# Implementing open science

## Jonas Schöley

@jschoeley

0000-0002-3340-8518

schoeley@demogr.mpg.de

MAX PLANCK INSTITUTE
FOR DEMOGRAPHIC RESEARCH

# Benefits of reproducibility and open science
## Some personal experiences

## Together with



Roland **Rau**
@Demographie
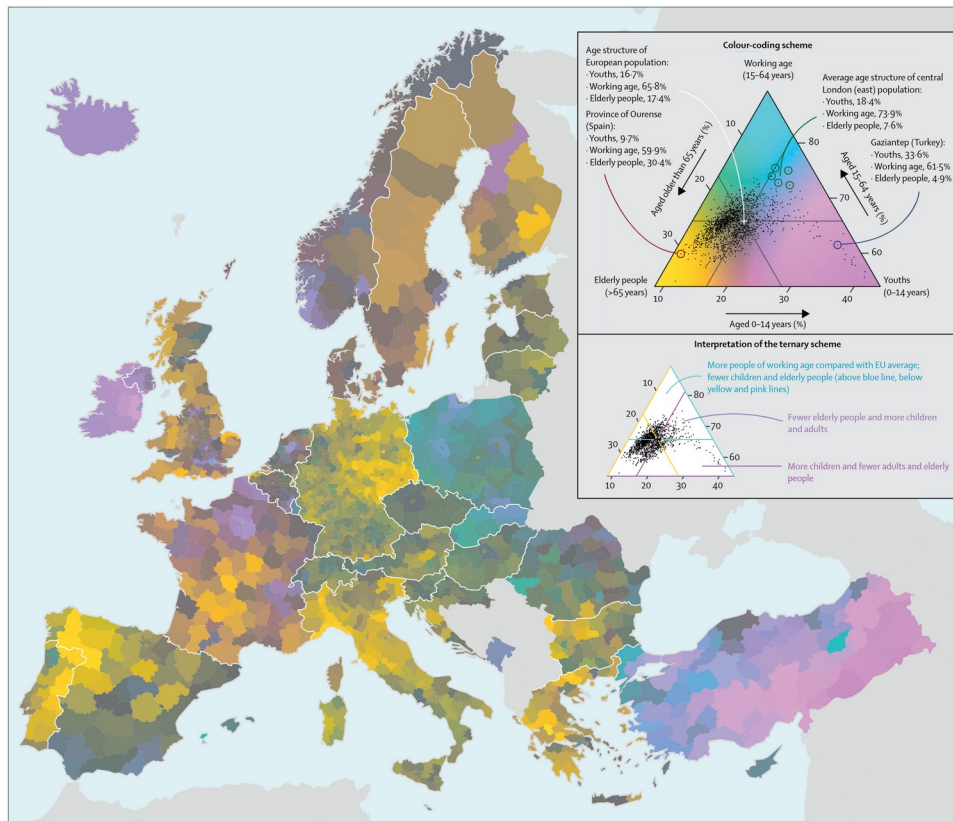
Testing

- The survdiff function in survival compares survival curves using the Fleming-Harrington G-rho family of test. NADA implements this class of tests for left-censored data.
- The maxcombo package compares survival curves using the max-combo test, which is often based on the Fleming-Harrington G-rho family of tests and is designed to have higher power than the logrank test in the scenario of non-proportional hazards such as those resulting from delayed treatment effects.
- clinfun implements a permutation version of the logrank test and a version of the logrank that adjusts for covariates.
- The exactRankTests implements the shift-algorithm by Streitberg and Roehmel for computing exact conditional p-values and quantiles, possibly for censored data.
- SurvTest in the coin package implements the logrank test reformulated as a linear rank test.
- The maxstat package performs tests using maximally selected rank statistics.
- The interval package implements logrank and Wilcoxon type tests for interval-censored data.
- Three generalised logrank tests and a score test for interval-censored data are implemented in the glrt (archived) package.
- survcomp compares 2 hazard ratios.
- The TSHRC implements a two stage procedure for comparing hazard functions.
- The FHtest package offers several tests based on the Fleming-Harrington class for comparing surival curves with right- and interval censored data.
- The LogrankA (archived) package provides a logrank test for which aggregated data can be used as input.
- The short term and long term hazard ratio model for two samples survival data can be found in the YPmodel package.
- The controlTest implements a nonparametric two-sample procedure for comparing the median survival time.
- The survRM2 package performs two-sample comparison of the restricted mean survival time
- The emplik2 package permits to compare two samples with censored data using empirical likelihood ratio tests.
- The KONPsurv package provides powerful nonparametric K-sample tests for right-censored data. The tests are consistent against any differences between the hazard functions of the groups.

# tricolore

## Together with

**Ilya Kashnitsky**
@ikashnitsky



Figure: Colour-coded map of population structures in European Nomenclature of Territorial Units for Statistics 3 regions in 2015
Each population composition in the European Nomenclature of Territorial Units for Statistics 3 regions is uniquely colour coded. Colours show direction and magnitude of deviations from the centrepoint, which represents the average age of the European population, and is dark grey. The hue component of a colour encodes the direction of deviation: yellow indicates an elderly population (>65 years), cyan indicates people of working age (15–64 years), and magenta indicates children (0–14 years). Chroma and lightness components signify the distance from the centre ranging from desaturated and dark colours near the centre to vivid and bright colours at the corners. We provide R code to fully reproduce this map.[10]

Kashnitsky & Schöley (2018). Regional population structures at a glance. 10.1016/S0140-6736(18)31194-2

# tricolore

## Together with
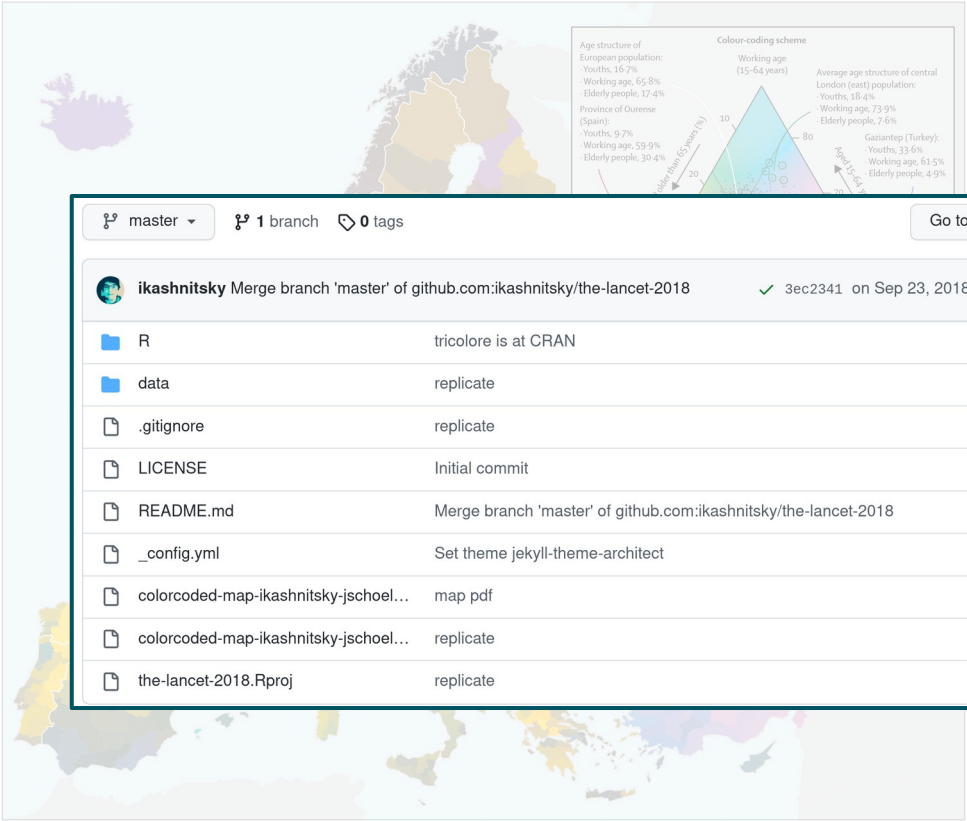
Ilya **Kashnitsky**
@ikashnitsky



Figure: Colour-coded map of population structures in European Nomenclature of Territorial Units for Statistics 3 regions in 2015
Each population composition in the European Nomenclature of Territorial Units for Statistics 3 regions is uniquely colour coded. Colours show direction and magnitude of deviations from the centrepoint, which represents the average age of the European population, and is dark grey. The hue component of a colour encodes the direction of deviation: yellow indicates an elderly population (>65 years), cyan indicates people of working age (15–64 years), and magenta indicates children (0–14 years). Chroma and lightness components signify the distance from the centre ranging from desaturated and dark colours near the centre to vivid and bright colours at the corners. We provide R code to fully reproduce this map.[3]

Kashnitsky & Schöley (2018). Regional population structures at a glance. 10.1016/S0140-6736(18)31194-2

## Together with



Ilya **Kashnitsky**
@ikashnitsky



**Figure 3:** Different representations of the color key for the (centered) ternary balance scheme showing the workforce composition by region in Europe, 2016. Data by Eurostat.

Kashnitsky & Schöley (2018). Regional population structures at a glance. 10.1016/S0140-6736(18)31194-2

# tricolore

Together with

Ilya **Kashnitsky**
@ikashnitsky

Kashnitsky & Schöley (2018). Regional population structures at a glance. 10.1016/S0140-6736(18)31194-2

# tricolore

### Income distribution in Canadian cities



### LMIC education disparity



### Vienna's population by origin



### French election results



### Regional age distribution of COVID deaths In Brazil



Kashnitsky & Schöley (2018). Regional population structures at a glance. 10.1016/S0140-6736(18)31194-2

# tricolore

## Agricultural and Forest Meteorology



Fig. 3. Relative importance (relative $R^2$) of the three hydrometeorological drivers of transpiration regulation. Relative $R^2$ were calculated at each cell dividing the projection of each three drivers partial $R^2$ from the FULL model by the sum of the three partial $R^2$ at the same cell. Partial $R^2$ were projected at the global scale using linear models with climate, soil and vegetation structural variables as explanatory variables (Table S6). Grid colours were calculated using the 'tricolore' package (Schöley and Kashnitsky, 2020) for each cell. Colour gradient indicate the relative importance of the three hydrometeorological constraints. Light grey colour are non-forested areas or with vegetation not taller than 0.5 m (Simard et al., 2011). % VPD: vapour pressure deficit relative importance. % SWC: soil water content relative importance. % PPFD: photosynthetic photon flux density relative importance. Points indicate locations of study sites.

## Soil composition



**Figure 8** **Maps of soil texture composition at two zoom levels, for a region at the southern edge of the Abitibi and James Bay Lowlands soil province.** Only productive forest land characterized by mineral soils was mapped. Agricultural and unproductive forest land, organic soils, anthropogenic infrastructures, and water areas were excluded. Maps were produced with QGIS software, version 3.4 (*QGIS, 2020*). Basemap credit: ©2021 TerraMetrics, ©2021 Google.

Full-size 🖼 DOI: 10.7717/peerj.11685/fig-8

*Figure: Colour-coded map of population structures in European Nomenclature of Territorial Units for Statistics 3 regions in 20... Each population composition in the European Nomenclature of Territorial Units for Statistics 3 regions is uniquely colour coded. Col... centrepoint, which represents the average age of the European population, and is dark grey. The hue component of a colour encode... (>65 years), cyan indicates people of working age (15–64 years), and magenta indicates children (0–14 years). Chroma and lightness... desaturated and dark colours near the centre to vivid and bright colours at the corners. We provide R code to fully reproduce this map.*

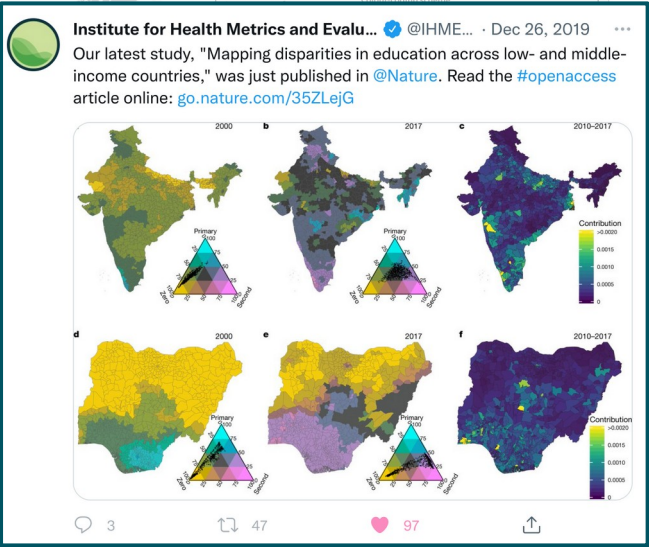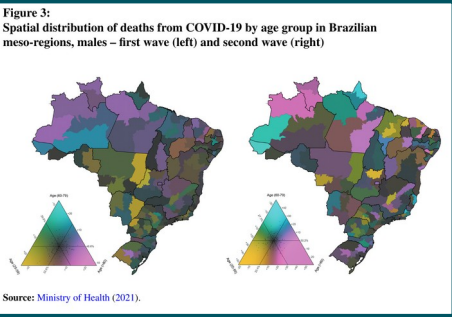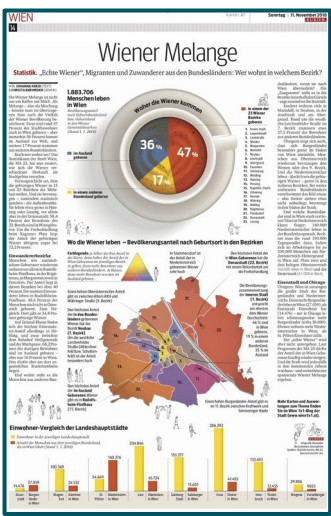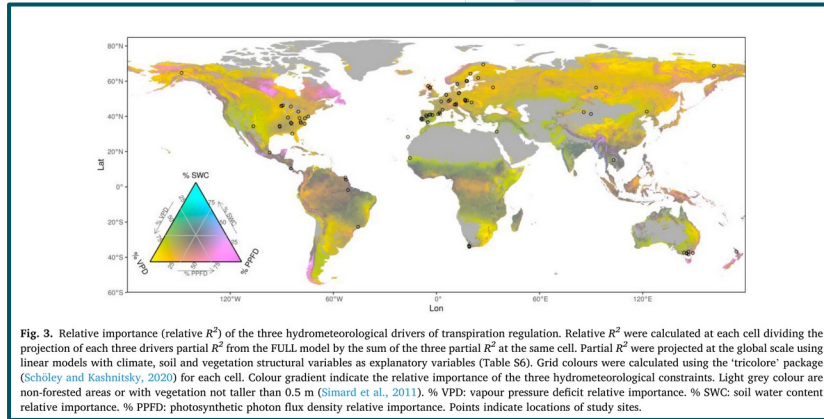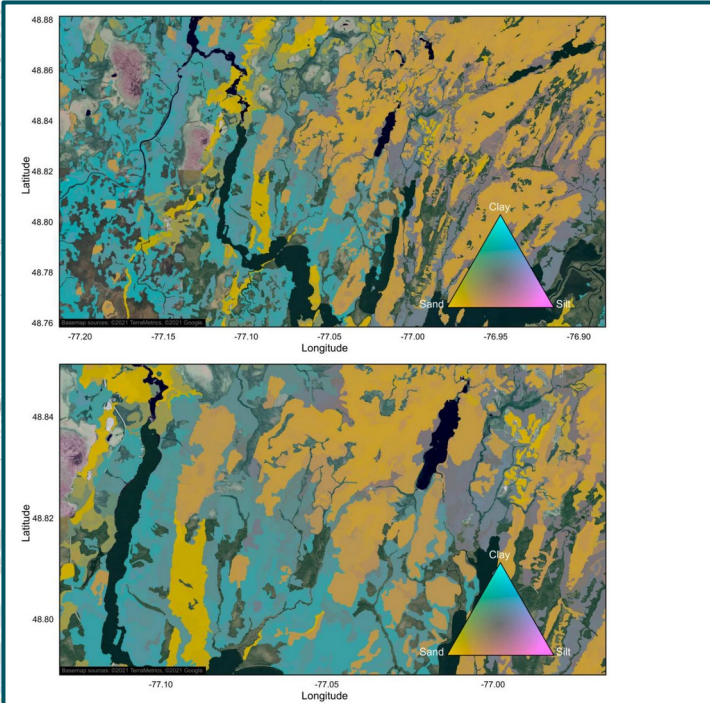Kashnitsky & Schöley (2018). Regional population structures at a glance. <u>10.1016/S0140-6736(18)31194-2</u>

# Reproducible Demographers

Tim **Riffe** @timriffe1 & Enrique **Acosta** @Acosta_Kike_ &
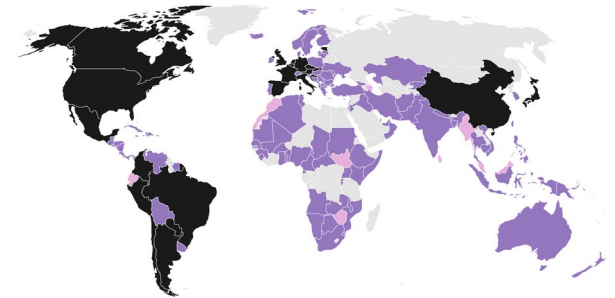Manal **Elzalabany** & Maxi **Kniffka** @MaxiKniffka & Jessica **Donzowa** @jdonzowa

Created a fully reproducible data base of age specific COVID-19 statistics.

**Data availability**

You can get the most up-to-date data at the `OSF` site that we mirror to: https://osf.io/mpwjq/.

Here's an overview of global coverage as of now. A country marked as *forthcoming* means we've identified a source, but that collection is pending for one reason or another. Are you from one of the countries not yet in the collection and want to pitch in? Are you interested in adopting collection for one of our time series that has fallen behind? We're in need of more support. Please reach out, if so by emailing us at `coverage-db < at > demogr.mpg.de`.

■ National and subnational ■ National ■ Forthcoming ■ Not included yet

# Reproducible Demographers

## Christina **Bohk**



| | |
|---|---|
| 📁 | Method00_FreezeRates |
| 📁 | Method01_Hadwiger1940 |
| 📁 | Method02_CoaleMcNeil1972 |
| 📁 | Method03_CoaleTrussell1974 |
| 📁 | Method04_Brass1974 |
| 📁 | Method05_Evans1986 |
| 📁 | Method06_Chandola1999 |
| 📁 | Method07_Schmertmann2003 |
| 📁 | Method08_PeristeraKostaki2007M1 |
| 📁 | Method09_PeristeraKostaki2007M2 |
| 📁 | Method10_MyrskylaGoldstein2013 |
| 📁 | Method11_Saboia1977 |
| 📁 | Method12_WillekensBaydar1984 |
| 📁 | Method13_deBeer1985and1989 |
| 📁 | Method14_Lee1993Log |
| 📁 | Method16_HyndmanUllah2007 |
| 📁 | Method17_ChengLin2010 |
| 📁 | Method18_Myrskyla2013 |
| 📁 | Method22_LiWu2003 |

Bohk et al. (2018). Forecast accuracy hardly improves with methods complexity when completing cohort fertility. 10.1109/5.771073

## Implemented, shared and compared 22 fertility forecasting methods

github.com/fertility-forecasting/validate-forecast-methods/tree/master/basic-scripts-forecast-methods

# Reproducible Demographers

## Rob **Hyndman**

Author of countless R packages widely applied in demography.

demography: Forecasting Mortality, Fertility, Migration and Population Data

Functions for demographic analysis including lifetable calculations; Lee-Carter modelling; functional data analysis of mortality rates, fertility rates, net migration numbers; and stochastic population forecasting.

| | |
|---|---|
| Version: | 1.22 |
| Depends: | R (≥ 3.4), forecast (≥ 8.5) |
| Imports: | ftsa (≥ 4.8), rainbow, cobs, mgcv, strucchange, RCurl |
| Published: | 2019-04-22 |
| Author: | Rob J Hyndman with contributions from Heather Booth, Leonie Tickle and John Maindonald. |
| Maintainer: | Rob J Hyndman <Rob.Hyndman at monash.edu> |
| BugReports: | https://github.com/robjhyndman/demography/issues |
| License: | GPL-2 | GPL-3 [expanded from: GPL (≥ 2)] |
| URL: | https://github.com/robjhyndman/demography |
| NeedsCompilation: | no |
| Materials: | README ChangeLog |
| CRAN checks: | demography results |

github.com/robjhyndman

# Consuming open science



Karlinsky & Kobak (2022). World Mortality Database. github.com/akarlinsky/world_mortality

# Consuming open science



Karlinsky & Kobak (2022). World Mortality Database. github.com/akarlinsky/world_mortality

# Consuming open science



Derived from Karlinsky & Kobak (2022). World Mortality Database.
github.com/akarlinsky/world_mortality

Given **your data** and **your analysis**
**I** arrive at **your results**

– – – – – – – – – – – – – – – – – – – – – –

Given **your research question**,
**my data** and **my analysis**
**I** arrive at **your results**

# Implementing reproducibility & replicability

**Ensure everyone can run your analysis**

**Share and version your analysis**

**Share your data Archive your data Get DOIs**

Schöley – Implementing open science

17

# The computational reproducibility stack

**Ensure everyone can run your analysis**

| |
|---|
| **Scripts & Data** |
| **R Libraries** |
| **System dependencies** |
| **Hardware** |

Schöley – Implementing open science

# The computational reproducibility stack

**Ensure everyone can run your analysis**

| Scripts & Data |
| R Libraries |
| System dependencies |
| Hardware |

Schöley – Implementing open science

# The computational reproducibility stack

**Ensure everyone can run your analysis**

Scripts & Data

```
Warning message:
"`funs()` was deprecated in dplyr 0.8.0.
i Please use a list of either functions or lambdas:

# Simple named list: list(mean = mean, median = median)

# Auto named with `tibble::lst()`: tibble::lst(mean, median)

# Using lambdas list(~ mean(., trim = .2), ~ median(., na.rm = TRUE
i The deprecated feature was likely used in the dataiku package.
  Please report the issue to the authors."
```
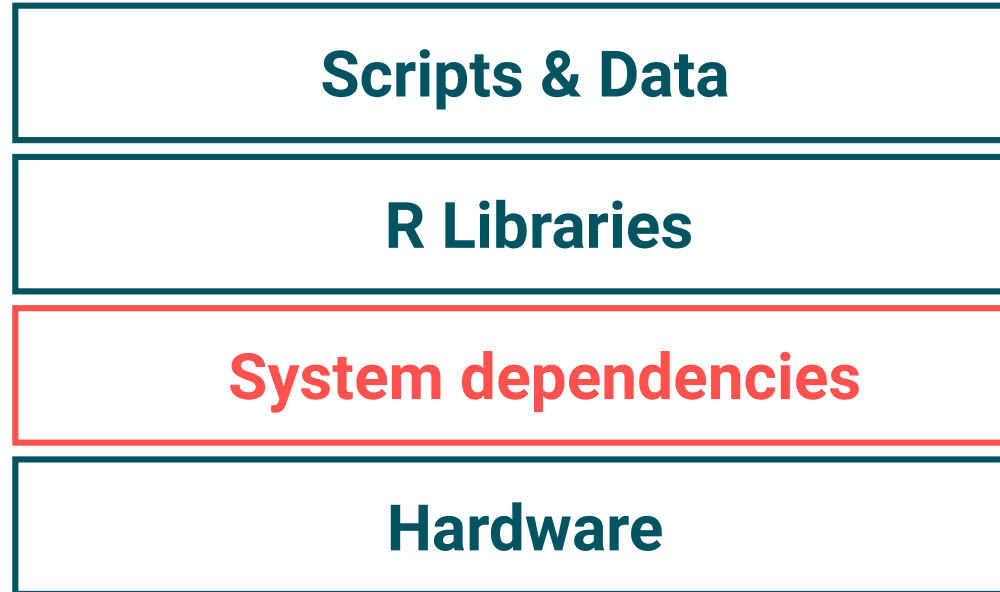
Schöley – Implementing open science

# The computational reproducibility stack

**Ensure everyone can run your analysis**



| Scripts & Data |
|---|
| R Libraries |
| **System dependencies** |
| Hardware |

Schöley – Implementing open science

# The computational reproducibility stack

**Ensure everyone can run your analysis**

Scripts & Data

git zenodo

```
Error in stop_no_virtualenv_starter(version = version, python =
ython) :
   Suitable Python installation for creating a venv not found.
   Requested Python: /usr/bin/python3.10
   Requested version constraint: 3.10
Please install Python with one of following methods:
- https://github.com/rstudio/python-builds/
- reticulate::install_python(version = '<version>')
- Install python3-venv and python3-pip using the system package
anager
```

docker

# The computational reproducibility stack

**Ensure everyone can run your analysis**

| Scripts & Data |
| :---: |
| R Libraries |
| System dependencies |
| Hardware |

**Ensure everyo~~ne~~
can run your
analysis**



R Session Aborted

R encountered a fatal error.

The session was terminated.
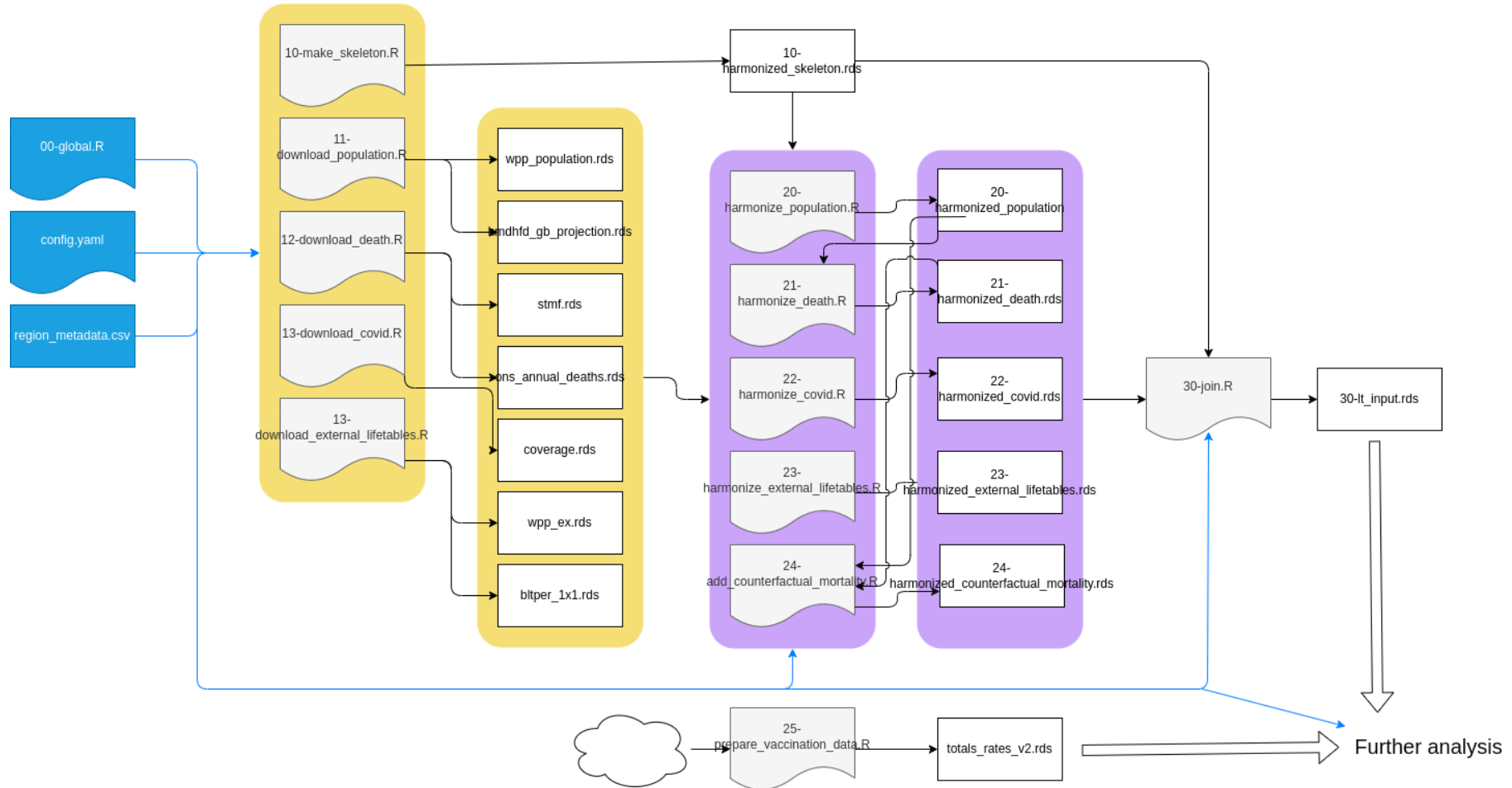
**Start New Session**

Schöley – Implementing open science

24

# Demonstrating the reproducible workflow
## Sharing a statistical model

# Demonstrating the reproducible workflow
## Life expectancy changes in 2021

**Roadblocks to open science**
I can't share my data

# Roadblocks to open science

My code sucks

**Roadblocks to open science**
Others will copy my stuff

# Roadblocks to open science
My co-authors are not on-board

# Reproducible analysis
## github.com/jschoeley

Jonas Schöley

@jschoeley

0000-0002-3340-8518

schoeley@demogr.mpg.de

**MAX PLANCK INSTITUTE**
FOR DEMOGRAPHIC RESEARCH