# Collecting Email Archives at the British Library

Jonathan Pledge

# Born Digital Archives at the British Library

- Born Digital research started at the BL in 2000.

- Pioneers in the born digital field, in particular forensic capture.

- Since April 2015 we have focused on curator-led processing of born digital archives.

- Work with BL Digital Preservation

# Contemporary British Born Digital Archives

- Usually hybrid archives; personal archives, primarily of writers and scientists.

- Smallest born digital collection is 71MB and the largest is 380GB (including system files).

- Many archives contain email particularly where we have acquired a hard drive with operating system.

- Largest email archive is of the poet Wendy Cope consisting of 40,000 emails

# Workflow: Forensic capture

- Hard drives/Flash drives acquired as .E01 files

- CDs/DVDs acquired as iso or .E01 files exported as raw files (.001 or .img)

- Floppy disks (3.5 inch) acquired via Kryoflux as .img or raw files.

# Workflow: Processing

- Extract folders and digital objects from captured forensic image.

- Create master copy of extracted material (with a minimum of two back up copies).

- Export master capture to create 'Extracted Captures' (replicating original structures and titles)

# Workflow: Processing

| Name | Date modified | Type | Size |
|------|---------------|------|------|
| Harwood-88881-Captures | 29/03/2017 14:50 | File folder | |
| Harwood-88881-Captures Extracted | 29/03/2017 14:51 | File folder | |
| Harwood-88881-Metadata | 03/07/2017 11:43 | File folder | |
| Harwood-88881-PDFA | 29/03/2017 14:51 | File folder | |
| Harwood-88881-PDFA Closed | 29/03/2017 14:51 | File folder | |

# Workflow: Processing

- DROID is used to extract metadata from the 'Extracted Captures'. The .csv metadata file is then saved to Excel, sorted and arranged. This information is then copied to a catalogue migration spreadsheet and published via an uploader to our Explore Archives and Manuscripts catalogue.

- PDF/As (PDF/A-b) are created from the 'Extracted Captures' and are renamed in sequential order with archive reference numbers (e.g. Add MS)

# Workflow: Processing

# Workflow: Processing

# Workflow: Delivery

- PDF/As are uploaded to a BL FTP Server in labelled directories. The server is hosted in the BL Manuscripts Reading Room.

- Born Digital Material is accessed, as PDF/As, through dedicated terminals in the Manuscripts Reading Room via a web browser.

- Digital objects can only be accessed in the Manuscripts Reading Room and cannot be edited, printed or saved.

# Email processing

- Available resources: Aid4Mail Forensic (Fookes Holding, Switzerland) and ePADD (Stanford University)

- Aid4Mail Forensic allows for processing and delivery as per our current workflow

- ePADD would seem to offer a holistic approach to email preservation – to be investigated

# Email processing

# Data Protection Act 1998

- UK Act of Parliament.

- In response to 1995 EU Data Protection Directive.

- Defines the law on the processing of data as 'Personal information' on identifiable living people.

- 'Personal information' means any information about an identifiable living individual in any format. It applies to paper, emails, photographs and sound recordings.

- Exemptions for journalism, literature and art.

# Challenges

- Email is a unique entity which often contains 'personal information'. The nature of the medium.

- Processing becomes a time consuming process due to DP checking of hundreds/thousands of emails.

- Complications arise where emails are saved as threads meaning that in the current BL DP model a single email can mean that the entire thread will be removed.

- Solution may be redaction of offending material at the risk of a loss of meaning.

# Thank you

jonathan.pledge@bl.uk
eleanor.dickens@bl.uk