

# The International Journal of Digital Curation

Issue 1, Volume 6 | 2011

## Use and Impact of UK Research Data Centres

Ellen Collins,  
Research Information Network

### Abstract

UK data centres are an important part of efforts to gain maximum value from research data. However, if they are to operate effectively, the services that they provide must be based upon an understanding of researchers' practices and needs. Furthermore, in order to build a case for ongoing funding, data centres must be able to demonstrate their value to researchers work and, increasingly, their contribution to wider political "impact" agendas. This paper presents the findings of a survey of users of five UK data centres. It suggests that research data centres are highly valued by their users. Benefits appear to be particularly strong around improving research efficiency, especially access to data. Data centres are less important in terms of stimulating novel research questions. Despite a few interesting cases of observable impact, in the main it remains difficult to understand the wider reach of research which draws upon data centre resources.<sup>1</sup>

---

<sup>1</sup> This paper is based on the paper given by the authors at the 6th International Digital Curation Conference, December 2010; received December 2010, published March 2011.

The *International Journal of Digital Curation* is an international journal committed to scholarly excellence and dedicated to the advancement of digital curation across a wide range of sectors. ISSN: 1746-8256 The IJDC is published by UKOLN at the University of Bath and is a publication of the Digital Curation Centre.





## Introduction

In recent years, there has been a growing understanding of the importance of data as a primary research output. This is demonstrated, for example, by the increasing interest in, and policies for, data preservation and management on the part of research councils. Such attention is due in part to recognition of the data deluge: the enormous quantities of data created as a result of changing research processes and, in particular, the growth of e-science. (Hey & Trefethen, [2003](#)). The production of large datasets is expected to continue to expand and data outputs from so-called “small science” are also being recognised as important resources for preservation and reuse (Lyon, [2007](#)).

The proliferation of data creates both opportunities and challenges for researchers. As the US National Science Board ([2005](#)) has pointed out, large amounts of data can be aggregated to permit new forms of scientific research, and data can be reused and aggregated to answer research questions beyond those for which it was originally gathered. However, such work depends upon data being readily discoverable and in a format which easily allows reuse. In practice, publication methods and locations are varied, which has tended to limit the ways in which datasets can be used by other researchers (UKRDS, [2008](#)).

It has long been recognised – by both policy bodies and by researchers themselves – that a suitable structure for the collection, management and access of research data is crucial if that data is to be useful for other researchers (Lievesley & Jones, [1998](#), Research Information Network, [2008a](#)). Data centres present one important attempt to ensure that the potential of research data is fully exploited. They do this in two main ways. First, they attempt to ensure that data is discoverable by providing a single location where researchers can deposit their work and tools which allow other researchers to find and access it. However, it is important to note that most have not yet achieved comprehensive coverage within their discipline, with Beagrie et al. finding that only 18% of researchers deposit their work in a data centre – although 43% use one to access data ([2009](#)). Furthermore, even data which is deposited may be governed by restrictions on access or reuse for ethical reasons or to enable the original researchers to get the maximum publication benefit from it before opening it up to the wider research community (Research Information Network, [2008b](#)).

The second important role played by research data centres is to provide a support structure for researchers who need to get their data – and metadata – into shape prior to deposit. Most researchers are not accustomed to preparing their data for use by others and few have the time to do so, or to learn to do so (Research Information Network, [2008b](#)). Indeed, this was identified as a major barrier to fully open research data by the Australian National Data Service Technical Working Group ([2007](#)). Repositories are an important resource for these researchers, providing advice, guidance and structures to ensure that data is ready for reuse.

The shape, size and number of datasets will continue to change over coming years. New research techniques and technologies may also allow novel uses of existing datasets (National Science Board, [2005](#)). Thus, it is necessary periodically to evaluate the usage of research data centres, to ensure that they are relevant and meeting the needs of researchers, both as depositors and as users of data. This can also help to inform future service development and ensure that the limited available funding is



## Findings

### *Patterns in Research Data Centre Usage*

The survey began by gathering some basic information about data centre users. Figure 2 shows the sectors in which survey respondents work. Sectors which account for at least 10% of respondents have been highlighted for ease of reference. The distribution of users for individual data centres varies considerably. The atypical result for the NGDC is probably due to the large number of respondents who work for the BGS, which is a public research organisation. Overall, while the majority of users appear to come from academic backgrounds, there is a relatively strong showing for public research organisations for the BADC, and for central and local government and community and charity organisations for the ADS. Use by researchers outside remains relatively small scale compared to use within academia.

Sector	Data centre				
	ADS	BADC	CDS	ESDS	NGDC
Academic	<b>51%</b>	<b>67%</b>	<b>95%</b>	<b>78%</b>	4%
Public research organisation	0%	<b>21%</b>	4%	4%	<b>86%</b>
Private research organisation	8%	2%	1%	4%	0%
Private / Independent researcher	8%	2%	1%	2%	2%
Central / local government	<b>10%</b>	5%	0%	7%	2%
Business	5%	1%	0%	1%	0%
Community / charity organisation	<b>11%</b>	1%	0%	2%	0%
Other (please specify)	8%	3%	1%	3%	6%
N	83	759	200	292	51

Figure 2. Sectors in which Survey Respondents Work, by Data Centre.

Figure 3 shows how survey respondents use research data: respondents could tick all categories which applied. The most common response in each category has been highlighted. The findings suggest that most data centres are used in one primary way by a large number of researchers, with other uses being less widespread. For the ADS and CDS, the data is primarily accessed for reference purposes, while most users of the BADC and the CDS use data for their research work. These differences may be due in part to the nature of the data centre's holdings; the CDS, for example, holds primarily experimental data which academic researchers reference routinely and frequently, while the ESDS holdings are much better suited to form the basis of a user's own research. The NGDC appears to be well-used for a range of purposes, although this may again reflect the small and homogeneous group of survey respondents.

Use of research data	Data centre				
	ADS	BADC	CDS	ESDS	NGDC
For research	51%	<b>75%</b>	48%	<b>88%</b>	72%
For combining with other data	46%	46%	39%	34%	<b>81%</b>
For reference	<b>79%</b>	22%	<b>82%</b>	21%	72%
As a basis for further data collection	51%	7%	24%	17%	58%
N	72	713	190	289	42

Figure 3. Use to which Research Data is Put, by Data Centre.

Figure 4 shows trends in use over time for each data centre. Again, the most common response in each category has been highlighted. For most centres, usage has broadly stayed the same over time, in some cases with fluctuations. Data centres' own research shows that overall usage has gone up over time. When combined with the results presented here, this suggests that data centres are gathering new users rather than seeing more intense usage from existing ones. That said, the ADS and NGDC have seen many users increase the frequency of their usage over time, while the BADC has seen over twice as many users decrease their usage over time as increase it.

Frequency of use	Data centre				
	ADS	BADC	CDS	ESDS	NGDC
Decreased over time	13%	26%	16%	21%	12%
Stayed the same / fluctuated	39%	<b>62%</b>	<b>54%</b>	<b>60%</b>	33%
Increased over time	<b>48%</b>	12%	30%	20%	<b>55%</b>
N	83	730	200	283	50

Figure 4. Frequency of Use of Research Data, by Data Centre.

Where data centre usage has changed, the survey indicated that this is often to do with changes to the researcher's circumstances rather than changes made by the data centre. Increases in usage were attributed to new research questions (28%) and changes in role or position (21%), while decreases were attributed to research questions being addressed or shifting emphasis (42%) or changes in role or position, including retirement (33%). However, 26% of respondents said that they increased their usage due to improvements to the range and quality of data available, suggesting that developments made by the data centre can have a positive impact upon usage.

Use is not limited, of course, to downloading data; data centres also provide a service for researchers who wish to share data that they have created. Accordingly, survey respondents were asked about their data sharing habits, and their perception of the impact that data centres have had upon data sharing and reuse within their disciplinary field. Figure 5 shows the proportion of data-creating researchers who submit content to data centres; the most common response for each data centre has been highlighted. For the BADC, this question was phrased slightly differently and so results are not comparable.

Submission of new data	Data centre			
	ADS	CDS	ESDS	NGDC
No, never	36%	32%	<b>70%</b>	8%
Yes, sometimes	<b>52%</b>	<b>41%</b>	20%	27%
Yes, always	11%	28%	11%	<b>65%</b>
N	44	111	82	26

Figure 5. Submission of New Research Data, by Data Centre.

The "N" figures for this information are themselves interesting; for most data centres, roughly half of the survey respondents created new data; for the ESDS it was roughly a third. Users of both the ADS and the CDS were most likely to submit some, but not all, of their data to a data centre. The low levels of submission – and low overall data creation – from ESDS users is probably related to the nature of research in the social sciences, where the high long-term value of data and ethical concerns about human subjects can inhibit sharing. The high submission levels by users of the NGDC

may well be explained by the large proportion of respondents who are based at the BGS and who are therefore bound by that organisation's data sharing policies.

Figure 6 shows how survey respondents perceive data centres to have improved the culture of data sharing and reuse within their own research communities. The most common response for each data centre has been highlighted. Users of all data centres seem to see a strong improvement in data sharing and reuse, which they consider attributable to the existence of the data centre. The slightly less enthusiastic response from users of the ESDS may be because this data centre has been in operation for more than 40 years, with a web presence for the last 15, meaning any behavioural changes have already filtered into the mainstream.

Extent of improvement in data sharing and reuse	Data centre				
	ADS	BADC	CDS	ESDS	NGDC
Not at all	0%	2%	1%	7%	3%
To a small extent	16%	29%	27%	40%	30%
To a large extent	<b>84%</b>	<b>69%</b>	<b>72%</b>	<b>54%</b>	<b>68%</b>
N	61	601	164	244	37

Figure 6. Extent of Improvement in Data Sharing and Reuse Due to Data Centre, by Data Centre.

### *Impact of Data Centres on Researchers and their Work*

Beyond understanding how research data centres are used, the survey also sought to establish the impact of data centres on researchers and their work. Figure 7 shows how important researchers consider data accessed via data centres to be to their research. For each centre, the most common response has been highlighted. For most data centres, researchers consider the data to be either "very important" or "essential". The NGDC represents a particularly extreme case, with 85% of researchers considering the data to be "essential"; this may again be due to the homogeneous nature of respondents to the survey for that particular data centre.

Importance	Data centre				
	ADS	BADC	CDS	ESDS	NGDC
Not at all important	0%	1%	1%	2%	0%
Not very important	6%	7%	7%	5%	0%
Quite important	<b>34%</b>	32%	28%	25%	5%
Very important	31%	<b>37%</b>	30%	27%	10%
Essential	29%	24%	<b>34%</b>	<b>41%</b>	<b>85%</b>
N	65	700	189	282	40

Figure 7. Importance of Accessed Data to Researchers' Work, by Data Centre.

Researchers were asked to gauge their level of agreement with several statements about the benefits of being able to access data from a data centre. These can be broadly grouped into four main areas: research efficiency, research practice and quality, research novelty, and researcher training. These areas are examined in greater detail, and by data centre, in Figures 8, 9, 10 and 11. Overall, however, the most widely-cited benefits fall into the research efficiency category. Most of the free text responses about data centre benefits also concerned research efficiency. Benefits relating to research novelty achieved the lowest overall rate of agreement. However, even the least widely-

supported statement, about new intellectual opportunities, was rejected by only 22% of researchers overall, while 36% agreed with it “to a large extent”.

Figure 8 shows the percentage of researchers, by data centre, who agreed “to a large extent” with the statements about research efficiency. The most common response has been highlighted for each data centre. The most widely-agreed benefit for all data centres was around saving time for data acquisition and processing, which was a primary aim of many data centres when they were established.

Research efficiency benefit	Data centre				
	ADS	BADC	CDS	ESDS	NGDC
It has reduced the time required for data acquisition / processing	<b>79%</b>	<b>68%</b>	<b>76%</b>	<b>80%</b>	<b>92%</b>
It has improved the efficiency of research	<b>79%</b>	62%	75%	67%	89%
It has reduced the financial cost of data acquisition / processing	65%	62%	61%	73%	78%
It has reduced duplication of effort (i.e. unnecessary recreation of data)	57%	57%	68%	62%	81%
It has enabled me to undertake a greater quantity of research	52%	42%	50%	54%	77%

Figure 8. Research Efficiency Benefits of Data Centres, by Data Centre.

Figure 9 shows the percentage of researchers, by data centre, who agreed “to a large extent” with statements about research practice and quality. Again, the most common response has been highlighted for each data centre. The rankings here are less clear than those for research efficiency; for most data centres, however, the most widely-agreed benefit appears to be an improvement in the evidence base for research. Few respondents for each data centre, other than the NGDC, agree that the centre has increased the use of data in their own research. This suggests that data centre benefits are concentrated around giving researchers access to core data, rather than encouraging them to undertake research which is more heavily data-focused.

Research practice and quality benefit	Data centre				
	ADS	BADC	CDS	ESDS	NGDC
It has improved the evidence base of my research	<b>58%</b>	46%	<b>60%</b>	<b>56%</b>	<b>77%</b>
It has helped to improve the quality of my research outputs	56%	<b>47%</b>	58%	<b>56%</b>	69%
It has improved the quality of the data I use within my research	55%	<b>47%</b>	48%	51%	72%
It has increased the use of data in my research	48%	40%	38%	46%	75%

Figure 9. Research Practice and Quality Benefits of Data Centres, by Data Centre.

Figure 10 shows the percentage of researchers, by data centre, who agreed ‘to a large extent’ with statements about research novelty. Again, the most common response has been highlighted for each data centre. Benefits here appear to be concentrated around enabling research that might not otherwise have happened. It is not entirely clear whether the research would not have happened because the

techniques used are only possible due to the aggregation of data through the data centre, or whether the research would not have happened because the data itself would have been inaccessible. Given the overall character of responses to this set of questions, the latter seems more likely: for most data centres there was relatively limited agreement with the statements about new types of research and new intellectual opportunities.

Research novelty benefit	Data centre				
	ADS	BADC	CDS	ESDS	NGDC
It has enabled research to go ahead that otherwise might not have done	62%	48%	41%	60%	89%
It has permitted more novel research questions to be answered / tackled	46%	38%	44%	51%	69%
It has enabled new types of research to be carried out	56%	34%	33%	44%	77%
It has created new intellectual opportunities (e.g. merging of several data sets to answer new questions)	51%	33%	27%	40%	69%

Figure 10. Research Novelty Benefits of Data Centres, by Data Centre.

Finally, a single question was asked about researcher training, the results of which are presented in Figure 11. This area revealed the greatest differences between data centres. The free text comments suggest that most of these benefits stem from the availability of a resource that can introduce researchers to the important data sets within their field and demonstrate best practice in collecting and handling data. This relates closely to the two founding aims of many data centres: to widen access to data sets and to improve researcher practices around curation and storage of data.

Researcher training benefits	Data centre				
	ADS	BADC	CDS	ESDS	NGDC
It has enabled me to improve researcher training	69%	48%	34%	47%	100%

Figure 11. Researcher Training Benefits of Data Centres, by Data Centre.

### *Wider Impact of Data Centres*

As set out in the introduction, it is important that data centres are able to demonstrate their impact and value beyond the academic community. Figure 2 above indicated the broad reach of some data centres, and some of the researcher benefits outlined in the previous section therefore relate to researchers working in non-academic settings. However, the aim of “impact” is to reach beyond the research community, and the project therefore sought evidence that this was happening.

Figure 12 shows the intended audiences for research produced using data acquired from data centres. Of course, an intended audience is by no means the same thing as an actual audience; surveying the supply side gives little concrete information about demand or usage. Nonetheless, this table gives an indication of the potential reach of research based upon data centres’ holdings. Furthermore, it seems unlikely that researchers would continue to produce work for a certain audience if that audience

showed no interest in it whatsoever.

Research audience	Data centre				
	ADS	BADC	CDS	ESDS	NGDC
Academics	69%	84%	90%	79%	68%
Policy makers	16%	28%	2%	47%	66%
Individuals within your organisation	43%	21%	26%	26%	82%
Own use only	31%	13%	26%	16%	16%
Business	19%	4%	7%	7%	75%
Unknown	4%	3%	2%	1%	7%
Other	24%	6%	8%	10%	25%
N	70	710	189	287	44

Figure 12. Intended Research Audiences, by Data Centre.

For all data centres except the NGDC, the primary audience is academics, while for the NGDC the primary audience is other individuals within the user's organisation. This is consistent with the profile of survey respondents which, as we have said above, was particularly inclined towards BGS staff members for the NGDC survey. Each data centre presents some interesting trends. NGDC users appear to be quite outward-facing, with high response rates for a number of different audiences, and a low number of users suggesting that their research was for their own use only. Again, this probably reflects the fact that respondents to this survey were, for the most part, BGS staff members and therefore part of an organisation with a strong cross-sectoral brief. CDS users, on the other hand, have a much more homogeneous target audience, with most of the attention focused on academics and very little on audiences outside the respondent's own organisation. This may well be due to the highly technical nature of the data within the CDS, which may require several degrees of analysis before it can be translated to non-specialist audiences.

Policy makers achieved a reasonably strong showing across all data centres except CDS and, to a lesser extent, ADS, while business was an important end audience for users of the ADS and NGDC. The low level of response to the 'unknown' category suggests that researchers take a strong interest in producing work which is useful for specific audiences. The "other" category, which scored particularly highly for ADS and NGDC, in many cases contained very specific sub-groups which could be considered part of "business" or "policy makers". However, "students" and "the general public" came up relatively frequently for all data centres; it may be that if these had been a prompted answer that they would have scored more highly still.

As mentioned above, although these intended audiences give a useful perspective on the reach that data centres may have, they are not themselves indicative of wider impact. We attempted to capture such impact through short case studies of some of the research projects mentioned by survey respondents. This was not an easy task. Beyond the well-rehearsed problems about connecting "impact" directly to academic activity, the great majority of survey respondents were professional researchers with no direct view of the ultimate use and outcomes of their data centre enabled research.

This difficulty is confirmed by the free text responses to the survey question about impact, most of which related to the impact of a data centre on a researcher's own work, rather than their work on society. Some respondents were overtly hostile to the



notion that research impact is worth measuring, while many others suggested areas where their research could add value, but no evidence to indicate that it had already done so. However, a small number of researchers were able to provide examples of instances where their work had influenced practice and policy outside academia.

For the most part, the subject matter of the data centre determines where research is likely to have impact. Thus, most of the impacts identified by researchers using BADC data were in environmental fields, while ESDS researchers influenced areas of social policy. Broken down by type of impact, most responses talked about either new models or tools which helped to support decision making by public or private bodies, new policies and regulatory controls, and development of new commercial materials, particularly drugs. Effects can be observed in the public, private and voluntary sectors, although given the small size of the sample for this question it would be unwise to attempt any estimation of the distribution between these three groups.

## Conclusions

This research suggests that UK data centres are playing a valuable role in the research community. They are making it easier and cheaper to access data, are supporting new ways of doing research and are helping researchers to manage and curate their own data more effectively. Overall, they are fulfilling many of the needs identified within the literature. However, it is important to emphasise that this survey only contacted existing users of research data centres: it is possible that there are researchers in these fields who are not accessing the benefits brought by data centres because they are not yet using them. Further study should focus on these non-users, and in particular any barriers which might prevent their use of data centres.

Most data centres have a fairly homogeneous user group consisting of researchers from academia or from public research institutions. ADS users represent the widest range of backgrounds, but overall there is relatively little usage from private researchers or business. This may be related to the nature of the research – it is possible that non-academic researchers have less time or interest in completing a user survey. However, data centre funders should consider whether they can encourage use from a more diverse community. Most users reported that they have maintained their level of usage (with some fluctuations) over their period of data centre use. Data centres themselves report increasing overall levels of usage, suggesting that they are attracting new users rather than encouraging existing ones to increase their intensity or frequency of usage by, for example, exploring new kinds of research question.

Indeed, there was a strong sense overall that, while data centres may have an effect on some elements of researcher behaviour – such as the propensity to share data – they are having a relatively weak effect on the types of research that are undertaken. This was the least widely-supported benefit of data centres; those to do with research efficiency and cost achieved much higher levels of agreement. It may be that research novelty is an important function of data centres for some researchers; for the majority, however, it is less important than the ability to access data quickly and cheaply. In developing further services for researchers, data centres should take into account the relative value of these benefits to researchers.

When considering possible future service developments, it is also important to note that on some issues the views of users were not homogeneous across data centres. For example, the ways in which researchers use the data they acquire varies by data centre, as does their view of the value of the data centre to researcher training. In many cases this will be determined by the content of the centre as well as the needs of the researchers. However, it raises interesting questions for the UKRDS, particularly in terms of researcher training and development, and the extent to which a national framework can be sensitive to the specific needs of researchers in different disciplines.

This research also confirmed the difficulty of tracing the impact of academic research. Several researchers suggested that their work *could* have an impact, in some cases suggesting very specific ways in which this could happen, but were not able to show that it had actually occurred. However, a few researchers were able to cite specific instances where their work had supported developments in public, private or voluntary sector organisations. The fact that researchers cannot always see the impact of their work suggests that such impact may be more widespread than this survey reveals. Future research could address this problem in more depth by contacting the eventual end users of the research, although this is bound to present new problems around traceability and access.

There are some other important questions that this research was not able to address. The value of data centres to small science was not covered explicitly within the survey and in the context of developments such as Dryad UK it would be useful to understand whether and how an established data centre might support researchers in these fields. The research also highlights the number of researchers that produce new data but do not submit it to a data centre. It is likely that in some cases this is because they do not think that the data has any value to other researchers; in others, however, potentially useful data will be going unshared. Funders should consider how they can encourage researchers to submit data to data centres, and in particular whether stronger guidelines about data citation might help.

## Acknowledgements

We would like to thank Technopolis, who undertook the survey work which forms the basis of this paper.

## References

- Beagrie, N., Beagrie, R. and Rowlands, I. (2009) Research data preservation and access: The views of researchers, *Ariadne 60*. Retrieved October 26, 2010, from <http://www.ariadne.ac.uk/issue60/beagrie-et-al/#7>.
- Grant, J., Brutscher, P.-B., Kirk, S., Butler, L. and Wooding, S. (2009). *Capturing Research Impacts: A review of international practice*. Rand Europe. Retrieved July 23, 2010, from [http://www.hefce.ac.uk/pubs/rdreports/2009/rd23\\_09/rd23\\_09.pdf](http://www.hefce.ac.uk/pubs/rdreports/2009/rd23_09/rd23_09.pdf).
- Hey, T. and Trefethen, A. (2003). The Data Deluge: An e-Science Perspective. In Berman, F., Fox, G. and Hey, A., Eds. *Grid Computing – Making the Global Infrastructure a Reality*. Chichester, John Wiley & Sons.

- 
- Lievesley, D. and Jones, S. (1998). *An Investigation into the Digital Preservation Needs of Universities and Research Funders*. UKOLN. Retrieved October 26, 2010 from <http://www.ukoln.ac.uk/services/papers/bl/blri109/datrep.html#Heading1>.
- Lyon, L. (2007). *Dealing with Data: Roles, Rights, Responsibilities and Relationships*. JISC. Retrieved July 7, 2010, from [http://www.jisc.ac.uk/media/documents/programmes/digitalrepositories/dealing\\_with\\_data\\_report-final.pdf](http://www.jisc.ac.uk/media/documents/programmes/digitalrepositories/dealing_with_data_report-final.pdf).
- National Science Board (2005). *Long-Lived Digital Data Collections: Enabling Research and Education in the 21<sup>st</sup> Century*. National Science Foundation. Retrieved July 7, 2010, from <http://www.nsf.gov/pubs/2005/nsb0540/>.
- Research Information Network (2008a). *Stewardship of digital research data: A framework of principles and guidelines*. Retrieved July 7, 2010, from <http://www.rin.ac.uk/our-work/data-management-and-curation/stewardship-digital-research-data-principles-and-guidelines>.
- Research Information Network (2008b). *To Share or not to Share: Publication and Quality Assurance of Research Data Outputs*. Retrieved July 7, 2010, from <http://www.rin.ac.uk/our-work/data-management-and-curation/share-or-not-share-research-data-outputs>.
- The ANDS Technical Working Group (2007). *Towards the Australian Data Commons: A proposal for an Australian National Data Service*. Australian Government Department of Education, Science and Training. Retrieved July 7, 2010, from <https://www.pfc.org.au/pub/Main/Data/TowardstheAustralianDataCommons.pdf>.
- UKRDS (2008). *The UK research data service feasibility study: Report and Recommendations to HEFCE*. UKRDS. Retrieved October 26, 2010, from <http://www.ukrds.ac.uk/resources/>.