# The Longitudinal Impossible Dataset: Helping Users Navigate the ONS Longitudinal Study

Andreas Mastrosavvas

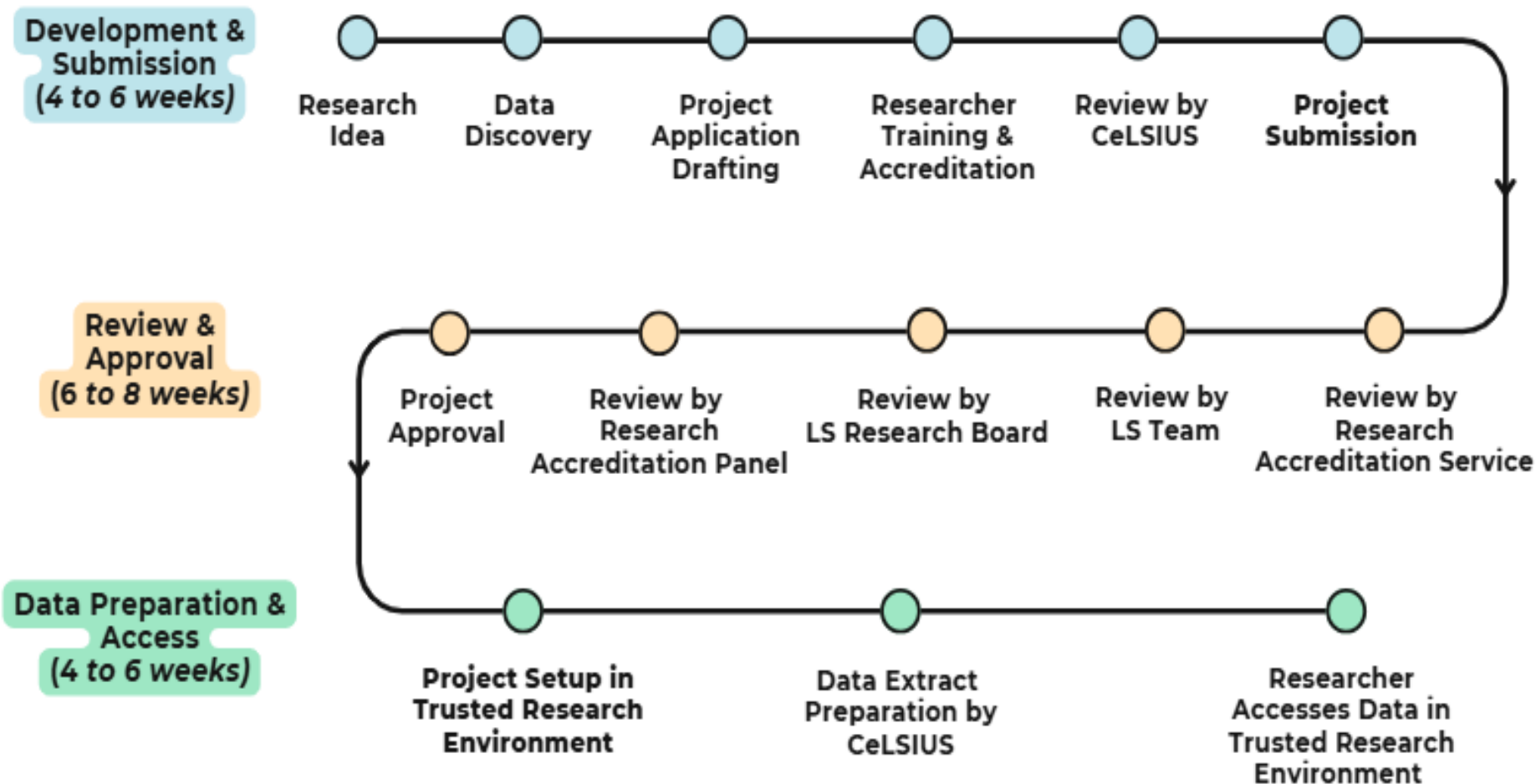Nicola Shelton

Epidemiology and Public Health, University College London (UCL)

# Accessing the ONS LS



**Development & Submission (4 to 6 weeks)**
- Research Idea
- Data Discovery
- Project Application Drafting
- Researcher Training & Accreditation
- Review by CeLSIUS
- Project Submission

**Review & Approval (6 to 8 weeks)**
- Project Approval
- Review by Research Accreditation Panel
- Review by LS Research Board
- Review by LS Team
- Review by Research Accreditation Service

**Data Preparation & Access (4 to 6 weeks)**
- Project Setup in Trusted Research Environment
- Data Extract Preparation by CeLSIUS
- Researcher Accesses Data in Trusted Research Environment

## Accessing the ONS LS – common sources of delay

- Data discovery is complex

  – Over 600 LS variables per Census

  – Researchers need to think carefully about variable selection pre-application

  – Researchers sometimes find they need to re-apply for additional variables

- Researchers unable to develop code from the start of their project

# The CeLSIUS data dictionary

# The CeLSIUS data dictionary

- Can be challenging to navigate, esp. for new users

- Variable metadata can be unstructured, esp. for linked data

| AHACC | |
|---|---|
| **Table** | CANC Cancer registrations |
| **Short description** | Area Health Authority within Regional Health Authority (1971 - 1981). |
| **Source** | Cancer registration data collected by ONS through the National Cancer Registration Scheme. |
| **Code notes** | 1971-1981: A to Q; Space . 1982-1986: Space. 1987 onwards: See 2nd character of NHSDHCC. See 1981 Census appendix 12c |

## Variable coding

The codelist for this variable is not available.

# The Longitudinal Impossible Dataset
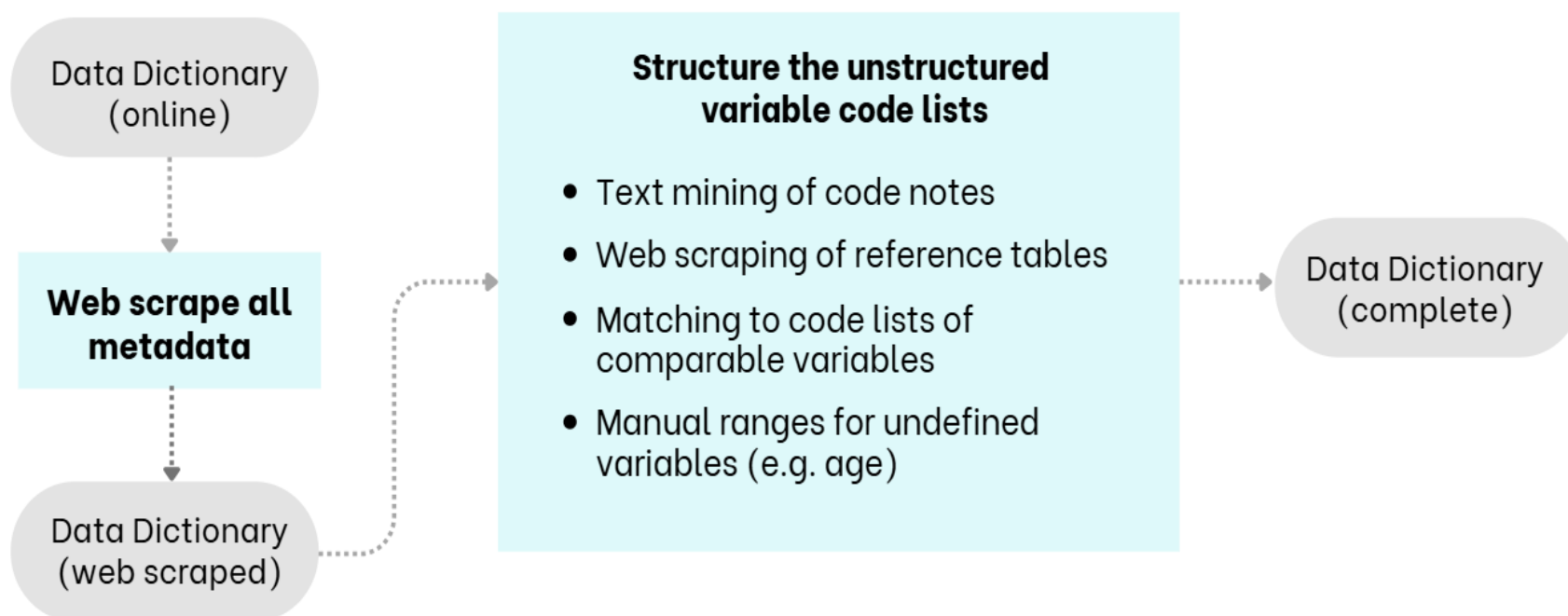
## Objectives

- Enable users to better understand what variables they need to request

- Enable users to develop code prior to data access

- Avoid real or perceived disclosure of ONS LS data

## Concept

- A fake, user-generated dataset that…

  … has the *same structure* as ONS LS data

  … is purely based on and representative of *public metadata*

  … is *impossible* – observations are clearly not real

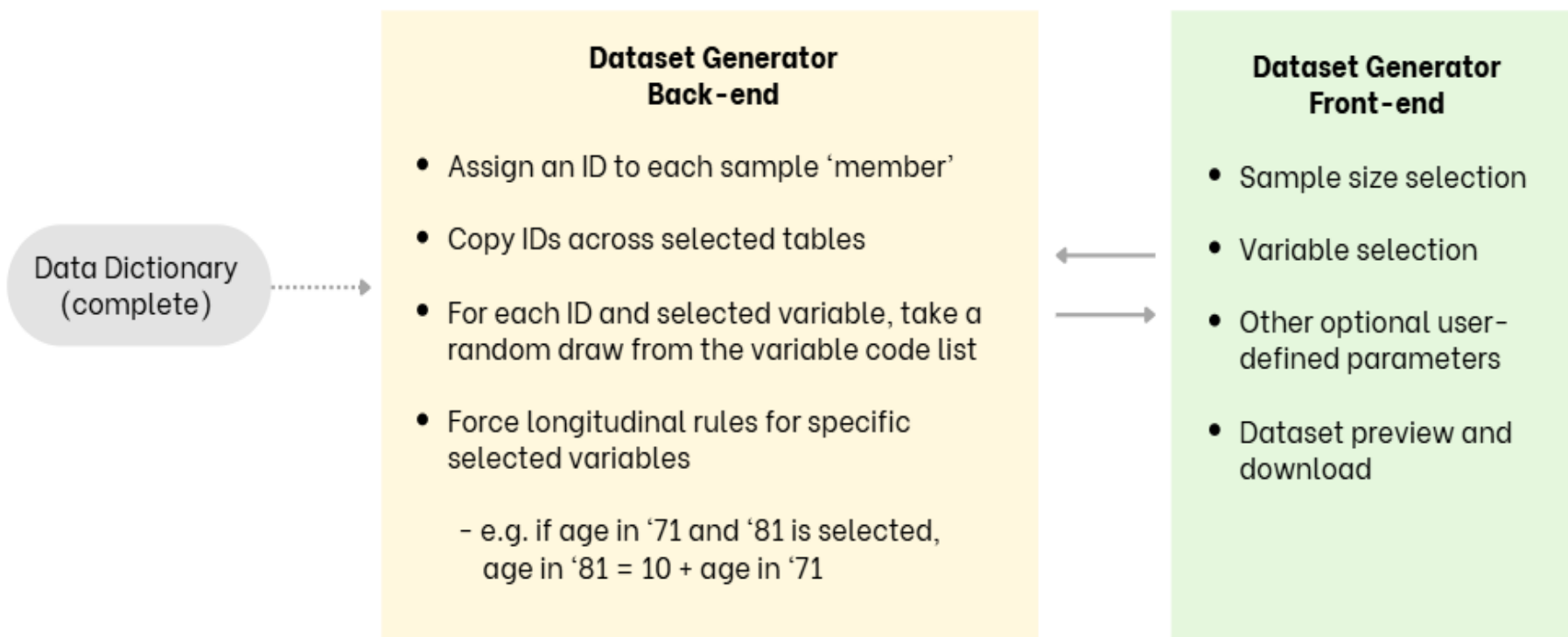  … is *longitudinal* – adheres to some realistic rules e.g. on ageing

# The Longitudinal Impossible Dataset

## Development – Input Metadata



**Data Dictionary (online)**

**Web scrape all metadata**

**Data Dictionary (web scraped)**

**Structure the unstructured variable code lists**

- Text mining of code notes
- Web scraping of reference tables
- Matching to code lists of comparable variables
- Manual ranges for undefined variables (e.g. age)

**Data Dictionary (complete)**

# The Longitudinal Impossible Dataset

## Development – Dataset Generator Application (RShiny)

Data Dictionary (complete)

### Dataset Generator Back-end

- Assign an ID to each sample 'member'

- Copy IDs across selected tables

- For each ID and selected variable, take a random draw from the variable code list

- Force longitudinal rules for specific selected variables

  - e.g. if age in '71 and '81 is selected, age in '81 = 10 + age in '71

### Dataset Generator Front-end

- Sample size selection

- Variable selection

- Other optional user-defined parameters

- Dataset preview and download

# The Longitudinal Impossible Dataset

**Key features**

- Purely based on public metadata – anyone could 'make their own'

- Purely random relationships between variables – anything is possible

- Customisable – helps users think about what variables they need

**What LIDS <u>can</u> be used for**

- Discovering relevant variables and their coding

- Familiarising with LS data extract relational table structure

- Developing mock code before accessing real data

**What LIDS <u>cannot</u> be used for**

- Estimating size of sub-populations of interest

- Identifying relationships between variables/events

# The Longitudinal Impossible Dataset

**Next steps**

- Sense-checking extracted code lists from unstructured metadata

- Explore additional longitudinal rules (e.g. births and deaths)

- Explore additional user-defined parameters

- Supporting documentation (interface & website)

- Launch – by end of 2025

# Census Innovation at CeLSIUS: other initiatives

**User support for secure Census flow data and microdata**

- CeLSIUS-ONS Steering Group

- Developing supporting guides

  - Census 2021 and the COVID-19 pandemic

  - Data catalogue

  - Data access process and restrictions (SRS/IDS)

- Exemplary research projects

- Other support and guidance

# Questions and comments

Thank you!

Andreas Mastrosavvas – amastrosavvas@ucl.ac.uk

Nicola Shelton – n.shelton@ucl.ac.uk

CeLSIUS – celsius@ucl.ac.uk