

Digital Forensics Innovation: Searching A Terabyte of Data in 10 minutes

Simson L. Garfinkel Associate Professor, Naval Postgraduate School Jan 17, 2013 http://simson.net/ https://domex.nps.edu/deep/

NPS is the Navy's Research University.

Monterey, CA - 1500 students

- US Military & Civilian (Scholarship for Service & SMART)
- Foreign Military (30 countries)

Graduate Schools of

Operational & Information Sciences (GSOIS)

- Computer Science
- Defense Analysis
- Information Sciences
- Operations Research
- Cyber Academic Group

National Capital Region (NCR) Office

• 900 N Glebe (Ballston)/Virginia Tech building







The Digital Evaluation and Exploitation (DEEP) Group: Research in "trusted" systems and exploitation.

"Evaluation"

- Trusted hardware and software
- Cloud computing

"Exploitation"

- MEDEX "Media" Hard drives, camera cards, GPS devices.
- CELEX Cell phone
- DOCEX Documents
- DOMEX Document & Media Exploitation

Current Partners:

- Law Enforcement (FBI & Local)
- DHS (HSARPA; Video Games & Insider Threat)
- NSF (Courseware development)
- DoD





Digital information is pervasive in today's society, but attorneys, judges and juries are not digital experts.

Many potential sources of digital evidence:

Laptops; Cell Phones; Email messages



Many possible goals:

- Establish possession of contraband information (child pornography, credit card #s)
- Recover stolen information
- Document a conspiracy (stock fraud; murder-for-hire)



The digital forensics process makes *digital evidence* available for [legal] decisions





Most work to date focuses on the first half.



PRAESTANTIA PER SCIENTIAM

My focus is developing better analysis approaches

Identification of high-value data.

- What is important?
 - -Contacts, calendar, documents?
 - -Software?
 - Geolocation information?
 - Temporal / time sequence?





Correlation — are there copies of the same or similar information?

- Identify previously unknown organizations or networks
- Identify data that is unusual or emerging

Presentation and Integration:

- Make the results understandable.
- Effect organizational change through adoption & integration





Three principles underly my research:

- 1. Automation is essential.
 - Today most forensic analysis is done manually.
 - We are developing techniques & tools to allow automation.

2. Concentrate on the invisible.

- It's *easy* to wipe a computer....
 - but targets don't erase what they can't see.
- So we look for:
 - Deleted and partially overwritten files.
 - Fragments of memory in swap & hibernation.
 - Tool marks.

3. Large amounts of data is essential.

- Most research is based on search & recognition
 - 10x the data produces 10x the false-positives
- We develop algorithms that work *better* with more data.







We do science with "real data."

The Real Data Corpus (30TB)

- Disks, camera cards, & cell phones purchased on the secondary market.
- Most contain data from previous users.
- Mostly acquire outside the US:
 - -Canada, China, England, Germany, France, India, Israel, Japan, Pakistan, Palestine, etc.
- Thousands of devices (HDs, CDs, DVDs, flash, etc.)

Mobile Phone Application Corpus

• Android Applications; Mobile Malware; etc.

The problems we encounter obtaining, curating and exploiting this data mirror those of national organizations

- Garfinkel, Farrell, Roussev and Dinolt, Bringing Science to Digital Forensics with Standardized Forensic Corpora, DFRWS 2009 <u>http://digitalcorpora.org</u>/



We manufacture data that can be freely redistributed.

Files from US Government Web Servers (500GB)

- \approx 1 million heterogeneous files
 - -Documents (Word, Excel, PDF, etc.); Images (JPEG, PNG, etc.)
 - -Database Files; HTML files; Log files; XML
- Freely redistributable; Many different file types
- This database was surprising difficulty to collect, curate, and distribute:
 - -Scale created data collection and management problems.
 - -Copyright, Privacy & Provenance issues.

Advantage over flickr & youtube: persistence & copyright



<abstract>NOAA's National Geophysical Data Center (NGDC) is building high-resolution digital elevation models (DEMs) for select U.S. coastal regions. ... </abstract>

<abstract>This data set contains data for birds caught
with mistnets and with other means for sampling Avian
Influenza (AI)....



This talk presents today's digital forensic challenges and presents a research project that helps address them.

Introducing digital forensics





Today's digital forensics challenges



Random sampling for high speed forensics









Challenges Facing Digital Forensics

Extracting digital evidence was simple five years ago

"Imaging tools" extracted data without modification.



"Write Blocker" prevents accidental overwriting.



stored on a storage array.

Analyzing digital evidence was simple five years ago

Commercial tools extracted *files* from disk images

- Display of *allocated* & *deleted* files.
- String search
- File extraction
- File "carving"
- Examining disk sectors

Job of analyst:

- Find interesting data
- Report on it.



		able Sallery Timeline Report			
6 .		Name	Filter	Logical Size	In Report
BOM See	55	C TTP		3	
Documents and Settings	₹ 56	C Northcode Inc		13	
Administrator	57	Pewlett-Packard		15	
Application Data	- 58	C Heidi Computers Ltd		19	
Bluetooth Software	59	C R-TT		4	
Cookies	₩ 60	LicenseManager		14	
B-D Cosktop	61	SCC		3	
😑 🕞 🔂 Favorites		Ahead		5	
installAnywhere	☑ 63	Crair		8	
Correction Corrections	1 64	local AppWittand-Generated Applications		38	
Or Counterts	2 65	lasterTech		10	
B-D NetHood	2 66	Nuch Tor		10	
D D Phone Browser Printhland	67	Autodayt		8	
	- 60 - 60	Autoutesk.		3	
	00	DAGE And Discourse		/	
Sant Menu	09	PALE ARCHITACY		10	
Loc Templates	N 10	Googe		6	
I NTUSER.DAT	71	Microsoft		9	_
n-D 🗸 🥧 NTRegistry	72	Adobe		5	
- C V Caster State	73	C Alas		5	
🗉 🕞 🏹 🚞 AppEvents	▼ 74	C Northwood Designs		17	
	<u>_ ای _</u>			~	•
	de				
D Text III Hex III Picture 🖉 Disk 🔲 Report 🖂 🤇	.onsole 🧡 P	iters [Queries V Look M 11/6/0[3068/4 : 0: PS 11/1	LS 11/1 CL 11/1 50/3	12 FOULED	
Name: AccessData					-
File Type: 2000000					
Description: Folder, Registry Entry					
Last Written: 11/08/05 09:30:10PM					
Logical Size: 10					
Physical Size: 10					



These approaches no longer work.

Hypothetical case



Dead Body



Cell Phone



All of the details matter

Operating system:

Android? iPhone? Blackberry? Feature Phone?

Access to the data:

- PIN lock?
- Encrypted Storage?
- Stored locally or in the cloud?

Applications:

- Built-in? Downloaded from "App Store"?
- Custom-written?
- Self-destruct / remote wipe?
- Malware?





Human Language: English? Korean? Chinese?



Opening up the phone doesn't guarantee the data.

A typical Android phone has:

- On-phone storage
- SD card
- SIM chip

Also:

Google SkyDrive facebook.





Digital forensics is fundamentally different from other kinds of scientific exploration...



There are five key challenges that we face...



1: Diversity — of systems, file and content

Our charter:

"Analyze any data that might be found on a computer."

Non-DF research is typically confined to a single area:



energy

literature

chemistry

DF must analyze any OS, application, protocol, encryption, etc...

math



2: Diversity over time

Today's DF tools must process:

- Today's computers / phones / cameras
 Because some criminals like to buy what's new!
- Yesterday's computers / phones / cameras
 - -Because criminals are using old devices too!



Implications for DF users and developers:

- Upgrade DF software as soon as possible.
- DF software will become geometrically more complicated over time....
 - -... or DF software will adapt on the fly to new data formats and representations.
 - -automated code analysis; pattern matching; hidden Markov models; etc.





Every year we have more data to analyze



Moore's law helps the adversary as much as us!

- We are using top-of-the-line system to analyze top-of-the-line systems
- We need to analyze in days what a subject spent weeks, months or years assembling
 - -We will never outpace the performance curve.

We must adopt "big data solutions"



4: Human capital challenges — especially in DF

Users (examiners, analysts):

- Overwhelmingly in law enforcement.
- Little or no background in CS or IS
- Deadline-driven; over-worked
- Knowledgable users tend to focus in just one particular area.
 - -Result: It takes two years to train most DF examiners.

Researchers and Developers:

- Data diversity means developers need to know the whole stack
 —opcodes & Unicode ⇒ OS & Apps ⇒ networking, encryption, etc.
- Scale issues means developers need to know HPC:
 - -threading, systems engineering, supercomputing, etc.
- Result:
 - It's hard to find qualified developers
 - Developers must be generalists





5: The "CSI Effect" — unrealistic expectations.

TV digital forensics:

- Every investigator is trained on every tool.
- Correlation is easy and instantaneous.
- There are no false positives.
- Overwritten data can be recovered.
- Encrypted data can usually be cracked.
- It is impossible to delete anything.

The reality:

- Overwritten data cannot be recovered
- Encrypted data usually can't be decrypted
- Forensics rarely answers questions or establishes guilt
- Tools crash a lot

Result:



-DF is a difficult process that looks easy





DF must respond with new science.

Current approaches don't scale.

- User spent years assembling email, documents, etc
- · Analysts have days or hours to process it
- Police analyze top-of-the-line systems
 - -with top-of-the-line systems
- National Labs have large-scale server farms
 - -to analyze huge collections

Our new algorithms must leverage our advantage: massive data

- Outlier detection and correlation
- Operate autonomously on incomplete, heterogeneous datasets
- Automatically calibrate; have no false positives







High speed forensic analysis with random sampling

Traditionally forensic analysis was leisurely. Today much analysis is under time pressure.

US agents encounter hard drives at border crossings...





US agents might have a need to search a room of computers:





What can we learn about a 1TB drive in five minutes?



Random sampling is a powerful tool for analyzing data

Simple random sampling can determine % free space



Data characterization can determine the kind of stored data



Capacity 148.87 GB	Audio 2.25 GB	Photos 8.18 GB	Other 20.25 GB	Free 118.20 GB
Audio Data reported by iTunes: 2.25 <i>GiB</i>			2.42~GB	
MP3 files reported by file	system:	2.39~GB		
Estimated MP3 usage with random sampling :			2.49 <i>GB</i> 2.71 <i>GB</i>	10,000 random samples 5,000 random samples

Sector hashing can identify specific target files





It takes 3.5 hours to read a 1TB hard drive.

In 5 minutes you can read:

- 36 GB in one strip
- 100,000 randomly chosen 64KiB strips (assuming 3 msec/seek)

		11/12	11/12
Minutes	208	5	5
Data	1 TB	36 GB	6.5 GB
# Seeks	1	1	100,000
% of data	100%	3.6%	0.65%



The statistics of a *randomly chosen sample* predict the *statistics of a population*.

US elections can be predicted by sampling thousands of households:



Hard drive contents can be predicted by sampling thousands of sectors:



The challenge is identifying *likely voters.*

The challenge is *identifying the sector* content that is sampled.



Challenge for political polls: interpreting each phone call

"On Tuesday, how will you vote for governor?"



PRAISTINTIA PER SCIENTIAN

Challenge for forensic sampling: interpreting each sector

"What data do you have?"

• Easy	/:
--------	----

000000:	ffd8	ffe0	0010	4a46	4946	0001	0201	0048	JFIFH
000010:	0048	0000	ffe1	1d17	4578	6966	0000	4 d 4d	.HExifMM
000020:	002a	0000	8000	0007	0112	0003	0000	0001	. *
000030:	0001	0000	011a	0005	0000	0001	0000	0062	b
0000040:	011b	0005	0000	0001	0000	006a	0128	0003	j.(
0000050:	0000	0001	0002	0000	0131	0002	0000	001b	
0000060:	0000	0072	0132	0002	0000	0014	0000	008d	r.2
0000070:	8769	0004	0000	0001	0000	00a4	0000	00d0	.i
0000080:	0000	0048	0000	0001	0000	0048	0000	0001	HH
0000090:	4164	6f62	6520	5068	6f74	6f73	686f	7020	Adobe Photoshop
00000a0:	4353	2057	696e	646f	7773	0032	3030	353a	CS Windows.2005:
00000b0:	3035	3a30	3920	3136	3a30	313a	3432	0000	05:09 16:01:42
00000c0:	0000	0003	a001	0003	0000	0001	0001	0000	
00000d0:	a002	0004	0000	0001	0000	00c8	a003	0004	
00000e0:	0000	0001	0000	0084	0000	0000	0000	0006	
00000f0:	0103	0003	0000	0001	0006	0000	011a	0005	•••••

• Hard:

f 00bd 15fb 5dfdqW.o	5dfdqW.o]	qW.o].
9 e154 97f4 efd5I	efd5I.T	I.	T
7 da9e d87f c12fGq	c12fGq	Gq	/
0 c57e ab7a ff00	ff00z.	• • • • • • • • • • •	~.z



We think of computers as devices with *files*.

1	🕌 Mobile Appl	ications				_	
•	GO - I	▼ Music ▼ í	Tunes - Mobile Applications	Search Mobile Applicati	Search Mobile Applications		
	Organize 🔻	New folder				= -	(?)
		Tony Birc	Name *	Date modified	Туре	Size	
		Tori Amo	iBooks 2.1 1.ipa	3/12/2012 4:04 AM	IPA File	50,045 KB	_
d		Trance C	iDisk 1.2.1.ipa	3/30/2011 4:13 AM	IPA File	3,956 KB	
1		TV Show	Keynote 1.6.ipa	3/30/2012 4:20 AM	IPA File	375,238 KB	
		Unknowr	Kindle 3.0.1.ipa	3/30/2012 4:20 AM	IPA File	20,719 KB	
		Various (MadPad 1.1.0.ipa	12/3/2011 5:59 PM	IPA File	20,096 KB	
		Vengabo	Magic Piano 4.0.2.ipa	3/30/2012 4:20 AM	IPA File	25,086 KB	
	l i	Vera Bet	MagicPlan 1.5.ipa	3/31/2012 12:35 PM	IPA File	19,026 KB	
		Vivaldi	MarbleMash 1.9.ipa	6/5/2011 1:01 PM	IPA File	7,550 KB	
		Vnv Nati	MarketDash 1.2.1.ipa	2/14/2012 8:37 PM	IPA File	4,022 KB	
		Voodoo (Memory Cards 4.3.0.ipa	6/5/2011 1:01 PM	IPA File	11,691 KB	
		Wanda L	Mobile News.ipa	7/13/2008 4:28 PM	IPA File	388 KB	
		White De	Molecules 2.02.ipa	6/5/2011 1:01 PM	IPA File	972 KB	
		William Fi	Molecules.ipa	8/3/2008 10:02 AM	IPA File	273 KB	
		Within Te	MustEatBirds.ipa	11/22/2010 2:54 PM	IPA File	9,177 KB	
:0		Wolfgan	MyFitnessPal 3.2.1.ipa	8/11/2011 7:45 PM	IPA File	21,949 KB	
		Yeah Yea	MyPad 2.5.4.ipa	2/9/2012 6:53 PM	IPA File	8,926 KB	
		yellowca	Nearby 1.ipa	9/17/2008 10:49 PM	IPA File	709 KB	
		Yes	Netflix 2.1.2.ipa	3/30/2012 3:02 PM	IPA File	17,896 KB	
-		Zack Bor	Night Stand 2.02.ipa	8/11/2011 7:31 PM	IPA File	153,568 KB	
		新居昭江	Nightstand 1.2.2.ipa	8/11/2011 7:32 PM	IPA File	14.746 KB	
c		梶浦由記	nook 3, 1, 2, 13, ipa	2/27/2012 8:40 PM	IPA File	19.344 KB	
		菅野よう	NPR 2.2.ipa	11/27/2011 7:33 AM	IPA File	3,740 KB	
	ה 📗	Tunes Plu <u>c</u>	NYTimes 2.2.ina	3/30/2012 4:20 AM	IPA File	6,715 KB	
	<u> </u>	1obile App		3/12/2012 4:04 AM		4 870 KB	
	📕 P	revious iT 🚽		5/12/2012 4:04 AM		1,070 ND	-



Data on computers is stored in fixed-sized sectors.

Data in a sector can be resident:



No Data

blank sectors

Resident data is the data you see from the root directory. e.g. "allocated" files.



Resident Data



"Deleted data" is on the disk, but can only be recovered with forensic tools.



Deleted Data



Some sectors are blank. They have "No data."



No Data


Sampling can't distinguish *allocated* from *deleted* data.





Sampling can tell us about the content of the data

Sampling can tell us the proportion of...

-blank sectors; video; HTML files; other data types...

—data with distinct signatures...



...provided we can identify it



Simplify the problem. Can we use statistical sampling to verify wiping?

Many organizations discard used computers.

Can we verify if a disk is properly wiped in 5 minutes?





Simple solution:

- 1. Read a random sector
 - If there is data, the drive is not wiped.
- 2. Repeat until satisfied.



A 1TB drive has 2 billion sectors. What if we read 10,000 and they are all blank?



A 1TB drive has 2 billion sectors. What if we read 10,000 and they are all blank?





A 1TB drive has 2 billion sectors. What if we read 10,000 and they are all blank?



Chances are good that they are all blank.



Random sampling won't find a single written sector.

If the disk has 1,999,999,999 blank sectors (1 with data)

• The sample is representative of the population.



We will only find that 1 sector with exhaustive search.



If half of the sectors are blank...

Sectors	Blank	Data
Sampled	5,000 (50%)	5,000 (50%)
Total:	1,000,000,000 (50%)	1,000,000,000 (50%)



The distribution of the data *does not matter* if sampling is random.



What if the the sampled sectors are the only blank sectors?

Sectors	Blank	Data
Sampled	10,000 (100%)	0 (0%)
Total:	10,000 (0.0005%)	1,999,990,000 (99%)



If the the only sectors read are blank...

- -We are incredibly unlucky.
- -Somebody has hacked our random number generator!



This is an example of the "urn" problem from statistics

Assume a 1TB disk has 10MB of data.

- 1TB = 2,000,000,000 = 2 Billion 512-byte sectors!
- 10MB = 20,000 sectors

Read just 1 sector; the odds that it is blank are:

$$\frac{2,000,000,000-20,000}{2,000,000,000} = .99999$$



This is an example of the "urn" problem from statistics

Assume a 1TB disk has 10MB of data.

- 1TB = 2,000,000,000 = 2 Billion 512-byte sectors!
- 10MB = 20,000 sectors

Read just 1 sector; the odds that it is blank are:

$$\frac{2,000,000,000-20,000}{2,000,000,000} = .99999$$

Read 2 sectors. The odds that both are blank are:

$$(\frac{2,000,000,000-20,000}{2,000,000})(\frac{1,999,999,999-20,000}{2,000,000,000}) = .99998$$

first pick second pick Odds we may have missed something



The more sectors picked, the less likely we are to miss the data....

$$P(X=0) = \prod_{i=1}^{n} \frac{\left(\left(N - (i-1) \right) - M \right)}{\left(N - (i-1) \right)}$$
(5)

		Non-nu	ll data	Probability of not finding data
Sampled sectors	Probability of not finding data	Sectors	Bytes	with 10,000 sampled sectors
1	0.99999	20,000	10 MB	0.90484
100	0.99900	100,000	50 MB	0.60652
1000	0.99005	200,000	100 MB	0.36786
10,000	0.90484	300,000	150 MB	0.22310
100,000	0.36787	400,000	200 MB	0.13531
200,000	0.13532	500,000	250 MB	0.08206
300,000	0.04978	600,000	300 MB	0.04976
400,000	0.01831	700,000	350 MB	0.03018
500,000	0.00673	1,000,000	500 MB	0.00673

Table 1: Probability of not finding any of 10MB of data ona 1TB hard drive for a given number of randomly sampledsectors. Smaller probabilities indicate higher accuracy.

- Pick 500,000 random sectors

- If are all NULL, the disk has p=(1-.00673) chance of having 10MB of non-NULL data



- The disk has a 99.3% chance of having less than 10MB of data

Table 2: Probability of not finding various amounts of data when sampling 10,000 disk sectors randomly. Smaller probabilities indicate higher accuracy.

In practice, we use a modified algorithm...

Sample with 64KiB "blocks" instead of 512-byte sectors.

- It takes the same amount of time to read 65,536 bytes as 512 bytes
- Analyze 64KiB block with a 4KiB sliding window
- On a 1TB drive, there are 15,258,789 64KiB sections

Identify data "type"

- Blank
- JPEG
- Video
- Encrypted

Update results in real-time

- Provides immediate feedback
- Catches important data faster
- Stop when analyst is satisfied.





We used this technique to calculate the size of the TrueCrypt volume on this iPod.

It takes 3+ hours to read all the data on a 160GB iPod.

• Apple bought very slow hard drives.





We got a statistically significant sample in two minutes.



The % of the sample will approach the % of the population.



The challenge: identifying a file "type" from a fragment.





One approach: hand-tuned discriminators based on a close reading of the specification.

For example, the JPEG format "stuffs" FF with a 00.

87654321	<mark>00</mark> 11	2233	4455	6677	8899	aabb	ccdd	eeff	0123456789abcdef 📃		
00006a20:	6b4c	cd62	54a0	b214	52ff	0074	ba4f	4622	kL.bTRt.OF"		
00006a30:	d1bf	bf4c	67c4	aa2a	4a91	036f	f3b3	7ddc	Lg*Jo}.		
00006a40:	98d5	f078	7f28	d327	340d	a2f2	c916	da4f	x.(.'40		
00006a50:	aefa	0cbc	e9a6	a580	4b20	952c	17d2	7a09	K .,z.		
00006a60:	377b	097c	7395	b7e4	c661	730c	447 f	9b5a	7{.lsas.DZ		
00006a70:	7675	e9d1	e14a	81a8	26a2	2948	93bc	4749	vuJ&.)HGI		
00006a80:	94fd	8d3f	fce2	4a13	e529	2b64	8f31	b961	?J)+d.1.a		
00006a90:	368b	827f	677e	7a64	9a62	60f9	9826	c4e0	6g~zd.b`&		
00006aa0:	b65e	bfa9	97fc	5aa9	6a94	626a	602e	4ac7	.^Z.j.bj`.J.		
00006ab0:	9cb1	0311	3d9d	3e33	e941	482e	caf2	8676	=.>3.AHv		
00006ac0:	240d	43ae	ce27	a39e	98d3	f14a	6a23	116a	\$.C'Jj#.j		
00006ad0:	af80	dffc	1867	58be	0eaa	a9a9	b29f	3331	gX31		
00006ae0:	20b1	9da6	46d3	eb6d	4846	774c	1870	4c98	FmHFwL.pL.		
00006af0:	60fd	0f7d	8382	2f04	e2a9	e314	d982	5947	`}/YG		
00006600:	<u>1</u> 1ef	bef1	7df3	9c6a	f0ab	289d	2d99	b6fb	<u>.</u> }j(
00006610:	f f00	9b6d	a903	35aa	8b3c	8014	9240	6006	m5<@`.		
00006b20:	cece	5c3b	9f4d	af7f	8934	44d8	bd10	4044	\;.M4D@D		
00006630:	0124	bd6e	b80d	61ff	001d	388c	8b74	aaef	.\$.na8t		
00006640:	3219	3010	c487	a6fa	681a	4a23	4a8a	5441	2.0h.J#J.TA		
00006650:	5600	3e19	7762	443b	1376	07a1	96c6	5553	[.>.wbD;.vUS		
00006660:	4bbc	285a	7e57	393d	e521	e8ce	b48a	c99a	K.(Z~₩9=.! M		
00006670:	69aa	9129	bdab	0361	ba5b	6636	418d	3e85	1)a.[l6A.>.		
00006680:	2c2b	5fc4	55c2	162e	0a60	1209	2144	5887	,+U`!DX.		
00006690:	20a4	3055	81c3	a566	799d	84b2	1493	28ac	.0Ufy(.		
-:F1	iStoc	k Priv	vacy.	jpg	8%	_1714	(He)	(L)	-8:37PM 🍸		
Mark saved	d whei	re sed	arch (starte	ed				×		



We built detectors to recognize the different parts of a JPEG file.



000107.jpg Bytes: 41,572



Sectors: 82



Nearly 50% of this 57K file identifies as "JPEG"





000897.jpg Bytes: 57596

Sectors: 113



Nearly 100% of this file identifies as "JPEG."





000512.jpg Bytes: 195,311

Sectors: 382

This is called the *file fragment classification problem*.

We can reliably classify JPEG, MPEG, Huffman, and other types.



Combine random sampling with sector ID to obtain the forensic contents of a storage device.

Our numbers from sampling are similar to those reported by iTunes.

Capacity 148.87 GB	Audio 2.25 GB	Photos 8.18 GB	Other 20.25 GB	Free 118.20 GB
Audio Data reported by iTu	nes:	2.25~GiB	2.42~GE	3
MP3 files reported by file sy	vstem:		2.39~GE	3
Estimated MP3 usage with 1	andom samplin	g :	2.49~GE	3 10,000 random samples
			2.71~GE	3 5,000 random samples

Figure 1: Usage of a 160GB iPod reported by iTunes 8.2.1 (6) (top), as reported by the file system (bottom center), and as computing with random sampling (bottom right). Note that iTunes usage actually in GiB, even though the program displays the "GB" label.

We accurately determined:

% of free space; % JPEG; % encrypted

-Simson Garfinkel, Vassil Roussev, Alex Nelson and Douglas White, Using purpose-built functions and block hashes to enable small block and sub-file forensics, DFRWS 2010, Portland, OR









Finding Known Content with Sector Hashing...

Most forensics processing tries to understand the internal structure of data files...





Files can also be viewed as a set of ordered blocks.



41,572 bytes

Block #	Byte Range	Values			
0	0- 511	ffd8 ffe0 0010 4a46 4946 0001 0201 0048			
1	512-1023	0c0c 0c0c ffc0 0011 0800 6a00 a003 0122			
2	1024-1535	4fa7 7567 ded2 cac5 8c82 2bf4 9e1c 23f9			
3	1536-2047	fafd 1527 e459 e934 c173 59ad 9234 f09f			
4	•••				



Compute the cryptographic hash of each block. These are "block hashes."



MD5*(block(N))	Byte Range	Block #
dc0c20abad42d487a74f308c69d18a5	0- 511	0
9e7bc64399ad87ae9c2b54506195977	512-1023	1
6e7f3577b100f9ec7fae18438fd5b04	1024-1535	2
4594899684d0565789ae9f364885e30	1536-2047	3
	• • •	4

Question: how often do *these* block hashes occur in other JPEGs?



Should these block hashes be in other files?



Specific byte sequences in high-entropy data are very rare.

• 512 bytes = $256^{512} = 10^{1,233}$ possible sectors

But metadata might be common:

- Specific headers
- Common color tables
- "all black"

You need to survey the datasphere to find out.

Header	[FF D8 FF E0] or [FF D8 FF E1]
Icons	MD5*(block(N))
EXIF	dc0c20abad42d487a74f308c69d18a5a
Color Table	9e7bc64399ad87ae9c2b545061959778
	6e7f3577b100f9ec7fae18438fd5b047
Huffman Encoded	4594899684d0565789ae9f364885e303
Data	
Footer	◀– [FF D9]



We examined sector hashes from \simeq 4 million files

- \simeq 1 million in GOVDOCS1 collection
- = 109,282 JPEGs (including 000107)
- \simeq 3 million samples of Windows malware

Our results:

- Most of the block hashes in 000107.jpg did not appear elsewhere in the corpus.
- Some of the block hashes appeared in other JPEGs.
- None of the block hashes appeared in files that were not JPEGs



The beginning of 000107.jpg contained distinct hashes...

dc0c20abad42d487a74f308c69d18a5a 9e7bc64399ad87ae9c2b545061959778 6e7f3577b100f9ec7fae18438fd5b047 4594899684d0565789ae9f364885e303 4d21b27ceec5618f94d7b62ad3861e9a 03b6a13453624f649bbf3e9cd83c48ae c996fe19c45bc19961d2301f47cabaa6 0691baa904933c9946bbda69c019be5f 1bd9960a3560b9420d6331c1f4d95fec 52ef8fe0a800c9410bb7a303abe35e64 b8d5c7c29da4188a4dcaa09e057d25ca 3d7679a976b91c6eb8acd1bfa3414f96 8649f180275e0b63253e7ee0e8fa4c1d 60ebc8acb8467045e9dcbe207f61a6c2 440c1c1318186ac0e42b2977779514a1 72686172f8c865231e2b30b2829e3dd9 fdff55c618d434416717e5ed45cb407e fcd89d71b5f728ba550a7bc017ea8ff1 2d733e47c5500d91cc896f99504e0a38 2152fdde0e0a62d2e10b4fecc369e4c6 692527fa35782db85924863436d45d7f 76dbb9b469273d0e0e467a55728b7883 171310e61a8e78364b4965b995f16ff5 6865477474f8a6011108c9cbf1fff0f9

offset	0-511	1	L
offset	512-1023	1	L
offset	1024-1535	1	L
offset	1536-2047	1	L
offset	2048-2559	1	L
offset	2560-3071	1	L
offset	3072-3583	1	L
offset	3584-4095	1	L
offset	4096-4607	1	L
offset	4608-5119	1	L
offset	5120-5631	1	L
offset	5632-6143	1	L
offset	6144-6655	1	L
offset	6656-7167	1	L
offset	7168-7679	1	L
offset	7680-8191	1	L
offset	8192-8703	1	L
offset	8704-9215	1	L
offset	9216-9727	1	L
offset	9728-10239	1	L
offset	10240-10751	1	L
offset	10752-11263	1	L
offset	11264-11775	1	L
offset	11776-12287	1	L



The middle of 000107.JPG had hash collisions...

offset 14848-15359 1 offset 15360-15871 1 offset 15872-16383 1 offset 16384-16895 1 offset 16896-17407 1 offset 17408-17919 1 offset 17920-18431 1 offset 18432-18943 14 offset 18944-19455 9198 offset 19456-19967 9076 offset 19968-20479 9118 offset 20480-20991 9237 offset 20992-21503 9708 offset 21504-22015 9615 offset 22016-22527 9564 offset 22528-23039 7 offset 23040-23551 83 83 offset 23552-24063 offset 24064-24575 84 offset 24576-25087 84 offset 25088-25599 84 offset 25600-26111 1 offset 26112-26623 1 offset 26624-27135 1

9df886fdfa6934cc7dcf10c04be3464a 95399e7ecc7ba1b38243069bdd5c263a ef1ffcdc11162ecdfedd2dde644ec8f2 7eb35c161e91b215e2a1d20c32f4477e 38f9b6f045db235a14b49c3fe7b1cec3 edceba3444b5551179c791ee3ec627a5 6bc8ed0ce3d49dc238774a2bdeb7eca7 5070e4021866a547aa37e5609e401268 13d33222848d5b25e26aefb87dbdf294 Odfcde85c648d20aed68068cc7b57c25 756f0bbe70642700aafb2557bf2c5649 c2c29016d3005f7a1df247168d34e673 42ff3d72b2b25f880be21fac46608cc9 b943cd0ea25e354d4ac22b886045650d a003ec2c4145b0bc871118842b74f385 1168c351f57aad14de135736c06665ea 51a50e6148d13111669218dc40940ce5 365b122f53075cb76b39ca1366418ff9 9ad9660e7c812e2568aaf063a1be7d05 67bd01c2878172e2853f0aef341563dc fc3e47d734d658559d1624c8b1cbf2c1 cb9aef5b7f32e2a983e67af38ce8ff87 531aea9e5b2987f923b0f0812bd5846e cef61251eb556fd095b3347dc87d8a24



Block 37 had 9198 collisions.. The sector is filled with blank lines 100 characters long...

13d33222	2848d	5b25e	26ae	fb87d	lbdf2	94	offs	et	18944-19455
\$ dd if=0	000107	7.jpg	skip=	=18944	l cour	nt=512	2 bs=1	L x :	ĸd
000000:	2020	2020	2020	2020	2020	2020	2020	202	20
000010:	2020	2020	2020	2020	2020	2020	0a20	202	20
000020:	2020	2020	2020	2020	2020	2020	2020	202	20
000030:	2020	2020	2020	2020	2020	2020	2020	202	20
0000040:	2020	2020	2020	2020	2020	2020	2020	202	20
0000050:	2020	2020	2020	2020	2020	2020	2020	202	20
0000060:	2020	2020	2020	2020	2020	2020	2020	202	20
0000070:	2020	2020	2020	2020	2020	2020	2020	202	20
0000080:	200a	2020	2020	2020	2020	2020	2020	202	20.
0000090:	2020	2020	2020	2020	2020	2020	2020	202	20
00000a0:	2020	2020	2020	2020	2020	2020	2020	202	20
0000b0:	2020	2020	2020	2020	2020	2020	2020	202	20
00000c0:	2020	2020	2020	2020	2020	2020	2020	202	20
00000d0:	2020	2020	2020	2020	2020	2020	2020	202	20
00000e0:	2020	2020	2020	0a20	2020	2020	2020	202	20.
00000f0:	2020	2020	2020	2020	2020	2020	2020	202	20
0000100:	2020	2020	2020	2020	2020	2020	2020	202	20
0000110:	2020	2020	2020	2020	2020	2020	2020	202	20
0000120:	2020	2020	2020	2020	2020	2020	2020	202	20
0000130:	2020	2020	2020	2020	2020	2020	2020	202	20
0000140:	2020	2020	2020	2020	2020	200a	2020	202	20
0000150:	2020	2020	2020	2020	2020	2020	2020	202	20
0000160:	2020	2020	2020	2020	2020	2020	2020	202	20



9198

Block 45 had 83 collisions.. It appears to contain EXIF metadata

5	1a5	50e6	148d1	31116	6921	8dc40)940c	e5	offs	et 23	040-23551
\$	dd	if=	000107	7.jpg	skip=	=23040) cour	nt=512	2 bs=1	l xxd	
0	000	000:	3936	362d	322e	3100	0000	0000	0000	0000	966-2.1
0	000	010:	0000	0000	0000	0000	0000	0000	0000	0000	• • • • • • • • • • • • • • • • •
0	000	020:	0000	0000	0000	0000	0000	0000	0000	0000	• • • • • • • • • • • • • • • • •
0	000	030:	0000	0000	0000	0000	0058	595a	2000	0000	XYZ
0	000	040:	0000	00f3	5100	0100	0000	0116	cc58	595a	QXYZ
0	000	050:	2000	0000	0000	0000	0000	0000	0000	0000	• • • • • • • • • • • • • • •
0	000	060:	0058	595a	2000	0000	0000	006f	a200	0038	.XYZ8
0	000	070:	£500	0003	9058	595a	2000	0000	0000	0062	b
0	000	080:	9900	00b7	8500	0018	da58	595a	2000	0000	XYZ
0	000	090:	0000	0024	a000	000f	8400	00b6	cf64	6573	\$des
0	000	0a0:	6300	0000	0000	0000	1649	4543	2068	7474	cIEC htt
0	000	0b0:	703a	2f2f	7777	772e	6965	632e	6368	0000	p://www.iec.ch
0	000	0c0:	0000	0000	0000	0000	0016	4945	4320	6874	IEC ht
0	000	0d0:	7470	3a2f	2£77	7777	2e69	6563	2e63	6800	tp://www.iec.ch.
0	000	0e0:	0000	0000	0000	0000	0000	0000	0000	0000	
0	000	OfO:	0000	0000	0000	0000	0000	0000	0000	0000	• • • • • • • • • • • • • • • •
0	000	100:	0000	0000	0000	0000	0000	0000	0064	6573	des
0	000	110:	6300	0000	0000	0000	2e49	4543	2036	3139	c IEC 619
0	000	120:	3636	2d32	2e31	2044	6566	6175	6c74	2052	66-2.1 Default R
0	000	130:	4742	2063	6f6c	6f75	7220	7370	6163	6520	GB colour space
0	000	140:	2d20	7352	4742	0000	0000	0000	0000	0000	- sRGB
0	000	150:	002e	4945	4320	3631	3936	362d	322e	3120	IEC 61966-2.1
0	000	160:	4465	6661	756c	7420	5247	4220	636f	6c6f	Default RGB colo
0	000	170:	7572	2073	7061	6365	202d	2073	5247	4200	ur space - sRGB.



Block 48 had 84 collisions.. It appears to contain part of a JPEG color table...

67bd0	1c28783	L72e28	853f0	aef34	1563	dc	offs	et 245	576-25087
\$ dd i	f=00010	7.jpg	skip	=24576	5 cour	nt=512	2 bs=1	l xxd	
000000	0: 7a27	ab27	dc28	0d28	3f28	7128	a228	d429	z'.'.(.(?(q(.(.)
000001	0: 0629	3829	6b29	9d29	d02a	022a	352a	682a	.)8)k).).*.*5*h*
000002	0: 9b2a	cf2b	022b	362b	692b	9d2b	d12c	052c	.*.+.+6+i+.+.,.,
000003	0: 392c	: 6e2c	a22c	d72d	0c2d	412d	762d	ab2d	9, n, .,A-v
000004	0: e12e	e 162e	4c2e	822e	b72e	ee2f	242f	5a2f	L/\$/Z/
000005	0: 912f	c72f	fe30	3530	6c30	a430	db31	1231	././.05010.0.1.1
000006	0: 4a31	8231	ba31	f232	2a32	6332	9b32	d433	J1.1.1.2*2c2.2.3
000007	0: 0d33	4633	7£33	b833	f134	2b34	6534	9e34	.3F3.3.3.4+4e4.4
000008	0: d835	5 1335	4d35	8735	c235	fd36	3736	7236	.5.5M5.5.5.676r6
000009	0: ae36	e937	2437	6037	9c37	d738	1438	5038	.6.7\$7`7.7.8.8P8
00000a	0: 8c38	c839	0539	4239	7£39	bc39	f93a	363a	.8.9.9B9.9.9.:6:
00000b	0: 743a	b23a	ef3b	2d3b	6b3b	aa3b	e83c	273c	t:.:;-;k;.;.<'<
00000c	:0: 653c	a43c	e33d	223d	613d	a13d	e03e	203e	e<.<.="=a=.=.> >
00000đ	0: 603e	a03e	e03f	213f	613f	a23f	e240	2340	`>.>.?!?a?.?.@#@
00000e	0: 6440	a640	e741	2941	6a41	ac41	ee42	3042	d@.@.A)AjA.A.BOB
00000f	0: 7242	b542	£743	3a43	7d43	c0 44	0344	4744	rB.B.C:C}C.D.DGD
000010	0: 8a44	ce45	1245	5545	9a45	de46	2246	6746	.D.E.EUE.E.F"FgF
000011	0: ab46	5 f047	3547	7b47	c048	0548	4b48	9148	$.F.G5G{G.H.HKH.H}$
000012	0: d749	1 d 49	6349	a949	f04a	374a	7d4a	c44b	.I.IcI.I.J7J}J.K
000013	0: 0c4b	534b	9a4b	e24c	2a4c	724c	ba4d	024d	.KSK.K.L*LrL.M.M
000014	0: 4a4d	l 934d	dc4e	254e	6e4e	b74f	004f	494f	JM.M.N%NnN.O.OIO
000015	0: 934f	dd50	2750	7150	bb51	0651	5051	9b51	.O.P'PqP.Q.QPQ.Q
000016	0: e652	3152	7c52	c753	1353	5f53	aa53	f654	.R1R R.S.S_S.S.T
000017	0: 4254	8f54	db55	2855	7555	c256	0f56	5c56	BT.T.U(UuU.V.V\V



84

With blocks of 512 bytes and 4KiB, the vast majority of sectors had distinct hashes.

Table 1. Incidence of singleton, paired, and common sectors in three file corpora.						
Govdocs		OpenMalv	ware 2012	2009 NSRL RDS		
Block size: 512 bytes						
911.4 M	(98.93%)	1,063.1 M	(88.69%)	Ν	I/A	
7.1 M	(.77%)	75.5 M	(6.30%)	Ν	I/A	
2.7 M	(.29%)	60.0 M	(5.01%)	N/A		
Block size: 4 kibibytes						
117.2 M	(99.46%)	143.8 M	(89.51%)	567.0 M	(96.00%)	
0.5 M	(.44%)	9.3 M	(5.79%)	16.4 M	(2.79%)	
	Incidenc Gov 911.4 M 7.1 M 2.7 M 117.2 M 0.5 M	Incidence of single in thre in three in	Incidence of singleton, paire in three file corporation in the file corporation in three file corporation in the file corporation in three file corporation in the file corporation in three file corporation in three file corporation in	Incidence of singleton, paired, and com in three file corpora. Govdocs OpenMalware 2012 Block size: 512 bytes 911.4 M (98.93%) 1,063.1 M (88.69%) 7.1 M (.77%) 75.5 M (6.30%) 2.7 M (.29%) 60.0 M (5.01%) Block size: 4 kibibytes 117.2 M (99.46%) 143.8 M (89.51%) 0.5 M (.44%) 9.3 M (5.79%)	Incidence of singleton, paired, and common second sin three file corpora. Govdocs OpenMalware 2012 2009 N Block size: 512 bytes 911.4 M (98.93%) 1,063.1 M (88.69%) N 7.1 M (.77%) 75.5 M (6.30%) N 2.7 M (.29%) 60.0 M (5.01%) N Block size: 4 kibibytes 117.2 M (99.46%) 143.8 M (89.51%) 567.0 M 0.5 M (.44%) 9.3 M (5.79%) 16.4 M	

Young, Foster, Garfinkel & Fairbanks, IEEE Computer, Dec. 2012

(.11%)

(4.71%)

7.6 M

7.1 M

(1.21%)

0.1 M



Common

File systems align large files on sector boundaries. We hash file blocks and identify sectors that match.





Using distinct sectors in media sampling and full media analysis to detect presence of documents from a corpus,

This means we can use distinct sectors to find known content.

Method #1 — Full media sampling

- Read & hash every disk sector.
- Lookup hash values in a database of block hashes.
- Distinct hash imply presence of files.
- Advantage: Can find a single sector of target content

Method #2 — Random sampling

- Read & hash randomly chosen sectors.
- Lookup hash values in a database of block hashes.
- Distinct hash implies presence of files.
- Advantage: Can find presence of target content very quickly





There are significant hash and database requirements.

1TB data in 208 minutes

- \simeq 80 Mbyte/sec
- \simeq 150,000 512-byte sectors/sec
- \simeq 150,000 database lookups/sec

		11/12	
Minutes	208	5	
Max Data	1 TB	36 GB	
Max Seeks		90,000	

Alignment uncertainty gives 4096-byte sectors same performance requirements:




By combining a Bloom filter & database, we can perform up to 2.7M TPS on low-cost hardware

Table 2. Total transactions per second (TPS) for best execution.									
Bloom filter			Database		TPS at 1 M lookups		TPS at 1,200 seconds		
k	м	Size	Strategy	Size	Present	Absent	Present	Absent	
100 million records									
3	31	257 MiBytes	B-tree (preload)	2.3 GiBytes	35.3 K	49.5 K	161.3 K	1.8 M	
3	31	257 MiBytes	B-tree	2.3 GiBytes	11.6 K	565.8 K	156.8 K	2.3 M	
3	31	257 MiBytes	Hash map	5.3 GiBytes	13.9 K	656.9 K	641.9 K	3.0 M	
3	31	257 MiBytes	Flat map	2.2 GiBytes	28.2 K	746.9 K	356.4 K	2.6 M	
3	31	257 MiBytes	Red/black tree	6.0 GiBytes	12.9 K	694.5 K	187.0 K	2.7 M	
1 billion records									
3	34	2.1 GiBytes	B-tree (preload)	23 GiBytes	2.2 K	6.1 K	3.6 K	23.1 K	
3	33	1.1 GiBytes	B-tree	23 GiBytes	2.6 K	85.8 K	3.7 K	114.9 K	
3	33	1.1 GiBytes	Hash map	57 GiBytes	-	-	0.3 K	3.1 K	
3	34	2.1 GiBytes	Flat map	22 GiBytes	-	-	0.4 K	4.0 K	
3	33	1.1 GiBytes	Red/black tree	60 GiBytes	-	-	0.1 K	1.4 K	

Hardware: 8GiB Laptop; 250GB external SSD.

- "Distinct sector hashes for target file detection," Young, Garfinkel, Foster & Fairbanks, IEEE Computer, Dec. 2012



Putting it all together, we have a significant innovation... field deployable on a single laptop.

Use Case #1: Rapidly search for known contraband:

- 1TB subject hard drive.
- 10 min x 60 min/sec x 1000 msec/sec / 3 msec/sample = 200,000 samples
- Searching for a sector from a corpus of 512GB
- 100% recognition of a single sector; 0% false positive rate

Amount of Contraband	p (prob of missing contraband)
5 MB	0.3654
10 MB	0.1335
15 MB	0.0488
20 MB	0.0178
25 MB	0.0065





Use Case #2: Find a single sector of known contraband:

Time to read data & search database: 208 minutes

Technique is file type and file system agnostic

-JPEG; Video; MSWord; Encrypted PDFs...

-provided data is not modified when copied or otherwise re-coded





Find out more!

For further reading...

Innovative Technology for Computer Professionals

SZ FORENSI

CS PRESIDENT'S MESSAGE, P. 6 COMPUTING CONVERSATIONS: VINT CERF, P. 10

IAAS CLOUD ARCHITECTURE, P. 65

IEEE



വ

-

ດ

œ

14

8

 \geq ш U

1-1

http://www.computer.org



oel Young, Kristina Foster, and Simson Garfinkel, Naval Postgraduate School Kevin Fairbanks, Johns Hopkins University

sing an alternative approach to traditional file hashing, digital forensic investigators can hash individually sampled subject drives on sector boundaries and then check these hashes against a prebuilt da-tabase, making it possible to process raw media without reference to the underlying file system

Grensic examiners frequently search disk drives, cell phones, and even network flows to determine if specific known content is present. For example, a corporate security officer might examine a sus-picious employee's laptop for unauthorized documents, we notoccenter officers might search a suspect's home computer for illegal pornography: and network analysts up to the search and the search and the search and any other case, examines ty typically identify files by computing their cryptographic hash—othen with MD5 or thet resulting hash value. Used hash values for file identification is pervasive age hash values for file identification is pervasive age hash values for file identification is pervasive age hash values for file identification is 25, released in March 2012, contains 25, 59/244 distinct file hashes towstoment of specific companies and to law enforcement

ustomers of specific companies and to law enforcement ganizations

COVER FEATURI

S

such as GPS, gyros

r hackers. Since the fir

28 COMPLITER Published by the IEEE Computer Society 0018-9162/12/531.00 © 2012 IEEE

There are many limitations when using file hashes to

identify known content. Because changing just a single bit of a file changes its hash, pornographers, malware au-thors, and other miscreants can evade detection simply by

thors, and other miscreants can evade detection simply by changing a comma to a period or appending a few random bytes to a file. Likewise, hash-based identification will not work if sections of the file are damaged or otherwise un-recoverable. This is especially a problem when large video files are deleted and the operating system reuses a few sectors for other purposes, most of the video is still present on the drive, but recovered video segments will not appear to a drabinet or <u>6.16</u> he hashes.

SECTOR HASHING We are developing alternative systems for detecting target files in large disk images using cryptographic bashes on sectors of data rather than entite files. Modern file sys-tems align the start of most files with the beginning of a disk sector. Thus, when a megaphyse-side video is stored on a modern hard drive, the first 4 kiblytes are stored in one disk sector, the second 4 kiblytes are stored in another disk sector, typically the adjacent one, and so on, (In our work, we distinguish thetwore power of-two based sizes of digital

sector, typically the adjacent one, and so on, (in our work, we distinguish the twene power-of-two-based sizes of digital artifacts, such as kiblytes, and power-of-ten-based sizes, such as kiblotytes. See the "Docimal versus Binary Prefixes" sidebar for more details.) Furthermore, by sampling ran-domly chosen sectors from the drive, it is only necessary to read a unit praction of the drive, to determine with high probability if a target file is present. This enables rapid triage of drive images. We compare drive sector hashes to a hash database of "two-dairoff fle frammers: which was call block. The terms

fixed-sized file fragments, which we call blocks. The terms "sector" and "block" are often used incorrectly as syn-

in a database of file hashes

SECTOR HASHING



SCADA Systems: **Challenges** for Forensic **Investigators**

stems consisted of simple I/O devices that transmi

the signals between master and remote terminal units. In recent years, SCADA systems have evolved to communi-

cate over public IP networks.2 Some are also connected to a corporate intranet or directly to the Internet to seamless

integrate SCADA data with external information such as

The integration of SCADA systems within a much wider

The integration of S-ADA systems within a much wider network brings threas that were unimagined at the time these systems were conceived. During the past decade, vendors, asset owners, and regulators recognized this growing concern and began to address it through new laws and various security mechanisms, processes, and

The discoveries in the wild of Stuxnet in June 2010

and Flame in May 2012 were additional eye-openers for SCADA owners and operators. Stuxnet, the first known

SCADA OWIER's and Operators. Subtriet, the first known malware designed to target automation systems, has infected 50,000 to 100,000 computers worldwide," while Flame is a cyberespionage tool an order of magnitude more sophisticated than Stuxnet.²

SCADA ARCHITECTURE As Figure 1 shows, a typical SCADA system for control-ling infrastructures for utilities such as power, gas, oil, or water generally consists of a control center and nu-

merous field sites. The sites are distributed over a wide geographical area and are connected to the control center

by different communication media such as satellites, wide

rate email or weather data.

standards.

Irfan Ahmed, University of New Orleans er and Martin Naedele, ABB Corporate Research Golden G. Richard III, University of New Orleans

hen security incidents occur, several challenges exist for conducting an effective forensic investigation of SCADA systems, which run 24/7 to control and monitor industrial and infrastructure processes.

n industrial automation and control system is a set of devices that regulate the behavior of physical processes. For example, a thermo-statis a simple control system that senses the temperature and turns a heater on or off to maintain the temperature at a set point. These systems are used to monitor and control industrial and infrastructure processes such as chemical plant and oil refinery opera tions, electricity generation and distribution, and water management. A control system that is spread over a wide area and

A control system that is spread over a wide area and can supervise its individual components is often called a supervisory control and data acquisition (SCADA) system.¹ However, here we use the term SCADA to refer to all kinds of control systems that share a common key characteristic: they are connected to physical processes and thus need to be continuously available and able to respond within a deterministic time bound. Early SCADA systems were intended to run as isolated networks, not connected to the Internet, and thus did not

require any specific cybersecurity mechanisms. These

44 COMPLETER

Published by the IFFE Computer Society

0018-9162/12/\$31 00 © 2012 IFFF

Smartphone Security **Challenges**

Yong Wang, Kevin Streff, and Sonell Raman, Dakota State University

Recause of their unique characteristics 2011, malware attacks on the Android platform inc seconds of their unique characteristics, smartphones present challenges requiring new business models that offer counter-neasures to help ensure their security. 2011, marware attacks on the Antorou platform increased 3,325 percent¹⁰ As the use of smartphones continues its rapid growth, subscribers must be assured that the services they offer are reliable, secure, and trustworthy.

SMARTPHONE THREATS AND ATTACKS

Sumarphones are quickly becoming the dominant device for accessing interent securces. Sales of marker in Q4 2010: "Shipments of samarphones surpassed those of feature phones in Western Europe ing Q2 2011-"According to a May 2011 Nielsen survey, smarphones outsold feature phones in the US in this same period." Compared to 3.9 Million workholds phone subactibers, smarphone usage G35 million is sull some priod." Compared to 3.9 Million estiphenes that an experiment of the same state of the same prior and the same state of the same state of the phone subactibers, smarphone usage G35 million is sull some state. The same state of the same state of the linear state of the same hone, which carries a large amount ensitive data. After infiltrating a smartphone, the malwar attempts to control its resources, collect data, or redirect th nartphone to a premium account or malicious website This model divides a smartphone into three layers:

 The application layer includes all of the smartphone's apps, such as social networking software, email, text messaging, and synchronization software.
The communication layer includes the carrier netsmartphones ouer many more uncloses una tradi-tional mobile phones. In addition to a perinstalled mobile operating system, such as iOS, Android, or Windows Mobile, most smartphones also typically support carrier networks; Wi-Fr connectivity, and Buetocoths on that users can access the Internet to download and run various third-The communication tayer includes the carrier net-works, Wi-Fi connectivity, Bluetooth network, Micro USB ports, and MicroSD slots. Malware can spread through any of these channels.

Through any or unact the flash memory, camera microphone, and sensors within a smartphone. Be cause smartphones contain sensitive data, malwan

sessage service (mins) and include entreduce sensors uch as GPS, gyroscopes, and accelerometers, as well as high-resolution camera, a microphone, and a speaker. Smartphones' increasing popularity raises many secu-ty concerns.⁶⁴ Their central data management makes An attack forms a loop starting with the launch of the nalware, moving through the smartph alayers, on to pro

Please try our tools!

bulk_extractor, a high-performance stream-based feature extractor

(dev tree)

(downloads)

- https://github.com/simsong/bulk_extractor
- http://digitalcorpora.org/downloads/bulk_extractor
- http://www.sciencedirect.com/science/article/pii/S0167404812001472 (paper)
 - -Computers & Security, 2013
 - -http://simson.net/clips/academic/2013.COSE.bulk_extractor.pdf

DFXML — An XML language for doing computer forensics

- provenance, file extraction, hashes and piecewise-hashes, registry values, etc.
- https://github.com/simsong/dfxml
- http://www.sciencedirect.com/science/article/pii/S1742287611000910
 - -Digital Investigation, 2012
 - -http://simson.net/clips/academic/2012.DI.dfxml.pdf

Data!



http://digitalcorpora.org/

In summary, there are many opportunities in digital forensics.

Math and Science:

- Algorithms tolerant of data that is dirty and damaged.
- New approaches for handling data that are compressed, encoded or encrypted
- Linguistics, Natural Language Processing & Machine Learning
- Visualization

Engineering:

- Reverse engineering & product development
- Approaches for dealing with large data volumes (100TB 10PB)
- Software that doesn't crash

Many of the techniques here are also applicable to:

- Social Network Analysis
- Personal Information Management
- Data mining unstructured information



Contact Information: Simson L. Garfinkel simsong@acm.org http://simson.net/