



# WK 37 - Advanced Tree Physiological Data Processing in R



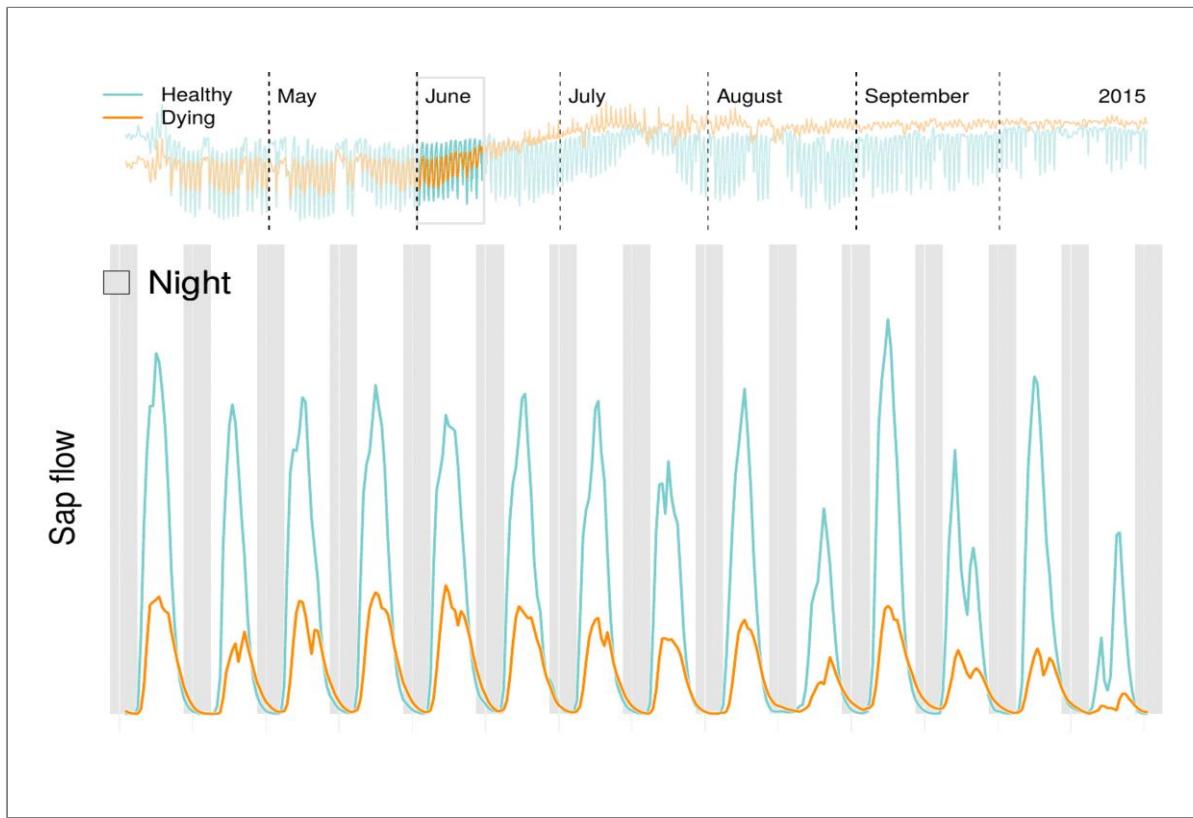
Friday, August 19, 2022

10:00 AM – 11:30 AM EDT

Location: 514C

## Goal?

*Facilitate the work with raw time series data and make methods for their processing accessible and reproducible*

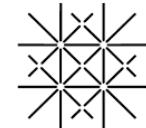


Christoforos Pappas



ΠΑΝΕΠΙΣΤΗΜΙΟ  
**ΠΑΤΡΩΝ**  
UNIVERSITY OF PATRAS

Richard L. Peters



University  
of Basel

Alexander G. Hurley



Roman Zweifel





## Outline

### (1) *Data handling & cleaning*

`datacleanr`

<https://the-hull.github.io/datacleanr>

### (2) *Sap flow data processing*

`TREX`

<https://the-hull.github.io/TREX>

### (3) *Dendrometer data processing*

`treenetproc`

<https://github.com/treenet/treenetproc>

### (4) *Hands-on & Q/A*



## Outline

### (1) *Data handling & cleaning*

<https://deep-tools.netlify.app/>



### (2) *Sap flow data processing*

#### Workshops

Teaching sessions and talks

#### Time series processing in R: from raw measurements to their analysis and interpretation

We show how to facilitate the work with raw time series data and make methods for their processing accessible and reproducible.

Jun 23, 2022 10:00 – 12:00 · Rapide-Danseur, Quebec, Canada  
Christoforos Pappas, Richard L. Peters, Alexander Hurley

[PDF](#) [View Session](#)



### (3) *Dendrometer data processing*

#### A Comprehensive Toolbox for Tree Physiological Data Processing in R

Reproducible cleaning, processing and assessment of high-frequency time series data sets.

Aug 6, 2021 13:00 – 14:30 · Virtual Meeting  
Christoforos Pappas, Alexander Hurley, Richard L. Peters, Roman Zweifel

[PDF](#) [View Materials](#) [View Files](#)



### (4) *Hands-on & Q/A*

#### Time series data processing course

Introduction to methods for processing time series in ecological research

Feb 9, 2021 08:00 – 12:00 · Virtual Meeting - University of Helsinki  
Richard L. Peters, Alexander Hurley

[PDF](#) [View Materials](#) [View Files](#)



## Relevance

Ecological research is becoming increasingly data-rich!

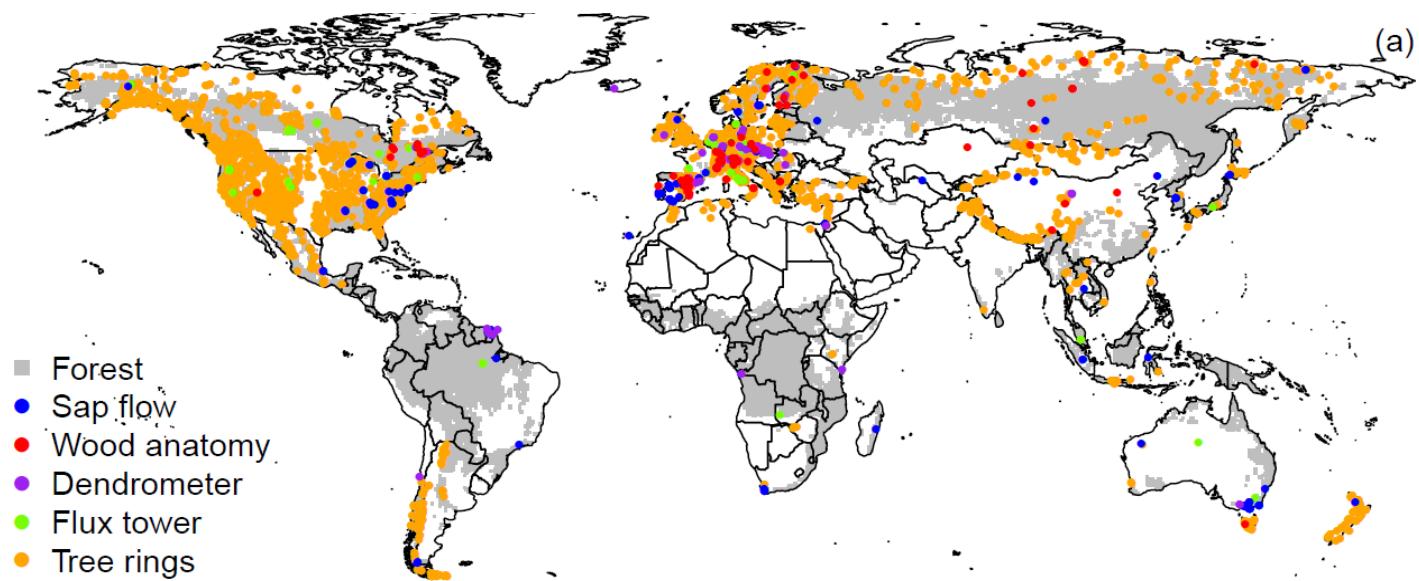
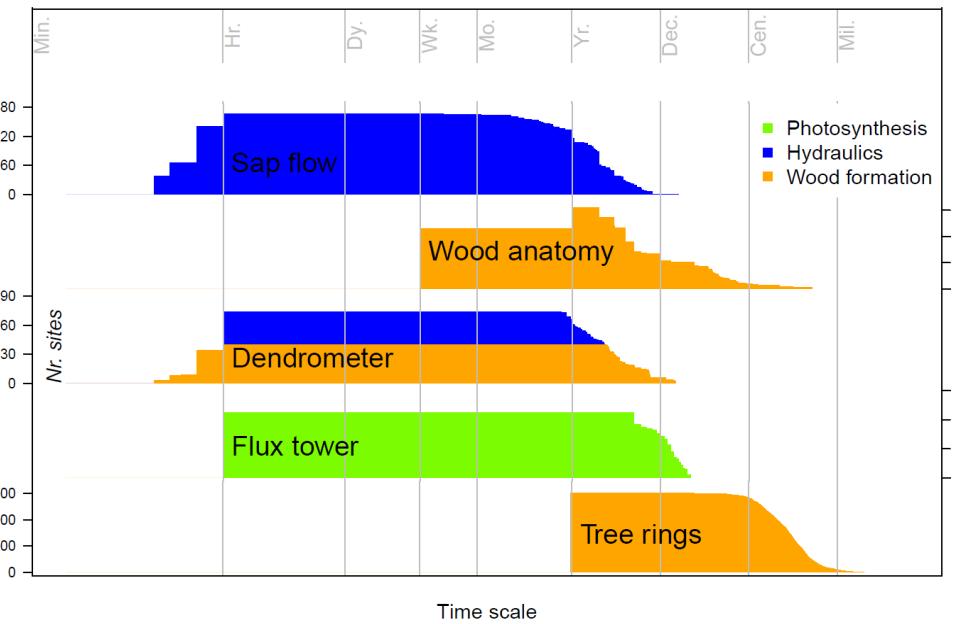
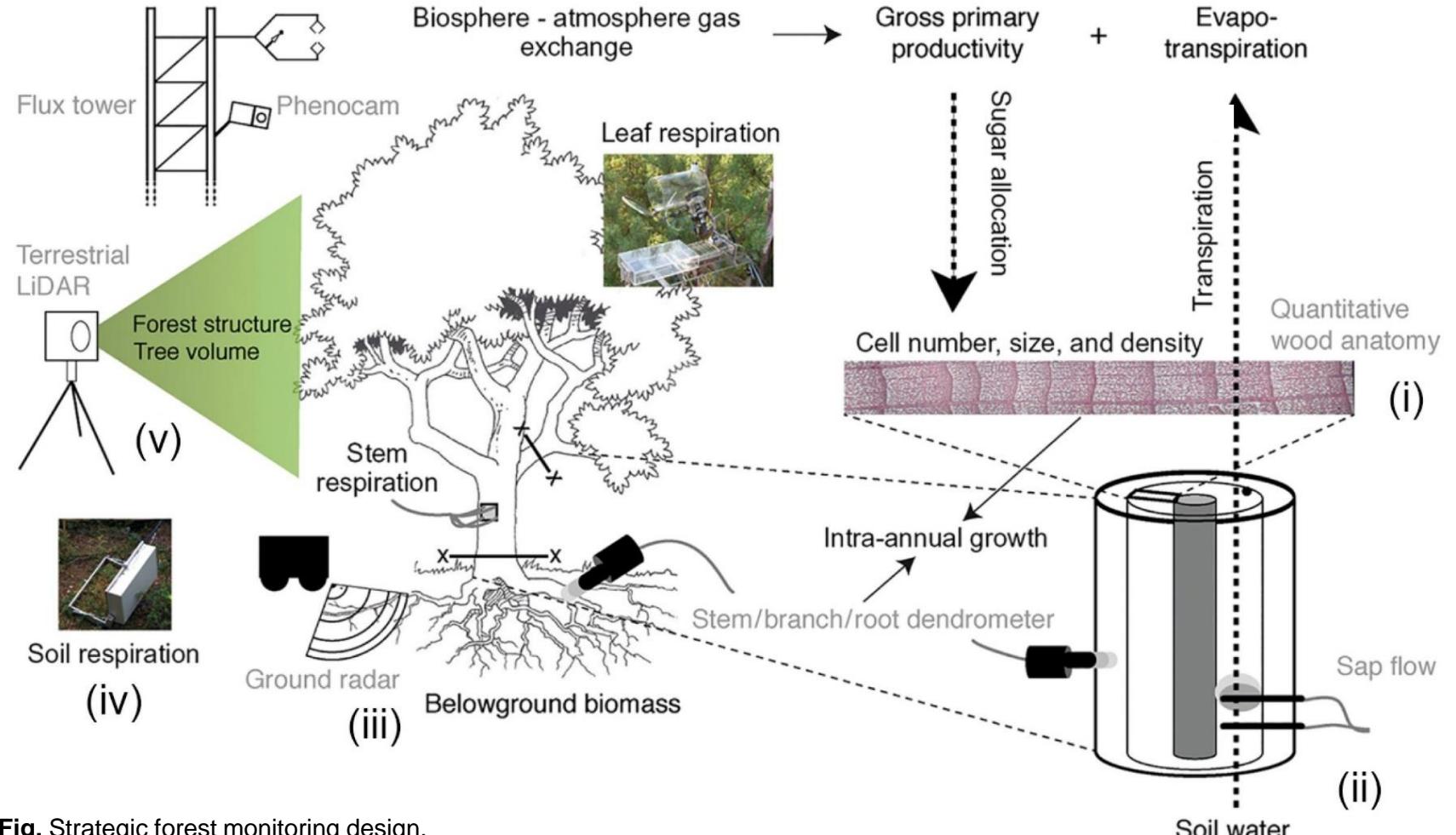


Fig. Global distribution of sites with relevant tree physiological measurements.

Source: Peters (2018) doi: <https://doi.org/10.5451/unibas-007085812>



## Relevance



**Fig.** Strategic forest monitoring design.

**Source:** Babst et al. (2021) doi: <https://doi.org/10.1016/j.tplants.2020.10.002>



## Data handling

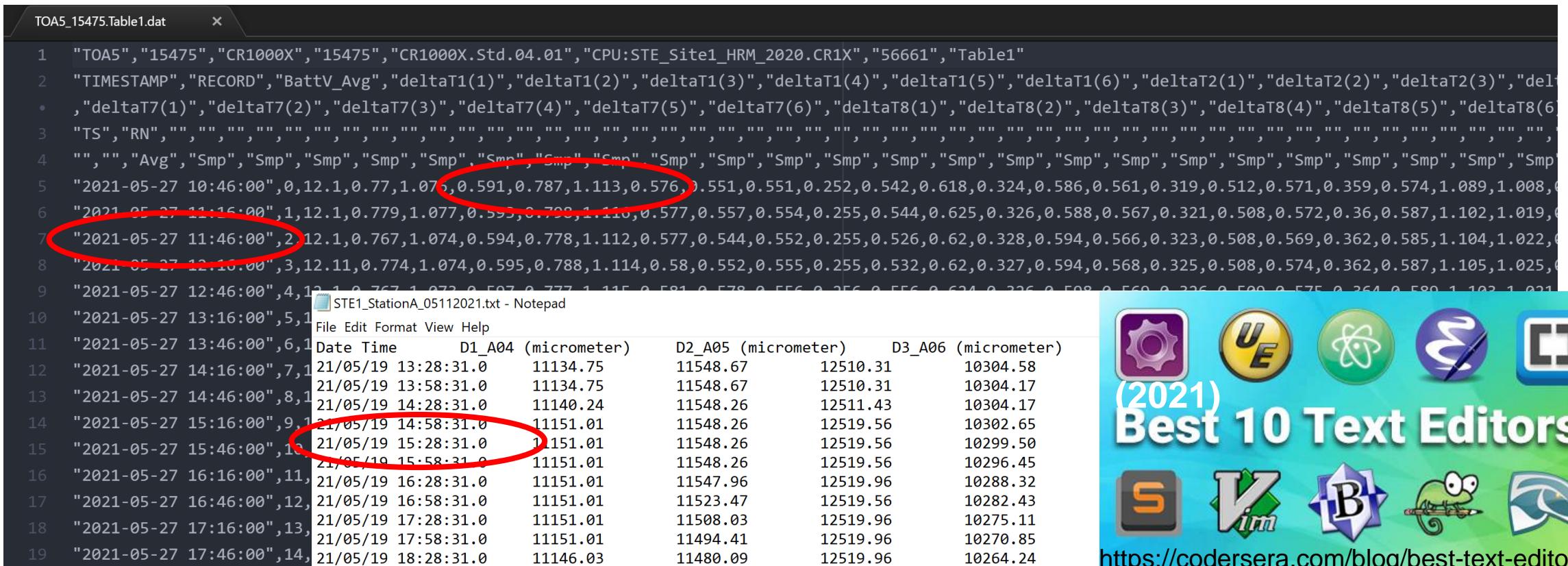
**Goal?** *Raw data import in R, data organization, aggregation, format homogenization.*



# Time series handling

## Raw data

- Typically, text files/spreadsheets
- Delimited formats (tab, space, ...), data types (numeric, strings, factors, ...)
- Quick view with a text editor



Date	Time	D1_A04 (micrometer)	D2_A05 (micrometer)	D3_A06 (micrometer)
2021-05-27	10:46:00	12.1, 0.77, 1.07	0.591, 0.787, 1.113	0.576, 0.551, 0.551
2021-05-27	11:16:00	12.1, 0.779, 1.077	0.593, 0.793, 1.118	0.577, 0.557, 0.554
2021-05-27	11:46:00	12.1, 0.767, 1.074	0.594, 0.778, 1.112	0.577, 0.544, 0.552
2021-05-27	12:16:00	12.11, 0.774, 1.074	0.595, 0.788, 1.114	0.58, 0.552, 0.555
2021-05-27	12:46:00	12.1, 0.767, 1.072	0.597, 0.777, 1.115	0.581, 0.578, 0.556
2021-05-27	13:16:00	12.1, 0.767, 1.072	0.597, 0.777, 1.115	0.581, 0.578, 0.556
2021-05-27	13:46:00	12.1, 0.767, 1.072	0.597, 0.777, 1.115	0.581, 0.578, 0.556
2021-05-27	14:16:00	12.1, 0.767, 1.072	0.597, 0.777, 1.115	0.581, 0.578, 0.556
2021-05-27	14:46:00	12.1, 0.767, 1.072	0.597, 0.777, 1.115	0.581, 0.578, 0.556
2021-05-27	15:16:00	12.1, 0.767, 1.072	0.597, 0.777, 1.115	0.581, 0.578, 0.556
2021-05-27	15:46:00	12.1, 0.767, 1.072	0.597, 0.777, 1.115	0.581, 0.578, 0.556
2021-05-27	16:16:00	12.1, 0.767, 1.072	0.597, 0.777, 1.115	0.581, 0.578, 0.556
2021-05-27	16:46:00	12.1, 0.767, 1.072	0.597, 0.777, 1.115	0.581, 0.578, 0.556
2021-05-27	17:16:00	12.1, 0.767, 1.072	0.597, 0.777, 1.115	0.581, 0.578, 0.556
2021-05-27	17:46:00	12.1, 0.767, 1.072	0.597, 0.777, 1.115	0.581, 0.578, 0.556
2021-05-27	18:16:00	12.1, 0.767, 1.072	0.597, 0.777, 1.115	0.581, 0.578, 0.556



## Data import in R

- base R
- tidyverse

```
read.table(file, header = FALSE, sep = "", quote = "\"\"",  
          dec = ".", numerals = c("allow.loss", "warn.loss", "no.loss"),  
          row.names, col.names, as.is = !stringsAsFactors,  
          na.strings = "NA", colClasses = NA, nrows = -1,  
          skip = 0, check.names = TRUE, fill = !blank.lines.skip,  
          strip.white = FALSE, blank.lines.skip = TRUE,  
          comment.char = "#",  
          allowEscapes = FALSE, flush = FALSE,  
          stringsAsFactors = default.stringsAsFactors(),  
          fileEncoding = "", encoding = "unknown", text, skipNul = FALSE)  
  
read.csv(file, header = TRUE, sep = ",", quote = "\"\"",  
         dec = ".", fill = TRUE, comment.char = "", ...)  
  
read.csv2(file, header = TRUE, sep = ";", quote = "\"\"",  
          dec = ",", fill = TRUE, comment.char = "", ...)  
  
read.delim(file, header = TRUE, sep = "\t", quote = "\"\"",  
           dec = ".", fill = TRUE, comment.char = "", ...)  
  
read.delim2(file, header = TRUE, sep = "\t", quote = "\"\"",  
            dec = ",", fill = TRUE, comment.char = "", ...)
```



readr

- [read\\_csv\(\)](#): comma-separated values (CSV) files
- [read\\_tsv\(\)](#): tab-separated values (TSV) files
- [read\\_delim\(\)](#): delimited files (CSV and TSV are important special cases)
- [read\\_fwf\(\)](#): fixed-width files
- [read\\_table\(\)](#): whitespace-separated files
- [read\\_log\(\)](#): web log files

<https://readr.tidyverse.org/>

<https://www.rdocumentation.org/packages/utils/versions/3.6.2/topics/read.table>



# ≡ Time series handling

## Dates

- time stamps

<https://rawgit.com/rstudio/cheatsheets/main/lubridate.pdf>

2017-11-28T14:02:00

**ymd\_hms()**, **ymd\_hm()**, **ymd\_h()**.  
ymd\_hms("2017-11-28T14:02:00")

2017-22-12 10:00:00

**ydm\_hms()**, **ydm\_hm()**, **ydm\_h()**.  
ydm\_hms("2017-22-12 10:00:00")

11/28/2017 1:02:03

**mdy\_hms()**, **mdy\_hm()**, **mdy\_h()**.  
mdy\_hms("11/28/2017 1:02:03")

1 Jan 2017 23:59:59

**dmy\_hms()**, **dmy\_hm()**, **dmy\_h()**.  
dmy\_hms("1 Jan 2017 23:59:59")

20170131

**ymd()**, **ydm()**. ymd(20170131)

July 4th, 2000

**mdy()**, **myd()**. mdy("July 4th, 2000")

4th of July '99

**dmy()**, **dym()**. dmy("4th of July '99")

2018-01-31 11:59:59

**day(x)** Day of month. day(dt)  
**wday(x, label, abbr)** Day of week.  
**qday(x)** Day of quarter.

2018-01-31 11:59:59

**hour(x)** Hour. hour(dt)

2018-01-31 11:59:59

**minute(x)** Minutes. minute(dt)

2018-01-31 11:59:59

**second(x)** Seconds. second(dt)



**lubridate**

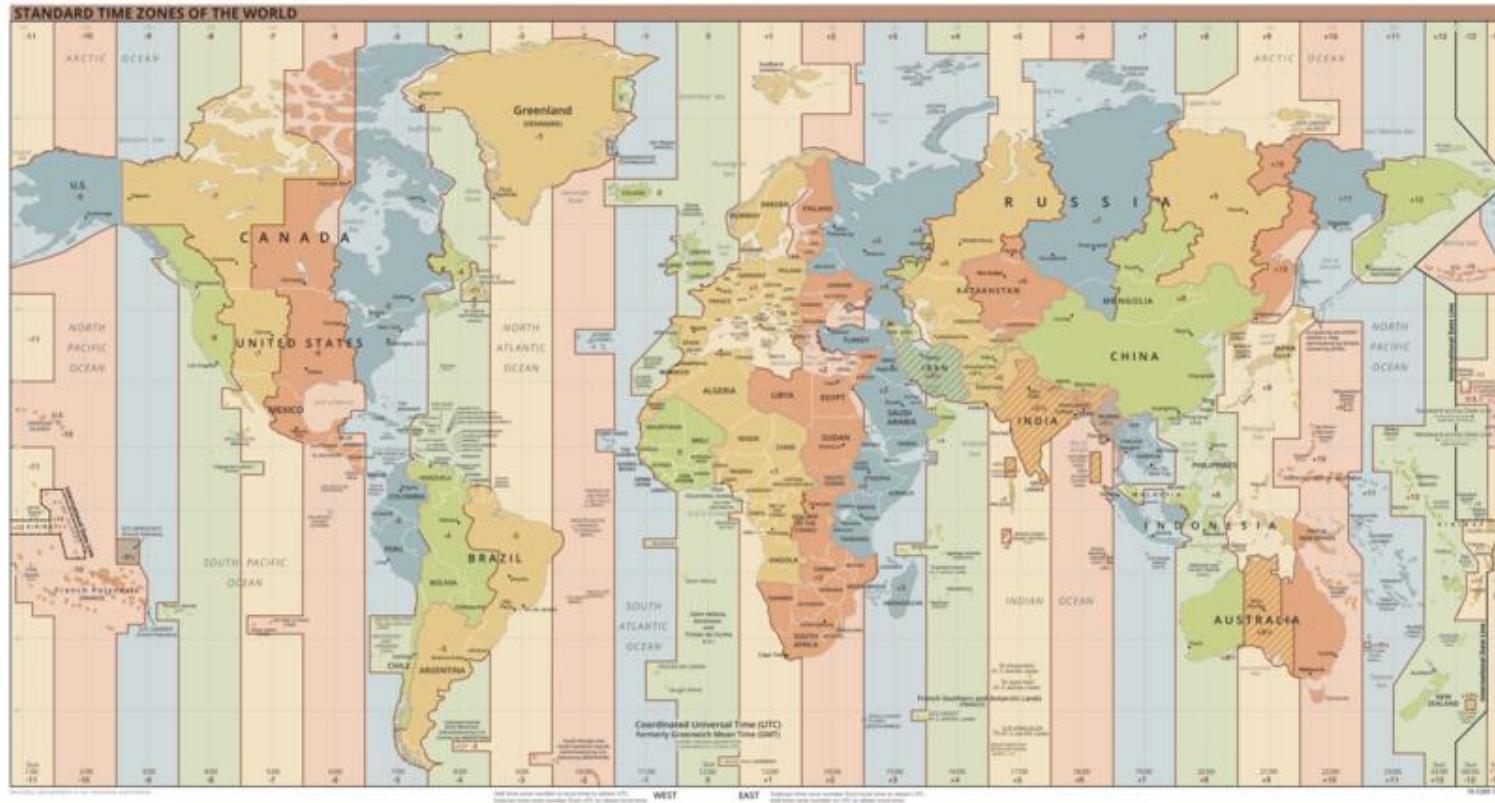
## More examples

- base R: [https://deep-tools.netlify.app/docs-workshops/uhelsinki-workshop2021/01\\_datacleanr/](https://deep-tools.netlify.app/docs-workshops/uhelsinki-workshop2021/01_datacleanr/)
- lubridate: <https://lubridate.tidyverse.org/>

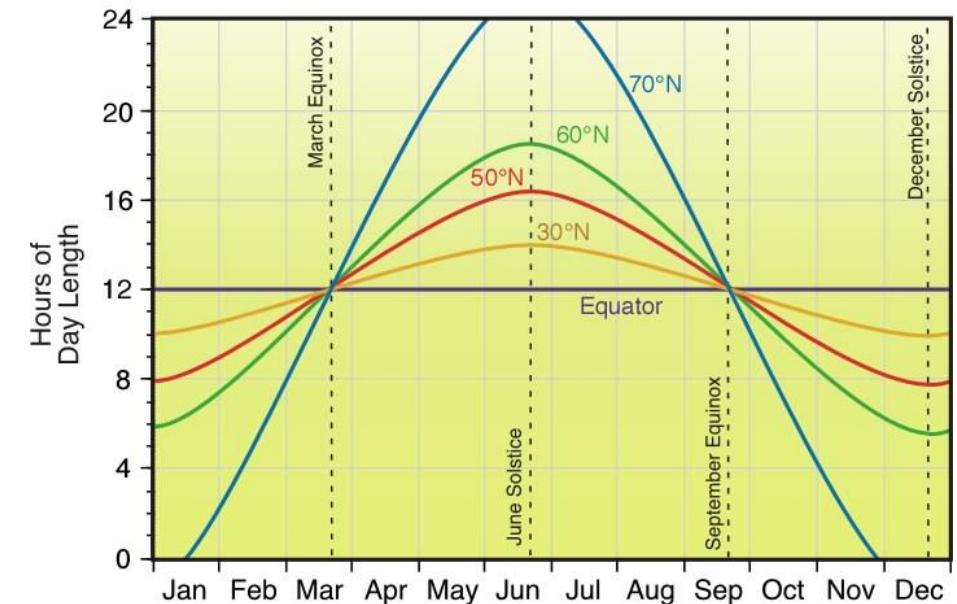
# ≡ Time series handling

## Dates

- time zones (daylight savings ?)
- solartime (<https://github.com/bgctw/solartime>; noon when sun is at zenith)



[https://commons.wikimedia.org/wiki/File:World\\_Time\\_Zones\\_Map.png](https://commons.wikimedia.org/wiki/File:World_Time_Zones_Map.png)

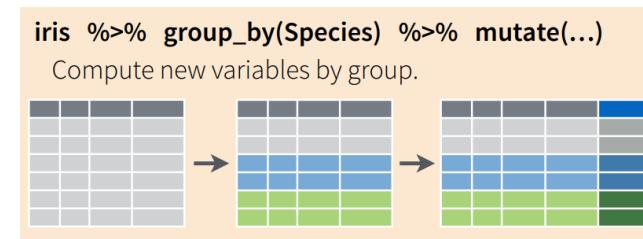


<http://www.physicalgeography.net/fundamentals/6i.html>

# ≡ Time series handling

## Temporal aggregation

- base R
- dplyr



### 30-min data

site	sensor	tree	depth	timestamp	time	tz	value
STE1	deltaT4_out	4	out	2020-08-24 18:46:00	2020-08-24 18:46:00	EDT	-0.025268594
STE1	deltaT4_out	4	out	2020-08-24 19:16:00	2020-08-24 19:16:00	EDT	-0.027292142
STE1	deltaT4_out	4	out	2020-08-24 19:46:00	2020-08-24 19:46:00	EDT	-0.025367214
STE1	deltaT4_out	4	out	2020-08-24 20:16:00	2020-08-24 20:16:00	EDT	-0.029156584
STE1	deltaT4_out	4	out	2020-08-24 20:46:00	2020-08-24 20:46:00	EDT	-0.029327615
STE1	deltaT4_out	4	out	2020-08-24 21:16:00	2020-08-24 21:16:00	EDT	-0.031191612
STE1	deltaT4_out	4	out	2020-08-24 21:46:00	2020-08-24 21:46:00	EDT	-0.031252544

### 1-h data

sensor	site	date	doy	hour	tree	depth	tz	value
deltaT4_out	STE1	2020-08-24 17:00:00	237	17	4	out	EDT	-0.0232645516
deltaT4_out	STE1	2020-08-24 18:00:00	237	18	4	out	EDT	-0.0223245229
deltaT4_out	STE1	2020-08-24 19:00:00	237	19	4	out	EDT	-0.0263296781
deltaT4_out	STE1	2020-08-24 20:00:00	237	20	4	out	EDT	-0.0292420997

```
# aggregate to hourly data
data_h = data %>%
  group_by(sensor, site,
    date=ymd_h(substr(data$timestamp, 1, 13)),
    doy=yday(data$timestamp),
    hour=hour(data$timestamp)) %>%
  summarise(tree = unique(tree),
            depth=unique(depth),
            tz=unique(tz),
            value=mean(value, na.rm=T))
```

## More examples

- aggregate(·) <https://deep-tools.netlify.app/talk/uhelsinki-2021-rpeters-ahurley/>
- dplyr | group\_by(·) <https://www.rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheatsheet.pdf>

## ≡ Data ‘cleaning’

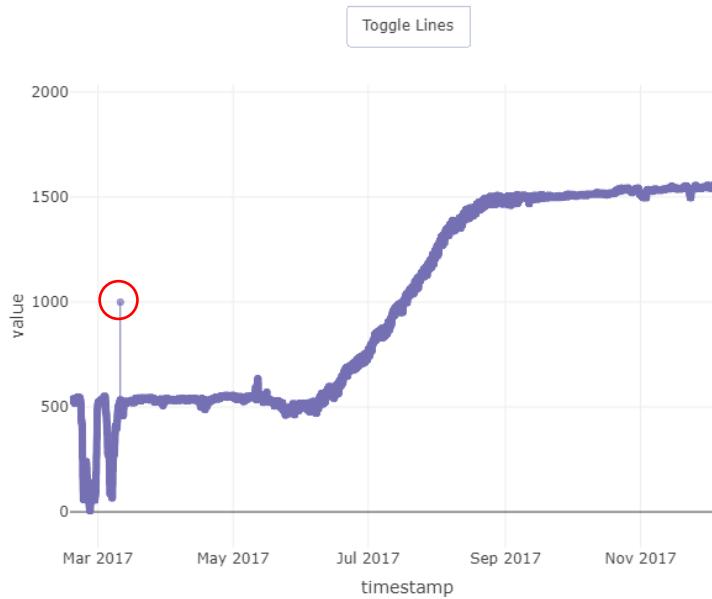
**Goal?** *Interactive and reproducible data inspection, visualization and cleaning.*



# ≡ Typical data issues

**Outlier and sensor failure issues** *Removing data should always be done with care!*

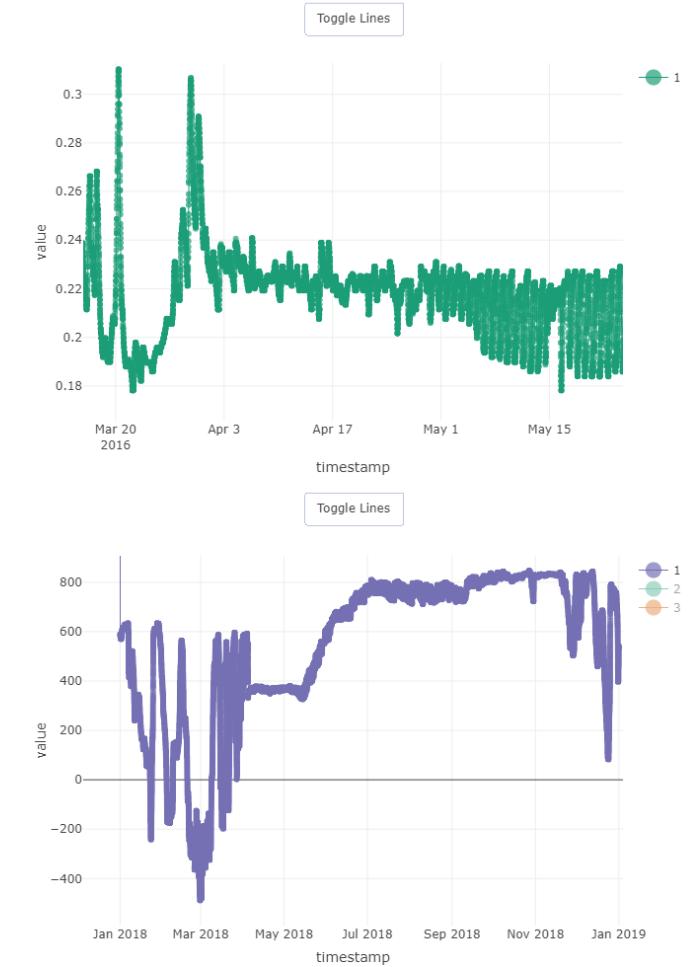
Outliers (dendrometer data)



Sensor failure (sap flow data)



What about spring data?



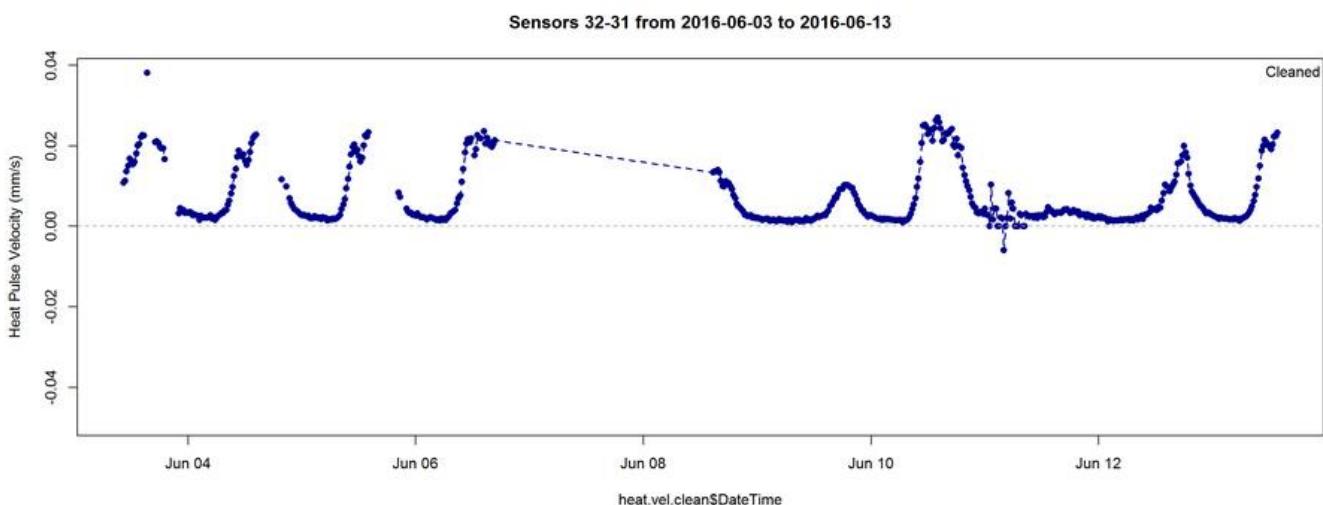
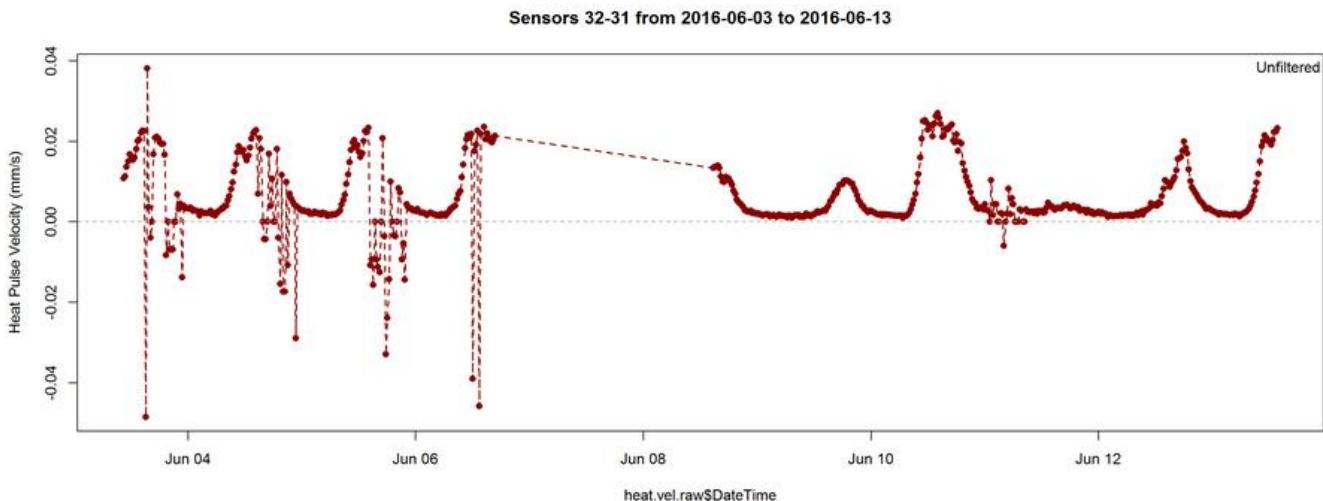


# Data ‘cleaning’ – ways forward?

## Manual processing:

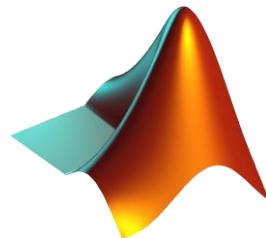
- *Quality Control with diagnostic plots*
  - *Time Series*
  - *X/Y (e.g., Sapflow vs. Environment)*
- *Quality Assurance with code:*
  - *Filters*
    - *Moving average*
    - *Threshold*
    - *Outlier detection*
  - *Adjustments*

- Long series & many sites = much time!
- No “one-size-fits-all” solutions
- Must be done reproducibly



## (Reproducible) processing and cleaning:

- Several excellent **programmatic** tools
  - Spreadsheet software (!)
  - Software & libraries
  - Bespoke code



## Focus on R and packages!

- Extensively used in *Environmental Sciences*
  - interoperable



Long series, many sites.. ?

## A new R-based package

*A flexible and efficient tool for interactive data cleaning*



## PLOS ONE

PUBLISH ABOUT BROWSE

OPEN ACCESS PEER-REVIEWED

RESEARCH ARTICLE

### Addressing the need for interactive, efficient, and reproducible data processing in ecology with the datacleanr R package

Alexander G. Hurley, Richard L. Peters, Christoforos Pappas, David N. Steger, Ingo Heinrich

Published: May 12, 2022 • <https://doi.org/10.1371/journal.pone.0268426>

```
# grab all data from ZHANG
zhang <- cosore::csr_table("data", c("d20190424_ZHANG_maple",
                                         "d20190424_ZHANG_oak")) %>%
  # adjust for grouping
  mutate(CSR_PORT = as.factor(CSR_PORT))

# group by CSR_DATASET and CSR_PORT
datacleanr::dcr_app(zhang)
```

### Properties:

- Uses R (links with other packages);
- Freely available (avoid license costs);
- Uses R shiny (interactive approach).
- **Reproducible**

### Structure of the tool:

- Set-up & overview;
- Filtering;
- Visual cleaning and annotating;
- Extract.



# Data ‘cleaning’ with

datacleanr



## Versatile

*Tested on multiple temporal-and spatial-specific data sets*

The screenshot shows the RStudio interface with the datacleanr package loaded. The left pane displays an R script named 'Untitled5.R' containing the following code:

```
1 saveRDS(iris, file = "./testiris.Rds")
2
3 library(datacleanr)
4 dcr_app("./testiris.Rds")
```

The right pane shows the package documentation for `dcr_app`. The documentation includes:

- package**: a character vector giving the package(s) to look in for data sets, or `NULL`. By default, all packages in the search path are used, then the 'data' subdirectory (if present) of the current working directory.
- lib.loc**: a character vector of directory names of `R` libraries, or `NULL`. The default value of `NULL` corresponds to all libraries currently known.
- verbose**: a logical. If `TRUE`, additional diagnostics are printed.
- envir**: the environment where the data should be loaded.
- overwrite**: logical: should existing objects of the same name in `envir` be replaced?

**Details**

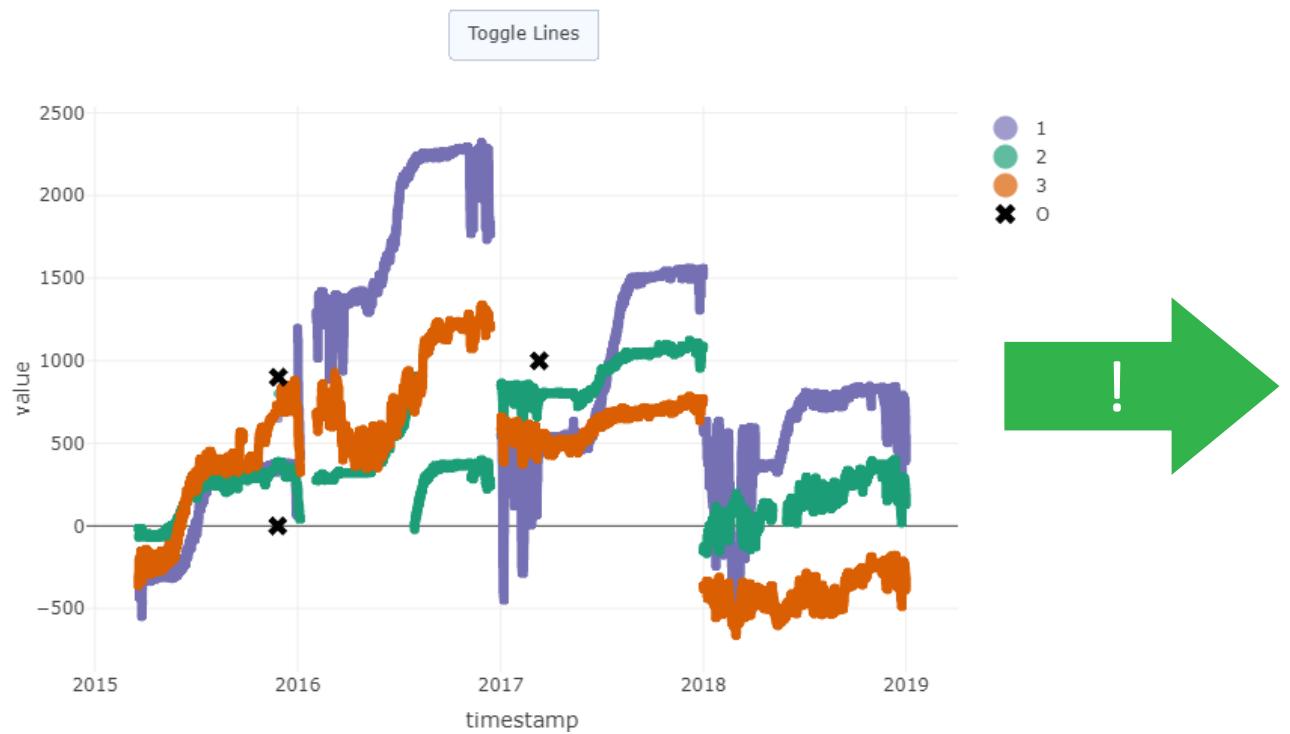
Currently, four formats of data files are supported:

1. files ending '`.r`' or '`.r`' are `source()`d in, with the `R` working directory changed temporarily to the directory containing the respective file. (`data` ensures that the `utils` package is attached, in case it had been run via `utils::data()`)
2. files ending '`.RData`' or '`.rda`' are `load()`ed.



## Extract

*Reproducible recipe to cook up some fresh data*



!

datacleanr Set-up & Overview Filtering Visual Cleaning & Annotating Extract Close Cancel

**Reproducible Recipe**  
[Click for Help](#)

All commands and operations in previous tabs are translated to code on the right, ensuring reproducibility.

Concise code?

[Send to RStudio](#) [Copy to clipboard](#)

**Set Output Locations**

Meta & Recipe  Same folder for cleaned data?

**Set and Save Outputs**

Suffix: Cleaned Data  
cleaned

Suffix: Filter + Outlier Data  
meta\_RAW

```
# datacleaning with datacleanr (1.0.1)
# ----- Sun Feb 07 10:13:50 2021 -----
library(dplyr)
library(datacleanr)

output_long <- readRDS("D:/Documents/UL - POSTDOC/02_communication/Education - Finland/Course -")

# adding column for unique IDs;
output_long$.dcrkey <- seq_len(nrow(output_long))

# observations from manual selection (Viz tab);
output_long_outlier_selection <- readRDS("D:/Documents/UL - POSTDOC/02_communication/Education -")

# create data set with annotation column (non-outliers are NA);
output_long <- dplyr::left_join(output_long, output_long_outlier_selection, by = ".dcrkey")

# remove comment below to drop manually selected obs in data set;
# output_long <- output_long %>% dplyr::filter(is.na(.annotation))

saveRDS(output_long, "D:/Documents/UL - POSTDOC/02_communication/Education - Finland/Course -")
```



## Additional Resources

Hurley AG, Peters RL, Pappas C, Steger DN, Heinrich I (2022)

**Addressing the need for interactive, efficient, and reproducible data processing in ecology with the datacleanr R package.** PLoS ONE 17(5): e0268426. <https://doi.org/10.1371/journal.pone.0268426>

<https://github.com/the-hull/datacleanr>

<https://deep-tools.netlify.app/#workshops>



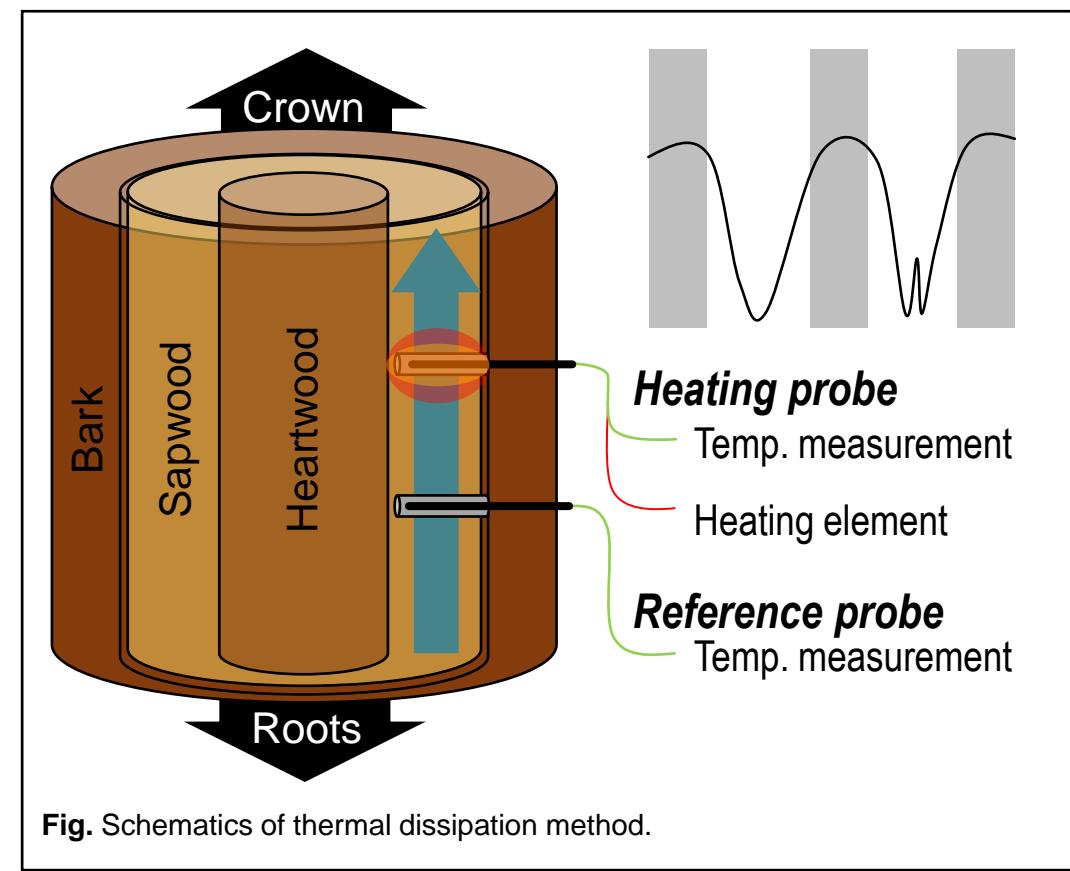
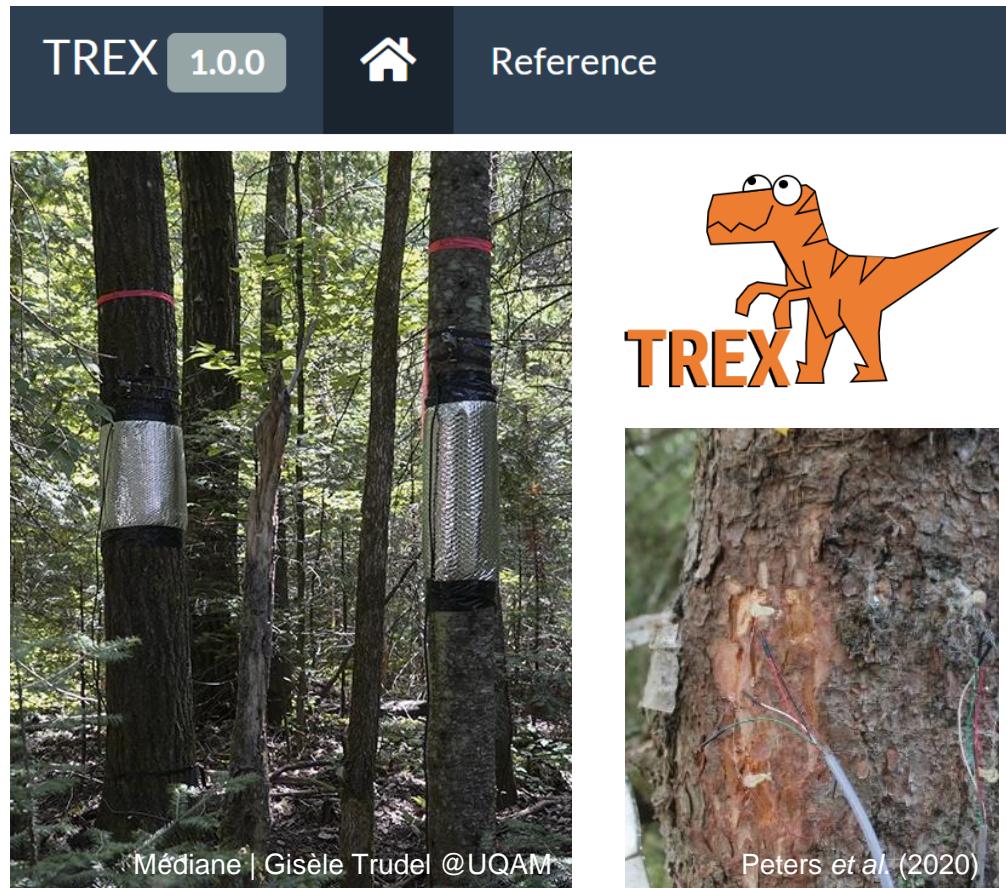
## Data processing

TREX

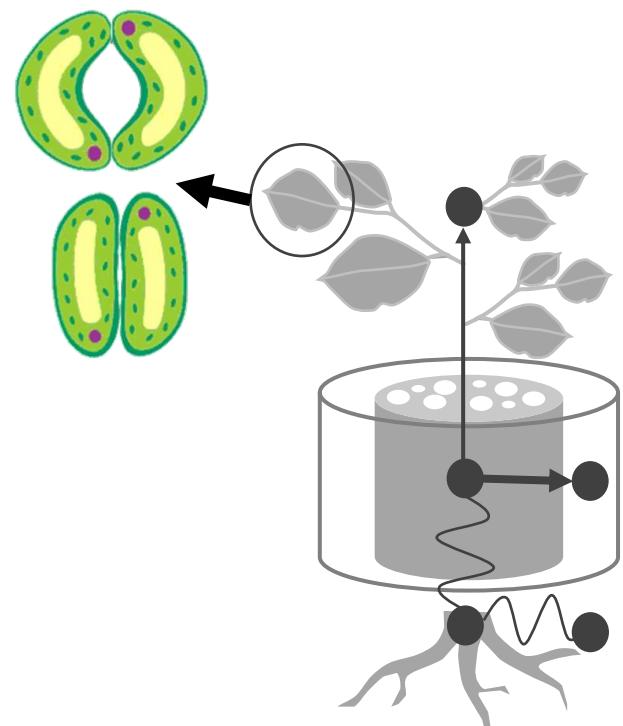
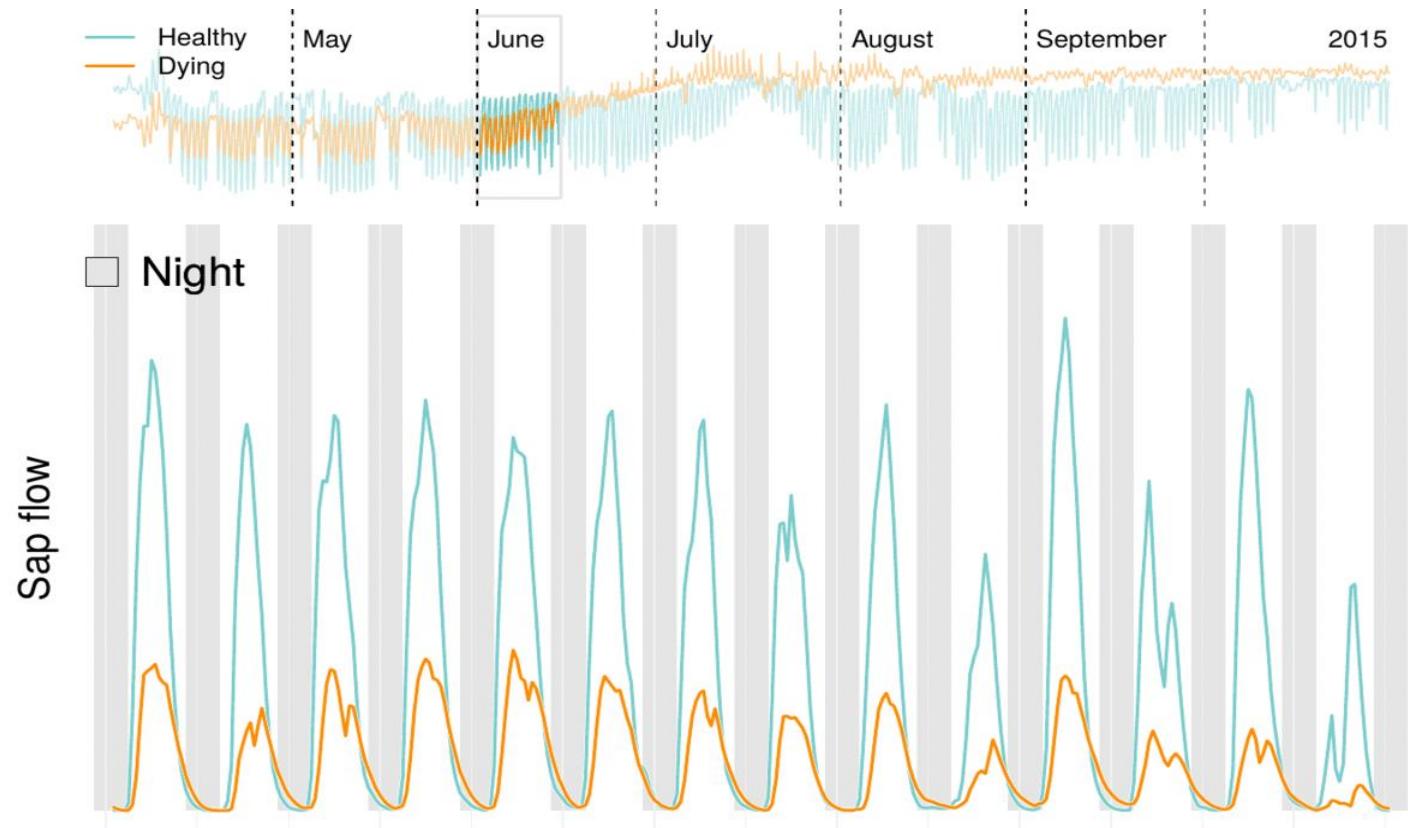
treenetproc

**Goal?** *Make standardized time series processing workflows  
for widely used monitoring techniques more accessible.*

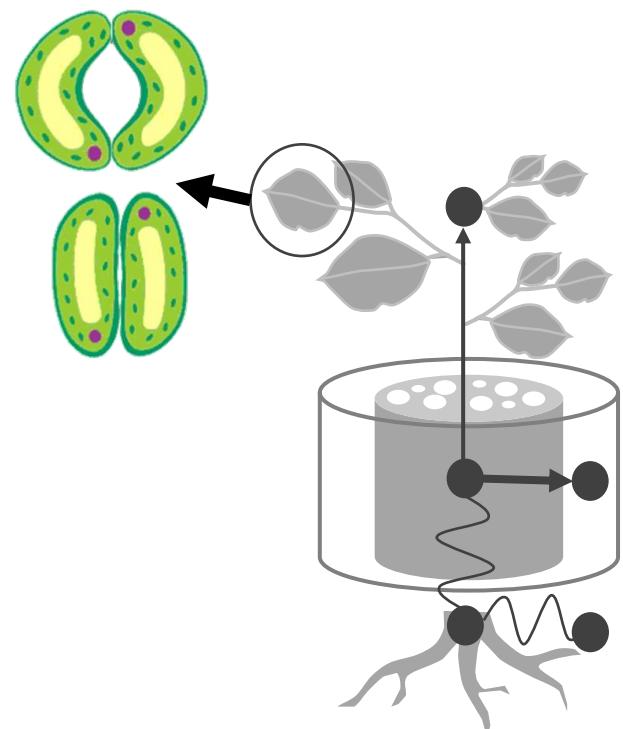
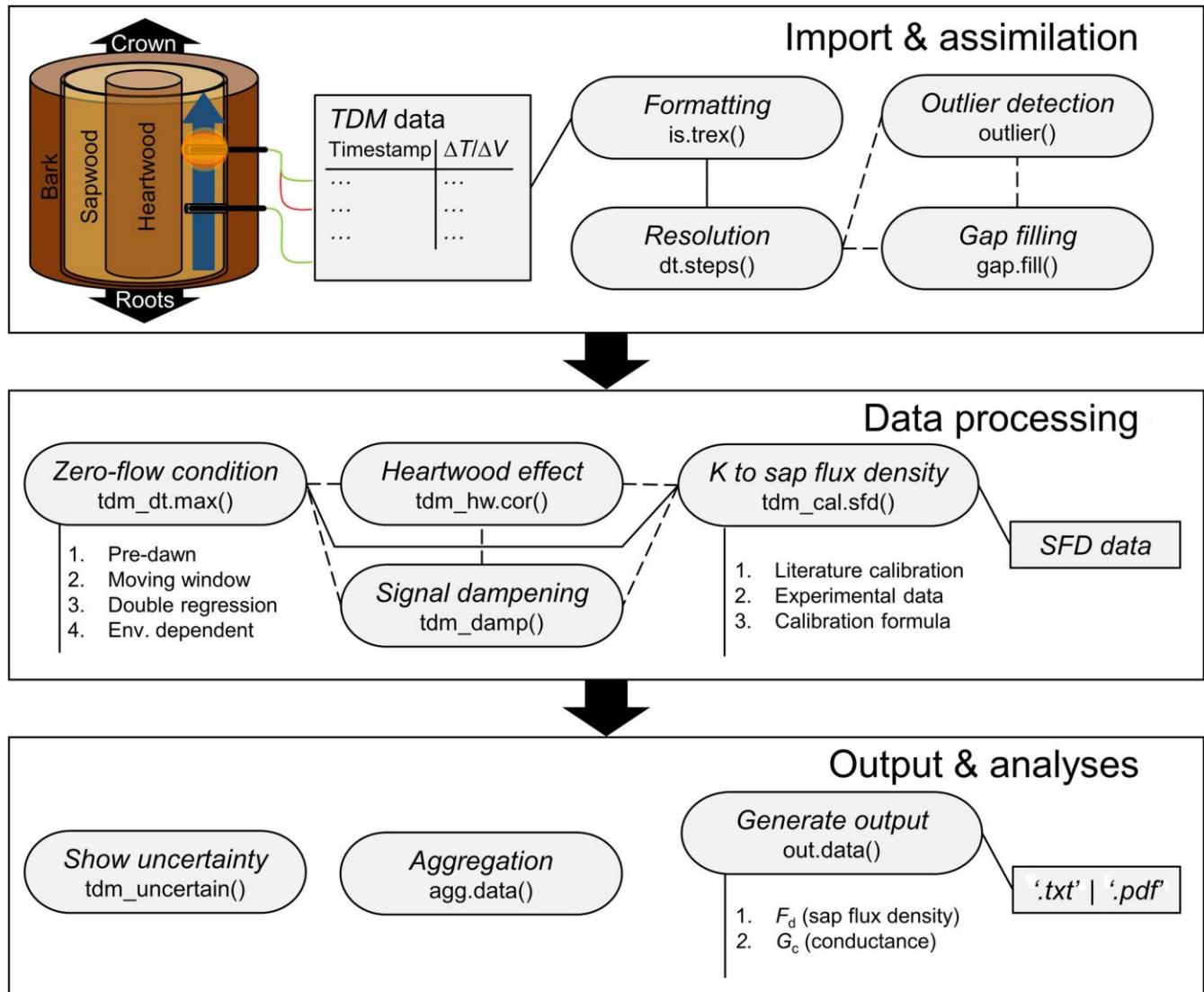
### Sap flow sensors - Utilizing the TREX R package (<https://the-hull.github.io/TREX>)



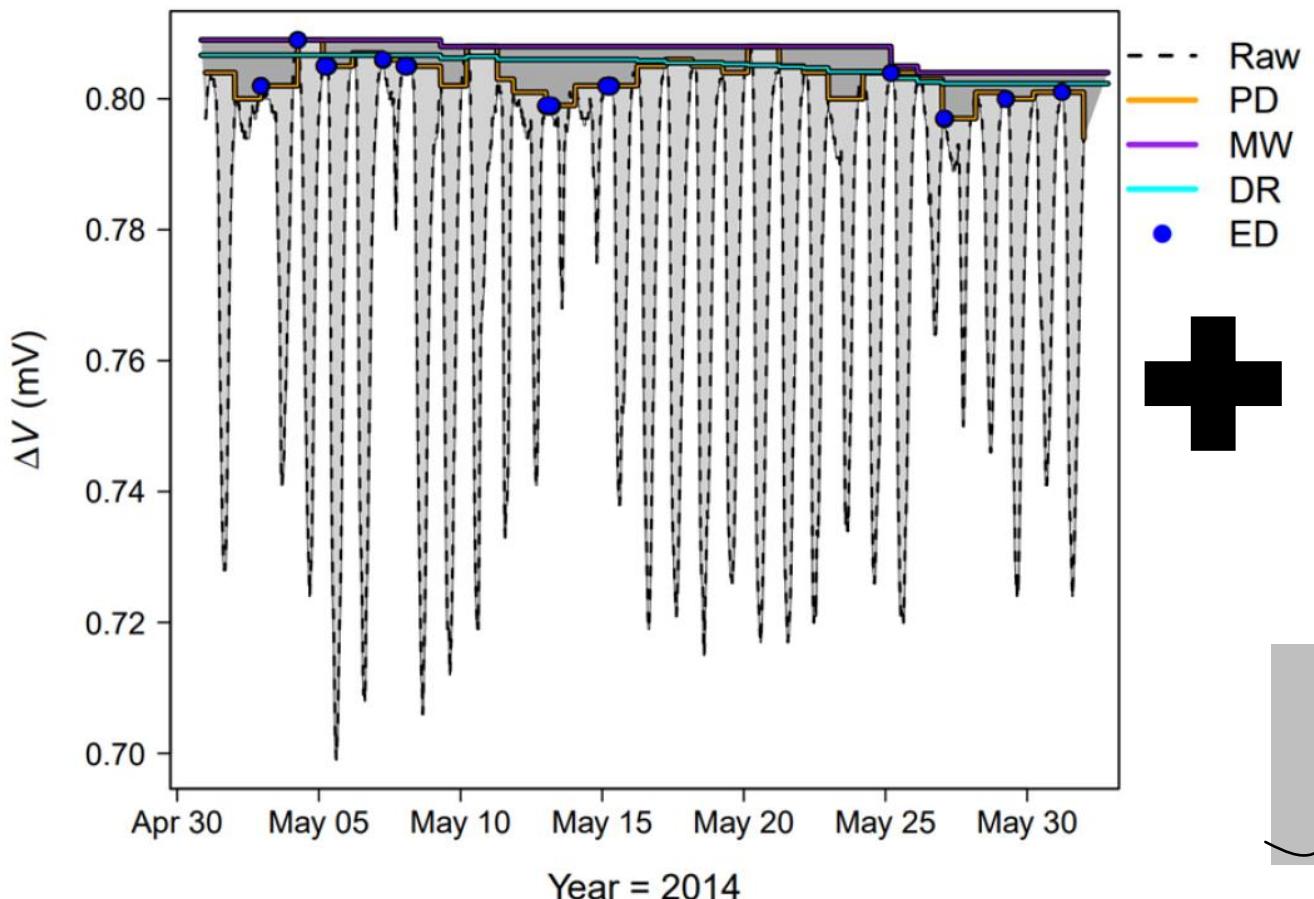
## Input – Thermal dissipation probe measurements *tdm.data()*



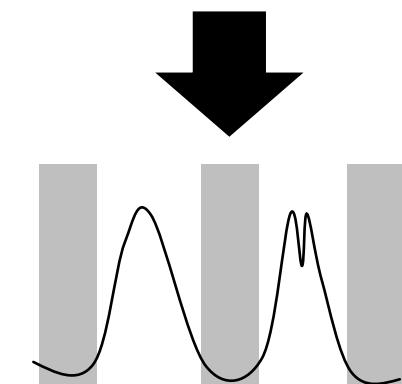
## Workflow



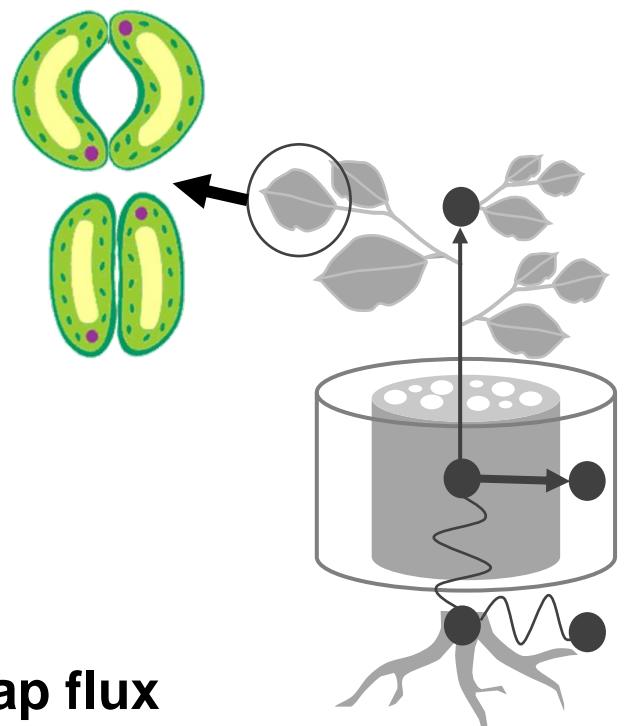
## Functionalities – Determining zero flow *tdm\_dt.max()*



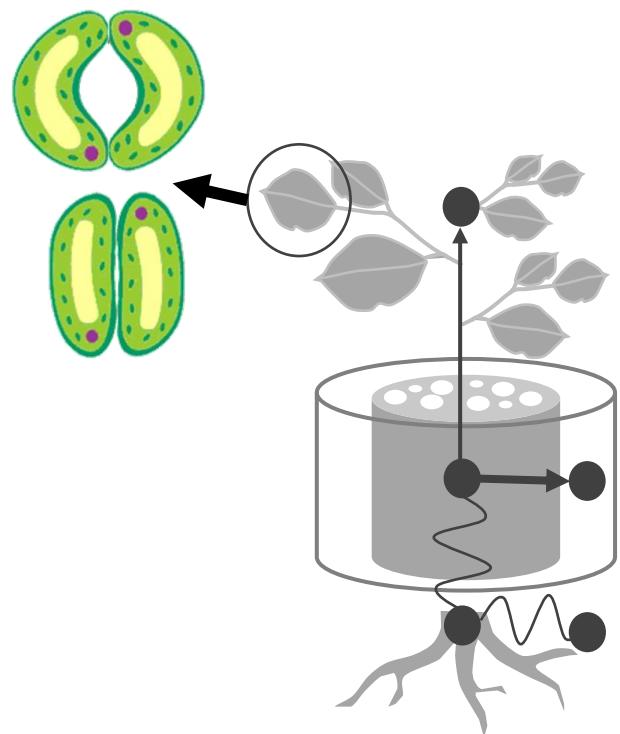
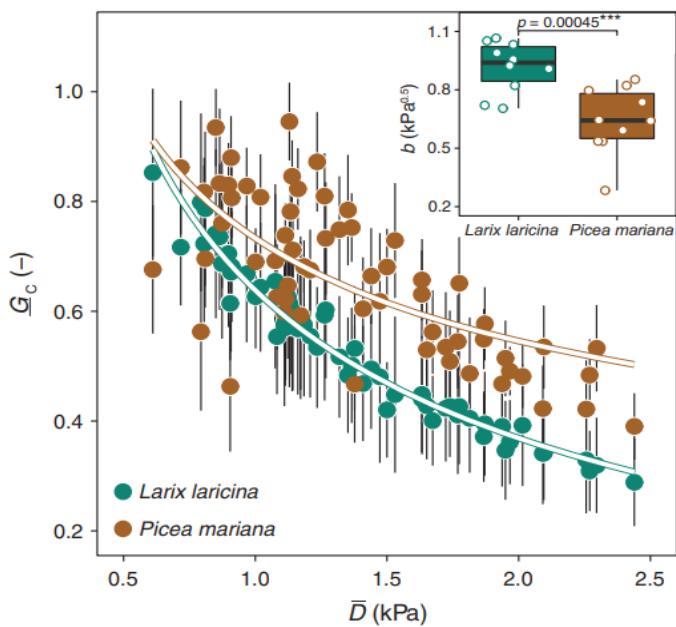
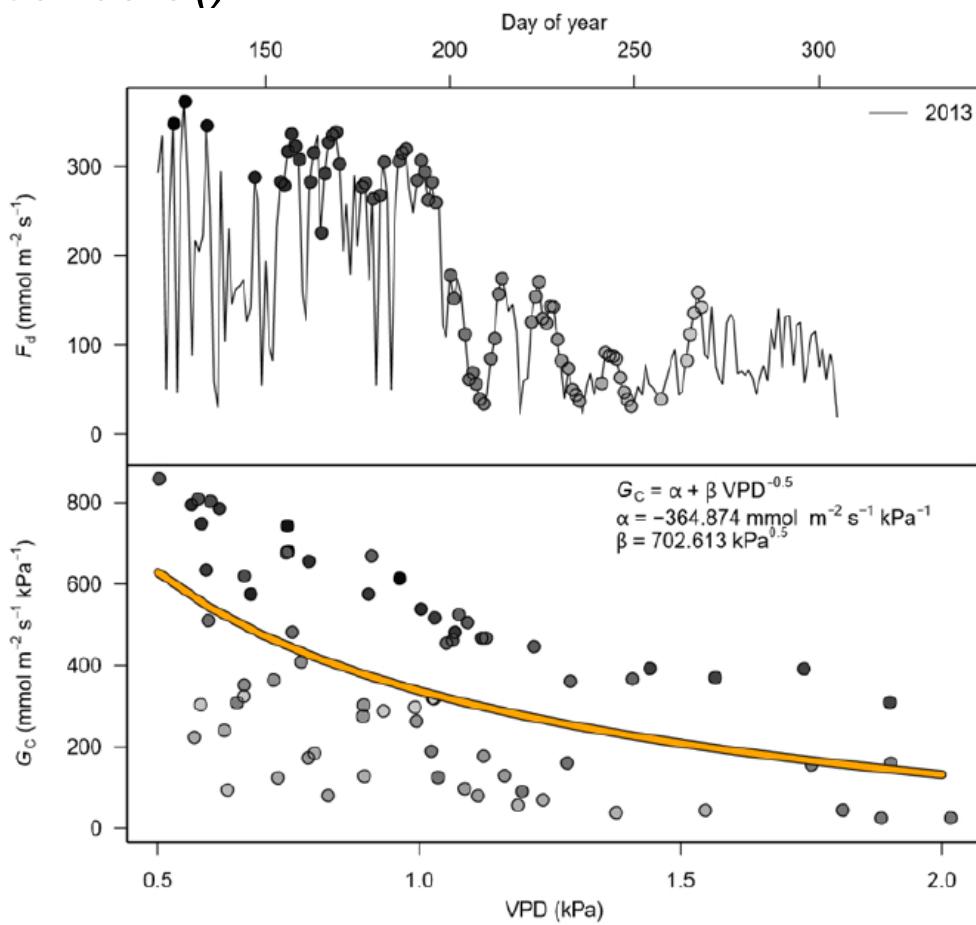
Calibration  
values  
*tdm\_cal.sfd()*



Sap flux  
density



## Output – Canopy conductance behavior *out.data()*



## Additional Resources

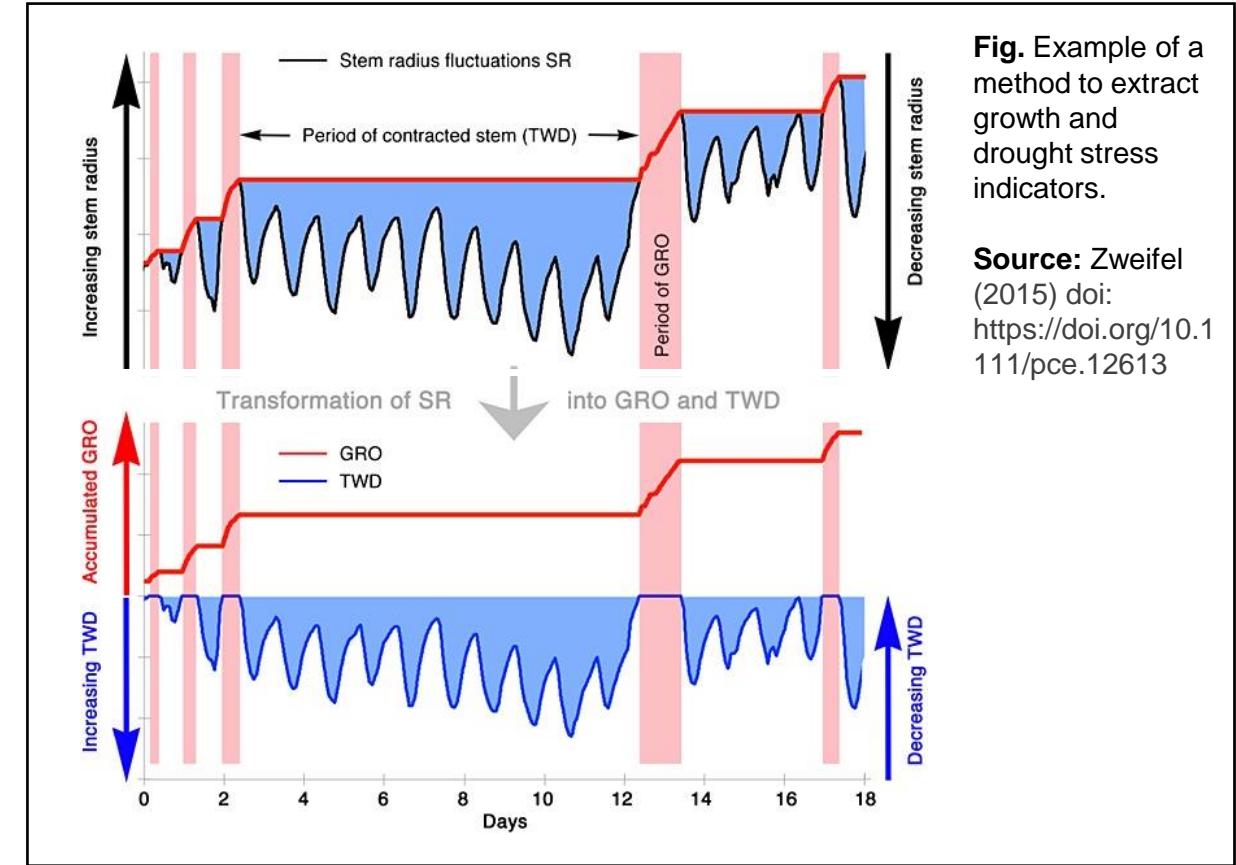
Peters, RL, Pappas, C, Hurley, AG, et al. (2021) **Assimilate, process and analyse thermal dissipation sap flow data using the TREX r package.** Methods Ecol Evol. <https://doi.org/10.1111/2041-210X.13524>

<https://the-hull.github.io/TREX/index.html>

<https://deep-tools.netlify.app/#workshops>



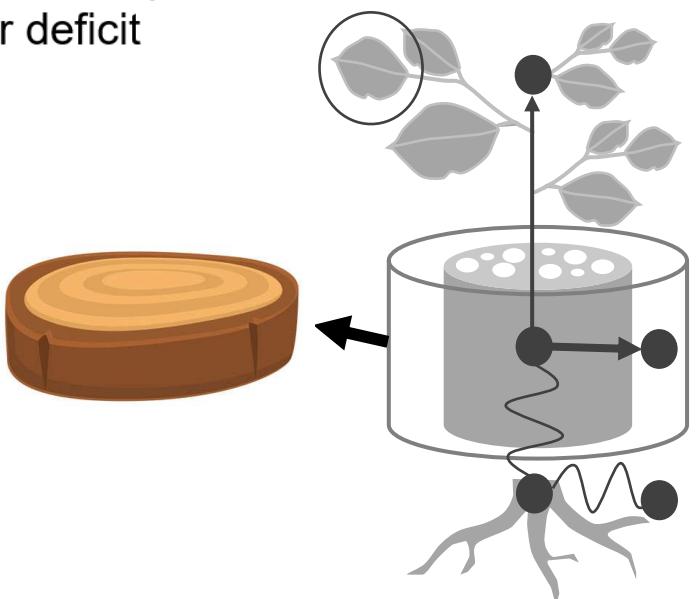
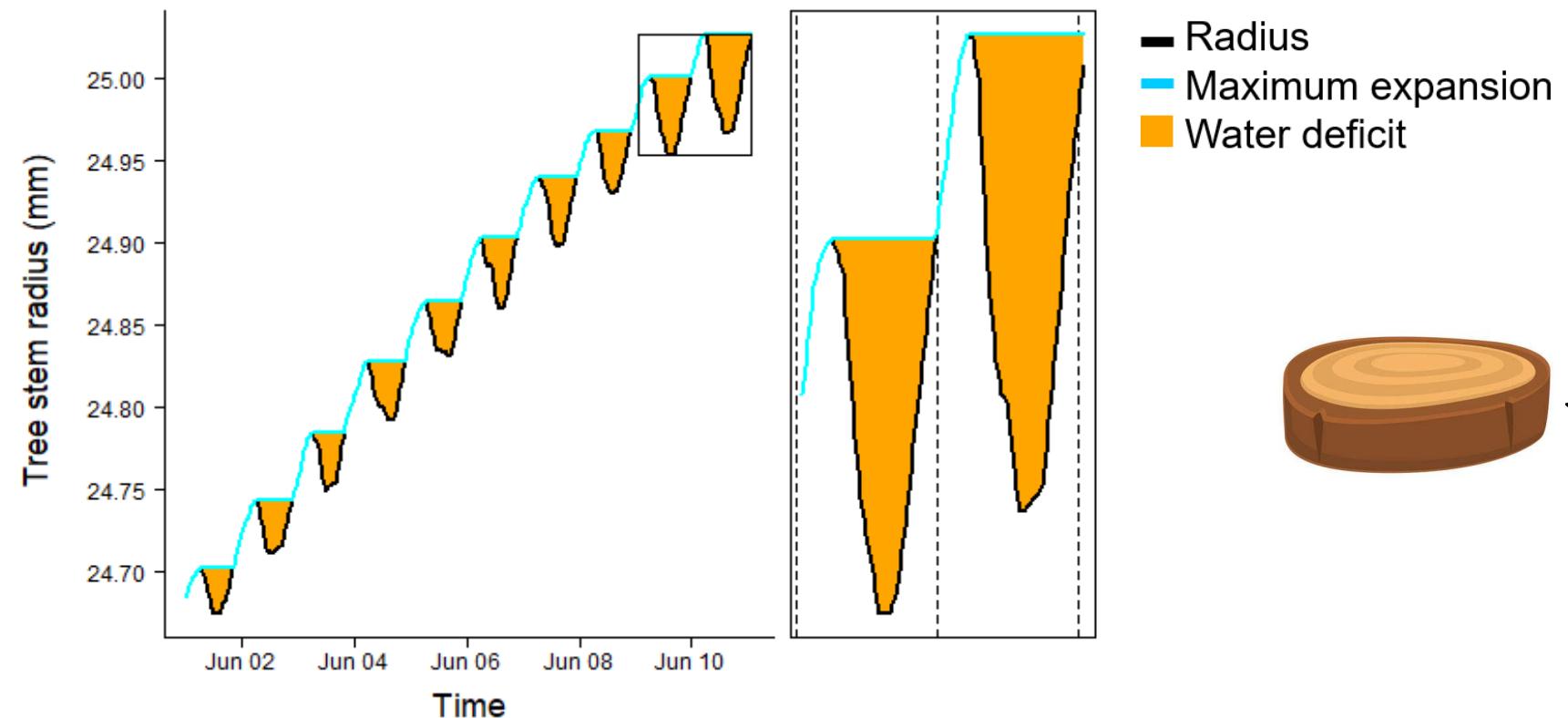
### Dendrometers - Utilizing the treenetproc R package (<https://github.com/treenet/treenetproc>)



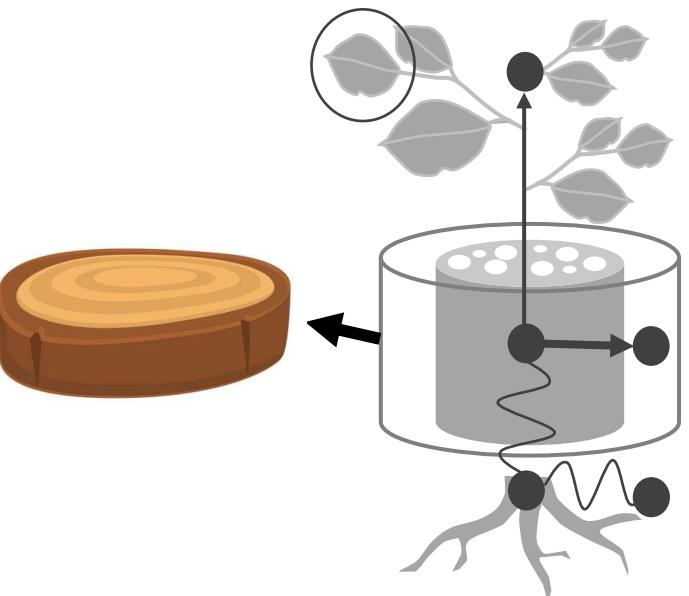
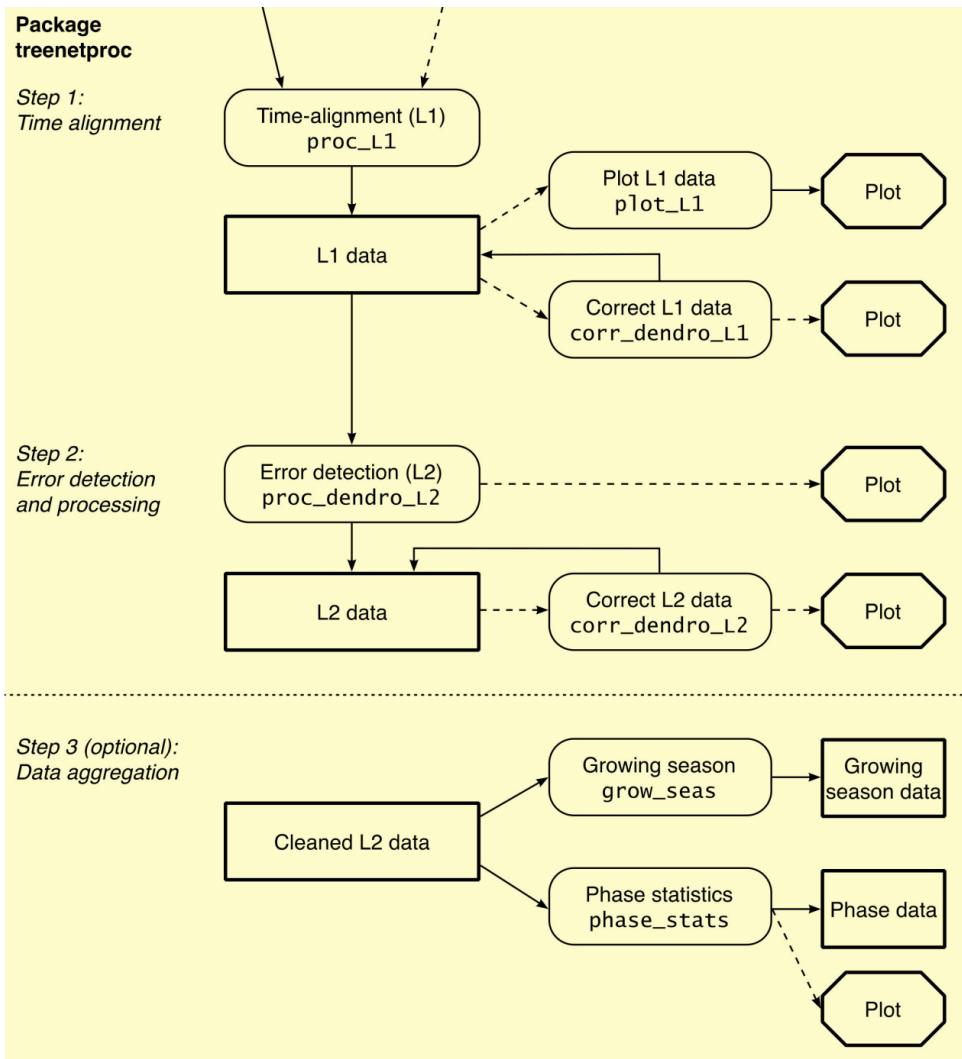
**Fig.** Example of a method to extract growth and drought stress indicators.

**Source:** Zweifel (2015) doi: <https://doi.org/10.111/pce.12613>

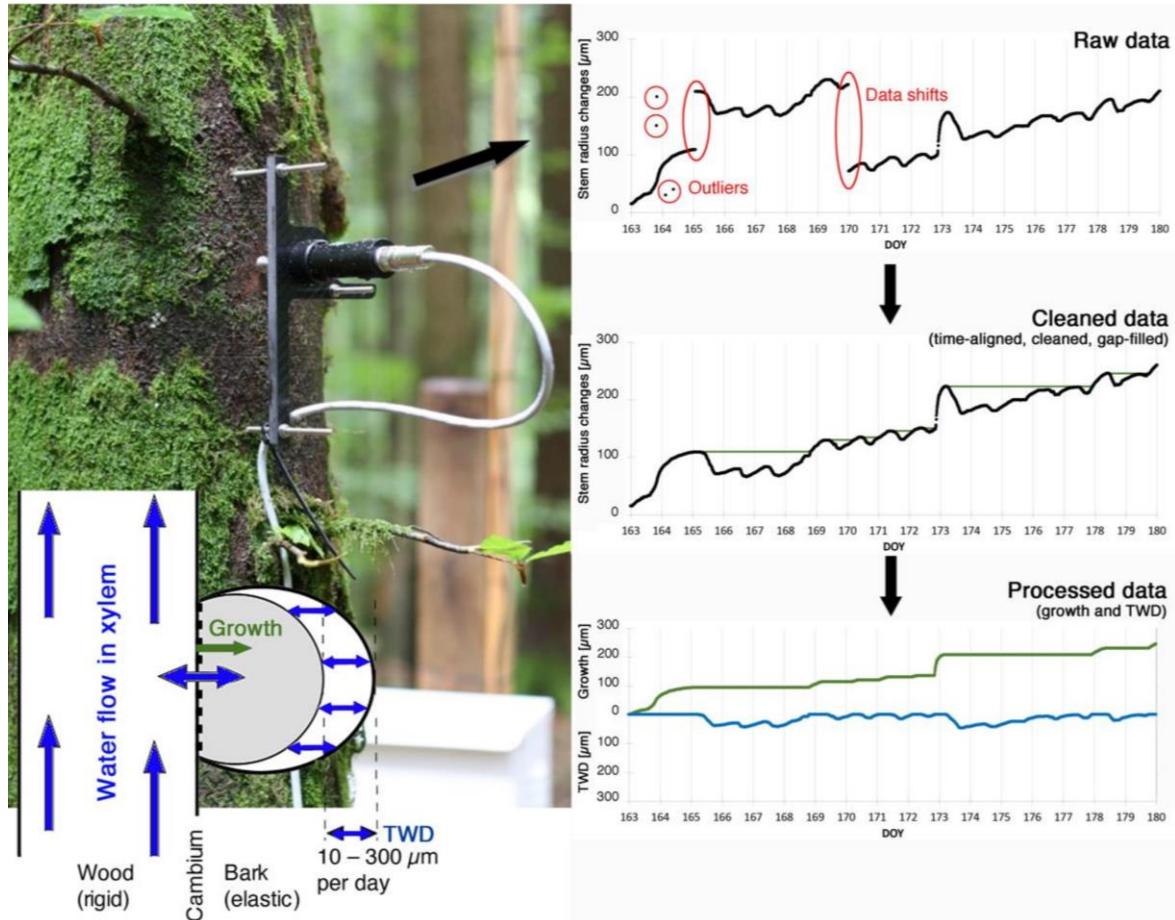
## Input – Point dendrometer measurements *dendro\_data\_L0()*



## Workflow



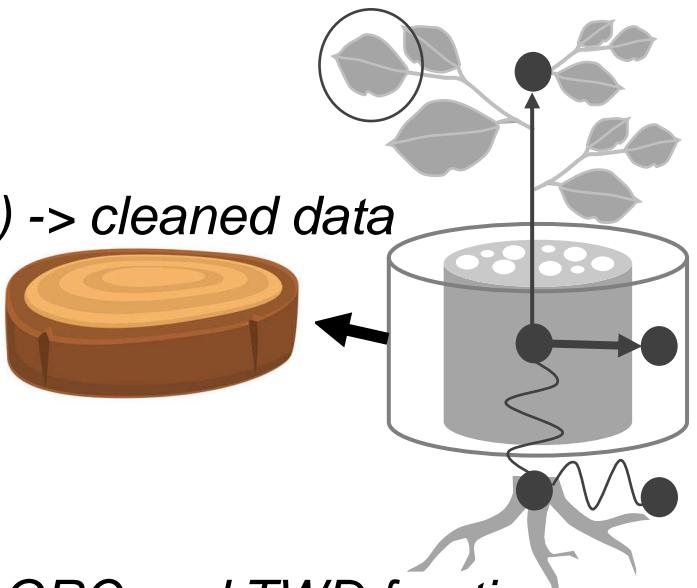
## Functionalities – Extracting growth and TWD *proc\_dendro\_L2()*



*proc\_dendro\_L1() -> Time aligned data*

*proc\_dendro\_L2() -> cleaned data*

*-> separated into GRO and TWD fractions*



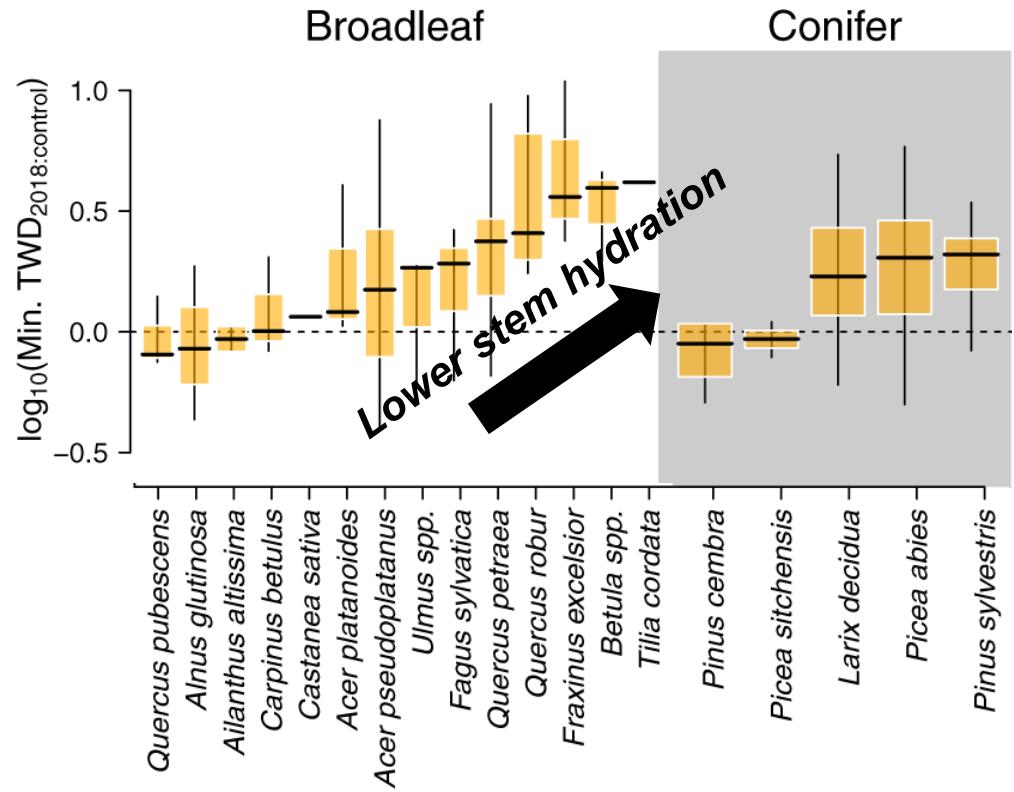


## Output – Night-time tree water deficit as indicator for drought stress

*proc\_dendro\_L2()*

*gro\_seas*

*gro\_stats*

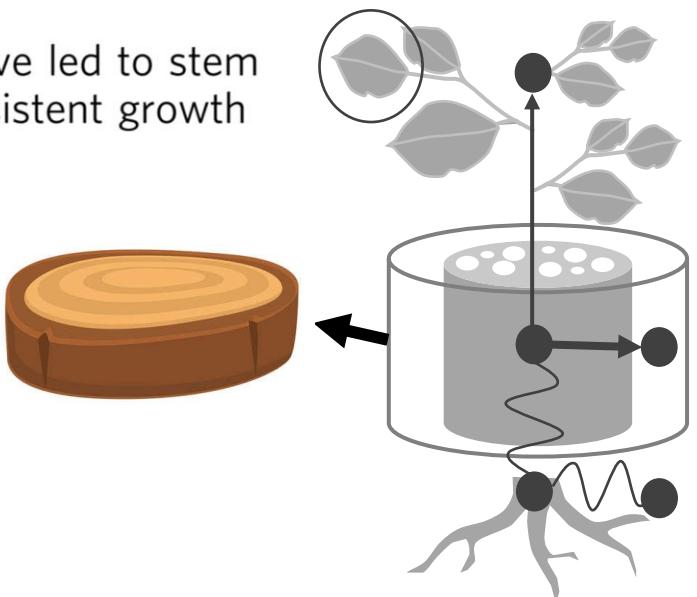


ARTICLE

<https://doi.org/10.1038/s41467-021-27579-9>

OPEN

The 2018 European heatwave led to stem dehydration but not to consistent growth reductions in forests



### Additional Resources

Knüsel S., Peters R.L., Haeni M., Wilhelm M., Zweifel R. (2021) **Processing and extraction of seasonal tree physiological parameters from stem radius time series.** Forests 12(6) <https://doi.org/10.3390/f12060765>

<https://github.com/treenet/treenetproc>

<https://deep-tools.netlify.app/#workshops>





## Hands-on & Q/A