protocols.io

NOV 09, 2022

⮜ SHARE

WORKS FOR ME 1

🌐 Formatting Data for Use With ebFRET Hidden Markov Modeling Software

**This protocol is published without a DOI.**

Clark Fritsch[1]

[1]University of Pennsylvania

Clark Fritsch
University of Pennsylvania

COMMENTS 0

## ABSTRACT

This protocol follows from the "Selecting a Region of Interest for Hidden Markov Modeling using ebFRET" protocol and is the second step towards analyzing your data using Hidden Markov Modeling. This protocol describes the steps you need to take to convert traces containing your region of interest (selected in the previous protocol) and formatting them for use with ebFRET, a software package that is capable of performing Hidden Markov Modeling on single-molecule FRET datasets.

## PROTOCOL CITATION

## LICENSE

1   This protocol is a continuation of the "Selecting a Region of Interest for Hidden Markov Modeling using ebFRET" practice protocol.

2   To begin this protocol, you should have your "SavedData.csv" file and traces in one directory, as shown below:



It is likely that when you are doing your real analysis that you don't include every trace that is included in your "GoodOnes.txt" file in your "SavedData.csv" file. For example, you might think that a trace is good at first and include it in your "GoodOnes.txt" file, but find upon closer inspection that the trace is not worthy of further analysis and therefore exclude it from your "SavedData.csv" file.

Because of this, you want to separate the traces that are included in your "SavedData.csv" file from the remainder of your traces.

3   To separate the traces included in your "SavedData.csv" file from the remainder of your traces, you can use the "SavedData_Copy_and_Transfer_v3.R" program that is attached below:

4    To use this program in R, simply copy the path to the directory that contains all of your traces and the "SavedData.csv" file and paste the path into the "path_var_old" variable in between the parentheses, as highlighted below:



Then press "Ctrl A" to highlight all of the code and press "Ctrl Enter" to run the code.

5    After running the code, you will see that the following information has been outputted to the console (boxed in yellow):



This output simply shows that each trace (TRUE) has been transferred successfully to the output directory. If a trace is not transferred successfully, it was output FALSE. Additionally, your "SavedData.csv" file is transferred to the output directory.

6    After running the program successfully, you will see that the "SavedData" output directory has been created at the path directory that you inputted into the program:

This "SavedData" directory will contain only your "SavedData.csv" file and the traces that were contained within your "SavedData.csv" file, as shown below:



7  Next, you want to cut out the parts of each of your traces that are not included within the range set by Click 1 and Click 2 that is included in your "SavedData.csv" file.

**KEEP THIS**



You can do this using the "Click1_Click2_Cut.R" program that is included below:

☐ Click1_Click2_Cut.R

8  The "Click1_Click2_Cut.R" program, as described above, simply takes the values for Click1 (leftmost boundary of your ROI) and Click2 (rightmost boundary of your ROI) and removes the frames from each trace that are not included between Click1 and Click2.

To use this program in R, simply copy the path to "SavedData" directory and paste the path into the "path_var" variable in between the parentheses, as highlighted below:

```r
# This program is meant to take the Click1 and Click2 values determined by the TTimes.py program and cut out everything that is not within these values.

# Note that this program assumes that you have the "SavedData_formatted_NoNeg.csv" file in a directory with only the files contained within the file. If you include
# files in the directory that are not in the "SavedData_formatted_NoNeg.csv" file, then the program will give you an error when it encounters a file that is not in the
# SavedData.csv file.

library(tidyverse)

# create the output folder
path_var <- r"(C:\Users\clark\Desktop\Example_Analysis\SavedData)"

dir.create(paste(toString(path_var), "T1_to_T2_Synchronized", sep = "///"), showWarnings = FALSE)
output_dir <- paste(path_var, "T1_to_T2_Synchronized", sep = "\\")


# get the names of the input files with their full path
files <- list.files(path_var, "pairProfile", full.name = TRUE)
savedData <- read.csv(list.files(path_var, "SavedData.csv", full.name = TRUE), header = FALSE, sep = ",")


# loop through all the input files
for (file in files) {

    # read the file, specify the correct separator
    data <- read.csv(file, header = FALSE, sep = ",")

    #use the select function in Dplyr to select only the V1 and V2 columns
    varSelected_data <- select(data, V1, V2)

    #use the slice function in Dplyr to cut coordinate row out of the dataframe
    cut_data <- slice(varSelected_data, 2:nrow(varSelected_data))

    #use the filter and grepl function in Dplyr to select the row that has the same filename as the current file
    formatted_Name <- strsplit(file, "/")[[1]][2]
    dummyTest2 <- dplyr::filter(savedData, grepl(formatted_Name, V13))

    #use the pull function in Dplyr to obtain the T1 and T2 values that correspond to the current file
    T1_dummyTest2 <- pull(dummyTest2, V3)
    T2_dummyTest2 <- pull(dummyTest2, V4)

    #Cut the current file according to the (T1 - certain number of frames) and (T2 + certain number of frames) values. This can be changed depending on how you want to display

    synch_data <- slice(cut_data, (T1_dummyTest2):(T2_dummyTest2))

    #Define the new file name

    outfile <- paste(strsplit(formatted_Name, ".csv")[[1]][1], "T1_to_T2_Synchronized", sep = "-")
    outfile <- paste(outfile, "csv", sep = ".")

    # write the file.
    write.table(synch_data, paste(output_dir, outfile, sep = "\\"), sep = ",", col.names = FALSE, row.names = FALSE, quote = FALSE)
}
```
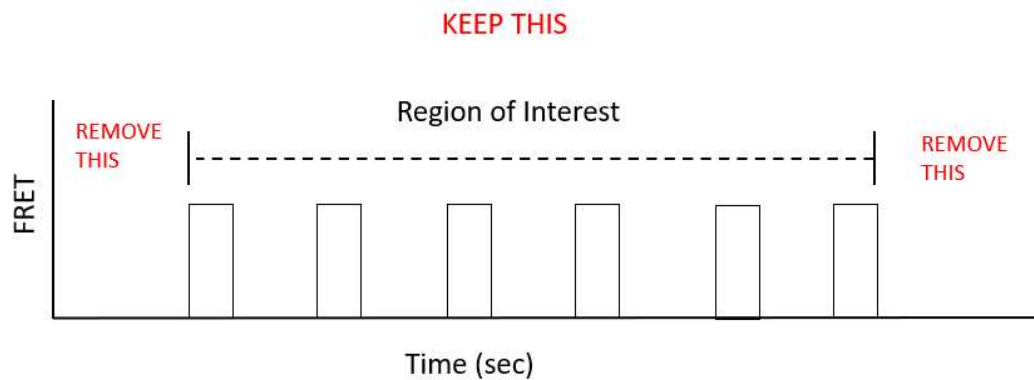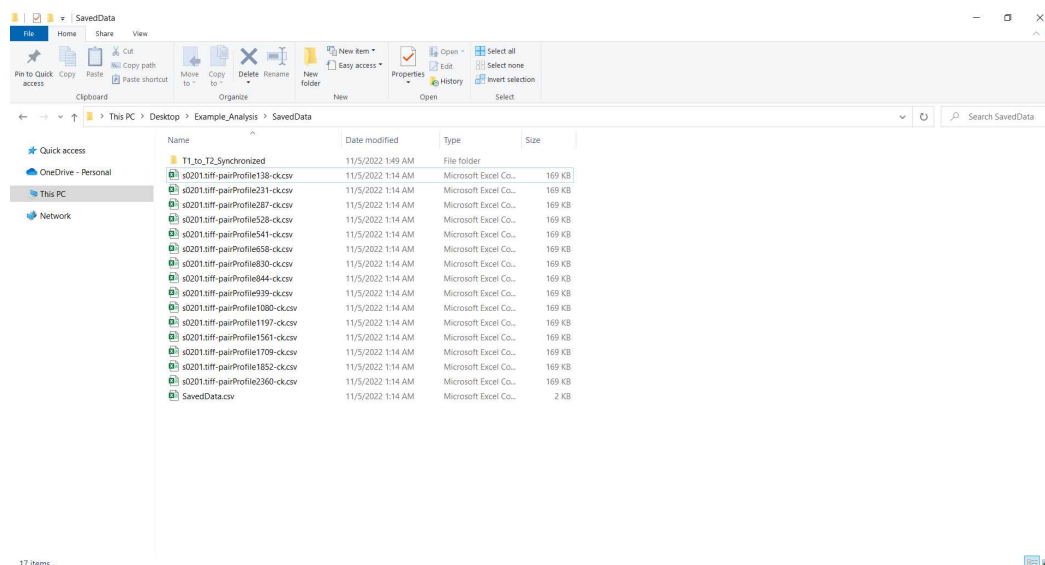
Then press "Ctrl A" to highlight all of the code and press "Ctrl Enter" to run the code.

9  Once you have run the "Click1_Click2_Cut.R" program, a new directory called "T1_to_T2_Synchronized" that will contain your new traces, each of which has been cut to contain

only the frames between your Click1 and Click2 values:



Each of these new "trimmed" traces can be distinguished from the original traces because the trimmed traces contain "-T1_to_T2_Synchronized" appended to the end of each filename:



10    Now that you have your trimmed traces prepared, you need to format the traces so that they can be run using the ebFRET Hidden Markov Modeling program.

To do this, you can use "RegionOfInterest_to_ebFRET.R" program that is attached below:

📄 RegionOfInterest_to_ebFRET.R

**11**    To use ebFRET, your traces of interest need to be stitched together with each trace being associated with a unique ID number. Additionally, the donor and acceptor fluorophore intensity for each timepoint in each trace needs to be included in the file.

To use the "RegionOfInterest_to_ebFRET.R" program to format your traces in this way, simply copy the path that contains your trimmed traces into the "path" variable in the program, as highlighted below:

```r
################################################################################
# This program is meant to take the .csv files created by the Click1_Click2_Cut.R program and format them for use in ebFRET.
################################################################################

#create a variable that contains the full names (full paths) for each text file in your directory of interest.

path = r"(C:\Users\clark\Desktop\Example_Analysis\SavedData\T1_to_T2_Synchronized)"
setwd(path)
files <- list.files(path, pattern="*.csv", full.names=TRUE, recursive=FALSE)

b <- 1
new_file_name <- "ebFRET_formatted_traces_NearCognate_25nMTC_5uMeEF2_CoInjection_SecondaryEvents_T1_to_T2_CKH_Filtered_Traces.txt"

for (i in files){

    # read the file, specify the correct separator
    data <- read.csv(i, header = FALSE, sep = ",")

    #Add a new column to each file that specifies the ID for the file
    ID_row <- b
    formatted_table <- cbind(ID_row, data)

    colnames(formatted_table) <- NULL

    #Save the formatted files to one file where each formatted trace appended to the last trace in the file.
    write.table(formatted_table, new_file_name, row.names = FALSE, col.names = FALSE, append = TRUE)

    #Iterate b so that the next trace has a unique identifier.
    b = b + 1

}
```

Additionally, you need to name the formatted file that your trace information will be outputted to by entering the new filename into the "new_file_name" variable, as shown below:

```r
################################################################################
# This program is meant to take the .csv files created by the Click1_Click2_Cut.R program and format them for use in ebFRET.
################################################################################

#create a variable that contains the full names (full paths) for each text file in your directory of interest.

path = r"(C:\Users\clark\Desktop\Example_Analysis\SavedData\T1_to_T2_Synchronized)"
setwd(path)
files <- list.files(path, pattern="*.csv", full.names=TRUE, recursive=FALSE)

b <- 1
new_file_name <- "practice_protocol_ebFRET.dat"

for (i in files){

    # read the file, specify the correct separator
    data <- read.csv(i, header = FALSE, sep = ",")

    #Add a new column to each file that specifies the ID for the file
    ID_row <- b
    formatted_table <- cbind(ID_row, data)

    colnames(formatted_table) <- NULL

    #Save the formatted files to one file where each formatted trace appended to the last trace in the file.
    write.table(formatted_table, new_file_name, row.names = FALSE, col.names = FALSE, append = TRUE)

    #Iterate b so that the next trace has a unique identifier.
    b = b + 1

}
```

Note that the new file name must have the ".dat" extension, as this is the only text file format that ebFRET will recognize.
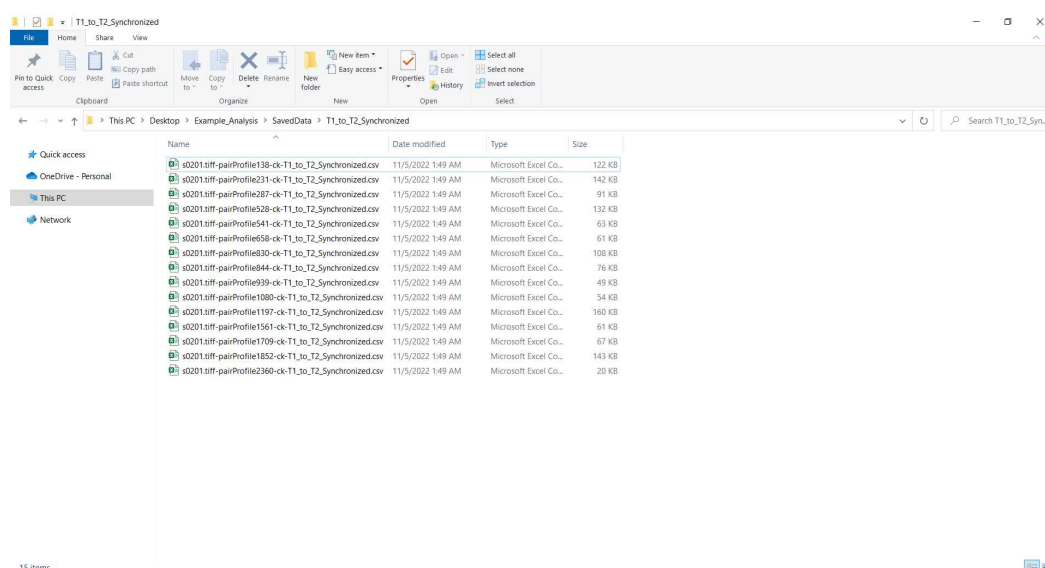
Press "Ctrl A" to highlight all of the code and press "Ctrl Enter" to run the code.

**12**    Once you have run the program, your ebFRET formatted file will appear in your "T1_to_T2_Synchronized" directory. It is now ready to input into ebFRET for Hidden Markov Analysis.