



Fluffy Guide

0. Introduction	2
What is deep learning?	2
What is my workflow doing?	2
The command line isn't a pick-up line for computers...	2
What should I do if things fail?	3
1. Setting up the environment	4
Code Repository	4
Fiji	4
Anaconda	4
2. Generating labeled data	7
File setup	7
File annotation	8
File post-processing	9
3. Training and testing a model	11
4. Using a pre-trained model	13
About	14

0. Introduction

This guide will provide you with enough information to use and train a fluffily certified model for biomedical image segmentation. Read everything carefully and if you find things to be unclear or outdated let me know by mail – [bastian.eichenberger\(at\)fmi.ch](mailto:bastian.eichenberger(at)fmi.ch).

What is deep learning?

I might add more detail here. For now, if you're interested you can read this [introduction](#), this [paper](#) or this [video](#).

What is my workflow doing?

The entirety of my workflow is divided into four main steps:

1. Setting up the environment

Initially, the programming and application environment must be set up. This is important as programming applications constantly change. A piece of code is only guaranteed to run if the same dependency versions (code written and published by other people) are used.

2. Generating labeled data

For a deep learning model to train it needs data. Ideally the more the better. As some wise old man once said – “garbage in, garbage out” – the higher quality your input is, the better your output will be. So please don't rush this step and take your time. The number of images depends on how variable the images to predict will be. For example, if you try to segment a pattern in a high throughput assay which always looks the same you will need significantly less images than if you try to predict DAPI stained nuclei across multiple microscopes, magnifications, and samples. A good starting point for simpler models is 20-30 images.

3. Training and testing a model

This step automatically trains and tests a model. There is a lot to be said here and I will add more information soon.

4. Using a pre-trained model

Quite often there might already be a model available to solve your problem. In that case one simply has to use the model on the data to be analyzed.

The command line isn't a pick-up line for computers...

All subsequent code has to be executed in the command line. Below is a short tutorial to get you going. Note the commands are only for Unix systems (Linux / Mac). All windows users... Please consider your choices and use a real operating system.

To access the command line open the application “Terminal”. Type the commands below in the text field and execute them by simply pressing [enter].

A directory is nothing but a folder. All you need to know is how to find which directory you are in and how to change the current directory. There are three commands you will use:

```
pwd
ls
cd
```

- “pwd” – small P, small W, and small D. This stands for “present working directory” and tells you the complete path of your directory. For windows – “echo %cd%”.
- “ls” – small L and small S. This stands for “list” and lists all files in the current directory. Use this to find out what directories are available for you to go to. For windows – “dir”.
- “cd” – small “C and small D. This stands for “change directory” and is used to navigate to other directories. There are two directions of travel:
 - “cd ..” – Upwards / out of the current directory
 - “cd folder_xy” – Into the specified folder, here “folder_xy”

Always use the complete path given by “pwd” in the commands below. This ensures that no mixups between two equally named directories can happen. Watch this [video](#) for a more visual experience.

Note – the easiest way to get to your desired folder simply use the “cd” command and drag and drop the folder you want to access into the terminal window. This will automatically generate the right path.

```
cd # Drag and drop the desired folder
```

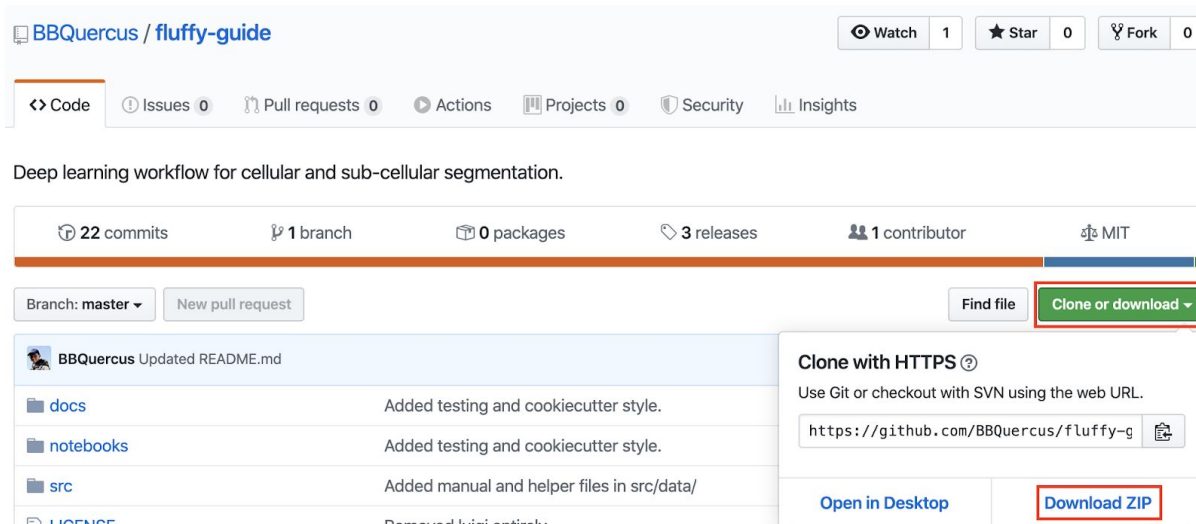
What should I do if things fail?

Don’t panic. There might still be some undetected bugs. Try to solve the problem **on your own** for at least 15 minutes by reading and understanding what went wrong. Also don’t forget to try turning it on and off. If things haven’t solved themselves feel free to contact me at the mail address above.

1. Setting up the environment

Code Repository

All code that will be used is found in the Github repository [here](#). Download its contents and unzip.



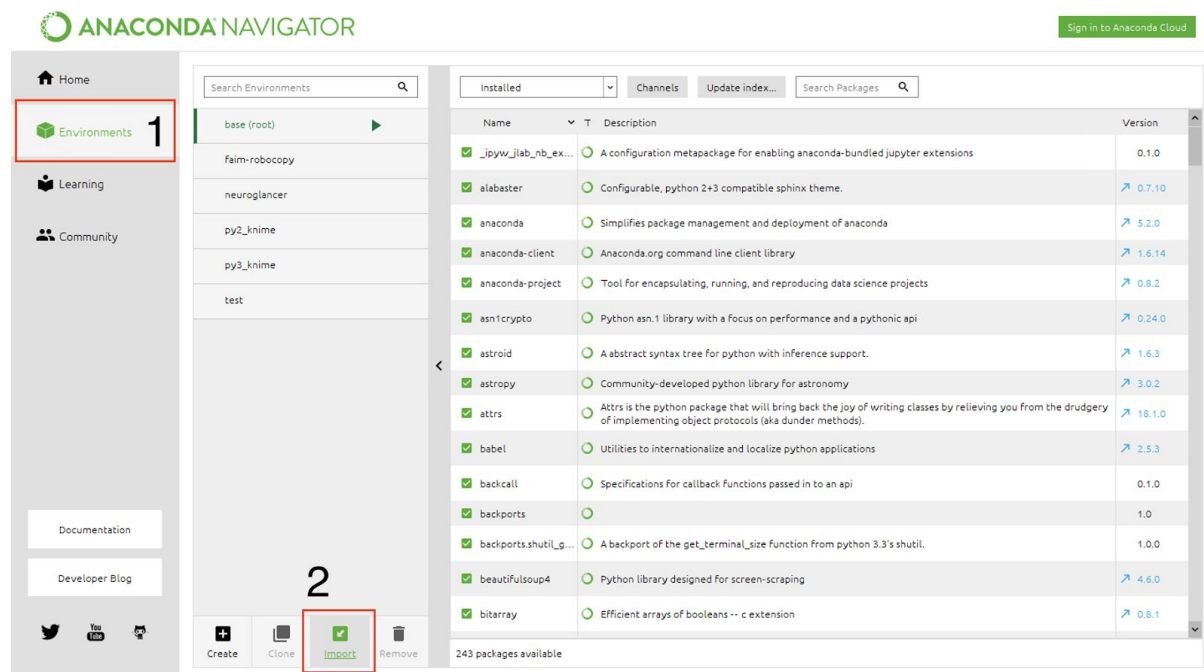
Practice your command line skills by navigating to the just downloaded repository.

Fiji

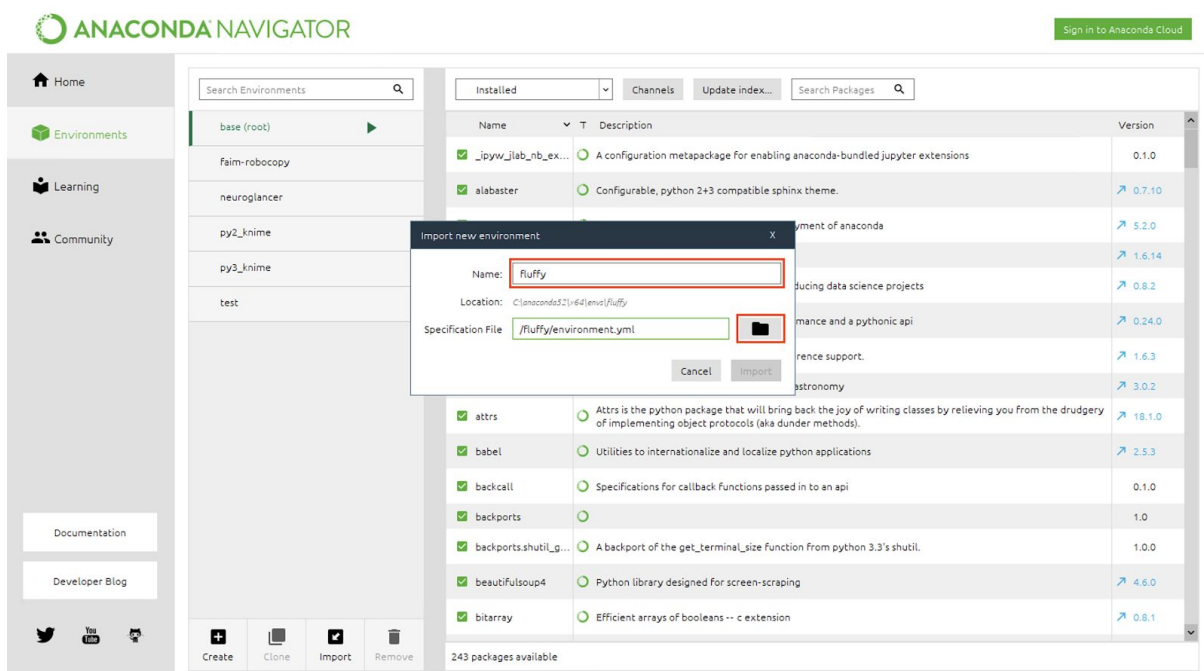
To annotate data we will need the free image processing program Fiji. Please download and install the latest version for your operating system [here](#).

Anaconda

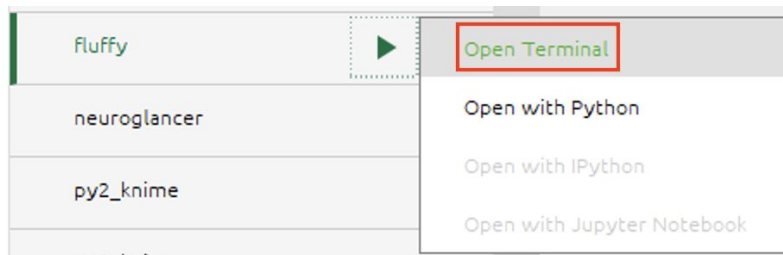
Anaconda is a package manager for python. We will use it to install all packages needed to run the workflow. Please follow the graphical installation manual for your operating system [here](#). Once installed, open up the Anaconda navigator. To create a new environment (set of python packages) go to the "Environments" tab and click on "Import".



In the popup window select a name for your environment such as “fluffy”. Select the “environment.yml” file which can be found in the code repository you downloaded above.



Follow the installation procedure and you should have a new “environment” (set of python packages) called “fluffy” (or your name of choice). Click on the arrow symbol next to the environment name and try to open a terminal from within the Anaconda navigator.



This completes your installation. Don't forget to always open a new terminal window through the Anaconda navigator.

2. Generating labeled data

Note – the ALL CAPS words used below must be changed to suit your specific situation. The names are selected to make sense if read in context.

File setup

Before we get to labeling images, we need to set up a proper file structure. Please run the following script in the command line after selecting a base directory (e.g. the Desktop) and the folder name (e.g. nucleus_labeling). Don't forget you can drag and drop folders to automatically get the path.

```
# From within the fluffy directory
cd src/data/
python make_directories.py

# Path to the base directory: PATH/TO/BASE/DIRECTORY
# Folder name: NAME_OF_FOLDER
```

Once run, you should have the basic file structure at your specified location. There should be three folders ("Raw", "Labeling", and "Processed") which will be described in more detail.

Add all images that you want to label to the "Raw" directory. Please only use one of the supported file types: ".stk", ".czi", ".tiff", ".tif", ".jpeg", ".jpg" or ".png". We will now convert all files into a uniform format which will provide the basis to start labeling. The trained model will only take 2D input. Therefore, images will be converted to 2D – specify whether the images in the "Raw" directory are Z-stacks (performs an automatic maximum projection) or time-lapses (only selects one time frame). Alternatively, if your images are already 2D you can simply say N (no) both times. Please run:

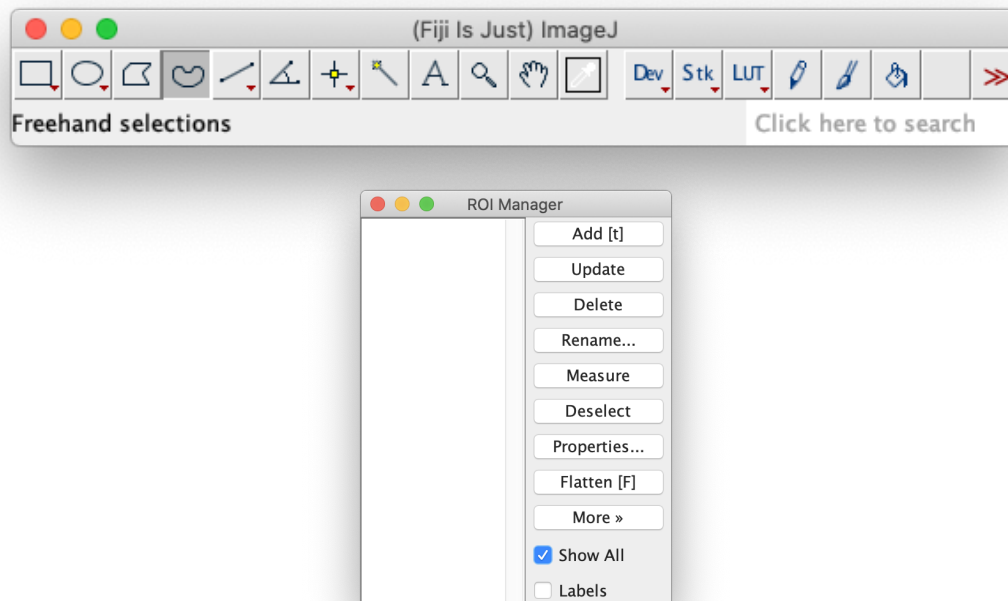
```
# From within the fluffy directory
cd src/data/
python make_labelling.py

# Path to the base directory: PATH/TO/BASE/DIRECTORY
# Timelapse [y/N]: CHOOSE
# ZStack [y/N]: CHOOSE
```

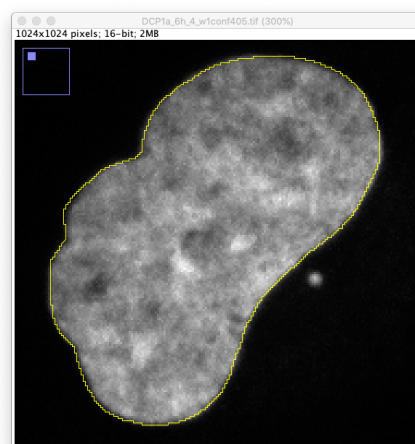
This just converted all files in the "Raw" directory to the 2D images of the ".tif" file format inside the "Labelling" directory.

File annotation

To annotate images, we will use Fiji. Select the “Freehand selections tool. Furthermore, open up the “ROI Manager” which can be found under “Analyze > Tools > ROI Manager”. Select the “Show All” option in the ROI Manager. Your setup should now look like this.



Open the first image from within the “Labeling/IMAGE_NAME/images/” directory and start drawing around the border of one object in your image. Here is an example of a nucleus.



After circelling one object click “Add” or [t] to add the annotated object to the ROI manager. An item with a numeric name should appear in the ROI manager. Continue circelling and adding the objects until all annotations have been added to the manager. Save all annotations by selecting “More > Save...” in the ROI manager. Navigate to the directory

"Labeling/IMAGE_NAME/" in the just opened popup window. Use the default name "RoiSet.zip" select "Save". To clear the cache of the ROI Manager select all annotations and press "Delete". Repeat the same procedure until all your images have been annotated.

File post-processing

If you followed the steps above you should now have "RoiSet.zip" files in all "Labeling/IMAGE_NAMES/" directories corresponding to that image. We will now use a Fiji macro to convert every set of annotations into label masks. Open up the "make_masks.ijm" file found in the "src/data/" directory inside fluffy by dragging it to the Fiji menu bar. At the very top of the document change "root" to your "Labeling" directory. Execute the macro by selecting "Run" right below the document.

All "Labeling/IMAGE_NAME/masks/" directories should now contain a bunch of files corresponding to each annotation labeled "mask_NUMBER.png". Check if your file structure looks the same:

```
# Inside the Processed directory
├── File_Name_1
│   ├── RoiSet.zip
│   ├── images
│   │   └── File_Name_1.tif
│   └── masks
│       ├── mask_0.png
│       ├── mask_....png
│       └── mask_n.png
├── File_Name_2
│   ├── RoiSet.zip
│   ├── images
│   │   └── File_Name_2.tif
│   └── masks
│       ├── mask_0.png
│       ├── mask_....png
│       └── mask_n.png
└── etc.
```

If this is the case, the last step is to merge all of these annotations into a single label map. You can do this by running:

```
# From within the fluffy directory
cd src/data/
python make_label_maps.py
```

```
# Path to base directory: PATH/TO/BASE/DIRECTORY
```

You should be greeted by a populated “Processed” directory. Inside, the image files and merged label maps are within the corresponding folders.

3. Training and testing a model

“Huston we have a problem” – Training a model, especially if the process is fully automated, takes a lot of computational power. To avoid waiting forever for training to finish one has to use GPUs. Unfortunately setting up GPUs isn’t as straightforward as one could hope for. Currently the easiest option is to speak to contact me once you have prepared enough labeled training data.

The automated training and testing process consists of the following steps:

1. Splitting of training and testing data

We split data into two main sets. One for training and one for testing. This step is important to evaluate the model performance. Our model learns from the training data set and is evaluated data it has never seen (testing data).

2. Bayesian probability based hyperparameter optimization

No set of training data is equal. Similarly, there is no one-size-fits-all neural network. Instead, a variety of parameters (hyperparameters) have to be tweaked and tuned to find suitable ones for the problem at hand. Bayesian probability improves this process by iteratively searching across a hyperparameter search space. At the end of this process one should have somewhat suitable hyperparameters for the following steps.

3. Cross validation

This step is not as important as the others but gives us a good feeling on how well the model generalizes (performs on other data). We try to validate the model by selecting different sets of training and validation data to train new models on. To learn more about this procedure you can read up [here](#).

4. Training one final model

Using the found hyperparameters and enough confidence provided by cross validation, we train a final model. This model is trained for longer to maximize the capability of the model to learn important features. This training step is performed on the entirety of training and validation data. The testing data is now used as validation to minimize overfitting.

5. Training multiple final models

Lastly, we have the option to train multiple models on the same dataset. Due to random weight initialization and floating point inaccuracies we will never get the exact same model. Therefore, we can pool together predictions from multiple models to increase performance.

Two main types are available – “binary” and “categorical”. The binary model outputs a True/False mask telling us if an object is there or not (e.g. for granules). The categorical model can be used to predict instances (individual objects) by also predicting object borders which can be used to separate nearby masks (e.g. for nuclei). Start the training regiment by running:

Note – this will change once the pipeline is complete.

```
# From within the fluffy directory
cd src/models/
```

```
python train_model.py \  
  --data_dir='PATH/TO/DATA/DIRECTORY' \  
  # Select if binary or categorical  
  --model_type='TYPE' \  
  --name='NAME_OF_MODEL'
```

4. Using a pre-trained model

Streamlit code will follow shortly. For now a bunch of pretrained models can be found [here](#). Use this code to test them out.

```
# From within the fluffy directory
cd src/models/
python predict_model.py

# Model file: PATH/TO/MODEL_FILE.h5
# Image file/folder: PATH/TO/FOLDER/OR/IMAGE.png
```

About

If you find this helpful for your research please cite:

```
@misc{Fluffy,  
  author = {Bastian Th., Eichenberger},  
  title = {Fluffy},  
  year = {2020},  
  publisher = {GitHub},  
  journal = {GitHub repository},  
  howpublished = {\url{https://github.com/bbquercus/fluffy}}
```

For assistance or to report bugs, please raise an issue on [GitHub](#). Title image designed by vector-pocket / [Freepik](#). A thank you to the Friedrich Miescher Institute's bioinformatics team that provides GPU access for easy model training.