

## **Wikiprint Book**

**Title:** udg/ecoms/RPackage/homogeneization

**Subject:** TracMeteo - udg/ecoms/RPackage/homogeneization

**Version:** 24

**Date:** 05/19/2022 01:52:01 PM

## Table of Contents

Vocabulary definition	3
Dictionary	3

The different nature of the different climate products, models and variables, and the idiosyncratic naming and storage conventions often applied by the various modelling centres, makes necessary a previous homogenization across datasets in order to implement a truly user-friendly toolbox for data access. The `ecomSUDG.Raccess` package achieves this aim by defining a common *vocabulary* to all climate datasets. The particular variables of each dataset are translated -and transformed if necessary- into the common vocabulary by means of a *dictionary*. Both features are described below:

## Vocabulary definition

In order to set a common framework with a precise definition of the variables, the `ecomSUDG.Raccess` package is based on the use of a vocabulary. In essence, the vocabulary is a table containing the standard names of a number of variables commonly used in impact studies and downscaling applications, subject to permanent revision or addition of new standard variables. The naming conventions and the units are based on the [?NetCDF Climate and Forecast Metadata Convention](#) and to the [?SPECS naming convention](#).

Note that not all the variables listed in the vocabulary are necessarily available at the ECOMS-UDG, although all variables available must be included in the dictionary, ensuring in all cases the retrieval of homogenized data by default. The vocabulary contains the following variable information:

- `identifier`: this is the standard name that the `loadECOMS` function will recognize automatically (see [this link](#) to see all available variables and their standard naming in R).
- `standard_name`: standard long name (description) of the variable as defined by the CF convention.
- `units`: units in which the standard variable is returned.

The vocabulary has been included as a built-in dataset of the `ecomSUDG.Raccess` package in order to provide the user with a reference of standard variables. This vocabulary is subject to permanent revision as new variables are required.

```
library(ecomSUDG.Raccess)
data(vocabulary)
print(vocabulary)
```

identifier	standard_name	units
1 hurs	2-meter relative humidity	%
2 hursmax	maximum 2-meter relative humidity	%
3 hursmin	minimum 2-meter relative humidity	%
4 hus	specific humidity	kg.kg-1
5 huss	2-meter specific humidity	kg.kg-1
6 prsn	total snowfall amount	mm
7 ps	air pressure at surface level	Pa
8 psl	air pressure at sea level	Pa
9 rlds	surface downwelling longwave radiation	W.m-2
10 rsds	surface downwelling shortwave radiation	W.m-2
11 snld	snow depth mm (water eq.)	
12 ta	air temperature	degrees Celsius
13 tas	2-meter air temperature	degrees Celsius
14 tasmax	maximum 2-m air temperature	degrees Celsius
15 tasmin	minimum 2-m air temperature	degrees Celsius
16 tdps	2-meter dewpoint temperature	degrees Celsius
17 tp	total precipitation amount	mm
18 ua	eastward wind	m.s-1
19 uas	eastward near-surface wind	m.s-1
20 va	northward wind	m.s-1
21 vas	northward near-surface wind	m.s-1
22 wss	near-surface wind speed	m.s-1
23 wssmax	maximum near-surface wind speed	m.s-1
24 z	geopotential height	m
25 zs	surface geopotential height	m

## Dictionary

The dictionary is a table whose aim is twofold:

1. On the one hand, the dictionary is intended for the translation of generic variables, as idiosyncratically defined in each particular dataset, to the standard variables defined in the vocabulary with their corresponding nomenclature and units. This is achieved by providing a correspondence

between the name of the variable as encoded in the dataset (this is the variable name returned by the `datasetInventory` function) and the corresponding name of the standard variable as defined in the vocabulary (i.e., the `identifier`), and by applying the corresponding transformation to the native variable in order to match the standard units by means of an `offset` and a `scale` factor. In some particular cases (e.g. the precipitation provided by the System4 model outputs), the variables are also deaccumulated.

- The dictionary also provides additional metadata often not explicitly declared in the datasets, regarding the *time* aggregation of the dataset (often referred to as the *cell method*). This includes the fields `time_step`, which is merely informative, and describes the time interval between two consecutive values, and the `lower_time_bound` and `upper_time_bound`, which are the values that should be summed to each verification time to unequivocally delimit the time span encompassed by each value.

In essence, the dictionary is a comma-separated text file (csv) that, by default, is identified with the same name than the argument `dataset` of the `loadSeasonalForecast` function. The dictionaries for the currently available datasets are included in the `ecomSUDG.Raccess` package, within the `dictionaries` folder. The dictionaries are read internally by `loadSeasonalForecast` to undertake the conversions needed for returning the standard variables, so by default, the user does not need to worry about them. In case an interested user wants to inspect a particular dictionary, he/she can proceed as follows:

```
# Path to the installed library
libPath <- find.package("ecomSUDG.Raccess")
# Path to the dictionaries folder
dicPath <- file.path(libPath, "dictionaries")
list.files(dicPath)
```

```
[1] "CFSv2_seasonal_16.dic"  "NCEP.dic"
[3] "readme.txt"            "System4_annual_15.dic"
[5] "System4_seasonal_15.dic" "System4_seasonal_51.dic"
[7] "WFDEI.dic"
```

In this case, we inspect the dictionary of the CFSv2 dataset:

```
dic.cfs <- read.csv(list.files(dicPath, full = TRUE)[grep("CFSv2", list.files(dicPath))])
str(dic.cfs)
```

```
'data.frame':      4 obs. of  11 variables:
 $ identifier      : Factor w/ 4 levels "tas","tasmax",...: 2 3 1 4
 $ short_name      : Factor w/ 4 levels "Maximum_temperature_height_above_ground",...: 1 2 4 3
 $ time_step       : Factor w/ 1 level "6h": 1 1 1 1
 $ lower_time_bound: int  0 0 0 0
 $ upper_time_bound: int  6 6 6 6
 $ cell_method     : Factor w/ 4 levels "max","mean","min",...: 1 3 2 4
 $ offset          : num  -273 -273 -273 0
 $ scale           : int  1 1 1 21600
 $ deaccum         : int  0 0 0 0
 $ derived         : int  0 0 0 0
 $ interface       : logi  NA NA NA NA
```

The latest version of the dictionaries can be checked-out in the development version of the package at the [?GitHub repository](#).

The columns of the dictionary are next described:

- `identifier`: this is the name of the standard variable, as defined in the vocabulary
- `short_name`: this is the name with which the original variable has been coded in the dataset
- `time_step`: the time interval between consecutive times in the time dimension axis (in hours). The standard nomenclature number of hours followed by the letter *h* (i.e.: 12h, 6h, 24h ...) is used in the default dictionaries included in the R package.
- `lower_time_bound`: lower time bound of the variable. This should be in hours. This is the number of hours to subtract from the verification times to define the lower time boundary of the verification period. See the example below.
- `upper_time_bound`: upper time bound of the variable. This should be in hours. This is the number of hours to add to the verification times to define the upper time boundary of the verification period. See the example below.
- `cell_method`: function of time aggregation between the lower and upper time bound. For instance, its value is "none" for instantaneous variables, "mean" for mean daily temperatures or "sum" for daily precipitation values. See the example below.
- `offset`: constant summed to the original variable for units conversion (e.g.: `offset = -273.15` for conversion from Kelvin to Celsius)

- `scale`: scale factor applied to the original variable for units conversion (e.g.: `scale = 0.001` for conversion from m to mm)
- `deaccum`. This is a logical flag (0 = FALSE, 1= TRUE), which indicates if the variable should be de-accumulated at each time step. Typically applied to precipitation in some forecast datasets.
- `derived`: this value is internally used by the loading functions to know if the variable is derived from any other variable(s) or can be directly read from the dataset.
- `interface`: this is a internal value used by the loading functions.

**Example:** Suppose we have daily temperature data, each record associated to a date (e.g. 19780420 for 20th April 1978). The lower time bound in this case is 0, and the upper time bound 24, indicating that the verification period of the value starts 20 April 1978 at 00:00 and ends 20 April 1978 23:59:59 (i.e., the time interval [19780420, 19780421]). Because it is mean temperature, the `cell_method` for this variable is "mean".

If the same record has a time associated at 12:00 rather than 00:00 (i.e. 20 April 1978 at 12:00 UTC), then the lower and the upper time bounds would be defined at -12 and +12 respectively. Similarly, if the variable is instantaneous rather than a daily mean (e.g. instantaneous temperature at 00:00 UTC) the `cell_method` would correspond to "none", and both the lower and upper time bounds would be equal to 0).