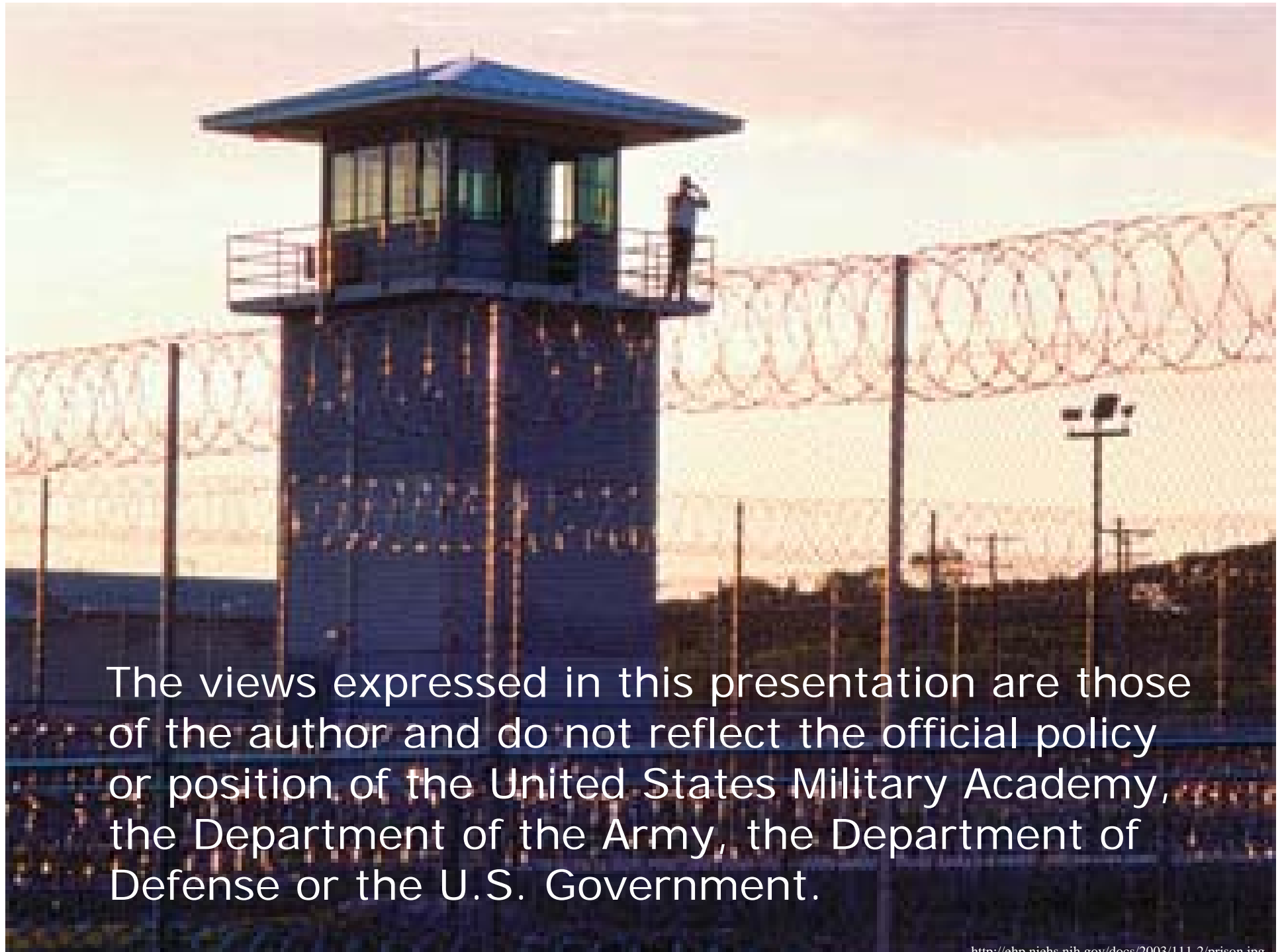# Googling

I'm Feeling (un)Lucky

*Gregory Conti*

*United States Military Academy*

*West Point, New York*

*gregory-conti@usma.edu*

The views expressed in this presentation are those of the author and do not reflect the official policy or position of the United States Military Academy, the Department of the Army, the Department of Defense or the U.S. Government.

# Scenarios

- Patent search
- Serious medical condition
- Failing company
- Anonymous blogger
- Death in the family

# Usama Fayyad
# Chief Data Officer
# 2005 Interview with ACM

**What are some of the biggest data mining challenges you face now at Yahoo?**
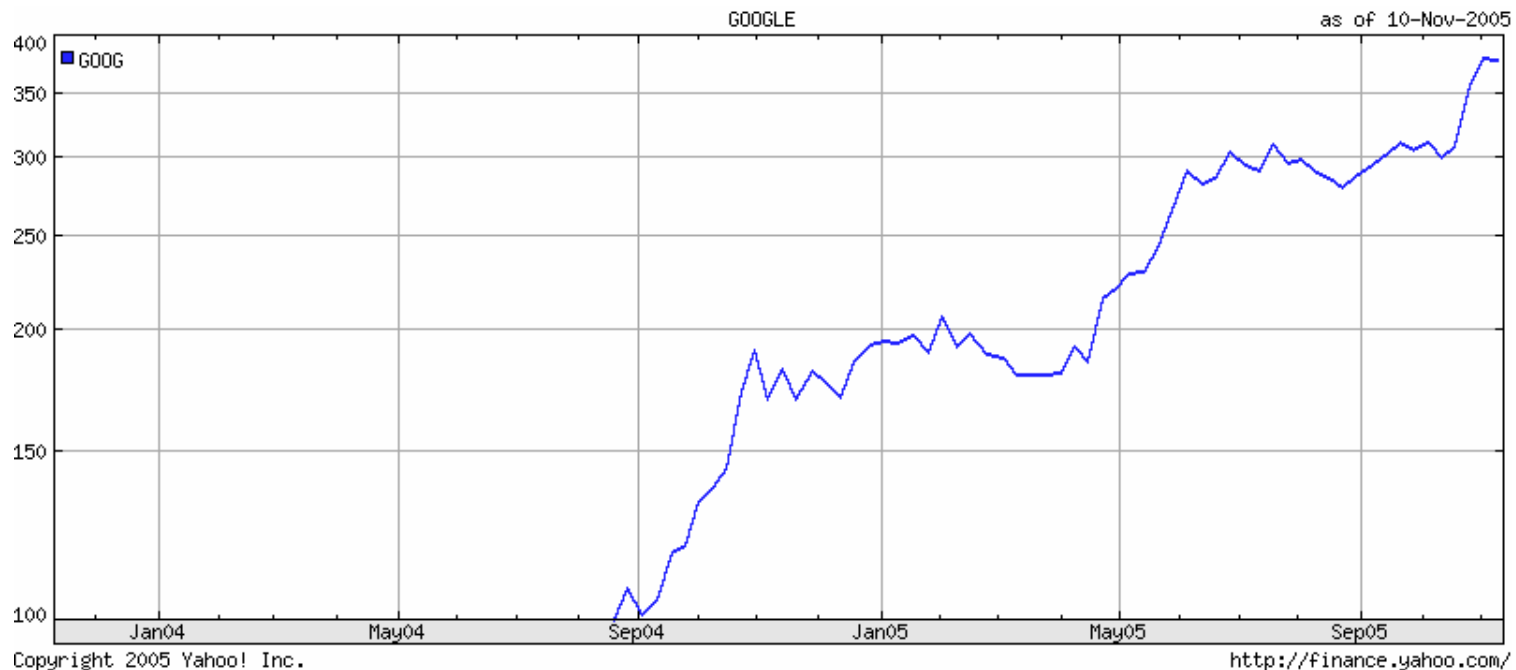
"Yahoo!'s users, through their use of our network of products, generate over <u>10 terabytes of data per day</u>. This is the equivalent of the entire text contents of the library of Congress. This is data that describes product usage, and does not include content, email, or images, etc."

"<u>The first and largest challenge is the ability to capture all of this data reliably</u>, process it, reduce it, and use it to feed the many, many [reports and applications]"

http://www.acm.org/sigs/sigkdd/explorations/issues/7-2-2005-12/fayyad.html

## What about Privacy and Data Mining at Yahoo?

"Yahoo! was built around consumer trust and it has been our number one priority since day one. <u>We would never do anything to compromise our users privacy</u>"

"Yahoo!, perhaps more than any other online consumer company, understands that <u>the long-term survival and health of the business is built on consumer trust</u>."

GOOGLE                                    as of 10-Nov-2005

■ GOOG

Copyright 2005 Yahoo! Inc.                http://finance.yahoo.com/

Organize the world's information and
make it universally accessible and useful. →

Evil _____ Good

← Generate revenue by providing advertisers
with the opportunity to deliver measurable,
cost-effective online advertising

# Related News...

"Google.cn granted license for operation in China"
censorship required to do "legal business" in China

"American Airlines subpoenas Google, YouTube"
AA demands Google reveal name of person who posted airline's training videos.

"Google Subpoena Woes Double"
DOJ, Child Pornography

"Yahoo helped Jail China writer"

"Phishing with Google Desktop"
(remote code execution via malicious website)

Sony Rootkit

Google balances privacy, reach
Elinor Mills (CNET)

http://www.doxpara.com.nyud.net:8090/planetsony_usa.JPG

# What if an entirely altruistic provider

- Insider threat
- Accidental information disclosure
- Merger / Sale of Company
- Change in leadership / corporate philosophy
- Impact of technical advances
- Legal compulsion

## What if not?

# Information Disclosed



URL
Shortening
Service

Online
Auction

Instant
Messaging
Provider

Web
Search

**Consolidated
Information
Service
Provider**

Online
Bookseller

Online
Mapping

Free
Email
Provider

File
Transfer
Service

Tiny URL
16 Million URLs
325M hits/month

I have searched for things I wouldn't want my grandmother* to know about.

# Case Study: Google

# Google Zeitgeist 2005

Google.com
Top Gainers of 2005

1. Myspace
2. Ares
3. Baidu
4. wikipedia
5. orkut
6. iTunes
7. Sky News
8. World of Warcraft
9. Green Day
10. Leonardo da Vinci

Google News
Top Searches in 2005

1. Janet Jackson
2. Hurricane Katrina
3. tsunami
4. xbox 360
5. Brad Pitt
6. Michael Jackson
7. American Idol
8. Britney Spears
9. Angelina Jolie
10. Harry Potter

Froogle
Top Searches in 2005

1. ipod
2. digital camera
3. mp3 player
4. ipod mini
5. psp
6. laptop
7. xbox
8. ipod shuffle
9. computer desk
10. ipod nano

January                                                      December

# Google Talent

"Passionate about these topics? You should work at Google."

- algorithms
- artificial intelligence
- compiler optimization
- computer architecture
- computer graphics
- data compression
- **data mining**
- file system design
- **genetic algorithms**
- **information retrieval**
- **machine learning**
- natural language processing
- operating systems
- **profiling**
- robotics
- text processing
- user interface design
- web information retrieval
- and more!

http://labs.google.com/

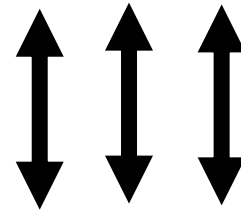# Disclosure Vectors

# "When Did You Lose Your Google Virginity?"

## Google Services

**Alerts**
Receive news and search results via email

**Answers**
Ask a question, set a price, get an answer

**Blog Search**
Find blogs on your favorite topics

**Book Search**
Search the full text of books

**Catalogs**
Search and browse mail-order catalogs

**Directory**
Browse the web by topic

**Froogle**
Shop smarter with Google

**Groups**
Create mailing lists and discussion groups

**Images**
Search for images on the web

**Labs**
Try out new Google products

**Local**
Find local businesses and services

**Maps**
View maps and get directions

**Mobile**
Use Google on your mobile phone

**News**
Search thousands of news stories

**Scholar**
Search scholarly papers

**SMS**
Use text messaging for quick info

**Special Searches**
Search within specific topics

**University Search**
Search a specific school's website

**Web Search**
Search over billions of web pages

**Web Search Features**
Do more with search

## Google Tools

**Blogger**
Express yourself online

**Code**
Download APIs and open source code

**Desktop**
Info when you want it, right on your desktop

**Earth**
Explore the world from your PC

**Gmail**
A Google approach to email

**Local for mobile**
View maps and get directions on your phone

**Picasa**
Find, edit and share your photos

**Talk**
IM and call your friends through your computer

**Toolbar**
Add a search box to your browser

**Translate**
View web pages in other languages

Information others provide or retrieve

Information
you provide
or retrieve

Google™

Information Google retrieves
(think Googlebot)

# Google Research…

google

google api

google autolink

google base

google blacklist

google chat

google contest

google hacked

google intervened

google isp

google keyhole

google labs

google maps

google scholar

google screen

google sets

google toolbar

google wallet

google watch

google wireless

googlebot

# I use Google as an Address Book.

# Google Alerts



- search terms
- email address
- frequency
- category (news, web, news and web, groups)

http://www.google.com/alerts

# Google Maps / Satellite Imagery



Locations
of interest
to you
(down to street
level)

## What have you looked at?

# Gmail



- email recipients
- email sources
- content
- N-order contact

# Gmail



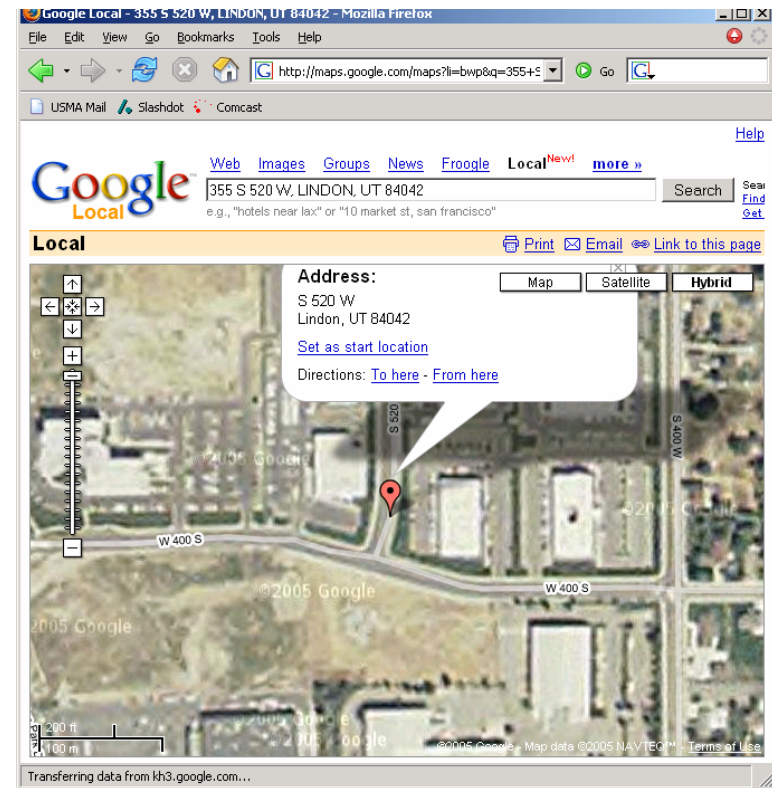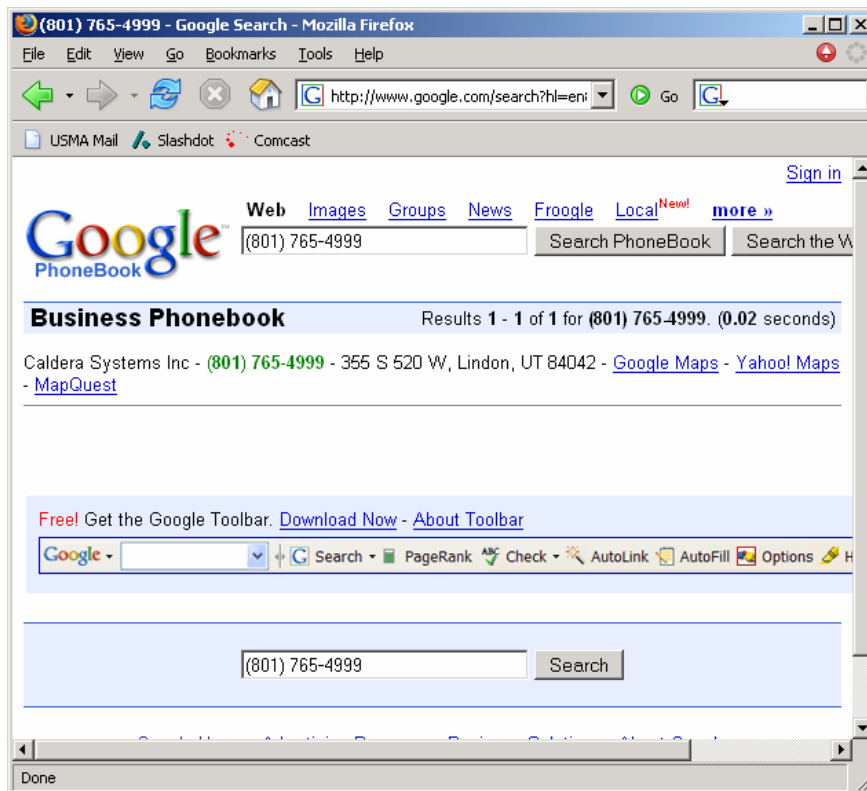- email recipients
- email sources
- content
- N-order contact

# Gmail



- email recipients
- email sources
- content
- N-order contact

# I use a Gmail account.

# I have sent an email to a Gmail account.

My Employees are registered Gmail users.

# Residential/Commercial
# Phone Number Lookup



## a search returns...

- name            -address
- phone number    -links to 3 mapping services

# Travel Support

# Collaborative Word Processing

# Google Calendar

# Local Information for Local Devices

- **Local for Mobile**
  - detailed directions
  - search results integrated with map
  - zoom in/out, drag maps
  - satellite imagery
- **free download**

http://www.google.com/glm/index.html

# Google Desktop
*"Info when you want it, right on your desktop"*

# Search Appliances



The Google Search Appliance makes the sea of lost data on your web servers, file systems and relational databases instantly available with one mouse click.

# Keyphrases
(from Google.com)

- greg conti
- gregory conti
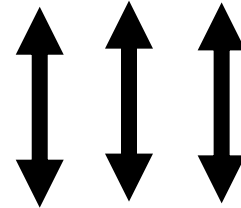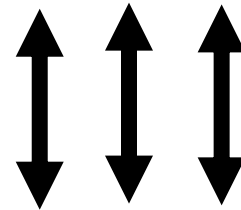- a framework for countering denial-of-information attacks
- ids rainstorm
- network security visualization
- 124th mi bn
- awareness teach information security filetype:pdf
- passive visual fingerprinting of network attack tools
- usma ia
- 92d information warfare aggressor squadron
- conti gatech
- greg conti rumint
- hacking and innovation cacm

Information others provide or retrieve

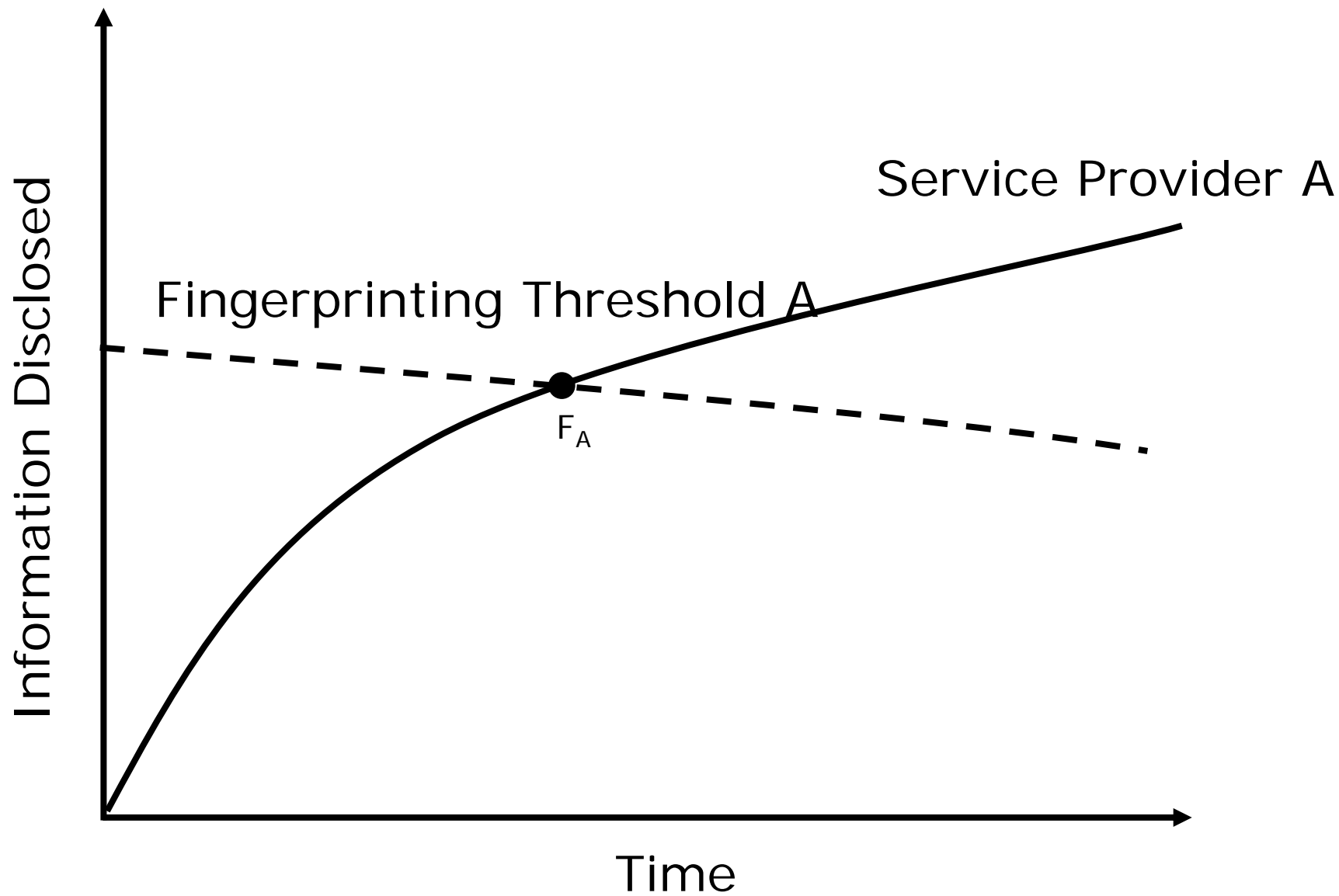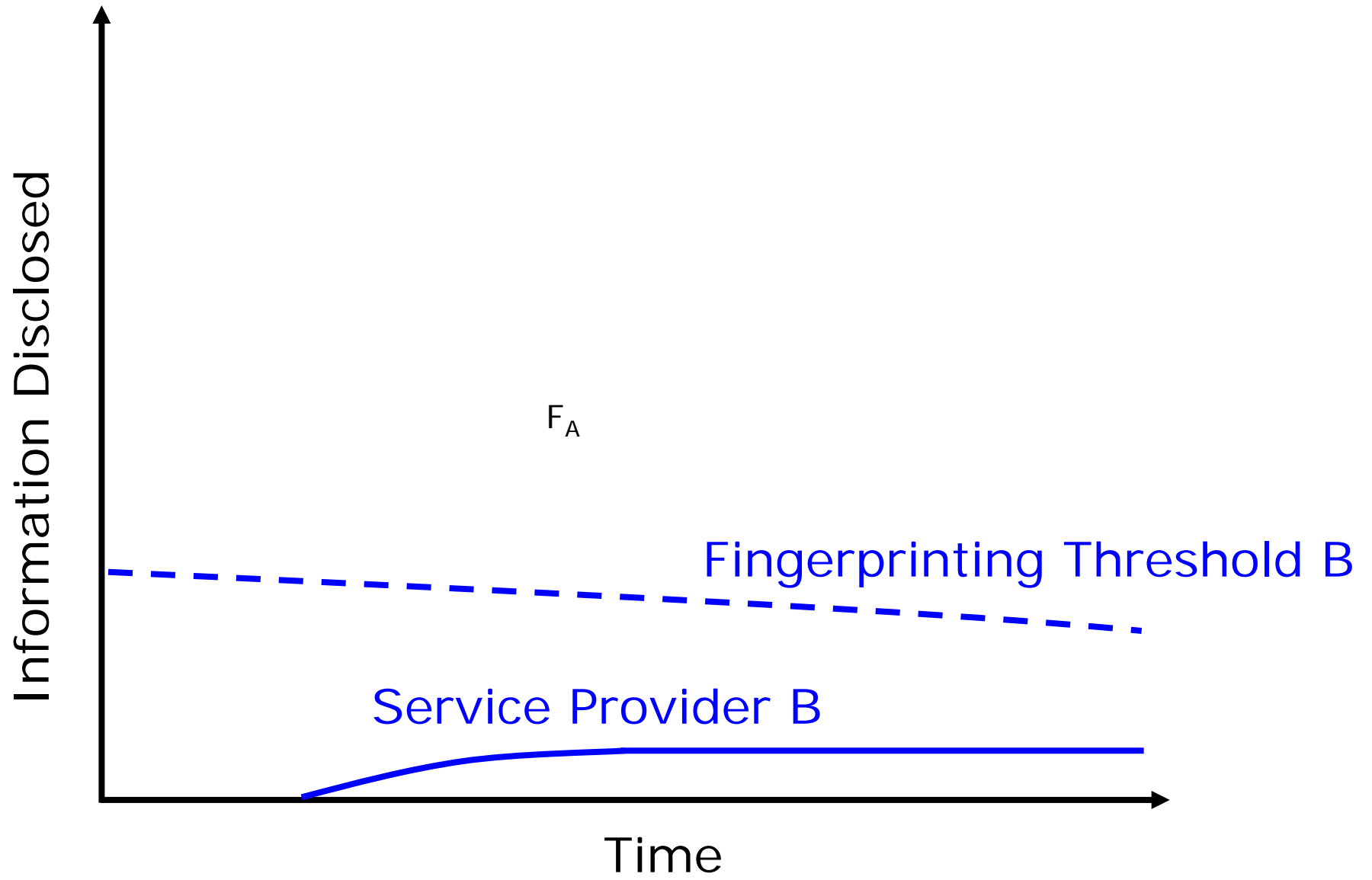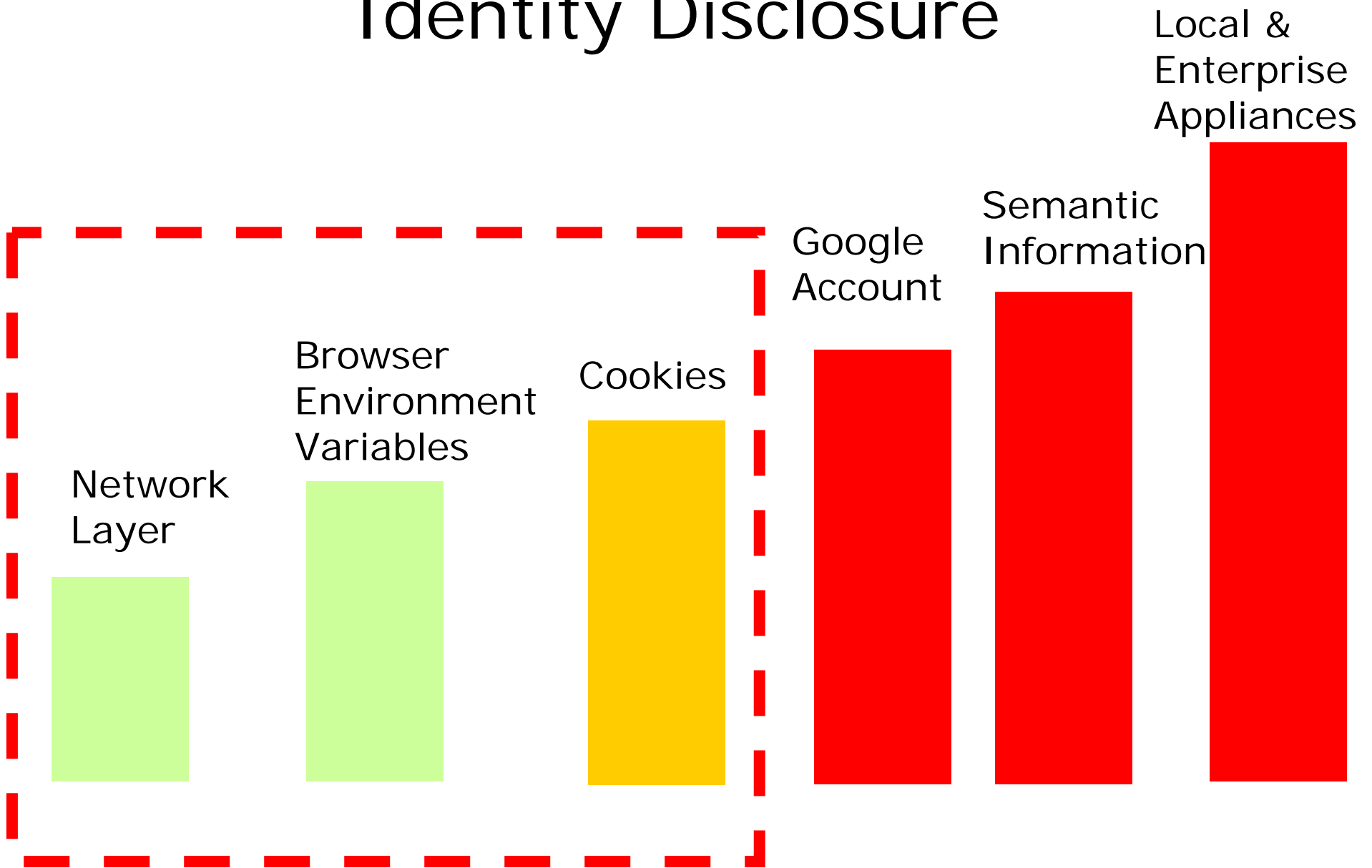Information you provide or retrieve

Google™

Information Google retrieves
(think Googlebot)

# Fingerprinting

# Identity Disclosure

Network Layer

Browser Environment Variables

Cookies

Google Account

Semantic Information

Local & Enterprise Appliances

Book Search

Finance

Email

Web Searches

Scholar Search

Instant Messaging

News

Groups

Calendar

Shopping
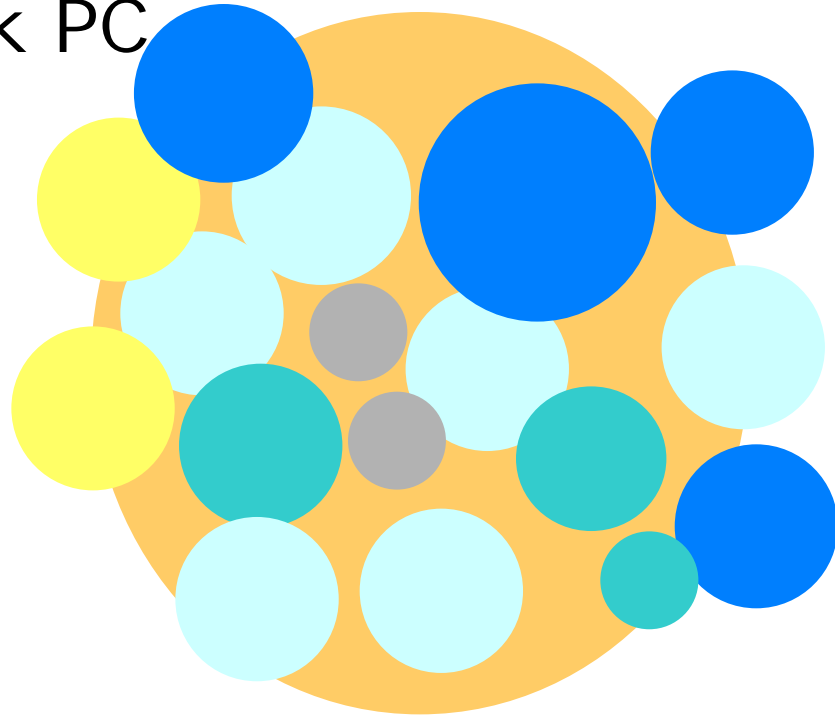
Mapping

Alerts

Language Translation

Blog

Word Processing

Local

Image Search

Network Layer
Cookies
Browser Environment Variables
Registration
Semantics

Work PC

Work PC

Mobile

Laptop

Home PC

Network Layer
Cookies
Browser Environment Variables
Registration
Semantics

**Network Layer**
Cookies
Browser Environment Variables
Registration
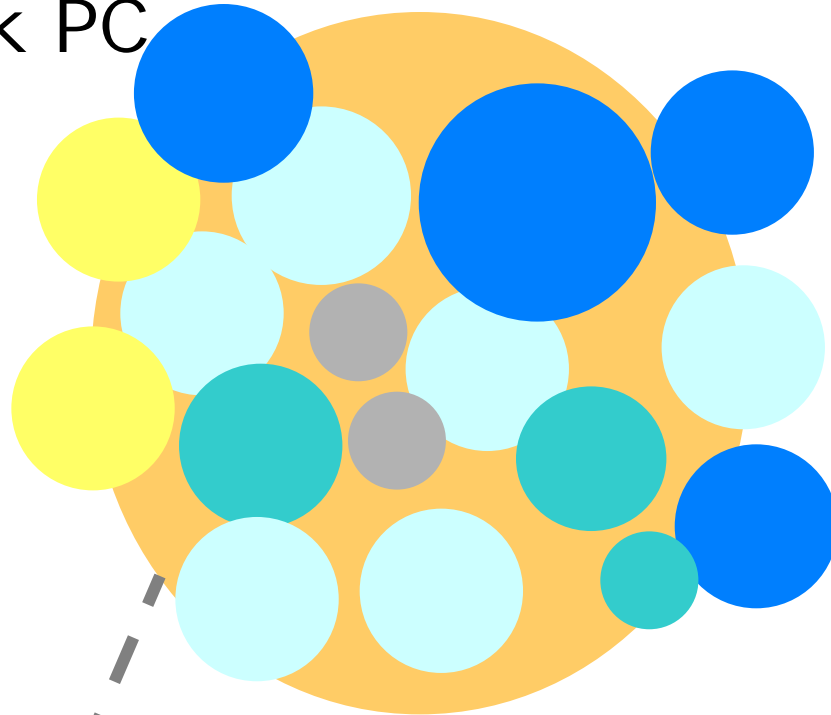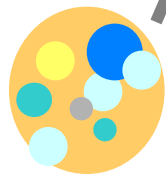Semantics

**Company**
**Organization**
**Social Network**

**Network Layer**
Cookies
Browser Environment Variables
Registration
Semantics

**Company**
**Organization**
**Social Network**

I have a Google cookie on my computer.

My Employees have Google cookies on their computers.

I have registered for a Google account.

# Countermeasures

# Countermeasures

- Anonymous browsing (Tor, anonymizer...)
- Diversity
- Policy/Law
- Go directly to the website you want
- Anonymizer-ish application for individuals and organizations but operates higher up the semantic scale
- Aggregators
- Encrypt content
- Switching Proxy Plug-ins
- Cookie rewriters

# Countermeasures

- Multiple email accounts, host own email
- Chaffing
- DHCP/NAT
- Retain only the "necessary" amount of information
- Time limits
- Transparency
- Opt-in / out

I actively manage/remove cookies.

# I regularly use web anonymization tools.

# Examples

# Co-worker example #1

- 696 entries

- 18 Feb 05 - 26 Jul 06
  - 1.33 unique searches / day

- Sanitization criteria
  - names of family, close friends
  - addresses
  - credit card numbers
  - SSN's

**5% sanitization**

# Co-worker example #2

- 1264 entries,
- 9 Nov 05 - 31 Jul 06
  - 4.8 unique searches / day
- criteria
  - names
  - locations close to home
  - phone numbers
  - friends blog

## 24% sanitization

# Co-worker example #3

- 3974 entries,
- 8 Dec 04 - 31 Jul 06
  - 6.6 unique searches / day
- criteria
  - names
  - friends
  - phone numbers
  - co-workers

## 31%
### sanitization

# Demo

# Discussion