

# **ENCORE. A practical implementation to improve reproducibility and transparency of computational research**

## **Supplementary Information**

### **Pre-defined sub-directories and files in the standardized file system structure (sFSS)**

The README Markdown files in the sFSS describe the content of the sub-directories and the content of the pre-defined files (including the README files) in more detail. The list below provides a summary of all sub-directories and pre-defined files.

ID_ProjectName	00_README-FIRST.{md, txt} 0_GETTINGSTARTED.{docx, tex, txt, html} 0_PROJECT.md 1_Step-by-Step-ENCORE-Guide.{docx,pdf} 2_CITATION.{md, txt} Navigate.py / Navigate_U.sh / Test_Navigate_Module.py Navigation.conf
o .navigate	
o Data	0_README.md
• NameOfDataset_1	
• Meta	0_README.md
• Processed	0_README.md
• Raw	0_README.md
o Processing	README.md, github.txt, gitignore-templates
• .git	
• 0_SoftwareEnvironment	0_README.md
• Anaconda	0_README-General.md, 0_README-ProjectSpecific.md
• C++	
• Matlab	
• Python	
• R	
• Data	
• NameOfDataset_1	
• Meta	
• Processed	
• Raw	
• NameOfComputation_1	
• Code	0_README.md
• CodeDocumentation	0_README.md
• Data	
• NameOfDataset_1	
• Meta	
• Processed	
• Raw	
• NoteBooks	0_README.md
• Results	0_README.md
• Settings	0_README.md
o ProjectDocumentation	LabJournal.{docx, tex, md, txt}
• BackgroundDocumentation	
• Literature	0_README.md
• MyPresentations	0_README.md
o Manuscript	0_README.md
o Sharing	0_README.md

14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26

**Supplementary Figure 1. The standardized File System Structure (FSS) and associated pre-defined files.** Standardized directory structure of the sFSS containing pre-defined files (brown), which include README files (in Markdown format) that provide a documentation template and instructions. Note that the pre-defined files in the 'Data' directories (red) and the '0\_SoftwareEnvironment' subdirectories are only shown once. The names of the directories 'NameOfDataset\_1' and 'NameOfComputation\_1' are placeholders and should be replaced with more descriptive names. These directories can be replicated if multiple datasets are used or if different computation procedures are performed. Subdirectories shown in light blue are under version control using Git/GitHub. The '0' prefix ensures that the corresponding files/directories are always on top of the file list when using lexicographic ordering. The README.md in 'Processing' is the default GitHub repository README file and therefore does not have the '0' prefix.

## Description of the sFSS sub-directories

### **\Manuscript**

This directory contains the (draft) manuscript(s) corresponding to this computational project including figures, tables, and supplementary information.

### **\Data**

**\Data\NameOfDataset\_1**

**\Data\NameOfDataset\_1\Raw**

**\Data\NameOfDataset\_1\Processed**

**\Data\NameOfDataset\_1\Meta**

These directories contain the raw, processed, and meta-data. Raw data comprises unprocessed data that come from the physical measurement device. Processed data comprises data obtained from collaborators or public databases and, consequently, not produced as part of the computational analyses. If the data (pre)processing is part of the computational analyses then, preferably, it should be placed in the \Results directory within \Processing. However, ENCORE allows the flexibility to store one's own processed data in the \Data\ NameOfDataset\_1\Processed directory. Meta-data is the description of the data including the data license, description of the samples, experimental design, content and format of the data files, etc.

### **\Processing**

Contains all sub-directories and pre-defined files related to the computational part of the project.

**\Processing\0\_SoftwareEnvironment**

**\Processing\0\_SoftwareEnvironment\Anaconda**

**\Processing\0\_SoftwareEnvironment\C++**

**\Processing\0\_SoftwareEnvironment\Matlab**

**\Processing\0\_SoftwareEnvironment\Python**

**\Processing\0\_SoftwareEnvironment\R**

One challenge that is only partially addressed by ENCORE concerns the preservation of the full computing environment. This environment is defined by (interdependencies of) the operating system, software tools, versions and dependencies, programming language libraries, etc. Gruning and co-workers proposed a software stack of interconnected technologies to preserve the computing environment (Gruning, Chilton, et al., 2018). This stack comprises (Bio)Conda (Anaconda Software Distribution, 2020; Gruning, Dale, et al., 2018) to provide virtual execution environments addressing software versions and dependencies, container platforms such as Docker (Nust et al., 2020) to preserve other aspects of the runtime environment, and virtual machines using cloud systems or dedicated applications such as VMware, to overcome the dependencies on the operating system and hardware. We are currently investigating how to best approach this within the ENCORE environment.

However, some basic information about the computing environment (e.g., export of Conda environments) can be stored in this directory.

The sub-directories provide basic information about different environments (e.g., R/Rstudio, Python/PyCharm, Anaconda) for peers not familiar with the used computing environment. In addition, one may find other files such as cheat sheets, tutorials, and exports of (Anaconda) environments.

## **\Processing\Data**

See \Data

## **\Processing \NameOfComputation\_1**

- This directory should also contain a conceptual description of applied methodology to improve transparency. For example,
  - Brief description of methods (and version) used including specification of the mathematical/statistical model, parameters, variables, references, etc.
  - If a new method is developed, this method should be described in full detail.
  - Description of why the selected or developed computational approach is valid for the research question that is addressed. This enables peers to make their own judgement about the approach and results.
  - Considered alternatives?
  - Detailed description of all data filtering, reduction, normalization etc. steps that are performed prior to the downstream analysis.
  - Avenues of exploration examined throughout development, including information about negative findings.

## **\Processing \NameOfComputation\_1\Code**

Contains the (in-house developed) software used for the computational analysis.

## **\Processing \NameOfComputation\_1\CodeDocumentation**

External (user) documentation of the code. Possibly automatically generated with documentation tools such as Sphinx.

## **\Processing \NameOfComputation\_1\Data**

See \Data

## **\Processing \NameOfComputation\_1\NoteBooks**

Notebooks ((web-based) interactive computing platform that combines live code, equations, narrative text, visualizations etc.) should be placed in this sub-directory.

## **\Processing \NameOfComputation\_1\Results**

(Intermediate) results (e.g., figures, tables) from the computational analysis. Record of intermediate results (preferably in a standardized format). It is strongly advised to generate hierarchical analysis output, allowing layers of increasing detail to be inspected. This can reveal discrepancies toward what is assumed and can in this way uncover bugs or faulty interpretation that is not apparent in the final results. It also allows any inconsistency to be tracked to the step where the problem occurs. Furthermore, it enables critical examination of the full process behind a result. It is

also advised to clearly document the intermediate/final results and the imposed hierarchy.

For any figure or table that ends up in a publication, report, or presentation at a meeting, the underlying data and a stand-alone piece of code should be available to regenerate the figure. This also allows easy modification of a figure and retrieval of the data corresponding to a figure (instead of having to redo a complete analysis). Equally important, the data corresponding to a figure can be further analyzed or inspected.

#### **\Processing \NameOfComputation\_1\Settings**

This file/sub-directory concerns settings/parameters for the algorithms that have been developed. For settings related to the computing environment, see \0\_SoftwareEnvironment for further instructions.

### **\ProjectDocumentation**

This subdirectory contains (background) information about any part of the project. However, as a rule, documentation should be close to the component (e.g., data, code) that is described. For example, the documentation about the data should be in the \Data directory. However, more general information can be placed in the sub-directories in \ProjectDocumentation.

#### **\ProjectDocumentation\BackgroundDocumentation**

Documents relevant as project background documentation. For example, the project proposal, presentations from collaborators or peers (thus, not from the project team), and relevant tutorials about applied methodology.

#### **\ProjectDocumentation\Literature**

This subdirectory should contain relevant scientific literature (e.g., PDF files) obtained from PubMed, bioRxiv, etc. The PDF files should follow the naming convention Author-Year-Journal (or similar) to allow easy retrieval during project discussions. In addition, it should contain the reference manager (e.g., EndNote, Mendeley, BibTeX) file and an export to the standardized RIS format. The README file in this sub-directory should briefly describe the relevance of each paper (e.g., specific information relevant for the computational analyses).

#### **\ProjectDocumentation\MyPresentations.**

This sub-directory contains oral or poster presentations (and abstracts) given by the project team during meetings (e.g., at progress meetings, seminars, conferences).

### **\Sharing**

*Rationale.* In general, the complete file system structure (FSS) and its contents should be shared unmodified with peers that aim at reproducing the computational analysis. However, in specific cases it might be desirable to share only parts of the sFSS and/or restructure the sFSS. The reduced and/or restructured sFSS is then stored (as a compressed file) in the \Sharing directory. This file should at least indicate what one did (not) share and how/why the sFSS was restructured. In addition, document when and with whom this directory was shared.

*Typical use.* A typical use of the \Sharing directory is for support projects in which computational analyses are performed for other researchers as a service (e.g., biomedical, clinical researchers). These researchers might only be interested in the final results (figures and tables) and not in the code that produced these results. That is, they will not aim to repeat the analyses. In such a situation, the results and tables can be shared in a (flat) structure that is more convenient for them to browse and use, and that leaves out all code and background documentation.

*Restrictions.* Sharing an sFSS with an (external) colleague may be restricted due to, for example, copyright on PDF files of papers, sensitive/private information (in LabJournal.docx), non-open-access of data, etc. Make sure to remove such information from \Sharing.

## Description of the sFSS pre-defined files

### **00\_README-FIRST.{md, txt}, 0\_README.md / README.md**

Throughout the sFSS there are README files that explain the content of the sub-directories and provide instructions and a template to guide documentation of the project. Most of these files are so-called Markdown files that can be opened in any text editor, but require a Markdown viewer (e.g., Notepad++, Typora) to show the resulting document.

To facilitate first-time users, the 00\_README-FIRST file is also provided as a text file that contains instructions w.r.t. the Markdown files.

Each \*\_README file starts with a '0\_' or '00\_' prefix to ensure it appears at the top of the file list. The only exception is the README.md file in the /Processing directory, which is the default GitHub README file that should not contain a prefix. Because the sFSS (and not GitHub) is the entry point for a project, the GitHub README.md file does not necessarily have to contain a project description. However, you may want to copy the information from 0\_PROJECT.md (see below) into this README.md file. More importantly, it should provide an explanation of the code in the Processing directory and instructions about its execution.

### **0\_PROJECT.md**

Short description of project, contact person, and project team. This file is used by the sFSS Navigator.

### **0\_GETTINGSTARTED.txt.**

Template document (plain text format). One can copy this file to one's favorite editor to add content. Examples:

- **0\_GETTINGSTARTED.docx.** Template document (Microsoft Word format containing an example of how to include links). The DOCX file can be saved as HTML (make sure to use UTF-8 encoding).
- **0\_GETTINGSTARTED.tex.** Template document (LaTeX format containing an example of how to include links). The LaTeX file can be converted to HTML with pandoc (<https://pandoc.org/index.html>).

- **0\_GETTINGSTARTED.html.** Example of exported HTML file used by the sFSS Navigator.

The main use of the GETTINGSTARTED files is to guide a first-time user of a finished project to the most important aspects (e.g., results, code) of the project before he/she manually explores other information contained in the sFSS. The GETTINGSTARTED files provide links to the relevant sub-directories and files. 'Getting started' templates are provided in different file formats, which can be converted to HTML once finished. The 0\_PROJECT.md and 0\_GETTINGSTARTED.html files are used by the sFSS Navigator

#### Help files

- **1\_Step-by-Step-ENCORE-Guide.{pdf, docx}.**  
User guide for ENCORE, describing the standardized File System Structure (sFSS) and how to set up a corresponding GitHub repository.

#### General

- **2\_CITATION.{md,txt}.** How to cite ENCORE and the sFSS Navigator
- **.FSSignore.** Currently not used but to be used with an application that selects all files needed for sharing.

#### sFSS Navigator

- **Navigate.py.** Standalone Python 3 script to generate Navigate.html used for navigating the sFSS. Can be run from the command line (Navigate.py -h)
- **Navigate\_U.sh.** Shell script to run Navigate on Unix/Linux systems. Change the first line (!/usr/bin/Python) if necessary. Make executable using chmod +x
- **Navigate.html.** Open in a browser to navigate the standardized file system. This file is created after running Navigate.py.

There are also executables available for Windows and macOS. These are available from the latest GitHub release and from Zenodo (DOI: <https://doi.org/10.5281/zenodo.7985655>):

- **Navigate\_W.exe.** Windows executable if one doesn't have Python installed (Navigate.exe -h).
- **Navigate\_M.** MacOS executable (macOS 13.3.1 (Ventura), Apple M1)
- **Navigate\_MacIntel.** MacOS executable (macOS 10.13.6 (High Sierra), Intel Core i5)

#### Test\_Navigate\_Module.py.

Python script to show how to use Navigate.py as a module in other Python scripts. This may help to keep Navigate.html up to date without manually executing Navigate.py.

#### Navigation.conf.

Configuration file for the sFSS Navigator.

#### \ProjectDocumentation\LabJournal.{docx, md, tex, txt}.

Templates in different file formats for the lab journal.

In general, any documentation should be kept in the sub-directory that it describes. Thus, use the 0\_README.md and/or additional (e.g., PowerPoint) files to document the data, software, and results in their respective directories.

The lab journal should contain more general documentation. For example,

- General information and concepts
- Summaries of project discussions
- Steps to be taken
- New (future) research ideas
- Pointers to the location of certain pieces of information

In addition, it may also contain integrated parts of the various README files whenever useful (but keep it consistent with the source files). Include figures and tables when necessary.

Although, strictly speaking, a lab journal is not intended for recording new/future ideas or providing summaries of discussions, it is important for a research group to also have a record of this. Therefore, these can also be included in the lab journal.

Parts of the lab journal that should not be shared with peers (e.g., new research ideas) should be clearly labelled with 'Not for sharing' such that one can easily remove these parts. Alternatively, one may maintain two separate documents.

If necessary, ENCORE allows to maintain lab journals in multiple sub-directories. For example, one may decide to have a separate lab journal in '\Processing\NameOfComputation\_1' for a specific part of the computational analysis.

#### **\Processing\github.txt.**

Provides the URL to the corresponding GitHub repository and any other relevant information about the repository. This file is used by the sFSS Navigator.

#### **\Processing\gitignore-FSS-template.txt and .gitignore**

The template file should be adapted (if needed) and then be renamed to .gitignore. It contains instructions for Git about files that should not be synchronized with the GitHub repository (e.g., data and results). This file should be modified depending on the contents of \Processing. Optionally, one can use other language-specific templates that are also found in \Processing:

- gitignore-C++-template.txt
- gitignore-FSS-template.txt
- gitignore-JetBrains-template.txt
- gitignore-Matlab-template.txt
- gitignore-Python-template.txt
- gitignore-R-template.txt

To use any of these templates, simply merge its content into .gitignore. It is considered good practice to keep a single .gitignore in the top-level directory and not in individual subdirectories, which would make debugging more troublesome.