# Use of Technology to Detect Collusive Practices

## Colin Ehren
## Managing Partner, C&SE

Tackling Corruption and Collusion in Public Procurement: Latin America and the Caribbean

Panama City, 2013

September 12, 2013

# Overview

- Challenges of Gathering Evidence from the Internet.

- Automating Collusion Detection in Public Procurement.

- Combining Manual and Automated Bid Analyses
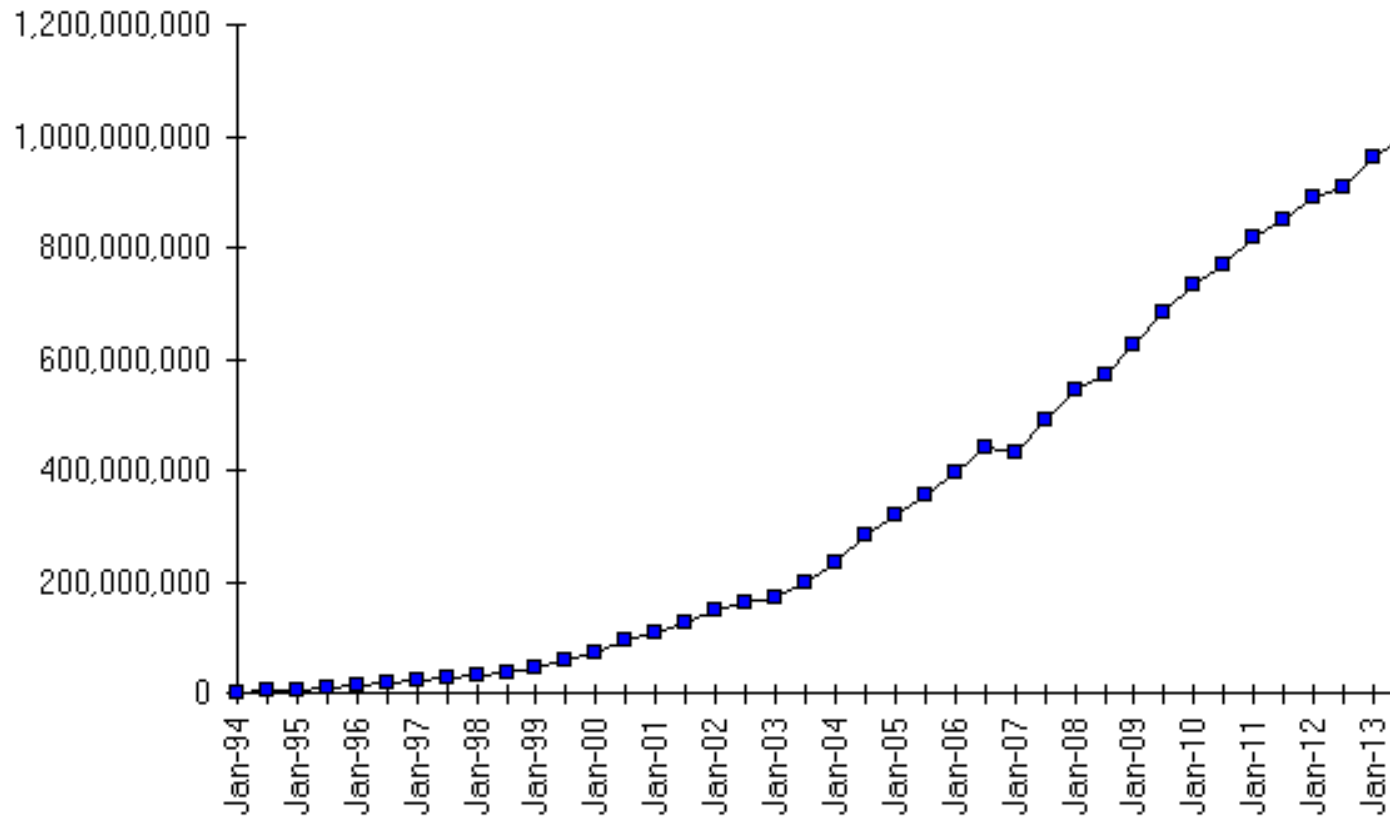
Colin Ehren
candse@gmail.com

# Challenges of Gathering Evidence from the Internet

- Size of the Internet

- Search Engines

- Invisible Web

- Social Media

- Gathering & Evidencing content.

- Compromise Issues & Internet Footprints

September 12, 2013

# Size of the Internet

## July 2013 - 996,230,757 Hosts



Internet Domain Survey Host Count

Source: Internet Systems Consortium (www.isc.org)

September 12, 2013

Colin Ehren
candse@gmail.com

# Size of the Internet

## August 2013

- **Google – 47+ bn**

- **Yahoo – 29.4 bn**

- **Bing – 30.9 bn**

- **Lycos – 30.4 bn**

- **Baidu – 4.5 bn (estimate)**

- **Yandex – 3.2 bn**

- **Terra – 2+ bn**

## Deep or Invisible Web

September 12, 2013

Colin Ehren
candse@gmail.com
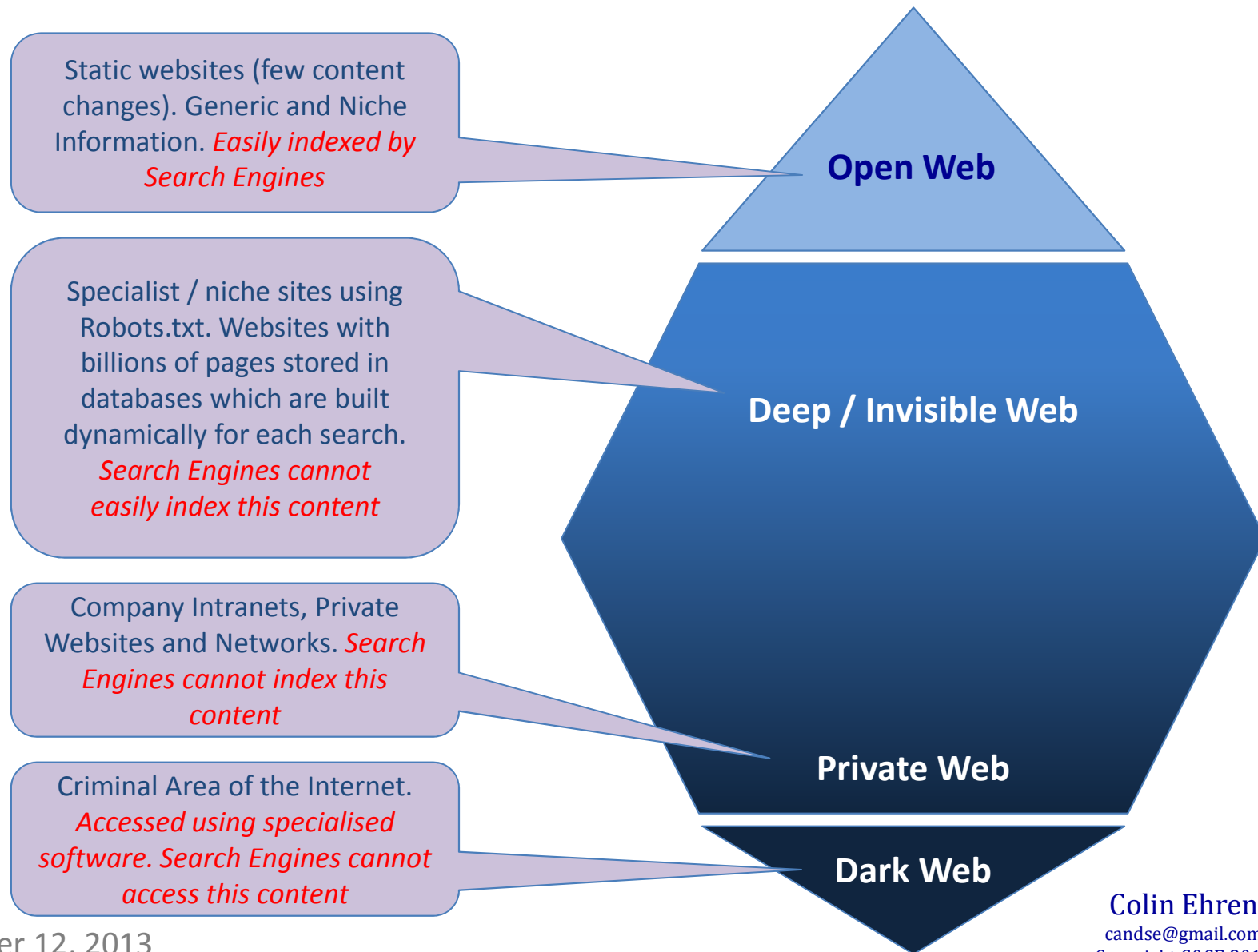
# Search Engines

## What does Google know?

2007 - Eric Schmidt (Google CEO) estimated  Google had indexed roughly 0.004% of the Internet.

July 2008 – Google had identified 1 trillion (1,000,000,000,000) unique URL's.

## Everything is on the Internet.

It is estimated that 80% of Open Source information exists in Books, Magazines, Literature and other Media.

September 12, 2013

Colin Ehren
candse@gmail.com

# Invisible Web

Static websites (few content changes). Generic and Niche Information. *Easily indexed by Search Engines*

Specialist / niche sites using Robots.txt. Websites with billions of pages stored in databases which are built dynamically for each search. *Search Engines cannot easily index this content*

Company Intranets, Private Websites and Networks. *Search Engines cannot index this content*

Criminal Area of the Internet. *Accessed using specialised software. Search Engines cannot access this content*

**Open Web**

**Deep / Invisible Web**

**Private Web**

**Dark Web**

September 12, 2013

Colin Ehren
candse@gmail.com

# Invisible Web

July 2001, Michael Bergman produced a white paper for www.brightplanet.com

- 400 to 550 times larger than the World Wide Web.

- 7,500 terabytes of information compared to 19 terabytes on WWW.

- 550 billion documents compared to 30 billion on the WWW.

- 200,000+ deep Web sites.

- 60 of the largest sites collectively contained over 40x the information on the WWW.

2004 Study - identified 330,000+ Deep Web sites.

Has grown almost exponentially since.

Colin Ehren
candse@gmail.com

September 12, 2013

# Social Media

- Austria, Belgium , Croatia, Cyprus, Czech Republic, Denmark, Finland, FYR of Macedonia, Germany, Greece, Iceland, Ireland, Italy, Luxembourg, Malta, Norway, Slovakia, Slovenia, Spain, Sweden, Switzerland, Turkey & UK - **Facebook**, **Youtube**

- Bulgaria – **Facebook**, **VBox7**, **Youtube**

- Estonia – **Facebook**,**Youtube**, **VKontakte**

- France – **Facebook**, **Skyrock**, **Youtube**

- Hungary – **Facebook**, **Iwiw**, **Youtube**

- Latvia – **Youtube**, **Draugiem**, **Facebook**

- Lithuania – **Facebook**, **Youtube**, **One**

- Netherlands – **Facebook**, **Hyves**, **Youtube**

- Poland –  **Facebook**, **Youtube**, **Chomikuj**

- Portugal – **Facebook**, **Youtube**, **Twitter**

- Romania – **Facebook**, **Youtube**, **Hi5**

September 12, 2013

Colin Ehren
candse@gmail.com

# Social Media

- China.
  - Qzone, Tencent Weibo, Sina Weibo, RenRen
- Russia.
  - Vkontakte, Youtube, Odnoklassniki, Facebook, Livejournal
- India.
  - Facebook, Youtube, Orkut, Ibibo
- Georgia
  - Facebook, Youtube, Odnoklassniki, VKontakte
- Brazil
  - Facebook, Youtube, Twitter, Orkut
- Panama
  - Facebook, Youtube, Twitter

September 12, 2013

Colin Ehren
candse@gmail.com

Create Journals
Update Journals
Journals

Find Users
Random
Read

Search
Create New
Communities

blurty

flickr GAMMA

myGamma myWorld | mobile Community

8 CarDomain years old!

facebook

myspace.com
a place for friends

listography, your life in lists
create. share. read.

sconex Your High School Online

care2 Make a Difference!

Bolt

eons

Bebo

photobucket
Video and Image Sharing

downelink BETA

hi5.

Kaneva beta

LiVEJOURNAL
Create an Account    Post to Journal

xanga.com
THE WEBLOG COMMUNITY

YouTube

Faceparty
THE BIGGEST PARTY ON EARTH

FOTOLOG

VampireFreaks.com

Gaia ONLINE

Music forte

friendster

Mashable! social networking 2.0
All | MySpace | YouTube | Bebo | Facebook | Xanga

last.fm
the social music revolution

MULTIPLY

Second LIFE

VOX

Windows Live Spaces

MOG  VALIDATING MUSIC GEEKS SINCE 2006  BETA

profileheaven v2.0
...FUN IN THE AFTERLIFE!

ODEO

myyearbook BETA
YOU'VE GOT FRIENDS!

www.mymidishare.com
sound
MIDIs and MP3s

imvu
beyond instant messaging

WAYN
WHERE ARE YOU NOW?

TagWorld

LibraryThing BETA

orkut beta

PICZO

OUTERWORLDS
Where Virtual Reality
And Real People Meet!

StumbleUpon  SU

# Traditional Social Media Tools

- Internet Relay Chat

- Usenet Talk Groups

- Google Groups / Yahoo Groups

- MSN/Skype, Yahoo, AOL Messengers and Chat Rooms

- E-mail

- Dedicated Discussion Forums

- Dating – Muslim Match, Uniform Match, Adult Friend

- Reunion Sites.

September 12, 2013

Colin Ehren
candse@gmail.com

# Gathering Data

## Common Issue – Large amounts of data

- Social Media Tools
- Data Extraction Tools.
- Visualisation Tools

Colin Ehren
candse@gmail.com

September 12, 2013

# Gathering Data

Colin Ehren
candse@gmail.com

September 12, 2013

# Evidencing Data

## Prove it or Lose it

All your efforts will be wasted unless you can prove what you found.

Colin Ehren
candse@gmail.com

# Evidencing Data

## General Rules

- Captured data should be given unique file names when saved.

- Data from different investigations kept separately.

- Data from different investigations must not be saved to the same CD/DVD disk.

- Separate CD/DVD disk should be used for each person subject to the investigation.

September 12, 2013

Colin Ehren
candse@gmail.com

# Evidencing Data

## Original Notes

- Full record of the Investigators actions

- Made contemporaneously at the time or as soon as possible after.

- Content – Acts as basis for Statement.

- May be many months until Court case.

- Importance cannot be over stated.

September 12, 2013

Colin Ehren
candse@gmail.com
Copyright C&SE 2013

# Evidencing Data

## Preserving the Evidential Chain

- When moving Data the Evidential Chain must be preserved.

- Proving the Integrity of the Data – MD5 / SHA1 Hash

- Moving Data – CD / DVD Disk.

- Exhibiting Data – Evidential Bag

- Notes

- Stored Securely

Colin Ehren
candse@gmail.com
Copyright C&SE 2013

# Compromise Issues

Multiple ways to access the Internet;

- Corporate networked PC's

- Corporate stand-alone PC's

- Re-claimed stand-alone PC's

- Covert / Unattributed stand-alone PC's

- Covert / Unattributed networked PC's

- Working from Home (Stand-alone or Networked)

- Mobile Devices

- Internet Café's

Colin Ehren
candse@gmail.com

September 12, 2013

# Compromise Issues

Criminals are aware that the Organisations use the Internet for their Research or Investigations. There are several issues that surround these methods of access.

**Recommendations**

- Corporate Networked and Stand-Alone workstations should only be used for generic research such as obtaining Crime Trend information or Research Publications from accredited websites.

- All detailed or sensitive Internet research or Open Source Investigations should be undertaken on a covert or unattributed and registered PC, using a covert or unattributed Internet connection.

September 12, 2013

Colin Ehren
candse@gmail.com

# Compromise Issues

# Why?

Your Internet Footprint could compromise yourself, your colleagues or an Investigation or Intelligence Operation that your Organisation or a Partner may be engaged in.

Colin Ehren
candse@gmail.com

# Compromise Issues

Every time you use the Internet you leave your footprints on the websites that you visit.

The size of your whole footprint depends on the 'environment variables' that your computer and internet browser pass and on your web activity.

Colin Ehren
candse@gmail.com

September 12, 2013

# Compromise Issues

**User**

**Web Server**

**Webmaster**

*Device ID*
*Browser/Software*

**Internet Access**

**Reports**

**Access logs**

- When you click on a hyperlink your browser sends your 'Internet footprint' to a Web Server so that it can find what you want and send it back to you.

- The Webmaster in control of the Web Server can view your 'footprint' to gain information about you or your organization (physical location, ISP, your interests, type of PC/Software, etc.)

September 12, 2013

Colin Ehren
candse@gmail.com

# Compromise Issues

- Footprints are left on every web server of every web page that you visit.

- Be aware of what footprints you are leaving before visiting websites which could implicate your organization.

- Should you visit: badguy.com from a Government / Police Internet Connection?

- Your footprints can include;

  - The name of your computer (or gateway).

  - The IP address of your computer, or proxy gateway.

  - The URL of the page you were previously viewing. *(Web masters use this to see what web pages lead visitors to their site.)*

September 12, 2013

Colin Ehren
candse@gmail.com

# Compromise Issues



**Google.com**

"search terms"

**User**

**webmaster**

http://www.google.com/keywords=searchterms

**Hackdiary.com**

**webmaster**

Footprint:
- @your.org OR
- you@yahoo.com

google.com webmaster knows your "search terms"

hackdiary.com webmaster knows what "search terms" you used to find them.

September 12, 2013

Colin Ehren
candse@gmail.com

# Compromise Issues

**In Summary** - Your internet footprint could compromise an investigation or intelligence operation that your organisation or Law Enforcement Agency may be engaged in.

```
Browser ID : Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.0;
AWARELockDown2001;   NET CLR 1.1.4322)
Browser:MicroSoft Internet Exploder Version 6.0
Operating System:Windows 2000
Referrer:  http://www.google.co.uk/search?
q=am+i+secure&hl=en&lr=&start=10&sa=N
Your IP Address : 195.173.172.10
IP Address Behind Firewall
Your hostname: mailgate.exitstrategy.co.uk
Your AS : 2529 ARIN ASN block
Your Country: GB
Your Abuse Contact  : abuse@demon.net
Additional Software :
```

September 12, 2013

Colin Ehren
candse@gmail.com
Copyright C&SE 2013

# Compromise Issues

Web  Images  Videos  Maps  News  Shopping  Gmail  more ▼

**Google**

awarelockdown2001

About 70 results (0.43 seconds)

**Everything**
**Images**
**Videos**
▼ **More**

**The web**
Pages from the

▼ More search to

**Home** / **Forums Index** / **Hardware and OS Related Technologies** / **Website Technology Issues**

Forum Library : Charter : Moderators: lammert

## Website Technology Issues

These terms have been highlighted:
**awarelockdown2001** [ remove highlighting]

**AWARELockDown2001**
What user agent is this?

**dunne**

#:672992

What user agent would "**awarelockdown2001**" be, anyone know? The host
IP traces back to ls4689-s0.metpolice.router.uk.quza.net, which is, ah,
"interesting".

Colin Ehren
candse@gmail.com

September 12, 2013

# Compromise Issues

GovRisk
The International Governance & Risk Institute

Browser ID : Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.0; AWARELockDown2001; .NET CLR 1.1.4322)
Browser:Micro$oft Internet Exploder Version 6.0
Operating System:Windows 2000
Referrer: http://www.google.co.uk/search?q=am+i+secure&hl=en&lr=&start=10&sa=N
Your IP Address : 195.173.172.10
IP Address Behind Firewall
Your hostname: mailgate.exitstrategy.co.uk
Your AS : 2529 ARIN ASN block
Your Country: GB
Your Abuse Contact : abuse@demon.net
Additional Software :

Colin Ehren
candse@gmail.com

September 12, 2013

# Compromise Issues

Colin Ehren
candse@gmail.com

# Compromise Issues

Browser ID : Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.0; AWARELockDown2001; .NET CLR 1.1.4322)

Browser:Micro$oft Internet Exploder Version 6.0

Operating System:Windows 2000

Referrer: http://www.google.co.uk/search?q=am+i+secure&hl=en&lr=&start=10&sa=N

Your IP Address : 195.173.172.10

IP Address Behind Firewall

Your hostname: mailgate.exitstrategy.co.uk

Your AS : 2529 ARIN ASN block

Your Country: GB

Your Abuse Contact : abuse@demon.net

Additional Software :

September 12, 2013

Colin Ehren
candse@gmail.com

# Compromise Issues

# Compromise Issues

GovRisk
The International Governance & Risk Institute

| | |
|---|---|
| Browser ID : Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.0; AWARELockDown2001; .NET CLR 1.1.4322) | |
| Browser:Micro$oft Internet Exploder Version 6.0 | |
| Operating System:Windows 2000 | |
| Referrer: http://www.google.co.uk/search? q=am+i+secure&hl=en&lr=&start=10&sa=N | |
| Your IP Address : 195.173.172.10 | |
| IP Address Behind Firewall | |
| Your hostname: mailgate.exitstrategy.co.uk | |
| Your AS : 2529 ARIN ASN block | |
| Your Country: GB | |
| Your Abuse Contact : abuse@demon.net | |
| Additional Software : | |

September 12, 2013

Colin Ehren
candse@gmail.com

## Network Whois record

Queried **whois.ripe.net** with "**-B 195.173.172.10**"...

```
% This is the RIPE Whois query server #2.
% The objects are in RPSL format.
%
% Note: the default output of the RIPE Whois server
% is changed. Your tools may need to be adjusted. See
% http://www.ripe.net/db/news/abuse-proposal-20050331.html
% for more details.
%
% Rights restricted by copyright.
% See http://www.ripe.net/db/copyright.html

% Information related to '195.173.172.0 - 195.173.172.15'

inetnum:        195.173.172.0 - 195.173.172.15
netname:        METROPOLICE
descr:          Metropolitan Police Service
descr:          London SW1H
country:        GB
admin-c:        AC2375-RIPE
tech-c:         AC2375-RIPE
status:         ASSIGNED PA
mnt-by:         AS2529-MNT
mnt-lower:      AS2529-MNT
mnt-routes:     AS2529-MNT
notify:         hostmaster@demon.net
changed:        hostmaster@demon.net 20030121
source:         RIPE

person:         Alan Cooper
address:        Metropolitan Police Service
address:        London SW1H
phone:          +44-20 8649 3658
notify:         hostmaster@demon.net
mnt-by:         AS2529-MNT
nic-hdl:        AC2375-RIPE
changed:        hostmaster@demon.net 20030121
source:         RIPE
```

September 12, 2013

Colin Ehren
candse@gmail.com
Copyright C&SE 2013

# Compromise Issues

Colin Ehren
candse@gmail.com
Copyright C&SE 2013

# Compromise Issues

## **Beware**

Parallel Surfing can associate a Covert or non-attributed PC to the your organisation. Parallel Surfing would occur when the same websites or Search Terms are researched on a Corporate workstation and then those same sites or Search Terms are researched on a Covert or non-attributed workstation.

Webmasters can run reports that identify who has been using the same Search Terms, etc.

Colin Ehren
candse@gmail.com
Copyright C&SE 2013

September 12, 2013

# Compromise Issues

## The "parallel surfing" Problem...

— **User #1: leaves "your.org" footprints whilst visiting "target.com"**

═ **User #2: leaves "Covert" footprints whilst visiting "target.com"**

User #1

User #2

**Your.org**

**1234@aol.com**

**target.com**

**The "Covert" User may now be recognized as an "your.org" visitor.**

Colin Ehren
candse@gmail.com

# Compromise Issues

## 195.173.172.10

195.173.172.10

antnicuk

antnicuk@msn.com

antnicuk nicola

192.com

September 12, 2013

Colin Ehren
candse@gmail.com

# Compromise Issues

lin Ehren
se@gmail.com

# Compromise Issues

## 195.173.172.10

195.173.172.10

antnicuk

antnicuk@msn.com

antnicuk nicola

192.com

September 12, 2013

# Compromise Issues

## **Beware**

Your footprints are left on every website that you visit. If those websites display adverts or images from 3$^{rd}$ parties, your footprints are sent automatically to them as well.

You do not need to have visited the 3$^{rd}$ party site for them to obtain your information.

Colin Ehren
candse@gmail.com

September 12, 2013

# Compromise Issues

Web pages can include images or adverts from third parties.

**Ad-Image.com**

**AdultFriend.com**
**Hot_stuff**
**JoeHot@webmail.com**
**Viewing history**

**Insurance.co.uk**
**Joseph Hotman**
**JHotman@isp.com**
**Address & phone**
**Viewing history**

**hacker.com**
**Black_hat**
**hatter@hushmail.com**
**Viewing history**

**Cookies on your PC**

**AdultFriend.com ID#_201**
**insurance.co.uk ID#_4873**
**hacker.com ID#_539**
**Ad-Image.com ID#_435349**

**Your Profile**

AdultFriend.com
Hot_stuff
JoeHot@webmail.com
Viewing history
Likes / Dislikes

Insurance.co.uk
Joseph Hotman
JHotman@isp.com
Address & phone
Car / Jewellry
Viewing history

hacker.com
Black_hat
hatter@hushmail.com
Viewing history
Contacts

Companies, such as "Ad-Image.com" are able to compile a significant profile on you and your surfing habits, which they are able to trade or sell to their partners or customers.

September 12, 2013

Colin Ehren
candse@gmail.com

# Compromise Issues



Ordering Pizza?

Colin Ehren
candse@gmail.com

September 12, 2013

# Automating Collusion Detection & Combining Manual and Automated Bid Analyses

September 12, 2013

Colin Ehren
candse@gmail.com

# Search Automation

## Query Servers

- Integrated Solution
- Live Queries
- Automated Queries
- Builds own Database
- Expensive

September 12, 2013

Colin Ehren
candse@gmail.com

# Search Engine Alerts

Colin Ehren
candse@gmail.com

# Data Extraction Tools



WebHarvy

Automation Anywhere

WebDataExtractor

Data Toolbar
The world easiest data scraping tool

Visual Web Ripper

mozenda™

Web-Harvest
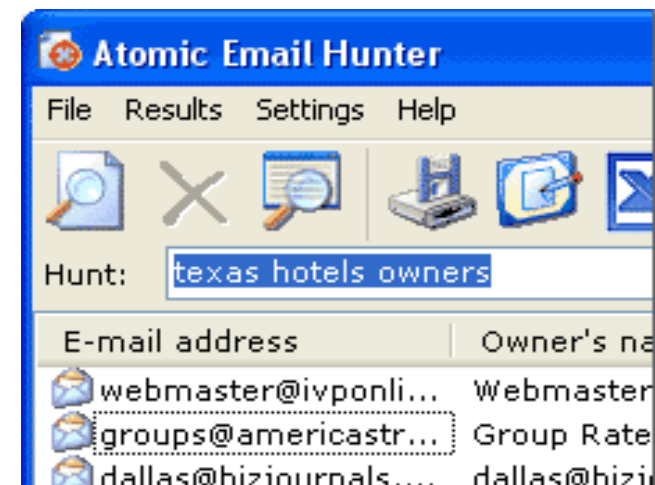version 2.0

Colin Ehren
candse@gmail.com

September 12, 2013

# Web Data Extractor

Colin Ehren
candse@gmail.com

# Data Extraction Tools



- Atomic E-Mail Hunter

- Atomic Web Spider

- Atomic Whois Explorer

- Atomic Newsgroup Explorer

- Atomic E-Mail Studio



September 12, 2013

Colin Ehren
candse@gmail.com

# Assimilation of Data

## Huge data volumes

- Petabytes of data  (1 Petabyte (1000 terabytes) = approx 3000 million documents
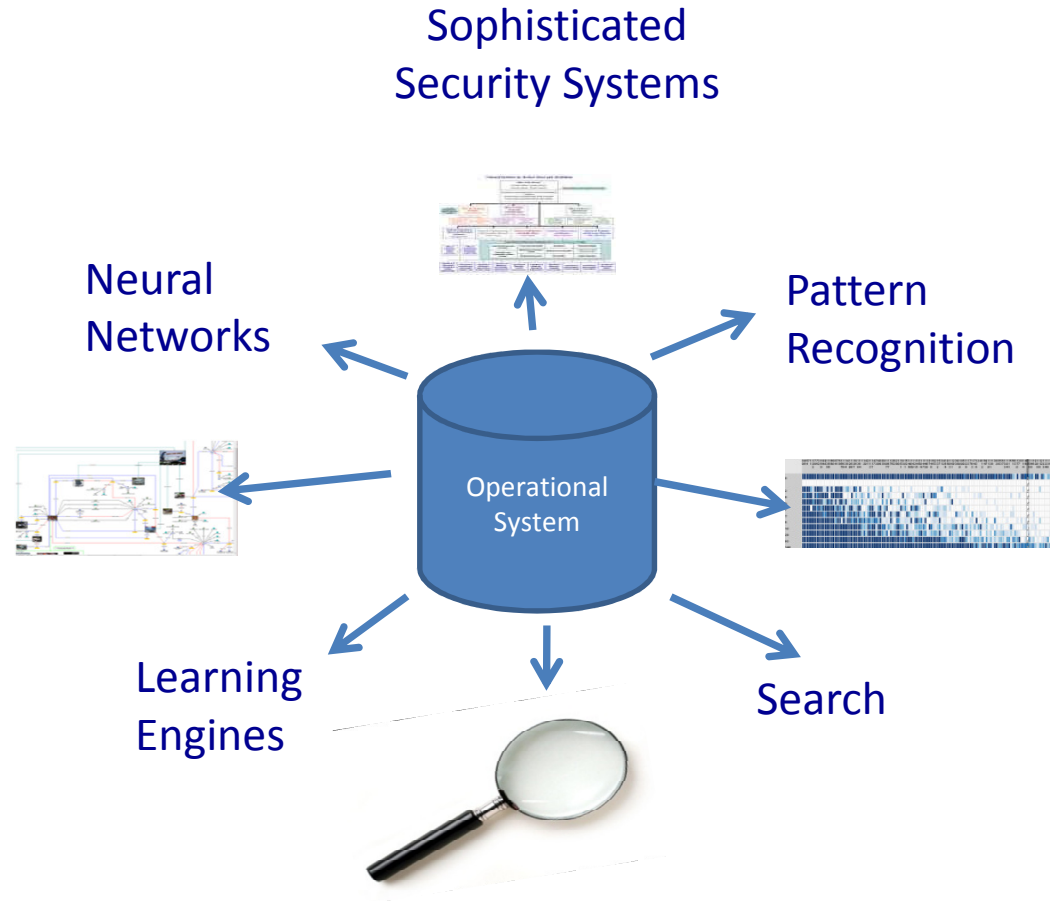
## Diverse sources

- The Internet – www, blogs, twitter, social networks, virtual worlds, chat-rooms
- Internal – E-mail, Office Systems, Knowledge Management and Analyst Reports
- Computers, storage devices, mobile phones

## Integration of data

- Multiple formats
  - Structured, unstructured (text)
  - Languages / alphabets
- Organisationally
  - Across departments
  - Across boundaries (different legislation)
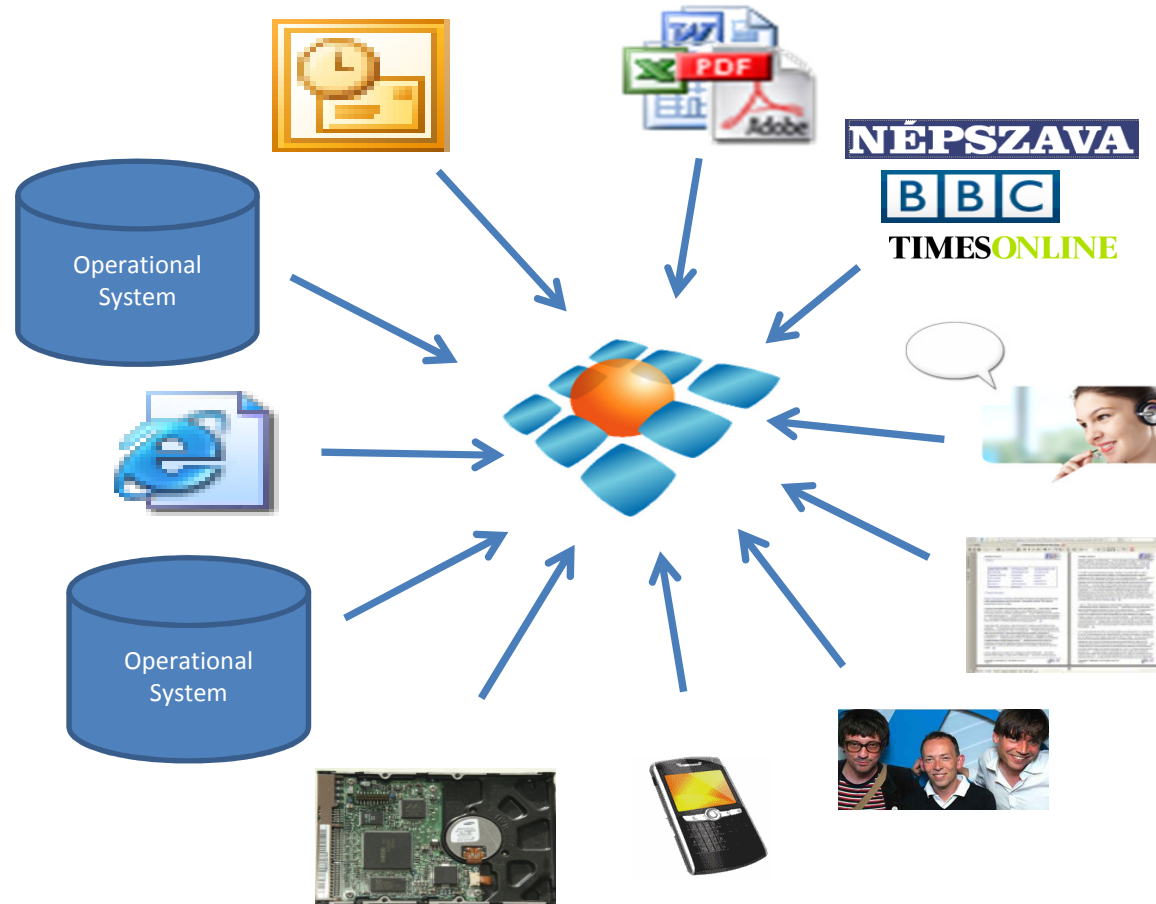  - Between organisations

September 12, 2013

Colin Ehren
candse@gmail.com

# Traditional Software

Sophisticated
Security Systems

Neural
Networks

Pattern
Recognition

Operational
System

Learning
Engines

Search

**Data traditionally held in isolated silos**

September 12, 2013

Colin Ehren
candse@gmail.com

# Next Generation Software



**Integration of all types of data**

Colin Ehren
candse@gmail.com
Copyright C&SE 2013

# Next Generation Software

## Entity Abstraction

Colin Ehren
candse@gmail.com

# Next Generation Software

## Links and Relationships

Group data and harmonise it

Connect data groups

Find definitive connections

Find implied connections

Get rid of noise

September 12, 2013

Colin Ehren
candse@gmail.com
Copyright C&SE 2013

# Next Generation Software



http://www.dexterintel.com/

Colin Ehren
candse@gmail.com

# Any Questions?

September 12, 2013

Colin Ehren
candse@gmail.com

# Use of Technology to Detect Collusive Practices

## Colin Ehren
## Managing Partner, C&SE

### candse@gmail.com
### +44 (0)7941 338 449

Tackling Corruption and Collusion in Public Procurement: Latin America and the Caribbean

Panama City, 2013

September 12, 2013

Colin Ehren
candse@gmail.com