# Identification and Redaction of Sensitive Information in Born-Digital Archival Materials: Research and Development Directions

**Christopher (Cal) Lee**
UNC School of Information and Library Science

Society of American Archivists Research Forum
Atlanta, GA
August 2, 2016

## BitCurator Access

**The Andrew W. Mellon Foundation**

UNC
SCHOOL OF INFORMATION
AND LIBRARY SCIENCE

# Sensitive Stuff

- Archivists must often identify information within their collections (either in their existing holdings or as part of new acquisitions) that is sensitive and should thus be removed, closed, restricted or filtered from public view.

- Manually identifying such information can be time-consuming and prone to error.

- Luckily, there are many automated methods that archivists can use to support this work*

*Lee, Christopher A., and Kam Woods. "Automated Redaction of Private and Personal Data in Collections: Toward Responsible Stewardship of Digital Heritage." In *Proceedings of Memory of the World in the Digital Age: Digitization and Preservation: An International Conference on Permanent Access to Digital Documentary Heritage, 26-28 September 2012, Vancouver, British Columbia, Canada*, edited by Luciana Duranti and Elizabeth Shaffer, 298-313: United Nations Educational, Scientific and Cultural Organization, 2013.
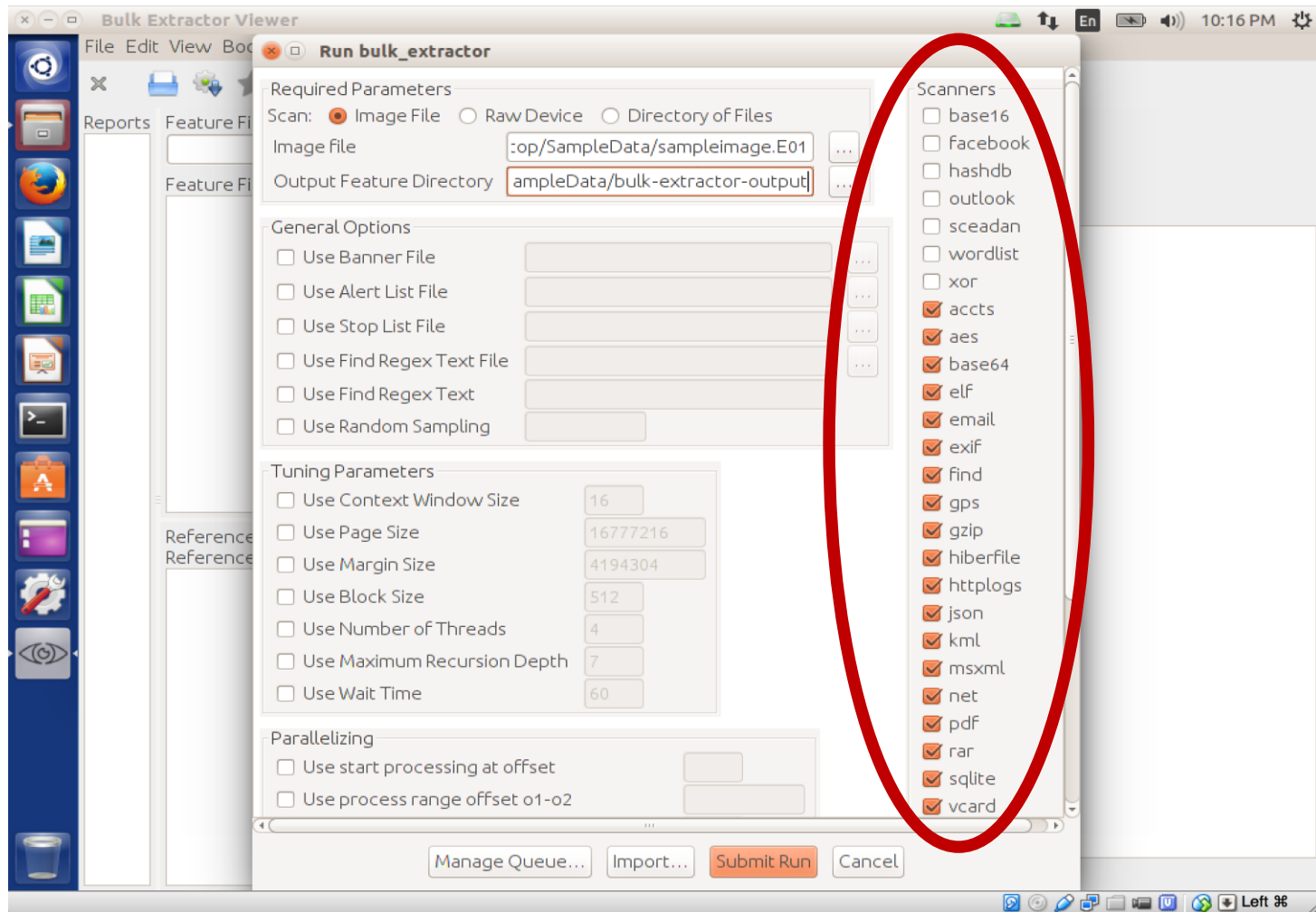
# Computer-Assisted Redaction

- Still a relatively undeveloped area of archival practice
- Good news for us: numerous opportunities for further research and development

# Examples of Potentially Private and Sensitive Information

- Personal identifiers (e.g. SSNs, DOBs, Drivers License #s, corp. and govt. IDs)
- Financial information (e.g. credit card numbers)
- Geolocation data
- Email messages, email addresses, attachments
- Traces of online activity (e.g. search histories, web caches, domain names, IP addresses)
- Recoverable data from deleted files
- Partially overwritten data

# We know how to look for these things…



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

*bulk_extractor, developed by Simson Garfinkel.  Available to run directly or within the BitCurator environment
http://wiki.bitcurator.net
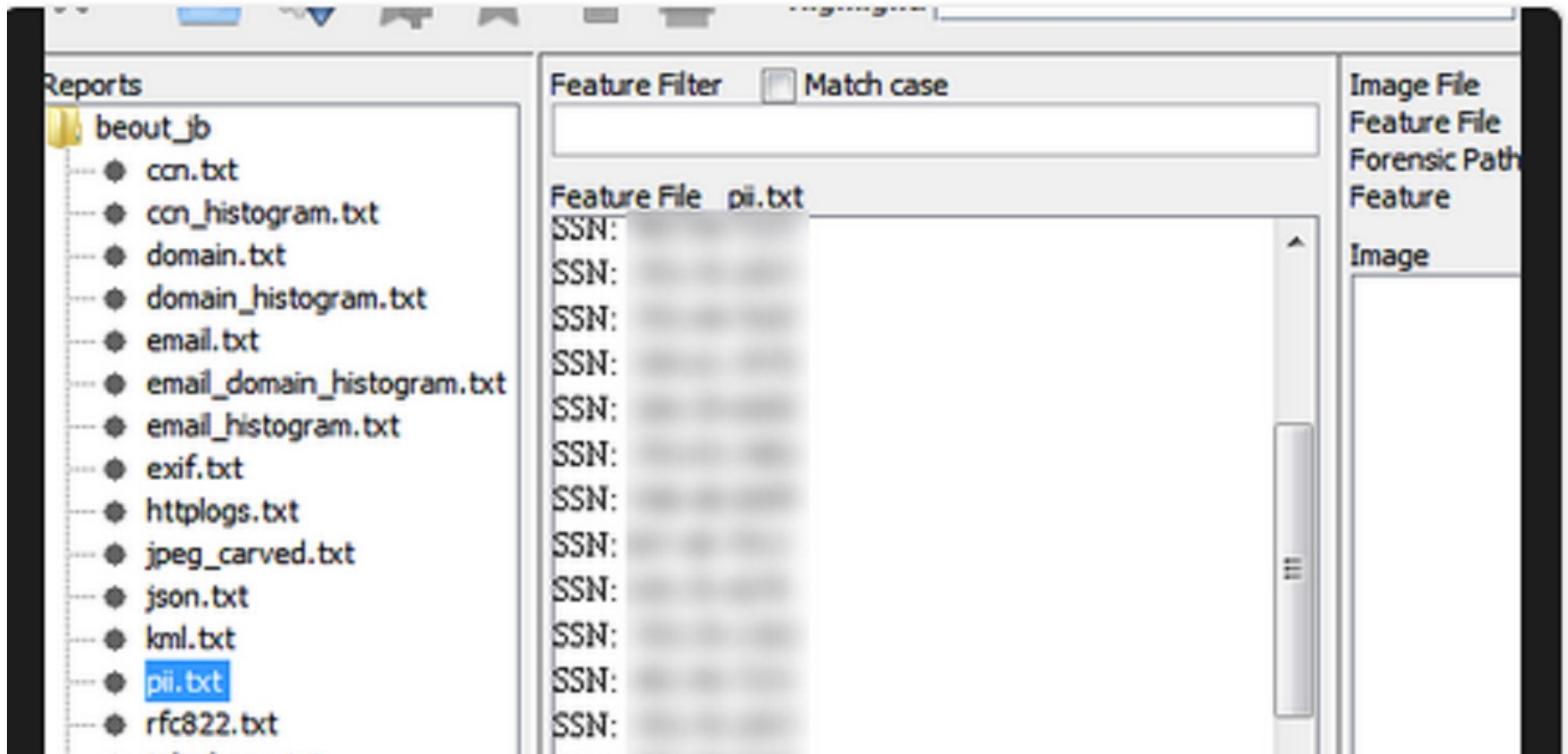
That news story about information leaked through email…

# 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  ✉ Email  🐦 @chillmage

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

# A Real World Example of Forensic Feature Extraction:

# 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

# Encryption may be a marker for sensitivity



PDFs from 2003 w/ 40-bit encryption

# Other PII and sensitive data may be harder to find, requiring e.g.:

- OCR
- named entity recognition
- partial file reconstruction
- format-specific tools
- visual inspection

## 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
```

# Geolocation data embedded in EXIF metadata from a smartphone photo

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

# BitCurator Access

- Two-year project (October 1, 2014 – September 30, 2016) at School of Information and Library Science, University of North Carolina at Chapel Hill

- Funded by Andrew W. Mellon Foundation

- Developing open-source software to support access to disk images. Three core areas of focus:

  – Tools and reusable libraries to support web access services for disk images

  – Analyzing contents of file systems and associated metadata

  – **Redacting complex born-digital objects (disk images) and emulated access to redacted images**
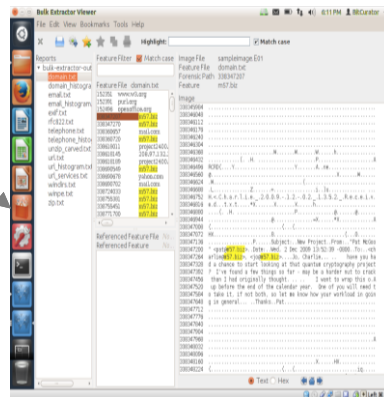
# Redacting Born-Digital Materials

- Locating relevant items can be problematic
  - compressed files
  - proprietary formats
  - formatting variations
  - Encryption
  - …
- Digital forensics tools can help
  - Scan file systems and disk images block-by-block
  - Extract features when possible
  - Report on materials that resist analysis

# Redaction: Automating search for redaction candidates and redaction actions is desirable, and few tools exist to support this search across heterogeneous materials

Acquire disk image from original media

Generate redacted disk image and/or files

Identify items to redact

Report on redacted items for preservation record

# Forensic disk imaging and metadata extraction provides clear provenance for redacted copies

File object identified in disk image and recorded in a forensic metadata format (DFXML)

Original file (unredacted in disk image)

PII identified at byte offsets

Redacted copy (alternate location)

```
-<fileobject>
 -<parent_object>
    <inode>98</inode>
  </parent_object>
    <filename>Papers6/58142.PKPpairs.Kenneth+Thompson.pdf</filename>
    <partition>1</partition>
    <id>734</id>
    <name_type>r</name_type>
    <filesize>100521</filesize>
    <alloc>1</alloc>
    <inode>6253</inode>
    <meta_type>1</meta_type>
    <mode>511</mode>
    <nlink>1</nlink>
    <uid>0</uid>
    <gid>0</gid>
    <mtime prec="2">2009-11-17T19:46:08</mtime>
    <atime prec="86400">2009-12-10T05:00:00</atime>
    <crtime prec="2">2009-12-10T19:33:28</crtime>
    <libmagic>PDF document, version 1.4 </libmagic>
  -<byte_runs>
      <byte_run file_offset="0" fs_offset="43913728" img_offset="43945984" len="100521"/>
    </byte_runs>
    <hashdigest type="md5">f5495bd2b5520984bb4c54a42485f9f0</hashdigest>
    <hashdigest type="sha1">2a005bffd8145c374427a929585ffa0e49ad79e3</hashdigest>
</fileobject>
```

# Automated Redaction and Access Options



EaaS = Emulation-as-a-Service. http://bw-fla.uni-freiburg.de/

# Automated Redaction and Access Options



EaaS = Emulation-as-a-Service. http://bw-fla.uni-freiburg.de/

# Emulation as a Service

See also: Woods, Kam, Christopher Lee, Oleg Stobbe, Thomas Liebetraut and Klaus Rechert. "Functional Access to Forensic Disk Images in a Web Service." In *Proceedings of the 12th International Conference on Digital Curation*, edited by Christopher A. Lee, Jonathan Crabtree, Leo Konstantelos, Nancy McGovern, Yukio Maeda, Maureen Pennock, Helen Tibbo, Kam Woods, and Eld Zierau, 191-195. Chapel Hill, NC: University of North Carolina, School of Information and Library Science, 2015.

# Automated Redaction and Access Options



**Option A: Redact from live image in EaaS via copy-on-write overlay**

**Option B: EaaS access to previously redacted image**

**Option C: Browse non-live file system with redaction mask**

EaaS = Emulation-as-a-Service. http://bw-fla.uni-freiburg.de/

BitCurator / bca-redtools

Watch  1    Star  0    Fork  1

<> Code     ⓘ Issues  0     Pull requests  0     Pulse     Graphs

Redaction Tools for Disk Images

⊙ 25 commits          ⑃ 2 branches          ⬡ 1 release          ⚇ 2 contributors

Branch: master ▾    New pull request                    Find file    Clone or download ▾

kamwoods Updated docs for release                     Latest commit 2a6fc79 2 days ago

| libredact | added support for FILE_SEQ_MATCH rules, redacting entire file based o... | 9 days ago |
| .gitignore | Initial commit | 9 months ago |
| ChangeLog.txt | Updated docs for release | 2 days ago |
| LICENSE | Initial commit | 9 months ago |
| README.md | Updated docs for release | 2 days ago |

📖 README.md

# bca-redtools: BitCurator Access Redaction Tools

Disk image and bitstream redaction tools for the BitCurator Access project.

https://github.com/BitCurator/bca-redtools
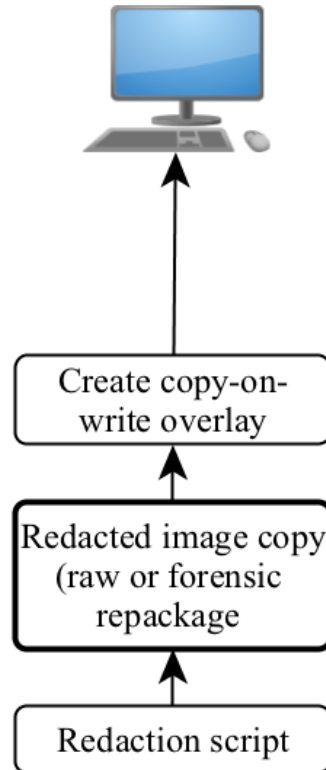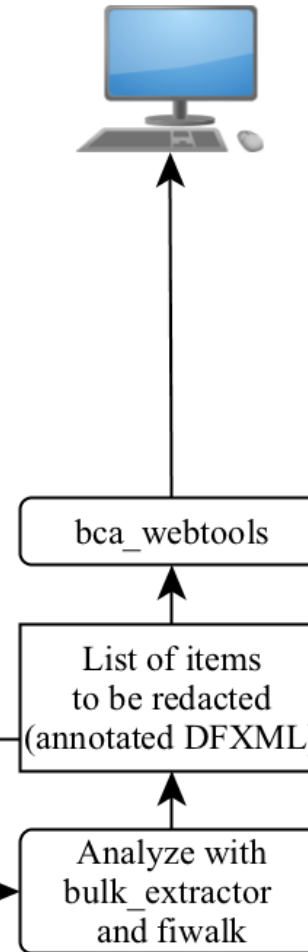
# Automated Redaction and Access Options
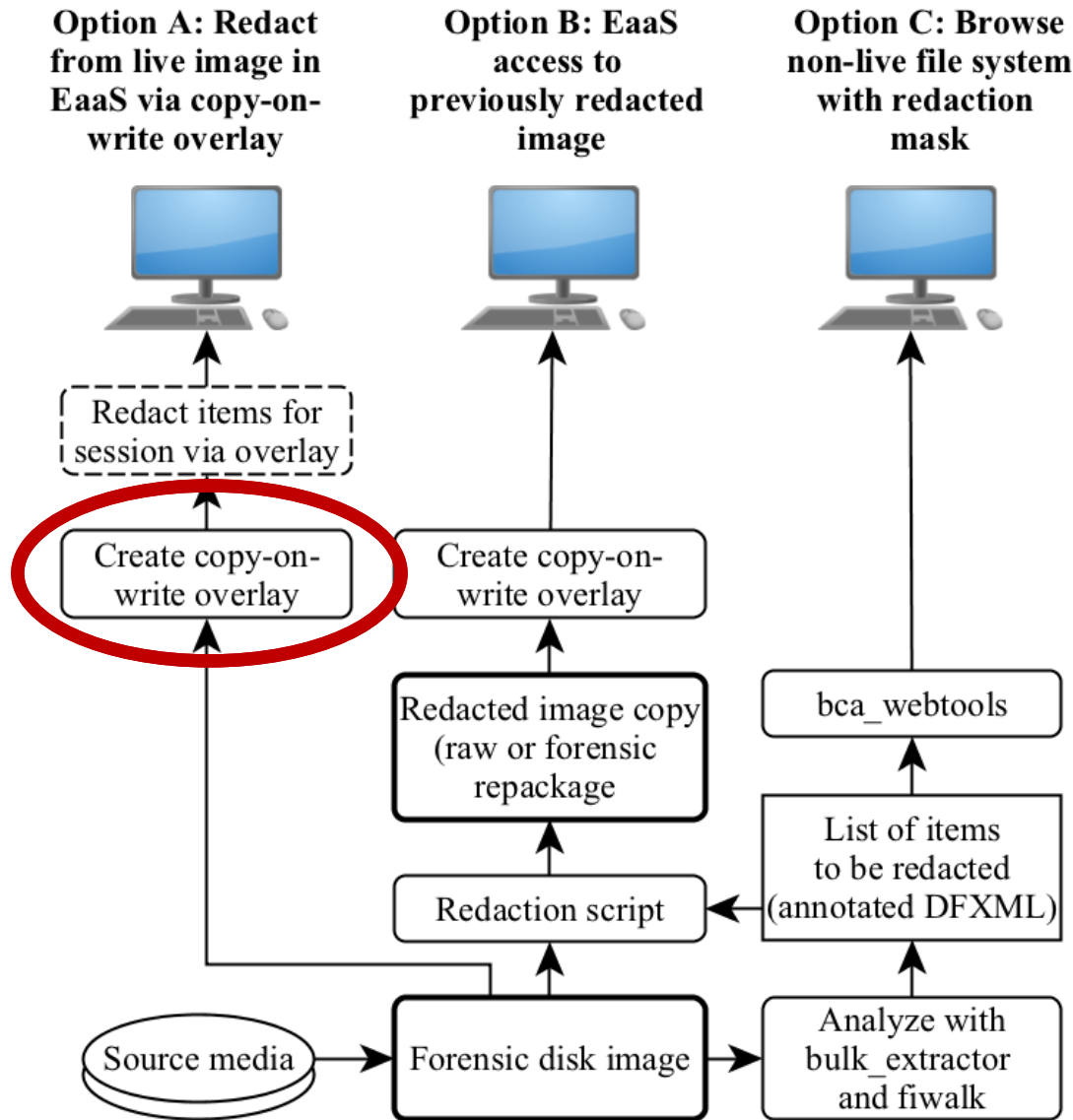


**Option A: Redact from live image in EaaS via copy-on-write overlay**
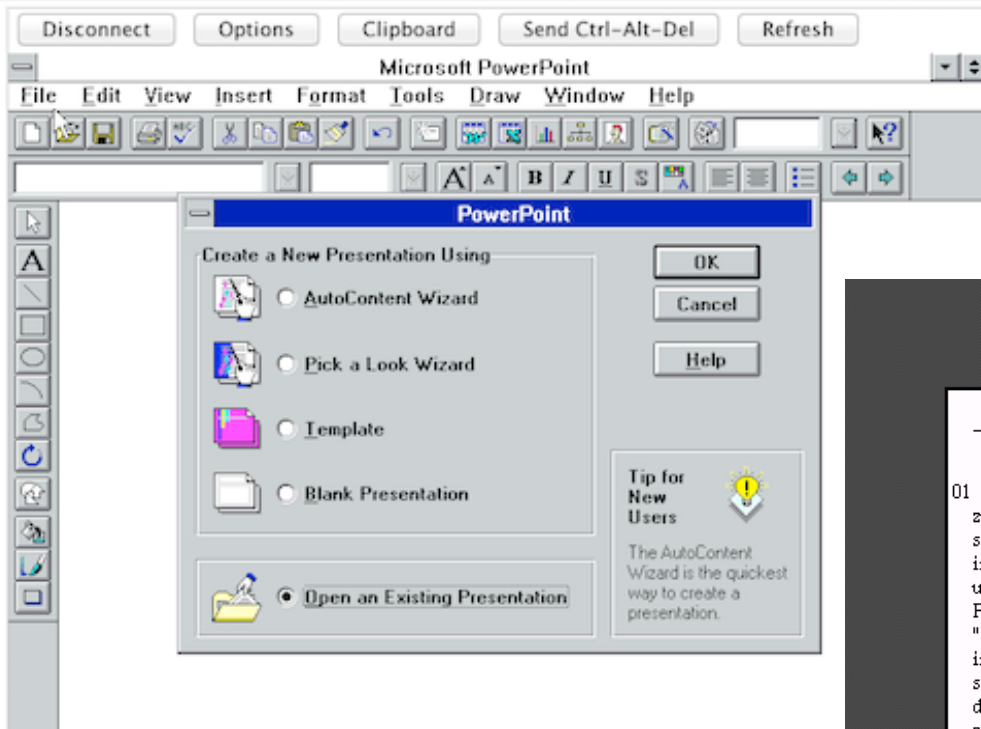
**Option B: EaaS access to previously redacted image**

**Option C: Browse non-live file system with redaction mask**

EaaS = Emulation-as-a-Service. http://bw-fla.uni-freiburg.de/

# BCA (BitCurator Access) Web Tools

- Integrates digital forensics software libraries and lightweight web-services tools
- Drop disk images in a local or network-accessible location, start up the service, and start browsing
- Most analysis runs server-side (via Sleuthkit and DFXML Python bindings, among others)
- Service is database-agnostic (we're using postgres)
- Automatic metadata production – Digital Forensics XML (DFXML), PREMIS, others)

**https://github.com/kamwoods/bca-webtools**

Sunitha Misra, Christopher A. Lee, and Kam Woods, "A Web Service for File-Level Access to Disk Images," *Code4Lib Journal 25* (2014), http://journal.code4lib.org/articles/9773

## Primary Server

### Application Stack

| The Sleuth Kit (TSK) | fiwalk |
|---|---|

pytsk3 → bca-webtools

DFXML tools → bca-webtools

**bca-webtools**
Flask (Python microframework) application

Jinja2 templating engine

Dynamically generated HTML file system views

PyLucene

Apache Lucene

SQLAlchemy

### Web Service Stack

uWSGI server

NGINX

### Search Index

Search terms extracted from text-rich document types (Word, ODT, PDT, TXT)

## Disk Image Store

Raw and forensically packaged disk images

## Database Server

PostgreSQL

### Application DB

Disk image metadata

Admin and access control data

Internet or Intranet

## Client

Web browser

Local / network storage for downloaded items

127.0.0.1:8080

# BCA
# Webtools

## Home

The bca-webtools application provides access to forensically-packaged and raw disk images. Supported file systems include FAT16, FAT32, NTFS, HFS+, and EXT2/3/4. Click on 'Browse' to navigate through the file system(s) within the disk image, or 'Download' to download the complete disk image.

| Image Name | Info | Browse | Download |
|---|---|---|---|
| charlie-work-usb-2009-12-11.E01 | 🔍 | 🗄 | 💿 |
| terry-work-usb-2009-12-11.E01 | 🔍 | 🗄 | 💿 |

Select an option below to search available disk images by filename or file contents. (Currently indexing all filenames, contents of .doc, .odt, .pdf, and .txt)

○ Search by filename
● Search by content

[                    ] Search

- Admin

# BCA Webtools

Home

## bca-webtools - Admin Tools

- ○ Build Image Table
- ○ Build DFXML Table
- ○ Build All Tables
- ○ Drop Image Table
- ○ Drop DFXML Table
- ○ Drop All Tables
- ○ Generate Index
- ○ Clear Index
- ● Show Image Matrix

Submit

### Image Matrix

| Index | Image name | Image DB? | DFXML DB? | Indexed? | Add Table | Delete Table |
|-------|-----------|-----------|-----------|----------|-----------|--------------|
| 0 | charlie-work-usb-2009-12-11.E01 | True | False | True | ☐ Add | ☐ Delete |
| 1 | terry-work-usb-2009-12-11.E01 | True | False | True | ☐ Add | ☐ Delete |

- Admin

127.0.0.1:8080/image/charlie-work-usb-2009-12-11.E01/1

**BCA Webtools**

Home

Browse directories and download files. Items marked "r" in the first column are regular files. Items marked "d" are directories.

| d/r | Filename | Size | Last Modified | Deleted? |
|-----|----------|------|---------------|----------|
| r | $AttrDef | 2560 | 2009-11-20T17:38:09Z | No |
| r | $BadClus | 0 | 2009-11-20T17:38:09Z | No |
| r | $Bitmap | 32320 | 2009-11-20T17:38:09Z | No |
| r | $Boot | 8192 | 2009-11-20T17:38:09Z | No |
| d | $Extend | 552 | 2009-11-20T17:38:09Z | No |
| r | $LogFile | 7405568 | 2009-11-20T17:38:09Z | No |
| r | $MFT | 262144 | 2009-11-20T17:38:09Z | No |
| r | $MFTMirr | 4096 | 2009-11-20T17:38:09Z | No |
| r | $Secure | 0 | 2009-11-20T17:38:09Z | No |
| r | $UpCase | 131072 | 2009-11-20T17:38:09Z | No |
| r | $Volume | 0 | 2009-11-20T17:38:09Z | No |
| d | . | 56 | 2009-12-03T21:17:01Z | No |
| r | 01.zip | 108438 | 2009-11-24T21:21:16Z | No |
| r | astronaut.jpg | 713418 | 2009-11-24T21:33:33Z | No |
| r | astronaut1.jpg | 722717 | 2009-11-24T21:43:42Z | No |
| d | Email | 56 | 2009-12-10T22:27:55Z | No |
| d | Immortality | 56 | 2009-11-24T21:55:45Z | No |
| r | invsecr2.exe | 1291720 | 2009-11-19T18:42:25Z | No |
| r | microscope.jpg | 136274 | 2009-11-24T21:27:51Z | No |

# Many Professional Decisions – One Example:

- Remote access over the Internet
- Direct access through the reading room
- More highly mediated access through selected surrogates

Research questions also abound.  A couple examples:

# Block-Level vs. File-Level Redaction

| Block-Level | File-Level |
|---|---|
| Finding the bits based on where they appear on the disk, without regard for what files or folders they're in | Identifying bits to remove based on where they appear in specific files or folders |
| Usually much faster and simpler to implement than file-level redaction and can find data in unallocated space (deleted content) | Requires more time and more knowledge of underlying data structures (filesystems, file formats), and often requires format-specific tools |
| Potential to hinder the mounting and navigation of file systems and ability to use/render files | Can require much more human effort, including identification, installation and running of format-specific redaction tools |

# More Product, Less Process (MPLP) Mean with Born-Digital Materials

- Options I've discussed represent significantly different levels of human and computational resources

- Different materials and use cases will warrant different strategies

- How do we decide which to adopt (i.e. what are the decision rules or heuristics)?

# BitCurator CONSORTIUM

About Us ▾    Why Digital Forensics ▾    Using BitCurator ▾    Get Involved ▾

## 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.

### 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⌐

### How our members are using BitCurator

### 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

https://bitcuratorconsortium.org/