

EcoTransLearn R-package

Version 1.0-3

USER MANUAL

G. WACQUET

*IFREMER
LABORATOIRE ENVIRONNEMENT LITTORAL & RESSOURCES AQUICOLES
UNITE LITTORAL
CENTRE MANCHE MER DU NORD
BOULOGNE-SUR-MER, FRANCE*

TABLE OF CONTENTS

INTRODUCTION	4
INSTALLATION AND EXECUTION	4
R-package installation	4
Python installation	4
Launching the Graphical User Interface	5
TRAINING and TEST SETS CREATION	5
USE OF THE GRAPHICAL USER INTERFACE	6
DATA SELECTION button	6
SETTINGS button	6
CLASSIFY button	12
VIEW button	13
MORE... (+) button.....	14

INTRODUCTION

In recent years, Deep Learning (DL) has been increasingly used in many fields, in particular in image recognition, due to its ability to solve problems where traditional machine learning algorithms fail. However, building an appropriate DL model from scratch especially in the context of ecological studies, such as monitoring marine ecosystems, is a difficult task due to the dynamic nature and morphological variability of living organisms, as well as the high cost in terms of time, human resources, and skills required to label a large number of training images. To overcome this problem, Transfer Learning (TL) can be used to improve a classifier by transferring information learnt from many domains thanks to a very large training set composed of various images, to another domain with a smaller amount of training data. To compensate the lack of “easy-to-use” software optimized for ecological studies, we propose the *EcoTransLearn* R-package, which allows greater automation in classification of images acquired with various devices, thanks to different TL methods pre-trained on the generic ImageNet dataset.

INSTALLATION AND EXECUTION

R-package installation

The version 1.0-0 of the *EcoTransLearn* package needs a recent version of R (version 4.0.x or upper). It can be directly downloaded on the CRAN website (<http://cran.r-project.org>).

By double-clicking on the R icon on the desktop, or by selecting R in the start menu, a window appears on the screen: this is the R console. This allows to control R directly by command lines. It also allows to display the main results and messages of the different actions performed with *EcoTransLearn*.

The R-packages needed by *EcoTransLearn* (*colorRamps*, *ggplot2*, *grid*, *jpeg*, *mapplots*, *maps*, *randomForest*, *reticulate*, *SDMTools*, *shapefiles*, *stringr*, *svDialogs*, *svMisc*, *tcltk2*, *tiff*, *zooimage*) can be installed directly from the R console, by typing:

```
install.packages(c("colorRamps","ggplot2","grid","jpeg","mapplots","maps","randomForest",  
"reticulate","SDMTools","shapefiles","stringr","svDialogs","svMisc","tcltk2","tiff","zooimage"))
```

Then choose a mirror (default: 0-cloud) to start downloads and installations.

Python installation

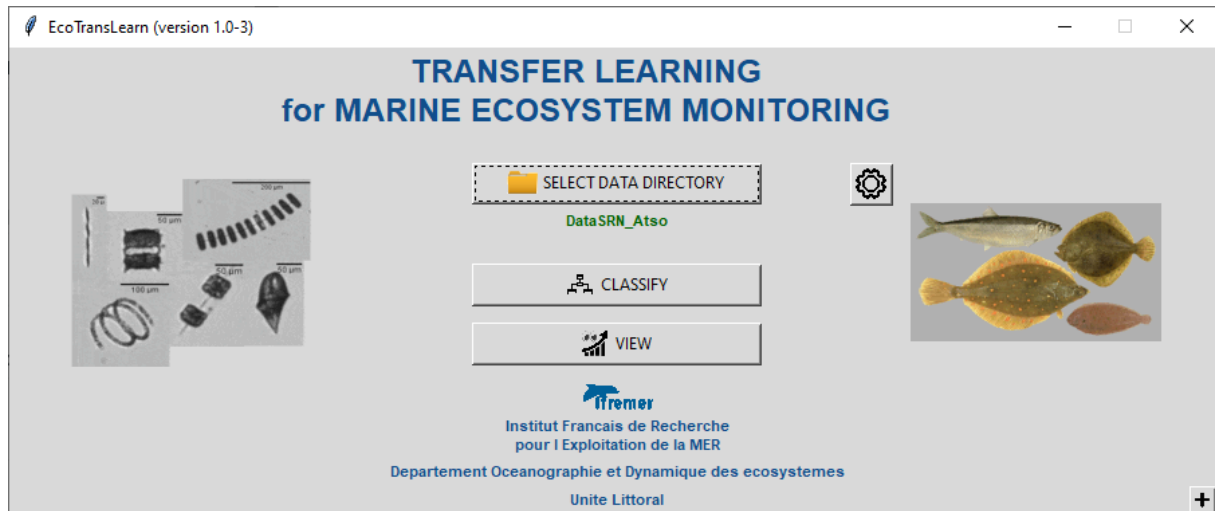
Anaconda is a scientific distribution of Python, which allows to use several applications (such as Spyder, Jupyter Notebook, ...) and to manage different libraries. It can be downloaded on the website: <https://docs.anaconda.com/anaconda/navigator/install/>. The *EcoTransLearn* package needs the version 3.7 (or upper) of Python.

Once the distribution installed, the Python libraries needed by *EcoTransLearn* (*matplotlib*, *numpy*, *pandas*, *sklearn*, *tensorflow*) can be installed directly by opening an Anaconda Prompt, and by entering the commands:

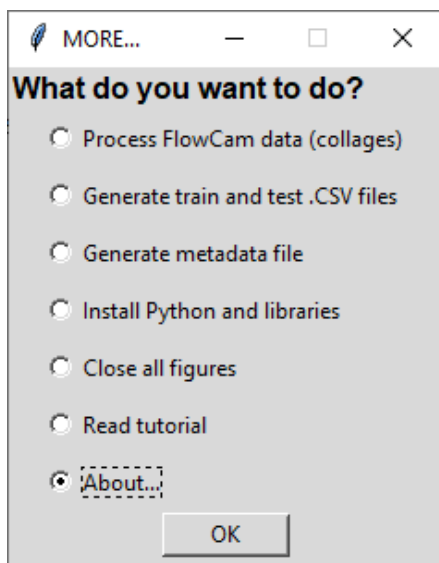
```
conda install Image  
conda install pillow  
conda install matplotlib  
conda install numpy  
conda install pandas  
conda install sklearn  
conda install tensorflow
```

Launching the Graphical User Interface

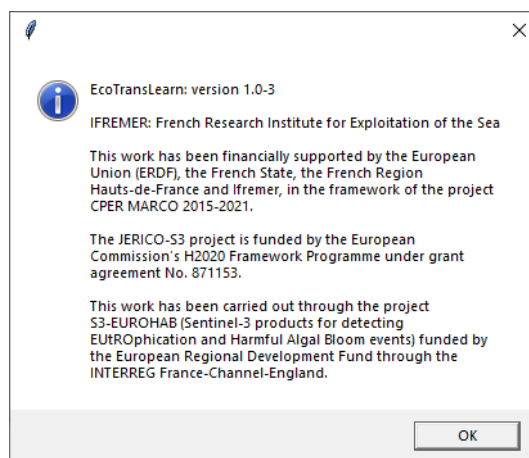
Once the installation of the packages is finished, it is possible to make sure that the previous steps run smoothly by checking that the installed version is 1.0-3. To do this, first type in the R console: **require(EcoTransLearn)**, to load the package, then: **EcoTransLearn()**, to launch the Graphical User Interface (GUI):



Click on the + button (bottom right). A new window appears:

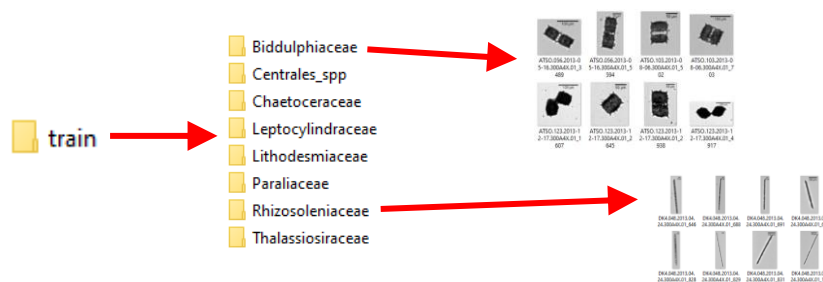


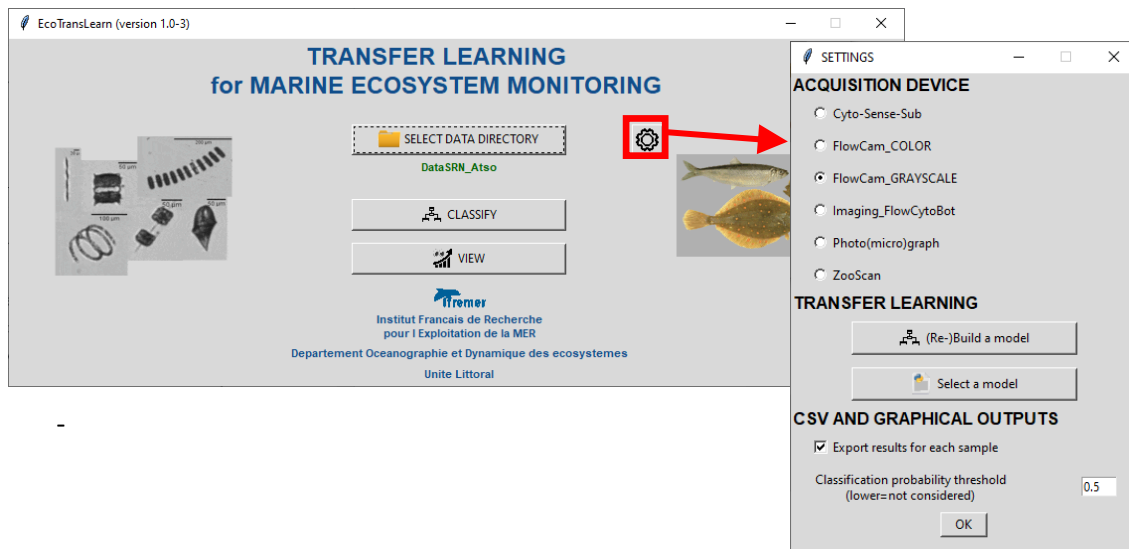
By selecting **About...**, a dialog box appears and informs the user of the *EcoTransLearn* version.



TRAINING and TEST SETS CREATION

It is important to note that training sets and test sets must contain individual images (and not collages), i.e. one image file per particle. These image files must have a unique filename. The images are then manually sorted into subdirectories (as many as necessary) bearing the names of the different groups.





❖ ACQUISITION parameters

Choose the kind of instrument used for image acquisition.

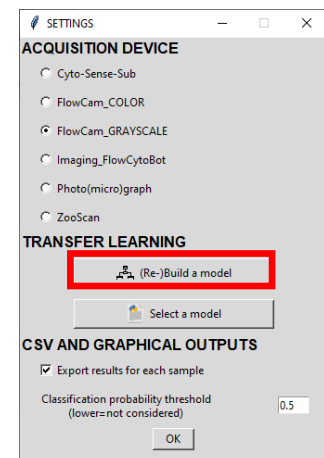
***Note:** for the FlowCam_COLOR and the FlowCam_GRAYSCALE, it is possible to directly use the raw data from the instruments (collages). The package allows to directly creating the vignettes (1 image per particle) from the collages and the 1st file. However, in the training sets and test sets, only the vignettes (1 image per particle) are considered. To create these sets, refer to “Training and test sets creation” section.*

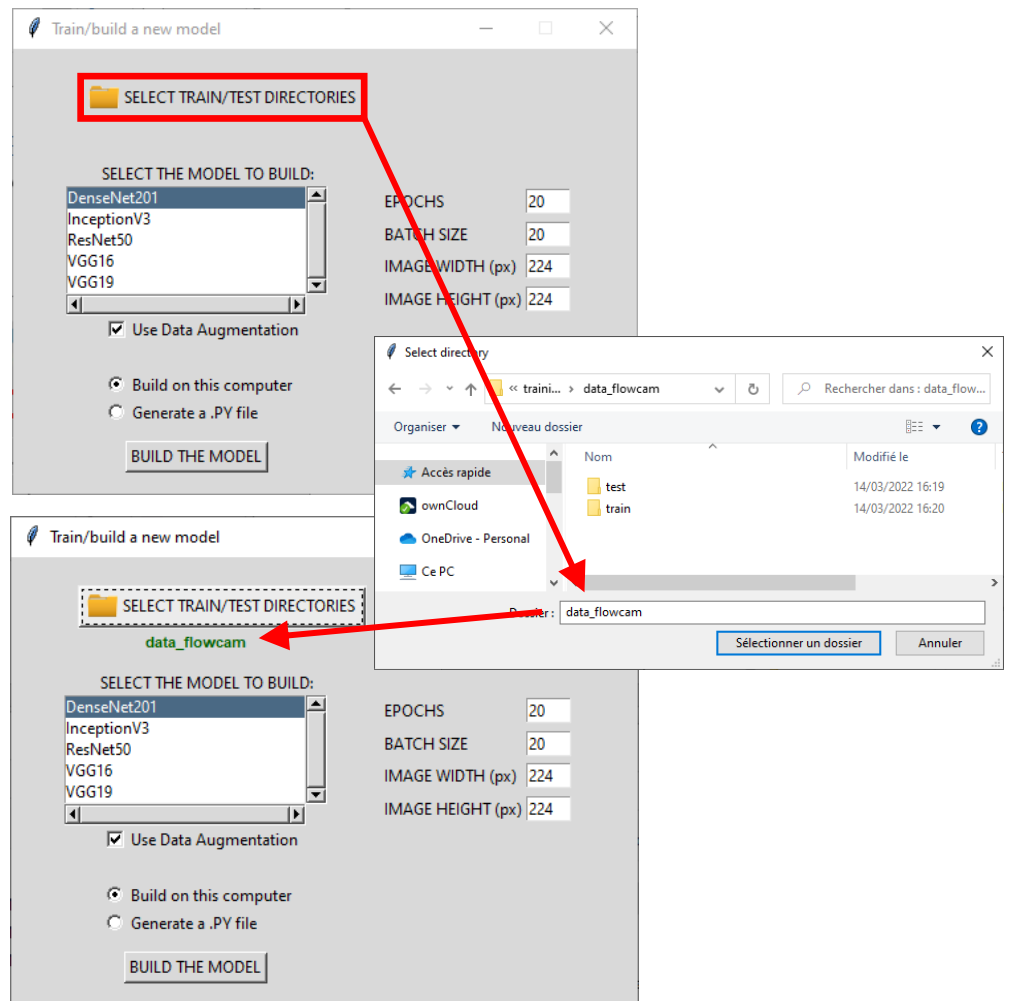
❖ CLASSIFICATION parameters

▪ (Re-)Build a model

This button allows to build (or rebuild) a classification model by Transfer Learning using a data set including a training set (directory named **train**) for learning, and a test set (directory named **test**) for validation and evaluation.

By clicking on this button, a new window appears.





Click on the **SELECT TRAIN/TEST DIRECTORIES** button, select the directory containing the two sub-folders train and test, then confirm by clicking **OK**: the name of the selected directory is then displayed below the selection button

It is then possible to select different Convolutional Neural Network (CNN) architectures in the **SELECT THE MODEL TO BUILD** list, among:

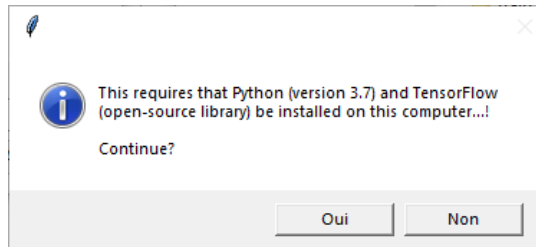
- DenseNet201
- InceptionV3
- ResNet50
- VGG16
- VGG19

and adjust parameters related to images and training step, by setting values for **EPOCH** (default=20), **BATCH SIZE** (default=20), **IMAGE WIDTH** (default=224) and **IMAGE WEIGHT** (default=224), but also choose the possibility of using the technique of data augmentation (**Use Data Augmentation**). In the case of a training set with few images, this option is used to automatically generate additional images from the original images in the training set, by applying geometric transformations such as rotations (by default, rotation_range=45) and horizontal and vertical flips (by default, horizontal_flip=True and vertical_flip=True).

The last step is to choose the material on which to build and adapt the classification model. Depending on the selected option, and after clicking on the **BUILD THE MODEL** button:

- **Build on this computer**

A dialog box appears:



***Warning:** the training duration can be long (several hours) depending on the number of images in the training set and according to the parameters defined in the previous step.*

- **Generate a .PY file**

A script is automatically created, and can be run on other hardware (dedicated computer, calculation server, etc.).

```
VGG16_script.py - Bloc-notes
Fichier Edition Format Affichage Aide
modelName = 'VGG16'

##### LIBRARIES IMPORTATION #####
import tensorflow as tf
from tensorflow.keras.layers import Flatten, Dense, Dropout, Input, BatchNormalization
from tensorflow.keras.preprocessing.image import ImageDataGenerator
from tensorflow.keras import applications
from tensorflow.keras.callbacks import EarlyStopping, ModelCheckpoint, ReduceLROnPlateau, CSVLogger
from tensorflow.keras import Model
from tensorflow.keras.optimizers import Adam
from tensorflow.keras import backend as K
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import label_binarize
import pandas as pd

if modelName == "DenseNet201":
    from tensorflow.keras.applications.densenet import DenseNet201
    from tensorflow.keras.applications.densenet import preprocess_input
if modelName == "InceptionV3":
    from tensorflow.keras.applications.inception_v3 import InceptionV3
    from tensorflow.keras.applications.inception_v3 import preprocess_input
if modelName == "ResNet50":
    from tensorflow.keras.applications.resnet50 import ResNet50
    from tensorflow.keras.applications.resnet50 import preprocess_input
if modelName == "VGG16":
    from tensorflow.keras.applications.vgg16 import VGG16
    from tensorflow.keras.applications.vgg16 import preprocess_input
if modelName == "VGG19":
    from tensorflow.keras.applications.vgg19 import VGG19
    from tensorflow.keras.applications.vgg19 import preprocess_input

##### GLOBAL VARIABLES #####
weightPath = './train'
trainPath = 'C:/Users/Administrateur/Desktop/ImgTLCClass/training_data/data_flowcam/train'
testPath = 'C:/Users/Administrateur/Desktop/ImgTLCClass/training_data/data_flowcam/test'
save_dir = 'C:/Users/Administrateur/Desktop/ImgTLCClass/training_data/data_flowcam/saved_models'
BATCH_SIZE = '20'
EPOCH = '20'
img_width = '224'
img_height = '224'
data_aug = 'TRUE'

Ln 1, Col 1    100%    Windows (CRLF)    UTF-8
```

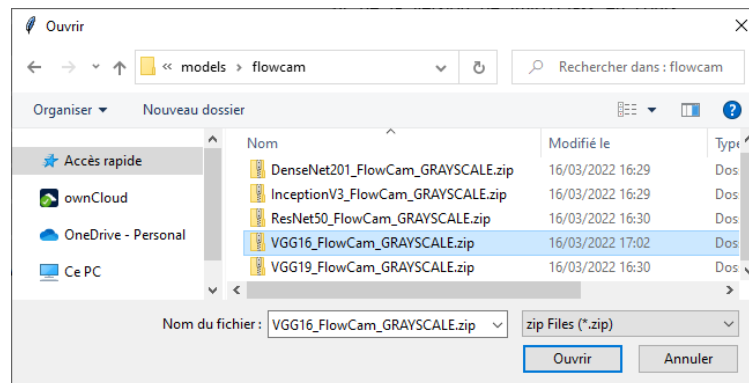
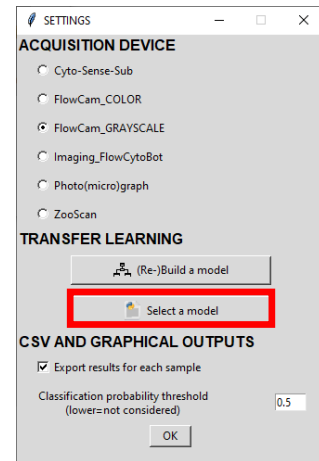
At the end of the execution of the script, a .ZIP file is created. It contains all the information necessary for the classification of a new set of images.

- VGG16_classnames.csv
- VGG16_history.csv
- VGG16_model.h5

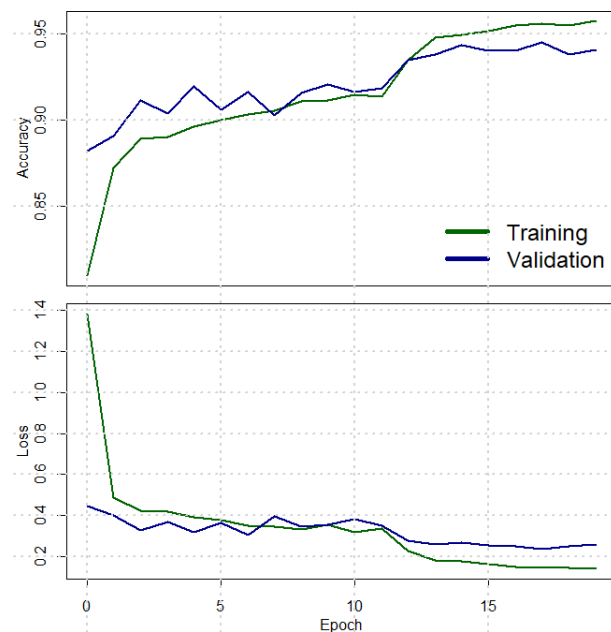
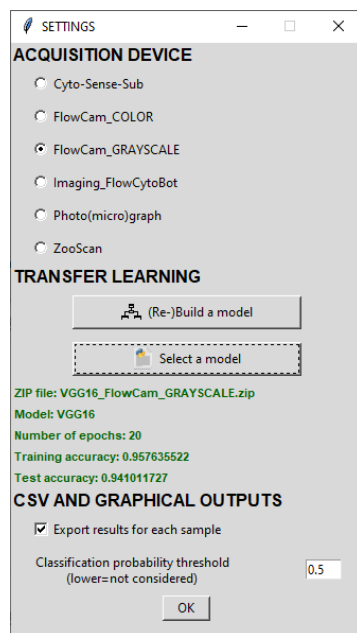
modelName_classNames.csv contains the class names in the training set.
modelName_history.csv contains the evaluation scores of the built model.
modelName_model.h5 is the model built with Python.

- **Select a model**

This button allows to choose a model for the automated classification of the images contained in the directory selected during the first step (see section “Selection of input data”).



Select the ZIP file generated during the previous step, then validate by clicking on **Open**: information on the model performance (calculated during the training and validation steps) is then displayed below the selection button, as well as **Accuracy** (\approx percentage of correctly classified data) and **Loss** (\approx distance between actual data and predicted data) curves.



❖ OUTPUTS parameters

▪ Export results for each sample

This option allows to export the results for each sub-folder of the directory selected during the first step (see section "Selection of input data").

For each processed sample, three CSV files are created:

- sampleName.csv

Sample	Group	Relative	Count	Date	Volume	Prob_threshold	Percent_used
flowcam_FCM.U0003.2016-01-21.300A4X.01	Asterionellopsis	0	0		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	C_curvisetus	8.474576271	5		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	C_danicus	15.25423729	9		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	C_socialis	1.694915254	1		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	Ciliophora	8.474576271	5		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	Dactyliosolen	3.389830508	2		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	Dytilum	6.779661017	4		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	G_flaccida	1.694915254	1		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	G_striata	1.694915254	1		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	Gymnodinium	3.389830508	2		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	Lauderia	1.694915254	1		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	Leptocylindrus	1.694915254	1		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	Odontella	3.389830508	2		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	P_globosa	0	0		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	Pleuro_Gyrosigma	3.389830508	2		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	Prorocentrum	11.86440678	7		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	PseudoNitzschia	3.389830508	2		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	Rhizosolenia	16.94915254	10		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	Thalassionema	3.389830508	2		1000	0.75	92.15500945
flowcam_FCM.U0003.2016-01-21.300A4X.01	Thalassiosira	3.389830508	2		1000	0.75	92.15500945

- sampleName_CLASSIF.csv

Filename	Class
unknown\flowcam_FCM.U0003.2016-01-21.300A4X.01_1.jpg	dark
unknown\flowcam_FCM.U0003.2016-01-21.300A4X.01_10.jpg	dark
unknown\flowcam_FCM.U0003.2016-01-21.300A4X.01_100.jpg	dark
unknown\flowcam_FCM.U0003.2016-01-21.300A4X.01_1000.jpg	dark
unknown\flowcam_FCM.U0003.2016-01-21.300A4X.01_1001.jpg	dark
unknown\flowcam_FCM.U0003.2016-01-21.300A4X.01_1002.jpg	dark
unknown\flowcam_FCM.U0003.2016-01-21.300A4X.01_1003.jpg	dark
unknown\flowcam_FCM.U0003.2016-01-21.300A4X.01_1004.jpg	dark
unknown\flowcam_FCM.U0003.2016-01-21.300A4X.01_1005.jpg	dark
unknown\flowcam_FCM.U0003.2016-01-21.300A4X.01_1006.jpg	Ciliophora
unknown\flowcam_FCM.U0003.2016-01-21.300A4X.01_1007.jpg	Odontella
unknown\flowcam_FCM.U0003.2016-01-21.300A4X.01_1008.jpg	dark
unknown\flowcam_FCM.U0003.2016-01-21.300A4X.01_1009.jpg	fiber

- sampleName_PRED.csv

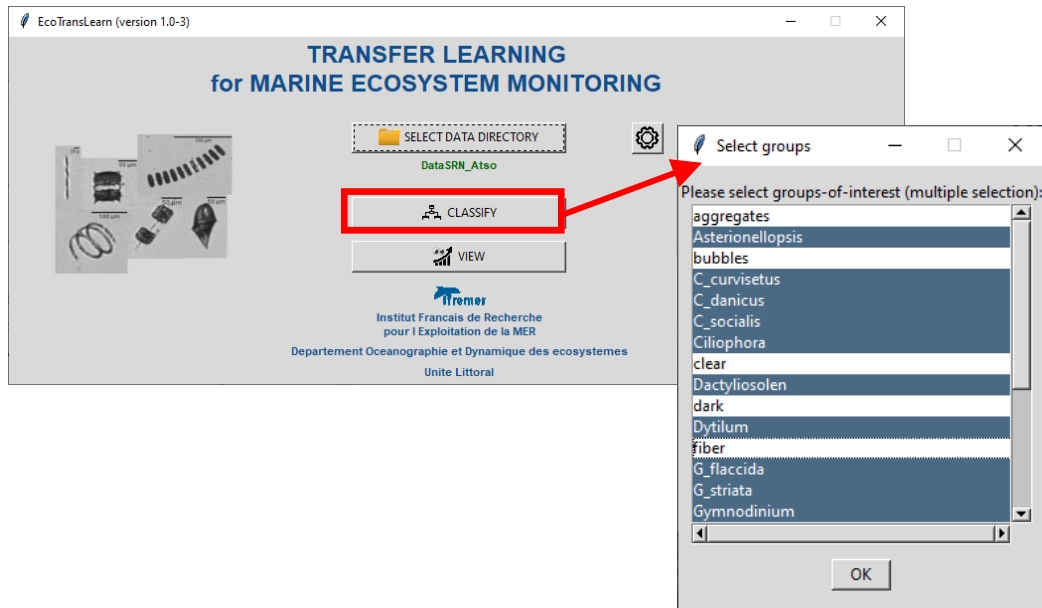
Asterionellopsis	C_curvisetus	C_danicus	C_socialis	Ciliophora	Dactyliosolen	Dytilum	G_flaccida
3.52E-09	1.24E-08	4.22E-10	3.65E-11	0.000285448	1.65E-06	9.95E-09	7.16E-10
1.34E-11	2.78E-11	7.02E-14	1.16E-16	4.84E-10	5.95E-14	3.10E-13	3.12E-16
1.86E-12	2.97E-09	1.78E-13	8.65E-14	8.47E-05	1.23E-07	4.80E-09	1.05E-11
9.69E-17	1.54E-11	3.29E-18	7.69E-20	3.97E-08	1.30E-10	1.50E-12	2.34E-17
6.27E-08	5.92E-07	3.55E-09	2.57E-09	7.62E-06	1.90E-06	1.73E-07	1.10E-07
1.42E-08	1.77E-08	5.07E-08	7.74E-11	1.00E-05	1.52E-07	1.58E-08	1.10E-11
1.01E-24	1.02E-15	9.25E-27	4.89E-27	3.53E-14	3.10E-17	2.96E-18	2.73E-22
3.96E-15	1.95E-12	1.79E-16	2.53E-18	6.86E-06	1.27E-08	2.07E-13	3.37E-14
6.43E-15	2.69E-09	8.91E-14	5.48E-17	0.001769048	2.66E-08	1.20E-12	4.52E-14

▪ Classification probability threshold

The threshold value defined at this step allows to take into account only the images having a probability of "good" classification greater than this threshold. To take into account all the images, this value must be set to 0.

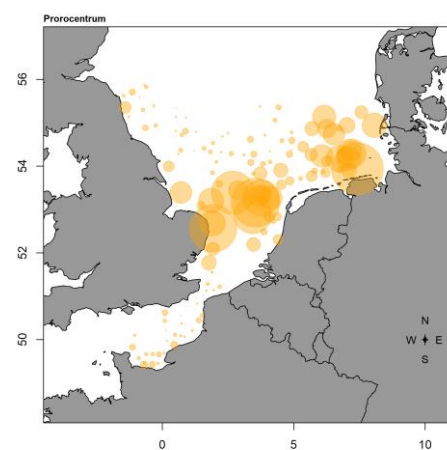
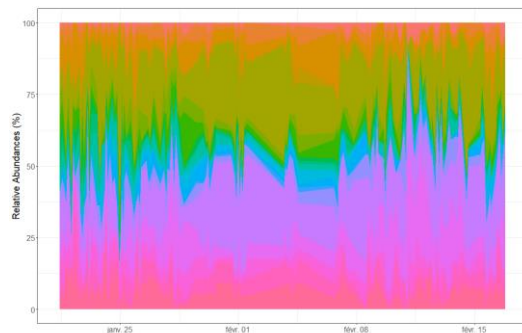
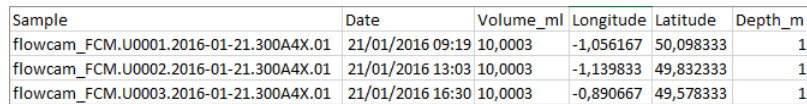
CLASSIFY button

To classify new images, click on the **CLASSIFY** button. A new window appears for group selection.



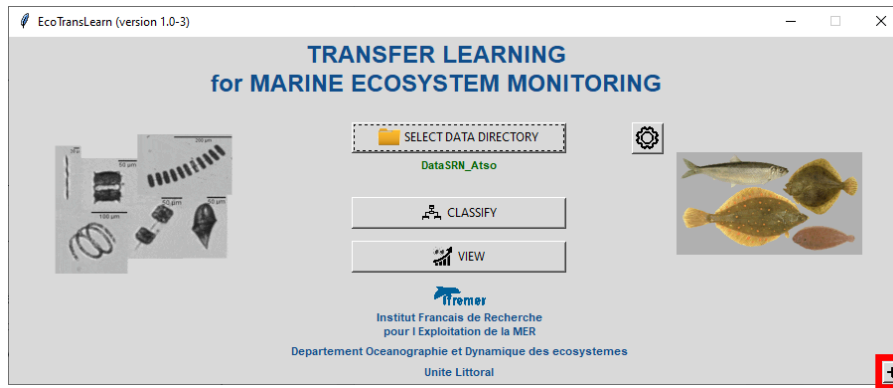
Sample	Group	Relative	Count	Prob_threshold	Percent_used
flowcam_FCM.U0001.2016-01-21.300A4X.01	Asterionellopsis	2.941176471	1	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	C_curvisetus	5.882352941	2	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	C_danicus	8.823529412	3	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	C_socialis	0	0	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	Ciliophora	5.882352941	2	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	Dactyliosolen	8.823529412	3	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	Dytilum	8.823529412	3	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	G_flaccida	8.823529412	3	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	G_striata	5.882352941	2	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	Gymnodinium	2.941176471	1	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	Lauderia	0	0	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	Leptocylindrus	0	0	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	Odontella	0	0	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	P_globosa	0	0	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	Pleuro_Gyrosigma	5.882352941	2	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	Prorocentrum	14.70588235	5	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	Pseudonitzschia	2.941176471	1	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	Rhizosolenia	11.76470588	4	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	Thalassionema	2.941176471	1	0.75	89.88476312
flowcam_FCM.U0001.2016-01-21.300A4X.01	Thalassiosira	2.941176471	1	0.75	89.88476312
flowcam_FCM.U0002.2016-01-21.300A4X.01	Asterionellopsis	1.388888889	1	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	C_curvisetus	1.388888889	1	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	C_danicus	12.5	9	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	C_socialis	0	0	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	Ciliophora	6.944444444	5	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	Dactyliosolen	5.555555556	4	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	Dytilum	6.944444444	5	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	G_flaccida	4.166666667	3	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	G_striata	5.555555556	4	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	Gymnodinium	5.555555556	4	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	Lauderia	0	0	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	Leptocylindrus	0	0	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	Odontella	4.166666667	3	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	P_globosa	0	0	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	Pleuro_Gyrosigma	4.166666667	3	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	Prorocentrum	15.27777778	11	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	Pseudonitzschia	8.333333333	6	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	Rhizosolenia	6.944444444	5	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	Thalassionema	5.555555556	4	0.75	90.75812274
flowcam_FCM.U0002.2016-01-21.300A4X.01	Thalassiosira	5.555555556	4	0.75	90.75812274

By clicking on this button, a new window appears for group selection.

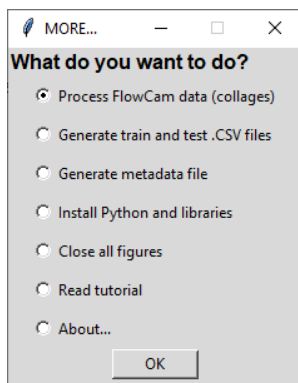


MORE... (+) button

To help the user in the different data formatting steps for image analysis, several options are available. To view the list of these additional tools, click on the + button (**MORE...**, located at the bottom right of the main window).

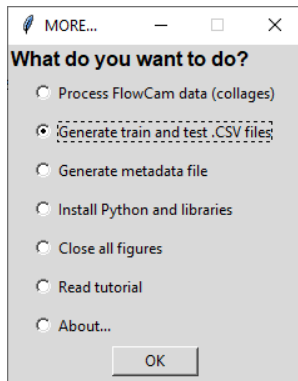


❖ Process FlowCam data (collages)

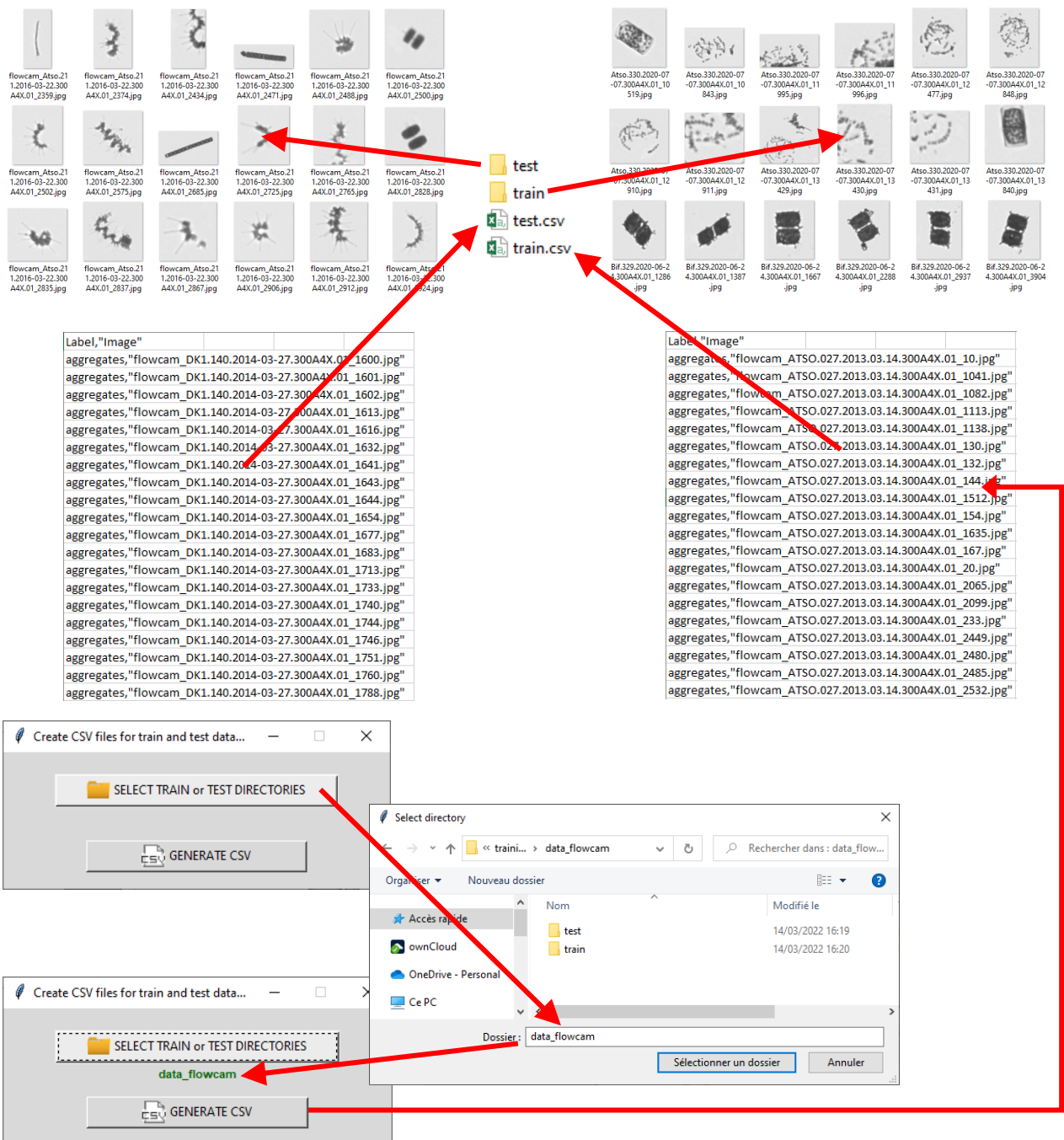


Raw data from the FlowCam is often presented as collages (one file with multiple images). It is then possible to cut and save vignettes (one image file per particle) from these collages. To do this, choose **Process FlowCam data (collages)**.

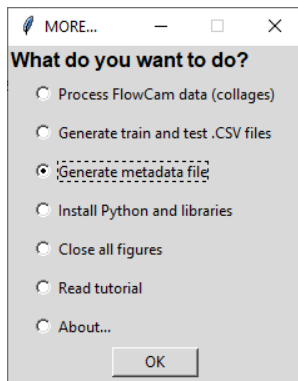
❖ Generate train and test .CSV files



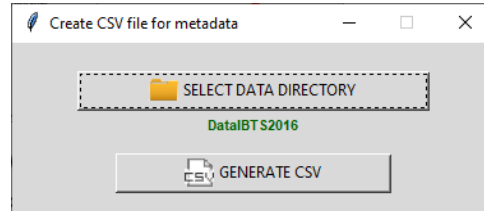
'train' and 'test' directories must contain individual images (=vignettes) sorted into different subdirectories. In order to generate summary files that can be easily used in R, it is possible to use the **Generate train and test .CSV files** option. These two .CSV files are then created and saved in the root directory, and all vignettes are then merged into a unique directory (for 'train' and 'test').



❖ Generate metadata file

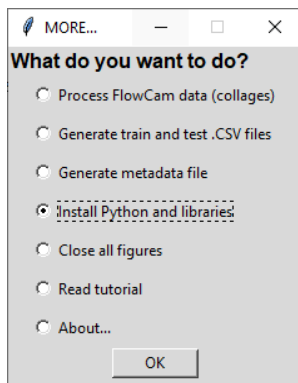


For the visualisation of classification results, a metadata file is required. This CSV file contains information on each sample, such as Date, Volume, GPS coordinates, and other metadata. A template can be generated thanks to this option.

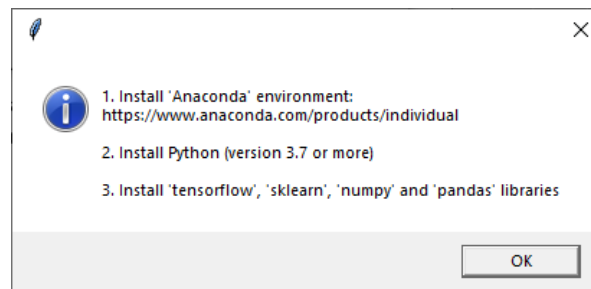


Sample	Date	Volume_ml	Longitude	Latitude	Depth_m
flowcam_FCM.U0001.2016-01-21.300A4X.01	21/01/2016 09:19	10,0003	-1,056167	50,098333	1
flowcam_FCM.U0002.2016-01-21.300A4X.01	21/01/2016 13:03	10,0003	-1,139833	49,832333	1
flowcam_FCM.U0003.2016-01-21.300A4X.01	21/01/2016 16:30	10,0003	-0,890667	49,578333	1

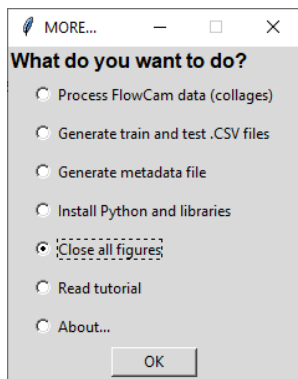
❖ Python and libraries requirements



This option allows having information on Python installation and the required libraries.



❖ Close all figures



This option allows closing all figures.

❖ Read tutorial

A user manual is available (PDF file).

Menu EcoTransLearn_pac... x + Créer Se connecter

Tous les outils Modifier Convertir Signer Rechercher du texte ou des o...

EcoTransLearn R-package
Version 1.0-3

USER MANUAL

G. WACQUET

IFREMER
LABORATOIRE ENVIRONNEMENT LITTORAL & RESSOURCES AQUICOLES
UNITE LITTORAL
CENTRE MANCHE MER DU NORD
BOULOGNE-SUR-MER, FRANCE

210 x 297 mm

MORE...

What do you want to do?

- ☐ Process FlowCam data (collages)
- ☐ Generate train and test .CSV files
- ☐ Generate metadata file
- ☐ Install Python and libraries
- ☐ Close all figures
- ☒ Read tutorial
- ☐ About...

OK