# Digital Forensics for Archivists: Fundamentals

**Instructor:**
Christopher (Cal) Lee
University of North Carolina at Chapel Hill

September 22, 2016
Greeley, CO

1

---

## Digital Archives Specialist (DAS)

Curriculum and Certification Program offered by SAA:

- <u>Foundational</u> Courses—*must pass 4*

- <u>Tactical and Strategic</u> Courses—*must pass 3*

- <u>Tools and Services</u> Courses—*must pass 1*

- <u>Transformational</u> Courses—*must pass 1*

- **Course examinations are administered online.**

2

## Agenda

- Welcome and introductions
- Motivation and scope
- Technical background
- Representation Information
- File systems and file management
- Extracting data from media
- Tools and methods
- Conclusions, questions, discussion

3

## Welcome and Introductions

4

## Motivation and Scope
### Applying Digital Forensics to Archival Work

## Many archivists know how to process this stuff:



Source: The Processing Table: Reflections on a manuscripts internship at the Lilly Library.
https://processingtable.wordpress.com/tag/archival-processing/

## How about processing this stuff?

Source: "Digital Forensics and creation of a narrative." *Da Blog: ULCC Digital Archives Blog.*
http://dablog.ulcc.ac.uk/2011/07/04/forensics/

7

---

## Same Goals as When Acquiring Analog Materials

- Ensure integrity of materials
- Allow users to make sense of materials and understand their context
- Prevent inadvertent disclosure of sensitive data

8

## Same Fundamental Archival Principles Apply

| | |
|---|---|
| Provenance | • Reflect "life history" of records<br>• Records from a common origin or source should be managed together as an aggregate unit |
| Original Order | Organize and manage records in ways that reflect their arrangement within the creation/use environment |
| Chain of Custody | • "Succession of offices or persons who have held materials from the moment they were created"[1]<br>• Ideal recordkeeping system would provide "an unblemished line of responsible custody"[2] |

1. Pearce-Moses, Richard. *A Glossary of Archival and Records Terminology.* Chicago, IL: Society of American Archivists, 2005.
2. Hilary Jenkinson, *A Manual of Archive Administration: Including the Problems of War Archives and Archive Making* (Oxford: Clarendon Press, 1922), 11.

**9**

## But you might need some of this stuff:



**10**

**Luckily, there are a lot of people with expertise in using such tools in places like this:**



El Paso County Sheriff's Office (Colorado)

http://shr.elpasoco.com/Law+Enforcement+Bureau/Investigations+Division/Computer+Crime+Lab.htm **11**

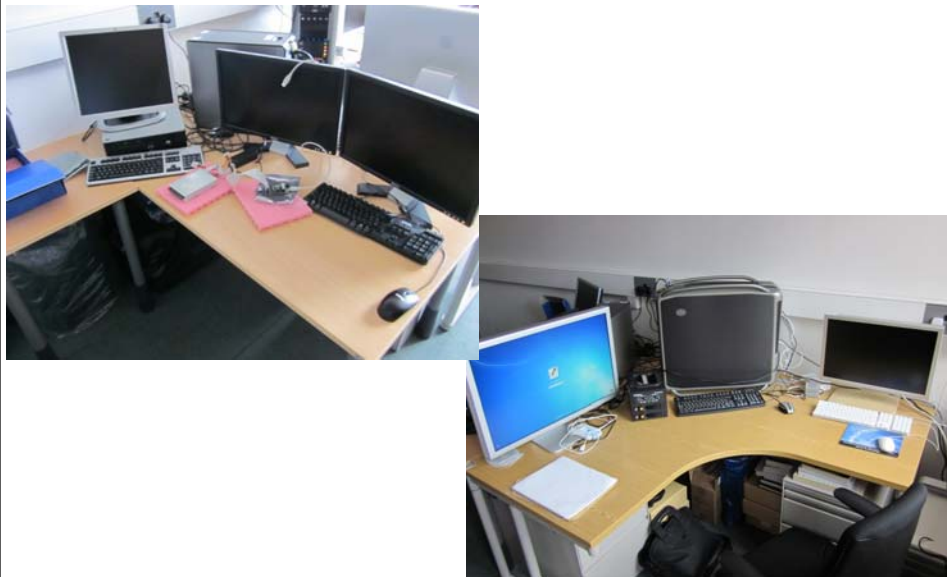**Here's what it looks like in libraries and archives:**

**12**

**Stanford University Libraries and Academic Information Resources (SULAIR)**



13

**British Library, London**



14

**UNC School of Information and Library Science**

15



USB 3.5" Floppy Disk Drive
Still available new from online retailers, look for a drive that can read both 1.44 MB(HD) and 800 KB (DD) 3.5" diskettes. Most drives support HD diskettes in both PC and Mac format, but only support PC formatted DD diskettes. New units are still available for around $20.

Device Side Data's FC5025
The FC5025 is a controller card for 5.25" floppy disk drives that can be used as an internal or external—as seen here—interface. Device Side Data charges $55.25 per controller.

5.25" Floppy Disk Drive
These units are no longer available new, but can still be purchased off of eBay for about $50. We recommend purchasing a number of drives as well as a floppy disk drive cleaning kit.

External USB 250MB Zip Drive
These units are available both new and used. We recommend the 250MB model as it is backwards compatible with the 100MB Zip disks. New units retail for around $200 and used units for around $50.

Wiebetech UltraDock
Hardware Write Protector
This unit serves as both an interface with IDE and Serial ATA type hard disk drives and as a write protector. Because it is common for the OS to overwrite metadata on a hard drive, write protection ensures that no interactions of the archivist or researcher affects the integrity of the original media. Wiebetech charges $250 for the UltraDock Hardware Write Protector.

http://www.bitcurator.net/building-a-digital-curation-workstation-with-bitcurator-update/

16

**Outfitting a Born-Digital Archives Program**
Ben Goldman, Penn State University



http://practicaltechnologyforarchives.org/issue2_goldman/          **17**

---

## Motivation

- Archivists are often responsible for acquiring or helping others access materials on removable storage media
- Information is often not packaged nor described as one would hope
- Information professionals must extract whatever useful information resides on the medium, while avoiding the accidental alteration of data or metadata

**18**

## Digital Forensics Can Help Archivists to Fulfill their Principles

| | |
|---|---|
| Provenance | • Identify, extract and save essential information about context of creation |
| Original Order | • Reflect original folder structures, files associations, related applications and user accounts |
| Chain of Custody | • Documentation of how records were acquired and any transformations to them<br>• Use well-established hardware and software mechanisms to ensure that data haven't been changed inadvertently |
| Identifying Sensitive Information | • Identify personally identifying information, regardless of where it appears<br>• Flag for removal, redaction, closure or restriction |

**19**

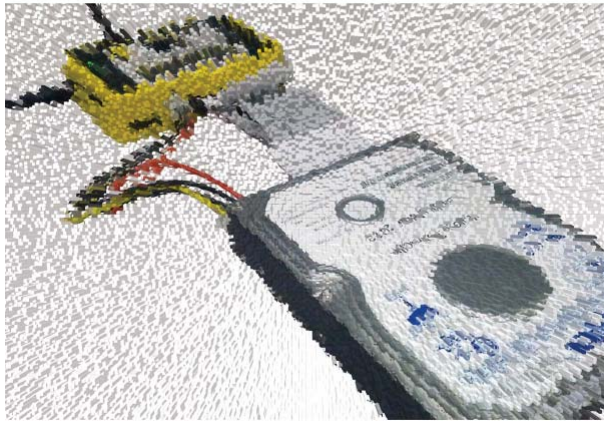## Applying Digital Forensics to Digital Collections – Previous Work*

- Ross and Gow (1999) - potential relevance of advances in data recovery and digital forensics to collecting institutions
- More recently - active stream of literature related to use of forensic tools and methods for digital collections, including activities at the British Library, National Library of Australia and Indiana University
- PERPOS (Georgia Tech) – has applied data capture and extraction to US presidential materials
- "Computer Forensics and Born-Digital Content in Cultural Heritage Collections" - symposium and report (2010)
- Born Digital Collections: An Inter-Institutional Model for Stewardship (AIMS) - framework for the stewardship of born-digital materials, including digital forensics methods
- Digital Records Forensics project - has articulated connections between the concepts of digital forensics and archival science

*See citations in:
http://ica2012.ica.org/files/pdf/Full%20papers%20upload/ica12Final00290.pdf  **20**

**From Bitstreams to Heritage:**

Putting Digital Forensics into Practice
in Collecting Institutions

Christopher A. Lee, Kam Woods, Matthew Kirschenbaum, and Alexandra Chassanoff

http://www.bitcurator.net/docs/bitstreams-to-heritage.pdf

21

---

## What this Course Covers

- Computational operations
- Layers of hardware and software that allow bitstreams on digital media to be read as files
- Roles and relationships of these layers
- Tools and techniques for ensuring completeness and evidential value of data

22

## Caveats and Such

- A vast space – we are only scratching the surface!

- Focus is on the foundational principles, methods and tasks that are applicable by a variety of tools

- This is a dynamic and evolving area, and these instructional materials evolve over time – your input is appreciated

**23**

## What is Digital Forensics (aka Forensic Computing)?

- "The process of identifying, preserving, analyzing and presenting digital evidence in a manner that is legally acceptable."*

- "Involves multiple methods of
  - □ **Discovering digital data (computer system, mobiles)**
  - □ **Recovering deleted, encrypted, or damaged file information**
  - □ Monitoring live activity
  - □ Detecting violations of corporate policy"**

*McKemmish, R. "What is Forensic Computing?" *Trends and Issues in Crime and Criminal Justice* 118 (1999).
**Brad Glisson, Introduction to Computer Forensics & E-discovery, University of Glasgow, Week 1 Lecture, September 2008.

**24**

## Why should we care about digital forensics

- **Not** because you're expected to solve crimes or catch malicious users
- Recognition of how data can be recovered when **layers** of technology fail or are no longer available
- **Capturing evidence** from places that are not always immediately visible
- Ensuring that actions taken on files **don't make irreversible changes** to essential characteristics (e.g. timestamps)
- Attending to the **order of volatility** – some types of data change much more quickly and often than others
- Learning about wide array of **tools and techniques** already available to deal with born-digital materials
- Established practices for **documenting** what we do, so others will know what we might have changed
- Considerable **overlap** between **technical knowledge** required to do digital forensics and ad hoc acquisition of digital materials by libraries/archives

25

## Digital Forensics vs. Intelligence Gathering vs. Electronic Discovery

Roughly in order of least to most targeted:

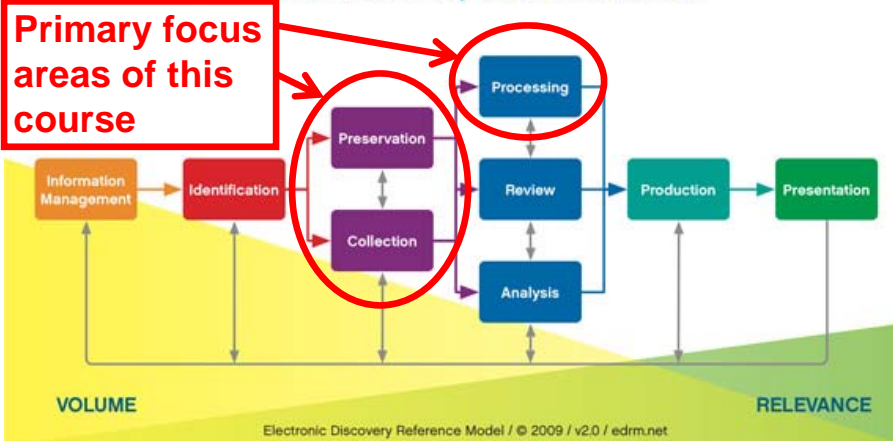| Activity | Main Emphasis | Common Scenario |
|---|---|---|
| Intelligence Gathering | Finding specific timely and relevant facts about target individuals or organizations | Seize whole physical medium or covertly collect data; systematically search and analyze for bits of data and interesting patterns |
| Digital Forensics | Obtaining evidence in order to solve or prove a specific crime | Seize whole physical medium or intervene into a live system to capture data; prove chain of custody and evidential value at bit level; search for offending or incriminating data, often within "hidden" areas |
| Electronic Discovery | Identifying and collecting documents relevant to a specific legal claim or dispute | Plaintiff makes explicit requests for specific types of information; issue queries that reflect the specific requests; prove chain of custody and evidential value at document and procedural level; parties share results |

26

## Common Digital Forensics Scenarios

- Evidence seized from home/office of "person of interest" in a criminal investigation (dead forensics)
- Response to system security breach, to determine what was done, by whom and how (live forensics)

27

---

**Electronic Discovery Reference Model**

**Primary focus areas of this course**

Information Management → Identification → Preservation → Processing

Collection → Review → Analysis → Production → Presentation

VOLUME

RELEVANCE

Electronic Discovery Reference Model / © 2009 / v2.0 / edrm.net

28

# Technical Background

# Nature of Digital Materials

# Digital objects are sets of **instructions for future interaction**

- Digital objects are useless (and don't even exist) if no one can interact with them
- Interactions depend on numerous technical components

---

"Errors typically occur at the juncture between analog and digital states, such as when a drive's magnetoresistive head assigns binary symbolic value to the voltage differentials it has registered, or when an e-mail message is reconstituted from independent data packets moving across the TCP/IP layer of the Internet, itself dependent on fiber-optic cables and other hardwired technologies. All forms of modern digital technology incorporate **hyper-redundant error-checking** routines that serve to sustain an **illusion of immateriality** by detecting error and correcting it, reviving the quality of the signal, like old-fashioned telegraph relays, such that any degradation suffered during a subsequent interval of transmission will not fall beyond whatever **tolerances of symbolic integrity** exist past which the original value of the signal (or identity of the symbol) cannot be reconstituted."

Kirschenbaum, Matthew G. *Mechanisms: New Media and the Forensic Imagination.* Cambridge, MA: MIT Press, 2008. p.12 (emphasis mine).

## Translation Across Layers

- Users view, read, write and click on things
- Programmers usually write & reuse source code
- Software & firmware manipulates data and instructions as bits (10100001110101)
- Physical equipment deals with magnetic charges, holes in optical disks, holes in punch cards

33

## Digital Resources - Levels of Representation

| Level | Label | Explanation |
|---|---|---|
| 8 | Aggregation of objects | Set of objects that form an aggregation that is meaningful encountered as an entity |
| 7 | Object or package | Object composed of multiple files, each of which could also be encountered as individual files |
| 6 | In-application rendering | As rendered and encountered within a specific application |
| 5 | File through filesystem | Files encountered as discrete set of items with associate paths and file names |
| 4 | File as "raw" bitstream | Bitstream encountered as a continuous series of binary values |
| 3 | Sub-file data structure | Discrete "chunk" of data that is part of a larger file |
| 2 | Bitstream through I/O equipment | Series of 1s and 0s as accessed from the storage media using input/output hardware and software (e.g. controllers, drivers, ports, connectors) |
| 1 | Raw signal stream through I/O equipment | Stream of magnetic flux transitions or other analog electronic output read from the drive without yet interpreting the signal stream as a set of discrete values (i.e. not treated as a digital bitstream that can be directly read by the host computer) |
| 0 | Bitstream on physical medium | Physical properties of the storage medium that are interpreted as bitstreams at Level 1 |

34

# Interaction Examples

**Level**

| Aggregation of objects |
| --- |
| Object or package |
| In-application rendering |
| File through filesystem |
| File as "raw" bitstream |
| Sub-file data structure |
| Bitstream through I/O equipment |
| Raw signal stream through I/O equipment |
| Bitstream on physical medium |



**35**

---

# Interaction Examples

**Level**

| Aggregation of objects |
| --- |
| **Object or package** |
| In-application rendering |
| File through filesystem |
| File as "raw" bitstream |
| Sub-file data structure |
| Bitstream through I/O equipment |
| Raw signal stream through I/O equipment |
| Bitstream on physical m... |



**36**

## Slide 37

### Interaction Examples

**Level**

| |
|---|
| Aggregation of objects |
| Object or package |
| **In-application rendering** |
| File through filesystem |
| File as "raw" bitstream |
| Sub-file data structure |
| Bitstream through I/O equipment |
| Raw signal stream through I/O equipment |
| Bitstream on physical medium |



37

## Slide 38



**Level**

| |
|---|
| Aggregation of objects |
| Object or package |
| In-application rendering |
| **File through filesystem** |
| File as "raw" bitstream |
| Sub-file data structure |
| Bitstream through I/O equipment |
| Raw signal stream through I/O equipment |
| Bitstream on physical medium |

38

# Interaction Examples

**Level**

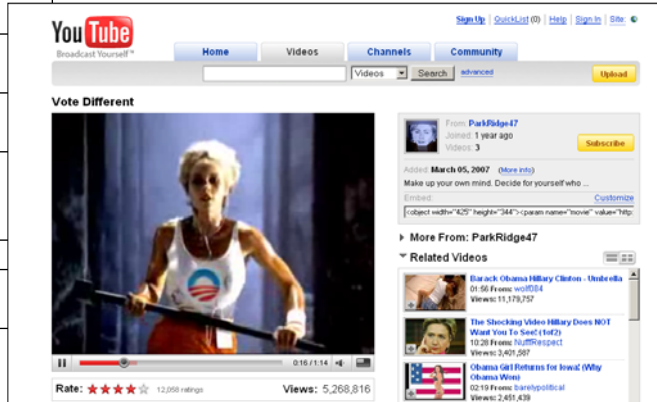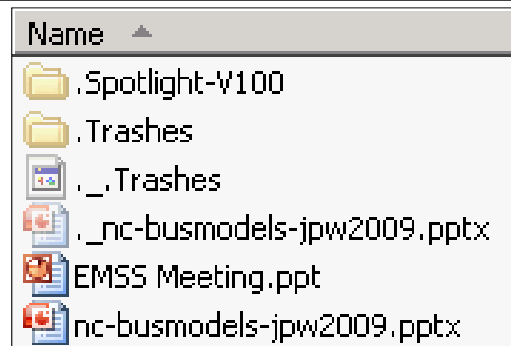| Level |
|---|
| Aggregation of objects |
| Object or package |
| In-application rendering |
| File through filesystem |
| **File as "raw" bitstream** |
| Sub-file data structure |
| Bitstream through I/O equipment |
| Raw signal stream through I/O equipment |
| Bitstream on physical medium |

HView 2000
File  Edit  Window  Help

G:\._nc-busmodels-jpw2009.pptx

```
00000000: 00 05 16 07 00 02 00 00 4D 61 63 20 4F 53 20 58 ........Mac OS X
00000010: 20 20 20 20 20 20 20 20 00 02 00 00 09 00 00 00  ........
00000020: 00 32 00 00 0E B0 00 00 00 02 00 00 0E E2 00 00 .2........
00000030: 01 1E 50 50 54 58 50 50 54 33 00 00 00 00 00 00 ..PPTXPPT3......
00000040: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
00000050: 00 00 00 00 00 41 54 54 52 3B 9A C9 FF 00 00 0E E2 ....ATTR;.......
00000060: 00 00 00 00 78 00 00 00 00 00 00 00 00 00 00 00 ....x........
00000070: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
00000080: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
00000090: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
000000A0: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
000000B0: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
000000C0: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
000000D0: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
000000E0: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
000000F0: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
00000100: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
00000110: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
00000120: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
00000130: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
00000140: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
00000150: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
00000160: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
00000170: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
00000180: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
00000190: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
000001A0: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 ................
000001B0: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 
```

DWord: 118883584   Word: 1280   Byte: 0   Position: 00000000   Size: 00001000

**39**

---

# Interaction Examples

**Level**

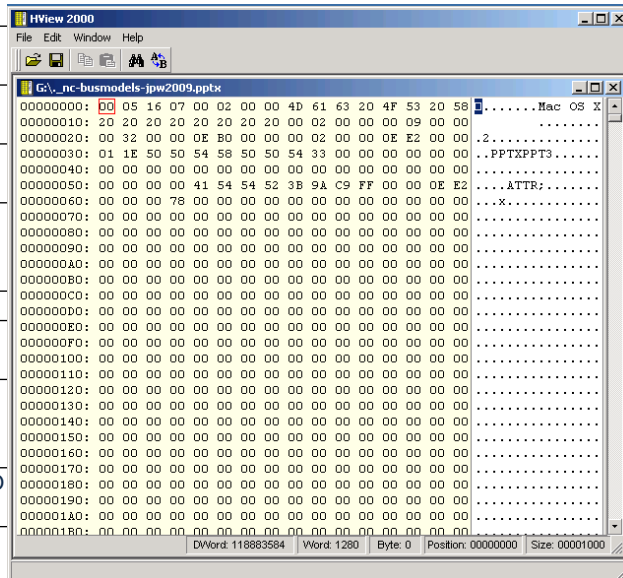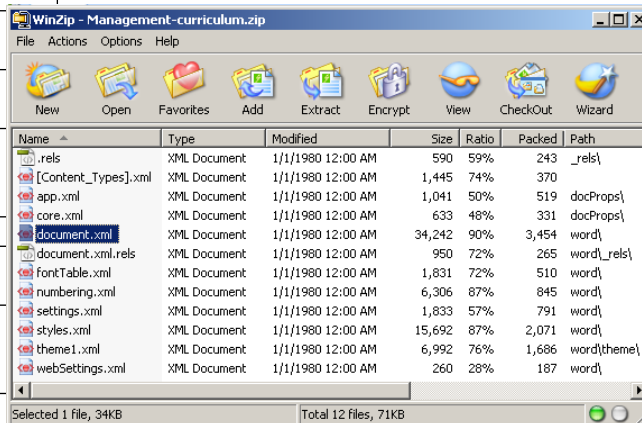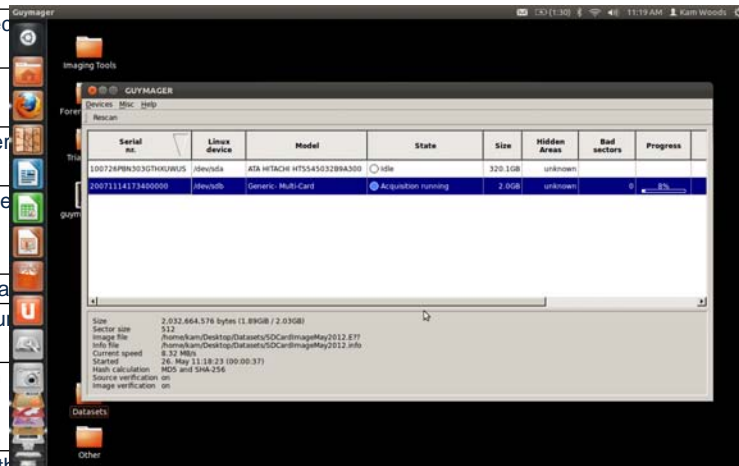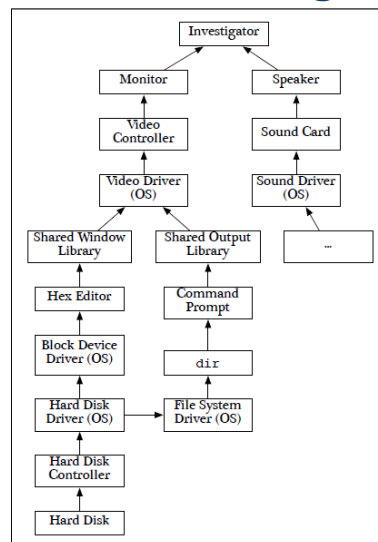| Level |
|---|
| Aggregation of objects |
| Object or package |
| In-application rendering |
| File through filesystem |
| File as "raw" bitstream |
| **Sub-file data structure** |
| Bitstream through I/O equipment |
| Raw signal stream through equipment |
| Bitstream on physical medium |

WinZip - Management-curriculum.zip
File  Actions  Options  Help

New   Open   Favorites   Add   Extract   Encrypt   View   CheckOut   Wizard

| Name | Type | Modified | Size | Ratio | Packed | Path |
|---|---|---|---|---|---|---|
| .rels | XML Document | 1/1/1980 12:00 AM | 590 | 59% | 243 | _rels\ |
| [Content_Types].xml | XML Document | 1/1/1980 12:00 AM | 1,445 | 74% | 370 | |
| app.xml | XML Document | 1/1/1980 12:00 AM | 1,041 | 50% | 519 | docProps\ |
| core.xml | XML Document | 1/1/1980 12:00 AM | 633 | 48% | 331 | docProps\ |
| document.xml | XML Document | 1/1/1980 12:00 AM | 34,242 | 90% | 3,454 | word\ |
| document.xml.rels | XML Document | 1/1/1980 12:00 AM | 950 | 72% | 265 | word\_rels\ |
| fontTable.xml | XML Document | 1/1/1980 12:00 AM | 1,831 | 72% | 510 | word\ |
| numbering.xml | XML Document | 1/1/1980 12:00 AM | 6,306 | 87% | 845 | word\ |
| settings.xml | XML Document | 1/1/1980 12:00 AM | 1,833 | 57% | 791 | word\ |
| styles.xml | XML Document | 1/1/1980 12:00 AM | 15,692 | 87% | 2,071 | word\ |
| theme1.xml | XML Document | 1/1/1980 12:00 AM | 6,992 | 76% | 1,686 | word\theme\ |
| webSettings.xml | XML Document | 1/1/1980 12:00 AM | 260 | 28% | 187 | word\ |

Selected 1 file, 34KB        Total 12 files, 71KB

**40**

## Interaction Examples

**Level**

| |
| --- |
| Aggregation of obje... |
| Object or package |
| In-application render... |
| File through filesyste... |
| File as "raw" bitstrea... |
| Sub-file data structu... |
| **Bitstream through I/O equipment** |
| Raw signal stream through I/O equipment |
| Bitstream on physical medium |

Guymager / GUYMAGER — Devices  Misc  Help — Rescan

| Serial nr. | Linux device | Model | State | Size | Hidden Areas | Bad sectors | Progress |
| --- | --- | --- | --- | --- | --- | --- | --- |
| 100726PBN303GTHXUWUS | /dev/sda | ATA HITACHI HTS54503289A300 | Idle | 320.1GB | unknown | | |
| 20071114173400000 | /dev/sdb | Generic- Multi-Card | Acquisition running | 2.0GB | unknown | 0 | 8% |

Size 2,032,664,576 bytes (1.89GiB / 2.03GB)
Sector size 512
Image file /home/kam/Desktop/Datasets/SDCardImageMay2012.E??
Info file /home/kam/Desktop/Datasets/SDCardImageMay2012.info
Current speed 8.32 MB/s
Started 26. May 11:18:23 (00:00:37)
Hash calculation MD5 and SHA-256
Source verification on
Image verification on

Imaging Tools — Datasets — Other

41

---

**Level**

| |
| --- |
| Aggregation of ... |
| Object or pack... |
| In-application r... |
| File through file... |
| File as "raw" bi... |
| Sub-file data st... |
| Bitstream throu... equipment |
| **Raw signal stream through I/O equipment** |
| Bitstream on physical medium |

KryoFlux

Tracks — 0 1 2 3 4 5 6 7 8 9     0 1 2 3 4 5 6 7 8 9 — Side 0 — Side 1

Information — 8.0 7.5 7.0 6.5 6.0 5.5 5.0 4.5 4.0 — Track  Summary  Hist  Scatter

Control — Motor  Stream  Error — Batman — AmigaDOS sector image — Start

42

21

## Interaction Examples

| Level |
| --- |
| Aggregation of objects |
| Object or package |
| In-application rendering |
| File through filesystem |
| File as "raw" bitstream |
| Sub-file data structure |
| Bitstream through I/O equipment |
| Raw signal stream through I/O equipment |
| **Bitstream on physical medium** |



Veeco Instruments. http://www.veeco.com/library/nanotheater_detail.php?type=application&id=78&app_id=34

**43**

---

## Multiple Paths for Viewing Bits



Carrier, Brian D. "A Hypothesis-Based Approach to Digital Forensic Investigations." Doctoral Dissertation, Purdue University, 2006. Figure 3-3 (p.60)

**44**

## Three Complicating Factors for Archivists:

1. **Medium Failure / Bit Rot**

2. **Obsolescence**

3. **Volatility**

45

## Bit Rot

- Preventing measures can help (proper storage and handling), but bits on a given medium will eventual flip or become unreadable
- In repositories
  - We maintain integrity of bit stream through security, checksums, periodic sampling and other validation
  - Bit rot and advantages of newer media both call for periodic refreshing and reformatting
- But:
  - The media we receive may not be so well maintained
  - Ensuring the **integrity of the bit stream** when transferring from one medium to another is extremely important

46

## Obsolescence

"Obsolete power corrupts obsoletely."

- Ted Nelson

The technology associated with interpreting the representation at each of the layers can change or become less available

47

## Order of Volatility

- Some types of data change much more quickly and often than others
- Important to recognize in order to recover data from a computer system or media, while ensuring that actions don't make irreversible changes to their record characteristics
- Example: If the contents of the browser cache are important to you, capture the cache before using the browser

48

## How and where does a computer store information?

## Bits – How Data are Conveyed in Computers

- Variable voltage electrical signals or pulses of light
- Bit represents a tiny "switch" with two possible states – on/off, true/false, 1/0
- Bit string or bitstream: a consecutive sequence of bits (e.g. 101000111010101)
- Rarely meaningful to humans – when looking at bitstream, usually use a hex editor (discussed later)

## Motivations for Storage Hierarchy

- Different forms of memory/storage have significantly different costs and performance
- Store recent data close by, in fast, expensive, volatile storage
- Store data that has not been used recently and is rarely used in slower, cheaper, less volatile storage

**51**

---

## Computer Memory Hierarchy

| | | |
|---|---|---|
| small size small capacity | | processor registers very fast, very expensive |
| | power on immediate term | |
| small size small capacity | | processor cache very fast, very expensive |
| medium size medium capacity | power on very short term | random access memory fast, affordable |
| small size large capacity | power off short term | flash / USB memory slower, cheap |
| large size very large capacity | power off mid term | hard drives slow, very cheap |
| large size very large capacity | power off long term | tape backup very slow, affordable |

Source: http://en.wikipedia.org/wiki/File:ComputerMemoryHierarchy.svg

**52**

## The Low-Level Building Blocks of Storage – Sectors and Clusters

- Your computer's processor manipulates data in the form of bitstreams, and data is stored on your computer's hard drive as bitstreams
- But moving the data from the hard drive to the processor depends on higher-level chunks: sectors and clusters
- Think of mail sent to a member of a family who all live in the same house – the envelope will indicate the house address but won't identify where that person's bedroom is located within the house

53

## Sectors

- Smallest unit of storage that can be assigned an address (i.e. can be directly identified & found by the computer system)
- Have specified size, depending on the type of storage, e.g.
  - CD-ROM = 2048 bytes (2,352 including error checking)
  - floppies (usually) = 512 bytes
  - modern hard drives = 4,096 (previously 512 bytes)
- Created when disk is low-level formatted (usually by manufacturer) with bad sectors identified by disk controller so data won't be written to them

54

## Clusters

- Groups of sectors
- Smallest unit of storage that can be tracked by the operating system
- Sizes depends on operating system, type & size of storage device – examples are 2048 bytes (4 sectors of 512 bytes) or 4096 bytes
- Defined during high-level formatting performed by operating system

55

## Magnetic Disk (e.g. Hard Drive or Floppy)

- Bits stored as magnetic fields of different polarity
- Magnetized surface of disk rotates under a read/write head
- Divided into tracks (like rings of a tree)
- Tracks divided into sectors and clusters
- Windows: File Allocation Table (FAT) or Master File Table (for NTFS) indicates, for given file, what clusters contain its content

Main spindle
Head 0
Side 0
Platter 1
(has sides 0-1)
Arm for head 1
Head stack
assembly
Head 2
Arm for
Tracking/Alignment head (head 3)

Image from : "Concepts." In Active UNDELETE v2.0 Documentation. Active Data Recovery Software. www.active-undelete.com/3tracks.htm

56

Hard Drive Structure:

A = track
B = sector
C = sector of a track
D = cluster

Source: http://en.wikipedia.org/wiki/File:Disk-structure2.svg

**57**

# Optical Media – CD-ROM as Example

Label layer
Protective layer
Reflective coating
Land
Pit
Laser beam
Substrate (PC)
Laser spot
Pit
Land
Light intensity

Source of Images: Compact Disk (CD). USByte.
http://www.usbyte.com/common/compact_disk_3.htm

**58**

## Solid-State Drives (SSDs)

Source:
http://www.tomshardware.com/gallery/Samsung-SSD-256-ToggleDDR,0101-260898-0-0-0-0-jpg-.html

- Uses integrated circuits to store data
- No moving parts
- Can be read using same I/O equipment as used for hard drives
- Increasingly common in laptops

Source:
http://www.tomshardware.com/gallery/ssd-controller-external-cache,0101-260900-0-0-0-0-jpg-.html

Host

SATA II | 3 or 6 Gb/s

SATA
Core          Cache
Ch-1  Ch-2  Ch-3  Ch-n

Flash   Flash   Flash   Flash

**59**

# Floppy Disks

- Physical storage is similar to hard drives described above (magnetic charges in a spinning disk)
- Various types and sizes, e.g. high density, double density, 3.5 inch, 5.25 inch, 8 inch
- 3.5 inch floppies are relatively easy to read using a USB drive, but older ones are more complicated…

**60**

## Floppy Controller Hardware

CatWeasel[1] (no longer available)     Disc Ferret[2]          Kryoflux[3]



FC 5025[4]          Disk2FDI[5]          SuperCard Pro[6]

1. http://lib.stanford.edu/digitial-forensics-stanford-university-libraries/catweasel-universal-floppy-drive-controller
2. http://discferret.com/wiki/DiscFerret
3. http://www.kryoflux.com/
4. http://www.deviceside.com/fc5025.html
5. http://disk2fdi.joguin.com/D2FCABLE.htm
6. http://www.cbmstuff.com/proddetail.php?prod=SCP

61

---

## Two Important Considerations for Internal Media that are Used as External Media

- Power - internal drive needs different connector (often Molex), not the kind that plugs into the wall



http://en.wikipedia.org/wiki/File:Molex_female_connector.jpg

- Cooling – when pulled from the computer, you've also separated the drive from the fan, so you should often add an external one to ensure cooling



http://www.tigerdirect.com/applications/SearchTools/item-details.asp?EdpNo=1648567

62

Kryoflux Running on a "Mini JukeBox"*



*Adapted from a Mini JukeBox setup designed by the National Library of Australia [63]

# Areas Designed to Store Temporary Data

- Files on disk used for virtual memory management – e.g. "swap files" in Windows 95/98, "page files" in Windows NT/2000/XP
- Temp files
- Various caches - e.g. browser cache, which includes copies of recently downloaded files
- "Recent documents" in Windows
- Cookies – "expires" attribute can indicate quick deletion or long-term retention
- History files – e.g. browsing & download history

[64]

# Caching

- Storing a copy of a subset of data from a slower data source to a faster (more readily available) data source
- Examples:
  - CPU cache from main memory
  - Main memory cache from hard disk
  - Hard disk cache from CD-ROM
  - Proxy server cache from web sites

65

# Configuration and Log Files

- Often contain information about where files are located, when last opened, user preferences, state of files when last used
- In Windows, much of this happens in the Registry
- On a Mac, much of this happens in property list (p-list) files
- Another examples:
  - Index.dat – RSS feeds, URLs visited, search queries and recently opened files in Internet Explorer

66

## Windows Registry

- Information about:
  - □ Applications installed
  - □ Application settings
  - □ Hardware installed
  - □ Hardware settings
  - □ User interface and system preferences
  - □ User accounts
  - □ Locations of files and recent activities, e.g. Most Recently Used (MRU)
  - □ Lots of online activities, e.g. user names and passwords, browsing and search query history

**67**

---

# Representation Information

**68**

# "No computation without representation"

Smith, Brian Cantwell. "Limits of Correctness in Computers." In *Computerization and Controversy: Value Conflicts and Social Choices*, edited by Rob Kling, 810-25. San Diego, CA: Academic Press, 1996. 815.

69



The 8 bits highlighted in the bit stream shown below can be interpreted in many ways, e.g., as an integer, a simple character code, a sound, a floating point number, an image, a logical bitmap, etc.

character
"U"

integer
"21"

sound

bit stream: 010111000000**00010101**000000000100000111101110

real number
"1.3125"

logical bitmap
"no, no, no, yes"
no, yes, no, yes

image

**Figure 4: A bit stream can represent anything at all**

Rothenberg, Jeff. "Ensuring the Longevity of Digital Information." Washington, DC: Council on Library and Information Resources, 1999.

70

35

# Representation Information

- "Information that maps a Data Object into more meaningful concepts" (OAIS) - makes humanly-perceptible properties happen
- Examples: file format, encoding scheme, data type



Reference Model for an Open Archival Information System (OAIS). CCSDS 650.0-M-2 (Magenta Book). Consultative Committee for Space Data Systems, 2012. [ISO 14721:2012]. Figure 2-2.

71



Reference Model for an Open Archival Information System (OAIS). CCSDS 650.0-M-2 (Magenta Book). Consultative Committee for Space Data Systems, 2012. [ISO 14721:2012]. Figure 4-10.

72

Reference Model for an Open Archival Information System (OAIS). CCSDS 650.0-M-2 (Magenta Book). Consultative Committee for Space Data Systems, 2012. [ISO 14721:2012]. Figure 4-10.

**73**

---

## Representation Information can Reside in Many Places

- Within digital object itself

- Stored separately as metadata

- Encoded within software required to read and parse the digital object

**74**

## Finding Representation Information Within a File

- Keys fields

- Headers

- Manifests

75

---

The 4 bits at the start of this bit stream are intended to be read as the "key" integer 7, meaning that the remaining bytes in the bit stream are each 7 bits long. However, there is no way to tell from the bit stream itself how long the key integer is; if we were to erroneously read the first 5 bits of the bit stream as the key (instead of the first 4), we would erroneously conclude that the remaining bytes were each 15 bits long.

Intended 4-bit key
(value of 0111 = 7)

Intended
7-bit data bytes

bit stream: 011111100000000010101000000000100000111101110

Unintended 5-bit key
(value of 01111 = 15)

Unintended
15-bit data byte

**Figure 6: Bit streams cannot be made self-explanatory**

Rothenberg, Jeff. "Ensuring the Longevity of Digital Information." Washington, DC: Council on Library and Information Resources, 1999.

76

## Not Just a Series of Bytes – Pointers and Offsets

- Pointer – reference within a file or programming code that leads from one place to another
  - Causes the data to be read out of serial order (i.e. a jump from one place to another place that does not immediately follow it within the data stream)
  - Ability to resolve the pointer is essential

- Offset – location that's some given distance from a starting point
  - Location calculated by adding offset to a base address (location)
  - Again, ability to resolve the offset to the precise location within a bitstream is essential

77

## Fonts and Character Encoding

78

## Font



PDF with missing fonts by prwheatly, on Flickr.

## Font

- Determines how characters will appear on screen (generation of glyphs)
- Same character can appear completely different in two different fonts
- Can be a major issue in digital preservation, when consistent rendering is important
- Not usually a focus of digital forensics or data recovery efforts, which focus on simply making sense of the characters within a bitstream
- However, changes of fonts within a document can provide some hints to versioning and authorship

**Character Encoding**

---

## ASCII – Major "installed base" of Character Encoding, Designed for the English-Speaking World

| Symbol | Decimal | Binary |
|--------|---------|----------|
| 7      | 55      | 00110111 |
| 8      | 56      | 00111000 |
| 9      | 57      | 00111001 |
| :      | 58      | 00111010 |
| ;      | 59      | 00111011 |
| <      | 60      | 00111100 |
| =      | 61      | 00111101 |
| >      | 62      | 00111110 |
| ?      | 63      | 00111111 |
| @      | 64      | 01000000 |
| A      | 65      | 01000001 |
| B      | 66      | 01000010 |
| C      | 67      | 01000011 |

## Unicode

- **Huge** number of possible characters – not limited to 8-bits for each

- Mapped to unique codes (numbers)

- Standard first published in 1991

- Current version is 8.0 (June 2015)

83

## UTF-8

- *Unicode Transformation Format (UTF)* = set of conventions for how specific Unicode code points are represented as unique byte sequences

- UTF 8 is widely used – including in email and web pages

- Codes 0 to 127 are backward compatible with ASCII

See: Frequently Asked Questions: UTF-8, UTF-16, UTF-32 & BOM. Unicode, Inc.
http://www.unicode.org/faq/utf_bom.html

84

## Escape Codes and Character Entities

- When a system doesn't allow use of certain characters (either because reserved for special uses or because not allowed at all), must do a translation to characters that it does allow

- Examples

    □ In programming languages

    □ In HTML – use & or % convention

    □ In URLs - Use of "%" + hexadecimal label

85

| Special-Use or Disallowed Character | Hex Replacement | Notes |
| --- | --- | --- |
| Space | %20 | Extremely common when posting to Web from OS that allows white spaces in file names |
| " | %22 | |
| # | %23 | Used as "anchor" within URLs (link to specific section of page) |
| $ | %24 | |
| % | %25 | Imagine what problems this might cause! |
| & | %26 | Used within URL to separate query parameters |
| + | %2B | |
| , | %2C | |
| / | %2F | Used as separator between parts of a URL or directory path |
| : | %3A | |
| ; | %3B | |
| < | %3C | Can appear when XML /HTML markup gets passed as part of URL |
| = | %3D | Used within URL to assign parameter value |
| > | %3E | Can appear when XML /HTML markup gets passed as part of URL |
| ? | %3F | Used within URL to indicate query parameters |
| @ | %40 | Often appears as part of email address |

86

43

# Compression

---

As a simple example of compressing a bit stream without loss, "run-length encoding" replaces each sequence of 0s (000...0) by a count, indicating how many 0 bits were present in the given "run" (similarly for 1s). This can reduce the size of a bit stream without losing any information. For example, each run in the original bit stream shown can be represented by a 5-bit byte whose first bit specifies whether the run is of 0s or 1s and whose remaining 4 bits specify the length of a run (of up to 15 bits). This scheme is most appropriate for data that contains long sequences of 0s or 1s, such as digital imagery.

original bit stream: 000000111111111111110000000000000111111111 (42 bits)

a run of 6 0-bits          a run of 14 1-bits

Representing each run in the original bit stream as a pair **b:n** (where b is 0 or 1 to indicate the contents of the run, and n is the length of the run) produces:

sequence of runs:          0:6, 1:14, 0:13, 1:9

resulting 5-bit bytes:     00110, 11110, 01101, 11001

compressed bit stream:     00110111100110111001   (20 bits)

**Figure 7: Compressing a bit stream**

Rothenberg, Jeff. "Ensuring the Longevity of Digital Information." Washington, DC: Council on Library and Information Resources, 1999.

## Three Levels of Compression*

- Format of file implements compression internally - e.g. body of JPEG file is compressed but not file header
- Application creates completely new, compressed copy of file(s) – e.g. WinZip, gzip
- File system compresses data units – e.g. not writing data to series of sectors that are all filled with zeros

*Carrier, Brian. *File System Forensic Analysis. Boston, MA: Addison-Wesley, 2005.*

**89**

---

## Encryption

**90**

# Encryption

- Special data ("keys") and algorithms used to transform data into a form that is purposely less easily readable
- Used for:
  - Confidentiality
  - Integrity
  - Non-repudiation
  - Authentication

**91**

# Encryption at Various Levels*

- Application that creates the file
- Application that reads an unencrypted file and creates an encrypted file
- Operating System – "Before a file is written to disk, the OS encrypts the file and saves the cipher text to the data units. The non-content data, such as the file name and last access time, are typically not encrypted. The application that wrote the data does not know the file is encrypted on the disk."
- Encrypt an entire volume – implemented in storage system below file system level

*Carrier, Brian. *File System Forensic Analysis. Boston, MA: Addison-Wesley, 2005.*

**92**

## Checksums – Compact Representations of Bitstreams

- A given bitstream, fed into an algorithm, will generate a short string of characters that is **<u>extremely</u>** unlikely to be generated by a different bistream fed into that same algorithm
- Most common = MD5, SHA-1
- Can determine:
  - ☐ If bits have changed after a transfer
  - ☐ If bits have flipped within a storage environment
  - ☐ Whether two different files are identical bitstreams
- A library of hash values can identify "known and notable" (EnCase terminology) files
  - ☐ Known – files that can be ignored (e.g. software listed in National Software Reference Library)
  - ☐ Notable – specific bitstreams that you're trying to find

**93**

## Checksums – Compact Representations of Bitstreams

- Tools for checksum generation
  - ☐ MD5Summer
  - ☐ HashDeep
  - ☐ MD5Deep
  - ☐ Fileverifier++
  - ☐ FF MD5Drop or command-line tool (Mac)
  - ☐ GtkHash (available in BitCurator environment)
  - ☐ Many others…

**94**

In BitCurator environment: Right Click on File or Directory and Calculate MD5



95



96

MD5 Hashes of an Entire Directory of Files

## Hexadecimal Notation

- A more compact and more humanly readable way of representing a stream of bits
  - Each character represents one of 16 possible values (0,1,2,3,4,5,6,7,8,9,A,B,C,D,E,F)
  - Conveniently, a series of two characters represented in hexadecimal can represent exactly one byte ($2^8$ = 256 possible values) of data, because $16^2 = 256$
- Hex dumps from computer's memory often used for debugging or reverse engineering software and for data recovery

**99**

## How to Generate a Hex Dump

- Many free or inexpensive tools available for download, e.g. Cygnus Hex Editor, Hex Workshop, HexAssistant, HxD, Hex Fiend (Mac), GHex (Linux), MiniDumper*
- To generate your own hex dump from a given file, try: http://www.fileformat.info/tool/hexdump.htm
- Hex viewing will usually include a separate view to the right that presents the ASCII equivalent of all bytes, which can help the human eye to detect patterns
- Hex viewing only necessary when a file includes either non-ASCII strings of bits or corrupted file elements
- If file is composed completely of ASCII-encoded data, using a simple text editor (e.g. Notepad) is simpler way to view file contents

* See http://en.wikipedia.org/wiki/Comparison_of_hex_editors          **100**

## Syllabus for a class (HTML File)



Beginning of file tells us what kind of file it is

**101**

## Slides from a lecture (PDF/A file)



Contents of this PDF file are not as easy to read within ASCII view as the contents of the HTML file were, but note that, again, beginning of file tells us what kind of file it is

**102**

In the BitCurator environment:

Let's corrupt a bitstream.

107



108

**113**



**114**

115



116

117



118

119



120

**That doesn't look right.**
**Let's compare it to our previous MD5 hash...**

| File Name | Change Made | MD5 Hash |
|---|---|---|
| powerpoint-to-break.ppt | ------ | a5feaf7d5b7d107e1805c2d223bce6e4 |
| powerpoint-to-break.ppt | **One Bit Different** From: Character "C" (Hex = 43) To: Character "D" (Hex = 44) | b6035b3e2048973a208666a519165288 |

| File Name | Change Made | MD5 Hash |
|---|---|---|
| powerpoint-to-break.ppt | ------ | a5feaf7d5b7d107e1805c2d223bce6e4 |
| powerpoint-to-break.ppt | **One Bit Different** From: Character "C" (Hex = 43) To: Character "D" (Hex = 44) | b6035b3e2048973a208666a519165288 |

**Note: A 1-byte change resulted in a completely different MD5 hash of the file.**

# Now do it yourself!

# Get yourself a Microsoft Word File:

- Create a new folder on your desktop called dfa-test
- Find a file by searching in Google on a topic of interest along with "filetype:doc" or "filetype:docx"
- Save the file to your dfa-test folder

## Generate a Hash

- Open FileVerifier++
- Click on the "Options" button
- If the "Default Algorithm" is not set to MD5, then change it to MD5 and select "Apply"
- Click "Ok"
- Close the program and then launch it again (ensures that the settings have been changed)
- Click on the "Dirs" button
- Navigate to the dfa-test folder on your desktop
- Click "Ok"
- You should now see the file path and associated hash for the file that you downloaded
- Leave this application running

**127**

## Breaking the File and then Fixing It at the Bitstream Level

- **Open the file in HxD or Hex Fiend** (use File > Open or drag the file onto the HxD icon on your desktop).
- **Change a byte** within the file in HxD or Hex Fiend and then save the changed file (note the **specific place** that you changed)
- **Re-Verify hash values**
    - Exit out of HxD or Hex Fiend
    - In FileVerifier++, again click on "Verify All"
    - What status do you now see in the Verification column?
- Use HxD or Hex Fiend to **change the byte back** to its earlier state and save it
- **Re-Verify hash values**
    - Go back to FileVerifier++
    - Again click on "Verify All" and "OK"
    - What status do you now see in the Verification column?

**128**

## Breaking the File and then "Fixing" It Using MS Word

- **Open the file in MS Word**
- **Change one word** within the file's content (note the **specific place** that you changed) and save it
- **Re-Verify hash values**
  - □ Exit out of MS Office
  - □ In FileVerifier++, again click on "Verify All"
  - □ What status do you now see in the Verification column?
- Use PowerPoint/Word to **change the text back** to its earlier state and save it
- **Re-Verify hash values**
  - □ Go back to FileVerifier++
  - □ Again click on "Verify All" and "OK"
  - □ What status do you now see in the Verification column?

**129**

## Note on MD5 – Potential Collisions

- From a security perspective, MD5 has been "broken" since 2005
- Someone with malicious intent can create two different bitstreams that result in the same MD5 hash (i.e. MD5 collisions)

**130**

## Implications of MD5 Being "Broken"

- Rarely a concern when MD5 is used for integrity checks on known items (e.g. verifying that a file was transferred correctly to a repository or that files in storage are still intact)
- Can be a concern if one is relying on a hash as proof of record authenticity – risks can include cases of internal tampering
- There are more robust hashes to address this (SHA-2)
- MD5 is still widely used, because it is fast to calculate and widely supported

**131**

## What have we covered so far?

- ✓ **Welcome and introductions**
- ✓ **Motivation and scope**
- ✓ **Technical background**
- ✓ **Representation Information**
4. File systems and file management
5. Extracting data from media
6. Tools and methods
7. Conclusions, questions, discussion

**132**

**File systems and file management**

**How do computers store and manage files?**

## Volumes and Partitions

- **Volume**
  - Storage area defined at the logical OS level, which has a single filesystem & usually resides on one disk partition
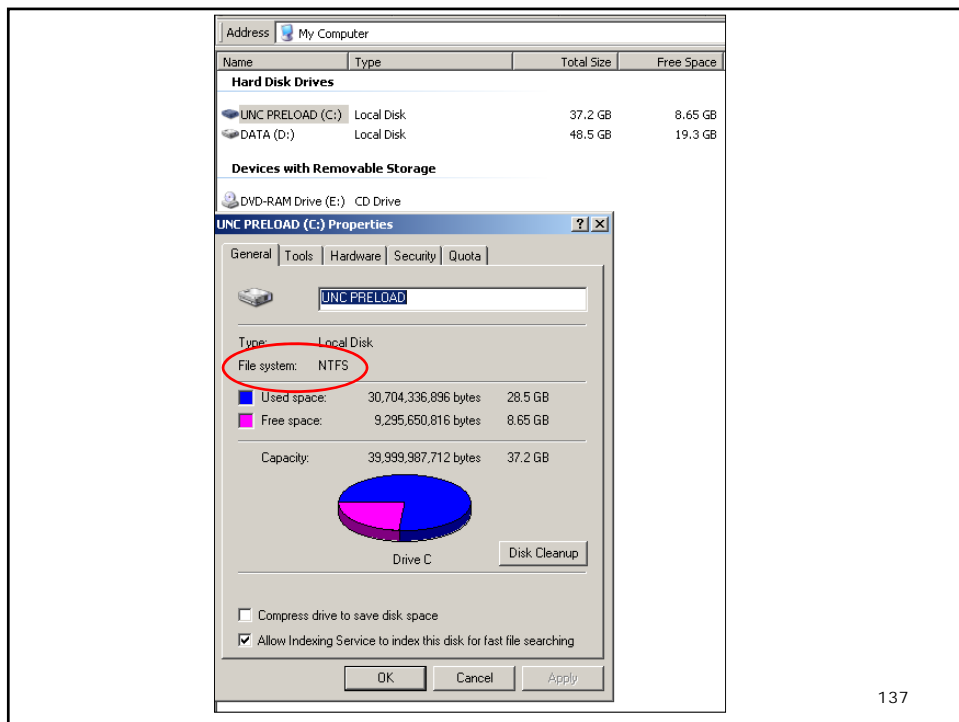- **Partition**
  - Exists at physical, media-specific level
  - May be used to set up multiple operating systems on same computer
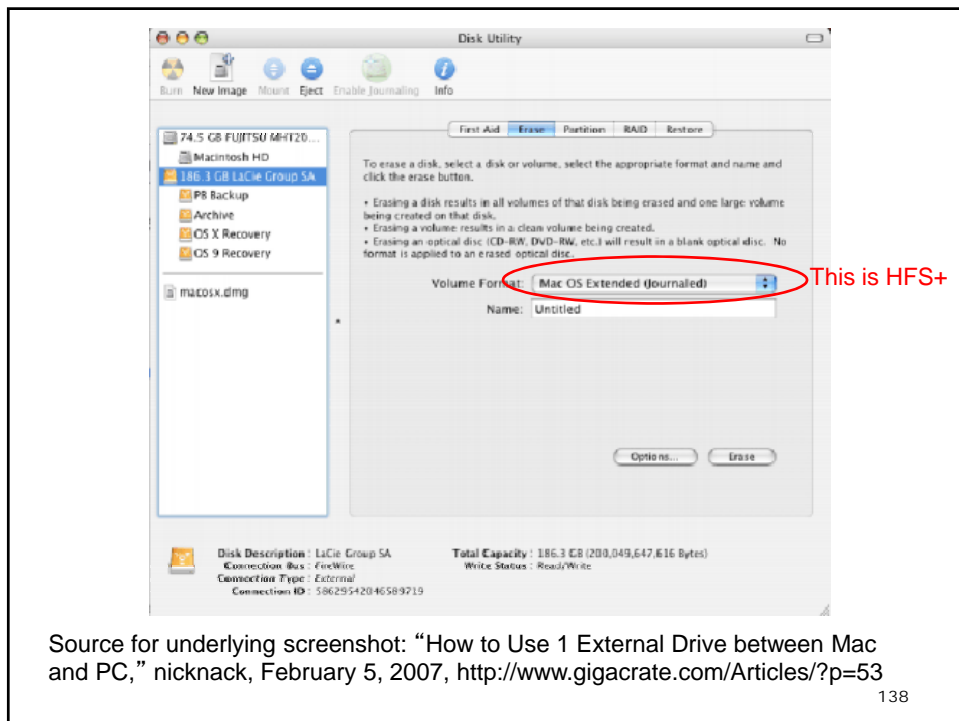
**135**

## File System

- Access controls
- File names & identifiers
- File size (length)
- Where to find files in storage (sectors and clusters)
- MAC times
  - Modified – when the content was last changed
  - Accessed – time file was last accessed (by person or software)
  - Changed – last time metadata changed
  - Created – (implemented inconsistently, if at all, across different file systems)

**136**

137



This is HFS+

Source for underlying screenshot: "How to Use 1 External Drive between Mac and PC," nicknack, February 5, 2007, http://www.gigacrate.com/Articles/?p=53

138

## File System Examples

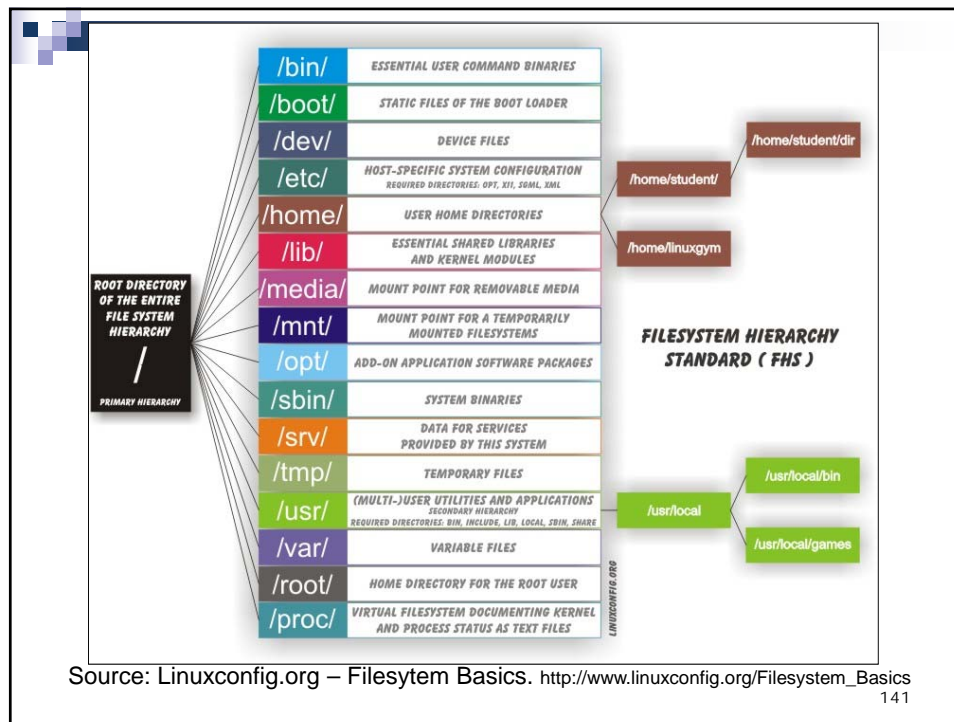| Name | Operating System(s) Using it as Native File System [often other OSs can also recognize it] |
| --- | --- |
| FAT12, FAT16 | MS-DOS |
| FAT32 (VFAT) | Windows 95, 98 |
| exFAT | Windows XP SP2 and later (primary use: USB drives, SD cards) |
| NTFS | Windows NT, 2000, XP, Server 2003, Server 2008, Vista |
| MFS | Macintosh System 1-3 |
| HFS (Hierarchical File System) | Macintosh System 4-8 |
| HFS+ | Macintosh System 8.1 – 9, OS X 10.0 – 10.11 |
| APFS | macOS 10.12 |
| ext, ext2, ext3, ext4 (Extended File System) | Linux |
| XFS | Linux, typically Enterprise variants (RHEL) |
| HPFS (High Performance File System) | OS/2 |
| ISOFS (ISO 9660) | Any OS that reads data from a CD |
| JFS1 (Journaled File System) | AIX (IBM) |
| ReiserFS | Several Linux distributions |
| UFS (Unix File System) aka FFS (Fast File System) | Various flavors of Unix |

139

## File System Examples

| Name | Operating System(s) Using it as Native File System [often other OSs can also recognize it] |
| --- | --- |
| FAT12, FAT16 | |
| FAT32 (VFAT) | |
| exFAT | s, SD cards) |
| NTFS | Vista |
| MFS | |
| HFS (Hierarchical File System) | |
| HFS+ | |
| APFS | |
| ext, ext2, ext3, ext4 (Extended File System) | |
| XFS | |
| HPFS (High Performance File System) | |
| ISOFS (ISO 9660) | |
| JFS1 (Journaled File System) | |
| ReiserFS | |
| UFS (Unix File System) aka FFS (Fast File System) | |

The filesystems you're most likely to encounter within archival collections

140

70

Source: Linuxconfig.org – Filesytem Basics. http://www.linuxconfig.org/Filesystem_Basics

**Microsoft File Systems:**

**FAT and NTFS**

## FAT

- Supported by all versions of Windows and most versions of UNIX

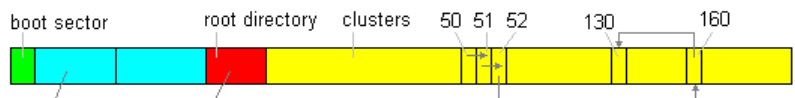- Common in many removable media, e.g. thumb drives, memory cards

http://mono-for-android.1047100.n5.nabble.com/detect-SD-Card-path-td5710218.html

http://www.ubergizmo.com/2008/09/usb-thumb-drive/

**143**

---

**FAT-16**

boot sector     root directory     clusters     50 51 52     130     160

FAT

```
50:    51
51:    52
52:    160
.........
130: EOF
........
160: 130
```

Root directory

`file.txt ... 20000 ... 50`

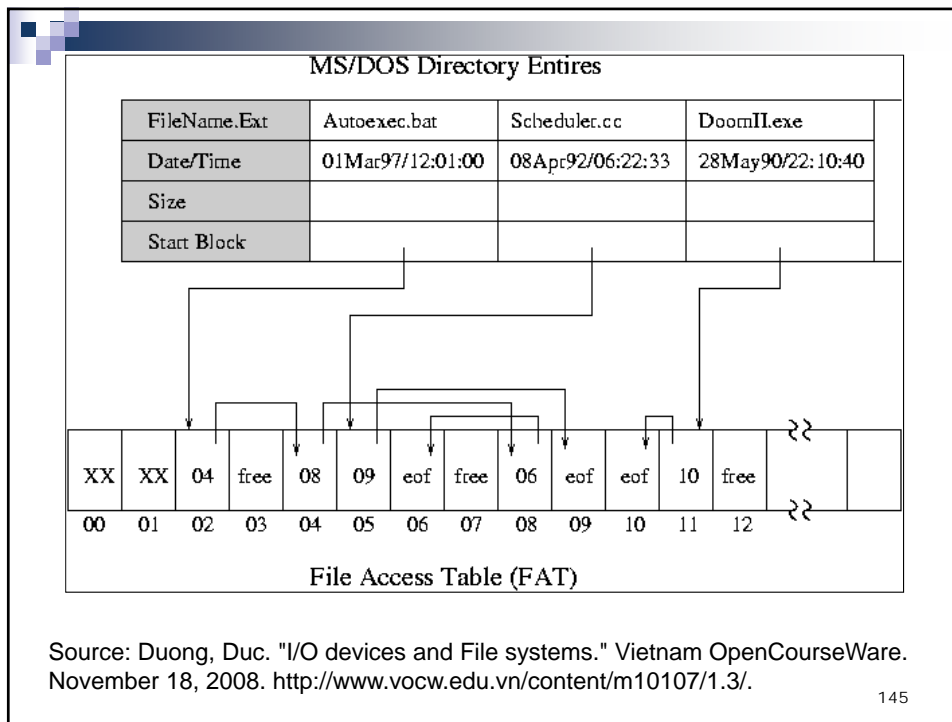Directory entry contains file size and number of the first cluster in the file.

Suppose there is a file "file.txt" on the disk which occupies clusters 50, 51, 52, 160, 130

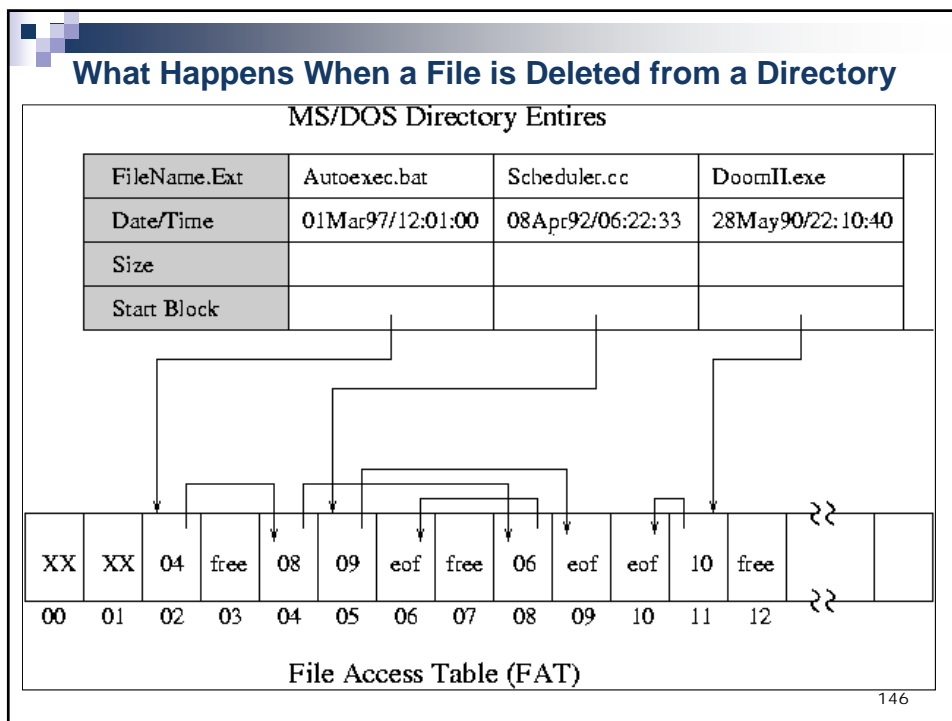Also we suppose disk is divided into 4k clusters

Source: Mikhail, Ranish. "Partitioning Primer." August 5, 1998.
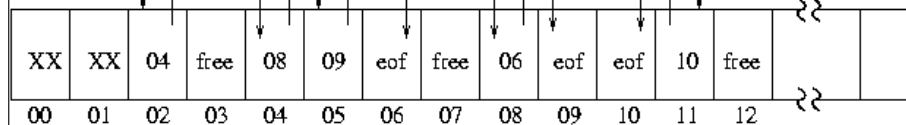http://www.ranish.com/part/primer.htm

**144**

## MS/DOS Directory Entires

| FileName.Ext | Autoexec.bat | Scheduler.cc | DoomII.exe |
|---|---|---|---|
| Date/Time | 01Mar97/12:01:00 | 08Apr92/06:22:33 | 28May90/22:10:40 |
| Size | | | |
| Start Block | | | |

| XX | XX | 04 | free | 08 | 09 | eof | free | 06 | eof | eof | 10 | free | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | |

### File Access Table (FAT)

Source: Duong, Duc. "I/O devices and File systems." Vietnam OpenCourseWare.
November 18, 2008. http://www.vocw.edu.vn/content/m10107/1.3/.

**145**

## What Happens When a File is Deleted from a Directory

### MS/DOS Directory Entires

| FileName.Ext | Autoexec.bat | Scheduler.cc | DoomII.exe |
|---|---|---|---|
| Date/Time | 01Mar97/12:01:00 | 08Apr92/06:22:33 | 28May90/22:10:40 |
| Size | | | |
| Start Block | | | |

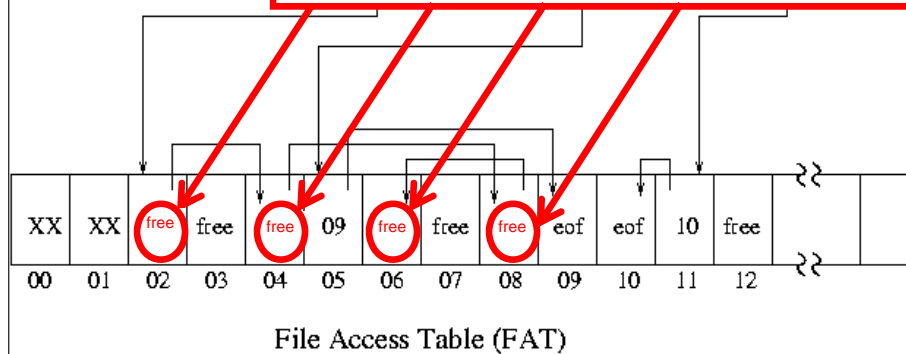| XX | XX | 04 | free | 08 | 09 | eof | free | 06 | eof | eof | 10 | free | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | |

### File Access Table (FAT)

**146**

## What Happens When a File is Deleted from a Directory

### MS/DOS Directory Entires

| FileName.Ext | E5utoExec.bat | Scheduler.cc | DoomII.exe |
|---|---|---|---|
| Date/Time | 01Mar97/12:01:00 | 08Apr92/06:22:33 | 28May90/22:10:40 |
| Size | | | |
| Start Block | | | |

Directory entry is marked as deleted, by replacing first character of the file name with a hex value of E5.

| XX | XX | 04 | free | 08 | 09 | eof | free | 06 | eof | eof | 10 | free | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | | |

File Access Table (FAT)

147

## What Happens When a File is Deleted from a Directory

### MS/DOS Directory Entires

| FileName.Ext | E5utoExec.bat | Scheduler.cc | DoomII.exe |
|---|---|---|---|
| Date/Time | 01Mar97/12:01:00 | 08Apr92/06:22:33 | 28May90/22:10:40 |
| Size | | | |
| Start Block | | | |

In the FAT, pointers to the clusters containing this file's data are marked as "free."

| XX | XX | free | free | free | 09 | free | free | free | eof | eof | 10 | free | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | | |

File Access Table (FAT)

148

74

**What Happens When a File is Deleted from a Directory**

MS/DOS Directory Entires

| FileName.Ext | E5utoExec.bat | Scheduler.cc | DoomII.exe |
|---|---|---|---|
| Date/Time | 01Mar97/12:01:00 | 08Apr92/06:22:33 | 28May90/22:10:40 |
| Size | | | |
| Start Block | | | |

Data from the deleted file will remain in these clusters, until they're over-written with data that is later allocated to those same clusters.

| XX | XX | free | free | free | 09 | free | free | free | eof | eof | 10 | free |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 |

File Access Table (FAT)

149

# NTFS

- Directory and FAT functions are combined in the Master File Table ($MFT)
- Each MFT record is assigned a unique number
- Good for forensic discovery:
  - ☐ For small files (< about 600 bytes), content is stored directly in the MFT itself & remains until overwritten by another MFT record
- Not so good for forensic discovery:
  - ☐ After deletion of a file, NTFS replaces (overwrites) the MFT record the next time a new file is created

150

## "Archive" Formats as Portable File Systems

- Most popular: zip and tar
- Retain important metadata that was in original file system, but does add a layer of representation information (packaging and possibly compression) that software needs to understand
- Compression also reduces robustness in the face of bit loss (any given bit flip is more likely to prevent recovery/rendering of content)
- Note that both zip and tar can be stored uncompressed (more common with tar)

## Bit-Level Treatment of Individual Files

**Example scenario:**
**A professor lets a couple of speakers use a USB thumb drive for their presentation slides.**

**Let's see what's on it...**

153

154

155



156

157



158

159

Note: The previous screenshots were from an older version of Windows. On a later version, you'll probably see something like this instead:



160

**161**



**162**

163



/a switch shows directories, read-only files, hidden files, system files

164

Period at beginning of file names indicates "hidden" files
•Not visible by default in Macintosh Finder
•Visible through Windows Explorer (with proper settings), DOS (using /a switch: dir /a) or Unix (-a switch: ls –a)

**165**



Spotlight is the desktop search utility on the Macintosh OS X.  The contents of this folder serve as an index of all files that were on the thumb drive that last time was last used on a Mac OS X computer.

**166**
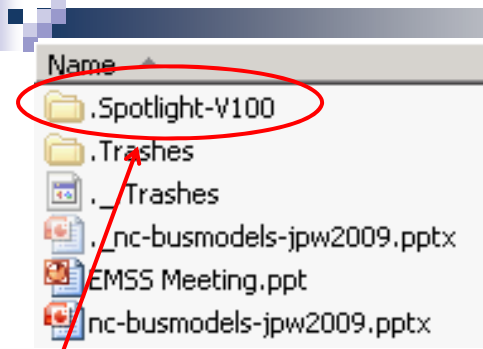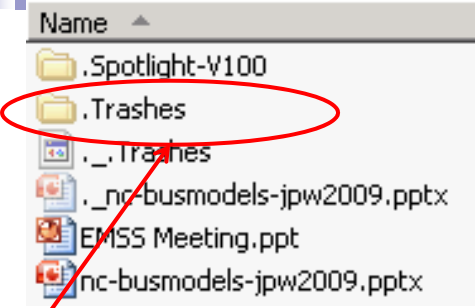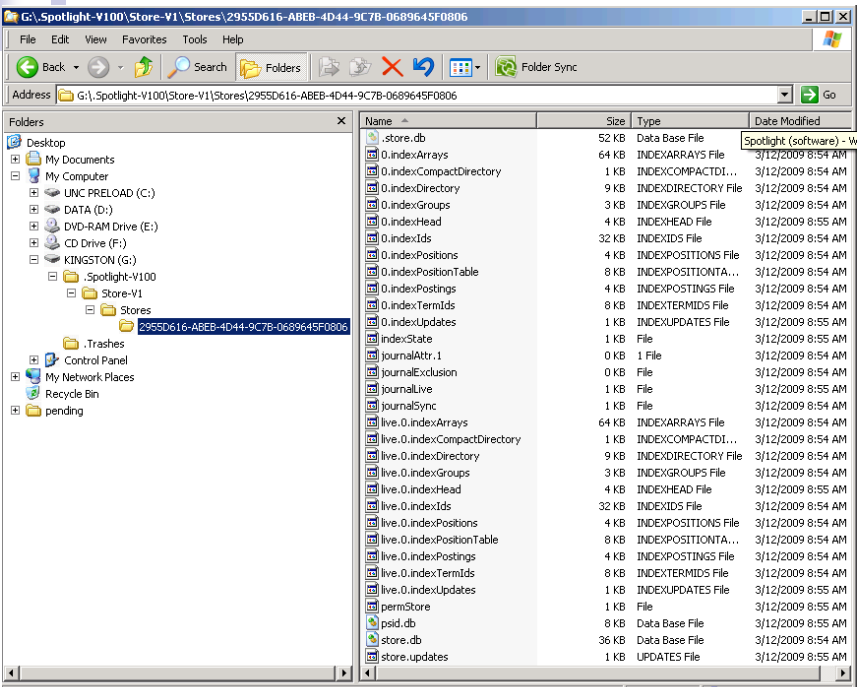
Name ▲

📁 .Spotlight-V100
📁 .Trashes
📄 ._.Trashes
📄 ._nc-busmodels-jpw2009.pptx
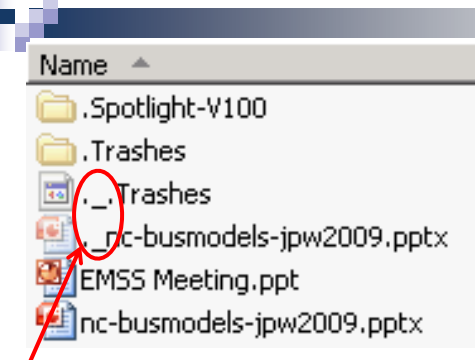📄 EMSS Meeting.ppt
📄 nc-busmodels-jpw2009.pptx

On a Mac, each user's home directory has a .Trash folder, and each volume has a .Trashes folder.   This is similar to the Windows Recycling Bin
.Trashes for each volume contains a separate sub-folder for the trash files contributed by each user (named by UID).

167

G:\.Spotlight-V100\Store-V1\Stores\2955D616-ABEB-4D44-9C7B-0689645F0806

File   Edit   View   Favorites   Tools   Help

Back ▼ ⊙ ▼ 🗁 | 🔍 Search 🗁 Folders | 🗐 🗐 ✕ 🖺 | ▥ ▼ | 🗐 Folder Sync

Address 🗁 G:\.Spotlight-V100\Store-V1\Stores\2955D616-ABEB-4D44-9C7B-0689645F0806                    ▼ → Go

| Folders                                          | ✕ | Name ▲ | Size | Type | Date Modified |
|--------------------------------------------------|---|--------|------|------|---------------|
| 🖳 Desktop                                        |   | .store.db | 52 KB | Data Base File | Spotlight (software) - W |
| ⊞ 📁 My Documents                                 |   | 0.indexArrays | 64 KB | INDEXARRAYS File | 3/12/2009 8:54 AM |
| ⊟ 🖳 My Computer                                   |   | 0.indexCompactDirectory | 1 KB | INDEXCOMPACTDI... | 3/12/2009 8:54 AM |
| ⊞ 🖴 UNC PRELOAD (C:)                              |   | 0.indexDirectory | 9 KB | INDEXDIRECTORY File | 3/12/2009 8:54 AM |
| ⊞ 🖴 DATA (D:)                                     |   | 0.indexGroups | 3 KB | INDEXGROUPS File | 3/12/2009 8:54 AM |
| ⊞ 💿 DVD-RAM Drive (E:)                            |   | 0.indexHead | 4 KB | INDEXHEAD File | 3/12/2009 8:55 AM |
| ⊞ 💿 CD Drive (F:)                                 |   | 0.indexIds | 32 KB | INDEXIDS File | 3/12/2009 8:54 AM |
| ⊟ 🖴 KINGSTON (G:)                                 |   | 0.indexPositions | 4 KB | INDEXPOSITIONS File | 3/12/2009 8:54 AM |
| ⊟ 📁 .Spotlight-V100                               |   | 0.indexPositionTable | 8 KB | INDEXPOSITIONTA... | 3/12/2009 8:54 AM |
| ⊟ 📁 Store-V1                                      |   | 0.indexPostings | 4 KB | INDEXPOSTINGS File | 3/12/2009 8:54 AM |
| ⊟ 📁 Stores                                        |   | 0.indexTermIds | 8 KB | INDEXTERMIDS File | 3/12/2009 8:54 AM |
| 📁 2955D616-ABEB-4D44-9C7B-0689645F0806           |   | 0.indexUpdates | 1 KB | INDEXUPDATES File | 3/12/2009 8:55 AM |
| 📁 .Trashes                                        |   | indexState | 1 KB | File | 3/12/2009 8:55 AM |
| ⊞ 🗐 Control Panel                                 |   | journalAttr.1 | 0 KB | 1 File | 3/12/2009 8:54 AM |
| ⊞ 📁 My Network Places                             |   | journalExclusion | 0 KB | File | 3/12/2009 8:54 AM |
| 🗑 Recycle Bin                                     |   | journalLive | 1 KB | File | 3/12/2009 8:54 AM |
| ⊞ 📁 pending                                       |   | journalSync | 1 KB | File | 3/12/2009 8:55 AM |
|                                                  |   | live.0.indexArrays | 64 KB | INDEXARRAYS File | 3/12/2009 8:54 AM |
|                                                  |   | live.0.indexCompactDirectory | 1 KB | INDEXCOMPACTDI... | 3/12/2009 8:54 AM |
|                                                  |   | live.0.indexDirectory | 9 KB | INDEXDIRECTORY File | 3/12/2009 8:54 AM |
|                                                  |   | live.0.indexGroups | 3 KB | INDEXGROUPS File | 3/12/2009 8:54 AM |
|                                                  |   | live.0.indexHead | 4 KB | INDEXHEAD File | 3/12/2009 8:55 AM |
|                                                  |   | live.0.indexIds | 32 KB | INDEXIDS File | 3/12/2009 8:54 AM |
|                                                  |   | live.0.indexPositions | 4 KB | INDEXPOSITIONS File | 3/12/2009 8:54 AM |
|                                                  |   | live.0.indexPositionTable | 8 KB | INDEXPOSITIONTA... | 3/12/2009 8:54 AM |
|                                                  |   | live.0.indexPostings | 4 KB | INDEXPOSTINGS File | 3/12/2009 8:54 AM |
|                                                  |   | live.0.indexTermIds | 8 KB | INDEXTERMIDS File | 3/12/2009 8:54 AM |
|                                                  |   | live.0.indexUpdates | 1 KB | INDEXUPDATES File | 3/12/2009 8:55 AM |
|                                                  |   | permStore | 1 KB | File | 3/12/2009 8:54 AM |
|                                                  |   | psid.db | 8 KB | Data Base File | 3/12/2009 8:55 AM |
|                                                  |   | store.db | 36 KB | Data Base File | 3/12/2009 8:54 AM |
|                                                  |   | store.updates | 1 KB | UPDATES File | 3/12/2009 8:55 AM |

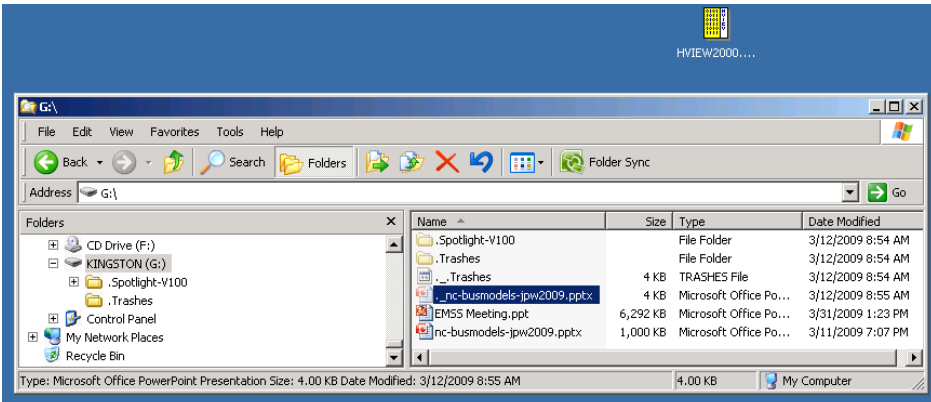32 objects                                                           367 KB        🖳 My Computer

168

Files beginning with ._ are resource fork files, which the Mac file system uses to store icon, data type, and some other metadata about the file (content being stored in the data fork)

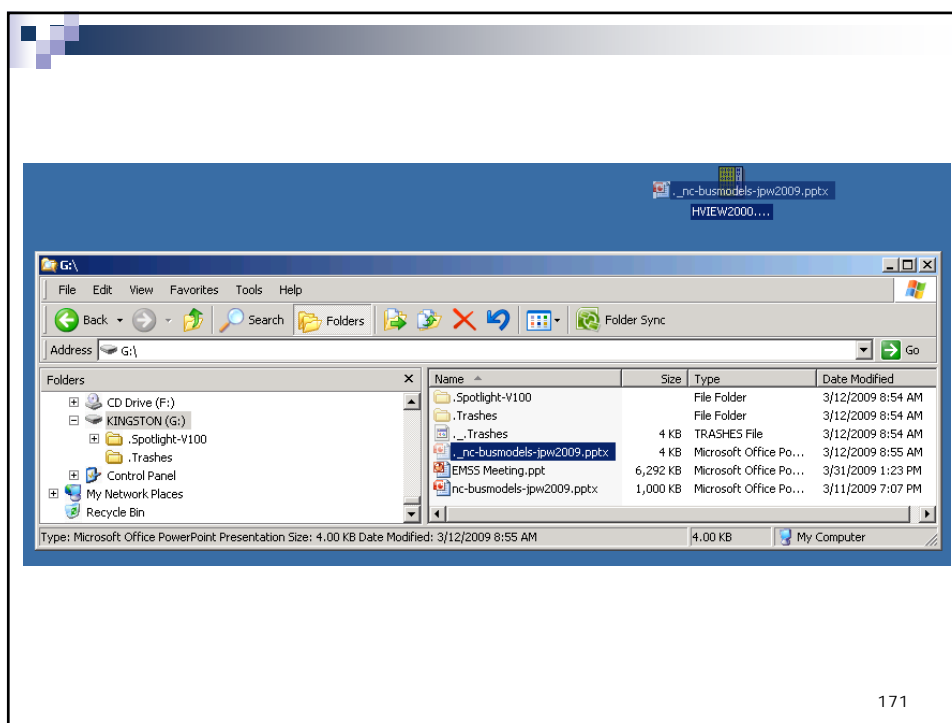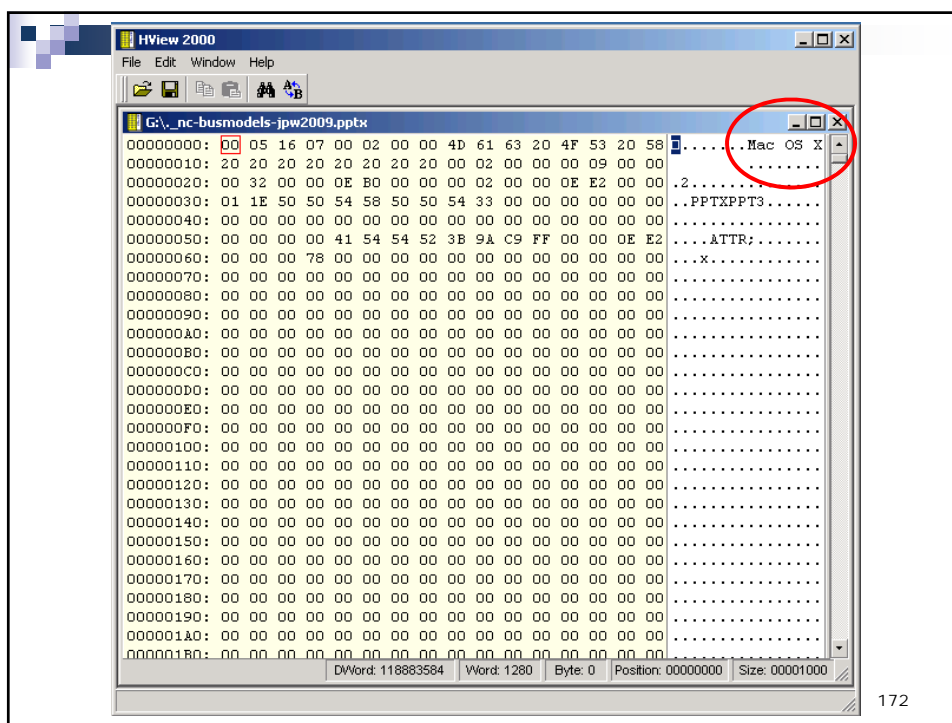Maybe I should look at a hex dump of one of them to be sure...
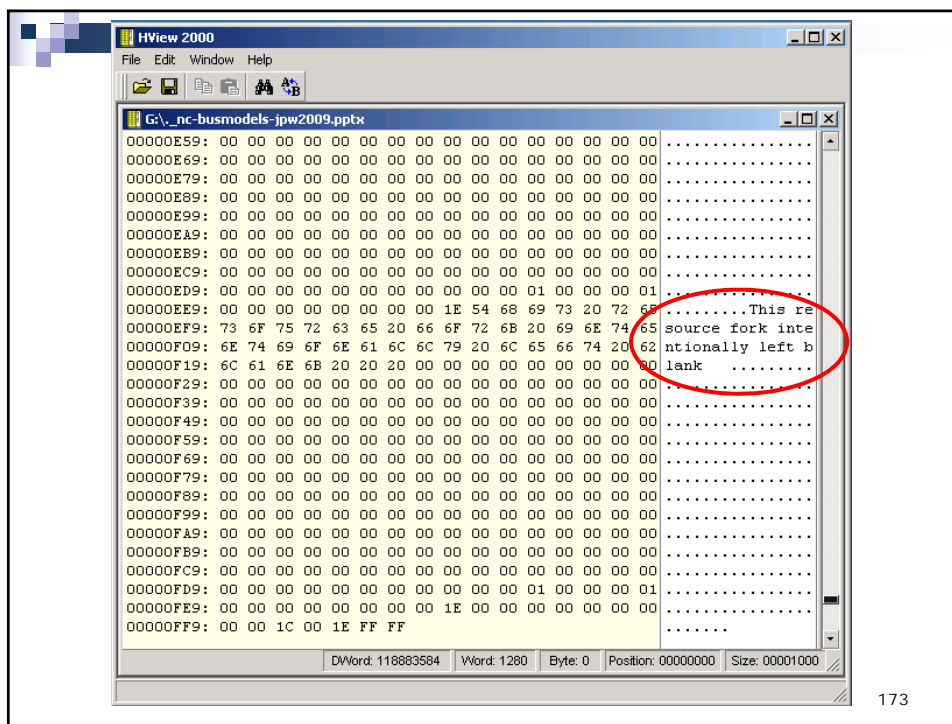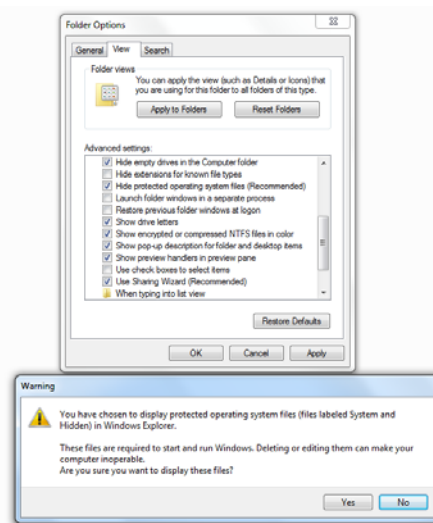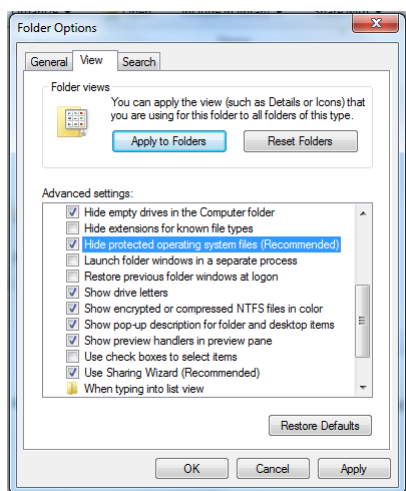
**169**



**170**

171



172

**173**

This wasn't necessary for seeing the hidden files (resource forks) in the previous example, but another system change you might want to make on your processing/forensics workstation to show system files that are normally hidden from the user:



**174**

# Forms of "Hidden Data"

## **Not** just what you see when you open a file in its native application

## Listed roughly in order of difficulty for identification and retrieval

175

---

# Sanitization Taxonomy

| LEVEL | WHERE FOUND | DESCRIPTION |
|---|---|---|
| Level 0 | Regular files | Information contained in the file system. Includes file names, file attributes, and file contents. By definition, no attempts are made to sanitize Level 0 files information. Level 0 also includes information that is written to the disk as part of any sanitization attempt. For example, if a copy of Windows 95 had been installed on a hard drive in an attempt to sanitize the drive, then the files installed into the C:\WINDOWS directory would be considered Level 0 files. No special tools are required to retrieve Level 0 data. |
| Level 1 | Temporary files | Temporary files, including print spooler files, browser cache files, files for "helper" applications, and recycle bin files. Most users either expect the system to automatically delete this data or are not even aware that it exists. Note: Level 0 files are a subset of Level 1 files. Experience has shown that it is useful to distinguish this subset, because many naive users will overlook Level 1 files when they are browsing a computer's hard drive to see if it contains sensitive information. No special tools are required to retrieve Level 1 data, although special training is required to teach the operator where to look. |
| Level 2 | Deleted files | When a file is deleted from a file system, most operating systems do not overwrite the blocks on the hard disk that the file is written on. Instead, they simply remove the file's reference from the containing directory. The file's blocks are then placed on the free list. These files can be recovered using traditional "undelete" tools, such as Norton Utilities. |
| Level 3 | Retained data blocks | Data that can be recovered from a disk, but which does not obviously belong to a named file. Level 3 data includes information in slack space, backing store for virtual memory, and Level 2 data that has been partially overwritten so that an entire file cannot be recovered. A common source of Level 3 data is disks that have been formatted with Windows Format command or the Unix newfs command. Even though the output of these commands might imply that they overwrite the entire hard drive, in fact they do not, and the vast majority of the formatted disk's information is recoverable with the proper tools. Level 3 data can be recovered using advanced data recovery tools that can "unformat" a disk drive or special-purpose forensics tools. |
| Level 4 | Vendor-hidden data | This level consists of data blocks that can only be accessed using vendor-specific commands. This level includes the drive's controlling program and blocks used for bad-block management. |
| Level 5 | Overwritten data | Many individuals maintain that information can be recovered from a hard drive even after it is overwritten. We reserve Level 5 for such information. |

Garfinkel, Simson L., and Abhi Shelat. "Remembrance of Data Passed: A Study of Disk Sanitization Practices." *IEEE Security and Privacy* 1 (2003): 17-27.

176

## Within the Files Themselves

- Lots of data in many files that you don't always see with the naked eye. For example:

    - ☐ Comments within the code

    - ☐ Stored rules & styles

    - ☐ Change tracking information

    - ☐ Metadata stored in file headers & elsewhere

    - ☐ Viruses!

**177**

## Examples of Hidden Data in MS Office Documents

- Application used to create document
- Authors, user names, organizational affiliations & author history
- Comments
- Custom properties
- Database queries
- Embedded objects (OLE) – elements not immediately visible (e.g. spreadsheet)
- Fast save – change history appended to end of file, rather than applied to body of document
- GUID – globally unique identifier for computer (see Leach et al, 2005)
- Hidden cells, slides, text – purposely hidden but then possibly forgotten

- Outlook (email) properties & routing slips
- Path information – audio & video paths, author history, linked objects, printers, hyperlinks, include fields, template
- Presentation notes
- Printer driver information
- RSID – Revision save ID (differentiates changes from different editing sessions)
- Tracked changes (added to PPT and Excel in Office XP)
- Versions
- Visual Basic code – including macros & viruses (and identity of code creators)
- Web server information
- White text (on white background)

**178**

## Jonathan Larson Fast Save Example

179

---

## Hidden Image Data

- Content outside crop area
- Layered objects (hidden from view)
- Pixel information in resized and embedded image
- Metadata:
  - GIF – comment extensions and application extensions
  - JPEG - camera use, date/time, distance settings, location, thumbnail image

180

90

## Example of EXIF Metadata from a JPEG File (Generated Using exiftool*)

---- ExifTool ----
ExifTool Version Number     : 9.38
---- System ----
File Name                   : IMG_20130823_151811.jpg
Directory                   : C:/Users/callee/Documents/images/digital-forensics-lab
File Size                   : 1785 kB
File Modification Date/Time : 2013:08:23 16:36:44-04:00
File Access Date/Time       : 2013:10:14 17:13:02-04:00
File Creation Date/Time     : 2013:08:23 16:36:44-04:00
File Permissions            : rw-rw-rw-
---- File ----
File Type                   : JPEG
MIME Type                   : image/jpeg
Exif Byte Order             : Big-endian (Motorola, MM)
Image Width                 : 2592
Image Height                : 1944
Encoding Process            : Baseline DCT, Huffman coding
Bits Per Sample             : 8
Color Components            : 3
Y Cb Cr Sub Sampling        : YCbCr4:2:0 (2 2)
---- GPS ----
GPS Img Direction           : 83
GPS Img Direction Ref       : Magnetic North
GPS Latitude Ref            : North
GPS Latitude                : 35 deg 55' 2.24"
GPS Longitude Ref           : West
GPS Longitude               : 79 deg 2' 57.55"
GPS Altitude Ref            : Above Sea Level
GPS Altitude                : 0 m
GPS Time Stamp              : 19:18:06
GPS Processing Method       : NETWORK
GPS Date Stamp              : 2013:08:23
---- IFD0 ----
Orientation                 : Unknown (0)
Camera Model Name           : Galaxy Nexus
Modify Date                 : 2013:08:23 15:18:11
Y Cb Cr Positioning         : Centered
Y Resolution                : 72
Resolution Unit             : inches
X Resolution                : 72
Make                        : Samsung
---- ExifIFD ----
Create Date                 : 2013:08:23 15:18:11
Date/Time Original          : 2013:08:23 15:18:11
Exif Version                : 0220
Flash Energy                : 0
Image Unique ID             : OAEL01
Exposure Time               : 1/17
ISO                         : 125, 0, 0

Scene Type                  : Directly photographed
Exposure Index              : undef
Components Configuration    : Y, Cb, Cr, -
F Number                    : 2.8
Compressed Bits Per Pixel   : 0
Sensing Method              : One-chip color area
Exposure Program            : Aperture-priority AE
Aperture Value              : 2.6
Brightness Value            : 0
Subject Distance Range      : Unknown
Shutter Speed Value         : 1/15
Subject Distance            : 0 m
Saturation                  : Normal
Color Space                 : sRGB
Contrast                    : Normal
Metering Mode               : Multi-spot
Flashpix Version            :
Exposure Compensation       : 0
Exif Image Height           : 1944
Max Aperture Value          : 2.6
Sharpness                   : Normal
Exif Image Width            : 2592
Focal Length                : 3.4 mm
Digital Zoom Ratio          : 1
Light Source                : Fluorescent
Scene Capture Type          : Standard
Flash                       : Off, Did not fire
Custom Rendered             : Custom
White Balance               : Auto
Exposure Mode               : Auto
---- IFD1 ----
Compression                 : JPEG (old-style)
Image Width                 : 160
Image Height                : 120
Thumbnail Offset            : 1239
Thumbnail Length            : 7164
---- Composite ----
Aperture                    : 2.8
GPS Altitude                : 0 m Above Sea Level
GPS Date/Time               : 2013:08:23 19:18:06Z
GPS Latitude                : 35 deg 55' 2.24" N
GPS Longitude               : 79 deg 2' 57.55" W
GPS Position                : 35 deg 55' 2.24" N, 79 deg 2' 57.55" W
Image Size                  : 2592x1944
Shutter Speed               : 1/17
Thumbnail Image             : (Binary data 7164 bytes, use -b option to extract)
Focal Length                : 3.4 mm
Light Value                 : 6.7

*http://www.sno.phy.queensu.ca/~phil/exiftool/ (Also available through the BitCurator environment)

181

---

## Example of EXIF Metadata from a JPEG File (Generated Using exiftool*)

---- ExifTool ----
ExifTool Version Number     : 9.38

**File Name                   : IMG_20130823_151811.jpg**
**Directory                   : C:/Users/callee/Documents/images/digital-forensics-lab**
**File Size                   : 1785 kB**
**File Modification Date/Time : 2013:08:23 16:36:44-04:00**
**File Access Date/Time       : 2013:10:14 17:13:02-04:00**
**File Creation Date/Time     : 2013:08:23 16:36:44-04:00**
**File Permissions            : rw-rw-rw-**

File Type                   : JPEG
MIME Type                   : image/jpeg
Exif Byte Order             : Big-endian (Motorola, MM)
Image Width                 : 2592
Image Height                : 1944
Encoding Process            : Baseline DCT, Huffman coding
Bits Per Sample             : 8
Color Components            : 3
Y Cb Cr Sub Sampling        : YCbCr4:2:0 (2 2)
---- GPS ----
GPS Img Direction           : 83
GPS Img Direction Ref       : Magnetic North
GPS Latitude Ref            : North
GPS Latitude                : 35 deg 55' 2.24"
GPS Longitude Ref           : West
GPS Longitude               : 79 deg 2' 57.55"
GPS Altitude Ref            : Above Sea Level
GPS Altitude                : 0 m
GPS Time Stamp              : 19:18:06
GPS Processing Method       : NETWORK
GPS Date Stamp              : 2013:08:23
---- IFD0 ----
Orientation                 : Unknown (0)
Camera Model Name           : Galaxy Nexus
Modify Date                 : 2013:08:23 15:18:11
Y Cb Cr Positioning         : Centered
Y Resolution                : 72
Resolution Unit             : inches
X Resolution                : 72
Make                        : Samsung
---- ExifIFD ----
Create Date                 : 2013:08:23 15:18:11
Date/Time Original          : 2013:08:23 15:18:11
Exif Version                : 0220
Flash Energy                : 0
Image Unique ID             : OAEL01
Exposure Time               : 1/17
ISO                         : 125, 0, 0

Scene Type                  : Directly photographed
Exposure Index              : undef
Components Configuration    : Y, Cb, Cr, -
F Number                    : 2.8
Compressed Bits Per Pixel   : 0
Sensing Method              : One-chip color area
Exposure Program            : Aperture-priority AE
Aperture Value              : 2.6
Brightness Value            : 0
Subject Distance Range      : Unknown
Shutter Speed Value         : 1/15
Subject Distance            : 0 m
Saturation                  : Normal
Color Space                 : sRGB
Contrast                    : Normal
Metering Mode               : Multi-spot
Flashpix Version            :
Exposure Compensation       : 0
Exif Image Height           : 1944
Max Aperture Value          : 2.6
Sharpness                   : Normal
Exif Image Width            : 2592
Focal Length                : 3.4 mm
Digital Zoom Ratio          : 1
Light Source                : Fluorescent
Scene Capture Type          : Standard
Flash                       : Off, Did not fire
Custom Rendered             : Custom
White Balance               : Auto
Exposure Mode               : Auto
---- IFD1 ----
Compression                 : JPEG (old-style)
Image Width                 : 160
Image Height                : 120
Thumbnail Offset            : 1239
Thumbnail Length            : 7164
---- Composite ----
Aperture                    : 2.8
GPS Altitude                : 0 m Above Sea Level
GPS Date/Time               : 2013:08:23 19:18:06Z
GPS Latitude                : 35 deg 55' 2.24" N
GPS Longitude               : 79 deg 2' 57.55" W
GPS Position                : 35 deg 55' 2.24" N, 79 deg 2' 57.55" W
Image Size                  : 2592x1944
Shutter Speed               : 1/17
Thumbnail Image             : (Binary data 7164 bytes, use -b option to extract)
Focal Length                : 3.4 mm
Light Value                 : 6.7

*http://www.sno.phy.queensu.ca/~phil/exiftool/ (Also available through the BitCurator environment)

182

## Example of EXIF Metadata from a JPEG File (Generated Using exiftool*)

```
---- ExifTool ----
ExifTool Version Number    : 9.38
---- System ----
File Name            : IMG_20130823_151811.jpg
Directory            : C:/Users/callee/Documents/images/digital-forensics-lab
File Size            : 1785 kB
File Modification Date/Time   : 2013:08:23 16:36:44-04:00
File Access Date/Time    : 2013:10:14 17:13:02-04:00
File Creation Date/Time    : 2013:08:23 16:36:44-04:00
File Permissions      : rw-rw-rw-
---- File ----
File Type            : JPEG
MIME Type            : image/jpeg
                       Big-endian (Motorola, MM)
Image Width          : 2592
Image Height         : 1944
Bits Per Sample      : 8
Color Components     : 3
Y Cb Cr Sub Sampling     : YCbCr4:2:0 (2 2)
---- GPS ----
GPS Img Direction        : 83
GPS Img Direction Ref      : Magnetic North
GPS Latitude Ref      : North
GPS Latitude         : 35 deg 55' 2.24"
GPS Longitude Ref     : West
GPS Longitude        : 79 deg 2' 57.55"
GPS Altitude Ref      : Above Sea Level
GPS Altitude         : 0 m
GPS Time Stamp       : 19:18:06
GPS Processing Method      : NETWORK
GPS Date Stamp       : 2013:08:23
---- IFD0 ----
Orientation          : Unknown (0)
Camera Model Name        : Galaxy Nexus
Modify Date          : 2013:08:23 15:18:11
Y Cb Cr Positioning      : Centered
Y Resolution         : 72
Resolution Unit       : inches
X Resolution          : 72
Make                 : Samsung
---- ExifIFD ----
Create Date          : 2013:08:23 15:18:11
Date/Time Original       : 2013:08:23 15:18:11
Exif Version          : 0220
Flash Energy          : 0
Image Unique ID          : OAEL01
Exposure Time         : 1/17
ISO                  : 125, 0, 0
```

```
Scene Type           : Directly photographed
Exposure Index        : undef
Components Configuration   : Y, Cb, Cr, -
F Number              : 2.8
Compressed Bits Per Pixel   : 0
Sensing Method        : One-chip color area
Exposure Program          : Aperture-priority AE
Aperture Value        : 2.6
Brightness Value      : 0
Subject Distance Range      : Unknown
Shutter Speed Value       : 1/15
Subject Distance      : 0 m
Saturation           : Normal
Color Space          : sRGB
Contrast             : Normal
Metering Mode         : Multi-spot
Flashpix Version      :
Exposure Compensation        : 0
Exif Image Height     : 1944
Max Aperture Value        : 2.6
Sharpness            : Normal
Exif Image Width      : 2592
Focal Length         : 3.4 mm
Digital Zoom Ratio        : 1
Light Source          : Fluorescent
Scene Capture Type        : Standard
Flash                : Off, Did not fire
Custom Rendered           : Custom
White Balance         : Auto
Exposure Mode         : Auto
---- IFD1 ----
Compression           : JPEG (old-style)
Image Width          : 160
Image Height         : 120
Thumbnail Offset      : 1239
Thumbnail Length      : 7164
---- Composite ----
Aperture             : 2.8
GPS Altitude          : 0 m Above Sea Level
GPS Date/Time         : 2013:08:23 19:18:06Z
GPS Latitude          : 35 deg 55' 2.24" N
GPS Longitude         : 79 deg 2' 57.55" W
GPS Position          : 35 deg 55' 2.24" N, 79 deg 2' 57.55" W
Image Size           : 2592x1944
Shutter Speed         : 1/17
Thumbnail Image       : (Binary data 7164 bytes, use -b option to extract)
Focal Length         : 3.4 mm
Light Value           : 6.7
```

*http://www.sno.phy.queensu.ca/~phil/exiftool/ (Also available through the BitCurator environment)

183

---

## Example of EXIF Metadata from a JPEG File (Generated Using exiftool*)

```
---- ExifTool ----
ExifTool Version Number    : 9.38
---- System ----
File Name            : IMG_20130823_151811.jpg
Directory            : C:/Users/callee/Documents/images/digital-forensics-lab
File Size            : 1785 kB
File Modification Date/Time    : 2013:08:23 16:36:44-04:00
File Access Date/Time    : 2013:10:14 17:13:02-04:00
File Creation Date/Time    : 2013:08:23 16:36:44-04:00
File Permissions      : rw-rw-rw-
---- File ----
File Type            : JPEG
MIME Type            : image/jpeg
Exif Byte Order       : Big-endian (Motorola, MM)
Image Width          : 2592
Image Height         : 1944
Encoding Process      : Baseline DCT, Huffman coding
Bits Per Sample      : 8
Color Components     : 3
Y Cb Cr Sub Sampling     : YCbCr4:2:0 (2 2)
GPS Img Direction        : 83
GPS Img Direction Ref      : Magnetic North
GPS Latitude Ref      : North
GPS Latitude         : 35 deg 55' 2.24"
GPS Longitude Ref     : West
GPS Longitude        : 79 deg 2' 57.55"
GPS Altitude Ref      : Above Sea Level
GPS Altitude         : 0 m
GPS Time Stamp       : 19:18:06
GPS Processing Method      : NETWORK
GPS Date Stamp       : 2013:08:23
Orientation          : Unknown (0)
Camera Model Name        : Galaxy Nexus
Modify Date          : 2013:08:23 15:18:11
Y Cb Cr Positioning      : Centered
Y Resolution         : 72
Resolution Unit       : inches
X Resolution          : 72
Make                 : Samsung
---- ExifIFD ----
Create Date          : 2013:08:23 15:18:11
Date/Time Original       : 2013:08:23 15:18:11
Exif Version          : 0220
Flash Energy          : 0
Image Unique ID          : OAEL01
Exposure Time         : 1/17
ISO                  : 125, 0, 0
```

```
Scene Type           : Directly photographed
Exposure Index        : undef
Components Configuration   : Y, Cb, Cr, -
F Number              : 2.8
Compressed Bits Per Pixel   : 0
Sensing Method        : One-chip color area
Exposure Program          : Aperture-priority AE
Aperture Value        : 2.6
Brightness Value      : 0
Subject Distance Range      : Unknown
Shutter Speed Value       : 1/15
Subject Distance      : 0 m
Saturation           : Normal
Color Space          : sRGB
Contrast             : Normal
Metering Mode         : Multi-spot
Flashpix Version      :
Exposure Compensation        : 0
Exif Image Height     : 1944
Max Aperture Value        : 2.6
Sharpness            : Normal
Exif Image Width      : 2592
Focal Length         : 3.4 mm
Digital Zoom Ratio        : 1
Light Source          : Fluorescent
Scene Capture Type        : Standard
Flash                : Off, Did not fire
Custom Rendered           : Custom
White Balance         : Auto
Exposure Mode         : Auto
---- IFD1 ----
Compression           : JPEG (old-style)
Image Width          : 160
Image Height         : 120
Thumbnail Offset      : 1239
Thumbnail Length      : 7164
---- Composite ----
Aperture             : 2.8
GPS Altitude          : 0 m Above Sea Level
GPS Date/Time         : 2013:08:23 19:18:06Z
GPS Latitude          : 35 deg 55' 2.24" N
GPS Longitude         : 79 deg 2' 57.55" W
GPS Position          : 35 deg 55' 2.24" N, 79 deg 2' 57.55" W
Image Size           : 2592x1944
Shutter Speed         : 1/17
Thumbnail Image       : (Binary data 7164 bytes, use -b option to extract)
Focal Length         : 3.4 mm
Light Value           : 6.7
```

*http://www.sno.phy.queensu.ca/~phil/exiftool/ (Also available through the BitCurator environment)

184

## Example of EXIF Metadata from a JPEG File (Generated Using exiftool*)

```
---- ExifTool ----
ExifTool Version Number    : 9.38
---- System ----
File Name                  : IMG_20130823_151811.jpg
Directory                  : C:/Users/callee/Documents/images/digital-forensics-lab
File Size                  : 1785 kB
File Modification Date/Time : 2013:08:23 16:36:44-04:00
File Access Date/Time      : 2013:10:14 17:13:02-04:00
File Creation Date/Time    : 2013:08:23 16:36:44-04:00
File Permissions           : rw-rw-rw-
---- File ----
File Type                  : JPEG
MIME Type                  : image/jpeg
Exif Byte Order            : Big-endian (Motorola, MM)
Image Width                : 2592
Image Height               : 1944
Encoding Process           : Baseline DCT, Huffman coding
Bits Per Sample            : 8
Color Components           : 3
Y Cb Cr Sub Sampling       : YCbCr4:2:0 (2 2)
---- GPS ----
GPS Img Direction          : 83
GPS Img Direction Ref      : Magnetic North
GPS Latitude Ref           : North
GPS Latitude               : 35 deg 55' 2.24"
GPS Longitude Ref          : West
GPS Longitude              : 79 deg 2' 57.55"
GPS Altitude Ref           : Above Sea Level
GPS Altitude               : 0 m
GPS Time Stamp             : 19:18:06
GPS Processing Method      : NETWORK
GPS Date Stamp             : 2013:08:23
---- IFD0 ----
Camera Model Name          : Galaxy Nexus
Modify Date                : 2013:08:23 15:18:11
Y Cb Cr Positioning        : Centered
Y Resolution               : 72
Resolution Unit            : inches
X Resolution               : 72
Make                       : Samsung
---- ExifIFD ----
Create Date                : 2013:08:23 15:18:11
Date/Time Original         : 2013:08:23 15:18:11
Exif Version               : 0220
Flash Energy               : 0
Image Unique ID            : OAEL01
Exposure Time              : 1/17
ISO                        : 125, 0, 0

Scene Type                 : Directly photographed
Exposure Index             : undef
Components Configuration   : Y, Cb, Cr, -
F Number                   : 2.8
Compressed Bits Per Pixel  : 0
Sensing Method             : One-chip color area
Exposure Program           : Aperture-priority AE
Aperture Value             : 2.6
Brightness Value           : 0
Subject Distance Range     : Unknown
Shutter Speed Value        : 1/15
Subject Distance           : 0 m
Saturation                 : Normal
Color Space                : sRGB
Contrast                   : Normal
Metering Mode              : Multi-spot
Flashpix Version           :
Exposure Compensation      : 0
Exif Image Height          : 1944
Max Aperture Value         : 2.6
Sharpness                  : Normal
Exif Image Width           : 2592
Focal Length               : 3.4 mm
Digital Zoom Ratio         : 1
Light Source               : Fluorescent
Scene Capture Type         : Standard
Flash                      : Off, Did not fire
Custom Rendered            : Custom
White Balance              : Auto
Exposure Mode              : Auto
---- IFD1 ----
Compression                : JPEG (old-style)
Image Width                : 160
Image Height               : 120
Thumbnail Offset           : 1239
Thumbnail Length           : 7164
---- Composite ----
Aperture                   : 2.8
GPS Altitude               : 0 m Above Sea Level
GPS Date/Time              : 2013:08:23 19:18:06Z
GPS Latitude               : 35 deg 55' 2.24" N
GPS Longitude              : 79 deg 2' 57.55" W
GPS Position               : 35 deg 55' 2.24" N, 79 deg 2' 57.55" W
Image Size                 : 2592x1944
Shutter Speed              : 1/17
Thumbnail Image            : (Binary data 7164 bytes, use -b option to extract)
Focal Length               : 3.4 mm
Light Value                : 6.7
```

*http://www.sno.phy.queensu.ca/~phil/exiftool/ (Also available through the BitCurator environment)          185

---

# EXIF Metadata From Header of a TIFF File*



*Using pyExifToolGUI in the BitCurator environment          186

## Identifying File Types

- Magic numbers and file signatures

- File extensions

- Metadata stored in file system

- MIME types

187

## Magic Numbers and File Signatures

- Distinct string or pattern that is found within files of a given type (most often in the header)
- Most effective searches for magic numbers often involve regular expressions (e.g. grep) in order to indicate multiple variations of a pattern
- Utilities that use this: file (Unix), TrID, DROID, FITS
- Examples:

| File Format | Hex | ASCII |
|---|---|---|
| DOC | D0 CF 11 E0 A1 B1 1A E1 | ÐÏà¡±á |
| JPG | FF D8 FF | ÿØÿ |
| PDF | 25 50 44 46 2D 31 2E | %PDF-1. |
| ZIP | 50 4B 03 04 | PK.. |

188

# Try it Yourself!

- Go back to the file that you downloaded earlier and open it again in HxD or Hex Fiend
- What do you see at the beginning?
- Is there a pattern that matches one of these?

| File Format | Hex | ASCII |
|---|---|---|
| DOC | D0 CF 11 E0 A1 B1 1A E1 | ÐÏà¡±á |
| JPG | FF D8 FF | ÿØÿ |
| PDF | 25 50 44 46 2D 31 2E | %PDF-1. |
| ZIP | 50 4B 03 04 | PK.. |

**189**

---

# File Information Tool Set (FITS)
## https://code.google.com/p/fits/

- FITS "identifies, validates, and extracts technical metadata for various file formats. It wraps several third-party open source tools, normalizes and consolidates their output, and reports any errors. FITS was created by the Harvard University Library Office for Information Systems for use in its Digital Repository Service (DRS)."
- Tools currently bundled into it:
    - Jhove
    - Exiftool
    - National Library of New Zealand Metadata Extractor
    - DROID
    - FFIdent
    - File Utility (windows)
- Note: you can find and run FITS from the command line in the BitCurator environment (found in "Additional Tools")

**190**

## File Extensions

- Changing file extension usually changes default application that OS uses to open (i.e. associates with) the file

- The "8.3" (eight characters, followed by three-character extension) limit in the past – based on FAT – resulted in many creative uses of the extension portion of file name (e.g. reports1.994, april-94.rpt)

- Convention is often still to use only three letters

- No authority for standardizing use, so three-letter extensions are often shared by many formats

- Security risks associated with trusting the file extension to be accurate – malicious code masquerading as another type of file (e.g. viruses sent as email attachments)

**191**



**192**

# MIME types

- Widely adopted and recognized by applications
- Based on two-level hierarchy (e.g. text/html, application/octet-stream, image/tiff)
- Major advantage is official registration of MIME types through a central authority

**193**

---

Reproducing In-Application Rendering (Sort of) for Archival Processing



*Viewing a WordPerfect document on a computer that doesn't have WordPerfect installed, using Quick View Plus*

**194**

**Properly Extracting Data from Media**

## Strategies for avoiding accidental manipulation of volatile data

- Use write-blocking equipment when first reading from a medium (hardware, if possible)
- Make bit-level image
- Create checksums before and after file transfers and transformations
- Pay special attention to irreversible changes...

## Examples of Irreversible Changes

- Lossy compression (e.g. JPEG)
- Lower-quality surrogate (e.g. thumbnail image, access copy of video)
- Format conversion (e.g. Word to PDF/A, Excel to CSV)
- Character encoding (e.g. EBCDIC to ASCII)
- Normalization of data values (e.g. date values in a database to a common date encoding)
- Rewriting pointers (e.g. links in a web site from absolute to relative or vice versa)
- Overwriting older versions files or values with newer versions
- Pulling files out of their native file system

**197**

## Write Blocking – One-Way Streets for Data

- Ensures that data can be read from the device, but no bits can be changed
- Doesn't just prevent changes conscious made by user but also changes made by the system
- Options for write blocking (in order of most to least certain to prevent writes to the drive):
  - Dedicated write blockers
  - Writing blocking tabs or settings on the device itself
  - Software-based write blocking



Image source: http://thinng.com/1555-one-way-sign-seat

**198**

# Dedicated Hardware Write Blockers



**199**

---

5.25 Inch Floppy – If light can get through, it's **not** write protected



These black Write Protect Tabs should be used when you do not want stored information on your disk to be changed or lost accidentally.

Simply cover the Write Enable Notch by taking 1 black tab and wrapping it around the notch.

Write Enable Notch

http://en.wikipedia.org/wiki/File:Floppy_tabs_3x2.jpg

**200**

## 3.5 Inch Floppy – If light can get through, it **<u>is</u>** write protected



http://www.techmint.info/2009/09/security-write-protecting-floppy-disks.html

**201**

---

Example of Software Write Blocking – Mounted Devices set to Read-Only by Default



**202**

## Getting below the File System – Low-Level Copying

- Getting an "image" of a storage medium involves working at a level below the file system
  - Can get at file attributes and deleted files not visible through higher-level copy operations
- Most commonly used tool is dd (or variant) - UNIX program for low-level copying and conversion of data from a storage device
- More specialized tools for creating forensic images include:
  - FTK Imager
  - Guymager
  - Imaging utilities in commercial applications (e.g. EnCase)

**203**

Main Acquisition Interface for Guymager



**204**

"The file system provides little if any help when you want to know about details of where and how it stores information; indeed, the entire **purpose of file systems is to hide such detail**. In order to look under the file system, you have to **bypass the file system code** and use tools that duplicate some of the file system's functionality."*

*Farmer, Dan, and Wietse Venema. *Forensic Discovery*. Upper Saddle River, NJ: Addison-Wesley, 2005. (emphasis added)

# Why Make a Bit-Level Image?

# 7 Reasons

**1. Make sure full set of bits is safe**

(allows you still to have the disk but not have to depend on fragile physical medium)

**2. There may be surprises within the structure of the file system**

(e.g. hidden files)

## 3. You could inadvertently change something in the act of examining or dealing with the files

- Byte order
- Character encoding
- File system information
    - □ MAC values
    - □ Access permissions
    - □ File typing
    - □ File sizes
- For example, when using a Windows machine to accession an HFS (Mac) disk, good idea to image the disk right away, so errors in translation across file systems can be noticed and corrected

**209**

## 4. Proof of file integrity and chain of custody

If there are questions about whether a given source was the basis for a given set of digital objects, one can go back to the original bits and compare hash values

**210**

## 5. Corrupted files and viruses

Having the whole bitstream available (in a controlled and safe staging area) makes it possible to determine what subset of the bitstream can actually be recovered in a useful way

## 6. There are likely to be changes in preservation strategy or access conditions over time.

Default ingest process is to create a normalized AIP from a given type of SIP (e.g. convert all Word documents to PDF).
- This is almost certain to lose some information in the process
- Future techniques or access scenarios might require access to the original Word files
- Possibly also information embedded in the file system

## 7. Embedded Contextual Information and User Artifacts

Depending on understanding of arrangement with the Producer, hidden data can also serve as important evidence for the curation of a collection, e.g. traces of data that indicate what application created the files, login or password information that's necessary for accessing various data sources*

*For further discussion of possibilities, see: Garfinkel, Simson, and David Cox. "Finding and Archiving the Internet Footprint." Paper presented at the First Digital Lives Research Conference: Personal Digital Archives for the 21st Century, London, UK, February 9-11, 2009.

**213**

A major rule of digital curation is to minimize irreversible transformations.  Copying files off of the original bitstream and then discarding that bistream is just such an irreversible change (no way to then derive the original bits from the files you have).

It's pretty easy to screw up files or file attributes in numerous ways, and the imaging step can dramatically increase the possibility that those screw ups are reversible.

**214**

The main point:

If you image first and ask questions later, you have a baseline data stream to which you can return if/when necessary.

## Examples of Disk Image formats

- RAW and Split RAW (RAW stored across multiple files)

- Advanced Forensics Format (AFF)

- EnCase Evidence File (.E01)

- ISO (for CD-ROM)

- IMG (floppy)

# RAW (dd)

- Copies of the raw media data. Often split into smaller chunks to make them more manageable and so that the resulting images can fit onto limited filesystems and media such as FAT or DVD/CDROM.
- Advantages:
  - Very simple, use simple tools to manipulate the image.
  - Image can be easily split for storage and transport on removable media
  - Output can be piped to other applications for immediate processing
- Disadvantages:
  - Can be very large (no compression). Zipped raw images cannot be operated on directly with regular tools (efficiently perform arbitrary seeks).
  - Often too large to store on FAT formatted media
  - No metadata other than filenames, no hashes.
  - No checksumming on files – not robust
    - ☐ Missing segments (for example from scratched CD/DVD – can sometimes be overwritten with 0's).
    - ☐ Overwritten data (unrecoverable – no checksums on small blocks in file).

**217**

# AFF

- Original AFF format: single file containing segments with drive data and metadata. Can be compressed.

- Large AFF files can be broken into multiple AFD format files. The smaller AFD files can be readily moved around a FAT32 file system which limits files to 2GB or stored on DVDs, which have similar size restrictions.

- AFM format stores the metadata in an AFF file, and the disk data in a separate raw file. This format allows analysis tools that support the raw format to access the data, but without losing the metadata.

Start of AFF file

AFF Header (8 bytes)

*named*
*metadata*
*segments*

page0 (first data segment)
page_sha1_0 (first segment's sha-1)
page_sha1_sig0 (segment PKCS7 signature)

page1 (second data segment)

page_sha1_1 (second segment's sha-1)
page_sha1_sig1 (segment PKCS7 signature)

...
...
...    *additional segments*
...

sha1 (data file sha-1)
sha1_sig (file PKCS7 signature)

End of AFF file

**AFF File
(Not to Scale)**

Image source: Simson Garfinkel, David Malan, Karl-Alexander Dubec, Christopher Stevens, and Cecile Pham, "AFF An Open Extensible Format for Disk Imaging," *Advances in Digital Forensics II*, edited by Martin S. Olivier and Sujeet Shenoi (New York: Springer, 2006), 13-28.      **218**

# Expert Witness Format – EWF (EnCase)

- Evidence file consists (in order) of: Acquisition information, Data Block, CRC (cyclic redundancy check), acquisition hash (MD5)
- Can be split for storage, transport
- CRC computed for every 32K block; balance between integrity and speed, also makes it very difficult to tamper with the evidence file (1 in 4 billion chance of collision)
- Cannot be manipulated with simple (open source UNIX) tools; support reverse engineered in libewf
- Previously limited to 2GB size
- Largely proprietary
- Has been reverse engineered by Joachim Metz in libewf (used in open source tools that read EWF) - http://sourceforge.net/projects/libewf/files/

Header
Contains: Case info, Acquisition info, Notes, etc.

Data Blocks

CRC's

Acquisition Hash

**219**

---

**220**

**EnCase Disk Image Info**

Computer
- Home
- Desktop
- Documents
- Downloads
- Music
- Pictures
- Videos
- File System
- Trash

Network
- Browse Net...

Home   Desktop

8:43 AM   BitCurator

charlie-work-usb-2009-12-11.E01

**EnCase Disk Image Info**

ewfinfo 20130416

Acquiry information
    Acquisition date: Wed Jan 19 12:09:18 2011
    System date:      Wed Jan 19 12:09:18 2011
    Operating system used:  Linux
    Software version used:  20100226
    Password:         N/A

EWF information
    File format:      EnCase 6
    Sectors per chunk:   64
    Error granularity: 64
    Compression method:   deflate
    Compression level:   best compression
    Set identifier:    4eb6701d-6cf0-2f4a-a0c6-0cb5d5e20959

Media information
    Media type:       fixed disk
    Is physical:      yes
    Bytes per sector: 512
    Number of sectors:   2068480
    Media size:       1010 MiB (1059061760 bytes)

Digest hash information
    MD5:          9c0de6c8532d7a66ddcf01861dfb6535

Cancel   OK

**221**

"charlie-work-usb-2009-12-11.E01" selected (9.3 MB)

---

# Four Ways to Interact with Disk Images

- Emulation
- Mount them like regular drives:
  - For ISO images - disk utilities in Mac OS or Windows 8/10
  - For forensically packaged disk images: ewfmount, OSFMount, BitCurator (mounting scripts built into the environment)
- Inspect them as forensic objects
  - FTK Imager
  - Autopsy
  - BitCurator (Disk Image Access tool)
- Dynamically navigate them from within a web browser (BCA Webtools)

**222**

Emulation as a Service

http://bw-fla.uni-freiburg.de/demos.html

223

Mounting a Forensically Packaged Disk Image in the BitCurator Environment



224

## Exporting Files from a Disk Image



**225**

## Viewing Contents of a Disk Image (.E01 format) in FTK Imager



**226**

**Forensic analysis meets the archives**

---

## Guidelines for Evidence Collection & Archiving (RFC 3227) – Main Lessons

- "Such collection represents a considerable efforts on the part of the System Administrator."
- "Keep detailed notes."
- "Minimize changes to the data as you are collecting it."
- "Do collection first and analysis later."
- "Proceed from the volatile to the less volatile."
- Computer evidence should be: *admissible*, *authentic*, *complete*, *reliable*, *believable*

## Digital Resources - Levels of Representation

| Level | Label | Explanation |
|---|---|---|
| 8 | Aggregation of objects | Set of objects that form an aggregation that is meaningful encountered as an entity |
| 7 | Object or package | Object composed of multiple files, each of which could also be encountered as individual files |
| 6 | In-application rendering | As rendered and encountered within a specific application |
| 5 | File through filesystem | Files encountered as discrete set of items with associate paths and file names |
| 4 | File as "raw" bitstream | Bitstream encountered as a continuous series of binary values |
| 3 | Sub-file data structure | Discrete "chunk" of data that is part of a larger file |
| 2 | Bitstream through I/O equipment | Series of 1s and 0s as accessed from the storage media using input/output hardware and software (e.g. controllers, drivers, ports, connectors) |
| 1 | Raw signal stream through I/O equipment | Stream of magnetic flux transitions or other analog electronic output read from the drive without yet interpreting the signal stream as a set of discrete values (i.e. not treated as a digital bitstream that can be directly read by the host computer) |
| 0 | Bitstream on physical medium | Physical properties of the storage medium that are interpreted as bitstreams at Level 1 |

229

## Digital Resources - Levels of Representation

| Level | Label | Explanation |
|---|---|---|
| 8 | Aggregation of objects | Set of objects that form an aggregation that is meaningful encountered as an entity |
| 7 | Object or package | Obje... also ... |
| 6 | In-application rendering | As re... |
| 5 | File through filesystem | Files... paths |
| 4 | File as "raw" bitstream | Bitstr... value |
| 3 | Sub-file data structure | Discr... |
| 2 | Bitstream through I/O equipment | Serie... using... contr... |
| 1 | Raw signal stream through I/O equipment | Stream of magnetic flux transitions or other analog electronic output read from the drive without yet interpreting the signal stream as a set of discrete values (i.e. not treated as a digital bitstream that can be directly read by the host computer) |
| 0 | Bitstream on physical medium | Physical properties of the storage medium that are interpreted as bitstreams at Level 1 |

**Levels where digital forensics methods and tools can provide a lot of assistance**

230

## Digital Forensics Industry and Tools

- Commercial products:
  - EnCase (Guidance Software)
  - FTK (AccessData)
- Open source tools – see:
  - http://www.sleuthkit.org/autopsy/
  - http://bitcurator.net
  - http://www.forensicswiki.org

231

## Commercial products

- Pros
  - GUI for non-specialists
  - Powerful search and indexing
  - Bookmarking, annotation, reporting
- Cons
  - Designed for law enforcement
  - High cost, low uptake in archival community

232

## Open source tools: Autopsy

- GUI for The Sleuth Kit (TSK) – software suite also used by BitCurator
- Available on Mac/Linux/Windows
- For "law enforcement, military, and corporate examiners"
- See http://www.sleuthkit.org/autopsy/

233

# BitCurator

- Funded by Andrew W. Mellon Foundation
  - Phase 1: October 1, 2011 – September 30, 2013
  - Phase 2 – October 1, 2013 – September 30, 2014
- Partners: School of Information and Library Science (SILS) at UNC and Maryland Institute for Technology in the Humanities (MITH)

234

# BitCurator Goals

- Develop a system for professionals in libraries, archives and museums that incorporates the functionality of open-source digital forensics tools
- Address two fundamental needs not usually addressed by the digital forensics industry:
  - ☐ incorporation into the workflow of archives/library ingest and collection management environments
  - ☐ provision of public access to the data

**235**

# BitCurator Environment*

- Bundles, integrates and extends functionality of open source software
- Can be run as:
  - ☐ Self-contained environment (based on Ubuntu Linux) running directly on a computer (download installation ISO)
  - ☐ Coming very soon: installation scripts to turn any Ubuntu Linux machine into a BitCurator Environment
  - ☐ Self-contained Linux environment in a virtual machine using e.g. Virtual Box or VMWare
  - ☐ As individual components run directly in your own Linux environment or (whenever possible) Windows environment

*To read about and download the environment, see: http://wiki.bitcurator.net/    **236**

## BitCurator-Supported Workflow



- **Acquisition**
- **Reporting**
- **Redaction**
- **Metadata Export**

See: http://bitcurator.net

237

---

# Digging Deeper – Unallocated Space and Data Carving

238

## Unallocated Data Segments in Storage and File Slack

- Recall that computers store data within storage segments of given size (clusters)
- Unallocated storage
  - When file is deleted, pointer to data in each segment is removed, but data will sit in that segment until overwritten by something else (i.e. clusters are freed up for further use, but data within them can persist for a while)
  - Special case is when file is only partially written to a medium (e.g. failed copy due to insufficient space) and then "deleted" but data remains on medium until overwritten
- Overwriting of storage across a computer system is not uniform – "most activity accesses the same data, programs, and other resources over and over again"*

*Farmer, Dan, and Wietse Venema. *Forensic Discovery*. Upper Saddle River, NJ: Addison-Wesley, 2005.

**239**

---

Consider an NTFS file system with a 2048-byte cluster and 512-byte sectors. Our file is 612 bytes, so it uses the entire first sector and 100 bytes of the second sector in the cluster. The remaining 412 bytes of the second sector are padded with data of the OSes choice. The third and fourth sectors may be wiped with zeros by the OS, or they might not be touched and might keep the data from a deleted file. We can see this in Figure 8.9, where the grayed areas are the file content and the white space is the slack space.

**Figure 8.9. Slack space of a 612-byte file in a 4096-byte cluster.**

Cluster 4910



Sector 1    Sector 2    Sector 3    Sector 4

Source: Carrier, Brian. *File System Forensic Analysis*. Boston, MA: Addison-Wesley, 2005.

**240**

## Data in "Slack" Space

- Storage space from the end of the file to the end of the last cluster used to store the file
  - **RAM slack** (in last sector of file) - If there is not enough data in file to fill the last sector, OS fills the space with whatever data is in memory at the moment (can be anything since computer was last booted) – no longer happens in recent versions of Windows (overwrites with zeros)
  - **Drive slack** (in last cluster of file) - If data from a file does not completely fill the cluster, then data remaining from previous files will remain at the end of the cluster

**241**

## Data Carving

- Extracting data from raw data blocks, rather than examining the contents of files through the file system.
- Can be done to reconstruct files, identify deleted information, and find data that was purposely hidden in various ways

**242**

# Carving Taxonomy

**Carving:** Extract data (files) from raw data (blocks).

**Block-Based Carving:** Analyze input on block-by-block basis to determine if a block is part of a possible output file.

**Statistical Carving:** Algorithm that analyzes the input on characteristic or statistic for example, entropy) to determine if the input is part of a possible output file.

**Header/Footer Carving:** A method for carving files out of raw data using a distinct header (start of file marker) and footer (end of file marker).

**Header/Maximum (file) size Carving:** A method for carving files out of raw data using a distinct header (start of file marker) and a maximum (file) size.

**Header/Embedded Length Carving:** A method for carving files out of raw data using a distinct header and a file length (size) which is embedded in the file format

**File structure based Carving:** A method for carving files out of raw data using a certain level of knowledge of the internal structure of file types.

**Semantic Carving:** A method for carving files based on a linguistic analysis of the file's content.

**Fragment Recovery Carving (split carving):** A carving method in which two or more fragments are reassembled to form the original file or object.

**Repackaging Carving:** A carving method that modifies the extracted data by adding new headers, footers, or other information so that it can be viewed with standard utilities.

*Source: http://www.forensicswiki.org/wiki/File_Carving*

**243**

---

# Digging as Deep as Possible – Physical Medium

**244**

"Because of the way magnetic media are written it is very difficult to lose everything. With sufficient resources much material that most of us would expect to be lost can be recovered."*

*Ross, Seamus, and Ann Gow. "Digital Archaeology: Rescuing Neglected and Damaged Data Resources." London: British Library, 1999.

245



Veeco Instruments. http://www.veeco.com/library/nanotheater_detail.php?type=application&id=78&app_id=34

246

Note: Finding over-written data from the surface of magnetic media is **<u>extremely</u>** slow and expensive – only applicable in cases of "heroic" recovery

Identifying "Features" of Interest in Disk Images

Bulk Extractor

## Identifying Potentially Sensitive Data using Bulk Extractor - Scanning Options



See: http://www.forensicswiki.org/wiki/Bulk_extractor

249

## Histogram of Email Addresses (Specific Instances in Context on Right)



250

## Bulk Extractor Output*

| File | Description |
|------|-------------|
| aes_keys.txt | AES encryption keys |
| alerts.txt | Processing errors |
| ccn.txt | Credit card numbers |
| ccn_track2.txt | Credit card "track 2" information, which has previously been found in some bank fraud cases |
| domain.txt | Internet domains found on the drive, including dotted-quad addresses found in text |
| email.txt | Email addresses |
| ether.txt | Ethernet MAC addresses found through IP packet carving of swap files and compressed system hibernation files and fragments |
| exif.txt | EXIF data from JPEG images and video segments |
| find.txt | Results of specific regular expression searches |
| gps.txt | Extracted GPS coordinates from Garmin XML and GPS-enabled JPEG files |
| ip.txt | IP addresses found through IP packet carving |
| json.txt | Extracted and validated JavaScript Object Notation fragments |
| kml.txt | Extracted KML files |

*See http://afflib.org/archives/tag/bulk_extractor

**251**

## Bulk Extractor Output (continued)*

| File | Description |
|------|-------------|
| report.txt | DFMXL file that explains what happened |
| rfc822.txt | Email message headers including Date:, Subject:, and Message-ID: fields |
| tcp.txt | TCP flow information found through IP packet carving |
| telephone.txt | Phone numbers (US and other countries) |
| url.txt | URLs, typically found in browser caches, email messages, and pre-compiled into executables |
| url_searches.txt | Histogram of terms used in Internet searches |
| url_services.txt | Histogram of the domain name portion of all URLs found on the media |
| winpefect.txt | Windows prefetch files and fragments, recorded as XML |
| wordlist.txt | A list of all "words" extracted from the disk, useful for password cracking |
| wordlist_*.txt | The wordlist with duplicates removed, formatted to be imported into a popular password-cracking program |
| zip.txt | Information about ZIP file components found on media (including compound files such as MS Office documents) |

*See http://afflib.org/archives/tag/bulk_extractor

**252**

Technical Metadata (about the System Used to do the Capture) in a Bulk Extractor Report



253

# A Real World Example of Forensic Feature Extraction:



Jeb Bush dumps emails including social security numbers of Florida residents online

*Florida man strikes again*

By T.C. Sottek on February 10, 2015 01:37 pm

254

https://www.theverge.com/2015/2/10/8013531/jeb-bush-florida-email-dump-privacy

## A Real World Example of Forensic Feature Extraction:

**kamwoods**
@kamwoods

🐦 Follow

Ran a few tools over the Jeb Bush emails. And...yeah. Pages of SSNs, DOBs, CCNs in the output.

1:30 PM - 10 Feb 2015

**144** RETWEETS **47** FAVORITES

↩ ⇄ ★

255

## A Real World Example of Forensic Feature Extraction:



256

## Conclusions: Implied Changes with the Archival Profession

- Professional vocabulary evolving to include terms such as disk image, hex viewer, cryptographic hash, and filesystem
- Gaining access to new professional communities and guidance
- Use of tools designed to treat data at a low level – as raw bitstreams off media – rather than at the file level
- Potential to shift "center of gravity" about electronic records from design of institutional recordkeeping systems toward acquisition and management of records from a more diverse and unpredictable set of sources

**257**

## For further Guidance

- See supplemental materials
- Digital forensics vendors offer workshops (e.g. AccessData, Digital Intelligence, Guidance Software)
- http://www.forensicswiki.org/

**258**

## BitCurator Resources

Get the software
Documentation and technical specifications
Screencasts
Google Group
**http://wiki.bitcurator.net/**

People
Project overview
Publications
News
**http://www.bitcurator.net/**

Twitter: @bitcurator

# BitCurat●r Access

- Web-based access to raw and forensically packaged disk images
- Transforming and using digital forensics metadata in collecting environments
- Redaction of file items, metadata and hidden data from disk images
- OS and executable virtualization for legacy disk images

http://bitcuratorconsortium.org

**261**



http://www2.archivists.org/groups/manuscript-repositories-section/jump-in-3 **262**

# To Learn More about Available Software

Forensics Wiki. http://forensicswiki.org

BitCurator environment. http://wiki.bitcurator.net

BitCurator Access software. http://access.bitcurator.net

Community Owned digital Preservation Tool Registry (COPTR)
http://coptr.digipres.org/

Information Guides on Tools for Electronic Records. Minnesota State Archives.
http://www.mnhs.org/preserve/records/tools.php

Lifecycle Tools for Archival Email Stewardship.
https://docs.google.com/spreadsheets/d/1V1N22xnr5e0EbDlZWx58bjYO6rkrMrYH9wGX9-CK8c4/

Tools for processing, managing, and preserving electronic records. University of Minnesota.
https://www.lib.umn.edu/dp/guides

**263**

---

# Online Forums

## BitCurator User Group
https://groups.google.com/forum/#!forum/bitcurator-users



## Digital Curation List
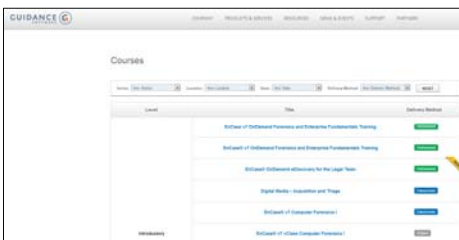https://groups.google.com/forum/#!forum/digital-curation

**264**

# Further Education

## Digital Forensics for Archivists: Advanced (SAA)

http://www2.archivists.org/prof-education/course-catalog/digital-forensics-for-archivists-advanced#.V1SSg-S0OzA



**265**

---

# Thank you!

Digital Archives Specialist (DAS) Questions?
http://www2.archivists.org/prof-education/das/FAQs
education@archivists.org

**266**