

20 years of archiving and sharing qualitative data in the UK

Louise Corti

Abstract

In this presentation I will provide an overview of archiving qualitative data in the UK. The UK Data Archive has been collecting and providing free access to high quality data for the social science research base since 1967, and since then has nurtured a close working relationship with data owners and producers across a range of sectors. Our major UK funder, the Economic and Social Research Council (ESRC), established its Data Policy in 1995, including qualitative research, which led to real success for changing the culture of data sharing in the UK. I briefly show how data are being promoted and re-used and compare the state of qualitative data archiving across Europe. I discuss the art of managing and sharing qualitative data, building on the procedures that we have developed, and illustrate an example of sharing emotionally sensitive data.

1. Introduction

This paper provides a short tour through the 20 plus years of archiving and sharing qualitative data in the UK. Disciplinary definitions of qualitative research and data are considered, and what this might mean for archiving. It outlines what users are doing with archived data and how data can be best prepared to reveal context in the data, and improve usability. The practicalities and challenges of sharing qualitative data are highlighted followed by an example of a ‘sensitive’ qualitative data collection that has been shared. The paper notes strategies used by the UK to successfully bring users to the data and encourage the practice of secondary analysis. It concludes by emphasising the power of collaboration and how important it is that data centres in one country work closely together to provide united services to users.

2. Who is sharing qualitative data?

If we consider whether experienced researchers are worried about sharing their qualitative research data, many will likely say yes. Over my 25-year period of working in this field, I see that this is often down to lack of appreciation of what this means and what is involved, and an assumption that it will expose them, their methods, and possibly their participants to some kind of potential damage or harm. Time and resources are also frequently noted as being especially precious and in competition with new research endeavours. However once we begin to unpick concerns and demonstrate the processes of how data can be, and have been successfully prepared and published, the worry often subsides.

If we look across the disciplines that routinely use qualitative methods, what is captured as *data* varies quite a lot. Disciplines like sociology, anthropology and education tend to prefer interviewing and ethnographic methods, generating data such as interview or focus group recordings and transcripts, diaries, other kind of observations, fieldnotes and kinship diagrams. Political science qualitative methods prefer to focus on text analysis using historical and legal documents, newspaper sources, political manifestoes and records of meetings, for example, in the policy environment or in human rights organisations. Humanities scholars such as historians or linguists often focus on a close examination of the source, using methods of text or audio-visual analysis to discover linguistic features of speech and on verbal communication and create transcribed and annotated text and audiovisual data. Figures 1a–d demonstrate the variety of transcriptions that could be created depending on the disciplinary approach taken. They vary from a basic summary of an interview, in Figure 1a, while Figure 1b shows a more detailed transcription of an interview. Figure 1c shows non-verbal communication that has been recorded, and Figure 1d shows a more detailed orthographic transcription, that captures some of the linguistic features of spoken word including pronunciation.

Figure 1a

Short Summary
This is a short passage from an interview conducted by Stef Scagliola on 15 April 2007 at the Imperial War Museum with major general Julian Thompson on the use of oral history in military history. As curator of the Dutch Veterans Interviewproject initiated by the Dutch Veterans Institute in the Netherlands, she was visiting the Imperial War Museum to do research on the requirements and best practices for a large scale oral history collections related to military topics. Julian Thompson agreed to be interviewed about the benefits and risks of using oral sources for research and on how he used the collections of the IWM for his publications. In the interview he stresses four major points. Firstly, that almost all history starts with oral accounts of events and that the trustworthiness of written accounts is often overestimated, secondly that the author of a book needs to know the subject well in order to distinguish facts from fiction, thirdly that despite these merits, oral history is not regarded as a serious discipline in the academic realm. An other important value, is it social character.

Transcript on the basis of light editing

J.T. - The problem which I recognize, and people here will be the first to admit it, with oral history, there are a number of problems, is that if you interview people who are old, their memories sometimes are bad. Quite often they have had a chance to discuss with other people what went on, so they have come to an agreed story.

Figure 1b

Algemene gegevens interview

| | |
|---------------------------|---|
| Titel | Erfgoed van de Oorlog, Bystander Memories, interview 01 |
| Geïnterviewde | [echte naam] |
| Interviewer | [echte naam] |
| Plaats | |
| Tijd | |
| Trefwoorden | |
| Samenvatting | |
| Transcriptie gemaakt door | |
| Eigenaar van het bestand | |

00:00:00

[image: 00:00:00]

IV: En kropen die mannen dan gewoon in de schuur, verstopten ze zich..

GI: Ja, bovenaan was een scheve kap, en daar kropen ze met een ladder naartoe, moesten wij een ladder daar neer zetten en dan kropen ze er naartoe, en dan werd dat dicht gepakt, en dan de ladder weer weg natuurlijk dan een paar uur zaten ze daar stil, der was er zelfs een uit Eindhoven, die was hier gekomen van Eindhoven naar zijn familie, om ja, hier veiliger te zitten, maar die moest hier ook onderduiken, die moest hier ook wegkruipen voor die razzia die er gehouden werd.

IV: Nelly, jullie hadden ook een appelboomgaard. Wat gebeurde daar mee?

GI: Ja, wij hadden een heel grote aard, appelboomgaard daar en peren natuurlijk. En die hadden wij, die moesten geplukt worden natuurlijk, en die hadden wij verkocht aan de zusters Ursuline van Venray, maar aangezien de mannen niet durfden te komen plukken vanwege die razzia allemaal, hebben de nonnen zelf, de zusters zelf, kwamen de appels plukken,

Figure 1c

I: *eh, eh*, ik wilde u eens vragen: wat vindt u van deze werkplek?

R: *ja*, wat vind ik van deze werkplek... dat is een moeilijke vraag, *da da da* daar moet ik even over nadenken. Ik vind het een hele aardige omgeving, *ja, een hele aardige omgeving*. Leuke collega's, inhoudelijk interessant werk, enzo. Het gebouw is wat minder, maar ja, je kunt niet alles hebben in het leven, hé?

I: kunt u daar wat dieper op in gaan, wat bedoelt u met leuke collega's?

R: wat ik met leuke collega's bedoel?

Figure 1d

la bo'rea e il fa'vònio ||
 una 'vòlta ala bo'rea 'vene 'vòlta di 'prender mari'θo || aŋ'do dal fa'vònio e 'li 'dis'e ||
 'vwo(i) 'eser(e) [il] 'mi(o) 'spozo || il fa'vònio 'er(a) uŋ 'ti'fò atakka:θo a(i) hwa'trimi
 e le 'dome non 'l aŋ'darvan(o) a 'dʒemio || le 'dis'e || 'no | per'ke non 'a'i ne'ank(e) uŋ 'sòldo
 di 'dòθe || la bo'rea | 'punta sul 'virvo | si 'miz(e) a so'fjare hon 'tute le sur(e) 'forse || so'fjo
 per 'tre 'dʒomi | e 'nevi'ho 'fatto 'fito || 'kwand(o) 'ebe finit'θo di 'stender(e) il su(o) ar'dʒento
 (i)ŋ'tomo | 'dis'e || 'ek:θi la 'mi(a) 'dote | 'tu kie di'f'evi he non 'tje l'v || e aŋ'do
 a ri'fò'sarsi delta fa'θi'ha || il fa'vònio skro'l:θe 'spate | e si 'miz(e) a so'fjare 'lui ||
 la ham'pagna e (i) 'monti restarono 'sot(o) uŋ 'fja'θo 'haldo he 'j:θse fin l' 'ult'imo 'fjoko
 di 'neve || la bo'rea | ri'fò'sa'fasi per 'beme | 'vide he 'di'eta 'dòθe non restarva 'fju
 'nula || dov' e aŋ'da'θa 'θuta la θur(a) 'dòθe | la han'son'o il fa'vònio || iŋ'soma | mi 'vwo(i)
 aŋ'kora per mari'θo || la bo'rea 'i(?) 'ris'pose || 'no | nom' v'ore(i) 'ma(i) 'esere [la] θur(a)
 'spozo | per'ke (i) n un 'dʒomo 'se(i) ha'fa:ʃe di man'darmi ŋ 'fumo 'θuta la 'dòθe ||

Source: Scagiola and Calamai (2017) from a workshop on tools for oral history, Arezzo

One thing is clear: a generally similar approach to the formal practices of archiving data can be taken, regardless of the flavour of data. This does not mean putting all research methods and their outputs into one box, but it means that we can utilise a common underlying approach. And, we can think about different levels and types of documentation depending on the likely needs of future users. Again, this does not imply we must fit the description of data into one box, but we can consider different ‘layers’ to help reveal the context of raw data. Finally, we can treat and control data in a variety of ways to make sure we are being legally and ethically compliant, and to ensure that ‘data’ are only shared, where they can be shared using appropriate pathways for access.

3. Qualitative data sharing in the UK: The UK Data Archive

My own organization, the UK Data Archive, is a department at the University of Essex and was established in 1967 by the Economic and Social Research Council (ESRC) as a data bank for social science. The Essex Archive brings with it 50 years of experience in curating and providing access to data and currently directs the national flagship service for data in the UK – the UK Data Service (UKDS). The Service supports the curation, long-term preservation, and access to data and supports the re-use of data for research and teaching. It specialises in social surveys, historical databases and qualitative research data and increasingly, bio-medical studies and data from real-time streaming devices.

The ESRC was an early adopter of data archiving and was the first major research funder to develop and implement an overarching research data policy across the spectrum of social sciences (ESRC 2015). Data Management Plans were later formally required as part of funding. The UKDS works as the organisation to operate the policy, often acting as the ‘Data Police’ when it comes to monitoring and liaising with researchers about data sharing on completion of their funded award. For qualitative researchers, discussions can be challenging, for example, when maybe they forgot to read their contracts or they were unaware that the policy also covers qualitative data! Or they did not seek consent to share outputs, despite setting out plans to share data in their submitted data management plans. Each time researchers are reminded by ESRC about their contractual obligations and to set out in advance in their application and DMP any problems foreseen in sharing data.

The data teams at UKDS have built up years of experience and awareness of the range of issues confronted by researchers and their data regarding the reality of sharing their outputs. Great emphasis is placed on advocacy, support and training for managing and sharing data, to ensure it becomes more familiar across all disciplines. The UK has had almost 25 years to embed principles and everyday practices of data sharing, with the methods literature embracing arguments for and against data sharing

and secondary analysis. For Germany, the state of enculturation for sharing and re-using qualitative data is still in its infancy, but it can learn from the UK, and feel relieved to hear that hostility, confusion and protest are not uncommon emotions to hear as research communities adjust to new practices. Indeed, UK anthropologists sought a formal exemption from the Research Council to be exempted from the ESRC Research Data Policy, but it was (thankfully) turned down.

At the UKDS, the 7000 plus data collections are available for research, learning and teaching and are used across many sectors and by many different disciplines across the world. But, first we must rewind back to the origins of UK qualitative data sharing. *Qualidata* was set up in 1994, following a small grant awarded to Paul Thompson, a Professor of Sociology at Essex University who set up the National Life Story Collection (NLSC) at the British Library, who decided that he wanted to find out the fate of the qualitative data created from well-known qualitative investigations in the UK; and how the investigators felt about sharing these data, if, indeed, they still existed! On conducting a survey and establishing that some great material did still exist, but some famous studies having been destroyed, he bid for a Search and Rescue grant, and was awarded 5 years of funding from the ESRC to catalogue the whereabouts of data arising from qualitative research grants, funded by the ESRC, and to try to rescue them and put them in accessible places.

Ironically, the UK Data Archive, also based at Essex since 1967, was housed literally 300 metres across the square, but was not interested in dealing with qualitative data at that time. The recruitment in *Qualidata* of a traditional archivist, a survey methodologist and a Historical Sociology Professor, enabled them to come together to develop a fused approach across the survey data publishing, oral history and traditional (paper based) archiving domains to create a new systematic approach to preparing, documenting and cataloguing qualitative data. Guinea pig studies were chosen, representing fieldwork from some of the UK Pioneers of Sociology (Thompson 2017). With *Qualidata* emerged a new culture of preserving and re-using qualitative data by defining approaches and methods for preparing and documenting data. In turn this work spawned a new (critical) literature on secondary analysis of qualitative data from the mid-1990s, encouraging those against sharing data to come out and demand that the qualitative research community also share its concerns. The debate slowly made its way into key methods literature, for example, Silverman (2016). The main concerns are briefly addressed in the following section, and have helped shape and refine robust methods for archiving and curating data. The *Qualidata* approach has been sought as a model for many embryonic archives across the world, and while some have flourished then faltered (often because the message has not been delivered in a convincing way, or that it has not been grasped as business-critical by funders), we should congratulate the Bremen QualiService, especially Andreas Witzel (and colleagues of the UKDS since 1998!), which will be applying for accreditation as a research data centre (RDC) by the RatSWD in the near future.

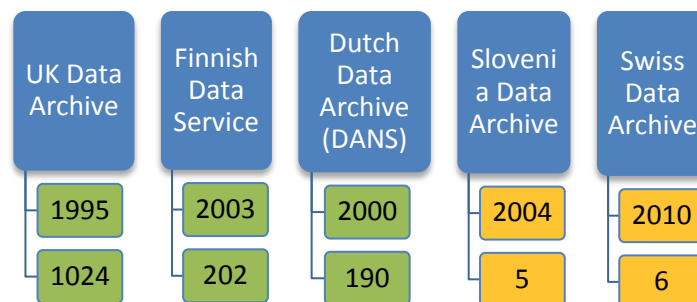
The centre, which reached a maximum of seven staff at its peak in 1998, gained funding until 2000, when money for smaller data initiatives ran out. It did continue its life by joining the UK Data Archive in 2000, when Louise Corti took a new job as Director of User Services, and brought the *Qualidata* unit she has been codirecting with her. It was fully incorporated into business in 2000, along with the separate (funding also running out) History Data Service. A year was spent enculturing these staff in new ways of working and, equally, helping staff in the UK Data Archive to accept this unfamiliar 'new kind' of data, which they did. Where integration is working well at home, then promoting positive messages to the research community is easier. Today, the UKDS has around 60 staff in total, with substantial benefits to sharing infrastructure. While many of the basic archiving processes are shared, such as acquisition and licensing, access pathways and cataloguing and promotion, a small

number of qualitative research specialists deal with evaluating data, ethical issues, preparing and documenting data, and training in re-use of data.

4. The European Landscape

While the UK is very well resourced with 1000+ published qualitative data collections, other European countries lack this resource. The UK funding bodies have continued to support infrastructure for qualitative data, which include the continuing support for the UK Data Service and one-off funding for a longitudinal data repository at the University of Leeds (University of Leeds 2017). In 2018, very few of the main social science data services are actively ingesting qualitative data; while some countries have had smaller pilots as noted above, others have built more lightweight self-deposit repositories that are able to take a greater volume of data to include qualitative data. Figure 2 shows a snapshot of the number of individual qualitative data collections in national data repositories at May 2018. All of these archives are part of a long running European Consortium of Social Science Data Archives that is now a more formally supported European infrastructure (CESSDA ERIC 2018). CESSDA aims to provide large-scale, integrated and sustainable data services to the social sciences and bring together archives to share expert guidance on matters like repository certification, data policy, data management and data discovery and persistent identifiers.

Figure 2: Qualitative data collection in national repositories



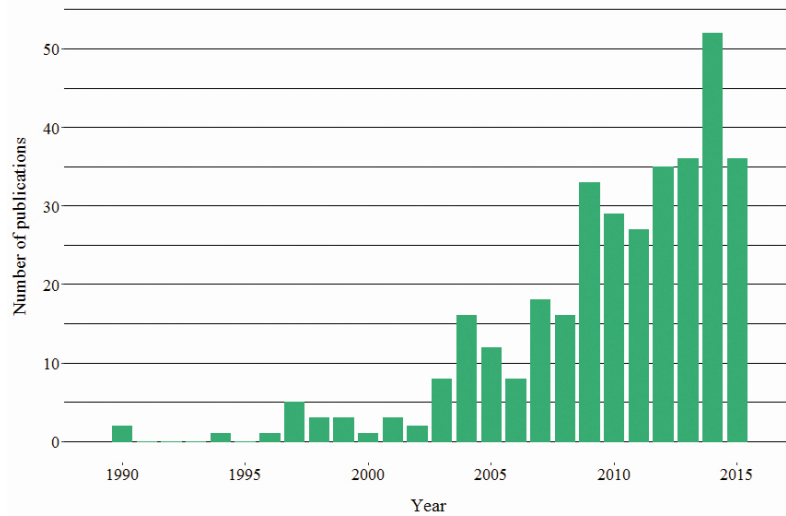
Source: Own illustration based on desk research

Legend: Numbers in the first row under the name of the archive indicate the respective founding year of the archive. The bottom row provides the number of individual qualitative data collections in the archive as of May 2018.

5. Reuse of qualitative data

To examine demand for archived data for re-use, Bishop and Kuula-Luumi (2017) looked at measure of popularity of data re-use over a 25-year period (1990–2015) using published articles in Thomson–Reuters Web of Science citation metrics portal (webofknowledge.com), across all disciplines, that had ‘re-used’ qualitative data in their analyses or discussed the technique. The number of publications identified was almost zero from the period 1990 until 1997 when five items were published (Figure 3) and rose to over 50 per year by 2014, and reaching reach a total of 347 over the 25-year period.

Figure 3: Mentions of re-use of qualitative data in the citation literature (1990–2015)

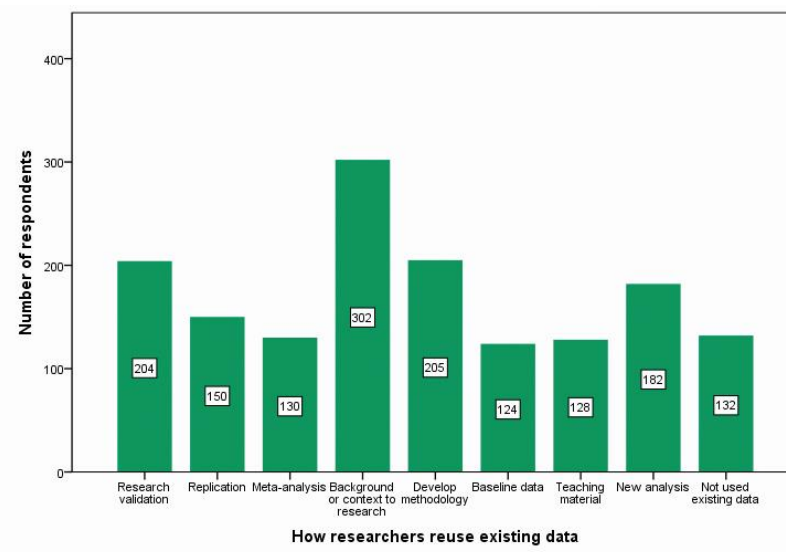


Source: Bishop and Kuula-Luumi (2017)

The case of two archives demonstrates how data are being re-used in practice: The UK Data Service and the Finnish Social Science Data Archive (FSD). The data are mostly text-based, arising from in-depth interviews, focus groups, essays, and observations; and the scale of studies varies hugely. Accompanying audio material is often not offered, due to disclosure risk, but there are sometimes related materials, such as images. Where raw data from a study has not been archived, there are often rich documentation about the conduct of the research itself and fieldwork.

Building on types of use, originally classified by Corti and Thompson (2004), they fall into eight broad categories: providing background description and historical context: comparative research, restudy or follow-up; secondary analysis; text mining and natural language processing; replication of published work; research design and methodological advancement; and teaching and learning. A survey carried out in 2016 of UK holders of Wellcome and ESRC research grants supports this pattern of use (van den Eynden et al. 2016), as shown in Figure 4.

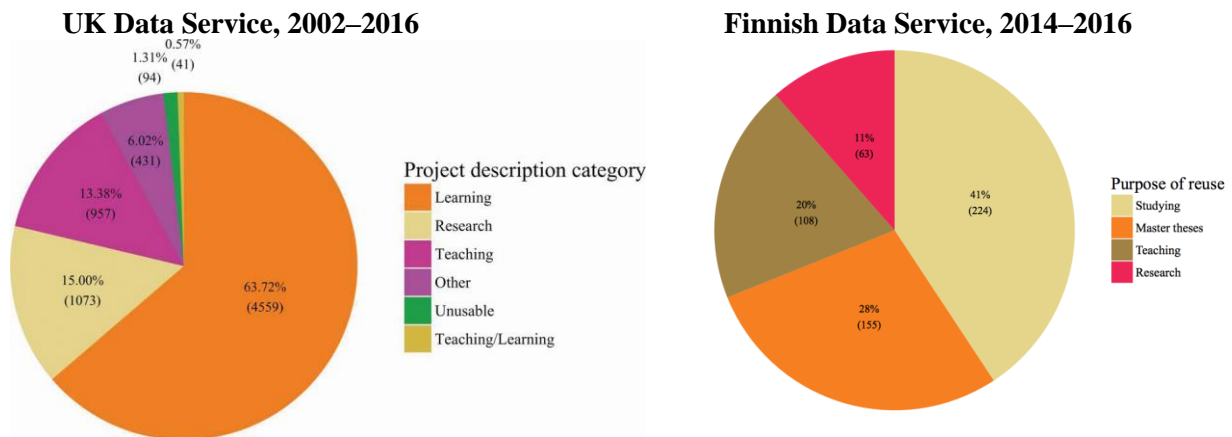
Figure 4: Reason for re-using data



Source: van den Eynden et al. (2016)

20 years of evidence, reflected in Bishop and Kuula-Luumi's (2017) paper on re-use of data at the UKDS and FSD suggests that people are undertaking secondary analysis of qualitative collections rather than replication, and there is a great deal of use for teaching and learning across all instructional levels, from undergraduate to postgraduate; indeed, over two thirds of total use (Figure 5).

Figure 5: Re-use purposes of qualitative data in two social science data archives



Source: Bishop and Kuula-Luumi (2017)

Capturing the methodological perspectives and details under which studies are undertaken, provides great value for new users who may be unfamiliar with the raw data, but also as teaching exemplars of showing various methods and fieldwork approaches. Evidence from Finland suggests that when data are especially rich, users do not demand much documentation about the study itself. Finally some topics seem to be popular, such as health, family practices and work. Crow and Edwards (2012) present a range of researchers' own methodological challenges when using archives. Further, it can be enlightening to read other's experiences of re-using data, for example, those captured in case studies of re-use in research and teaching published by the UK Data Service (2018).

There isn't room here to debate some of the more challenging epistemological and practical problems when re-using qualitative data. The sociology literature in the UK was the first to encounter critical debate around the role of research fit and context. Some of the real and perceived problems of re-using data, are uncovered by Corti and Thompson in their original early 2004 chapter on Secondary Analysis of Qualitative Data, and relate to: concerns about ethical re-use of data; not enough data available to analyse; unfamiliarity with the secondary analysis method; 'enough' contextual information; biased selection and selective sampling in secondary analysis.

6. The practicalities: preparing qualitative data

The role of a data archive includes the acquisition, preservation, documentation, cataloguing, anonymisation, and dissemination of data, to prepare a collection for future use. The current 'Qualidata' systemic approach to combining traditional archiving practice with survey data archiving protocols seeks to anticipate what future users are likely to need. Whatever one presents, a user will still want to review and assess the provenance of the data and identify any limitations or assumptions for their proposed analysis. What ends up in an archive, is in part driven by what the donor has to or is prepared to offer and what the archive wishes to take. Thus it is hard to assess what might be missing as part of a donation, or what an archivist might have weeded out. Archiving can never be a 'scientifically-neutral activity', nor are data passive resources to be mined by researchers. Indeed, archiving can be even viewed as a value-laden political practice, with what and how material gets archived being influenced by power.

One of the key roles then of the data archivist is to appraise data using criteria of research value, usability, formats, and of course, ethical and legal. Where collections contain detailed documentation about the research process such as administrative documents, grant proposals and final reports, this can help appraise the completeness of the material. Once a collection is accepted, and legal issues around ownership and suitable access conditions are agreed, various ‘data processing’ activities can begin. The UKDS assigns priorities to collections, so that the inflow of data meets the resources available for preparation.

Despite which data formats come in, users typically want data in a user-friendly format, such as a common word-processing package for text; and long-term preservation requires formats that can be accessed by researchers, now and in the future. Data archives undertake various data curation activities which include checking and validation of the collection’s contents, ensuring that consent and confidentiality agreements are met; anonymising direct identifiers where needed; and storing data for the longer term. In term of format, clear speech demarcation and the use of speaker tags are crucial when transcribing interviews and focus groups (e.g. transcription template, UK Data Service 2017).

Research should conform to ethical and legal requirements especially with respect to data protection. Informed consent is required to participate in research and to specify what kinds of uses are to be made of the information is also necessary. Many qualitative researchers have been successful in adding clauses for future data sharing into their consent forms. Where personal information might be passed on, data protection issues considerations must be built in – including that identifying information be removed or pseudonymised, or the participant must agree (e.g. oral history testimony). Where (legally protected) personal information needs to be anonymised, pseudonyms or generic descriptors should be used to edit identifying information, rather than blanking-out information, follow any article publication strategy, and be indicated in the text. Pre-planning and agreeing with participants during the consent process, on what may and may not be recorded or transcribed, can be a much more effective way of creating data that accurately represents the research process and the contribution of participants. Consideration should be given to the level of anonymity required to meet the needs agreed during the informed consent process. Obtaining informed consent for data sharing or regulating access to data should also be considered together with anonymization.

At the UKDS, documenting the data collection to enable informed use requires that information is collated about the study, methods, questionnaires and data is compiled into a User Guide. It is useful to think about levels of context that can help us document data. For example, in addition to situational features, context can also include factors resulting from everyday interaction and interpretive processes. ‘Data-level’ documentation includes information about participants/settings and the data files collected. File-level attributes can be detailed in a number of ways, including providing a summary at the top of a text file (e.g. an interview transcript), attaching descriptive information to a separate file, or setting out file-level information for all data files in a collection in one document. A data list, in an MS Excel spreadsheet can be used to identify: descriptive attributes or biographical characteristics of participants or entities studied, such as: age, gender, occupation or location; and identifying details of the data items, such as file name, description, file format and size; connections between related files e.g. audio, images and so on; and indicate where parts of the data might be missing, such as partial transcripts or those completely missing from a collection. See for example a data list from the study, *Being a Doctor* (Nettleton 2009).

What an archive can receive and store in terms of annotated data depends on the software in which annotations/coding have been created. CAQDAS packages, like NVivo and Atlas-ti, are proprietary and not all the value-added work can be exported. But they can aid with data description tasks as they

contain features to help create attributes about sources used, interview settings and interviewees, pseudonyms used, and the final coding list and analytic memos that can also be exported and archived with the data. Finally, it is essential to provide an audit trail of what was done at each stage to provide full transparency. Many repositories use a read file for users, for example: <http://doc.ukdataservice.ac.uk/doc/2000/read2000.htm>.¹ This should include information on Privacy Impact Assessment, legal agreements on condition of use, and enhancements by the depositor or the archive.

A systematic catalogue or ‘metadata’ record is created for studies, providing an overview of the study, the size and content of the data files, availability and terms and conditions of access. The Data Documentation Initiative (DDI), is seen as the de facto standard for archiving and providing systematic resource discovery for social science collections. The UKDS uses DataCite DOIs (Digital Object Identifiers) for persistent identification of their data collections. Many data archives provide web-based delivery of data via access facilities in a secure and managed environment, where data files can be downloaded, after the appropriate level of user authentication and authorisation that has been agreed for that data collection. Under certain circumstances, sensitive and confidential data can be safeguarded by regulating or restricting access to and use of data.

Data presented digitally via user-friendly tools and interfaces, enables more efficient discovery and manipulations of large data. The UK Data Service developed a fully searchable qualitative data platform, QualiBank (UK Data Service, 2014) where users can search for key terms and explore the entirety of published texts. QualiBank also allows users to select extracts of text, such as a paragraph from an interview transcript, highlight these and cite that extract in an article with a persistent web address. The reader can view and use the quoted extract in context as a highlighted paragraph (e.g. extract: <https://discover.ukdataservice.ac.uk/QualiBank/Document/?cid=q-1dba72b1-d148-40e7-b3dc-a81ae230ca80>).² Fielding and Corti (2016) discuss how being able to persistently point to quote or extracts of qualitative data from publications is a useful feature for the reader to engage themselves ‘within’ the primary data source.

7. Sharing data from an emotionally sensitive topic

Seymour’s 2012 study of *Managing suffering at the end of life: a study of continuous deep sedation until death*³ covers discussion about delicate issues which prompt heightened ethical and legal concerns about how to ‘treat’ such data beyond the immediate needs of the research team. A tension arises between the desire to share valuable data and the need to protect these assets. While the ‘topic’ is undeniably sensitive, the data themselves need not be, at least not by the definition of common data protection legislation. Seymour discussed data sharing with relevant governance committees where a two-stage consent process was agreed: first, for participation in the research and then second, afterwards to allow for data archiving/sharing. Seymour identified the danger of clinicians as qualitative researchers being tempted to switch to their clinical roles in order to help vulnerable participants, and thus leading to a reluctance to add additional “burden” of asking for permission to share data; yet missing an opportunity for participants to feel further empowered by the increased reach of their voices. Preparing the data for the Archive was straightforward; a suitable anonymisation strategy was devised and carried out at the transcription stage, ensuring that the transcriber had signed a confidentiality agreement. A safeguarded access pathway enabled 32 anonymised transcripts from

¹ Accessed on 08.08.2018.

² Accessed on 08.08.2018.

³ See <https://www.researchcatalogue.esrc.ac.uk/grants/RES-062-23-2078/read> (accessed on 08.08.2018).