# Technology Watch

# Digital Archiving Innovation and Research (DAIR)

- Leading research and innovation in digital archiving policy, methodologies, techniques, and practice
- Promoting and supporting staff digital archiving capability to support the National Archives vision, priorities and strategies
- Contributing subject matter expertise to NAA projects and operations
- Technology watch is a component of our brief

# NAA

**NATIONAL ARCHIVES OF AUSTRALIA**

# Digital·Archives·Strategic· Research·Priorities· Framework¶

### Digital·Archives·Innovation·&·Research¶

February·2021¶

Australian Government

National Archives of Australia

# Types of Research

## Horizon scanning

- Archival strategies & emerging practice, including indigenous collections
- 'Stretch' topics less based in current operations
- May not have a clear implementation pathway
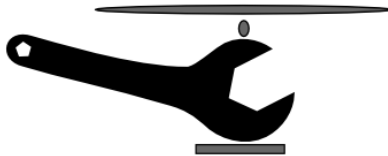- Often with national / international partners

## Commissioned archival practices and techniques

- More attuned to operational needs, pain points and challenges
- Must have an identified business owner within NAA
- Focus on pragmatic outcomes for business improvement

DCC

DigitalPreservationCoalition

E-ARK

ICA
International Council on Archives
Conseil International des Archives

Software Preservation Network

Open Preservation Foundation

pasig
preserving and archiving
special interest group

CAARA
Council of Australasian Archives
and Records Authorities

WE MISS iPRES

Jisc

NDSA

International Association of Sound and Audiovisual Archives
Internationale Vereinigung der Schall- und audiovisuellen Archive
Association Internationale d'Archives Sonores et Audiovisuelles
Asociación Internacional de Archivos Sonoros y Audiovisuales

iasa

FA DGI Federal Agencies Digital Guidelines Initiative

dcc

# Horizon Scanning and Technology Watch Register

| Category | Sub Category | Topic/Technology | Description | Sourced from | Date added | Added by | Referred to DAIR by | Action Required | Staff member | Completed | Date last updated | Updated by | Outcome | RkS Reference | Comments | external l |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Technology | digital preservation | Emulation as a Service Infrastructure | A project to develop a scalable emulation to support a shared emulation capability | Yale University, Software Preservation Network | 6-Jan-21 | JD | | Test | James Doig, Carey Garvie, Tim Mifsud, Bridget Dexter | N | | | | 2019/3387 | Tim Mifsud attended workshop on EaaSI at IDCC in Dublin 2020 | https://w |
| Technology | digital preservation | Virtual Research Environment (VRE) - DDHN & OPF | virtual platform containing key digital preservation tools | DP JISC listserv | 5-Jan-21 | CG | Yaso Arumugam | Test | James Doig, Tim Mifsud, Carey Garvie | Y | 05-Jan-21 | CG | Response provided to ADG | 2020/4159 | | Virtual Re |
| Technology | email | Review, Appraisal and Triage of Mail (RATOM) - UNC, Chapel Hill | software to assist with email analysis, selection and appraisal tasks | | 5-Jan-21 | CG | | Test | James Doig | N | | | | | | Review, A |
| Technology | digital forensics | BitCurator environment - UNC, Chapel Hill | Suite of digital forensics tools | Cal Lee | 5-Jan-21 | CG | | Test | James Doig, Carey Garvie | N | | | | | Carey Garvie attended workshop on BitCurator post IDCC 2019. Initial testing on café laptop with 3.5" floppy disks. | https://bi |
| Technology | risk management | Digital Preservation Framework - NARA | NARA risk management approach to digital preservation of file formats. Includes risk matrix and preservation action plans for over 500 file formats | NARA | 5-Jan-21 | CG | | | | | | | | 2020/984 | Nathan Andrews looked at initial version as part of DAT work. No review undertaken as yet on the finalised product. | GitHub - u |
| Technology | risk management | DiAGRAM - Digital Archives Graphical Risk Assessment Model - TNA & Uni of Warwick | Tool to assess level of risk to digital collections utilising Bayesian network | DPC listserv | 5-Jan-21 | CG | | Watch | James Doig | Y | | | | 2020/2573 | Participated in TNA workshop in July 2020 | Safeguard |
| Technology | file format | PRONOM - TNA | File Format Registry developed by TNA (supports DROID - file format identification tool). TNA proposing key updates as outlined in R4102021 - includes new data model (graph data, linked data) | OPF Conference documentation 2020 | 5-Jan-21 | CG | | Watch | | | | | | 2021/76 | | |
| Technology | web archiving | Web Curator Tool (WCT) - NLNZ & NL Netherlands | open source workflow management tool for selective web harvesting. It supports selecting, crawling websites, performing QA and preparing websites for ingest to archival storage. | OPF Conference documentation 2020 | 5-Jan-21 | CG | | Watch | | | | | | 2021/76 | Web Archiving Tool, may be useful if NAA decides to undertake this work in future | |

# Product Assessment Template



National Archives of Australia

# Communication Plan

# Thank You!

Australian Government

National Archives of Australia

**naa.gov.au**