

Exploring digital preservation requirements: A case study from the National Geoscience Data Centre (NGDC)

Jaana Pinnick

British Geological Survey, Keyworth, Nottingham

Abstract

Purpose - This case study is based on an MSc dissertation research undertaken at Northumbria University. The aim was to explore digital preservation requirements within the wider NGDC organisational framework in preparation for developing a preservation policy and integrating associated preservation workflows throughout the existing research data management processes.

Design/methodology/approach - This mixed methods case study used quantitative and qualitative data to explore the preservation requirements and triangulation to strengthen the design validity. Corporate and the wider scientific priorities were identified through literature and a stakeholder survey. Organisational preparedness was investigated through staff interviews.

Findings - Stakeholders expect data to be reliable, reusable, and available in preferred formats. To ensure digital continuity, the creation of high quality metadata is critical, and data depositors need data management training to achieve this. Recommendations include completing a risk assessment, creating a digital asset register, and a technology watch to mitigate against risks.

Research limitations/ implications - The main constraint in this study is the lack of generalisability of results. As the NGDC is a unique organisation, it may not be possible to generalise the organisational findings although those relating to research data management may be transferrable.

Originality/value - This research examines the specific nature of geoscience data retention requirements and looks at existing NGDC procedures in terms of enhancing digital continuity, providing new knowledge on the preservation requirements for a number of national datasets.

Keywords Digital preservation, digital continuity, geoscience data management, data centre, digital repository

Article type Case study

Introduction and background

This article explores the requirements of the National Geoscience Data Centre (NGDC) to ensure that the long-term preservation and usability of its digital data are supported and aligned to the corporate aims. It examines the specific characteristics of geoscience and geospatial data and looks at the efficiency of the existing data management procedures in terms of digital continuity within the current challenging funding climate.

The NGDC is the designated repository for the Natural Environment Research Council (NERC) grant-funded geoscience research data and the guardian for many commercially funded datasets. It is hosted by the British Geological Survey (BGS) and responsible for ensuring the availability of the data as one of the NERC Environmental Data Centres. BGS corporate budgets and staffing levels have decreased during 2010-2015, whilst the volume of digital data has more than doubled. Data management needs to consider the existing organisational framework under the Research Councils UK and the Department of Business, Energy and Industrial Strategy (BEIS).

As a public sector organisation BGS is committed, on behalf of NERC as the legal entity, to look after certain geoscience data in its care in perpetuity (Bowie, 2010) and to make most of it openly available to a wide range of stakeholders, who in turn use the data to develop products and services as well as to inform their decision-making. This requires the organisation to monitor the ongoing condition of digital data and to take appropriate actions in collaboration with its stakeholders to ensure the usability, trustworthiness, and future interoperability of those data. These attributes can only be achieved if data remain both accessible and understandable for future users.

The geoscience data held at the NGDC includes a wide range of data types including but not limited to borehole, bedrock, hydrogeology, geochemistry, seismic, marine geoscience, oil and gas, airborne geophysical, and geohazards data. They have been collected and accumulated over long periods of time and are used by industry, manufacturing, construction, and transport sections, as well as the public sector and academia and researchers, to build UK infrastructure, develop insurance and other data products, innovate, build risk models, answer science questions e.g.

1
2
3 in climate change research, and to support many geoscience applications. Many
4 stakeholders (40%) have been using the NGDC datasets for over 10 years. The use
5 of numerous proprietary software packages (Vulcan, MicroStation) over the years,
6 the lack of restrictions in used file formats in the past, and the occasionally
7 incomplete contextual metadata means older digital data is not always easily
8 accessible to current users if preservation actions are not taken at appropriate times.
9 Past decisions made – or not made – by data creators and guardians at the
10 ingestion phase have a direct impact on the data quality today.

11
12
13 An additional strategic driver for building a digital preservation programme is the plan
14 to apply for a Trusted Repository Status under the Data Seal of Approval (DSA) and
15 the International Council for Science World Data System (ICSU-WDS) certification
16 (RDA, 2016), requiring the NGDC to provide evidence for its long-term preservation
17 capability as a data repository. This includes having a continuity plan in place to
18 ensure the ongoing preservation of data holdings, ensuring the integrity and
19 authenticity of the data, and managing the long-term preservation in a planned and
20 documented way. Although a lot of work has already been done, some of it is not
21 documented consistently. The data centre currently holds over 275TB of data and,
22 although it has considered the digital preservation aspect before, has no formal
23 workflows in place to incorporate preservation actions within data management
24 processes outside the Oracle relational database management systems. Maintaining
25 the value and usability of the data means introducing these workflows has to be a
26 priority in going forward.

40 **Methodology**

41
42 This exploratory case study started by reviewing literature in order to place the
43 findings in the context of long-term preservation of digital geoscience data at a digital
44 repository. The field work phase used a mixed methods approach employing both
45 quantitative and qualitative data. The sequential design included quantitative and
46 qualitative data to provide a comprehensive analysis of the case 'by combining
47 information from complementary kinds of data or sources' (Denscombe, 2008). The
48 iterative approach used ethnographic methods acknowledging the role of the
49 researcher as part of the organisation studied (O'Reilly, 2012) and constant
50 comparative analysis to create categories from raw data (Pickard, 2013).
51
52
53
54
55
56
57
58
59
60

1
2
3 Manipulation and interpretation of dataset access statistics, a literature review, and a
4 macroanalysis of corporate documentation provided a framework on which an
5 external stakeholder survey was based. The datasets and products were selected
6 based on them being widely accessed via the BGS OpenGeoScience web service
7 (<http://www.bgs.ac.uk/opengeoscience>) and included borehole geology, DiGMapGB,
8 groundwater flooding aquifer designation, radon, GeoSure, and mining hazards data.
9 These data are used by industry, academia and general public alike. Corporate and
10 academic users were invited to participate in a small online survey. Transcripts from
11 a purposive sample of staff interviews were compared with the findings. The aim was
12 to cover the research data lifecycle from data management planning and data
13 creation through to accession, archiving, preservation, and reuse.
14
15
16
17
18
19
20
21

22 Establishing trustworthiness of the qualitative research strands can be judged by
23 demonstrating credibility, transferability, and dependability (Pickard, 2013). To
24 increase credibility, this study used triangulation in the form of multiple data
25 collection techniques and sources, which also 'reflects an attempt to secure an in-
26 depth understanding of the phenomenon in question' (Denzin, 2012). Transferability
27 was made possible by providing rich contextual data. To strengthen dependability,
28 member checking was employed during the transcript phase.
29
30
31
32
33

34 **Literature review**

35 *Digital preservation*

36 The early report of the Task Force on Archiving of Digital Information 'laid
37 foundations for most subsequent work in the field, and continues to shape the
38 agenda even today' (Brown, 2013). Its first conclusion, still relevant, is one of the
39 cornerstones of this research: 'The first line of defence against loss of valuable
40 digital information rests with the creators, providers and owners of digital information'
41 (Garrett and Waters, 1996).
42
43
44
45
46
47

48 Key areas in digital preservation research are wide-ranging and include, but are not
49 limited to:
50

- 51 • Preservation planning, policy and strategy development (*Farquhar and Hockx-*
52 *Yu, 2007; Becker et al., 2009*)
 - 53 • Preservation metadata and standards (*BS ISO 15836, 2009; Lavoie and*
54 *Gartner, 2013; Library of Congress, 2016*)
- 55
56
57
58
59
60

- Digital preservation infrastructure, models and toolkits (*Jones, 2006; Ruusalepp and Dobрева, 2012; Lavoie, 2014*)
- Trusted/institutional repositories (*Hockx-Yu, 2006, Ball 2010*)

The Digital Preservation Coalition (DPC) Technology Watch series has captured many salient topics (DPC, 2016). As recently as in 2012, Ross expressed 'an urgent need for a theory of digital preservation and curation' (Ross, 2012), listing nine themes (e.g. repository management, preservation as risk management, and preserving the context) agreed by the Digital Preservation Europe (DPE) and providing a framework for digital preservation research.

To create a digital preservation programme, a solid understanding of corporate drivers is required. This may include collection development, use access, information reuse, legal and regulatory compliance, and efficiencies and savings (Brown, 2013). A recent survey by the Information Governance Initiative (IGI) preferred the term 'long-term protection and access' to 'digital preservation' and suggested there is no distinction between permanent retention and keeping digital content for at least ten years (IGI, 2016).

A UK report on research data management suggests greatest benefits are created by developing ingest and access activities (Beagrie *et al.*, 2010). The UK Data Archive argues that 'maximizing ingest processing efforts as early as possible in the process of digital curation ensures that long-term access is available at a lower total cost' (Woollard and Corti, 2014) and emphasises the importance of the researcher creating good documentation to support data preservation (Van den Eynden *et al.*, 2011).

The NGDC has recently put a lot of effort into developing its digital data ingestion processes and offers guidance to depositors on its data portal website. It employs a NERC-wide Data Value Check List to support the appraisal of data in the pre-ingest phase, captures metadata as part of the ingest process, and offers a list of preferred file formats. All this contributes towards creating better data documentation early in the data lifecycle, but more focused preservation planning and process automation is required as data volumes increase.

1
2
3 The draft BGS digital preservation policy indicates it will use the Open Archival
4 Information System (OAIS) reference model covering ingestion, archival storage,
5 data management, administration, preservation planning, and access functions (BSI,
6 2012a). The other high-level standard, BS ISO 16363 (BSI, 2012b), whilst extremely
7 competent, is acknowledged to require too many resources to achieve at this point in
8 time.
9

13 *Repositories and research data management*

14 Ruusalepp and Dobрева recommend the analysis of preservation 'from the point of
15 view of the organisation's business processes and stakeholders' and state
16 'understanding the *specific* requirements and their implication on the preservation
17 infrastructure is important' (Ruusalepp and Dobрева, 2012, *italics in original*). This
18 chimes with McGovern and McKay who state '[o]rganizations cannot acquire ready-
19 made, out-of-the-box digital preservation programs' (McGovern and McKay, 2008).
20

21
22
23
24
25
26 Kenney and Buckley studied the relationship between repositories and digital
27 preservation and concluded 'insufficient attention was being paid to the
28 organizational context of digital preservation programs' (2005, quoted in McGovern
29 and McKay, 2008).
30
31

32
33 In the UK context the expectation is to preserve research data for ten years. The
34 RCUK Data Policy states: 'Data with acknowledged long-term value should be
35 preserved and remain accessible and usable for future research' (RCUK, 2015b).
36 The accompanying guidance notes add: 'To maximise the research benefit which
37 can be gained from limited budgets, the mechanisms for these activities should be
38 both efficient and cost-effective in the use of public funds' (RCUK, 2015a). The
39 NGDC, which receives its funding from NERC, considers the retention of most of its
40 geoscience data to be longer than ten years due to its 'national good' value,
41 unrepeatability, and the high expense of collecting/creating some of the data. It
42 therefore carries even more responsibility for maintaining the long-term usability of
43 its digital assets.
44
45
46
47
48
49
50

51
52 To oversee long-term research data management, several non-profit and open
53 source digital repository solutions and frameworks, including Dryad
54 (<http://datadryad.org/>) and Fedora (<http://fedorarepository.org/>), offer repositories the
55 benefit of wider user communities but require an extra development layer, and a
56
57
58
59
60

1
2
3 resource need, to integrate the technology with existing infrastructure and to capture
4 sufficient metadata. They are also generic data repositories for any types of research
5 data, whereas the NGDC has the added benefit of in-house geoscientists to e.g.
6 support the creation of descriptive metadata for even externally generated data and
7 to provide guidance in the future use of data.
8
9

10
11 Risk assessment and preservation toolkits, such as DRAMBORA (Digital Repository
12 Audit Method Based on Risk Assessment) (DCC and DPE, 2015) and SCIDIP-ES
13 (SCience Data Infrastructure for Preservation with focus on Earth Science) (EC,
14 2014), assist repositories in identifying and documenting risks associated with digital
15 research data. The collaboration of repository staff with geoscientists is essential to
16 capture the preservation metadata and representation information.
17
18
19
20
21

22 Pryor suggests 'lower long-term costs for the preservation of datasets are perhaps
23 the largest shared benefit' (Pryor, 2014) of a shared infrastructure. He also points out
24 that the allocation of public funds in the RCUK data policy does not say explicitly how
25 the costs are covered. NERC has allocated its environmental data centres, including
26 the NGDC, a small percentage from a central grant budget for long-term
27 management of standard grants data. For larger grants this is costed in the project
28 plans.
29
30
31
32
33
34

35 Data repositories can be divided into institutional repositories based at universities;
36 discipline-specific repositories such as the NERC data centres; and repositories
37 attached to publishing houses, which are increasingly storing data underpinning
38 articles in their scientific journals (Campbell, 2015; Jones, 2014).
39
40
41

42 In the UK '[m]uch could be done to consider digital preservation from the outset, to
43 involve the authors and to embed digital preservation into repository workflow, which
44 will ease the later preservation tasks' (Hockx-Yu, 2006). It is an ongoing challenge
45 for repositories 'to balance the need for fixity in the datasets they offer with the
46 fluidity of changing practices in their designated communities' (Daniels *et al.*, 2012).
47
48
49
50

51 A key challenge in the data lifecycle is the transfer of data from the researcher(s) to
52 the repository, which ought to be clear about its role and responsibilities to ensure
53 high-quality metadata is created and ingested with the data for the purposes of
54 discovery, accessibility, restrictions, user terms and conditions, and preservation,
55
56
57
58
59
60

1
2
3 and to provide tools to aid the researchers in the transfer process (Jones, 2014).
4 Data preservation forms the end of the data lifecycle, but efforts to strengthen
5 usability are best made earlier using effective data management planning and risk
6 mitigation as tools to organise the data and capture preservation metadata.
7
8

9
10 A proposed set of criteria to measure metadata quality at repositories includes
11 completeness, accuracy, and consistency (Park, 2009). In another study 39 tools for
12 semi-automated metadata generation were evaluated, suggesting these are
13 'important considering the fast development of digital repositories and the recent
14 explosion of data and information' (Park and Brenza, 2015). The NGDC is planning
15 to identify and evaluate some of the available digital preservation tools to help
16 automate its processes, ideally by integrating open source tools within its in-house
17 ingestion and data management systems. Of interest are for example tools for
18 extracting preservation metadata information and file format identification or
19 validation.
20
21

22
23 The DCC provides resources and guidance for appraising research data and
24 evaluating repositories (DCC, 2016; Whyte, 2015). Whyte's checklist contains three
25 levels of proficiency, from basic to more advanced, and relates to the three-level
26 framework endorsed by the European Commission (EC, 2016).
27
28

29 *Geoscience and geospatial data*

30 According to the NERC Data Policy, 'NERC requires that all environmental data of
31 long-term value generated through NERC-funded activities must be submitted to
32 NERC for long-term management and dissemination' (NERC, 2010). There is a
33 pressing need to evaluate the preservation requirements of all NGDC geoscience
34 data, a large part of which are geospatial.
35
36

37
38 Geoscience data are of numerous types and cover research areas such as climate
39 change, earth characteristics, rocks, sediments and soils, seismology, marine
40 geology, land contamination, geological processes including erosion and volcanic
41 activity, natural resources, and many more. A specific characteristic of all geoscience
42 data is their long validity, which considerably extends the usable lifespan of digital
43 data assuming actions are taken to preserve them. Geoscience is an interpretive
44 discipline working on previous data and hypotheses which must remain available for
45 future research and interpretation. In addition, not all uses for scientific data are
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 known at the time of data creation and capture, as has been the case e.g. in the use
4 of historical environmental and palaeontological data in modern climate change
5 research or of basic borehole data in creating 3D-models.
6
7

8
9 As future research questions cannot usually be predicted, it is essential to document
10 geoscience data well at the earliest opportunity to enhance their future usability and
11 interoperability potential. This includes developing vocabularies relating to the data
12 bearing the risk of changing semantics over time. For example, the BGS Lexicon of
13 Named Rock Units dictionary (<http://www.bgs.ac.uk/lexicon>) provides previous and
14 alternative names for terms currently in use to facilitate the interpretation of historical
15 data.
16
17
18
19

20
21 Some geoscience data acquisition projects are too expensive to repeat, such as
22 drilling deep boreholes to the depth of many kilometres costing tens of millions of
23 pounds. These data must be managed particularly well in the long-term for future
24 users to get the benefit of previous researchers having made the investment to
25 capture the data in the first place. Seismic data originating from earthquakes is
26 another example of unique and unrepeatable data which is useful for e.g. building
27 hazard models.
28
29
30
31
32

33 Geospatial data is defined as 'conveying information about the Earth, the location of
34 specific features, and attributes and properties of those geo-located features'
35 (Library of Congress, 2014). They are data represented in various geographic
36 coordinate systems, and to use the data it is essential to know which system was
37 used to create them. Not all geospatial data is geoscientific; it may also refer to data
38 used in town planning e.g. street information and post codes.
39
40
41
42
43

44 Geographic Information Systems (GIS) is often used as an alternative term for
45 geospatial data, but strictly speaking it refers to geographic vector or raster data
46 stored as layers and used in mapping. GIS data may also include remote sensing
47 and geo-referenced satellite imagery data. The main GIS used at BGS is the
48 proprietary Esri software, which is considered by staff to be an industry *de facto*
49 standard and as such fairly stable and secure in terms of continuity. GIS data are
50 combined and integrated to create maps and derived datasets, and therefore having
51 accurate and sufficient metadata, such as the scale of the data, is a key attribute to
52 their trustworthiness. In product development, it is important to maintain previous
53
54
55
56
57
58
59
60

1
2
3 versions of the master datasets which may get updated, in case of possible litigation
4 action by third parties, should issues arise due to them using inaccurate out-of-date
5 version of data.
6
7

8
9 Key preservation challenges include 'data versioning, file size, proprietary data
10 formats, copyright, and the complexity of file formats' (Sweetkind-Singer *et al.*,
11 2006). The DPC Technology Watch Report, stating 'geospatial data inherits the
12 preservation challenges inherent to all digital information', points out additional risks
13 including the variety of data structures and the granularity at which the data are
14 processed (McGarva *et al.*, 2009). The National Digital Stewardship Alliance report
15 adds to the mix frequently changing data and issues with scale and resolution of the
16 datasets (Morris, 2013).
17
18
19
20
21

22
23 The US National Research Council (NRC) report 'Geoscience Data and Collections:
24 National Resources in Peril' (NRC, 2002) discussed mainly physical data, but its
25 recommendations apply equally on digital data (prioritisation of data difficult or
26 impossible to replace; funding the collection of information about the data). Its criteria
27 of inaccessible data also pertain: data thought lost, residing elsewhere, in proprietary
28 formats, or not properly indexed and curated.
29
30
31
32
33

34 The US Geological Survey (USGS) has provided preservation guidelines for digital
35 scientific data containing descriptions 'of different levels of increasing assurance that
36 digital data will be preserved' (USGS, 2014). Lessons learned from the development
37 of a digital geospatial archive include collecting data on behalf of organisations
38 without a mandate to preserve them; not knowing future priorities and uses of data;
39 the difference between storing data and preserving it; and preserving only the most
40 important data (Erwin *et al.*, 2009).
41
42
43
44
45

46 BGS has contributed to a more recent strand of geoscience data preservation, the
47 SCIDIP-ES, a European Commission 7th Framework earth science project, as a
48 geoscience expert. One of the outputs was an online survey of long-term data
49 preservation practices (EC, 2012). A preservation toolkit is available on the project
50 website, and NERC holds a copy at the Science and Technology Facilities Council
51 (STFC) as a result of BGS's and STFC's joint contribution to the project.
52
53
54
55
56
57
58
59
60

1
2
3 The use of metadata and standards for archiving geospatial data (Hoebelheinrich
4 and Banning, 2008) is vital because the multiplicity of data types and the variety of
5 GIS (geospatial information systems), used to manipulate the data and possibly
6 employing different coordinate systems, mean that combining data from various
7 sources can be problematical depending upon the number of translation/projection
8 stages.

9
10
11
12
13
14 The EU INSPIRE (Infrastructure for Spatial Information in the European Community)
15 Directive came into force in 2007 and will be fully implemented by 2019. It is the
16 main European framework for sharing spatial data across public sector organisations
17 and applies on all environmental data. In the UK it is implemented via a pan-
18 government UK Location initiative, in which BGS has participated by contributing
19 towards the development of Geoscience Markup Language (GeoSciML), a data
20 transfer standard for geological data. The Directive harmonises the vocabulary and
21 aims to provide discovery metadata and increase interoperability of spatial datasets,
22 so it will have a strong impact on preservation planning once implemented. BGS has
23 made its own discovery metadata service available on its website (BGS, 2016a) and
24 complies with the UK GEo-spatial Metadata INteroperability initiative (GEMINI2), the
25 UK Government specification for spatial metadata.

26
27
28
29
30
31
32
33
34
35 The Open Geospatial Consortium (OGC) provides open standards aimed at the
36 global geospatial data community. 'A goal of open standards is to ensure that
37 "interoperability" (the ability to integrate datasets and related services of different
38 types and from different sources) will minimize such costs and problems.' (OGC *et*
39 *al.*, 2015).

40 41 42 43 44 **Findings**

45
46 In this part of the paper, the first section discusses findings from the corporate
47 documentation. The subsequent sections describe the outcomes from the
48 stakeholder survey and staff interviews.

49 50 51 *Organisational framework and resources*

52
53 The corporate vision of BGS is 'to be a global geological survey, working with new
54 technology and data to understand and predict the geological processes that matter
55 to people's lives and livelihoods' (BGS, 2014). The strategy points to the national
56
57
58
59
60

1
2
3 geological database as one of the core strengths of the organisation and states it
4 'will play a critical role in BGS in the next decade' (BGS, 2014).
5
6

7 The documentary analysis suggested a group of aspects which have an impact on
8 research data management:
9

- 10 1) priority science areas
- 11 2) priority data and datasets
- 12 3) stakeholders, partnerships, user communities
- 13 4) corporate aims and drivers, change drivers
- 14 5) data management drivers and requirements
- 15 6) long-term outlook, data reuse, digital continuity
- 16 7) challenges and opportunities

17
18 The key science areas are Sustainable Natural Resources; Environmental or Climate
19 Change; and Environmental Hazards. A fourth area, entitled Discovery Science, is a
20 funding stream supporting 'excellent environmental research that is driven by
21 curiosity rather than by NERC's wider strategic priorities' (BGS, 2016b). As some
22 research projects and programmes are still in development, related corporate
23 datasets may not be fully available to wider public yet.
24
25

26
27 The BGS has legislative obligations to manage some types of data, e.g. borehole
28 data collected under the Petroleum Operations Notice 9 (PON9), or the Mining
29 Industry Act of 1926. It accepts voluntary donations of other earth science data types
30 and makes data available under the Public Records Act (PRA) 1958/1967, the
31 Freedom of Information Act (FOIA) 2000, and the Environmental Information
32 Regulations (EIR) 2004 (Bowie, 2010).
33
34

35
36 Corporate drivers affecting data preservation are diverse: the need to create
37 efficiency savings; the drive for openness of publicly funded research data; the need
38 to address societal and scientific challenges; increased collaboration between
39 scientific disciplines; the need to update the skills and technologies employed by the
40 organisation; the financial drivers to innovate and to look for new opportunities to
41 contribute towards economic growth.
42
43

44
45 With fewer resources to maintain datasets, they risk eventually becoming unusable
46 and not being further developed into data products. If data are not accessible, a loss
47 of potential new business and product development, and of investment and
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 credibility in favour of other companies, may ensue. Building services using funds
4 from various sources poses a complex task.
5

6
7 Scientists have little time available to condition data for archiving after project
8 closure, as they are expected to quickly move on to the next project. Reducing
9 central funds tensioned against strategic/science needs means limited funds are
10 available to maintain existing datasets, so they may not be as well looked after as
11 they could, reflecting one of the modern challenges of data and information
12 management. However, the funders expect to receive high-quality data as a return
13 for their investment, so a culture change is required of project managers to include
14 the required resources as part of their project lifecycle and management.
15
16
17
18
19

20
21 The NGDC is interested in process automation for metadata creation, appraisal and
22 ingestion, such as the tools and functions analysed by Ruusalepp and Dobrova
23 (2012). Using tools created by the SCIDIP-ES project would support the NGDC in
24 building its capabilities with minimal investment. The SCIDIP-ES Interactive Platform
25 offers a range of services including Preservation Strategy and Certification Toolkits,
26 whereas DRAMBORA provides a framework for assessing digital repository risk in
27 eight categories, including acquisition and ingest, and metadata management. In
28 addition, The National Archives provide a Risk Assessment Handbook and related
29 tools on their website (TNA, 2011).
30
31
32
33
34
35

36 *Development of data management*

37
38 There has been a clear move from analogue to digital data at the NGDC over the
39 last 20 years leading to considerable developments in data management. The
40 introduction of the OGL has encouraged BGS to rethink the way it provides access
41 to data and to consider data quality issues, whilst keeping in mind the legislation
42 giving the tax-payers more access to public sector information.
43
44
45
46

47
48 The changes within BGS over the last two decades include a strong centralisation,
49 standardisation, and deduplication of data. BGS integrated individual datasets into
50 one Oracle database (Geoscience Integrated Database System or GeIDS (Baker
51 and Giles, 2000) and created a discovery metadata schema (Digital Geoscience
52 Spatial Model or DGSM (Smith, 2000)) to support data discovery. Corporate
53 scientific vocabularies and dictionaries completed the creation of BGS integrated
54 data model. This improved the consistency of data but required enhancements in
55
56
57
58
59
60

1
2
3 pre-ingestion and ingestion processes and more knowledge from data centre staff,
4 used to handling a more limited range of data types.
5
6

7 Being the custodian of over 180 years' worth of data makes decision-making about
8 data challenging. Converting large amounts of data to new formats is a labour-
9 intensive, expensive, and often a highly-skilled task. Data may need converting to
10 modern languages before reuse, with a risk of losing original attributes or nuances in
11 interpretation in the process. The organisation may no longer hold software to read
12 all legacy formats, and the annual cost of available licences may be too prohibitive to
13 justify their use.
14
15
16
17
18

19 It may be possible to use crowdsourcing to improve the availability of legacy data
20 e.g. in digitising field notebooks and field slips, as these contain historical information
21 that interests people and may encourage them to contribute. A data amnesty
22 campaign might help capture data from staff who could provide metadata and
23 descriptive information. Insufficiently documented data, with no reuse rights, should
24 consequently be disposed of.
25
26
27
28
29

30 File format obsolescence and technological changes may lead to the loss of access
31 to data. Old media such as DAT (digital audio tape) tapes and floppy disks have
32 become unreadable, and recovery is a specialised task and not guaranteed.
33 Deploying a technology watch and promoting preferred, open, and described formats
34 would help mitigate against obsolescence.
35
36
37
38

39 As for creating full and accurate metadata, its importance and value cannot be
40 overstated. Insufficient metadata hampers the discovery, interpretation, and
41 preservation of data, decreases their value and threatens to make them unusable,
42 and insufficient description may ultimately lead to false scientific conclusions.
43
44
45

46 The BGS discovery metadata, a mature and interoperable schema, provides a
47 powerful tool for data discovery, preservation, interpretation, and validation across
48 different levels (NERC, Data.gov.uk, EC INSPIRE Geoportal). Work is currently
49 underway to build on this by developing a preservation metadata extension
50 containing elements based on the PREMIS (PREservation Metadata:
51 Implementation Strategies) Data Dictionary (Library of Congress, 2016a) and by
52 configuring it to include key geoscience data preservation requirements.
53
54
55
56
57
58
59
60

External stakeholder requirements

BGS data is utilised by the national and local government, academia, industry and commerce. Amalgamating the contradictory requirements of all stakeholders needs to be considered in the business and strategic planning of the organisation.

The small snapshot stakeholder survey (38 respondents, 27% response rate) explored the NGDC user requirements. The results reflected the access statistics with core borehole data being accessed by most. Groundwater and geohazard datasets were used by several participants, indicating that apart from general geoscience, BGS data are used to investigate environmental geohazards. This aligns with and supports the strategic priorities in the list of current key science areas.

The main access to data was via the website, although data are also licensed or shared within collaborative projects. Data were also found using the Discovery Metadata service or by contacting BGS staff directly. 40% have used the data over a 10-year-period, and 80% of users describe the data quality as good, as needed by their own internal systems. Over 70% use raw data, and 25% transform them into other formats including .las (Log ASCII Standard), Oracle spatial, GIS shapefiles, or AGS (Association of Geotechnical & Geoenvironmental Specialists Standard) data.

The main uses of NGDC data are business re-use, decision- or policymaking, research, and science projects. The main benefit of the NGDC is perceived to be the provision of a centralised access point to geoscience data, with reliable long-term availability of data and the provision of an opportunity for data reuse being other advantages. All participants believe the NGDC is a trustworthy repository.

Suggestions to improve long-term data usability range from providing richer extended metadata to complying with industry standards. Standardised web access and data downloads, and data management and preservation training for users, are also suggested. The main concerns highlighted are file format obsolescence and the lack of required preservation skills and resources, followed by the lack of metadata, the creation of digital silos, and increasing volumes of digital data.

The data centre and research data management

The repository requirements include the need to minimise the cost of ingestion and to prioritise the data to be made immediately available, and the creation of sufficient

1
2
3 metadata which complies with industry standards and facilitates data discovery and
4 reuse. The NGDC has the ambition to pull together community user data both within
5 Europe and globally, providing services to other geological surveys and aligned or
6 commercial organisations. This requires an infrastructure resilient enough to cover
7 the expected upturn in data volumes and sufficient resources to manage that
8 potential increase.
9

10
11
12
13 The capability to answer research questions quickly, fully, and accurately using
14 reliable and up-to-date data is essential if BGS is to maintain its reputation as a
15 scientific and archival organisation. A possible legal liability, caused by old versions
16 of datasets or their provenance not being available to customers, can be mitigated
17 against by maintaining a comprehensive, up-to-date data collection, by providing full
18 and accurate metadata record and version control for all datasets, and by using
19 digital object identifiers (DOIs) to provide persistent links to different versions of
20 datasets.
21
22
23
24
25
26

27
28 Clear lines of responsibility govern BGS data management: the Informatics
29 Directorate leads on strategy and implementation; data providers ensure the data
30 quality by selecting appropriate formats and providing sufficient metadata when
31 depositing their data and providing quality science data; NGDC staff ingest, store,
32 help deliver and preserve the data, and provide guidance to users. Science projects
33 use additional measures, such as in-project quality assurance procedures and
34 metadata reports. Information professionals educate scientists, provide guidance on
35 data formats and documentation, and advice on file naming conventions.
36
37
38
39
40

41
42 Ingestion processes have been digitised and standardised resulting in the creation of
43 an online Digital Data Deposit Application (<http://transfer.bgs.ac.uk/ingestion>).
44 Updated procedures are already creating efficiency savings, and automation of
45 metadata creation and an increased use of open formats would further facilitate data
46 collection. The archival storage is good and the cost per unit has been decreasing.
47 Data retrieval and discovery facilities include good external search tools. Web
48 services providing an interface for users to copy raw data to other formats or to
49 embed it within their own systems are increasingly popular.
50
51
52
53
54

55
56 The management of digital data faces challenges, such as predicting future societal
57 and scientific requirements, rapid technological changes leading to increase in real-
58
59
60

1
2
3 time monitoring and sensor data, and the dependence on new technologies vis-à-vis
4 creation, management, storage and retrieval of the data. Corti *et al.* list some
5 advantages specialist data centres offer here, including having appropriate access
6 controls in place and acting as a point of contact between the data and its users
7 (Corti *et al.*, 2014). The NGDC offers limited embargo periods for data creators,
8 metadata-only publication for commercial data, and data licences for businesses
9 who wish to use data in their own products. For this category, metadata may also be
10 hidden.
11

12
13 The long validity of geoscience data means permanent retention is often expected.
14 Adding descriptive and contextual metadata on datasets to ensure their longevity
15 becomes essential, but convincing scientists to provide them remains challenging.
16 With growing volumes of data coming in (175% increase between January 2014 and
17 February 2017), more automated processes, infrastructure and resources are
18 required to manage them.
19

20
21 Drivers influencing data management are the cost of data creation and collection, the
22 support of data publication agenda, the need to create operational efficiencies and to
23 maximise the cost-benefit ratio. As a national repository, the NGDC has a unique
24 selling point in its position of providing open and centralised access to the nation's
25 geoscience data while securing its validity and reliability. It benefits from its
26 connection with the internationally recognised BGS brand, associated with its
27 scientists' ability to answer complex science questions using the available data and
28 knowledge.
29

30
31 The monetary value of data is difficult to measure, although it has been attempted
32 (BGS, 2003). Geologists add value to BGS data e.g. by creating derived datasets
33 and national good science, available to all under the Open Government Licence
34 (OGL) and accessible via the OpenGeoScience website. Links to the underlying data
35 are increasingly a requirement when submitting research papers for publication, and
36 the NGDC provides the DOI service for this purpose.
37

38
39 As data users are increasingly taking the easiest option in sourcing data for their
40 research, the NGDC will need to increase its visibility. Increasing demands from
41 research funders for data to be stored in trustworthy repositories have led the NGDC
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 to capitalise on its current reputation through collecting evidence for the Data Seal of
4 Approval repository certification.
5

6
7 The corporate information culture has evolved considerably. Initially data
8 management procedures were enforced through the appraisal procedure, but a shift
9 towards corporate data holdings and a general acknowledgement of the value of
10 metadata has occurred. More recently, attitudes towards data sharing and reuse and
11 data collection processes have also changed.
12
13
14

15
16 The NGDC needs to continue to enhance its present practices and implement best
17 practice and comply with the latest standards and legislation; create robust
18 metadata; appraise and ingest data efficiently, creating added value; minimise the
19 costs and maximise the long-term usability, accessibility and availability of
20 information. All earth science data must be backed up securely, stored and delivered
21 to a diverse group of users.
22
23
24
25

26
27 If data are to be reused and retained for longer periods, preservation planning
28 becomes critical. If unique or irreplaceable geoscience data are stored in the long-
29 term, a technology watch needs to be in place monitoring risks and changes and
30 keeping research data safe. The value of data can only be realised if the data remain
31 usable and accessible to people who can achieve that value and create data
32 products. All future uses of data are not known today, and even 'useless' data may
33 become important.
34
35
36
37

38
39 The pressure on resources are encouraging BGS to think more strategically. It needs
40 to consider its data management strategy and the skills required to implement it. The
41 data ingested and created need to be fit for multiple purposes and of sufficiently high
42 quality to justify the use of public funds on their preservation. Maintaining data
43 consistency faced with growing data volumes and heterogeneity of data types and
44 formats needs to be tackled to make the best use of the resources.
45
46
47
48

49
50 On a strategic level, the critical role of data management at BGS has been
51 acknowledged by the senior management, who aim at a wider reach of BGS data,
52 forming global partnerships. This necessitates data interoperability and the use of
53 common standards to enhance the longevity and applicability of the data. For this
54 purpose, even scientists agree that basic data management must be seen as critical.
55
56
57
58 This requires high-level corporate support, raising awareness of digital preservation,
59
60

1
2
3 and providing appropriate guidance. As for prioritising data types for preservation,
4 those more expensive to create and collect or completely unique are a priority.
5 Preferred and open formats and consistent file naming should be considered on the
6 outset to enhance data longevity for born-digital data.
7
8

9
10 It is impossible to recollect or re-create long time-series datasets and one-off
11 observations, and some data are too expensive to re-create, e.g. deep boreholes
12 costing tens of millions of pounds to drill. Historical datasets are used to monitor
13 long-term trends and to inform Government public safety decisions in areas such as
14 groundwater flooding and radioactive waste storage. Knowledge on archived
15 datasets may be lost before they are reused. If data are needed urgently, as is the
16 case in emergencies, properly archived and preserved data are available for
17 decision-making immediately.
18
19

20
21 The impact and benefit of using corporate resources on making data more widely
22 available needs to be demonstrated. The cost of the investment must be offset by a
23 future financial or scientific gain as in any other business, possibly even more so
24 when using tax-payers' money. Recreating lost data bears a huge cost and is not
25 always feasible. Using scientists' time to search for data is not a cost-effective
26 option. Decreased funding may lead to weaker quality data and science; however,
27 this very much depends on the strategy basis of using supplied funds. Lack of
28 resources for data management and digital preservation can eventually diminish the
29 value of data as a corporate asset. Raising the awareness about digital continuity to
30 ensure sufficient resources are available may go some way when deciding on
31 corporate priorities and budgets.
32
33

34
35 Data centres have a key role in providing a facility for accessing and reusing
36 scientific data. The NGDC is well placed within the earth science community to
37 provide this service – as a national repository it has little UK competition and can
38 guarantee the quality, validity, and provenance of the geodata in its care.
39
40

41 **Conclusions**

42
43 The users of NGDC geoscience data come from a variety of backgrounds and have
44 varying requirements. This has wider implications on the need to maintain
45 accessibility and usability of borehole, groundwater, and geohazards data, which are
46 often used for ten years or longer by the same stakeholders. The interpretation of
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 what digital preservation means depends on the role and experience of the users.
4 People may know and understand the data very well without fully grasping what
5 actions they could take to improve its long-term usability. It is essential first to
6 develop a common understanding of the concept across the user community, to
7 agree on key messages, and to communicate them to all stakeholders in order to
8 develop the digital preservation culture.
9

10
11
12
13 If the aim is to provide high-quality data and efficiency savings with fewer resources,
14 it will be critical to integrate research activities better with data management
15 procedures. This will ensure some progress is made in enhancing the long-term
16 usability of the UK geoscience data archive for future generations and in maintaining
17 the BGS's ability to answer research questions not yet known. The geoscience data
18 archive is a national asset, which support, BGS in its core mission of providing
19 objective and authoritative geoscientific data, expert services, and impartial advice to
20 all sectors of the UK society. Merely having and storing the data is not sufficient for
21 this aim but a long-term strategy and active, working preservation procedures are
22 required to maintain the value and continuity of the archive.
23
24
25
26
27
28
29

30
31 The development of the NGDC data management procedures and discovery
32 metadata schema, as well as the updated digital ingestion processes, have
33 enhanced the digital continuity of data. However, by making some preservation
34 activities mandatory would further strengthen their reusability. The inspection of
35 current issues with legacy data suggests that by documenting data more extensively
36 at the ingestion stage some of these issues could be avoided or at least be mitigated
37 against – this is a key component of the data deposit portal.
38
39
40
41
42

43 A starting point to implement preservation methods would be to survey all different
44 data types and risks relating to their long-term availability. This could be followed by
45 the creation of a technology watch and a priority list to monitor risks for mitigation
46 and response as and if they are realised. The experience of staff shows that such a
47 function, implemented earlier, would have alerted them to migrating at-risk file
48 formats before them becoming obsolete.
49
50
51
52

53 The stakeholder survey indicated that the user priorities with regard to BGS data are
54 well aligned with the current key science areas and could be used to inform the
55 preservation priorities to some extent. It also showed that in general the users
56
57
58
59
60

1
2
3 consider the NGDC to be a trusted geoscience data repository, a characteristic
4 which should be further enriched by the acquisition of the DSA accreditation and a
5 widening use of the DOIs.
6
7

8
9 Training and raising awareness amongst all the stakeholders is another key aspect
10 of digital preservation and will be addressed in the coming months and years. This
11 point is validated by a comment from a survey participant:
12

13
14 *'The biggest threat to the preservation of digital research data is the lack of*
15 *understanding by those outside the field'.*
16

17
18 In the digital era, there should be nobody left 'outside the field'. The NGDC has
19 always had an interest in the longevity of its data, but in this digital age there are new
20 challenges, risks and opportunities for the data archive which spans over 180 years.
21 This research looked at the current organisational and data management framework
22 of the NGDC, and the findings indicate that the groundwork to manage digital data to
23 a high standard has been done over the last 20 years. Now is the time to take a
24 long-term outlook for data reuse and digital continuity and to develop a functional
25 preservation work plan to make the geoscience data work for the good of the nation.
26
27

28
29 Slightly paraphrasing William Kilbride, the Director of DPC, digital preservation
30 doesn't do itself (Kilbride, 2013). It requires a broad range of skills, a strategic
31 outlook, far-sightedness, long-term planning, financial resources, and above all,
32 perseverance. To conclude in the words of one of the interviewees, in an
33 organisation such as BGS, the custodian of data assets covering almost two
34 centuries:
35
36

37
38 *'We will be judged in the future by how good we tackled data preservation and*
39 *continuity, it's not something that we can ignore. It will become another pillar that*
40 *underpins our well-founded data centre.'*
41
42

43 44 45 46 47 48 **Notes**

49 The article is based on the author's MSc dissertation 'Exploring digital preservation
50 requirements: a case study from the National Geoscience Data Centre (NGDC)'
51 (2016) for Northumbria University.
52
53

54 55 56 **Acknowledgements**

57 The MSc research was supported by the British Geological Survey.
58
59
60

References

- 1
2
3
4
5 Baker, G. and Giles, J. (2000) 'BGS Geoscience Integrated Database System: A
6 repository for corporate data'. In: *Earthwise*, 16, pp.12-13. Available at:
7
8 [http://www.bgs.ac.uk/discoveringGeology/newsAndEvents/earthwise/downloadSearch](http://www.bgs.ac.uk/discoveringGeology/newsAndEvents/earthwise/downloadSearch.cfc?method=viewIssues)
9 [h.cfc?method=viewIssues](http://www.bgs.ac.uk/discoveringGeology/newsAndEvents/earthwise/downloadSearch.cfc?method=viewIssues) (Accessed: 10 August 2016)
10
11
12 Ball, A. (2010) *Preservation and Curation in Institutional Repositories*. Digital
13 Curation Centre Technology Watch Report. Available at:
14 <http://www.dpconline.org/advice/technology-watch-reports> (Accessed: 7 February
15 2016)
16
17
18
19 Beagrie N., Lavoie, B. and Woollard, M. (2010) *Keeping Research Data Safe 2*.
20 Available at: <http://www.beagrie.com/publications/> (Accessed: 7 February 2016)
21
22
23
24 Becker, C., Kulovits, H., Guttenbrunner, M., Strodl, S., Rauber, A., and Hofman, H.
25 (2009) 'Systematic planning for digital preservation: evaluating potential strategies
26 and building preservation plans', *International Journal of Digital Libraries* 10, pp.133-
27 157.
28
29
30
31 Bowie, R. (2010) *BGS Geological Survey – the Legislative Framework*. Available at:
32 <http://www.bgs.ac.uk/services/NGDC/records/policy.html>
33
34
35
36 British Geological Survey (BGS) (2003) *The economic benefit of BGS*. Available at:
37 <http://www.bgs.ac.uk/about/economicbenefits.html> (Accessed: 9 August 2016)
38
39
40 British Geological Survey (BGS) (2014) *Gateway to the Earth: Science for the next*
41 *decade*. Available at: <http://www.bgs.ac.uk/about/strategy.html> (Accessed: 26
42 January 2016)
43
44
45 British Geological Survey (BGS) (2016a) *OGC Catalogue Service for the Web*
46 *(CSW)*. Available at: <http://www.bgs.ac.uk/data/services/csw.html> (Accessed: 2 June
47 2016)
48
49
50
51 British Geological Survey (BGS) (2016b) *Summary of BGS highlights for November*
52 *to January 2016*. Internal BGS report. Unpublished.
53
54
55
56 British Standards Institution (BSI) (2009) *BS ISO 15836 Information and*
57 *documentation – The Dublin Core Metadata Element Set*. London: British Standards
58 Institution.
59
60

1
2
3 British Standards Institution (BSI) (2012a) *BSI ISO 14721:2012 Space data and*
4 *information transfer systems. Open archival information system (OAIS) Reference*
5 *model*. London: British Standards Institution.
6
7

8
9 British Standards Institution (BSI) (2012b) *BS ISO 16363:2012. Space data and*
10 *information transfer systems. Audit and certification of trustworthy digital repositories*.
11 London: British Standards Institution.
12

13
14 Brown A. (2013) *Practical Digital Preservation: A How-to Guide for Organisations of*
15 *Any Size*. London: Facet.
16

17
18 Campbell, J. (2015) 'Access to scientific data in the 21st Century: rationale and
19 illustrative usage rights review', *Data Science Journal*, Vol.13, pp.203-230.
20

21
22 Corti, L., Van den Eynden, V., Bishop, L. and Woollard, M. (2014) *Managing and*
23 *sharing research data: A guide to good practice*. London: Sage.
24

25
26 Daniels, M., Faniel, I., Fear, K. and Yakel, E. (2012) 'Managing fixity and fluidity in
27 data repositories', *Proceedings of the 2012 iConference*, Toronto, ON, Canada,
28 February 7-10, 2012. New York: ACM, pp. 279-286. doi: 10.1145/2132176.2132212.
29

30
31 Denscombe, M. (2008) 'Communities of practice: a research paradigm for the mixed
32 methods approach', *Journal of Mixed Methods Research* 2(3), pp.270-283.
33

34
35 Denzin, N.K. (2012) 'Triangulation 2.0', *Journal of Mixed Methods Research* 6(2),
36 pp.80-88.
37

38
39 Digital Curation Centre (DCC) (2016). *How-to Guides and Checklists*. Available at:
40 <http://www.dcc.ac.uk/resources/how-guides> (Accessed: 16 August 2016)
41

42
43 Digital Curation Centre (DCC) and DigitalPreservationEurope (DPE) (2015) *Digital*
44 *Repository Audit Method Based on Risk Assessment (DRAMBORA)*. Available at:
45 <http://www.repositoryaudit.eu/> (Accessed: 8 August 2016)
46

47
48 Digital Preservation Coalition (DPC) (2016) *Technology Watch Reports*. Available at:
49 <http://www.dpconline.org/advice/technology-watch-reports> (Accessed: 16 August
50 2016)
51
52
53
54
55
56
57
58
59
60

1
2
3 Erwin, T., Sweet-Kinder, J. and Larsgaard, M.L. (2009) 'The National Geospatial
4 Digital Archives – Collection Development: Lessons Learned', *Library Trends* 57(3),
5 pp.490-515.
6
7

8
9 European Commission (EC) (2012) *SCIDIP-ES D15.1 Report on the survey of*
10 *technologies, policies, metadata, semantics and ontologies*. Available at:
11 <http://www.scidip-es.eu/wp-content/uploads/sites/6/2014/11/SCIDIP-ES-DEL-WP15->
12 [D15-1.pdf](http://www.scidip-es.eu/wp-content/uploads/sites/6/2014/11/SCIDIP-ES-DEL-WP15-) (Accessed: 16 August 2016)
13
14

15
16 European Commission (EC) (2016) *European Framework for Audit and Certification*
17 *of Digital Repositories*. Available at:
18 <http://www.trusteddigitalrepository.eu/Trusted%20Digital%20Repository.html>
19
20 (Accessed: 3 June 2016)
21
22

23 Farquhar, A. and Hockx-Yu, H. (2007) 'Planets: Integrated services for digital
24 preservation', *International Journal of Digital Curation* 2(2), pp.88-99.
25
26

27 Garrett, J. and Waters, D. (1996) *Preserving Digital Information: Report of the Task*
28 *Force on Archiving of Digital Information*. The Final Report of the Commission on
29 Preservation and Access and the Research Libraries Group. Available at:
30 <http://www.clir.org/pubs/reports/pub63watersgarrett.pdf> (Accessed: 9 February 2016)
31
32
33

34 Hockx-Yu, H. (2006) 'Digital preservation in the context of institutional repositories',
35 *Program* 40(3), pp.232-243.
36
37

38 Hoebelheinrich, N., and Banning, J. (2008) *An investigation into metadata for long-*
39 *lived geospatial data formats*. National Geospatial Digital Archive project, available
40 at: http://www.ngda.org/reports/InvestigateGeoDataFinal_v2.pdf (Accessed: 28 April
41 2016)
42
43
44

45 Information Governance Initiative (IGI) (2016) *The governance of long-term digital*
46 *information: IGI 2016 benchmark*. Available at: <http://iginitiative.com/reports/>
47
48 (Accessed: 6 June 2016)
49
50

51 Jones, M. (2006) 'The Digital Preservation Coalition', *The Serials Librarian* 49(3), pp.
52 95-104.
53
54
55
56
57
58
59
60

1
2
3 Jones, S. (2014) 'The range and components of RDM infrastructure and services', in
4 Pryor, G., Jones, S. and Whyte, A. (Eds.) *Delivering research data management*
5 *services: Fundamentals of good practice*. London: Facet, pp.89-114.
6
7

8
9 Kilbride, W. (2013) 'Getting Started in Digital Preservation: what do I need to know?'
10 [PowerPoint presentation.] *Getting Started in Digital Preservation: Extra Stop in*
11 *London*. Available at: [http://www.dpconline.org/events/previous-events/1148-getting-](http://www.dpconline.org/events/previous-events/1148-getting-started-in-digital-preservation-extra-stop-in-london)
12 [started-in-digital-preservation-extra-stop-in-london](http://www.dpconline.org/events/previous-events/1148-getting-started-in-digital-preservation-extra-stop-in-london) (Accessed: 9 August 2016)
13
14

15
16 Lavoie, B. (2014) *The Open Archival Information System (OAIS) Reference Model:*
17 *Introductory Guide* (2nd edn.). DPC Technology Watch Report. Available at:
18 <http://www.dpconline.org/advice/technology-watch-reports> (Accessed: 7 February
19
20
21 2016)
22

23
24 Lavoie, B. and Gartner, R. (2013) *Preservation Metadata* (2nd edn.). DPC
25 Technology Watch Report. Available at: [http://www.dpconline.org/advice/technology-](http://www.dpconline.org/advice/technology-watch-reports)
26 [watch-reports](http://www.dpconline.org/advice/technology-watch-reports) (Accessed: 28 February 2016)
27
28

29
30 Library of Congress (2014) *Introduction to Geospatial Resources and Formats*.
31 Available at: http://www.digitalpreservation.gov/formats/content/gis_intro.shtml
32 (Accessed: 24 May 2016)
33

34
35 Library of Congress (2016a) PREMIS Data Dictionary for Preservation Metadata,
36 Version 3.0. Available at: <http://www.loc.gov/standards/premis/v3/index.html>
37 (Accessed: 20 March 2017)
38
39

40
41 Library of Congress (2016b) *Recommended Formats Statement*. Available at:
42 <https://www.loc.gov/preservation/resources/rfs/RFS%202016-2017.pdf> (Accessed:
43
44 22 August 2016)
45

46
47 McGarva, G., Morris, S. and Janée, G. (2009) *Preserving Geospatial Data*. DPC
48 Technology Watch Report. Available at: [http://www.dpconline.org/advice/technology-](http://www.dpconline.org/advice/technology-watch-reports)
49 [watch-reports](http://www.dpconline.org/advice/technology-watch-reports) (Accessed: 28 February 2016)
50

51
52 McGovern, N.Y. and McKay, A.C. (2008) 'Leveraging short-term opportunities to
53 address long-term obligations: a perspective on institutional repositories and digital
54 preservation programs', *Library Trends* 57(2), pp.262-279.
55
56
57
58
59
60

- 1
2
3 Morris, S. (2013) *Issues in the appraisal and selection of geospatial data. An NDSA*
4 *Report*. Available at:
5
6 [http://www.digitalpreservation.gov/ndsa/working_groups/documents/NDSA_Appraisal](http://www.digitalpreservation.gov/ndsa/working_groups/documents/NDSA_AppraisalSelection_report_final102413.pdf)
7 [Selection_report_final102413.pdf](http://www.digitalpreservation.gov/ndsa/working_groups/documents/NDSA_AppraisalSelection_report_final102413.pdf) (Accessed: 28 February 2016)
8
9
10 National Geoscience Data Centre (NGDC) (2016). Available at:
11 (<http://www.bgs.ac.uk/services/ngdc>) (Accessed: 12 August 2016)
12
13
14 National Research Council (NRC) (2002) *Geoscience data and collections: National*
15 *resources in peril*. Available at: [http://www.nap.edu/catalog/10348/geoscience-data-](http://www.nap.edu/catalog/10348/geoscience-data-and-collections-national-resources-in-peril)
16 [and-collections-national-resources-in-peril](http://www.nap.edu/catalog/10348/geoscience-data-and-collections-national-resources-in-peril) (Accessed: 28 February 2016)
17
18
19
20 Natural Environment Research Centre (NERC) (2010) *NERC data policy*. Available
21 at: <http://www.nerc.ac.uk/research/sites/data/policy/> (Accessed: 21 February 2016)
22
23
24 Open Geospatial Consortium (OGC), The International Organization for Standards
25 (ISO) Technical Committee 211 Geographic information/Geomatics, and the
26 International Hydrographic Organization (IHO) (2015) 'A Guide to the Role of
27 Standards in Geospatial Information Management'. Available at:
28 <http://ggim.un.org/docs/Standards%20Guide%20for%20UNGGIM%20-%20Final.pdf>
29 (Accessed: 1st June 2016)
30
31
32
33
34 O'Reilly, K. (2012) *Ethnographic methods*. (2nd edn.) London: Routledge.
35
36
37 Park, J.-R. (2009) 'Metadata quality in digital repositories: a survey of the current
38 state of the art', *Cataloging & Classification Quarterly* 47(3-4), pp.213-228. DOI:
39 10.1080/01639370902737240
40
41
42 Park, J.-R. and Brenza, A. (2015) 'Evaluation of semi-automatic metadata generation
43 tools', *Information Technology and Libraries*, September 2015, pp.22-42.
44
45
46 Pickard, A.J. (2013) *Research Methods in Information* (2nd edn). London: Facet
47 Publishing.
48
49
50 Pryor, G., Jones, S. and Whyte, A. (Eds.) (2014) *Delivering research data*
51 *management services: Fundamentals of good practice*. London: Facet.
52
53
54 Pryor, G. (2014) 'Options and approaches to RDM service provision', in Pryor, G.,
55 Jones, S. and Whyte, A. (eds.) *Delivering research data management services:*
56 *Fundamentals of good practice*. London: Facet, pp.21-40.
57
58
59
60

1
2
3 Research Councils UK (RCUK) (2015a) *Guidance on best practice in the*
4 *management of research data*. Available at: [http://www.rcuk.ac.uk/RCUK-](http://www.rcuk.ac.uk/RCUK-prod/assets/documents/documents/RCUKCommonPrinciplesonDataPolicy.pdf)
5 [prod/assets/documents/documents/RCUKCommonPrinciplesonDataPolicy.pdf](http://www.rcuk.ac.uk/RCUK-prod/assets/documents/documents/RCUKCommonPrinciplesonDataPolicy.pdf)
6
7

8 (Accessed: 17 February 2016)
9

10 Research Councils UK (RCUK) (2015b) *RCUK Common Principles on Data Policy*.
11 Available at: <http://www.rcuk.ac.uk/research/datapolicy/> (Accessed: 17 February
12 2016)
13
14

15 Research Data Alliance (RDA) (2016) *DSA–WDS Partnership Working Group*
16 *Catalogue of Common Requirements*. Available at: [https://rd-](https://rd-alliance.org/system/files/DSA%E2%80%93WDS%20Catalogue%20of%20Common%20Requirements%20V2.2.pdf)
17 [alliance.org/system/files/DSA%E2%80%93WDS%20Catalogue%20of%20Common](https://rd-alliance.org/system/files/DSA%E2%80%93WDS%20Catalogue%20of%20Common%20Requirements%20V2.2.pdf)
18 [%20Requirements%20V2.2.pdf](https://rd-alliance.org/system/files/DSA%E2%80%93WDS%20Catalogue%20of%20Common%20Requirements%20V2.2.pdf) (Accessed: 12 May 2016)
19
20
21
22

23 Ross, S. (2012) 'Digital preservation, archival science and methodological
24 foundations for digital libraries', *New Review of Information Networking* 17(1), pp.43-
25 68.
26
27

28 Ruusalepp, R. and Dobрева, M. (2012) *Digital Preservation Services: State of the Art*
29 *Analysis*. Available at: www.dc-net.org/getFile.php?id=467 (Accessed: 8 February
30 2016)
31
32
33

34 Smith, I. (2000) 'The Digital Geoscience Spatial Model: The shape of the BGS of the
35 future'. *Earthwise*, 15, p.8. Available at:
36 [http://www.bgs.ac.uk/discoveringGeology/newsAndEvents/earthwise/downloadSearch](http://www.bgs.ac.uk/discoveringGeology/newsAndEvents/earthwise/downloadSearch.h.cfc?method=viewIssues)
37 [h.cfc?method=viewIssues](http://www.bgs.ac.uk/discoveringGeology/newsAndEvents/earthwise/downloadSearch.h.cfc?method=viewIssues) (Accessed: 10 August 2016)
38
39
40
41

42 Sweetkind-Singer, J., Laarsgaard, M.L. and Erwin, T. (2006) 'Digital preservation of
43 geospatial data', *Library Trends* 55(2), pp. 304-214.
44
45

46 The National Archives (TNA) (2011) *Risk Assessment Handbook*. Available at:
47 [http://www.nationalarchives.gov.uk/information-management/manage-](http://www.nationalarchives.gov.uk/information-management/manage-information/policy-process/digital-continuity/risk-assessment/)
48 [information/policy-process/digital-continuity/risk-assessment/](http://www.nationalarchives.gov.uk/information-management/manage-information/policy-process/digital-continuity/risk-assessment/) (Accessed: 10 August
49 2016)
50
51
52

53 United States Geological Survey (USGS) (2014) '*USGS Guidelines for the*
54 *Preservation of Digital Scientific Data*'. Available at:
55 <http://www2.usgs.gov/datamanagement/documents/USGS%20Guidelines%20for%2>
56 [0](http://www2.usgs.gov/datamanagement/documents/USGS%20Guidelines%20for%2)
57
58
59
60

1
2
3 [0the%20Preservation%20of%20Digital%20Scientific%20Data%20Final.pdf](#)

4
5 (Accessed: 4 May 2016)

6
7 Van den Eynden, V., Corti, L., Woollard, M., Bishop, L. and Horton, L. (2011)

8 *Managing and sharing data: Best practice for researchers*. UK Data Archive.

9 Available at: <http://www.data-archive.ac.uk/media/2894/managingsharing.pdf>

10
11 (Accessed: 7 June 2016)

12
13
14 Whyte, A. (2015) *'Where to keep research data: DCC checklist for evaluating data*

15 *repositories' v. 1.1* Edinburgh: Digital Curation Centre. Available on:

16
17 <http://www.dcc.ac.uk/resources/how-guides-checklists/where-keep-research-data>

18
19 (Accessed: 4 May 2016)

20 21 22 23 24 **About the author**

25
26 Jaana Pinnick works in geoscience data management and digital preservation at the

27 British Geological Survey. She is also the Chair of the Midlands Branch of the

28 Information and Records Management Society (IRMS). She can be contacted on

29
30 jpak@bgs.ac.uk.