# Current ecosystem of R package for outbreak analytics

## CURRENT LANDSCAPE



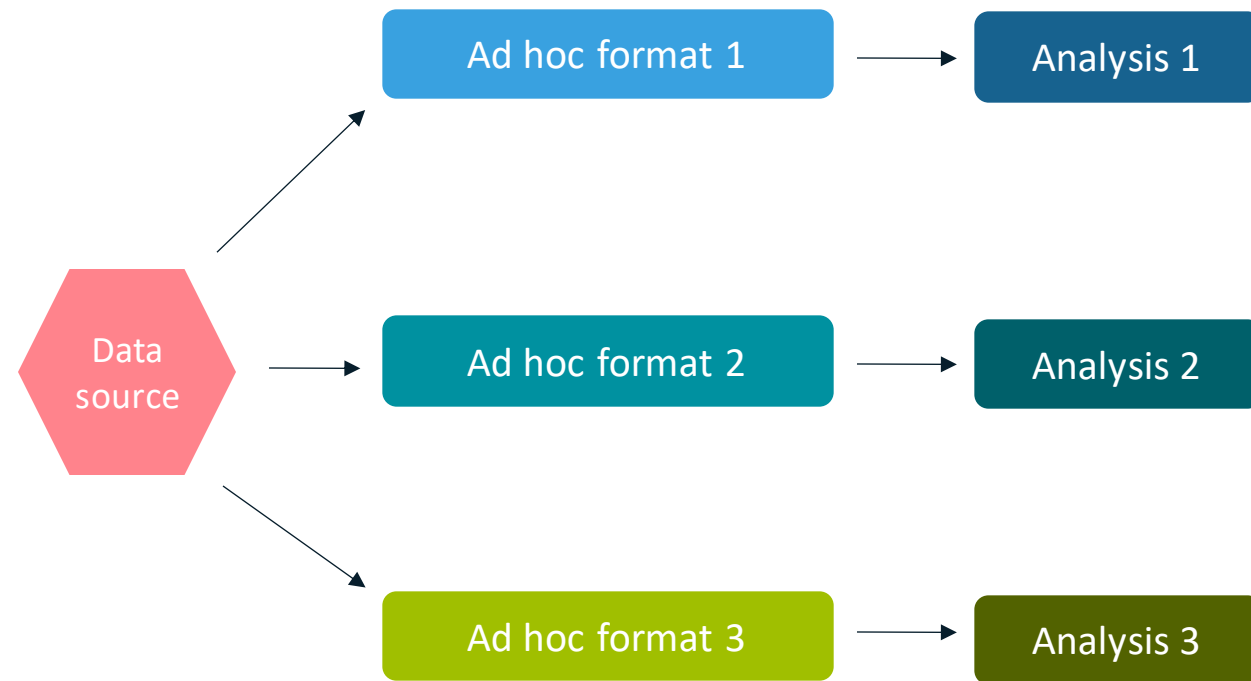## CHALLENGES

- Limited interoperability

- Lack of common data structures

- Data pipelines harder to build

# *Ad hoc* versus canonical data representations

## CURRENTLY COMMON MODEL



## ISSUES

- Pre-processing overhead
  - Requires ad-hoc scripts
  - Sometimes more complex than analysis itself
  - May require dedicated import functions
  - Discourages multiple analyses and methods comparisons

- Prone to data discrepancies
  - Mistakes can be introduced during pre-processing
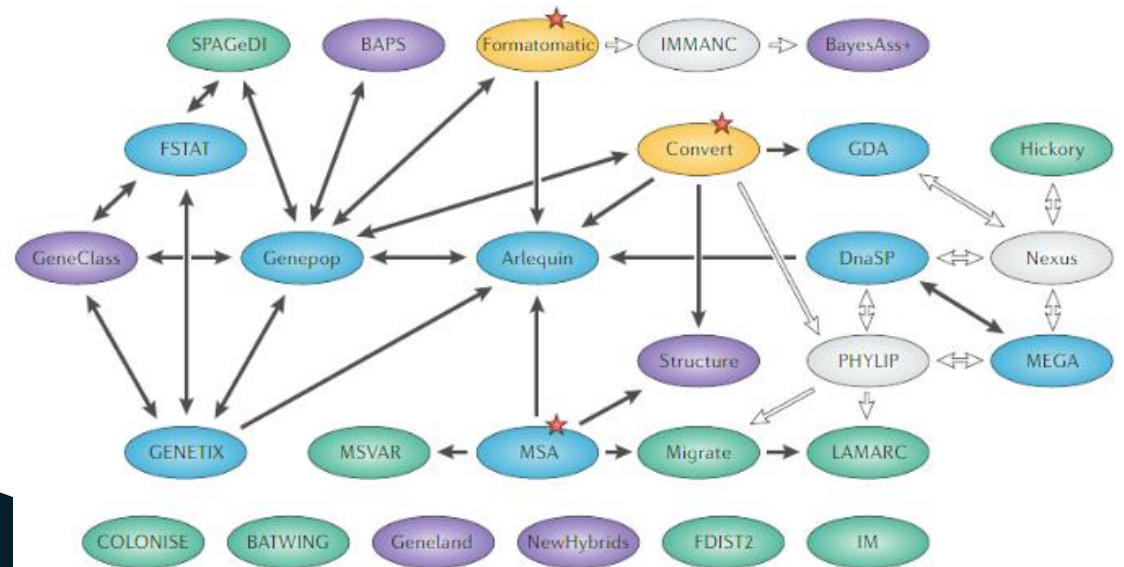  - Hinders comparison of methods and results

data.org

# Why *ad hoc* formats are a problem: lessons from population genetics



Computer programs for population genetics data analysis: a survival guide
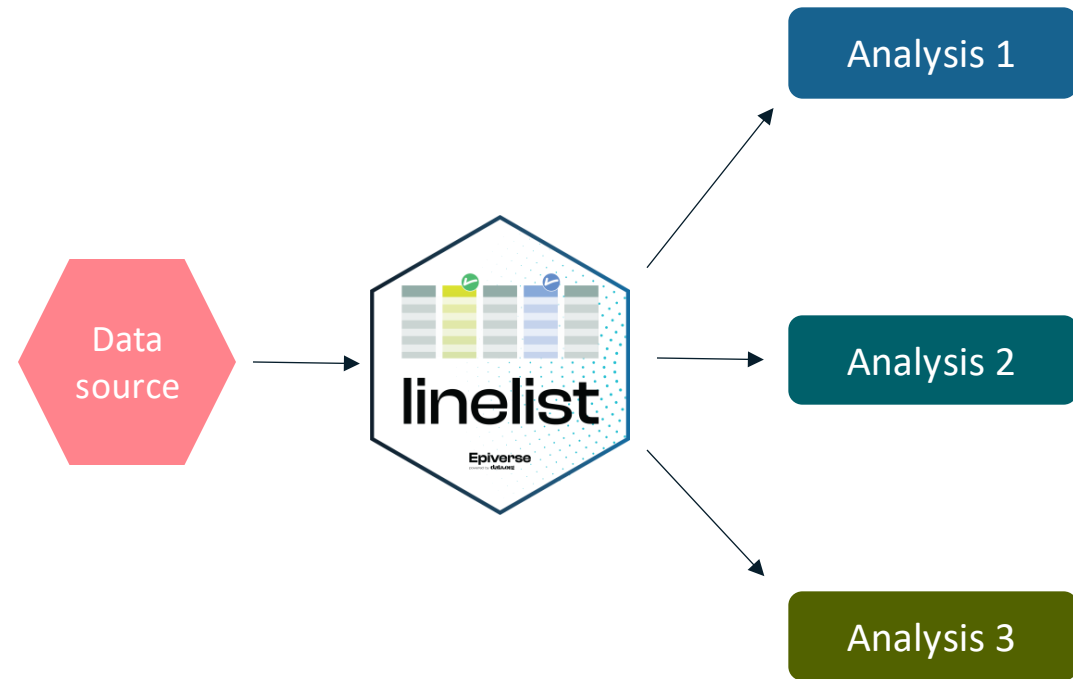
*Laurent Excoffier and Gerald Heckel*

Abstract | The analysis of genetic diversity within species is vital for understanding evolutionary processes at the population level and at the genomic level. A large quantity of data can now be produced at an unprecedented rate, requiring the use of dedicated computer programs to extract all embedded information. Several statistical packages have been recently developed, which offer a panel of standard and more sophisticated analyses. We describe here the functionalities, special features and assumptions of more than 20 such programs, indicate how they can interoperate, and discuss new directions that could lead to improved software and analyses.

- Interoperability but…
  - Requires many conversion steps
  - Complicated data pathways
  - Prone to mistakes/errors

- No coding, testing, or documentation standards

Excoffier & Heckel 2006 *Nature Reviews Genetics*

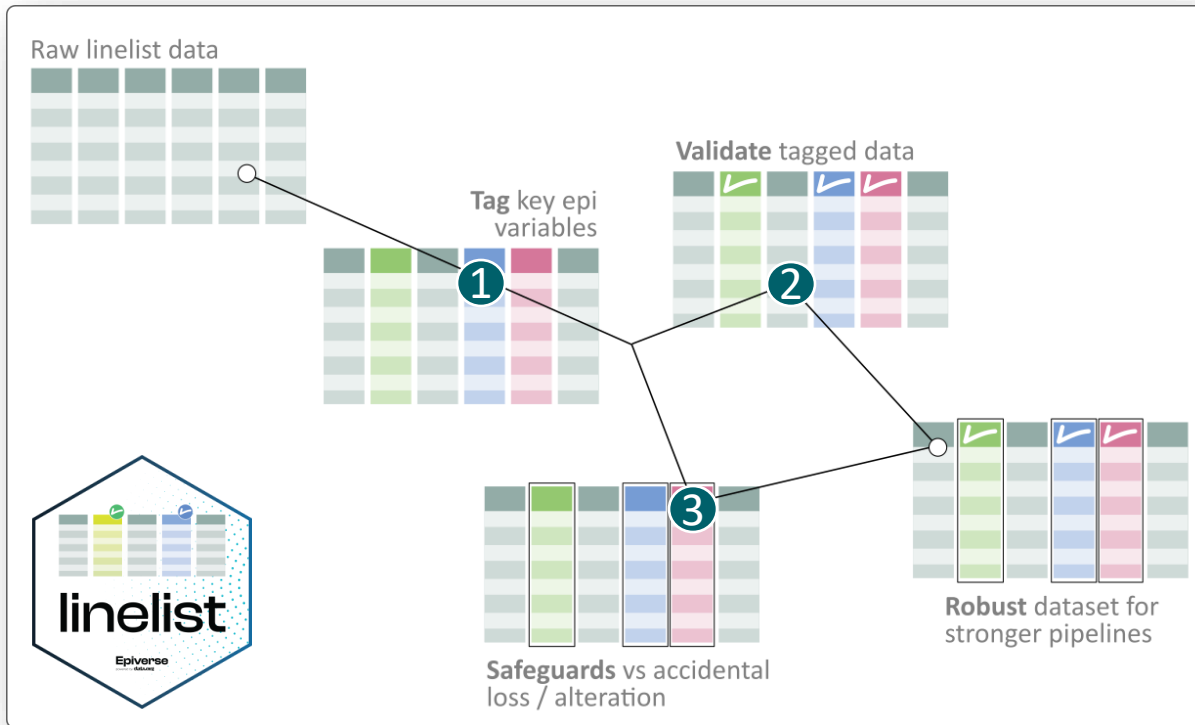# *linelist*: a package to handle case line list data



- Case line list data
  - Spreadsheet-like dataset
  - Rows = different patients
  - Columns = variables, including:
    - Dates of events: symptom onset, case reporting, …
    - Case data: epi case definition, disease outcome, symptoms
    - Patient info: age, gender, occupation, location, …

- Common issues
  - No standardized names for key epi variables
  - Data have the wrong type / class *e.g.:*
    - dates stored as `factors` rather than `Date`,
    - Epi case definition coded as `integer` rather explicit character strings
  - Key data may be altered or lost accidentally

# *linelist* in a nutshell



Raw linelist data

Tag key epi variables

Validate tagged data

Safeguards vs accidental loss / alteration

Robust dataset for stronger pipelines

https://github.com/epiverse-trace/linelist

**❶ Tagging system**

- `linelist`: S3 class extending `data.frame` or `tibble`
- Tags identify columns storing key epi variables in a dataset
- Tags are stored as named `list` in `attr(., "tags")`

**❷ Data validation**

- Validate tagged variables against their expected types
- Easy to extend to more advanced validations (e.g. consistency in timing of events)

**❸ Safeguards**

- Protections against alteration / loss of tagged variables
- Done through S3 dispatching with dedicated methods for 'dangerous' operations
- Implementation for non-generic more cumbersome

data.org

# Coding practices

```r
103  make_linelist <- function(x,
104                            ...,
105                            allow_extra = FALSE) {
106    # assert inputs
107    checkmate::assertDataFrame(x, min.cols = 1)
108    checkmate::assertLogical(allow_extra)
109
110    # The approach is to replace default values with user-provided ones, and then
111    # tag each variable in turn. Validation the tagged variables is done
112    # elsewhere.
113    tags <- tags_defaults()
114
115    args <- list(...)
116    if (length(args) && is.list(args[[1]])) {
117      args <- args[[1]]
118    }
119
120    tags <- modify_defaults(tags, args, strict = !allow_extra)
121
122    out <- x
123    for (i in seq_along(tags)) {
124      out <- tag_variable(out, var_type = names(tags)[i], var_name = tags[[i]])
125    }
126
127    # shape output and return object
128    class(out) <- c("linelist", class(out))
129    out
130
131  }
```

- **Modular** code, short functions

- Short argument list

- **Assertion** of inputs

- Comments
  - mark code structure
  - explain global strategy
  - avoid mere code reformulation

- **Document** as you code

- **Test** as you code (or before)

data.org

# *linelist*

# Documentation



- All functions documented using **roxygen2**

- All exported functions have examples

- README.Rmd

- Vignette

- Website (**pkgdown**)

- Infographics

- Hex sticker!!

- Cheat-sheet underway

## data.org

# Testing and continuous integration

`CRAN OK` `R-CMD-check passing` `codecov 100%`

```
1    test_that("tests for make_linelist", {
2
3      # test errors
4      msg <- "Assertion on 'x' failed: Must be of type 'data.frame', not 'NULL'."
5      expect_error(make_linelist(NULL), msg)
6
7      msg <- "Assertion on 'x' failed: Must have at least 1 cols, but has 0 cols."
8      expect_error(make_linelist(data.frame()), msg)
9
10     msg <- "Assertion on 'var_name' failed: Must be element of set \\{'speed','dist'\\}, but is 'bar'."
11     expect_error(make_linelist(cars, outcome = "bar"), msg)
12
13     msg <- "Unknown variable types: foo\n  Use only tags listed in `tags_names()`, or set `allow_extra = TRUE`"
14     expect_error(make_linelist(cars, foo = "speed", allow_extra = FALSE), msg, fixed = TRUE)
15
16     # test functionalities
17     expect_identical(tags_defaults(), tags(make_linelist(cars), TRUE))
18
19     x <- make_linelist(cars, date_onset = "dist", date_outcome = "speed")
20     expect_identical(tags(x)$date_onset, "dist")
21     expect_identical(tags(x)$date_outcome, "speed")
22     expect_null(tags(x)$outcome)
23     expect_null(tags(x)$date_reporting)
24
25     x <- make_linelist(cars, foo = "speed", bar = "dist", allow_extra = TRUE)
26     expect_identical(tags(x, TRUE), c(tags_defaults(), foo = "speed", bar = "dist"))
27
28     x <- make_linelist(cars, date_onset = "dist", date_outcome = "speed")
29     y <- make_linelist(cars, list(date_onset = "dist", date_outcome = "speed"))
30     expect_identical(x, y)
31
32   })
```

- **Full test coverage** for all functions

- Each function has its own test file

- Testing guides coding

- ~ 50% tests on errors and warnings; ~50% tests on functionalities

- **Github actions** for CI (`usethis`)

data.org

9

# Code reviews, collaboration workflows, communities

```
`[.linelist` <- function(x, i, j, drop = FALSE) {
  # Strategy for subsetting
  #
  # Subsetting is done using the next method in line, for which we drop the
  # linelist class (we cannot use NextMethod because of the extra argument
  # `lost_action`). Then we need to check two things:
  #
  # 1. that the subsetted object is still a `data.frame` or a `tibble`; if not,
  # we automatically drop the `linelist` class and tags
  # 2. if the output is going to be a `linelist` we need to restore previous
  # tags with the appropriate behaviour in case of missing tagged variables

  lost_action <- get_lost_tags_action()

  # Case 1
  out <- drop_linelist(x)[i, j, drop = drop]
  if (is.null(ncol(out))) {
    return(out)
  }
}
```

```
`[<-.linelist` <- function(x, i, j, value) {
  lost_action <- get_lost_tags_action()
  out <- NextMethod()
  old_tags <- tags(x, TRUE)
  out <- restore_tags(out, old_tags, lost_action)
  out
}
```

Solo project is a bad example but...

- **Code reviews**
  - Useful, but non-trivial
  - Pair-programming better?
  - Part of Pull Requests, or separate exercise?

- **Lean / Agile development**
  - Important to incorporate user feedback early and often
  - Short cycles with regular production of MVPs better for user engagement?

- **Dev and user community**
  - How to engage end-users?
  - Metrics of community health or performance?

data.org

# Blueprints: how many boxes to ticks?

goodpractice::gp() ?

- ❑ Coding standards

- ❑ Documentation standards

- ❑ Testing and continuous integration

- ❑ Collaboration framework
  - ❑ Code reviews
  - ❑ Pair programming
  - ❑ Agile/Lean methodologies

- ❑ Community engagement and monitoring

- ❑ …

data.org

# THANK
# YOU

Senior RSE position for *Epiverse* open until Friday at data.org!
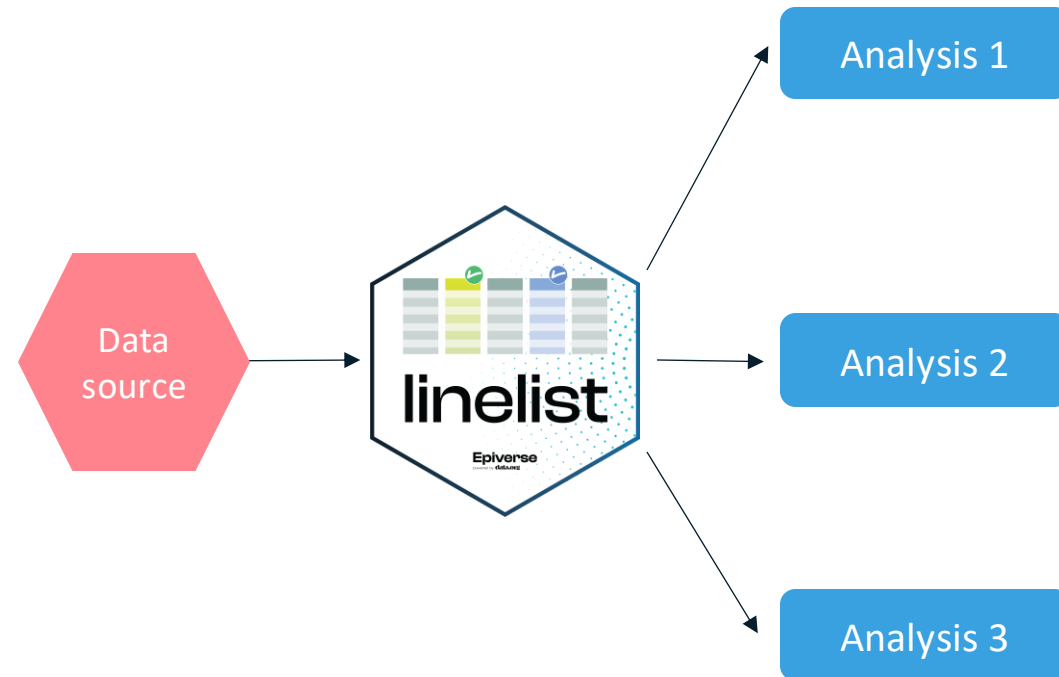
Contact:

👤 Thibaut Jombart

✉ thibaut@data.org

# linelist: tagging, validating, and safeguarding key epi data



**Linelist = S3 class extending `data.frame`**

- Tagging system: identify key epi data
  - `set_tags()`: identify columns storing key epi variables (e.g. date of onset, disease outcome, age)
  - `select()`, `select_tags()`, `tags_df()`: access tagged variables
  - Tags are stored in the attributes of the object

- Data validation
  - Key epi variables have pre-defined acceptable types
  - `validate_linelist()`: check tags consistency and types

- Safeguards
  - Protections against alteration / loss of tagged variables
  - `lost_tags_action()`: defines behaviour to adopt
  - `rename()`, `names() <-`: update tags as needed
  - `select(), [, [[, …`: issues warnings/errors if variable lost

data.org

# Epiverse-TRACE: building a coherent ecosystem for outbreak analytics in R

- Simplest workflows rely on:

```
import_data() %>%
    preprocess_data() %>%
    use_analysis() %>%
    summarise_or_plot()
```

- Complications:
  - Lack of canonical representations
  - Multiple inputs (*e.g.* case counts + serial interval)
  - Multiple layers of analysis
    (*e.g.* estimate R -> make projections)



data.org