
Multiscale Genomics Documentation

Release 0.1

Mark McDowall

Nov 05, 2018

Contents:

1	VRE	3
2	Development APIs	5
3	Workflows	7
4	RESTful APIs	9
4.1	HOWTO	9
4.2	MuG Coding Guidelines	36
4.3	License	41
5	Indices and tables	45

Virtual Research Environment- Supporting the 3D/4D genomics community with tools to integrate navigation from sequence to 3D/4D chromatin dynamics data.

Full details of the project can be found on the [MuG website](#).

From this site there are links to all of the documentation for each of the repositories within the [MuG GitHub](#). Each set of documents contains details about the installation and usage both for an end user and for developers creating workflows.

If you want to develop tools and workflows that can run within the VRE please checkout the [HOWTO](#) section of the site. This lists how to write a workflow, tool, the configuration files and testing that the workflow works. Please also read about how to apply the Apache2.0 license to your code and a document on Coding Standards that you should adhere to to ensure the best chance of your workflow being integrated with the VRE as smoothly as possible.

CHAPTER 1

VRE

VRE MuG Virtual Research Environment. MuG VRE is a web application aimed to allow MuG users to access MuG data and explore and exploit it together with its own data via a selection of tools and visualizers. It is written in PHP, HTML and Javascript.

[GitHub VRE](#)

CHAPTER 2

Development APIs

mg-dm-api Data Management API. This API tracks files within the VRE and contains meta data about how the file was generated with access to the file geneology.

[GitHub mg-dm-api](#)

mg-tool-api Tool API. This API provides the interface between the pyCOMPSs architecture and the tool. It provides a standard way for all tools to be wrapped to allow for a common interface layer.

[GitHub mg-tool-api](#)

mg-process-fastq Workflows for processing FASTQ data. These workflows can handle ChIP-seq, MNase-Seq, RNA-Seq and Whole Genome BiSulphate Sequencing (WGBS). There are also scripts for generating the initial set of indexes for given genome assemblies. There are also workflows for processing Hi-C data to generate adjacency matrices and calculate TAD regions

[GitHub mg-process-fastq](#)

mg-process-files Workflows for processing results files into an indexed form for use in a RESTful interface.

[GitHub mg-process-files](#)

mg-rest-service The root RESTful server. This provides links to the main root end points. Each end point is a provides a unique function within the defined URL so that the service as a whole appears seamless to the end user.

[GitHub mg-rest-service](#)

mg-rest-dm RESTful interface to the DM API along with end-points to manage the stored files and track the relevant metadata.

[GitHub mg-rest-dm](#)

mg-rest-file RESTful interface to the DM API along with end-points to servicing out regions from basic file based data, such as Bed, Wig and TSV files.

[GitHub mg-rest-file](#)

mg-rest-adjacency Interface for RESTfully querying adjacency matrices generated by the TADbit workflows developed in mg-process-fastq.

[GitHub mg-rest-adjacency](#)

mg-rest-3d Interface for RESTfully querying 3D models generated by the TADbit workflows developed in mg-process-files.

[GitHub mg-rest-3d](#)

mg-rest-util Set of common functions that are required by the RESTful interfaces for interacting with the DM API.

[GitHub mg-rest-util](#)

4.1 HOWTO

4.1.1 Development Checklist

This document describes the standard work flow to help developers when creating a new tool or pipeline. The purpose is to aid the developer in the most efficient way for integrating a new tool or pipeline and ensure that all steps have

been addressed so that they have a ready to deploy Tool and Pipeline within the MuG VRE.

Note: If you are adding a new tool and pipeline to an already existing repository then you can skip ahead and concentrate on **steps 1 to 6**.

Note: If you are adding just a new pipeline that just integrates already existing tools then you need to look at **steps 3 to 6**.

0 - Copy mg-process-test from GitHub

0.0 - Create an empty repository

In GitHub create a blank repository with no README, license or *.gitignore* file. These files will be inherited from the mg-process-test file. For this example it will be called *mg-process-test1.git*.

0.1 - Copy the mg-process-test repository

From GitHub take a copy of the mg-process-test repository:

```
git clone --depth 1 -b master https://github.com/Multiscale-Genomics/mg-process-test

rm -rf mg-process-test/.git
mv mg-process-test mg-process-test1
cd mg-process-test1

git init
git add .
git commit -m 'Initial commit'

git remote add origin https://github.com/<USERNAME>/mg-process-test1.git
git remote -v

git push origin master
```

From here you can then customise the following files to match your new repository:

- *README.md*
- *NOTICE*
- *setup.py*
- *__init__.py*
- *docs/conf.py*

The files in *docs* contain boilerplate data that matches the processes and tools already in the repository, so should be updated as you add new pipelines and tools.

1 - Setup Your Python environment

1.1 - pyenv

This is required for managing the version of Python and the installation environment for the Python modules so that they can be installed in the user space.

```

1 git clone https://github.com/pyenv/pyenv.git ~/.pyenv
2 echo 'export PYENV_ROOT="$HOME/.pyenv"' >> ~/.bash_profile
3 echo 'export PATH="$PYENV_ROOT/bin:$PATH"' >> ~/.bash_profile
4 echo 'eval "$(pyenv init -)"' >> ~/.bash_profile
5
6 # Add the .bash_profile to your .bashrc file
7 echo 'source ~/.bash_profile' >> ~/.bashrc
8
9 git clone https://github.com/pyenv/pyenv-virtualenv.git ${PYENV_ROOT}/plugins/pyenv-
  ↳virtualenv
10
11 pyenv install 2.7.14
12 pyenv virtualenv 2.7.14 mg-process-test

```

1.2 - Install Tool API

```

1 pyenv activate mg-process-test
2 pip install git+https://github.com/Multiscale-Genomics/mg-tool-api.git

```

2 - Create a Tool

See the [HOWTO - Tools](#) for details about writing a tool and [HOWTO - Test Your Code](#) about how to write relevant tests

2.1 - Tool Development

Using the *testTool.py* script as a template, create you new tool.

Checklist

1. There is a license at the header of the script
2. Documentation for each function.
3. Code matches the PEP8 standard (by running pylint).
4. Tool has been added to *docs/tool.rst*

3 - Create a Test to run the Tool

3.1 - Test Dataset

Create a small test dataset that can be used when testing the code. This should match the input file type required by the Tool.

When the tool has been run the output for the test datasets should provide a valid result. For example if wrapping a peak caller there should be enough of the genome selected and matching reads that when aligned and the peak caller analyses the alignments it should generate results similar to the original for that region.

Once the datasets have been generated the procedure for how the test sets were created should be documented in a new “NNN.rst” file. This should contain the source of the data, publications, where the files were downloaded from and how the data was handled so that this can be repeated if the datasets need to be regenerated or changed at a later stage. This file should then be linked into the rest of the documentation, this is usually done by linking the file in the table of contents block in the index page.

3.2 - Test Scripts

Create a script that uses pytest to check that the required output files have been generated and are not empty. Other tests can be added here if there are other aspects that should be tested. Examples could include testing if a JSON object has the expected parameters.

Checklist

1. There is a test to run each single tool
2. There is a license header in each test script
3. All functions in the test script are fully documented with details about how to run the test or if other tests need to be run first
4. Test dataset generation has been fully documented and linked to the index.rst file
5. Any scripts developed to create the datasets are stored in *scripts/* and have matching license headers and documentation
6. All code matches the PEP8 standard (by running pylint).
7. All new tests have been added to TravisCI
8. All tests are passing
9. Ensure that the output of running the tests matches what you would expect

4 - Create a Pipeline

See the [HOWTO - Pipelines](#) for details about writing a pipeline and [HOWTO - Test Your Code](#) about how to write relevant tests.

4.1 - Pipeline Development

Using the *process_test.py* script as a template, create a pipeline to accept the configuration and input JSON files that describe the parameters and files to get passed into the pipeline. The pipeline should manage the passing of file locations and parameters to each of the tools.

4.2 - Create a Test to run the Pipeline

Create a script that uses pytest to check that the required input files and configuration parameters are accepted by the pipeline and the relevant output files have been generated and are not empty. Other tests can be added to be more comprehensive.

The pipeline is running tools developed as part of part 1, so there should be no need for creating new datasets.

4.3 - Create test config and input JSON files

JSON files need to be created that duplicate what would be the expected input coming from the VRE and saved in the *tests/json/* directory of the repository. Example files can be found in the [HOWTO on Configuration](#). There are also examples of these files in *mg-process-test* in the *test/json/*. These files allow a user to run the sample datasets from the command line either on their own computer or on one with (py)COMPSs installed.

Checklist

1. There is a license in the header of all pipelines and tests
2. There is a test to run each pipeline
3. There is documentation for all functions in the pipeline script and test script
4. Update docs/pipelines.rst to include documentation and links to the new pipeline to import all function documentation
5. All code matches the PEP8 standard (by running pylint).
6. All new tests have been added to TravisCI
7. All tests are passing
8. Ensure that the output of running the tests matches what you would expect
9. The script can be run from the command line

5 - VRE JSON Configuration

See the [HOWTO - Configuration Files](#) for details about writing a MuG VRE JSON configuration files.

Checklist

1. Ensure that there is a JSON configuration file present in the *tool_json* for each pipeline.

6 - Installation Documentation

Checklist

1. Make sure that *setup.py*, *setup.cfg* and *requirements.txt* are updated with any new packages required for installation
2. Update *docs/install.rst* if there is any external software that is required by tool or pipeline along with the required command to install that software

7 - COMPSs testing

Now that you have a functional pipeline and tool it now needs to be tested within a COMPSs environment. Download the latest version of the [COMPSs virtual machine](#) from the BSC website.

Checklist

1. Was it possible to install everything based on the installation scripts and documentation?
2. Do all the test scripts pass when they are run?
3. When the test scripts have run do you get the expected results?
4. Can the pipeline be run using the “runcomps” command?

8 - Hook up your repository for continuous integration

Now that you have a fully documented pipeline, with tests it is possible to hook up your GitHub repository with ReadTheDocs.org, TravisCI.org and Landscape.io. These services will automatically build you documentation, run the tests and check the compliance of the code with that of PEP8 respectively.

It is possible to login to each service using your GitHub account and link the repository.

Checklist

1. You have your documentation building on ReadTheDocs.org
2. You have your test scripts running on TravisCI and passing
3. Your code is being continually analysed by Landscape.io

9 - Congratulations

You now have a pipeline that could be integrated into the MuG VRE.

4.1.2 HOWTO - Tools

This document provides a tutorial for the creation of a tool that can be used within a pipeline within the MuG VRE. All functions should be wrapped up as a tool, this then allows for the tools to be easily reused by other pipelines and also deployed onto the compute cluster.

The Tool is the core element when it comes to running a program of function within the COMPSs environment. It defines the procedures that need to happen to prepare the data along with the function that is parallelised to run over the chunks of data provided. A function can be either a piece of code that is written in python or an external package that is run with given chunks of data or defined parameters. The results are then returned to the calling function for merging.

All functions contain at least a *run(self)* function which is called by the pipeline. The run function takes the input files (list), defined output files (list) and relevant metadata (dict). Returned by the run function is a list containing a list of the output files as the first object and a list of metadata dict objects as the second element.

Repository Structure

All tools should be placed within the *tools* directory within the package.

Basic Tool

This is a test tool that takes an input file, then counts the number of characters in that file and then prints the result to a second file. The matching code can be found in the GitHub repository [mg-process-test](#). The file is called *testTool.py*.

```

1  """
2  .. License and copyright agreement statement
3  """
4  from __future__ import print_function
5
6  from utils import logger
7
8  try:
9      from pycompss.api.parameter import FILE_IN, FILE_OUT
10     from pycompss.api.task import task
11     from pycompss.api.api import compss_wait_on
12 except ImportError:
13     logger.warn("[Warning] Cannot import \"pycompss\" API packages.")
14     logger.warn("          Using mock decorators.")
15
16     from utils.dummy_pycompss import FILE_IN, FILE_OUT # pylint: disable=ungrouped-
↪imports
17     from utils.dummy_pycompss import task # pylint: disable=ungrouped-imports
18     from utils.dummy_pycompss import compss_wait_on # pylint: disable=ungrouped-
↪imports
19
20 from basic_modules.tool import Tool
21 from basic_modules.metadata import Metadata
22
23 # -----
24
25 class testTool(Tool):
26     """
27     Tool for writing to a file
28     """
29
30     def __init__(self, configuration=None):
31         """
32         Init function
33         """
34         print("Test writer")
35         Tool.__init__(self)
36
37         if configuration is None:
38             configuration = {}
39
40         self.configuration.update(configuration)
41
42     @task(returns=bool, file_in_loc=FILE_IN, file_out_loc=FILE_OUT, isModifier=False)
43     def test_writer(self, file_loc):
44         """
45         Count the number of characters in a file and return a file with the count
46
47         Parameters
48         -----
49         file_in_loc : str
50             Location of the input file
51         file_out_loc : str

```

(continues on next page)

(continued from previous page)

```

52         Location of an output file
53
54     Returns
55     -----
56     bool
57         Writes to the file, which is returned by pyCOMPSSs to the defined location
58     """
59     try:
60         with open(file_loc, "w") as file_handle:
61             file_handle.write("This is the test writer")
62     except IOError as error:
63         logger.fatal("I/O error({0}): {1}".format(error.errno, error.strerror))
64         return False
65
66     return True
67
68 def run(self, input_files, input_metadata, output_files):
69     """
70     The main function to run the test_writer tool
71
72     Parameters
73     -----
74     input_files : dict
75         List of input files - In this case there are no input files required
76     input_metadata: dict
77         Matching metadata for each of the files, plus any additional data
78     output_files : dict
79         List of the output files that are to be generated
80
81     Returns
82     -----
83     output_files : dict
84         List of files with a single entry.
85     output_metadata : dict
86         List of matching metadata for the returned files
87     """
88
89     results = self.test_writer(
90         input_files["input_file_location"],
91         output_files["output_file_location"]
92     )
93
94     results = compss_wait_on(results)
95
96     if results is False:
97         logger.fatal("Test Writer: run failed")
98         return {}, {}
99
100     output_metadata = {
101         "test": Metadata(
102             data_type="<data_type>",
103             file_type="txt",
104             file_path=output_files["test"],
105             sources=[input_metadata["input_file_location"].file_path],
106             taxon_id=input_metadata["input_file_location"].taxon_id,
107             meta_data={
108                 "tool": "testTool"

```

(continues on next page)

(continued from previous page)

```

109         }
110     )
111 }
112
113     return (output_files, output_metadata)

```

This is the simplest case of a Tool that will run a function within the COMPSS environment. The run function takes the input files, if the output files are defined it can use those as the output locations and any relevant metadata. The locations of the output files can also be defined within the run function as sometimes functions can generate a large number of files that are not always easy to define up front if the Tool is being run as part of the VRE or as part of a larger pipeline.

The run function then calls the *test_writer* function. This uses the python decorator syntax to highlight that it is a function that can be run in parallel to pyCOMPSs library. The *task* decorator is used to define the list of files and parameters that need to be passed to the function. It also requires a list of the files that are to be returned. As such the most common types will be *FILE_IN*, *FILE_OUT*, *FILE_INOUT*.

The *__init__* function is important as it loads the configuration parameters into the class from the VRE. In this case there are no parameters used, but these can be parameters required for the tool that has been wrapped by the code.

Decorators can also be used to define the resources that are required by function. They can be used to define a set of machines that the task should be run on, required CPU capacity or the amount of RAM that is required by the task. Defining these parameters helps the COMPSS infrastructure correctly allocate jobs so that they are able to run as soon as the resources allow and prevent the job failing by being run on a machine that does not have the correct resources.

Further details about COMPSS and pyCOMPSs can be found at the BSC website along with specific tutorials about how to write functions that can utilise the full power of COMPSS.

pyCOMPSs within the Tool

When importing the pyCOMPSs modules it is important to provide access to the dummy_pycompss decorators as well. This will allow scripts to be run on computers where COMPSs has not been installed.

Practical Example

Now that we know the basics it is possible to apply this to writing a tool that can run and perform a real operation within the cluster.

Here is a tool that uses BWA to index a genome sequence file that has been saved in FASTA format.

The run function takes the input FASTA file, from this it generates a list of the locations of the output files. The input file and output files are passed to the *bwa_indexer* function. The files do not need to be listed in the return call so True is fine. COMPSS handles the passing back of the files to the run function. The run function then returns the output files to the pipeline or the VRE.

```

1  from __future__ import print_function
2
3  import os
4  import shlex
5  import shutil
6  import subprocess
7  import sys
8  import tarfile
9
10 from utils import logger

```

(continues on next page)

(continued from previous page)

```

11
12 try:
13     if hasattr(sys, '_run_from_cmdl') is True:
14         raise ImportError
15     from pycompss.api.parameter import FILE_IN, FILE_OUT
16     from pycompss.api.task import task
17     from pycompss.api.api import compss_wait_on
18 except ImportError:
19     logger.warn("[Warning] Cannot import \"pycompss\" API packages.")
20     logger.warn("        Using mock decorators.")
21
22     from utils.dummy_pycompss import FILE_IN, FILE_OUT # pylint: disable=ungrouped-
↳ imports
23     from utils.dummy_pycompss import task # pylint: disable=ungrouped-imports
24     from utils.dummy_pycompss import compss_wait_on # pylint: disable=ungrouped-imports
25
26 from basic_modules.tool import Tool
27 from basic_modules.metadata import Metadata
28
29 # -----
30
31 class bwaIndexerTool(Tool):
32     """
33     Tool for running indexers over a genome FASTA file
34     """
35
36     def __init__(self, configuration=None):
37         """
38         Init function
39         """
40         print("BWA Indexer")
41         Tool.__init__(self)
42
43         if configuration is None:
44             configuration = {}
45
46         self.configuration.update(configuration)
47
48     def bwa_index_genome(self, genome_file):
49         """
50         Create an index of the genome FASTA file with BWA. These are saved
51         alongside the assembly file. If the index has already been generated
52         then the locations of the files are returned
53
54         Parameters
55         -----
56         genome_file : str
57             Location of the assembly file in the file system
58
59         Returns
60         -----
61         amb_file : str
62             Location of the amb file
63         ann_file : str
64             Location of the ann file
65         bwt_file : str
66             Location of the bwt file

```

(continues on next page)

(continued from previous page)

```

67     pac_file : str
68         Location of the pac file
69     sa_file : str
70         Location of the sa file
71
72     """
73     command_line = 'bwa index ' + genome_file
74
75     amb_name = genome_file + '.amb'
76     ann_name = genome_file + '.ann'
77     bwt_name = genome_file + '.bwt'
78     pac_name = genome_file + '.pac'
79     sa_name = genome_file + '.sa'
80
81     if os.path.isfile(bwt_name) is False:
82         args = shlex.split(command_line)
83         process = subprocess.Popen(args)
84         process.wait()
85
86     return (amb_name, ann_name, bwt_name, pac_name, sa_name)
87
88 @task(file_loc=FILE_IN, idx_out=FILE_OUT)
89 def bwa_indexer(self, file_loc, idx_out): # pylint: disable=unused-argument
90     """
91     BWA Indexer
92
93     Parameters
94     -----
95     file_loc : str
96         Location of the genome assembly FASTA file
97     idx_out : str
98         Location of the output index file
99
100    Returns
101    -----
102    bool
103    """
104    amb_loc, ann_loc, bwt_loc, pac_loc, sa_loc = self.bwa_index_genome(file_loc)
105
106    # tar.gz the index
107    print("BS - idx_out", idx_out, idx_out.replace('.tar.gz', ''))
108    idx_out_pregz = idx_out.replace('.tar.gz', '.tar')
109
110    index_dir = idx_out.replace('.tar.gz', '')
111    os.mkdir(index_dir)
112
113    idx_split = index_dir.split("/")
114
115    shutil.move(amb_loc, index_dir)
116    shutil.move(ann_loc, index_dir)
117    shutil.move(bwt_loc, index_dir)
118    shutil.move(pac_loc, index_dir)
119    shutil.move(sa_loc, index_dir)
120
121    index_folder = idx_split[-1]
122
123    tar = tarfile.open(idx_out_pregz, "w")

```

(continues on next page)

(continued from previous page)

```
124     tar.add(index_dir, arcname=index_folder)
125     tar.close()
126
127     command_line = 'pigz ' + idx_out_pregz
128     args = shlex.split(command_line)
129     process = subprocess.Popen(args)
130     process.wait()
131
132     return True
133
134     def run(self, input_files, metadata, output_files):
135         """
136         Function to run the BWA over a genome assembly FASTA file to generate
137         the matching index for use with the aligner
138
139         Parameters
140         -----
141         input_files : dict
142             List containing the location of the genome assembly FASTA file
143         meta_data : dict
144         output_files : dict
145             List of output files generated
146
147         Returns
148         -----
149         output_files : dict
150             index : str
151                 Location of the index file defined in the input parameters
152         output_metadata : dict
153             index : Metadata
154                 Metadata relating to the index file
155         """
156         results = self.bwa_indexer(
157             input_files["genome"],
158             output_files["index"]
159         )
160         results = compss_wait_on(results)
161
162         if results is False:
163             logger.fatal("BWA Indexer: run failed")
164             return {}, {}
165
166         output_metadata = {
167             "index": Metadata(
168                 data_type="sequence_mapping_index_bwa",
169                 file_type="TAR",
170                 file_path=output_files["index"],
171                 sources=[metadata["genome"].file_path],
172                 taxon_id=metadata["genome"].taxon_id,
173                 meta_data={
174                     "assembly": metadata["genome"].meta_data["assembly"],
175                     "tool": "bwa_indexer"
176                 }
177             )
178         }
179
180         return (output_files, output_metadata)
```

(continues on next page)

(continued from previous page)

181
182

```
# -----
```

Troubleshooting Common Issues

Program is installed but fails to run

There are several points that need to be checked in this instance:

1. **Is the program available on your \$PATH?** - If not either add it, or place a symlink in a directory that is.
2. **Check that the command that you are running matches the command run by subprocess** - Use the `logger.info()` to print the command and check that it works.
3. **Subprocess runs commands in a sandbox** - The normal way to run `subprocess()` is to use `subprocess.Popen(args)` and pass it a list of arguments that represent the command to be run (as shown in the practical example above). Sometimes this fails as extra environment parameters may be required by the program, in this case it is possible to run the whole command as a single string and tell the subprocess to use a shell:

```
1 command_line = "python --version"
2 process = subprocess.Popen(command_line, shell=True)
3 process.wait()
```

4.1.3 HOWTO - Pipelines

This document is a tutorial about creating pipelines that can be easily integrated into the MuG VRE. The aim of a pipeline is to bring together a number of tools (see [Creating a Tool](#)) and running them as part of a workflow for end to end processing of data.

Each pipeline consists of the main class for the pipeline, a main function for running the class and a section of global code to catch if the pipeline has been run from the command line. All functions should have full documentation describing the function, inputs and outputs. For details about the coding style please consult the [coding style documentation](#).

Example Pipeline

This example code uses the `testTool.py` from the [Creating a Tool](#) tutorial. The matching code can be found in the GitHub repository [mg-process-test](#).

There are 2 ways of calling this function, either directly from another program or via the command line.

```
1 #!/usr/bin/env python
2
3 """
4 .. License and copyright agreement statement
5 """
6 from __future__ import print_function
7
8 # Required for ReadTheDocs
9 from functools import wraps # pylint: disable=unused-import
10
11 import argparse
12
```

(continues on next page)

(continued from previous page)

```

13 from basic_modules.workflow import Workflow
14 from utils import logger
15
16 from mg_process_test.tools.testTool import testTool
17
18 # -----
19
20 class process_test(Workflow):
21     """
22     Functions for demonstrating the pipeline set up.
23     """
24
25     configuration = {}
26
27     def __init__(self, configuration=None):
28         """
29         Initialise the tool with its configuration.
30
31         Parameters
32         -----
33         configuration : dict
34             a dictionary containing parameters that define how the operation
35             should be carried out, which are specific to each Tool.
36         """
37         logger.info("Processing Test")
38         if configuration is None:
39             configuration = {}
40
41         self.configuration.update(configuration)
42
43     def run(self, input_files, metadata, output_files):
44         """
45         Main run function for processing a test file.
46
47         Parameters
48         -----
49         input_files : dict
50             Dictionary of file locations
51         metadata : list
52             Required meta data
53         output_files : dict
54             Locations of the output files to be returned by the pipeline
55
56         Returns
57         -----
58         output_files : dict
59             Locations for the output txt
60         output_metadata : dict
61             Matching metadata for each of the files
62         """
63
64         # Initialise the test tool
65         tt_handle = testTool(self.configuration)
66         tt_files, tt_meta = tt_handle.run(input_files, metadata, output_files)
67
68         return (tt_files, tt_meta)
69

```

(continues on next page)

(continued from previous page)

```

70
71 # -----
72
73 def main_json(config, in_metadata, out_metadata):
74     """
75     Alternative main function
76     -----
77
78     This function launches the app using configuration written in
79     two json files: config.json and input_metadata.json.
80     """
81     # 1. Instantiate and launch the App
82     logger.info("1. Instantiate and launch the App")
83     from apps.jsonapp import JSONApp
84     app = JSONApp()
85     result = app.launch(process_test,
86                        config,
87                        in_metadata,
88                        out_metadata)
89
90     # 2. The App has finished
91     logger.info("2. Execution finished; see " + out_metadata)
92
93     return result
94
95 # -----
96
97 if __name__ == "__main__":
98
99     # Set up the command line parameters
100     PARSER = argparse.ArgumentParser(description="Index the genome file")
101     PARSER.add_argument("--config", help="Configuration file")
102     PARSER.add_argument("--in_metadata", help="Location of input metadata file")
103     PARSER.add_argument("--out_metadata", help="Location of output metadata file")
104     PARSER.add_argument("--local", action="store_const", const=True, default=False)
105
106     # Get the matching parameters from the command line
107     ARGS = PARSER.parse_args()
108
109     CONFIG = ARGS.config
110     IN_METADATA = ARGS.in_metadata
111     OUT_METADATA = ARGS.out_metadata
112     LOCAL = ARGS.local
113
114     if LOCAL:
115         import sys
116         sys._run_from_cmdl = True # pylint: disable=protected-access
117
118     RESULTS = main_json(CONFIG, IN_METADATA, OUT_METADATA)
119     print(RESULTS)

```

Code Walk Through

I'll step through each of the sections of the example code describing what is happening at each point.

Header

This section defines the license and any modules that need to be loaded for the code to run correctly. As a bare minimum is shown in the example with the license, import of the Workflow and Metadata basic_tools and the Data Management (DM) API. Theoretically the pipeline does not have to call a tool, but for completeness this uses the Tool generated as part of the [HOWTO - Tools](#) tutorial.

def main_json()

This is the main entry point into the pipeline. It allows the pipeline to be run either locally or as part of a series of function calls within the VRE.

The *main_json()* function is the primary function of the script and is what initiates running the pipeline. It is from here that the VRE or locally run function will call to with any matching input file, defined output files (is required) and any necessary meta data.

At the bottom of the script the `__main__` is triggered when being run from the command line. It can take in parameters from the command line and pass them to the *main_json()* function. As the VRE is responsible for loading of files into the Data Management (DM) API, if files that are used locally are to be tracked then they should also be loaded into the DM API at this point. For clarity of creating a pipeline this has not been included within the example.

Once *main_json()* has been called it launches the *WorkflowApp()* with the name of the pipeline (*process_test* in this case) along with the input files, output files (if known) and relevant meta data for running the application.

process_test - __init__

Instantiates the pipeline and passes on any configuration data to the WorkFlowApp.

process_test - run

This is a required function which is called by the *main_json()* function. It is responsible for orchestrating the flow of data within the pipeline. The run function ensures that the Tools are initiated correctly and are passed the correct variables. If there are multiple Tools in the pipeline each relying on the output from the previous then the *run()* function is responsible for handing the output files from one tool to the next. At this point the handling of files is managed by the pyCOMPSs API and files only become accessible from the final location once the *run()* function has returned to *main_json()*. If you require the output of a tool locally for launching the next then you need to stream the file out of compss, this can be done with the following snippet:

```
1 if hasattr(sys, '_run_from_cmdl') is True:
2     pass
3 else:
4     with compss_open(intermediate_file_in_compss, "rb") as f_in:
5         with open(local_loc_for_file, "wb") as f_out:
6             f_out.write(f_in.read())
```

This will only work within the COMPSs environment so you will need to test for how your code is getting run.

4.1.4 HOWTO - Documentation

As part of the development of sustainable software it is important that code is well documented to inform developers that need to implement, extend or replace the code about what it does, the inputs, outputs and any dependencies on other software or code. All classes and functions should have matching documentation.

There are 2 key parts of the documentation, the first is for the classes and functions. The documentation should match the PEP8 standard, an example of this is in the [MuG Coding Guidelines](#). The second part is the Architectural Design Record. The ADR should record why key choices have been made, this is especially true if the choices do not match the norm or there has been a major change in a function (addition, removal or completely rewritten). The ADR provides the reasoning behind the code and the documentation string in the functions describe the code. Between them they provide a log of the development of the project.

An example function description should therefore match the following:

```

1  """
2  Assembly Index Manager
3
4  Manages the creation of indexes for a given genome assembly file. If the
5  downloaded file has not been unzipped then it will get unzipped here.
6  There are then 3 indexers that are available including BWA, Bowtie2 and
7  GEM. If the indexes already exist for the given file then the indexing
8  is not rerun.
9
10 Parameters
11 -----
12 file_name : str
13     Location of the assembly FASTA file
14
15 Returns
16 -----
17 dict
18     bowtie : str
19         Location of the Bowtie index file
20     bwa : str
21         Location of the BWA index file
22     gem : str
23         Location of the gem index file
24
25 Example
26 -----
27 .. code-block:: python
28 :linenos:
29
30 from tool.common import common
31 cf = common()
32
33 indexes = cf.run_indexers('/<data_dir>/human_GRCh38.fa.gz')
34 print(indexes)
35
36
37 """

```

Building the Documentation

Full documentation for a repository can be built using [Sphinx](#). If the pipeline has been developed based on a fork of the mg-process-test repository it can be done by:

```

1  cd ${mg-process-test}
2  pip install sphinx
3
4  cd docs/
5  make html

```

Updating the documentation

If new pipelines or tools are added to the repository then it is important that they are included in the documentation.

Updates for a new tool - docs/tool.rst

A new section can be added to the docs/tool.rst file to reflect the new tool.

Before:

```
1 .. automodule:: tool
2
3     Test Tool
4     -----
5     .. autoclass:: tool.testTool.testTool
6         :members:
```

After:

```
1 .. automodule:: tool
2
3     Test Tool
4     -----
5     .. autoclass:: tool.testTool.testTool
6         :members:
7
8     Test Tool 2
9     -----
10    .. autoclass:: tool.testTool.testTool2
11        :members:
```

Updates for a new pipeline - docs/pipelines.rst

A new section can be added to the docs/pipelines.rst file to reflect the new pipeline. This requires providing a larger description about the input required for running the pipeline, what it returns and examples about how to run the code locally and within the COMPSs environment.

An example of a pipeline block is as follows:

```
1 Test Tool
2 -----
3 .. automodule:: process_test
4
5     This is a demonstration pipeline using the testTool.
6
7     Running from the command line
8     =====
9
10    Parameters
11    -----
12    config : file
13        Location of the config file for the workflow
14    in_metadata : file
15        Location of the input list of files required by the process
16    out_metadata : file
```

(continues on next page)

(continued from previous page)

```

17     Location of the output results.json file for returned files
18
19     Returns
20     -----
21     output : file
22             Text file with a single entry
23
24     Example
25     -----
26     To run the script locally this can be done as follows:
27
28     .. code-block:: none
29         :linenos:
30
31         cd ${mg-process-test}
32         python mg_process_test/process_test.py --config mg_process_test/tests/json/
↪process_test.json --in_metadata mg_process_test/tests/json/input_test.json --out_
↪metadata mg_process_test/tests/results.json --local
33
34     The `--local` parameter should be used if the script is being run within an
↪environment where (py)COMPSS is not installed. It can also be used in an
↪environment where (py)COMPSS is installed, but the script needs to be run locally
↪for testing purposes.
35
36     When using a local version of the [COMPS virtual machine](http://www.bsc.es/
↪computer-sciences/grid-computing/comp-superscalar/downloads-and-documentation):
37
38     .. code-block:: none
39         :linenos:
40
41         cd /home/compss/code/mg-process-test
42         runcompss --lang=python mg_process_test/process_test.py --config /home/compss/
↪code/mg-process-test/mg_process_test/tests/json/process_test.json --in_metadata /
↪home/compss/code/mg-process-test/mg_process_test/tests/json/input_test.json --out_
↪metadata /home/compss/code/mg-process-test/mg_process_test/tests/results.json
43
44     Methods
45     =====
46     .. autoclass:: process_test.process_test
47         :members:

```

4.1.5 HOWTO - Licensing

All software developed as part of the VRE by the MuG consortium should be openly licensed using the Apache 2.0 software license. This should encompass the APIs, Tool wrappers and pipelines that have been developed.

Implementing the Apache 2.0 license

There are 3 parts to the license

LICENSE file

This is the full Apache license. This should be an unmodified version of the Apache LICENSE file which can be downloaded from:

```
wget http://www.apache.org/licenses/LICENSE-2.0.txt -O LICENSE
```

Often when starting a new project on GitHub this is automatically generated and included in the repository by default.

File headers

At the top of all code and documentation there needs to be a header including the license agreement:

```
See the NOTICE file distributed with this work for additional information
regarding copyright ownership.
```

```
Licensed under the Apache License, Version 2.0 (the "License");
you may not use this file except in compliance with the License.
You may obtain a copy of the License at
```

```
    http://www.apache.org/licenses/LICENSE-2.0
```

```
Unless required by applicable law or agreed to in writing, software
distributed under the License is distributed on an "AS IS" BASIS,
WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
See the License for the specific language governing permissions and
limitations under the License.
```

This should be surrounded by the appropriate commenting depending on the language

NOTICE file

Lists those that own the Copyright on the software in the repository and the dates that they have been involved with the development of the software. There is then a second list of institutes that have contributed. Often these will be the same.

```
Multiscale Genomics (MuG)
Copyright 2015-2017 Institute A
Copyright 2015-2016 Institute B
Copyright 2016-2017 Institute C
```

```
This product includes software developed at:
```

- Institute A
- Institute B
- Institute C

Benefits for developers

For those that are developing software this format means that there is only a single file that needs to be updated at the start of each year to reflect the involvement of those in the development of the software.

If there are additional licenses that need to match specific sections of code then these can be added at the end of the LICENSE file along with the files that they relate to and who owns the Copyright.

There is a single header for all files referring the reader to the NOTICE file for details about the developers and those that have contributed to the code.

4.1.6 HOWTO - Configuration Files

Tool Description

This configuration file is used to describe the Tool and inform the VRE about what arguments are required by the tool and a list of the file types that can be used as inputs and the matching names that should be used as input parameters.

Below is the example Tool config file for the process_test workflow. It is located in the tool_config directory within the repository. For a full description of all of the parameters please consult the [Tool Integration Document](#).

```
{
  "_id": "process_test",
  "name": "Process Test",
  "title": "Test Workflow",
  "short_description": "Generates file with some text",
  "owner": {
    "institution": "EMBL-EBI",
    "author": "Mark McDowall",
    "contact": "mcdowall@ebi.ac.uk",
    "url": "https://github.com/Multiscale-Genomics/mg-process-test"
  },
  "external": true,
  "has_custom_viewer": false,
  "keywords": [
    "dna"
  ],
  "infrastructure": {
    "memory": 12,
    "cpus": 4,
    "executable": "/home/pmes/code/mg-process-test/process_test.py",
    "clouds": {}
  },
  "input_files": [
    {
      "name": "input",
      "description": "Input file",
      "help": "path to the input file",
      "file_type": ["TXT"],
      "data_type": [
        "text"
      ],
      "required": true,
      "allow_multiple": false
    }
  ],
  "input_files_combinations": [
    [
      "input"
    ]
  ],
  "arguments": [],
}
```

(continues on next page)

(continued from previous page)

```
"output_files": [  
  {  
    "name": "output",  
    "required": true,  
    "allow_multiple": false,  
    "file": {  
      "file_type": "TXT",  
      "meta_data": {  
        "visible": true,  
        "tool": "process_test",  
        "description": "Output"  
      },  
      "file_path": "test.txt",  
      "data_type": "text",  
      "compressed": ""  
    },  
  }  
]
```

Input Files

The `input_files` section defines the types of files that are able to be processed. This can be one or many files. Each file object within the list needs to have the following key-pairs:

- name
- description
- help
- file_type
- data_type
- required
- allow_multiple

`file_type` and `data_type` can have multiple values in. For example in the case of a DNA sequence this can have the type of “sequence_genomic” or “sequence_dna”, so a tool that is able to accept both can have both in the list.

The `input_files_combinations` is a list of lists of the valid permutations of files that can be accepted by the tool. For example with aligners that are able to handle single or paired-end alignments would need to be able to accept 1 or 2 FASTQ files. These lists use the name value from the `input_files` file objects.

Arguments

If extra arguments are required by a tool to perform its functions these are defined in the arguments section of the JSON. The arguments section is a list of key-value objects consisting of the following keys:

- name
- description
- help
- type

- required
- default

Examples that can be used within the list include:

```
{
  "name": "test_example_bool_param",
  "description": "Example boolean parameter",
  "help": "Example of a boolean selector",
  "type": "boolean",
  "required": false,
  "default": false
},
{
  "name": "test_example_integer_param",
  "description": "Example integer parameter",
  "help": "Example of an integer input",
  "type": "integer",
  "required": false,
  "default": 5
},
{
  "name": "test_example_string_param",
  "description": "Example string parameter",
  "help": "Example of a string input",
  "type": "string",
  "required": false,
  "default": "default_string_value"
},
{
  "name": "test_example_selector_param",
  "description": "Example selector parameter",
  "help": "Example of a selector input",
  "type": {
    "type": "string",
    "enum": ["abc", "def", "xyz"]
  }
  "required": false,
  "default": "xyz"
}
```

Examples

For larger examples of VRE JSON configuration files have a look at the [mg-process-fastq configuration files](#) on GitHub.

Test Configuration Files

There are 2 configuration JSON files as inputs for the test instance. These describe the input and output files and an required arguments that need to get passed to the workflow. These configuration files are those that would get passed to the workflow by the VRE.

config.json

Defines the configurations required for by the pipeline including parameters that need to be passed from the VRE submission form, file and the related metadata as well as the output files that need to be produced by the pipeline.

```

1  {
2    "input_files": [
3      {
4        "required": true,
5        "allow_multiple": false,
6        "name": "input",
7        "value": "<unique_file_id>"
8      }
9    ],
10   "arguments": [
11     {
12       "name": "project",
13       "value": "run001"
14     },
15     {
16       "name": "execution",
17       "value": "../run001"
18     },
19     {
20       "name": "description",
21       "value": null
22     },
23     {
24       "name": "<tool_argument>"
25       "value": "<value_from_form>"
26     }
27   ],
28   "output_files": [
29     {
30       "required": true,
31       "allow_multiple": false,
32       "name": "output",
33       "file": {
34         "file_type": "TXT",
35         "meta_data": {
36           "visible": true,
37           "tool": "testTool",
38           "description": "Output"
39         },
40         "file_path": "tests/data/test.txt",
41         "data_type": "text",
42         "compressed": ""
43       }
44     }
45   ]
46 }

```

In the arguments there are 2 sets (project and execution) that will always be present and are provided by the VRE at the point of submission of the to the tool. These are the name of the project that has been given in the VRE and is defined by the user. The second is the execution path, this is the location for where the input files are located and can be used as the working directory for the tool. The other parameters in the arguments list are from form elements based on what parameters the tool requires from the user at run time.

input_file_metadata.json

Lists the file location that are used as input. The configuration names should match those that are in the config.json file defined above.

```

1  [
2    {
3      "_id": "<unique_file_id>",
4      "data_type": "text",
5      "file_type": "TXT",
6      "file_path": "tests/data/test_input.txt",
7      "compressed": 0,
8      "sources": [],
9      "taxon_id": "0",
10     "meta_data": {
11       "visible": true,
12       "validated": 1
13     }
14   }
15 ]

```

Examples

For larger examples of JSON configuration files that can be used to test pipelines have a look at the [mg-process-fastq test configuration files](#) on GitHub.

4.1.7 HOWTO - Logging

As the pipelines and tools with the MuG VRE environment run without the terminal returning to the user, it is important to have a way to communicate to the user that there is an error with the pipeline. As the code is run within a cluster, the text that is printed to screen won't be returned to the user. Within the Tool API a logging interface has been implemented.

Levels of Logging

When there is an issue it can be passed back to the VRE. These are tracked and passed back to the VRE as the application finishes. There is the option to raise errors in 1 of 6 states:

- **INFO** Confirmation that Tool execution is working as expected.
- **DEBUG** Detailed information, typically of interest only when diagnosing problems.
- **WARNING** An indication that something unexpected happened, but that the Tool can continue working successfully.
- **ERROR** A more serious problem has occurred, and the Tool will not be able to perform some function.
- **FATAL** A serious error, indicating that the Tool may be unable to continue running.
- **PROGRESS** Provide the VRE with information about Tool execution progress, in the form of a percentage (0-100)

Using Logging

The code is present within the Tools API, so adding it into a tool or pipeline requires minimal effort. Improving the logging functions requires the following code:

```
1 from utils import logger
```

To add elements to the log can be implemented by:

```
1 logger.info("Processing Text")
```

This logging has been implemented within the `mg-process-test` repository within the `process_test.py` and within the `testTool.py` scripts. There is no logging within the `@task` as from this it is possible to return an actual object that can then be checked by the `run()` function to determine the correct error to return to the main pipeline.

4.1.8 HOWTO - Testing Your Code

Running the Code

To run the code it needs a `config.json` file and an `input_metadata.json` file to provide the input.

Running the pipeline manually

```
1 python mg_process_test/process_test.py --config config.json --in_metadata input_files.  
  ↪ json --out_metadata output_metadata.json
```

Testing a Tools and Pipelines

As defined in the coding standards documentation, it is important to generate scripts for testing the functionality of the tools and workflows. If there are then changes to the code, if it raises errors this is identified sooner rather than later. Within python the use of `pytest` provides the relevant framework around testing code functionality.

Scripts should be placed in the `<repo>/tests` directory.

An example `pytest` for the `test_writer` tool:

```
1 """  
2 .. See the NOTICE file distributed with this work for additional information  
3 regarding copyright ownership.  
4  
5 Licensed under the Apache License, Version 2.0 (the "License");  
6 you may not use this file except in compliance with the License.  
7 You may obtain a copy of the License at  
8  
9     http://www.apache.org/licenses/LICENSE-2.0  
10  
11 Unless required by applicable law or agreed to in writing, software  
12 distributed under the License is distributed on an "AS IS" BASIS,  
13 WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.  
14 See the License for the specific language governing permissions and  
15 limitations under the License.  
16 """  
17
```

(continues on next page)

(continued from previous page)

```

18 from __future__ import print_function
19
20 import os.path
21
22 from mg_process_test.tool.testTool import testTool
23
24 def test_testTool():
25     """
26     Test case to ensure that the GEM indexer works.
27     """
28     resource_path = os.path.dirname(__file__)
29     text_file = resource_path + "/test.txt"
30
31     input_files = {}
32
33     output_files = {
34         "output": text_file
35     }
36
37     metadata = {}
38
39     print(input_files, output_files)
40
41     tt_handle = testTool()
42     tt_files, tt_meta = tt_handle.run(input_files, metadata, output_files)
43
44     assert output_files['output'] == tt_files['output']
45     assert os.path.isfile(text_file) is True
46     assert os.path.getsize(text_file) > 0

```

Automated Testing

Once you have defined your test functions it is handle to then hook up the repository with an automated testing framework that can notify you if there are unexpected changes to the behaviour of your code. This is often triggered whenever there is a push to the repository.

Running in COMPSs

It is possible to use a local version of the [COMPS virtual machine](#) as used by the MuG VRE. Within the VM is possible to install any required software. To run the application the following command can then be used:

```

runcompss                                     \\  

  --lang=python                               \\  

  --library_path=${HOME}/bin                  \\  

  --pythonpath=<pyenv_virtenv_dir>/lib/python2.7/site-packages/ \\  

  --log_level=debug                           \\  

  mg_process_test/process_test.py             \\  

  --config <repo>/tool_config/process_test.json \\  

  --in_metadata <repo>/tests/json/input_process_test.json \\  

  --out_metadata <repo>/tests/json/output_process_test.json

```

The following is a walk through of developing a tool and pipeline wrapper to include new functionality within the MuG VRE. There are several stages covering the Tool development, using the tool within a pipeline and defining the configuration files required so that the final product can be smoothly integrated into the MuG VRE.

4.1.9 Common Coding Standards

When it comes to developing the code all the code should stick to a common standard. This has been defined within the [Coding Standards](#) documentation as well as how to set up the licenses correctly so that your package can be integrated.

4.1.10 Adding a new function

All of the examples in the following sections describe code that has been incorporated into a functional pipeline and tool within a demonstration VRE Tool that is ready for deployment within the VRE. The code can be found in GitHub repository [mg-process-test](#).

[GitHub mg-process-test](#)

In the test process there are example pipelines, tools, documentation, setup scripts unit tests and config files. This repository can be forked and used as the base for developing new pipelines and tools.

The following documents will help guide with the creation of all the components required for creating a tool ready to be integrated into the VRE. To help with the development and [Development Checklist](#) has been created to provide a generic guide and checklist to help make sure that nothing has been forgotten.

Wrapping a Tool This section guides you through how to wrap an external tool, or create a tool that utilises the pyCOMPSs framework and should be capable of running within the MuG VRE environment.

Creating a Pipeline Once you have created a tool you can now incorporate one or multiple tools into a pipeline. This will handle the passing of variables from the VRE to the tool and the tracking of outputs ready for handing back to the VRE. This document will also help in creating test input metadata and file location JSON files that are required to run the pipeline.

Documentation This provides a overview of the documentation requirements as described by the [MuG Coding Standards](#).

Logging Takes you through adding logging to your pipelines and tools to return messages to the user via the MuG VRE.

Testing Your Code A important part of making sure that a pipeline or tool is ready for integration is ensuring that the code has been tested. This covers testing the code is functional and that it is capable of running within the infrastructure used by the VRE.

Licensung Your Code The Apache 2.0 license is required for pipelines and tools to to be integrated into the MuG VRE.

VRE Configuration This takes you through creating JSON configuration files for your tool. This should define all the inputs, outputs and any arguments that are required by the pipelines and tools.

4.1.11 Integrating a new tool into the VRE

The next step is the integration of the pipeline/tool into the MuG VRE. The [Configuration](#) should provide guide to the initial JSON files. The full JSON specification is located in this [GoogleDoc](#). The details the requirements for correctly creating the Tool description JSON file and the requirements for parameters needed for an application.

4.2 MuG Coding Guidelines

The purpose of this document is to provide a description of the standards that code should conform to so that everything can be share and developed with ease.

4.2.1 Language and Versions

- Python 2.7
- Python 3.6

4.2.2 Installation Method

- PIP

4.2.3 Environment Management

- pyenv
- pyenv-virtualenv

4.2.4 Style

This should follow the PEP8 standard defined by the [Python community](#). Also check out the [Google Python Style Guide](#), a quick and easy reference.

This should be enforced with the use of [pylint](#) to ensure that we are matching the PEP8 coding standard.

In addition for every script that is written at the top of ALL python scripts/modules should be the stub license agreement

4.2.5 Header

At the top of all scripts and modules there should be the minified license version for the code. There should also be a full copy of the licence in with the repo as part of the root _dir and a reStructuredText version as part of the documentation.

As part of the head section is also the shebang (`#!/`) line. This should only be included if the script is an executable and refer to the form:

```
#!/usr/bin/env python
```

If the file just contains classes and functions then no shebang is required.

4.2.6 Repository Structure

This is based on python coding standards (PEP8) and the requirements for installation (pip) and documentation (as detailed below). The base contents of a git repository should include:

```
1 <repo_name>/
2   docs/
3     conf.py
4     index.rst
5     install.rst
6     license.rst
7     ...
8 <module>/
9   __init__.py
10  ...
```

(continues on next page)

(continued from previous page)

```
11  scripts/  
12     travis/  
13         <travis_test_scripts>.sh  
14  tests/  
15     data/  
16         test_<function_name>.py  
17  .travis.yml  
18  LICENSE  
19  README.md  
20  requirements.txt  
21  setup.cfg  
22  setup.py
```

4.2.7 Documentation

For this we use ReadTheDocs. This is based on the Sphinx annotation servers and the reStructuredText format ([Primer and RTD related docs](#))

The code for a basic setup within a repo is as follows:

```
1  cd <repo_root>  
2  
3  pip install sphinx  
4  
5  mkdir docs  
6  cd docs  
7  sphinx-quickstart
```

Once the *docs* folder has been generated the documentation can be built with:

```
1  cd <repo_root>/docs  
2  make html
```

It is advisable to build the repo locally to remove the majority of the bugs before submitting to GitHub and letting the docs build on RTD.

Common extensions include:

```
1  extensions = [  
2     'sphinx.ext.autodoc',  
3     'sphinx.ext.napoleon',  
4     'sphinx.ext.viewcode',  
5  ]
```

The current theme across all projects is *default*. This can be set like so:

```
1  html_theme = 'default'
```

There is an issue with the display of code blocks, so there needs to be 2 extra style files:

[_static/style.css](#)

```

1 .rst-content .highlight > pre {
2     line-height: 1.5;
3 }

```

_templates/layout.html

```

1 {% extends "!layout.html" %}
2 {% block extrahead %}
3     <link href="{{ pathto("_static/style.css", True) }}" rel="stylesheet" type="text/
4     ↪css">
5 {% endblock %}

```

4.2.8 Classes and Functions

All functions should have matching documentation describing the purpose of the function, the inputs, outputs and where relevant an example piece of code showing how to call the function:

```

1 """
2 Assembly Index Manager
3
4 Manges the creation of indexes for a given genome assembly file. If the
5 downloaded file has not been unzipped then it will get unzipped here.
6 There are then 3 indexers that are available including BWA, Bowtie2 and
7 GEM. If the indexes already exist for the given file then the indexing
8 is not rerun.
9
10 Parameters
11 -----
12 file_name : str
13     Location of the assembly FASTA file
14
15 Returns
16 -----
17 dict
18     bowtie : str
19         Location of the Bowtie index file
20     bwa : str
21         Location of the BWA index file
22     gem : str
23         Location of the gem index file
24
25 Example
26 -----
27 .. code-block:: python
28     :linenos:
29
30     from tool.common import common
31     cf = common()
32
33     indexes = cf.run_indexers('/<data_dir>/human_GRCh38.fa.gz')
34     print(indexes)
35
36
37 """

```

4.2.9 Architectural Design Record (ADR)

For all repositories there should be a document called `adr.rst`. This should record choices that have been made and summaries the reason for those decisions. This is to provide an in-code record of the design process and reasoning behind why technologies have been selected. In the case of python, pytest, pyenv and pyenv-virtualenv this is the standard setup for use within the pyCOMPSs environment. It is the selection of the key technology that is important for the most part, but there will be times that one technology was chosen over another due to the libraries that are used.

4.2.10 Testing

pytest is the standard in the Python community and has been adopted for testing within the MuG WP4 related code.

As with all python scripts these should have the licence stub and documentation for all functions.

Runs of tests should also tidy up after themselves once they have completed so that the environment is clean ready for the next test case to run. This could mean that some files will get generated multiple times, but these should be small sample datasets.

To avoid the use of too many datasets and provide function level testing Mock should be used.

The following options should be used to test code:

```
1 # Run only the tests
2 pytest
3
4 # Run only pylint as a test
5 pytest --pylint --pylint-rcfile=pylintrc -m pylint
6
7 # Run both
8 pytest --pylint --pylint-rcfile=pylintrc
```

There will also be times when there are sections of code that are under development or when a test needs to not be included as it is long running or has a bug. To handle this pytest has decorators for this. If a test is to not be used within the TravisCI environment then the following decorator should be used:

```
1 @pytest.mark.underdevelopment
```

pytest can then be run in the following manner:

```
1 # Runs all tests
2 pytest
3
4 # Runs only those marked as underdevelopment
5 pytest -m "underdevelopment"
6
7 # Runs all tests except those underdevelopment
8 pytest -m "not underdevelopment"
```

Sample Data

For all test cases there should be matching datasets that are packaged within the repo.

All datasets should be in the directory `<repo>/tests/data` with a name patching the pattern `<script_name>.<species>.<assembly>.fasta` for genome files and `<script_name>.<accession>.fastq` for read files.

Only the raw files should be stored. For testing these should be small files (~100kB).

Large files can be store, but in cases like that it might be best to have a generation script that can calculate the relevant file with the data structure. If this is part of a reader then it should be part of the DM API and stored within the *dm_generator* directory. The script should be runnable from the command line but should also be able to be run by the reader when the *user_id* is *test*. The generated file should be saved to the */tmp/* folder as *sample_<reader-tag>.<file-tag>*.

4.3 License

Apache License Version 2.0, January 2004 <http://www.apache.org/licenses>

1. Definitions.

“License” shall mean the terms and conditions for use, reproduction, and distribution as defined by Sections 1 through 9 of this document.

“Licensor” shall mean the copyright owner or entity authorized by the copyright owner that is granting the License.

“Legal Entity” shall mean the union of the acting entity and all other entities that control, are controlled by, or are under common control with that entity. For the purposes of this definition, “control” means (i) the power, direct or indirect, to cause the direction or management of such entity, whether by contract or otherwise, or (ii) ownership of fifty percent (50%) or more of the outstanding shares, or (iii) beneficial ownership of such entity.

“You” (or “Your”) shall mean an individual or Legal Entity exercising permissions granted by this License.

“Source” form shall mean the preferred form for making modifications, including but not limited to software source code, documentation source, and configuration files.

“Object” form shall mean any form resulting from mechanical transformation or translation of a Source form, including but not limited to compiled object code, generated documentation, and conversions to other media types.

“Work” shall mean the work of authorship, whether in Source or Object form, made available under the License, as indicated by a copyright notice that is included in or attached to the work (an example is provided in the Appendix below).

“Derivative Works” shall mean any work, whether in Source or Object form, that is based on (or derived from) the Work and for which the editorial revisions, annotations, elaborations, or other modifications represent, as a whole, an original work of authorship. For the purposes of this License, Derivative Works shall not include works that remain separable from, or merely link (or bind by name) to the interfaces of, the Work and Derivative Works thereof.

“Contribution” shall mean any work of authorship, including the original version of the Work and any modifications or additions to that Work or Derivative Works thereof, that is intentionally submitted to Licensor for inclusion in the Work by the copyright owner or by an individual or Legal Entity authorized to submit on behalf of the copyright owner. For the purposes of this definition, “submitted” means any form of electronic, verbal, or written communication sent to the Licensor or its representatives, including but not limited to communication on electronic mailing lists, source code control systems, and issue tracking systems that are managed by, or on behalf of, the Licensor for the purpose of discussing and improving the Work, but excluding communication that is conspicuously marked or otherwise designated in writing by the copyright owner as “Not a Contribution.”

“Contributor” shall mean Licensor and any individual or Legal Entity on behalf of whom a Contribution has been received by Licensor and subsequently incorporated within the Work.

2. **Grant of Copyright License.** Subject to the terms and conditions of this License, each Contributor hereby grants to You a perpetual, worldwide, non-exclusive, no-charge, royalty-free, irrevocable copyright license to reproduce, prepare Derivative Works of, publicly display, publicly perform, sublicense, and distribute the Work and such Derivative Works in Source or Object form.
3. **Grant of Patent License.** Subject to the terms and conditions of this License, each Contributor hereby grants to You a perpetual, worldwide, non-exclusive, no-charge, royalty-free, irrevocable (except as stated in this section) patent license to make, have made, use, offer to sell, sell, import, and otherwise transfer the Work, where such license applies only to those patent claims licensable by such Contributor that are necessarily infringed by their Contribution(s) alone or by combination of their Contribution(s) with the Work to which such Contribution(s) was submitted. If You institute patent litigation against any entity (including a cross-claim or counterclaim in a lawsuit) alleging that the Work or a Contribution incorporated within the Work constitutes direct or contributory patent infringement, then any patent licenses granted to You under this License for that Work shall terminate as of the date such litigation is filed.
4. **Redistribution.** You may reproduce and distribute copies of the Work or Derivative Works thereof in any medium, with or without modifications, and in Source or Object form, provided that You meet the following conditions:
 - (a) You must give any other recipients of the Work or Derivative Works a copy of this License; and
 - (b) You must cause any modified files to carry prominent notices stating that You changed the files; and
 - (c) You must retain, in the Source form of any Derivative Works that You distribute, all copyright, patent, trademark, and attribution notices from the Source form of the Work, excluding those notices that do not pertain to any part of the Derivative Works; and
 - (d) If the Work includes a “NOTICE” text file as part of its distribution, then any Derivative Works that You distribute must include a readable copy of the attribution notices contained within such NOTICE file, excluding those notices that do not pertain to any part of the Derivative Works, in at least one of the following places: within a NOTICE text file distributed as part of the Derivative Works; within the Source form or documentation, if provided along with the Derivative Works; or, within a display generated by the Derivative Works, if and wherever such third-party notices normally appear. The contents of the NOTICE file are for informational purposes only and do not modify the License. You may add Your own attribution notices within Derivative Works that You distribute, alongside or as an addendum to the NOTICE text from the Work, provided that such additional attribution notices cannot be construed as modifying the License.

You may add Your own copyright statement to Your modifications and may provide additional or different license terms and conditions for use, reproduction, or distribution of Your modifications, or for any such Derivative Works as a whole, provided Your use, reproduction, and distribution of the Work otherwise complies with the conditions stated in this License.

5. **Submission of Contributions.** Unless You explicitly state otherwise, any Contribution intentionally submitted for inclusion in the Work by You to the Licensor shall be under the terms and conditions of this License, without any additional terms or conditions. Notwithstanding the above, nothing herein shall supersede or modify the terms of any separate license agreement you may have executed with Licensor regarding such Contributions.
6. **Trademarks.** This License does not grant permission to use the trade names, trademarks, service marks, or product names of the Licensor, except as required for reasonable and customary use in describing the origin of the Work and reproducing the content of the NOTICE file.
7. **Disclaimer of Warranty.** Unless required by applicable law or agreed to in writing, Licensor provides the Work (and each Contributor provides its Contributions) on an “AS IS” BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied, including, without

limitation, any warranties or conditions of TITLE, NON-INFRINGEMENT, MERCHANTABILITY, or FITNESS FOR A PARTICULAR PURPOSE. You are solely responsible for determining the appropriateness of using or redistributing the Work and assume any risks associated with Your exercise of permissions under this License.

8. **Limitation of Liability.** In no event and under no legal theory, whether in tort (including negligence), contract, or otherwise, unless required by applicable law (such as deliberate and grossly negligent acts) or agreed to in writing, shall any Contributor be liable to You for damages, including any direct, indirect, special, incidental, or consequential damages of any character arising as a result of this License or out of the use or inability to use the Work (including but not limited to damages for loss of goodwill, work stoppage, computer failure or malfunction, or any and all other commercial damages or losses), even if such Contributor has been advised of the possibility of such damages.
9. **Accepting Warranty or Additional Liability.** While redistributing the Work or Derivative Works thereof, You may choose to offer, and charge a fee for, acceptance of support, warranty, indemnity, or other liability obligations and/or rights consistent with this License. However, in accepting such obligations, You may act only on Your own behalf and on Your sole responsibility, not on behalf of any other Contributor, and only if You agree to indemnify, defend, and hold each Contributor harmless for any liability incurred by, or claims asserted against, such Contributor by reason of your accepting any such warranty or additional liability.

END OF TERMS AND CONDITIONS

APPENDIX: How to apply the Apache License to your work.

To apply the Apache License to your work, attach the following boilerplate notice, with the fields enclosed by brackets “{}” replaced with your own identifying information. (Don’t include the brackets!) The text should be enclosed in the appropriate comment syntax for the file format. We also recommend that a file or class name and description of purpose be included on the same “printed page” as the copyright notice for easier identification within third-party archives.

Copyright [yyyy] [name of copyright owner]

Licensed under the Apache License, Version 2.0 (the “License”); you may not use this file except in compliance with the License. You may obtain a copy of the License at

<http://www.apache.org/licenses/LICENSE-2.0>

Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an “AS IS” BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

CHAPTER 5

Indices and tables

- `genindex`
- `modindex`
- `search`