

Yi-Hung Kuo<sup>a</sup> Dani Coleman<sup>b</sup> Thomas Jackson<sup>c</sup> Chih-Chieh (Jack) Chen<sup>b</sup> Andrew Gettelman<sup>b</sup>  
J. David Neelin<sup>a</sup> Eric Maloney<sup>d</sup> John Krasting<sup>c</sup>  
(a: UCLA; b: NCAR; c: GFDL; d: CSU)

Aug 01, 2020

---

Developer information

### Developer's quickstart guide

This walkthrough contains information for developers wanting to contribute a process-oriented diagnostic (POD) module to the MDTF framework. We assume that you have read the [Getting Started](#)<sup>1</sup>, and have followed the instructions therein for installing and testing the MDTF package, thus having some idea about the package structure and how it works. We further recommend running the framework on the sample model data with both `save_ps` and `save_nc` in the configuration input `src/default_tests.jsonc` set to `true`. This will preserve directories and files created by individual PODs, and help you understand how a POD is expected to write output.

For developers already familiar with version 2.0 of the framework, [section 2](#) (page 2) concisely summarizes changes from v2.0 to facilitate migration to v3.0. New developers can skip this section, as the rest of this walkthrough is self-contained.

For new developers, [section 3](#) (page 3) provides a to-do list of steps for implementing and integrating a POD into the framework, with more technical details in subsequent sections. [Section 4](#) (page 6) discusses the choice of programming languages, managing language and library dependencies through Conda, how to make use of and extend an existing Conda environment for POD development, and create a new Conda environment if necessary. In [section 5](#) (page 13), we walk the developers through the workflow of the framework, focusing on aspects that are relevant for the operation of individual PODs, and using the [Example Diagnostic POD](#)<sup>2</sup> as a concrete example to illustrate how a POD works under the framework.

We require developers to manage POD codes and submit them through [GitHub](#)<sup>3</sup>. See [section 8](#) (page 20) for how to manage code through the GitHub website and, for motivated developers, how to manage using the `git` command.

[@@@Moved from instruct:

### Scope of the analysis your POD conducts

See the [BAMS article](#)<sup>4</sup> describing version 2.0 of the framework for a description of the project's scientific goals and what we mean by a "process oriented diagnostic" (POD). We encourage PODs to have a specific, focused scope.

PODs should be relatively lightweight in terms of computation and memory requirements (eg, run time measured in minutes, not hours): this is to enable rapid feedback and iteration cycles to assist users in model development. Bear in mind that your POD may be run on model output of potentially any date range and spatial resolution. Your POD should not require strong assumptions about these quantities, or other details of the model's operation. @@@]

---

[https://mdtf-diagnostics.readthedocs.io/en/latest/\\_static/MDTF\\_getting\\_started.pdf](https://mdtf-diagnostics.readthedocs.io/en/latest/_static/MDTF_getting_started.pdf)  
<https://github.com/NOAA-GFDL/MDTF-diagnostics/tree/main/diagnostics/example>  
<https://github.com/NOAA-GFDL/MDTF-diagnostics>  
<https://doi.org/10.1175/BAMS-D-18-0042.1>

## Migration from framework v2.0

In this section we summarize issues to be aware of for developers familiar with the organization of version 2.0 of the framework. New developers can skip this section, as the rest of this documentation is self-contained.

### Getting Started and Developer's Walkthrough

A main source of documentation for v2.0 of the framework were the “Getting Started” and “Developer’s Walkthrough” documents. Updated versions of these documents are:

[Getting Started v3.0 \(PDF\)](#)<sup>5</sup>

[Developer’s Walkthrough v3.0 \(PDF\)](#)<sup>6</sup>

Note: these documents contain a subset of information available on this website, rather than new material: the text is reorganized to be placed in the same order as the v2.0 documents, for ease of comparison.

### Checklist for migrating a POD from v2.0

Here we list the broad set of tasks needed to update a diagnostic written for v2.0 of the framework to v3.0.

Update settings and varlist files: In v3.0 these have been combined into a single `settings.jsonc` file. See the settings file [format guide](#) (page 17), example POD, or [reference documentation](#) (page 22) for a description of the new format.

Update references to framework environment variables: See the table below for an overview, and the [reference documentation](#) (page 30) for complete information on what environment variables the framework sets. Note that your diagnostic should not use any hard-coded paths or variable names, but should read this information in from the framework’s environment variables.

Resubmit digested observational data: To minimize the size of supporting data users need to download, we ask that you only supply observational data specifically needed for plotting, as well as any code used to perform that data reduction from raw sources.

Remove HTML templating code: Version 2.0 of the framework required that your POD’s top-level driver script take particular steps to assemble its HTML file. In v3.0 these tasks are done by the framework: all that your POD needs to do is generate figures of the appropriate names in the specified folders, and the framework will convert and link them appropriately.

### Conversion from v2.0 environment variables

In v3.0, the paths referred to by the framework’s environment variables have been changed to be specific to your POD. The variables themselves have been renamed to avoid possible confusion. Here’s a table of the appropriate substitutions to make:

Table 1 Environment variable name conversion

Path Description	v2.0 environment variable expression	Equivalent v3.0 variable
Top-level code repository	<code>\$DIAG_HOME</code>	No variable set: PODs should not access files outside of their own source code directory within <code>\$POD_HOME</code>
POD’s source code	<code>\$VARCODE/&lt;pod name&gt;</code>	<code>\$POD_HOME</code>
POD’s observational/supporting data	<code>\$VARDATA/&lt;pod name&gt;</code>	<code>\$OBS_DATA</code>
POD’s working directory	<code>\$variab_dir/&lt;pod name&gt;</code>	<code>\$WK_DIR</code>
Path to requested netcdf data file	Currently unchanged:	
for <variable name> at date frequency <freq>	<code>\$DATADIR/&lt;freq&gt;/&lt;variable name&gt;.&lt;freq&gt;.nc</code>	<code>\$CASENAME.&lt;variable name&gt;.&lt;freq&gt;.nc</code>
Other v2.0 paths	<code>\$DATA_IN</code> , <code>\$WKDIR</code>	<code>\$DIAG_ROOT</code> , No equivalent variable set. PODs shouldn’t access files outside of their own directories; instead use one of the quantities above.

[https://mdtf-diagnostics.readthedocs.io/en/latest/\\_static/MDTF\\_getting\\_started.pdf](https://mdtf-diagnostics.readthedocs.io/en/latest/_static/MDTF_getting_started.pdf)  
[https://mdtf-diagnostics.readthedocs.io/en/latest/\\_static/MDTF\\_walkthrough.pdf](https://mdtf-diagnostics.readthedocs.io/en/latest/_static/MDTF_walkthrough.pdf)

## POD Development Checklist

In this section, we compile a to-do list summarizing necessary steps for POD implementation, as well as a checklist for mandatory POD documentation and testing before submitting your POD.

We recommend running the framework on the sample model data again with both `save_ps` and `save_nc` in the configuration input `src/default_tests.jsonc` set to `true`. This will preserve directories and files created by individual PODs in the output directory, which could come in handy when you go through the instructions below, and help understand how a POD is expected to write output.

### Preparation for POD implementation

We assume that, at this point, you have a set of scripts, written in [languages](#) (page 6) consistent with the framework's open source policy, that a) read in model data, b) perform analysis, and c) output figures. Here are 3 steps to prepare your scripts for POD implementation.

Give your POD an official name (e.g., Convective Transition; referred to as `long_name`) and a short name (e.g., `convective_transition_diag`). The latter will be used consistently to name the directories and files associated with your POD, so it should (1) loosely resemble the `long_name`, (2) avoid space bar and special characters (`!@#%&^*`), and (3) not repeat existing PODs' name (i.e., the directory names under `diagnostics/`). Try to make your PODs name specific enough that it will be distinct from PODs contributed now or in the future by other groups working on similar phenomena.

If you have multiple scripts, organize them so that there is a main driver script calling the other scripts, i.e., a user only needs to execute the driver script to perform all read-in data, analysis, and plotting tasks. This driver script should be named after the POD's short name (e.g., `convective_transition_diag.py`).

You should have no problem getting scripts working as long as you have (1) the location and filenames of model data, (2) the model variable naming convention, and (3) where to output files/figures. The framework will provide these as environment variables that you can access (e.g., using `os.environ` in Python, or `getenv` in NCL). DO NOT hard code these paths/filenames/variable naming convention, etc., into your scripts. See the [complete list](#) of environment variables supplied by the framework.

Your scripts should not access the internet or other networked resources.

### An example of using framework-provided environment variables

The framework provides a collection of environment variables, mostly in the format of strings but also some numbers, so that you can and MUST use in your code and make your POD portable and reusable.

For instance, using 3 of the environment variables provided by the framework, `CASENAME`, `DATADIR`, and `pr_var`, the full path to the hourly precipitation file can be expressed as

```
MODEL_OUTPUT_DIR = os.environ["DATADIR"]+"/1hr/"
pr_filename = os.environ["CASENAME"]+*. "+os.environ["pr_var"]+".1hr.nc"
pr_filepath = MODEL_OUTPUT_DIR + pr_filename
```

You can then use `pr_filepath` in your code to load the precipitation data.

Note that in Linux shell or NCL, the values of environment variables are accessed via a `$` sign, e.g., `os.environ["CASENAME"]` in Python is equivalent to `$CASENAME` in Linux shell/NCL.

### Relevant environment variables

The environment variables most relevant for a POD's operation are:

`POD_HOME`: Path to directory containing POD's scripts, e.g., `diagnostics/convective_transition_diag/`.

`OBS_DATA`: Path to directory containing POD's supporting/digested observation data, e.g., `inputdata/obs_data/convective_transition_diag/`.

`DATADIR`: Path to directory containing model data files for one case/experiment, e.g., `inputdata/model/QB0i.EXP1.AMIP.001/`.

`WK_DIR`: Path to directory for POD to output files. Note that this is the only directory a POD is allowed to write its output. E.g., `wkdir/MDTF_QB0i.EXP1.AMIP.001_1977_1981/convective_transition_diag/`.

1. Output figures to `$WK_DIR/obs/` and `$WK_DIR/model/` respectively.
2. `$WK_DIR/obs/PS` and `$WK_DIR/model/PS`: If a POD chooses to save vector-format figures, save them as EPS under these two directories. Files in these locations will be converted by the framework to PNG for HTML output. Caution: avoid using PS because of potential bugs in recent matplotlib and converting to PNG.
3. `$WK_DIR/obs/netCDF` and `$WK_DIR/model/netCDF`: If a POD chooses to save any digested data for later analysis/plotting, save them in two directories in NetCDF.

Note that (1) values of `POD_HOME`, `OBS_DATA`, and `WKDIR` change when the framework executes different PODs; (2) the `WKDIR` directory and subdirectories therein are automatically created by the framework. Each POD should output files as described here so that the framework knows where to find what, and also for the ease of code maintenance.

More environment variables for specifying model variable naming convention can be found in the `src/filedlist_$convention.jsonc` files. Also see [the comprehensive list](#) of environment variables supplied by the framework.

### To-do list for POD implementation

The following are the necessary steps for the POD module implementation and integration into the framework. You can use the PODs currently included in the code package under `diagnostics/` as concrete examples since they all have the same structure as described below:

1. Create your POD directory under `diagnostics/` and put all scripts in. Among the scripts, there should be 1) a driver script written in Python, 2) a template html, and 3) a `settings.jsonc` file. The POD directory, driver script, and html template should all be named after your POD's short name.

For instance, `diagnostics/convective_transition_diag/` contains its driver script `convective_transition_diag.py`, `convective_transition_diag.html`, and `settings.jsonc`, etc.

The framework will call the driver script, which calls the other scripts in the same POD directory.

The html template will be copied by the framework into the output directory to display the figures generated by the POD. You can create a new html template by simply copying and modifying the example templates from existing PODs without prior knowledge about html syntax.

`settings.jsonc` contains a POD's information. The framework will read this setting file to find out the driver script's name, verify the required environment and model data files are available, and prepare the necessary environment variables before executing the driver script.

2. Create a directory under `inputdata/obs_data/` named after the short name, and put all your digested observation data in (or more generally, any quantities that are independent of the model being analyzed).

Digested data should be in the form of numerical data, not figures.

Raw data, e.g., undigested reanalysis data will be rejected.

The data files should be small (preferably a few MB) and just enough for producing figures for model comparison.

If you really cannot reduce the data size or require GB of space, consult with the lead team.

3. Provide the Conda environment your POD requires. Either you can use one of the Conda environments currently supplied with the framework, defined by the YAML (.yml) files in `src/conda/`, or submit a .yml file for a new environment.

We recommend using existing Conda environments as much as possible. Consult with the lead team if you would like to submit a new one.

If you need a new Conda environment, add a new .yml file to `src/conda/`, and install the environment using the `conda_env_setup.sh` script as described in the Getting Started.

4. If your POD requires model data not included in the samples, prepare your own data files following instructions given in the Getting Started, and create a new configuration input from the template `src/default_tests.jsonc`.

Update `case_list` and `pod_list` in the configuration input file for your POD. Now you can try to run the framework following the Getting Started and start debugging. Good luck!

## Checklist before submitting your POD

After getting your POD working under the framework, there are 2 additional steps regarding the mandatory POD documentation and testing before you can submit your work to the lead team.

### 4. Provide documentation following the templates:

- A. Provide a comprehensive POD documentation in reStructuredText (.rst) format. This should include a one-paragraph synopsis of the POD, developers' contact information, required programming language and libraries, and model output variables, a brief summary of the presented diagnostics as well as references in which more in-depth discussions can be found.

Create a doc directory under your POD directory (e.g., `diagnostics/convective_transition_diag/doc/`) and put the .rst file and figures inside. It should be easy to copy and modify the .rst examples from existing PODs.

- B. All scripts should be self-documenting by including in-line comments. The main driver script (e.g., `convective_transition_diag.py`) should contain a comprehensive header providing information that contains the same items as in the POD documentation, except for the "More about this diagnostic" section.

- C. The one-paragraph POD synopsis (in the POD documentation) as well as a link to the Full Documentation should be placed at the top of the html template (e.g., `convective_transition_diag.html`).

### 5. Test before distribution. It is important that you test your POD before sending it to the lead team contact. Please take the time to go through the following procedures:

- A. Test how the POD fails. Does it stop with clear errors if it doesn't find the files it needs? How about if the dates requested are not presented in the model data? Can developers run it on data from another model? Have you added any code to scripts outside your own POD directory. Here are some simple tests you should try:

Move the `inputdata` directory around. Your POD should still work by simply updating the values of `OBS_DATA_ROOT` and `MODEL_DATA_ROOT` in the configuration input file.

Try to run your POD with a different set of model data. For POD development and testing, the MDTF-1 team produced the Timeslice Experiments output from the [NCAR CAM5<sup>7</sup>](#) and [GFDL AM4 \(contact the lead team programmer for password\)<sup>8</sup>](#).

If you have problems getting another set of data, try changing the files' `CASENAME` and variable naming convention. The POD should work by updating `CASENAME` and `convention` in the configuration input.

Try your POD on a different machine. Check that your POD can work with reasonable machine configuration and computation power, e.g., can run on a machine with 32 GB memory, and can finish computation in 10 min. Will memory and run time become a problem if one tries your POD on model output of high spatial resolution and temporal frequency (e.g., avoid memory problem by reading in data in segments)? Does it depend on a particular version of a certain library? Consult the lead team if there's any unsolvable problems.

- B. After you have tested your POD thoroughly, make clean tar files for distribution. Make a tar file of your digested observational data (preserving the `inputdata/obs_data/` structure). Do the same for model data used for testing (if different from what is provided by the MDTF page). Upload your POD code to your [GitHub repo](#) (page 20). The tar files (and your GitHub repo) should not include any extraneous files (backups, `pyc`, `*~`, or `#` files).

Use `tar -tf` to see what is in the tar file.

- C.  $\beta$ -test before distribution. Find people ( $\beta$ -testers) who are not involved in your POD's implementation and are willing to help. Give the tar files and point your GitHub repo to them. Ask them to try running the framework with your POD following the Getting Started instructions. Ask for comments on whether they can understand the documentation.

Possible  $\beta$ -tester candidates include nearby postdocs/grads and members from other POD-developing groups.

6. Submit your POD code through [GitHub pull request](#) (page 20), and share the tar files of digested observation (and model data if any) with the lead-team contact. Please also provide a list of tests you've conducted along with the machine configurations (e.g., memory size).

<https://www.earthsystemgrid.org/dataset/ucar.cgd.cesm4.NOAA-MDTF.html>  
<http://data1.gfdl.noaa.gov/MDTF/>

## Language choice and managing library dependencies

In this section, we discuss restrictions on coding languages and how to manage library dependencies. These are important points to be aware of when developing your POD, and may require you to modify existing code.

You must manage your POD's language/library dependencies through [Conda](#)<sup>9</sup>, since the dependencies of the framework are so managed by design, and this is also how the end-users are instructed to set up and manage their own environments for the framework. Note that Conda is not Python-specific, but allows coexisting versioned environments of most scripting languages, including, [R](#)<sup>10</sup>, [NCL](#)<sup>11</sup>, [Ruby](#)<sup>12</sup>, [PyFerret](#)<sup>13</sup>, and more.

To prevent the proliferation of dependencies, we suggest that new POD development use existing Conda environments whenever possible, e.g., [python3\\_base](#)<sup>14</sup>, [NCL\\_base](#)<sup>15</sup>, and [R\\_base](#)<sup>16</sup> for Python, NCL, and R, respectively.

### Choice of language(s)

The framework itself is written in Python, and can call PODs written in any scripting language. However, Python support by the lead team will be “first among equals” in terms of priority for allocating developer resources, etc.

To achieve portability, the MDTF cannot accept PODs written in closed-source languages (e.g., MATLAB and IDL; try [Octave](#)<sup>17</sup> and [GDL](#)<sup>18</sup> if possible). We also cannot accept PODs written in compiled languages (e.g., C or Fortran): installation would rapidly become impractical if users had to check compilation options for each POD.

Python is strongly encouraged for new PODs; PODs funded through the CPO grant are requested to be developed in Python. Python version  $\geq 3.6$  is required (official support for Python 2 was discontinued as of January 2020).

If your POD was previously developed in NCL or R (and development is not funded through a CPO grant), you do not need to re-write existing scripts in Python 3 if doing so is likely to introduce new bugs into stable code, especially if you're unfamiliar with Python.

If scripts were written in closed-source languages, translation to Python 3.6 or above is required.

We do not allow new PODs using Python 2 in principle. However, for a POD primarily coded in NCL and R, and uses Python only for the main driver script, an exception can be made on the basis of better managing existing Conda environments, after consulting with the lead team.

### POD development using exiting Conda environment

We assume that you've followed the instructions in the Getting Started to set up the Conda environments for the framework. We recommend developing POD and managing POD's dependencies following the same approach.

### Developers working with Python

The framework provides the [\\_MDTF\\_python3\\_base](#)<sup>19</sup> Conda environment (recall the [\\_MDTF](#) prefix for framework-specific environment) as the generic Python environment, which you can install following the instructions. You can then activate this environment by running in a terminal:

```
% source activate $CONDA_ENV_DIR/_MDTF_python3_base
```

where `$CONDA_ENV_DIR` is the path you used to install the Conda environments.

```
https://docs.conda.io/en/latest/
https://anaconda.org/conda-forge/r-base
https://anaconda.org/conda-forge/ncl
https://anaconda.org/conda-forge/ruby
https://anaconda.org/conda-forge/pyferret
https://github.com/NOAA-GFDL/MDTF-diagnostics/blob/develop/src/conda/env_python3_base.yml
https://github.com/NOAA-GFDL/MDTF-diagnostics/blob/develop/src/conda/env_NCL_base.yml
https://github.com/NOAA-GFDL/MDTF-diagnostics/blob/develop/src/conda/env_R_base.yml
https://www.gnu.org/software/octave/
https://github.com/gnudatalanguage/gdl
https://github.com/NOAA-GFDL/MDTF-diagnostics/blob/develop/src/conda/env_python3_base.yml
```



## MDTF Developer's Walkthrough, Release 3.0 beta 2

For developers' convenience, [JupyterLab](#)<sup>20</sup> (including [Jupyter Notebook](#)<sup>21</sup>) has been included in `python3_base`. Run `% jupyter lab` or `% jupyter notebook`, and you can start working on development.

If there are any [commonly used Python libraries](#)<sup>22</sup> that you'd like to add to `python3_base`, e.g., `jupyterlab`, run `% conda install -c conda-forge jupyterlab`.

- a. Only add libraries when necessary. We'd like to keep the environment small.
- b. Include the `-c` flag and specify using the [conda-forge](#)<sup>23</sup> channel as the library source. Combining packages from different channels (in particular, `conda-forge` and `anaconda`'s channel) may create incompatibilities. Consult with the lead team if encounter any problem.
- c. After installation, run `% conda clean --a` to clear cache.
- d. Don't forget to update `src/conda/env_python3_base.yml` accordingly.

After you've finished working under this environment, run `% conda deactivate` or simply close the terminal.

In case you need any exotic third-party libraries, e.g., a storm tracker, consult with the lead team and create your own Conda environment following [instructions](#) (page 7) below.

## Developers working with NCL or R

The framework also provides the `_MDTF_NCL_base`<sup>24</sup> and `_MDTF_R_base`<sup>25</sup> Conda environments as the generic NCL and R environments. You can install, activate/deactivate or add common NCL-/R-related libraries (or `jupyterlab`) to them using commands similar to those listed above.

### Create a new Conda environment

If your POD requires languages that aren't available in an existing environment or third-party libraries unavailable through the common [conda-forge](#)<sup>26</sup> and [anaconda](#)<sup>27</sup> channels, we ask that you notify us (since this situation may be relevant to other developers) and submit a [YAML \(.yml\) file](#)<sup>28</sup> that creates the environment needed for your POD.

The new YAML file should be added to `src/conda/`, where you can find templates for existing environments from which you can create your own.

The YAML filename should be `env_$your_POD_short_name.yml`.

The first entry of the YAML file, name of the environment, should be `_MDTF_$your_POD_short_name`.

We recommend listing `conda-forge` as the first channel to search, as it's entirely open source and has the largest range of packages. Note that combining packages from different channels (in particular, `conda-forge` and `anaconda` channels) may create incompatibilities.

We recommend constructing the list of packages manually, by simply searching your POD's code for `import` statements referencing third-party libraries. Please do not exporting your development environment with `% conda env export`, which gives platform-specific version information and will not be fully portable in all cases; it also does so for every package in the environment, not just the "top-level" ones you directly requested.

We recommend specifying versions as little as possible, out of consideration for end-users: if each POD specifies exact versions of all its dependencies, `conda` will need to install multiple versions of the same libraries. In general, specifying a version should only be needed in cases where backward compatibility was broken (e.g., Python 2 vs. 3) or a bug affecting your POD was fixed (e.g., postscript font rendering on MacOS with older NCL). `Conda` installs the latest version of each package that's consistent with all other dependencies.

<https://jupyterlab.readthedocs.io/en/stable/>  
<https://jupyter-notebook.readthedocs.io/en/stable/>  
<https://conda-forge.org/feedstocks/>  
<https://anaconda.org/conda-forge>  
[https://github.com/NOAA-GFDL/MDTF-diagnostics/blob/develop/src/conda/env\\_NCL\\_base.yml](https://github.com/NOAA-GFDL/MDTF-diagnostics/blob/develop/src/conda/env_NCL_base.yml)  
[https://github.com/NOAA-GFDL/MDTF-diagnostics/blob/develop/src/conda/env\\_R\\_base.yml](https://github.com/NOAA-GFDL/MDTF-diagnostics/blob/develop/src/conda/env_R_base.yml)  
<https://conda-forge.org/feedstocks/>  
<https://docs.anaconda.com/anaconda/packages/pkg-docs/>  
<https://docs.conda.io/projects/conda/en/latest/user-guide/tasks/manage-environments.html#creating-an-environment-file-manually>

## Testing with new Conda environment

If you've updated an existing environment or created a new environment (with corresponding changes to the YAML file), verify that your POD works.

Recall how the framework finds a proper Conda environment for a POD. First, it searches for an environment matching the POD's short name. If this fails, it then looks into the POD's `settings.jsonc` and prepares a generic environment depending on the language(s). Therefore, no additional steps are needed to specify the environment if your new YAML file follows the naming conventions above (in case of a new environment) or your `settings.jsonc` correctly lists the language(s) (in case of updating an existing environment).

For an updated environment, first, uninstall it by deleting the corresponding directory under `$CONDA_ENV_DIR`.

Re-install the environment using the `conda_env_setup.sh` script as described in the installation instructions, or create the new environment for you POD:

```
% cd $CODE_ROOT
% ./src/conda/conda_env_setup.sh --env $your_POD_short_name --conda_root $CONDA_ROOT --
  ↪ env_dir $CONDA_ENV_DIR
```

Have the framework run your POD on suitable test data.

1. Add your POD's short name to the `pod_list` section of the configuration input file (template: `src/default_tests.jsonc`).
2. Prepare the test data as described in `start_config`.

## Coding best practices: avoiding common issues

In this section we describe issues we've seen in POD code that have caused problems in the form of bugs, inefficiencies, or unintended consequences.

### All languages

PS vs. EPS figures: Save vector plots as `.eps` (Encapsulated PostScript), not `.ps` (regular PostScript).

Why: Postscript (`.ps`) is perhaps the most common vector graphics format, and almost all plotting packages are able to output postscript files. [Encapsulated Postscript](#)<sup>29</sup> (`.eps`) includes bounding box information that describes the physical extent of the plot's contents. This is used by the framework to generate bitmap versions of the plots correctly: the framework calls [ghostscript](#)<sup>30</sup> for the conversion, and if not provided with a bounding box ghostscript assumes the graphics use an entire sheet of (letter or A4) paper. This can cause plots to be cut off if they extend outside of this region.

Note that many plotting libraries will set the format of the output file automatically from the filename extension. The framework will process both `*.ps` and `*.eps` files.

### Python: General

Whitespace: Indent python code with four spaces per indent level.

Why: Python uses indentation to delineate nesting and scope within a program, and indentation that's not done consistently is a syntax error. Using four spaces is not required, but is the generally accepted standard.

Indentation can be configured in most text editors, or fixed with scripts such as `reindent.py` described [here](#)<sup>31</sup>. We recommend using a [linter](#)<sup>32</sup> such as `pylint` to find common bugs and syntax errors.

Beyond this, we don't impose requirements on how your code is formatted, but voluntarily following standard best practices (such as described in [PEP8](#)<sup>33</sup> or the Google [style guide](#)<sup>34</sup>) will make it easier for you and others to understand your code, find bugs, etc.

[https://en.wikipedia.org/wiki/Encapsulated\\_PostScript](https://en.wikipedia.org/wiki/Encapsulated_PostScript)  
<https://www.ghostscript.com/>  
<https://stackoverflow.com/q/1024435>  
<https://books.agiliq.com/projects/essential-python-tools/en/latest/linters.html>  
<https://www.python.org/dev/peps/pep-0008/>  
<https://github.com/google/styleguide/blob/gh-pages/pyguide.md>



Filesystem commands: Use commands in the `os`<sup>35</sup> and `shutil`<sup>36</sup> modules to interact with the filesystem, instead of running unix commands using `os.system()`, `commands` (which is deprecated), or `subprocess`.

Why: Hard-coding unix commands makes code less portable. Calling out to a subprocess introduces overhead and makes error handling and logging more difficult. The main reason, however, is that Python already provides these tools in a portable way. Please see the documentation for the `os`<sup>37</sup> and `shutil`<sup>38</sup> modules, summarized in this table:

Table 2 Recommended python functions for filesystem interaction

Task	Recommended function
Construct a path from <code>dir1</code> , <code>dir2</code> , ..., <code>filename</code>	<code>os.path.join</code> <sup>39</sup> ( <code>dir1</code> , <code>dir2</code> , ..., <code>filename</code> )
Split a path into directory and filename	<code>os.path.split</code> <sup>40</sup> ( <code>path</code> ) and related functions in <code>os.path</code> <sup>41</sup>
List files in directory <code>dir</code>	<code>os.scandir</code> <sup>42</sup> ( <code>dir</code> )
Move or rename a file or directory from <code>old_path</code> to <code>new_path</code>	<code>shutil.move</code> <sup>43</sup> ( <code>old_path</code> , <code>new_path</code> )
Create a directory or sequence of directories <code>dir</code>	<code>os.makedirs</code> <sup>44</sup> ( <code>dir</code> )
Copy a file from <code>path</code> to <code>new_path</code>	<code>shutil.copy2</code> <sup>45</sup> ( <code>path</code> , <code>new_path</code> )
Copy a directory <code>dir</code> , and everything inside it, to <code>new_dir</code>	<code>shutil.copytree</code> <sup>46</sup> ( <code>dir</code> , <code>new_dir</code> )
Delete a single file at <code>path</code>	<code>os.remove</code> <sup>47</sup> ( <code>path</code> )
Delete a directory <code>dir</code> and everything inside it	<code>shutil.rmtree</code> <sup>48</sup> ( <code>dir</code> )

In particular, using `os.path.join`<sup>49</sup> is more verbose than joining strings but eliminates bugs arising from missing or redundant directory separators.

## Python: Arrays

To obtain acceptable performance for numerical computation, people use Python interfaces to optimized, compiled code. `NumPy`<sup>50</sup> is the standard module for manipulating numerical arrays in Python. `xarray`<sup>51</sup> sits on top of `NumPy` and provides a higher-level interface to its functionality; any advice about `NumPy` applies to it as well.

`NumPy` and `xarray` both have extensive documentation and many tutorials, such as:

`NumPy`'s own `basic`<sup>52</sup> and `intermediate`<sup>53</sup> tutorials; `xarray`'s `overview`<sup>54</sup> and climate and weather `examples`<sup>55</sup>;

A `demonstration`<sup>56</sup> of the features of `xarray` using earth science data;

The 2020 SciPy conference has open-source, interactive `tutorials`<sup>57</sup> you can work through on your own machine or fully online using `Binder`<sup>58</sup>. In particular, there are tutorials for `NumPy`<sup>59</sup> and `xarray`<sup>60</sup>.

<https://docs.python.org/3.7/library/os.html>  
<https://docs.python.org/3.7/library/shutil.html>  
<https://docs.python.org/3.7/library/os.html>  
<https://docs.python.org/3.7/library/shutil.html>  
<https://docs.python.org/3.7/library/os.path.html?highlight=os%20path#os.path.join>  
<https://docs.python.org/3.7/library/os.path.html?highlight=os%20path#os.path.split>  
<https://docs.python.org/3.7/library/os.path.html?highlight=os%20path>  
<https://docs.python.org/3.7/library/os.html#os.scandir>  
<https://docs.python.org/3.7/library/shutil.html#shutil.move>  
<https://docs.python.org/3.7/library/os.html#os.makedirs>  
<https://docs.python.org/3.7/library/shutil.html#shutil.copy2>  
<https://docs.python.org/3.7/library/shutil.html#shutil.copytree>  
<https://docs.python.org/3.7/library/os.html#os.remove>  
<https://docs.python.org/3.7/library/shutil.html#shutil.rmtree>  
<https://docs.python.org/3.7/library/os.path.html?highlight=os%20path#os.path.join>  
<https://numpy.org/doc/stable/index.html>  
<http://xarray.pydata.org/en/stable/index.html>  
[https://numpy.org/doc/stable/user/absolute\\_beginners.html](https://numpy.org/doc/stable/user/absolute_beginners.html)  
<https://numpy.org/doc/stable/user/quickstart.html>  
<http://xarray.pydata.org/en/stable/quick-overview.html>  
<http://xarray.pydata.org/en/stable/examples.html>  
[https://rabernat.github.io/research\\_computing/xarray.html](https://rabernat.github.io/research_computing/xarray.html)  
<https://www.scipy2020.scipy.org/tutorial-information>  
<https://mybinder.org/>  
<https://github.com/enthought/Numpy-Tutorial-SciPyConf-2020>  
<https://xarray-contrib.github.io/xarray-tutorial/index.html>

## MDTF Developer's Walkthrough, Release 3.0 beta 2

Eliminate explicit for loops: Use NumPy/xarray functions instead of writing for loops in Python that loop over the indices of your data array. In particular, nested for loops on multidimensional data should never need to be used.

Why: For loops in Python are very slow compared to C or Fortran, because Python is an interpreted language. You can think of the NumPy functions as someone writing those for-loops for you in C, and giving you a way to call it as a Python function.

It's beyond the scope of this document to cover all possible situations, since this is the main use case for NumPy. We refer to the tutorials above for instructions, and to the following blog posts that discuss this specific issue:

“Look Ma, no for-loops<sup>61</sup>,” by Brad Solomon;

“Turn your conditional loops to Numpy vectors<sup>62</sup>,” by Tirthajyoti Sarkar;

“‘Vectorized’ Operations: Optimized Computations on NumPy Arrays<sup>63</sup>”, part of “Python like you mean it<sup>64</sup>,” a free resource by Ryan Soklaski.

Use xarray with netCDF data:

Why: This is xarray's use case. You can think of NumPy as implementing multidimensional matrices in the fully general, mathematical sense, and xarray providing the specialization to the case where the matrix contains data on a lat-lon-time-(etc.) grid.

xarray lets you refer to your data with human-readable labels such as 'latitude,' rather than having to remember that that's the second dimension of your array. This bookkeeping is essential when writing code for the MDTF framework, when your POD will be run on data from models you haven't been able to test on.

In particular, xarray provides seamless support for [time axes](#)<sup>65</sup>, with [support](#)<sup>66</sup> for all CF convention calendars through the `cftime` library. You can, eg, subset a range of data between two dates without having to manually convert those dates to array indices.

Again, please see the xarray tutorials linked above.

Memory use and views vs. copies: Use scalar indexing and [slices](#)<sup>67</sup> (index specifications of the form `start_index:stop_index:stride`) to get subsets of arrays whenever possible, and only use [advanced indexing](#)<sup>68</sup> features (indexing arrays with other arrays) when necessary.

Why: When advanced indexing is used, NumPy will need to create a new copy of the array in memory, which can hurt performance if the array contains a large amount of data. By contrast, slicing or basic indexing is done in-place, without allocating a new array: the NumPy documentation calls this a “view.”

Note that array slices are native [Python objects](#)<sup>69</sup>, so you can define a slice in a different place from the array you intend to use it on. Both NumPy and xarray arrays recognize slice objects.

This is easier to understand if you think about NumPy as a wrapper around C-like functions: array indexing in C is implemented with pointer arithmetic, since the array is implemented as a contiguous block of memory. An array slice is just a pointer to the same block of memory, but with different offsets. More complex indexing isn't guaranteed to follow a regular pattern, so NumPy needs to copy the requested data in that case.

See the following references for more information:

The [numpy documentation](#)<sup>70</sup> on indexing;

“Numpy Views vs Copies: Avoiding Costly Mistakes<sup>71</sup>,” by Jessica Yung;

“How can I tell if NumPy creates a view or a copy?<sup>72</sup>” on [stackoverflow](#).

MaskedArrays instead of NaNs or sentinel values: Use NumPy's [MaskedArrays](#)<sup>73</sup> for data that may contain missing or invalid values, instead of setting those entries to NaN or a sentinel value.

<https://realpython.com/numpy-array-programming/>  
<https://towardsdatascience.com/data-science-with-python-turn-your-conditional-loops-to-numpy-vectors-9484ff9c622e>  
[https://www.pythonlikeyoumeanit.com/Module3\\_IntroducingNumpy/VectorizedOperations.html](https://www.pythonlikeyoumeanit.com/Module3_IntroducingNumpy/VectorizedOperations.html)  
<https://www.pythonlikeyoumeanit.com/>  
<http://xarray.pydata.org/en/stable/time-series.html>  
<http://xarray.pydata.org/en/stable/weather-climate.html>  
<https://numpy.org/doc/stable/reference/arrays.indexing.html#basic-slicing-and-indexing>  
<https://numpy.org/doc/stable/reference/arrays.indexing.html#advanced-indexing>  
<https://docs.python.org/3.7/library/functions.html?highlight=slice#slice>  
<https://numpy.org/doc/stable/reference/arrays.indexing.html>  
<https://www.jessicayung.com/numpy-views-vs-copies-avoiding-costly-mistakes/>  
<https://stackoverflow.com/questions/11524664/how-can-i-tell-if-numpy-creates-a-view-or-a-copy>  
<https://numpy.org/doc/stable/reference/maskedarray.generic.html>

Why: One sometimes encounters code which sets array entries to fixed “sentinel values” (such as 1.0e+20 or NaN<sup>74</sup>) to indicate missing or invalid data. This is a dangerous and error-prone practice, since it’s frequently not possible to detect if the invalid entries are being used by mistake. For example, computing the variance of a timeseries with missing elements set to 1e+20 will either result in a floating-point overflow, or return zero.

NumPy provides a better solution in the form of [MaskedArrays](#)<sup>75</sup>, which behave identically to regular arrays but carry an extra boolean mask to indicate valid/invalid status. All the NumPy mathematical functions will automatically use this mask for error propagation. For [example](#)<sup>76</sup>, trying to an array element by zero or taking the square root of a negative element will mask it off, indicating that the value is invalid: you don’t need to remember to do these sorts of checks explicitly.

### Python: Plotting

Use the ‘Agg’ backend when testing your POD: For reproducibility, set the shell environment variable MPLBACKEND to Agg when testing your POD outside of the framework.

Why: Matplotlib can use a variety of [backends](#)<sup>77</sup>: interfaces to low-level graphics libraries. Some of these are platform-dependent, or require additional libraries that the MDTF framework doesn’t install. In order to achieve cross-platform portability and reproducibility, the framework specifies the ‘Agg’ non-interactive (ie, writing files only) backend for all PODs, by setting the MPLBACKEND environment variable.

When developing your POD, you’ll want an interactive backend – for example, this is automatically set up for you in a Jupyter notebook. When it comes to testing your POD outside of the framework, however, you should be aware of this backend difference.

### NCL

Deprecated calendar functions: Check the [function reference](#)<sup>78</sup> to verify that the functions you use are not deprecated in the current version of NCL<sup>79</sup>. This is especially necessary for [date/calendar functions](#)<sup>80</sup>.

Why: The framework uses a current version of NCL<sup>81</sup> (6.6.x), to avoid plotting bugs that were present in earlier versions. This is especially relevant for calendar functions: the `ut_*` set of functions have been deprecated in favor of counterparts beginning with `cd_` that take identical arguments (so code can be updated using find/replace). For example, use `cd_calendar`<sup>82</sup> instead of the deprecated `ut_calendar`<sup>83</sup>.

This change is necessary because only the `cd_*` functions support all calendars defined in the CF conventions, which is needed to process data from some models (eg, weather or seasonal models are typically run with a Julian calendar.)

### Extra tips for POD implementation

#### Scope of your POD’s code

As described above, your POD should accept model data as input and express the results of its analysis in a series of figures, which are presented to the user in a web page. Input model data will be in the form of one netCDF file (with accompanying dimension information) per variable, as requested in your POD’s [settings file](#) (page 17). Because your POD may be run on the output of any model, you should be careful about the assumptions your code makes about the layout of these files. Supporting data may be in any format and will not be modified by the framework.

The above data sources are your POD’s only input: you may provide options in the settings file for the user to configure when the POD is installed, but these cannot be changed each time the POD is run. Furthermore, your POD should not access the internet or other networked resources.

<https://en.wikipedia.org/wiki/NaN>

<https://numpy.org/doc/stable/reference/maskedarray.html>

<https://numpy.org/doc/stable/reference/maskedarray.generic.html#numerical-operations>

<https://matplotlib.org/tutorials/introductory/usage.html#backends>

<https://www.ncl.ucar.edu/Document/Functions/index.shtml>

<https://www.ncl.ucar.edu/>

<https://www.ncl.ucar.edu/Document/Functions/date.shtml>

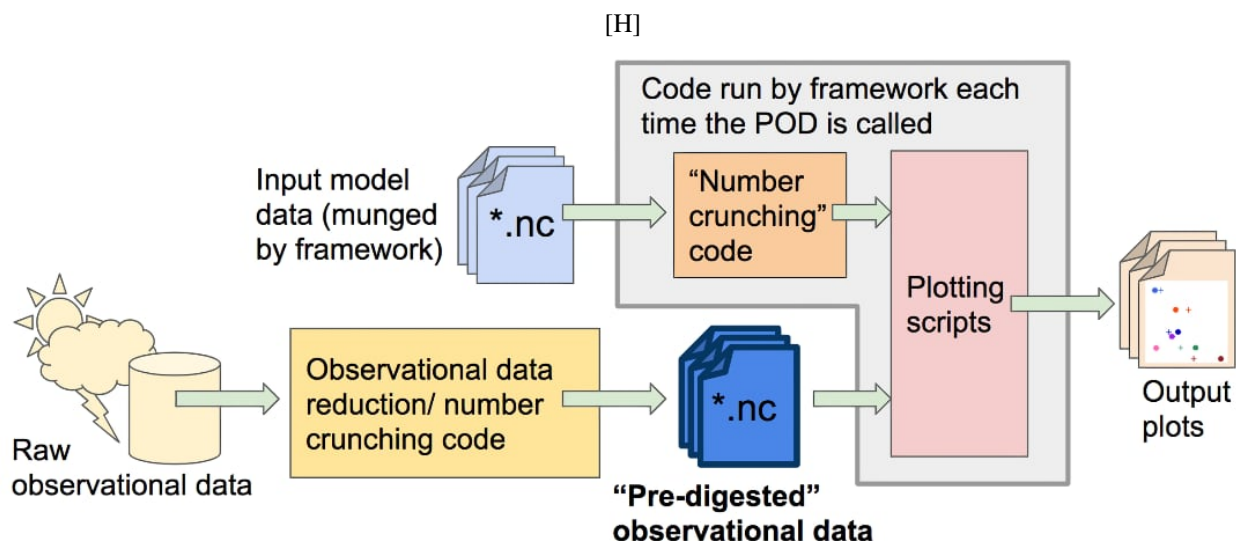
<https://www.ncl.ucar.edu/>

[https://www.ncl.ucar.edu/Document/Functions/Built-in/cd\\_calendar.shtml](https://www.ncl.ucar.edu/Document/Functions/Built-in/cd_calendar.shtml)

[https://www.ncl.ucar.edu/Document/Functions/Built-in/ut\\_calendar.shtml](https://www.ncl.ucar.edu/Document/Functions/Built-in/ut_calendar.shtml)

The output of your POD should be a series of figures in vector format (.eps or .ps), written to a specific working directory (described below). Optionally, we encourage POD developers to also save relevant output data (eg, the output data being plotted) as netcdf files, to give users the ability to take the POD's output and perform further analysis on it.

### Observational and supporting data; code organization.



In order to make your code run faster for the users, we request that you separate any calculations that don't depend on the model data (eg, pre-processing of observational data), and instead save the end result of these calculations in data files for your POD to read when it is run. We refer to this as "digested observational data," but it refers to any quantities that are independent of the model being analyzed. For purposes of data provenance, reproducibility, and code maintenance, we request that you include all the pre-processing/data reduction scripts used to create the digested data in your POD's code base, along with references to the sources of raw data these scripts take as input (yellow box in the figure).

Digested data should be in the form of numerical data, not figures, even if the only thing the POD does with the data is produce an unchanging reference plot. We encourage developers to separate their "number-crunching code" and plotting code in order to give end users the ability to customize output plots if needed. In order to keep the amount of supporting data needed by the framework manageable, we request that you limit the total amount of digested data you supply to no more than a few gigabytes.

In collaboration with PCMDI, a framework is being advanced that can help systematize the provenance of observational data used for POD development. Some frequently used datasets have been prepared with this framework, known as PCMDIobs. Please check to see if the data you require is available via PCMDIobs. If it is, we encourage you to use it, otherwise proceed as described above.

### Other tips on implementation:

1. Structure of the code package: Implementing the constituent PODs in accordance with the structure described in sections 2 and 3 makes it easy to pass the package (or just part of it) to other groups.
2. Robustness to model file/variable names: Each POD should be robust to modest changes in the file/variable names of the model output; see section 5 regarding the model output filename structure, and section 6 regarding using the environment variables and robustness tests. Also, it would be easier to apply the code package to a broader range of model output.
3. Save intermediate output: Can be used, e.g. to save time when there is a substantial computation that can be re-used when re-running or re-plotting diagnostics. See section 3.I regarding where to save the output.
4. Self-documenting: For maintenance and adaptation, to provide references on the scientific underpinnings, and for the code package to work out of the box without support. See step 5 in section 2.
5. Handle large model data: The spatial resolution and temporal frequency of climate model output have increased in recent years. As such, developers should take into account the size of model data compared with the available memory.

## MDTF Developer's Walkthrough, Release 3.0 beta 2

For instance, the example POD `precip_diurnal_cycle` and `Wheeler_Kiladis` only analyze part of the available model output for a period specified by the environment variables `FIRSTYR` and `LASTYR`, and the `convective_transition_diag` module reads in data in segments.

6. Basic vs. advanced diagnostics (within a POD): Separate parts of diagnostics, e.g., those might need adjustment when model performance out of obs range.

7. Avoid special characters (!@#\$\$%^&\*) in file/script names.

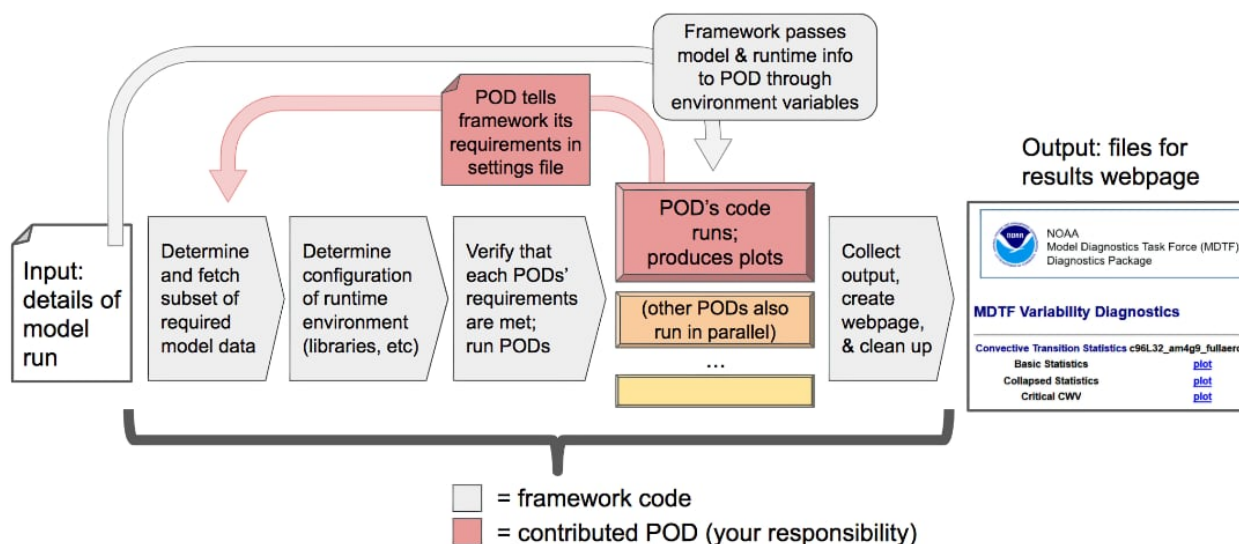
See section 3 of the Getting Started for more details on how the package is called. See the command line reference for documentation on command line options (or run `mdtf --help`).

Avoid making assumptions about the machine on which the framework will run beyond what's listed here; a development priority is to interface the framework with cluster and cloud job schedulers to enable individual PODs to run in a concurrent, distributed manner.

### Walkthrough of framework operation

We now describe in greater detail the actions that are taken when the framework is run, focusing on aspects that are relevant for the operation of individual PODs. The [Example Diagnostic POD](#)<sup>84</sup> (short name: `example`) is used as a concrete example here to illustrate how a POD is implemented and integrated into the framework.

[H]



We begin with a reminder that there are 2 essential files for the operation of the framework and POD:

`src/default_tests.jsonc`: configuration input for the framework.

`diagnostics/example/settings.jsonc`: settings file for the example POD.

To setup for running the example POD, (1) download the necessary supporting and NCAR-CAM5.timeslice sample data @@@hyperlinks required@@@ and unzip them under `inputdata/`, and (2) open `default_tests.jsonc`, uncomment the whole `NCAR-CAM5.timeslice` section in `case_list`, and comment out the other cases in the list. We also recommend setting both `save_ps` and `save_nc` to `true`.

<https://github.com/NOAA-GFDL/MDTF-diagnostics/tree/main/diagnostics/example>



## Step 1: Framework invocation

The user runs the framework by executing the framework's main driver script `$CODE_ROOT/mdtf`, rather than executing the PODs directly. This is where the user specifies the model run to be analyzed, and chooses which PODs to run via the `pod_list` section in `default_tests.jsonc`.

Some of the configuration options can be input through command line, see the command line reference or run `%$CODE_ROOT/mdtf --help`.

At this stage, the framework also creates the directory `$OUTPUT_DIR/` (default: `mdtf/wkdir/`) and all subdirectories therein for hosting the output files by the framework and PODs from each run.

If you've run the framework with both `save_ps` and `save_nc` in `default_tests.jsonc` set to `true`, check the output directory structure and files therein.

Note that when running, the framework will keep collecting the messages relevant to individual PODs, including (1) the status of required data and environment, and (2) texts printed out by PODs during execution, and will save them as log files under each POD's output directory. These log files can be viewed via the top-level results page `index.html` and, together with messages printed in the terminal, are useful for debugging.

## Example diagnostic

Run the framework using the `NCAR-CAM5.timeslice` case. After successful execution, open the `index.html` under the output directory in a web browser. The `plots` links to the webpage produced by the example POD with figures, and `log to example.log` including all example-related messages collected by the framework. The messages displayed in the terminal are not identical to those in the log files, but also provide a status update on the framework-POD operation.

## Step 2: Data request

Each POD describes the model data it requires as input in the `varlist` section of its `settings.jsonc`, with each entry in `varlist` corresponding to one model data file used by the POD. The framework goes through all the PODs to be run in `pod_list` and assembles a list of required model data from their `varlist`. It then queries the source of the model data (`$DATADIR/`) for the presence of each requested variable with the requested characteristics (e.g., frequency, units, etc.).

The most important features of `settings.jsonc` are described in the [settings documentation](#) (page 17) and full detail on the [reference page](#) (page 22).

Variables are specified in `varlist` following [CF convention](#)<sup>85</sup> wherever possible. If your POD requires derived quantities that are not part of the standard model output (e.g., column weighted averages), incorporate necessary pre-processings for computing these from standard output variables into your code. PODs are allowed to request variables outside of the CF conventions (by requiring an exact match on the variable name), but this will severely limit the POD's application.

Some of the requested variables may be unavailable or without the requested characteristics (e.g., frequency). You can specify a backup plan for this situation by designating sets of variables as alternates if feasible: when the framework is unable to obtain a variable that has the `alternates` attribute in `varlist`, it will then (and only then) query the model data source for the variables named as alternates.

If no alternates are defined or the alternate variables are also unavailable, the framework will skip executing your POD, and an error log will be presented in `index.html`.

Once the framework has determined which PODs are able to run given the model data, it prepares the necessary environment variables, including directory paths and the requested variable names (as defined in `src/filedlist_$convention.jsonc`) for PODs' operation.

At this step, the framework also checks the PODs' observational/supporting data under `inputdata/obs_data/`. If the directory of any of the PODs in `pod_list` is missing, the framework would terminate with error messages showing on the terminal. Note that the framework only checks the presence of the directory, but not the files therein.

---

<http://cfconventions.org/>



## Example diagnostic

The example POD uses only one model variable in its `varlist`<sup>86</sup>: surface air temperature, recorded at monthly frequency. In the beginning of `example.log`, the framework reports finding the requested model data file under `Found files`. If the framework could not locate the file, the log would instead record `Skipping execution` with the reason being missing data.

### Step 3: Runtime environment configuration

The framework reads the other parts of your POD's `settings.jsonc`, e.g., , and generates the additional environment variables accordingly (on top of those being defined through `default_tests.jsonc`).

Furthermore, in the `runtime_requirements` section of `settings.jsonc`, we request that you provide a list of languages and third-party libraries your POD uses. The framework will check that all these requirements are met by one of the Conda environments under `$CONDA_ENV_DIR/`.

The requirements should be satisfied by one of the existing generic Conda environments (updated by you if necessary), or a new environment you created specifically for your POD.

If there isn't a suitable environment, the POD will be skipped.

## Example diagnostic

In its `settings.jsonc`, the example POD lists its `requirements`<sup>87</sup>: Python 3, and the `matplotlib`, `xarray` and `netCDF4` third-party libraries for Python. In this case, the framework assigns the POD to run in the generic `python3_base`<sup>88</sup> environment provided by the framework.

In `example.log`, under `Env vars:` is a comprehensive list of environment variables prepared for the POD by the framework. A great part of them are defined as in `src/filedlist_$convention.jsonc` via convention in `default_tests`. Some of the environment variables are POD-specific as defined under `'pod_env_vars'` in the POD's `settings.jsonc`, e.g., `EXAMPLE_FAV_COLOR`.

In `example.log`, after `--- MDTF.py calling POD example`, the framework verifies the Conda-related paths, and makes sure that the `runtime_requirements` in `settings.jsonc` are met by the Conda environment assigned to the POD.

### Step 4: POD execution

At this point, your POD's requirements have been met, and the environment variables are set. The framework then activates the right Conda environment, and begins execution of your POD's code by calling the top-level driver script listed in its `settings.jsonc`.

See [Relevant environment variables](#) (page 3) for most relevant environment variables, and how your POD is expected to output results.

All information passed from the framework to your POD is in the form of Unix/Linux shell environment variables; see [reference](#) for a complete list of environment variables (another good source is the log files for individual PODs).

For debugging, we encourage that your POD print out messages of its progress as it runs. All text written to `stdout` or `stderr` (i.e., displayed in a terminal) will be captured by the framework and added to a log file available to the users via `index.html`.

Properly structure your code/scripts and include error and exception handling mechanisms so that simple issues will not completely shut down the POD's operation. Here are a few suggestions:

- A. Separate basic and advanced diagnostics. Certain computations (e.g., fitting) may need adjustment or are more likely to fail when model performance out of observed range. Organize your POD scripts so that the basic part can produce results even when the advanced part fails.

<sup>86</sup><https://github.com/NOAA-GFDL/MDTF-diagnostics/blob/d8d9f951d2c887b9a30fc496298815ab7ee68569/diagnostics/example/settings.jsonc#L46>  
<sup>87</sup><https://github.com/NOAA-GFDL/MDTF-diagnostics/blob/d8d9f951d2c887b9a30fc496298815ab7ee68569/diagnostics/example/settings.jsonc#L38>  
<sup>88</sup>[https://github.com/NOAA-GFDL/MDTF-diagnostics/blob/main/src/conda/env\\_python3\\_base.yml](https://github.com/NOAA-GFDL/MDTF-diagnostics/blob/main/src/conda/env_python3_base.yml)

B.If some of the observational data files are missing by accident, the POD should still be able to run analysis and produce figures for model data regardless.

C.Say a POD reads in multiple variable files and computes statistics for individual variables. If some of the files are missing or corrupted, the POD should still produce results for the rest (note that in this case, the framework would skip this POD due to missing data, but PODs should have this robustness property for ease of workarounds or running outside the framework).

The framework contains additional exception handling so that if a POD experiences a fatal or unrecoverable error, the rest of the tasks and POD-calls by the framework can continue. The error messages, if any, will be included in the POD's log file.

### Example diagnostic

The framework activates the `_MDTF_python3_base` Conda environment, and calls the driver script `example-diag.py`<sup>89</sup> listed in `settings.jsonc`. Take a look at the script and the comments therein.

`example-diag.py` performs tasks roughly in the following order:

- 1)It reads the model surface air temperature data at `input_path`,
- 2)computes the model time average,
- 3)saves the model time averages to `$WK_DIR/model/netCDF/temp_means.nc` for later use,
- 4)plots model figure `$WK_DIR/model/PS/example_model_plot.eps`,
- 5)reads the digested data in time-averaged form at `$OBS_DATA/example_tas_means.nc`, and plots the figure to `$WK_DIR/obs/PS/example_obs_plot.eps`.

Note that these tasks correspond to the code blocks 1) through 5) in the script.

When the script is called and running, it prints out messages which are saved in `example.log`. These are helpful to determine when and how the POD execution is interrupted if there's a problem.

The script is organized to deal with model data first, and then to process digested observations. Thus if something goes wrong with the digested data, the script is still able to produce the html page with model figures. This won't happen if code block 5) is moved before 4), i.e., well-organized code is more robust and may be able to produce partial results even when it encounters problems.

In code block 7) of `example-diag.py`, we include an example of exception handling by trying to access a non-existent file (the final block is just to confirm that the error would not interrupt the script's execution because of exception-handling).

The last few lines of `example.log` demonstrate the script is able to finish execution despite an error having occurred. Exception handling makes code robust.

### Step 5: Output and cleanup

At this point, your POD has successfully finished running, and all remaining tasks are handled by the framework. The framework converts the postscript plots to bitmaps according to the following rule:

`$WK_DIR/model/PS/filename.eps` → `$WK_DIR/model/filename.png`

`$WK_DIR/obs/PS/filename.eps` → `$WK_DIR/obs/filename.png`

The html template for each POD is then copied to `$WK_DIR` by the framework.

In writing the template file all plots should be referenced as relative links to this location, e.g., "`<A href=model/filename.png>`". See templates from existing PODs.

Values of all environment variables referenced in the html template are substituted by the framework, allowing you to show the run's CASENAME, date range, etc. Text you'd like to change at runtime must be changed through environment variables: the v3 framework does not allow other ways to alter the text of your POD's output webpage at run-time.

If `save_ps` and `save_nc` are set to `false`, the `.eps` and `.nc` files will be deleted.

Finally, the framework links your POD's html page to the top-level `index.html`, and copies all files to the specified output location.

[https://github.com/NOAA-GFDL/MDTF-diagnostics/blob/main/diagnostics/example/example\\_diag.py](https://github.com/NOAA-GFDL/MDTF-diagnostics/blob/main/diagnostics/example/example_diag.py)

If `make_variab_tar` in `default_tests.jsonc` is set to `true`, the framework will create a tar file for the output directory, in case you're working on a server, and have to move the file to a local machine before viewing it.

### Diagnostic settings file quickstart

This page gives a quick introduction to how to write the settings file for your diagnostic. See the full [documentation](#) (page 22) on this file format for a complete list of all the options you can specify.

### Overview

The MDTF framework can be viewed as a “wrapper” for your code that handles data fetching and munging. Your code communicates with this wrapper in two ways:

The settings file is where your code talks to the framework: when you write your code, you document what model data your code uses and what format it expects it in. When the framework is run, it will fulfill the requests you make here (or tell the user what went wrong).

When your code is run, the framework talks to it by setting [environment variables](#) (page 30) containing paths to the data files and other information specific to the run (not covered on this page, follow the link for details).

In the settings file, you specify what model data your diagnostic uses in a vocabulary you're already familiar with:

The [CF conventions](#)<sup>90</sup> for standardized variable names and units.

The netCDF4 (classic) data model, in particular the notions of [variables](#)<sup>91</sup> and [dimensions](#)<sup>92</sup> as they're used in a netCDF file.

### Example

```
// Any text to the right of a '//' is a comment
{
  "settings" : {
    "long_name": "My example diagnostic",
    "driver": "example_diagnostic.py",
    "realm": "atmos",
    "runtime_requirements": {
      "python": ["numpy", "matplotlib", "netCDF4"]
    }
  },
  "data" : {
    "frequency": "day"
  },
  "dimensions": {
    "lat": {
      "standard_name": "latitude"
    },
    "lon": {
      "standard_name": "longitude"
    },
    "plev": {
      "standard_name": "air_pressure",
      "units": "hPa",
      "positive": "down"
    },
    "time": {
      "standard_name": "time",
      "units": "day"
    }
  }
}
```

(continues on next page)

<http://cfconventions.org/>

<https://www.unidata.ucar.edu/software/netcdf/workshops/2010/datamodels/NcVars.html>

<https://www.unidata.ucar.edu/software/netcdf/workshops/2010/datamodels/NcDims.html>

```

},
"varlist" : {
  "my_precip_data": {
    "standard_name": "precipitation_flux",
    "path_variable": "PATH_TO_PR_FILE",
    "units": "kg m-2 s-1",
    "dimensions" : ["time", "lat", "lon"]
  },
  "my_3d_u_data": {
    "standard_name": "eastward_wind",
    "path_variable": "PATH_TO_UA_FILE",
    "units": "m s-1",
    "dimensions" : ["time", "plev", "lat", "lon"]
  }
}
}

```

## Settings section

This is where you describe your diagnostic and list the programs it needs to run.

Display name of your diagnostic, used to describe your diagnostic on the top-level index.html page. Can contain spaces.

Filename of the driver script the framework should call to run your diagnostic.

One or more of the eight CMIP6 modeling realms (aerosol, atmos, atmosChem, land, landIce, ocean, ocnBgchem, seaIce) describing what data your diagnostic uses. This is give the user an easy way to, eg, run only ocean diagnostics on data from an ocean model.

This is a list of key-value pairs describing the programs your diagnostic needs to run, and any third-party libraries used by those programs.

The key is program's name, eg. languages such as “python<sup>93</sup>” or “ncl<sup>94</sup>” etc. but also any utilities such as “ncks<sup>95</sup>”, “cdo<sup>96</sup>”, etc.

The value for each program is a list of third-party libraries in that language that your diagnostic needs. You do not need to list built-in libraries: eg, in python, you should to list `numpy`<sup>97</sup> but not `math`<sup>98</sup>. If no third-party libraries are needed, the value should be an empty list.

## Data section

This section contains settings that apply to all the data your diagnostic uses. Most of them are optional.

The time frequency the model data should be provided at, eg. “1hr”, “6hr”, “day”, “mon”, ...

---

<https://www.python.org/>  
<https://www.ncl.ucar.edu/>  
<http://nco.sourceforge.net/>  
<https://code.mpimet.mpg.de/projects/cdo>  
<https://numpy.org/>  
<https://docs.python.org/3/library/math.html>

## Dimensions section

This section is where you list the dimensions (coordinate axes) your variables are provided on. Each entry should be a key-value pair, where the key is the name your diagnostic uses for that dimension internally, and the value is a list of settings describing that dimension. In order to be unambiguous, all dimensions must specify at least:

The CF [standard name](#)<sup>99</sup> for that coordinate.

The units the diagnostic expects that coordinate to be in (using the syntax of the [UDUnits library](#)<sup>100</sup>). This is optional: if not given, the framework will assume you want CF convention [canonical units](#)<sup>101</sup>.

In addition, any vertical (Z axis) dimension must specify:

Either "up" or "down", according to the [CF conventions](#)<sup>102</sup>. A pressure axis is always "down" (increasing values are closer to the center of the earth).

## Varlist section

This section is where you list the variables your diagnostic uses. Each entry should be a key-value pair, where the key is the name your diagnostic uses for that variable internally, and the value is a list of settings describing that variable. Most settings here are optional, but the main ones are:

The CF [standard name](#)<sup>103</sup> for that variable.

Name of the shell environment variable the framework will use to pass the location of the file containing this variable to your diagnostic when it's run. See the environment variable [documentation](#) (page 30) for details.

The units the diagnostic expects the variable to be in (using the syntax of the [UDUnits library](#)<sup>104</sup>). This is optional: if not given, the framework will assume you want CF convention [canonical units](#)<sup>105</sup>.

List of names of dimensions specified in the "dimensions" section, to specify the coordinate dependence of each variable.

## General developer resources

The following links to third-party pages contain information that may be helpful.

### Git tutorials/references

The official git [tutorial](#)<sup>106</sup>.

A more verbose [introduction](#)<sup>107</sup> to the ideas behind git and version control.

A still more detailed [walkthrough](#)<sup>108</sup> which assumes no prior knowledge.

<http://cfconventions.org/Data/cf-standard-names/72/build/cf-standard-name-table.html>  
<https://www.unidata.ucar.edu/software/udunits/udunits-2.0.4/udunits2lib.html#Syntax>  
<http://cfconventions.org/Data/cf-standard-names/current/build/cf-standard-name-table.html>  
[http://cfconventions.org/faq.html#vertical\\_coords\\_positive\\_attribute](http://cfconventions.org/faq.html#vertical_coords_positive_attribute)  
<http://cfconventions.org/Data/cf-standard-names/72/build/cf-standard-name-table.html>  
<https://www.unidata.ucar.edu/software/udunits/udunits-2.0.4/udunits2lib.html#Syntax>  
<http://cfconventions.org/Data/cf-standard-names/current/build/cf-standard-name-table.html>  
<https://git-scm.com/docs/gittutorial>  
<https://www.atlassian.com/git/tutorials/what-is-version-control>  
<http://swcarpentry.github.io/git-novice/>

## Python coding style

PEP8<sup>109</sup>, the officially recognized Python style guide.

Google's Python style guide<sup>110</sup>.

## Code documentation

Documentation for the framework's code is managed using sphinx<sup>111</sup>, which works with files in reStructured text<sup>112</sup> (reST, .rst) format. The framework uses Google style conventions for python docstrings.

Sphinx quickstart<sup>113</sup>.

reStructured text introduction<sup>114</sup>, quick reference<sup>115</sup> and in-depth guide<sup>116</sup>.

reST syntax comparison<sup>117</sup> to other text formats you may be familiar with.

Style guide for google-style python docstrings<sup>118</sup> and quick examples<sup>119</sup>.

## Developing for MDTF Diagnostics through GitHub website and git command

We recommend developers to manage the MDTF package using the GitHub webpage interface:

If you have no prior experience with GitHub<sup>120</sup>, create an account first.

Create a fork of the project by clicking the Fork button in the upper-right corner of NOAA's MDTF GitHub page<sup>121</sup>. This will create a copy (also known as repository, or simply repo) in your own GitHub account which you have full control over.

Before you start working on the code, remember to switch to the develop branch (instead of main) as expected from a POD developer.

It should be easy to figure out how to add/edit files through your repo webpage interface.

After updating the code in your repo, submit a Pull request so that the changes you have made can be incorporated into the official NOAA's repo.

Your changes will not affect the official NOAA's repo until the pull request is accepted by the lead-team programmer.

Note that if any buttons are missing, try CTRL`~+`~+ or CTRL`~+`~` to adjust the webpage font size so the missing buttons may magically appear.

Managing through the webpage interface as described above is quick and easy. Another approach, unfortunately with a steeper learning curve, is to create a local repo on your machine and manage the code using the git command in a terminal. In the interests of making things self-contained, the rest of this section gives brief step-by-step instructions on git for interested developers.

Before following the instructions below, make sure that a) you've created a fork of the project, and b) the git command is available on your machine (installation instructions<sup>122</sup>).

<https://www.python.org/dev/peps/pep-0008/>  
<https://github.com/google/styleguide/blob/gh-pages/pyguide.md>  
<https://www.sphinx-doc.org/en/master/index.html>  
<https://docutils.sourceforge.io/rst.html>  
<http://www.sphinx-doc.org/en/master/usage/quickstart.html>  
<http://docutils.sourceforge.net/docs/user/rst/quickstart.html>  
<http://docutils.sourceforge.net/docs/user/rst/quickref.html>  
<http://docutils.sourceforge.net/docs/ref/rst/restructuredtext.html>  
<http://hyperpolyglot.org/lightweight-markup>  
<https://github.com/google/styleguide/blob/gh-pages/pyguide.md#38-comments-and-docstrings>  
[https://sphinxcontrib-napoleon.readthedocs.io/en/latest/example\\_google.html](https://sphinxcontrib-napoleon.readthedocs.io/en/latest/example_google.html)  
<https://github.com/>  
<https://github.com/NOAA-GFDL/MDTF-diagnostics>  
<https://git-scm.com/download/>



### Some online git resources

If you are new to git and unfamiliar with many of the terminologies, [Dangit, Git?!<sup>123</sup>](#) provides solutions in plain English to many common mistakes people have made.

There are many comprehensive online git tutorials, such as:

The official [git tutorial<sup>124</sup>](#).

A more verbose [introduction<sup>125</sup>](#) to the ideas behind git and version control.

A still more detailed [walkthrough<sup>126</sup>](#), assuming no prior knowledge.

### Set up SSH with GitHub

You have to generate an [SSH key<sup>127</sup>](#) and [add it<sup>128</sup>](#) to your GitHub account. This will save you from having to re-enter your GitHub username and password every time you interact with their servers.

When generating the SSH key, you'll be asked to pick a passphrase (i.e., password).

The following instructions assume you've generated an SSH key. If you're using manual authentication instead, replace the "git@github.com:" addresses in what follows with "https://github.com/".

### Clone a local repository onto your machine

Clone your fork onto your computer: `git clone git@github.com:<your_github_account>/MDTF-diagnostics.git`. This not only downloads the files, but due to the magic of git also gives you the full commit history of all branches.

Enter the project directory: `cd MDTF-diagnostics`.

Clone additional dependencies of the code: `git submodule update --recursive --init`.

Git knows about your fork, but you need to tell it about NOAA's repo if you wish to contribute changes back to the code base. To do this, type `git remote add upstream git@github.com:NOAA-GFDL/MDTF-diagnostics.git`. Now you have two remote repos: origin, your GitHub fork which you can read and write to, and upstream, NOAA's code base which you can only read from.

### Start coding

Switch to the develop branch: `git checkout develop`.

If it's been a while since you created your fork, other people may have updated NOAA's develop branch. To make sure you're up-to-date, get these changes with `git pull upstream develop` and `git submodule update --recursive --remote`.

That command updates the working copy on your computer, but you also need to tell your fork on GitHub about the changes: `git push origin develop`.

Now you're up-to-date and ready to start working on a new feature. `git checkout -b feature/<my_feature_name>` will create a new branch (-b flag) off of develop and switch you to working on that branch.

If you are unfamiliar with git and want to practice with the commands listed here, we recommend you to create an additional feature branch just for this. Remember: your changes will not affect NOAA's repo until you've submitted a pull request through the GitHub webpage and accepted by the lead-team programmer.

If you encounter problems during practice, you can first try looking for plain English instructions to unmess the situation at [Dangit, Git?!<sup>129</sup>](#).

---

<https://dangitgit.com/>  
<https://git-scm.com/docs/gittutorial>  
<https://www.atlassian.com/git/tutorials/what-is-version-control>  
<http://swcarpentry.github.io/git-novice/>  
<https://help.github.com/en/articles/generating-a-new-ssh-key-and-adding-it-to-the-ssh-agent>  
<https://help.github.com/en/articles/adding-a-new-ssh-key-to-your-github-account>  
<https://dangitgit.com/>

## MDTF Developer's Walkthrough, Release 3.0 beta 2

Write your code! A useful command is `git status` to remind you what branch you're on and changes you've made (but have not committed yet). `git branch -a` lists all branches with `*` indicating the branch you're on.

If you've added new files, `git add --all` before commit the changes.

Commit changes with `git commit -a`. This creates a snapshot of the code into the history in your local repo.

The snapshot will exist until you intentionally delete it (after confirming a warning message). You can always revert to a previous snapshot.

Don't commit code that you know is buggy or non-functional!

You'll be asked to enter a commit message. Good commit messages are key to making the project's history useful.

Write in present tense describing what the commit, when applied, does to the code – not what you did to the code.

Messages should start with a brief, one-line summary, less than 80 characters. If this is too short, you may want to consider entering your changes as multiple commits.

When finish updating your feature, merge it back into develop: first `git checkout develop` then `git merge --no-ff feature/<my_feature_name>`. The `'--no-ff'` flag is important: it tells git not to compress ("fast-forward") your commit history onto the develop branch.

`git push origin` so that the changes to your local repo is incorporated to the your GitHub fork (displayed on the webpage).

Enter the SSH key passphrase when asked for password.

If you haven't finished working on your feature, you can checkout and update the local feature branch again by repeating the above commands.

When your feature is ready, submit a pull request by going to the GitHub page of your fork and clicking on the Pull request button. This is your proposal to the maintainers to incorporate your feature into NOAA's repo.

When the feature branch is no longer needed, delete the branch locally with `git branch -d feature/<my_feature_name>`. If you pushed it to your fork, you can delete it remotely with `git push --delete origin feature/<my_feature_name>`.

Remember that branches in git are just pointers to a particular commit, so by deleting a branch you don't lose any history.

If you want to let others work on your feature, push its branch to your GitHub fork with `git push -u origin feature/<my_feature_name>`. The `-u` flag is for creating a new branch remotely and only needs to be used the first time.

Framework reference

### Diagnostic settings file format

The settings file is how your diagnostic tells the framework what it needs to run, in terms of software and model data.

Each diagnostic must contain a text file named `settings.jsonc` in the JSON<sup>130</sup> format, with the addition that any text to the right of `//` is treated as a comment and ignored (sometimes called the "JSONC" format).

### Brief summary of JSON

We'll briefly summarize subset of JSON syntax used in this configuration file. The file's JSON expressions are built up out of items, which may be either

- 1.a boolean, taking one of the values `true` or `false` (lower-case, with no quotes).
- 2.a number (integer or floating-point).
- 3.a case-sensitive string, which must be delimited by double quotes.

In addition, for the purposes of the configuration file we define

- 4.a "unit-ful quantity": this is a string containing a number followed by a unit, eg. `"6hr"`. In addition, the string `"any"` may be used to signify that any value is acceptable.

Items are combined in compound expressions of two types:

[https://en.wikipedia.org/wiki/JSON#Data\\_types\\_and\\_syntax](https://en.wikipedia.org/wiki/JSON#Data_types_and_syntax)

5. arrays, which are one-dimensional ordered lists delimited with square brackets. Entries can be of any type, eg [true, 1, "two"].
6. objects, which are un-ordered lists of key:value pairs separated by colons and delimited with curly brackets. Keys must be strings and must all be unique within the object, while values may be any expression, eg. {"red": 0, "green": false, "blue": "bagels"}.
- Compound expressions may be nested within each other to an arbitrary depth.

## File organization

```
{
  "settings" : {
    <...properties describing the diagnostic...>
  },
  "data" : {
    <...properties for all requested model data...>
  },
  "dimensions" : {
    "my_first_dimension": {
      <...properties describing this dimension...>
    },
    "my_second_dimension": {
      <...properties describing this dimension...>
    },
    ...
  },
  "varlist" : {
    "my_first_variable": {
      <...properties describing this variable...>
    },
    "my_second_variable": {
      <...properties describing this variable...>
    },
    ...
  }
}
```

At the top level, the settings file is an **object** (page 23) containing four required entries, described in detail below.

**settings** (page 23): properties that label the diagnostic and describe its runtime requirements.

**data** (page 25): properties that apply to all the data your diagnostic is requesting.

**dimensions** (page 26): properties that apply to the dimensions (in **netCDF**<sup>131</sup> terminology) of the model data. Each distinct dimension (coordinate axis) of the data being requested should be listed as a separate entry here.

**varlist** (page 28): properties that describe the individual variables your diagnostic operates on. Each variable should be listed as a separate entry here.

## Settings section

This section is an **object** (page 23) containing properties that label the diagnostic and describe its runtime requirements.

<https://www.unidata.ucar.edu/software/netcdf/workshops/2010/datamodels/NcDims.html>

## Example

```

"settings" : {
  "long_name" : "Effect of X on Y diagnostic",
  "driver" : "my_script.py",
  "realm" : ["atmos", "ocean"],
  "runtime_requirements": {
    "python": ["numpy", "matplotlib", "netCDF4", "cartopy"],
    "ncl": ["contributed", "gsn_code", "gsn_csm"]
  },
  "pod_env_vars" : {
    // RES: Spatial Resolution (degree) for Obs Data (0.25, 0.50, 1.00).
    "RES": "1.00"
  }
}

```

## Diagnostic description

String, required. Human-readable display name of your diagnostic. This is the text used to describe your diagnostic on the top-level index.html page. It should be in sentence case (capitalize first word and proper nouns only) and omit any punctuation at the end.

String, required. Filename of the top-level driver script the framework should call to run your diagnostic's analysis.

String or [array](#) (page 22) (list) of strings, required. One of the eight CMIP6 modeling realms (aerosol, atmos, atmosChem, land, landIce, ocean, ocnBgchem, sealce) describing what data your diagnostic uses. If your diagnostic uses data from multiple realms, list them in an array (eg. ["atmos", "ocean"]). This information doesn't affect how the framework fetches model data for your diagnostic: it's provided to give the user a shortcut to say, eg., "run all the atmos diagnostics on this output."

## Diagnostic runtime

[object](#) (page 23), required. Programs your diagnostic needs to run (for example, scripting language interpreters) and any third-party libraries needed in those languages. Each executable should be listed in a separate key-value pair:

The key is the name of the required executable, eg. languages such as "python"<sup>132</sup> or "ncl"<sup>133</sup> etc. but also any utilities such as "ncks"<sup>134</sup>, "cdo"<sup>135</sup>, etc.

The value corresponding to each key is an [array](#) (page 22) (list) of strings, which are names of third-party libraries in that language that your diagnostic needs. You do not need to list standard libraries or scripts that are provided in a standard installation of your language: eg, in python, you need to list [numpy](#)<sup>136</sup> but not [math](#)<sup>137</sup>. If no third-party libraries are needed, the value should be an empty list.

In the future we plan to offer the capability to request specific [versions](#)<sup>138</sup>. For now, please communicate your diagnostic's version requirements to the MDTF organizers.

[object](#) (page 23), optional. Names and values of shell environment variables used by your diagnostic, in addition to those supplied by the framework. The user can't change these at runtime, but this can be used to set site-specific installation settings for your diagnostic (eg, switching between low- and high-resolution observational data depending on what the user has chosen to download). Note that environment variable values must be provided as strings.

---

<https://www.python.org/>  
<https://www.ncl.ucar.edu/>  
<http://nco.sourceforge.net/>  
<https://code.mpimet.mpg.de/projects/cdo>  
<https://numpy.org/>  
<https://docs.python.org/3/library/math.html>  
<https://docs.conda.io/projects/conda/en/latest/user-guide/concepts/pkg-specs.html#package-match-specifications>

## Data section

This section is an [object](#) (page 23) containing properties that apply to all the data your diagnostic is requesting.

### Example

```
"data": {
  "format": "netcdf4_classic",
  "rename_dimensions": false,
  "rename_variables": false,
  "multi_file_ok": true,
  "frequency": "3hr",
  "min_frequency": "1hr",
  "max_frequency": "6hr",
  "min_duration": "5yr",
  "max_duration": "any"
}
```

### Example

String. Optional: assumed "any\_netcdf\_classic" if not specified. Specifies the format(s) of model data your diagnostic is able to read. As of this writing, the framework only supports retrieval of netCDF formats, so only the following values are allowed:

"any\_netcdf" includes all of:

"any\_netcdf3" includes all of:

"netcdf3\_classic" (CDF-1, files restricted to < 2 Gb)

"netcdf3\_64bit\_offset" (CDF-2)

"netcdf3\_64bit\_data" (CDF-5)

"any\_netcdf4" includes all of:

"netcdf4\_classic"

"netcdf4"

"any\_netcdf\_classic" includes all the above except "netcdf4" (classic data model only).

See the [netCDF FAQ<sup>139</sup>](#) (under "Formats, Data Models, and Software Releases") for information on the distinctions. Any recent version of a supported language for diagnostics with netCDF support will be able to read all of these. However, the extended features of the "netcdf4" data model are not commonly used in practice and currently only supported at a beta level in NCL, which is why we've chosen "any\_netcdf\_classic" as the default.

Boolean. Optional: assumed false if not specified. If set to true, the framework will change the name of all [dimensions](#) (page 26) in the model data from the model's native value to the string specified in the name property for that dimension. If set to false, the diagnostic is responsible for reading dimension names from the environment variable. See the environment variable [documentation](#) (page 30) for details on how these names are provided.

Boolean. Optional: assumed false if not specified. If set to true, the framework will change the name of all [variables](#) (page 28) in the model data from the model's native value to the string specified in the name property for that variable. If set to false, the diagnostic is responsible for reading dimension names from the environment variable. See the environment variable [documentation](#) (page 30) for details on how these names are provided.

Boolean. Optional: assumed false if not specified. If set to true, the diagnostic can handle datasets for a single variable spread across multiple files, eg [xarray<sup>140</sup>](#).

[Unit-ful quantities](#) (page 22). Optional: assumed "any" if not specified. Set minimum and maximum length of the analysis period for which the diagnostic should be run: this overrides any choices the user makes at runtime. Some example uses of this setting are:

<https://www.unidata.ucar.edu/software/netcdf/docs/faq.html>  
[http://xarray.pydata.org/en/stable/generated/xarray.open\\_mfdataset.html](http://xarray.pydata.org/en/stable/generated/xarray.open_mfdataset.html)

If your diagnostic uses low-frequency (eg seasonal) data, you may want to set `min_duration` to ensure the sample size will be large enough for your results to be statistically meaningful.

On the other hand, if your diagnostic uses high-frequency (eg hourly) data, you may want to set `max_duration` to prevent the framework from attempting to download a large volume of data for your diagnostic if the framework is called with a multi-decadal analysis period.

The following properties can optionally be set individually for each variable in the `varlist` section (page 28). If so, they will override the global settings given here.

Boolean. Optional: assumed `false` if not specified. If set to `true`, the framework will ensure that the dimensions of each variable's array are given in the same order as listed in `dimensions`. If set to `false`, your diagnostic is responsible for handling arbitrary dimension orders: eg. it should not assume that 3D data will be presented as (time, lat, lon).

**Unit-ful quantities** (page 22). Time frequency at which the data is provided. Either `frequency` or the min/max pair, or both, is required:

If only `frequency` is provided, the framework will attempt to obtain data at that frequency. If that's not available from the data source, your diagnostic will not run.

If the min/max pair is provided, the diagnostic must be capable of using data at any frequency within that range (inclusive). The diagnostic is responsible for determining the frequency if this option is used.

If all three properties are set, the framework will first attempt to find data at `frequency`. If that's not available, it will try data within the min/max range, so your code must be able to handle this possibility.

## Dimensions section

This section is an `object` (page 23) contains properties that apply to the dimensions of model data. "Dimensions" are meant in the sense of the netCDF `data model`<sup>141</sup>: informally, they are "coordinate axes" holding the values of independent variables that the dependent variable is sampled at.

All `dimensions` (page 29) and `scalar coordinates` (page 29) referenced by variables in the `varlist` section must have an entry in this section. If two variables reference the same dimension, they will be sampled on the same set of values.

Note that the framework only supports the (simplest and most common) "independent axes" case of the `CF conventions`<sup>142</sup>. In particular, the framework only deals with data on lat-lon grids.

## Example

```
"dimensions": {
  "lat": {
    "standard_name": "latitude",
    "units": "degrees_N",
    "range": [-90, 90],
    "need_bounds": false
  },
  "lon": {
    "standard_name": "longitude",
    "units": "degrees_E",
    "range": [-180, 180],
    "need_bounds": false
  },
  "plev": {
    "standard_name": "air_pressure",
    "units": "hPa",
    "positive": "down",
    "need_bounds": false
  },
  "time": {
    "standard_name": "time",
```

(continues on next page)

<https://www.unidata.ucar.edu/software/netcdf/workshops/2010/datamodels/NcDims.html>  
[http://cfconventions.org/Data/cf-conventions/cf-conventions-1.8/cf-conventions.html#\\_independent\\_latitude\\_longitude\\_vertical\\_and\\_time\\_axes](http://cfconventions.org/Data/cf-conventions/cf-conventions-1.8/cf-conventions.html#_independent_latitude_longitude_vertical_and_time_axes)



```

    "units": "days",
    "calendar": "noleap",
    "need_bounds": false
  }
}

```

## Latitude and Longitude

Required, string. Must be "latitude" and "longitude", respectively.

Optional. String, following syntax of the [UDUnits library](#)<sup>143</sup>. Units the diagnostic expects the dimension to be in. Currently the framework only supports decimal degrees\_north and degrees\_east, respectively.

Array (page 22) (list) of two numbers. Optional. If given, specifies the range of values the diagnostic expects this dimension to take. For example, "range": [-180, 180] for longitude will have the first entry of the longitude variable in each data file be near -180 degrees (not exactly -180, because dimension values are cell midpoints), and the last entry near +180 degrees.

Boolean. Optional: assumed false if not specified. If true, the framework will ensure that bounds are supplied for this dimension, in addition to its midpoint values, following the [CF conventions](#)<sup>144</sup>: the bounds attribute of this dimension will be set to the name of another netCDF variable containing the bounds information.

## Time

Required. Must be "time".

String. Optional, defaults to "day". Units the diagnostic expects the dimension to be in. Currently the diagnostic only supports time axes of the form "<units> since <reference data>", and the value given here is interpreted in this sense (eg. settings this to "day" would accommodate a dimension of the form "[decimal] days since 1850-01-01".)

String, Optional. One of the CF convention [calendars](#)<sup>145</sup> or the string "any". Defaults to "any" if not given. Calendar convention used by your diagnostic. Only affects the number of days per month.

Boolean. Optional: assumed false if not specified. If true, the framework will ensure that bounds are supplied for this dimension, in addition to its midpoint values, following the [CF conventions](#)<sup>146</sup>: the bounds attribute of this dimension will be set to the name of another netCDF variable containing the bounds information.

## Z axis (height/depth, pressure, ...)

Required, string. [Standard name](#)<sup>147</sup> of the variable as defined by the [CF conventions](#)<sup>148</sup>, or a commonly used synonym as employed in the CMIP6 MIP tables.

Optional. String, following syntax of the [UDUnits library](#)<sup>149</sup>. Units the diagnostic expects the dimension to be in. If not provided, the framework will assume CF convention [canonical units](#)<sup>150</sup>.

String, required. Must be "up" or "down", according to the [CF conventions](#)<sup>151</sup>. A pressure axis is always "down" (increasing values are closer to the center of the earth), but this is not set automatically.

Boolean. Optional: assumed false if not specified. If true, the framework will ensure that bounds are supplied for this dimension, in addition to its midpoint values, following the [CF conventions](#)<sup>152</sup>: the bounds attribute of this dimension will be set to the name of another netCDF variable containing the bounds information.

<https://www.unidata.ucar.edu/software/udunits/udunits-2.0.4/udunits2lib.html#Syntax>  
<http://cfconventions.org/Data/cf-conventions/cf-conventions-1.8/cf-conventions.html#cell-boundaries>  
<http://cfconventions.org/Data/cf-conventions/cf-conventions-1.8/cf-conventions.html#calendar>  
<http://cfconventions.org/Data/cf-conventions/cf-conventions-1.8/cf-conventions.html#cell-boundaries>  
<http://cfconventions.org/Data/cf-standard-names/72/build/cf-standard-name-table.html>  
<http://cfconventions.org/>  
<https://www.unidata.ucar.edu/software/udunits/udunits-2.0.4/udunits2lib.html#Syntax>  
<http://cfconventions.org/Data/cf-standard-names/current/build/cf-standard-name-table.html>  
[http://cfconventions.org/faq.html#vertical\\_coords\\_positive\\_attribute](http://cfconventions.org/faq.html#vertical_coords_positive_attribute)  
<http://cfconventions.org/Data/cf-conventions/cf-conventions-1.8/cf-conventions.html#cell-boundaries>

## Other dimensions (wavelength, ...)

Required, string. [Standard name](#)<sup>153</sup> of the variable as defined by the [CF conventions](#)<sup>154</sup>, or a commonly used synonym as employed in the CMIP6 MIP tables.

Optional. String, following syntax of the [UDUnits library](#)<sup>155</sup>. Units the diagnostic expects the dimension to be in. If not provided, the framework will assume CF convention [canonical units](#)<sup>156</sup>.

Boolean. Optional: assumed `false` if not specified. If `true`, the framework will ensure that bounds are supplied for this dimension, in addition to its midpoint values, following the [CF conventions](#)<sup>157</sup>; the `bounds` attribute of this dimension will be set to the name of another netCDF variable containing the bounds information.

## Varlist section

This section is an [object](#) (page 23) contains properties that apply to the model variables your diagnostic needs for its analysis. “Dimensions” are meant in the sense of the netCDF [data model](#)<sup>158</sup>; informally, they are the “independent variables” whose values are being computed as a function of the values stored in the dimensions.

Each entry corresponds to a distinct data file (or set of files, if `multi_file_ok` is `true`) downloaded by the framework. If your framework needs the same physical quantity sampled with different properties (eg. slices of a variable at multiple pressure levels), specify them as multiple entries.

## Varlist entry example

```
"u500": {
  "standard_name": "eastward_wind",
  "path_variable": "U500_FILE",
  "units": "m s-1",
  "dimensions" : ["time", "lat", "lon"],
  "dimensions_ordered": true,
  "scalar_coordinates": {"pressure": 500},
  "requirement": "optional",
  "alternates": ["another_variable_name", "a_third_variable_name"]
}
```

## Varlist entry properties

The key in a varlist key-value pair is the name your diagnostic uses to refer to this variable (and must be unique). The value of the key-value pair is an [object](#) (page 23) containing properties specific to that variable:

String, required. [Standard name](#)<sup>159</sup> of the variable as defined by the [CF conventions](#)<sup>160</sup>, or a commonly used synonym as employed in the CMIP6 MIP tables (eg. “ua” instead of “eastward\_wind”).

String, required. Name of the shell environment variable the framework will set with the location of this data. See the environment variable [documentation](#) (page 30) for details.

If `multi_file_ok` is `false`, `<path_variable>` will be set to the absolute path to the netcdf file containing this variable’s data.

If `multi_file_ok` is `true`, `<path_variable>` will be a single path or a colon-separated list of paths to the files containing this data. Files will be listed in chronological order.

<http://cfconventions.org/Data/cf-standard-names/72/build/cf-standard-name-table.html>  
<http://cfconventions.org/>  
<https://www.unidata.ucar.edu/software/udunits/udunits-2.0.4/udunits2lib.html#Syntax>  
<http://cfconventions.org/Data/cf-standard-names/current/build/cf-standard-name-table.html>  
<http://cfconventions.org/Data/cf-conventions/cf-conventions-1.8/cf-conventions.html#cell-boundaries>  
<https://www.unidata.ucar.edu/software/netcdf/workshops/2010/datamodels/NcVars.html>  
<http://cfconventions.org/Data/cf-standard-names/72/build/cf-standard-name-table.html>  
<http://cfconventions.org/>

If the variable is listed as "optional" or "alternate" or has alternate variables listed, <path\_variable> will be defined but set to the empty string if the framework couldn't obtain this data from the data source. Your diagnostic should test for this possibility.

Boolean. Optional: assumed false if not specified. If true, the framework will ignore the model's naming conventions and only look for a variable with a name matching the key of this entry, regardless of what model or data source the framework is using. The only use case for this setting is to give diagnostics the ability to request data that falls outside the CF conventions: in general, you should rely on the framework to translate CF standard names to the native field names of the model being analyzed.

Optional. String, following syntax of the [UDUnits library](#)<sup>161</sup>. Units the diagnostic expects the variable to be in. If not provided, the framework will assume CF convention [canonical units](#)<sup>162</sup>.

Required. List of strings, which must be selected the keys of entries in the [dimensions](#) (page 26) section. Dimensions of the array containing the variable's data. Note that the framework will not reorder dimensions (transpose) unless `dimensions_ordered` is additionally set to true.

Boolean. Optional: assumed false if not specified. If true, the framework will ensure that the dimensions of this variable's array are given in the same order as listed in `dimensions`. If set to false, your diagnostic is responsible for handling arbitrary dimension orders: eg. it should not assume that 3D data will be presented as (time, lat, lon). If given here, overrides the values set globally in the data section (see [description](#) (page 26) there).

`object` (page 23), optional. This implements what the CF conventions refer to as "[scalar coordinates](#)<sup>163</sup>", with the use case here being the ability to request slices of higher-dimensional data. For example, the snippet at the beginning of this section shows how to request the u component of wind velocity on a 500 mb pressure level.

keys are the key (name) of an entry in the [dimensions](#) (page 26) section.

values are a single number (integer or floating-point) corresponding to the value of the slice to extract. Units of this number are taken to be the `units` property of the dimension named as the key.

In order to request multiple slices (eg. wind velocity on multiple pressure levels, with each level saved to a different file), create one varlist entry per slice.

[Unit-ful quantities](#) (page 22). Optional. Time frequency at which the variable's data is provided. If given here, overrides the values set globally in the data section (see [description](#) (page 26) there).

String. Optional: assumed "required" if not specified. One of three values:

"required": variable is necessary for the diagnostic's calculations. If the data source doesn't provide the variable (at the requested frequency, etc., for the user-specified analysis period) the framework will not run the diagnostic, but will instead log an error message explaining that the lack of this data was at fault.

"optional": variable will be supplied to the diagnostic if provided by the data source. If not available, the diagnostic will still run, and the `path_variable` for this variable will be set to the empty string. The diagnostic is responsible for testing the environment variable for the existence of all optional variables.

"alternate": variable is specified as an alternate source of data for some other variable (see next property). The framework will only query the data source for this variable if it's unable to obtain one of the other variables that list it as an alternate.

[Array](#) (page 22) (list) of strings, which must be keys (names) of other variables. Optional: if provided, specifies an alternative method for obtaining needed data if this variable isn't provided by the data source.

If the data source provides this variable (at the requested frequency, etc., for the user-specified analysis period), this property is ignored.

If this variable isn't available as requested, the framework will query the data source for all of the variables listed in this property. If all of the alternate variables are available, the diagnostic will be run; if any are missing it will be skipped. Note that, as currently implemented, only one set of alternates may be given (no "plan B", "plan C", etc.)

---

<https://www.unidata.ucar.edu/software/udunits/udunits-2.0.4/udunits2lib.html#Syntax>  
<http://cfconventions.org/Data/cf-standard-names/current/build/cf-standard-name-table.html>  
<http://cfconventions.org/Data/cf-conventions/cf-conventions-1.8/cf-conventions.html#scalar-coordinate-variables>

## MDTF Environment variables

This page describes the environment variables that the framework will set for your diagnostic when it's run.

### Overview

The MDTF framework can be viewed as a “wrapper” for your code that handles data fetching and munging. Your code communicates with this wrapper in two ways:

The [settings file](#) (page 17) is where your code talks to the framework: when you write your code, you document what model data your code uses (not covered on this page, follow the link for details).

When your code is run, the framework talks to it by setting shell [environment variables](#)<sup>164</sup> containing paths to the data files and other information specific to the run. The framework communicates all runtime information this way: this is in order to 1) pass information in a language-independent way, and 2) to make writing diagnostics easier (you don't need to parse command-line settings).

Note that environment variables are always strings. Your script will need to cast non-text data to the appropriate type (eg. the bounds of the analysis time period, FIRSTYR, LASTYR, will need to be converted to integers.)

Also note that names of environment variables are case-sensitive.

### Paths

Path to the top-level directory containing any observational or reference data you've provided as the author of your diagnostic. Any data your diagnostic uses that doesn't come from the model being analyzed should go here (ie, you supply it to the framework maintainers, they host it, and the user downloads it when they install the framework). The framework will ensure this is copied to a local filesystem when your diagnostic is run, but this directory should be treated as read-only.

Path to the top-level directory containing your diagnostic's source code. This will be of the form `.../MDTF-diagnostics/diagnostics/<your POD's name>`. This can be used to call sub-scripts from your diagnostic's driver script. This directory should be treated as read-only.

Path to your diagnostic's working directory, which is where all output data should be written (as well as any temporary files).

The framework creates the following subdirectories within this directory:

`$WK_DIR/obs/PS` and `$WK_DIR/model/PS`: All output plots produced by your diagnostic should be written to one of these two directories. Only files in these locations will be converted to bitmaps for HTML output.

`$WK_DIR/obs/netCDF` and `$WK_DIR/model/netCDF`: Any output data files your diagnostic wants to make available to the user should be saved to one of these two directories.

### Model run information

User-provided label describing the run of model data being analyzed.

Four-digit years describing the analysis period.

---

[https://en.wikipedia.org/wiki/Environment\\_variable](https://en.wikipedia.org/wiki/Environment_variable)

## Locations of model data files

These are set depending on the data your diagnostic requests in its [settings file](#) (page 17). Refer to the examples below if you're unfamiliar with how that file is organized.

Each variable listed in the `varlist` section of the settings file must specify a `path_variable` property. The value you enter there will be used as the name of an environment variable, and the framework will set the value of that environment variable to the absolute path to the file containing data for that variable.

From a diagnostic writer's point of view, this means all you need to do here is replace paths to input data that are hard-coded or passed from the command line with calls to read the value of the corresponding environment variable.

If the framework was not able to obtain the variable from the data source (at the requested frequency, etc., for the user-specified analysis period), this variable will be set equal to the empty string. Your diagnostic is responsible for testing for this possibility for all variables that are listed as `optional` or have alternates listed (if a required variable without alternates isn't found, your diagnostic won't be run.)

If `multi_file_ok` is set to `true` in the settings file, this environment variable may be a list of paths to multiple files in chronological order, separated by colons. For example, `/dir/precip_1980_1989.nc:/dir/precip_1990_1999.nc:/dir/precip_2000_2009.nc` for an analysis period of 1980-2009.

## Names of variables and dimensions

These are set depending on the data your diagnostic requests in its [settings file](#) (page 17). Refer to the examples below if you're unfamiliar with how that file is organized.

If `<key>` is the name of the key labeling the key:value entry for this dimension, the framework will set an environment variable named `<key>_dim` equal to the name that dimension has in the data files it's providing.

If `rename_dimensions` is set to `true` in the settings file, this will always be equal to `<key>`. If `rename_dimensions` is `false`, this will be whatever the model or data source's native name for this dimension is, and your diagnostic should read the name from this variable. Your diagnostic should only use hard-coded names for dimensions if `rename_dimensions` is set to `true` in its [settings file](#) (page 22).

If `<key>` be the name of the key labeling the key:value entry for this variable in the `varlist` section, the framework will set an environment variable named `<key>_var` equal to the name that variable has in the data files it's providing.

If `rename_variables` is set to `true` in the settings file, this will always be equal to `<key>`. If `rename_variables` is `false`, this will be whatever the model or data source's native name for this variable is, and your diagnostic should read the name from this variable. Your diagnostic should only use hard-coded names for dimensions if `rename_dimensions` is set to `true` in its [settings file](#) (page 22).

## Simple example

We only give the relevant parts of the [settings file](#) (page 22) below.

```
"data": {
  "rename_dimensions": false,
  "rename_variables": false,
  "multi_file_ok": false,
  ...
},
"dimensions": {
  "lat": {
    "standard_name": "latitude",
    ...
  },
  "lon": {
    "standard_name": "longitude",
    ...
  },
  "time": {
    "standard_name": "time",
    ...
  }
}
```

(continues on next page)

```

},
"varlist": {
  "pr": {
    "standard_name": "precipitation_flux",
    "path_variable": "PR_FILE"
  }
}

```

The framework will set the following environment variables:

1. lat\_dim: Name of the latitude dimension in the model's native format (because rename\_dimensions is false).
2. lon\_dim: Name of the longitude dimension in the model's native format (because rename\_dimensions is false).
3. time\_dim: Name of the time dimension in the model's native format (because rename\_dimensions is false).
4. pr\_var: Name of the precipitation variable in the model's native format (because rename\_variables is false).
5. PR\_FILE: Absolute path to the file containing pr data, eg. /dir/precip.nc.

### More complex example

Let's elaborate on the previous example, and assume that the diagnostic is being called on model that provides precipitation\_flux but not convective\_precipitation\_flux.

```

"data": {
  "rename_dimensions": true,
  "rename_variables": false,
  "multi_file_ok": true,
  ...
},
"dimensions": {
  "lat": {
    "standard_name": "latitude",
    ...
  },
  "lon": {
    "standard_name": "longitude",
    ...
  },
  "time": {
    "standard_name": "time",
    ...
  }
},
"varlist": {
  "prc": {
    "standard_name": "convective_precipitation_flux",
    "path_variable": "PRC_FILE",
    "alternates": ["pr"]
  },
  "pr": {
    "standard_name": "precipitation_flux",
    "path_variable": "PR_FILE"
  }
}

```

Comparing this with the previous example:

lat\_dim, lon\_dim and time\_dim will be set to "lat", "lon" and "time", respectively, because rename\_dimensions is true. The framework will have renamed these dimensions to have these names in all data files provided to the diagnostic.

prc\_var and pr\_var will be set to the model's native names for these variables. Names for all variables are always set, regardless of which variables are available from the data source.

In this example, PRC\_FILE will be set to '', the empty string, because it wasn't found.



## MDTF Developer's Walkthrough, Release 3.0 beta 2

PR\_FILE will be set to /dir/precip\_1980\_1989.nc:/dir/precip\_1990\_1999.nc:/dir/precip\_2000\_2009.nc, because multi\_file\_ok was set to true.

**Acknowledgements** Development of this code framework for process-oriented diagnostics was supported by the National Oceanic and Atmospheric Administration<sup>165</sup> (NOAA) Climate Program Office Modeling, Analysis, Predictions and Projections<sup>166</sup> (MAPP) Program (grant # NA18OAR4310280). Additional support was provided by University of California Los Angeles<sup>167</sup>, the Geophysical Fluid Dynamics Laboratory<sup>168</sup>, the National Center for Atmospheric Research<sup>169</sup>, Colorado State University<sup>170</sup>, Lawrence Livermore National Laboratory<sup>171</sup> and the US Department of Energy<sup>172</sup>.

Many of the process-oriented diagnostics modules (PODs) were contributed by members of the NOAA Model Diagnostics Task Force<sup>173</sup> under MAPP support. Statements, findings or recommendations in these documents do not necessarily reflect the views of NOAA or the US Department of Commerce.

### Disclaimer

This repository is a scientific product and is not an official communication of the National Oceanic and Atmospheric Administration, or the United States Department of Commerce. All NOAA GitHub project code is provided on an 'as is' basis and the user assumes responsibility for its use. Any claims against the Department of Commerce or Department of Commerce bureaus stemming from the use of this GitHub project will be governed by all applicable Federal law. Any reference to specific commercial products, processes, or services by service mark, trademark, manufacturer, or otherwise, does not constitute or imply their endorsement, recommendation or favoring by the Department of Commerce. The Department of Commerce seal and logo, or the seal and logo of a DOC bureau, shall not be used in any manner to imply endorsement of any commercial product or activity by DOC or the United States Government.

---

<https://www.noaa.gov/>  
<https://cpo.noaa.gov/Meet-the-Divisions/Earth-System-Science-and-Modeling/MAPP>  
<https://www.ucla.edu/>  
<https://www.gfdl.noaa.gov/>  
<https://ncar.ucar.edu/>  
<https://www.colostate.edu/>  
<https://www.llnl.gov/>  
<https://www.energy.gov/>  
<https://cpo.noaa.gov/Meet-the-Divisions/Earth-System-Science-and-Modeling/MAPP/MAPP-Task-Forces/Model-Diagnostics-Task-Force>