

Introduction to privacy protection

Part 2 - Design and operation principles

Guillaume Piolle
guillaume.piolle@supelec.fr
<http://guillaume.piolle.fr/>

Master “Machine Learning and Data Mining”, Saint-Étienne

December 16th 2013

Design and operation principles

- 1 Computer security and privacy
- 2 General principles
- 3 Anonymization and deanonymization
- 4 Privacy by Design
- 5 An example use-case: Social networks

Computer security and privacy

Is privacy protection
a component of computer security?

Classical dimensions of computer security:

- Confidentiality ;
- Integrity ;
- Availability ;
- + authentication, non-repudiation, access control, flow control. . .

Dimensions of personal data protection:

- User information ;
- User consent ;
- Access/correction/deletion rights;
- Finality and proportionality ;
- Retention time ;
- Forwarding to third parties.

Computer security and privacy

- Privacy (or personal data protection) can be considered **from the point of view of computer security**;
- Some operational privacy requirements **can be met** thanks to classical computer security tools;
- Some operational privacy requirements **cannot be met** thanks to classical computer security tools;
- Some operational privacy requirements may be deemed **incompatible** with some computer security requirements.

Sometimes presented as a **sub-discipline**, sometimes as a **connex** or **transversal** discipline, sometimes as a **rival** discipline.

Computer security and privacy

Auditability

A computer security requirement: make sure that we are able to spot malicious or faulty behaviours and to blame the entities responsible for them.

Main tool: maintaining **logs** tracing the activity of the system (piece of software, web server, etc.)

Example: authentication log file (/var/log/auth.log)

```
Sep 29 20:59:01 vpsxxx CRON[14089]: pam_unix(cron:session): session opened for user root
by (uid=0)
Sep 29 20:59:01 vpsxxx CRON[14090]: pam_unix(cron:session): session opened for user root
by (uid=0)
Sep 29 20:59:01 vpsxxx CRON[14090]: pam_unix(cron:session): session closed for user root
Sep 29 20:59:01 vpsxxx CRON[14089]: pam_unix(cron:session): session closed for user root
Sep 29 20:59:46 vpsxxx sshd[14140]: Did not receive identification string
from 161.139.xxx.xxx
Sep 29 21:00:01 vpsxxx CRON[14141]: pam_unix(cron:session): session opened for user root
by (uid=0)
Sep 30 11:53:18 vpsxxx sshd[6842]: Authentication tried for root with correct key
but not from a permitted host
(host=nat-profs.rennes.supelec.fr, ip=193.54.192.3).
Sep 30 11:53:18 vpsxxx sshd[6842]: Authentication tried for root with correct key
but not from a permitted host
(host=nat-profs.rennes.supelec.fr, ip=193.54.192.3).
Sep 30 11:53:21 vpsxxx sshd[6842]: Accepted password for root from 193.54.192.3
port 55130 ssh2
Sep 30 11:53:21 vpsxxx sshd[6842]: pam_unix(sshd:session): session opened
for user root by (uid=0)
```

Example: firewall log (in /var/log/messages)

```
Sep 26 10:33:55 vpsxxx kernel: netfilter-input IN=eth0 OUT= MAC=[masqué]  
SRC=90.84.xxx.xxx DST=46.105.yyy.yyy LEN=60 TOS=0x00 PREC=0x00 TTL=48  
ID=25922 DF PROTO=TCP SPT=59766 DPT=8080 WINDOW=5840 RES=0x00 SYN URGP=0
```

```
Sep 26 10:34:47 vpsxxx kernel: netfilter-input IN=eth0 OUT= MAC=[masqué]  
SRC=90.84.xxx.xxx DST=46.105.yyy.yyy LEN=60 TOS=0x00 PREC=0x00 TTL=48  
ID=20314 DF PROTO=TCP SPT=55315 DPT=8080 WINDOW=5840 RES=0x00 SYN URGP=0
```

```
Sep 26 10:34:48 vpsxxx kernel: netfilter-input IN=eth0 OUT= MAC=[masqué]  
SRC=90.84.xxx.xxx DST=46.105.yyy.yyy LEN=60 TOS=0x00 PREC=0x00 TTL=48  
ID=20315 DF PROTO=TCP SPT=55315 DPT=8080 WINDOW=5840 RES=0x00 SYN URGP=0
```

Example: Apache access log file (/var/log/apache2/access.log)

```
212.113.aaa.aaa - [28/Sep/2011:15:53:07 +0200] "GET / HTTP/1.1" 200 460 "-" "Mozilla/5.0
(compatible; MJ12bot/v1.4.0; http://www.majestic12.co.uk/bot.php?+)"
188.165.bbb.bbb - [28/Sep/2011:20:07:05 +0200] "GET /phpmyadmin/main.php HTTP/1.0" 200 977
 "-" "-"
188.165.bbb.bbb - [28/Sep/2011:20:07:07 +0200] "GET /phpmyadmin/libraries/select_lang.lib.p
HTTP/1.0" 403 522 "-" "-"
188.40.ccc.ccc - [28/Sep/2011:20:21:55 +0200] "GET / HTTP/1.0" 200 453 "-" "-"
188.40.ccc.ccc - [28/Sep/2011:20:21:55 +0200] "GET / HTTP/1.0" 200 453 "-" "-"
188.40.ccc.ccc - [28/Sep/2011:20:21:56 +0200] "HEAD / HTTP/1.1" 200 276 "-"
"Visited by http://tools.geek-tools.org"
188.165.ddd.ddd - [28/Sep/2011:22:49:36 +0200] "GET /w00tw00t.at.ISC.SANS.test0:)
HTTP/1.1"400 513 "-" "-"
213.91.eee.eee - [29/Sep/2011:00:37:47 +0200] "GET /w00tw00t.at.ISC.SANS.DFind:)
HTTP/1.1" 400 513 "-" "-"
85.214.fff . fff - [29/Sep/2011:10:03:39 +0200] "GET /w00tw00t.at.ISC.SANS.DFind:)
HTTP/1.1" 400 513 "-" "-"
209.160.ggg.ggg - [29/Sep/2011:15:30:17 +0200] "GET / HTTP/1.0" 206 498 "-"
"Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1)"
```


Computer security and privacy

Logging obligations in France

- 2001: *Loi sur la Sécurité Quotidienne (LSQ)*: ISPs must keep connection data for one year (temporary measure, extended *ad vitam*);
- 2004: *Loi sur la Confiance dans l'Économie Numérique (LCEN)*: retention of data to identify users uploading content (extended to all service providers);
- 2011: application decree of LCEN: retention of logins, pseudonyms, passwords, payment info, contact info.

What if I don't log enough?

Fine up to 375 k€ for a company, 75 k€ and one year's imprisonment for its CEO.

Computer security and privacy

Who may access log files ?

- Justice (rogatory commission, emergency or standard court ruling);
- The police, with a simple requisition (without judiciary authorization/oversight), since the January 23rd 2006 law against terrorism;
- The network/system administrator, who has an **obligation of confidentiality** (even towards his employer) and can only access the data in the context of his network security mission (*Cour de Cassation*, June 17th 2009).

An aggravated operational risk

With a security excuse (anti-terrorism), the risk of loss/harm in case of intrusion has been increased, and potential attackers have been given an incentive (they know the kind of valuable data they are supposed to find).

Computer security and privacy

Trade-offs with other concepts

- Freedom of speech ;
- Right to information ;
- Data dissemination and availability (e.g. in social network systems);
- Usability (procedures, privacy policies).

Computer security and privacy

Complementarity between security and privacy

Actors are often the same...

...but not always.

System designers, administrators, information system security officers, ISPs, service providers... What are the imperatives and constraints of each?

Similar issues

“Hard” problems in information security and in privacy protection are very close: information flow control, and efficient remote usage control.

Computer security and privacy

Auditability is not only a threat

- Article 39 of the I&L law ;
- Obligation to notify privacy breaches;
- Importance of the capability to demonstrate one's commitment to privacy.

Computer security tools can be useful to privacy protection

- Classical encryption tools;
- More “exotic” cryptography (IBE/PBE, repudiability guarantees, plausible deniability, group signature. . .);
- Converging research in usage and flow control.

Privacy according to ISO

Common Criteria for Information Technology Security Evaluation

Norm ISO/IEC 15408, successor of the DoD *Orange Book*.

Section 14 : privacy.

Technical requirements for privacy protection

- **Anonymity**: inability of observers to determine the identity of a user;
- **Pseudonymity**: same thing, except that the user must remain accountable for his actions;
- **Unlinkability**: inability of observers to determine whether two actions have been performed by the same user;
- **Unobservability**: Unability of observers to determine whether an action is being performed.

General principles

Data sovereignty

Ensure that **users remain in control** of the data relating to them:

- Store data and keys in priority on their own personal devices;
- Enforce a tight control over usage and dissemination, imposing obligations (security obligations, obligations to notify, to delete...).

General principles

Data minimization

cf. proportionality principle

- Collect only data absolutely necessary to the finality;
- Keep/share them only if absolutely necessary;
- Destroy, as soon as possible, all data non absolutely necessary;

All within the limits of system auditability obligations.

Anonymous and anonymized databases

- Absence of data allowing the identification of a physical person:
 - Deletion of name and surname;
 - Replacement by a random number;
 - Replacement by arbitrary pseudonyms;
- Anonymous polls, official or other;
- Polls and forms whose identifying section is afterwards separated from the rest.

Legal framework for anonymous databases

If we consider that there are no “personal data” because there are no identifying data, then the I&L law (or the 95/46 directive) **does not apply!**

Consequences:

- **No rights** for data subjects;
- **No obligations** for data processors;
- **No limitations** to retention, publication, exploitation, correlation with other databases.

But... no problem, since everything is anonymous?

Perfect anonymization is impossible

**The anonymization of a database in the eye of the law
IS NOT ENOUGH TO FORBID
IDENTIFICATION OF INDIVIDUALS
in most cases.**

Deanonimization: the AOL case

July 2006

AOL publishes a history of 20 million queries (3 months, 658000 users).

Good intentions

- Communication to the scientific community, for research purposes;
- “Anonymized” queries (but a unique ID for each user).

Deanonymization: the AOL case

Data dissemination

- Data gets republished by many mirror sites;
- Many specialized analysis tools are developed.

Deanonymization

Many users are re-identified.

The user n° 4417749 queried AOL about:

- 60 year old singles;
- Hand tremor;
- Effects of nicotine;
- Dogs urinating everywhere;
- ...

Deanonymization: the AOL case

User n° 4417749

Thelma Arnold, 62 ans, Lilburn, Georgia.

Identified by the New York Times, identity published with the user's authorization.



Deanonymization: the AOL case

Consequences

Many query compilations and derived works are made, some of them exotic or artistic.

Catastrophic commercial impact for AOL, two lay-offs, resignation of the Chief Technology Officer, one class action in California.

Educational impact: stress is put on the risk related to search engines, not on data deanonymization (not very technical in this case).

Deanonimization : the Netflix case

A harmless initiative

Netflix : collaborative website including movie evaluations and recommendations.

2010 : Netflix publishes anonymous evaluation data, in the context of a competition (*Netflix prize*, 1M\$) aiming at improving its recommendation algorithm.

A researcher correlates anonymous data with IMDb's info and "deanonimizes" the database. The users' cinematographic preferences become personally identifying information!

Knowing two evaluations is enough to identify 68 % of users.

Federal complaint filed, Netflix retracts and puts an end to the competition, saying it bears a risk to privacy.

Deanonimization : The adventures of Latanya Sweeney, episode 1

The GIC case

Mid 90's: Massachusetts's *Group Insurance Commission* decides to publish "anonimized" data about hospitalizations of state employees.

L. Sweeney, master student at Carnegie Mellon, correlates the data with local electoral rolls and sends his medical file to the governor.

After that, the governor made GIC back off. . .

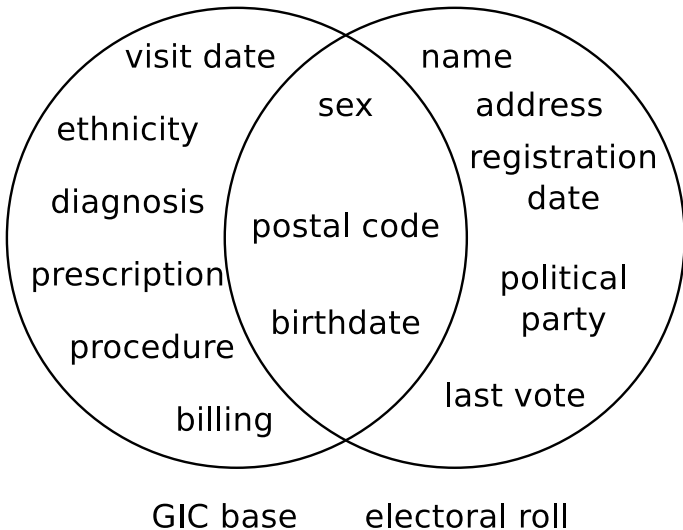
Deanonimization : The adventures of Latanya Sweeney, episode 2

Introduction of k -anonymity

2000: L. Sweeney shows that 87 % of U.S. citizens can be uniquely identified by their sex, birthdate and postal code (through an easy correspondence with public registers).

2002 paper: she introduces the concept of k -anonymity, which quantifies the degree of anonymity of an “anonymous” database.

Deanonimization : the database interconnection issue



Deanonimization : an “anonymous” database

(Fictional) anonymous poll on a university campus

Sex	Height	Sexual orientation
...
M	170-180	heterosexual
M	170-180	heterosexual
M	170-180	heterosexual
M	170-180	bisexual
M	170-180	heterosexual
M	170-180	heterosexual
M	170-180	heterosexual
M	170-180	heterosexual
M	180-190	heterosexual
M	180-190	homosexual
M	180-190	heterosexual
M	180-190	other
M	180-190	heterosexual
M	190-200	heterosexual
M	200-210	homosexual

Sexual orientation : **sensitive**
in the meaning of I&L art. 8.

HOWEVER: completely
anonymous poll... therefore
outside the scope of the law!

- How much anonymous is it really?
- Are all students equally “anonymous”?
- To which questions do you answer during “anonymous” polls?

Deanonimization : principles of k-anonymity

Quasi-identifier

Set of database attributes allowing, in at least one case, the identification of a tuple with the help of external information.

Any attribute may belong to a quasi-identifier!

Sheds a new light on the concept of “personal data” (art. 2 of I&L).

k-Anonymity

A database is said to be **k-anonymous** if any tuple is indistinguishable from at least $k - 1$ other tuples in the database, projected on any quasi-identifier.

Deanonymisation : back to the example database

(Fictional) anonymous poll on a university campus

Sex	Height	Sexual orientation
...
M	170-180	heterosexual
M	170-180	heterosexual
M	170-180	heterosexual
M	170-180	bisexual
M	170-180	heterosexual
M	170-180	heterosexual
M	170-180	heterosexual
M	170-180	heterosexual
M	180-190	heterosexual
M	180-190	homosexual
M	180-190	heterosexual
M	180-190	other
M	180-190	heterosexual
M	190-200	heterosexual
M	200-210	homosexual

La base est **1-anonyme** : the worst case !

At least one person (two here) is re-identifiable with the help of an (easy to make) external database.

It can be said that the individual in bold is 8-anonymous within the base.

Deanonymisation : k-anonymization

Starting from a 1-anonymous base, how to obtain a “publishable” k-anonymous database?

- Falsification or obfuscation of some attributes;
- Aggregation;
- Relational projection;
- Introduction of artificial tuples. . .

Other anonymization (but not necessarily k-anonymization) techniques:

- Noise injection in some attributes;
- Relational selection. . .

Warning!

- k-anonymizing a database also reduces its utility;
- The size and nature of the quasi-identifier depends on every single available external database.

Deanonymisation : k-anonymization

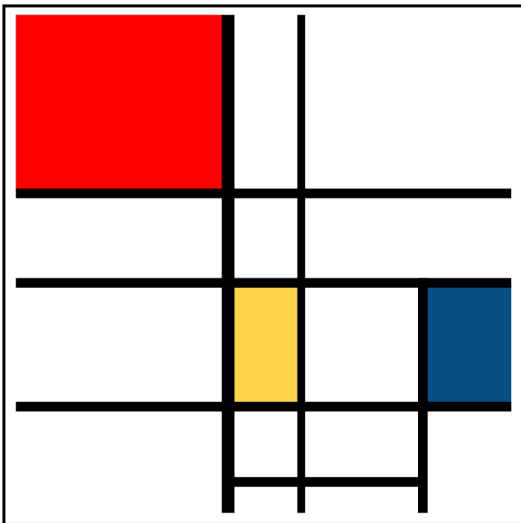
k-anonymization example: Mondrian algorithms

- **Hypothesis:** a database attribute is either part of a QI, or a sensitive attribute;
- **Goal:** make as equivocal as possible the relationship between a QI value and the corresponding sensitive information;
- **Mechanism:** partitioning of the QI space, in order to form clusters of at most k elements, then replacement in the database of the QI values by partition identifiers.

- Start from the set of all QIs;
- Partition this set step by step (by dichotomy for instance), until a maximum number of zones is obtained, each containing at least k elements.

This kind of algorithm can be parametrized by the distribution characteristics in each QI zone, for instance.

Deanonymisation : k-anonymization



Deanonimization : attacks against k-anonymous databases

Unsorted matching attack

Valid only for databases where tuple order is meaningful.

Principle: searching correspondences between two different k-anonymized versions of a same database (the tuples being in the same order in both).

Complementary release attack

Correlation of two excerpts of the same database, k-anonymized in different ways.

Temporal attacks

Comparison between the results of queries made at different times, correlated with external information.

Deanonimization: how to avoid data-mining attacks?

Restriction on multiple queries

Example database: physician \times patient \times prescription

- Patient/prescription queries are very sensitive, and therefore prohibited;
- Physician/patient queries are moderately sensitive;
- Physician/prescription queries are not very sensitive.

Problem

Possible correlation between the results of the last two queries, possibility of a partial deanonimization.

Is it necessary to forbid any query about prescriptions?

Deanonimization: other metrics of anonymity

- I-Diversity (Machanavajjhala 2006) ;
- **Differential privacy** (Dwork 2006) ;
- Closeness (Li 2007) ;
- (c, k) -Safety (Martin 2007) ;
- 3D-Privacy (Chen 2007) ;
- (d, γ) -Privacy (Rastogi 2007) ;
- ϵ -Privacy (Machanavajjhala 2009) ;
- Towards a unified theory of privacy and utility (Kifer 2010) ;
- ...

Deanonimization: in a nutshell

- An anonymous or anonymized database is **never really anonymous** ;
- **Correlation** between several databases may allow **deanonimization** ;
- One may be identified by **any kind of information**.

Privacy by Design

Principle

Privacy protection, just like security, cannot be efficient if it is not thought and integrated from the very beginning of the system design. Later additions will never be able to perfectly seal off privacy breaches stemming from existing design flaws.

The principle, more and more present in the literature and in regulation texts, must involve both technical and non-technical participants, in a coordinated effort.

Implementation examples:

- Specification phase including engineers, lawyers, users and deciders;
- Application of formal conception methods;
- Generalization of *privacy impact assessments*;
- Policy-constrained systems;
- ...

Privacy by Design

The seven principles of PbD

- Proactive rather than reactive, preventive rather than remedial;
- Put privacy as the default setting;
- Embed privacy into design;
- Ensure full functionality: positive sum, not zero sum;
- Ensure end-to-end security, full lifecycle protection;
- Visibility and transparency – keep it open;
- Respect for user privacy – keep it user-centric.

Social Network Systems



Risks associated to social network systems

- **Security risks:** identity theft, phishing, predation, blackmailing, various scams;
- **Profiling risks:** data collection by the SNS platform or by third-party applications, sale of contact lists and social data;
- **e-Reputation risks:** exploration of the network by recruiters, employers, clients. . .

Social Network Systems

Textbook example: Kevin Colvin, 2007

Intern in a British bank in 2007, maybe the first person to be fired because of Facebook, very famous case.

A party picture posted on FB by one of his contacts was used to prove that he had lied on the motivation of an absence.

Nathalie Blanchard, 2009

During a long-term sick leave because of a serious nervous breakdown, her allowance was cancelled because she had posted a picture of herself in a bar, smiling.

A lot of people complain that pictures of them broadcasted on social networks have a negative impact on their employability.

Social Network Systems

Which responsibility for Facebook?

- Their privacy policy is longer (5830 words) than the U.S. constitution;
- 50 different settings related to privacy, more than 170 options;
- Frequent (and sometimes aggressive) updates to the terms of service.

(May 2010)

Consequences: Most users do not bother to tune these settings nor to ask themselves the (complex!) corresponding questions.

Social Network Systems

Third-party applications, a source of risks

- Data access and sharing conditions are often abusive (possibly leading to behavioural profiling);
- Some apps demand that users stop using SSL/TLS encryption;
- Some application organize the spam of Facebook walls and pages (by forcing users to “like” some content without them being aware of it);
- This spam often hides a scam (in which the bait may be access to a video, or the knowledge of “who visits your profile”), which may in the worst case end up in a wire transfer or a malware infection.

Facebook has worked a lot on these threats since 2010, a lot of them have practically disappeared.

Facebook vs Europe

The Max Schrems case

Austrian law student, he tries to exert his access right with Facebook, as defined in Irish law.

He receives a CD corresponding to a 1200-page document.

After a cautious analysis of these data, he files 22 complaints against Facebook with the Irish Data Protection Commissioner.

Details of the ongoing case:

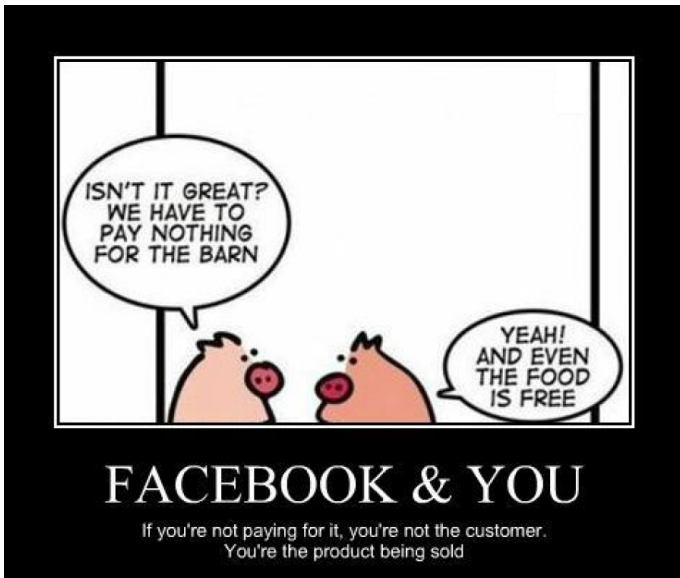
<http://www.europe-v-facebook.org/>

Facebook vs Europe

The Max Shrems case: a few complaint motives

- “Pokes”, posted messages and private messages were kept after their deletion by user;
- FB collected information on non-users, in order to create replacement profiles;
- Information collected through the “friend finder” were used without the users’ consent;
- Users were not aware of re-publication settings for messages posted on other people’s walls;
- The access request was not fully met by Facebook;
- ...

Facebook: the "free" model



Alternatives to Facebook?

Diaspora*

Social networking software in development:

- Free software;
- Listed in the 10 best open source pieces of software of 2010 (although development had barely started);
- Decentralized architecture and management;
- Focused on privacy protection;
- Currently unfinished and hard to deploy/use, development was almost halted by the death of one of the founders.

Many research proposals exist.

Several trade-offs must be resolved with data availability, performance, security, computational complexity, policy system complexity. . .

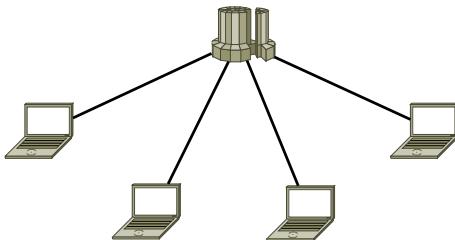
Alternatives to Facebook?

Centralized architectures

Many examples: Facebook, LinkedIn, Twitter, FourSquare...

Relatively easy to design, to deploy, to manage...

A central authority may exert a unilateral control on user information.



Alternatives to Facebook?

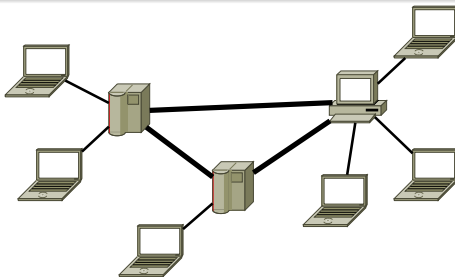
Decentralized architectures

Diaspora*, SuperNova, PeerSon...

Users connect to a “superpeer” of their choice, which provides part of the services.

Repartition of the computational burden, while allowing users to remain “clients”, possibility of a collective management system...

Superpeers remain a control centralization point and may be sensitive to collusion.



Alternatives to Facebook?

Fully distributed architectures

PrivacyWatch, Safebook, FOAF...

All peers are “equal” and responsible for part of the services provided by the system.

No entity gets more control or power (no designated point of failure), potentially better control of their users over their data.

Issues with data availability and complex policy enforcement.

