

CLOSER and the Challenge of Historical Metadata

Gemma Seabrook

CLOSER

- 5 year collaborative project
- Maximise the use, value and impact of cohort and longitudinal studies
- Based at UCL Institute of Education

CLOSER Studies

- Studies:
 - Avon Longitudinal Study of Parents and Children (ALSPAC)
 - British Cohort Study (BCS70)
 - Hertfordshire Cohort Study (HCS)
 - Life Study
 - Millennium Cohort Study (MCS)
 - National Child Development Study (NCDS 1958)
 - National Study of Health and Development (NSHD 1946)
 - Southampton Women's Survey (SWS)
 - Understanding Society

CLOSER Partners

- Other partners:
 - UK Data Archive
 - British Library
- Funders:
 - Economic and Social Research Council
 - Medical Research Council
- Also working with Colectica and Metadata Technologies

Work Streams

- Data Harmonisation demonstration projects
- Data Linkage demonstration projects
- Research Impact
- Training and Capacity Building
- **Uniform Search Platform**

Current situation

- I want to know about smoking in pregnancy and effects on income in later life
- Want to look at ALSPAC, NCDS, NSHD, BCS and predictively at MCS
- NCDS, BCS, MCS – UKDA
- ALSPAC and NSHD – contact studies



The (near) future

- Search for smoking on CLOSER
- Filter by 'pregnancy & birth' life stage and required studies
- Get questions and variables
- Suggested items – similar, nearby, in and cross study

Why isn't this already there?

- Started with a dream: the search platform
- Many have tried
- Limited success
- Technology supports it
- Not about technology – all about content
- Metadata not in a standard format
- Lack of managed processes for consistency

Content needs to be ...

- Consistent
- Comparable
- Support discovery
- Re-usable
- Support processing
- Sound familiar? Metadata standards
- DDI Lifecycle

Why DDI-Lifecycle?

- Maintained and developed
- Mature
- International reach
- Used by major archives (DDI-Codebook)
- Fits longitudinal studies well

DDI is not easy

- Large and complex
- Difficult to explain to non-experts
- Even harder to demonstrate to non-experts
- Not commonly used in medical circles

Historical DDI is even harder

- Data goes back to 1930s
- Projects go back to 1946
- “The Schleswig-Holstein question is so complicated, only three men in Europe have ever understood it. One was Prince Albert, who is dead. The second was a German professor who became mad. I am the third and I have forgotten all about it.”
— [Lord Palmerston](#)

Getting started

- Assess the scale of the problem
- Try stuff
- Don't panic
- Assess the risks and challenges

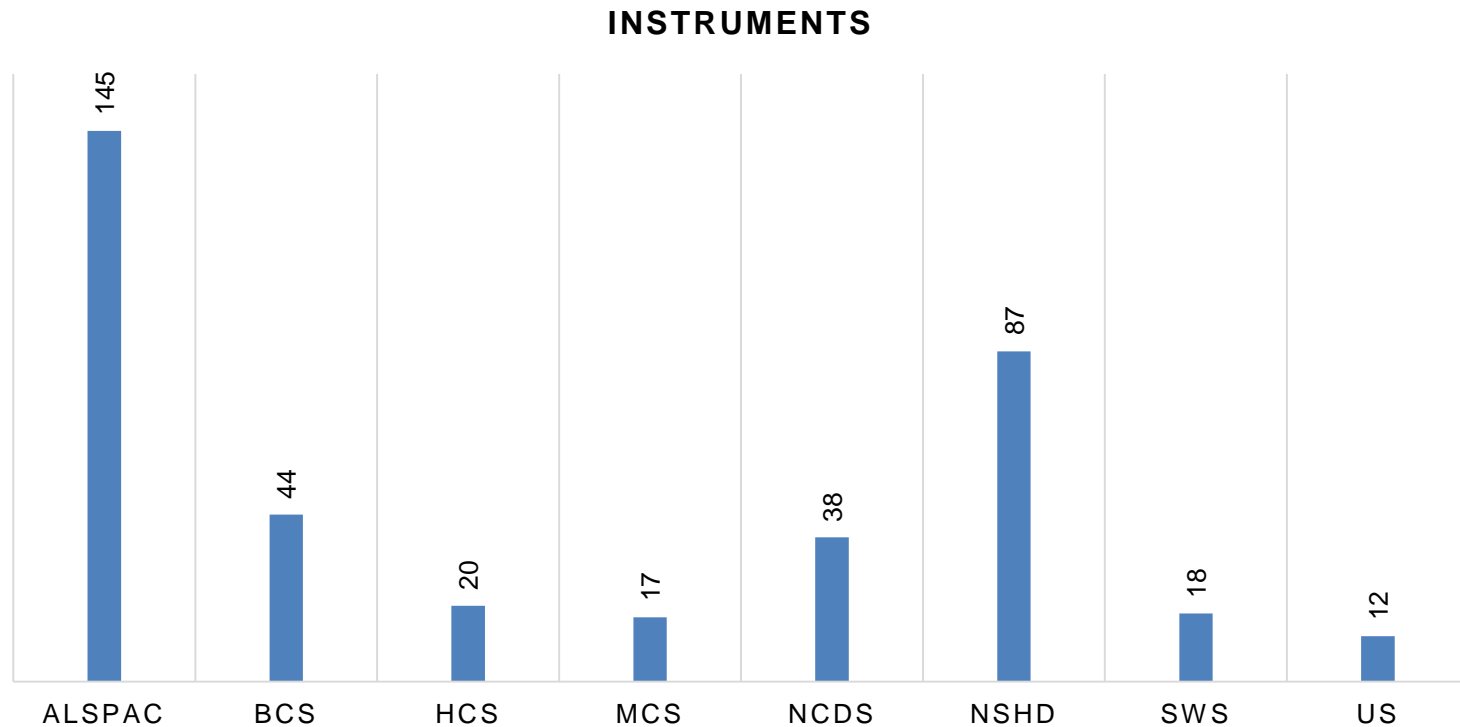


BilibioArchives/Library Archives @Flickr

Metadata Project Cake Recipe

- 500g Expert knowledge
- 250g Project management
- 250g Patience (you may need to add more as required)
- 100g Software
- 500g Funding
- **1 kg of Commitment**

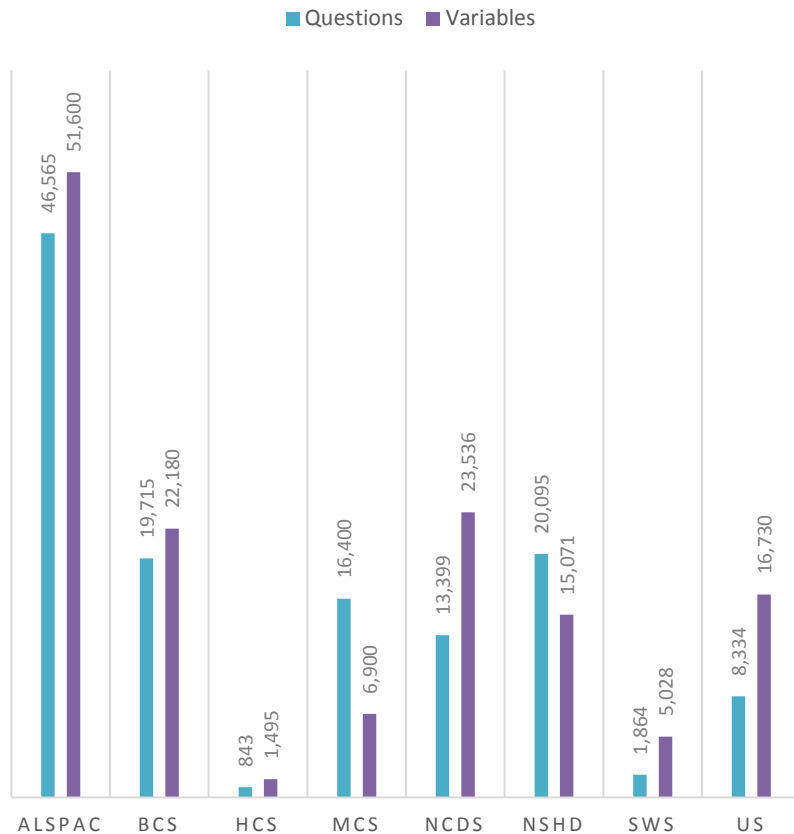
Challenges: Scale



- Instruments: 381

Challenges: Scale

QUESTIONS AND VARIABLES



- Huge:
 - >200k questions and variables
 - >300 questionnaires
- Ambitious – no one has attempted this scale before

Solution: Reduce, Reuse, Recycle

- Limited DDI profile
- Define allowable entries
- Determine content to enter
- Import from older DDI
- Import from databases, spreadsheets, etc.

Solution: Manageable chunks

- Unit of measurement: the questionnaire or the dataset
- Estimation:
 - Effort per questionnaire
 - Effort per question
 - Bonus for complexity
 - Controlled runs
 - Flexible
- Workflow

Pilot

- Had developed a tool for questionnaire entry (Caddies)
- Piloted this
- Outcomes:
 - Consistency
 - Study capacity (not about money)
 - Timelines

1. Ingest

- Manual entry
- Metadata Team
- Manual and standard operating procedures
- Covers all types:
 - Questionnaire
 - Clinical tool
 - Data extract

Principles

1. Maintain and do not alter the semantic meaning of the questionnaire
2. Do not correct the questionnaire
3. Only record what is contained within the questionnaire
4. Do not allow the data recorded (i.e. the variables) to inform the metadata archiving

Nothing is ever easy ...

b) **If no, or no partner:**

Does this older child have (please tick):

you as the natural mother (but his/her
natural father is not present)

1

answer (c) below
and then go to (e)

your partner as the natural father
(but his/her natural mother not present)

2

answer (d) on page 37
and then go to (e)

neither of his/her natural parents present

3

answer (c), (d)
and (e)

... but its still possible

If no, or no partner:

- if-true branch:

qc_D15_b

Does this older child have (please tick):

- 1 you as the natural mother (but his/her natural father is not present)
- 2 your partner as the natural father (but his/her natural mother not present)
- 3 neither of his/her natural parents present

If you as the natural mother (but his/her natural father is not present) to question D15b or neither of his/her natural parents present to question D15b

- if-true branch:

qc_D15_c

How often do you or your partner talk to the child's natural father about this older child?

- 1 once a month or more
- 2 less than once a month
- 3 once a year or less
- 4 never
- 9 don't know
- 7 natural father is dead

If your partner as the natural father (but his/her natural mother not present) to question D15b or neither of his/her natural parents present to question D15b

- if-true branch:

qc_D15_d

How often do you or your partner talk to this older child's natural mother about the child?

- 1 once a month or more
- 2 less than once a month
- 3 once a year or less
- 4 never
- 9 don't know
- 7 natural mother is dead

qc_D15_e

Are your relations with this older child's other parent(s) :

- 1 generally warm and friendly
- 2 sometimes friendly
- 3 polite
- 4 distant
- 5 usually unfriendly
- 6 no relationship
- 7 child's other parent is dead

2. Verification

- Still within the Metadata Team
- Checks for accuracy
- Also checks for consistency

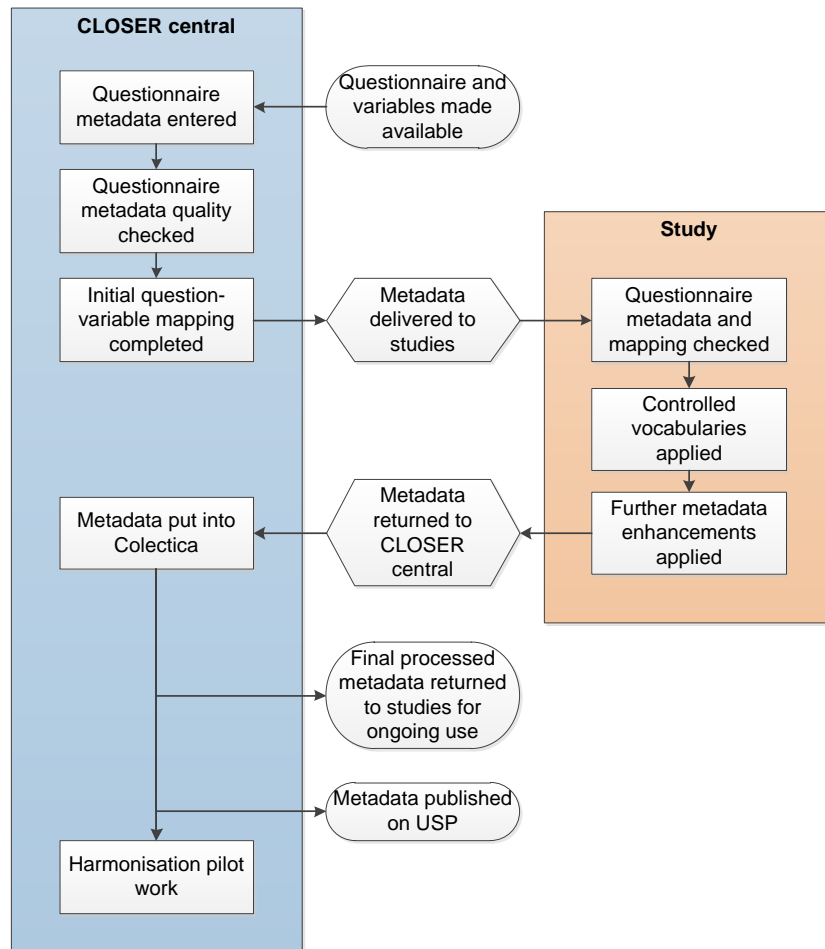
3. Mappings

- Map all questionnaires against variables
- E.g. q1 is mapped to variable N123
- Process also accommodates derived variables
- Topic mappings

4. Repository ingest

- Returned to CLOSER
- Put through systematic checks before entry
- Queries from checks addressed
- Entered into Colectica repository
- Made available in USP at this point

Process



- Process bound – each step relies on those prior
- Efficiency/consistency of CLOSER central and subject expertise of studies

Challenges: varied material

- Age of questionnaires
- Condition of questionnaires and data
- Variation in methodology – PAPI, CAPI, CATI, etc.
- Medical data, images, videos, biosamples, etc.

Challenges: Intent

- What was the intent of a question?
- Layout
- Machine logic
- Limitations of technology
- Assumptions & implied logic

BCS 1970

11 Smoking: *AC043B*

a) Does the mother smoke now?

Yes1

No2

Not known.....0

(66)

b) If **NO** to 11a, did she ever smoke?

Yes1

No2

Not known.....0

(67)

c) If **YES** to 11b how long ago did she stop?

..... years months

(68-70)

d) How much does/did she smoke?

Cigarettes per day

1-4.....1

5-14.....2

15-24.....3

25 or more.....4

Not known.....0

(71)

e) Has she smoked during this pregnancy?

Yes1

No2

Not known.....0

(72)

Challenges: Computers

- No machine logic
- No machine readable copies
- Limited variable creation

- Then, when computers arrive ...
- **MILLIONS** of questions

MCS Birth

PDOB (PDBD, PDBM, PDBY)³

^And now yourself. What is your date of birth?

DATE

CHECK HH9

^And now yourself. What is your date of birth? = 'And now yourself. What is your date of birth?' if referring to first person in household (i.e. the person answering household module)

'What is ^name given at PNAME's date of birth?' if referring to second and subsequent members of household.

IF date of birth is don't know and person being asked about is not respondent [PDOB=DK AND PNO>1]

|

| **PAGE**

| Do you know ^name given at PNAME's age last birthday?

| ENTER AGE OR <Ctrl+K> FOR DONT KNOW

| Range: 0..997

|

ELSEIF date of birth is given

|

| **PAGE** [computed using INTD and PDOB]

| Range: 0..997

|

ENDIF

Challenges: Documentation

- Studies changed hands
- Formats
- What was considered important
- Changing ideas of data management

NCDS Guide from 1969

<u>Question Number</u>	<u>Field Description</u>	<u>Computer Item</u>	<u>Explanation, and Computer Codes Associated with Permitted Punches</u>
Q 6	Col 26 Polylog SIC 1671 Type of school attended by child	825 Coded	For schools not maintained by a Local Education Authority: Types: (1) = 1 Independent School (including grant-aided) catering wholly or mainly for children who are not handicapped (2) = 2 Day Special School for Handicapped Children (3) = 3 Residential Special School (4) = 4 Other (5) = 5 Approved School (6) = 6 Controlled School (7) = 7 Junior Training Centre, or SSN provision (DNA) = 9 Not applicable - see Col 25 (NA) = b No answer

Challenges: Versions

- Gender
- Country
- In-field fixes
- Attempts to increase response rate
- Change of policy

Challenges: consistency

- Larger the group, harder consistency is
- Agree standards
- Manual
- How long do you debate a point?

Extracts from the issue log

- When to use a grid?
- When is a scale not a scale?
- Images as questions
- Images as response domains
- Character set encoding
- Are corner labels on grids cosmetic or semantically significant?



The Madness of IDs

- What is the purpose of an ID?
 - Does an ID have meaning?
 - How do you ensure IDs are unique?
 - Using IDs as references
-
- When we got to the portal interface, it all changed!

Solution: Communication

- Meet regularly
- Group size
 - Who needs to contribute
 - Consequences for other parts of the project

What happens after CLOSER?

- Documented the metadata
- Created the portal
- Made it available
- What next?
- Need sustainable processes to produce metadata
- Business as usual

Sustainable futures

- Flexible
 - Software independent
- Sustainable
 - Involvement with standard development
 - Survey agencies
 - Produce once, use multiple times
- Documentation and training

More than just metadata

- Publishing the protocols
- Publishing estimate methods
- Making CADDIES available
- Other training materials?
- Enhancing Colectica search UI
- Demonstration projects
- Future improvements

Good practices

- Controlled vocabularies
- Tracking software
- Methods for ensuring consistency
- Approval process
- Flexibility – one size will not fit all
- Rigidity – where you can't be flexible

The integrated future

- Survey agencies
 - Make DDI a contractual requirement
 - Improved exchange of information
- Data managers
 - Use DDI as the backbone for your data
 - Metadata APIs