

The CORESIDENCE Database:

National and Subnational Data on Household Size and Composition Around the World, 1964-2021

Juan Galeano¹, Albert Esteve¹, Anna Turu², Joan García-Roman², Federica Becca², Huifen Fang², Maria Louise Christine Pohl², Rita Trias Prats²

1: Universitat Autònoma de Barcelona / Centre d'Estudis Demogràfics

2: Centre d'Estudis Demogràfics

Corresponding autor: Juan Galeano (jgaleano@ced.uab.es)

Abstract

The CORESIDENCE Database (CoDB) represents a significant advancement in the field of family studies, addressing existing data gaps and facilitating comprehensive analysis of households' composition and living arrangements at the national and subnational levels. This article introduces the CoDB, developed for the ERC project Intergenerational Coresidence in Global Perspective: Dimensions of Change. The database draws on global-scale individual microdata from four main repositories and national household surveys, encompassing over 150 million individual records representing more than 98% of the world's population. The CoDB provides datasets at the national, subnational, and subnational-harmonized levels, covering 156 countries, 3950 regions, and 1511 harmonized regions. It includes 146 indicators on household composition and family arrangements, allowing researchers to explore intergenerational co-residence patterns, gender dynamics within households, and longitudinal trends in living arrangements. The CoDB fills an important gap in comparative household studies, enabling researchers to undertake ground breaking research at both macro and micro levels, ultimately fostering a deeper understanding of the complex dynamics of family structures and living arrangements worldwide.

Background & Summary

Households represent the most fundamental unit of human organization. They play a crucial role in child-rearing, elderly care, resource allocation, and shaping gender roles. While the composition of households is primarily based on familial bonds, practices vary significantly across societies. Factors such as demographics, economics, and social norms influence variations in household size and composition. Consequently, households have significant implications for social reproduction, urbanization, housing demands, and consumption. Despite their significance, the availability of household level data at the global scale is underdeveloped and could be substantially expanded thanks to the increasing availability of household level microdata. To bridge this information gap, the Coresidence database (CoDB) provides access to 146 harmonized indicators on household size and composition for 156 countries, 3950 subnational areas, and 58 data points in time. Compared to the United Nations database on Household Size and Composition (<https://www.un.org/development/desa/pd/data/household-size-and-composition>), CoDB complements, updates, and introduces new features. Firstly, CoDB exclusively includes data from countries where microdata is accessible to researchers.

46 While this slightly limits the number of countries compared to the United Nations database, it
47 significantly expands analytical possibilities. Secondly, by leveraging microdata, CoDB
48 broadens the number of indicators on household size and composition. A total of 146 indicators
49 have been calculated. Thirdly, CoDB offers open-source code in R, allowing users to observe
50 how the microdata has been processed and indicators have been built. This ensures replicability
51 and empowers users to create new indicators. Finally, CoDB provides detail at the subnational
52 level.

53 CoDB has been developed within the project “Intergenerational Coresidence in Global
54 Perspective: Dimensions of Change (CORESIDENCE)¹”, funded by the European Research
55 Council. The available indicators in CoDB have been calculated from individual microdata
56 samples from four large data repositories of international microdata, supplemented by national
57 household surveys. All included samples allow grouping individuals into households and
58 examining the relationships established among their members. Additionally, they provide basic
59 sociodemographic information about household members, including age, sex, and marital
60 status. With all the samples combined, the original microdata database contains more than 150
61 million individual records, representing more than 98% of the world's population and spanning
62 from the 1960's to the present. The 146 indicators contained in CoDB represent different
63 aggregations of the original microdata, both by country and subnational areas. Within each
64 country, subnational areas have been harmonized to facilitate the study of change over time. As
65 a final output, CoDB consists of three datasets: The National dataset contains 156 countries, the
66 Subnational dataset contains 3950 subnational areas, and the Subnational harmonized dataset
67 contains 1511 subnational areas for the period 1964 to 2021, and it provides 146 indicators on
68 household composition and family arrangements across the world.

69

70 **1. Methods**

71 **1.1 Overview**

72 Figure 1 provides a schematic overview of the entire process of creating CoDB, starting with
73 data acquisition, and followed by data processing, harmonization, indicator's construction,
74 output datasets, and external validation. CoDB draws on four main repositories of global-scale
75 individual microdata (Fig.1): The International Integrated Public Use Microdata Series
76 (IPUMS-I), the Demographic Health Surveys (DHS), the Multiple Indicator Cluster Surveys
77 (MICS), and the European Union Labor Force Survey (EU-LFS). Additionally, CoDB includes
78 country-specific surveys and censuses not available in any of the previous repositories, such as
79 the EU Statistics on Income and Living Conditions (EU-SILC) surveys, the Income and Labour
80 Dynamics in Australia (HILDA) surveys, the Household Income and Expenditure Survey
81 (HIES) for South Korea and the China Family Panel Studies (CFPS). Contextual indicators

¹ HE-ERC-2021-AdG-GA No 101052787-CORESIDENCE

82 come directly from various UN datasets, specifically the United Nations World Population
83 Prospects¹ (UNWPP), the United Nations Development Programme² (UNDP) and from gridded
84 data of the Human Development Index from Kumm et al. (2008)³.

85 All data cleaning, processing, harmonization and aggregation were performed in R⁴. For the
86 Subnational harmonized dataset, the use of QGIS⁵ was additionally required. All the coded used
87 in the construction of CoDB is available in the GitHub repository of this project (see section
88 **Code Availability**).

89 The output data of CoDB includes three datasets: National, Subnational, and Subnational
90 harmonized.

91 The National dataset includes 791 country-year samples from 155 countries (Fig 2). Figure 2
92 provides an overview of the number of countries included in the database and the data available
93 for each of them. For each sample, Figure 2 informs about the source of reference and about
94 what type of subnational data is available per sample. The National dataset contains 146
95 indicators on household size and composition worldwide for over 60 years. The selected
96 indicators provide information on the size and composition of households. Regarding
97 composition, details are provided on the age, relationship to the reference person, type of
98 household (e.g. unipersonal, nuclear, extended), and sex of the reference person in the
99 households (see section 1.4). These are standard measures in household research using similar
100 data sources⁶. This dataset incorporates an additional set of 20 contextual indicators obtained
101 from the UNWPP and the UNDP. These additional indicators provide information on
102 population size, life expectancy by sex, fertility rates, and the human development index for
103 each country in a given year.

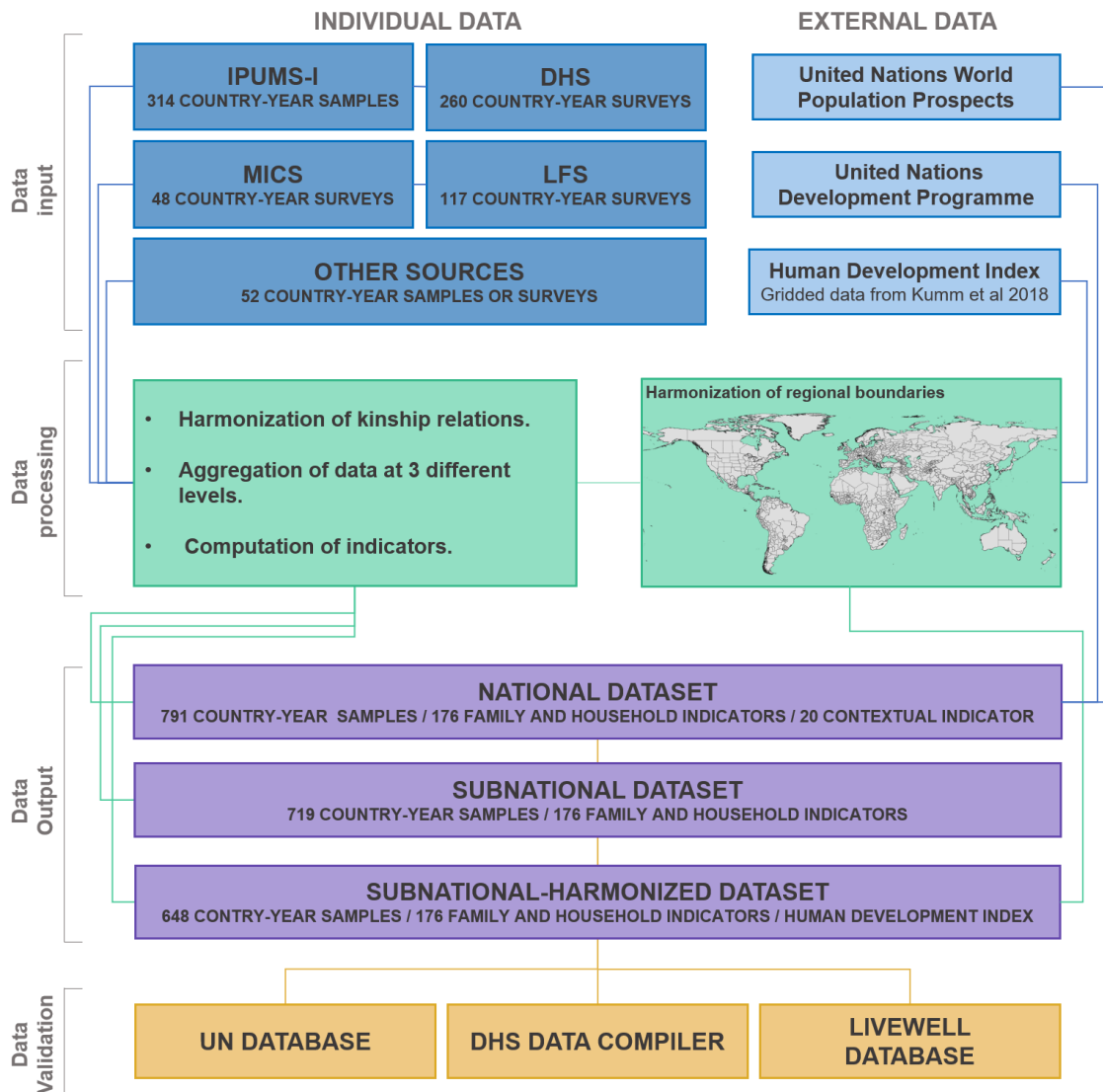
104 The Subnational dataset includes 719 country-year samples covering 149 countries and 3,950
105 unique regions. The 146 indicators were calculated based on the major administrative unit in
106 which households were enumerated in each of the primary data sources. Out of the original 791
107 samples, 72 were not included in this dataset due to the absence of territorial disaggregation
108 information (see Figure 2).

109 Last, the Subnational Harmonized dataset consists of 648 country-year samples from 138
110 countries and 1,511 unique regions. To ensure consistency and minimize repetition, only
111 countries for which we had more than one sample and regions could be harmonized over time
112 were included. As a result, the Subnational Harmonized dataset covers 82% of the original
113 samples. Figure 3 shows the regional breakdown available in the Subnational Harmonized
114 dataset (Figure 3 in green). The regions marked in green are present in this database. For
115 countries with only on sample (e.g. Canada), it is necessary to retrieve the data from the
116 Subnational dataset. Regarding the indicators, the same 146 are available for all these regions.
117 The harmonization of geographic boundaries is explained in sub-section **Harmonization of**
118 **regional subnational boundaries**. Regarding contextual data, indicators such as life

119 expectancy or fertility are not available at this scale. However, data from the Human
 120 Development Index, extracted through from Kumm et al. (2008)³, has been included (see section
 121 **Contextual indicators**).

122 In addition to the three datasets, the CoDB also provides a spatial file with the boundaries of
 123 the subnational harmonized regions either as *sf* object or a multi-polygon geopackage (see
 124 section **Data Records**). For the production of the spatial file, we relied on the already
 125 harmonized geographies provided by the IPUMS international, the DHS Spatial Repository⁷,
 126 the work done by the LiveWell project⁸ for harmonizing DHS boundaries and the Database of
 127 Global Administrative Areas (GADM)⁹.

128 To ensure the accuracy and reliability of the CoDB, we validated our database by comparing
 129 the results of a selected set of indicators from the three datasets with corresponding data from
 130 reputable sources such as the UN database on Household Size and Composition¹⁰, the DHS
 131 STAT compiler¹¹ and the LiveWell project (see section 3 on **Technical Data Validation**).



132
 133

Fig.1: Flowchart representing the different stages to build the CoDB

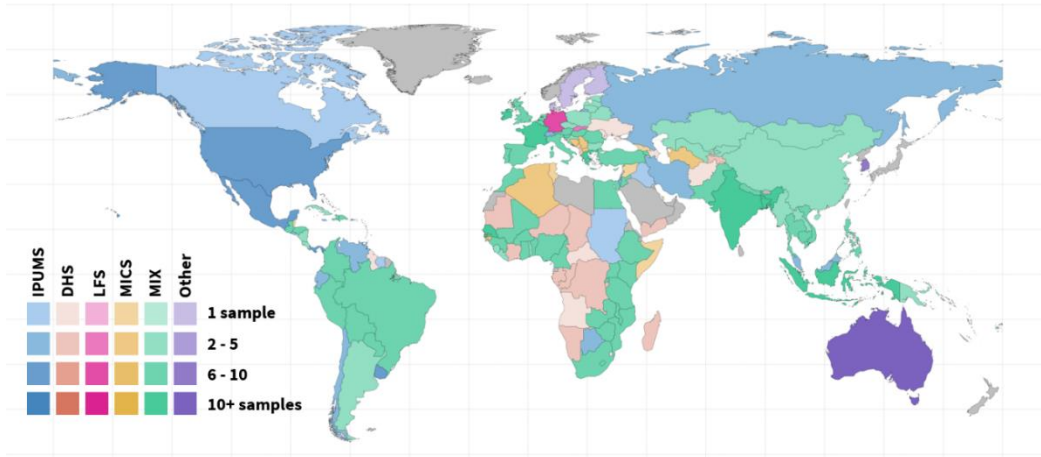


Fig. 2: Country coverage by number of samples available of the CoDB



Fig. 3: Availability of samples by country, year, and source in the CoDB

141 1.2 Data Sources

142 The CoDB is a comprehensive source of information on household structure and family
143 composition at national and subnational levels. The database draws on four major repositories
144 of individual microdata on a global scale, along with country-specific surveys and censuses
145 from countries not included in those repositories. Additionally, we employed external data
146 sources to provide a set of contextual (demographic and socioeconomic) indicators.

147 The first source of individual microdata for the CoDB is the International Integrated Public Use
148 Microdata Series (IPUMS-I)¹², consisting of 314 census samples from 94 countries
149 (<https://international.ipums.org>). The IPUMS International project is a global initiative that
150 aims to collect, preserve, harmonize, and distribute census microdata from countries worldwide.
151 In all cases, except for countries with fewer households in a specific year, a sample of 20,000
152 households from the original microdata was randomly selected to build the CoDB indicators.
153 This was done to minimize data storage and speed processing, but users can rebuild these
154 indicators with larger samples using the same source. In the validation process (see section 3),
155 we show that our estimates are consistent with those of the United Nations based on indicators
156 that are available in both UN and the CoDB sources.

157 The second source of individual microdata for the CoDB is the Demographic Health Surveys
158 (DHS)¹³ (<https://dhsprogram.com/data/>), which have been collecting demographic and health
159 information for low- and middle-income countries since 1986. A total of 260 samples from 75
160 countries were retrieved. DHS surveys rely on a two-stage cluster sampling design that ensures
161 the representativeness of the data at the national and subnational level.

162 To expand the coverage of the CoDB beyond the countries and years included in the two
163 previous repositories, 49 additional samples from 33 countries were included from the Multiple
164 Indicator Cluster Surveys (MICS) program¹⁴ (<https://mics.unicef.org/surveys>), which collects
165 data related to key indicators of health, education, child protection, and water and sanitation.
166 MICS surveys are designed to collect data at both national and subnational levels. The data is
167 publicly available and has been widely utilized for studying family structures and change in a
168 variety of countries.

169 Microdata from 117 samples of the European Labour Force Survey (EU-LFS)¹⁵ were used to
170 complement the information available on European countries from IPUMS
171 (<https://ec.europa.eu/eurostat/web/microdata/european-union-labour-force-survey>). The EU-
172 LFS is a large household sample survey on the labour force participation of the 15-year and
173 older population, also collecting information on all members of the household surveyed, as well
174 as the kinship relations among them. As LFS collects data on a quarterly basis, samples included
175 in the CoDB correspond to the yearly samples to ensure consistency with the specific time frame
176 for which the data was downloaded.

177 The CoDB includes information from country-specific surveys and censuses for countries
178 and/or years not present in the previous repositories. This includes: 22 samples of the EU
179 Statistics on Income and Living Conditions (EU-SILC) survey¹⁶ for the year 2021
180 (<https://ec.europa.eu/eurostat/web/microdata/eu-silc>), 21 samples from the Household, Income
181 and Labour Dynamics in Australia (HILDA)¹⁷ survey between 2001 and 2021
182 (<https://melbourneinstitute.unimelb.edu.au/hilda>), 9 samples from the South Korean Census
183 (<http://kosis.kr/eng/>) covering the period 1970-2010 and 2 sample from the China Family Panel
184 Studies (CFPS)¹⁸ for the years 2010 and 2018 (<https://www.issp.pku.edu.cn/cfps/en/>).
185 Last, for the set of contextual socio-demographic indicators provided in the National dataset of
186 the CoDB we used data from the United Nations World Population Prospects (UNWPP)
187 (<https://population.un.org/wpp/>) and the United Nations Development Programme (UNDP)
188 (<https://hdr.undp.org/data-center>). The UNWPP provides information on global population
189 trends, projections, and demographic indicators, whereas the UNDP focuses on promoting
190 human development globally. To get subnational estimates of the Human Development Index
191 for the Subnational Harmonized dataset, we utilized the HDI gridded dataset developed by
192 Kummu *et al.*¹⁹ (<https://datadryad.org/stash/dataset/doi:10.5061/dryad.dk1j0>).
193 The CoDB has been designed with a forward-looking perspective, poised to accommodate the
194 ongoing growth of its constituent data repositories. As the aforementioned data sources continue
195 to release new samples, the CoDB is primed to seamlessly integrate these additions, ensuring
196 its comprehensiveness over time.

197

198 **1.3 Harmonization processes**

199 **1.3.1 Harmonization of household interrelationship variables**

200 In the construction of the CoDB, a crucial step was the harmonization of relationships among
201 household members from diverse data sources. Most of these relationships involve a certain
202 degree of kinship, but the amount of detail varies widely. We followed the IPUMS-I
203 harmonization coding scheme to harmonize intrahousehold relationships for the other sources.
204 The IPUMS-I samples include a harmonized variable called "relate" which captures 75 distinct
205 types of relationships (or their absence) with respect to the reference person of the household,
206 often named the household head. Not all types of relationships are present in every sample. The
207 detailed classification of types is grouped into six categories: Head, Spouse/Partner, Child,
208 Other relative, Other non-relative, and Other relative or non-relative.

209 In the case of DHS and MICS surveys, before establishing equivalences with the IPUMS-I
210 categories, an additional step was necessary to harmonize the data internally, as the same
211 kinship category was recorded in slightly different ways across different surveys. For instance,
212 variations like "brother-in-law or sister-in-law," "brother-in-law/sister-in-law," and "brother-in-
213 law/sister-in-law" were observed. Through the internal harmonization process, these variations

214 were consolidated into 24 distinct categories for DHSs and 39 categories for MICSS. These
215 categories were then aligned with the corresponding ones from the IPUMS-I samples.
216 The EU Labour Force Surveys (LFS) only capture 6 types of relations, but crucially for the
217 purpose of this project the type: ‘Ascendant relative of reference person (or of his/her spouse or
218 cohabiting partner)’ is included. The EU statistics on income and living conditions (EU-SILC)
219 offer a broader perspective, encompassing 19 different types of relations to the head. In the case
220 of South Korea, the census samples provided by the National Office of Statistics provide a wider
221 range of recorded relations to the head, varying between 13 and 38 depending on the specific
222 year.

223 In the case of the Australian data, the absence of a designated head or reference person made
224 the procedure more complex. However, leveraging the available information on the total income
225 of household members, we employed a specific criterion to define the head of the household.
226 The person with the highest total income was identified as the head, ensuring consistency
227 between the surveys provided by the National Statistical Institute of Australia. In the rare
228 instances where two members had exactly the same income, the older person was designated as
229 the head. Additionally, we had to re-code all the relations within the household as they were
230 originally recorded from the perspective of the individual (ego) to all other members of the
231 household, ensuring a consistent and standardized representation of kinship relations.

232 The Chinese Family Panel Survey (CFPS) also provides the relations between household
233 members as a matrix of "all versus all" type. The source code for the re-coding of relations can
234 be accessed and downloaded from the CORESIDENCE project's GitHub repository (see section
235 **Code Availability**). In total, 17 types of relations to the head were defined and aligned with
236 IPUMS-I.

237

238 **1.3.2 Harmonization of regional subnational boundaries**

239 One of the key and original features of the CoDB with respect to other databases is the provision
240 of subnational data in the Subnational and the Subnational Harmonized datasets. For this latter
241 one, geographical boundaries were harmonized to facilitate the study of change over time.

242 For harmonizing the subnational regions we relied on four major sources of spatial data
243 information: the spatially harmonized first-level geography from IPUMS International
244 (https://international.ipums.org/international/gis_harmonized_1st.shtml), the work done by the
245 LiveWell project²⁰ for the harmonization of DHS boundaries, the DHS Spatial Repository
246 (<http://spatialdata.dhsprogram.com/home/>), and the Database of Global Administrative Areas,
247 GADM (<https://gadm.org/>)

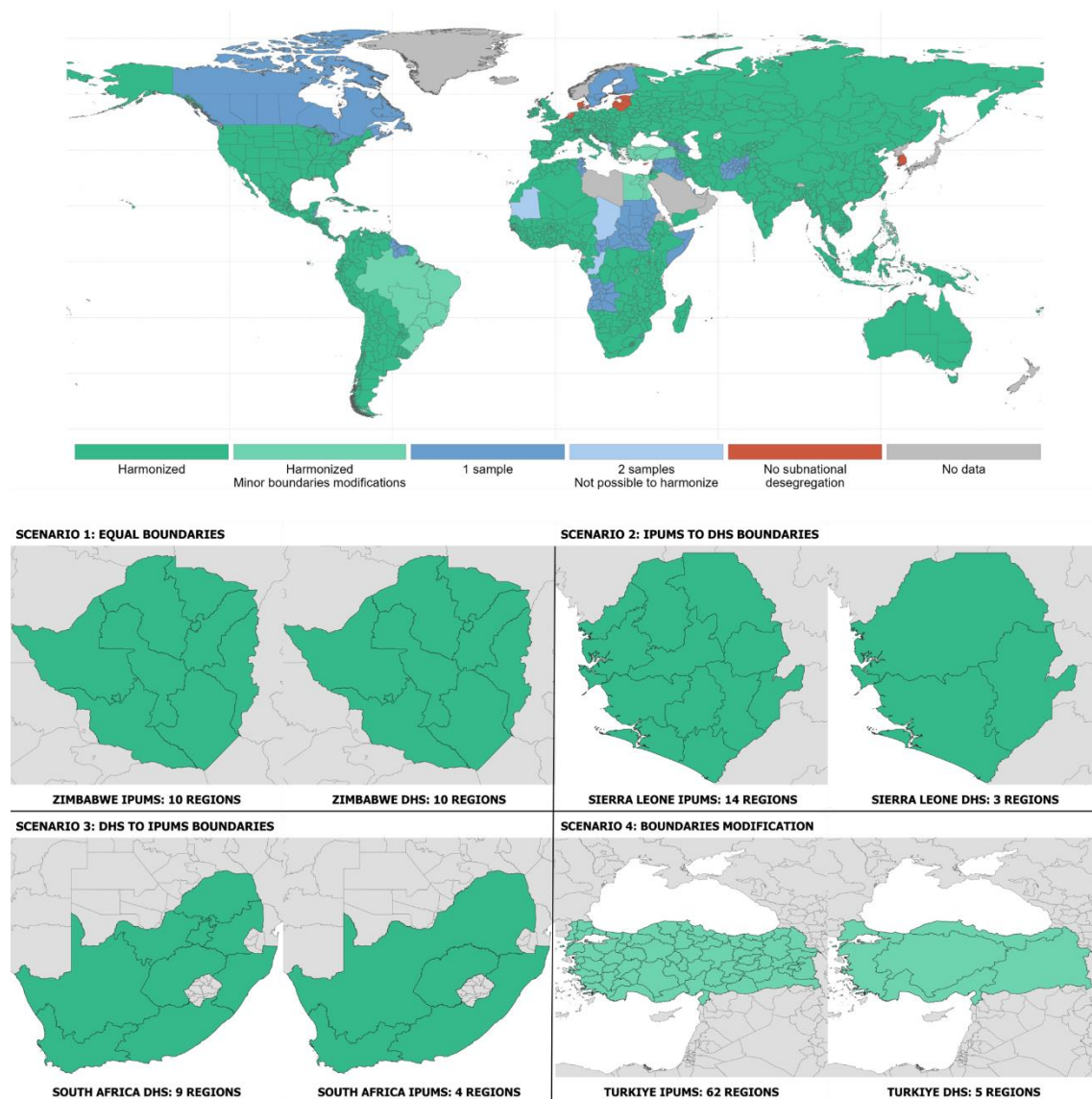
248 The harmonization process involved multiple steps. First, we selected countries with at least
249 two data samples. Second, we identified the smallest common spatial denominator to allow for
250 comparisons over time. Third, we categorized the selected countries based on whether all the

251 data samples originated from the same data source or not. When all samples originated from the
252 same source, we encountered two distinct scenarios. Firstly, if the samples were obtained from
253 IPUMS, which already had a pre-existing harmonized subnational division and identification
254 system, no further harmonization was needed. Secondly, when the data comes from sources
255 other than IPUMS, it has been necessary to harmonize administrative boundaries in some
256 countries. When all samples from a country come from the DHS, we assigned an IPUMS-like
257 ID to each of the harmonized regions (6 digits where the 3 first digits are the ISO numeric code
258 of the country), following the process developed by the LiveWell project. The same process
259 was applied when all samples were obtained from the LFS or SILC. For the 21 samples from
260 Australian HILDA data, the subnational regions were already harmonized and we only assigned
261 a new ID to each region.

262 When dealing with samples from different data repositories for a given country, the
263 harmonization process became more complex. Where there was a perfect match between
264 sources, such as Zimbabwe, the harmonization process was straightforward (Fig 4, scenario 1).
265 In this country, both the IPUMS-I and DHS samples used the same regional breakdown of the
266 country. In these instances, we used the GEOLEVEL1 IDs from IPUMS to harmonize the DHS
267 data. In other cases, for instance that of South Africa or Sierra Leone, the harmonization process
268 involved the aggregation of regions (Fig 4, scenario 2 and 3). When aggregating data from DHS
269 to IPUMS, we retained the region IDs provided by IPUMS. Conversely, when aggregating data
270 from IPUMS to DHS, we created new IDs for both sources, as it was the case for countries with
271 samples from the LFS and SILC repositories. The last scenario we encountered involved
272 making slight modifications to regional boundaries between sources (Fig 4, scenario 4). This
273 was the case for samples from Turkey, Philippines, Egypt, and Brazil, and the affected regions
274 are listed in the harmonization table provided within the CoDB.

275 The R code for re-coding the individual data of each sample to the harmonized regions can be
276 found, consulted and downloaded from the GITHUB repository of the CORESIDENCE project
277 (see section **Code availability**).

278



279
280 **Fig. 4: Harmonized Subnational coverage of the CoDB**

281
282 **1.4 Construction of indicators on household composition and living arrangements.**

283 From the original microdata presented in section 1.2., we have calculated 146 indicators on
284 household composition and family arrangements globally over a span of 60 years
285 (Supplementary Table 1). To generate these indicators, we aggregated the individual data to the
286 national, subnational, and harmonized subnational levels across our three datasets. The data was
287 weighted using the individual weights provided by each sample.

288 The indicators provided in the CoDB can be grouped into four categories: (i) indicators related
289 to size and age composition of households, (ii) indicators derived from the relation of family
290 members to the person defined as the head of the household, (iii) indicators related to
291 household's typology and (iv) indicators related to household headship. Before computing the
292 indicators included in the CoDB, the population was weighted by the relevant survey weight
293 (household or individual weight), ensuring representativeness with respect to the underlying

294 population. This comprehensive set of indicators offers a rich resource for studying household
295 composition and living arrangements across different contexts and time periods.

296 *(i) Indicators related to size and age composition of households:* The first set of indicators
297 (HS01 to HS11) focuses on the relative distribution of households by size (ranging from 1 to
298 10 persons and 11 or more persons). To further explore the composition of households,
299 indicators HS12 to HS14 provide information on the proportion of households with 2-3, 4-5, or
300 6 or more persons as computed in the UN database on Household Size and Composition; thus,
301 enabling external validation of our own computations (see section on **Technical Validation**).
302 Indicators HS15 and HS16 compute the proportion of households with at least one person aged
303 0-4 or 65 or more years old respectively. In the CoDB, this is presented as an average (HS17)
304 at the national, subnational, or subnational harmonized level. Indicators HS18 to HS21 provide
305 additional insights into the average number of persons in households, categorized by age
306 groups: 0-4 years, below 18 years, above 18 years, and 65 years or older. These indicators shed
307 light on the age distribution within households. Moreover, indicators HS22 to HS30 provide
308 information on the average number of persons in households within 10-year age intervals. This
309 allows for a more detailed understanding of the age composition of households.

310 *(ii) Indicators derived for the relation of family members to the person defined as the head or*
311 *of the household:* These indicators offer insights into the structure and dynamics of family
312 relationships. The first group of indicators (HR01 to HR06) provides information about the
313 average number of heads, spouses, children, other relatives, and non-relatives in the household.
314 These indicators help us understand the composition of the household in terms of these specific
315 family relationships. Moreover, this information is further disaggregated based on the size of
316 the household, specifically households with 2 to 5 people. Indicators HR07 to HR30 present the
317 average number of heads, spouses, children, other relatives, and non-relatives in households of
318 this size range. This allows for a more detailed analysis of the relationship dynamics within
319 different household configurations. By examining these indicators, we can improve our
320 understanding of the social structure and interdependencies among family members within
321 households of various sizes. This information contributes to a deeper understanding of family
322 dynamics and relationships within different contexts.

323 *(iii) Indicators related to household's typology:* Indicators related to household typology in the
324 CoDB offer valuable insights into the diverse forms and compositions of households across
325 different contexts. To ensure comparability and overcome variations in the types of kinship
326 relations recorded in the different data sources, we computed indicators based on two distinct
327 typologies.

328 The first typology, developed by the CORESIDENCE team, consists of eight categories:

329 1. Unipersonal households.

- 330 2. Nuclear households: consisting of a head, a spouse, and their children, or a head and
331 their children.
- 332 3. Nuclear households with additional relatives.
- 333 4. Nuclear households with non-relatives.
- 334 5. Nuclear households with both relatives and non-relatives.
- 335 6. Other relative households.
- 336 7. Other non-relative households.
- 337 8. Other households with a combination of relatives and non-relatives.

338 The second typology, based on the work of John Bongaarts⁶ for developing countries in the
339 1990s, comprises five categories:

- 340 1. Unipersonal households.
- 341 2. Nuclear households: consisting of a head, a spouse, and their children, or a head and
342 their children.
- 343 3. Stem family additions: including parents or grandchildren of the head.
- 344 4. Other family households: encompassing other relatives of the head.
- 345 5. Other non-family households: comprising individuals not related to the head.

346 Using these two sets of typologies, the indicators (HT01 to HT31) provide information on the
347 proportion and average size of each household type. These indicators shed light on the
348 prevalence and characteristics of various household types, contributing to a deeper
349 understanding of household structures and arrangements across different populations and time
350 periods within the CoDB.

351 *(iv) Indicators related to household headship:* Indicators related to household headship in the
352 CoDB capture important dimensions covered in the previous sets of indicators, such as
353 proportions of n persons households, average sizes, and typologies. However, they specifically
354 consider the gender dimension in relation to the household head (HH01 to HH56). These
355 indicators provide key information on the roles and dynamics of gender within households.
356 They shed light on the distribution of male-headed and female-headed households, offering a
357 deeper understanding of how gender influences household structures and arrangements. By
358 examining proportions, average sizes, and typologies of male-headed and female-headed
359 households, these indicators contribute to a comprehensive analysis of household composition
360 and dynamics, while considering gender dimension.

361

362 **1.4.1 Contextual indicators**

363 In addition to the household level indicators, the National and Subnational harmonized datasets
364 included in the CoDB provide contextual indicators. Within the National dataset, for each
365 country-year sample included in the CoDB, we provide population counts, total fertility rates

366 (TFR), and life expectancy by sex. In addition to demographic indicators, CoDB includes socio-
367 economic measures, such as the Human Development Index (HDI) and its components. The
368 HDI is a composite index that assesses the overall development and well-being of a country,
369 considering factors such as life expectancy, education, and income. The components of HDI
370 included in CoDB are: expected years of schooling, mean years of schooling, Gross Domestic
371 Income (GDI), and Gross National Income (GNI) per capita. The socio-economic indicators are
372 also divided by sex.

373 In the case of the Subnational Harmonized (SH) dataset, we utilized the HDI gridded dataset
374 developed by Kummu *et al.* (2008)³ to provide the Human Development Index (HDI) at the
375 subnational level for all the harmonized samples between 1990 and 2015 included in the CoDB.
376 This allowed us to capture the variations in development within countries at a more detailed
377 geographical level.

378 To calculate the average HDI values for each Subnational Harmonized region in our dataset,
379 we proceeded as follows. First, we transformed the gridded HDI data for each year into a spatial
380 points layer using the "*raster pixels to points*" function from the processing toolbox of QGIS.
381 Next, we clipped the spatial boundaries of our SH dataset with the points shapefile by joining
382 their attributes based on location. Finally, we summarized the joined data by the harmonized ID
383 and year and computed the mean HDI values using R. This process allowed us to provide the
384 HDI at the subnational level for 72.8% of the region-year entries of the SH dataset.

385 By including total fertility rates, life expectancy, and socio-economic indicators like the HDI,
386 the CoDB empowers researchers and policymakers to explore the demographic and socio-
387 economic landscapes of different countries and time periods in relation to changes in family
388 arrangements and households' composition. These indicators facilitate a deeper understanding
389 of population dynamics, thereby supporting evidence-based decision-making and policy
390 formulation.

391

392 **2. Data records**

393 The CoDB is hosted in Zenodo, an open-access digital repository that allows researchers,
394 scientists, and scholars from various disciplines to share and preserve their research outputs.
395 Zenodo is operated by CERN (European Organization for Nuclear Research) and supported by
396 various organizations, including the European Commission's OpenAIRE project. The CoDB is
397 hosted at the permanent DOI: <https://doi.org/10.5281/zenodo.8142652>. The repository is
398 composed of the following elements: a RData file named CORESIDENDE_DB containing the
399 CoDB in the form of a List. In R, a List object is a versatile data structure that can contain a
400 collection of different data types, including vectors, matrices, data frames, other lists, spatial
401 objects or even functions. It allows to store and organize heterogeneous data elements within a
402 single object. The CORESIDENDE_DB R-list object is composed of six elements:

403

- 404 1. NATIONAL: a data frame with the household composition and living arrangements
405 indicators at the national level.
- 406 2. SUBNATIONAL: a data frame with the household composition and living
407 arrangements indicators at the subnational level computed over the original subnational
408 division provided in each sample and data source.
- 409 3. SUBNATIONAL_HARMONIZED: a data frame with the household composition and
410 living arrangements indicators computed over the harmonized subnational regions.
- 411 4. SUBNATIONAL_BOUNDARIES_CORESIDENCE: a spatial data frame (a sf object)
412 with the boundary's delimitation of the subnational harmonized regions created for this
413 project.
- 414 5. CODEBOOK: a data frame with the complete list of indicators, their code names and
415 description.
- 416 6. HARMONIZATION_TABLE: a data frame with the full list of individual country-year
417 samples employed in this project and their state of inclusion in the 3 datasets composing
418 the CoDB.

419 Elements 1, 2, 3, 5 and 6 of the R-list are also provided as *csv* files under the same names.

420 Element 4, the harmonized boundaries, is at disposal as *gpkg* (Geopackage) file.

421

422 **3. Technical Validation**

423 To ensure the accuracy and reliability of the CoDB, we employ a two-stage validation process.

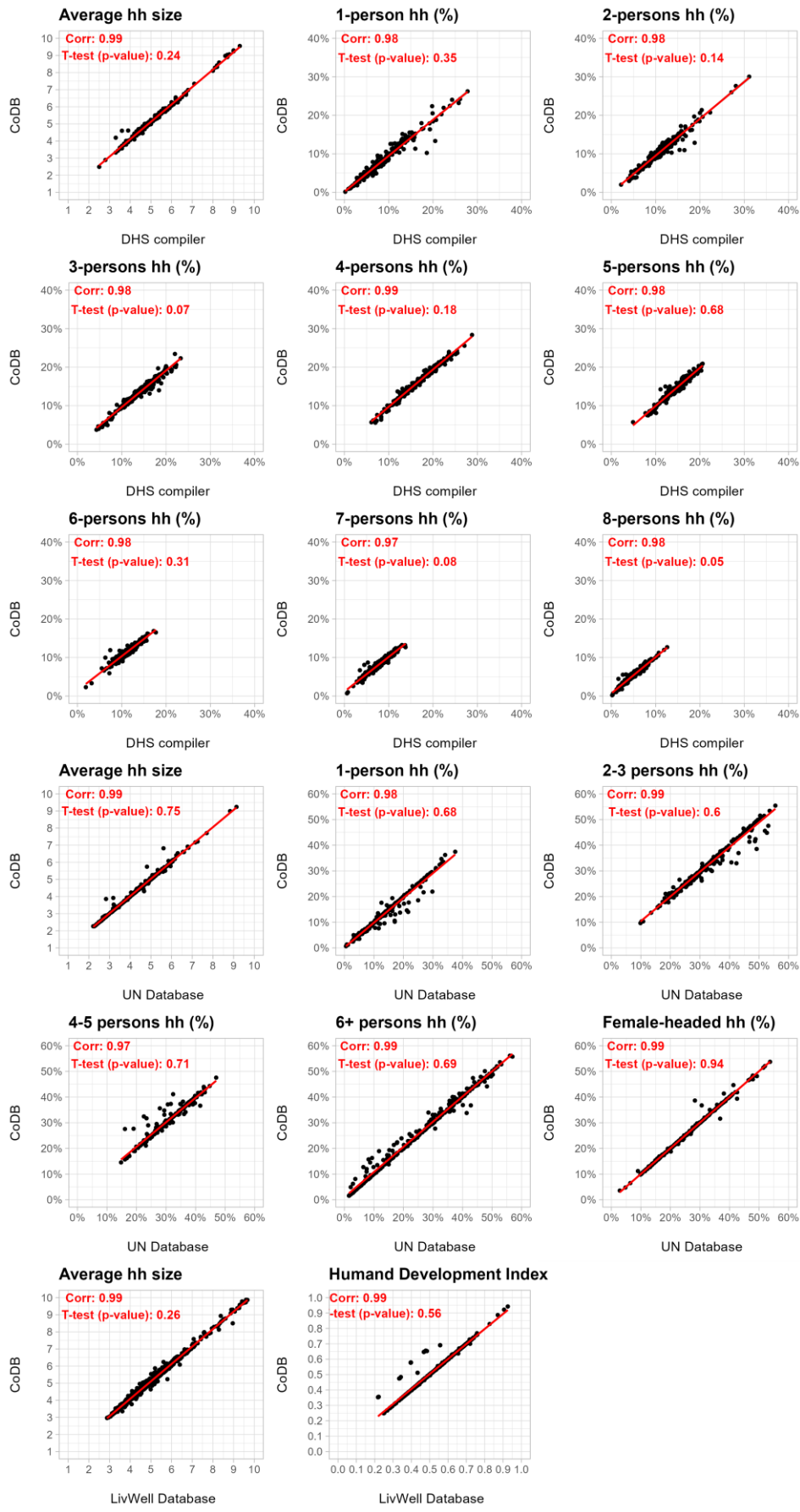
424 In the first stage, we validate our National dataset by comparing some of our indicators to those
425 from the DHS STAT compiler²¹ and the UN database on Household Size and Composition.

426 The DHS STAT compiler, developed by the DHS Program, is a user-friendly interface that
427 facilitates the exploration and visualization of indicators derived from DHS survey data at the
428 national and subnational levels. Complementing this, the United Nations (UN) database on
429 Household Size and Composition serves as a comprehensive repository that gathers data from
430 diverse sources to offer insights into the worldwide size and composition of households at the
431 national level. By harmonizing and standardizing the measurement and classification of
432 household characteristics, it enables comparisons and analysis across countries.

433 Among the indicators provided by the STAT compiler at the national level, there is a specific
434 set of nine indicators providing information on the average number of people per household and
435 the relative distribution of them by size, which allow us to compare 255 surveys from 74
436 countries. Using the UN database, we compared 269 samples from IPUMS and 14 surveys from
437 MICS, encompassing data from 91 countries, over a set of six indicators connected with the
438 same dimensions plus the share of female-headed households. Leveraging these tools, we assess
439 the consistency and alignment of our National dataset indicators with these two reputable

440 sources, ensuring the reliability and validity of our data. Overall, the correlation between the
441 country-level indicator of the CoDB and the ones from the STAT compiler and the UN database
442 is highly linear, suggesting a good fit of our computations (Figure 3). Additionally, we
443 computed an equal variance T-test for each of the selected indicators. The p-values, greater than
444 the common significance level of 0.05, suggest that the observed difference in means is likely
445 due to random variation, primarily associated with the data cleaning and processing steps. This
446 indicates that the disparities between the compared databases are more likely a result of data
447 handling rather than genuine differences in means.

448 In the second stage, we validate the Subnational Harmonized dataset using data from the
449 LiveWell project and the subnational human development database²². These additional sources
450 of data enable us to cross-reference and corroborate the harmonized indicators at the subnational
451 level. To validate the Subnational Harmonized dataset, we conducted the same analysis as for
452 the National Database using three directly comparable indicators sourced from the LiveWell
453 database. This validation process encompassed 1485 region-year entries, accounting for
454 approximately 20.4% of our dataset. This validation process is crucial to ensure the robustness,
455 accuracy, and overall quality of our subnational harmonized dataset, as well as to support its
456 usefulness for demographic analysis and/or to inform policy decision-making.



457

458

Fig. 5: External validation of the CoDB

459 **4. Code availability**

460 The processing steps to build the three datasets composing the CoDB were carried out in R,
461 utilizing the libraries tidyverse²³, haven²⁴, labelled²⁵, and tibble²⁶. All the code is available on the
462 GitHub repository of this project: <https://github.com/JuanGaleano/CORESIDENCE>

463

464 **5. Acknowledgements**

465 We thank PhD. Ginevra Floridi for her thoughtful comments and suggestion on the first version
466 of this manuscript.

467

468 **6. Funding**

469 European Research Council (ERC). Advanced Grant. Reference: HE-ERC-2021-AdG-GA No
470 101052787-CORESIDENCE

471

472 **7. Author contributions**

473 Albert Esteve conceived the project and is the principal investigator of the CORESIDENCE
474 project.

475 Juan Galeano, designed the analytic strategy and processed the data, wrote the R code for
476 building the CoDB, produced the subnational harmonized spatial file, wrote the initial
477 manuscript and prepared the figures for it.

478 Joan García-Roman compiled the row data necessary for building the CoDB.

479 Anna Turu compiled and processed the data.

480 Federica Becca tested the outcome data of the CoDB for Latin-American countries

481 Huifen Fang tested the outcome data of the CoDB for Asian countries and contributed in the
482 production of the subnational harmonized spatial file.

483 Maria Louise Christine Pohl tested the outcome data of the CoDB for African countries.

484 Rita Trias Prats tested the outcome data of the CoDB for European countries.

485

486 **8. Competing interests**

487 The authors declare no competing interests.

488

489

490

491

492

493

494

495

-
- ¹ United Nations, Department of Economic and Social Affairs, Population Division. World Population Prospects 2022. <https://population.un.org/wpp/2022/> (2022).
- ² United Nations Development Programme. Human Development Index (HDI). <http://hdr.undp.org/en/content/human-development-index-hdi> (2022)
- ³ Kummu, M., Taka, M. & Guillaume, J. H. A. Gridded global datasets for gross domestic product and human development index over 1990–2015. *Sci. Data* 5, 180004 (2018). <https://doi.org/10.1038/sdata.2018.4>
- ⁴ R Core Team. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>. (2021).
- ⁵ QGIS Development Team. QGIS Geographic Information System. Open Source Geospatial Foundation. <https://qgis.org> (2023).
- ⁶ Bongaarts, John. “Household Size and Composition in the Developing World in the 1990s.” *Population Studies* 55, no. 3 (2001): 263–79. <http://www.jstor.org/stable/3092865>.
- ⁷ ICF. The DHS program spatial data repository. Funded by USAID. (2020).
- ⁸ Belmin, C., Hoffmann, R., Elkasabi, M. et al. LivWell: a sub-national Dataset on the Living Conditions of Women and their Well-being for 52 Countries. *Sci Data* 9, 719 (2022). <https://doi.org/10.1038/s41597-022-01824-2>
- ⁹ Hijmans, R, GADM, the Database of Global Administrative Areas, <https://gadm.org/> (2022)
- ¹⁰ United Nations, Department of Economic and Social Affairs, Population Division (2022). Database on Household Size and Composition (2022).
- ¹¹ ICF. STATcompiler. (2012).
- ¹² Minnesota Population Center. IPUMS International [multiple datasets]. Minneapolis, MN: University of Minnesota. (2023).
- ¹³ ICF. Demographic and health surveys [multiple datasets]. Funded by USAID. (2023).
- ¹⁴ UNICEF. Multiple Indicator Cluster Surveys (MICS) [multiple datasets]. New York: UNICEF (2023).
- ¹⁵ Eurostat. European Union Labour Force Survey [multiple datasets]. European Union. (2023)
- ¹⁶ Eurostat. EU Statistics on Income and Living Conditions (EU-SILC) survey. [multiple datasets]. European Union. (2023)
- ¹⁷ Wilkins, R. HILDA Survey. HILDA Project, Melbourne Institute of Applied Economic and Social Research, University of Melbourne. [multiple datasets]. (2021).
- ¹⁸ Li, J., & Zhu, H. China Family Panel Studies, 2010-2018 [Data set]. Peking University. <https://doi.org/10.18170/CNCFM.2020.001> (2020).

-
- ¹⁹ Kummu, M., Taka, M. & Guillaume, J. H. A. Gridded global datasets for gross domestic product and human development index over 1990–2015. *Sci. Data* 5, 180004 (2018). <https://doi.org/10.1038/sdata.2018.4>
- ²⁰ Belmin, C., Hoffmann, R., Elkasabi, M. et al. LivWell: a sub-national Dataset on the Living Conditions of Women and their Well-being for 52 Countries. *Sci Data* 9, 719 (2022). <https://doi.org/10.1038/s41597-022-01824-2>
- ²¹ ICF. STATcompiler. (2012).
- ²² Smits, J. & Permanyer, I. The subnational human development database. *Sci. Data* 6, 190038 (2019).
- ²³ Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., ... & Hester, J. Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), 1686. <https://doi.org/10.21105/joss.01686> (2019)
- ²⁴ Wickham, H., & Miller, E. haven: Import and Export ‘SPSS’, ‘Stata’ and ‘SAS’ Files. R package version 2.3.1. <https://CRAN.R-project.org/package=haven> (2020).
- ²⁵ Christoph, J., & Lazarevic, M. labelled: Manipulating Labelled Data. R package version 2.8.0. <https://CRAN.R-project.org/package=labelled> (2020).
- ²⁶ Müller, K., & Wickham, H. tibble: Simple Data Frames. R package version 3.1.4. <https://CRAN.R-project.org/package=tibble> (2021).