# Cleaning and analysis of the SCTS database

Graeme L Hickey[1,2]; Stuart W Grant[2]; Kate McAllister[1]; Norman Stein[1]; Iain Buchan[1]; Ben Bridgewater[2]

[1]*Northwest Institute of BioHealth Informatics*
[2]*University Hospital South Manchester*
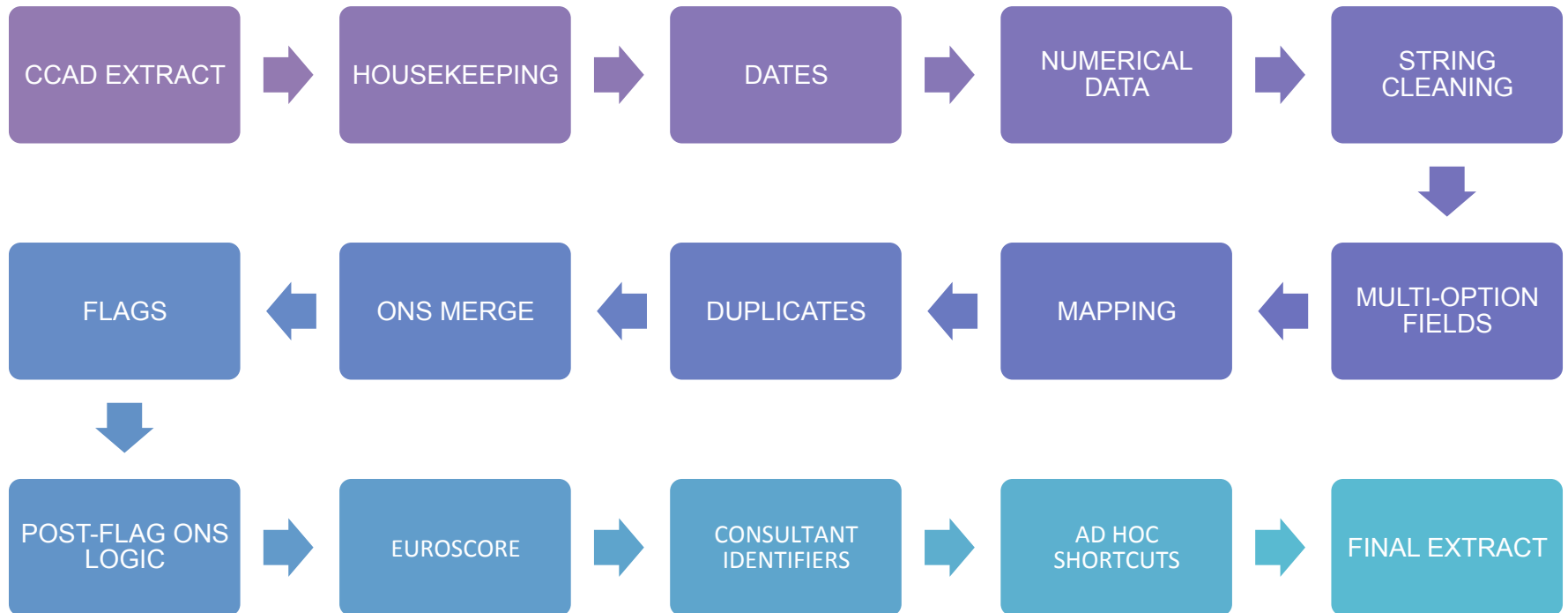
# Structure: 20 March 2012

- 444,289 records pre-cleaning
- 422,493 records post-cleaning
- 181 fields made available
- 45 hospitals in UK and Ireland

- Real world data is messy:
  - missingness
  - measurement error
  - conflicts / miscoding

requires cleaning

NIBHI

# Cleaning schema

CCAD EXTRACT → HOUSEKEEPING → DATES → NUMERICAL DATA → STRING CLEANING ↓

FLAGS ← ONS MERGE ← DUPLICATES ← MAPPING ← MULTI-OPTION FIELDS

POST-FLAG ONS LOGIC → EUROSCORE → CONSULTANT IDENTIFIERS → AD HOC SHORTCUTS → FINAL EXTRACT

# Implementation

-  : a language and environment for statistical computing and graphics

- Transparent (common S language and open source)

- Sharable (free software);

- Reproducible (tweak and re-run)

- Programmable reports (data organisation, cleaning, analysis, presentation)

- Seamless transition from cleaning to analysis

# Database in action

# Cleaning schema

# Housekeeping

- Remove identifiable fields

- Delete free text and low-importance fields

- Tidy-up field names (spelling, whitespace, etc.)

# Cleaning schema

CCAD EXTRACT → HOUSEKEEPING → DATES → NUMERICAL DATA → STRING CLEANING

FLAGS ← ONS MERGE ← DUPLICATES ← MAPPING ← MULTI-OPTION FIELDS

POST-FLAG ONS LOGIC → EUROSCORE → CONSULTANT IDENTIFIERS → AD HOC SHORTCUTS → FINAL EXTRACT

# Dates

- Formatting – time discarded except for procedure

- Delete records < 1$^{st}$ Jan 1998

- Delete dates (pre-67 and future)

- Delete records not satisfying sensible logic:

$$admission \leq procedure \leq discharge$$

# Cleaning schema

CCAD EXTRACT → HOUSEKEEPING → DATES → **NUMERICAL DATA** → STRING CLEANING

FLAGS ← ONS MERGE ← DUPLICATES ← MAPPING ← MULTI-OPTION FIELDS

POST-FLAG ONS LOGIC → EUROSCORE → CONSULTANT IDENTIFIERS → AD HOC SHORTCUTS → FINAL EXTRACT

# Numerical data

- Delete free text and symbols
- Delete impossible values (e.g. 5 valves operated on)
- Delete [clinically] unlikely values (e.g. > 11 grafts)
- Resolve 'obvious' serial imputation errors (e.g. height recorded in mm and not cm)
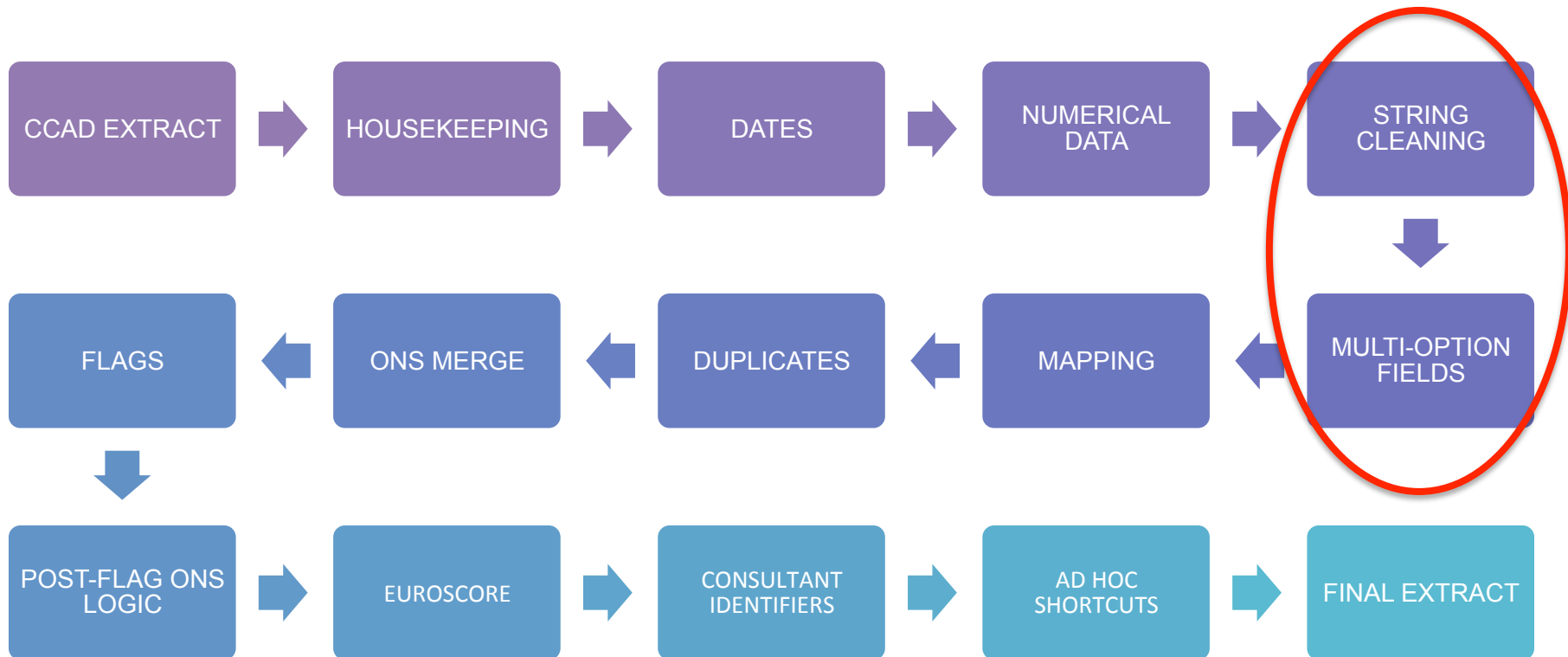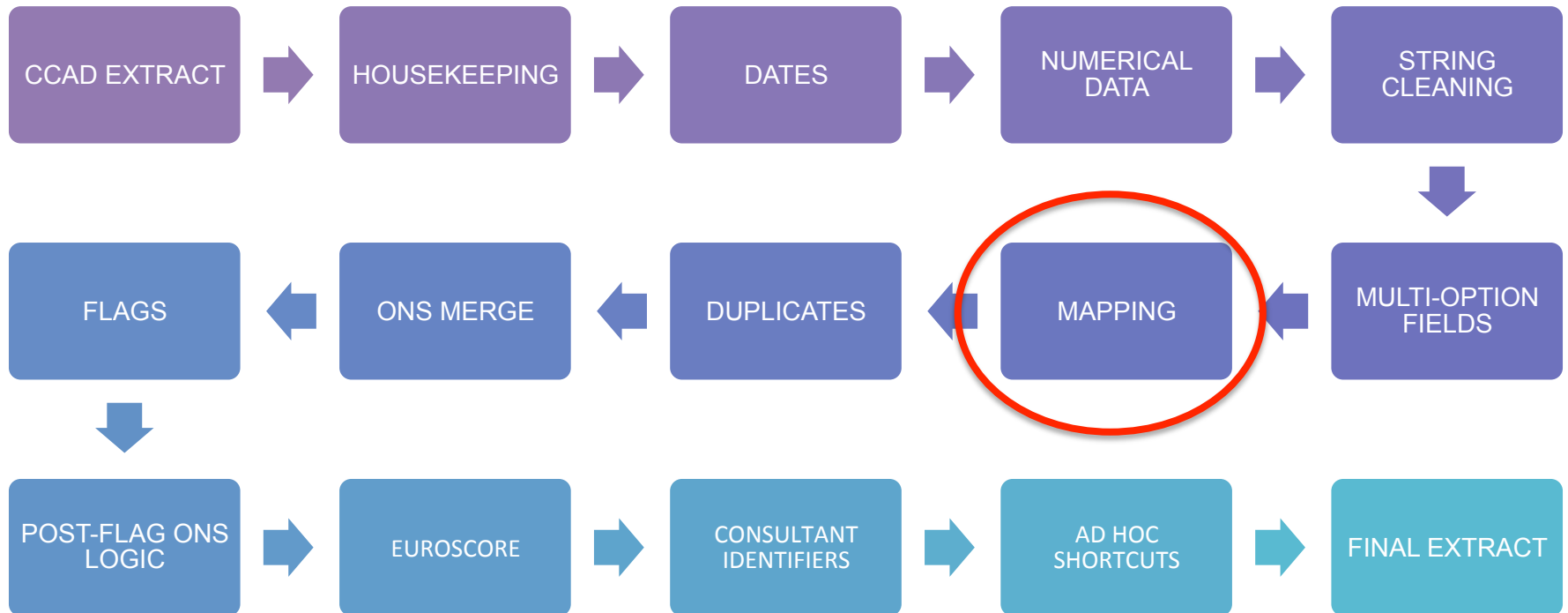
# Cleaning schema

# String cleaning

- Transcriptional errors harmonized (e.g. 'female' ➔ '2. Female')
  - manual
  - automated macros
- Invalid inputs (e.g. free text) assigned to [clinically] appropriate options
- Multi-option fields (ordered + unordered) – structure retained
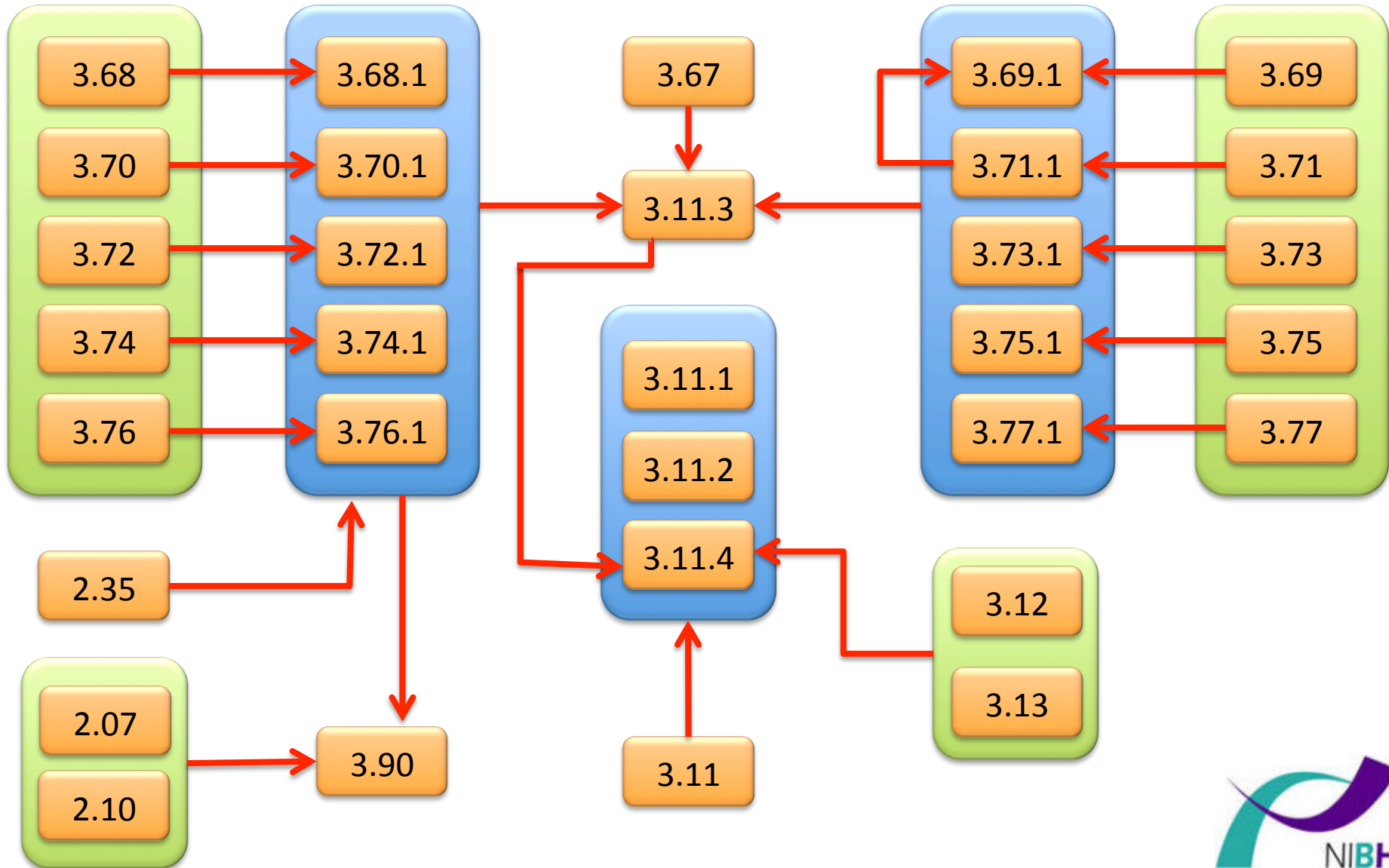- Small number of conflicts and mappings handled
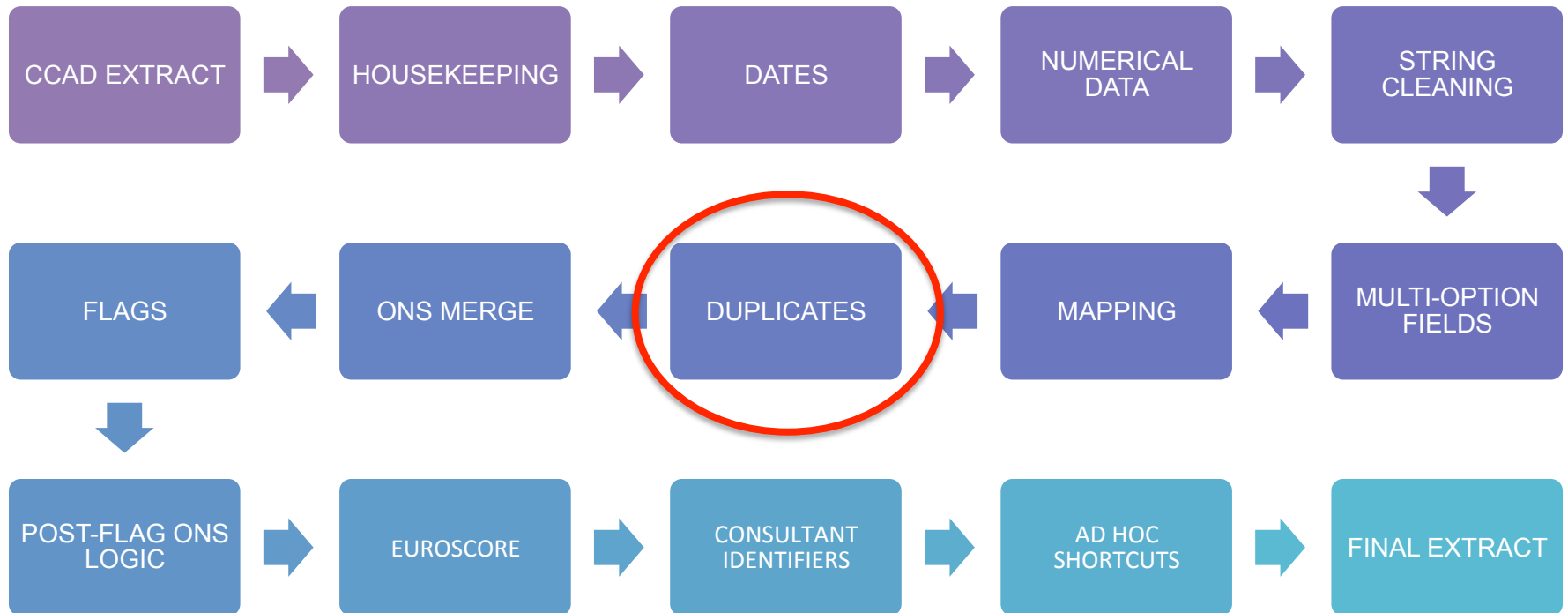
# Cleaning schema

# Mapping

- Partially fragmented about March 2010: Version 3 & 4.

- Scripts written to <span style="color:red">map</span> V3.8 into V4.1.2

- Simultaneous pre- and post-mapping cleaning

- Retrospectively deleted isolated abdominal procedure records

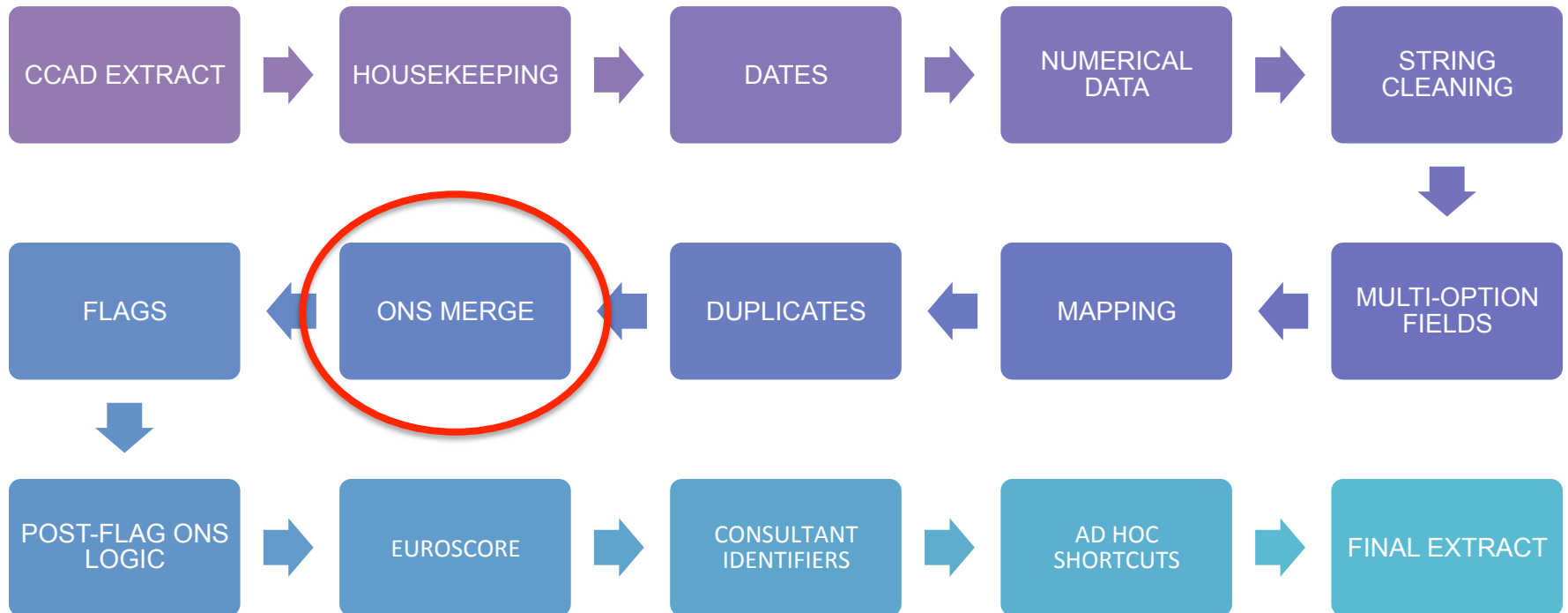# Example: major aortic fields

# Cleaning schema

# Duplicate records

- A record is classed as a duplicate if it matches on a subset.

- The most recent record created is kept; others deleted

- Records inspected after removal to 'confirm' duplicates and not re-dos

Match criteria

- ✓ hospital
- ✓ gender
- ✓ age (decimal precision)
- ✓ Apollo number (where available)
- ✓ number of previous heart operations
- ✓ procedure indicators (CABG, valve, major aortic, other)
- ✓ admission, procedure (incl. time) and discharge date
- ✓ elective (true/false)
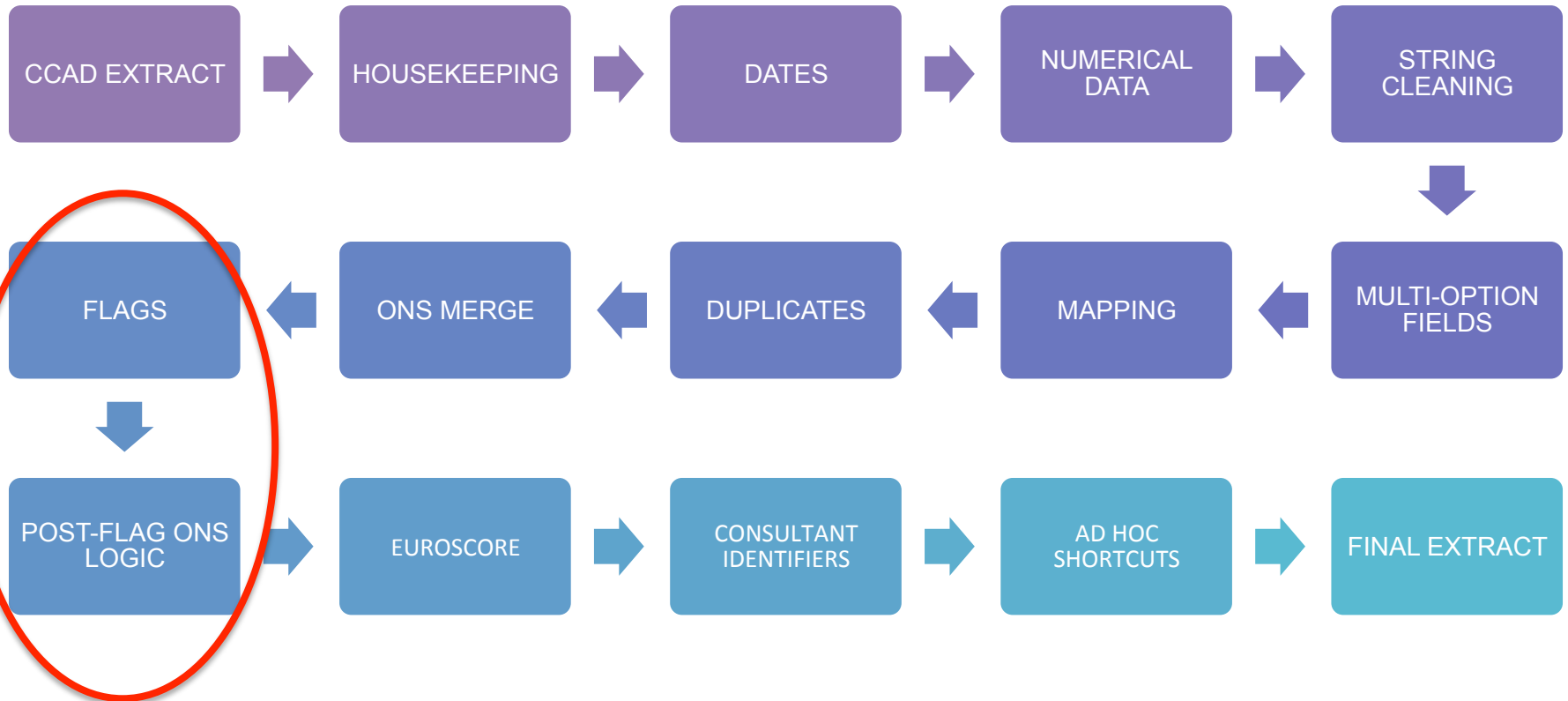
NIBHI

# Cleaning schema

# ONS data linkage



- Life status data extracted from the Office for National Statistics (ONS)

- ONS data removed if precedes procedure date

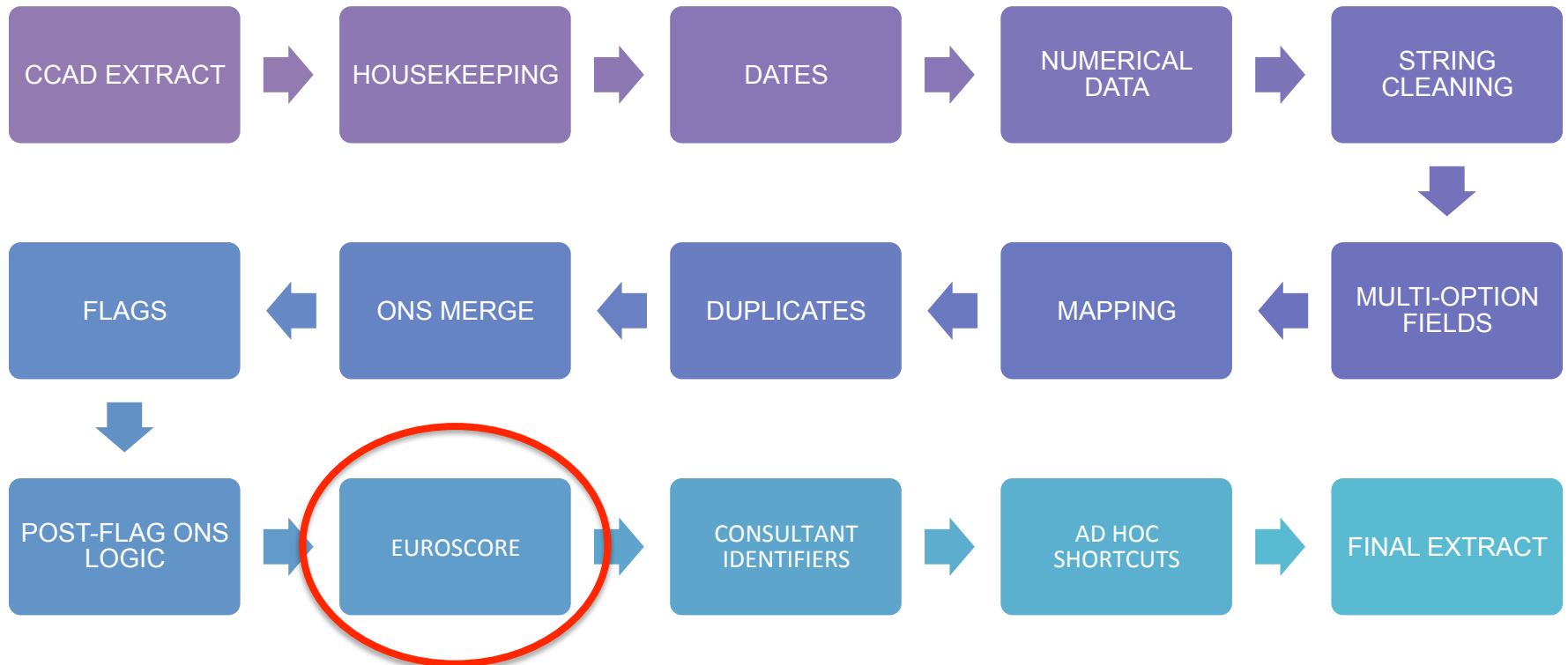- Records deleted if patient deceased prior to a first-time cardiac procedure

# Cleaning schema

# Flags

- Resolve conflicts
  - in-hospital mortality (e.g. deceased but sent home)
  - back-fill missing mortality from ONS
- Evidence based indicators (incl. resolving conflicts):
  - (individual) valve procedures
  - first operation in a single admission spell
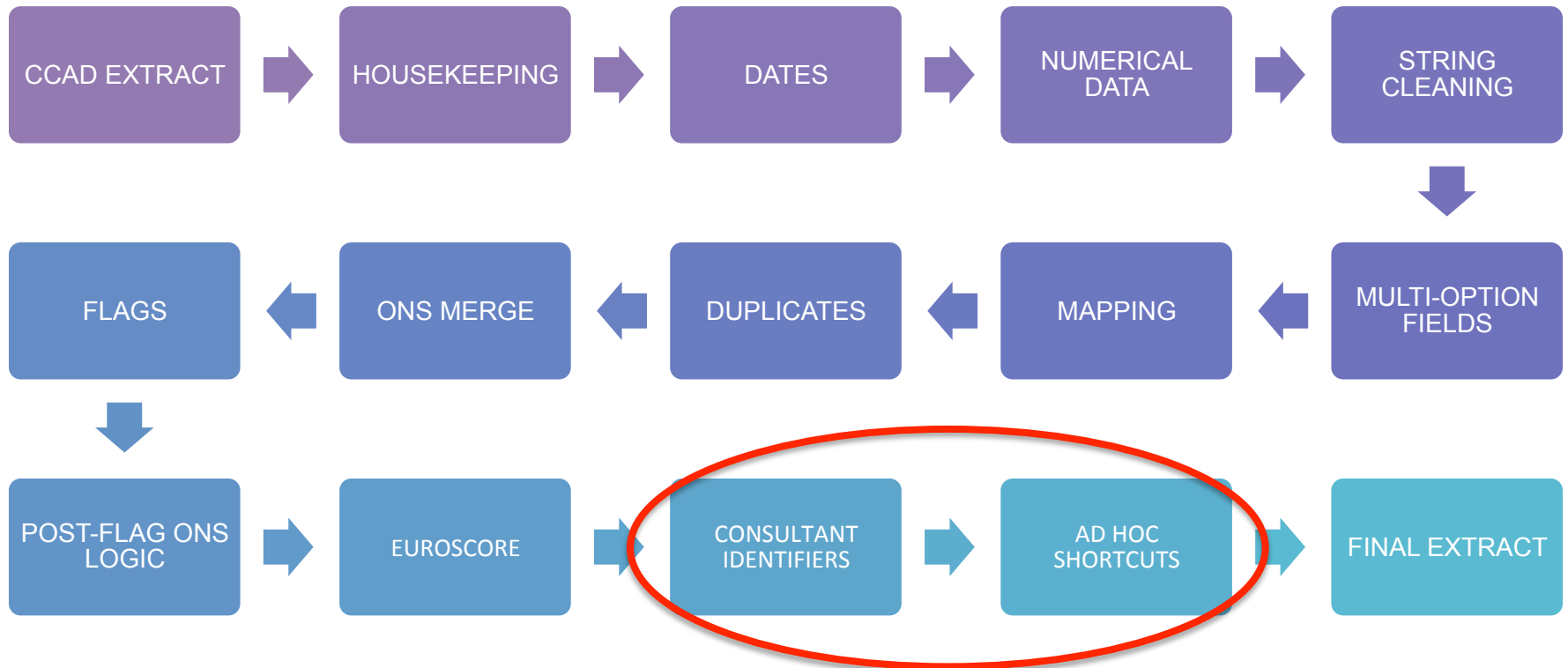  - first-time cardiac surgery

# Cleaning schema

# EuroSCORE

- 3 predictions calculated: logistic, mEuroSCORE & EuroSCORE II

- Emphasis on identifying true missing values:
  - data quality measure
  - future analysis of consequences of SCTS imputation

- Database not developed with EuroSCORE II in mind
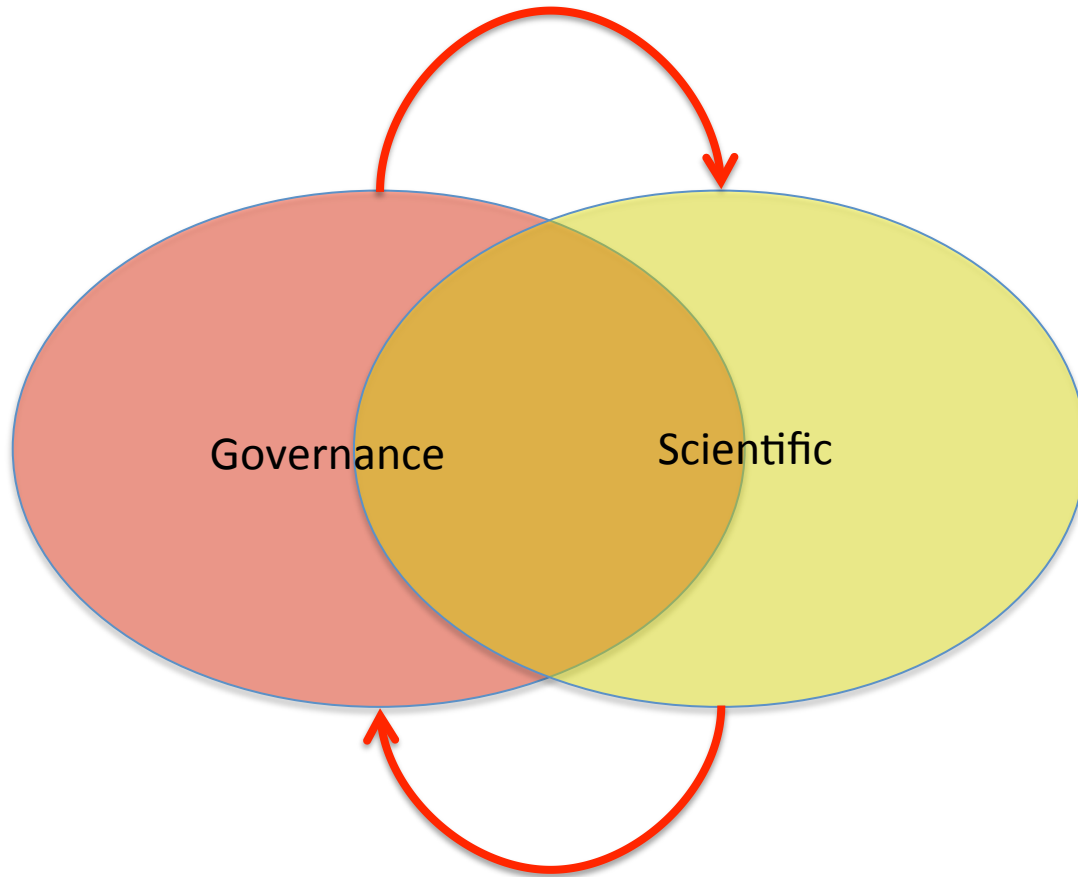
# Cleaning schema

```
CCAD EXTRACT  →  HOUSEKEEPING  →  DATES  →  NUMERICAL DATA  →  STRING CLEANING
                                                                      ↓
FLAGS  ←  ONS MERGE  ←  DUPLICATES  ←  MAPPING  ←  MULTI-OPTION FIELDS
  ↓
POST-FLAG ONS LOGIC  →  EUROSCORE  →  CONSULTANT IDENTIFIERS  →  AD HOC SHORTCUTS  →  FINAL EXTRACT
```

# Additional modules

- Consultant identifiers coded to GMC numbers
  - GMC database; hospital webpage; Dr. Forster
- Records deleted for serious ONS date discrepancies
- Expanding list of shortcut fields (e.g. country, financial year)

# Future cleaning

- Trust-level publication of deleted records

- Tweaks based on validation feedback

- Revisit assumptions + 'quick-fixes' of numerical values

- Refinement of the aortic field mappings

- Centralized cleaning / mapping by NICOR

# Analyzing the data

# Governance



EuroSCORE II: all cardiac surgery

# Informing our members

# Responding to contemporary questions

# Measuring data quality



Distribution of ranks of EuroSCORE risk factor prevalence might be expected to homogenous across hospital

Further investigation required

# Scientific

- Mitral valve prosthesis: mechanical vs. biological

- Model validation (➜ ensure current governance)

- Calibration drift detection methodology (➜ inform future governance)

# Further information

- SCTS website
  - [www.scts.org/](www.scts.org/)
- SCTS-NIBHI project website (incl. contacts)
  - [personalpages.manchester.ac.uk/staff/graeme.hickey/scts/](personalpages.manchester.ac.uk/staff/graeme.hickey/scts/)
- NICOR website
  - [www.ucl.ac.uk/nicor](www.ucl.ac.uk/nicor)

# Acknowledgements

- **Heart Research UK** – funding

- **Sue Manuel (NICOR)** – database extracts

- **All hospital audit leads and database managers** – validating audit summaries

- **UK cardiac surgeons** – ensuring the validity and accuracy of the data inputted

- **The SCTS and all its members** – for supporting the audit project