

Vigicaen: A 'vigibase®' Pharmacovigilance Database Toolbox.

Charles Dolladille¹ and Basile Chrétien²

¹ University of Caen Normandy, Pharmacology Department, Centre Hospitalier Universitaire de Caen, Caen, France ² University of Nagoya, Department of biostatistics, Nagoya University Hospital, Nagoya, Japan

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [↗](#)

Submitted: 25 August 2025

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

Advanced methodologies are essential when conducting disproportionality analyses using pharmacovigilance data, as traditional approaches are susceptible to various biases such as reporting bias and confounding. The aim of vigicaen is to provide a toolbox for the VigiBase® Extract Case Level database, resolving technical challenges related to the database large size, and providing easier and reproducible access to advanced features. The package is built on top of the parquet file format. Functions related to drug and adverse event identification, descriptive features such as time to onset, dechallenge and rechallenge outcomes are provided. Command line side-effect outputs aim at fast resolving of common issues related to drug and adverse event identification. The package is intended for pharmacovigilance practitioners, clinicians and researchers with or without advanced biostatistical skills. A graphical output can be produced for routine use, to support daily assessment of causality.

Statement of need

Disproportionality analysis represents an essential component in the domain of drug safety signal detection. Advanced methodologies are required to address common biases within pharmacovigilance databases. These analyses necessitate expertise in biostatistical software, such as R, which may present substantial challenges in terms of acquiring and maintaining the requisite skills—in addition to a solid understanding of pharmacovigilance principles and reporting systems.

For decades, the World Health Organization (WHO) has been collecting adverse drug reaction reports, called Individual Case Safety Reports (ICSRs), from its member countries, populating more than 40 millions reports to date. This pharmacovigilance database is called VigiBase® and is managed by the Uppsala Monitoring Centre in Sweden. (Centre, n.d.) These ICSRs describe the course of patients who experienced an adverse event (a medical condition) after taking a drug. The burning question is whether this adverse event was actually related to the drug intake, e.g. if it is an adverse drug *reaction* (ADR). The pharmacovigilance database aims at uncovering the very first potential signals of association between drugs and ADRs. (Montastruc et al., 2011)

It relies on disproportionality analysis, a statistical method that produces estimators of how unlikely the number of observed ICSRs reporting on a specific drug and adverse event is to be attributable to chance alone. Together with an incertitude margin, these estimators are used to raise safety signals on drugs. (Montastruc et al., 2011)

The Uppsala Monitoring Centre grants access to VigiBase® to researchers, either academic or industrial, under a licence contract. The most extensive available version is called Extract Case Level: It contains all the ICSRs, with information such as the patient demographics, the

42 drug intake, the adverse events, the outcome, the dechallenge and rechallenge outcome, and
43 the time to onset. However, this version is provided as large text files, and requires a lot of
44 processing before being usable for analysis. Those text files might be particularly challenging
45 to use in R, as they would often exceed the size of the available Random Access Memory, thus
46 requiring advanced knowledge of R computing techniques. Clinicians and pharmacovigilance
47 practitioners typically lack these skills, and therefore struggle to use the VigiBase® data for
48 their research. As a result, they would often rely in partial data, with limited statistical
49 modelling options.

50 The vigicaen package aims at providing a toolbox for the VigiBase® Extract Case Level
51 database, tackling a few technical challenges to run on low-specification computers, and
52 provide easy and reproducible access to advanced features.(Dolladille & Chrétien, 2025)
53 This article will explain the technical choices and data management logic underlying the
54 package, and provide some examples of its main features. Additional examples and use
55 cases are treated in the package vignettes, which can be found on the package website
56 at <https://pharmacologie-caen.github.io/vigicaen/>. Of important note, the package is not
57 supported nor reflects the opinion of the WHO. The Uppsala Monitoring Centre, in charge
58 of maintaining VigiBase®, was informed of the package development and kindly allowed its
59 publication, acknowledging the potential benefit to promote the use of VigiBase®.

60 Processing vigibase® source files.

61 Clinicians and pharmacovigilance researchers are used to work with low-specification computers.
62 The typical available Random Access Memory rarely exceeds 16GB, which is one of the key
63 resources to deal with large data files in R.(22 Arrow – r for Data Science (2e), n.d.) VigiBase®
64 Extract Case Level files currently exceed 30GB once unpacked, which is way too large to be
65 loaded in-memory for mainstream readers like `read.table()`.

66 Vigicaen relies on parquet files a recent format based on open standards.(Parquet, n.d.) Arrow
67 is a cross-language development platform that allows for manipulation of large datasets.(Apache
68 Arrow, n.d.) It is implemented in R via the arrow package.(Richardson et al., 2025) Datasets
69 remain out of memory, allowing for processing of large files on low-specification computers.
70 Various tests of vigicaen on 16GB RAM computers succeeded in processing the source files.
71 This, in combination with an as close as possible alignment with the tidyverse style guide, is also
72 aimed at providing a modern and more rigorous approach as compared to base R.(Wickham &
73 RStudio, 2023)

74 Sourcing VigiBase® Extract Case Level files is done with the `tb_*` family functions.

75 First, we define paths to the source folders.

```
library(vigicaen)

path_base <- paste0(tempdir(), "/main/")
path_sub  <- paste0(tempdir(), "/sub/")

dir.create(path_base)
dir.create(path_sub)
```

76 Example files can be put in these folders.

```
create_ex_main_txt(path_base)
create_ex_sub_txt(path_sub)
```

77 Then, we run the related `tb_*` function, `tb_vigibase()`.

```
tb_vigibase(path_base, path_sub)
```

78 `##`

```

79  ## -- tb_vigibase() -----
80  ## i Checking for existing tables.
81  ## i Creating vigibase tables.
82  ## This process must only be done once per database version.
83  ## It can take up to 30minutes.
84  ## =====>----- 33% | 1s | Remove duplicates
85  ##
86  With an average computer, the real running time is around 20-30minutes on current database
87  version.
88  If the dictionaries for drugs and adverse events are also required, tb_who() and tb_meddra()
89  can be used.

```

Identifying drugs and adverse events

Exposure to drugs and occurrence of adverse events are located in the drug and adr tables, respectively. They connect together through the demo table, in a many-to-one relationship, via the UMCReportId key variable. Drugs and adverse events themselves are identified by codes (or IDs) from the WHO Drug Dictionary and the Medical Dictionary for Regulatory Activities (MedDRA), respectively. Disproportionality analysis requires a dataset with one row per ICSR, with the corresponding drugs and adverse events.

The following logic is implemented in vigicaen:

- 1 Use drug and adverse event names to collect their IDs.
- 2 Match the IDs in drug and adr tables to identify the cases.
- 3 Report this information in demo (or any other Vigibase® table).

This is done with the get_* functions (step 1), and the add_* functions (steps 2 and 3). The overall process requires the sequential use of both. Below is an example to identify the drugs. The same principle is applied to adverse events.

```

# load vigibase tables and drug dictionary
demo <- dt_parquet(path_base, "demo")
drug <- dt_parquet(path_base, "drug")

# for the demonstration, we will use built-in example files
demo <- demo_
drug <- drug_
mp <- mp_

# Select drug names
d_sel <-
  list(ipilimumab = "ipilimumab")

# Get the drug IDs
d_drecno <-
  get_drecno(
    d_sel,
    mp = mp
  )

```

```

104 ##
105 ## -- get_drecno() -----

```

```

106 ##
107 ## -- `d_sel`: Matching drugs --
108 ##
109 ## -- v Matched drugs
110 ##
111 ## > `ipilimumab`: "ipilimumab" and "ipilimumab;nivolumab"
112 ##
113 ##
114 ## i Set `verbose` to FALSE to suppress this section.
115 ##
116 ##
117 ##
118 ## -----
# report into demo
demo <-
  demo |>
  add_drug(
    d_drecno,
    drug_data = drug
  )
119 ## i `.data` detected as `demo` table.

120 Displaying information at the command line

121 As seen in the output above, the get_* functions do 2 things: They return drug or adverse
122 event IDs (stored in d_drecno in the example), and they display command line information
123 about the matching process. This is especially useful since drugs and adverse events name
124 may vary in their spelling and case, while the underlying dictionary only accepts exact matches.
125 Matched and un-matched names are displayed, along with some hints for the unmatching
126 reasons.

meddra <- meddra_

a_sel <-
  list(colitis_term = c("Colitis", "Autoimmune colitis"),
       pneumonitis_term = "pneumonitis")

a_llt <- get_llt_soc(a_sel, term_level = "pt", meddra = meddra)

127 ##
128 ## -- get_llt_soc() -----
129 ##
130 ## -- v Matched reactions at `pt` level (number of codes) --
131 ##
132 ## > `colitis_term`: "Autoimmune colitis (1)" and "Colitis (25)"
133 ## > `pneumonitis_term`: x No match
134 ##
135 ##
136 ## i Set `verbose` to FALSE to suppress this section.
137 ##

```

```

138 ##
139 ##
140 ## -- x Unmatched reactions --
141 ##
142 ##
143 ##
144 ## -- ! Some reactions did not start with a Capital letter
145 ##
146 ##
147 ##
148 ## * In `pneumonitis_term`: x "pneumonitis"

```

149 The named list for inputting drug and adverse event names

150 The `get_*` and `add_*` functions are built on top of named list as first argument. This structure
 151 may seem a bit busy, especially for new comers, but it allows for genuine flexibility when analyses
 152 plan increment. As an example, one may create `list(drug_group_1 = c("ipilimumab",`
 153 `"nivolumab"))` to automatically gather all ICSRs reporting one of these two drugs, through
 154 `get_drecno()` and `add_drug()`.

155 Descriptive features

156 Descriptive features often take an important place in pharmacovigilance studies. They may
 157 be as important as producing statistical estimands, to assess the liability of a given drug.
 158 Among them, the time to onset is rather challenging to compute. The main reasons are the
 159 incertitude around the exact reported time to onset, and the potential multiple reports for a
 160 given drug-adverse event pair in a single ICSR. The first is tackled by the Uppsala Monitoring
 161 Centre, which recommends in internal documentation to analyze ICSR where the incertitude
 162 interval is no more than a day. The second is addressed in `extract_tto()` or `desc_tto()`,
 163 which only extracts the longest time to onset reported for a given drug-adverse event pair in a
 164 given ICSR. This variable is called `tto_max`. Admittedly, this is a simplification that might not
 165 cover all potential use cases, for example if the question is the time since last infusion of a
 166 drug.

167 A similar simplifying approach is applied to drug dechallenge (`desc_dch()`) and rechallenge
 168 (`desc_rch()`) outcomes, as well as adverse event outcome (`desc_outcome()`).

169 Disproportionality estimates

170 Although the aim of the package is to prepare readily available datasets for users to compute
 171 disproportionality on their own via advanced modelling techniques, it also provides basic esti-
 172 mates through the `compute_dispro()` and `compute_interaction()` functions. The underlying
 173 computations rely on the Norén et al methodology, for both point estimates, confidence and
 174 credibility intervals. (Norén et al., 2013)

175 Routine use

176 As a routine pharmacovigilance practitioner, key information on a drug - adverse event pair
 177 may be needed out-of-the-box, without further need for manipulating the underlying tables. To
 178 adress the typical needs (disproportionality estimand, time to onset, dechallenge and rechallenge
 179 outcomes), `vigi_routine()` creates a graphical output for a given pair. It is intended as
 180 a daily practice tool, to support routine assessment of causality. The graph can easily be
 181 exported to an external file with the `export_to` argument.

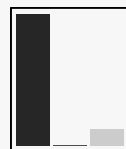
```
vigi_routine(
  demo,
  drug,
  adr_,
  link_,
  d_code = d_dreco,
  a_code = a_llt[1],
  vigibase_version = "Current"
)
```

VIGIBASE ANALYSIS

Drug: ipilimumab
Adverse event: colitis_term
Setting: All reports
VigiBase version: Current

N° of cases: 9

Rechallenge



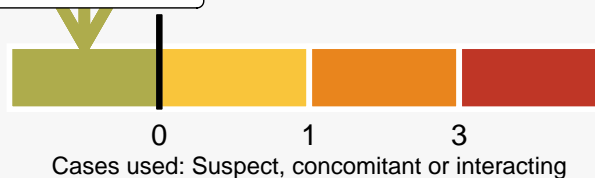
Suspected: 8
Concomitant: 0
Interacting: 1

Total	3
Positive	0
Rate	0%

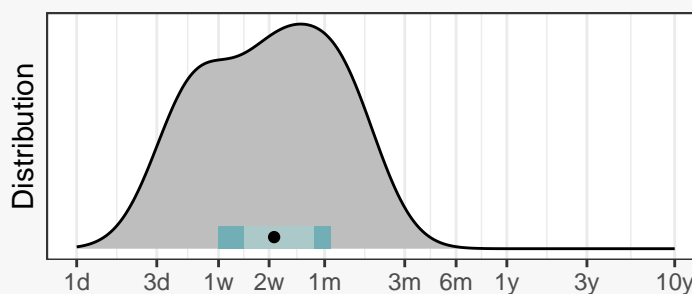
Informative
rechallenges only

Disproportionality Analysis

IC025 = -0.1



Time to onset



d: day, w: week, m: month, y: year
x axis capped at 1 day (min) and 10 years (max)

Created with vigicaen, the R package for VigiBase®

Conclusion

Easier, reproducible research in pharmacovigilance databases is key to appropriate safety signal detection. Vigicaen proposes a set of tools based on popular open standards to facilitate pharmacovigilance analysis in R.

Acknowledgements

The information presented in this study does not represent the opinion of the Uppsala Monitoring Centre or the World Health Organization. We thank the research team at the Uppsala Monitoring Centre (Uppsala, Sweden) who provided case-level data from VigiBase®.

References

- 22 *arrow – r for data science (2e)*. (n.d.). <https://r4ds.hadley.nz/arrow.html>
- Apache arrow*. (n.d.). <https://arrow.apache.org/>
- Centre, U. M. (n.d.). *About VigiBase*. <https://who-umc.org/vigibase/>
- Dolladille, C., & Chrétien, B. (2025). *vigicaen: 'VigiBase' Pharmacovigilance Database Toolbox*. <https://github.com/pharmacologie-caen/vigicaen>
- Montastruc, J.-L., Sommet, A., Bagheri, H., & Lapeyre-Mestre, M. (2011). Benefits and strengths of the disproportionality analysis for identification of adverse drug reactions in a pharmacovigilance database. *British Journal of Clinical Pharmacology*, 72(6), 905–908. <https://doi.org/10.1111/j.1365-2125.2011.04037.x>
- Norén, G. N., Hopstadius, J., & Bate, A. (2013). Shrinkage observed-to-expected ratios for robust and transparent large-scale pattern discovery. *Statistical Methods in Medical Research*, 22(1), 57–69. <https://doi.org/10.1177/0962280211403604>
- Parquet*. (n.d.). <https://parquet.apache.org/>
- Richardson, N., Cook, I., Crane, N., Dunnington, D., François, R., Keane, J., Moldovan-Grünfeld, D., Ooms, J., Wujciak-Jens, J., Luraschi, J., Werner, K. D., Wong, J., & Arrow, A. (2025). *Arrow: Integration to 'apache' 'arrow'*. <https://cran.r-project.org/web/packages/arrow/index.html>
- Wickham, H., & RStudio. (2023). *Tidyverse: Easily install and load the 'tidyverse'*. <https://cran.r-project.org/web/packages/tidyverse/index.html>