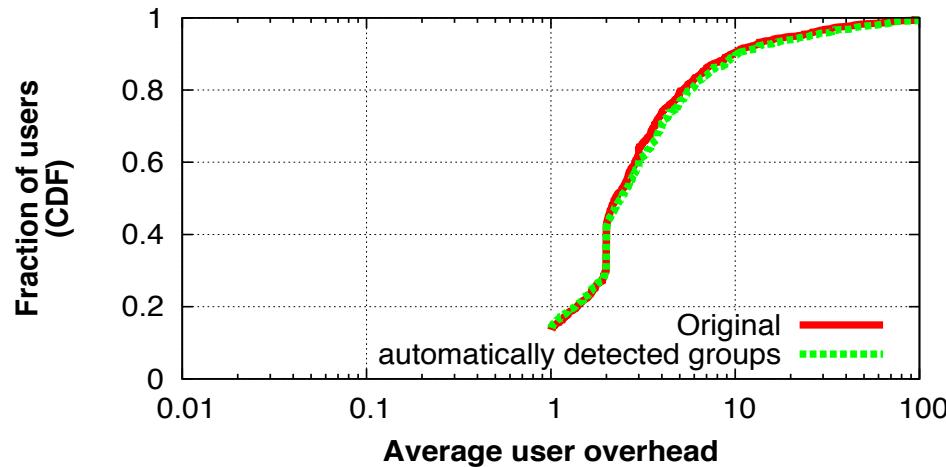
We tried all!

Can we reduce overhead using automatically detected groups?

In the existing work there are three types of group detection for SACLs:

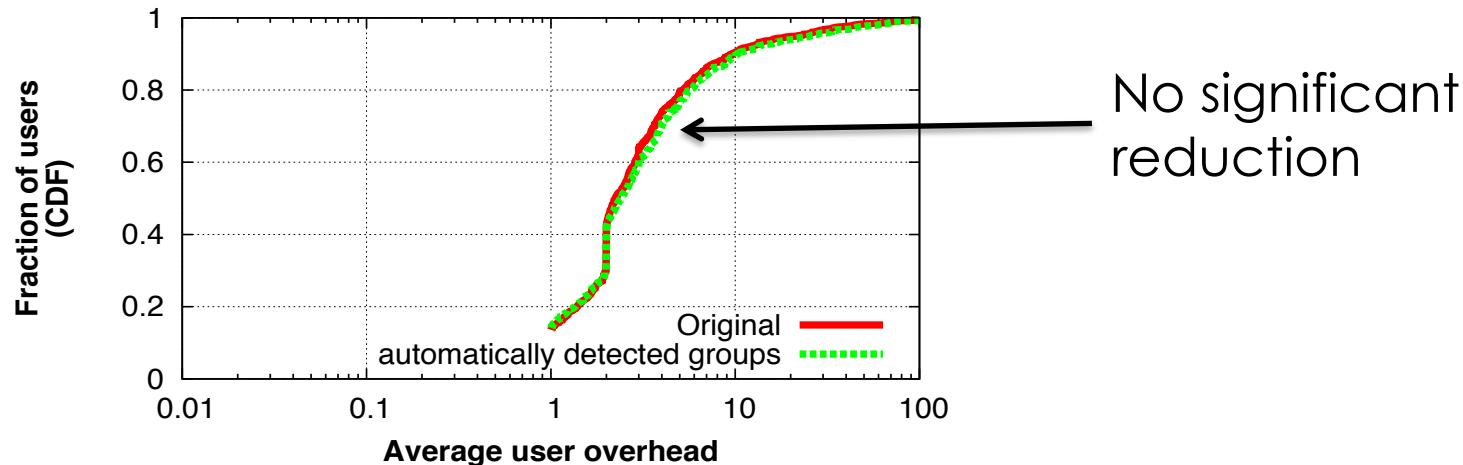
Based on **network structure**, **user profile attributes**, **user activity**

We tried all!



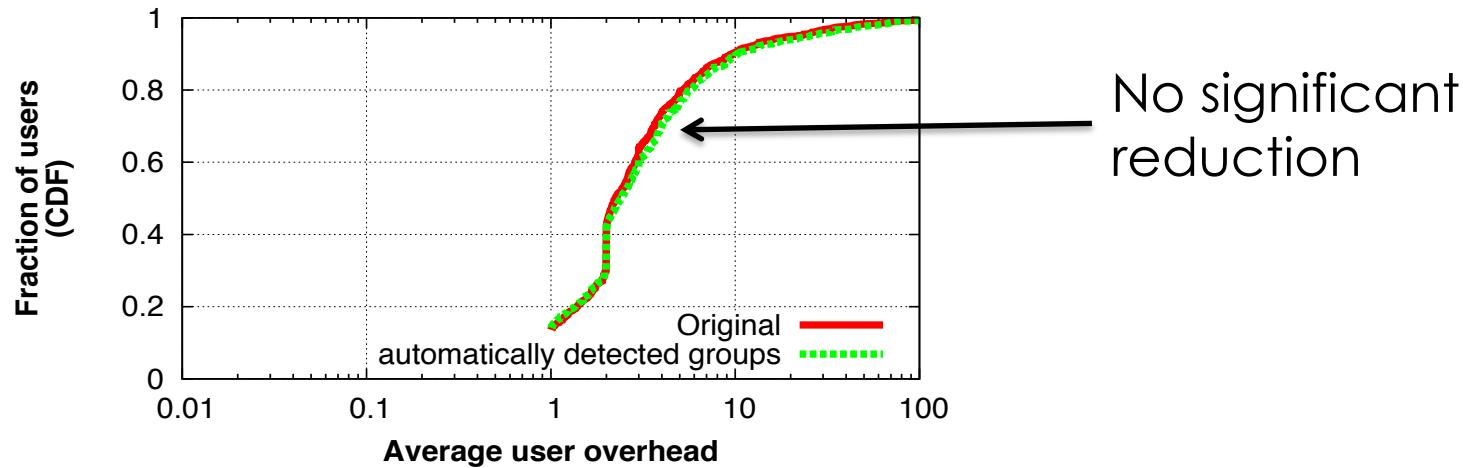
Can we reduce overhead using automatically detected groups?

In the existing work there are three types of group detection for SACLs:
Based on network structure, user profile attributes, user activity
We tried all!



Can we reduce overhead using automatically detected groups?

In the existing work there are three types of group detection for SACLs:
Based on network structure, user profile attributes, user activity
We tried all!

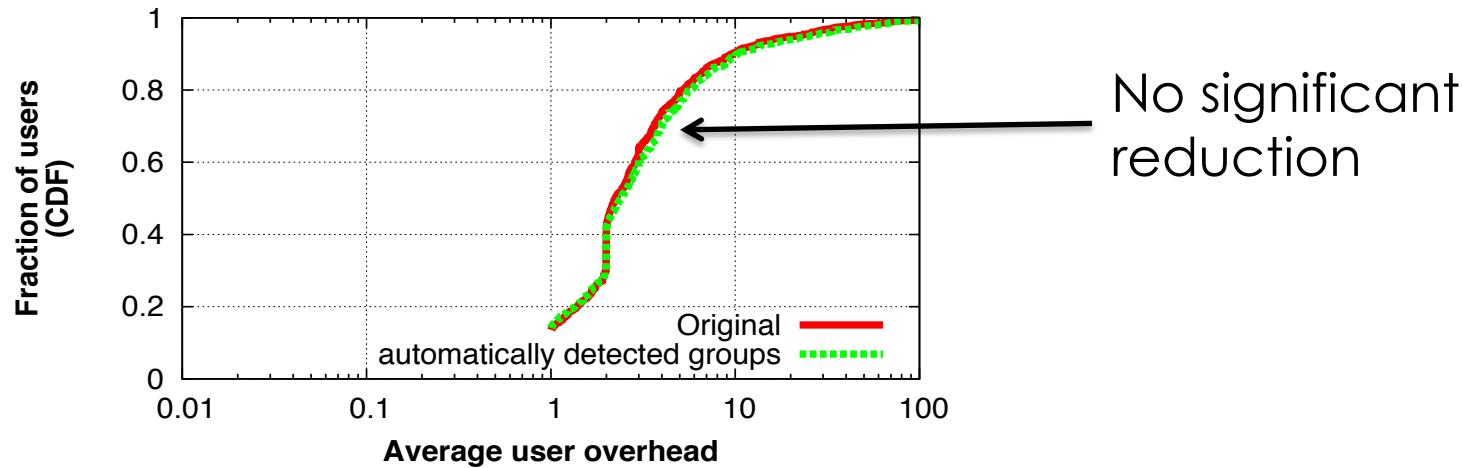


Automated groups **do not** significantly reduce overhead

Reason: Most SACLs are **not highly correlated** with automatically detected groups (used F-score, ARI, normalized entropy)

Can we reduce overhead using automatically detected groups?

In the existing work there are three types of group detection for SACLs:
Based on network structure, user profile attributes, user activity
We tried all!

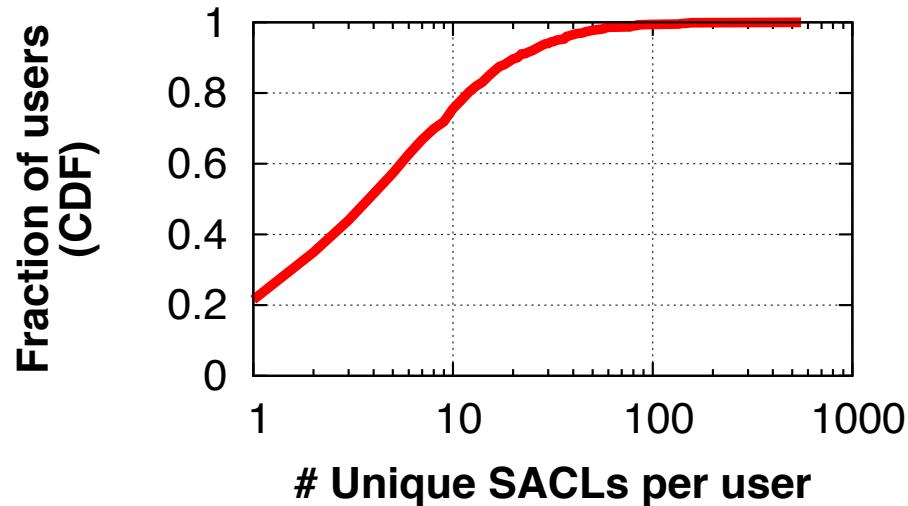


Automated groups **do not** significantly reduce overhead

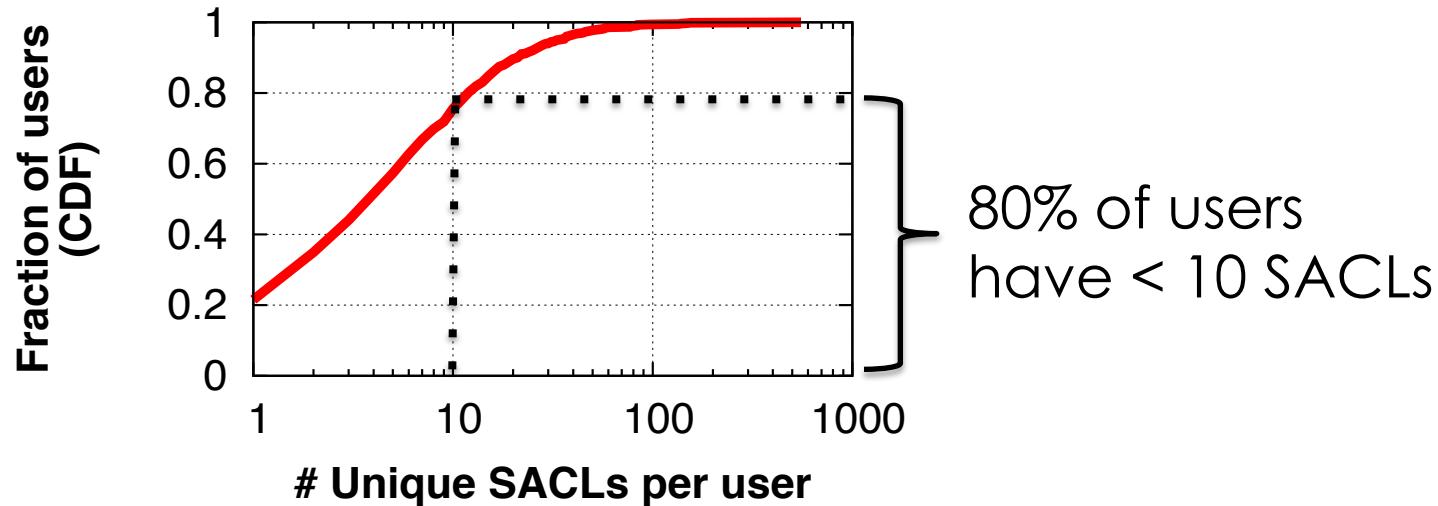
Reason: Most SACLs are **not highly correlated** with automatically detected groups (used F-score, ARI, normalized entropy)

Is there any other way?

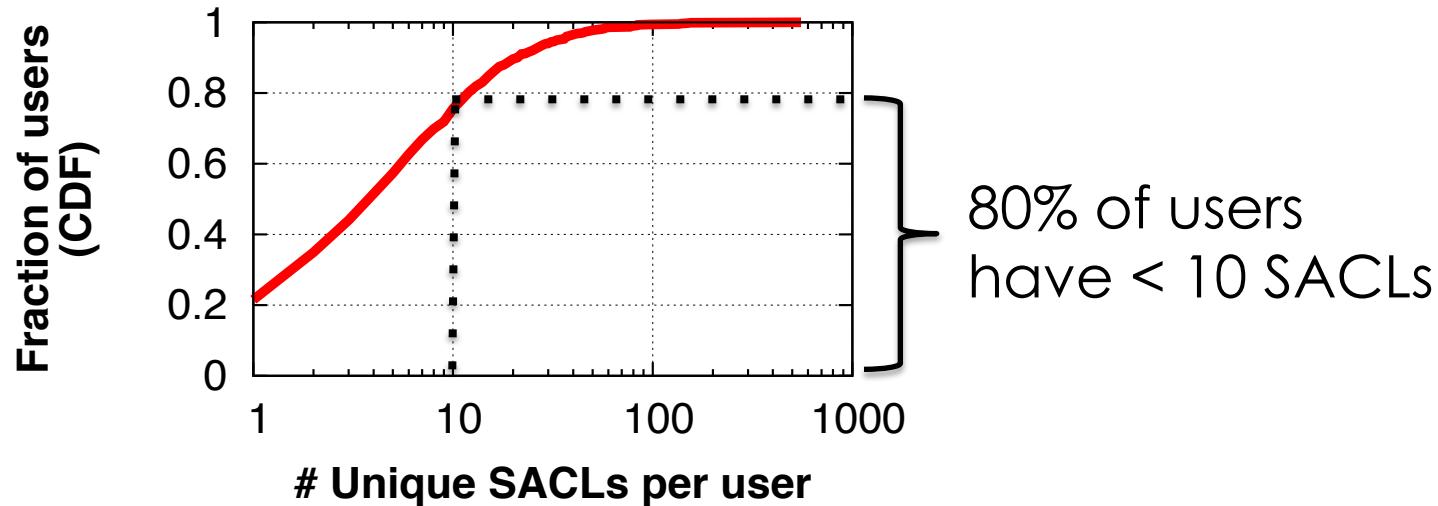
Insight via modelling user behavior: Users reuse SACLs repeatedly



Insight via modelling user behavior: Users reuse SACLs repeatedly

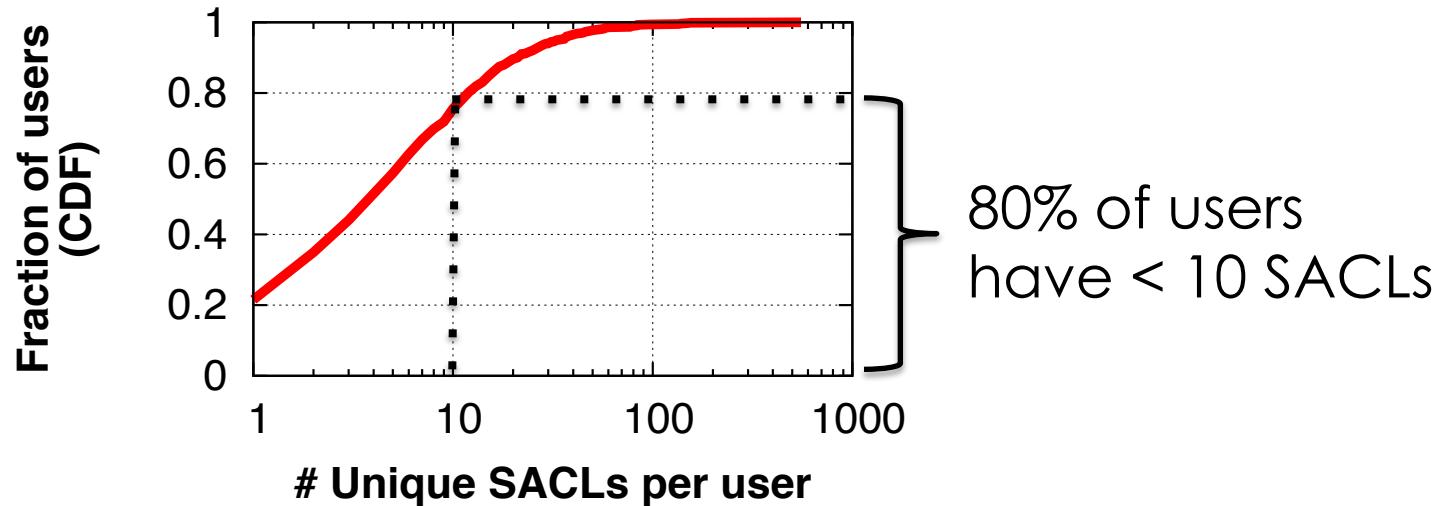


Insight via modelling user behavior: Users reuse SACLs repeatedly



Moreover on average a SACL is reused for **28 contents**

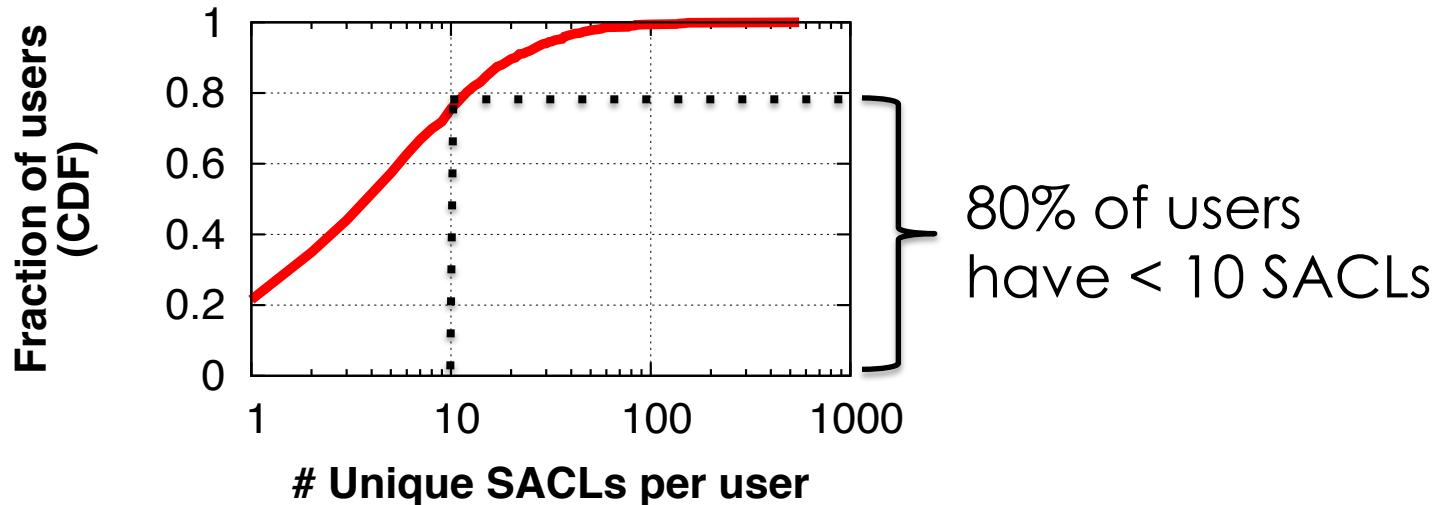
Insight via modelling user behavior: Users reuse SACLs repeatedly



Moreover on average a SACL is reused for **28 contents**

Most users **reuse** their **SACLs repeatedly**

Insight via modelling user behavior: Users reuse SACLs repeatedly

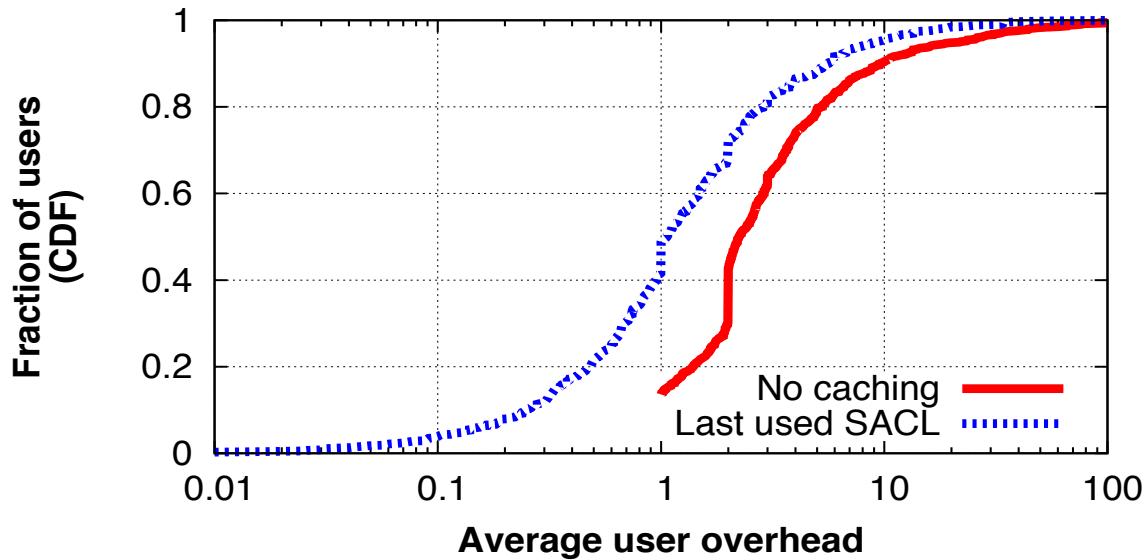


Moreover on average a SACL is reused for **28 contents**

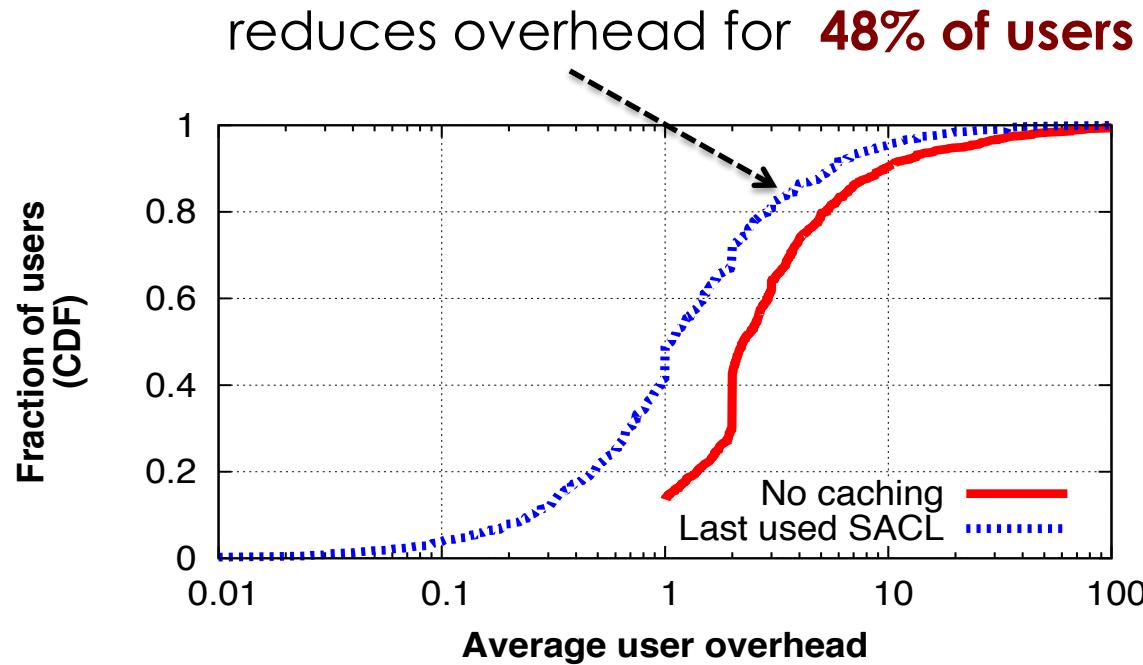
Most users **reuse** their **SACLs repeatedly**

Idea: How about **caching** a few of the past SACLs to reduce overhead?

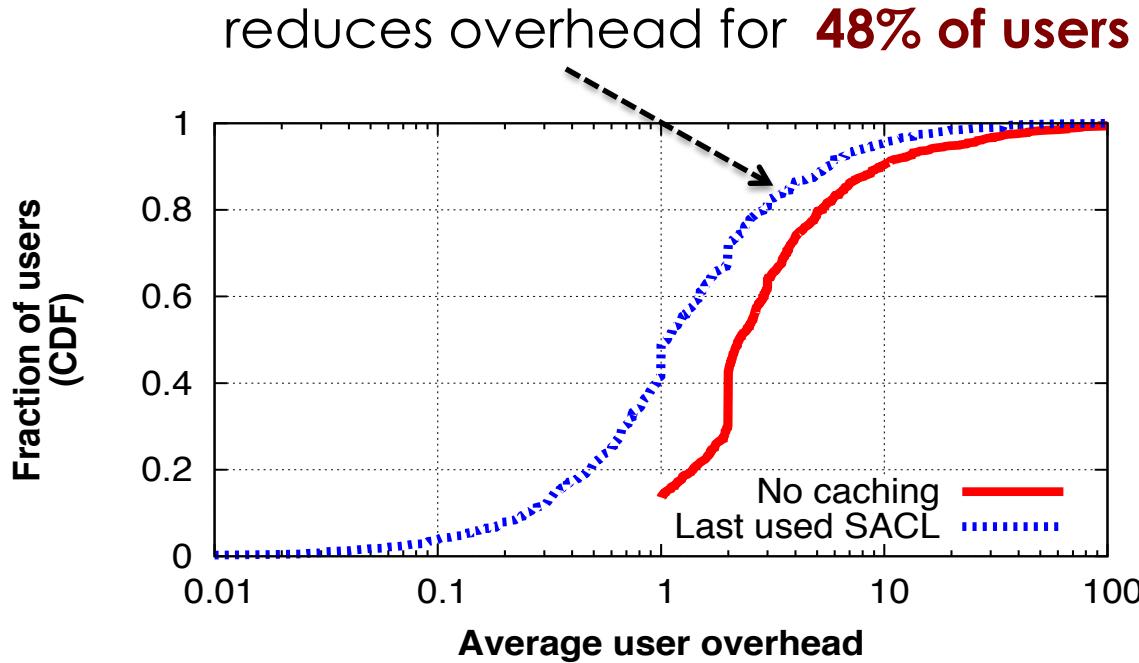
Does caching reduce SACL specification overhead?



Does caching reduce SACL specification overhead?

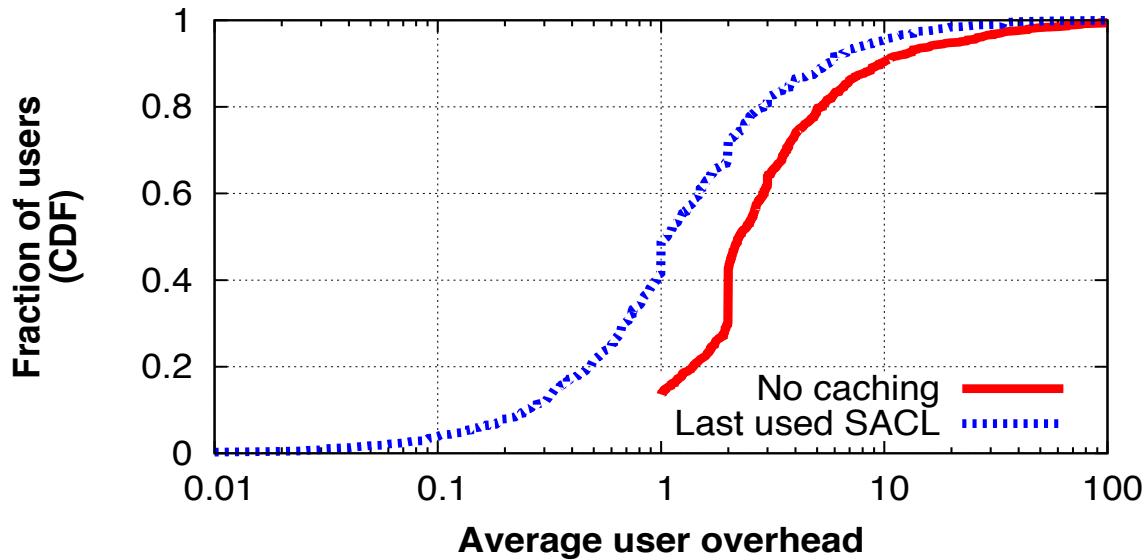


Does caching reduce SACL specification overhead?



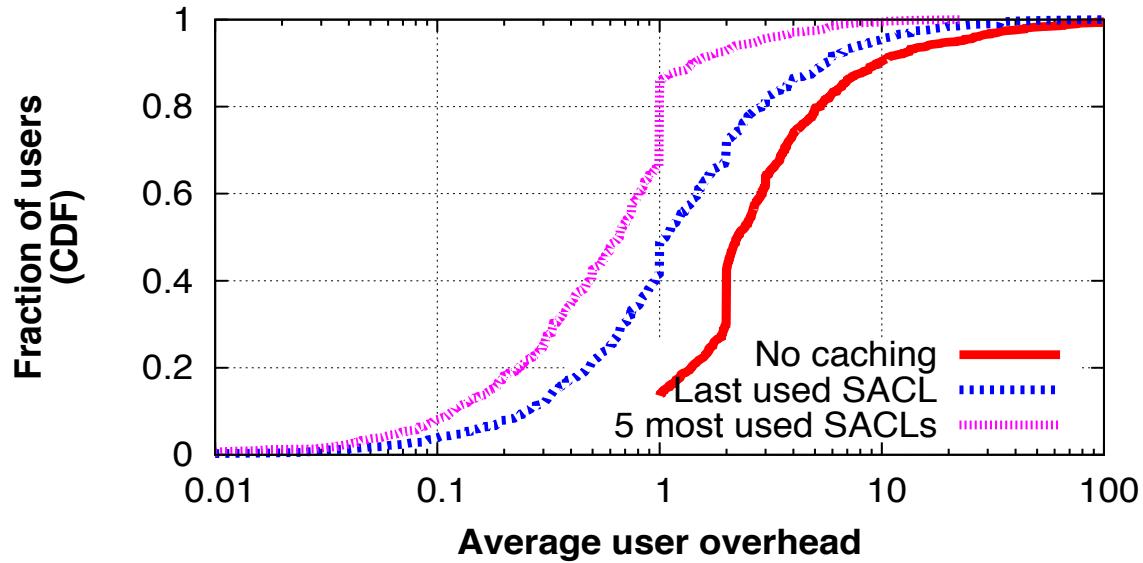
Can we improve by caching a few SACLs instead of one SACL?

Does caching reduce SACL specification overhead?



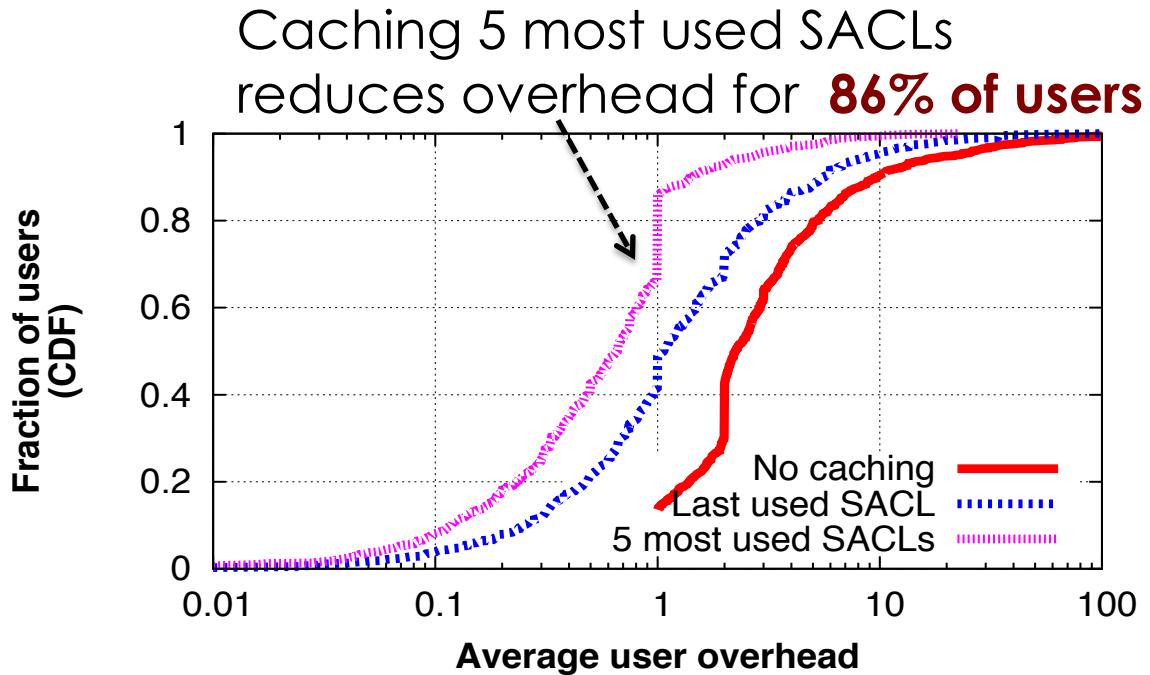
Can we improve by caching a few SACLs instead of one SACL?

Does caching reduce SACL specification overhead?



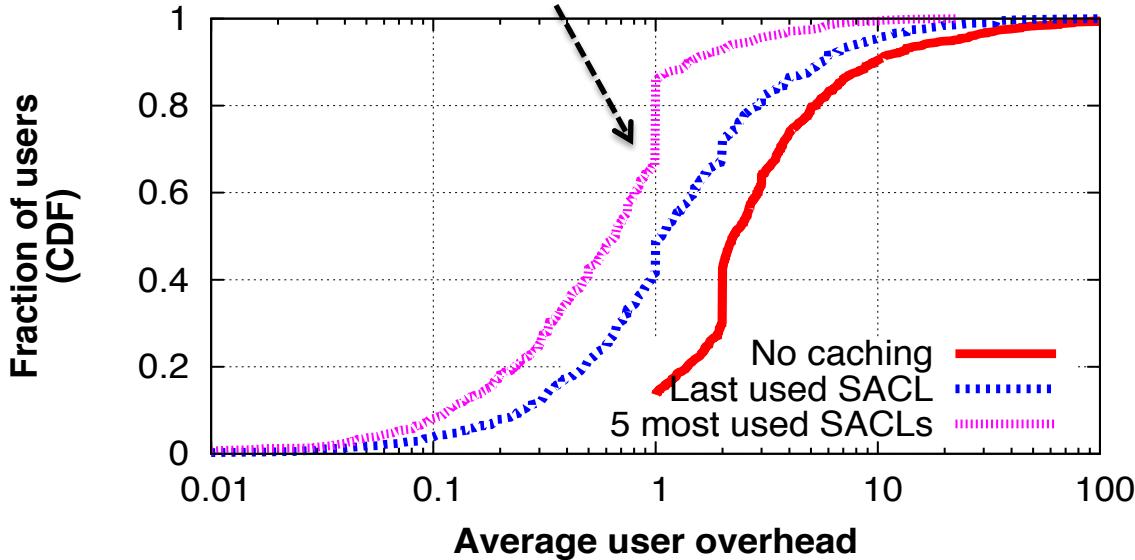
Can we improve by caching a few SACLs instead of one SACL?

Does caching reduce SACL specification overhead?



Does caching reduce SACL specification overhead?

Caching 5 most used SACLs reduces overhead for **86% of users**



SACL specification overhead significantly reduces if we cache just a few SACLs!

Limiting third party crawlers (data aggregators) in social media

- Questions
 - Why are third party crawlers bad?
 - What do they do with the data?
 - How to protect against these data aggregators?

Limiting third party crawlers (data aggregators) in social media

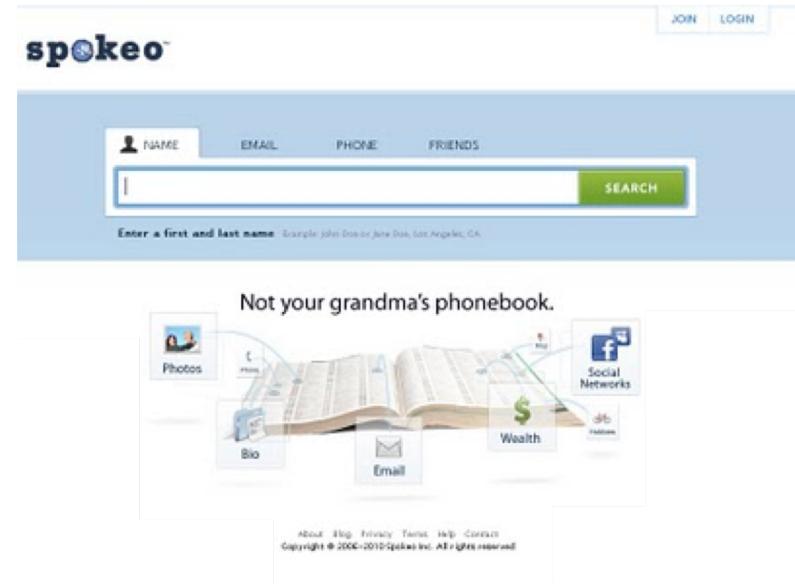
- Questions
 - Why are third party crawlers bad?
 - What do they do with the data?
 - How to protect against these data aggregators?

Why are third party crawlers bad?

Recap: Example of a third party crawler -- Spokeo

Service aggregating “public” data
from web

Others get all of this data
by searching Spokeo



After aggregation: Inferring non public data become easier
Users complained of privacy violation



Problem: Third party crawlers **violate exposure**

Problem: Third party crawlers **violate** exposure



Problem: Third party crawlers **violate** exposure



Crawlers like **Spokeo.com**,
pipl.com are **not** in the
user's **expected exposure**
set

Problem: Third party crawlers **violate** exposure



Crawlers like **Spokeo.com**, **pipl.com** are **not** in the user's **expected exposure set**

Crawlers can **republish/sell** user data in easily accessible form

Problem: Third party crawlers **violate** exposure



Crawlers like **Spokeo.com**, **pipl.com** are **not** in the user's **expected exposure set**

Crawlers can **republish/sell** user data in easily accessible form

In 2010, 171 M Facebook user's data published in torrent

Problem: Third party crawlers **violate** exposure



Crawlers like **Spokeo.com**, **pipl.com** are **not** in the user's **expected exposure set**

Crawlers can **republish/sell** user data in easily accessible form

Problem: Third party crawlers **violate** exposure



Crawlers like **Spokeo.com**, **pipl.com** are **not** in the user's **expected exposure set**

Crawlers can **republish/sell** user data in easily accessible form
Problem for **OSM operators**

User data is valuable asset to OSM operators

OSM operators are blamed for misuse of user data [NYTimes '10]

Problem: Third party crawlers **violate** exposure



Crawlers like **Spokeo.com**, **pipl.com** are **not** in the user's **expected exposure set**

Crawlers can **republish/sell** user data in easily accessible form
Problem for **OSM operators**

User data is valuable asset to OSM operators

OSM operators are blamed for misuse of user data [NYTimes '10]

OSMs need to limit large-scale third party crawlers

Limiting third party crawlers (data aggregators) in social media

- Questions
 - Why are third party crawlers bad?
 - **What do they do with the data?**
 - How to protect against these data aggregators?

**What can be done with
aggregated data?**

Sale the data to databrokers

- What are data brokers [from wiki]
 - ...collects information about individuals from public records and private sources...

Sale the data to databrokers

- What are data brokers [from wiki]
 - ...collects information about individuals from public records and private sources...
 - ...**census and change of address records**, motor vehicle and driving records, user-contributed material to **social networking sites**...

Sale the data to databrokers

- What are data brokers [from wiki]
 - ...collects information about individuals from public records and private sources...
 - ...**census and change of address records**, motor vehicle and driving records, user-contributed material to **social networking sites**...
 - The data are aggregated to create **individual profiles**

Sale the data to databrokers

- What are data brokers [from wiki]
 - ...collects information about individuals from public records and private sources...
 - ...**census and change of address records**, motor vehicle and driving records, user-contributed material to **social networking sites**...
 - The data are aggregated to create **individual profiles**
 - **Profiles:** age, race, gender, height, weight, marital status, religious affiliation, ..., net worth, home ownership status, investment habits, product preferences, health-related interests

What do data brokers do with this data?

- They sell it (Duh!)
 - To marketers and ad companies
 - To government agencies
 - To background checkers
 - To anybody who can afford...

Facebook used to partner with data brokers

- Why? : For getting more data about users
 - Study: **Auditing Offline Data Brokers via Facebook's Advertising Platform**
 - Venkatadri et al., WWW 2019
 - <http://www.ccs.neu.edu/home/amislove/publications/DataBrokers-WWW.pdf>

Facebook used to partner with data brokers

- Why? For getting more data about users
 - Study: **Auditing Offline Data Brokers via Facebook's Advertising Platform**
 - Venkatadri et al., WWW 2019
 - <http://www.ccs.neu.edu/home/amislove/publications/DataBrokers-WWW.pdf>
- In a nutshell:
 - They targeted participants with Transparency Enhancing Ads (Treads)
 - Treads: embed the targeting parameter in the image so that users know what Facebook thinks their attributes are

Lastly: Some data brokers

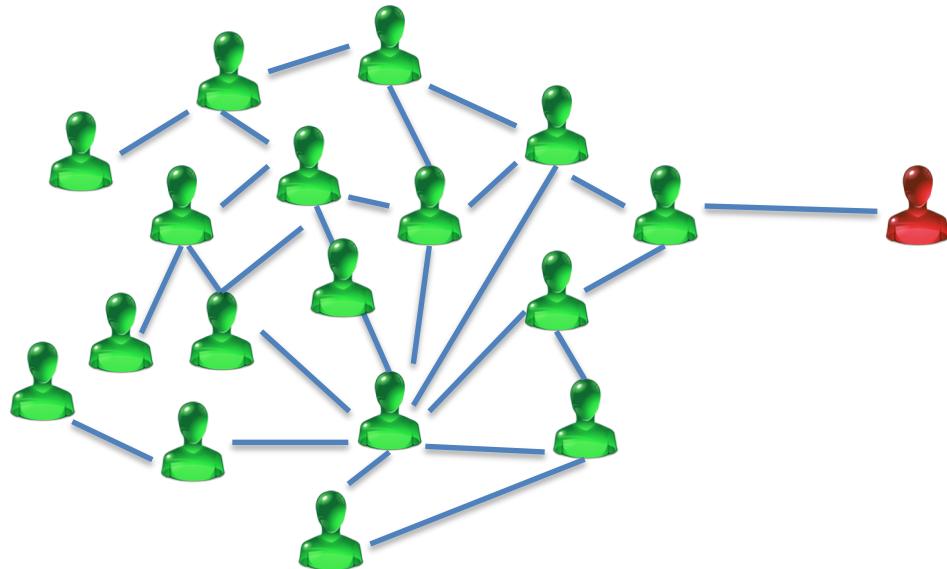
Country	Partner	Attribute count	Targetable		Percent
			Overall	Partner	
U.S.	All	507	210M	190M	90.5%
	Acxiom	128	210M	160M	76.2%
	Datalogix	350	210M	160M	76.2%
	Others ⁵	10	210M	150M	71.4%
	Experian	5	210M	140M	66.7%
	Epsilon	14	210M	130M	61.9%
Australia	All	58	16M	13M	81.3%
	Experian	34	16M	12M	75.0%
	Acxiom	24	16M	9.1M	56.9%
U.K.	All	139	39M	29M	74.4%
	Acxiom	103	39M	22M	56.4%
	Datalogix	19	39M	17M	43.6%
	Experian	17	39M	15M	38.5%
Germany	Acxiom	60	31M	20M	64.5%
France	Acxiom	21	32M	18M	56.3%
Brazil	Experian	20	120M	61M	50.8%
Japan	Acxiom	17	25M	12M	48.0%

Table 1: Coverage of different data brokers across countries with partner categories. We show the total number of broker attributes, the number of Facebook identities that are targetable (Overall), the number of these identities that have at least one attribute from that broker (Partner), and the resulting coverage. Countries with more than one broker have a row indicating the coverage of all the brokers together (All).

**How to stop the third party crawlers
(who help data brokers)?**

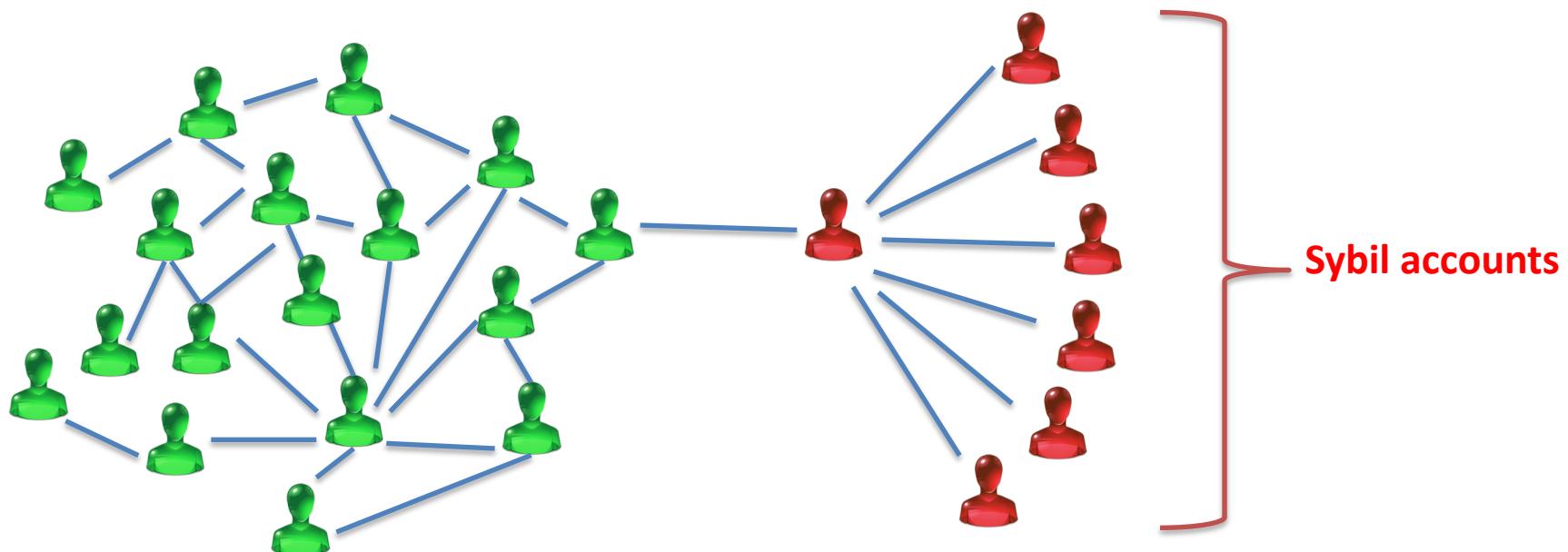
Existing solution: Simple rate-limiting

OSMs rate-limit on per-account or per IP address basis
Crawlers can **defeat rate-limit** using multiple accounts



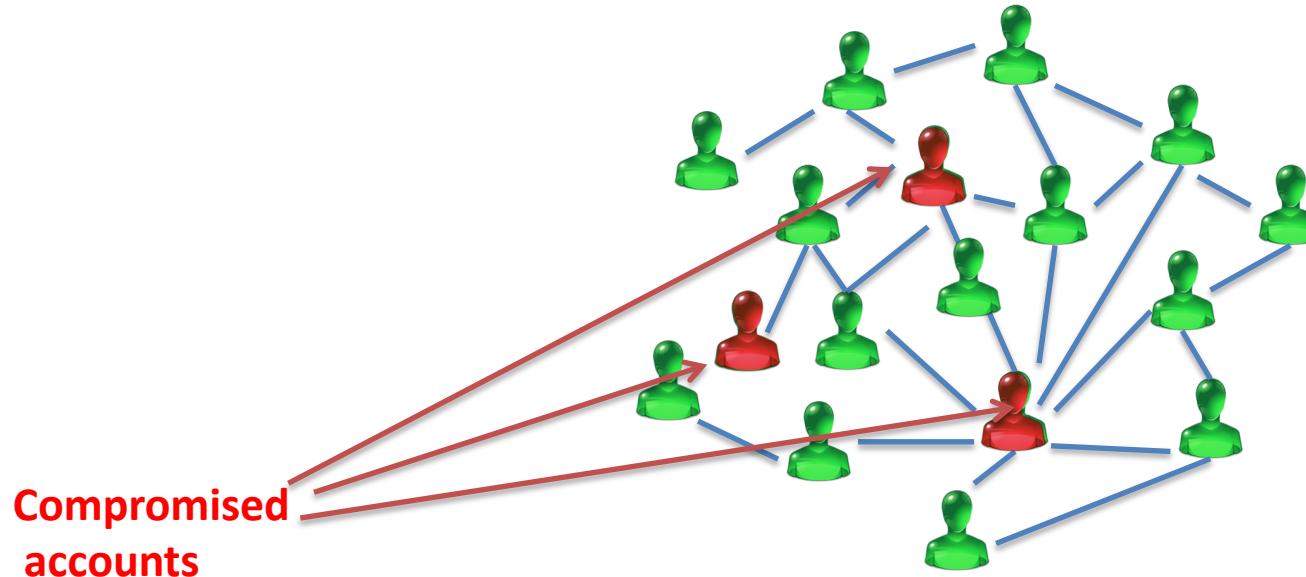
Existing solution: Simple rate-limiting

OSMs rate-limit on per-account or per IP address basis
Crawlers can **defeat rate-limit** using multiple accounts



Existing solution: Simple rate-limiting

OSMs rate-limit on per-account or per IP address basis
Crawlers can **defeat rate-limit** using multiple accounts



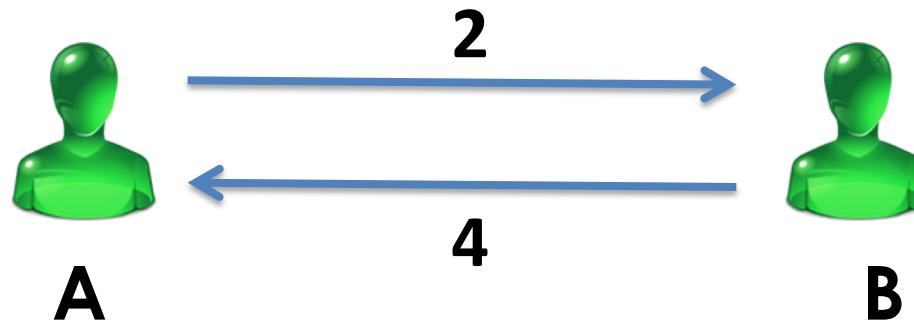
Our solution: Genie

Intuition: Links to good users are harder to get than accounts

Replace user account based rate limiting with link based rate limiting

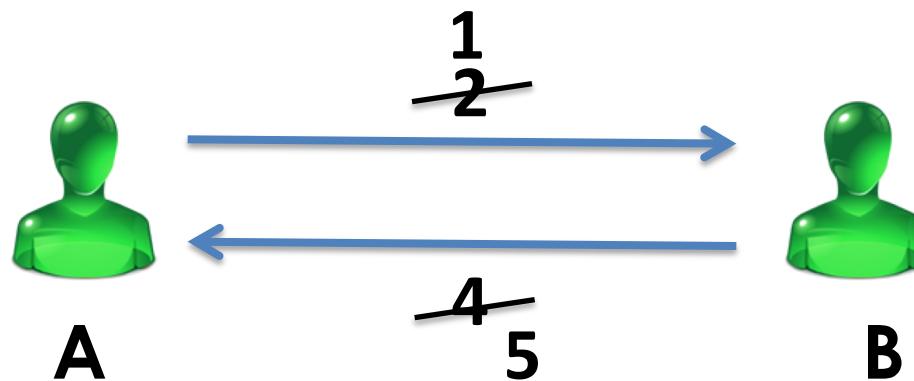
Credit Networks [EC '11]

Nodes trust each other by providing pair-wise credit
Credit is used to pay for the services received



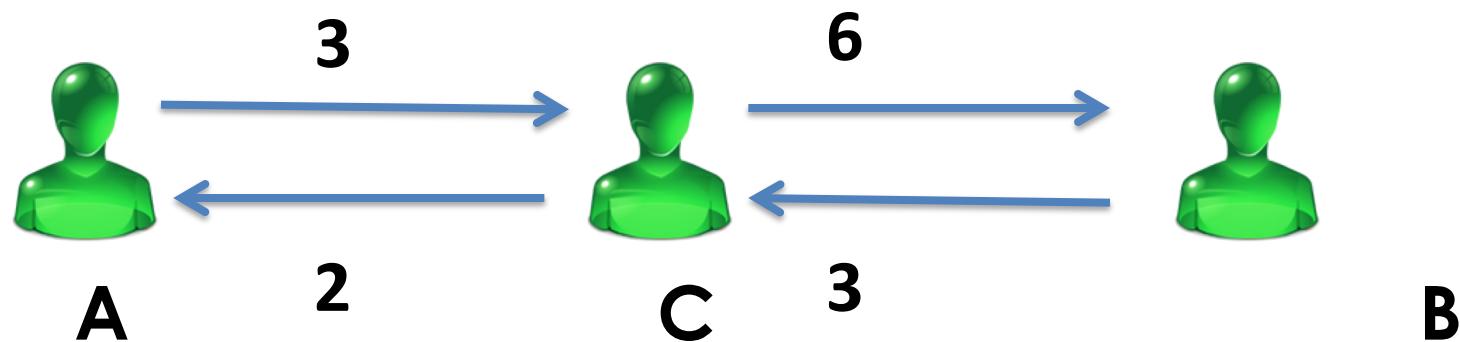
Credit Networks [EC '11]

Nodes trust each other by providing pair-wise credit
Credit is used to pay for the services received



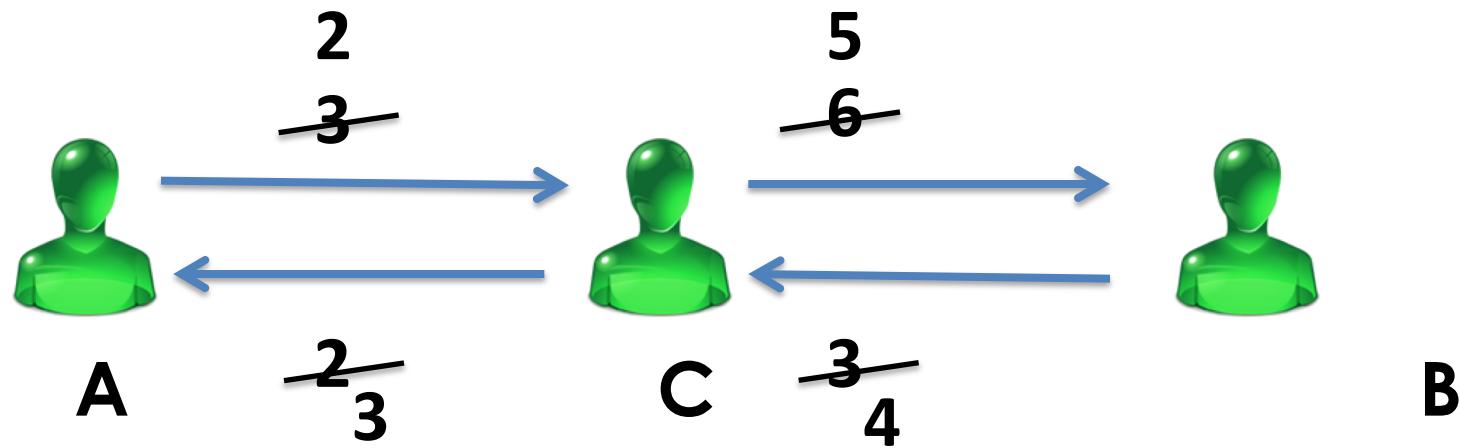
Credit Networks [EC '11]

Nodes trust each other by providing pair-wise credit
Credit is used to pay for the services received



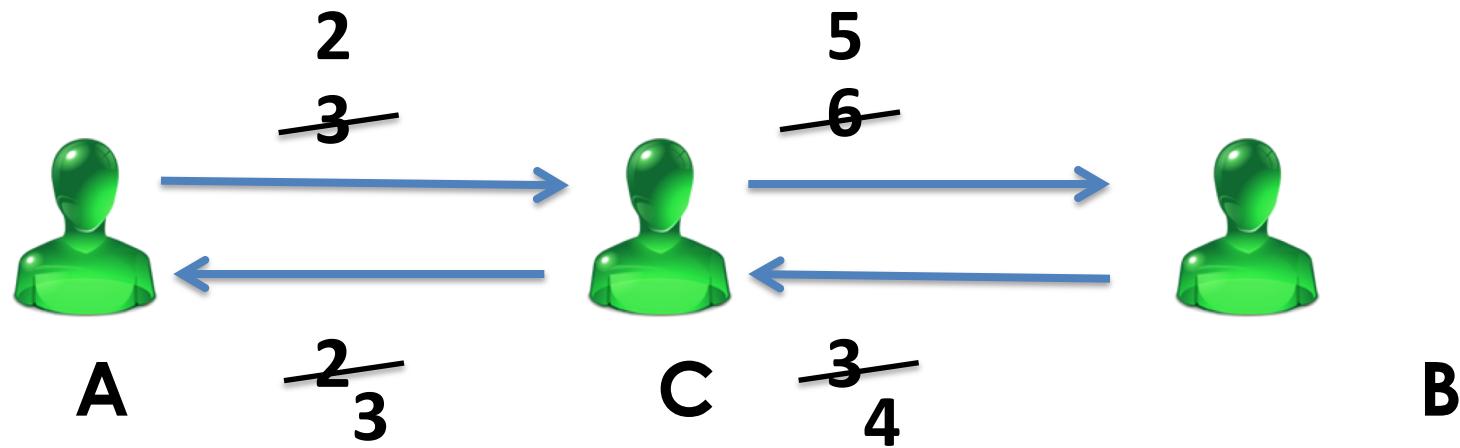
Credit Networks [EC '11]

Nodes trust each other by providing pair-wise credit
Credit is used to pay for the services received



Credit Networks [EC '11]

Nodes trust each other by providing pair-wise credit
Credit is used to pay for the services received



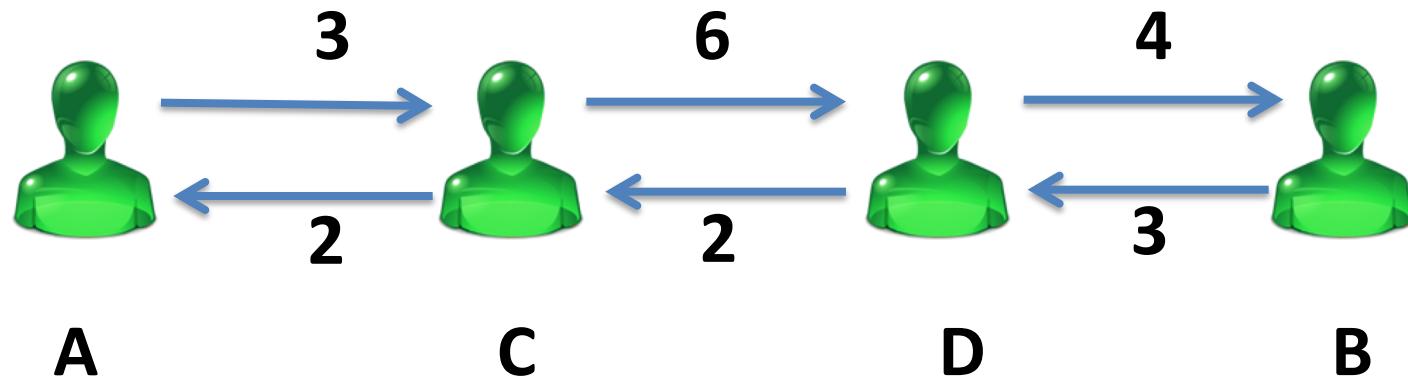
To obtain a service, find path(s) with sufficient credits

How can we map OSM to credit networks ?

OSM operator forms credit network from the social network

Operator assigns initial credit to the links

Credit deducted from links to view another user's profile

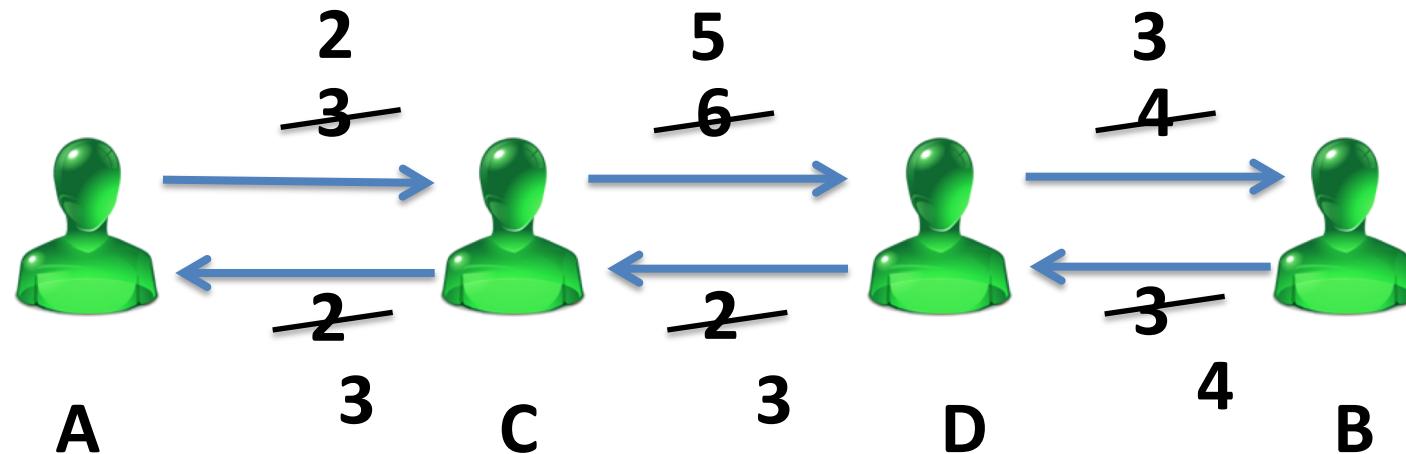


How can we map OSM to credit networks ?

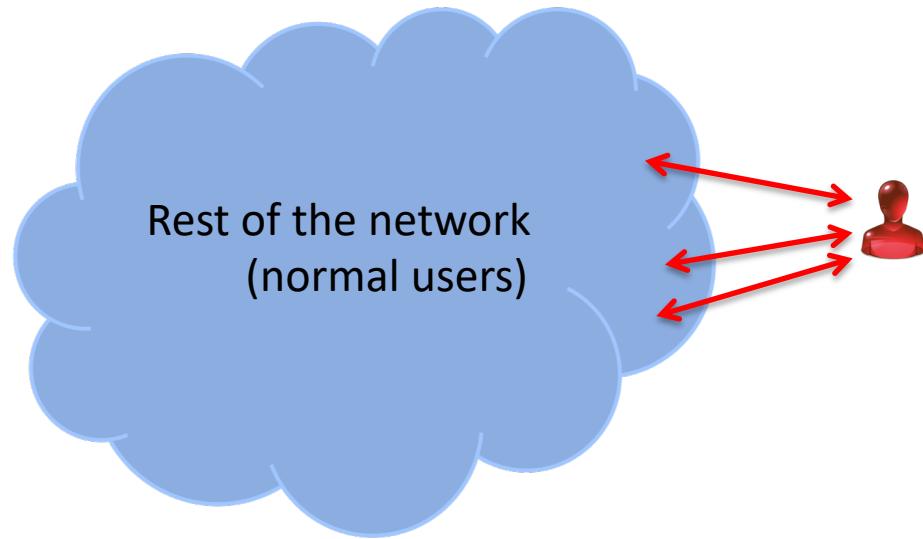
OSM operator forms credit network from the social network

Operator assigns initial credit to the links

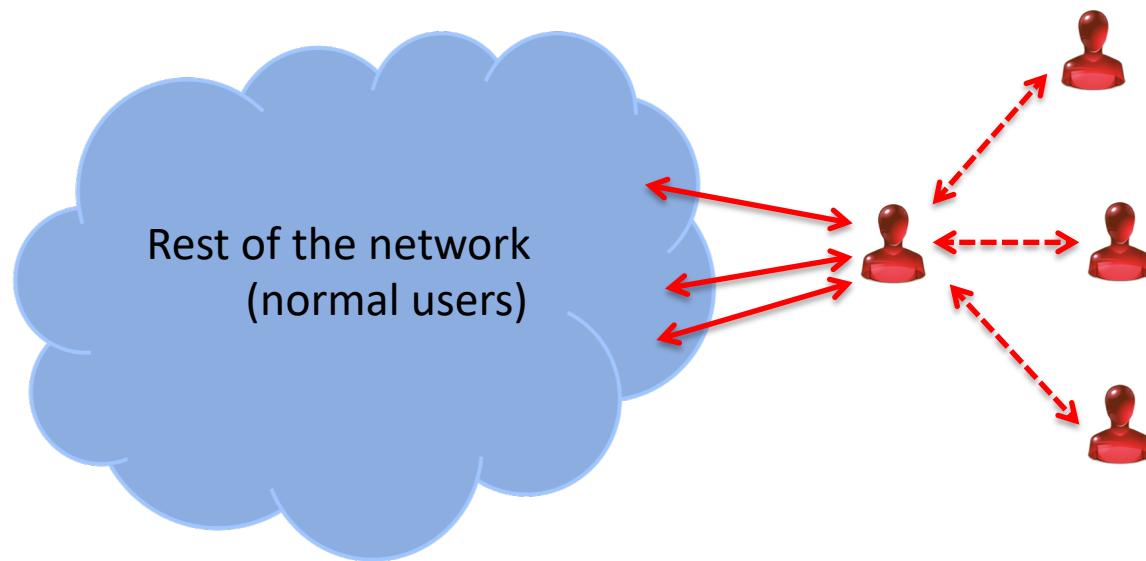
Credit deducted from links to view another user's profile



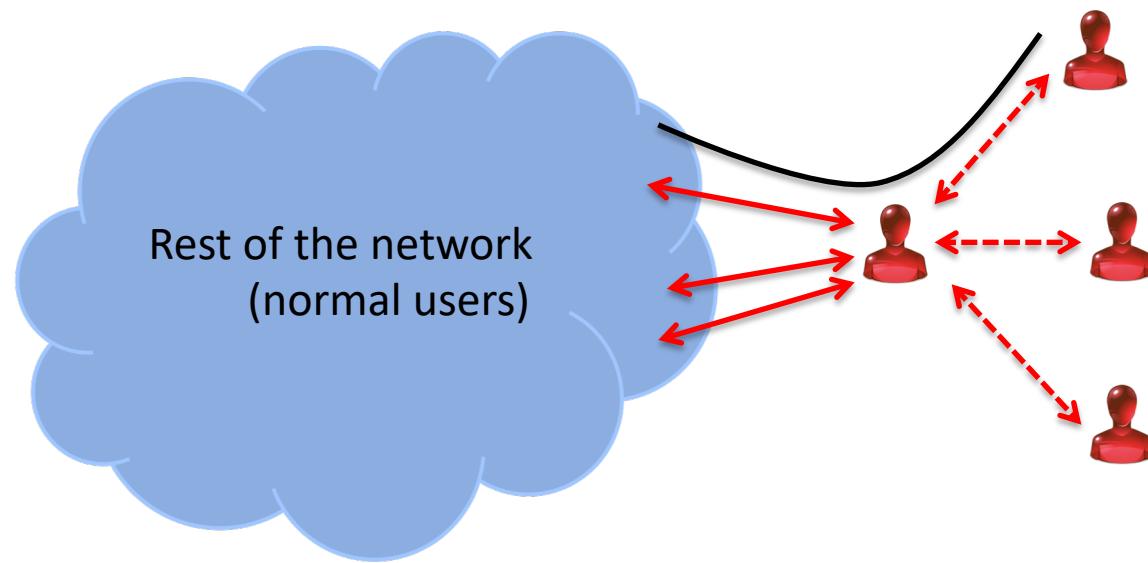
How do credit network defend against crawlers?



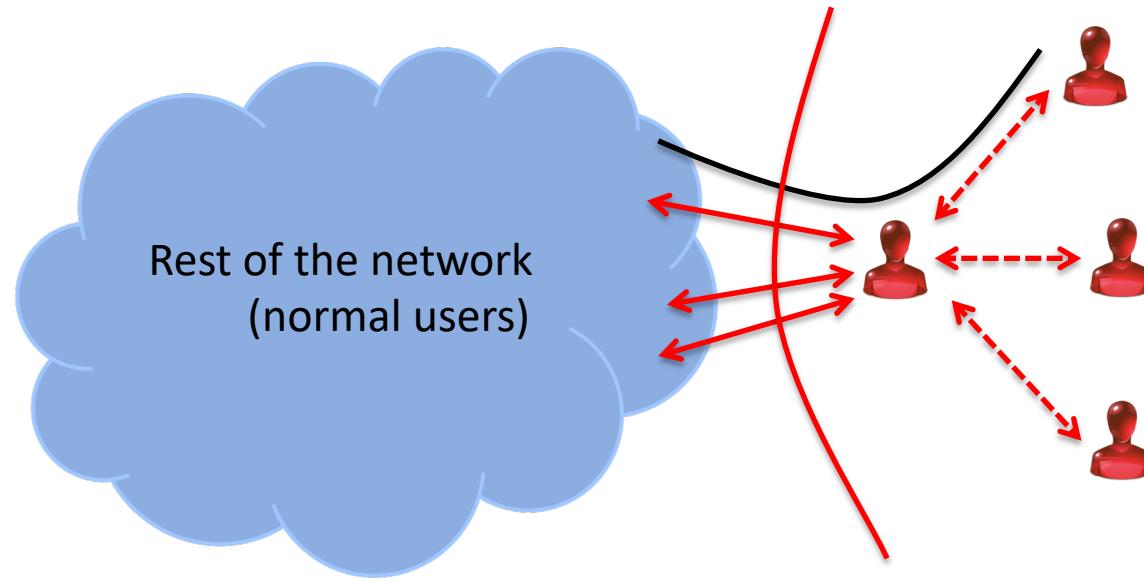
How do credit network defend against crawlers?



How do credit network defend against crawlers?

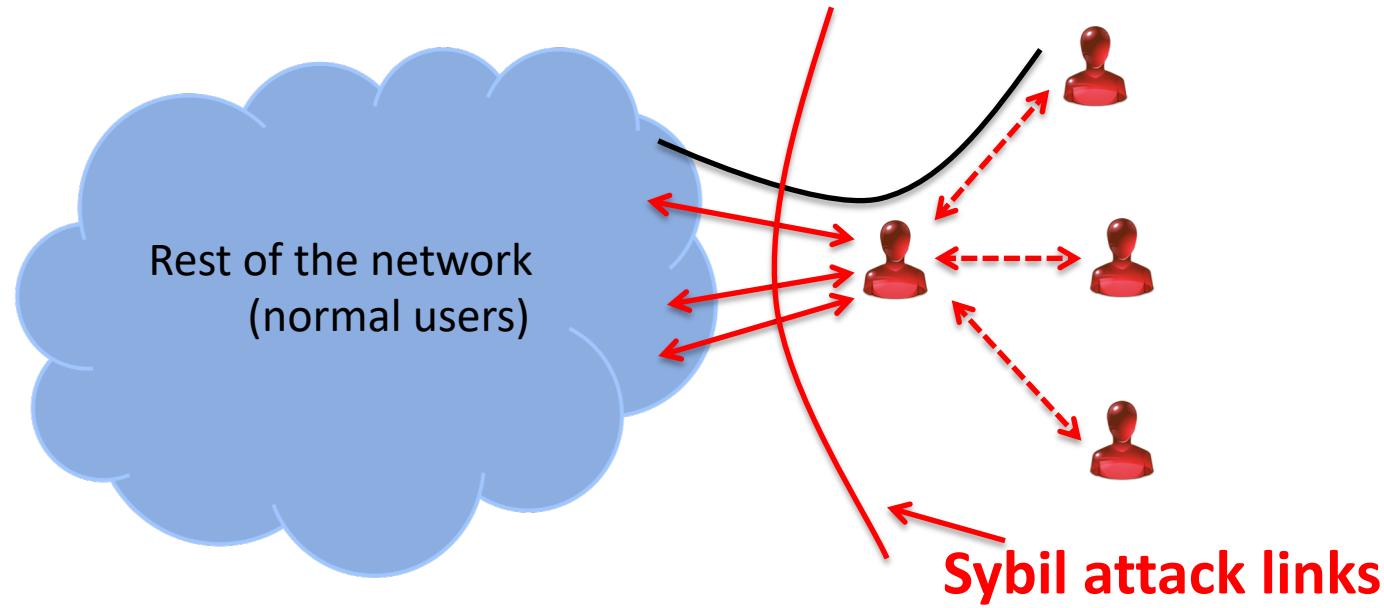


How do credit network defend against crawlers?



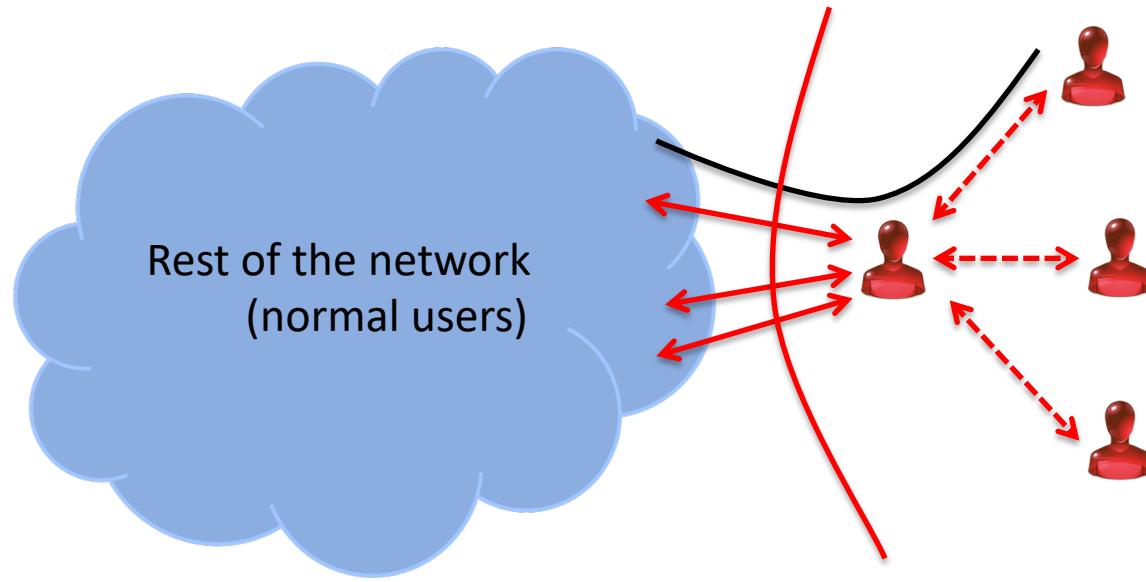
Amount of crawling proportional to attack links

How do credit network defend against crawlers?



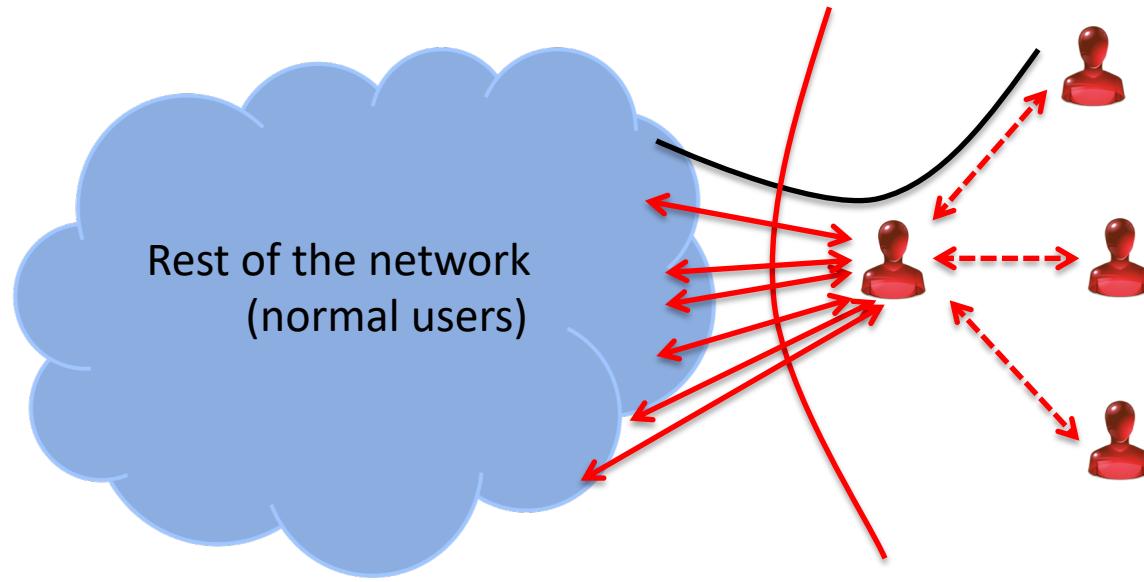
Amount of crawling proportional to attack links

How do credit network defend against crawlers?



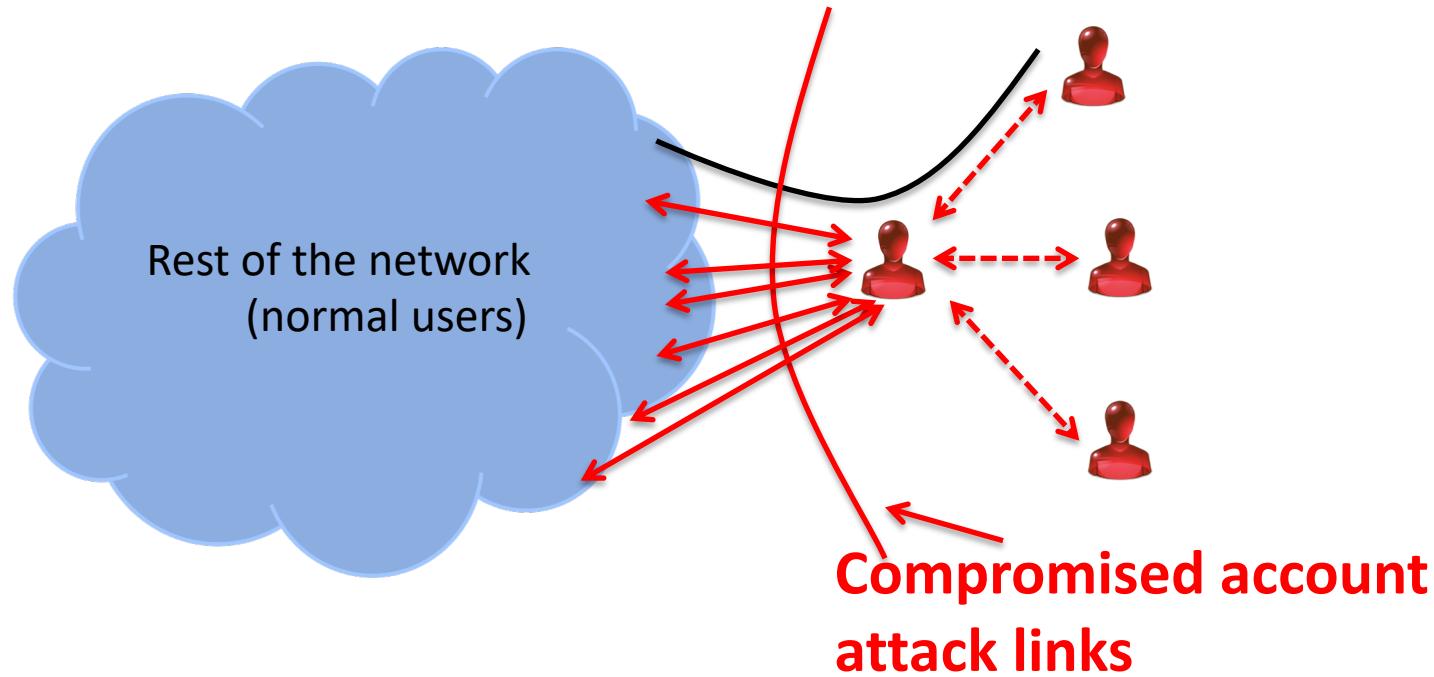
Amount of crawling proportional to attack links

How do credit network defend against crawlers?



Amount of crawling proportional to attack links

How do credit network defend against crawlers?



Amount of crawling proportional to attack links

Strengthening Genie against compromised accounts

Real-world data driven finding:

More than 92% of the normal views are local (1 or 2 hop away)

Strengthening Genie against compromised accounts

Real-world data driven finding:

More than 92% of the normal views are local (1 or 2 hop away)

New charging model: Pay more to view profiles far away

Credit charged per link = Distance between two nodes - 1

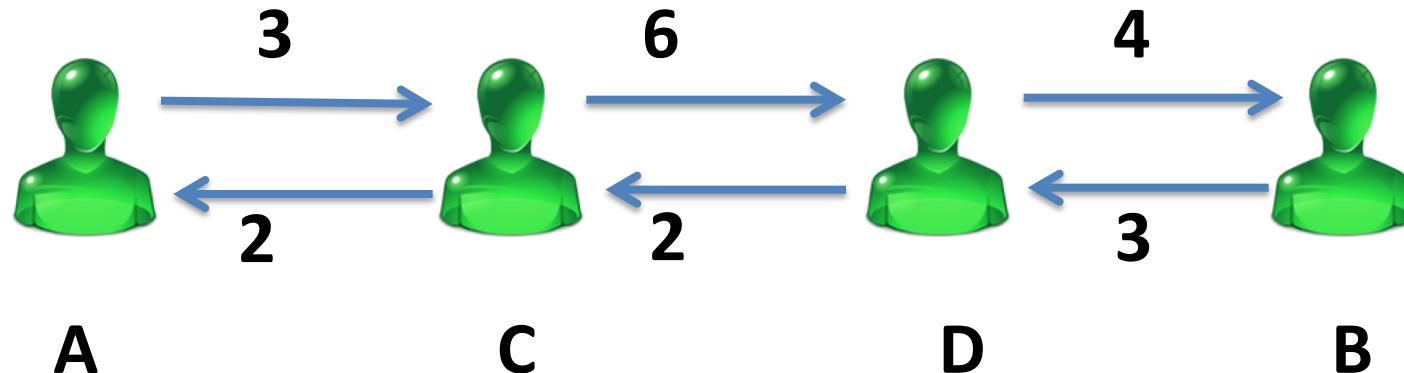
Strengthening Genie against compromised accounts

Real-world data driven finding:

More than 92% of the normal views are local (1 or 2 hop away)

New charging model: Pay more to view profiles far away

Credit charged per link = Distance between two nodes - 1



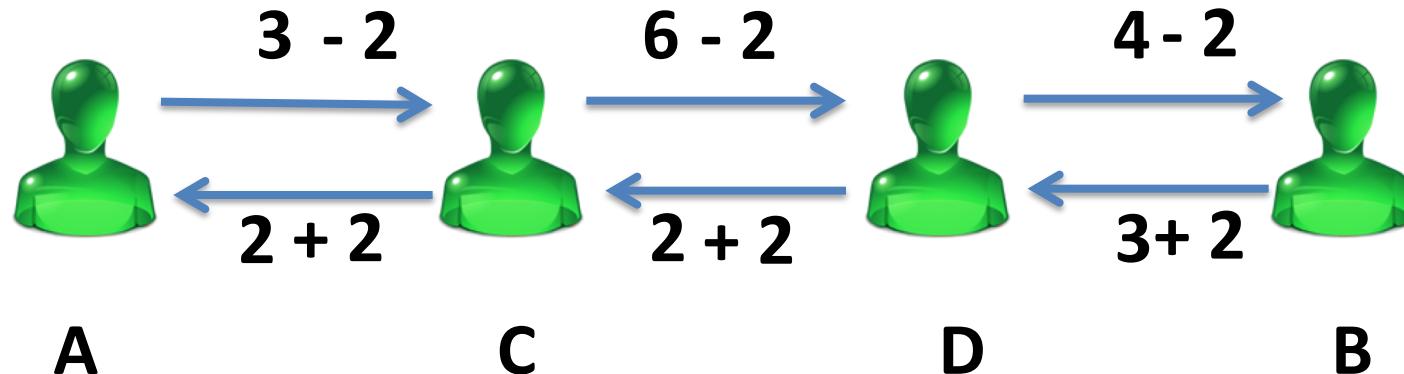
Strengthening Genie against compromised accounts

Real-world data driven finding:

More than 92% of the normal views are local (1 or 2 hop away)

New charging model: Pay more to view profiles far away

Credit charged per link = Distance between two nodes - 1



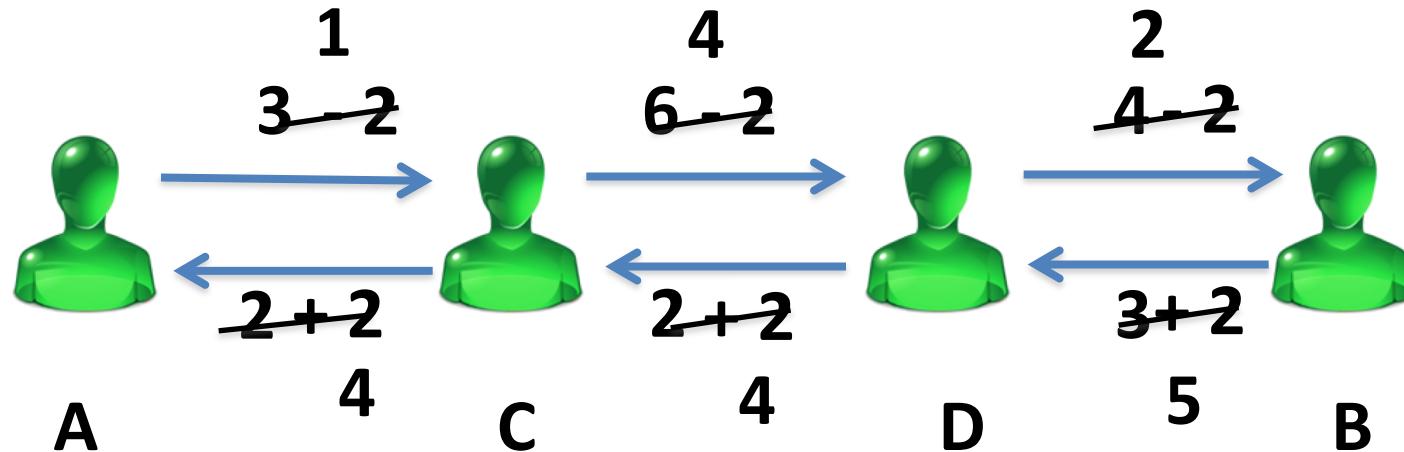
Strengthening Genie against compromised accounts

Real-world data driven finding:

More than 92% of the normal views are local (1 or 2 hop away)

New charging model: Pay more to view profiles far away

Credit charged per link = Distance between two nodes - 1



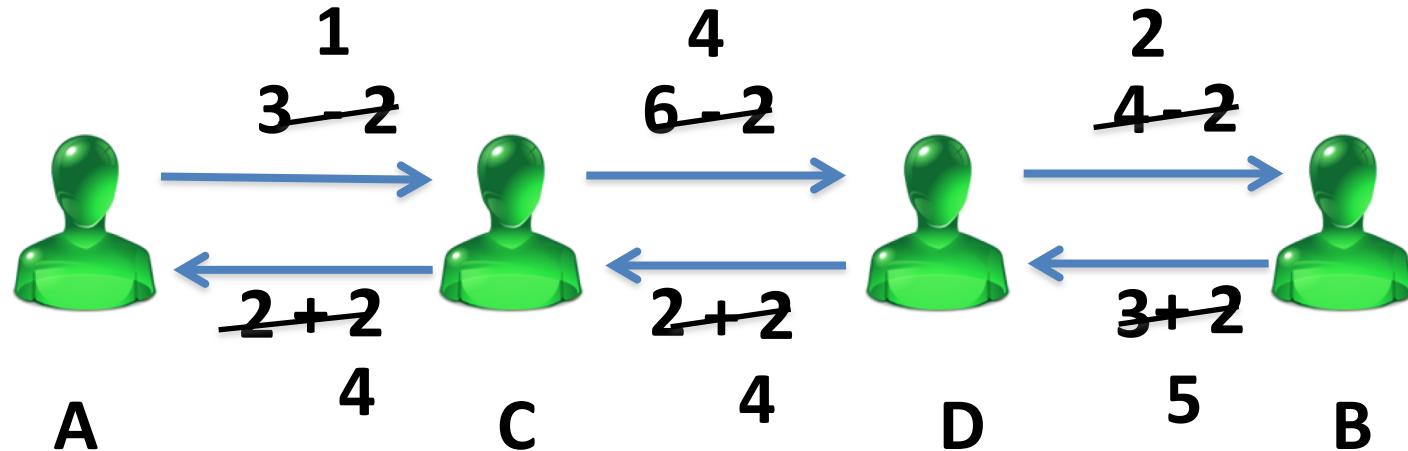
Strengthening Genie against compromised accounts

Real-world data driven finding:

More than 92% of the normal views are local (1 or 2 hop away)

New charging model: Pay more to view profiles far away

Credit charged per link = Distance between two nodes - 1



Rate of crawling decreases with increased path length

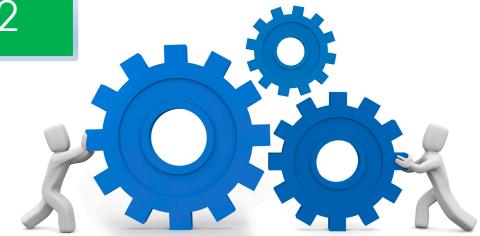
Genie evaluation

Genie simulator using **Canal** Library

EuroSys'12

Input: Social graph and user activity trace

Output: allowed/flagged user activity



Genie evaluation

Genie simulator using **Canal** Library

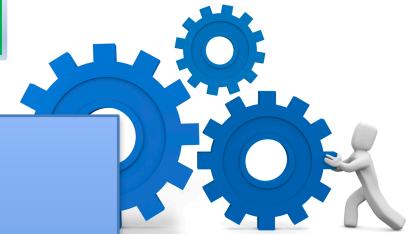
EuroSys'12

Input: Social graph and user activity trace

Output: allowed/flagged

Make graph operations fast!

Less than millisecond latency for
million node graphs



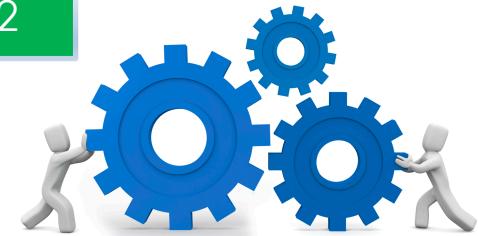
Genie evaluation

Genie simulator using **Canal** Library

EuroSys'12

Input: Social graph and user activity trace

Output: allowed flagged user activity



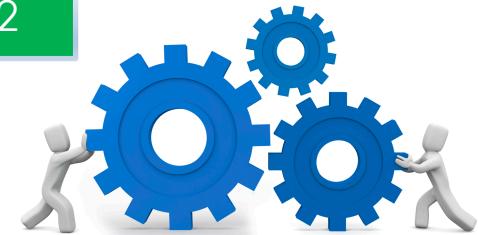
Genie evaluation

Genie simulator using **Canal** Library

EuroSys'12

Input: Social graph and user activity trace

Output: allowed/flagged user activity



We evaluated Genie on viewing activity trace from real world profile viewing data collected by earlier work

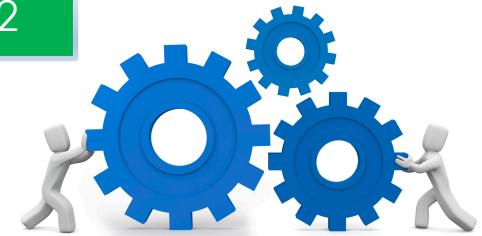
Genie evaluation

Genie simulator using **Canal** Library

EuroSys'12

Input: Social graph and user activity trace

Output: allowed flagged user activity



We evaluated Genie on viewing activity trace from real world profile viewing data collected by earlier work

Genie **heavily slows down** crawlers

Allows majority of the **normal user activity**

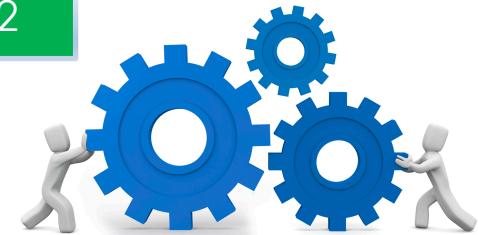
Genie evaluation

Genie simulator using **Canal** Library

EuroSys'12

Input: Social graph and user activity trace

Output: allowed flagged user activity



We evaluated Genie on viewing activity trace from real world profile viewing data collected by earlier work

Genie **heavily slows down** crawlers

Allows majority of the **normal user activity**

Genie controls exposure by limiting third party crawlers