

# BIG Data, BIG responsibility

Maneage: *Managing data lineage for long-term and archivable reproducibility*

(Published in CiSE 23 (3), pp 82-91: DOI:10.1109/MCSE.2021.3072860, arXiv:2006.03018)

Mohammad Akhlaghi

Centro de Estudios de Física del Cosmos de Aragón (CEFCA), Teruel, Spain

SoftwareHeritage 5th Anniversary  
November 30th, 2021 (Inria, Paris)

Most recent slides available in link below (this PDF is built from Git commit fc30dd7):

<https://maneage.org/pdf/slides-intro-short.pdf>



Financiado por la Unión Europea-NextGenerationEU



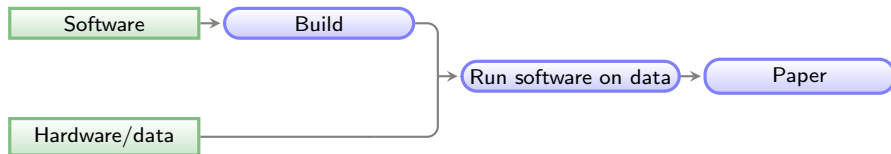
"Una manera de hacer Europa"



This project received financial support from the European Union's Horizon 2020 research and innovation grant agreement No. 717388 to the RDA EU 4.0 project.



## General outline of a project (after data collection)

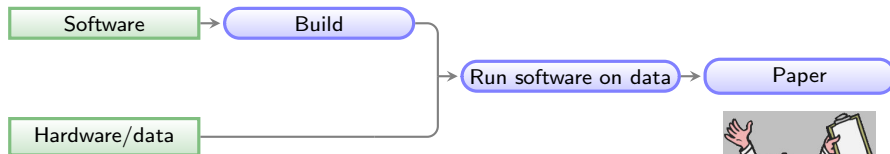


Green boxes with sharp corners: *source*/input components/files.

Blue boxes with rounded corners: *built* components.

Red boxes with dashed borders: questions that must be clarified for each phase.

## General outline of a project (after data collection)



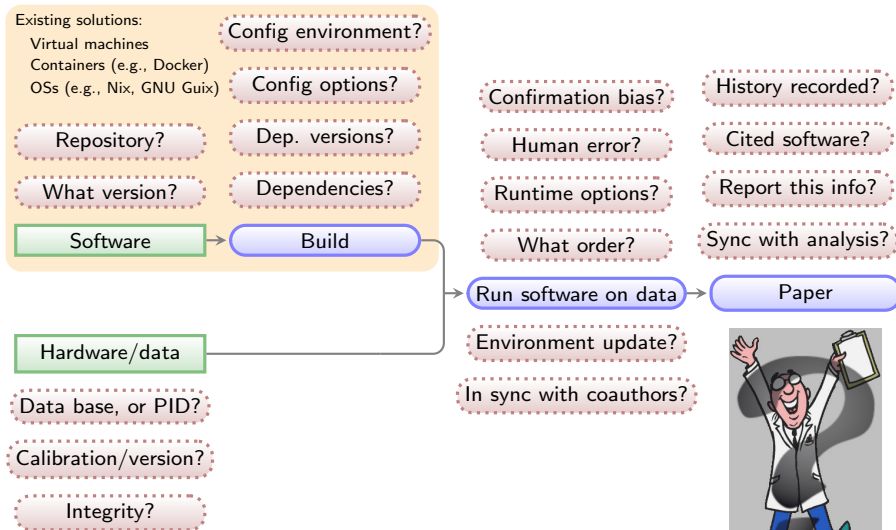
Green boxes with sharp corners: *source*/input components/files.

Blue boxes with rounded corners: *built* components.

Red boxes with dashed borders: questions that must be clarified for each phase.

<https://heywhatwhatdidyousay.wordpress.com>

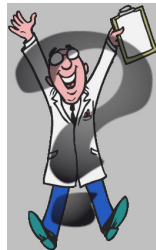
# General outline of a project (after data collection)



Green boxes with sharp corners: *source*/input components/files.

Blue boxes with rounded corners: *built* components.

Red boxes with dashed borders: questions that must be clarified for each phase.



<https://heywhatwhatdidyousay.wordpress.com>

<http://pngimages.net>



## Science is a tricky business

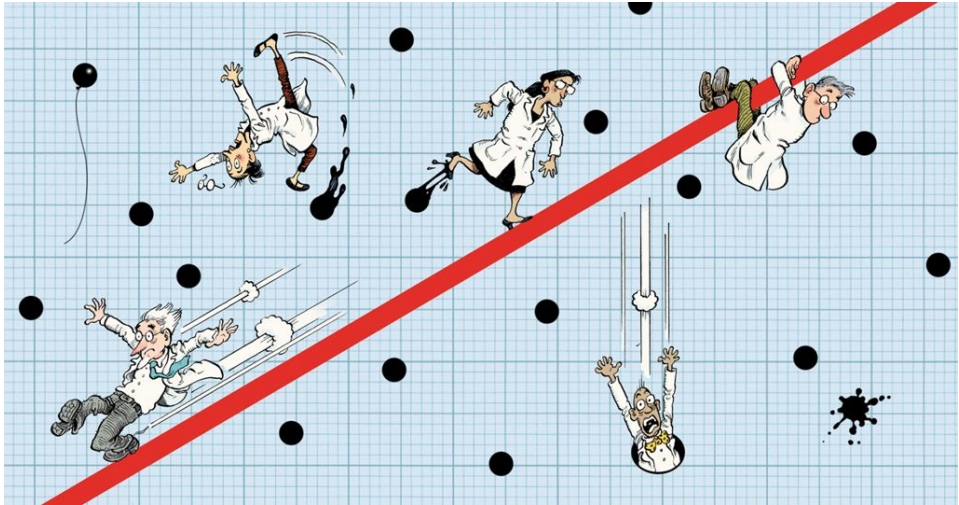
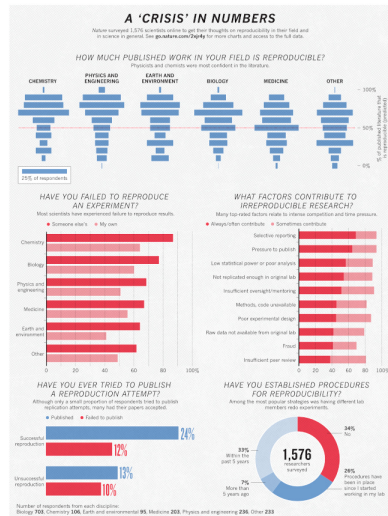
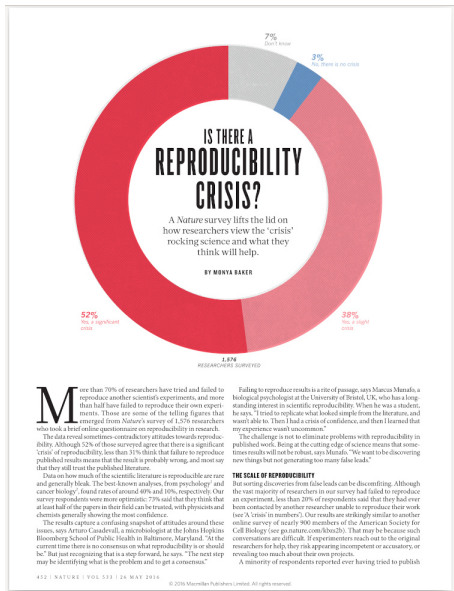


Image from nature.com ("Five ways to fix statistics", Nov 2017)

Data analysis [...] is a **human behavior**. Researchers who hunt hard enough will turn up a result that fits statistical criteria, but their **discovery** will probably be a **false positive**.

Five ways to fix statistics, Nature, 551, Nov 2017.

# "Reproducibility crisis" in the sciences? (Baker 2016, Nature 533, 452)



## Founding criteria

Basic/simple principle:

Science is defined by its METHOD, **not** its result.

## Founding criteria

Basic/simple principle:

Science is defined by its METHOD, **not** its result.

### ► Complete/self-contained:

- **Only dependency** should be **POSIX** tools (discards Conda or Jupyter which need Python).
- Must **not require root** permissions (discards tools like Docker or Nix/Guix).
- Should be **non-interactive** or runnable in batch (user interaction is an incompleteness).
- Should be usable **without internet** connection.

## Founding criteria

Basic/simple principle:

Science is defined by its METHOD, **not** its result.

▶ **Complete/self-contained:**

- ▶ **Only dependency** should be **POSIX** tools (discards Conda or Jupyter which need Python).
- ▶ Must **not require root** permissions (discards tools like Docker or Nix/Guix).
- ▶ Should be **non-interactive** or runnable in batch (user interaction is an incompleteness).
- ▶ Should be usable **without internet** connection.

▶ **Modularity:** Parts of the project should be **re-usable** in other projects.

## Founding criteria

Basic/simple principle:

Science is defined by its METHOD, **not** its result.

- ▶ **Complete/self-contained:**
  - ▶ **Only dependency** should be **POSIX** tools (discards Conda or Jupyter which need Python).
  - ▶ Must **not require root** permissions (discards tools like Docker or Nix/Guix).
  - ▶ Should be **non-interactive** or runnable in batch (user interaction is an incompleteness).
  - ▶ Should be usable **without internet** connection.
- ▶ **Modularity:** Parts of the project should be **re-usable** in other projects.
- ▶ **Plain text:** Project's source should be in **plain-text** (binary formats need special software)
  - ▶ This includes high-level analysis.
  - ▶ It is easily publishable (very low volume of  $\times 100\text{KB}$ ), archivable, and parse-able.
  - ▶ **Version control** (e.g., with Git) can track project's history.

# Founding criteria

Basic/simple principle:

Science is defined by its METHOD, **not** its result.

## ▶ Complete/self-contained:

- ▶ **Only dependency** should be **POSIX** tools (discards Conda or Jupyter which need Python).
- ▶ Must **not require root** permissions (discards tools like Docker or Nix/Guix).
- ▶ Should be **non-interactive** or runnable in batch (user interaction is an incompleteness).
- ▶ Should be usable **without internet** connection.

## ▶ Modularity: Parts of the project should be **re-usable** in other projects.

## ▶ Plain text: Project's source should be in **plain-text** (binary formats need special software)

- ▶ This includes high-level analysis.
- ▶ It is easily publishable (very low volume of  $\times 100\text{KB}$ ), archivable, and parse-able.
- ▶ **Version control** (e.g., with Git) can track project's history.

## ▶ Minimal complexity: Occum's razor: "Never posit pluralities without necessity".

- ▶ Avoiding the **fashionable** tool of the day: tomorrow another tool will take its place!
- ▶ Easier **learning curve**, also doesn't create a **generational gap**.
- ▶ Is **compatible** and **extensible**.

# Founding criteria

Basic/simple principle:

Science is defined by its METHOD, **not** its result.

- ▶ **Complete/self-contained:**
  - ▶ **Only dependency** should be **POSIX** tools (discards Conda or Jupyter which need Python).
  - ▶ Must **not require root** permissions (discards tools like Docker or Nix/Guix).
  - ▶ Should be **non-interactive** or runnable in batch (user interaction is an incompleteness).
  - ▶ Should be usable **without internet** connection.
- ▶ **Modularity:** Parts of the project should be **re-usable** in other projects.
- ▶ **Plain text:** Project's source should be in **plain-text** (binary formats need special software)
  - ▶ This includes high-level analysis.
  - ▶ It is easily publishable (very low volume of  $\times 100\text{KB}$ ), archivable, and parse-able.
  - ▶ **Version control** (e.g., with Git) can track project's history.
- ▶ **Minimal complexity:** Occum's razor: "Never posit pluralities without necessity".
  - ▶ Avoiding the **fashionable** tool of the day: tomorrow another tool will take its place!
  - ▶ Easier **learning curve**, also doesn't create a **generational gap**.
  - ▶ Is **compatible** and **extensible**.
- ▶ **Verifiable inputs and outputs:** Inputs and Outputs must be **automatically verified**.

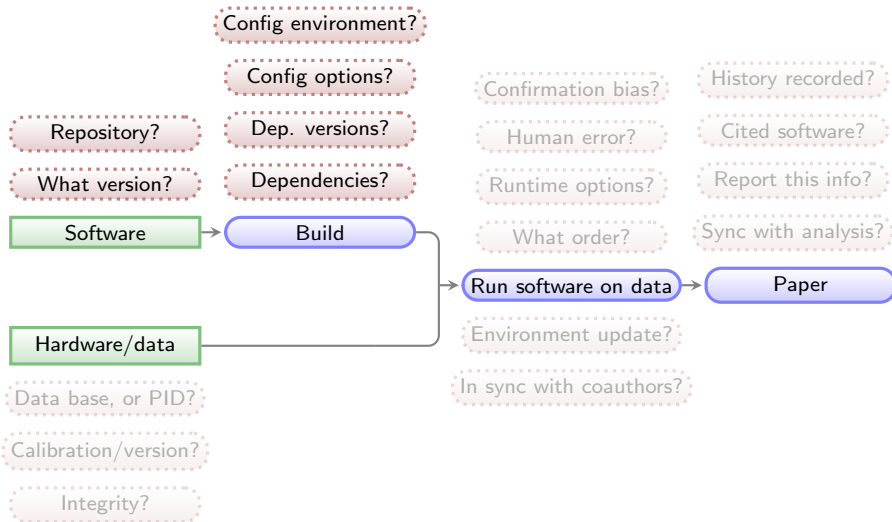


Basic/simple principle:

Science is defined by its METHOD, **not** its result.

- ▶ **Complete/self-contained:**
  - ▶ **Only dependency** should be **POSIX** tools (discards Conda or Jupyter which need Python).
  - ▶ Must **not require root** permissions (discards tools like Docker or Nix/Guix).
  - ▶ Should be **non-interactive** or runnable in batch (user interaction is an incompleteness).
  - ▶ Should be usable **without internet** connection.
- ▶ **Modularity:** Parts of the project should be **re-usable** in other projects.
- ▶ **Plain text:** Project's source should be in **plain-text** (binary formats need special software)
  - ▶ This includes high-level analysis.
  - ▶ It is easily publishable (very low volume of  $\times 100\text{KB}$ ), archivable, and parse-able.
  - ▶ **Version control** (e.g., with Git) can track project's history.
- ▶ **Minimal complexity:** Occum's razor: "Never posit pluralities without necessity".
  - ▶ Avoiding the **fashionable** tool of the day: tomorrow another tool will take its place!
  - ▶ Easier **learning curve**, also doesn't create a **generational gap**.
  - ▶ Is **compatible** and **extensible**.
- ▶ **Verifiable inputs and outputs:** Inputs and Outputs must be **automatically verified**.
- ▶ **Free and open source software:** **Free software** is essential: non-free software is not configurable, not distributable, and dependent on non-free provider (which may discontinue it in N years).

## General outline of a project (after data collection)



Green boxes with sharp corners: *source*/input components/files.

Blue boxes with rounded corners: *built* components.

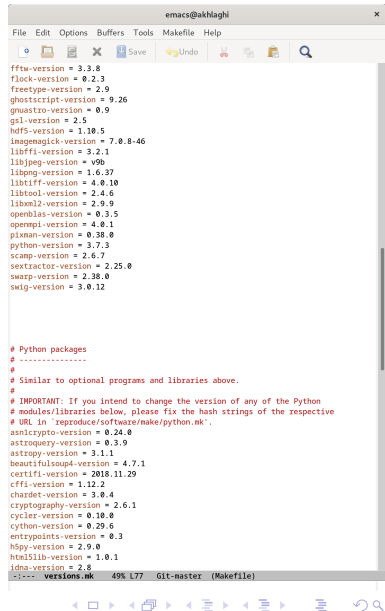
Red boxes with dashed borders: questions that must be clarified for each phase.

# Predefined/exact software tools

## Reproducibility & software

Reproducing the environment (specific **software versions**, **build instructions** and **dependencies**) is also critically important for reproducibility.

- ▶ *Containers or Virtual Machines* are a **binary black box**.
  - ▶ e.g., with 'FROM ubuntu:16.04' (released in April 2016),
  - ▶ in a Dockerfile, the OS image will come from (updated monthly!): <https://partner-images.canonical.com/core/xenial>
- ▶ Maneage **installs fixed versions** of all necessary research software.
  - ▶ Including their dependencies.
  - ▶ All the way down to the C compiler.
- ▶ Installs similar environment on **GNU/Linux**, or **macOS** systems.
- ▶ Works like a package manager (e.g., **apt**, **brew** or Conda).
  - ▶ ... **but (!)**, its not a third party package manager.
  - ▶ Build instructions are within same analysis project.
  - ▶ e.g., see Conda's build of Gnuastro (its gets updated behind your back): <https://anaconda.org/conda-forge/gnuastro/files>
- ▶ Source code of all software in Maneage is archived on [zenodo.3883409](https://zenodo.org/record/3883409).



```
emacs@akhlaghi
File Edit Options Buffers Tools Makefile Help
[Icons] Save Undo [Icons]

fftw-version = 3.3.8
flock-version = 0.2.3
freetype-version = 2.9
ghostscript-version = 9.26
gnuastro-version = 0.9
gsl-version = 2.5
hdf5-version = 1.10.5
imagemagick-version = 7.0.8-46
libffi-version = 3.2.1
libjpeg-version = v9b
libpng-version = 1.6.37
libtiff-version = 4.0.10
libtool-version = 2.4.6
libxml2-version = 2.9.9
openblas-version = 0.3.5
openmpi-version = 4.0.1
pixman-version = 0.38.0
python-version = 3.7.3
scamp-version = 2.6.7
sexttractor-version = 2.25.0
swarp-version = 2.38.0
swig-version = 3.0.12

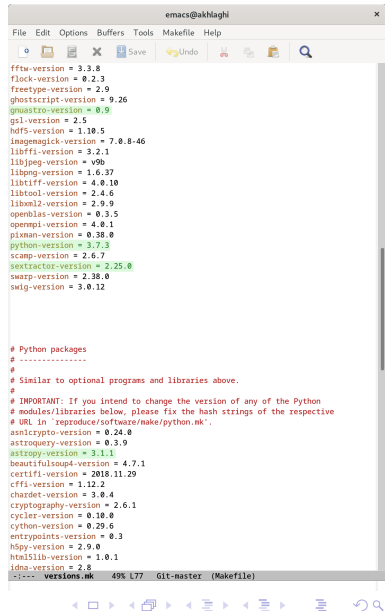
# Python packages
# -----
#
# Similar to optional programs and libraries above.
#
# IMPORTANT: If you intend to change the version of any of the Python
# modules/libraries below, please fix the hash strings of the respective
# URL in 'reproduce/software/make/python.mk'.
asnycrypto-version = 0.24.0
astroquery-version = 0.3.9
astropy-version = 3.1.1
beautifulsoup4-version = 4.7.1
certifi-version = 2018.11.29
cffi-version = 1.12.2
chardet-version = 3.0.4
cryptography-version = 2.6.1
cycler-version = 0.10.0
cython-version = 0.29.6
entrypoints-version = 0.3
h5py-version = 2.9.0
html5lib-version = 1.0.1
idna-version = 2.8
-i--- versions.mk 49% L77 Git-master (Makefile)
```

# Predefined/exact software tools

## Reproducibility & software

Reproducing the environment (specific **software versions**, **build instructions** and **dependencies**) is also critically important for reproducibility.

- ▶ *Containers or Virtual Machines* are a **binary black box**.
  - ▶ e.g., with 'FROM ubuntu:16.04' (released in April 2016),
  - ▶ in a Dockerfile, the OS image will come from (updated monthly!): <https://partner-images.canonical.com/core/xenial>
- ▶ Maneage **installs fixed versions** of all necessary research software.
  - ▶ Including their dependencies.
  - ▶ All the way down to the C compiler.
- ▶ Installs similar environment on **GNU/Linux**, or **macOS** systems.
- ▶ Works like a package manager (e.g., **apt**, **brew** or Conda).
  - ▶ ... **but (!)**, its not a third party package manager.
  - ▶ Build instructions are within same analysis project.
  - ▶ e.g., see Conda's build of Gnuastro (its gets updated behind your back): <https://anaconda.org/conda-forge/gnuastro/files>
- ▶ Source code of all software in Maneage is archived on [zenodo.3883409](https://zenodo.org/record/3883409).



```
fftw-version = 3.3.8
flock-version = 0.2.3
freetype-version = 2.9
ghostscript-version = 9.26
gnuastro-version = 0.9
gsl-version = 2.5
hdff5-version = 1.10.5
imagemagick-version = 7.0.8-46
libffi-version = 3.2.1
libjpeg-version = v9b
libpng-version = 1.6.37
libtiff-version = 4.0.10
libtool-version = 2.4.6
libxml2-version = 2.9.9
openblas-version = 0.3.5
openmpi-version = 4.0.1
pixman-version = 0.38.0
python-version = 3.7.3
scamp-version = 2.6.7
sextractor-version = 2.25.0
swarp-version = 2.38.0
swig-version = 3.0.12

# Python packages
# -----
#
# Similar to optional programs and libraries above.
#
# IMPORTANT: If you intend to change the version of any of the Python
# modules/libraries below, please fix the hash strings of the respective
# URL in 'reproduce/software/make/python.mk'.
asn1crypto-version = 0.24.0
astroquery-version = 0.3.9
astropy-version = 3.1.1
beautifulsoup4-version = 4.7.1
certifi-version = 2018.11.29
cffi-version = 1.12.2
chardet-version = 3.0.4
cryptography-version = 2.6.1
cycler-version = 0.10.0
cython-version = 0.29.6
entrypoints-version = 0.3
h5py-version = 2.9.0
html5lib-version = 1.0.1
idna-version = 2.8
-i-- versions.mk 49% L77 Git-master (Makefile)
```

# Example: Matplotlib (a Python visualization library) build dependencies

## Matplotlib library

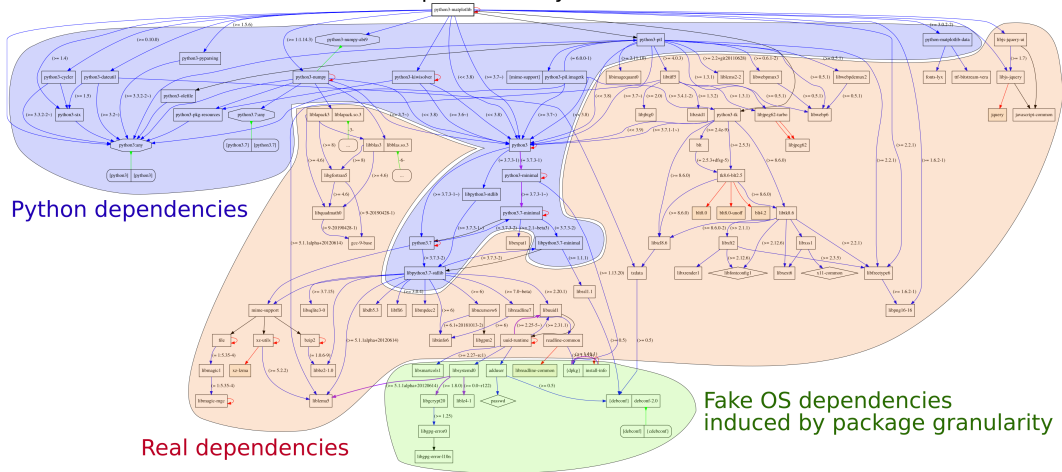


Fig. 1. Transitive dependencies of the software environment required by a simple "import matplotlib" command in the Python 3 interpreter.

## Advantages of this build system

- ▶ Project runs in fixed/controlled environment: custom build of **Bash**, **Make**, GNU Coreutils (**ls**, **cp**, **mkdir** and etc), **AWK**, or **SED**, **L<sup>A</sup>T<sub>E</sub>X**, etc.
- ▶ No need for **root**/administrator **permissions** (on servers or super computers).
- ▶ Whole system is built **automatically** on any Unix-like operating system (less 2 hours).
- ▶ Dependencies of different projects will **not conflict**.
- ▶ Everything in **plain text** (human & computer readable/archivable).



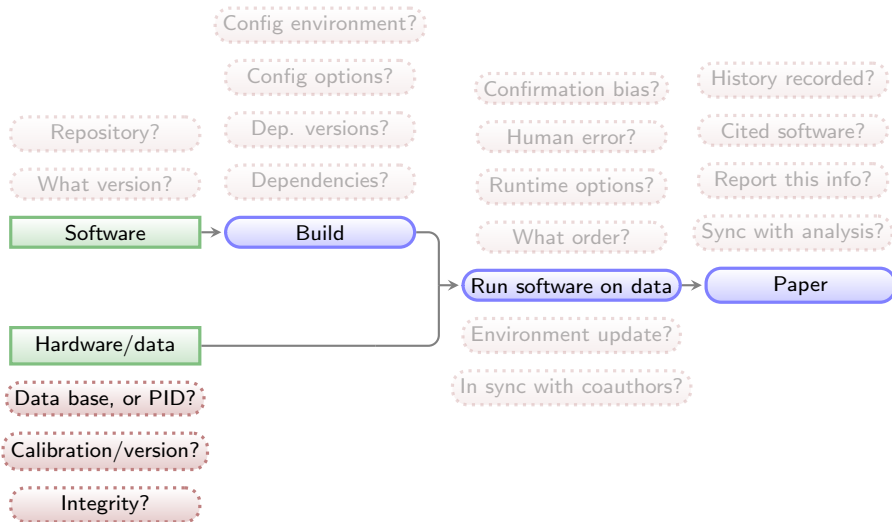
<https://natemowry2.wordpress.com>

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡

## YOUTH NAME: 000000



## General outline of a project (after data collection)



Green boxes with sharp corners: *source*/input components/files.

Blue boxes with rounded corners: *built* components.

Red boxes with dashed borders: questions that must be clarified for each phase.

# Input data source and integrity is documented and checked

Stored information about each input file:

- ▶ **PID** (where available).
- ▶ Download **URL**.
- ▶ **MD5**-sum to check integrity.

All inputs are **downloaded** from the given PID/URL when necessary (during the analysis).

MD5-sums are **checked** to make sure the download was done properly or the file is the same (hasn't changed on the server/source).

Example from the reproducible paper [arXiv:1909.11230](https://arxiv.org/abs/1909.11230).  
This paper needs three input files (two images, one catalog).



```
emacs@akhlaghi
File Edit Options Buffers Tools Makefile Help
[Icons] Save Undo [Icons] Search

# Input files necessary for this project.
#
# This file is read by the configure script and running Makefiles.
#
# Copyright (C) 2018-2019 Mohammad Akhlaghi <mohammad@akhlaghi.org>
#
# Copying and distribution of this file, with or without modification, are
# permitted in any medium without royalty provided the copyright notice and
# this notice are preserved. This file is offered as-is, without any
# warranty.

W51S05SRURL = https://dr12.sdss.org/sas/dr12/boos/photoObj/frames/301/3716/6
W51S05SRIMAGE = frame-r-003716-6-0117.fits.bz2
W51S05SRMDS = 965da8bd861e94a9701521a11b2d88aa
W51S05SRSIZE = 2.8M

XDF775WURL = http://archive.stsci.edu/pub/hlsp/xdff
XDF775WIMGE = hlsp_xdf_hst_acswfc-60mas_hudf_f775w_v1_sc1.fits
XDF775WMDS = 81408ed0949bd3a93c4bfe7e229472e6
XDF775WSIZE = 106M

UVUDFSEGURL = https://asd.gsfc.nasa.gov/UVUDF
UVUDFSEGIMAGE = segmentation_map_rafelski_2015.fits.gz
UVUDFSEGMDS = 29d5b3e5311b77512ba727db6ad0e11b
UVUDFSEGSIZE = 1.3M

-:--- INPUTS.mk All L1 Git-master (GNUmakefile)
For information about GNU Emacs and the GNU system, type C-h C-a.
```

# Input data source and integrity is documented and checked

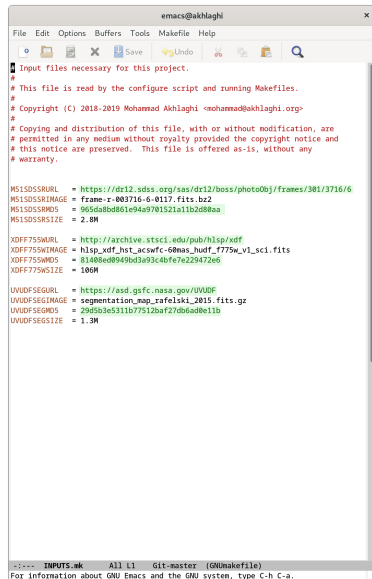
Stored information about each input file:

- ▶ **PID** (where available).
- ▶ Download **URL**.
- ▶ **MD5**-sum to check integrity.

All inputs are **downloaded** from the given PID/URL when necessary (during the analysis).

MD5-sums are **checked** to make sure the download was done properly or the file is the same (hasn't changed on the server/source).

Example from the reproducible paper [arXiv:1909.11230](https://arxiv.org/abs/1909.11230).  
This paper needs three input files (two images, one catalog).



```
# Input files necessary for this project.
# This file is read by the configure script and running Makefiles.
# Copyright (C) 2018-2019 Mohammad Akhlaghi <mohammad@akhlaghi.org>
# Copying and distribution of this file, with or without modification, are
# permitted in any medium without royalty provided the copyright notice and
# this notice are preserved. This file is offered as-is, without any
# warranty.

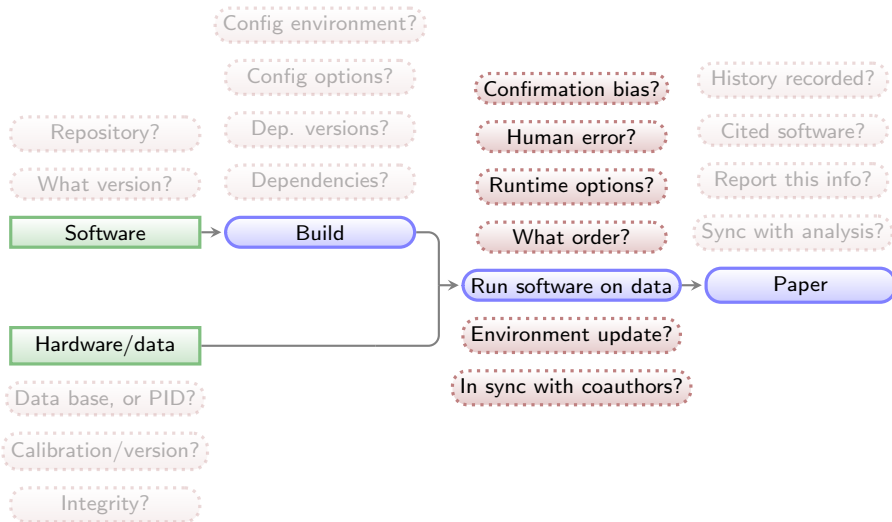
MS1S0SSRURL = https://dr12.sdss.org/sas/dr12/boos/photoObj/frames/301/3716/6
MS1S0SSRIMAGE = frame-r-003716-6-0117.fits.bz2
MS1S0SSRMDS = 965da8bd861e94a9701521a11b2d88aa
MS1S0SSRSIZE = 2.8M

XDFF75SWURL = http://archive.stsci.edu/pub/hlsp/xdff
XDFF75SWIMAGE = hlsp_xdff_hst_acswfc-60mas_hudf_f775w_v1_sc1.fits
XDFF75SWMDS = 81408ed0949bd3a93c4bfe7e229472e6
XDFF75WSIZE = 106M

UVUDFSEGURL = https://asd.gsfc.nasa.gov/UVUDF
UVUDFSEGIMAGE = segmentation_map_rafelski_2015.fits.gz
UVUDFSEGMDS = 29d5b3e5311b77512ba727db6ad0e11b
UVUDFSEGSIZE = 1.3M

-:- INPUTS.mk All L1 Git-master (GNUmakefile)
For information about GNU Emacs and the GNU system, type C-h C-a.
```

## General outline of a project (after data collection)



Green boxes with sharp corners: *source*/input components/files.

Blue boxes with rounded corners: *built* components.

Red boxes with dashed borders: questions that must be clarified for each phase.

# Reproducible science: Maneage is managed through a Makefile

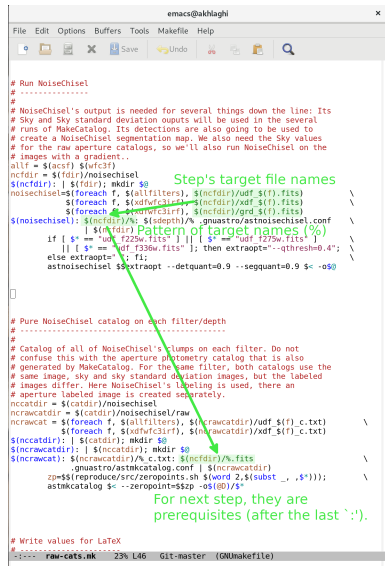
All steps (downloading and analysis) are managed by Makefiles (example from [zenodo.1164774](https://zenodo.org/record/1164774)):

- ▶ Unlike a script which always starts from the top, a Makefile **starts from the end** and steps that don't change will be left untouched (not remade).
- ▶ A single *rule* can **manage any number of files**.
- ▶ Make can identify independent steps internally and do them in **parallel**.
- ▶ Make was **designed for complex projects** with thousands of files (all major Unix-like components), so it is highly evolved and efficient.
- ▶ Make is a very **simple** and **small** language, thus easy to learn with great and free documentation (for example [GNU Make's manual](#)).

# Reproducible science: Manage is managed through a Makefile

All steps (downloading and analysis) are managed by Makefiles (example from [zenodo.1164774](https://zenodo.org/record/1164774)):

- ▶ Unlike a script which always starts from the top, a Makefile **starts from the end** and steps that don't change will be left untouched (not remade).
- ▶ A single *rule* can **manage any number of files**.
- ▶ Make can identify independent steps internally and do them in **parallel**.
- ▶ Make was **designed for complex projects** with thousands of files (all major Unix-like components), so it is highly evolved and efficient.
- ▶ Make is a very **simple** and **small** language, thus easy to learn with great and free documentation (for example [GNU Make's manual](#)).



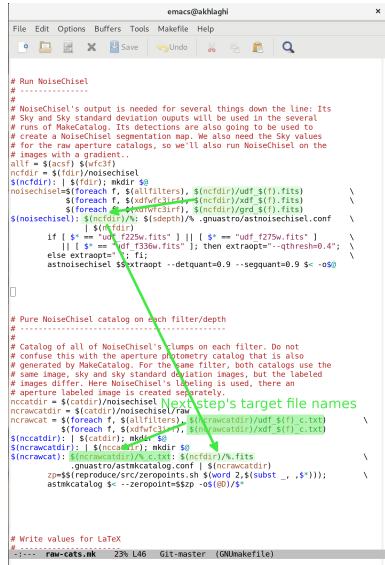
```
# Run NoiseChisel
# -----
#
# NoiseChisel's output is needed for several things down the line: Its
# Sky and Sky standard deviation outputs will be used in the several
# runs of MakeCatalog. Its detections are also going to be used to
# create a NoiseChisel segmentation map. We also need the Sky values
# for the raw aperture catalogs, so we'll also run NoiseChisel on the
# images with a gradient..
allf = $(acff) $(wfc3f)
ncfdir = $(fdir)/noisechisel
noisechisel: | $(fdir); mkdir $@
noisechisel=$(foreach f, $(allfilters), $(ncfdir)/udf $(f).fits) \
$(foreach f, $(xdfsfc3irf), $(ncfdir)/xdf $(f).fits) \
$(foreach f, $(xdfsfc3irf), $(ncfdir)/grd $(f).fits)
$(noisechisel): $(ncfdir)/%. $(sdepth)/%. gnuastro/astnoisechisel.conf \
| $(fdir) Pattern of target names (%)
if [ $* == "udf_f235w.fits" ]; then extraopt="--qthresh=0.4"; \
else extraopt=""; fi;
astnoisechisel $$extraopt --detquant=0.9 --segquant=0.9 $< -o$@

# Pure NoiseChisel catalog on each filter/depth
# -----
#
# Catalog of all of NoiseChisel's clumps on each filter. Do not
# confuse this with the aperture photometry catalog that is also
# generated by MakeCatalog. For the same filter, both catalogs use the
# same image, sky and sky standard deviation images, but the labeled
# images differ. Here NoiseChisel's labeling is used, there an
# aperture labeled image is created separately.
ncatdir = $(catdir)/noisechisel
ncrawcatdir = $(catdir)/noisechisel/raw
ncrawcat = $(foreach f, $(allfilters), $(ncrawcatdir)/udf $(f)_c.txt) \
$(foreach f, $(xdfsfc3irf), $(ncrawcatdir)/xdf $(f)_c.txt)
$(ncatdir): | $(catdir); mkdir $@
$(ncrawcatdir): | $(ncatdir); mkdir $@
$(ncrawcat): $(ncrawcatdir)/%. c.txt: $(ncfdir)/%.fits \
, gnuastro/astmkcatalog.conf | $(ncrawcatdir)
zp=$(reproduce/src/zeropoints.sh $(word 2,$(subst _,,$*)) \
astmkcatalog $< --zeropoint=$zp -o$(@)/$*
```

# Reproducible science: Maneage is managed through a Makefile

All steps (downloading and analysis) are managed by Makefiles (example from [zenodo.1164774](https://zenodo.org/record/1164774)):

- ▶ Unlike a script which always starts from the top, a Makefile **starts from the end** and steps that don't change will be left untouched (not remade).
- ▶ A single *rule* can **manage any number of files**.
- ▶ Make can identify independent steps internally and do them in **parallel**.
- ▶ Make was **designed for complex projects** with thousands of files (all major Unix-like components), so it is highly evolved and efficient.
- ▶ Make is a very **simple** and **small** language, thus easy to learn with great and free documentation (for example [GNU Make's manual](#)).

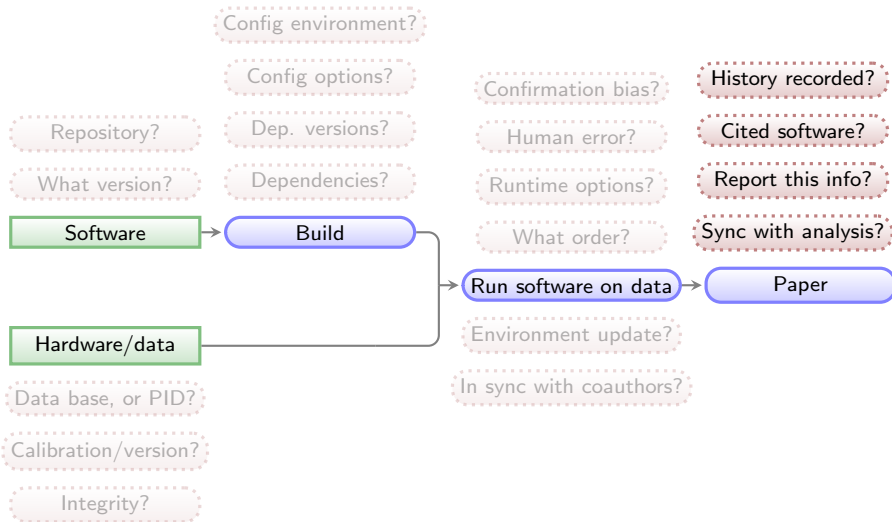


```
# Run NoiseChisel
# -----
#
# NoiseChisel's output is needed for several things down the line: Its
# Sky and Sky standard deviation outputs will be used in the several
# runs of MakeCatalog. Its detections are also going to be used to
# create a NoiseChisel segmentation map. We also need the Sky values
# for the raw aperture catalogs, so we'll also run NoiseChisel on the
# images with a gradient..
allf = $(acsf) $(wfc3f)
ncfdir = $(fdir)/noisechisel
$(ncfdir): | $(fdir); mkdir $@
noisechisel=$(foreach f, $(allfilters), $(ncfdir)/udf $(f).fits) \
$(foreach f, $(xdfs3irf), $(ncfdir)/xdf $(f).fits) \
$(foreach f, $(xdfs3irf), $(ncfdir)/grd $(f).fits)
$(noisechisel): $(ncfdir)/%. $(sdepth)/%. gnuastro/astnoisechisel.conf \
| $(fdir)
if [ $* == "udf_f225w.fits" ] || [ $* == "udf_f275w.fits" ] \
|| [ $* == "udf_f336w.fits" ]; then extraopt="--qthresh=0.4"; \
else extraopt=""; fi;
astnoisechisel $$extraopt --detquant=0.9 --segquant=0.9 $< -o$@

# Pure NoiseChisel catalog on each filter/depth
# -----
#
# Catalog of all of NoiseChisel's clumps on each filter. Do not
# confuse this with the aperture photometry catalog that is also
# generated by MakeCatalog. For the same filter, both catalogs use the
# same image, sky and sky standard deviation images, but the labeled
# images differ. Here NoiseChisel's labeling is used, there an
# aperture labeled image is created separately.
ncatdir = $(catdir)/noisechisel
ncrcatdir = $(catdir)/noisechisel/raw
ncrcat = $(foreach f, $(allfilters), $(ncrcatdir)/udf $(f)_c.txt) \
$(foreach f, $(xdfs3irf), $(ncrcatdir)/xdf $(f)_c.txt)
$(ncatdir): | $(catdir); mkdir $@
$(ncrcatdir): | $(ncatdir); mkdir $@
$(ncrcat): $(ncrcatdir)/%. c.txt: $(ncfdir)/%.fits \
, gnuastro/astnmkcatalog.conf | $(ncrcatdir)
zp=$(reproduce/src/zeropoints.sh $(word 2,$(subst _,, $*))); \
astnmkcatalog $< --zeropoint=$zp -o$(@D)/$*

# Write values for LaTeX
# -----
raw-cats.mk 23% L46 Git-master (GNUmakefile)
```

## General outline of a project (after data collection)



Green boxes with sharp corners: *source*/input components/files.

Blue boxes with rounded corners: *built* components.

Red boxes with dashed borders: questions that must be clarified for each phase.



## Values in final report/paper

All analysis **results** (numbers, plots, tables) written in paper's PDF as **L<sup>A</sup>T<sub>E</sub>X macros**. They are thus **updated automatically** on any change.

Shown here is a portion of the NoiseChisel paper and its L<sup>A</sup>T<sub>E</sub>X source ([arXiv:1505.01664](https://arxiv.org/abs/1505.01664)).

```
\begin{equation}
  \label{tSNeg}
  \mathrm{S/N}_{\mathrm{T}} = \frac{NF - NS_a}{\sqrt{NF + N\sigma_s^2}}
  = \frac{\sqrt{N}(F - S_a)}{\sqrt{F + \sigma_s^2}}.
\end{equation}
```

\noindent

See Section `\ref{SNegmodif}` for the modifications required when the input image is not in units of counts or has already been Sky subtracted. The distribution of `\small S/N`<sub>T</sub> from the objects in `$R_s$` for the three examples in Figure `\ref{dettf}` can be seen in column 5 (top) of that figure. Image processing effects, mainly due to shifting, rotating, and re-sampling the images for co-adding, on the real data further increase the size and count, and hence, the `\small S/N` of false detections in real, reduced/co-added images. A comparison of scales on the `\small S/N` histograms between the mock ((a.5.1) and (b.5.1)) and real (c.5.1) examples in Figure `\ref{dettf}` shows the effect quantitatively. In the histograms of Figure `\ref{dettf}`, the bin with the largest number of false pseudo-detections respectively has an `\small S/N` of `\onelargedettfmax`, `\sensitivedettfmax`, and `\fourdettfmax`.<sup>□</sup>

smaller than `--detsnminarea` are removed from the analysis in both  $R_s$  and  $R_d$ . In the examples in this section, it is set to 15. Note that since a threshold approximately equal to the Sky value is used, this is a very weak constraint. For each pseudo-detection,  $S/N_T$  can be written as,

$$S/N_T = \frac{NF - NS_a}{\sqrt{NF + N\sigma_s^2}} = \frac{\sqrt{N}(F - S_a)}{\sqrt{F + \sigma_s^2}}. \quad (3)$$

See Section 3.3 for the modifications required when the input image is not in units of counts or has already been Sky subtracted. The distribution of  $S/N_T$  from the objects in  $R_s$  for the three examples in Figure 7 can be seen in column 5 (top) of that figure. Image processing effects, mainly due to shifting, rotating, and re-sampling the images for co-adding, on the real data further increase the size and count, and hence, the  $S/N$  of false detections in real, reduced/co-added images. A comparison of scales on the  $S/N$  histograms between the mock ((a.5.1) and (b.5.1)) and real (c.5.1) examples in Figure 7 shows the effect quantitatively. In the histograms of Figure 7, the bin with the largest number of false pseudo-detections respectively has an  $S/N$  of 1.89, 2.37, and 4.77.

The  $S/N_T$  distribution of detections in  $R_s$  provides a very ro-

## Values in final report/paper

All analysis **results** (numbers, plots, tables) written in paper's PDF as **L<sup>A</sup>T<sub>E</sub>X macros**. They are thus **updated automatically** on any change.

Shown here is a portion of the NoiseChisel paper and its L<sup>A</sup>T<sub>E</sub>X source ([arXiv:1505.01664](https://arxiv.org/abs/1505.01664)).

```
\begin{equation}
  \label{tSNeg}
  \mathrm{S/N}_{\mathrm{T}} = \frac{NF - NS_a}{\sqrt{NF + N\sigma_s^2}}
  = \frac{\sqrt{N}(F - S_a)}{\sqrt{F + \sigma_s^2}}.
\end{equation}
```

\noindent

See Section `\ref{SNegmodif}` for the modifications required when the input image is not in units of counts or has already been Sky subtracted. The distribution of `\small S/N`<sub>T</sub> from the objects in `$R_s$` for the three examples in Figure `\ref{dettf}` can be seen in column 5 (top) of that figure. Image processing effects, mainly due to shifting, rotating, and re-sampling the images for co-adding, on the real data further increase the size and count, and hence, the `\small S/N` of false detections in real, reduced/co-added images. A comparison of scales on the `\small S/N` histograms between the mock ((a.5.1) and (b.5.1)) and real (c.5.1) examples in Figure `\ref{dettf}` shows the effect quantitatively. In the histograms of Figure `\ref{dettf}`, the bin with the largest number of false pseudo-detections respectively has an `\small S/N` of `\onelargedettfmax$`, `\sensitivitycdettfmax$`, and `\fourdettfmax$`.<sup>□</sup>

smaller than `--detsnminarea` are removed from the analysis in both  $R_s$  and  $R_d$ . In the examples in this section, it is set to 15. Note that since a threshold approximately equal to the Sky value is used, this is a very weak constraint. For each pseudo-detection,  $S/N_T$  can be written as,

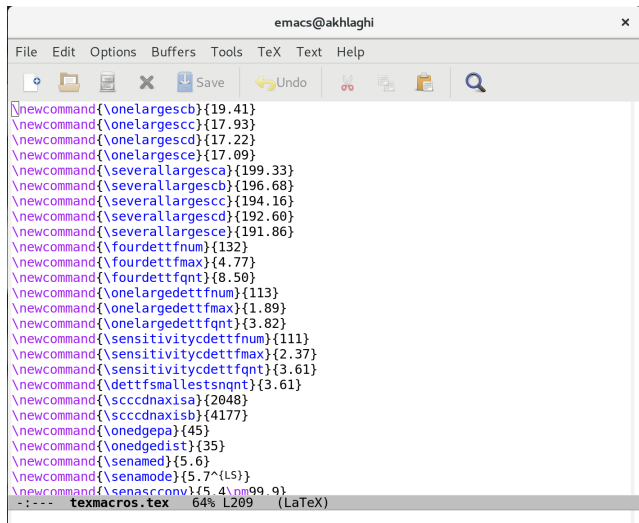
$$S/N_T = \frac{NF - NS_a}{\sqrt{NF + N\sigma_s^2}} = \frac{\sqrt{N}(F - S_a)}{\sqrt{F + \sigma_s^2}}. \quad (3)$$

See Section 3.3 for the modifications required when the input image is not in units of counts or has already been Sky subtracted. The distribution of  $S/N_T$  from the objects in  $R_s$  for the three examples in Figure 7 can be seen in column 5 (top) of that figure. Image processing effects, mainly due to shifting, rotating, and re-sampling the images for co-adding, on the real data further increase the size and count, and hence, the  $S/N$  of false detections in real, reduced/co-added images. A comparison of scales on the  $S/N$  histograms between the mock ((a.5.1) and (b.5.1)) and real (c.5.1) examples in Figure 7 shows the effect quantitatively. In the histograms of Figure 7, the bin with the largest number of false pseudo-detections respectively has an  $S/N$  of 1.89, 2.37, and 4.77.

The  $S/N_T$  distribution of detections in  $R_s$  provides a very ro-

Analysis step results/values concatenated into a single file.

All  $\text{\LaTeX}$  macros come from a **single file**.



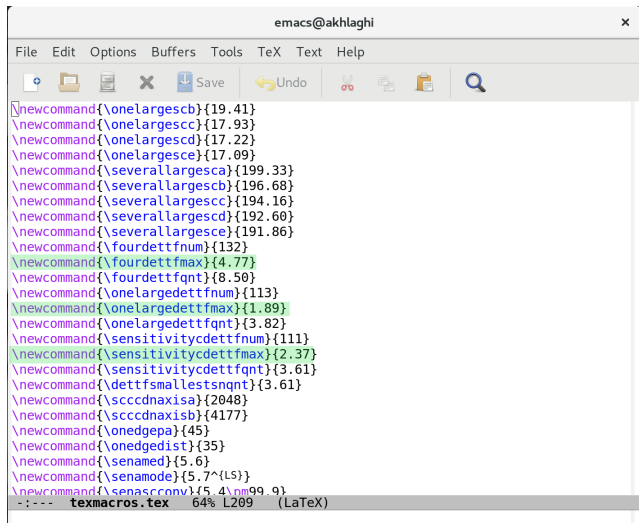
The screenshot shows an Emacs editor window titled "emacs@akhlaghi". The menu bar includes "File", "Edit", "Options", "Buffers", "Tools", "TeX", "Text", and "Help". The toolbar contains icons for opening a file, saving, undo, redo, and search. The main text area displays a list of LaTeX macro definitions, each starting with `\newcommand`. The macros are defined with a name and a numerical value in curly braces. The status bar at the bottom shows the file name "texmacros.tex", the cursor position "64%", and the page number "L209" followed by "(LaTeX)".

```
\newcommand{\onelargescb}{19.41}
\newcommand{\onelargescb}{17.93}
\newcommand{\onelargescd}{17.22}
\newcommand{\onelargescd}{17.09}
\newcommand{\severallargescb}{199.33}
\newcommand{\severallargescb}{196.68}
\newcommand{\severallargescb}{194.16}
\newcommand{\severallargescd}{192.60}
\newcommand{\severallargescd}{191.86}
\newcommand{\fourdettfnm}{132}
\newcommand{\fourdettfmax}{4.77}
\newcommand{\fourdettfqnt}{8.50}
\newcommand{\onelargedettfnm}{113}
\newcommand{\onelargedettfmax}{1.89}
\newcommand{\onelargedettfqnt}{3.82}
\newcommand{\sensitivitycdettfnm}{111}
\newcommand{\sensitivitycdettfmax}{2.37}
\newcommand{\sensitivitycdettfqnt}{3.61}
\newcommand{\dettfsmallestsnqnt}{3.61}
\newcommand{\scccdnaxisa}{2048}
\newcommand{\scccdnaxisb}{4177}
\newcommand{\onedgepa}{45}
\newcommand{\onedgedist}{35}
\newcommand{\senamed}{5.6}
\newcommand{\senamode}{5.7^{LS}}
\newcommand{\senascconv}{5.4^{nm}99.9}
```

--- texmacros.tex 64% L209 (LaTeX)

Analysis step results/values concatenated into a single file.

All  $\text{\LaTeX}$  macros come from a **single file**.



The screenshot shows an Emacs editor window titled "emacs@akhlaghi". The menu bar includes "File", "Edit", "Options", "Buffers", "Tools", "TeX", "Text", and "Help". The toolbar contains icons for opening a file, saving, undo, redo, and search. The main text area displays a list of LaTeX macro definitions, each starting with `\newcommand`. The macros are defined with a name and a value in curly braces. The values are numerical or symbolic expressions. The status bar at the bottom shows the file name "texmacros.tex", the cursor position "64%", and the page number "L209".

```
\newcommand{\onelargescb}{19.41}
\newcommand{\onelargescd}{17.93}
\newcommand{\onelargescd}{17.22}
\newcommand{\onelargescd}{17.09}
\newcommand{\severallargescb}{199.33}
\newcommand{\severallargescb}{196.68}
\newcommand{\severallargescd}{194.16}
\newcommand{\severallargescd}{192.60}
\newcommand{\severallargescd}{191.86}
\newcommand{\fourdettfnum}{132}
\newcommand{\fourdettfmax}{4.77}
\newcommand{\fourdettfqnt}{8.50}
\newcommand{\onelargedettfnum}{113}
\newcommand{\onelargedettfmax}{1.89}
\newcommand{\onelargedettfqnt}{3.82}
\newcommand{\sensitivitycdettfnum}{111}
\newcommand{\sensitivitycdettfmax}{2.37}
\newcommand{\sensitivitycdettfqnt}{3.61}
\newcommand{\dettfsmallestsnqnt}{3.61}
\newcommand{\scccdnaxisa}{2048}
\newcommand{\scccdnaxisb}{4177}
\newcommand{\onedgepa}{45}
\newcommand{\onedgedist}{35}
\newcommand{\senamed}{5.6}
\newcommand{\senamode}{5.7^{LS}}
\newcommand{\senascconv}{5.4^{nm}99.9}
-:--- texmacros.tex 64% L209 (LaTeX)
```

## Analysis results stored as $\text{\LaTeX}$ macros

The analysis scripts write/update the  $\text{\LaTeX}$  macro values automatically.

```
# Numbers for dettf.tex:
sqnt=9999999
function dettfhist
{
  # Set the file name.
  if [ $2 == 4 ]; then          obase=four;
  elif [ $2 = sensitivity3 ]; then obase=sensitivityc;
  else                          obase=$2;
  fi
  if [ $2 == onelarge ]; then ind="_7"; else ind="_12"; fi
  name=$1$2$ind"_detsn"$txt

  dettfnum=$(awk '/points binned in/{print $4; exit(0)}' $name)
  dettfqnt=$(awk '/quantile has a value of/{
    printf("%.2f", $9); exit(0);}' $name)
  dettfmax=$(awk 'BEGIN { max=-999999 }
    !/^#/ { if($2>max){max=$2; mv=$1} }
    END { printf("%.2f", mv) }' $name)
  addtexmacro $obase"dettfnum" $dettfnum
  addtexmacro $obase"dettfmax" $dettfmax
  addtexmacro $obase"dettfqnt" $dettfqnt

  # Find the smallest S/N quantile:
  sqnt=$(echo " " | awk '{if('$dettfqnt'<'$sqnt') print '$dettfqnt'}}')
}
for base in 4 onelarge sensitivity3
do dettfhist $texdir/dettf/ $base; done
addtexmacro dettfsmallestsqnt $sqnt
```

## Analysis results stored as $\text{\LaTeX}$ macros

The analysis scripts write/update the  $\text{\LaTeX}$  macro values automatically.

```
# Numbers for dettf.tex:
sqnt=9999999
function dettfhist
{
  # Set the file name.
  if [ $2 == 4 ]; then          obase=four;
  elif [ $2 = sensitivity3 ]; then obase=sensitivityc;
  else                          obase=$2;
  fi
  if [ $2 == onelarge ]; then ind="_7"; else ind="_12"; fi
  name=$1$2$ind"_detsn"$txt

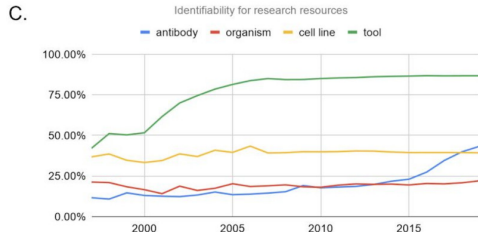
  dettfnum=$(awk '/points binned in/{print $4; exit(0)}' $name)
  dettfqnt=$(awk '/quantile has a value of/{
    printf("%.2f", $9); exit(0);}' $name)
  dettfmax=$(awk 'BEGIN { max=-999999 }
    !/^#/ { if($2>max){max=$2; mv=$1} }
    END { printf("%.2f", mv) }' $name)
  addtexmacro $obase"dettfnum" $dettfnum
  addtexmacro $obase"dettfmax" $dettfmax
  addtexmacro $obase"dettfqnt" $dettfqnt

  # Find the smallest S/N quantile:
  sqnt=$(echo " " | awk '{if('$dettfqnt'<'$sqnt') print '$dettfqnt'}}')
}
for base in 4 onelarge sensitivity3
do dettfhist $texdir/dettf/ $base; done
addtexmacro dettfsmallestsqnt $sqnt
```

Let's look at the data lineage to replicate Figure 1C (green/tool) of Menke+2020 (DOI:10.1101/2020.01.15.908111), as done in arXiv:2006.03018 for a demo.

## ORIGINAL PLOT

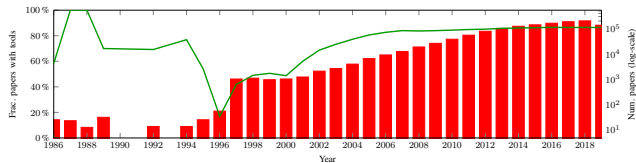
The Green plot shows the fraction of papers mentioning software tools from 1997 to 2019.



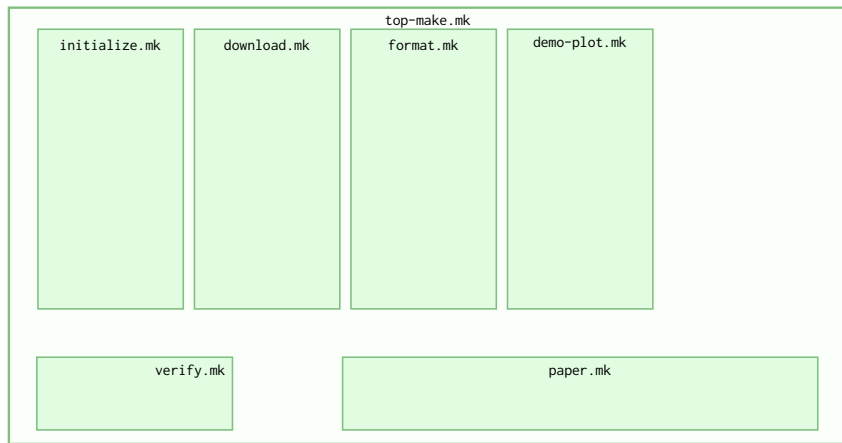
## OUR enhanced REPLICATION

The green line is same as above but over their full historical range.

Red histogram is the number of papers studied in each year



Makefiles (`.mk`) keep contextually separate parts of the project, all imported into `top-make.mk`



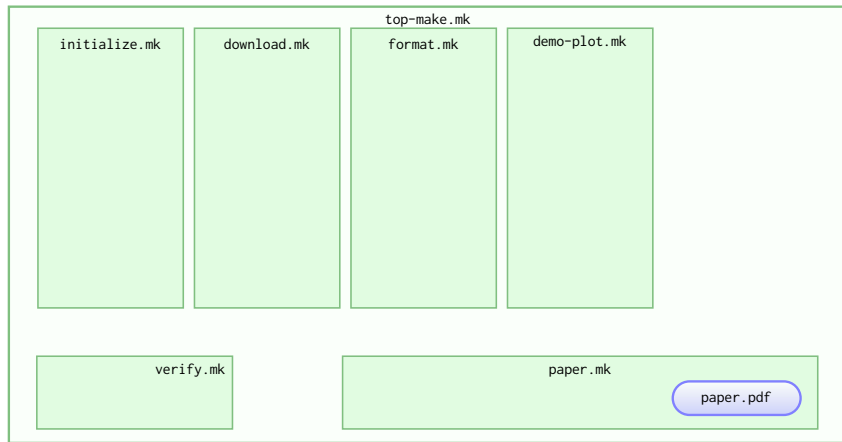
Green boxes with sharp corners: *source* files (hand written).

Blue boxes with rounded corners: *built* files (automatically generated),

built files are shown in the Makefile that contains their build instructions.



The ultimate purpose of the project is to produce a paper/report (in PDF).

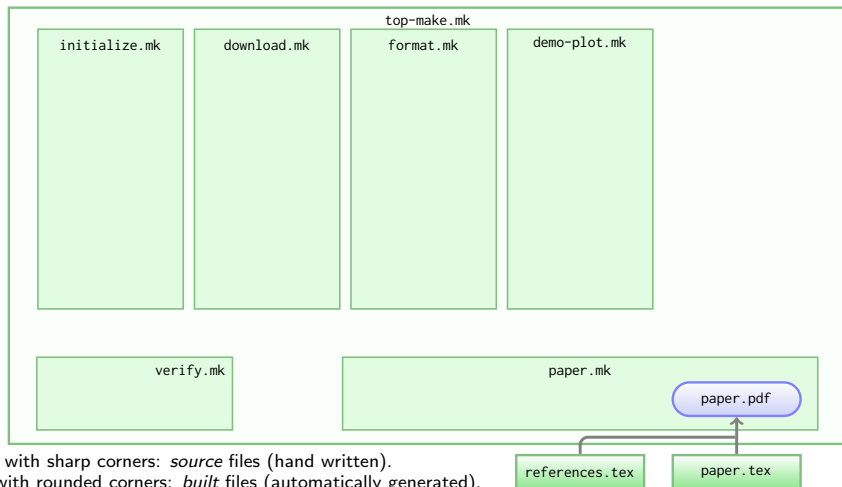


Green boxes with sharp corners: *source* files (hand written).

Blue boxes with rounded corners: *built* files (automatically generated),

built files are shown in the Makefile that contains their build instructions.

The narrative description, typography and references are in `paper.tex` & `references.tex`.

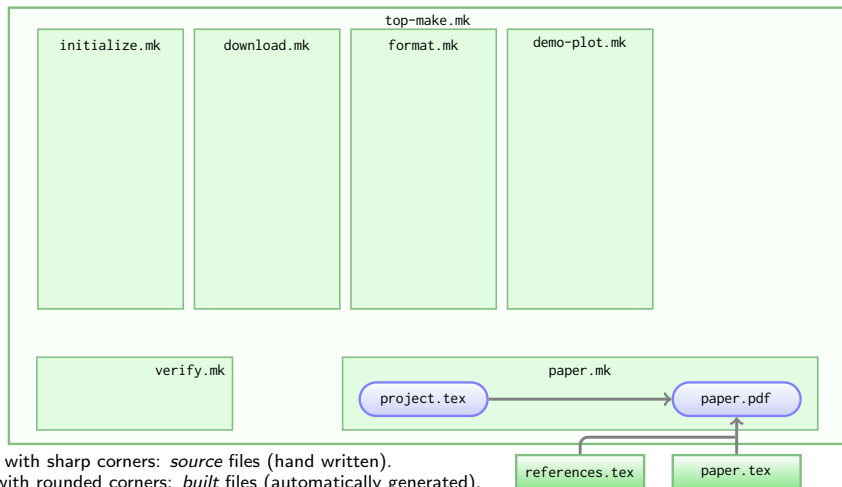


Green boxes with sharp corners: *source* files (hand written).

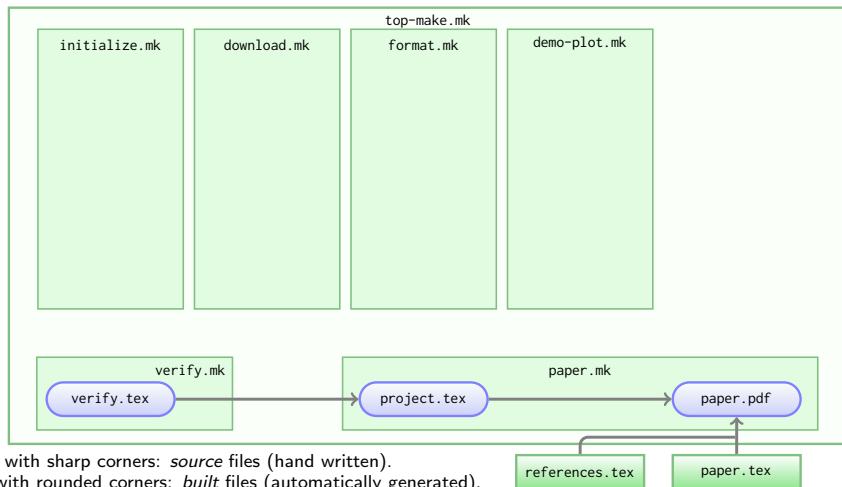
Blue boxes with rounded corners: *built* files (automatically generated),

built files are shown in the Makefile that contains their build instructions.

Analysis outputs (blended into the PDF as  $\LaTeX$  macros) come from `project.tex`.



But analysis outputs must first be *verified* (with checksums) before entering the report/paper.

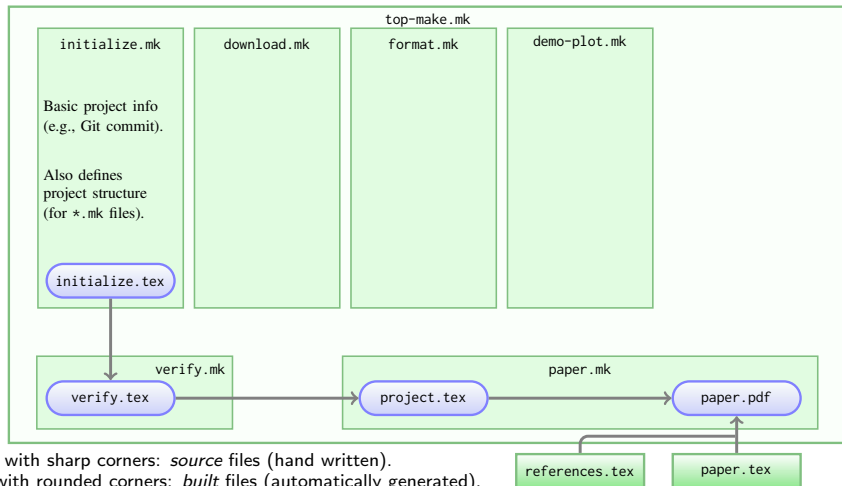


Green boxes with sharp corners: *source* files (hand written).

Blue boxes with rounded corners: *built* files (automatically generated),

built files are shown in the Makefile that contains their build instructions.

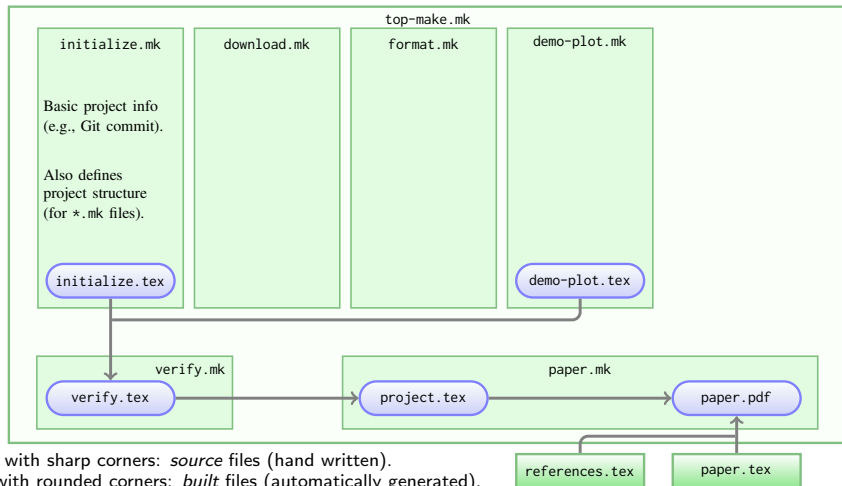
Basic project info comes from `initialize.tex`.



Green boxes with sharp corners: *source* files (hand written).

Blue boxes with rounded corners: *built* files (automatically generated),  
built files are shown in the Makefile that contains their build instructions.

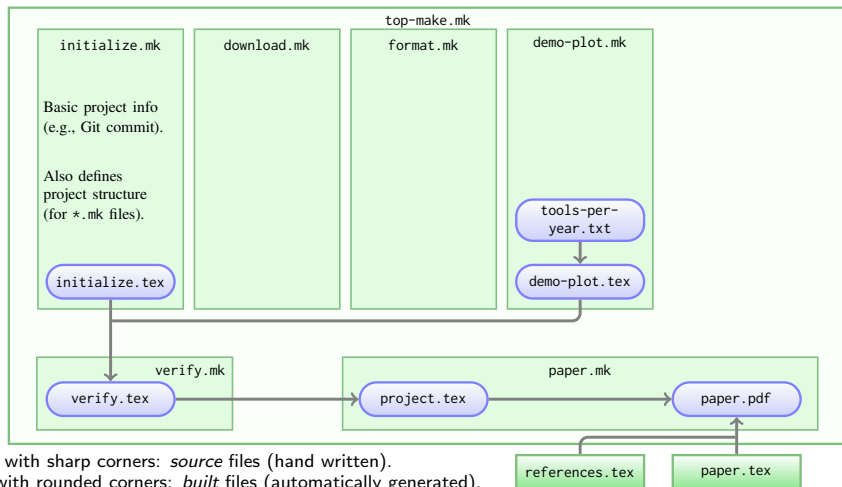
The paper includes some information about the plot.



Green boxes with sharp corners: *source* files (hand written).

Blue boxes with rounded corners: *built* files (automatically generated),  
built files are shown in the Makefile that contains their build instructions.

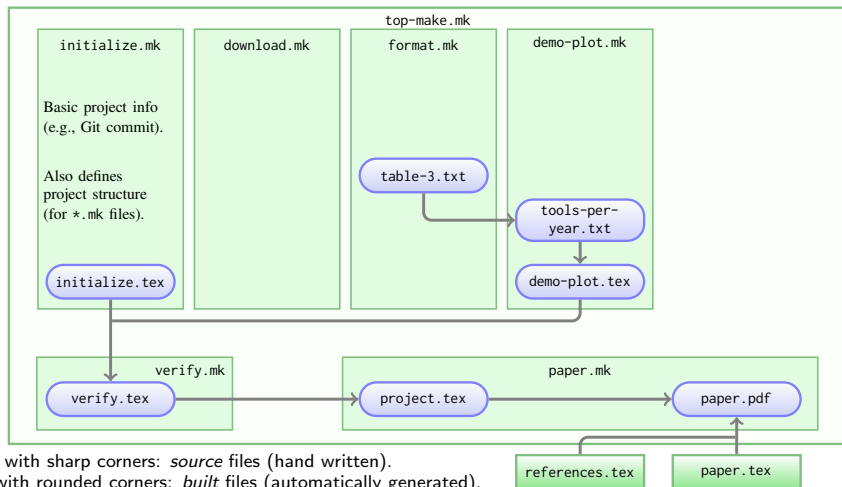
The final plotted data are calculated and stored in `tools-per-year.txt`.



Green boxes with sharp corners: *source* files (hand written).

Blue boxes with rounded corners: *built* files (automatically generated),  
built files are shown in the Makefile that contains their build instructions.

The plot's calculation is done on a formatted sub-set of the raw input data.

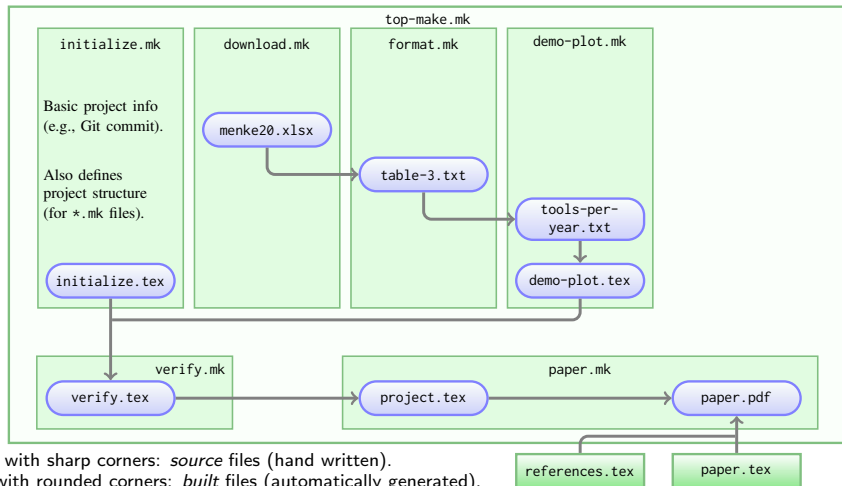


Green boxes with sharp corners: *source* files (hand written).

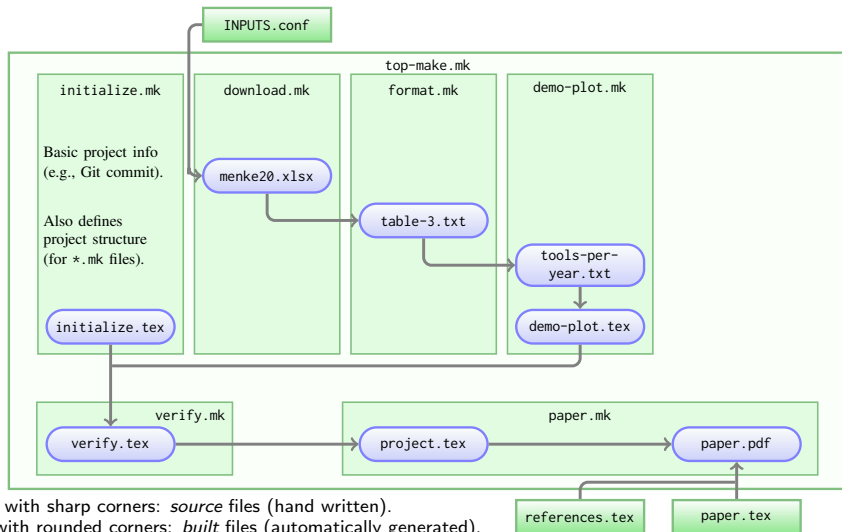
Blue boxes with rounded corners: *built* files (automatically generated),  
built files are shown in the Makefile that contains their build instructions.



The raw data that were downloaded are stored in XLSX format.



The download URL *and* a checksum to validate the raw inputs, are stored in `INPUTS.conf`.

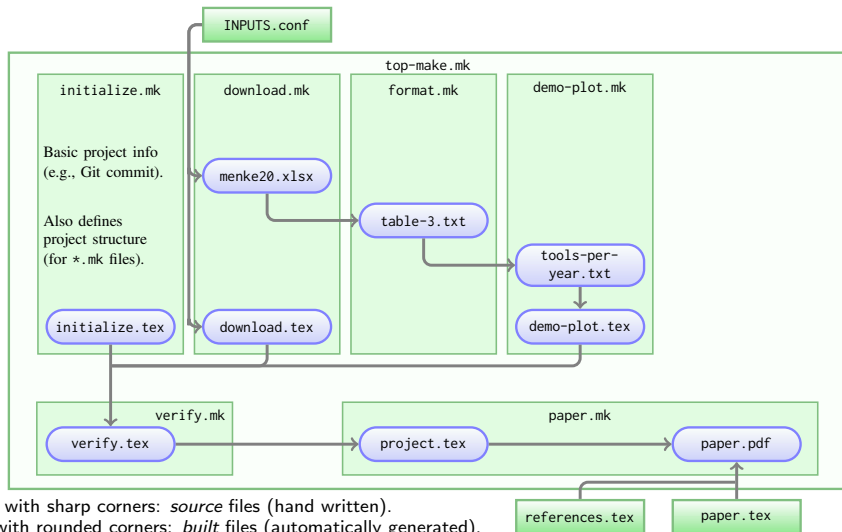


Green boxes with sharp corners: *source* files (hand written).

Blue boxes with rounded corners: *built* files (automatically generated),

built files are shown in the Makefile that contains their build instructions.

We also need to report the URL in the paper...

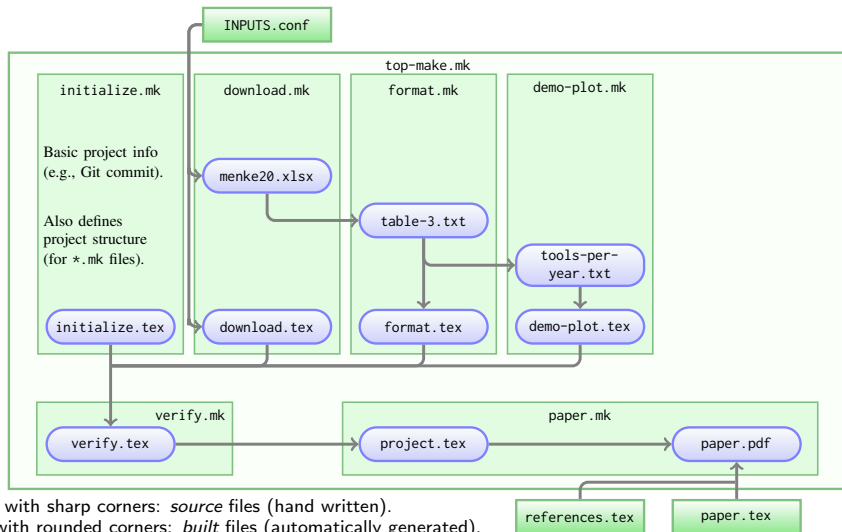


Green boxes with sharp corners: *source* files (hand written).

Blue boxes with rounded corners: *built* files (automatically generated),

built files are shown in the Makefile that contains their build instructions.

Some general info about the full dataset may also be reported.

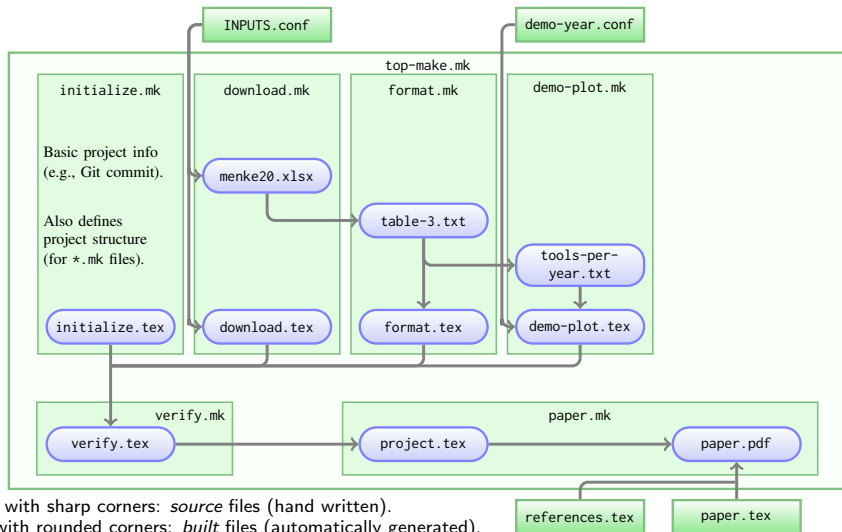


Green boxes with sharp corners: *source* files (hand written).

Blue boxes with rounded corners: *built* files (automatically generated),

built files are shown in the Makefile that contains their build instructions.

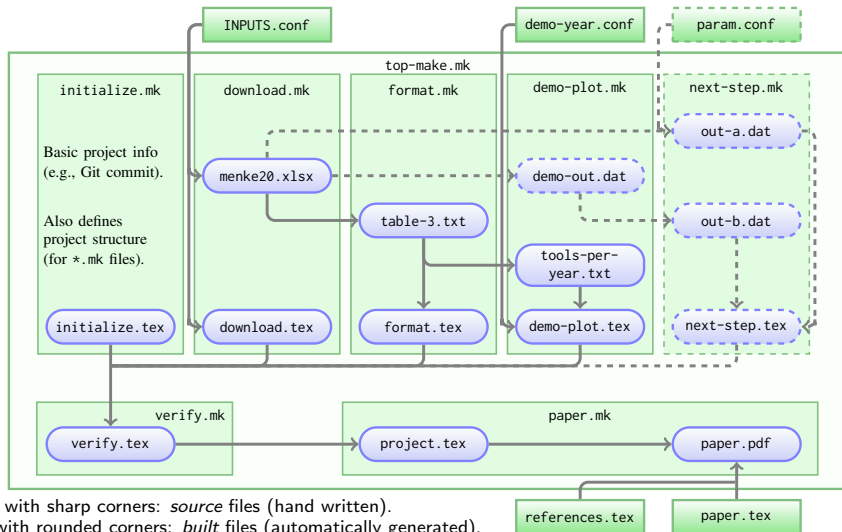
We report the number of papers studied in a special year, desired year is stored in `.conf` file.



Green boxes with sharp corners: *source* files (hand written).

Blue boxes with rounded corners: *built* files (automatically generated),  
built files are shown in the Makefile that contains their build instructions.

It is very easy to expand the project and add new analysis steps (this solution is scalable)

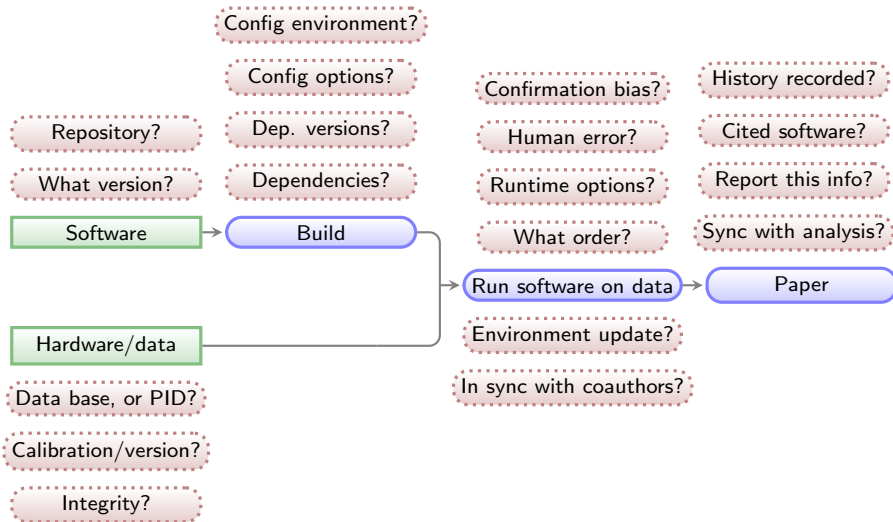


Green boxes with sharp corners: *source* files (hand written).

Blue boxes with rounded corners: *built* files (automatically generated),

built files are shown in the Makefile that contains their build instructions.

All questions have an answer now (in **plain text**: human & computer readable/archivable).

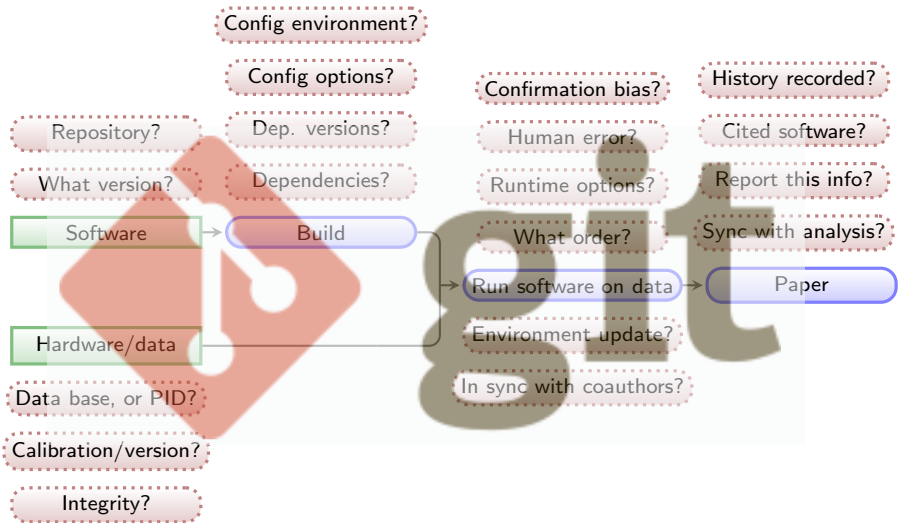


Green boxes with sharp corners: *source*/input components/files.

Blue boxes with rounded corners: *built* components.

Red boxes with dashed borders: questions that must be clarified for each phase.

All questions have an answer now (in plain text: so we can use Git to keep its history).



Green boxes with sharp corners: *source*/input components/files.  
Blue boxes with rounded corners: *built* components.  
Red boxes with dashed borders: questions that must be clarified for each phase.



## New projects branch from Maneage

- ▶ The project (answers to questions above) will evolve.



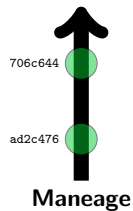
## New projects branch from Maneage

- The project (answers to questions above) will evolve.



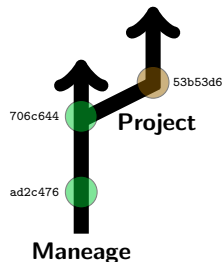
## New projects branch from Maneage

- ▶ Each point of project's history is recorded with Git.

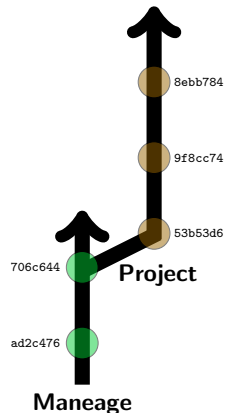


## New projects branch from Maneage

- ▶ Each point of project's history is recorded with Git.
- ▶ New project: a branch from the template.  
Recall that **every commit** contains the following:
  - ▶ Instructions to download, verify and build **software**.
  - ▶ Instructions to download and verify input **data**.
  - ▶ Instructions to run software on data (do the **analysis**).
  - ▶ Narrative description of project's purpose/**context**.

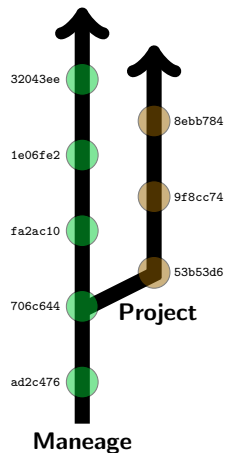


## New projects branch from Maneage



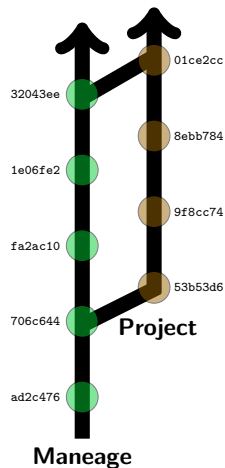
- ▶ Each point of project's history is recorded with Git.
- ▶ New project: a branch from the template.  
Recall that **every commit** contains the following:
  - ▶ Instructions to download, verify and build **software**.
  - ▶ Instructions to download and verify input **data**.
  - ▶ Instructions to run software on data (do the **analysis**).
  - ▶ Narrative description of project's purpose/**context**.
- ▶ Research progresses in the project branch.

## New projects branch from Maneage



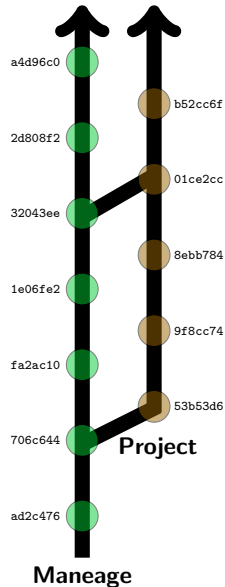
- ▶ Each point of project's history is recorded with Git.
- ▶ New project: a branch from the template.  
Recall that **every commit** contains the following:
  - ▶ Instructions to download, verify and build **software**.
  - ▶ Instructions to download and verify input **data**.
  - ▶ Instructions to run software on data (do the **analysis**).
  - ▶ Narrative description of project's purpose/**context**.
- ▶ Research progresses in the project branch.
- ▶ Template will evolve (improved infrastructure).

## New projects branch from Maneage



- ▶ Each point of project's history is recorded with Git.
- ▶ New project: a branch from the template.  
Recall that **every commit** contains the following:
  - ▶ Instructions to download, verify and build **software**.
  - ▶ Instructions to download and verify input **data**.
  - ▶ Instructions to run software on data (do the **analysis**).
  - ▶ Narrative description of project's purpose/**context**.
- ▶ Research progresses in the project branch.
- ▶ Template will evolve (improved infrastructure).
- ▶ Template can be imported/merged back into project.

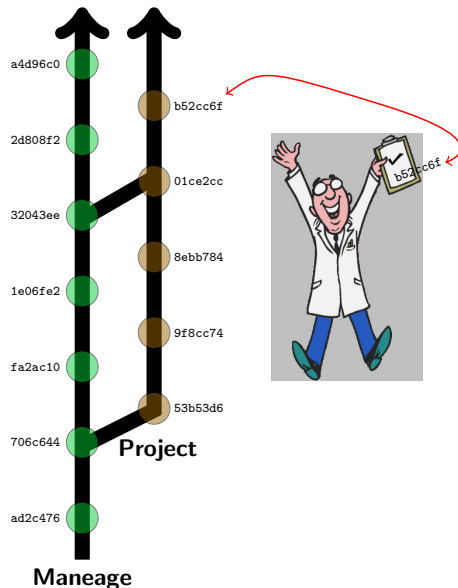
## New projects branch from Maneage



- ▶ Each point of project's history is recorded with Git.
- ▶ New project: a branch from the template.  
Recall that **every commit** contains the following:
  - ▶ Instructions to download, verify and build **software**.
  - ▶ Instructions to download and verify input **data**.
  - ▶ Instructions to run software on data (do the **analysis**).
  - ▶ Narrative description of project's purpose/**context**.
- ▶ Research progresses in the project branch.
- ▶ Template will evolve (improved infrastructure).
- ▶ Template can be imported/merged back into project.
- ▶ The template and project will **evolve**.
- ▶ During research this **encourages creative tests** (previous research states can easily be retrieved).
- ▶ **Coauthors** can work on same project in parallel (separate project branches).

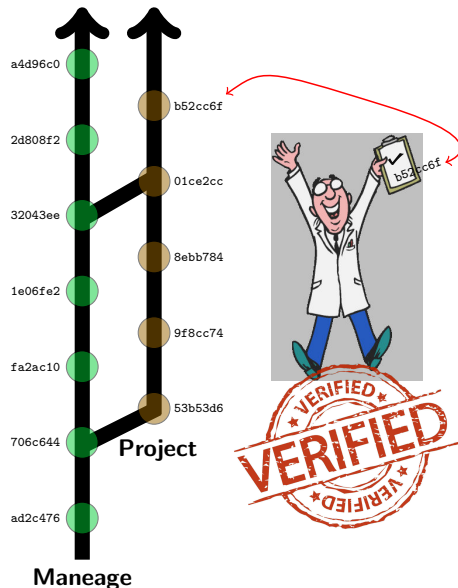


## New projects branch from Maneage



- ▶ Each point of project's history is recorded with Git.
- ▶ New project: a branch from the template.  
Recall that **every commit** contains the following:
  - ▶ Instructions to download, verify and build **software**.
  - ▶ Instructions to download and verify input **data**.
  - ▶ Instructions to run software on data (do the **analysis**).
  - ▶ Narrative description of project's purpose/**context**.
- ▶ Research progresses in the project branch.
- ▶ Template will evolve (improved infrastructure).
- ▶ Template can be imported/merged back into project.
- ▶ The template and project will **evolve**.
- ▶ During research this **encourages creative tests** (previous research states can easily be retrieved).
- ▶ **Coauthors** can work on same project in parallel (separate project branches).
- ▶ Upon **publication**, the **Git checksum** is enough to verify the integrity of the result.

## New projects branch from Maneage



- ▶ Each point of project's history is recorded with Git.
- ▶ New project: a branch from the template.  
Recall that **every commit** contains the following:
  - ▶ Instructions to download, verify and build **software**.
  - ▶ Instructions to download and verify input **data**.
  - ▶ Instructions to run software on data (do the **analysis**).
  - ▶ Narrative description of project's purpose/**context**.
- ▶ Research progresses in the project branch.
- ▶ Template will evolve (improved infrastructure).
- ▶ Template can be imported/merged back into project.
- ▶ The template and project will **evolve**.
- ▶ During research this **encourages creative tests** (previous research states can easily be retrieved).
- ▶ **Coauthors** can work on same project in parallel (separate project branches).
- ▶ Upon **publication**, the **Git checksum** is enough to verify the integrity of the result.



# Two recent examples (publishing Git checksum in abstract)

arXiv:1909.11230v1 [astro-ph.IM] 24 Sep 2019

The Results of the Low-Surface-Brightness Universe  
Proceedings IAU Symposium No. 355, 2019  
D. Valle-Gabaud, J. Trujillo & S. Okamoto, eds.

© 2019 International Astronomical Union  
DOI: 00.0000/X000000000000000X

## Carving out the low surface brightness universe with NoiseChisel

Mohammad Akhlaghi<sup>1,2</sup>

<sup>1</sup>Instituto de Astrofísica de Canarias, C/ Vía Láctea, 38200 La Laguna, Tenerife, Spain.  
email: mohammad@akhlaghi.org

<sup>2</sup>Facultad de Física, Universidad de La Laguna, Avda. Astrofísico Fco. Sánchez s/n, 38200 La Laguna, Tenerife, Spain.

**Abstract.** NoiseChisel is a program to detect very low signal-to-noise ratio (S/N) features with minimal assumptions on their morphology. It was introduced in 2015 and released within a collection of data analysis programs and libraries known as GNU Astronomy Utilities (Gnuastro). Over the last ten stable releases of Gnuastro, NoiseChisel has significantly improved: detecting even fainter signal, enabling better user control over its inner workings, and many bug fixes. The most important change may be that NoiseChisel's segmentation features have been moved into a new program called Segment. Another major change is the final growth strategy of its true detections, for example NoiseChisel is able to detect the outer wings of M51 down to S/N of 0.25, or 28.27 mag/arcsec<sup>2</sup> on a single-exposure SDSS image (r-band). Segment is also able to detect the localized HII regions as “clumps”, which are non-spheroidal. Finally, to orchestrate a controlled analysis, the concept of a “reproducible analysis” is proposed: this paper itself is exactly reproducible (snapshot v4.4.0-g8505dc6).

**Keywords.** galaxies: halos, galaxies: photometry, galaxies: structure, methods: data analysis, methods: reproducible, techniques: image processing, techniques: photometric

### 1. Introduction

Signal from the low surface brightness universe is buried deep in the datasets noise and thus requires accurate detection methods. In Akhlaghi and Ichikawa (2015) (henceforth AI15) a new method was introduced to detect such very low signal-to-noise ratio (S/N) signal from the images in a non-parametric manner. It allows accurate detection of the diffuse outer features of galaxies (that often have a different morphology from the centers). The software implementation of this method (NoiseChisel) is released as part of a larger collection of data analysis software known as GNU Astronomy Utilities (Gnuastro). It was the first professional astronomical software to be independently refereed by an independent panel (GNU Evaluation committee) and fully conforms with the GNU Coding Standards<sup>†</sup>.

Since its release, NoiseChisel has been used in many studies. For example Bacon et al. (2017) used it to identify objects that were missed by Rafelski et al. (2015) (henceforth R15), who used a combination of six SExtractor (Bertin and Arnouts 1996) runs with different configurations to avoid deblending problems, but still missed many sources with significant signal, see Figure 1. Borlaff et al. (2019), Miller et al. (2019), and Trujillo et al. (2019) used it for accurate flat field and Sky subtraction to create deeper co-added images in galaxy fields for optimal detection of the low surface brightness features. Calvi et al. (2019) used it to find Lyman- $\alpha$  emitters in spectra. For future studies, Laine et al.

<sup>†</sup> <https://www.gnu.org/s/gnuastro>  
<sup>‡</sup> <https://www.gnu.org/prep/standards>

### Monthly Notices

mnras.000.0000-0000000000000000

MNRAS **00**, 000–000 (2020)  
Advance Access publication 2019 November 14

doi:10.1093/mnras/mtz111

## The Sloan Digital Sky Survey extended point spread functions

Raúl Infante-Sainz<sup>1,2\*</sup>, Ignacio Trujillo<sup>0,1,2</sup> and Javier Román<sup>0,1,2,3</sup>

<sup>0</sup>Instituto de Astrofísica de Canarias, C/ Vía Láctea s/n, E-38205 La Laguna, Tenerife, Spain

<sup>1</sup>Departamento de Astrofísica, Universidad de La Laguna, E-38205 La Laguna, Tenerife, Spain

<sup>2</sup>Instituto de Astrofísica de Andalucía (CSIC), Glorieta de la Astronomía, E-18008 Granada, Spain

Accepted 2019 October 30. Received 2019 October 29; in original form 2019 September 10

### ABSTRACT

A robust and extended characterization of the point spread function (PSF) is crucial to extract the photometric information produced by deep imaging surveys. Here, we present the extended PSFs of the Sloan Digital Sky Survey (SDSS), one of the most productive astronomical surveys of all time. By stacking ~1000 images of individual stars with different brightness, we obtain the bidimensional SDSS PSFs extending over 9 arcmin in radius for all the SDSS filters (i.e.  $u$ ,  $r$ ,  $i$ ,  $z$ ). This new characterization of the SDSS PSFs is near a factor of 10 larger in extension than previous PSFs characterizations of the same survey. We found asymmetries in the shape of the PSFs caused by the drift scanning observing mode. The flux of the PSFs is larger along the drift scanning direction. Finally, we illustrate with an example how the PSF models can be used to remove the scattered light field produced by the brightest stars in the central region of the Coma cluster field. This particular example shows the huge importance of PSFs in the study of the low-surface brightness Universe, especially with the upcoming of ultradep surveys, such as the Large Synoptic Survey Telescope (LSST). Following a reproducible science philosophy, we make all the PSF models and the scripts used to do the analysis of this paper publicly available (snapshot v0.4.0-gd966ad0).

**Key words:** instrumentation: detectors – methods: data analysis – techniques: image processing – techniques: photometric – galaxies: halos.

### 1 INTRODUCTION

The point spread function (PSF) describes the response of an imaging system to the light produced by a point source. Real PSFs have complex structures as their shapes depend on the optical path that light takes as it travels through the atmosphere and multiple optical elements, mirrors, lenses, detectors, etc. For the vast majority of astronomical works, only a tiny portion of the PSF (i.e. essentially a few inner arcscales; see e.g. Trujillo et al. 2004a, b) is characterized. In practice, however, the light of both point and extended sources are spread over the entire detector due to the effect of the PSF at large radii. Therefore, it is necessary to have a good understanding of its structure along the entire detector (typically extending over arcminutes or more).

Extended PSFs have become a vital tool to obtain precise photometric information in modern astronomical surveys. For instance, Stare, Harding & Mibow (2009) modelled the extended PSF and the internal reflections produced by the stars of the Burrell Schmidt telescope and showed that virtually all the pixels of the image are dominated by the scattered light by both stars and galaxies at 20.5 mag/arcsec<sup>2</sup> (i-band; Trujillo & Pflanz 2016)

also characterized and used the extended PSF of the 10.4 m Gran Telescopio Canarias (GTC) telescope to model and remove the scattered light in ultradep observations of the UGC 00100 galaxy. Even more troublesome for low-surface brightness studies is the finding (see e.g. Trujillo & Baker 2013; Sandin 2014, 2015) that the outer regions of astronomical objects are severely affected by their own scattered light produced by the convolution with the PSF. In order to correct this effect, Karabal et al. (2017) generated the PSF models and the internal reflections from images of the Canada-France-Hawaii Telescope (CFHT) to deconvolve a sample of three galaxies and correct them from instrumental scattered light. More recently, Román, Trujillo & Montes (2019) characterized the PSFs of the Stripe 82 survey and used them to model and correct the scattered light field produced by stars to study the optical properties of the Galactic rim. All the above works have shown that having an extended PSF is crucial when accurate photometric and structure properties of astronomical objects at low-surface brightness levels are required.

One of the most commonly used surveys for measuring photometric properties of astronomical objects is the Sloan Sky Digital Survey (SDSS; York et al. 2000, covering 14 555 deg<sup>2</sup> on the sky (just over 35 per cent of the full sky) in five photometric bands (i.e.  $u$ ,  $r$ ,  $i$ , and  $z$ ). Although SDSS is a relatively shallow survey compared

\*E-mail: rinfante@gaia.es

## Publication of the project

A reproducible project using Maneage will have the following (**plain text**) components:

- ▶ Makefiles.
- ▶  $\LaTeX$  source files.
- ▶ Configuration files for software used in analysis.
- ▶ Scripts/programming files (e.g., Python, Shell, AWK, C).

The **volume** of the project's source will thus be **negligible** compared to a single figure in a paper (usually  $\sim 100$  kilo-bytes).

The project's pipeline (customized Maneage) can be **published** in

- ▶ **arXiv**: uploaded with the  $\LaTeX$  source to always stay with the paper (for example [arXiv:1505.01664](#) or [arXiv:2006.03018](#)).
- ▶ **Zenodo**: Along with all the input datasets (many Gigabytes) and software (for example [zenodo.3872248](#)) and given a unique DOI.
- ▶ **Software Heritage**: to archive the full version-controlled history of the project. (for example [swh:1:dir:33fea87068c1612daf011f161b97787b9a0df39fk](#))

## Executing a Maneaged project (for example arXiv:2006.03018)

```
$ git clone https://gitlab.com/makhlaghi/maneage-paper    # Import the project.
```

## Executing a Maneaged project (for example arXiv:2006.03018)

```
$ git clone https://gitlab.com/makhlaghi/maneage-paper    # Import the project.
```

```
$ ./project configure    # You will specify the build directory on your system,  
                          # and it will build all software (about 1.5 hours).
```

## Executing a Maneaged project (for example arXiv:2006.03018)

```
$ git clone https://gitlab.com/makhlaghi/maneage-paper    # Import the project.
```

```
$ ./project configure    # You will specify the build directory on your system,  
# and it will build all software (about 1.5 hours).
```

```
$ ./project make    # Does all the analysis and makes final PDF.
```



## Future prospects...

Adoption of reproducibility by many researchers will enable the following:

- ▶ A repository for education/training (PhD students, or researchers in other fields).
- ▶ Easy **verification/understanding** of other research projects (when necessary).
- ▶ Trivially **test** different steps of others' work (different configurations, software and etc).
- ▶ Science can progress **incrementally** (shorter papers actually building on each other!).
- ▶ **Extract meta-data** after the publication of a dataset (for future ontologies or vocabularies).
- ▶ Applying **machine learning** on reproducible research projects will allow us to solve some Big Data Challenges:
  - ▶ *Extract the relevant parameters automatically.*
  - ▶ *Translate the science to enormous samples.*
  - ▶ *Believe the results when no one will have time to reproduce.*
  - ▶ *Have confidence in results derived using machine learning or AI.*

Achievements: RDA adoption grant (2019) to IAC for Maneage



A wide, horizontal banner with a blue gradient background. In the center is a globe showing the Earth. The word 'HORIZON' is written in large, white, sans-serif capital letters on the left, and '2020' is on the right, separated by the central globe.

For Maneage, the **IAC** is selected as a **Top European organization** funded to adopt RDA Recommendations and Outputs.

- ▶ Research Data Alliance was launched by the **European Commission**, NSF, National Institute of Standards and Technology, and the Australian Government's Department of Innovation.
- ▶ RDA Outputs are the technical and social infrastructure solutions developed by RDA Working Groups or Interest Groups that enable data sharing, exchange, and interoperability.

# Achievements: “News and Views” in Nature Astronomy (DOI:10.1038/s41550-021-01402-3)

## REPRODUCIBILITY

### No expiration date

The short lifespan of software puts a time limit on the reproducibility of computational research. To extend software longevity, guidelines and tools to preserve scientific workflows and analysis are helpful, but the challenge is to get researchers to use them.

Michelle M. Kitzel

A software has an expiration date: there is no guarantee that you will be able to run today's popular software tool at some date in the future, or even that its stored data will be usable. This problem is annoying when the software is your favorite fitness app, but for more serious users for scientific software. Short software lifespans not only complicate the sustainability of computational research projects, but also impede the long-term reproducibility of the research results. The ability to demonstrate reproducibility — that is, that consistent results are obtained for the same input data, code and conditions of analysis — is ever more important. Given reports of a 'reproducibility crisis' in science<sup>1</sup>, computational scientists are increasingly expected to include data and software together with publications. In a recent issue of *Computing in Science and Engineering*, Michael and Abigail and colleagues propose a set of key guidelines to ensure extended software longevity and so the reproducibility of published research.

Ensuring software longevity (an extended time span wherein software is functional) is a key requirement for long-term computational reproducibility. Software longevity is particularly difficult to preserve for research workflows that link a series of software tools together to process data in an analysis pipeline. To ensure the longevity of an entire workflow, each component of the pipeline must be specified (including the precise version of the software and the packages it depends on) and maintained. The more steps (or tools) there are in a workflow, the broader the web of software dependencies and the more difficult it is to ensure the longevity, and thus reproducibility, of the pipeline as a whole. Worse, scientific software pipelines are continually increasing in complexity as a direct result of the rising complexity and volume of data produced by ever-developing instruments. For example, the development of the Square Kilometer Array (which will



Fig. 1 | A reproducible document. Manage supports reproducibility by tying precise software workflows to the data graphics in an academic paper, thus creating a dynamic document that links narrative, data and results.

be the largest radio telescope in the world) has necessitated concomitant efforts to design effective software pipelines to process the astronomical flood of data.

Unfortunately, software longevity is still often ignored by researchers: many are not formally trained in software management, which can be seen as an unnecessary burden, distraction and delay when their focus is on the 'real work' — producing scientific results. However, there are obvious benefits to ensuring that the next generation of astronomers is equipped with a basic understanding of computer science and current best practice in software development (such as modularization, version control, and so on). In the absence of such training, a single set of rules for researchers to follow can be very helpful. The requirements set out by Abigail et al. for scientific workflows are rigorous — the software components must be freely available, completely self-contained, minimally complex and modular, with simple command-line execution (because GUIs [graphical user interfaces] are quickly a source of central software dependencies and data storage primarily in plain text. The specifications for the data and the project documentation are similarly rigorous: the data inputs and outputs must be automatically verified, storage should be

primarily as text, the project must have version control and a fully documented history (including unsuccessful attempts) with any analysis (such as the generation of a plot for the paper) linked with a complete narrative.

Adhering to these guidelines would seem like a Herculean task, were it not for the increasing range of online tools to support reproducible research and software longevity. Data and software repositories (for example, *Conda*) provide version control and enable long-term accessibility. Visual image registries (for example, *DockerHub*) store images of a virtual machine or container, a useful method for distributing and preserving research software (excepting when images are impractically large). Frameworks such as *Occure* and others allow users to create, store and execute a complete scientific software workflow, including precise versions of all software and data used, the dependencies between them, and detailed information on the computational environment. At the highest level, tools for creating dynamic documents allow researchers to tie a workflow to the data analysis within a narrative. In this category, Abigail et al. announce their tool *Manage* (managing data lineage), a general-purpose reproducible research platform.

## news & views

Manage enables a precise software workflow to be stored and linked to a research paper, which then becomes a dynamic, reproducible document combining a scientific narrative with data and code (Fig. 1). This approach allows users to produce, and then archive, a paper together with its software and data, so that the computations necessary to regenerate the data behind the figures and tables in the paper can be reproduced at any time. Abigail et al. point out that it is best to plan for longevity when starting a software project, as trying to impose suitable rules afterwards is much trickier. To help with this, the *Manage* system provides users with a working template of a single paper containing an example calculation that they can then customize, with full version control. The *Manage* system has already been used to produce some research publications, including (pleasantly) the Abigail et al. paper itself.

Manage is a valuable proof-of-concept demonstrating that long-term reproducibility of scientific results is

achievable. However, the biggest difficulty here will be to get researchers to adhere to the authors' guidelines for longevity, which will seem rather laudible at first glance. The small set of tools used in *Manage* have been around for decades. Those with a training in computer science in the 1990s (or earlier) will have no problem in seeing the value of the robust Unix command line, the job manager *Make*, version control with *Git* and the typesetting language *LaTeX*, and computer science thus fulfils Dijkstra's would-be highly approved of their sensible recommendations to remove unnecessary 'bells and whistles'. However, to a younger generation of researchers raised on the digital luxuries of smart devices with GUIs and web software, these retrograde ideas will be a hard sell. A further stumbling block is that the *Manage* implementation is for Unix/Linux operating systems, and there are no clear alternatives for Windows provided. The authors claim that once researchers try the tools they will see the value of them, but it is getting them to try that will be

the hard task. Anyone convinced of the superiority of *LaTeX* with counters attached to Microsoft Office will have had a taste of the experience. Given eggs and hens, anyone?

Michelle M. Kitzel<sup>1</sup>  
<sup>1</sup>Department of Computer Science, University of Cape Town, Rondebosch, South Africa  
My email: m.kitzel@uct.ac.za

Published online: 21 June 2021  
<https://doi.org/10.1038/s41550-021-01402-3>

### References

1. Nature 575, 461–464 (2020).
2. Abigail, M. et al. *Comput. Sci. Eng.* 16, 40–42 (2021).
3. Kitzel, M. et al. *Nature* 594, 386–387 (2021).
4. Williams, C. *Front. Astron. Space Sci.* 7, 624647 (2020).
5. Dijkstra, E. W. *Communications of the ACM* 17, 773–774 (1974).
6. Dijkstra, E. W. *Communications of the ACM* 17, 773–774 (1974).
7. Dijkstra, E. W. *Communications of the ACM* 17, 773–774 (1974).
8. Dijkstra, E. W. *Communications of the ACM* 17, 773–774 (1974).
9. Dijkstra, E. W. *Communications of the ACM* 17, 773–774 (1974).
10. Dijkstra, E. W. *Communications of the ACM* 17, 773–774 (1974).

### Competing interests

The author declares no competing interests.

## Summary:

Maneage and its principles are described in [arXiv:2006.03018](https://arxiv.org/abs/2006.03018). It is a customizable template that will do the following steps/instructions (all in simple plain text files).

- ▶ **Automatically downloads** the necessary *software* and *data*.
- ▶ **Builds** the software in a **closed environment**.
- ▶ Runs the software on data to **generate** the final **research results**.
- ▶ Modification of part of the analysis will only result in re-doing that part, not the whole project.
- ▶ Using LaTeX macros, paper's figures, tables and numbers will be **Automatically updated** after a change in analysis. Allowing the scientist to focus on the scientific interpretation.
- ▶ The whole project is under **version control** (Git) to allow easy reversion to a previous state. This **encourages tests/experimentation** in the analysis.
- ▶ The **Git commit hash** of the project source, is **printed** in the published paper and **saved on output** data products. Ensuring the integrity/reproducibility of the result.
- ▶ These slides are available at <https://maneage.org/pdf/slides-intro-short.pdf>.
- ▶ Longer slides are available at <https://maneage.org/pdf/slides-intro.pdf>.

For a technical description of Maneage's implementation, as well as a checklist to customize it, and tips on good practices, please see this page:

<https://gitlab.com/maneage/project/-/blob/maneage/README-hacking.md>