# Introduction to BitCurator

BitCuratorEdu
Last Updated: August 11, 2020

# Many information professionals 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/
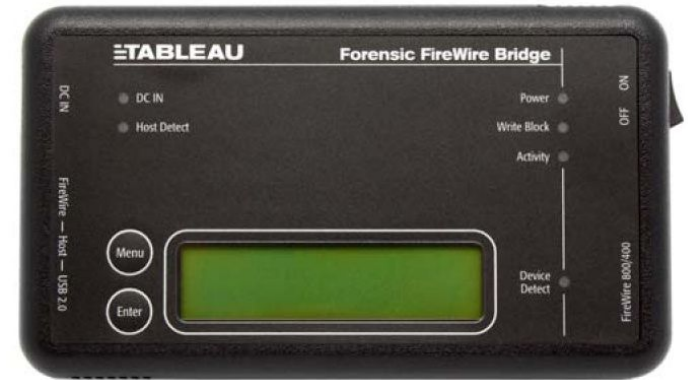
# Same Goals as When Acquiring Analog Materials

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

# 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.

# But you might need some of this stuff:

# 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
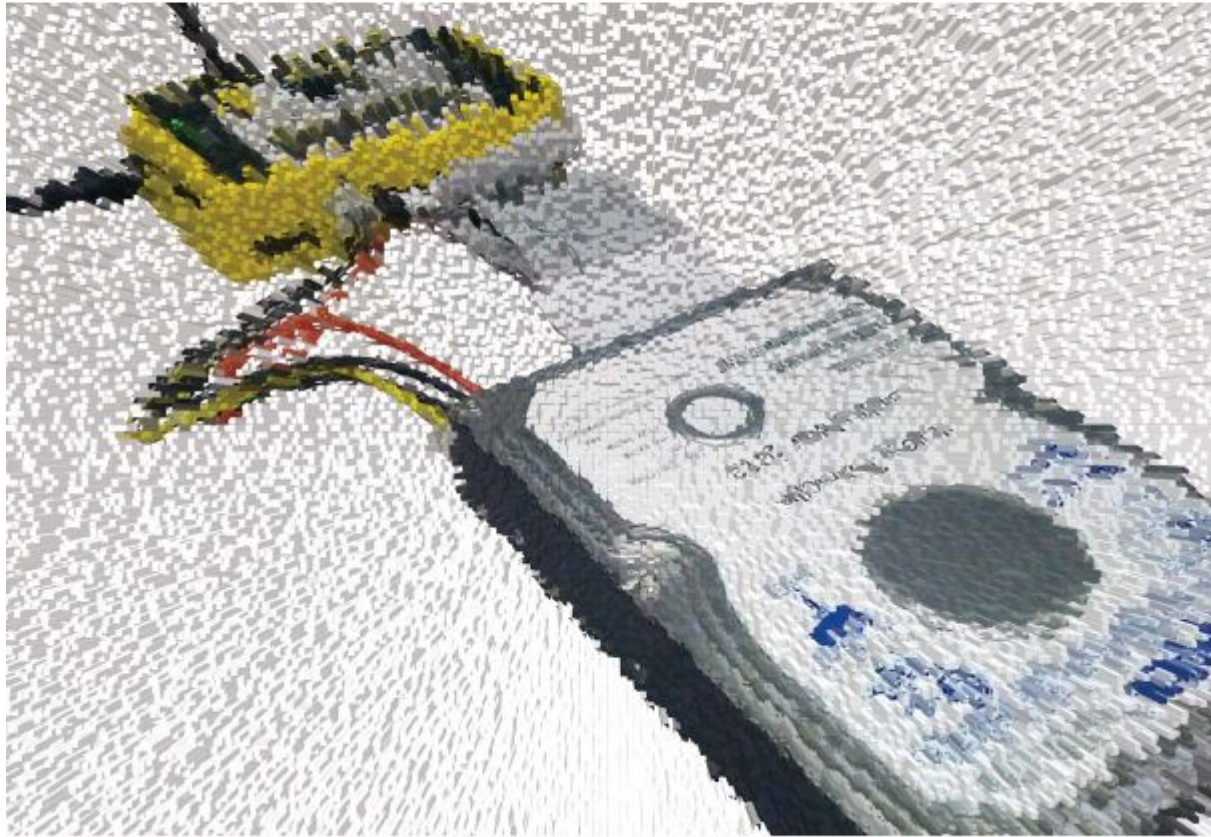
# Digital Forensics Can Help Archivists to Fulfill their Principles

| | |
|---|---|
| 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 |

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

# 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 |

# 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 |

ContextMiner Alpha 3.0

[Home][Publications][Reports][Add][View][Search][Profile][Visualize][Monitor][Tools][Developer]

This page lists all the seed queries that are used for monitoring videos related to elections on YouTube. Clicking on a query will show all the results collected over several crawls. Total number of these results are also listed here for each query. The last column in the following table shows how many total results YouTube had for a given query during our latest crawl. Clicking on 'Setup' associated with a query will bring up an interface where the curator can specify what constitutes as a "significant" change for a video of that query.

| # | Query | Setup | Total results so far | Max results on last crawl |
|---|---|---|---|---|
| 1 | election 2008 | Setup | 574 | 6150 |
| 2 | US election 2008 | Setup | 349 | 795 |
| 3 | United States election 2008 | Setup | 216 | 257 |
| 4 | presidential election 2008 | Setup | 206 | 1820 |
| 5 | campaign 2008 | Setup | 273 | 2530 |
| 6 | decision 2008 | Setup | 168 | 142 |
| 7 | Joe Biden | Setup | 209 | 1080 |
| 8 | Hillary Rodham Clinton | Setup | 193 | 353 |
| 9 | Christopher Dodd | Setup | 267 | 815 |
| 10 | John Edwards | Setup | 902 | 7540 |
| 11 | Mike Gravel | Setup | 301 | 1210 |
| 12 | Dennis Kucinich | Setup | 229 | 1600 |
| 13 | Barack Obama | Setup | 861 | 9140 |
| 14 | Bill Richardson | Setup | 287 | 1100 |
| 15 | Wesley Clark | Setup | 191 | 375 |
| 16 | Al Gore | Setup | 613 | 4910 |
| 17 | Tom Vilsack | Setup | 89 | 68 |
| 18 | Sam Brownback | Setup | 254 | 404 |

# 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 throu... equipment |
| Bitstream on physical m... |

ContextMiner Alpha 3.0

[Home][Publications][Reports][Add][View][Search][Profile][Visualize][Monitor][Tools][Developer]

This page presents contextual information for a video captured over a number of days. Contextual information is defined as the information about a video that may change with time. Usually this information is contributed by the visitors of the video page. See the metadata information for this video. Description of various attributes displayed is given here.

Query: *Rudy Giuliani*
I Got A Crush On.... Giuliani
Collaboration with the very talented JackDanyells, who came up with the concept for this video. Check out his channel at: http://www.youtube.com/jackdanyells -Lyrics by JackDanyells -Vocal melody composed and sung by me -Royalty free background music from sounddogs.com
Comedy
Crawling since 2007-07-19

Color coding for % changes
< | 0.05 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1.0 | 5.0 | >

| Crawl # | Crawl date | Rank | Views | Ratings | Avg Rating | Comments | Links | Favorited | Honors | Change |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2007-07-31 | 5 | 27357 | 301 | 3.74 | 288 | 5 | 44 | 0 | -- |
| 2 | 2007-08-01 | 5 | 27452 | 303 | 3.73 | 290 | 5 | 44 | 0 | -- |
| 3 | 2007-08-02 | 5 | 27780 | 307 | 3.72 | 291 | 5 | 45 | 0 | -- |
| 4 | 2007-08-03 | 5 | 28048 | 309 | 3.71 | 291 | 5 | 45 | 0 | -- |
| 5 | 2007-08-04 | 2 | 28398 | 310 | 3.71 | 291 | 5 | 45 | 0 | -- |
| 6 | 2007-08-05 | 2 | 28443 | 314 | 3.69 | 294 | 5 | 45 | 0 | -- |
| 7 | 2007-08-06 | 3 | 28980 | 314 | 3.69 | 296 | 5 | 45 | 0 | -- |
| 8 | 2007-08-07 | 3 | 29265 | 318 | 3.65 | 298 | 5 | 45 | 0 | -- |
| 9 | 2007-08-08 | 3 | 29551 | 319 | 3.65 | 299 | 5 | 46 | 0 | -- |
| 10 | 2007-08-09 | 3 | 30094 | 320 | 3.64 | 300 | 5 | 47 | 0 | -- |
| 11 | 2007-08-10 | 3 | 30384 | 323 | 3.61 | 302 | 5 | 47 | 0 | -- |
| 12 | 2007-08-10 | 5 | 30419 | 324 | 3.62 | 303 | 5 | 48 | 0 | -- |
| 13 | 2007-08-11 | 3 | 30540 | 324 | 3.62 | 305 | 5 | 49 | 0 | -- |
| 14 | 2007-08-12 | 3 | 30697 | 326 | 3.61 | 306 | 5 | 49 | 0 | -- |
| 15 | 2007-08-13 | 3 | 30848 | 326 | 3.61 | 306 | 5 | 49 | 0 | -- |
| 16 | 2007-08-14 | 3 | 31036 | 326 | 3.61 | 306 | 5 | 49 | 0 | -- |
| 17 | 2007-08-15 | 2 | 31181 | 326 | 3.61 | 306 | 5 | 49 | 0 | -- |
| 18 | 2007-08-16 | 2 | 31321 | 326 | 3.61 | 307 | 5 | 51 | 0 | -- |
| 19 | 2007-08-17 | 2 | 31459 | 327 | 3.61 | 307 | 5 | 51 | 0 | -- |
| 20 | 2007-08-18 | 2 | 31662 | 331 | 3.59 | 308 | 5 | 51 | 0 | -- |
| 21 | 2007-08-19 | 2 | 31792 | 332 | 3.58 | 308 | 5 | 51 | 0 | -- |
| 22 | 2007-08-20 | 2 | 31937 | 335 | 3.57 | 310 | 5 | 51 | 0 | -- |
| 23 | 2007-08-21 | 2 | 32135 | 335 | 3.57 | 311 | 5 | 52 | 0 | -- |

# Interaction Examples

**Level**

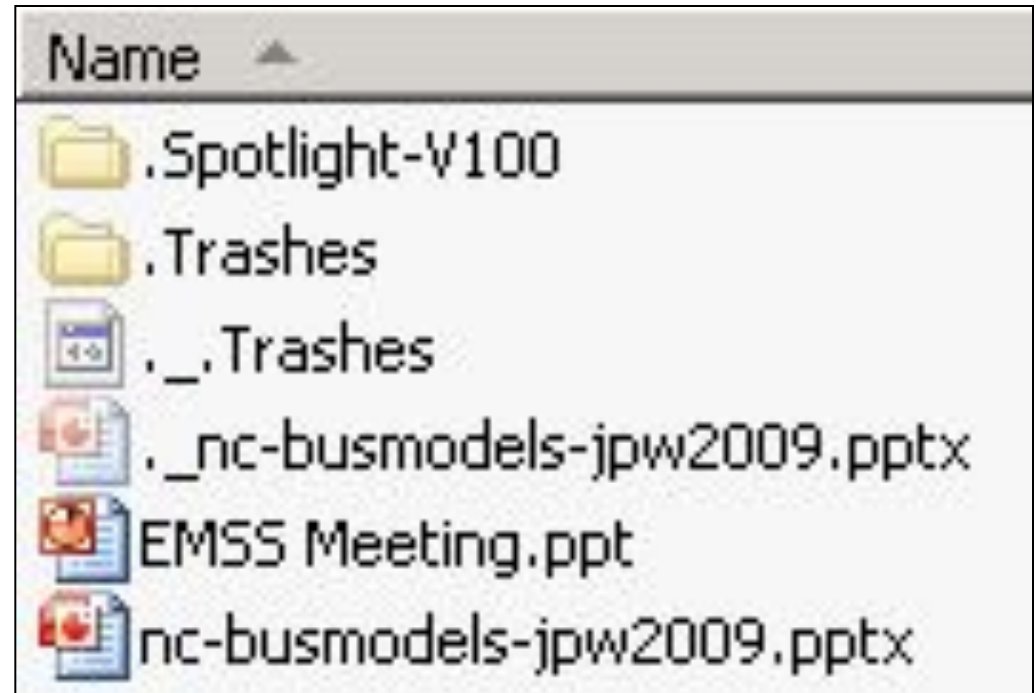| |
|---|
| Aggregation of objects |
| Object or package |
| **In-application rendering** |
| File through filesystem |
| File as "raw" bitstream |
| Sub-file data structure |
| Bitstream through I/O equipment |
| Raw signal stream through I/O equipment |
| Bitstream on physical medium |

**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 |
| |





14

# Interaction Examples

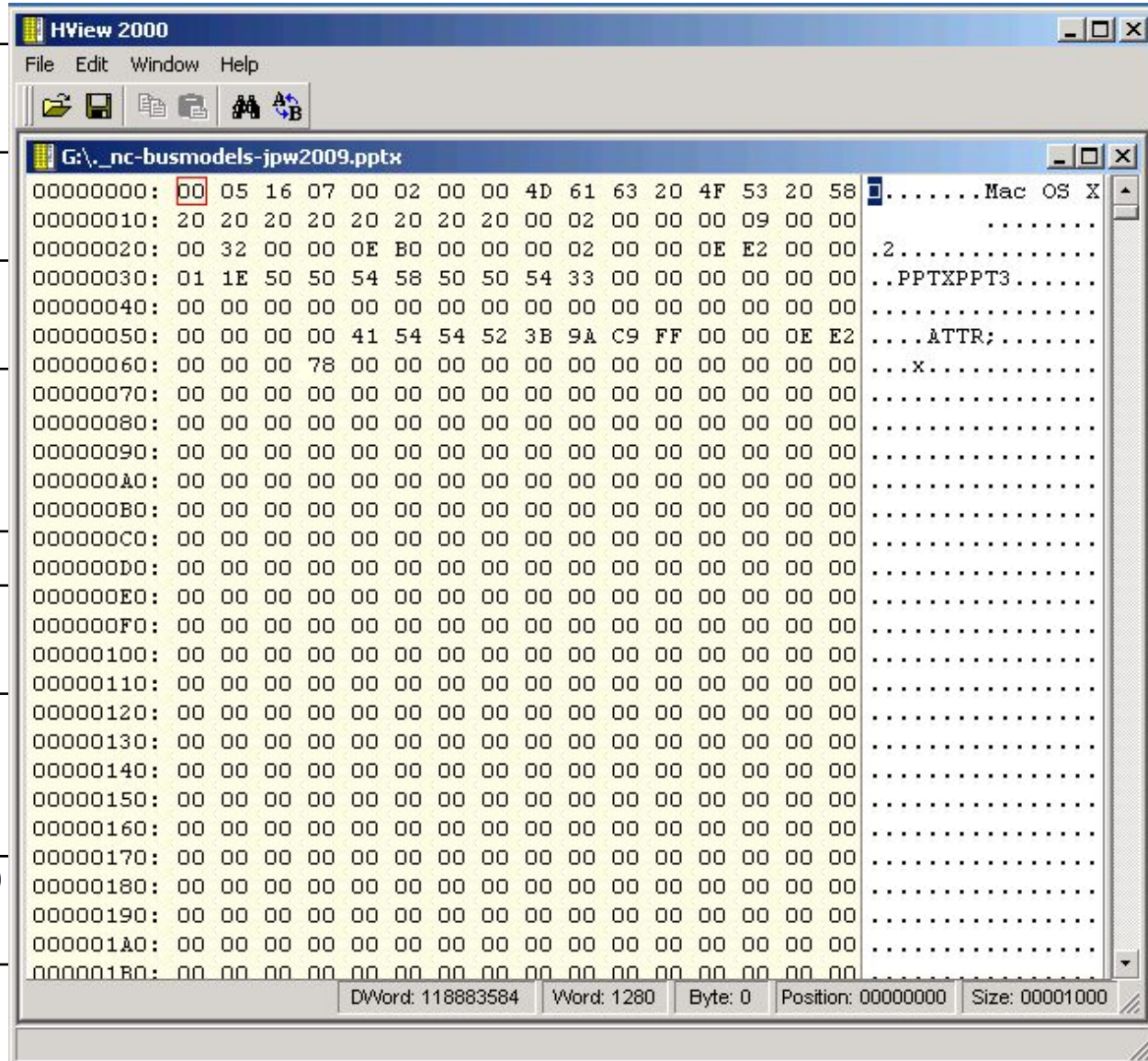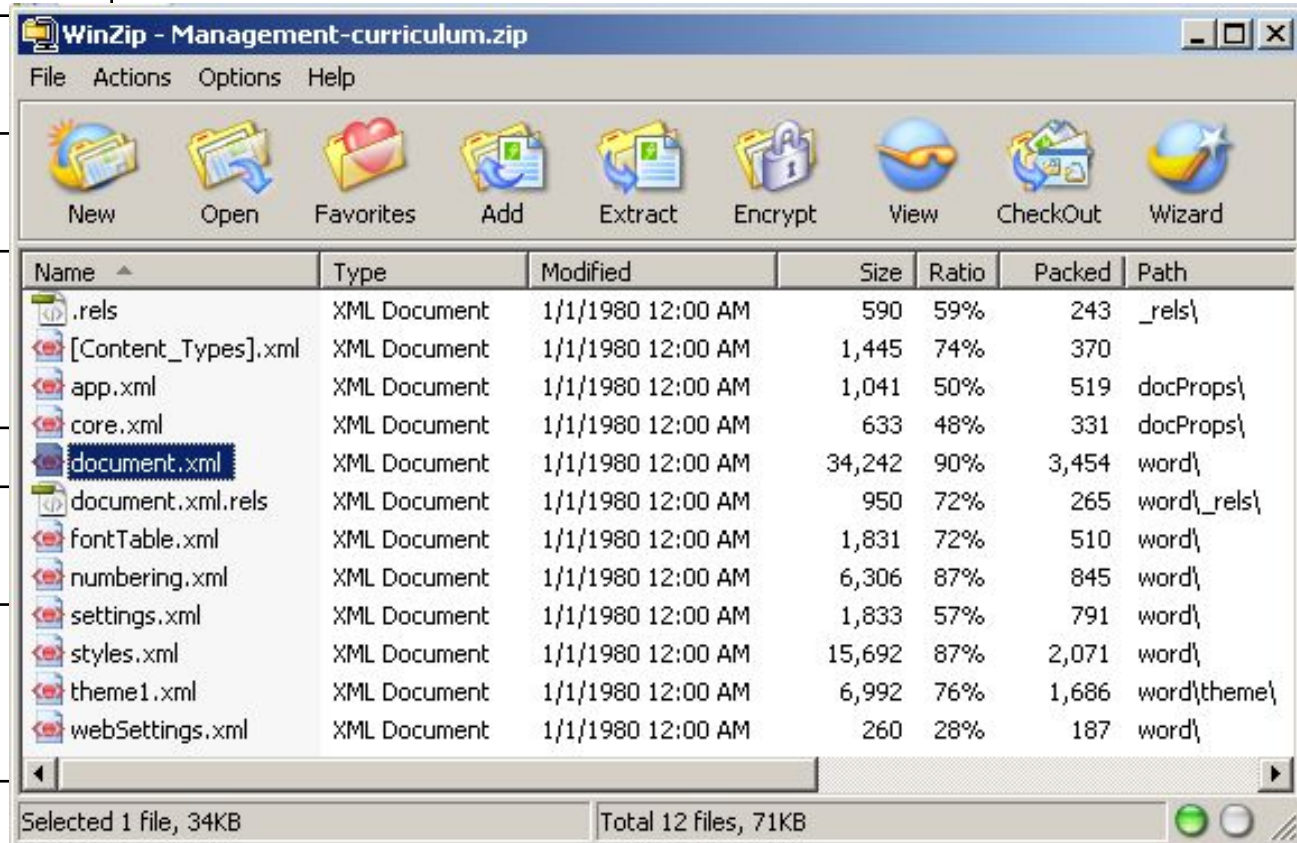| Level |
|---|
| Aggregation of objects |
| Object or package |
| In-application rendering |
| File through filesystem |
| **File as "raw" bitstream** |
| Sub-file data structure |
| Bitstream through I/O equipment |
| Raw signal stream through I/O equipment |
| Bitstream on physical medium |

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 20 02 00 00 00 09 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 41 54 54 52 3B 9A C9 FF 00 00 0E E2  ....ATTR;.......
00000060: 00 00 00 78 00 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

15

# Interaction Examples

**Level**

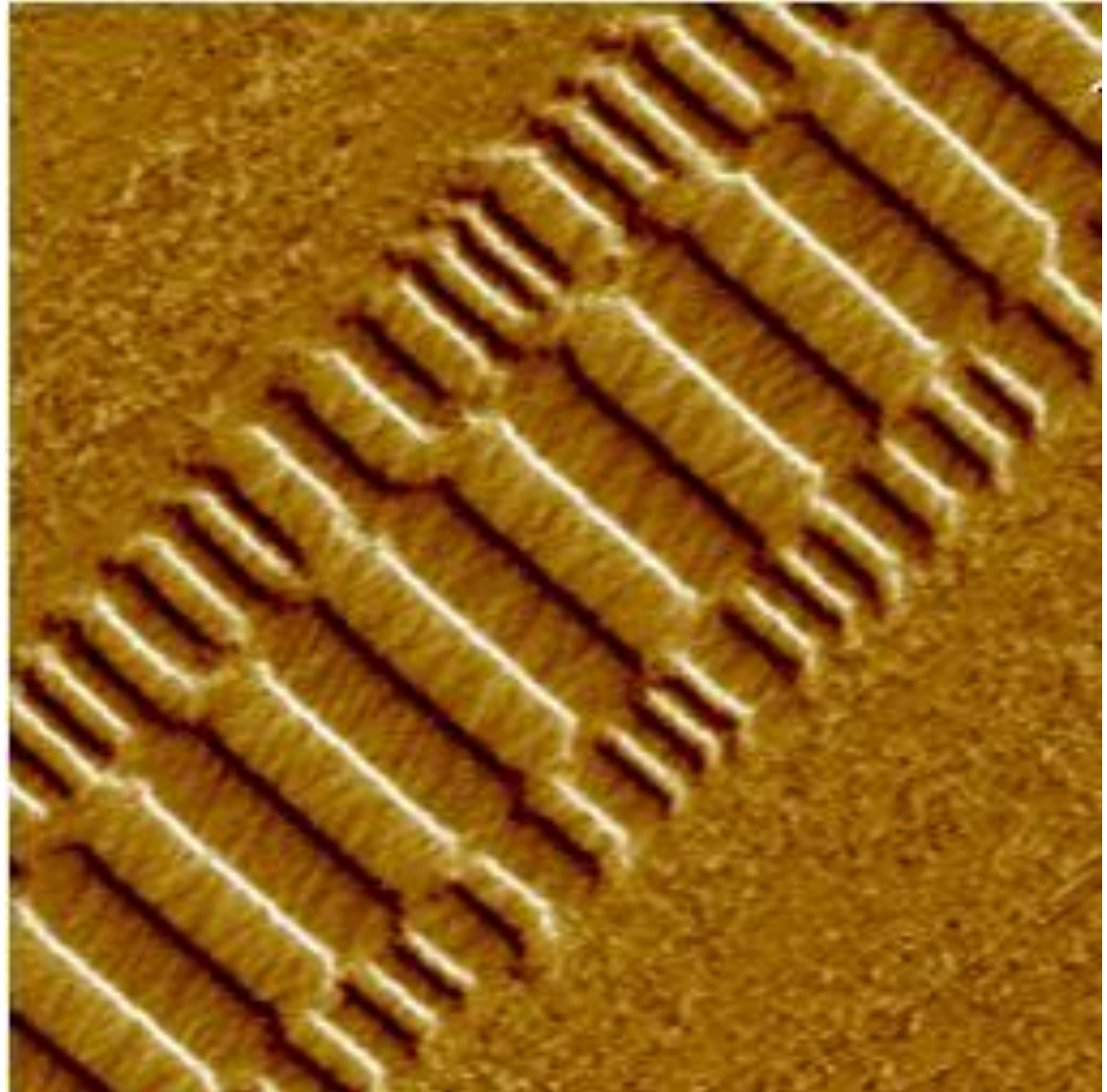| |
|---|
| Aggregation of objects |
| Object or package |
| In-application rendering |
| File through filesystem |
| File as "raw" bitstream |
| **Sub-file data structure** |
| Bitstream through I/O equipment |
| Raw signal stream through equipment |
| Bitstream on physical medium |



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

# Interaction Examples

**Level**

| |
|---|
| Aggregation of objec |
| 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 |

**Level**

| 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 |



18

| | **Examples** |
|---|---|
| | Browsing the contents of an archival collection using a finding aid |
| | Viewing a web page that contains several files, including HTML, a style sheet and several images |
| In-application rendering | Using Microsoft Excel to view an .xls file, watching an online |
| File through filesystem | Windows Explorer, typing to show the contents of a |
| File as "raw" bitstream | Opening an individual file in a hex editor |
| Sub-file data structure | Extra value |
| Bitstream through I/O equipment | Conn gene comm |
| **Raw signal stream through I/O equipment** | Connecting a floppy drive to a host computer and then generating a magnetic flux transition image of the disk |
| Bitstream on physical medium | Using a high-power microscope and camera to take a picture arges on the surface of a hard drive or pits and lands on an optical disk |

http://www.pagetable.com/?p=32

# 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

# 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)

# BitCurator Goals

- Develop a system for collecting professionals 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

# BitCurator Environment*

- Bundles, integrates and extends functionality of open source software

- Can be run as:

  - Self-contained environment running directly on a computer (download installation ISO)

  - Using "bootstrapping" 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:
https://github.com/BitCurator/bitcurator-distro

# BitCurator Consortium

- Continuing home for hosting, stewardship and support of BitCurator tools and associated user engagement

- Administrative home: Educopia Institute

- Funding based on membership dues

- Software and documentation are free and open source, but membership provides benefits (e.g. support, training, consulting)

*https://bitcuratorconsortium.org/*

## A Growing Community

The BitCurator Consortium provides spaces for members to share documentation, develop their skills, and improve the BitCurator environment.

Membership is open >

---

Membership is open to libraries, archives, museums, and other institutions worldwide that seek a collaborative community within which they may explore and apply forensics approaches and solutions to their digital collections.

Become a member now >

## How to Use BitCurator

- Acquire and process digital collections.
- Maintain the original order of digital materials.
- Survey the extent and composition of digital collections.
- Redact personally identifiable information.
- Extract technical and preservation metadata.
- Package digital materials for archival storage.

Learn more about getting started.

## How our members are using BitCurator

## Member Benefits

- Use of the members-only BCC mailing list and help desk
- Access to the members-only videos and documentation
- Prioritized requests for BitCurator feature development
- Opportunities to serve on the BCC committees
- Voting rights for community governance
- Professional development opportunities
- Discounts for events including the BitCurator User Forum

## Members

McMaster University

Penn State University

Massachusetts Institute of Technology

Duke University

The University of Maryland, MITH

Stanford University

Yale University

The University of Manchester Library

University of

# BitCurator Consortium: Fostering Community

- ## Communication
  - Monthly community calls
  - Listserv
  - Maintains documentation feat. community scripts and data set libraries

- ## Active Subgroups
  - Software Development
  - Program
  - Membership Working Group
  - Executive Council

- ## Events
  - Mixers at various professional conferences
  - Annual User Forum

# BitCuratorEdu (2018-2021)

- **Partners:** University of North Carolina at Chapel Hill School of Information and Library Science (UNC SILS), Educopia Institute, BitCurator Consortium, and the Council of State Archivists (CoSA)

- **Purpose:** study and advance adoption of digital forensics tools and methods in libraries and archives through professional education

- **Research Questions**:
  - What are the primary **institutional and technological factors** that influence adoption of digital forensics tools and methods in LIS classes in different educational settings?
  - What are the most viable mechanisms for **sustaining collaboration** among LIS programs on the adoption of digital forensics tools and methods?

- **Objectives:**
  - **produce and disseminate** learning materials
  - **investigate and report** on institutional factors to facilitate, hinder and shape adoption of educational offerings
  - **advance** community of practice around digital forensics education

# Advisory Board

| | |
|---|---|
| Catholic University | Jane Zhang, Associate Professor |
| Indiana University | Devan Donaldson, Assistant Professor |
| New York University | Howard Besser, Professor, Associate Director of MIAP, and Senior Scientist for Digital Library Initiatives for NYU Library |
| San Jose State University | Sandra Hirsh, Professor and Director of the School of Information; Alyce Scott, Lecturer |
| University of Illinois | Rhiannon Bettivia, Postdoctoral Research Associate |
| University of Maryland | Ricky Punzalan, Assistant Professor at iSchool, Affiliate Assistant Professor in Anthropology, and Co-Director of Museum Scholarship and Material Culture Program |
| University of Michigan | Paul Conway, Associate Professor |
| University of Texas | Patricia Galloway, Professor |
| Wayne State University | Kimberly Schroeder, Lecturer |

# BitCurator-Supported Workflow

## Acquisition

Source media

Additional supported output formats
- Split raw
- E01

Imaging (aimage, guymager)

AFF packaged image → Capture and image metadata → Log, device info

## Reporting

Filesystem metadata extraction (fiwalk) → Filesystem report, DFXML

File <-> disk block map

Image analysis (private, sensitive info) → Bulk Extractor -- Human-readable reporting scripts / plug-in

Bulk Extractor → Annotated feature file

Reports: file distribution, sensitive info, hashes, etc.
Exported as PDF, .xlsx, plain text as appropriate

Analysis (accounts, filesystem activity) → reg2xml (and similar) → XML dump of registry

Reports: user accounts, device usage, environment
Exported as PDF, .xlsx, plain text as appropriate

Analysis (file similarity, deduplication) → sdhash

Reports: locations, status of similar files
Exported as PDF, .xlsx, plain text as appropriate

## Redaction

Redaction of disk image → Ruleset (patterns, hashes, etc) describing what to redact → Python redaction scripts

## Metadata export

Metadata export → Python (lxml) export scripts

METS, MODS, EAD as required by user

---

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

See: http://bitcurator.net

# For Further Information



https://bitcurator.github.io/

Most of the tasks we cover in this class are explained in the Quick Start Guide. The most recent version is always available at:
https://github.com/BitCurator/bitcurator-distro/wiki/Releases

# BitCurator
## Quick Start Guide

Last updated: August 1, 2018
Release(s): 2.0.4 and later

UNC
SCHOOL OF INFORMATION
AND LIBRARY SCIENCE

BitCurator
CONSORTIUM

En  2:41 PM

home

Imaging Tools

Forensics Tools

System Settings

Additional Tools

Shared Folders and
Media

sf_bc_share

Trash

Documentation
and Help

Network Servers

BitCurator

Right Ctrl

# Creating and Extracting Forensic Metadata

# High-Level view of Metadata Generation and Reporting



See: Woods, Kam, Christopher Lee, and Sunitha Misra. "Automated Analysis and Visualization of Disk Images and File Systems for Preservation." In *Proceedings of Archiving 2013* (Springfield, VA: Society for Imaging Science and Technology, 2013), 239-244.

BitCurator-0.3.0 [Running]

Bulk Extractor Viewer

2:00 AM  BitCurator

File Edit View Tools He...

**Run bulk_extractor**

Highlight:

Reports  Feature Filter

Feature File  *None*

Referenced Featu...
Referenced Featu...

Required Parameters

Scan:  ● Image File  ○ Raw Device  ○ Directory of Files

Image file                    :op/SampleData/sampleimage.E01   ...

Output Feature Directory   ampleData/bulk-extractor-output  ...

General Options

☐ Use Banner File

☐ Use Alert List File

☐ Use Stop List File

☐ Use Find Regex Text File

☐ Use Find Regex Text

Tuning Parameters

☐ Use Context Window Size   16

☐ Use Page Size             16777216

☐ Use Margin Size           1048576

☐ Use Min Word Size         6

☐ Use Max Word Size         14

☐ Use Block Size            512

☐ Use Number of Threads     1

Scanner Controls

☐ Use Plugin Directory

☐ Use Scan Option Name

Scanners

☐ bulk
☐ wordlist
☑ accts
☑ aes
☑ base16
☑ base64
☑ elf
☑ email
☑ exif
☑ gps
☑ gzip
☑ hiber
☑ json
☑ kml
☑ net
☑ pdf
☑ vcard
☑ windirs
☑ winpe
☑ winprefetch
☑ zip

Restore Defaults    Start bulk_extractor    Cancel

Left ⌘

# Bulk Extractor* – Identifying Potentially Sensitive Information



**Bulk Extractor Viewer**

File Edit View Boo

**Run bulk_extractor**

Reports

Required Parameters

Scan: ● Image File  ○ Raw Device  ○ Directory of Files

Image file

Output Feature Directory

General Options

☐ Use Banner File
☐ Use Alert List File
☐ Use Stop List File
☐ Use Find Regex Text File
☐ Use Find Regex Text
☐ Use Random Sampling

Tuning Parameters

☐ Use Context Window Size        16
☐ Use Page Size                  16777216
☐ Use Margin Size                4194304
☐ Use Block Size                 512
☐ Use Number of Threads          4
☐ Use Maximum Recursion Depth    7
☐ Use Wait Time                  60

Parallelizing

☐ Use start processing at offset

Scanners

☐ bulk
☐ wordlist
☐ xor
☑ accts
☑ aes
☑ base16
☑ base64
☑ elf
☑ email
☑ exif
☑ find
☑ gps
☑ gzip
☑ hiber
☑ json
☑ kml
☑ net
☑ pdf
☑ rar
☑ vcard
☑ windirs
☑ winpe
☑ winprefetc
☑ zip

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

*Developed by Simson Garfinkel

**Bulk Extractor Viewer**

🖼 ✉ 🔋 ↕ ◀)) 2:02 AM 👤 BitCurator ☼

File Edit View Tools Help

✖ 🗁 🗇 🗐 🗋 🖨

✖ Highlight: [                    ]  ☑ **Match case**

Reports  Feature Filter ☐ Match case          Navigation

**✕ bulk_extractor Scan**

Image File          sampleimage.E01
Feature Directory  bulk-extractor-output

Progress  Done
bulk_extractor scan completed. See Status, below, for details.

Options
'bulk_extractor'
'-o'
'/home/bcadmin/Desk
'/home/bcadmin/Desk

**⊗ Report is Ready**

ⓘ  bulk_extractor has completed.
Report bulk-extractor-output has been opened and is ready for viewing.

✔ OK

Status
Elapsed time: 0.4985 sec.
Overall performance: 2.958 MBytes/sec.
Total email features found: 0
Done.

Close

● Text ○ Hex    ⬅ 🏠 ➡

🖴 ⊙ ✏ 🖳 ▭ Ⓤ 🌐 ⬇ Left ⌘

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

# 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 |

*https://www.forensicswiki.org/wiki/Bulk_extractor#Output_Feature_Files

# 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) |

*https://www.forensicswiki.org/wiki/Bulk_extractor#Output_Feature_Files

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

# BitCurator Reporting Tool

# Provenance – DFXML Output from fiwalk

# Capturing Original Order - Filesystem Metadata Output from fiwalk*

```xml
-<fileobject>
  -<parent_object>
      <inode>102</inode>
   </parent_object>
   <filename>Papers8/37638.BrannyPhyle.Joseph+Moore.pdf</filename>
   <partition>1</partition>
   <id>901</id>
   <name_type>r</name_type>
   <filesize>100857</filesize>
   <alloc>1</alloc>
   <used>1</used>
   <inode>6783</inode>
   <meta_type>1</meta_type>
   <mode>511</mode>
   <nlink>1</nlink>
   <uid>0</uid>
   <gid>0</gid>
   <mtime prec="2">2009-11-17T19:35:10</mtime>
   <atime prec="86400">2009-12-10T05:00:00</atime>
   <crtime prec="2">2009-12-10T19:34:11</crtime>
   <libmagic>PDF document, version 1.4 </libmagic>
  -<byte_runs>
      <byte_run file_offset="0" fs_offset="56621568" img_offset="56653824" len="100857"/>
   </byte_runs>
   <hashdigest type="md5">eb60256dabffa67cef7211bcba659815</hashdigest>
   <hashdigest type="sha1">e56f606877f10daf91dc0304ea120b35452bd36e</hashdigest>
</fileobject>
```

*Developed by Simson Garfinkel

XML Schema for Digital Forensics XML

| ⊙ 43 commits | ⑂ 1 branch | ⬘ 9 releases | ⬚ 1 contributor |
|---|---|---|---|

⟨⟩ Code

① Issues    8

⑃ Pull requests    0

↯ Pulse

⬚ Graphs

⟲  ⑂ branch: master ▾   dfxml_schema / +    ☰

Document an XML validation step  ⋯

ajnelson authored on Dec 4, 2014    latest commit 4c8aab566e ⬚

| 📁 ref | Allow offline validation with local XSD cache | 2 years ago |
| 📄 LICENSE.txt | Add public domain license text | 2 years ago |
| 📄 README.md | Document an XML validation step | 6 months ago |
| 📄 dfxml.xsd | Document an XML validation step | 6 months ago |

**HTTPS** clone URL

https://github.com/c  ⬚

You can clone with HTTPS or Subversion. ⑦

🖥 Clone in Desktop

⬇ Download ZIP

📖 **README.md**

This is the schema repository for Digital Forensics XML, version 1.1.1.

If you intend to use the dfxml.xsd file as a DFXML document validator, note that you will also need to download two accompanying .xsd files under the "ref" directory. The easiest way to do this is by downloading the repository as a Git clone, or by downloading the zip archive from the Github page.

To report issues, questions, or feature requests, please either:

- File a Github issue, seeing first if it is already filed, here.
- Email the dfxml@nist.gov mailing list. If you wish to join the mailing list, send an email to dfxml-subscribe@nist.gov (no subject or message body is necessary), and a moderator will grant access.

https://github.com/dfxml-working-group/dfxml_schema

# Various Specialized BitCurator Reports

# Specialized BitCurator Reports

| File | Content |
|---|---|
| bc_format_bargraph.pdf | histogram of file formats found on the volume |
| bulk_extractor_report.pdf | high-level overview of feature locations on disk |
| fiwalk_deleted_files.pdf | shows paths to any deleted materials found in a given partition |
| fiwalk-output.xml.xlsx | Excel converted DFXML output (file system metadata) |
| fiwalk_report.pdf | high-level overview of file system characteristics |
| format_table.pdf | long-form file format names for formats shown in bar graph |
| premis.xml | PREMIS preservation metadata |

# PREMIS (Preservation) Metadata Generated from Running BitCurator Tools – Recorded as PREMIS Events

# BitCurator PDF Redaction Tool

bcadmin@ubuntu: ~/bitcurator-redact-pdf/build/libs

**BitCurator PDF Redact**

PDF Files   Entity Recognition   Text Patterns   Help

| Filesme | Path | Output |
|---------|------|--------|
| Abstract.... | /home/bc... | /home/bc... |

**Named Entities**   **Text Patterns**

Named entities are people, places, and organizations detected in the text of PDF files you have added.

| Entity Text | Type | # | Files | Action |
|-------------|------|---|-------|--------|
| Archive Analytics | ORGANIZATION | 1 | 1 | Ignore |
| Cassandra | PERSON | 1 | 1 | Ignore |
| Digital Curation Innovation... | ORGANIZATION | 1 | 1 | Ignore |
| Maryland | LOCATION | 1 | 1 | Ignore |
| NCSA | ORGANIZATION | 1 | 1 | Ignore |
| NLG for Libraries FY17 Nati... | ORGANIZATION | 1 | 1 | Ignore |
| University of Maryland | ORGANIZATION | 1 | 1 | Ignore |
| University of Maryland's Co... | ORGANIZATION | 1 | 1 | Ignore |

output folder: none

Trash

Documenta tion and Help

Network Servers

bcadmin@ubuntu: ~/bitcurator-redact-pdf/build/libs

**BitCurator PDF Redact**

**PDF Files**  **Entity Recognition**  **Text Patterns**  **Help**

| Files | me | Path | Output | New Pattern | terns |
|-------|-----|------|--------|-------------|-------|
| Abstract.... | /home/bc... | /home/bc | Open File(s).. | | |

**New Pattern**
**Open File(s)..**
**Save As..**
**Reset to Defaults**
**Save as Defaults**
**Clear All**
**Import Bulk Extractor features..**

ces, and organizations detected in the ed.

| | Type | # | Files | Action |
|---|------|---|-------|--------|
| | ORGANIZATION | 1 | 1 | Ignore | ▲ |
| | PERSON | 1 | 1 | Ignore |
| | ORGANIZATION | 1 | 1 | Ignore |
| Maryland | LOCATION | 1 | 1 | Ignore |
| NCSA | ORGANIZATION | 1 | 1 | Ignore |
| NLG for Libraries FY17 Nati... | ORGANIZATION | 1 | 1 | Ignore |
| University of Maryland | ORGANIZATION | 1 | 1 | Ignore |
| University of Maryland's Co... | ORGANIZATION | 1 | 1 | Ignore |

Documenta tion and Help

Network Servers

output folder: none

Trash

Activities  bca-redact-RedactionApp ▾                    Sun 15:42 ●

bcadmin@ubuntu: ~/bitcurator-redact-pdf/build/libs

**BitCurator PDF Redact**

PDF Files  Entity Recognition  Text Patterns  Help

| Filename | Path | Output |
|---|---|---|
| Abstract.... | /home/bc... | /home/bc... |

**Named Entities**  **Text Patterns**

Patterns are regular expressions used to redact matching text in PDFs.
Add new patterns by clicking in the empty first row.

| Name | Expression | Action |
|---|---|---|
| Social Security Num... | \d{3}-\d{2}-\d{4} | Redact |
| gross.joshua.b+job... | \Qgross.joshua.b+jo... | Ask |
| Glenn.Gunzelmann... | \QGlenn.Gunzelmann... | Ask |
| gross.joshua.b@gm... | \Qgross.joshua.b@g... | Ask |
| mathbio@math.pitt.... | \Qmathbio@math.pi... | Ask |
| cnbc-all@cnbc.cmu.... | \Qcnbc-all@cnbc.cm... | Ask |
| bard@math.pitt.edu | \Qbard@math.pitt.e... | Ask |
| mathbio@math.pitt... | \Qmathbio@math.pi... | Ask |
| cnbc-all@cnbc.cmu... | \Qcnbc-all@cnbc.cm... | Ask |
| leonardochiesi@gma... | \Qleonardochiesi@g... | Ask |
| gross.joshua.b@gm... | \Qgross.joshua.b@g... | Ask |
| gross.joshua.b@gm... | \Qgross.joshua.b@g... | Ask |
| gross.joshua.b@gm... | \Qgross.joshua.b@g... | Ask |
| buy.com_offers@en... | \Qbuy.com_offers@e... | Ask |
| gross.joshua.b@gm... | \Qgross.joshua.b@g... | Ask |
| 3C4A527E0E.40006... | \Q3C4A527E0E.400... | Ask |
| 3C4A527E0E.40006... | \Q3C4A527E0E.400... | Ask |
| leonardochiesi@gma... | \Qleonardochiesi@g... | Ask |
| 3C2acb011c090706... | \Q3C2acb011c0907... | Ask |
| daughtry@psu.edu | \Qdaughtry@psu.ed... | Ask |
| amsuich@nps.edu | \Qamsuich@nps.edu\E | Ask |
| 3C8AB9A1F305571... | \Q3C8AB9A1F30557... | Ask |
| amsuich@nps.edu | \Qamsuich@nps.edu\E | Ask |
| 3C8AB9A1F305571... | \Q3C8AB9A1F30557... | Ask |
| amsuich@nps.edu | \Qamsuich@nps.edu\E | Ask |
| 3C8AB9A1F305571... | \Q3C8AB9A1F30557... | Ask |
| hous-daccq-136905... | \Qhous-daccq-1369... | Ask |
| cherylseekingforoom... | \Qcherylseekingforo... | Ask |
| hous-daccq-136905... | \Qhous-daccq-1369... | Ask |
| hous-daccq-136905... | \Qhous-daccq-1369... | Ask |
| bw3maggers@gmail... | \Qbw3maggers@gm... | Ask |
| hous-daccq-136905... | \Qhous-daccq-1369... | Ask |
| gross.joshua.b+job... | \Qgross.joshua.b+jo... | Ask |

output folder: none

Trash

Documenta
tion and
Help

Network
Servers

Activities    bca-redact-RedactionApp ▾     Sun 15:48 ●

**Redact Document**

| Page | Text | Type | Action |
|---|---|---|---|
| 1 | DRASTIC | REGEX | Ask ▼ |
| 1 | DRASTIC | REGEX | Ignore |
| 1 | DRASTIC | REGEX | **Ask** |
| 1 | DRASTIC | REGEX | Redact |
| 1 | DRASTIC | REGEX | Ask |
| 1 | DRASTIC | REGEX | Ask |
| 1 | DRASTIC | REGEX | Ask |

NLG for Libraries FY17 National Digital Platform Research Grant full proposal narrative -- University of Maryl

**Abstract**

**Improving Fedora to Work with Web-scale Storage and Serv**

Memory institutions around the world face a rapidly expanding need for storage and acce and metadata. The Fedora Repository has long been at the forefront of their efforts, develop the challenge, including four major versions of the Fedora Repository software. Now th have put forward a bold call to the community to create new implementations of Fedor needs, publishing a formal API that specifies the expectations of a Fedora repository. Throu computational archives and through prior Fedora involvements, we have learned that scalability, by which we mean the ability to expand storage capacity without losing perfo that institutions must be able to incrementally grow a fully-functional repository as collec the need for expensive enterprise storage plans, massive data migrations, and performance the vertical storage strategy of previous repository implementations.

The Digital Curation Innovation Center (DCIC) at the University of Maryland's College of (Maryland's iSchool) intends to conduct a 2-year project to research, develop, and test soft improve the performance and scalability of the Fedora Repository for the Fedora communi this project will apply the new Fedora 5 application programming interface (API) to the stack called DRAS-TIC to create a new Fedora implementation we are calling *DRAS-TIC* which stands for Digital Repository at Scale that Invites Computation, was developed ov through a collaboration between UK-based storage company, Archive Analytics, and funding from an NSF DIBBs (Data Infrastructure Building Blocks) grant (NCSA "Brown leverages NoSQL industry standard distributed database technology, in the form of A provide near limitless scaling of storage without performance degradation. With Cassandra can also hold redundant copies of data in datacenters around the world. Even if an enti access can remain uninterrupted, and data re-replicated to a new datacenter. Beyond instit think this creates the possibility for new reciprocal storage arrangements between Fedora in

To meet with this potential, DRAS-TIC will first need to be adapted to the new Fedora API and tested to meet the performance expectations of our Fedora community partners. We ha of institutional partners in the Fedora community that will work with us to develop use ca expectations. As we develop and test *DRAS-TIC Fedora*, their institutional needs will g become our measure of success. The proposal has received the endorsement of the Fedor http://fedorarepository.org/leadership-group.

The proposed project will produce open-source software, tested cluster configurations, docu practice guides that will enable institutions to manage Fedora repositories with Petabyte-s

< Prev   Next >     Redact     Close

# Other Functionality to Meet Identified User Needs:

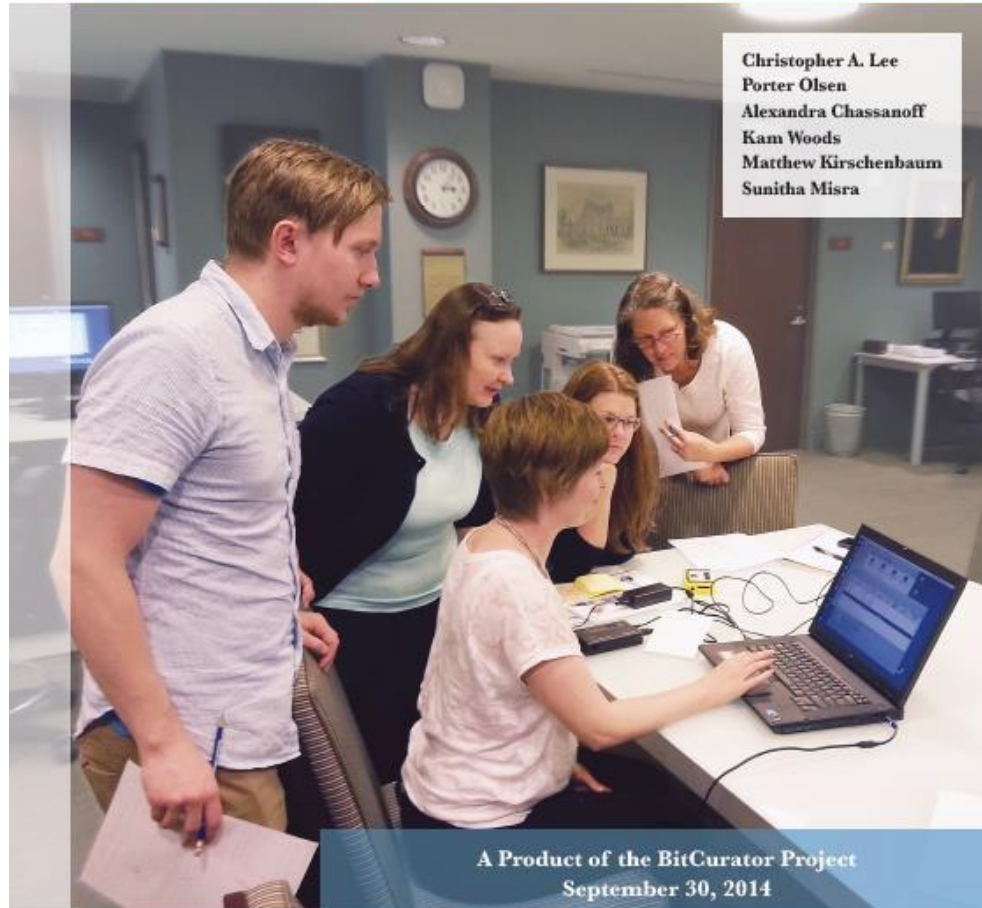| Function | Tool(s) |
| --- | --- |
| Identify duplicate files | FSLint |
| Characterize files | FIDO, Siegfried, Brunnhilde |
| Scan for viruses | ClamTK |
| Examine, copy and extract information from old Mac disks | HFS Utilities (including HFS Explorer) |
| Capture AV file metadata | MediaInfo, FFProbe |
| Extract text from older binary (.doc) Word files | antiword |
| Read contents of Mircosoft Outlook PST files | readpst |
| Examine embedded header information in images | pyExifToolGUI |
| Generate images of problematic disks or particular disk types (I addition to Guymager | dd, dcfldd, ddrescue, cdrdao (for audio CDs) |
| Extract and analyze data from Windows Registry files | regripper |
| Identify files that are partially similar but not identical | sdhash, ssdeep |
| Package files for storage and/or transfer | BagIt (Java) library, Bagger |
| File preview (left-click on file then hit space bar) | gnome-sushi |

# Other Functionality to Meet Identified User Needs (Continued):

| Function | Tool(s) |
|---|---|
| Play and examine metadata from AV media files | VLC media player |
| Damaged/lost partition recovery | TestDisk |
| Damaged/lost file recovery | PhotoRec |
| Identify the filesystem on a disk | disktype |
| Index and search for keywords in documents | recoll |
| Find blacklist data by using hashes calculated from hash blocks | hashdb |
| Generate hashes of files and blocks | GTK Hash, md5deep, md5sum |
| Compare hashes of files to hashes in the National Software Reference Library (NSRL) of known system files | nsrllookup |
| View and edit bytestreams (hex editor) | Bless Hex Editor, GHex |

# From Code to Community:
## Building and Sustaining BitCurator through Community Engagement

Christopher A. Lee
Porter Olsen
Alexandra Chassanoff
Kam Woods
Matthew Kirschenbaum
Sunitha Misra

A Product of the BitCurator Project
September 30, 2014

http://www.bitcurator.net/wp-content/uploads/2014/11/code-to-community.pdf

# BitCuratorEdu

Advancing the adoption of digital forensics tools and methods in libraries and archives through professional education efforts

EDUCOPIA INSTITUTE
Community Cultivators

INSTITUTE of Museum and Library SERVICES

Most resources from the BitCuratorEdu project are intentionally left with basic formatting and without project branding. We encourage educators, practitioners, and students to adapt these materials as much as needed and share them widely.

*The BitCuratorEdu project is a three-year effort funded by the Institute of Museum and Library Services (IMLS) to study and advance the adoption of digital forensics tools and methods in libraries and archives through professional education efforts. This project is a partnership between Educopia Institute and the School of Information and Library Science at the University of North Carolina at Chapel Hill, along with the Council of State Archivists (CoSA) and several Masters-level programs in library and information science.*