# CLS Linked Data Policy

| Document Information | |
|---|---|
| **Document Name** | CLS Linked Data Policy |
| **Author(s)** | Karen Dennison, Danielle Gomes, Aida Sanchez, Emla Fitzsimons, Lisa Calderwood |
| **Owner** | Records Linkage Managers |
| **Issue Date** | February 2023 |
| **Approved By** | CLS Records Linkage Project Board |
| **Next review** | 08/02/2024 |

This document includes data that is **PUBLIC** and can be disclosed outside UCL IOE CLS and shall not be duplicated, used, or disclosed in whole or in part for any purpose other than to evaluate and implement procedures defined within this document.

## Contents

## Scope

This document covers the acquisition, management and provision of access of Linked Data within the Centre for Longitudinal Studies (CLS). Linked Data within CLS are any data linked to CLS cohort data that are not collected or derived from a CLS data collection such as a survey question, physical sample, device from a participant or other person who has been contacted as part of the survey process. This includes both individual and non-individual data. The CLS cohort data are from the 1958 National Child Development Study (NCDS), 1970 Birth Cohort Study (BCS), Next Steps and Millennium Cohort Study (MCS)

## Rationale for Data Linkage

CLS seeks to obtain Linked Data for:

- Research and methodological purposes
    a. adding administrative and other relevant records to participants' study responses and making them available to the research community helps researchers and policymakers to build a more complete picture of participants' lives and use of healthcare, education and other services, enabling research for public benefit.
    b. CLS also has an internal programme of work to use administrative data to improve understanding of the quality of linked data, to validate participants' responses, to aid the handling of missing data and to provide the outcomes to the wider research community, to aid their use of the data.
- Operational purposes i.e. for tracing and contacting study participants, updating their records and flagging deaths and embarkations.

## Linked Data Programme

The CLS Linked Data programme comprises:

1. **Linkage of data for research and methodological purposes**
   - Individual administrative records including health, mortality (fact and cause of death), education, economic and crime
   - Non-individual geographic data (e.g. proximity to green space)
   - Non-individual administrative records (e.g. schools)
2. **Linkage of data for operational purposes only**
   - Individual records e.g. contact details, mortality (fact of death) and embarkation data from administrative records

Information about consents obtained and linkages enacted for each of the cohort studies that CLS runs can be found on the CLS website (see Appendix 5).

## Legal Basis for Processing and Sharing Linked Data

a. UCL's legal basis for processing and sharing linked personal data is GDPR Article 6(1)(d) **-** 'processing is necessary for the performance of a task carried out in the public interest or in the exercise of official authority vested in the controller'. The processing of special categories of personal data, such as health data, follows the additional condition from Article 9(2)(j), paraphrased: 'processing is necessary for scientific research purposes, subject to appropriate safeguards'. For operational (tracing) purposes, an ethics application is submitted by the Records Linkage Team to a Research Ethics Committee and a separate application is submitted

to the Confidentiality Advisory Group (CAG) for support under Section 251 of the NHS Act 2006 to set aside the common law duty of confidentiality.

b. A Data Sharing Agreement is agreed and signed with the external data provider. This must include the scope and content of the data being transferred to accurately link the cohort participant to the administrative records. The agreement also includes a list of those data items which the data holder will supply and the period of retention of personal identifiers supplied from CLS for linkage. CLS applies (jointly with a co-applicant if appropriate) to the data provider for acquisition of data, for those participants that have consented (where applicable).

## Informed consent

CLS's record linkage programme has ethical approval from HRA ethics committees, including the CLS approach to participant informed consent and for operational purposes.

1. **Data for research purposes**

a. Individual administrative data
For ethical and common law duty of confidentiality purposes, CLS obtains informed consent on an opt-in basis from participants to link their individual level administrative records to their survey responses. This involves providing participants with information about proposed data linkages, including participant information leaflets sent in advance of the interview sweep. Consent for each linkage type (e.g. health, education, economic, crime) is collected during an interview sweep. The only exception is fact and cause of death, for which consent is not used, though participants are able to object to this via a GDPR information rights request. Participant websites provide information about consents obtained and linkages enacted, extensive FAQs and privacy notices (see Appendix 5). Participants can withdraw their consent for any or all of the linkages at any time without giving a reason. Participant consents and any consent withdrawals are recorded in the CLS Cohort Maintenance Team database. Further details on consent checking during the linkage process, consent withdrawals and loss of contact can be found in Appendix 1 (items 1 and 11).

b. Non-individual data
This refers to data e.g. school level, and linkages to geographic data at dwelling or area level. Participant websites, including privacy notices and FAQs, include information about these linkages. CLS approach is that no consent is required for these linkages. Participants are able to object to this via a GDPR information rights request.

2. **Data for operational purposes**

This relates to tracing and contacting study participants, including establishing whether they have moved abroad (hence usually no longer contacted for the study) and to avoid attempting contact where they have died. CLS's approach is that no consent is required for these linkages as there is a reasonable expectation that we would make best efforts to stay in touch with study participants who have not permanently withdrawn from the study. Participant websites include information about how the studies apply for access to contact information from administrative records in order to keep in touch with participants. Participants are able to object to this via a GDPR information rights request.

## Individual-level Administrative Data Linkage for Research Purposes

CLS applies to government departments and agencies, such as the Department for Education and NHS

Digital for administrative data linkages, ensuring that appropriate Data Sharing Agreements and secure data transfer arrangements are in place. CLS securely sends uniquely identifiable information together with a bespoke identifier to the Data Provider or their Trusted Third Party for those participants who have consented. The Data Providers identify the participants in their systems, match their records and remove any direct identifiers before transferring the linked data. The Data Provider sends the linked data to CLS to store and process the data internally and, if they have agreed to onward sharing, to deposit with the UK Data Service (or other similar service) – this linkage process and data flow is detailed in the Appendices 1 and 2. In some cases the Data Provider may send the linked data direct to the relevant data service. Similar secure processes are in place for other data services - see *Data access and governance*.

## Linkage to Data at Small-Area Level for Research Purposes

CLS has a programme of linkages for data that are not about an individual but refer to an area or dwelling. These data are usually publicly available and do not require CLS to send any personal data to an external organisation. Small-area level data that CLS has linked to its cohort data include air pollution data, distance to green spaces, distance to fast food outlets (amongst others). There is the potential for further geographical linkage to a wide range of other external data such as Ordnance Survey, housing, environmental, energy efficiency of the property, broadband, and weather data. Where the data are freely and publicly available they can be linked without the requirement for specific agreements with the data owners, subject to the terms and conditions of use of the data being linked. There is a programme of work to explore which other data it would be most useful to link to in future and to plan for such linkages. Any new linkages must be approved by the DAC and are subject to disclosure control reviews before being disseminated. Information on these linkages are provided on participant websites (see Appendix 5).

## Linkage for Operational Purposes

Where CLS has lost contact with participants who took part in any of its four studies (NCDS, BCS70, Next Steps, MCS) an attempt will be made to trace these cohort members so that they can be re-contacted and invited to take part in the next survey. One of the ways CLS attempts to trace participants is by requesting contact information from other data providers such as NHS-Digital. CLS applies to these data providers to request new addresses for these participants; CLS also regularly applies to NHS-Digital to receive fact of death and embarkations (emigrations) data. The linkage process and data flow are detailed in Appendices 3 and 4.

## UCL Data Management Environment

CLS holds and manages its potentially identifiable personal data within the UCL Data Safe Haven[1] (DSH). The DSH is a computing environment separated from the rest of the UCL network and has been set up with additional technical measures aimed at increasing security by reducing risks. The DSH is certified to the international information security standard ISO27001:20131. UCL holds NHS Data Security and Protection (DSP) toolkit accreditation, reference EE133902-SLMS. All CLS staff complete mandatory IG training and all staff processing personal data must have basic disclosure checks. In addition, CLS aims to comply with any specific information governance requirements of Data Providers.

 The data is managed by the following two CLS teams:

---

[1] https://www.ucl.ac.uk/isd/services/file-storage-sharing/data-safe-haven-dsh

1. The Cohort Maintenance team who are responsible for, and have access to, study participants' contact details and other personal identifiers and also those of their relations (mainly, partners, parents), such as names, addresses, National Insurance Numbers and NHS numbers
2. The Research Data Management team who have access to the research data i.e. de-identified and pseudonymised research data collected during surveys and external data linked through administrative systems or geographical databases.

Each team has access to a logically separated DSH area. This is managed by access control, so that the two teams reside in separate access groups.

## Research Data Access and Governance

CLS Linked Data are categorised and shared in accordance with the *CLS Data Classification Scheme*, as described in the *CLS Data Classification Policy*[2]. This scheme sets out the classification of data and the appropriate levels of data security and data segmentation. The dissemination and access arrangements for data are described in the *CLS Data Access Framework*[3].

CLS is committed to ensuring that as much linked research data as possible are made available to the research community through the existing mechanisms of data dissemination.  This is a key aim of the Economic and Social Research Council (ESRC), which is the main funder of CLS, which has an international reputation for supporting high quality data for the research community, enacted primarily through the UK Data Service (UKDS). The linked data are pseudonymised, de-identified, curated and documented by the CLS data Management team for research purposes. In addition to making data available through the UKDS, CLS also deposits survey data with other data service providers, and securely transfers personal identifiers to the Trusted Third Parties used by those services, for the purposes of linkage to administrative data records held by those data services. These services include the UK Longitudinal Linkage Collaboration, (UKLLC) the Secure Anonymised Information Linkage Databank (SAIL) and the Office for National Statistics Secure Research Service (ONS SRS). UKDS, UKLLC, SAIL and ONS SRS are all Trusted Research Environments (TREs) with robust policies, protocols and safeguards for secure remote desktop access to data, accreditation to relevant international standards and adoption of the Five Safes Framework. The latter is a set of principles adopted by a range of TREs, which researchers and their organisations must adhere to and which provide complete assurance for data providers.[4]

For data not currently available through the above-mentioned data services, applicants need to apply for approval from the CLS Data Access Committee (DAC). The CLS DAC determines the most suitable onward sharing arrangement, which will usually be through the UCL Data Safe Haven. These applications may include methodological research projects to access linked data prior to their deposit through the UKDS. Researchers accessing data through the UCL DSH need to complete NHS Digital's Data Security Awareness course and a UCL DSH Acceptable Use Statement. Researchers are required to follow Statisical Disclosure Control guidelines and a member of the CLS RDM team check outputs for disclosure control.

## Data Protection Impact Assessments

CLS has completed separate UCL Data Protection Impact Assessments (DPIAs) for its administrative data

---

linkage programme for research. A DPIA is in development for geographical linkages and operational linkages. These DPIAs are periodically reviewed, including for any new data linkages, to ensure that any new risks are identified and mitigated. CLS may also complete separate DPIAs for specific linkages and the requirement for this is assessed on a case-by-case basis.

## Appendix 1 Linkage Process for Individual Administrative Data for Research Purposes

1. **Consent checking**

   CLS has sought informed consent for each type of data linkage from the participants of its studies and participants can withdraw this consent at any time. Any proposed data linkage is only conducted on those participants for whom consent has been given. Full details of the consent collection procedure and the numbers of individuals who have agreed are the basis upon which identifiable information is passed to the Data Provider to enact participant linkage.

2. **Participant matching**

   a. The CLS Cohort Maintenance Team securely send uniquely identifiable information to the Data Provider or their Trusted Third Party for those participants who have consented. This information enables the Data Provider to identify the participants in the Data Providers' systems. In addition to names, addresses, date of birth, and sex, CLS might send information such as National Insurance Number, NHS Number or institutional identifiers such as a school or GP. This information does not include any survey responses.

   b. The identifiable information is sent with a proxy serial ID for the purposes of the individual matching only. This proxy serial ID is generated by the Research Data Management Team and passed to the Cohort Maintenance Team together with an internal identifier shared by both teams. The identifiable information is sent to the Data Provider with ONLY the proxy serial ID. Please see Appendix 2- *Proxy Serial ID generation data flow diagram* for further details.

   c. Resolution of data matching queries between the Data Provider and the CLS Cohort Maintenance Team.

   d. Agreement on what constitutes a valid match, e.g. name, NHS number, date of birth, etc.

3. **Data linkage and extraction**

   a. Following the agreed criteria for matching individuals, the Data Provider is responsible for linking and extracting the relevant data related to that individual.

   b. The content of what needs to be extracted should be agreed in advance of the linkage process between CLS and the Data Provider and specified in the Data Sharing Agreement between the two parties.

   c. The Data Provider utilises the proxy serial ID to identify the linked data and removes all of the personal identifiers used for linkage (names, addresses, National Insurance Numbers and NHS numbers, etc), ensuring that the minimum amount of identifiable information is included in the matched dataset.

   d. The Data Provider sends the linked de-identified data to the CLS Research Data Management team in a secure manner.

e. Once CLS has checked the linked data, the Data Provider securely deletes the uniquely identifiable information provided to them by CLS. Please see Appendix 1 *CLS External Data Linkage data flow* for detailed information on the data linkage process.

## 4. Data storage

The CLS Research Data Management team stores the linked data returned from the Data Provider in the UCL DSH for de-identification, curation and data documentation.

## 5. Data enhancement and validation

Subject to the terms of the agreement with the Data Provider, CLS securely stores and manages access; carries out validation of the data received and combines the linked data with that collected from the surveys or other linked data (where appropriate).

## 6. Evaluation of disclosure risk

CLS evaluates the linked data in terms of whether the data might contain information that could re-identify an individual (disclosivity) and how damaging re-identification might be to an individual (sensitivity). Data may be modified to reduce the risk and are classified according to the *CLS Data Classification Scheme*.

## 7. Pseudonymisation

Versions of the data prepared for research purposes are 'pseudonymised' with 'Research Identifiers' used for the research survey data available from the UK Data Service. These 'Research Identifiers' are shared across different sweeps of the data to enable longitudinal research analysis. In some cases, a dataset is prepared with its own unique set of Research Identifiers so that the data cannot be linked to other data from the same study for reasons of data confidentiality.

## 8. Data deposit

Pseudonymised research data are deposited with the UK Data Service, or other similar service, under a licence agreement[5].

## 9. Data transfers
Data transfers are always encrypted in transit (where possible) and logged. For transfers between CLS and the Data Provider, this may be done using the UCL DSH secure transfer system or the Data Provider 's own secure FTP. For transfers of data to UKDS, these are securely uploaded as an encrypted archive via the University of Essex ZendTo portal.

## 10. Data access
a. Researchers apply to access linked data via the UKDS through a rigorous application process, including completion of application forms, user agreements that require institutional countersignatures, and relevant training. Applicants need to detail their intended use of the data and the public benefit of their research project (or benefit to health and social care for data from NHS Digital). Applicants must abide by the terms and conditions of the UKDS. In certain circumstances, where required by the Data Provider, researchers may also need to sign an agreement with CLS.

---

[5] https://ukdataservice.ac.uk/app/uploads/licenceform.pdf

b. The application must be approved by CLS and, by the Data Provider if they require this.

c. Researchers access the data securely and remotely via their own institutional or home computer (researchers can apply for temporary home access) or at a Safe Room or Safe Pod. Their access to data is limited to the data they have requested and removed upon the expiry date of their research project.

d. Data Providers may have different requirements around data minimisation which CLS adheres to. For example, researchers requesting NHS Digital health data are required to specify the particular linked administrative variables for their research project. In that case, CLS generates the necessary syntax and provides this to the UKDS. The UKDS then uses this syntax to give the researchers secure remote desktop access to the subset of required linked variables.

e. Researchers inform UKDS of any outputs for publication who check them for Statistical Disclosure Control (SDC) within the secure environment to ensure that the published results do not reveal the identity, or contain any confidential information, about a data subject. If they are safe for publication, UKDS removes the outputs from the secure environment and sends them to the researcher. Some Data Providers may require sight of draft publications and / or may require the UKDS to apply additional SDC requirements.

## 11. Consent and study withdrawals and loss of contact

If a respondent withdraws their consent for data linkage, or asks for their contact details to be deleted, if linkage has not already happened no linkage will be conducted. If the data have already been linked and made available for research purposes, no linked records will be updated. Any existing linked data will continue to be shared securely for research purposes, unless the agreement with the Data Provider does not permit this or the participant objects. If a respondent requests that all of their study data be deleted, data (including linked data) will deleted from the UKDS (or other similar service) within six months of receipt of a request. If a respondent withdraws from the study or we lose contact with them, but does not explicitly withdraw their consent for linkages or request that all of their study data to be deleted, CLS will, in the course of updating linked records of study participants, continue to link their data. More details can be found in the *CLS guidance on study consent withdrawal and right to erasure (data deletion) requests*.

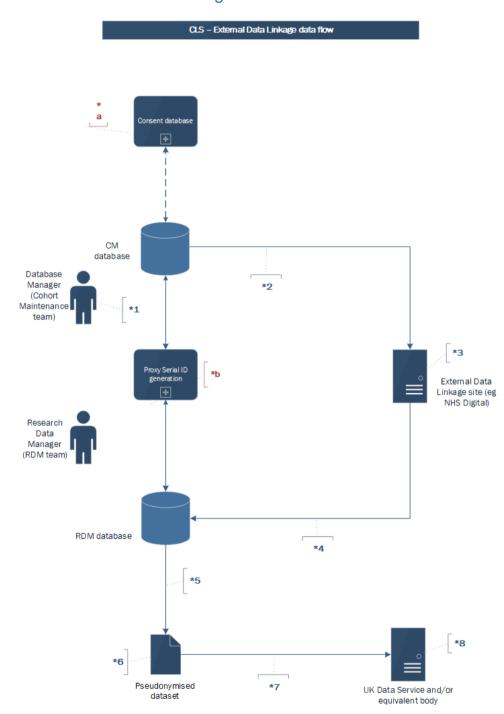## 12. Subject access requests

If study participants request access to their linked administrative records, such requests will be handled according to the agreement in place with the Data Provider. Depending upon what this agreement is, CLS may ask the participants to contact the relevant data controller to request their data directly or CLS may provide the participant with a copy of their linked data.

## Appendix 2 Data Flow Diagrams - Linkage of Individual Administrative Data for Research Purposes

1. ### External Data Linkage data flow



**CLS – External Data Linkage data flow**

**Process Notes**

**\*1** - The Database Manager in the Cohort Maintenance (CM) team sends the CLS IDs of the cohort members to be included in the data linkage to the RDM team. These are limited to those cohort members for whom CLS has consent. The RDM team generates Proxy Serial IDs to send to the External Data Linkage site in place of the CLS IDs.

**\*2** – The CM team compiles the personal data of the cohort members (eg. name, address, DOB, etc.) plus the RDM team-generated Proxy Serial IDs.

**\*3** – Data transfers are encrypted in transit and logged. This is done using the UCL Data Safe Haven secure transfer system or the Data Provider's own secure FTP.

**\*4** – Matched data file returned securely with relevant requested data (eg. NHS Digital variables) plus Proxy Serial IDs.

**\*5** – Assigned staff solely within the RDM team ingest the returned data into the RDM database and prepare the data for deposit with the UK Data Service (or other similar service).

**\*6** – The pseudo-anonymised dataset for deposit contains Research IDs which are different from the CLS IDs and the Proxy Serial IDs.

**\*7** - Transfer of data from CLS to the UKDS has to comply with their security policies.

**\*8** – The UKDS classifies the data by its level of sensitivity : either End User Licence for low risk, Special Licence or Secure Lab, the highest level of security.

**Misc info**

**\*a** – The consent database records which cohort members have consented to administrative data linkages and is kept up-to-date to reflect consent withdrawals.

**\*b** See Proxy Serial ID generation diagram for details.

**Legal basis**

CLS's legal basis for processing and sharing linked personal data is for a public task under GDPR article 6(1)(e). CLS also processes special categories of personal data for research under GDPR article 9(2)(j).

## 2. Proxy Serial ID generation



**Proxy Serial ID generation**

Database Manager (Cohort Maintenance team)

CM database

*1

*3

RDM database (SiR)

*2

Research Data Manager (RDM team)

**Business Rule**

There is an information barrier in place between the two areas of CLS.
The Cohort Maintenance team will only have access to the cohort members personal identifiers (e.g name, address, DoB, gender, etc).
The Research Data Management team will only have access to the research data.

**Process Rules**

*1 – The Database Manager in the Cohort Maintenance (CM) team will send the *'CLS ID'* to the Database Manager in the Research Data Management (RDM) team.
*2 – The Database Manager in the RDM team will generate a *'Proxy Serial ID'* for each of the *'CLS ID'* and send both IDs back to the CM team.
*3 –The Database Manager in the CM Team compiles the personal data of the cohort (forename, middle name, surname, DoB, gender, latest oldest NHS number, latest new NHS number, address,) plus the *'Proxy Serial ID'* and sends the data to the data provider for linkage.

## Appendix 3 Linkage Process for Operational Purposes

1. **Data Flow for tracing purposes**

Where the data request is for tracing purposes, CLS requests information for those participants who have not permanently withdrawn from the study or who are not known to be ineligible (e.g. died).

- CLS prepares a file containing the CLS member IDs, names, addresses, dates of birth, sex, institutional identifier such as NHS numbers (depending upon the data provider) and sends it to the data provider.

The file sent to data providers for linkage contains either data on the full cohort, or only a subset of the cohort. This will depend on CLS's needs and the data providers' requirements and working procedures.

a) CLS sends a matching file containing participants' information

- The file contains all those cohort members which CLS intends to trace except those who have withdrawn from the study or who have died or CLS know have emigrated.
- The data provider matches the participants against the information available in their database and prepares a file containing the latest addresses, names, date of birth, sex, institutional identifier and all other variables requested, for each participant and sends it back to CLS.
- Once information is provided to CLS, the data provider deletes the file received. For regular delivery of data, some providers such as NHS Digital, may flag the cohort member in their database to facilitate the regular delivery of information and to minimise the risk of mismatches.

On receipt of the data, CLS will:
- Update their database with new addresses and flag cohort members who have died or emigrated.
- Provide new information to agency field workers where necessary.
- Provide new information to third party mailing house of the new address where necessary.
- Please see Appendix 3 **-** *External data request for tracing purposes data flow diagram* for detailed information.

2. **Data flow for Notifications of Deaths or Embarkations**

Where a data request is to receive notifications of deaths and embarkations (people reported to NHS as having moved abroad)), CLS requests information for the full cohort, except for those cohort members who have requested to withdraw from the study along with a request for their data to be deleted

- CLS prepares a file containing the CLS member IDs, names, addresses, dates of birth, sex, NHS numbers for the full cohort and sends it to NHS-Digital.
- NHS-Digital flags the full cohort in their database.
- NHS-Digital notifies CLS on a regular basis, as agreed in the contract with NHS-Digital.
- Please see *External data requests for notifications of deaths data flow diagram* below.

# Appendix 4 Data Flow Diagram for Operational Purposes

**External data request for Notifications of Deaths and Embarkations and Tracing purposes**

*1

Database Manager (Cohort Maintenance team)

CM database

*2/*2a

*3

External Data Provider site (eg NHS Digital)

*4

Interview sub-contractor/ Mailing agency

**Business Rule**

The External Data Provider site contains a secure FTP portal area to drop off data, requiring a login/ user ID and password to gain access.

**Process Rules**

***1** – The Database Manager in the Cohort Maintenance (CM) Team compiles the personal data of the cohort (forename, middle, surname, gender, DoB, address, postcode, country, 'CLS ID' and the External provider organisation ID (where available) , such as NHS number  for those cohort members who CLS would like to trace.

***2** – **Tracing** - The external data provider finds the cohort in their database and returns to CLS the available data for the cohort member (forename, middle name, surname, DoB, gender, latest oldest NHS number, latest new NHS number, address, postcode and 'CLS ID'). The data provider may also flag the cohort in their database for regular delivery of data.

**\*2a**- **Notifications of deaths and embarkations**- The external data provider flags the cohort in their database and periodically notifies CLS of any deaths or moves. Data file is returned via secure FTP portal.

**\*3** – The Database Manager in the CM team updates the CM database where the External Data Provider provided a new address or a notification of death or embarkation.

**\*4 -** The Database Manager in the CM team updates new contact details with the sub-contractor agency and mailing agency  where needed.

**Misc info**
**\*a** – The consent database records which cohort members have consented to administrative data linkages and is kept up-to-date to reflect consent withdrawals. For the tracing exercise,  the CM team will exclude from the file sent to the external data provider, all participants who have requested to withdraw from the study. For the deaths and embarkations, participants who have requested CLS to stop using their data will be excluded from the File.
**\*b -** Where cohort members are flagged in the external provider database, CLS will submit an updated file should there be new withdrawal requests.
**Legal basis**
CLS's legal basis for processing and sharing linked personal data is for a public task under GDPR article 6(1)(e). CLS also processes special categories of personal data for research under GDPR article 9(2)(j).  Section 251 approval was obtained for processing data received for tracing and notifications of deaths and embarkations.

## Appendix 5 Data linkage information on CLS and participant facing websites

Information about administrative data linkages can be found on the CLS website at the following links, under 'Study features' in the 'Linked administrative data' box -

- 1958 National Child Development Study https://cls.ucl.ac.uk/cls-studies/1958-national-child-development-study/
- 1970 British Cohort Study https://cls.ucl.ac.uk/cls-studies/1970-british-cohort-study/
- Next Steps https://cls.ucl.ac.uk/cls-studies/next-steps/
- Millennium Cohort Study https://cls.ucl.ac.uk/cls-studies/millennium-cohort-study/

Information about data linkages can be found on the participant websites within privacy notices and study FAQs –

- 1958 National Child Development Study https://ncds.info/
- 1970 British Cohort Study https://bcs70.info/
- Next Steps https://nextstepsstudy.org.uk/
- Child of the New Century https://childnc.net/