# Effecting rigorous data harmonization and documentation to understand data heterogeneity and quality

Tina W. Wey & Isabel Fortier

Maelstrom Research
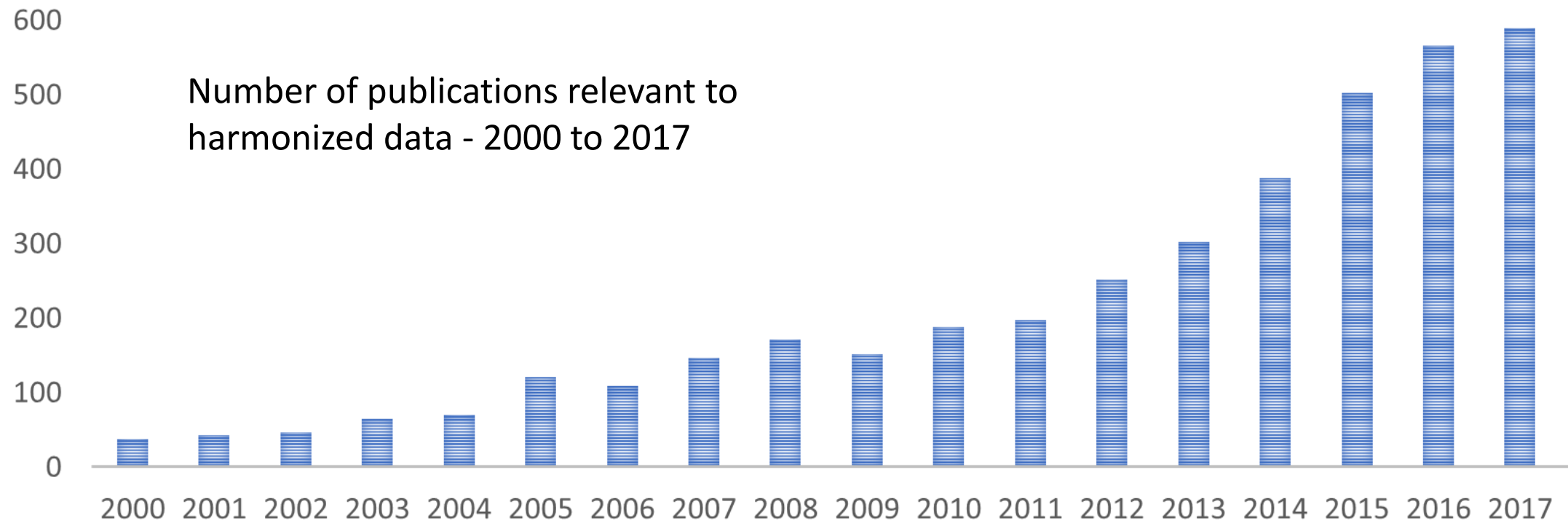
Building Multi-Source Databases for Comparative Analyses
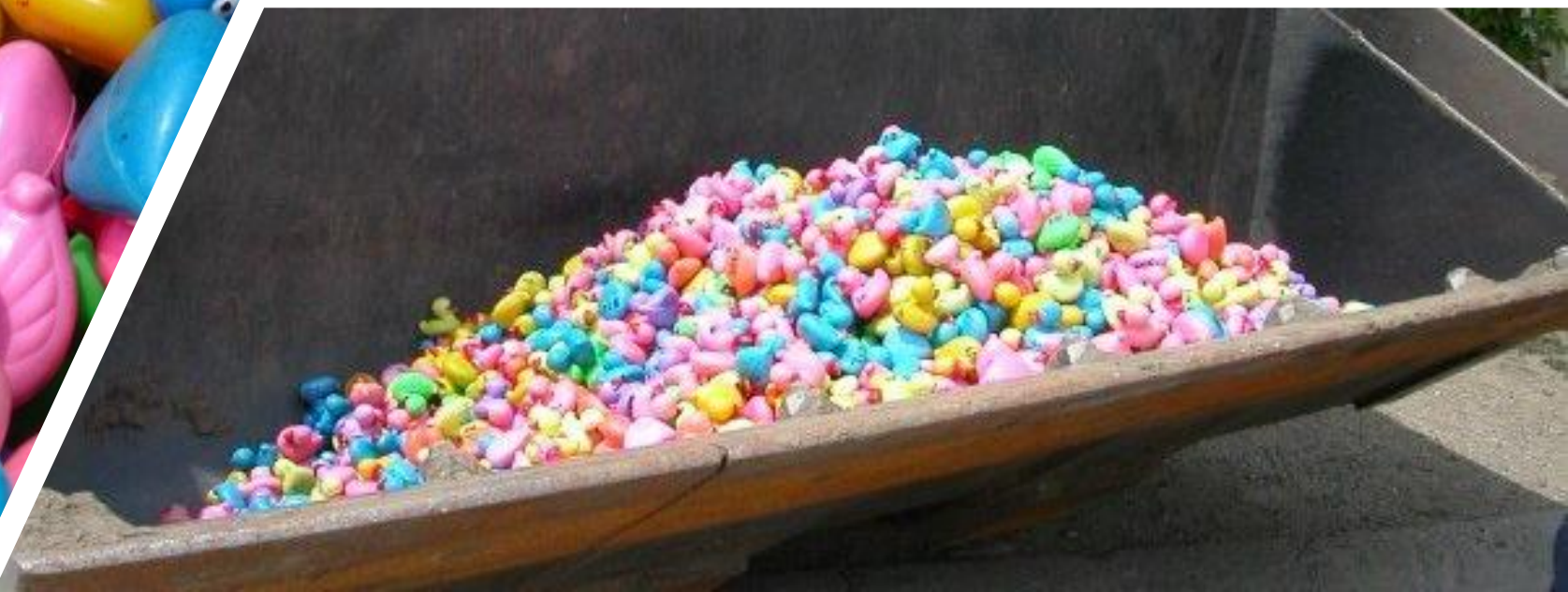
Warsaw, Poland, 17 December 2019

# Increasing need for harmonized data in epidemiological research

Driven by need to obtain **larger sample sizes** and statistical power; conduct **comparative research** across studies/populations; **extend** the scientific **impact** of individual studies/data sources.

Offers benefits: enabling **timely access** to available data and samples, increasing **potential to share** data across studies, and promoting a **collaborative approach**

Number of publications relevant to harmonized data - 2000 to 2017

Study A

Study B

Study C

# Maelstrom Research

**Facilitate collaborative epidemiological research through rigorous data documentation, harmonization, integration, and co-analysis**

## Who we are:

**Hosted at**
Research Institute of the McGill University Health Centre in Montreal, Canada

**International research program**
partnering with over 15 international networks and research consortia

**Multi-disciplinary team**
epidemiologists, data analysts, and computer scientists

## Activities:

**Methodological guidelines/support**
for data cataloguing, harmonization, integration, and co-analysis

**Web-based catalogues and harmonization platforms**
searchable metadata catalogues and platforms to generate common-format variables for co-analysis

**Open-source software**
for data cataloguing, harmonization, integration, and co-analysis

# mælstrom

Methodological guidelines and open-source software to support data collection, management, dissemination, harmonization and co-analysis

onyx    opal    mica    DataSHIELD Secure Bioscience Collaboration

DataSHaPER

A central study catalogue to foster usage of available data

National and international platforms harmonizing, integrating and co-analysing data across studies

# Maelstrom harmonization guidelines

**0** **Define the research question(s)**

**1** **Assemble information and select studies**
1. Document individual study designs, methods and content
2. Select participating studies

**2** **Define variables and evaluate harmonization potential**
1. Select and define the core variables to be harmonized
2. Determine the potential to create the core variables using the study-specific data items

**3** **Process data**
1. Ensure access to adequate study-specific data items and establish the overall data processing infrastructure
2. Process study-specific data under a common format to generate the harmonized datasets

**4** **Estimate quality of the harmonization dataset(s) generated**

**5** **Disseminate and preserve final harmonization products**

IEA

Original Article

## Maelstrom Research guidelines for rigorous retrospective data harmonization

Isabel Fortier,[1]* Parminder Raina,[2] Edwin R Van den Heuvel,[3] Lauren E Griffith,[2] Camille Craig,[1] Matilda Saliba,[1] Dany Doiron,[1] Ronald P Stolk,[4] Bartha M Knoppers,[5] Vincent Ferretti,[6] Peter Granda[7] and Paul Burton[8]

# A systematic but adaptable process

Iterative, dynamic process of consideration, evaluation, discussion, validation

Documentation and assessment of source data heterogeneity to understand harmonized output



Iterative Harmonization Steps

Step 0: Define the research questions, objectives and protocol

Step 1: Assemble information and select studies

Step 2: Define variables and evaluate harmonization potential

Step 3: Process data

Step 4: Estimate quality of the harmonized dataset(s) generated

Step 5: Disseminate and preserve final harmonization products

# Assemble information and select studies: Cohort metadata catalogue



**Study description**
(e.g., design, participant selection criteria, data collection events)

**Areas of information**
(e.g., smoking, cancer, anthropometrics)

**Variable metadata**
(e.g.,variable name/label, categories, units)

**Specific data**
(individual participants data collected)

Step 0 · Step 1 · Step 2 · Step 3 · Step 4 · Step 5
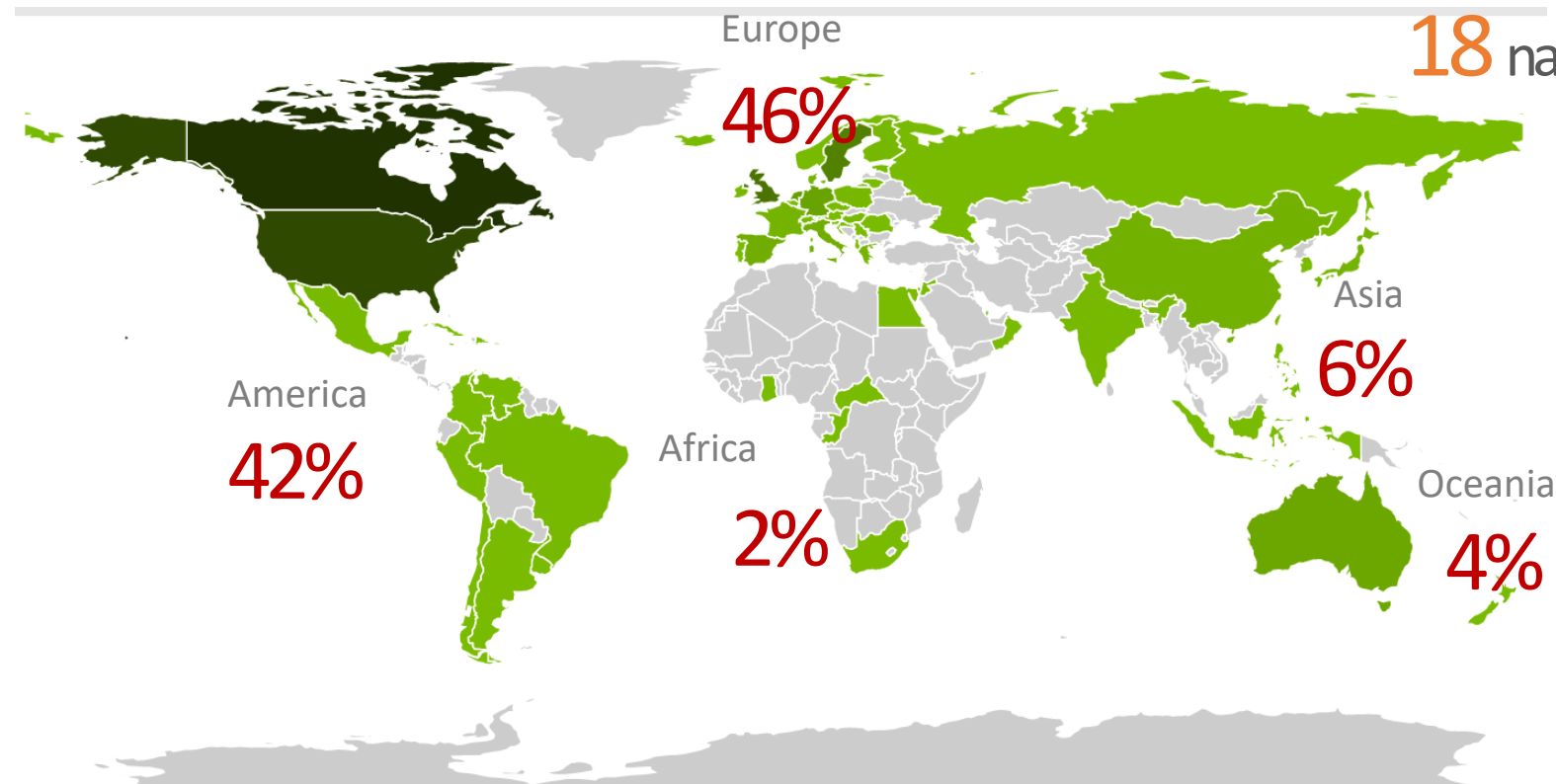
Iterative Harmonization Steps

PLOS ONE

RESEARCH ARTICLE

Fostering population-based cohort data discovery: The Maelstrom Research cataloguing toolkit

Julie Bergeron[1], Dany Doiron[1,2,3], Yannick Marcon[1], Vincent Ferretti[4], Isabel Fortier[1]*

1 Research Institute of the McGill University Health Centre, Montreal, Quebec, Canada, 2 Swiss Tropical and Public Health Institute, Basel, Switzerland, 3 University of Basel, Basel, Switzerland, 4 Research Center of the Sainte-Justine University Hospital, Montreal, Quebec, Canada

# The Maelstrom Research metadata catalogue



Europe
46%

America
42%

Africa
2%

Asia
6%

Oceania
4%

**18** national and international networks

**204** studies (**122** with variables)

**933,144** variables

**6,349,772** cohort participants

**Studies, including...**

# Illustrative harmonization projects

Canadian Partnership for Tomorrow Project

Environmental, lifestyle and genetic factors related to the development and progression of cancer and chronic diseases; Prospective design; 5 Canadian provinces

MINDMAP
Promoting mental well-being and healthy ageing in cities

Urban environments and promotion of mental wellbeing and cognitive function of older individuals;
Retrospective design; 7 European countries, Russia, and Canada

ReACH
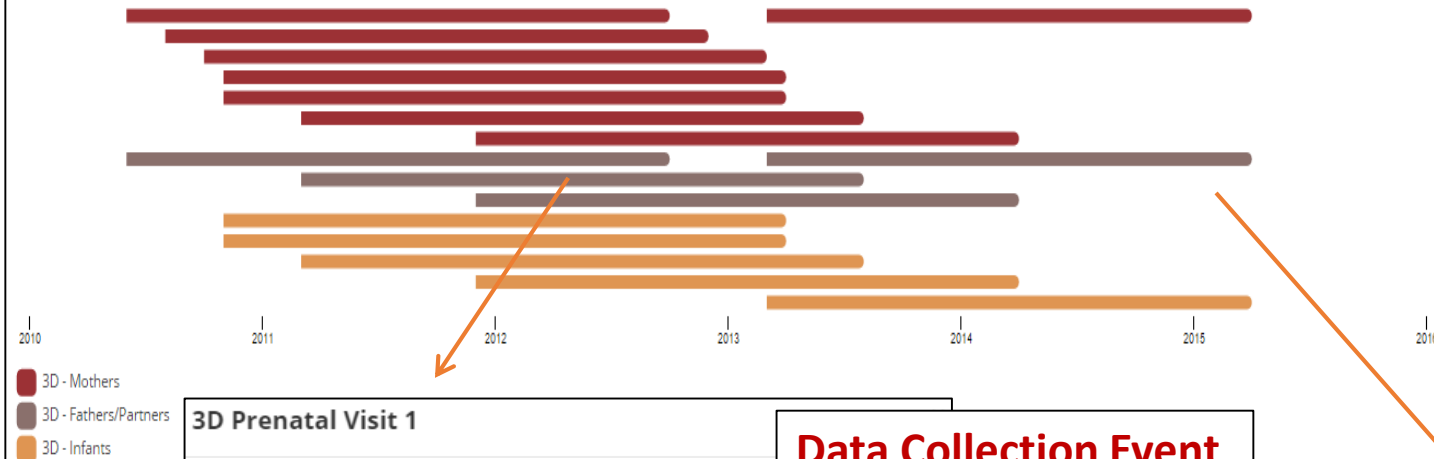Research Advancement through Cohort Cataloguing and Harmonization

Canadian pregnancy and birth cohorts data and biological samples to study Developmental Origins of Health and Disease (DOHaD);
Retrospective design; Canadian

# Study Description

## Timeline

Each colour in the timeline graph below represents a separate Study Population, while each segment in the graph represents a separate Data Collection Event. Clicking on a segment gives more detailed information on a Data Collection Event.

- 3D - Mothers
- 3D - Fathers/Partners
- 3D - Infants

2010  2011  2012  2013  2014  2015  2016

## Data Collection Event

### 3D Prenatal Visit 1

3D first trimester (8 to 14 weeks gestation) visit with the mother to be and her partner, including questionnaires administered by trained staff and self-administered, anthropometric measurements, and biospecimen collection.

Paternal blood and urine collection occurred as soon as possible post-conception, ideally at visit 1, but to optimize collection, the actual window of collection was extended until visit 5. If father was recruited postnatally, only blood for DNA or saliva for DNA was collected.

| Start Year | 2010 (June) |
| --- | --- |
| End Year | 2012 (September) |
| Data Sources | • Questionnaires<br>• Physical Measures<br>• Biological Samples |
| Biological Samples | • Blood<br>• Urine |

Description of the event
Start and end years
Data sources
Biological samples

**ReACH**
Research Advancement through Cohort Cataloguing and Harmonization

## 3D - 3D Study - Design, Develop, Discover

### Overview

| Acronym | 3D |
| --- | --- |
| Website | 3D website |
| Investigators | Dr. William D. Fraser (University of Sherbrooke)<br>Dr. Lise Dubois (University of Ottawa)<br>Dr. Zhong-Cheng Luo (University of Montreal)<br>Dr. Jacques Michaud (University of Montreal)<br>Dr. Jean-Marie Moutquin (University of Sherbrooke)<br>Dr. Gina Muckle (Laval University)<br>Dr. Jean Séguin (University of Montreal)<br>Dr. Margaret Somerville (Université McGill)<br>Dr. Jacquetta Trasler (McGill University)<br>Dr. Richard E. Tremblay (University of Montreal)<br>Dr. François Audibert (University of Montreal)<br>Dr. Pierre Julien (Laval University) |
| Contacts | Martine Fournier (CHU Sainte-Justine Research Centre)<br>Josée Poirier (CHU Sainte-Justine Research Centre)<br>Isabelle Krauss (CHU Sainte-Justine Research Centre) |
| Study Start Year | 2010 |

### Access

Access to external researchers or third parties provided or foreseen for:

| Data (questionnaire-derived, measured...) | ✔ |
| --- | --- |
| Biological samples | ✔ |

### Design

| Study Design | Cohort Study |
| --- | --- |
| General Information on Follow Up (profile and frequency) | Pregnant women and their partners were recruited during the first trimester of pregnancy and were followed throughout pregnancy and birth, and along with their children up to 2 years of age, with a total of 8 visits. |
| Recruitment Target | Families |
| Target number of participants | 2456 |
| Target number of participants with biological samples | 2357 |
| Supplementary information about target number of participants | 2456 participants were originally recruited. (The study is still ongoing, with completion in spring of 2015. There was some attrition throughout the study and therefore, each participant has a different number of visits completed).<br>There are 2357 mothers with at least one biological sample and 2333 fathers with at least one biological sample. |

### Design

Objectives
Study design
Start and end years
General information on follow-up
Recruitment target
Number of participants

### Populations

3D - Mothers
3D - Fathers/Partners
3D - Infants

**3D - Mothers**

A total of 2456 pregnant women from the general population who were attending prenatal clinics (ultrasound, midwife and/or doctor's clinics) during the first trimester of pregnancy were recruited for the study. The recruited women had to be between 18 and 45 years of age, 6 to 13 weeks pregnant at the time of recruitment, fluent in French or English, and plan to deliver in a study hospital to be eligible for the study.

**Sample Size**

| Number of participants | 2456 |
| --- | --- |
| Number of participants with biological samples | 2357 |

**Sources of Recruitment**

| Specific population | Clinic Patients |
| --- | --- |
| Supplementary information | The women were recruited from prenatal clinics. |

**Selection Criteria**

| Gender | Women only |
| --- | --- |
| Age | Minimum 18, Maximum 45 |
| Country | Canada |
| Territory | Quebec and Eastern Ontario |
| Other | Exclusion criteria:<br>Multiple pregnancies, intention to donate or bank cord blood, intravenous drug users, or if the woman has any one or more of the following conditions: HIV+ status, renal disease with altered renal function, any collagen vascular disease requiring active treatment (e.g. lupus, scleroderma), epilepsy, cardiovascular disease, serious pulmonary disease, cancer, or severe hematologic disorder. |

### Sub-population

Description of the population
Sources of the recruitment
Selection criteria
Number of participants

# Detailed study-specific source variable information

## Overview

| | |
|---|---|
| **Label** | 1a. lifetime: Smoke a total of 100 or more cigarettes |
| **Description** | 1.1. In your lifetime, have you smoked a total of 100 or more cigarettes (about 4 packs)? |
| **Individual Study** | 3D |
| **Dataset** | 3D_Prenatal_Visit1_Mother |
| **Value Type** | Integer |
| **Variable Type** | Collected |

## Classifications

| Additional information | |
|---|---|
| **Source** | Questionnaire |
| **Target** | Participant |

| Areas of information | |
|---|---|
| **Lifestyle and behaviours** | Tobacco |

## Categories

| Name | Label | Missing |
|---|---|---|
| 0 | not at all | |
| 1 | yes | |
| 88 | no data | |
| 98 | refuse to answer | |
| 99 | don't know | |

Clear | ❖ | ⓘ ReACH ▾ | x | ᵗⓐ | ⓘ Selection criteria - Pregnant ... ▾ | x | ⓘ Start year:[2010,2018] ▾ | x | ᵗ 🎨 | ⓘ Nutrition | Physical activity ▾ | x | ⓘ Age/birthdate | Education ▾ | x

ⓘ Anthropometry ▾ | x    Advanced

## ▾ Variables

| Areas of Information | › |
| Scales/Measures | › |
| Source & target | › |
| Properties | › |

## ▾ Studies

| Properties | › |

## ▾ Networks

| Properties | › |

List | **Comparison Table** | Summary Statistics

⬕ **Download**

All | Individual | Harmonization                    ☐ Population/Data Collection Event (DCE)

| ☐▾ | Study | Socio-demographic and economic characteristics ✕ | | Lifestyle and health behaviours ✕ | | Physical measures ✕ |
|---|---|---|---|---|---|---|
| | | Age/birthdate ✕ | Education ✕ | Nutrition ✕ | Physical activity ✕ | Anthropometry |
| ☐ | 3D | 24 | 5 | 191 | 151 | 104 |
| ☐ | ABC | 28 | 14 | 2,321 | 101 | 286 |
| ☐ | OBS | 53 | 6 | 33 | 47 | 233 |
| ☐ | START | 33 | 33 | 6,021 | 195 | 632 |
| | **All** | **138** | **58** | **8,566** | **494** | **1,255** |

# Define core variables (DataSchema)

***Quantity*** = number of studies to include
***Quality*** = scientific relevance/precision

# Evaluate harmonization potential



**Variable: Number of red wine drinks**

| Study A | Study B | Study C |
|---|---|---|
| **Period:**<br>Week (7 days) | **Period:**<br>Weekdays (Sunday to Thursday)<br>Weekend days (Friday to Saturday) | **Period :**<br>Weekday day (working day)<br>Weekend day (non-working day) |
| **Unit:**<br>Drinks/week | **Unit:**<br>Per weekdays<br>Per weekend days | **Unit:**<br>Per day |

status: **Complete**

status: **Complete**

status: **Impossible**

+

**Number of red wine drinks per week**  X

DataSchema variable

**Comment:** The information on alcohol quantity is collected differently. The number of drinks of alcohol is asked in separate questions for working days and non working days without specifying the number of days of each period.

# Target variable: Frequency of Binge Drinking During Pregnancy

## Study A

**3 collections**

**1st collection : 8 – 14 weeks**

Question: Since you have become pregnant, how often did you have 5 or more drinks on one occasion?

Response: #days of week OR #days of month OR #days since beginning of your pregnancy

**2nd collection: 20 - 24 weeks**

Question: Since your last visit, how often did you have 5 or more drinks on one occasion?

Response: #days of week OR #days of month OR #days since beginning of your pregnancy

**3rd collection: 32 - 35 weeks**

Question: Since your last visit, how often did you have 5 or more drinks on one occasion?

Response: #days of week OR #days of month OR #days since beginning of your pregnancy

## Study B

**2 collections**

**1st collection : 12 – 16 weeks**

Question: Please specify the number of times per month you have four or more drinks at the same sitting or occasion (during this pregnancy)?

Response: >= 1times/month. Please specify number:___ | < 1/month | None

**2nd collection : 28 – 32 weeks**

Question: Over the past 3 months, how often did you have four or more drinks at the same sitting or occasion?

Response: 6 to 7 times a week | 4 to 5 times a week | 2 to 3 times a week | once a week | 2 to 3 times a month | about once a month | 6 to 11 times a year | 1 to 5 times a year | never

## Study C

**1 collection**

**1st collection : 21 – 39 weeks**

Question: During this pregnancy, how many times have you consumed at least 5 or more drinks of alcohol in a day?

Response: continuous

### Challenges
- Timing
- Wording of questions
- Wording of categories
- Responses options
- Data collection events

**ReACH**
Research Advancement through
Cohort Cataloguing and Harmonization

# Process data: Access to data



**Step 0** | **Step 1** | **Step 2** | **Step 5**
**Step 4** | **Step 3**

Iterative Harmonization Steps

Data repository application integrating and storing data from multiple sources

R Studio® Open-source software for data analysis



**Opal clients**

**Opal Admin client**
To administrate Opal server, import, manipulate, and export study data

Primary user:
Study coordinator, Data manager

**R client**
To conduct statistical analyses using R or Rstudio.

Primary user:
Researcher, statistician

**Python client**
To conduct automated tasks and batch processing

Primary user:
Data manager

**Mica clients**

**Mica Admin client**
To administrate Mica server, add/edit information to be displayed on website

Primary user:
Study or network coordinator

**Mica Portal client**
To render and display information on website

Primary user:
Research community, Public

**R server**
Temporarily hosts data used for R analyses

**Opal server**
Hosts data and data dictionaries (i.e. codebooks)

**Mica server**
Hosts metadata and links to data on Opal

# opal

## Variable

# Tables

V CPTP / Coreqx_final / A_SDC_ADOPTED_CHILD ☆

| Dictionary | Summary | Values | Permissions |

⚙ Derive ▾     🛒     🗑     ⌃     ⌄

## Properties ✎

| Name | A_SDC_ADOPTED_CHILD | | Unit | |
|---|---|---|---|---|
| Entity Type | Participant | | Referenced Entity Type | |
| Value Type | integer | | Mime Type | |
| Repeatable | No | | Occurrence Group | |

## Categories

✎ Edit Categories

Total 2

| Name | Label | Missing |
|---|---|---|
| 0 | Not adopted | |
| 1 | Adopted | |

## Attributes

| Standard | Raw |

### Label ✎

Adopted

### Description ✎

Indicator of whether the participant was adopted.

## Annotations

✎ Edit Annotation ▾     🔍 Search similar variables

Areas of Information

| **Socio-demographic and economic characteristics** | Other socio-demographic and economic characteristics |
|---|---|
| Refers to sociodemographic and economic characteristics of an individual. | Information about other socio-demographic and economic characteristics (e.g. being adopted). |

# Assess study-specific source data


Iterative Harmonization Steps

- Verify:
  - Data format and compatibility with Opal
  - Entity IDs, duplicate IDs, IDs missing values
  - Inclusion criteria
  - Variable list, metadata, format
  - Univariate checks
  - Multivariate checks for cross-variable coherence
  - Document issues, summary reports, communication with cohorts
  - …

# Generate core variables

opal



Iterative Harmonization Steps

**Study specific variables**

Case 1:    Ever had sigmoidoscopy or colonoscopy

Case 2:    Ever had sigmoidoscopy
           Ever had colonoscopy

**DataSchema variable**

Ever had sigmoidoscopy or
colonoscopy

| Case | Rule | Script |
|------|------|--------|
| 1 | Direct mapping from source variable | ```$('uhlq_hc_3').map({`<br>`    '1': '1',`<br>`    '2': '0'`<br>`},`<br>`null,`<br>`null);``` |
| 2 | Sourced from DataSchema variables<br><br>If HS_SIG_EVER = 1  OR HS_COL_EVER = 1 --> code to 1<br><br>If HS_SIG_EVER = 0  AND HS_COL_EVER = 0 --> code to 0 | ```var sig_ever = $this('HS_SIG_EVER');`<br>`var col_ever = $this('HS_COL_EVER');`<br><br>`if (sig_ever.eq(1).or(col_ever.eq(1)).value()) { //if either is ever --> ever`<br>`  1;`<br>`} else if (sig_ever.eq(0).and(col_ever.eq(0)).value()) { //if both are never --> never`<br>`  0;`<br>`} else { //if either is null--> null`<br>`  null;`<br>`}``` |

Variable summary statistics

# Generate core variables



R Studio®



Iterative Harmonization Steps

```
44
45   # Table of Contents
46   1. [Alcohol consumption subdomain](#alcohol-consumption-subdomain)
47   2. [Tobacco consumption subdomain](#tobacco-consumption-subdomain)
48   3. [Physical activity subdomain](#physical-activity-subdomain)
49   4. [Diet and nutrition subdomain](#diet-and-nutrition-subdomain)
50   5. [Sleep quality subdomain](#sleep-quality-subdomain)
51
52   # Alcohol consumption subdomain
53
54   ### **Variable label**:  Current consumption of alcohol
55   **Variable name**: lsb_alc_cur_0
56   **Variable description**:  Indicator of whether the participant currently consumes alcohol
57   **Value type**: integer
58   **Variable unit**: N/A
59   **Category coding**:
60
61   **Code** | **Category Label**
62   ------------- | -------------
63   0 |  Does not currently consume alcohol
64   1 |  Currently consumes alcohol
65
66   **Harmonization status**:  complete
67   **Harmonization comment**:
68   **R script**:
69   ```{r, echo=TRUE}
70   lsb_GLOBE_0$lsb_alc_cur_0<-ifelse(GLOBE1991$v135<6,1L,
71                            ifelse(GLOBE1991$v135==6,0L,NA))
72
```

## MINDMAP

MINDMAP is a multi-cohort research project exploring the urban environment and mental well-being. This space is used to manage MINDMAP data harmonization work.

🔗 http://www.mindmap-cities.eu/

📖 Repositories 10

**lifestyle_behaviours**
Lifestyles and behaviours domain data harmonization work repository // Lead - Marielle Beenackers (EMC)
⑂2  ★0  ①0  ⑃1  Updated 2 days ago

**sociodem_characteristics**
Sociodemographic characteristics domain data harmonization work repository // Lead - Rita Wissa (RI-MUHC)
⑂4  ★0  ①0  ⑃1  Updated 2 days ago

**mental_health_outcomes**
Mental health outcomes domain data harmonization work repository // Lead - Milagros Ruiz (UCL)
⑂3  ★0  ①0  ⑃0  Updated 3 days ago

**other_outcomes**
Other outcomes domain data harmonization work repository // Lead - Marielle Beenackers (EMC)
⑂2  ★0  ①0  ⑃0  Updated 3 days ago

**Harmonized-Datasets**
●R  ⑂0  ★0  ①0  ⑃0  Updated 5 days ago

Top languages
● R

People ›
This organization has no public members. You must be a member to see who's a part of this organization.

GitHub

# Estimate quality of harmonized dataset

- For each study-specific harmonized dataset:
  - Validate harmonization process (algorithms, scripts)

  - Validate data content and consistency
  - Distributions and missing values
  - Consistency with DataSchema (format, categories)
  - Harmonization completion statuses
  - Multivariate checks for cross-variable coherence
  - Document issues, summary reports, communication with cohorts
  - …

# Alcohol consumption of mother 1 year prior to pregnancy : Y/N ?

## Study A

| Category | Freq. |
|---|---|
| Every day | 43 |
| 4-6 / week | 185 |
| 2-3 / week | 503 |
| 1/ week | 380 |
| 2-3 /month | 278 |
| 1 / month | 185 |
| < 1 /month | 325 |
| Never | 461 |
| Missing | 5 |
| Total | 2 365 |

### Study variable(s)

[Alcohol frequency during year before pregnancy]

### DataSchema variable values

| Value | Condition |
|---|---|
| 0 | Mapping from study variable if<br>• [Alcohol frequency during year before pregnancy] = *None* |
| 1 | Mapping from study variable if<br>• [Alcohol frequency during year before pregnancy] ≥ 1 |
| | Missing |

**Study-specific harmonized variable**

| Category | Freq. |
|---|---|
| Yes | 1 899 |
| No | 461 |
| Missing | 5 |
| Total | 2 365 |

## Study B

| Category | Freq. |
|---|---|
| Yes | 2 728 |
| No | 590 |
| Missing | 23 |
| Total | 3 341 |

### Study variable(s)

[Alcohol use 12 months before pregnancy]

### Dataschema variable values

| Value | Condition |
|---|---|
| 0,1 | Direct mapping from study variable |
| | Missing |

**Study-specific harmonized variable**

| Category | Freq. |
|---|---|
| Yes | 2 728 |
| No | 590 |
| Missing | 23 |
| Total | 3 341 |

## Study C

| Category | Freq. |
|---|---|
| Yes | 261 |
| No | 1835 |
| Missing | 90 |
| Total | 2 187 |

### Study variable(s)

[Never consumed alcohol]

### Dataschema variable values

| Value | Condition |
|---|---|
| 0,1 | Impossible |
| | Missing |

**Study-specific harmonized variable**

| Category | Freq. |
|---|---|
| Yes | 0 |
| No | 0 |
| Missing | 2 187 |
| Total | 2 187 |

**ReACH**
Research Advancement through
Cohort Cataloguing and Harmonization

## Harmonized variable



28%
59%
13%

- Yes
- No
- Missing

# Harmonization potential across studies: CPTP Health and Risk Factor Questionnaire



Health and Risk Factor Questionnaire – Harmonized variables

Variables status by datasets (N=935)

# Estimate quality of harmonized dataset

## Understand the potential and limitations of the harmonized dataset

|  | N | Mean |
|---|---|---|
| ATL_Measurements | 4872 | 130,88 |
| ATL_Measurements | 22703 | 88,01 |
| ATP_Measurements | 29347 | 89,8 |
| ATP_Measurements | 1149 | 90,2 |
| BCGP_Measurements | 16363 | 133,21 |
| CAG_Measurements | 19992 | 126,83 |
| LSC_Measurements | 649 | 129,78 |
| Pilot_Measurements | 7970 | 131,3 |

## Sitting height



Canadian Partnership for Tomorrow Project

# Missing values: CPTP data collection mode



Proportion of participants with missing values for more than 5% according to data collection mode and basic demographic characteristics

N

**Data collection mode**
Paper form (manual entry) — 18614
Paper form (teleform) — 56263
Electronic form (onsite computer) — 41518
Electronic form (web-based) — 190622

**Sex**
Female — 189204
Male — 117813

**Age Category**
<45 — 80025
45-54 — 94324
55-64 — 92583
>64 — 40085

**Overall** — 307017

0.05    0.10    0.15    0.20

Proportion missing values

# Harmonization process: MINDMAP variable profiles

# Disseminate and preserve harmonization products

Ensure transparency and leverage usage of harmonized data

# A data portal application used to describe central data and manage data access requests

## Data portal

# Document the overall harmonization process

**cmaj**OPEN

## Harmonization of the Health and Risk Factor Questionnaire data of the Canadian Partnership for Tomorrow Project: a descriptive analysis

Isabel Fortier PhD, Nataliya Dragieva MSc, Matilda Saliba PhD, Camille Craig MSc, Paula J. Robson PhD; with the Canadian Partnership for Tomorrow Project's scientific directors and the Harmonization Standing Committee*

### Abstract

**Background:** The Canadian Partnership for Tomorrow Project is a multistudy platform integrating the British Columbia Generations Project, Alberta's Tomorrow Project, the Ontario Health Study, CARTaGENE (Quebec) and the Atlantic Partnership for Tomorrow's Health. This paper describes the process used to harmonize the Health and Risk Factor Questionnaire data and provides an overview of the key information required to properly use the core data set generated.

**Methods:** This is a descriptive analysis of the harmonization process that was developed on the basis of the Maelstrom Research guidelines for retrospective harmonization. Core variables (DataSchema) to be generated across cohorts were defined and the potential for cohort-specific data sets to generate the DataSchema variables was assessed. Where relevant, algorithms were developed and applied to process cohort-specific data into the DataSchema format, and information to be provided to data users was documented.

**Results:** The Health and Risk Factor Questionnaire DataSchema (version 2.0, October 2017) comprised 694 variables. The assessment of harmonization potential for the variables over 12 cohort-specific data sets resulted in 6799 (81.6%) of the variables being considered as harmonizable. A total of 307 017 participants were included in the harmonized data set. Through the cohort data portal, researchers can find information about the definitions of variables, harmonization potential, algorithms applied to generate harmonized variables and participant distributions.

**Interpretation:** The harmonization process enabled the creation of a unique data set including data on health and risk factors from over 307 000 Canadians. These data, in combination with complementary data sets, can be used to investigate the impact of biological, environmental and behavioural factors on cancer and chronic diseases.

# Maelstrom Research process for rigorous data harmonization



Iterative Harmonization Steps
- Step 0: Define the research questions, objectives and protocol
- Step 1: Assemble information and select studies
- Step 2: Define variables and evaluate harmonization potential
- Step 3: Process data
- Step 4: Estimate quality of the harmonized dataset(s) generated
- Step 5: Disseminate and preserve final harmonization products

Retrospective data harmonization offers many benefits but is necessarily challenging.

Need a general systematic process that can be adapted to each initiative.

Applying systematic approach to ensure proper quality checks and documentation throughout is critical for assessing and interpreting results.

**F**indable **A**ccessible **I**nteroperable **R**eusable

# THANK YOU !

**www.maelstrom-research.org**

## Funding and support:

Centre universitaire de santé McGill
Institut de recherche

McGill University Health Centre
Research Institute

MUHC
McGILL UNIVERSITY HEALTH CENTRE
FOUNDATION

Économie, Innovation et Exportations
Québec

INNOVATION.CA
CANADA FOUNDATION FOR INNOVATION | FONDATION CANADIENNE POUR L'INNOVATION

CANADIAN PARTNERSHIP AGAINST CANCER
PARTENARIAT CANADIEN CONTRE LE CANCER

SEVENTH FRAMEWORK PROGRAMME

EUROPEAN COMMISSION

CIHR IRSC
Canadian Institutes of Health Research | Instituts de recherche en santé du Canada

National Institute on Aging

canarie @25

Canadian Partnership for Tomorrow Project

ReACH
Research Advancement through Cohort Cataloguing and Harmonization

MINDMAP
Promoting mental well-being and healthy ageing in cities

## Our numbers continue to grow

| | Networks | 18 |
| Individual Studies | 204 |
| Individual Studies with Variables | 122 |
| Individual Study Variables | 933,144 |

# Population-based cohort studies



**Study Timeline**

Each colour in the timeline graph below represents a separate Study Population, while each segment in the graph represents a separate Data Collection Event. Clicking on a segment gives more detailed information on a Data Collection Event.

- HAPIEE - Russia
- HAPIEE - Poland
- HAPIEE - Czech Republic
- HAPIEE - Lithuania

# A_DIS_ARTHRITIS_EVER

**Variable description**

## Overview

| Label | Lifetime occurrence of arthritis |
|---|---|
| Description | Occurrence of arthritis at any point during the life of the participant. |
| Dataset | Health and Risk Factor Questionnaire |
| Value Type | integer |

## Classification

| Areas of Information | |
|---|---|
| Diseases | Musculoskeletal system and connective tissue (M00-M99) |

## Categories

| Name | Label | Mis... |
|---|---|---|
| 0 | Never had arthritis | |
| 1 | Ever had arthritis | |
| 2 | Presumed - Never had arthritis | |

## Statistics

Cumulative summary of all studies:

| Value | Frequency |
|---|---|
| **Valid Values** | |

## Statistics

Cumulative summary of all studies:

| Value | Frequency |
|---|---|
| **Valid Values** | |
| 0 Never had arthritis | 145717 72.0% (73.8%) |
| 1 Ever had arthritis | 49809 24.6% (25.2%) |
| 2 Presumed - Never had arthritis | 1863 0.9% (0.9%) |
| Subtotal | 197389 97.5% |
| **Other Values** | |
| Missing | 5013 2.5% (100.0%) |
| Subtotal | 5013 2.5% |
| Total | 202402 |

Valid values frequencies

- 0 (Never had arthritis)
- 1 (Ever had arthritis)
- 2 (Presumed - Never had arthritis)

*Summary statistics (real time)*

Variable search