



MAFw: Modular Analysis Framework

Release 1.2.0

Antonio Bulgheroni et al.

May 06, 2026

1	First steps	3
1.1	Installation	3
1.2	Contributing	3
2	Introduction	5
2.1	Statement of need	5
2.2	MAFw conceptual design	5
2.3	The way ahead	6
2.3.1	A G- / T- UI for MAFw	6
2.3.2	Parallel processing	6
2.3.3	Introduce interactivity	6
3	Processor: The core of MAFw	7
3.1	Execution workflow	7
3.1.1	The single loop execution	7
3.1.2	The for loop execution	8
3.1.3	The while loop execution	8
3.1.4	How to switch from one loop type to another	9
3.2	Subclassing	9
3.2.1	List of methods to be overloaded	9
3.2.2	Customize and document you processor	10
3.3	Processor parameters	10
3.3.1	Public <code>ActiveParameter</code> and private <code>PassiveParameter</code>	11
3.3.2	Parameter configuration	12
3.3.3	The limitation of <code>ActiveParameter</code>	13
3.3.4	Saving parameter configuration	14
3.3.5	Parameter typing	17
3.4	What's next	19
4	Processor Examples	21
4.1	Simple and looping processors	21
4.2	Modify the <i>for loop</i> cycle using the <code>LoopingStatus</code>	23
4.3	For and while loop execution workflow	25
4.4	Importing elements to the database	32
4.4.1	Retrieving information from filenames	32
4.4.2	The basic importer	35
4.4.3	The <code>ImportExample</code> processor	35
5	ProcessorList: Combine your processors in one go	39
5.1	The <code>ProcessorExitStatus</code>	40
5.2	Resources acquisition and distribution	41
5.3	What's next	41

6	Plugins: import your processors	43
6.1	The quick and dirty: code inside MAFw	43
6.2	The clean and effective: turn your library into a plugin	43
6.3	What's next	45
7	Run your processors from the command line	47
7.1	List the processors	47
7.2	Prepare a generic steering file	48
7.3	Run the steering file	49
7.4	Debugging your processors	49
7.5	The exit code	50
7.6	Commands abbreviation	50
7.7	What's next	50
8	Database: your new buddy!	51
8.1	Database: one name, many implementations	51
8.2	Peewee: a simple and small ORM	51
8.2.1	Database drivers	52
8.3	One class for each table	52
8.4	Ready, go!	55
8.4.1	Configuring other types of databases	56
8.5	Triggers: when the database works on its own	58
8.5.1	The signal approach	58
8.5.2	The trigger approach	60
8.5.3	Triggers on different databases	65
8.6	Standard tables	66
8.7	Custom fields	68
8.7.1	Removing widow rows	68
8.7.2	Pruning orphan files	69
8.7.3	Keeping the entries updated	72
8.8	Filters: let us do only what is needed to be done!	72
8.8.1	How to configure a filter	73
8.8.2	How to use a filter	73
8.9	Multi-primary key columns	77
8.10	Importing an existing DB	80
8.11	What's next	81
9	Display your results	83
9.1	Add diagnostic information to your processors	83
9.2	Generate high level graphs	83
9.2.1	Getting the data	84
9.2.2	Plotting the data	85
9.2.3	Mixin approach	86
10	A step by step tutorial for a real experimental case study	91
10.1	The experiment scenario	91
10.1.1	Task 0. Generating the data	91
10.1.2	Task 1. Building your data bank	92
10.1.3	Task 2. Do the analysis	92
10.1.4	Task 3. Prepare a relation plot	92
10.2	The 'code'	92
10.2.1	The database definition	93
10.2.2	The processor library	96
10.2.3	The plugin module	100
10.3	Run the code!	100
11	API	105
11.1	maf	105
11.1.1	maf.active	106

11.1.2	mafwb.db	107
11.1.3	mafwb.decorators	125
11.1.4	mafwb.enumerators	129
11.1.5	mafwb.examples	131
11.1.6	mafwb.hookspecs	141
11.1.7	mafwb.mafwb_errors	141
11.1.8	mafwb.plugin_manager	143
11.1.9	mafwb.plugins	144
11.1.10	mafwb.processor	144
11.1.11	mafwb.processor_library	154
11.1.12	mafwb.runner	172
11.1.13	mafwb.scripts	174
11.1.14	mafwb.timer	182
11.1.15	mafwb.tools	184
11.1.16	mafwb.ui	190
Python Module Index		195
Index		197

1.1 Installation

MAFw can be installed using pip in a separated virtual environment.

Windows

```
c:\> python -m venv mafw-env
c:\> cd mafw-env
c:\mafw-env> Scripts\activate
(mafw-env)c:\mafw-env> pip install mafw
```

linux

```
$ python -m venv mafw-env
$ cd mafw-env
$ source bin/activate
(mafw-env) $ pip install mafw
```

By default mafw comes with an abstract plotting interface. If you want to use [seaborn](#), then just install the optional dependency `pip install mafw[seaborn]`

All MAFw dependencies will be automatically installed by pip.

1.2 Contributing

Contributions to the software development are very much welcome.

If you want to join the developer efforts, the best way is to clone/fork this repository on your system and start working.

The development team has adopted [hatch](#) for basic tasks. So, once you have downloaded the git repository to your system, open a shell there and type:

```
D:\mafw> hatch env create dev
D:\mafw> hatch env find dev
C:\path\to\.venv\mafw\KVhWIDtq\dev.py3.11
C:\path\to\.venv\mafw\KVhWIDtq\dev.py3.12
C:\path\to\.venv\mafw\KVhWIDtq\dev.py3.13
```

to generate the python environments for the development. This command will actually create the whole environment matrix, that means one environment for each supported python version. If you intend to work primarily with one single python version, simply specify it in the create command, for example:

```
D:\maf> hatch env create dev.py3.13
D:\maf> hatch env find dev.py3.13
C:\path\to\.venv\maf\KVhWIDtq\dev.py3.13
```

hatch will take care of installing MAFw in development mode with all the required dependencies. Use the output of the find command, if you want to add the same virtual environment to your favorite IDE. Once done, you can spawn a shell in the development environment just by typing:

```
D:\maf> hatch shell dev.py3.13
(dev.py3-13) D:\maf>
```

and from there you can simply run mafw and all other scripts.

MAFw uses `pre-commit` to assure a high quality code style. The pre-commit package will be automatically installed into your environment, but it must be initialised before its first use. To do this, simply run the following command:

```
(dev.py3-13) D:\maf> pre-commit install
```

And now you are really ready to go with your coding!

Before pushing all your commits to the remote branch, we encourage you to run the pre-push tests to be sure that everything still works as expected. You can do this by typing:

```
D:\maf> hatch run dev.py3-13:pre-push
```

if you are not in an activated development shell, or

```
(dev.py3-13) D:\maf> hatch run pre-push
```

if you are already in the dev environment.

These pre-push checks will include some cosmetic aspects (ruff check and format) and more relevant points like static type checking with mypy, documentation generation with sphinx and functionality tests with pytest.

2.1 Statement of need

MAFw addresses the need for a flexible and modular framework that enables data scientists to implement complex analytical tasks in a well-defined environment. Currently, data analysis workflows often require scientists to handle multiple tasks, such as data ingestion, processing, and visualization, which can be time-consuming and prone to errors. Moreover, the lack of standardization in data analysis pipelines can lead to difficulties in reproducing and sharing results.

MAFw aims to fill this gap by providing a Python-based tool that allows data scientists to focus on the analysis itself, rather than on the ancillary tasks. The framework is designed to be highly customizable, enabling users to create their own processors and integrate them into the workflow. A key feature of MAFw is its strong collaboration with a relational database structure, which simplifies the analysis workflow by providing a centralized location for storing and retrieving data. This database integration enables seamless data exchange between different processors, making it easier to manage complex data pipelines.

2.2 MAFw conceptual design

The concept behind MAFw is certainly not novel. Its functionality is so prevalent in data analysis that numerous developers, particularly data scientists, have attempted to create libraries with similar capabilities. MAFw's developers got inspired by [MARLIN](#): this object C++ oriented application framework, no longer being maintained, was offering a modular environment where particle physicists were developing their code in the form of shared libraries that could be loaded at run time in a plugin-like manner¹. One of MARLIN strengths was the strong connection with the serial I/O persistency data model offered by [LCIO](#).

Starting from those solid foundations, MAFw moved from C++ to python in order to facilitate the on-boarding of data scientists and to profit from the vast availability of analytical tools, replacing the obsolete [LCIO](#) backend with a more flexible database supported input/output able to deal with large amount of data with categorical variables without severely impacting on I/O performance.

The general concept behind MAFw has been developed by the authors to perform image analysis on [autoradiography images](#), featuring an ultra simplified database interface ([sqlite](#) only) along with some dedicated processors targeting autoradiography specific tasks.

Having understood the potentiality of this scheme, the authors decided to extract the core functionalities of the framework itself, expand the database interface making use of an ORM approach ([peewee](#)), include a plugin system to simplify the integration of processors developed for different purposes in external projects and supply an extensive general and API documentation, before releasing the code to the public domain as open source.

¹ One of MAFw developer was in charge of the original coding of [EUTelescope](#), a set of [MARLIN](#) processors for the reconstruction of particle trajectories recorded with beam telescopes.

2.3 The way ahead

The future development of MAFw is driven by code usability. The authors are trying their best to make the framework as functional as possible offering colleague scientists a platform where to perform their analyses. At the time of writing there are already three targets envisaged: improved user-friendliness via a GUI or at least a TUI, improved performance via parallel processing and improved interactivity.

2.3.1 A G- / T- UI for MAFw

For the time being, the authors focused their efforts in coding a functional framework leaving usability aspects for a later development stage. MAFw is able to *execute* one or more processors one after the other from the command line following the instructions given in a human- and machine- readable (`toml`) steering file. While users can generate examples of steering files to serve as a starting point, implementing a graphical or at least a textual user interface would significantly simplify the process of creating and executing these files.

Following the plugin approach already used elsewhere in the library, the authors are considering to implement a TUI based on `textual` where the user is guided in the creation of steering files.

2.3.2 Parallel processing

Python is often recognized for its slower performance compared to some other programming languages available today. Analytical tasks are often I/O and/or CPU demanding, making the performance of a single threaded single processor program somehow limited.

MAFw uses `pandas` and `numpy` when dealing with data and those libraries are already capable to perform concurrent operations under some specific circumstances.

The well-known `python GIL` is actually preventing multi-threaded applications to improve overall performance for CPU-bound tasks. In this respect, an important revolution is taking place with the release of the experimental `Free-threaded CPython` implementation of python (added in version 3.13) and MAFw authors are closely observing those developments to identify the more convenient approach to improve computational performance.

2.3.3 Introduce interactivity

Even though in the original implementation of MAFw precursor, interactive processors were already existing, they were temporary removed from the current implementation. The authors recognize that many data scientists prefer to conduct interactive analysis using `jupyter` or `marimo` notebooks. Therefore, they are actively exploring ways to seamlessly integrate interactivity into the processor workflow through these notebook environments.

The Processor is responsible to carry out a specific analytical task, in other simple words, it takes some input data, it does some calculations or manipulations, and it produces some output data.

The input and output data can be of any type: a simple list of numbers, a structured data frame, a path to a file where data are stored, an output graph or a link to a web resource, a database table, and so on.

In ultra simplified words, a Processor does three sequential actions:

1. Prepare the conditions to operate (*start()*)
2. Process the data (*process()*)
3. Clean up what is left (*finish()*).

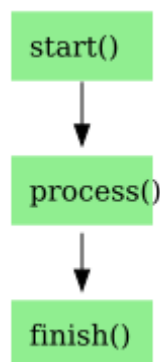
3.1 Execution workflow

There are instances where tasks can be executed in a single operation across your entire dataset, such as generating a graph from the data. However, there are also situations where you need to iterate over a list of similar items, applying the same process repeatedly. Additionally, you may need to continue a process until a specific condition is met. The *Processor* can actually accomplish all those different execution workflows just by changing a variable: *the loop type*.

You will be presented a more detailed description of the different workflows in the following sections.

3.1.1 The single loop execution

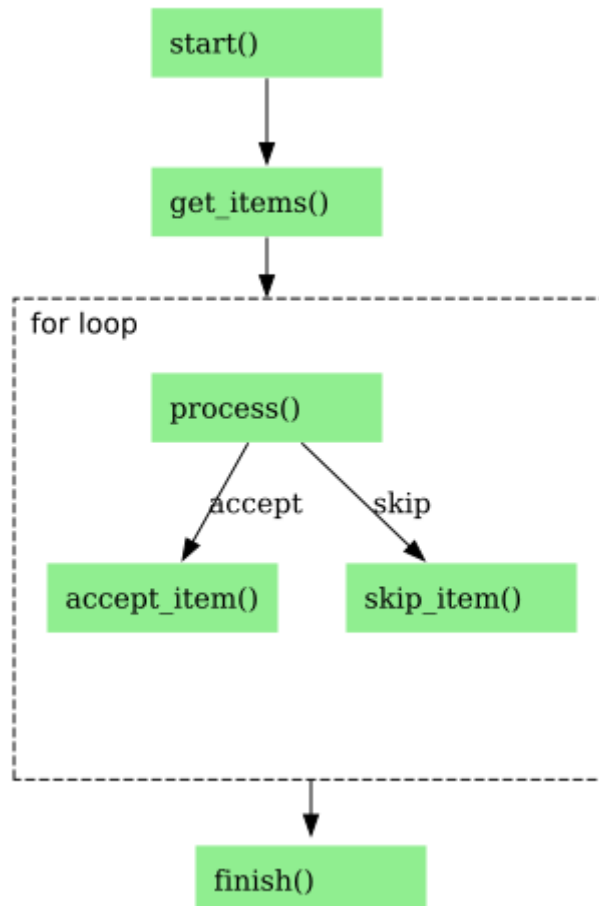
The first family has a simplified execution workflow schematically shown here below:



The distinction of roles between the three methods is purely academic, one could implement all the preparation, the real calculation and the clean up in one method and the processor will work in the same way. The methods are anyhow preserved to offer a similar execution scheme also for the other looping scheme.

3.1.2 The for loop execution

In python a for loop is generally performed on a list of items and MAFw is actually following the same strategy. Here below is the schematic workflow.

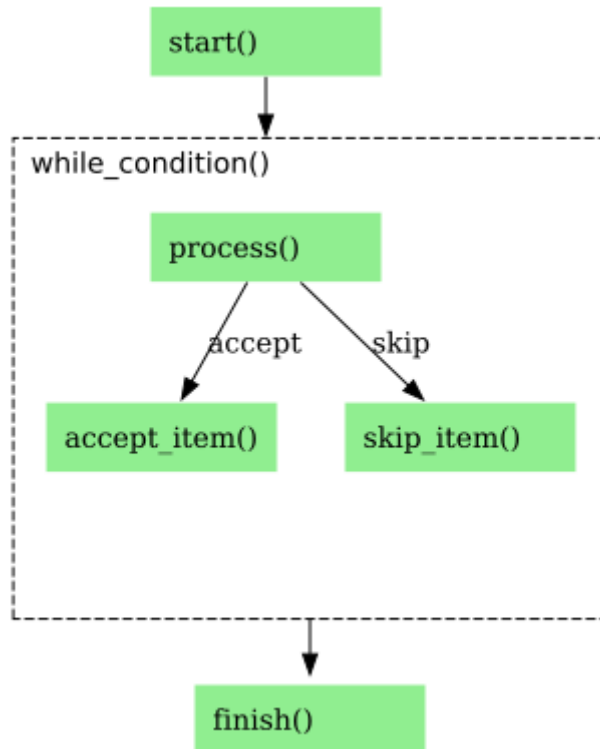


As you can see, after having called the `start()`, the user must provide a list of items to be processed implementing the `get_items()`. This can be the list of files in a directory, or the rows in a DB table or even a list of simple numbers; whatever best suit the user's needs. You will soon learn how to deal with *database* entries and how to *filter* them.

Now everything is ready to start the loop and call the `process()` as many times as the items in the input list. In the process implementation, the user can realize that something went wrong with a particular item and can modify what is executed next (`accept_item()` or `skip_item()`). See an example of such a possibility [here](#). At the end of the loop it is time to clean up everything, saving files, updating DBs and so one. This is again the task of the `finish()` method.

3.1.3 The while loop execution

From a programming point of view, the for loop and the while loop execution are rather similar. Between the execution of the start and the finish method, the process is repeated for a certain number of times until a certain condition is met.



In this case, there is no list of items to loop over, but a condition that should be checked. Thus the user has to overload the `while_condition()` method to return a boolean value: True if the loop has to continue or False if it has to stop and go to the finish.

3.1.4 How to switch from one loop type to another

It is extremely simple to switch from one execution scheme to another. The `Processor` class takes the argument `loop_type`, just change this value in the processor init and the rest will come automatically.

Remember that, by default a processor is set to work with a `for loop` workflow, and thus you have to implement the `get_items()`. If you switch to `while loop`, then you need to implement `while_condition()` for the system to work.

The last thing, you need to know is how to run a processor. First you create an instance of your processor, then you call the `execute()` method. There are other more practical ways, not involving any coding that we will discuss *later on*.

A comparison between the **for loop** and the **while loop** execution workflow is described in the *example page*.

3.2 Subclassing

The basic processor provided in the library does nothing specifically, it is only used as a skeleton for the execution. In order to perform some real analysis, the user has to subclass it and overload some methods. Having gained a clear understanding of the role of each steps in the execution workflow, read the list of methods that can be overloaded here below and then you are ready to see some simple *Processor Examples*.

3.2.1 List of methods to be overloaded

The `process()` method.

This is the central part of the processor, and it must contain all the calculations. If the processor is looping on a list of input items, the user can access the current item via the `item` attribute.

The `start()` method.

This is useful to prepare the condition to operate. For example, it is a good idea to open files and do all the preparatory work. In a looping processor, this method is called just before the cycle starts.

The `finish()` method.

This is useful to clean up after the work. For example, to save files, update DB tables and so on. In a looping processor this is performed just after the last iteration. Here is also the place where you can set the value of the `ProcessorExitStatus`. This attribute is particularly relevant when several processors are executed in a daisy-chained manner. See more about this [here](#) and [here](#).

Specifically for processors with a looping workflow, these methods need to be subclassed.

The `get_items()` method.

This is where the input collection of items is generated. Very likely this can be a list of files resulting from a glob or the rows of an input table or similar. This is required only for **for_loop** processors and must return an iterable object.

The `while_condition()` method.

This is the point at which the decision is made to either continue or terminate the loop. This is required only for **while_loop** processors and must return a boolean value.

The `accept_item()` and `skip_item()` methods (optional).

The execution of the for loop can be slightly modified using the `LoopingStatus` (see this [example](#)). If the current iteration was successful, then the user can decide to perform some actions, otherwise if the current iteration was failing, then some other actions can be taken.

The `format_progress_message()` (optional).

During the process execution, your console may look frozen because your CPU is working out your analysis, thus it may be relevant to have every now and then an update on the progress. The `Processor` will automatically display regular messages via the logging system about the progress (more or less every 10% of the total number of items), but the message it is using is rather generic and does not contain any information about the current item.

By overloading this method, you can include information about the current item and customize the content of the message. You can use `item` to refer to the current item being processed. Here below is an example:

```
def format_progress_message(self):
    self.progress_message = f'{self.name} is processing {self.item}'
```

Be aware that if your items are objects without a good `__repr__` or string conversion, the output may be a little messy.

3.2.2 Customize and document you processor

You are a scientist, we know, the only text you like to write is the manuscript with your last research results. Nevertheless, documenting your code, in particular your processor is a very helpful approach. We strongly recommend to use *docstring* to give a description of what a processor does. If you do so, you will get a surprise when you will *generate your first steering file*.

It is also a very good practice to provide help/doc information to your parameters using the `help_doc` argument of the `ActiveParameter` and `PassiveParameter`. If you do so, your *first steering file* will be much more readable.

One more point, each `Processor` has a class attribute named `description`, this is a short string that is used by some user interfaces (like `RichInterface`) to make the progress bars more meaningful.

3.3 Processor parameters

One super power of MAFw is its capability to re-use code, that means less work, less bugs and more efficiency.

In order to boost code re-usability, one should implement Processor accomplishing one single task and possibly doing it with a lot of general freedom. If you have a processor to calculate the Threshold of B/W images and you have hard coded the threshold algorithm and tomorrow you decide to give a try to another algorithm, then you have to recode a processor that actually already exists.

The solution is to have processor parameters, kind of variable that can be changed in order to make your processor more general and more useful. A note of caution, if you opt for too many parameters, then it may become too difficult to configure your processor. As always, the optimal solution often lies in finding a balance.

You can have parameters implemented in the processor subclass as normal attributes, but then you would need to modify the code in order to change them and this is far from practical. You can have them as specific variables passed to the processor init, but then you would need to code the command line to pass this value or to implement a way to read a configuration file. MAFw has already done all this for you as long as you use the right way to declare processor parameters.

Let us start with a bit of code.

```
class MyProcessor(Processor):
    """
    This is my wonderful processor.

    It needs to know the folder where the input files are stored.
    Let us put this in a processor parameter.

    :param input_folder: The input folder.
    """
    input_folder = ActiveParameter(name='input_folder', default=Path.cwd(),
                                   help_doc='This is the input folder for the file')

    def __init__(self, *args, **kwargs):
        super().__init__(*args, **kwargs)

        # change the input folder to something else
        self.input_folder = Path(r'D:\data')

        # get the value of the parameter
        print(self.input_folder)
```

Some important notes: `input_folder` is defined as a class attribute rather than as an instance attribute, it is to say it is outside the `init` method. **This must always be the case for processor parameters.** The fact that the variable and the parameter names are the same is not compulsory, but why do you want to make your life more complicated!

We have declared `input_folder` as an `ActiveParameter`, but if you put a break point in the code execution and you inspect the type of `self.input_folder` you will see that it is a `Path` object. This is because `ActiveParameter` is a `descriptor`, a very pythonic way to have a double interface for one single variable.

3.3.1 Public `ActiveParameter` and private `PassiveParameter`

When you define `input_folder`, you use up to four parameters in the `ActiveParameter` `init`, but if you try to access those attributes from `input_folder` you will get an error message, because class `Path` does not have the attributes you are looking for. Where are then those values gone?

When you create an `ActiveParameter`, there is a lot of work behind the scene: first of all your class is dynamically modified and an additional attribute named `param_input_file` is created and it is an instance of `PassiveParameter`. This private attribute¹ is actually playing the role of the container, storing the value of the parameter, plus a bunch of other interesting stuff (default, help...).

The set and get dunder methods of the `input_file` are overloaded to operate directly on the private interface, more or less in the same way like when you define `property` setter and getter, but here it is done all automatically.

¹ Python does not have public and private methods in a strict sense like other programming language. Usually variables starting with `_` are considered private. In this specific case, it does not start with an underscore, but we refer to it as private because the user is not aware of having created it.

If you want to access the private interface, you can still do it. And if you have forgotten the name that is automatically assigned (`param_param_name`), you can always use the processor `get_parameter()` using the parameter name as a key. Theoretically you can even set the value of the parameter using the private interface (`set_parameter_value()`), but this is equivalent to set public interface directly.

Let us summarize all this with an example:

```
from mafw.processor import Processor, ActiveParameter

class MyFavProcessor(Processor):
    useful_param = ActiveParameter('useful_param', default=0, help_doc='Important_
↪parameter')

# create your processor, initialize the parameter with a keyword argument.
my_fav = MyFavProcessor(useful_param=10)

# print the value of useful_param in all possible ways
print(my_fav.useful_param)
print(my_fav.param_useful_param.value)
print(my_fav.get_parameter('useful_param').value)

# change the value of useful_param in all possible ways
my_fav.useful_param += 1
my_fav.param_useful_param.value += 1
my_fav.get_parameter('useful_param').value += 1

print(my_fav.useful_param)

# access other fields of the parameter
print(my_fav.get_parameter('useful_param').doc)
```

This is the output that will be generated:

```
10
10
10
13
Important parameter
```

3.3.2 Parameter configuration

We have seen how to add flexibility to a processor including parameters, but how do you configure the parameters?

You have probably noticed that for both `ActiveParameter` and `PassiveParameter` you have the possibility to pass a default value, that is a very good practice, especially for very advanced parameters that will remain untouched most of the time.

If you want to set a value for a parameter, the easiest way is via the processor `__init__` method. The basic `Processor` accepts any number of keyword arguments that can be used exactly for this purpose. Just add a keyword argument named after the parameter and the processor will take care of the rest.

Have a look at the example below covering both the case of Active and Passive parameters:

Listing 3.1: Parameter setting via kwargs

```
1 class MyProcessor(Processor):
2     active_param = ActiveParameter('active', default=0, help_doc='An active parameter
↪')
3
4     def __init__(self, *args, **kwargs):
```

(continues on next page)

(continued from previous page)

```

5         super().__init__(*args, **kwargs)
6         self.passive_param = PassiveParameter('passive', default='I am a string',
↪ help_doc='A string')
7
8 my_p = MyProcessor(active=100, passive='a better string', looper='single')
9
10 print(my_p.active_param) # we get 100
11 assert my_p.active_param == 100
12
13 print(my_p.passive_param.value) # we get 'a better string'
14 assert my_p.passive_param.value == 'a better string'

```

Note that the best way is to avoid to explicitly include the parameter names in the init signature. They will be collected anyhow through keyword arguments and registered automatically.

The second approach is to use a configuration object, it is to say a dictionary containing all the parameters key and value pairs. This is particularly handy when using a configuration file. Exactly for this reason, the configuration object can have one of the two following structures. In both cases the configuration object has to be passed to the class using the keyword *config*

Listing 3.2: Parameter setting via configuration object

```

1 @single_loop
2 class ConfigProcessor(Processor):
3     def __init__(self, *args, **kwargs):
4         super().__init__(*args, **kwargs)
5         self.p1 = PassiveParameter('param1', default='value')
6         self.p2 = PassiveParameter('param2', default='another_value')
7
8 cp = ConfigProcessor(config=dict(param1='new_value', param2='better_value', param3=
↪ 'do not exists'))
9 assert cp.get_parameter('param1').value == 'new_value'
10 assert cp.get_parameter('param2').value == 'better_value'
11 dumped_config = cp.dump_parameter_configuration(option=2)
12 assert dumped_config == dict(param1='new_value', param2='better_value')
13 dumped_config = cp.dump_parameter_configuration(option=20)
14 assert dumped_config == dict(param1='new_value', param2='better_value')
15
16 config = {'ConfigProcessor': {'param1': '1st', 'param2': '2nd'}}
17 cp = ConfigProcessor(config=config)
18 assert cp.get_parameter('param1').value == '1st'
19 assert cp.get_parameter('param2').value == '2nd'
20 dumped_config = cp.dump_parameter_configuration()
21 assert config == dumped_config

```

3.3.3 The limitation of ActiveParameter

With great power comes great responsibility. The use of ActiveParameters has a lot of advantages as we have seen, we can make our processor easily reusable and customizable. The configuration process is straightforward, whether using keyword arguments or a configuration file. But there is at least one disadvantage the user should be aware of and it is connected with the fact that all *ActiveParameter* must be declared as class attributes (read about the difference between class and instance attributes [here](#)). The consequence of this is that two instances of the same class will share the same value for all parameters.

Have a look at the code below:

Listing 3.3: Active parameter usage

```

1 from mafw.processor import ActiveParameter, Processor
2
3 @single_loop
4 class MyProcess(Processor):
5     my_param = ActiveParameter('my_param', default=10)
6
7 my_proc = MyProcess(my_param=12)
8 print(my_proc.my_param) # we expect 12
9 assert my_proc.my_param == 12
10
11 second_proc = MyProcess(my_param=15)
12 print(second_proc.my_param) # we expect 15
13 assert second_proc.my_param == 15
14
15 print(my_proc.my_param) # we would expect 12, but we get 15
16 assert my_proc.my_param == 15

```

Having more than one instance of the same processor type in the same execution run is not really a common situation, but nevertheless there is a workaround. The user can move the definition of the parameter from the class to the instance (moving it inside the processor `__init__` method), and using a *PassiveParameter* instead of an active one. In this way, the parameter will remain bound to the instance (and not to the class), but the user will have to access attributes of the passive parameter using the dot notation. All other super-powers, like parameter registration and configuration, remain unchanged. See the following snippet for a demonstration.

Listing 3.4: Passive parameter usage

```

1 from mafw.processor import PassiveParameter, Processor
2
3 class MyProcess(Processor):
4     def __init__(self, *args, **kwargs):
5         super().__init__(*args, looper='single', **kwargs)
6         self.my_param = PassiveParameter('my_param', default=10)
7
8 my_proc = MyProcess(my_param=12)
9 print(my_proc.my_param.value) # we expect 12
10 assert my_proc.my_param.value == 12
11
12 second_proc = MyProcess(my_param=15)
13 assert second_proc.my_param.value == 15
14 print(second_proc.my_param.value) # we expect 15
15
16 print(my_proc.my_param.value) # we expect 12 and we get 12!
17 assert my_proc.my_param.value == 12

```

The emphasized lines are showing the difference between the two approaches. There is no `__init__` method in the *first snippet* because the *ActiveParameter* is defined at the class level. In the *second snippet*, we moved the definition of `my_param` in the `__init__` method to make it an instance attribute, but we had to change from *Active* to *PassiveParameter*. To access the value of the passive parameter we need to use the `value()` method.

3.3.4 Saving parameter configuration

We have seen that we can configure processor using a dictionary with parameter / value pairs. This is very handy, because we can load toml file with the configuration to be used for all the processors we want to execute.

We don't want you to write toml file hand, for this we have a function `dump_processor_parameters_to_toml()` that will generate an output file with all the parameter values.

But which value is stored? The default or the actual one? Good question!

Let us start from the basics: we have just seen that there are two types of parameters: the *ActiveParameter* and the *PassiveParameter* with the former being class attributes and the latter instance attributes.

This last detail (class or instance attributes) makes a lot of *difference*. If you change a class attribute, this impacts all instances of that class and of all its subclasses as well. If you change an instance attribute, this will only affect that specific instance.

Let us have a look at the following example.

Listing 3.5: Processors definition

```
from mafw.processor import ActiveParameter, PassiveParameter

class ActiveParameterProcessor(Processor):
    """A processor with one active parameter."""

    active_param = ActiveParameter('active_param', default=-1,
    ↪ help_doc='An active parameter with default value -1')

class AnotherActiveParameterProcessor(Processor):
    """Another processor with one active parameter."""

    active_param = ActiveParameter('active_param', default=-1,
    ↪ help_doc='An active parameter with default value -1')

class PassiveParameterProcessor(Processor):
    """A processor with one passive parameter."""

    def __init__(self, *args, **kwargs):
        super().__init__(*args, **kwargs)
        self.passive_param = PassiveParameter('passive_param', default=-1,
    ↪ help_doc='A passive parameter with
    ↪ default value -1')

class AnotherPassiveParameterProcessor(Processor):
    """Another processor with one passive parameter."""

    def __init__(self, *args, **kwargs):
        super().__init__(*args, **kwargs)
        self.passive_param = PassiveParameter('passive_param', default=-1,
    ↪ help_doc='A passive parameter with
    ↪ default value -1')
```

We have defined four different processors, two with active parameters and two with passive parameters.

Now let us have a look at what happens when we attempt to dump the configuration of these processors.

Listing 3.6: Case 1: all processor types

```
processor_list = [ActiveParameterProcessor, AnotherPassiveParameterProcessor,
                 PassiveParameterProcessor, AnotherActiveParameterProcessor]

dump_processor_parameters_to_toml(processor_list, 'test1.toml')
```

We have not created any instances of none of the four processors, and the `processor_list` consists of four classes (no instance). When passed to the `dump_processor_parameters_to_toml()` function, an instance of each class

type is created and the parameter dictionary is retrieved. The creation of an instance is required because only after `__init__` all parameters are registered.

Inside the function, the instances are created without passing any parameters, so it means that for all parameters the dumped value will be the initial value of the parameter, if specified, or the default value if a value is not given in the parameter definition.

The output of the *first example* will be with no surprise, all four processors will have their parameters listed and connected to their default values.

Listing 3.7: Output of case 1: test1.toml

```
[ActiveParameterProcessor] # A processor with one active parameter.
active_param = -1 # An active parameter with default value -1

[AnotherPassiveParameterProcessor] # Another processor with one passive parameter.
passive_param = -1 # A passive parameter with default value -1

[PassiveParameterProcessor] # A processor with one passive parameter.
passive_param = -1 # A passive parameter with default value -1

[AnotherActiveParameterProcessor] # Another processor with one active parameter.
active_param = -1 # An active parameter with default value -1
```

Now, let us take a step forward and add a bit more complexity to the situation. See the next snippet:

Listing 3.8: Case 2: mixed instances and classes

```
# create an instance of ActiveParameterProcessor with a specific value of the
↳parameter
# but we will include the class in the processor list.
active_processor_instance = ActiveParameterProcessor(active_param=100)

# create an instance of AnotherActiveParameterProcessor with a specific value of the
↳parameter
# and we will submit the instance
another_active_processor_instance = AnotherActiveParameterProcessor(active_param=101)

# create an instance of PassiveParameterProcessor with a specific value of the
↳parameter
# but we will submit the class via the use of type.
passive_processor_instance = PassiveParameterProcessor(passive_param=102)

# create an instance of AnotherPassiveParameterProcessor with a specific value of the
↳parameter
# and we will submit the instance.
another_passive_processor_instance = AnotherPassiveParameterProcessor(passive_
↳param=103)

processor_list = [
    ActiveParameterProcessor, # a class
    another_active_processor_instance, # an instance
    type(passive_processor_instance), # a class
    another_passive_processor_instance # an instance
]
dump_processor_parameters_to_toml(processor_list, 'test2.toml')
```

This time we have a mixed list, the first item is a class, so the function will have to create an instance of this. The same is happening with the third element (the `type` function is actually turning back the instance into a class). For these two elements, the generated instances will use no constructor arguments. You might assume that the the

default values are automatically written to the output file, but surprisingly, that is not the case. When you have created the `active_processor_instance` you set the value of `active_param` to 100 and since it is a class attribute, this will be applied to all instances of this class. In the case of the third one, since the parameter is passive, the default value will be actually dumped.

The second element is an instance, so it will not be regenerated inside the function, but since `active_param` is indeed a class attribute, the modified value will be stored in the file. The last one is an instance thus a new instance will not be generated internally and the actual value of the parameter will be stored.

Here is the produced configuration file:

Listing 3.9: Output of case 2: test2.toml

```
[ActiveParameterProcessor] # A processor with one active parameter.
active_param = 100 # An active parameter with default value -1

[AnotherActiveParameterProcessor] # Another processor with one active parameter.
active_param = 101 # An active parameter with default value -1

[PassiveParameterProcessor] # A processor with one passive parameter.
passive_param = -1 # A passive parameter with default value -1

[AnotherPassiveParameterProcessor] # Another processor with one passive parameter.
passive_param = 103 # A passive parameter with default value -1
```

Ultimately, the values that will be output depends on the type of parameter used (Active / Passive) and in case of Passive parameters also on the fact that you pass an instance or a class.

In the end, this is not really important because the user should dump TOML configuration files as a sort of template to be checked, modified and adapted by the user for a specific run.

3.3.5 Parameter typing

When creating a Passive or Active parameter, you have the option to directly specify the parameter's type using the typing template, but you can also do it, and probably in a simpler way assigning a reasonable default value. While this is not really important for numbers, it is extremely important if you want to interpret string as Path object.

If you declare the default value as a Path, for example `Path.cwd()`, then the string read from the configuration file will be automatically converted in a Path.

Note

If you intend to have a float parameter, use a decimal number, for example `0.`, as default, otherwise the interpreter will assume it is an integer and convert to int the parameter being read from the configuration file.

One more note about parameters. Theoretically speaking you could also have custom objects / classes as parameters, but this will become a problem when you will be loading the parameters from a TOML file. Actually two problems:

1. The TOML writer is not necessarily able to convert your custom type to a valid TOML type (number, string...). If your custom type has a relatively easy string representation then you can add an encoder to the TOML writer and teach it how to write your object. See for example the encoder for the Path object.

```
1 def path_encoder(obj: Any) -> Item:
2     """Encoder for PathItem."""
3     if isinstance(obj, PosixPath):
4         return PathItem.from_raw(str(obj), type_=StringType.SLB,
↳escape=False)
5     elif isinstance(obj, WindowsPath):
6         return PathItem.from_raw(str(obj), type_=StringType.SLL,
↳
```

(continues on next page)

(continued from previous page)

```

↪escape=False)
7     else:
8         raise ConvertError

```

2. Even though you managed to write your class to the TOML steering file, you have now the problem of reading back the steering file information and build your custom type with that.

One way to overcome this limitation might be to write to the steering file the `__repr__` of your custom class and at read back time to use `eval` to transform it back to your class. This below would be a more concrete implementation:

python

```

1 import tomlkit
2 from tomlkit.items import String, Item, ConvertError
3 from tomlkit.toml_file import TOMLFile
4
5 class MyClass:
6     """A custom class to be used as a processor parameter"""
7     def __init__(self, a, b):
8         self.a = a
9         self.b = b
10
11     def __repr__(self):
12         """IT MUST BE IMPLEMENTED"""
13         return f"{self.__class__.__name__}({self.a}, {self.b})"
14
15 class MyClassItem(String):
16     """TOML item representing a MyClass"""
17     def unwrap(self) -> MyClass:
18         return MyClass(*super().unwrap())
19
20 def my_class_encoder(obj: MyClassItem) -> Item:
21     """Encoder for MyClassItem."""
22     if isinstance(obj, MyClass):
23         # we write the class as a string using the class repr.
24         return MyClassItem.from_raw(repr(obj))
25     else:
26         raise ConvertError
27
28 # register the encoder
29 tomlkit.register_encoder(my_class_encoder)
30
31 # -----
32 # write to TOML file
33 # -----
34
35 my_class = MyClass(10,24)
36 doc = tomlkit.document()
37 doc.add('my_class', my_class)
38 doc.add('simple_i', 15)
39
40 with open('test.toml', 'w') as fd:
41     tomlkit.dump(doc, fd)
42
43 # -----
44 # read back from TOML file
45 # -----

```

(continues on next page)

(continued from previous page)

```
46 doc = TOMLFile('test.toml').read()
47 read_back_class = eval(doc['my_class'])
48
49 try:
50     simple_i = eval(doc['simple_i'])
51 except TypeError:
52     simple_i = doc['simple_i']
53
54 assert isinstance(read_back_class, MyClass)
55 assert read_back_class.a == my_class.a
56
57 assert isinstance(simple_i, int)
58 assert simple_i == 15
```

TOML

```
my_class = "MyClass(10, 24)"
simple_i = 15
```

This approach, even though possible, is rather risky, we all know how dangerous `eval` can be especially when using it directly with information coming from external files. Furthermore, it is worth considering whether having a custom class as a parameter in a processor is truly necessary. Often, there are simpler and safer alternatives available.

3.4 What's next

You reached the end of the first part. It means that by now you have understood what a processor is and how you can subclass the basic class to implement your analytical tasks, both with looping or single shot workflow. You learned a lot about parameters and how you can configure it.

The next section will be about chaining more processors one after the other using the *ProcessorList*.

From this page, you can see a few example of processors in order to simplify the creation of your first processor sub class.

4.1 Simple and looping processors

The first two examples of this library are demonstrating how you can implement a simple processor that execute all calculations in one go and a looping processor where you need to loop over a list of items to get either a cumulative results.

AccumulatorProcessor is calculating the sum of the first N integer numbers in a loop. The processor takes the *last_number* as an input to include and put the output in the *accumulated_value* parameter. This process is very inefficient, but it is here to demonstrate how to subclass a looping processor.

```

1 class AccumulatorProcessor(Processor):
2     r"""
3     A processor to calculate the sum of the first n values via a looping approach.
4
5     In mathematical terms, this processor solves this easy equation:
6
7     .. math::
8
9         N = \sum_{i=0}^n \{i\}
10
11     by looping. It is a terribly inefficient approach, but it works as a
12     ↪ demonstration of the looping structure.
13
14     The user can get the results by retrieving the `accumulated_value` parameter at
15     ↪ the end of the processor
16     execution.
17     """
18
19     last_value = ActiveParameter('last_value', default=100, help_doc='Last value of
20     ↪ the series')
21
22     def __init__(self, *args, **kwargs):
23         """Constructor parameters:
24
25         :param last_value: The `n` in the equation above. Defaults to 100
26         :type last_value: int

```

(continues on next page)

(continued from previous page)

```

24     :param accumulated_value: The `N` in the equation above at the end of the
↪process.
25     :type accumulated_value: int
26     """
27     super().__init__(*args, **kwargs)
28     self.accumulated_value: int = 0
29
30     def start(self):
31         """Resets the accumulated value to 0 before starting."""
32         super().start()
33         self.accumulated_value = 0
34
35     def get_items(self) -> list[int]:
36         """Returns the list of the first `last_value` integers."""
37         return list(range(self.last_value))
38
39     def process(self):
40         """Increase the accumulated value by the current item."""
41         self.accumulated_value += self.item

```

GaussAdder is calculating exactly the same result using the Gauss formula, eliminating the need for any looping. Indeed the looping is disabled and the output is the same.

```

1  class GaussAdder(Processor):
2      r"""
3      A processor to calculate the sum of the first n values via the so called *Gauss
↪formula*.
4
5      In mathematical terms, this processor solves this easy equation:
6
7      .. math::
8
9          N = \frac{n * (n - 1)}{2}
10
11     without any looping
12
13     The user can get the results by retrieving the `sum_value` parameter at the end of
↪the processor
14     execution.
15     """
16
17     last_value = ActiveParameter('last_value', default=100, help_doc='Last value of
↪the series.')
18
19     def __init__(self, *args, **kwargs):
20         """
21         Constructor parameters:
22
23         :param last_value: The `n` in the equation above. Defaults to 100
24         :type last_value: int
25         :param sum_value: The `N` in the equation above.
26         :type sum_value: int
27         """
28         super().__init__(looper=LoopType.SingleLoop, *args, **kwargs)
29         self.sum_value: int = 0
30

```

(continues on next page)

(continued from previous page)

```

31 def start(self):
32     """Sets the sum value to 0."""
33     super().start()
34     self.sum_value = 0
35
36 def process(self):
37     """Compute the sum using the Gauss formula."""
38     self.sum_value = int(self.last_value * (self.last_value - 1) / 2)

```

If you carefully look at line 28, you will notice that in the GaussAdder constructor, the `looper` option is set to `SingleLoop` and as we have *seen*, it means that that the processor will follow the single loop execution workflow.

The definition of the `looper` parameter in the `init` method can be sometimes hard to remember and unpractical especially if you have to overload the `init` method just to set the value of the `looper`. In such circumstances the use of a class decorator can be very handy. MAFw makes you available three class decorators for this purpose, to transform a processor in a *single loop*, a *for loop* or a *while loop*.

Using the decorator approach the GaussAdder above can be re-written in this way:

```

@single_loop
class GaussAdder(Processor):
    # the rest of the implementation remains the same

```

And here below is an example of execution of the two.

```

from mafw.examples.sum_processor import GaussAdder, AccumulatorProcessor

n = 35

# create the two processors
accumulator = AccumulatorProcessor(last_value=n)
gauss = GaussAdder(last_value=n)

# execute them
accumulator.execute()
gauss.execute()

# print the calculated results
print(accumulator.accumulated_value)
print(gauss.sum_value)

```

This will generate the following output:

```

595
595

```

4.2 Modify the *for loop* cycle using the `LoopingStatus`

In a looping processor, the `process()` method is invoked inside a loop, but the user can decide to skip a certain item and even to interrupt the (abort or quit) the loop.

The tool to achieve this is the `LoopingStatus`. This is set to `Continue` at the beginning of each iteration, but the user can turn it `Skip`, `Abort` or `Quit` inside the implementation of `process()`.

When set to `Skip`, a special callback is invoked `skip_item()` where the user can do actions accordingly. When set to `Abort` or `Quit`, the loop is broken and the user can decide what to do in the `finish()` method. Those two statuses seems to be redundant, but this gives the user the freedom to decide if everything was wasted (`Abort`) or if what done so far was still acceptable (`Quit`).

Here below is the implementation of a simple processor demonstrating such a functionality.

```

1 class ModifyLoopProcessor(Processor):
2     """
3     Example processor demonstrating how it is possible to change the looping
4     ↪ structure.
5
6     It is a looping processor where some events will be skipped, and at some point
7     ↪ one event will trigger an abort.
8     """
9
10    total_item: ActiveParameter[int] = ActiveParameter('total_item', default=100,
11    ↪ help_doc='Total item in the loop.')
12    items_to_skip: ActiveParameter[list[int]] = ActiveParameter(
13    ↪ 'items_to_skip', default=[12, 16, 25], help_doc='List of items to be skipped.'
14    )
15    item_to_abort: ActiveParameter[int] = ActiveParameter('item_to_abort', default=65,
16    ↪ help_doc='Item to abort')
17
18    def __init__(self, *args, **kwargs):
19        """
20        Processor Parameters:
21
22        :param total_item: The total number of items
23        :type total_item: int
24        :param items_to_skip: A list of items to skip.
25        :type items_to_skip: list[int]
26        :param item_to_abort: The item where to trigger an abort.
27        :type item_to_abort: int
28
29        """
30        super().__init__(*args, **kwargs)
31        self.skipped_items: [list[int]] = []
32        """A list with the skipped items."""
33
34    def start(self):
35        """Resets the skipped item container."""
36        super().start()
37        self.skipped_items = []
38
39    def get_items(self) -> list[int]:
40        """Returns the list of items, the range from 0 to total_item."""
41        return list(range(self.total_item))
42
43    def process(self):
44        """Processes the item"""
45        if self.item in self.items_to_skip:
46            self.looping_status = LoopingStatus.Skip
47            return
48        if self.item == self.item_to_abort:
49            self.looping_status = LoopingStatus.Abort
50            return
51
52    def skip_item(self):
53        """Add skipped item to the skipped item list."""
54        self.skipped_items.append(self.item)

```

And here below is how the processor can be used.

```

1 from mafw.examples.loop_modifier import ModifyLoopProcessor
2
3 # generate a random number corresponding to the last item
4 last_value = random.randint(10, 1000)
5
6 # get a sample with event to be skipped
7 skip_items = random.sample(range(last_value), k=4)
8
9 # find an event to abort after the last skipped one
10 max_skip = max(skip_items)
11 if max_skip + 1 < last_value:
12     abort_item = max_skip + 1
13 else:
14     abort_item = last_value - 1
15
16 # create the processor and execute it
17 mlp = ModifyLoopProcessor(total_item=last_value, items_to_skip=skip_items, item_to_
18     ↪ abort=abort_item)
19 mlp.execute()
20
21 # compare the recorded skipped items with the list we provided.
22 assert mlp.skipped_items == list(sorted(skip_items))
23
24 # check that the last item was the abort item.
25 assert mlp.item == abort_item

```

4.3 For and while loop execution workflow

We have seen in the previous chapter that there are different type of loopers and in the previous section we have seen in practice the execution workflow of a single loop and a while loop processor.

In this example, we will explore the difference between the **for loop** and the **while loop** execution workflow. Both processors will run the `Processor.process()` method inside a loop, but for the former we will loop over a pre-established list of items, while for the latter we will continue repeating the process until a certain condition is valid.

Both processors will work with prime number and we will use this *helper function* to check if an integer number is prime or not.

```

1 def is_prime(n: int) -> bool:
2     """
3     Check if n is a prime number.
4
5     :param n: The integer number to be checked.
6     :type n: int
7     :return: True if n is a prime number. False, otherwise.
8     :rtype: bool
9     """
10    prime = True
11    if n < 2:
12        prime = False
13    elif n == 2:
14        prime = True
15    elif n % 2 == 0:
16        prime = False
17    else:
18        sqrt_n = int(math.floor(math.sqrt(n)))

```

(continues on next page)

(continued from previous page)

```

19     for i in range(3, sqrt_n + 1, 2):
20         if n % i == 0:
21             prime = False
22
23     return prime

```

The task of the **for loop** processor is to find all prime numbers included in a given user defined range of integer numbers. In other words, we want to find all prime numbers between 1000 and 2000, for example. The brute force approach is to start a loop on 1000, check if it is prime and if not check the next one until you get to 2000. If a number is actually prime, then store it in a list for further use.

For the sake of clarity, along with the *API documentation*, we are copying here also the processor source code.

```

1 @for_loop
2 class FindPrimeNumberInRange(Processor):
3     """
4     An example processor to find prime numbers in the defined interval from ``start_
5     ↪from`` to ``stop_at``.
6
7     This processor is meant to demonstrate the use of a for_loop execution workflow.
8
9     Let us say we want to select only the prime numbers in a user defined range. One
10    ↪possible brute force approach is
11    ↪to generate the list of integers between the range extremes and check if it is
12    ↪prime or not. If yes,
13    ↪then add it to the list of prime numbers, if not continue with the next element.
14
15    This is a perfect application for a loop execution workflow.
16    """
17
18    start_from = ActiveParameter('start_from', default=50, help_doc='From which
19    ↪number to start the search')
20    stop_at = ActiveParameter('stop_at', default=100, help_doc='At which number to
21    ↪stop the search')
22
23    def __init__(self, *args: Any, **kwargs: Any):
24        """
25        Processor parameters:
26
27        :param start_from: First element of the range under investigation.
28        :type start_from: int
29        :param stop_at: Last element of the range under investigation.
30        :type stop_at: int
31        """
32        super().__init__(*args, **kwargs)
33        self.prime_num_found: list[int] = []
34        """The list with the found prime numbers"""

```

This is the class definition with its constructor. As you can see, we have decorated the class with the *for loop decorator* even though it is not strictly required because the for loop is the default execution workflow.

We have added two processor parameters, the `start_from` and the `stop_at` to allow the user to specify a range on interest where to look for prime numbers.

In the `init` method, we create a list of integer to store all the prime numbers that we will finding during the process.

Now let us overload all compulsory methods for a **for loop** processor.

```

1 def get_items(self) -> Collection[Any]:
2     """
3     Overload of the get_items method.
4
5     This method must be overloaded when you select a for loop workflow.
6
7     Here we generate the list of odd numbers between the start and stop that we need.
8     ↪to check.
9     We also check that the stop is actually larger than the start, otherwise we print.
10    ↪an error message, and we
11    return an empty list of items.
12
13    :return: A list of odd integer numbers between start_from and stop_at.
14    :rtype: list[int]
15    """
16    if self.start_from >= self.stop_at:
17        log.critical('%s must be smaller than %s' % (self.start_from, self.stop_at))
18        return []
19
20    if self.start_from != 2 and self.start_from % 2 == 0:
21        self.start_from += 1
22
23    if self.stop_at != 2 and self.stop_at % 2 == 0:
24        self.stop_at -= 1
25
26    return list(range(self.start_from, self.stop_at, 2))

```

The get items method is expected to return a list of items, that will be processed by the `Processor.process()` method. It is absolutely compulsory to overload this method, otherwise the whole loop structure will not have a list to loop over.

And now, let us have a look at the three stages: start, process and finish.

```

1 def start(self) -> None:
2     """
3     Overload of the start method.
4
5     **Remember:** to call the super method when you overload the start.
6
7     In this specific case, we just make sure that the list of found prime numbers is.
8     ↪empty.
9     """
10    super().start()
11    self.prime_num_found = []
12
13 def process(self) -> None:
14     """
15     The process method.
16
17     In this case, it is very simple. We check if :attr:`.Processor.item` is a prime.
18     ↪number, if so we added to the list,
19     otherwise we let the loop continue.
20     """
21    if is_prime(self.item):
22        self.prime_num_found.append(self.item)
23
24 def finish(self) -> None:

```

(continues on next page)

(continued from previous page)

```

23     """
24     Overload of the finish method.
25
26     **Remember:** to call the super method when you overload the finish method.
27
28     In this case, we just print out some information about the prime number found in_
↪ the range.
29     """
30     super().finish()
31     log.info(
32         'Found %s prime numbers in the range from %s to %s'
33         % (len(self.prime_num_found), self.start_from, self.stop_at)
34     )
35     if len(self.prime_num_found):
36         log.info('The smallest is %s', self.prime_num_found[0])
37         log.info('The largest is %s', self.prime_num_found[-1])

```

These three methods are the core of the execution workflow, so it is obvious that you have to overload them. Keep in mind to always include a call to the super method when you overload the start and finish because they perform some tasks also in the basic processor implementation. The code is written in a straightforward manner and includes clear, thorough explanations in the docstring.

The looping parameters: *Processor.i_item*, *Processor.n_item* and *Processor.item* can be used while implementing the *process()* and *finish()*. The *n_item* is calculated soon after the list of items is returned, while *item*, *i_item* are assigned in the for loop as the current item and its enumeration.

Optionally, one can overload the *format_progress_message()* in order to generate a nice progress message informing the user that something is happening. This is an example:

```

1 def format_progress_message(self) -> None:
2     self.progress_message = (
3         f'Checking integer number: {self.item}, already found {len(self.prime_num_}
↪ found)} prime numbers'
4     )

```

The task for the **while loop** processor is again about prime number finding but different. We want to find a certain number of prime numbers starting from an initial value. We cannot generate a list of integer number and loop over that in the *FindPrimeNumberInRange*, but we need to reorganize our workflow in order to loop until the number of found primes is equal to the requested one.

This is how such a task can be implemented using the while loop execution framework. You can find the example in the *API documentation* and an explanation of the here below.

Let us start again from the class definition.

```

1 @while_loop
2 class FindNPrimeNumber(Processor):
3     """
4     An example of Processor to search for N prime numbers starting from a given_
↪ starting integer.
5
6     This processor is meant to demonstrate the use of a while_loop execution workflow.
7
8     Let us say we need to find 1000 prime numbers starting from 12347. One possible_
↪ brute force approach to solve this
9     problem is to start checking if the initial value is a prime number. If this is_
↪ not the case, then check the next
10    odd number. If it is the case, then add the current number to the list of found_
↪ prime numbers and continue until

```

(continues on next page)

(continued from previous page)

```

11     the size of this list is 1000.
12
13     This is a perfect application for a while loop execution workflow.
14     """
15
16     prime_num_to_find = ActiveParameter(
17         'prime_num_to_find', default=100, help_doc='How many prime number we have to
↪find'
18     )
19     start_from = ActiveParameter('start_from', default=50, help_doc='From which
↪number to start the search')
20
21     def __init__(self, *args: Any, **kwargs: Any):
22         """
23         Processor parameters:
24
25         :param prime_num_to_find: The number of prime numbers to be found.
26         :type prime_num_to_find: int
27         :param start_from: The initial integer number from where to start the search.
28         :type start_from: int
29         """
30         super().__init__(*args, **kwargs)
31         self.prime_num_found: list[int] = []
32         """The list with the found prime numbers"""

```

The first difference compared to the previous case is the use of the `while_loop()` decorator, this time it is really necessary to specify the processor `LoopType` because the while loop is not the default strategy.

The processor has two parameters, the number of prime number to find and from where to start. Similarly as before, in the init method, we define a list of integer to store all the prime numbers that we have found.

For while loop processor, we don't have a list of items, but we need to have a condition either to continue or to stop the loop. For this reason we need to overload the `while_condition()` method, keeping in mind that we return True if we want the cycle to continue for another iteration and False otherwise.

Here is the implementation of the `while_condition()` for the `FindNPrimeNumber`.

```

1  def while_condition(self) -> bool:
2      """
3      Define the while condition.
4
5      First, it checks if the prime_num_to_find is positive. Otherwise, it does not
↪make sense to start.
6      Then it will check if the length of the list with the already found prime numbers
↪is enough. If so, then we can
7      stop the loop return False, otherwise, it will return True and continue the loop.
8
9      Differently from the for_loop execution, we are responsible to assign the value
↪to the looping variables
10     :attr:`.Processor.i_item`, :attr:`.Processor.item` and :attr:`.Processor.n_item`.
11
12     In this case, we will use the :attr:`.Processor.i_item` to count how many prime
↪numbers we have found and :attr:`.Processor.n_item`
13     will be our target. In this way, the progress bar will work as expected.
14
15     In the while condition, we set the :attr:`.Processor.i_item` to the current length
↪of the found prime number list.
16

```

(continues on next page)

(continued from previous page)

```

17     :return: True if the loop has to continue, False otherwise
18     """
19     if self.prime_num_to_find <= 0:
20         log.warning('You requested to find a negative number of prime numbers. It
↳makes no sense.')
21         return False
22
23     self.i_item = len(self.prime_num_found)
24     return self.i_item < self.prime_num_to_find

```

For a while loop, it is not easy to define an enumeration parameter and also the total number of items might be misleading. It is left to the user to decide if they want to use them or not. If yes, their definition and incrementation is under their responsibility. For this processor, it was natural to consider the requested number of primes as the *n_item* and consequently the value of *i_item* can be utilized to keep track of the quantity of prime numbers that have already been discovered. This choice is very convenient because then progress bar that uses *i_item* and *n_item* to calculate the progress will show the actual progress. In case, you do not have any way to assign a value to *n_item*, do not do it, or set it to None. In this way, the progress bar will display an indeterminate progress. You can set the value of *n_item* either in the *start()* or in the *while_condition()*, with a performance preference with the first option because it is executed only once before the start of the loop.

Here below is the implementation of the three stages.

```

1  def start(self) -> None:
2      """
3      The overload of the start method.
4
5      **Remember:** The start method is called just before the while loop is started.
↳So all instructions in this
6      method will be executed only once at the beginning of the process execution.
↳Always put a call to its `super`
7      when you overload start.
8
9      First, we empty the list of found prime numbers. It should not be necessary, but
↳it makes the code more readable.
10     Then set the :attr:`Processor.n_item` to the total number of prime numbers we
↳need to find. In this way, the progress bar
11     will display useful progress.
12
13     If the start value is smaller than 2, then let's add 2 to the list of found prime
↳number and set our first
14     item to check at 3. In principle, we could already add 3 as well, but maybe the
↳user wanted to find only 1
15     prime number, and we are returning a list with two, that is not what he was
↳expecting.
16
17     Since prime numbers different from 2 can only be odd, if the starting number is
↳even, increment it already by
18     1 unit.
19     """
20     super().start()
21     self.prime_num_found = []
22     self.n_item = self.prime_num_to_find
23     if self.start_from < 2:
24         self.prime_num_found.append(2)
25         self.start_from = 3
26
27     if self.start_from % 2 == 0:

```

(continues on next page)

(continued from previous page)

```

28     self.item = self.start_from + 1
29     else:
30         self.item = self.start_from
31
32 def process(self) -> None:
33     """
34     The overload of the process method.
35
36     **Remember:** The process method is called inside the while loop. It has access
↪ to the looping parameters:
37     :attr:'.Processor.i_item', :attr:'.Processor.item' and :attr:'.Processor.n_item'.
38
39     In our specific case, the process contains another while loop. We start by
↪ checking if the current
40     :attr:'.Processor.item' is a prime number or not. If so, then we have found the
↪ next prime number, we add it to the list,
41     we increment by two units the value of :attr:'.Processor.item' and we leave the
↪ process method ready for the next iteration.
42
43     If :attr:'.Processor.item' is not prime, then increment it by 2 and check it again.
44     """
45     while not is_prime(self.item):
46         self.item += 2
47     self.prime_num_found.append(self.item)
48     self.item += 2
49
50 def finish(self) -> None:
51     """
52     Overload of the finish method.
53
54     **Remember:** The finish method is called only once just after the last loop
↪ interaction.
55     Always put a call to its `super` when you overload finish.
56
57     The loop is over, it means that the while condition was returning false, and now
↪ we can do something with our
58     list of prime numbers.
59     """
60     super().finish()
61     log.info('Found the requested %s prime numbers' % len(self.prime_num_found))
62     log.info('The smallest is %s', self.prime_num_found[0])
63     log.info('The largest is %s', self.prime_num_found[-1])

```

Let us have a look at the `FindNPrimeNumber.start()`. First of all we set the value of `Processor.n_item` to our target value of primes. We use the `Processor.item` to store the current integer number being tested, so we initialize it to `start_from` or the first not prime odd number following it. In the `FindNPrimeNumber.process()` we need to include another while loop, this time we need to check the current value of `Processor.item` if it is a prime number. If yes, then we add it to the storage list, we increment it by two units (*remember that for while loop processors it is your responsibility to increment the loop parameters*) and we get ready for the next loop iteration. As for the other processor, we `FindNPrimeNumber.finish()` printing some statistics.

4.4 Importing elements to the database

Note

This example is using concepts that have not yet been introduced, in particular the database. So in a first instance, you can simply skip it and come back later.

Importing elements in the database is a very common task, that is required in all analytical projects. To accomplish this task, mafw is providing a dedicated base class (the *Importer*) that heavily relies on the use of the *FilenameParser* to extract parameters from the filenames.

The *ImporterExample* is a concrete implementation of the base *Importer* that can be used by the user to get inspiration in the development of their importer subclass.

Before diving into the *ImporterExample* code analysis, we should understand the role and the functionality of other two helper classes: the *FilenameElement* and the *FilenameParser*.

4.4.1 Retrieving information from filenames

When setting up an experimental plan involving the acquisition of several data files, there are different approaches.

1. The **descriptive** approach, where the filename is used to store information about the measurement itself,
2. the **metadata** approach, where the same information are stored inside the file in a metadata section,
3. or the **logbook** approach, where the filename is just a unique identifier and the measurement information are stored in a logbook (another file, database, piece of paper...) using the same unique identifier.

The **descriptive** approach, despite being sometime a bit messy because it may end up with very long filenames, it is actually very practical. You do not need to be a hacker including the metadata in the file itself and you do not risk to forget to add the parameters to the logbook.

The tricky part is to include those information to the database containing all your experiments, and you do not want to do this by hand to avoid errors.

The best way is to use **regular expression** that is a subject in which python is performing excellently and MAFw is helping you with two helpers.

The first helper is the *FilenameElement*. This represents one single piece of information that is stored in the filename.

Let us assuming that you have a file named as `sample_12_energy_10_repetition_2.dat`. You can immediately spot that there are three different pieces of information stored in the filename. The sample name, the value of the energy in some unit that you should know, and the value of the repetition. Very likely there is also a `repetition_1` file saved on disc.

In order to properly interpret the information stored in the filename, we need to define three *FilenameElement*s, one for each of them!

If you look at the documentation of the *FilenameElement*, you will see that you need four arguments to build it:

- its **name**, this is easy. Take one, and use it to name a named group in the regular expression.
- its **regular expression**, this is tricky. This is the pattern that python is using to read and parse the actual element.
- its **type**, this is the expected type for the element. It can be a string, an integer or a floating point number.
- its **default** value, this is used to make the element optional. It means that if the element is not found, then the default value is returned. If no default value is provided and the element is not found then an error is raised.

Let us see how you could use *FilenameElement* class to parse the example filename.

```

filename = 'sample_12_energy_10_repetition_2.dat'

sample = FilenameElement('sample', r'[_]*(?P<sample>sample_\d+)[_]*', value_type=str)
energy = FilenameElement('energy', r'[_]*energy_(?P<energy>\d+\.\d*)[_]*', value_
↳type=float)
repetition = FilenameElement(
    'repetition', r'[_]*repetition_(?P<repetition>\d+)[_]*', value_type=int, default_
↳value=1
)

sample.search(filename)
assert sample.value == 'sample_12'

energy.search(filename)
assert energy.value == 10

repetition.search(filename)
assert repetition.value == 2

```

The interesting thing is that you can swap the position of the elements in the filename, for example starting with the energy, and it will still be working absolutely fine.

Just open a python interpreter, import the `FilenameElement` class and give it a try yourself to familiarize with the regular expression. Be careful, when you write the regular expression pattern, since it usually contains a lot of `'\'`, it may be useful to prefix the string with a `r`, in order to inform python that what is coming must be interpreted as a raw string.

If you want to gain confidence with regular expressions, make some tests and understand their power, we recommend to play around with one of the many online tools available on the web, like [pythex](#).

The `FilenameElement` is already very helpful, but if you have several elements in the filename, the readability of your code will quickly degrade. To help you further, you can enjoy the `FilenameParser`.

This is actually a combination of filename elements and when you will try to interpret the filename by invoking `interpret()` all of the filename elements will be parsed and thus you can retrieve all parameters in a much easier way.

If you look at the `FilenameParser` documentation, you will see that you need a configuration file to build an instance of it. This configuration file is actually containing the information to build all the filename element.

In the two tabs here below you can see the configuration file and the python code.

Parser configuration

```

# FilenameParser configuration file
#
# General idea:
#
# The file contains the information required to build all the FilenameElement_
↳requested by the importer.
#
# Prepare a table for each element and in each table add the regexp, the type and_
↳optionally the default.
# Adding the default field, will make the element optional.
#
# Add the table name in the elements array. The order is irrelevant. The division in_
↳compulsory and optional elements
# is also irrelevant. It is provided here just for the sake of clarity.
#
# You can have as many element tables as you like, but only the one listed in the_

```

(continues on next page)

(continued from previous page)

```

→elements array will be used to
# configure the Importer.
#
elements = [
    # compulsory elements:
    'sample', 'energy',
    # optional elements:
    'repetition'
]

[sample]
regexp = '[_]*(?P<sample>sample_\d+)[_]*'
type='str'

[energy]
regexp = '[_]*energy_(?P<energy>\d+\.\.*\d*)[_]*'
type='float'

[repetition]
regexp = '[_]*repetition_(?P<repetition>\d+)[_]*'
type='int'
default = 1

```

Python test code

```

filename = 'energy_10.3_sample_12.dat'

parser = FilenameParser('example_conf.toml')
parser.interpret(filename)

assert parser.get_element_value('sample') == 'sample_12'
assert parser.get_element_value('energy') == 10.3
assert parser.get_element_value('repetition') == 1

```

The configuration file must contain a top level `elements` array with the name of all the filename elements that are included into the filename. For each value in `elements`, there must be a dedicated table with the same name containing the definition of the regular expression, the type and optionally the default value.

Important

In TOML configuration files, the use of single quotation marks allows to treat a string as a raw string, that is very important when passing expression containing backslashes. If you prefer to use double quotation marks, then you have to escape all backslashes.

The order of the elements in the `elements` array is irrelevant and also the fact we have divided them in compulsory and optional is just for the sake of clarity.

In the python tab, you can see how the use of `FilenameParser` makes your code looking much tidier and easier to read. In this second example, we have removed the optional specification of the `repetition` element and you can see that the parser is returning the default value of 1 for such element and we have swapped the energy field with the sample name. Moreover, now the energy field is actually a floating number with a decimal figure.

4.4.2 The basic importer

With the power of these two helper classes, building a processor for parsing all our measurement filenames is a piece of a cake. In the `processor_library` package, you can find a basic implementation of a generic `Importer` processor, that you can use as a base class for your specific importer.

The idea behind this importer is that you are interested in files inside an `input_folder` and possibly all its subfolders. You can force the processor to look recursively in all subfolder by turning the processor parameter `recursive` to `True`. The last parameter of this processor is the `parser_configuration` that is the path to the `FilenameParser` configuration file.

This configuration file is used during the `start()` method of `Importer` (or any of its subclasses) to configure its `FilenameParser`, so that you do not have to worry of this step. In your subclass process method, the filename parser will be straight away ready to use.

Let us have a loop and the `ImporterExample` processor (available in the `examples` package) for a concrete implementation of an importer processor.

4.4.3 The ImportExample processor

We will build a subclass of the `Importer` processor following the `for_loop` execution workflow.

In the `start()` method, we will assure that the target table in the database is existing. The definition of the target database Model (`InputElement` in this example) should be done in a separate database model module to facilitate import statements from other modules as well.

```
class InputElement(MAFwBaseModel):
    """A model to store the input elements"""

    element_id = AutoField(primary_key=True, help_text='Primary key for the input_
    ↪element table')
    filename = FileNameField(unique=True, checksum_field='checksum', help_text='The_
    ↪filename of the element')
    checksum = FileChecksumField(help_text='The checksum of the element file')
    sample = TextField(help_text='The sample name')
    exposure = FloatField(help_text='The exposure time in hours')
    resolution = IntegerField(default=25, help_text='The readout resolution in μm')
```

```
def start(self) -> None:
    """
    The start method.

    The filename parser is ready to use because it has been already configured in the_
    ↪super method.
    We need to be sure that the input table exists, otherwise we create it from_
    ↪scratch.
    """
    super().start()
    self.database.create_tables([InputElement])
```

In the `get_items()`, we create a list of all files, in this case matching the fact that the extension is `.tif`, included in the `input_folder`. We use the recursive flag to decide if we want to include also all subfolders.

The `steering` file may contain a `GlobalFilter` section (see *the Filter section*) and we use the `new_only` flag of the `filter_register`, to further filter the input list from all files that have been already included in the database. It is also important to check that the table is update because you may have an entry pointing to the same filename that in the mean time has been modified. For this purpose the `verify_checksum()` can be very useful. A more detailed explanation of this function will be presented in a *subsequent section*.

```
def get_items(self) -> Collection[Any]:
    r"""Retrieves the list of element to be imported.
```

(continues on next page)

(continued from previous page)

```

    The base folder is provided in the configuration file, along with the recursive_
    ↪ flags and all the filter options.

    :return: The list of items full file names to be processed.
    :rtype: list[Path]
    """
    pattern = '**/*tif' if self.recursive else '*tif'
    input_folder_path = Path(self.input_folder)

    file_list = [file for file in input_folder_path.glob(pattern) if file.is_file()]

    # verify the checksum of the elements in the input table. if they are not up to_
    ↪ date, then remove the row.
    verify_checksum(InputElement)

    if self.filter_register.new_only:
        # get the filenames that are already present in the input table
        existing_rows = InputElement.select(InputElement.filename).namedtuples()
        # create a set with the filenames
        existing_files = {row.filename for row in existing_rows}
        # filter out the file list from filenames that are already in the database.
        file_list = [file for file in file_list if file not in existing_files]

    return file_list

```

The *ImporterExample* follows an implementation approach that tries to maximise the efficiency of the database transaction. It means that instead of making one transaction for each element to be added to the database, all elements are collected inside a list and then transferred to the database with a cumulative transaction at the end of the process itself. This approach, as said, is very efficient from the database point of view, but it can be a bit more demanding from the memory point of view. The best approach depends on the typical number of items to be added for each run and the size of each element.

The implementation of the *process()* is rather simple and as you can see from the source code it is retrieving the parameter values encoded in the filename via the *FilenameParser*. If you are wondering why we have assigned the filename to the filename and to the checksum field, have a look at the section about *custom fields*.

```

def process(self):
    """
    The process method overload.

    This is where the whole list of files is scanned.

    The current item is a filename, so we can feed it directly to the FilenameParser_
    ↪ interpret command, to have it
    parsed. To maximise the efficiency of the database transaction, instead of_
    ↪ inserting each file
    singularly, we are collecting them all in a list and then insert all of them in_
    ↪ the :meth:`~.finish` method.

    In case the parsing is failing, then the element is skipped and an error message_
    ↪ is printed.
    """
    try:
        new_element = {}
        self._filename_parser.interpret(self.item.name)
        new_element['sample'] = self._filename_parser.get_element_value('sample_name')

```

(continues on next page)

(continued from previous page)

```

new_element['exposure'] = self._filename_parser.get_element_value('exposure')
new_element['resolution'] = self._filename_parser.get_element_value(
↪ 'resolution')
new_element['filename'] = self.item
new_element['checksum'] = self.item
self._data_list.append(new_element)
except ParsingError:
    log.critical('Problem parsing %s' % self.item.name)
    self.looping_status = LoopingStatus.Skip

```

The `finish()` is where the real database transaction is occurring. All the elements have been collected into a list, so we can use an `insert_many` statement to transfer them all to the corresponding model in the database. Since we have declared the filename field as unique (this was our implementation decision, but the user is free to relax this requirement), we have added a `on_conflict` clause to deal with the case the user is updating an entry with the same filename.

Since the `super method` is printing the execution statistics, we are leaving its call at the end of the implementation.

```

def finish(self) -> None:
    """
    The finish method overload.

    Here is where we do the database insert with a on_conflict_replace to cope with
↪ the unique constraint.
    """
    # we are ready to insert the lines in the database
    InputElement.insert_many(self._data_list).on_conflict_replace(replace=True).
↪ execute()

    # the super is printing the statistics, so we call it after the implementation
    super().finish()

```


CHAPTER 5

PROCESSORLIST: COMBINE YOUR PROCESSORS IN ONE GO

So far we have seen how a *Processor* can be coded to perform a simple task with a certain degree of generality offered by the configurable parameters. But analytical tasks are normally rather complex and coding the whole task in a single processor will actually go against the mantra of simplicity and code reusability of MAFw.

To tackle your complex analytical task, MAFw proposes a solution that involves chaining multiple processors together. The following processors can start where the previous one stopped so that like in a building game, scientists can put together their analytical solution with simple blocks.

From a practical point of view, this is achieved via the *ProcessorList* that is an evolution of the basic python list, which can contain only instances of processor subclasses or other ProcessorLists.

Once you have appended the processors in the order you want them to be executed, just call the *execute()* method of the list and it will take care of running all the processors.

As simple as that:

```
1 def run_simple_processor_list():
2     """Simplest way to run several processors in a go."""
3     from mafw.examples.sum_processor import AccumulatorProcessor, GaussAdder
4     from mafw.processor import ProcessorList
5
6     # create the list. name and description are optional
7     new_list = ProcessorList(name='AddingProcessor', description='Summing up numbers')
8
9     # append the processors. you can pass parameters to the processors in the
10    ↪ standard way
11    max_value = 120
12    new_list.append(AccumulatorProcessor(last_value=max_value))
13    new_list.append(GaussAdder(last_value=max_value))
14
15    # execute the list. This will execute all the processors in the list
16    new_list.execute()
17
18    # you can access single processors in the list, in the standard way.
19    # remember that the ProcessorList is actually a list!
20    assert new_list[0].accumulated_value == new_list[1].sum_value
```

5.1 The ProcessorExitStatus

We have seen in a *previous section* that the user can modify the looping behavior of a processor by using the *LoopingStatus* enumerator. In a similar manner, the execution loop of a processor list can be modified looking at the *ProcessorExitStatus* of each processors.

When one processor in the list is finishing its task, the *ProcessorList* is checking for its exit status before moving to the next item. If a processor is finishing with an Abort status, then the processor list will raise a *AbortProcessorException* that will cause the loop to be interrupted.

Let us have a look at the snippet here below:

```

1 def run_processor_list_with_loop_modifier():
2     """Example on deal with processors inside a processor list changing the loop
   ↳ structure.
3
4     In this example there are two processors, one that will run until the end and the
   ↳ other that will set the looping
5     status to abort half way. The user can see what happens when the :class:`~maf.
   ↳ processor.ProcessorList` is executed.
6     """
7     import time
8
9     from mafw.enumerators import LoopingStatus, ProcessorExitStatus, ProcessorStatus
10    from mafw.mafw_errors import AbortProcessorException
11    from mafw.processor import ActiveParameter, Processor, ProcessorList
12
13    class GoodProcessor(Processor):
14        n_loop = ActiveParameter('n_loop', default=100, help_doc='The n of the loop')
15        sleep_time = ActiveParameter('sleep_time', default=0.01, help_doc='So much
   ↳ work')
16
17        def get_items(self) -> list[int]:
18            return list(range(self.n_loop))
19
20        def process(self):
21            # pretend to do something, but actually sleep
22            time.sleep(self.sleep_time)
23
24        def finish(self):
25            super().finish()
26            print(f'{self.name} just finished with status: {self.processor_exit_
   ↳ status.name}')
27
28    class BadProcessor(Processor):
29        n_loop = ActiveParameter('n_loop', default=100, help_doc='The n of the loop')
30        sleep_time = ActiveParameter('sleep_time', default=0.01, help_doc='So much
   ↳ work')
31        im_bad = ActiveParameter('im_bad', default=50, help_doc='I will crash it!')
32
33        def get_items(self) -> list[int]:
34            return list(range(self.n_loop))
35
36        def process(self):
37            if self.item == self.im_bad:
38                self.looping_status = LoopingStatus.Abort
39                return
40            # let me do my job

```

(continues on next page)

(continued from previous page)

```

41         time.sleep(self.sleep_time)
42
43     def finish(self):
44         super().finish()
45         print(f'{self.name} just finished with status: {self.processor_exit_
↪status.name}')
46
47     proc_list = ProcessorList(name='with exception')
48     proc_list.extend([GoodProcessor(), BadProcessor(), GoodProcessor()])
49     try:
50         proc_list.execute()
51     except AbortProcessorException:
52         print('I know you were a bad guy')
53     assert proc_list.processor_exit_status == ProcessorExitStatus.Aborted
54     assert proc_list[0].processor_exit_status == ProcessorExitStatus.Successful
55     assert proc_list[1].processor_exit_status == ProcessorExitStatus.Aborted
56     assert proc_list[2].processor_status == ProcessorStatus.Init

```

We created two processors, a good and a bad one. The good one is doing nothing, but getting till the end of its job. The bad one is also doing nothing but giving up before the end of the item list. In the process method, the bad processor is setting the looping status to abort, causing the for loop to break immediately and to call finish right away. In the processor finish method, we check if the status was aborted and in such a case we set the exit status of the processor to Aborted.

At line 47, we create a list and we populate it with three elements, a good, a bad and another good processor and we execute it inside a try/except block. The execution of the first good processor finished properly as you can see from the print out and also from the fact that its status (line 54) is Successful. The second processor did not behave, the exception was caught by the except clause and this is confirmed at line 55 by its exit status. The third processor was not even started because the whole processor list got stopped in the middle of processor 2.

5.2 Resources acquisition and distribution

While it may seem somewhat technical for this for this tutorial, it is worth highlighting an intriguing implementation detail. If you look at the constructor of the *Processor* class, you will notice that you can provide some resources, like the Timer and UserInterface even though we have never done this so far. The idea is that when you execute a single processor, it is fully responsible of creating the required resources by itself, using them during the execution and then closing them when finishing.

Just as an example, consider the case of the use of a database interface. The processor is opening the connection, doing all the needed transactions and finally closing the connection. This approach is also very practical because it is much easier to keep expectations under control.

If you run a *ProcessorList*, you may want to move the responsibility of handling the resources from the single Processor to the output ProcessorList. This approach allows all processors to share the same resources efficiently, eliminating the need to repeatedly open and close the database connection each time the ProcessorList advances to the next item.

You do not have to care about this shift in responsibility, it is automatically done behind the scene when you add a processor to the processor lists.

5.3 What's next

In this part we have seen how we can chain the execution of any number of processors all sharing the same resources. Moreover, we have seen how we can change the looping among different processors using the exit status.

Now it is time to move forward and see how you can add your own processor library!

CHAPTER 6

PLUGINS: IMPORT YOUR PROCESSORS

We are almost half-way through our tutorial on MAFw. We have learned what a *Processor* is and how we can create our own processors just by subclassing the base class. With the *ProcessorList* we have seen an easy way to chain the execution of many processors. So by now, you might be tempted to open your IDE and start coding stuff... But please hold on your horses for another minute, MAFw has much more to offer.

In the previous pages of this tutorial, we always coded our *processors* and then executed them manually, either from the python console or via an ad-hoc script. This is already helpful, but not really practical, because everytime you want to change which *processors* are executed or one of their parameters, we need to change the script code. This seems very likely the job for a configuration file, doesn't it?

MAFw is providing you with exactly this capability, a generic script that will read your configuration file (**steering file** in the MAFw language) and execute it. Here comes the problem, *how can the execution framework load your [processors] since they are not living in the same library?*

To this problem there are two solutions: a quick and dirty and a clean and effective! Of course you will have the tendency to prefer the first one, but really consider the benefits of the second one before giving up reading!

6.1 The quick and dirty: code inside MAFw

When you install MAFw on your system, you can always do it in development mode (for pip this corresponds to the *-e* option). In this way, you will be allowed to code your *processors* library directly inside your locally installed MAFw.

It is very quick, because you can start coding right away, but it is also rather dirty, because it will be much harder to install MAFw updates and after a while the project will become too big and messy.

For these and many other reasons, we are convinced that this is not the right way to go.

6.2 The clean and effective: turn your library into a plugin

It may seem more challenging than it actually is, but in reality, it is quite simple. MAFw uses the plugin system developed by *pluggy*, that is very powerful and at the same time relatively easy to deploy.

After you have installed MAFw in your system (you can do it in development mode if you want to contribute!), create a new project to store your processor. How to create a project is well described in details [here](#) with a step by step guide. You don't need to upload your package to PyPI, you can stop as soon as you are able to build a wheel of your package and install it in a virtual environment.

To achieve this, you would need to have a *pyproject.toml* file, with the project metadata, the list of dependencies and other things. Here below is an example of what you should have:

Listing 6.1: An example of pyproject.toml for a plugin

```

1 [build-system]
2 requires = ["hatchling"]
3 build-backend = "hatchling.build"
4
5 [project]
6 name = "fantastic_analysis"
7 dynamic = ["version"]
8 description = 'My processor library'
9 readme = "README.md"
10 requires-python = ">=3.8"
11 license = "MIT"
12 keywords = []
13 authors = [
14     { name = "Surname Name", email = "this.is.me@my.domain.com" },
15 ]
16 classifiers = [
17     "Development Status :: 4 - Beta",
18     "Programming Language :: Python",
19     "Programming Language :: Python :: 3.8",
20     "Programming Language :: Python :: 3.9",
21     "Programming Language :: Python :: 3.10",
22     "Programming Language :: Python :: 3.11",
23     "Programming Language :: Python :: 3.12",
24     "Programming Language :: Python :: Implementation :: CPython",
25     "Programming Language :: Python :: Implementation :: PyPy",
26 ]
27
28 dependencies = ['mafw'] # plus all your other dependencies
29
30 [project.urls]
31 Documentation = "https://github.com/..."
32 Issues = "https://github.com/..."
33 Source = "https://github.com/."
34
35 ## THIS IS THE KEY PART
36 ## -----
37 # you can add as many lines you want to the table.
38 # always use unique names for the entry.
39
40 [project.entry-points.'mafw']
41 fantastic_analysis_plugin = 'fantastic_analysis.plugins'
42
43 [tool.hatch.version]
44 path = "src/fantastic_analysis/__about__.py"

```

Particularly important are lines 40 and 41. There you declare that when installed in an environment, your package is providing a plugin for MAFw, in particular this plugin that you named *fantastic_analysis_plugin* is located inside your package in a file named *plugins.py*.

So now, let us have a look at what do you have to have inside this file:

Listing 6.2: Example of plugins.py

```

1 """
2 This is the module that will be exposed via the entry point declaration.
3

```

(continues on next page)

(continued from previous page)

```
4 Make sure to have all processors that you need to export in the list.
5 """
6
7 import mafw
8 from mafw.processor import Processor
9 from fantastic_analysis import my_processor_lib
10
11 @mafw.mafw_hookimpl
12 def register_processors() -> list[mafw.processor.Processor]:
13     return [my_processor_lib.Processor1, my_processor_lib.Processor2]
```

In this python file you need to import mafw in order to have access to the `mafw_hookimpl` decorator. This is a marker for pluggy so that it knows that this function should return something for the host application. The second import statement is needed only for the typing of the return value. The third one is to import your library with the processors. It can be one or as many as you have.

Inside the returned list, just put all the processors you want to export and the trick is done!

Now install in development mode your package and MAFw will be able to access all your processor library, without you having to do anything else!

6.3 What's next

If you are still reading this tutorial, it means that you are aware of the benefits that MAFw can bring to your daily analysis task. And the best part is, there is even more to discover! On the following section, you will discover how to efficiently manage your complex analytical process by simply listing all the sequential steps in a steering file.

CHAPTER 7

RUN YOUR PROCESSORS FROM THE COMMAND LINE

Now that you have a wonderful library of *processors* that are exported as plugins to MAFw, you can benefit from an additional *bonus*, the possibility to execute them without having to code any other line!

When you have installed *MAFw* in your environment, an executable file *maf*w has been put in your path, so that if you open a console and you activate the python environment, you can run it.

Just try something like this!

```
(env) C:\>maf -v
maf, version 0.0.5
```

If you get this message (the version number may vary), you are good to go. Otherwise, check that everything was installed correctly and try again.

This executable is rather powerful and can launch several commands (see the full documentation [here](#)), but for the time we will focus on three.

7.1 List the processors

Before doing anything else, it is important to check that your processor library has been properly loaded in MAFw. To check that, run the following command from the console.

```
(env) C:\>maf list
```

Available processors (External Plugins = True)

Processor Name	Package Name	Module
AccumulatorProcessor	maf	examples.sum_processor
GaussAdder	maf	examples.sum_processor
ModifyLoopProcessor	maf	examples.loop_modifier
PrimeFinder	maf_plugin_test	prime_finder

Total processors = 4, internal = 3, external = 1

The output should be a table similar to the one above, containing a list of known processors in the first column, the package where they come from in the second and the module in the third one.

If you have exported your processors correctly the total external processors should be different from zero.

7.2 Prepare a generic steering file

The next step is to prepare a steering file. We do not want you to open your text editor and start preparing the TOML file from scratch. MAFw can prepare a generic template for you that you have to adapt to your needs before running.

```
(env) C:\>maf w steering steering-file.toml
A generic steering file has been saved in steering-file.toml.
Open it in your favourite text editor, change the processors_to_run list and save it.

To execute it launch: maf w run steering-file.toml.
```

This command is generating a generic steering file with a list of all processors available (internal to MAFw and imported from your library) with all their configurable parameters. You can pass the option `-show` to display the generated file on the console, or even more practical the `-open-editor` to open the steering file in your default editor so that you can immediately customize it.

Here below is an example of what you will get. The content may vary, but the ideas are the same.

Listing 7.1: An example of steering file.

```
1 # MAFw steering file generated on 2024-11-21 16:54:22.455297
2
3 # uncomment the line below and insert the processors you want to run from the
4   ↪available processor list
5 # processors_to_run = []
6
7 # customise the name of the analysis
8 analysis_name = "maf w analysis"
9 analysis_description = "Summing up numbers"
10 available_processors = ["AccumulatorProcessor", "GaussAdder", "ModifyLoopProcessor"]
11
12 [AccumulatorProcessor] # A processor to calculate the sum of the first n values via a
13   ↪looping approach.
14 last_value = 100 # Last value of the series
15
16 [GaussAdder] # A processor to calculate the sum of the first n values via the so
17   ↪called *Gauss formula*.
18 last_value = 100 # Last value of the series.
19
20 [ModifyLoopProcessor] # Example processor demonstrating how it is possible to change
21   ↪the looping structure.
22 item_to_abort = 65 # Item to abort
23 items_to_skip = [12, 16, 25] # List of items to be skipped.
24 total_item = 100 # Total item in the loop.
25
26 [UserInterface] # Specify UI options
27 interface = "rich" # Default "rich", backup "console"
```

The TOML file contains some tables and arrays that should be rather easy to understand. There is an important array at line 4 that you need to uncomment and to fill in with the *processors* that you want to run. For all known processors, there is table with all the configurable parameters, you can change them to your wishes. You do not need to remove the configuration information of the processors that you do not intend to use. They will be simply ignored.

If you have commented your *Processor* classes properly, the short description will appear next to the table.

You can customize the analysis name and short description to make the execution of the steering file a bit nicer. In the *UserInterface* section, you can change the way the processor will interact with you. If you are running the steering file on a headless HPC, then you do not need to have any fancy output. If you like to see progress bars and

spinners moving, then select rich as your interface.

7.3 Run the steering file

Here we are, finally! Now we can run our analysis and enjoy the first results.

From the console type the following command:

```
(env) C:\>maf run steering-file.toml
```

and enjoy the thrill of seeing your data being processed. From now on, you have a new buddy helping you with the analysis of your data, giving you more time to generate nice and interesting plots because MAFw will do all the boring stuff for you!

7.4 Debugging your processors

When you execute your processors from the command line using *maf* and something does not work as expected, you may want to debug your processors to fix the problem. Normally you would turn back to your IDE and click on the debug run button to follow your code line by line and see where it fails.

The problem here is that the module where you have coded your processor is not *running*, it just contains the definition of the processor! You would need to create an instance of your processor and to execute it in a standalone way, but then you would have to set the parameters manually. In other words, it is going to be slow and problematic and MAFw wants to alleviate your problems, not to create new ones!

The best way is to run the mafw executable in debug mode from your IDE inserting a breakpoint in your processor code.

In PyCharm, you would create a new python configuration launching a module instead of a script. See the screenshot below.

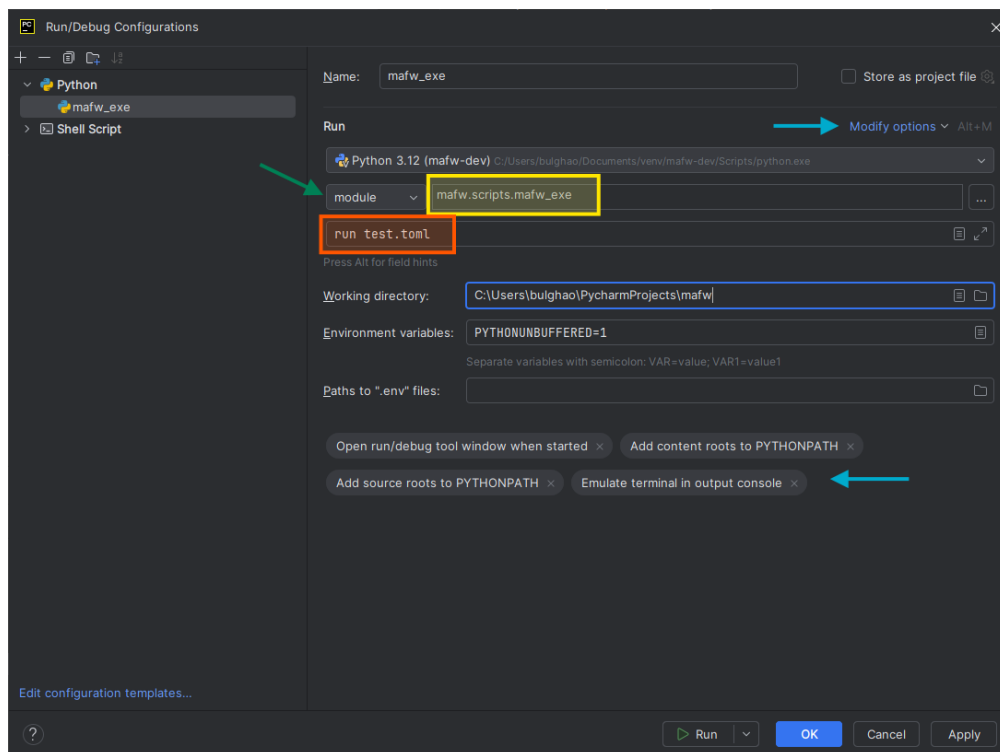


Fig. 7.1: The PyCharm configuration for running *mafw_exe* from your plugin module. It can be particularly helpful to execute your processor in debug mode.

The three most important points have been highlighted. First of all, you need to select module and not script from the drop down menu indicated by the green arrow. In the text field next to it (yellow in the screenshot), write `mafwscrip`. MAFw is installed in the environment where your plugins are living, so PyCharm will complete the module name while writing. In the text field below (orange in the screenshot) write the command line options as you were doing from the console. Optionally, if you want to enjoy the colorful output of the *RichInterface*, click on the *Modify options* (top blue arrow) and add *Emulate terminal in output console*. The corresponding label will appear in the bottom part as indicated by the blue arrow.

7.5 The exit code

The MAFw executable is always releasing an exit code when its execution is terminated. This exit code follows the standard convention of being 0 if the process was successful. Any other number is a symptom that something went wrong during the execution.

You will get an exit code different from 0 whenever, for example, you have not provided all the required parameters or they are invalid. You can also get an exit code different from 0 if your processor failed its execution.

Keep in mind that the *MAFwApplication* is actually dynamically generating and processing a *ProcessorList* according to the information stored in the *steering file*. The exit code for the whole app is linked to the output of the last executed processor in the list.

The reason for this is clearly explained. We have seen in a previous *section* that the execution of a *ProcessorList* can be modified using the *ProcessorExitStatus*. This parameter can assume three different values: **Successful**, **Failed** and **Aborted**. In MAFw convention, when a processor is aborted, then the execution of the processor list is immediately interrupted, thus resulting in the MAFw executable to return a no zero exit code.

On the contrary, when a processor exit status is Failed, the execution of the processor list is continued until the end of the list or the first processor terminating with an Abort status. The MAFw executable exit code is affected by the the processor exit status of the last executed processor, so it might be zero even if there were several processors with Failed exit status in the middle of the list.

The use of the exit code is particularly relevant if you want to execute MAFw on a headless cluster. You can, for example, trigger retry or logging actions when the execution fails.

7.6 Commands abbreviation

The MAFw executable is based on the very powerful *click* and allows the execution of complex nested commands with options and arguments.

You can get the list of available commands (for each level) using the `-h` option and keep in mind that MAFw allows you to shorten the command as long as the command selection will be unique, so actually `mafw list` and `mafw l` are both accepted and produce the same effect.

7.7 What's next

You believed we were done? Not really. We have just started, this was just the appetizer! Get ready for the next course, when we will enjoy the power of relational databases!

CHAPTER 8

DATABASE: YOUR NEW BUDDY!

As a scientist you are used to work with many software tools and often you need to write your own programs because the ones available do not really match your needs. Databases are not so common among scientists, we do not really understand why, but nevertheless their strength is beyond question.

The demonstration of the power and usefulness of a database assisted analysis framework will become clear and evident during our *tutorial* where we will build a small analytical experiment from scratch one step at the time.

For the time being, let us concentrate a bit on the technical aspect without delving too deeply into the specifics.

8.1 Database: one name, many implementations

Database is really a generic term and, from a technical point of view, one should try to be more specific and define better what kind of database we are talking about. We are not just referring to the brand or the producing software house: there are indeed different database architectures, but the one best suited for our application is the **relational database**, where each entity can be related directly to one or more other entities. You can read about relational databases on [wikipedia](#) for example, and if you find it too complicated, have a look at [this shorter and easier version](#).

The software market is full of relational databases, from the very common **MySQL** to the very famous **ORACLE**, passing through the open source **PostgreSQL** to finish with **SQLite** the simplest and most used database in the world. As you may have guessed from their names, they are all sharing the same query language (*SQL*: structured query language), making it rather simple to have an abstract interface, that is to say a layer in between the user and the actual database that allows your code to work in the same way irrespectively of the specific implementation.

8.2 Peewee: a simple and small ORM

Of course there are also lots of different abstract interfaces, some more performing than the others. We have selected: [link:peewee](#), because it is lightweight, easy to understand and to use, and it works with several different implementations.

[peewee](#) is a ORM (promised this is the last acronym for this page!), it is to say an object relational mappers, or in simpler words a smart way to connect the tables in your database with python classes in your code. Have a look at [this interesting article](#) for a more detailed explanation.

[peewee](#) offers at least three different backends: **SQLite**, **MySQL** and **PostgreSQL**. If the size of your project is small to medium and the analysis is mainly performed on a single computer, then we recommend **SQLite**: the entire database will be living in a single file on your disc, eliminating the need for IT experts to set up a database server. If you are aiming for a much bigger project with distributed computing power, then the other two choices are probably equally good and your local IT helpdesk may suggest you what is best for your configuration and the available IT infrastructure. As you see, **MAFw** is always offering you a tailored solution!

Now, take a short break from this page, move to the [peewee](#) documentation and read the *Quickstart* section before coming back here.

8.2.1 Database drivers

If you have installed MAFw via pip without specifying the *all-db* optional features, then your python environment is very likely missing the python drivers to connect to MySQL and PostgreSQL. This is not a bug, but more a feature, because MAFw gives you the freedom to select the database implementation that fits your needs. Sqlite is natively supported by python, so you do not need to install anything extra, but if you want to use MySQL, PostgreSQL or any other DB supported by [peewee](#) than it is your responsibility to install in your environment the proper driver. [Here](#) you can find a list of DB driver compatible with [peewee](#).

If you want, you can install MAFw adding the *all-db* optional feature and in this way the standard MySQL and PostgreSQL drivers will also be installed.

8.3 One class for each table

As mentioned before, the goal of an ORM is to establish a link between a database table and a python class. You can use the class to retrieve existing rows or to add new ones, and as always, you do not need take care of the boring parts, like establishing the connection, creating the table and so on, because this is the task of MAFw and we do it gladly for you!

Let us have a look together to the following example. We want a processor that lists recursively all files starting from a given directory and adds the filenames and the file hash digests to a table in the database.

Let us start with some imports

```

1 import datetime
2 from pathlib import Path
3
4 from peewee import DateTimeField, IntegerField, TextField
5
6 import mafw.processor
7 from mafw.db.db_configurations import db_scheme, default_conf
8 from mafw.db.db_model import MAFwBaseModel
9 from mafw.decorators import database_required
10 from mafw.tools.file_tools import file_checksum

```

The crucial one is at line 8, where we import *MAFwBaseModel* that is the base model for all the tables we want to handle with MAFw. Your tables **must inherit** from that one, if you want the *Processor* to take care of handling the link between the class and the table in the database. At line 4, we import some classes from peewee, that will define the columns in our model class and consequently in our DB table.

Now let us create the model class for our table of files.

```

1 class File(MAFwBaseModel):
2     """The Model class representing the table in the database"""
3
4     filename = TextField(primary_key=True)
5     digest = TextField()
6     creation_date = DateTimeField()
7     file_size = IntegerField()

```

Note

A much better implementation of a similar class will be given *later on* demonstrating the power of custom defined fields.

As you can see the class definition is extremely simple. We define a class attribute for each of the columns we want to have in the table and we can choose the field type from a long list of [available ones](#) or we can even easily implement our own. The role of a field is to adapt the type of the column from the native python type to the native database type and vice versa.

Our table will have just four columns, but you can have as many as you like. We will have one text field with the full filename, another text containing the hexadecimal hashlib digest of the file, the creation date for which we will use a datetime field, and finally a `file_size` field of integer type. We will be using the filename column as a primary key, because there cannot be two files with the same filename. On the contrary, there might be two identical files having the same hash but different filenames. According to many good database experts, using a not numeric primary key is not a good practice, but for our small example it is very practical.

If you do not specify any primary key, the ORM will add an additional number auto-incrementing column for this purpose. If you want to specify multiple primary keys, [this](#) is what you should do. If you want to create a model without a primary key, [here](#) is what you need to do.

The ORM will define the actual name of the table in the database and the column names. You do not need to worry about this!

i What is the name of my table?

We have created a Model class so far, that is the link between a python object and a database table. But what is the name of the table in the database? And the name of the columns?

All those information belong to the so called model metadata and are stored in a inner class of the `MAFwBaseModel`. `Peewee` uses a standard way to name tables and columns and as long as your are happy with this you do not have to worry about.

In general, MAFw is following this convention: model classes are written with CamelCase, while database tables are in snake_case. In other words, if you have a model class named `RawData` this will be linked to a database table named `raw_data`.

If you want to customize this naming convention, you can assign a specific name to a field as [explained here](#) and for the table name take a look at [this page](#)

And now comes the processor that will be doing the real work, it is to say, filling the File table.

```

1 @database_required
2 class FillFileTableProcessor(mafw.processor.Processor):
3     """Processor to fill a table with the content of a directory"""
4
5     root_folder: Path = mafw.processor.ActiveParameter(
6         'root_folder', default=Path.cwd(), help_doc='The root folder for the file_
7     listing'
8     )
9
10    def __init__(self, *args, **kwargs):
11        """
12        Constructor parameter:
13
14        :param root_folder: ActiveParameter corresponding to the directory from
15        where to start the recursive search
16        :type root_folder: Path, Optional
17        """
18        super().__init__(*args, **kwargs)
19        self.data: list[dict] = []
20
21    def format_progress_message(self):
22        self.progress_message = f'Upserting {self.item.name}'

```

(continues on next page)

(continued from previous page)

```

23 def start(self):
24     """Starts the execution.
25
26     Be sure that the table corresponding to the File model exists.
27     It it does already exists, it is not a problem.
28     """
29     super().start()
30     self.database.create_tables((File,))
31
32 def get_items(self) -> list[Path]:
33     """Retrieves the list of files.
34
35     Insert or update the files from the root folder to the database
36
37     :return: The list of full filename
38     :rtype: list[Path]
39     """
40     file_list = []
41     if self.root_folder.is_file():
42         return [self.root_folder]
43     elif self.root_folder.is_dir():
44         for f in self.root_folder.glob('* */*'):
45             if f.is_file():
46                 file_list.append(f)
47     else: # root_file is a glob
48         for f in self.root_folder.parent.glob(self.root_folder.name):
49             file_list.append(f)
50     return file_list
51
52 def process(self):
53     """Add all information to the data list"""
54     self.data.append(
55         dict(
56             filename=str(self.item),
57             digest=file_checksum(self.item),
58             file_size=self.item.stat().st_size,
59             creation_date=datetime.datetime.fromtimestamp(self.item.stat().st_
60 ←mtime),
61         )
62     )
63
64 def finish(self):
65     """Transfers all the data to the File table via an atomic transaction."""
66     with self.database.atomic():
67         File.insert_many(self.data).on_conflict_replace().execute()
68     super().finish()

```

The first thing to notice is at line 1, where we used the decorator `database_required()`. The use of this decorator is actually not compulsory, its goal is to raise an exception if the user tries to execute the processor without having a properly initialized database.

At line 30, in the `start` method we ask the database to create the table corresponding to our `File` model. If the table already exists, then nothing will happen.

In the `process` method we will store all the information we have collected from the files into a list and we interact with the database only in the `finish` method. At line 65, we use a context manager to create an `atomic` transaction and then, at line 66, we insert in the `File` all our entries and in case a row with the same primary key exists, then it is replaced.

We could have used several different insert approaches, here below are few examples:

```
# create an instance of File with all fields initialized
new_file = File(filename=str(self.item),
                digest=file_checksum(self.item),
                file_size=self.item.stat().st_size,
                creation_date=datetime.datetime.fromtimestamp(self.item.stat().st_
↳mtime))
new_file.save() # new_file is now stored in the database

# create an instance of File and add the fields later
new_file = File()
new_file.filename = str(self.item)
new_file.digest = file_checksum(self.item)
new_file.file_size = self.item.stat().st_size
new_file.creation_data = datetime.datetime.fromtimestamp(self.item.stat().st_mtime)
new_file.save()

# create and insert directly
new_file = File.create(filename=str(self.item),
                      digest=file_checksum(self.item),
                      file_size=self.item.stat().st_size,
                      creation_date=datetime.datetime.fromtimestamp(self.item.stat().st_
↳mtime))
```

To choice of approach to follow depends on various factor. Keep in mind that `peewee` operates by default in `auto commit mode`, meaning that for each database interaction, it creates a transaction to do the operation and it closes afterwards.

To be more performant from the database point of view, especially when you have several operations that can be grouped together, you can create an `atomic transaction` where the ORM will open one transaction only to perform all the required operations.

What we have done in the `finish` method is actually known as an `upsert`. It means that we will be inserting new items or updating them if they exist already.

8.4 Ready, go!

We have prepared the code, now we can try to run it. We can do it directly from a script

```
if __name__ == '__main__':
    database_conf = default_conf['sqlite']
    database_conf['URL'] = db_scheme['sqlite'] + str( Path.cwd() / Path('test.db'))

    db_proc = FillFileTableProcessor(root_folder =r'C:\Users\bulghao\Documents\
↳autorad-analysis\EdgeTrimming',
                                    database_conf=database_conf)

    db_proc.execute()
```

or in a more elegant way we can use the mafw app to run, but first we need to generate the proper steering file.

Console

```
c:\> mafw steering db-processor.toml
A generic steering file has been saved in db-processor.toml.
Open it in your favourite text editor, change the processors_to_run list and save it.

To execute it launch: mafw run db-processor.toml.
```

TOML

```

1 # MAFw steering file generated on 2024-11-24 22:13:38.248423
2
3 # uncomment the line below and insert the processors you want to run from the
4 ↪available processor list
5 processors_to_run = ["FillFileTableProcessor"]
6
7 # customise the name of the analysis
8 analysis_name = "maf_w_analysis"
9 analysis_description = "Using the DB"
10 available_processors = ["AccumulatorProcessor", "GaussAdder", "ModifyLoopProcessor",
11 ↪"FillFileTableProcessor", "PrimeFinder"]
12
13 [DBConfiguration]
14 URL = "sqlite:///file-db.db" # Change the protocol depending on the DB type. Update
15 ↪this file to the path of your DB.
16
17 [DBConfiguration.pragmas] # Leave these default values, unless you know what you are
18 ↪doing!
19 journal_mode = "wal"
20 cache_size = -64000
21 foreign_keys = 1
22 synchronous = 0
23
24 [FillFileTableProcessor] # Processor to fill a table with the content of a directory
25 root_folder = 'C:\Users\bulghao\PycharmProjects\maf_w' # The root folder for the file
26 ↪listing
27
28 [UserInterface] # Specify UI options
29 interface = "rich" # Default "rich", backup "console"

```

If you look at the steering file, you will notice that there is a `DBConfiguration` section, where we define the most important variable, it is to say the DB URL. This is not only specifying where the database can be found, but also the actual implementation of the database. In this case, it will be a sqlite database located in the file `file-db.db` inside the current directory.

There is also an additional sub table, named `pragmas`, containing advanced options for the sqlite DB. Unless you really know what you are doing, otherwise, you should leave them as they are.

In the following *Configuring other types of databases*, we will cover the case you want to use another DB implementation different from SQLite.

In the `FillFileTableProcessor` you can find the standard configuration of its processor parameters.

Now we are really ready to run our first DB processor and with a bit of luck, you should get your DB created and filled.

How to check the content of a DB?

There are several tools serving this purpose. One of those is `dbeaver` that works with all kind of databases offering an open source community version that you can download and install.

8.4.1 Configuring other types of databases

In the previous example, we have seen how to configure a simple SQLite database. For this database, you just need to indicate in the URL field the path on the local disc where the database file is stored.

SQLite does not require any user name nor password and there are no other fields to be provided. Nevertheless, it is worth adding the previously mentioned `pragmas` section to assure the best functionality of peewee.

In the case of MySQL and PostgreSQL, the URL should point to the server where the database is running. This could be the localhost but also any other network destination. Along with the server destination, you need also to specify the port, the database name, the user name and the password to establish the connection.

Of course, it is not a good idea to write your database password as plain text in a steering file that might be shared among colleagues or even worse included in a Git repository. To avoid this security issue, it is recommended to follow some other authentication approach.

Both MySQL and PostgreSQL offers the possibility to store the password in a separate file, that, at least in linux, should have a very limited access right. Have a look at the exemplary steering files with the corresponding password files here below.

SQLite

```
[DBConfiguration]
URL = "sqlite:///file-db.db" # change the filename to the absolute path of the db

[DBConfiguration.pragmas] # Leave these default values, unless you know what you are
↳doing!
journal_mode = "wal"
cache_size = -64000
foreign_keys = 1
synchronous = 0
```

PostgreSQL

```
[DBConfiguration]
# Update the database server and the database name to reflect your configuration
URL = "postgresql://database-server:5432/database-name"

# change it to your username
user = 'username'

# if you want, you can leave the pragmas section from the SQLite default
↳configuration because it
# want be used.
```

Instruction on how to create a PostgreSQL password file are provided [here](#). This is an example:

```
database-server:5432:database-name:username:password
```

MySQL

```
[DBConfiguration]
# Update the database server and the database name to reflect your configuration
URL = "mysql://database-server:3306/database-name"

# update to specify your username
user = 'username'

# update to specify the password file
read_default_file = '~/my.cnf'

# if you want, you can leave the pragmas section from the SQLite default
↳configuration because it
# want be used.
```

Instruction on how to create a MySQL password file are provided [here](#). This is an example:

```
[client]
user=username
password=password
host=database-server
```

8.5 Triggers: when the database works on its own

In the next paragraphs we will spend a few minutes understanding the roles of Triggers. Those are database entities performing some actions in response of specific events. You can have, for example, a trigger that is inserting a row in TableB whenever a row is inserted in TableA. If you are not really familiar with triggers, this is a [brief introduction](#).

Triggers are very handy for many applications, and in our *tutorial* we will see an interesting case, but they tend to struggle with ORM in general. In fact, no ORM system is natively supporting triggers. The reason is very simple. In an ORM, the application (the python code, if you wish) is the main actor and the database is just playing the role of the passive buddy. From the point of view of an ORM based application, if you want to have a trigger, then just write the needed lines of python code to have the actions performed in the other tables. It makes totally sense, you have only one active player and it simplifies the debugging because if something is going wrong, it can only be a fault of the application.

The standard implementation of trigger-like functions with ORM is to use [signals](#), where you can have callbacks called before and after high level ORM APIs calls to the underlying database. Signals are good, but they are not free from disadvantages: at a first glance, they look like a neat solution, but as soon as the number of callbacks is growing, it may become difficult to follow a linear path in the application debugging. Second, if you do a change in the database from another application, like the [dbeaver](#) browser, then none of your codified triggers will be executed. Moreover in the case of [Peewee](#), signals work only on Model instances, so all bulk inserts and updates are excluded.

Having triggers in the database would assure that irrespectively of the source of the change, they will always be executed, but as mentioned above, the user will have to be more careful in the debug because also the database is now playing an active role.

We let you decide what is the best solution. If you want to follow the pure ORM approach, then all models inheriting from `MAFwBaseModel` have the possibility to use signals. If you want to have triggers, you can also do so. An example for both approaches is shown here below.

8.5.1 The signal approach

As mentioned above, the signal approach is the favourite one if you plan to make all changes to the database only via your python code. If you are considering making changes also from other applications, then you should better use the trigger approach.

Another limitation is that only model instances emit signals. Everytime you use a *classmethod* of a Model, then no signals will be emitted.

The signal dispatching pattern functionality is achieved by linking the signal emitted by a sender in some specific circumstances to a handler that is receiving this signal and performing some additional operations (not necessarily database operations).

Every model class has five different signals:

1. **pre_save**: emitted just before that a model instance is saved;
2. **post_save**: emitted just after the saving of a model instance in the DB;
3. **pre_delete**: emitted just before deleting a model instance in the DB;
4. **post_delete**: emitted just after deleting a model instance from the DB;
5. **pre_init**: emitted just after the init method of the class is invoked. Note that the *pre* is actually a *post* in the case of init.

Let us try to understand how this works with the next example.

Listing 8.1: A test with signal

```

1 class MyTable(MAFwBaseModel):
2     id_ = AutoField(primary_key=True)
3     integer = IntegerField()
4     float_num = FloatField()
5
6 class TargetTable(MAFwBaseModel):
7     id_ = ForeignKeyField(MyTable, on_delete='cascade', primary_key=True, backref=
↳ 'half')
8     another_float_num = FloatField()
9
10    @post_save(sender=MyTable, name='my_table_after_save_handler')
11    def post_save_of_my_table(sender: type(MAFwBaseModel), instance: MAFwBaseModel,
↳ created: bool):
12        """
13        Handler for the post save signal.
14
15        The post_save decorator is taking care of making the connection.
16        The sender specified in the decorator argument is assuring that only signals
↳ generated from MyClass will be
17        dispatched to this handler.
18
19        The name in the decorator is optional and can be use if we want to disconnect
↳ the signal from the handler.
20
21        :param sender: The Model class sending this signal.
22        :type sender: type(Model)
23        :param instance: The actual instance sending the signal.
24        :type instance: Model
25        :param created: Bool flag if the instance has been created.
26        :type created: bool
27        """
28        TargetTable.insert({'id_id': instance.id, 'another_float_num': instance.
↳ float_num / 2}).execute()
29
30 database: Database = SqliteDatabase(':memory:', pragmas=default_conf['sqlite'][
↳ 'pragmas'])
31 database.connect()
32 database_proxy.initialize(database)
33 database.create_tables([MyTable, TargetTable], safe=True)
34
35 MyTable.delete().execute()
36 TargetTable.delete().execute()
37
38 # insert a single row in MyTable with the save method.
39 my_table = MyTable()
40 my_table.integer = 20
41 my_table.float_num = 32.16
42 my_table.save()
43 # after the save query is done, the signal mechanism will call the
44 # post_save_trigger_of_my_table and perform an insert on the target
45 # table as well.
46 assert MyTable.select().count() == 1
47 assert TargetTable.select().count() == 1
48
49 # add some bulk data to MyTable

```

(continues on next page)

(continued from previous page)

```
50 data = []
51 for i in range(100):
52     data.append(dict(integer=random.randint(i, 10 * i), float_num=random.gauss(i, 2 *
↪ i)))
53 MyTable.insert_many(data).execute()
54 # this is done via the Model class and not via a concrete instance of the Model, so
↪ no signals will be emitted.
55
56 assert MyTable.select().count() == 101
57 assert TargetTable.select().count() == 1
```

We created two tables linked via a foreign key. The goal is that everytime we fill in a row in MyTable, a row is also added to TargetTable with the same id but where the value of another_float_num is just half of the original float_num. The example is stupid, but it is good enough for our demonstration.

The signal part is coded in the lines from 11 to 28 (mainly doc strings). We use the `post_save` decorator to connect MyTable to the `post_save_of_my_table` function where an insert in the TargetTable will be made.

The code is rather simple to follow. Just to be sure, we empty the two tables, then we create an instance of the MyTable model, to set the integer and the float_num column. When we save the new row, the `post_save` signal of MyTable is emitted and the handler is reacting by creating an entry in the TargetTable as well. In fact the number of rows of both tables are equal to 1.

What happens later is to demonstrate the weak point of signals. At line 53, we insert several rows via a `insert_many`. It must be noted that the `insert_many` is a classmethod applied directly to the model class. The consequence is that the signal handler will not be invoked and no extra rows will be added to the TargetTable.

8.5.2 The trigger approach

In order to use a trigger you need to create it. This is an entity that lives in the database, so you would need the database itself to create it.

MAFw is providing a *Trigger* class that helps you in creating the required SQL query that needed to be issued in order to create the trigger. Once it is created it will operate continuously.

If you have a look at the [CREATE TRIGGER SQL command](#) you will see that it starts with the definition of when the trigger is entering into play (BEFORE/AFTER) and which operation (INSERT/DELETE/UPDATE) of which table. Then there is a section enclosed by the BEGIN and END keywords, where you can have as many SQL queries as you like.

The same structure is reproduced in the *Trigger* class. In the constructor, we will pass the arguments related to the configuration of the trigger itself. Then you can add as many SQL statement as you wish.

Python

Listing 8.2: python Trigger class

```

1 from mafw.db.db_model import Trigger
2
3 new_trigger = Trigger('trigger_after_update', (TriggerWhen.After,
4     TriggerAction.Update), 'my_table', safe=False, for_each_row=False)
5 new_trigger.add_sql('INSERT INTO another_table (col1, col2) VALUES (1, 2)')
6 new_trigger.add_sql('INSERT INTO third_table (col1, col2) VALUES (2, 3)')
7 new_trigger.create()

```

SQL

Listing 8.3: emitted SQL

```

1 CREATE TRIGGER trigger_after_update
2 AFTER UPDATE ON my_table
3
4 BEGIN
5     INSERT INTO another_table (col1, col2) VALUES (1, 2);
6     INSERT INTO third_table (col1, col2) VALUES (2, 3);
7 END;

```

Now let us have a look at how you can use this, following one of our test benches.

Standalone triggers

Listing 8.4: A test Trigger created manually

```

1 def test_manually_created_trigger():
2     class MyTable(MAFwBaseModel):
3         id_ = AutoField(primary_key=True)
4         integer = IntegerField()
5         float_num = FloatField()
6
7     class TargetTable(MAFwBaseModel):
8         id_ = ForeignKeyField(MyTable, on_delete='cascade', primary_key=True, backref=
9     ↪ 'half')
10        half_float_num = FloatField()
11
12        database: Database = SQLiteDatabase(':memory:', pragmas=default_conf['sqlite']
13     ↪ 'pragmas'])
14        database.connect()
15        database_proxy.initialize(database)
16        database.create_tables([MyTable, TargetTable], safe=True)
17
18        MyTable.delete().execute()
19        TargetTable.delete().execute()
20
21        # manually create a trigger
22        trig = Trigger('mytable_after_insert', (TriggerWhen.After, TriggerAction.Insert),
23     ↪ MyTable, safe=True)
24        trig.add_sql('INSERT INTO target_table (id_id, half_float_num) VALUES (NEW.id_,
25     ↪ NEW.float_num / 2)')
26        database.execute_sql(trig.create())
27
28        # add some data for testing to the first table
29        data = []

```

(continues on next page)

(continued from previous page)

```

26     for i in range(100):
27         data.append(dict(integer=random.randint(i, 10 * i), float_num=random.gauss(i,
↪ 2 * i)))
28     MyTable.insert_many(data).execute()
29
30     # check that the target table got the right entries
31     for row in MyTable.select(MyTable.float_num, TargetTable.half_float_num).
↪ join(TargetTable).namedtuples():
32         assert row.float_num == 2 * row.half_float_num
33
34     assert MyTable.select().count() == TargetTable.select().count()

```

In lines 20 - 22, we create a trigger and we ask the database to execute the generated SQL statement.

We insert 100 rows using the insert many class method and the trigger is doing its job in the background filling the other table. We can check that the values in the two tables are matching our expectations.

The drawback of this approach is that you may have triggers created all around your code, making your code a bit messy.

Model embedded triggers

An alternative approach is to define the trigger within the Model class, allowing it to be created simultaneously with model table. This is demonstrated in the code example below.

Listing 8.5: A test Trigger created within the Model

```

1  # the trigger is directly defined in the class.
2  class MyTable(MAFwBaseModel):
3      id_ = AutoField(primary_key=True)
4      integer = IntegerField()
5      float_num = FloatField()
6
7      @classmethod
8      def triggers(cls):
9          return [
10             Trigger('mytable_after_insert', (TriggerWhen.After, TriggerAction.Insert),
↪ cls, safe=True).add_sql(
11                 'INSERT INTO target_table (id_id, half_float_num) VALUES (NEW.id_,
↪ NEW.float_num / 2)'
12             )
13         ]
14
15  class TargetTable(MAFwBaseModel):
16      id_ = ForeignKeyField(MyTable, on_delete='cascade', primary_key=True, backref=
↪ 'half')
17      half_float_num = FloatField()
18
19  database: Database = SqliteDatabase(':memory:', pragmas=default_conf['sqlite']
↪ ['pragmas'])
20  database.connect()
21  database_proxy.initialize(database)
22  database.create_tables([MyTable, TargetTable], safe=True)

```

This approach is much cleaner. The Trigger is stored directly in the Model (lines 8 - 13). In the specific case, the triggers method returned one trigger only, but you can return as many as you like. When the tables are created (line 22), all the triggers will also be created.

In the example above, you have written the SQL statement directly, but nobody is preventing you to use `peewee`

queries for this purpose. See below, how exactly the same trigger might be re-written, using an insert statement:

Listing 8.6: A test Trigger created within the Model using an Insert statement

```

1 class MyTable(MAFwBaseModel):
2     id_ = AutoField(primary_key=True)
3     integer = IntegerField()
4     float_num = FloatField()
5
6     @classmethod
7     def triggers(cls):
8         trigger = Trigger('mytable_after_insert', (TriggerWhen.After, TriggerAction.
↪Insert), cls, safe=True)
9         sql = TargetTable.insert(id_=SQL('NEW.id_'), half_float_num=SQL('NEW.float_
↪num/2'))
10        trigger.add_sql(sql)
11        return [trigger]
12
13    class Meta:
14        depends_on = [TargetTable]

```

The key point here is at line 9, where the actual insert statement is generated by `peewee` (just for your information, you have generated the statement, but you have not *execute it*) and added to the existing trigger.

In the last two highlighted lines, we are overloading the Meta class, specifying that `MyTable`, depends on `TargetTable`, so that when the `create_tables` is issued, they are built in the right order. This is not necessary if you follow the previous approach because the trigger will be very likely executed only after that the tables have been created.

Warning

Even though starting from MAFw release v1.1.0, triggers are now properly generated for the three main *database backends*, its use has been deeply tested only with SQLite. For this reason, we (MAFw developers) encourage the user community to work also with other DBs and, in case, submit bugs or feature request.

Disabling triggers

Not all database implementations provide the same option to temporarily disable one or more triggers. In order to cope with this limitation, MAFw is providing a general solution that is always working independently of the concrete implementation of the database.

The standard SQL trigger definition allows to have one or more WHEN clauses¹, meaning that the firing of a trigger script might be limited to the case in which some other external conditions are met.

In order to achieve that, we use one of our *standard tables*, that are automatically created in every MAFw database.

This is the `TriggerStatus` table as you can see it in the snippet below:

Listing 8.7: `TriggerStatus` model

```

class TriggerStatus(StandardTable):
    """A Model for the trigger status"""

    trigger_type_id = AutoField(primary_key=True, help_text='Primary key')
    trigger_type = TextField(
        help_text='You can use it to specify the type (DELETE/INSERT/UPDATE) or the_
↪name of a specific trigger'

```

(continues on next page)

¹ MySQL does not directly support adding WHEN conditions to the trigger, but a similar behaviour is obtainable using an IF statement in the trigger SQL body. This adaptation is automatically implemented by the *MySQLDialect*.

(continued from previous page)

```

)
    status = BooleanField(default=True, help_text='False (0) = disable / True (1) =
↳enable')

    # noinspection PyProtectedMember
    @classmethod
    def init(cls) -> None:
        """Resets all triggers to enable when the database connection is opened."""
        data = [
            dict(trigger_type_id=1, trigger_type='DELETE', status=True),
            dict(trigger_type_id=2, trigger_type='INSERT', status=True),
            dict(trigger_type_id=3, trigger_type='UPDATE', status=True),
            dict(trigger_type_id=4, trigger_type='DELETE_FILES', status=True),
        ]

        # this is used just to make mypy happy
        # cls and meta_cls are exactly the same thing
        meta_cls = cast(PeeweeModelWithMeta, cls)

        db_proxy = meta_cls._meta.database
        if isinstance(db_proxy, peewee.DatabaseProxy):
            db = cast(peewee.Database, db_proxy.obj)
        else:
            db = cast(peewee.Database, db_proxy)

        if isinstance(db, peewee.PostgresqlDatabase):
            cls.insert_many(data).on_conflict(
                'update', conflict_target=[cls.trigger_type_id], update={cls.status:
↳True}
            ).execute()
        else:
            cls.insert_many(data).on_conflict_replace().execute()

```

You can use the `trigger_type` column to specify a generic family of triggers (DELETE/INSERT/UPDATE) or the name of a specific trigger. By default a trigger is active (`status = 1`), but you can easily disable it by changing its status to 0.

To use this functionality, the Trigger definition should include a WHEN clause as described in this modified model definition.

Listing 8.8: Trigger definition with when conditions.

```

class MyTable(MAFwBaseModel):
    id_ = AutoField(primary_key=True)
    integer = IntegerField()
    float_num = FloatField()

    @classmethod
    def triggers(cls):
        return [
            Trigger('mytable_after_insert', (TriggerWhen.After, TriggerAction.Insert),
↳ cls, safe=True)
                .add_sql('INSERT INTO target_table (id_id, half_float_num) VALUES (NEW.
↳id_, NEW.float_num / 2)')
                .add_when('1 == (SELECT status FROM trigger_status WHERE trigger_type ==
↳"INSERT")')
        ]

```

To facilitate the temporary disabling of a specific trigger family, MAFw provides a special class *TriggerDisabler* that can be easily used as a context manager in your code. This is an ultra simplified snippet.

Listing 8.9: Use of a context manager to disable a trigger

```
with TriggerDisabler(trigger_type_id = 1):
    # do something without triggering the execution of any trigger of type 1
    # in case of exceptions thrown within the block, the context manager is restoring
    # the trigger status to 1.
```

8.5.3 Triggers on different databases

We have seen that Peewee provides an abstract interface that allows interaction with various SQL databases like MySQL, PostgreSQL, and SQLite.

This abstraction simplifies database operations by enabling the same codebase to work across different database backends, thanks to the common SQL language they all support. However, while these databases share SQL as their query language, they differ in how they handle certain features, such as triggers. Each database has its own peculiarities and syntax for defining and managing triggers, which can lead to inconsistencies when using a single approach across all databases.

To address this challenge, the MAFw introduced the *TriggerDialect* abstract class and three specific implementations for the main databases. Relying on the use of the *TriggerDialect* class, a syntactically correct SQL statement for the creation or removal of triggers is generated. But, MAFw cannot read the mind of the user (yet!) and given the very different behaviour of the databases, the operation of the triggers will be different.

Have a look at the table below for an illustrative comparison on how triggers are handled by the different databases.

Feature	MySQL	PostgreSQL	SQLite
Trigger Event	<ul style="list-style-type: none"> • INSERT • UPDATE • DELETE 	<ul style="list-style-type: none"> • INSERT • UPDATE • DELETE • TRUNCATE 	<ul style="list-style-type: none"> • INSERT • UPDATE • DELETE
Trigger Time	<ul style="list-style-type: none"> • BEFORE • AFTER 	<ul style="list-style-type: none"> • BEFORE • AFTER • INSTEAD OF 	<ul style="list-style-type: none"> • BEFORE • AFTER • INSTEAD OF
Activation	Row-level only	Row-level and statement-level	Row-level and statement-level
Implementation	BEGIN-END block with SQL statements (supports non-standard SQL like SET statements)	Functions written in PL/pgSQL, PL/Perl, PL/Python, etc.	BEGIN-END block with SQL statements
Trigger Per Event	Multiple triggers allowed ordered by creation time	Multiple triggers allowed ordered alphabetically by default, can be specified	Multiple triggers allowed but unspecified execution order
Privileges required	TRIGGER privilege on the table and SUPER or SYSTEM_VARIABLES_ADMIN for DEFINER	CREATE TRIGGER privilege on schema and TRIGGER privilege on table	No specific privilege model
Foreign Key Cascading	Cascaded foreign key actions do not activate triggers	Triggers are activated by cascaded foreign key actions	Triggers are activated by cascaded foreign key actions
Disabled/Enabled Trigger	Yes, using ALTER TABLE ... DISABLE/ENABLE TRIGGER	Yes, using ALTER TABLE ... DISABLE/ENABLE TRIGGER	No direct mechanism to disable

PostgreSQL offers the most comprehensive trigger functionality, with built-in support for both row-level and statement-level triggers, INSTEAD OF triggers for views, and the widest range of programming languages for implementation. Its trigger functions can be written in any supported procedural language, providing considerable flexibility.

MySQL implements triggers using SQL statements within BEGIN-END blocks and only supports row-level triggers. It allows non-standard SQL statements like SET within trigger bodies, making it somewhat more flexible for certain operations. A critical limitation is that MySQL triggers are not activated by cascaded foreign key actions, unlike the other databases. This is a strong limiting factor and the user should consider it when designing the database model to store their data. In this case, it might be convenient to not rely at all on the cascading operations, but to have dedicated triggers for this purpose.

SQLite provides more trigger capabilities than it might initially appear. While its implementation is simpler than PostgreSQL's, it supports both row-level and statement-level triggers (statement-level being the default if FOR EACH ROW is not specified). Like PostgreSQL, SQLite triggers are activated by cascaded foreign key actions, which creates an important behavioral difference compared to MySQL.

When designing database applications that may need to work across different database systems, these implementation differences can lead to subtle bugs, especially around foreign key cascading behavior. MySQL applications that rely on triggers not firing during cascaded operations might behave differently when migrated to PostgreSQL or SQLite. Similarly, applications that depend on statement-level triggers will need to be redesigned when moving from PostgreSQL or SQLite to MySQL.

All so said, even though MAFw provides a way to handle triggers creation and removal in the same way across all the databases, the user who wants to move from one DB implementation to the other should carefully review the content of the trigger body to ensure that the resulting behavior is what is expected.

8.6 Standard tables

In the previous section, we discussed a workaround implemented by MAFw to address the limitations of database backends that cannot temporarily disable trigger execution. This is achieved querying a table where the status of a specific trigger or a family of triggers can be toggled from active to inactive and vice-versa.

This *TriggerStatus* model is one of the so-called MAFw standard tables, it is to say models that will be silently created by the execution of a *Processor*. In other words, as soon as you connect to a database using the MAFw infrastructure, all those tables will be created so that all processors can benefit from them.

Along with the *TriggerStatus*, there are two other relevant standard tables: the *OrphanFile* and the *PlotterOutput*.

OrphanFile: the house of files without a row

This table can be profitably used in conjunction with Triggers. The user can define a trigger fired when a row in a table is deleted. The trigger will then insert all file references contained in the deleted row into the *OrphanFile* model.

The next time a processor (it does not matter which one) having access to the database is executed, it will query the full list of files from the *OrphanFile* and remove them.

This procedure is needed to avoid having files on your disc without a reference in the database. It is kind of a complementary cleaning up with respect to *another function* you will discover in a moment.

Additional details about this function are provided directly in the *API*.

PlotterOutput: where all figures go.

Plotters are special *Processor* subclasses with the goal of generating a graphical representation of some data you have produced in a previous step.

The output of a plotter is in many cases one or more figure files and instead of having to define a specific table to store just one line, MAFw is providing a common table for all plotters where they can store the reference to their output files linked to the plotter name.

It is very useful because it allows the user to skip the execution of a plotter if its output file already exists on disc.

Triggers are again very relevant, because when a change is made in the data used to generate a plotter output, then the corresponding rows in this table should also be removed, in order to force the regeneration of the output figures with the updated data.

The role of those tables is to support the functionality of the general infrastructure and not of a single processor. If your processor needs a specific table, then it is its responsibility to create it. If all your processors need to have access to certain tables, then you may consider having them added to the standard tables.

In this case, you can follow the same plugin approach you used for sharing your processors.

First, you need to create the model classes, possibly inheriting from `StandardTable` and then you need to export it in the plugin.py module.

Here below is an example:

my_std_tables.py

```
from peewee import AutoField, TextField
from mafw.db.std_tables import StandardTable

class AnotherExtraTable(StandardTable):
    main_id = AutoField(primary_key=True)
    interesting_field = TextField()

    @classmethod
    def init(cls)
        # implement here some optional initialisation code.
        # it will be performed everytime the connection to the database is opened.
    -> from MAFw.
        pass
```

plugins.py

```
import mafw
import my_package.my_std_table

@mafw.mafw_hookimpl
def register_standard_tables() -> list:
    return [my_package.my_std_table.AnotherExtraTable]
```

pyproject.toml

```
# your standard pyproject goes here
# be sure to add this entry point
[project.entry-points.'mafw']
my_package = 'my_package.plugins'
```

Note

Keep in mind, that your extra standard tables will be added only when your processor will be executed by the `mafw executable`, that means that if you write your own main script to execute a processor or a processor list, this will not load the plugin imported tables.

8.7 Custom fields

We have seen in a previous section that there are plenty of field types for you to build up your model classes and that it is also possible to add additional ones. We have made a few for you that are very useful from the point of view of MAFw. The full list is available [here](#).

The role of the database in MAFw is to support the input / output operation. You do not need to worry about specifying filenames or paths. Simply instruct the database to retrieve a list of items, and it will automatically provide the various processors with the necessary file paths for analysis.

With this in mind, we have created a `FileNameField`, that is the evolution of a text field accepting a Path object as a python type and converting it into a string for database storage. On top of `FileNameField`, we have made `FileNameListField` that can contain a list of filenames. This second one is more appropriate when your processor is generating a group of files as output. The filenames are stored in the database as a ';' separated string, and they are seen by the python application as a list of Path objects.

Similarly, we have also a `FileChecksumField` to store the string of hexadecimal characters corresponding to the checksum of a file (or a list of files). From the python side, you can assign either the checksum directly, as generated for example by `file_checksum()` or the path to the file, and the field will calculate the checksum automatically.

The `FileNameField` and `FileNameListField` accept an additional argument in their constructor, called `checksum_field`. If you set it to the name of a `FileChecksumField` in the same table, then you do not even have to set the value of the checksum field because this will be automatically calculated when the row is saved.

With these custom fields in mind, our initial definition of a `File` table, can be re-factored as follows:

```

1 from peewee import AutoField
2
3 from mafw.db.db_model import MAFwBaseModel
4 from mafw.db.fields import FileNameField, FileChecksumField
5
6 class File(MAFwBaseModel):
7     file_id = AutoField(primary_key=True, help_text='The primary key')
8     file_name = FileNameField(checksum_field='file_digest', help_text='The full_
↳ filename')
9     file_digest = FileChecksumField(help_text='The hex digest of the file')
```

Pay attention at the definition of the `file_name` field. The `FileNameField` constructor takes an optional parameter `checksum_field` that is actually pointing to the variable of the `FileChecksumField`.

You can use the two custom fields as normal, for example you can do:

```

1 new_file = File()
2 # you can assign a Path object.
3 new_file.file_name = Path('/path/to/some/file')
4 # the real checksum will be calculated automatically.
5 # this next line is totally optional, you can leave it out and it will work in the_
↳ same way.
6 new_file.file_digest = Path('/path/to/some/file')
```

The super power of these two custom fields is that you can remove useless rows from the database, just issuing one command.

8.7.1 Removing widow rows

Due to its I/O support, the database content should always remain aligned with the files on your disc. If you have a row in your database pointing to a missing file, this may cause troubles, because sooner or later, you will try to access this missing file causing an application crash.

In MAFw nomenclature, those rows are called *widows*, following a similar concept in [typesetting](#), because they are a fully valid database entry, but their data counter part on disc disappeared.

To avoid any problem with widow rows, MAFw is supplying a *function* that the processor can invoke in the start method on the Model classes used as input:

```
class MyProcessor(Processor):

    def start():
        super().start()
        remove_widow_db_rows(InputData)
```

The `remove_widow_db_rows()` will check that all the `FileNameField` fields in the table are pointing to existing files on disc. If not, then the row is removed from the database.

The function is not automatically called by any of the Processor super methods. It is up to the user to decide if and when to use it. Its recommended use is in the overload of the `start()` method or as a first action in the `get_items()` in the case of a *for loop* workflow, so that you are sure to re-generate the rows that have been removed.

8.7.2 Pruning orphan files

The opposite situation is when you have a file on disc that is not linked to an entry in the database anymore. This situation could be even more perilous than the previous one and may occur more frequently than you realize. The consequences of this mismatch can be severe, imagine that during the *testing / development phase* of your *Processor* you generate an output figure saved on disc. You then realize that the plot is wrong and you fix the bug and update the DB, but for some reasons you have forgotten to delete the figure file from the disc. Afterwards, while looking for the processor output, you find this file and believe it is a valid result and you use it for your publication. In order to prevent this to happen, you just have to follow some simple rules, and then the reliable machinery of MAFw will do the rest.

The key point is to use a specific trigger in every table that has a file name field. This trigger has to react before any delete query on such a table and inserting all `FileNameFields` or `FileNameListFields` in the `OrphanFile` table. You will see an example of such a trigger in the next paragraphs. This standard tables will be queried by the next processor being executed and during the start super method, all files in the `Orphan` table will be removed from the disc.

Let us try to understand this better with a step-by-step example. For simplicity, we have removed the import statements from the code snippet, but it should not be too difficult to understand the code anyway.

We begin with the declaration of our input model:

Listing 8.10: File model definition with trigger

```
class File(MAFwBaseModel):
    file_id = AutoField(primary_key=True, help_text='primary key')
    file_name = FileNameField(checksum_field='check_sum', help_text='the file name')
    check_sum = FileChecksumField(help_text='checksum')

    @classmethod
    def triggers(cls) -> list[Trigger]:
        file_delete_file = Trigger(
            'file_delete_file',
            (TriggerWhen.Before, TriggerAction.Delete),
            source_table=cls,
            safe=True,
            for_each_row=True,
        )
        file_delete_file.add_when('1 == (SELECT status FROM trigger_status WHERE_
->trigger_type = "DELETE_FILES")')
        file_delete_file.add_sql(OrphanFile.insert(filenamees=SQL('OLD.file_name'),_
->checksum=SQL('OLD.file_name'))))
        return [file_delete_file]
```

(continues on next page)

(continued from previous page)

```
class Meta:
    depends_on = [OrphanFile]
```

Here you see the trigger definition: it is a before delete type and when triggered it is adding the filename field to the OrphanFile table. It is important to notice that this trigger has a when condition and will only be executed when the trigger type DELETE_FILES is enabled. This is necessary for the pruning mechanism to work, just copy this line in your trigger definition.

And now let us define some fake processors. First we import some files into our model, then we remove some rows from the file table and finally other two processors, doing nothing but useful to demonstrate the effect of the orphan removal.

Listing 8.11: Some example processors

```
@database_required
class FileImporter(Processor):
    input_folder = ActiveParameter('input_folder', default=Path.cwd(), help_doc='From_
↪where to import')

    def __init__(self, *args, **kwargs):
        super().__init__(*args, looper=LoopType.SingleLoop, **kwargs)
        self.n_files: int = -1

    def start(self):
        super().start()
        self.database.create_tables([File])
        File.delete().execute()

    def process(self):
        data = [(f, f) for f in self.input_folder.glob('**/*.dat')]
        File.insert_many(data, fields=['file_name', 'check_sum']).execute()
        self.n_files = len(data)

    def finish(self):
        super().finish()
        if File.select().count() != self.n_files:
            self.processor_exit_status = ProcessorExitStatus.Failed

@database_required
class RowRemover(Processor):
    n_rows = ActiveParameter('n_rows', default=0, help_doc='How many rows to be_
↪removed')

    def __init__(self, *args, **kwargs):
        super().__init__(*args, looper=LoopType.SingleLoop, **kwargs)
        self.n_initial = 0

    def start(self):
        super().start()
        self.database.create_tables([File])

    def process(self):
        self.n_initial = File.select().count()
        query = File.select().order_by(fn.Random()).limit(self.n_rows).execute()
        ids = [q.file_id for q in query]
        File.delete().where(File.file_id.in_(ids)).execute()
```

(continues on next page)

(continued from previous page)

```

def finish(self):
    super().finish()
    if File.select().count() != self.n_initial - self.n_rows or OrphanFile.
↪select().count() != self.n_rows:
        self.processor_exit_status = ProcessorExitStatus.Failed

@orphan_protector
@database_required
class OrphanProtector(Processor):
    def __init__(self, *args, **kwargs):
        super().__init__(looper=LoopType.SingleLoop, *args, **kwargs)
        self.n_orphan = 0

    def start(self):
        self.n_orphan = OrphanFile.select().count()
        super().start()

    def finish(self):
        super().finish()
        if OrphanFile.select().count() != self.n_orphan:
            self.processor_exit_status = ProcessorExitStatus.Failed

@single_loop
class LazyProcessor(Processor):
    def finish(self):
        super().finish()
        if OrphanFile.select().count() != 0:
            self.processor_exit_status = ProcessorExitStatus.Failed

```

The **FileImporter**² is very simple, it reads all dat files in a directory and loads them in the File model along with their checksum. Just to be sure we empty the File model in the start and in the finish we check that the number of rows in File is the same as the number of files in the folder.

The **RowRemover** is getting an integer number of rows to be removed. Even though the File model is already created, it is a good practice to repeat the statement again in the start method. Then we select some random rows from File and we delete them. At this point, we have created some orphan files on disc without related rows in the DB. Finally (in the finish method), we verify that the number of remaining rows in the database aligns with our expectations. Additionally, we ensure that the trigger functioned correctly, resulting in the appropriate rows being added to the OrphanFile model.

The **OrphanProtector** does even less than the others. But if you look carefully, you will see that along with the `database_required()` there is also the `orphan_protector()` decorator. This will prevent the processor to perform the check on the OrphanFile model and deleting the unrelated files. In the start method, we record the number of orphan files in the OrphanFile model and we confirm that they are still there in the finish. Since the actual removal of the orphan files happens in the processor start method, we need to count the number of orphans before calling the super start.

The **LazyProcessor** is responsible to check that there are no rows left in the OrphanFile, meaning that the removal was successful.

And now let us put everything together and run it.

Listing 8.12: ProcessorList execution

```

db_conf = default_conf['sqlite']
db_conf['URL'] = 'sqlite:///memory:'

```

(continues on next page)

² A much better implementation of an Importer could be achieved using a subclass of the `Importer`. See, for example, the `ImporterExample` class and its `documentation`.

(continued from previous page)

```

plist = ProcessorList(name='Orphan test', description='dealing with orphan files',
↳database_conf=db_conf)
importer = FileImporter(input_folder=tmp_path)
remover = RowRemover(n_rows=n_delete)
protector = OrphanProtector()
lazy = LazyProcessor()
plist.extend([importer, remover, protector, lazy])
plist.execute()

assert importer.processor_exit_status == ProcessorExitStatus.Successful
assert remover.processor_exit_status == ProcessorExitStatus.Successful
assert protector.processor_exit_status == ProcessorExitStatus.Successful
assert lazy.processor_exit_status == ProcessorExitStatus.Successful

```

In practice, the only thing you have to take care of is to add a dedicated trigger to each of your tables having at least a file field and then the rest will be automatically performed by MAFw.

8.7.3 Keeping the entries updated

One aspect is that the file exists; another is that the file content remains unchanged. You may have replaced an input file with a newer one and the database will not know it. If your processors are only executed on items for which there is still no output generated, then this replaced file may go unnoticed causing issues to your analysis.

For this reason, we are strongly recommending to always add a checksum field for each file field in your table. Calculating a checksum is just a matter of a split second on modern CPU while the time for the debugging your analysis code is for sure longer.

The function `verify_checksum()` takes a Model as argument and will verify that all checksums are still valid. In other words, for each `FileNameField` (or `FileNameListField`) with a link to a checksum field in the table, the function will compare the actual digest with the stored one. If it is different, then the DB row will be removed.

Also this function is not automatically called by any processor super methods. It is ultimately the user's responsibility to decide whether to proceed, bearing in mind that working with long tables and large files may result in delays in processor execution.

The implementation is very similar to the previous one, just change the function name. Keep in mind that the `verify_checksum()` will implicitly check for the existence of files and warn you if some items are missing, so you can avoid the `remove_widow_db_rows()`, if you perform the checksum verification.

```

class MyProcessor(Processor):

    def start():
        super().start()
        verify_checksum(InputData)

```

8.8 Filters: let us do only what is needed to be done!

As mentioned already several times, the main role of the database is to support the processor execution in providing the input items and in storing the output products. Not everything can be efficiently stored in a database, for example large chunk of binary data are better saved to the disc, in this case, the database will know the path where these data are stored.

One advantage of the database is that you can apply selection rules and you do not have to process the whole dataset if you do not need it. To help you in this, MAFw is offering you a ready to use solution, the `Filter` class. This is an object that you can configure via the steering file allowing you to run the processors only over the items passing the criteria you set.

8.8.1 How to configure a filter

In a steering file, there is a table for each processor where you can set the configurable active and passive parameters. To this table, you can add a sub-table named Filter, containing other tables, one for each input Model. This is how it will look like:

```
[GlobalFilter]
new_only = true

[MyProcessor]
param1 = 15

[MyProcessor.Filter.InputTable1]
resolution_value = 25

[MyProcessor.Filter.InputTable2]
sample_name = 'sample_1'
```

In the example above, MyProcessor has two *.Filter.* tables, one for InputTable1 and one for InputTable2. When the steering file will be parsed, the processor constructor will automatically generate two filters: for InputTable1 it will put a condition that the resolution field must be 25 and for InputTable2, the sample_name column should be 'sample_1'. If MyProcessor is using other tables for generating the input items, you could add them in the same way.

There is also an additional table named **GlobalFilter** where you can specify conditions that will be applied by default to all processors being executed. If the same field is specified in both the GlobalFilter and the Processor/Model specific one, the one in the Processor/Model will overrule the global one.

For the global filter, we can specify a magic boolean keyword `new_only == true`³. This will allow to implement a different input sequence considering only items for which an output does not exist.

8.8.2 How to use a filter

Let us put in practice what we have seen so far. The filter native playground is in the implementation of the `get_items()` method.

Let us assume that our *Processor*, named AdvProcessor, is using three models to obtain the item lists. Everything is nicely described in the ERD below. The three models are interconnected via foreign key relations. There is a fourth model, that is where the output data will be saved.

The real core of our database is the image table, where our data are first introduced, the other two on the right, are kind of helper tables storing references to the samples and to the resolution of our images. The processed_image table is where the output of our AdvProcessor will be stored.

To realize this database with our ORM we need to code the corresponding model classes as follows:

```
1 class Sample(MAFwBaseModel):
2     sample_id = AutoField(primary_key=True, help_text='The sample id primary key')
3     sample_name = TextField(help_text='The sample name')
4
5 class Resolution(MAFwBaseModel):
6     resolution_id = AutoField(primary_key=True, help_text='The resolution id primary_
7     ↪key')
8     resolution_value = FloatField(help_text='The resolution in μm')
9
10 class Image(MAFwBaseModel):
11     image_id = AutoField(primary_key=True, help_text='The image id primary key')
12     filename = FileNameField(help_text='The filename of the image', checksum_field=
13     ↪'checksum')
```

(continues on next page)

³ Remember that in TOML, the two booleans are **NOT** capitalized as in Python. Moreover, you can specify `new_only` also in Processor/Model but it will not be taken into account unless that model has a column named `new_only`.

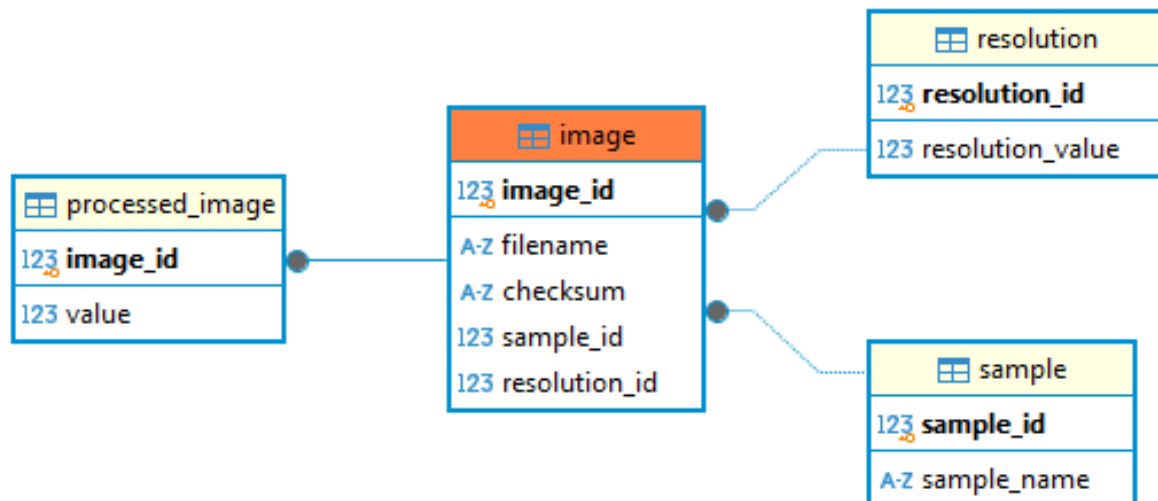


Fig. 8.1: The ERD of the AdvProcessor.

(continued from previous page)

```

12     checksum = FileChecksumField(help_text='The checksum of the input file')
13     sample = ForeignKeyField(
14         Sample, Sample.sample_id, on_delete='CASCADE', backref='sample', column_name=
↪ 'sample_id'
15     )
16     resolution = ForeignKeyField(
17         Resolution, Resolution.resolution_id, on_delete='CASCADE', backref='resolution
↪ ', column_name='resolution_id'
18     )
19
20 class ProcessedImage(MAFwBaseModel):
21     image = ForeignKeyField(
22         Image,
23         Image.image_id,
24         primary_key=True,
25         column_name='image_id',
26         backref='raw',
27         help_text='The image id, foreign key and primary',
28         on_delete='CASCADE',
29     )
30     value = FloatField(default=0)

```

By now, you should be an expert in ORM and everything there should be absolutely clear, otherwise, take your chance now to go back to the previous sections or to the `peewee` documentation to find an explanation. Note how the `Image` class is making use of our `FileNameField` and `FileChecksumField`. We added also a bit of help text to each field, in order to make even more evident what they are. For each foreign key field, we have specified a `backref` field, so that you could get access to the related models. Pay attention also at the highlighted lines, where we define foreign key fields to other tables. `Peewee` follows Django style references, so actually the field object is named with the noun of the object you are referring to. This will allow the following:

```

image.sample # Resolve the related object returning a Sample instance, it costs an
↪ additional query
image.sample_id # Return the corresponding Sample's id number

```

The primary source of input is the `Image`; however, you may wish to process only images that meet specific criteria, such as belonging to a particular sample or being captured at a certain resolution. Unfortunately, this information

is not explicitly included in the Image model. Only the `resolution_id` and the `sample_id` are included in the image table: those primary keys are excellent for a computer, but for a human being it is better to use sample names and resolution values. The solution is to use a `join query` in order to have all fields available and then we will be able to apply the configurable filters from the TOML steering file to limit the selection to what we want. We have three input models, so we may have up to three filters defined in the configuration file, along with the `GlobalFilter`. The processor, during its construction, will build the filters with the configuration information and store them in the processor `filter register`.

Let us have a look at a possible implementation of the `get_items()` for our advanced processor.

python

```

1 def get_items(self):
2
3     # first of all, let us be sure that the tables exist
4     # the order is irrelevant, the ORM will find the best creation strategy.
5     # if the table already exists, nothing will happen.
6     self.database.create_tables([Sample, Resolution, Image, ProcessedImage])
7
8     # if you want to remove widow rows from the output table or verify the checksum.
9     ↪do it now!
10    remove_widow_rows([Image, ProcessedImage])
11
12    # bind all input table filters
13    # this will establish a physical connection between the filter and
14    # the corresponding ORM model.
15    # the order of the model in the list is irrelevant.
16    # list all input models.
17    self.filter_register.bind_all([Sample, Resolution, Image])
18
19    # let us get a list of all already processed items.
20    # since the item in the output table are stored using the same primary key,
21    # this will simplify the discrimination between new and already processed items.
22    existing_entries = ProcessedImage.select(ProcessedImage.image_id).execute()
23
24    # did we select new_only in the global filter?
25    # if yes, prepare an additional condition in which we specify that the
26    # Image.image_id should not be among the image_id of the ProcessedImage.
27    # if no, then just accept everything.
28    if self.filter_register.new_only:
29        existing = ~Image.image_id.in_([i.image_id for i in existing_entries])
30    else:
31        existing = True
32
33    # finally let us make the query.
34    query = (Image.select(Image, Sample, Resolution)
35              .join(Sample, on=(Image.sample_id == Sample.sample_id), attr='s')
36              .switch(Image)
37              .join(Resolution, on=(Image.resolution_id == Resolution.resolution_
38    ↪id), attr='r')
39              .where(self.filter_register.filter_all())
40              .where(existing)
41              )
42    return query

```

TOML

```

1 [GlobalFilter]
2 new_only = true
3
4 [AdvProcessor.Filter.Sample]
5 sample_name = 'Sample_0000[123]'
6
7 [AdvProcessor.Filter.Resolution]
8 resolution_value = 50
9
10 [AdvProcessor.Filter.Image]
11 filename = '\*file_0[12345]0\*'

```

The comments in the code make it rather self explanatory, but let us have a look together at the key points.

- At line 16, we are binding the filters in the filter register with the model classes. The filter register contains all the filters that were found in the configuration file, and since their presence in the steering file is optional, you may have decided not to include all model specific filters. In such a case, the `FilterRegister()` will create an empty filter for the missing models using the global filter as default values.
- At line 21, we query the output table to get the list of already existing output. Since the input and output images are sharing the same primary key, this will make the identification of new items very easy. If the DB structure was not built in this way, you will have to work out a solution to find which are the new items.
- The list of items is finally retrieved in the lines 33 to 39. We select not only from Image, but we join also Sample and Resolution. In the join statement we specify the target Model and the relation to align the two models. Theoretically, if the two models are linked by a foreign key relation, then the `on` attribute could be skipped, but we like when everything is clearly stated. We will come back to the `attr` in a second.
- Lines 37 and 38 is where the filter gets into action, exactly in the `where` clause. First we ask at the `filter_register` to generate the filter statement for each of the filters it is holding and then we select according to the existence criteria.

In this specific situation (*just click on the TOML configuration tab to see the filter configuration*), the `filter_all` condition will be equivalent of writing:

```
(Sample.sample_name % 'Sample_0000[123]') &
(Resolution.resolution_value == 50) &
(Image.filename % '\*file_0[12345]0\')
```

where the `%` operator⁴ is the SQL GLOB.

In the last line, we return the query, that is an iterable object, so actually our looping processor will assign to the `item` attribute the value of one selected and filtered row after the other. The item will then be a modified version of the Image class. You will be able to access the image table columns directly using the dot notation, and, at the same time, you can access the two joined tables via the two `attr` fields, that you have specified in the join statement.

So for example, this will be a valid syntax:

```
print(self.item.filename.name) # print the file name (without the path) of the image
print(self.item.s.sample_name) # print the sample name, look at the .s in between
print(self.item.r.resolution_value) # print the value of the resolution. look at the .
↪r in between
```

Avoiding the N+1 problem

Since the Image table has foreign key dependency from Sample and Resolution, you could theoretically access the referenced row via the `backref` attribute, but this will trigger an additional query every time. The impact

⁴ The list of all overloaded operator and special method is available in the `peewee` doc.

of those additional queries can be excessive especially if the tables are large. This is the so-called N+1 problem in ORM. On the contrary, if you use the s and r attributes, no additional queries will be performed.

8.9 Multi-primary key columns

Special attention is required when you need to have a primary key that is spanning over two or more columns of your model. So far we have seen how we can identify one column in the model as the primary key and now we will see what to do if you want to use more than one column as primary key and, even more important, how you can use this composite primary key as a foreign key in another model.

To describe this topic, we will make use of an example that you can also find in the examples modules of MAFw named *multi_primary*.

Let us start with the model definition.

```

1 class Sample(MAFwBaseModel):
2     sample_id = AutoField(primary_key=True, help_text='The sample id primary key')
3     sample_name = TextField(help_text='The sample name')
4
5 class Resolution(MAFwBaseModel):
6     resolution_id = AutoField(primary_key=True, help_text='The resolution id primary_
↳key')
7     resolution_value = FloatField(help_text='The resolution in µm')
8
9 class Image(MAFwBaseModel):
10    image_id = AutoField(primary_key=True, help_text='The image id primary key')
11    filename = FileNameField(help_text='The filename of the image', checksum_field=
↳'checksum')
12    checksum = FileChecksumField()
13    sample = ForeignKeyField(Sample, on_delete='CASCADE', backref='+', lazy_
↳load=False, column_name='sample_id')
14    resolution = ForeignKeyField(
15        Resolution, on_delete='CASCADE', backref='+', lazy_load=False, column_name=
↳'resolution_id'
16    )
17
18 class ProcessedImage(MAFwBaseModel):
19    image = ForeignKeyField(
20        Image,
21        primary_key=True,
22        backref='+',
23        help_text='The image id, foreign key and primary',
24        on_delete='CASCADE',
25        lazy_load=False,
26        column_name='image_id',
27    )
28    value = FloatField(default=0)
29
30 class CalibrationMethod(MAFwBaseModel):
31    method_id = AutoField(primary_key=True, help_text='The primary key for the_
↳calculation method')
32    multiplier = FloatField(default=1.0, help_text='The multiplication factor of this_
↳method')
33
34 class CalibratedImage(MAFwBaseModel):
35    image = ForeignKeyField(
36        ProcessedImage,

```

(continues on next page)

(continued from previous page)

```

37     on_delete='CASCADE',
38     help_text='The reference to the processed image',
39     backref='+',
40     lazy_load=False,
41     column_name='image_id',
42 )
43 method = ForeignKeyField(
44     CalibrationMethod,
45     on_delete='CASCADE',
46     help_text='The reference to the calibration method',
47     backref='+',
48     lazy_load=False,
49     column_name='method_id',
50 )
51 calibrated_value = FloatField(default=0.0, help_text='The calibrated value')
52
53 @property
54 def primary_key(self) -> Iterable:
55     return self.image_id, self.method_id
56
57 class Meta:
58     primary_key = CompositeKey('image', 'method')
59
60 class ColoredImage(MAFwBaseModel):
61     image_id = IntegerField(help_text='The reference to the processed image. Combined_
62     ↪FK with method_id')
63     method_id = IntegerField(help_text='The reference to the calibration method.
64     ↪Combined FK with method_id')
65     red = FloatField(default=0, help_text='Fake red. Only for testing')
66     green = FloatField(default=0, help_text='Fake green. Only for testing')
67     blue = FloatField(default=0, help_text='Fake blue. Only for testing')
68
69 @property
70 def primary_key(self) -> Iterable:
71     return self.image_id, self.method_id
72
73 class Meta:
74     constraints = [
75         SQL(
76             'FOREIGN KEY (image_id, method_id) REFERENCES '
77             'calibrated_image(image_id, method_id) ON DELETE CASCADE'
78         )
79     ]
80     primary_key = CompositeKey('image_id', 'method_id')

```

As always, one single picture can convey more than a thousand lines of code. Here below the ERDs of Image and of CalibratedImage.

In the diagrams, the fields with bold font represent primary keys, also highlighted by the separation line, while the arrow are the standard relation.

As in the examples above, we have images of different samples acquired with different resolutions entering the Image model. We use those lines to make some calculations and we obtain the rows in the ProcessedImage model. These two tables are in 1 to 1 relation and this relation is enforced with a delete cascade, meaning that if we delete an element in the Image model, the corresponding one in the ProcessedImage will also be deleted.

The CalibrationMethod model contains different sets of calibration constants to bring each row from the Processed-

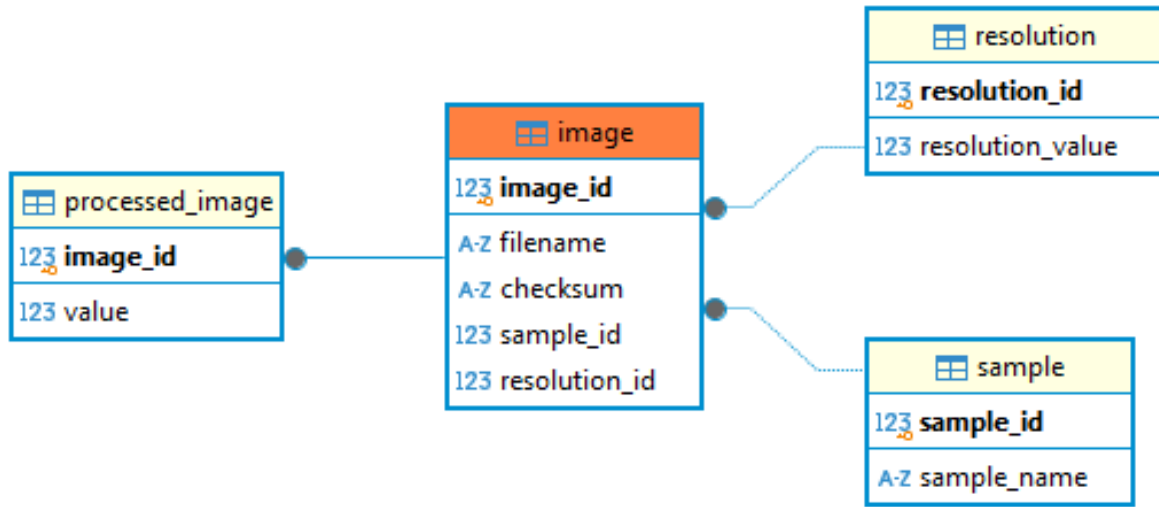


Fig. 8.2: The ERD of the Image Model

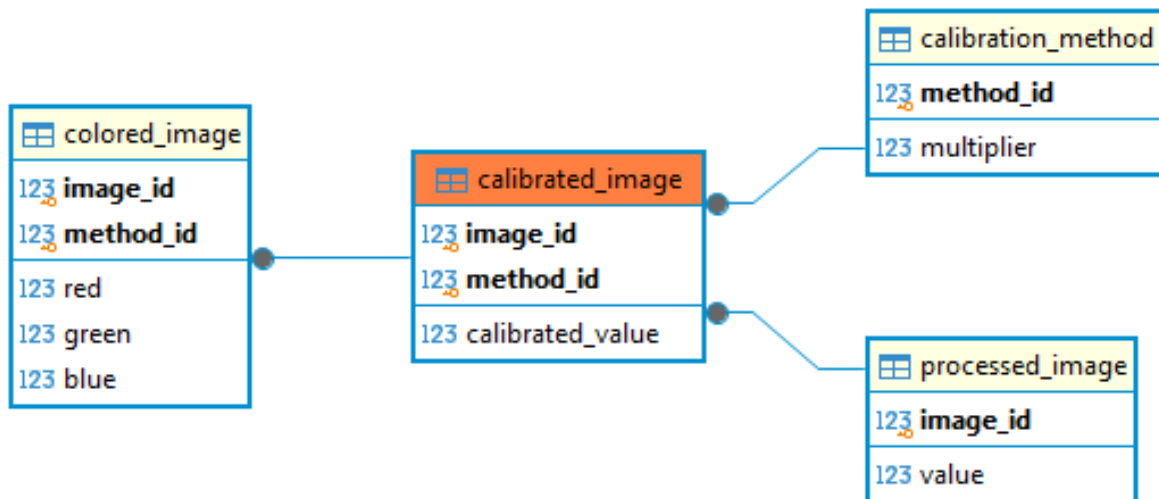


Fig. 8.3: The ERD of the CalibratedImage Model

Image model to the CalibratedImage one. It is natural to assume that the `image_id` and the `method_id` are the best candidates to be a combined primary key. To achieve this, in the CalibratedImage model, we need to add (line 57-58) an overload of the Meta class, where we specify our CompositeKey. Pay attention to an important detail: the CompositeKey constructor takes the name of the fields and not the name of the columns, that in the case of foreign keys differ of `'_id'`. Optionally we can also define a `primary_key` property (line 53-55) to quickly retrieve the values of our keys.

From the application point of view, we want all the processed images to be calibrated with all possible calibration methods, that means we need to make a cross join as described below:

```

1 with database.atomic() as txn:
2     cross_join = (
3         ProcessedImage.select(ProcessedImage, CalibrationMethod).
4         ↪join(CalibrationMethod, JOIN.CROSS).execute()
5     )
6     for row in cross_join:
7         calibrated_image = CalibratedImage()
8         calibrated_image.image_id = row.image_id
9         calibrated_image.method_id = row.method_id
10        calibrated_image.calibrated_value = row.value * row.multiplier
        calibrated_image.save(force_insert=True)

```

Up to this point we have seen what we have to do to specify a composite primary key, we cannot use the AutoField or the `primary_key` parameter, but we need to go through the Meta class in the way shown in the example.

The next step is to have another table (ColoredImage in our imaginary case) that is in relation with CalibratedImage. We would need to have again a composite primary key that is also a composite foreign key. Peewee does not support composite foreign keys, but we can use the workaround shown at lines 72-77. Along with the CompositeKey definition, we need to add a Constraint as well using the SQL function to convert a string into a valid SQL statement. This time, since we are using low level SQL directives, we have to use the column names (additional `'_id'`) instead of the field name.

And in a similar way we can insert items in the ColoredImage model.

```

1 with database.atomic() as txn:
2     full_list = CalibratedImage.select().execute()
3
4     for row in full_list:
5         colored_image = ColoredImage(image_id=row.image_id, method_id=row.method_id)
6         colored_image.red = row.calibrated_value / 3
7         colored_image.green = row.calibrated_value / 3
8         colored_image.blue = row.calibrated_value / 3
9         colored_image.save(force_insert=True)

```

Now, with all the tables linked to each other, try to delete one from a table, and guess what will happen to all other tables.

This tutorial might be a bit more complex than the examples we have seen so far, but we believed you have appreciated the power of such a relational tool.

8.10 Importing an existing DB

The last section of this long chapter on database will show you how to deal with an existing DB. It is possible that before you have adopted MAFw for your analysis tasks, you were already employing a relational database to store your dataset and results. So far we have seen how to create tables in a database starting from an object oriented description (a model) in a python library. But what do we have to do if the database already exists? Can we create the classes starting from a database? This process goes under the name of **reflection** and it is the subject of this section.

The reflection of tables in python classes cannot be performed automatically at 100% by definition. A typical case

is the use of application specific fields. Consider, for example, the `FileNameField` that we have discussed earlier. This field corresponds to a `Path` object when you look at it from the application point of view, but the actual path is saved as a text field in the concrete database implementation. If you now read the metadata of this table from the database point of view, you will see that the field will contain a text variable and thus the reflected class will not have any `FileNameField`.

Let us try to understand the process looking at the picture below. If we create the model in python, then we can assign special field descriptors to the table columns, but their concrete implementation in the database must be done using types that are available in the database itself. So when we perform the reverse process, we get only a good approximation of the initial definition.

```
class Image(MAFwBaseModel): 14 usages (2 dynamic)  ▸ BULGHERONI Antonio
    image_id = AutoField(primary_key=True, help_text='The image id primary key')
    filename = FileNameField(help_text='The filename of the image', checksum_field='checksum')
    checksum = FileChecksumField()
    sample = ForeignKeyField(Sample, on_delete='CASCADE', backref='+', lazy_load=False, column_name='sample_id')
    resolution_id = ForeignKeyField(
        Resolution, on_delete='CASCADE', backref='+', lazy_load=False, column_name='resolution_id'
    )
)
```

Fig. 8.4: This is the model implementation as you would code it making use of the specific field definitions.

Column Name	#	Data Type
I23 image_id	1	INTEGER
A2 filename	2	TEXT
A2 checksum	3	TEXT
I23 sample_id	4	INTEGER
I23 resolution_id	5	INTEGER

Fig. 8.5: During the actual implementation of the model as a database table, python column definitions will be translated into database types.

```
class Image(MAFwBaseModel): 3 usages (2 dynamic)
    image_id = AutoField()
    filename = TextField()
    checksum = TextField()
    sample = ForeignKeyField(column_name='sample_id', field='sample_id', model=Sample)
    resolution = ForeignKeyField(column_name='resolution_id', field='resolution_id', model=Resolution)

    class Meta:
        table_name = 'image'
```

Fig. 8.6: The reflection process will translate the database implementation in a generic model implementation, not necessarily including all the specific field definition.

Nevertheless the process is rather efficient and can generate an excellent starting point that we can use to customize the model classes to make them more useful in our application.

From a practical point of view, you just have to open a console and type the command `mafwb wizard --help` to get some help on the tool and also read its [documentation](#). You need to provide the name of the database and how to connect to it, in the case of `Sqlite DB`, it is enough to provide the filename, and you have to specify the name of the output python file that will contain all the model classes. This module is ready to go, you could theoretically import it into your project and use it, but it is strongly recommended to accurately check that everything is really the way you want it to be.

The reflection process is absolutely safe for your existing database, so it is worth to give it a try!

8.11 What's next

Congratulations! You reached the end of the most difficult chapter in this tutorial. It is difficult because as a scientist you might not be used to deal with databases everyday, but their power is incredible, isn't it?

The next chapter is about the library of advanced processors that MAFw is sharing to simplify your job even further. In particular, we are sure you will like a lot our plotting processors!

Scientists love to present their results graphically, because like a picture is worth a thousand words, a good plot can replace a lot of tables. The python scene is full of libraries for data visualization, from the most basic like `matplotlib` offering a low level interface to higher level tools like `seaborn` where the user can profit from predefined classes to display complex relationships among different variables in the dataset.

In a spirit of complete freedom, MAFw is not deciding for you what plotting or data frame library you should use for your analysis. For this reason, we have prepared a skeleton of the *generic plotter processor* that you can use as a template and implement it with your favorite libraries.

Nevertheless, if you are ok with the use of `seaborn` as a data visualization platform (using `matplotlib` under the hood for the low level plotting routines) and `pandas` and `numpy` for data handling, then MAFw can offer you a pre-cooked solution implemented in the *SNSPlotter*.

If you want to use *SNSPlotter*, then you need to install mafw including the *seaborn* optional feature, in this way:

```
c:\> pip install mafw[seaborn]
```

In the next section of this unit, we will assume that you are using the `seaborn` functionalities.

9.1 Add diagnostic information to your processors

Processors are responsible for performing simple analysis tasks; however, it is essential to verify and assess the accuracy of their output. We strongly recommend to include the generation of some diagnostic plots or storing some key results in a database table while implementing the `process()` method in order to verify the correct execution of the analysis task.

This aspect is totally in your hand and MAFw has little to offer; you will have to code your *Processor* so that it can create meaningful figures.

9.2 Generate high level graphs

Along with graphs aimed to verify that your processor was actually behaving as expected, you might have the need to have processors whose main task is indeed the generation of high level figures to be used for your publications or presentations starting from a data frame.

In this case, MAFw is providing you a powerful tool: the *generic plotter processor*.

The *SNSPlotter* is a member of the `processor_library` package and represents a specialized subclass of a *Processor*¹ with a predefined `process()` method.

¹ Actually *SNSPlotter* parent class is *GenericPlotter* that is a subclass of *Processor*. This intermediate inheritance step is needed to allow the user to define their preferred interface to graphical and data frame library.

When you are subclassing a *Processor*, you are normally required to implement some methods and the most important of all is clearly the *process()* one, where your analytical task is actually coded.

For a *SNSPlotter*, the *process method* is already implemented and it looks like the following code snippet:

```
def process(self) -> None:
    """
    Process method overload.

    In the case of a plotter subclass, the process method is already implemented and
    ↪the user should not overload
    ↪it. On the contrary, the user must overload the other implementation methods.
    ↪described in the general
    :class:`class description <.SNSPlotter>`.
    """
    if self.filter_register.new_only:
        if self.is_output_existing():
            return

    self.in_loop_customization()
    self.get_data_frame()
    self.patch_data_frame()
    self.slice_data_frame()
    self.group_and_aggregate_data_frame()
    if not self.is_data_frame_empty():
        self.plot()
        self.customize_plot()
        self.save()
        self.update_db()
```

In its basic form, all the methods included in the process have an empty implementation, and it is your task to make them doing something.

We can divide these methods into two blocks:

- **The data handling methods: where the data frame is acquired and prepared for the plotting step.**
 1. *get_data_frame()*
 2. *patch_data_frame()*
 3. *slice_data_frame()*
 4. *group_and_aggregate_data_frame()*
- **The plotting methods: where you will have to generate the figures from the data frame and save them.**
 1. *plot()*
 2. *customize_plot()*
 3. *save()*
 4. *update_db()*

9.2.1 Getting the data

The first thing to do when you want to plot some data is to have the data! The *SNSPlotter* has a *data_frame* attribute to store a pandas data frame exactly for this purpose. The data frame can be read in from a file or retrieved from a database table or even built directly in the processor itself. The *get_data_frame()* is the place where this is happening. It is important that you assign the data frame to *self.data_frame* attribute so that the following methods can operate on it.

The following three methods are giving you the possibility to add new columns based on the existing ones (*patch_data_frame()*), selecting a sub set of rows (*slice_data_frame()*) and grouping / aggregating the rows (*group_and_aggregate_data_frame()*).

Data frame patching

Use this method to add or modify data frame columns. One typical example is the application of unit of measurement conversion: you may want to store the data using the SI unit, but it is more convenient to display them using another more practical unit. Another interesting application is to replace numerical categorical values with textual one: for example you may want to add a column with the sample name to replace the column with the sample identification number. You can implement all these operations overloading the *patch_data_frame()* method. Below is an example of how to add an additional column that represents the exposure time of an image acquisition in hours, converted from the stored value in seconds.

```
self.data_frame['ExposureTimeH'] = self.data_frame['ExposureTime'] / 3600.
```

Slicing the data frame

Use this method to select a specific group of rows, corresponding to a certain value for a variable. You do not need to overload the *slice_data_frame()*, but just define the *slicing_dict* attribute and the basic implementation will do the magic. Have a look also at the documentation of the function *slice_data_frame()* from the *pandas_tools* module. In a few words, if you have a data frame with four columns named A, B, C and D and you want to select only the rows where column A is equal to 10, then set `self.slicing_dict = { 'A' : 10 }` and the *slice_data_frame()* method will do the job for you.

Group and aggregate the data frame

In a similar manner, you can group and aggregate the columns of the data frame by defining the *grouping_columns* and *aggregation_functions* attributes. Also for this method, have a look at the *group_and_aggregate_data_frame()* function from the *pandas_tools*. If, for example, you have a data frame containing two columns named SampleName and Signal and you want to get the average and standard deviation of the signal column for each of the different samples, then you can do:

```
self.grouping_columns = ['SampleName']
self.aggregation_functions = ['mean', 'std']
```

The *group_and_aggregate_data_frame()* will perform the grouping and aggregation and assign to the *self.data_frame* a modified data frame with three columns named SampleName, Signal_mean and Signal_std.

Per cycle customization

You may wonder, when and where in your code you should be setting the slicing and grouping / aggregating parameters. Those values can be ‘statically’ provided directly in the constructor of the generic plotter processor subclass, but you can also have them dynamically assigned at run time overloading the *in_loop_customization()* method.

In the implementation of this method you can also set other parameters relevant for your plotting.

9.2.2 Plotting the data

Now that you have retrieved and prepared your data frame, you can start the real plotting. First of all, just be sure that after all your modifications, the data frame is still containing some rows, otherwise, just move on to the next iteration.

The generation of the real plot is done in the implementation of the *plot()*. The *SNSPlotter* is meant to operate with *seaborn* so it is quite reasonable to anticipate that the results of the plotting will be assigned to a *facet grid attribute*, in order to allow the following methods to refer to the plotting output.

If your *plot()* implementation is not actually generating a facet grid, but something else, it is still ok. If you want to pass a static typing, then define a attribute in your *SNSPlotter* subclass with the proper typing; if you are not doing any static typing you can still use the facet grid attribute for storing your object.

Customising the plot

When producing high quality graphs, you want to have everything looking perfect. You can implement this method to set the figure and plot titles, the axis labels and to add specific annotations. If your *SNSPlotter* implementation is following a looping execution workflow, you can use the looping parameters if you need them. Theoretically speaking you could code this customisation directly in the plot method, but the reason for this split will become clear later on when explaining the *mixin approach*.

Save the plot

This is where your figure must be saved to a file. The very typical implementation of this method is the following:

```
def save(self) -> None:
    self.facet_grid.figure.savefig(self.output_filename)
    self.output_filename_list.append(self.output_filename)
```

This is assuming that your *SNSPlotter* subclass has an attribute (it can also be a processor parameter) named `output_filename`. The second statement is also very common: the *SNSPlotter* has a list type attribute to store the names of all generated files. This list, along with the cumulative checksum will be stored in the dedicated *PlotterOutput* standard table, to avoid the regeneration of identical plots. You do not have to worry about updating this standard table, this operation will be performed automatically by the private method `_update_plotter_db()`.

Update the db

Along with the *PlotterOutput* standard table, you may want to update also other objects in the database. For this purpose, you can implement the `update_db()` method.

9.2.3 Mixin approach

You might have noticed that the structure of the `process()` is a bit *over-granular*. The reason behind this fine distinction is due to the fact that the code can potentially be reused to a large extent thanks to a clever use of mixin classes.

Explanation

If you are not a seasoned programmer, you might not be familiar with the concept of mixins. Those are small classes implementing a specific subtask of a bigger class, finding extensive usage in developing interfaces. If you wish, you can have a general explanation following this [webpage](#).

If you are aware of what a mixin is and how to use it, then just go to the *next section*.

In the case of the *SNSPlotter*, we have seen that its methods can be divided into two categories: the data retriever and the actual plotter. You can have several different data retriever options and also different plotting strategies, thus making the creation of a matrix of subclasses of the *SNSPlotter* covering all possible cases really unpractical.

The mixin approach allows to have a flexible subclass scheme and at the same time an optimal reuse of code units.

To make things even more evident, here is a logical scheme of what happens when you mix a *DataRetriever* mixin with a *SNSPlotter*.

Let us start from the definition of the subclass itself:

```
class MyPlotterSNS(MyDataRetrieverMixin, MyPlotMixin, SNSPlotter):
```

`MyPlotterSNS` is mixing two mixins, one *DataRetriever* and one *FigurePlotter*, with the *SNSPlotter*. The order **matters**: always put mixins before the base class, because this is affecting the method resolution order (MRO) that follows a depth-first left-to-right traversal.

`MyDataRetrieverMixin` being a subclass of *DataRetriever* is implementing these two methods `get_data_frame()` and `patch_data_frame()` that are also defined in the *SNSPlotter*.

During the execution of the `process()`, the `get_data_frame` method will be called and, thanks to the MRO, the `MyDataRetrieverMixin` will provide its implementation of such method. The same will apply for the `plot`, where the `MyPlotMixin` will kick in and provide its implementation of the method.

Mixins can come with other class parameters along with the ones shared with the base class. Those can be accessed and assigned using the standard dot notation on the instances of the mixed class, or they can also be defined in the mixed class constructor directly.

Data-retriever mixins

Let us start considering the data retrieval part. Very likely, you will get those data either from a database table or from a HDF5 file. If this is the case, then you do not need to code the `get_data_frame` method everytime, but you can use one of the available mixin to grab the data.

For the moment, the following data retriever mixins exist²:

- `HDFPdDataRetriever`.

This data retriever is opening the provided HDF5 file name and obtaining the data frame identified by the provided key.

- `SQLPdDataRetriever`.

This data retriever is performing a SQL query on a table (`table_name`) using as columns the provided `required_columns` list and fulfilling the `where_clause`. A valid database connection must exist for this retriever to work, it means that the processor with which this class is mixed must have a working database instance.

- `FromDatasetDataRetriever`.

This data retriever is provided mainly for test and debug purposes. It allows to obtain a data frame from the standard data set library of `seaborn`.

If you plan to retrieve data several times from a common source, you may consider to add your own mixin class. Just make it inheriting from the `DataRetriever` (or `PdDataRetriever` if you are using Pandas) class and follow one of the provided case as an example.

Here below is an example of a Plotter retrieving the data from a HDF file:

```
@single_loop
class HDFPlotterSNS(HDFDataRetriever, SNSPlotter):
    pass

hdf = HDFPlotterSNS(hdf_filename=your_input_file, key=your_key)
hdf.execute()
```

The `hdf_filename` and the `key` can be provided directly in the class constructor, but this might not doable if you are planning to execute the processor using the `maf w steering file` approach.

In this case, you can add two parameters to your processor that can be read from the steering file and follow this approach:

```
@single_loop
class HDFPlotterSNS(HDFDataRetriever, SNSPlotter):
    hdf_filename_from_steering = ActiveParameter('hdf_filename_from_steering',
↪ default='my_datafile.h5')
    hdf_key_from_steering = ActiveParameter('hdf_key_from_steering', default='my_key')

    def start(self):
        super().start()
        # you can assign the parameter values to the mixin attributes in the start.
```

(continues on next page)

² The extra `Pd` in the name of these classes stands for Pandas. In fact the user is free to select any other data frame library and define their own data retriever mixin.

(continued from previous page)

```

↪method
    # but also in the in_loop_customization method if this is better suiting your
↪needs
    self.hdf_filename = self.hdf_filename_from_steering
    self.key = self.hdf_key_from_steering

# let us execute the processor from within the mafw runner.

```

Figure plot mixins

In a similar way, if what you want is to generate a `seaborn` figure level plot, then you do not need to define the plot method every time. Just mix the `SNSPlotter` with one of the three figure plotter mixins available:

- `RelPlot`: to plot relational plots.
- `DisPlot`: to plot distribution graphs.
- `CatPlot`: to produce categorical figures.

From the coding point of it, the approach is identical to the one shown before for the data retrieving mixin.

```

@single_loop
@processor_depends_on_optional(module_name='pandas;seaborn')
class DataSetRelPlotPlotterSNS(FromDatasetDataRetriever, RelPlot, SNSPlotter):
    def __init__(self, output_png: str | Path | None = None, *args: Any, **kwargs:
↪Any) -> None:
        super().__init__(*args, **kwargs)
        self.output_png = Path(output_png) if output_png is not None else Path()

    def save(self) -> None:
        self.facet_grid.figure.savefig(self.output_png)
        self.output_filename_list.append(self.output_png)

output_file = tmp_path / Path('relplot_from_dataset.png')

dsrp = DataSetRelPlotPlotterSNS(
    output_file, dataset_name='penguins', x='flipper_length_mm', y='bill_length_mm',
↪col='sex', hue='species'
)
dsrp.execute()

```

The code above will generate this plot:

As you can see, the plot method is taken from the `RelPlot` mixin and the output can be further customized (axis labels and similar things) implementing the `customize_plot` method.

A large fraction of the parameters that can be passed to the `seaborn` figure level functions corresponds to attributes in the mixin class, but not all of them. For the missing ones, you can use the dictionary like parameter `plot_kws`, to have them passed to the underlying plotting function. See the documentation of the mixins for more details (`RelPlot`, `DisPlot` and `CatPlot`).

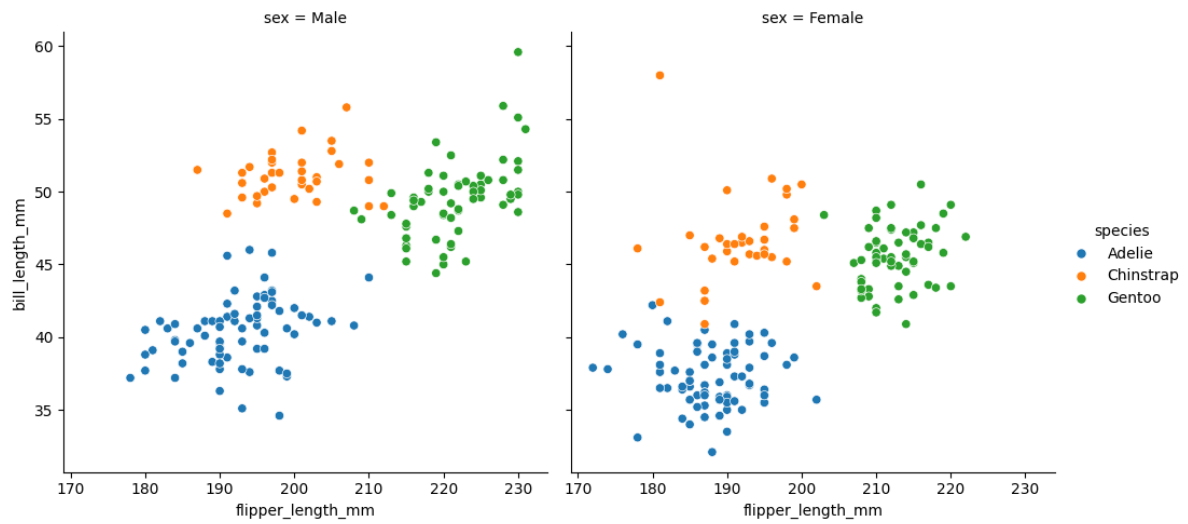


Fig. 9.1: This is the exemplary output of a SNSPlotter mixed with a data retriever and a relational plot mixin.

CHAPTER 10

A STEP BY STEP TUTORIAL FOR A REAL EXPERIMENTAL CASE STUDY

Congratulations! You have decided to give a try to MAFw for your next analysis, you have read all the documentation provided, but you got overwhelmed by information and now you do not know from where to start.

Do not worry, we have a simple but nevertheless complete example that will guide you through each and every steps of your analysis.

Before start typing code, let's introduce how experimental scenarion.

10.1 The experiment scenario

Let's imagine, that we have a measurement set that is integrating the amount of UV radiation that is reaching a sensor in a give amount of time.

The experimental data acquisition system is saving one file for each exposure containing the value read by the sensor. The DAQ is encoding the duration of the acquisition in the file name and let's assume we acquired 25 different exposures, starting from 0 up to 24 hours. The experimental procedure is repeated for three different detectors having different type of response and the final goal is to compare their performance.

It is an ultra simplified experimental case and you can easily make it as complex as you wish, just by adding other variables (different UV lamps, detector operating conditions, background illumination...). Nevertheless this simple scenario can be straightforward expanded to any real experimental campaign.

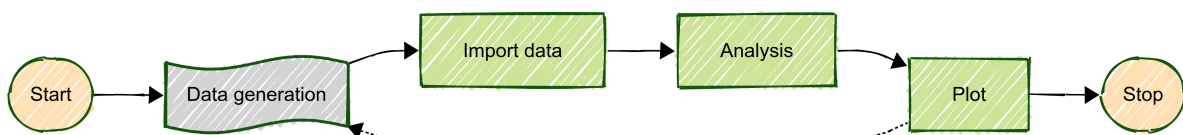


Fig. 10.1: The simplified pipeline for the tutorial example.

10.1.1 Task 0. Generating the data

This is not really an analysis task, rather the real acquisition of the data, that's why it is represented with a different color in schema above. Nevertheless it is important to include that in our planning, because new data might be generated during or after some data have been already analyzed. In this case, it is important that our analysis pipelines will be able to process only the new (or modified) data, without wasting time and resources re-doing what has been done already. This is the reason why there is a dashed line looping back from the plot step to the data generation. Since this is a simulated experiment, instead of collecting real data, we will use a *Processor* to generate some data.

10.1.2 Task 1. Building your data bank

From a conceptual point of view, the first thing you should do when using MAFw is to import all your data (in this case the raw data files) into a relational database. You do not need to store the content of the file in the database, otherwise it will soon explode in size, you can simply insert the full path from where the file can be retrieved and its checksum, so that we can keep an eye on possible modifications.

10.1.3 Task 2. Do the analysis

In the specific scenario the analysis is very easy. Each file contains only one number, so there is very little to be done, but your case can be as complicated as needed. We will open each file, read the content and then put in a second database table containing the results of the analysis of each file.

10.1.4 Task 3. Prepare a relation plot

Using the data stored in the data, we can generate a plot representing the integral flux versus the exposure time for the three different detectors using a relation plot.

10.2 The ‘code’

In the previous section, we have defined what we want to achieve with our analysis (it is always a good idea to have a plan before start coding!) Now we are ready to start with setting up the project containing the required processors to achieve the analytical goal described above.

If you want to use MAFw plugin mechanism, then you need to build your project in as a proper python package. Let's start then with the project specification contained in the `pyproject.toml` file.

```

1 [build-system]
2 requires = ["hatchling"]
3 build-backend = "hatchling.build"
4
5 [project]
6 name = "plug_test"
7 dynamic = ["version"]
8 description = 'A MAFw processor library plugin'
9 readme = "README.md"
10 requires-python = ">=3.11"
11 license = "EUPL-1.2"
12 authors = [
13     { name = "Antonio Bulgheroni", email = "antonio.bulgheroni@ec.europa.eu" },
14 ]
15 classifiers = [
16     "Development Status :: 4 - Beta",
17     "Programming Language :: Python",
18     "Programming Language :: Python :: 3.11",
19     "Programming Language :: Python :: 3.12",
20     "Programming Language :: Python :: 3.13",
21     "Programming Language :: Python :: Implementation :: CPython",
22     "Programming Language :: Python :: Implementation :: PyPy",
23 ]
24
25 dependencies = ['mafw', 'pluggy']
26
27 [project.urls]
28 Documentation = "https://github.com/..."
29 Issues = "https://github.com/..."
30 Source = "https://github.com/."
31

```

(continues on next page)

(continued from previous page)

```

32 [project.entry-points.mafw]
33 plug_test = 'plug.plugins'
34
35 [tool.hatch.version]
36 path = "src/plug/___about__.py"
37
38 [tool.hatch.build.targets.wheel]
39 packages = ['src/plug']

```

The specificity of this configuration file is the highlighted lines where we define an entry point where your processors are made available to MAFw.

Now before start creating the rest of the project, prepare the directory structure according to the python packaging specification. Here below is the expected directory structure.

```

plug
├── README.md
├── pyproject.toml
├── src
│   └── plug
│       ├── ___about__.py
│       ├── db_model.py
│       ├── plug_processor.py
│       └── plugins.py

```

You can divide the code of your analysis in as many python modules as you wish, but for this simple project we will keep all processors in one single module *plug_processor.py*. We will use a second module for the definition of database model classes (*db_model.py*). Additionally we will have to include a *plugins.py* module (this is the one declared in the *pyproject.toml* entry point definition) where we will list the processors to be exported along with our standard table.

10.2.1 The database definition

Let's start from the database definition.

Our database will contain tables corresponding to the three model classes: *InputFile* and *Data*, one helper table for the detectors along with all the *standard tables* that are automatically created by MAFw.

Before analysing the code let's visualize the database with the ERD.

The *InputFile* is the model where we will be storing all the data files that are generated in our experiment while the *Data* model is where we will be storing the results of the analysis processor, in our specific case, the value contained in the input file.

The rows of these two models are linked by a 1-to-1 relation defined by the primary key.

Remember that is always a good idea to add a checksum field every time you have a filename field, so that we can check if the file has changed or not.

The *InputFile* model is also linked with the *Detector* model to be sure that only known detectors are added to the analysis.

Let's have a look at the way we have defined the three models.

```

1 class Detector(StandardTable):
2     detector_id = AutoField(primary_key=True, help_text='Primary key for the detector_
   ↪table')
3     name = TextField(help_text='The name of the detector')
4     description = TextField(help_text='A longer description for the detector')
5
6     @classmethod

```

(continues on next page)

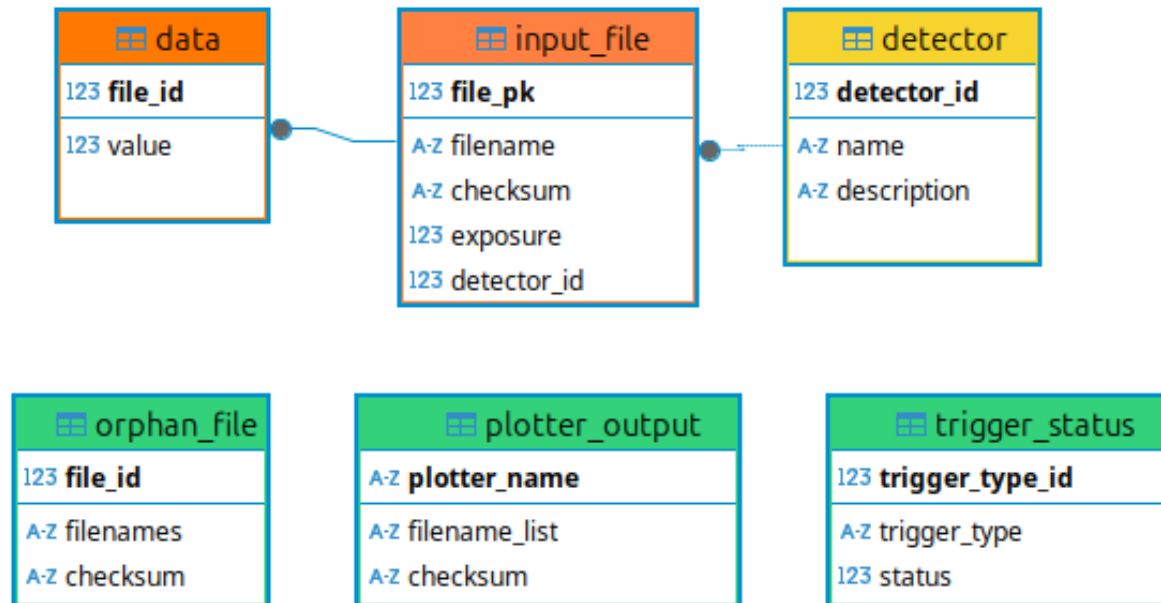


Fig. 10.2: The schematic representation of our database. The standard tables, automatically created are in green. The detector table (yellow) is exported as a standard table and its content is automatically restored every time mafw is executed.

(continued from previous page)

```

7  def init(cls) -> None:
8      data = [
9          dict(detector_id=1, name='Normal', description='Standard detector'),
10         dict(detector_id=2, name='HighGain', description='High gain detector'),
11         dict(detector_id=3, name='NoDark', description='Low dark current detector
↪ '),
12     ]
13
14     cls.insert_many(data).on_conflict(
15         conflict_target=[cls.detector_id],
16         update={'name': SQL('EXCLUDED.name'), 'description': SQL('EXCLUDED.
↪ description')},
17     ).execute()

```

The detector table is derived from the StandardTable, because we want the possibility to initialize the content of this table every time the application is executed. This is obtained in the init method. The use of the on_conflict clause assure that the three detectors are for sure present in the table with the value given in the data object. This means that if the user manually changes the name of one these three detectors, the next time the application is executed, the original name will be restored.

```

1  class InputFile(MAFwBaseModel):
2      @classmethod
3      def triggers(cls) -> list[Trigger]:
4          update_file_trigger = Trigger(
5              trigger_name='input_file_after_update',
6              trigger_type=(TriggerWhen.After, TriggerAction.Update),
7              source_table=cls,
8              safe=True,
9              for_each_row=True,
10         )
11         update_file_trigger.add_when(or_('NEW.exposure != OLD.exposure', 'NEW.

```

(continues on next page)

(continued from previous page)

```

12 ↪checksum != OLD.checksum'))
13     update_file_trigger.add_sql('DELETE FROM data WHERE file_pk = OLD.file_pk;')
14
15     return [update_file_trigger]
16
17     file_pk = AutoField(primary_key=True, help_text='Primary key for the input file.
18 ↪table')
19     filename = FileNameField(unique=True, checksum_field='checksum', help_text='The
20 ↪filename of the element')
21     checksum = FileChecksumField(help_text='The checksum of the element file')
22     exposure = FloatField(help_text='The duration of the exposure in h')
23     detector = ForeignKeyField(
24         Detector, Detector.detector_id, on_delete='CASCADE', backref='detector',
25 ↪column_name='detector_id'
26     )

```

The InputFile has five columns, one of which is a foreign key linking it to the Detector model. Note that we have used the `FileNameField` and `FileChecksumField` to take advantage of the `verify_checksum()` function. InputFile has a trigger that is executed after each update that is changing either the exposure or the file content (checksum). When one of these conditions is verified, then the corresponding row in the Data file will be removed, because we want to force the reprocessing of this file since it has changed. A similar trigger on delete is actually not needed because the Data model is linked to this model with an `on_delete` cascade option.

```

1 class Data(MAFwBaseModel):
2     @classmethod
3     def triggers(cls) -> list[Trigger]:
4         delete_plotter_sql = 'DELETE FROM plotter_output WHERE plotter_name =
5 ↪"PlugPlotter";'
6
7         insert_data_trigger = Trigger(
8             trigger_name='data_after_insert',
9             trigger_type=(TriggerWhen.After, TriggerAction.Insert),
10            source_table=cls,
11            safe=True,
12            for_each_row=False,
13        )
14        insert_data_trigger.add_sql(delete_plotter_sql)
15
16        update_data_trigger = Trigger(
17            trigger_name='data_after_update',
18            trigger_type=(TriggerWhen.After, TriggerAction.Update),
19            source_table=cls,
20            safe=True,
21            for_each_row=False,
22        )
23        update_data_trigger.add_when('NEW.value != OLD.value')
24        update_data_trigger.add_sql(delete_plotter_sql)
25
26        delete_data_trigger = Trigger(
27            trigger_name='data_after_delete',
28            trigger_type=(TriggerWhen.After, TriggerAction.Delete),
29            source_table=cls,
30            safe=True,
31            for_each_row=False,
32        )
33        delete_data_trigger.add_sql(delete_plotter_sql)

```

(continues on next page)

(continued from previous page)

```

33         return [insert_data_trigger, delete_data_trigger, update_data_trigger]
34
35     file_pk = ForeignKeyField(InputFile, on_delete='cascade', backref='file', primary_
36     ↪key=True, column_name='file_id')
37     value = FloatField(help_text='The result of the measurement')

```

The Data model has only two columns, one foreign key linking to the InputFile and one with the value calculated by the Analysis processor. It is featuring three triggers executed on INSERT, UPDATE and DELETE. In all these cases, we want to be sure that the output of the PlugPlotter is removed so that a new one is generated. Keep in mind that when a row is removed from the PlotterOutput model, the corresponding files are automatically added to the OrphanFile model for removal from the filesystem the next time a processor is executed.

Via the use of the foreign key it is possible to associate a detector and the exposure for this specific value.

10.2.2 The processor library

Let's now prepare one processor for each of the tasks that we have identified in our planning. We will create a processor also for the data generation.

GenerateDataFiles

This processor will accomplish Task 0 and it is very simple. It will generate a given number of files containing one single number calculated given the exposure, the slope and the intercept. The detector parameter is used to differentiate the output file name. As you see here below, the code is very simple.

```

1 class GenerateDataFiles(Processor):
2     n_files = ActiveParameter('n_files', default=25, help_doc='The number of 1-h
3     ↪increasing exposure')
4     output_path = ActiveParameter(
5     ↪'output_path', default=Path.cwd(), help_doc='The path where the data files
6     ↪are stored.'
7     )
8     slope = ActiveParameter(
9     ↪'slope', default=1.0, help_doc='The multiplication constant for the data
10    ↪stored in the files.'
11    )
12    intercept = ActiveParameter(
13    ↪'intercept', default=5.0, help_doc='The additive constant for the data stored
14    ↪in the files.'
15    )
16    detector = ActiveParameter(
17    ↪'detector', default=1, help_doc='The detector id being used. See the detector
18    ↪table for more info.'
19    )
20
21    def __init__(self, *args, **kwargs):
22        super().__init__(*args, **kwargs)
23        self.n_digits = len(str(self.n_files))
24
25    def start(self) -> None:
26        super().start()
27        self.output_path.mkdir(parents=True, exist_ok=True)
28
29    def get_items(self) -> Collection[Any]:
30        return list(range(self.n_files))

```

(continues on next page)

(continued from previous page)

```

27     def process(self) -> None:
28         current_filename = self.output_path / f'rawfile_exp{self.i_item:0{self.n_
↳digits}}_det{self.detector}.dat'
29         value = self.i_item * self.slope + self.intercept
30         with open(current_filename, 'wt') as f:
31             f.write(str(value))
32
33     def format_progress_message(self) -> None:
34         self.progress_message = f'Generating exposure {self.i_item} for detector
↳{self.detector}'

```

In order to generate the different detectors, you run the same processor with different values for the parameters.

PlugImporter

This processor will accomplish Task 1, i.e. import the raw data file into our database. This processor is inheriting from the basic *Importer* so that we can use the functionalities of the *FilenameParser*.

```

1 @database_required
2 class PlugImporter(Importer):
3     def __init__(self, *args: Any, **kwargs: Any) -> None:
4         super().__init__(*args, **kwargs)
5         self._data_list: list[dict[str, Any]] = []
6
7     def start(self) -> None:
8         super().start()
9         self.database.create_tables([InputFile])
10
11    def get_items(self) -> Collection[Any]:
12        pattern = '**/*dat' if self.recursive else '*dat'
13        input_folder_path = Path(self.input_folder)
14
15        file_list = [file for file in input_folder_path.glob(pattern) if file.is_
↳file()]
16
17        # verify the checksum of the elements in the input table. if they are not up_
↳to date, then remove the row.
18        verify_checksum(InputFile)
19
20        if self.filter_register.new_only:
21            # get the filenames that are already present in the input table
22            existing_rows = InputFile.select(InputFile.filename).namedtuples()
23            # create a set with the filenames
24            existing_files = {row.filename for row in existing_rows}
25            # filter out the file list from filenames that are already in the_
↳database.
26            file_list = [file for file in file_list if file not in existing_files]
27
28        return file_list
29
30    def process(self) -> None:
31        try:
32            new_file = {}
33            self._filename_parser.interpret(self.item.name)
34            new_file['filename'] = self.item
35            new_file['checksum'] = self.item
36            new_file['exposure'] = self._filename_parser.get_element_value('exposure')

```

(continues on next page)

(continued from previous page)

```

37         new_file['detector'] = self._filename_parser.get_element_value('detector')
38         self._data_list.append(new_file)
39     except ParsingError:
40         log.critical('Problem parsing %s' % self.item.name)
41         self.looping_status = LoopingStatus.Skip
42
43     def finish(self) -> None:
44         InputFile.insert_many(self._data_list).on_conflict_replace(replace=True).
↪execute()
45     super().finish()

```

The `get_items` is using the `verify_checksum()` to verify that the table is still actual and we apply the filter to be sure to process only new or modified files. The process and finish are very standard. In this specific case, we preferred to add all the relevant information in a list and insert them all in one single call to the database. But also the opposite approach (no storing, multiple insert) is possible.

Analyser

This processor will accomplish Task 2, i.e. the analysis of the files. In our case, we just need to open the file, read the value and put it in the database.

```

1 @database_required
2 class Analyser(Processor):
3     def start(self) -> None:
4         super().start()
5         self.database.create_tables([InputFile, Data])
6
7     def get_items(self) -> Collection[Any]:
8         self.filter_register.bind_all([InputFile])
9
10        existing_entries = Data.select(Data.file_pk).execute()
11
12        if self.filter_register.new_only:
13            existing = ~InputFile.file_pk.in_([i.file_pk for i in existing_entries])
14        else:
15            existing = True
16
17        query = InputFile.select().where(self.filter_register.filter_all()).
↪where(existing)
18        return query
19
20    def process(self) -> None:
21        with open(self.item.filename, 'rt') as fp:
22            value = float(fp.read())
23
24            Data.create(file_pk=self.item.file_pk, value=value)
25
26    def format_progress_message(self) -> None:
27        self.progress_message = f'Analysing {self.item.filename.name}'

```

Also in this case, the generation of the item list is done keeping in mind the possible filters the user is applying in the steering file. In the process, we are inserting the data directly to the database, so we will have one query for each item.

PlugPlotter

This processor will accomplish the last task, i.e. the generation of a relation plot where the performance of the three detectors is compared.

```

1 @database_required
2 @processor_depends_on_optional(module_name='seaborn')
3 @single_loop
4 class PlugPlotter(SQLPdDataRetriever, RelPlot, SNSPlotter):
5     output_plot_path = ActiveParameter(
6         'output_plot_path', default=Path('output.png'), help_doc='The filename of the
↳output plot'
7     )
8
9     def __init__(self, *args, **kwargs):
10        super().__init__(
11            *args,
12            table_name='data_view',
13            required_cols=['exposure', 'value', 'detector_name'],
14            x='exposure',
15            y='value',
16            hue='detector_name',
17            facet_kws=dict(legend_out=False, despine=False),
18            **kwargs,
19        )
20
21    def start(self) -> None:
22        super().start()
23
24        sql = """
25        CREATE TEMP VIEW IF NOT EXISTS data_view AS
26        SELECT
27            file_id, detector.detector_id, detector.name as detector_name, exposure,
↳value
28        FROM
29            data
30            JOIN input_file ON data.file_id = input_file.file_pk
31            JOIN detector USING (detector_id)
32        ORDER BY
33            detector.detector_id ASC,
34            input_file.exposure ASC
35        ;
36        """
37        self.database.execute_sql(sql)
38
39    def customize_plot(self):
40        self.facet_grid.set_axis_labels('Exposure', 'Value')
41        self.facet_grid.figure.subplots_adjust(top=0.9)
42        self.facet_grid.figure.suptitle('Data analysis results')
43        self.facet_grid._legend.set_title('Detector type')
44
45    def save(self) -> None:
46        self.facet_grid.figure.savefig(self.output_plot_path)
47        self.output_filename_list.append(self.output_plot_path)

```

This processor is a mixture of *SQLPdDataRetriever*, *RelPlot* and *SNSPlotter*.

Looking at the init method, you might notice a strange thing, the `table_name` variable is set to `data_view`, that does not corresponding to any of our tables. The reason for this strangeness is quickly explained.

The `SQLPdDataRetriever` is generating a `pandas` Dataframe from a SQL query. In our database the data table contains only two columns: the file reference and the measured value, but we have no direct access to the exposure nor to the detector. To get these other fields we need to join the data table with the `input_file` and the `detector` ones. The solution for this problem is the creation of a temporary view containing this join query. Have a look at the `start` method. This view will be deleted as soon as the connection will be closed.

10.2.3 The plugin module

Everything is ready, we just have to make MAFw aware of our processors and our standard tables. We are missing just a few lines of code in the `plugins` module

```
1 from plug.db_model import Detector
2 from plug.plugin_processor import Analyser, GenerateDataFiles, PlugImporter, PlugPlotter
3
4 import mafw
5
6
7 @mafw.mafw_hookimpl
8 def register_processors() -> list[mafw.processor.Processor]:
9     return [GenerateDataFiles, PlugImporter, Analyser, PlugPlotter]
10
11
12 @mafw.mafw_hookimpl
13 def register_standard_tables() -> list[mafw.db.std_tables.StandardTable]:
14     return [Detector]
```

The code is self-explaining. We need to invoke the two hooks and return the list of processors and standard tables that we want to export.

10.3 Run the code!

We are done with coding and we are ready to run our analysis.

First thing, we need to install our package in a separated python environment.

```
$ python -m venv my_venv
$ source my_venv/bin/activate
(my_venv) $ pip install -e /path/to/plugin
```

Now verify that the installation was successful. If you run `mafw list` command you should get the list of all available processors including the three that you have created.

```
(my_venv) $ mafw list
```

One last step, before running the analysis. We need to make the two steering files, one for the generation of the synthetic data and one for the real analysis and also the configuration file for the importer.

File generation

```
1 # file: generate-file.toml
2 processors_to_run = ["GenerateDataFiles"]
3
4 # customise the name of the analysis
5 analysis_name = "integration-test-p1"
6 analysis_description = """Generating data files"""
7 available_processors = ["GenerateDataFiles"]
8
9 [GenerateDataFiles]
```

(continues on next page)

(continued from previous page)

```

10 intercept = 5.0 # The additive constant for the data stored in the files.
11 n_files = 25 # The number of files being generated.
12 output_path = "/tmp/full-int/data" # The path where the data files are stored.
13 slope = 1.0 # The multiplication constant for the data stored in the files.
14 detector = 1 # The detector id being used. See the detector table for more info.
15
16 [UserInterface] # Specify UI options
17 interface = "rich" # Default "rich", backup "console"

```

Analysis

```

# file: analysis.toml
processors_to_run = ["PlugImporter", "Analyser", "PlugPlotter"]

# customise the name of the analysis
analysis_name = "integration-test-p2"
analysis_description = """"Analysing data""""
available_processors = ["PlugImporter", "Analyser", "PlugPlotter"]

[GlobalFilter]
new_only = true

[DBConfiguration]
URL = "sqlite:///tmp/full-int/plugin.db" # Change the protocol depending on the DB.
↳type. Update this file to the path of your DB.

[DBConfiguration.pragmas] # Leave these default values, unless you know what you are.
↳doing!
journal_mode = "wal"
cache_size = -64000
foreign_keys = 1
synchronous = 0

[PlugImporter]
input_folder = "/tmp/full-int/raw_data" # The input folder from where the images have.
↳to be imported.
parser_configuration = "/tmp/full-int/importer_config.toml" # The path to the TOML.
↳file with the filename parser configuration
recursive = true

[Analyser]

[PlugPlotter]
output_plot_path = "/tmp/full-int/output.png" # The filename of the output plot

[UserInterface] # Specify UI options
interface = "rich" # Default "rich", backup "console"

```

Importer configuration

```

1 # file: importer_config.toml
2 elements = ['exposure', 'detector']
3
4 [exposure]
5 regexp = '[_-]*exp(?P<exposure>\d+\.\d*)[_-]*'

```

(continues on next page)

(continued from previous page)

```
6 type='float'  
7  
8 [detector]  
9 regexp = '[_-]*det(?P<detector>\d+)[_-]*'  
10 type='int'
```

Adapt the steering files, in particular the paths and you are ready to run!

```
(my_venv) $ mafw run generate-file.toml
```

This will generate one set of files for the detector 1, now open the steering file in your favourite editor, change the detector to 2 and set the slope to 5 and the intercept to 15. Save the file and run it once again. Repeat it once more with detector 3, slope 0.2 and intercept 0.1 and you should now have 75 raw data files ready to be analysed.

```
(my_venv) $ mafw run analysis.toml
```

And here comes the magic! The three processors will be executed one after the other, the database is created and filled with all the provided data and the comparison plot is generated (default file name output.png).

This is just the beginning, now you can try all the benefits to use a clever database to drive your analysis pipeline. Try, for example, to remove one file and re-run the analysis, you will see a warning message informing you that a file was not found and that the corresponding row in the database has been removed as well. The rest of the analysis will remain the same, but the output plot will be regenerated with hole.

Try to manually modify the content of a file and re-run the analysis. The `verify_checksum()` will immediately detect the change and force the re-analysis of that file and the generation of the output plot.

Try to rename one file changing the exposure value. You will see that mafw will detect that one file is gone missing in action but a new one has been found. The output file will be update.

Try to repeat the analysis without any change and mafw won't do anything! Try to delete the output plot and during the next run mafw will regenerate it.

You can also play with the database. Open it in DBeaver (be sure that the `foreign_check` is enforced) and remove one line from the `input_file` table. Run the analysis and you will see that the output plot file is immediately removed because it is no more actual and a new one is generated at the end of chain.

It's not magic even if it really looks like so, it's just a powerful library for scientists written by scientists!

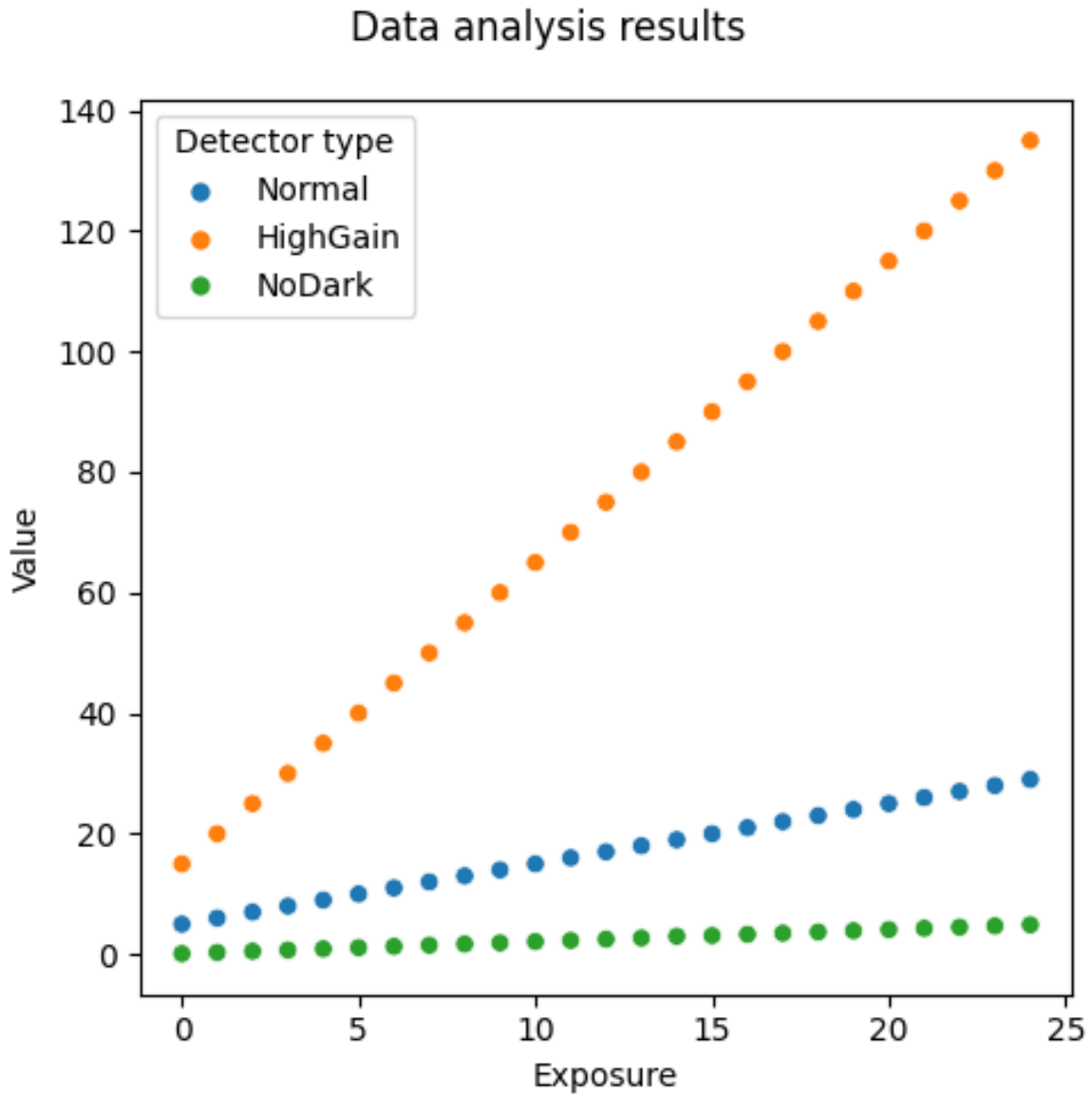


Fig. 10.3: The comparison plot of the three detectors.

CHAPTER 11

API

This page is particularly useful for the developers because it offers a quick view at the implemented API.

mafw

The Modular Analysis Framework

11.1 mafw

The Modular Analysis Framework

A software tool for scientists written by scientists!

Module Attributes

mafw_hookimpl

Marker to be imported and used in plugins loading

`mafw.mafw_hookimpl = <pluggy._hooks.HookimplMarker object>`

Marker to be imported and used in plugins loading

Modules

<i>active</i>	Module implements active variable for classes.
<i>db</i>	Defines the DB interface
<i>decorators</i>	The module provides some general decorator utilities that are used in several parts of the code, and that can be reused by the user community.
<i>enumerators</i>	Module provides a set of enumerators for dealing with standard tasks.
<i>examples</i>	A library of examples.
<i>hookspeccs</i>	Defines the hook specification decorator bound the MAFw library.
<i>mafw_errors</i>	Module defines MAFw exceptions
<i>plugin_manager</i>	Provides utilities to retrieve internal and external plugins.
<i>plugins</i>	Exports Processor classes to the execution script.

continues on next page

Table 11.3 – continued from previous page

<code>processor</code>	Module implements the basic Processor class, the ProcessorList and all helper classes to achieve the core functionality of the MAFw.
<code>processor_library</code>	Library of generic processors.
<code>runner</code>	Provides a container to run configurable and modular analytical tasks.
<code>scripts</code>	Executables.
<code>timer</code>	Module implements a simple timer to measure the execution duration.
<code>tools</code>	The package provides a set of tools for different range of applications.
<code>ui</code>	User interface modules.

11.1.1 mafw.active

Module implements active variable for classes.

Module Attributes

<code>ActiveType</code>	A type for templating the Active class.
<code>ActivableType</code>	A type for all classes that can include an Active.

Classes

<code>Active</code> ([default])	A descriptor class to make class variable active .
---------------------------------	-----------------------------------------------------------

class mafw.active.ActivableType

A type for all classes that can include an Active.

alias of `TypeVar('ActivableType')`

class mafw.active.Active(*default: ActiveType | None = None*)

Bases: `Generic[ActiveType]`

A descriptor class to make class variable **active**.

When assigned to a class variable, any change of this value will trigger a call back to a specific function.

Here is a clarifying example.

```
class Person:
    age = Active()

    def __init__(self, age):
        self.age = age

    def on_age_change(self, old_value, new_value):
        # callback invoked every time the value of age is changed
        # do something with the old and the new age
        pass

    def on_age_set(self, value):
        # callback invoked every time the value on age is set to the same
        # value as before.
        pass
```

(continues on next page)

(continued from previous page)

```
def on_age_get(self, value):
    # callback invoked every time the value of age is asked.
    # not really useful, but...
    pass
```

Once you have assigned an Active to a class member, you need to implement the callback in your class. If you do not implement them, the code will run without problems.

The three callbacks have the signature described in the example.

- `on_[var_name]_change(self, old, new)`
- `on_[var_name]_set(self, value)`
- `on_[var_name]_get(self, value)`

Constructor parameter:

Parameters

default (`ActiveType`) – Initial value of the Active value. Defaults to None.

class `maf.active.ActiveType`

A type for templating the Active class.

alias of `TypeVar('ActiveType')`

11.1.2 mafw.db

Defines the DB interface

Modules

<code>db_configurations</code>	Module provides default configurations for different database engines.
<code>db_filter</code>	The module provides the capability to filter model classes with information taken from the steering file.
<code>db_model</code>	The module provides functionality to MAFw to interface to a DB.
<code>db_types</code>	Database Type Definitions
<code>db_wizard</code>	The module allows to generate a DB model structure starting from an existing DB.
<code>fields</code>	Module provides customised model fields specific for MAFw.
<code>std_tables</code>	Module provides standard tables that are included in all database created by MAFw processor.
<code>trigger</code>	Module provides a Trigger class and related tools to create triggers in the database via the ORM.

maf.db.db_configurations

Module provides default configurations for different database engines.

Module Attributes

<code>default_conf</code>	default configuration dictionary used to generate steering files
<code>db_scheme</code>	default database scheme

```
mafwb.db.db_configurations.db_scheme = {'mysql': 'mysql://', 'postgresql':
'postgresql://', 'sqlite': 'sqlite:///'}

```

default database scheme

```
mafwb.db.db_configurations.default_conf = {'mysql': {'URL':
'mysql://user:passwd@ip:port/my_db'}, 'postgresql': {'URL':
'postgresql://postgres:my_password@localhost:5432/my_database'}, 'sqlite': {'URL':
'sqlite:///my_database.db', 'pragmas': {'cache_size': -64000, 'foreign_keys': 1,
'journal_mode': 'wal', 'synchronous': 0}}}

```

default configuration dictionary used to generate steering files

mafwb.db.db_filter

The module provides the capability to filter model classes with information taken from the steering file.

Classes

<code>Filter(name_, **kwargs)</code>	Class to filter rows from a model.
<code>FilterRegister([data])</code>	A special dictionary to store all <i>Filters</i> in a processors.

```
class mafwb.db.db_filter.Filter(name_: str, **kwargs: Any)

```

Bases: object

Class to filter rows from a model.

The filter object can be used to generate a where clause to be applied to `Model.select()`.

The construction of a `Filter` is normally done via a configuration file using the `from_conf()` class method. The name of the filter is playing a key role in this. If it follows a dot structure like:

ProcessorName.Filter.ModelName

then the corresponding table from the TOML configuration object will be used.

For each processor, there might be many `Filters`, up to one for each `Model` used to get the input list. If a processor is joining together three `Models` when performing the input select, there will be up to three `Filters` collaborating on making the selection.

The filter configuration can contain the following key, value pair:

- key / string pairs, where the key is the name of a field in the corresponding `Model`
- key / numeric pairs
- key / arrays

All fields from the configuration file will be added to the instance namespace, thus accessible with the dot notation. Moreover, the field names and their filter value will be added to a private dictionary to simplify the generation of the filter SQL code.

The user can use the filter object to store selection criteria. He can construct queries using the filter contents in the same way as he could use processor parameters.

If he wants to automatically generate valid filtering expression, he can use the `filter()` method. In order for this to work, the `Filter` object be *bound* to a `Model`. Without this binding the `Filter` will not be able to automatically generate expressions.

For each field in the filter, one condition will be generated according to the following scheme:

Filter field type	Logical operation	Example
Numeric, boolean	==	Field == 3.14
String	GLOB	Field GLOB '*rec'
List	IN	Field IN [1, 2, 3]

All conditions will be joined with a AND logic.

Consider the following example:

```

1 class MeasModel(MAFwBaseModel):
2     meas_id = AutoField(primary_key=True)
3     sample_name = TextField()
4     successful = BooleanField()
5
6    flt = Filter('MyProcessor.Filter.MyModel', sample_name='sample_00*',
7               meas_id=[1,2,3], successful=True)
8     flt.bind(MeasModel)
9
10    filter_select_manual = MeasModel.select().where((MeasModel.meas_id.in_(flt.meas_
11    ↪id) &
12    ↪sample_name) &
13    ↪successful))
14    filter_select_auto = MeasModel.select().where(flt.filter())

```

At line 6, we created a Filter and at 8 we bind it to our Model class. The two select at 10 and 14 are actually identical.

Constructor parameters:

Parameters

- **name_** (*str*) – The name of the filter. It should be in dotted format to facilitate the configuration via the steering file. The `_` is used to allow the user to have a keyword argument named name.
- **kwargs** – Keyword parameters corresponding to fields and filter values.

Changed in version 1.2.0: The parameter *name* has been renamed as *name_*.

classmethod from_conf(*name: str, conf: dict[str, Any], default: dict[str, Any] | None = None*) → Self
Builds a Filter object from a steering file dictionary.

If the name is in dotted notation, then this should be corresponding to the table in the configuration file. If a default configuration is provided, this will be used as a starting point for the filter, and it will be updated by the actual configuration in *conf*.

In normal use, you would provide the specific configuration via the *conf* parameter and the global filter configuration as default.

See details in the [class documentation](#)

Parameters

- **name** (*str*) – The name of the filter in dotted notation.
- **conf** (*dict*) – The configuration dictionary.
- **default** (*dict*) – Default configuration dictionary

Returns

A Filter object

Return type*Filter***bind**(*model*: *type[Model] | None = None*) → None

Connects a filter to a Model class.

If no model is provided, the method will try to bind a class from with global dictionary with a name matching the model name used during the Filter configuration. It only works when the Model is defined as global.

Parameters**model** (*Model*, *Optional*) – Model to be bound. Defaults to None**filter**() → Expression | bool

Generates a filtering expression logically ANDing all filtering fields.

See details in the [class documentation](#)**Returns**

The filtering expression.

Return type

peewee.Expression

Raises**TypeError** – when the field value type is not supported.**get_field**(*key*: *str*) → Any

Gets a field by name.

Parameters**key** (*str*) – The name of the field.**Returns**

The value of the field.

Return type

Any

Raises**KeyError** – if the requested field does not exist.**set_field**(*key*: *str*, *value*: *Any*) → None

Sets the value of a field by name

Parameters

- **key** (*str*) – The name of the field.
- **value** (*Any*) – The value of the field.

property field_names: **list[str]**

The list of field names.

property is_bound: **bool**

Returns true if the Filter has been bound to a Model

class mafw.db.db_filter.**FilterRegister**(*data*: *dict[str, Filter] | None = None*, / (*Positional-only parameter separator (PEP 570)*), ***kwargs*: *Any*)Bases: UserDict[str, *Filter*]A special dictionary to store all *Filters* in a processors.

It contains a publicly accessible dictionary with the configuration of each Filter using the Model name as keyword.

It contains a private dictionary with the global filter configuration as well. The global filter is not directly accessible, but only some of its members will be exposed via properties. In particular, the `new_only` flag

that is relevant only at the Processor level can be accessed directly using the *new_only*. If not specified in the configuration file, the *new_only* is by default True.

Constructor parameters:

Parameters

- **data** (*dict*) – Initial data
- **kwargs** – Keywords arguments

bind_all(*models: list[type[Model]] | dict[str, type[Model]]*) → None

Binds all filters to their models.

The *models* list or dictionary should contain a valid model for all the Filter in the registry. In the case of a dictionary, the key value should be the model name.

If the user provides a model for which there is no corresponding filter in the register, then a new filter for that model is created using the GlobalFilter default.

This can happen if the user did not provide a specific table for the Processor/Model, but simply put all filtering conditions in the GlobalFilter table. Even though, this behaviour is allowed and working, it may result in unexpected results. Also listing more than needed models in the input list can be dangerous because they will anyhow use the default filters.

Parameters

models (*list[type(Model)] | dict[str, type(Model)]*) – List or dictionary of a databank of Models from which the Filter can be bound.

filter_all()

Generates a where clause ANDing all filters.

If one Filter is not bound, then True is returned.

property new_only: bool

The new only flag.

Returns

True, if only new items, not already in the output database table must be processed.

Return type

bool

mafw.db.db_model

The module provides functionality to MAFw to interface to a DB.

Module Attributes

<i>database_proxy</i>	This is a placeholder for the real database object that will be known only at run time
-----------------------	----------------------------------------------------------------------------------------

Classes

<i>MAFwBaseModel</i> (*args, **kwargs)	The base model for the MAFw library.
----------------------------------------	--------------------------------------

Exceptions

<i>MAFwBaseModelDoesNotExist</i>	Raised when the base model class is not existing.
----------------------------------	---------------------------------------------------

exception `mafwb.db.db_model.MAFwBaseModelDoesNotExist`

Bases: `MAFwException`

Raised when the base model class is not existing.

class `mafwb.db.db_model.MAFwBaseModel(*args, **kwargs)`

Bases: `Model`

The base model for the MAFw library.

Every model class (table) that the user wants to interface must inherit from this base.

DoesNotExist

alias of `MAFwBaseModelDoesNotExist`

classmethod `create_table(safe: bool = True, **options) → None`

Create the table in the underlying DB and all the related trigger as well.

If the creation of a trigger fails, then the whole table dropped, and the original exception is re-raised.

Warning

Trigger creation has been extensively tested with `SQLite`, but not with the other database implementation. Please report any malfunction.

Parameters

- **safe** (*bool*, *Optional*) – Flag to add an IF NOT EXISTS to the creation statement. Defaults to True.
- **options** – Additional options passed to the super method.

classmethod `std_upsert(_MAFwBaseModel__data: dict[str, Any] | None = None, **mapping: Any) → ModelInsert`

Perform a so-called standard upsert.

An upsert statement is not part of the standard SQL and different databases have different ways to implement it. This method will work for modern versions of `sqlite` and `postgresql`. Here is a [detailed explanation for SQLite](#).

An upsert is a statement in which we try to insert some data in a table where there are some constraints. If one constraint is failing, then instead of inserting a new row, we will try to update the existing row causing the constraint violation.

A standard upsert, in the naming convention of MAFw, is setting the conflict cause to the primary key with all other fields being updated. In other words, the database will try to insert the data provided in the table, but if the primary key already exists, then all other fields will be updated.

This method is equivalent to the following:

```
class Sample(MAFwBaseModel):
    sample_id = AutoField(
        primary_key=True,
        help_text='The sample id primary key',
    )
    sample_name = TextField(help_text='The sample name')

(
    Sample.insert(sample_id=1, sample_name='my_sample')
    .on_conflict(
        preserve=[Sample.sample_name]
```

(continues on next page)

(continued from previous page)

```

) # use the value we would have inserted
.execute()
)

```

Parameters

- **__data** (*dict*, *Optional*) – A dictionary containing the key/value pair for the insert. The key is the column name. Defaults to None
- **mapping** – Keyword arguments representing the value to be inserted.

classmethod `std_upsert_many`(*rows: Iterable[Any]*, *fields: list[str] | None = None*) → `ModelInsert`
 Perform a standard upsert with many rows.

See also

Read the `std_upsert()` documentation for an explanation of this method.

Parameters

- **rows** (*Iterable*) – A list with the rows to be inserted. Each item can be a dictionary or a tuple of values. If a tuple is provided, then the *fields* must be provided.
- **fields** (*list[str]*, *Optional*) – A list of field names. Defaults to None.

classmethod `triggers`() → `list[Trigger]`

Returns an iterable of `Trigger` objects to create upon table creation.

The user must overload this returning all the triggers that must be created along with this class.

`maf.db.db_model.database_proxy = <peewee.DatabaseProxy object>`

This is a placeholder for the real database object that will be known only at run time

maf.db.db_types

Database Type Definitions

This module provides type definitions for database models and related components.

It defines Protocol classes that represent the expected interfaces for database models, helping with static type checking when working with ORM frameworks like Peewee.

Classes

<code>PeeweeModelWithMeta(*args, **kwargs)</code>	Protocol defining the interface for Peewee model classes with metadata.
---------------------------------------------------	-------------------------------------------------------------------------

class `maf.db.db_types.PeeweeModelWithMeta`(*args, **kwargs)

Bases: `Protocol`

Protocol defining the interface for Peewee model classes with metadata.

This Protocol helps with static type checking for Peewee ORM models, ensuring that objects passed to functions expecting Peewee models have the necessary methods and attributes.

Attributes

`_meta` : Any

The metadata object associated with the Peewee model class. Contains information about the table, fields, and other model properties.

Methods

`select(*args: Any, **kwargs: Any) -> Any`

Select records from the database.

`delete(*args: Any, **kwargs: Any) -> Any`

Delete records from the database.

maf.db.db_wizard

The module allows to generate a DB model structure starting from an existing DB.

It is strongly based on the `peewee` playhouse module `pwiz`.

Functions

<code>dump_models(output_file, introspector[, ...])</code>	Dumps all the ORM models in the output file.
------------------------------------------------------------	----------------------------------------------

Classes

<code>UnknownField(*_, **_)</code>	Placeholder class for an Unknown Field.
------------------------------------	-----------------------------------------

class `maf.db.db_wizard.UnknownField(*_: Any, **_: Any)`

Bases: `object`

Placeholder class for an Unknown Field.

`maf.db.db_wizard.dump_models(output_file: TextIO, introspector: Introspector, tables: list[str] | tuple[str, ...] | None = None, preserve_order: bool = True, include_views: bool = False, ignore_unknown: bool = False, snake_case: bool = True) -> None`

Dumps all the ORM models in the output file.

This function will write to the output stream a fully functional python module with all the models that the introspector class can access. The user has the possibility to pre-select a subset of tables to be dumped and also to optionally include views.

See also

This function is the core of the `db wizard`.

Parameters

- **output_file** (`TextIO`) – The output file. It must be open in text/write mode.
- **introspector** (`Introspector`) – A database introspector from the reflection module.
- **tables** (`list | None, Optional`) – A list of table names to be dumped. If `None`, then all tables in the DB will be dumped. Defaults to `None`.
- **preserve_order** (`bool, Optional`) – Preserve the order of the columns in the model. Defaults to `True`.
- **include_views** (`bool, Optional`) – Include saved views to be dumped. Views can be obtained joining tables. Defaults to `False`.

- **ignore_unknown** (*bool, Optional*) – Ignore unknown fields. If True, then an UnknownField type will be used as a placeholder. Defaults to False.
- **snake_case** (*bool, Optional*) – Use snake case for column and table name. Defaults to True.

mafwb.db.fields

Module provides customised model fields specific for MAFw.

Classes

<code>FileChecksumField</code> ([null, index, unique, ...])	A field to be used for file checksum.
<code>FileChecksumFieldAccessor</code> (model, field, name)	Field accessor specialized for file checksum fields.
<code>FileNameField</code> ([checksum_field])	Field to be used for filenames.
<code>FileNameFieldAccessor</code> (model, field, name)	A field accessor specialized for filename fields.
<code>FileNameListField</code> ([checksum_field])	A field for a list of file names.

```
class mafwb.db.fields.FileChecksumField(null=False, index=False, unique=False,
column_name=None, default=None, primary_key=False,
constraints=None, sequence=None, collation=None,
unindexed=False, choices=None, help_text=None,
verbose_name=None, index_type=None, db_column=None,
_hidden=False)
```

Bases: `TextField`

A field to be used for file checksum.

It is the evolution of the `TextField` for storing file checksum hexadecimal digest.

If linked to a `FileNameField` or `FileNameListField`, then it will be automatically filled when the primary file name (or list of file names) field is set.

If the user decides to set its value manually, then he can provide either the string with the hexadecimal characters as calculated by `file_checksum()`, or simply the filename (or filename list) and the field will perform the calculation automatically.

accessor_class

The specific field accessor class

alias of `FileChecksumFieldAccessor`

db_value(*value: str | Path | list[str | Path]*) → str

Converts the python assigned value to the DB type.

The checksum will be stored in the DB as a string containing only hexadecimal characters (see `hexdigest`).

The user can provide the checksum directly or the path to the file or a list of path to files. If a `Path`, or list of `Path`, is provided, then the checksum will be calculated, if a `str` (a path converted into a string) is provided, the function will try to see if a file with that path exists. If so, the checksum will be calculated, if not, the original string is assumed to be the checksum.

Parameters

value (*str | Path | list[str | Path]*) – The checksum or path to the file, or list of path to files for which the checksum has to be stored.

Returns

The checksum string for the database storage.

Return type

str

class mafw.db.fields.**FileChecksumFieldAccessor**(*model, field, name*)

Bases: `FieldAccessor`

Field accessor specialized for file checksum fields.

When the field is directly set, then an initialization flag in the model instance is turned to `True` to avoid that the linked primary field will overrule this value again.

For each checksum field named `my_checksum`, the model instance will get an attribute: `init_my_checksum` to be used as an initialization flag.

Once the field is manually set, to re-establish the automatic mechanism, the user has to manually toggle the initialization flag.

class mafw.db.fields.**FileNameField**(*checksum_field: str | None = None, *args: Any, **kwargs: Any*)

Bases: `TextField`

Field to be used for filenames.

It is just an overload of `TextField`, that allows to apply filters and python functions specific to filenames.

If the user specifies the name of a file checksum field, then when this field is updated, the checksum one will also be automatically updated.

See also

- [FileNameListField](#) for a field able to store a list of filenames.
- [remove_widow_db_rows\(\)](#) for a function removing entries from a database table where the corresponding files on disk are missing.
- [verify_checksum\(\)](#) for a function comparing the actual checksum with the stored one and in case removing outdated entries from the DB.

Constructor parameter:

Parameters

checksum_field (*str, Optional*) – The name of the checksum field linked to this filename. Defaults to `None`.

accessor_class

The specific accessor class

alias of [FileNameFieldAccessor](#)

db_value (*value: str | Path*) → `str`

Converts the input python value into a string for the DB.

python_value (*value: str*) → `Path | None`

Converts the db value from `str` to `Path`

The return value might also be `None`, if the user set the field value to null.

Parameters

value (*str*) – The value of the field as stored in the database.

Returns

The converted value as a path. It can be `None`, if value was stored as null.

Return type

`Path | None`

class mafw.db.fields.**FileNameFieldAccessor**(*model, field, name*)

Bases: `FieldAccessor`

A field accessor specialized for filename fields.

In the constructor of the `FileNameField` and subclasses, the user can specify the name of a checksum field linked to this filename. This is very useful because in this way, the user does not have to manually assign any value to this field that will simply be automatically updated when the filename field is updated.

The user can disable this automatic feature either removing the link in the `FileNameField` or simply assigning a value to the `FileChecksumField`.

class mafw.db.fields.**FileNameListField**(*checksum_field: str | None = None, *args: Any, **kwargs: Any*)

Bases: `FileNameField`

A field for a list of file names.

The evolution of the `FileNameField`, this field is able to store a list of filenames as a ';' separated string of full paths.

It is meant to be used when a processor is saving a bunch of correlated files that are to be used all together.

In a similar way as its parent class, it can be link to a checksum field, in this case, the checksum of the whole file list will be calculated.

Constructor parameter:

Parameters

checksum_field (*str, Optional*) – The name of the checksum field linked to this filename. Defaults to None.

db_value (*value: list[str | Path] | str | Path*) → str

Converts the list of paths in a ';' separated string

python_value (*value: str*) → list[Path]

Converts the ';' separated string in a list of paths

mafw.db.std_tables

Module provides standard tables that are included in all database created by MAFw processor.

Standard tables are automatically created and initialized by a `Processor` or a `ProcessorList` when opening a database connection.

This means that if a processor receives a valid database object, then it will suppose that the connection was already opened somewhere else (either from a `ProcessorList` or a third party) and thus it is not creating the standard tables.

If a processor is constructed using a database configuration dictionary, then it will first try to open a connection to the DB, then creating all standard tables and finally executing their `StandardTable.init` method. The same apply for the `Processor list`.

In other words, object responsible to open the database connection is taking care also of creating the standard tables and of initializing them. If the user opens the connection and passes it to a `Processor` or `ProcessorList`, then the user is responsible to create the standard tables and to initialize them.

All standard tables must derive from the `StandardTable` to have the same interface for the initialization.

Users can create their own standard tables and export them as a pluggable object using the same approach as for custom processors. An example is provided in the [documentation](#).

Module Attributes

<i>standard_tables</i>	A dictionary containing the standard tables being exported to the rest of the framework.
------------------------	------------------------------------------------------------------------------------------

Classes

<i>OrphanFile</i> (*args, **kwargs)	A Model for the files to be removed from disc
<i>PlotterOutput</i> (*args, **kwargs)	A model for the output of the plotter processors.
<i>StandardTable</i> (*args, **kwargs)	A base class for tables that are generated automatically by the MAFw processor.
<i>TriggerDisabler</i> (trigger_type_id)	A helper tool to disable a specific type of triggers.
<i>TriggerStatus</i> (*args, **kwargs)	A Model for the trigger status

Exceptions

<i>OrphanFileDoesNotExist</i>	An exception raised when trying to access a not existing table.
<i>PlotterOutputDoesNotExist</i>	An exception raised when trying to access a not existing table.
<i>StandardTableDoesNotExist</i>	An exception raised when trying to access a not existing table.
<i>TriggerStatusDoesNotExist</i>	An exception raised when trying to access a not existing table.

exception mafw.db.std_tables.**OrphanFileDoesNotExist**

Bases: `DoesNotExist`

An exception raised when trying to access a not existing table.

exception mafw.db.std_tables.**PlotterOutputDoesNotExist**

Bases: `DoesNotExist`

An exception raised when trying to access a not existing table.

exception mafw.db.std_tables.**StandardTableDoesNotExist**

Bases: `Exception`

An exception raised when trying to access a not existing table.

exception mafw.db.std_tables.**TriggerStatusDoesNotExist**

Bases: `Exception`

An exception raised when trying to access a not existing table.

class mafw.db.std_tables.**OrphanFile**(*args, **kwargs)

Bases: *StandardTable*

A Model for the files to be removed from disc

DoesNotExist

alias of *OrphanFileDoesNotExist*

class mafw.db.std_tables.**PlotterOutput**(*args, **kwargs)

Bases: *StandardTable*

A model for the output of the plotter processors.

The model has a trigger activated on delete queries to insert filenames and checksum in the OrphanFile model.

DoesNotExist

alias of *PlotterOutputDoesNotExist*

classmethod `triggers()` → list[*Trigger*]

Returns an iterable of *Trigger* objects to create upon table creation.

The user must overload this returning all the triggers that must be created along with this class.

class `mafwb.db.std_tables.StandardTable(*args, **kwargs)`

Bases: *MAFwBaseModel*

A base class for tables that are generated automatically by the MAFw processor.

DoesNotExist

alias of *StandardTableDoesNotExist*

classmethod `init()` → None

The user must overload this method, if he wants some specific operations to be performed on the model everytime the database is connected.

class `mafwb.db.std_tables.TriggerDisabler(trigger_type_id: int)`

Bases: object

A helper tool to disable a specific type of triggers.

Not all SQL dialects allow to temporarily disable trigger execution.

In order overcome this limitation, MAFw has introduced a practical workaround. All types of triggers are active by default but they can be temporarily disabled, by changing their status in the *TriggerStatus* table.

In order to disable the trigger execution, the user has to set the status of the corresponding status to 0 and also add a when condition to the trigger definition.

Here is an example code:

```
class MyTable(MAFwBaseModel):
    id_ = AutoField(primary_key=True)
    integer = IntegerField()
    float_num = FloatField()

    @classmethod
    def triggers(cls):
        return [
            Trigger(
                'mytable_after_insert',
                (TriggerWhen.After, TriggerAction.Insert),
                cls,
                safe=True,
            )
            .add_sql(
                'INSERT INTO target_table (id_id, half_float_num) VALUES (NEW.
→ id_, NEW.float_num / 2)'
            )
            .add_when(
                '1 == (SELECT status FROM trigger_status WHERE trigger_type_id_
→ == 1)'
            )
        ]
```

When you want to perform a database action with the trigger disabled, you can either use this class as context manager or call the *disable()* and *enable()* methods.

```
# as a context manager
with TriggerDisabler(trigger_type_id = 1):
    # do something without triggering any trigger of type 1.

# with the explicit methods
disabler = TriggerDisabler(1)
disabler.disable()
# do something without triggering any trigger of type 1.
disabler.enable()
```

When using the two explicit methods, the responsibility to assure that the triggers are re-enabled in on the user.

Constructor parameters:

Parameters

trigger_type_id (*int*) – the id of the trigger to be temporary disabled.

disable() → None

Disable the trigger

enable() → None

Enable the trigger

class mafw.db.std_tables.**TriggerStatus**(*args, **kwargs)

Bases: *StandardTable*

A Model for the trigger status

DoesNotExist

alias of *TriggerStatusDoesNotExist*

classmethod **init()** → None

Resets all triggers to enable when the database connection is opened.

```
mafw.db.std_tables.standard_tables: dict[str, type[StandardTable]] = {'OrphanFile':
<Model: OrphanFile>, 'PlotterOutput': <Model: PlotterOutput>, 'TriggerStatus':
<Model: TriggerStatus>}
```

A dictionary containing the standard tables being exported to the rest of the framework.

The key is the name of the model and the value is the model class itself.

mafw.db.trigger

Module provides a Trigger class and related tools to create triggers in the database via the ORM.

It supports SQLite, MySQL and PostgreSQL with dialect-specific SQL generation.

Functions

<i>and_</i> (*conditions)	Concatenates conditions with logical AND.
<i>or_</i> (*conditions)	Concatenates conditions with logical OR.

Classes

<i>MySQLDialect</i> ()	MySQL-specific trigger SQL generation.
<i>PostgreSQLDialect</i> ()	PostgreSQL-specific trigger SQL generation.
<i>SQLiteDatabase</i> ()	SQLite-specific trigger SQL generation.
<i>Trigger</i> (trigger_name, trigger_type, source_table)	Trigger template wrapper for use with peewee ORM.

continues on next page

Table 11.20 – continued from previous page

<i>TriggerAction</i> (*values)	String enumerator for the trigger action (Delete, Insert, Update)
<i>TriggerDialect</i> ()	Abstract base class for database-specific trigger SQL generation.
<i>TriggerWhen</i> (*values)	String enumerator for the trigger execution time (Before, After or Instead Of)

class mafw.db.trigger.MySQLDialectBases: *TriggerDialect*

MySQL-specific trigger SQL generation.

create_trigger_sql(trigger: *Trigger*) → str

Generate MySQL trigger SQL.

drop_trigger_sql(trigger_name: str, safe: bool = True) → str

Generate MySQL drop trigger SQL.

supports_safe_create() → bool

MySQL supports IF NOT EXISTS for triggers.

supports_trigger_type(when: *TriggerWhen*, action: *TriggerAction*, on_view: bool = False) → bool

MySQL doesn't support INSTEAD OF triggers.

supports_update_of_columns() → bool

MySQL doesn't support column-specific UPDATE triggers.

supports_when_clause() → bool

MySQL supports conditions but through WHERE instead of WHEN.

class mafw.db.trigger.PostgreSQLDialectBases: *TriggerDialect*

PostgreSQL-specific trigger SQL generation.

_clean_sql(sql: str) → str

Remove RETURNING clauses from SQL statements for PostgreSQL trigger functions.

Parameters**sql** – The SQL statement**Returns**

SQL statement without RETURNING clause

create_trigger_sql(trigger: *Trigger*) → str

Generate PostgreSQL trigger SQL.

drop_trigger_sql(trigger_name: str, safe: bool = True) → str

Generate PostgreSQL drop trigger SQL.

supports_safe_create() → bool

PostgreSQL doesn't support IF NOT EXISTS for triggers before v14, but we implement safety differently.

supports_trigger_type(when: *TriggerWhen*, action: *TriggerAction*, on_view: bool = False) → bool

PostgreSQL supports INSTEAD OF only on views.

supports_update_of_columns() → bool

PostgreSQL supports column-specific UPDATE triggers.

supports_when_clause() → bool

PostgreSQL supports WHEN conditions.

class mafw.db.trigger.SQLiteDialect

Bases: *TriggerDialect*

SQLite-specific trigger SQL generation.

create_trigger_sql(*trigger*: *Trigger*) → str

Generate SQLite trigger SQL.

drop_trigger_sql(*trigger_name*: str, *safe*: bool = True) → str

Generate SQLite drop trigger SQL.

supports_safe_create() → bool

SQLite supports IF NOT EXISTS for triggers.

supports_trigger_type(*when*: *TriggerWhen*, *action*: *TriggerAction*, *on_view*: bool = False) → bool

SQLite supports all trigger types except INSTEAD OF on tables (only on views).

supports_update_of_columns() → bool

SQLite supports column-specific UPDATE triggers.

supports_when_clause() → bool

SQLite supports WHEN conditions.

class mafw.db.trigger.Trigger(*trigger_name*: str, *trigger_type*: tuple[*TriggerWhen*, *TriggerAction*],
source_table: type[*Model*] | *Model* | str, *safe*: bool = False,
for_each_row: bool = False, *update_columns*: list[str] | None = None,
on_view: bool = False)

Bases: object

Trigger template wrapper for use with peewee ORM.

Constructor parameters:

Parameters

- **trigger_name** (str) – The name of this trigger. It needs to be unique!
- **trigger_type** (tuple[*TriggerWhen*, *TriggerAction*]) – A tuple with *TriggerWhen* and *TriggerAction* to specify on which action the trigger should be invoked and if before, after or instead of.
- **source_table** (type[*Model*] | *Model* | str) – The table originating the trigger. It can be a model class, instance, or also the name of the table.
- **safe** (bool, *Optional*) – A boolean flag to define if in the trigger creation statement a ‘IF NOT EXISTS’ clause should be included. Defaults to False
- **for_each_row** (bool, *Optional*) – A boolean flag to repeat the script content for each modified row in the table. Defaults to False.
- **update_columns** (list[str], *Optional*) – A list of column names. When defining a trigger on a table update, it is possible to restrict the firing of the trigger to the cases when a subset of all columns have been updated. An column is updated also when the new value is equal to the old one. If you want to discriminate this case, use the *add_when()* method. Defaults to None.
- **on_view** (bool, *Optional*) – A boolean flag to indicate if the target is a view. This affects the support for INSTEAD OF. Defaults to False.

_get_dialect() → *TriggerDialect*

Get the appropriate dialect based on the database type.

Returns

A dialect instance

add_sql(*sql: str | Query*) → Self

Add an SQL statement to be executed by the trigger.

The sql can be either a string containing the sql statement, or it can be any other peewee Query.

For example:

```
# assuming you have created a trigger ...

sql = AnotherTable.insert(
    field1=some_value, field2=another_value
)
trigger.add_sql(sql)
```

In this way the SQL code is generated with parametric placeholder if needed.

Parameters

sql (*str | peewee.Query*) – The SQL statement.

Returns

self for easy chaining

Return type

Trigger

add_when(**conditions: str*) → Self

Add conditions to the *when* statements.

Conditions are logically ANDed. To have mixed *OR* and *AND* logic, use the functions *and_()* and *or_()*.

Parameters

conditions (*str*) – Conditions to be added with logical AND

Returns

self for easy chaining

Return type

Trigger

create() → str

Generates the SQL create statement.

Returns

The trigger creation statement.

Raises

- *MissingSQLStatement* – if no SQL statements are provided.
- *UnsupportedDatabaseError* – if the trigger type is not supported by the database.

drop(*safe: bool = True*) → str

Generates the SQL drop statement.

Parameters

safe (*bool, Optional*) – If True, add an IF EXIST. Defaults to True.

Returns

The drop statement

Return type

str

set_database(*database: Database | DatabaseProxy*) → Self

Set the database to use for this trigger.

Parameters

database – The database instance

Returns

self for easy chaining

class mafw.db.trigger.**TriggerAction**(*values)

Bases: StrEnum

String enumerator for the trigger action (Delete, Insert, Update)

static **_generate_next_value_**(name, start, count, last_values)

Return the lower-cased version of the member name.

class mafw.db.trigger.**TriggerDialect**

Bases: ABC

Abstract base class for database-specific trigger SQL generation.

abstractmethod **create_trigger_sql**(trigger: Trigger) → str

Generate the SQL to create a trigger for a specific database dialect.

Parameters

trigger – The trigger object

Returns

SQL string to create the trigger

abstractmethod **drop_trigger_sql**(trigger_name: str, safe: bool = True) → str

Generate the SQL to drop a trigger for a specific database dialect.

Parameters

- **trigger_name** – The name of the trigger to drop
- **safe** – If True, add an IF EXISTS clause. Defaults to True.

Returns

SQL string to drop the trigger

abstractmethod **supports_safe_create**() → bool

Check if the database supports IF NOT EXISTS for triggers.

Returns

True if supported, False otherwise

abstractmethod **supports_trigger_type**(when: TriggerWhen, action: TriggerAction, on_view: bool = False) → bool

Check if the database supports the specified trigger type.

Parameters

- **when** – When the trigger should fire (BEFORE, AFTER, INSTEAD OF)
- **action** – The action that triggers the trigger (INSERT, UPDATE, DELETE)
- **on_view** – Whether the trigger is on a view

Returns

True if supported, False otherwise

abstractmethod **supports_update_of_columns**() → bool

Check if the database supports column-specific UPDATE triggers.

Returns

True if supported, False otherwise

abstractmethod supports_when_clause() → bool

Check if the database supports WHEN conditions.

Returns

True if supported, False otherwise

class mafw.db.trigger.TriggerWhen(*values)

Bases: StrEnum

String enumerator for the trigger execution time (Before, After or Instead Of)

static _generate_next_value_(name, start, count, last_values)

Return the lower-cased version of the member name.

mafw.db.trigger.and_(*conditions: str) → str

Concatenates conditions with logical AND.

Parameters

conditions (*str*) – The condition to join.

Returns

The and-concatenated string of conditions

Return type

str

mafw.db.trigger.or_(*conditions: str) → str

Concatenates conditions with logical OR.

Parameters

conditions (*str*) – The condition to join.

Returns

The or-concatenated string of conditions.

Return type

str

11.1.3 mafw.decorators

The module provides some general decorator utilities that are used in several parts of the code, and that can be reused by the user community.

Module Attributes

<i>F</i>	TypeVar for generic function.
<i>P</i>	TypeVar for generic processor.
<i>single_loop</i> (cls)	A decorator shortcut to define a single execution processor.
<i>for_loop</i> (cls)	A decorator shortcut to define a for loop execution processor.
<i>while_loop</i> (cls)	A decorator shortcut to define a while loop execution processor.

Functions

<i>class_depends_on_optional</i> (module_name[, ...])	Class decorator factory to check if module module_name is available.
<i>database_required</i> (cls)	Modify the processor start method to check if a database object exists.

continues on next page

Table 11.22 – continued from previous page

<code>depends_on_optional</code> (<i>module_name</i> [, <i>raise_ex</i> , ...])	Function decorator to check if <i>module_name</i> is available.
<code>execution_workflow</code> (<i>loop_type</i>)	A decorator factory for the definition of the looping strategy.
<code>for_loop</code> (<i>cls</i>)	A decorator shortcut to define a for loop execution processor.
<code>orphan_protector</code> (<i>cls</i>)	A class decorator to modify the <code>init</code> method of a Processor so that the <code>remove_orphan_files</code> is set to <code>False</code> and no orphan files will be removed.
<code>processor_depends_on_optional</code> (<i>module_name</i> [, ...])	Class decorator factory to check if module <i>module_name</i> is available.
<code>single_loop</code> (<i>cls</i>)	A decorator shortcut to define a single execution processor.
<code>singleton</code> (<i>cls</i>)	Make a class a Singleton class (only one instance)
<code>suppress_warnings</code> (<i>func</i>)	Decorator to suppress warnings during the execution of a test function.
<code>while_loop</code> (<i>cls</i>)	A decorator shortcut to define a while loop execution processor.

class `maf.decorators.F`

TypeVar for generic function.

alias of `TypeVar('F', bound=Callable[[...], object])`

class `maf.decorators.P`

TypeVar for generic processor.

alias of `TypeVar('P', bound=Processor)`

`maf.decorators.class_depends_on_optional`(*module_name*: *str*, *raise_ex*: *bool* = *False*, *warn*: *bool* = *True*) → `Callable[[Type[Any]], Type[Any]]`

Class decorator factory to check if module *module_name* is available.

It checks if all the optional modules listed in *module_name* separated by a ';' can be found.

If all modules are found, then the class is returned as it is.

If at least one module is not found:

- and *raise_ex* is `True`, an `ImportError` exception is raised and the user is responsible to deal with it.
- if *raise_ex* is `False`, instead of returning the class, a new empty class is returned.
- depending on the value of *warn*, the user will be informed with a warning message or not.

Parameters

- **module_name** (*str*) – The optional module(s) from which the class depends on. A “;” separated list of modules can also be provided.
- **raise_ex** (*bool*, *Optional*) – Flag to raise an exception if *module_name* not found, defaults to `False`.
- **warn** (*bool*, *Optional*) – Flag to display a warning message if *module_name* is not found, defaults to `True`.

Returns

The wrapped class.

Return type

`type(object)`

Raises

ImportError – if *module_name* is not found and *raise_ex* is `True`.

`mafw.decorators.database_required(cls)`

Modify the processor start method to check if a database object exists.

This decorator must be applied to processors requiring a database connection.

Parameters

cls – A Processor class.

`mafw.decorators.depends_on_optional(module_name: str, raise_ex: bool = False, warn: bool = True)`
 → Callable[[F], Callable[[...], Any]]

Function decorator to check if `module_name` is available.

If `module_name` is found, then returns the wrapped function. If not, its behavior depends on the `raise_ex` and `warn_only` values. If `raise_ex` is `True`, then an `ImportError` exception is raised. If it is `False` and `warn` is `True`, then a warning message is displayed but no exception is raised. If they are both `False`, then function is silently skipped.

If `raise_ex` is `True`, the value of `warn` is not taken into account.

Typical usage

The user should decorate functions or class methods when they cannot be executed without the optional library. In the specific case of Processor subclass, where the class itself can be created also without the missing library, but it is required somewhere in the processor execution, then the user is suggested to decorate the execute method with this decorator.

Parameters

- **module_name** (*str*) – The optional module(s) from which the function depends on. A “;” separated list of modules can also be provided.
- **raise_ex** (*bool*, *Optional*) – Flag to raise an exception if `module_name` is not found, defaults to `False`.
- **warn** (*bool*, *Optional*) – Flag to display a warning message if `module_name` is not found, default to `True`.

Returns

The wrapped function

Return type

Callable

Raises

ImportError – if `module_name` is not found and `raise_ex` is `True`.

`mafw.decorators.execution_workflow(loop_type: LoopType | str = LoopType.ForLoop)`

A decorator factory for the definition of the looping strategy.

This decorator factory must be applied to Processor subclasses to modify their value of `loop_type` in order to change the execution workflow.

See `single_loop()`, `for_loop()` and `while_loop()` decorator shortcuts.

Parameters

loop_type (*LoopType* | *str*, *Optional*) – The type of execution workflow requested for the decorated class. Defaults to `LoopType.ForLoop`.

`mafw.decorators.for_loop(cls)`

A decorator shortcut to define a for loop execution processor.

`mafw.decorators.orphan_protector(cls)`

A class decorator to modify the init method of a Processor so that the `remove_orphan_files` is set to `False` and no orphan files will be removed.

```
mafw.decorators.processor_depends_on_optional(module_name: str, raise_ex: bool = False, warn:  
                                             bool = True) → Callable[[Type[P]],  
                                             Type[Processor]]
```

Class decorator factory to check if module *module_name* is available.

It checks if all the optional modules listed in *module_name* separated by a ‘;’ can be found.

If all modules are found, then the class is returned as it is.

If at least one module is not found:

- and *raise_ex* is True, an ImportError exception is raised and the user is responsible to deal with it.
- if *raise_ex* is False, instead of returning the class, the *Processor* is returned.
- depending on the value of *warn*, the user will be informed with a warning message or not.

Typical usage

The user should decorate Processor subclasses everytime the optional module is required in their `__init__` method. Should the check on the optional module have a positive outcome, then the Processor subclass is returned. Otherwise, if *raise_ex* is False, an instance of the base *Processor* is returned. In this way, the returned class can still be executed without breaking the execution scheme but of course, without producing any output.

Should be possible to run the `__init__` method of the class without the missing library, then the user can also follow the approach described in this other [example](#).

Parameters

- **module_name** (*str*) – The optional module(s) from which the class depends on. A “;” separated list of modules can also be provided.
- **raise_ex** (*bool*, *Optional*) – Flag to raise an exception if *module_name* not found, defaults to False.
- **warn** (*bool*, *Optional*) – Flag to display a warning message if *module_name* is not found, defaults to True.

Returns

The wrapped processor.

Return type

type(*Processor*)

Raises

ImportError – if *module_name* is not found and *raise_ex* is True.

```
mafw.decorators.single_loop(cls)
```

A decorator shortcut to define a single execution processor.

```
mafw.decorators.singleton(cls)
```

Make a class a Singleton class (only one instance)

```
mafw.decorators.suppress_warnings(func: F) → F
```

Decorator to suppress warnings during the execution of a test function.

This decorator uses the `warnings.catch_warnings()` context manager to temporarily change the warning filter to ignore all warnings. It is useful when you want to run a test without having warnings clutter the output.

Usage:

```
@suppress_warnings  
def test_function():  
    # Your test code that might emit warnings
```

Parameters

func (*Callable*) – The test function to be decorated.

Returns

The wrapped function with suppressed warnings.

Return type

Callable

`mafwr.decorators.while_loop(cls)`

A decorator shortcut to define a while loop execution processor.

11.1.4 mafw.enumerators

Module provides a set of enumerators for dealing with standard tasks.

Classes

<i>LoopType</i> (*values)	The loop strategy for the processor.
<i>LoopingStatus</i> (*values)	Enumerator to modify the looping cycle.
<i>ProcessorExitStatus</i> (*values)	The processor exit status enumerator class
<i>ProcessorStatus</i> (*values)	Enumerator to describe the status of a processor.

class `mafwr.enumerators.LoopType(*values)`

Bases: `StrEnum`

The loop strategy for the processor.

Each processor can be executed in one of the following modes:

1. **Single mode.** The process method is executed only once.
2. **For loop mode.** The process method is executed inside a for loop after the start and before the finish. The loop is based on a list of elements, the user **must** overload the `get_items()` method to define the list of items for the loop.
3. **While loop mode.** The process method is executed inside a while loop after the start and before the finish. The user **must** overload the `while_condition()` to define when to stop the loop.

i Future development

Implement concurrent loop. Depending on the development of the [free-threading capabilities](#) of future python releases, this concurrent looping strategy might be based on threads or a porting of the autorad multi-processor approach.

static `_generate_next_value_(name, start, count, last_values)`

Return the lower-cased version of the member name.

ForLoop = `'for_loop'`

Value for the for loop on item list execution.

SingleLoop = `'single'`

Value for the single mode execution.

WhileLoop = `'while_loop'`

Value for the while loop execution.

class mafw.enumerators.**LoopingStatus**(*values)

Bases: IntEnum

Enumerator to modify the looping cycle.

In the case of a looping Processor, the user has the ability to slightly modify the looping structure using this enumerator.

In the *process()* the user can set the variable *looping_status* to one of the following values:

- *LoopingStatus.Continue*. It means that everything is working well and the loop cycle must go ahead as foreseen and the *accept_item()* will be invoked.
- *LoopingStatus.Skip*. The *skip_item()* will be called soon after the *process()* is finished. The status will be reset to Continue and the next item will be processed.
- *LoopingStatus.Abort*. The cycle is broken immediately.
- *LoopingStatus.Quit*. The cycle is broken immediately.

The last two options are apparently identical, but they offer the possibility to implement a different behaviour in the *finish()* method. When abort is used, then the *AbortProcessorException* will be raised. For example, the user can decide to rollback all changes if an abort as occurred or to save what done so far in case of a quit.

Abort = 3

Break the loop and force the outside container (*mafw.processor.ProcessorList*) to quit.

Continue = 1

The loop can continue

Quit = 4

Break the loop but let the outside container (*mafw.processor.ProcessorList*) to continue.

Skip = 2

Skip this item.

class mafw.enumerators.**ProcessorExitStatus**(*values)

Bases: IntEnum

The processor exit status enumerator class

- Successful: means that the processor reached the end with success
- Failed: means that the processor did not reach the end with success
- Aborted: means that the user aborted the processor execution

Aborted = 3

The processor execution was aborted by the user

Failed = 2

The processor execution was failed

Successful = 1

The processor execution was successfully concluded

class mafw.enumerators.**ProcessorStatus**(*values)

Bases: StrEnum

Enumerator to describe the status of a processor.

static *_generate_next_value_(name, start, count, last_values)*

Return the lower-cased version of the member name.

Finish = 'finishing'

Finished

```

Init = 'initializing'
    Initialized
Run = 'processing'
    Running
Start = 'starting'
    Started
Unknown = 'unknown'
    Unknown status

```

11.1.5 mafw.examples

A library of examples.

Modules

<i>db_processors</i>	The module provides some basic processors with DB interaction for demonstrating the basic functionalities.
<i>importer_example</i>	The module provides one concrete implementation of an Importer, as it was used in the autorad paper2.
<i>loop_modifier</i>	The module provides examples on how the user can change the looping structure in a looping processor using the looping status.
<i>multi_primary</i>	Module demonstrates how to use multicolumn primary keys and foreign keys.
<i>processor_list</i>	The module provides some examples on how to use ProcessorList to combine several processors.
<i>sum_processor</i>	The module provides some examples for the user to develop their own processors.

mafw.examples.db_processors

The module provides some basic processors with DB interaction for demonstrating the basic functionalities.

Classes

<i>CountStandardTables</i> (*args, **kwargs)	A processor to count the number of standard tables
<i>File</i> (*args, **kwargs)	The Model class representing the table in the database
<i>FillFileTableProcessor</i> (*args, **kwargs)	Processor to fill a table with the content of a directory

Exceptions

<i>FileDoesNotExist</i>	Exception raised if the table corresponding to the File model does not exist.
-------------------------	-------------------------------------------------------------------------------

exception mafw.examples.db_processors.**FileDoesNotExist**

Bases: **Exception**

Exception raised if the table corresponding to the File model does not exist.

class mafw.examples.db_processors.**CountStandardTables**(*args: Any, **kwargs: Any)

Bases: *Processor*

A processor to count the number of standard tables

Processor parameters

- **n_tables**: The number of standard tables (default: -1)

Constructor parameters

Parameters

- **name** (*str*, *Optional*) – The name of the processor. If None is provided, the class name is used instead. Defaults to None.
- **description** (*str*, *Optional*) – A short description of the processor task. Defaults to the processor name.
- **config** (*dict*, *Optional*) – A configuration dictionary for this processor. Defaults to None.
- **looper** (*LoopType*, *Optional*) – Enumerator to define the looping type. Defaults to `LoopType.ForLoop`
- **user_interface** (*UserInterfaceBase*, *Optional*) – A user interface instance to be used by the processor to interact with the user.
- **timer** (*Timer*, *Optional*) – A timer object to measure process duration.
- **timer_params** (*dict*, *Optional*) – Parameters for the timer object.
- **database** (*Database*, *Optional*) – A database instance. Defaults to None.
- **database_conf** (*dict*, *Optional*) – Configuration for the database. Default to None.
- **remove_orphan_files** (*bool*, *Optional*) – Boolean flag to remove files on disc without a reference to the database. See *Standard tables* and *_remove_orphan_files()*. Defaults to True
- **kwargs** – Keyword arguments that can be used to set processor parameters.

process()

Processes the current item.

This is the core of the Processor, where the user has to define the calculations required.

start() → None

Start method.

The user can overload this method, including all steps that should be performed at the beginning of the operation.

If the user decides to overload it, it should include a call to the super method.

class `maf.w.examples.db_processors.File(*args, **kwargs)`

Bases: *MAFwBaseModel*

The Model class representing the table in the database

DoesNotExist

alias of *FileDoesNotExist*

class `maf.w.examples.db_processors.FillFileTableProcessor(*args: Any, **kwargs: Any)`

Bases: *Processor*

Processor to fill a table with the content of a directory

Processor parameters

- **root_folder**: The root folder for the file listing (default: `PosixPath('/tmp/maf-w-docs-_mot&xgc/v1.2.0_pdf')`)

Constructor parameter:

Parameters

root_folder (*Path, Optional*) – ActiveParameter corresponding to the directory from where to start the recursive search

finish()

Transfers all the data to the File table via an atomic transaction.

format_progress_message()

Customizes the progress message with information about the current item.

The user can overload this method in order to modify the message being displayed during the process loop with information about the current item.

The user can access the current value, its position in the looping cycle and the total number of items using *Processor.item*, *Processor.i_item* and *Processor.n_item*.

get_items() → list[Path]

Retrieves the list of files.

Insert or update the files from the root folder to the database

Returns

The list of full filename

Return type

list[Path]

process()

Add all information to the data list

start()

Starts the execution.

Be sure that the table corresponding to the File model exists. If it does already exist, it is not a problem.

mafw.examples.importer_example

The module provides one concrete implementation of an Importer, as it was used in the autorad paper2.

Classes

<i>ImporterExample</i> (*args, **kwargs)	An exemplary implementation of an importer processor.
<i>InputElement</i> (*args, **kwargs)	A model to store the input elements

Exceptions

<i>InputElementDoesNotExist</i>	Exception raised if the InputElement does not exist
---------------------------------	-----------------------------------------------------

exception mafw.examples.importer_example.**InputElementDoesNotExist**

Bases: *DoesNotExist*

Exception raised if the InputElement does not exist

class mafw.examples.importer_example.**ImporterExample**(*args: Any, **kwargs: Any)

Bases: *Importer*

An exemplary implementation of an importer processor.

This importer subclass is looking for tif files in the `input_folder` and using the information stored in the filename, all required database fields will be obtained.

For this importer, we will use a filename parser including two compulsory elements and one optional ones. Those are:

- Sample name, in the form of `sample_xyz`, where `xyz` are three numerical digits (**compulsory**).
- Exposure time, in the form of `expYh`, where `Y` is the exposure time in hours. It can be an integer number or a floating number using `.` as decimal separator (**compulsory**).
- Resolution, in the form of `rXYu`, where `XY` is the readout granularity in micrometer. It is optional and its default is 25 μm if not provided.

This processor is subclassing:

1. The `start()`: in order to be sure that the database table is existing.
2. The `get_items()`: where the list of input files is retrieved. There we also verify that the table is still updated using the `verify_checksum()` and in case the user wants to process only new files, we will have to filter out from the list all items already stored in the database.
3. The `process()`: where we create a dictionary with the values to be stored in the database. This approach allows a much more efficient database transaction.
4. The `finish()`: where we actually do the database insert in a single go.
5. The `format_progress_message()`: to keep the user informed about the progress (optional).

Processor parameters

- **input_folder**: The input folder from where the images have to be imported. (default: `‘/tmp/mafwdocs-_mot8xgc/v1.2.0_pdf’`)
- **parser_configuration**: The path to the TOML file with the filename parser configuration (default: `‘parser_configuration.toml’`)
- **recursive**: Extend the search to sub-folder (default: `True`)

Constructor parameters

Parameters

- **name** (*str*, *Optional*) – The name of the processor. If `None` is provided, the class name is used instead. Defaults to `None`.
- **description** (*str*, *Optional*) – A short description of the processor task. Defaults to the processor name.
- **config** (*dict*, *Optional*) – A configuration dictionary for this processor. Defaults to `None`.
- **looper** (*LoopType*, *Optional*) – Enumerator to define the looping type. Defaults to `LoopType.ForLoop`
- **user_interface** (*UserInterfaceBase*, *Optional*) – A user interface instance to be used by the processor to interact with the user.
- **timer** (*Timer*, *Optional*) – A timer object to measure process duration.
- **timer_params** (*dict*, *Optional*) – Parameters for the timer object.
- **database** (*Database*, *Optional*) – A database instance. Defaults to `None`.
- **database_conf** (*dict*, *Optional*) – Configuration for the database. Default to `None`.
- **remove_orphan_files** (*bool*, *Optional*) – Boolean flag to remove files on disc without a reference to the database. See *Standard tables* and `_remove_orphan_files()`. Defaults to `True`
- **kwargs** – Keyword arguments that can be used to set processor parameters.

finish() → None

The finish method overload.

Here is where we do the database insert with a `on_conflict_replace` to cope with the unique constraint.

format_progress_message() → None

Customizes the progress message with information about the current item.

The user can overload this method in order to modify the message being displayed during the process loop with information about the current item.

The user can access the current value, its position in the looping cycle and the total number of items using `Processor.item`, `Processor.i_item` and `Processor.n_item`.

get_items() → Collection[Any]

Retrieves the list of element to be imported.

The base folder is provided in the configuration file, along with the recursive flags and all the filter options.

Returns

The list of items full file names to be processed.

Return type

list[Path]

process()

The process method overload.

This is where the whole list of files is scanned.

The current item is a filename, so we can feed it directly to the `FilenameParser` interpret command, to have it parsed. To maximise the efficiency of the database transaction, instead of inserting each file singularly, we are collecting them all in a list and then insert all of them in the `finish()` method.

In case the parsing is failing, then the element is skipped and an error message is printed.

start() → None

The start method.

The filename parser is ready to use because it has been already configured in the super method. We need to be sure that the input table exists, otherwise we create it from scratch.

class `mafw.examples.importer_example.InputElement(*args, **kwargs)`

Bases: `MAFwBaseModel`

A model to store the input elements

DoesNotExist

alias of `InputElementDoesNotExist`

mafw.examples.loop_modifier

The module provides examples on how the user can change the looping structure in a looping processor using the looping status.

Functions

`is_prime(n)`

Check if n is a prime number.

Classes

<code>FindNPrimeNumber(*args, **kwargs)</code>	An example of Processor to search for N prime numbers starting from a given starting integer.
<code>FindPrimeNumberInRange(*args, **kwargs)</code>	An example processor to find prime numbers in the defined interval from <code>start_from</code> to <code>stop_at</code> .
<code>ModifyLoopProcessor(*args, **kwargs)</code>	Example processor demonstrating how it is possible to change the looping structure.

```
class mafw.examples.loop_modifier.FindNPrimeNumber(*args: Any, **kwargs: Any)
```

Bases: *Processor*

An example of Processor to search for N prime numbers starting from a given starting integer.

This processor is meant to demonstrate the use of a `while_loop` execution workflow.

Let us say we need to find 1000 prime numbers starting from 12347. One possible brute force approach to solve this problem is to start checking if the initial value is a prime number. If this is not the case, then check the next odd number. If it is the case, then add the current number to the list of found prime numbers and continue until the size of this list is 1000.

This is a perfect application for a while loop execution workflow.

Processor parameters

- **prime_num_to_find**: How many prime number we have to find (default: 100)
- **start_from**: From which number to start the search (default: 50)

Processor parameters:

Parameters

- **prime_num_to_find** (*int*) – The number of prime numbers to be found.
- **start_from** (*int*) – The initial integer number from where to start the search.

finish() → None

Overload of the finish method.

Remember: The finish method is called only once just after the last loop interaction. Always put a call to its *super* when you overload finish.

The loop is over, it means that the while condition was returning false, and now we can do something with our list of prime numbers.

format_progress_message() → None

Customizes the progress message with information about the current item.

The user can overload this method in order to modify the message being displayed during the process loop with information about the current item.

The user can access the current value, its position in the looping cycle and the total number of items using *Processor.item*, *Processor.i_item* and *Processor.n_item*.

process() → None

The overload of the process method.

Remember: The process method is called inside the while loop. It has access to the looping parameters: *Processor.i_item*, *Processor.item* and *Processor.n_item*.

In our specific case, the process contains another while loop. We start by checking if the current *Processor.item* is a prime number or not. If so, then we have found the next prime number, we add it to the list, we increment by two units the value of *Processor.item* and we leave the process method ready for the next iteration.

If *Processor.item* is not prime, then increment it by 2 and check it again.

start() → None

The overload of the start method.

Remember: The start method is called just before the while loop is started. So all instructions in this method will be executed only once at the beginning of the process execution. Always put a call to its *super* when you overload start.

First, we empty the list of found prime numbers. It should not be necessary, but it makes the code more readable. Then set the *Processor.n_item* to the total number of prime numbers we need to find. In this way, the progress bar will display useful progress.

If the start value is smaller than 2, then let's add 2 to the list of found prime number and set our first item to check at 3. In principle, we could already add 3 as well, but maybe the user wanted to find only 1 prime number, and we are returning a list with two, that is not what he was expecting.

Since prime numbers different from 2 can only be odd, if the starting number is even, increment it already by 1 unit.

while_condition() → bool

Define the while condition.

First, it checks if the *prime_num_to_find* is positive. Otherwise, it does not make sense to start. Then it will check if the length of the list with the already found prime numbers is enough. If so, then we can stop the loop return False, otherwise, it will return True and continue the loop.

Differently from the *for_loop* execution, we are responsible to assign the value to the looping variables *Processor.i_item*, *Processor.item* and *Processor.n_item*.

In this case, we will use the *Processor.i_item* to count how many prime numbers we have found and *Processor.n_item* will be our target. In this way, the progress bar will work as expected.

In the while condition, we set the *Processor.i_item* to the current length of the found prime number list.

Returns

True if the loop has to continue, False otherwise

prime_num_found: list[int]

The list with the found prime numbers

class mafw.examples.loop_modifier.FindPrimeNumberInRange(*args: Any, **kwargs: Any)

Bases: *Processor*

An example processor to find prime numbers in the defined interval from *start_from* to *stop_at*.

This processor is meant to demonstrate the use of a *for_loop* execution workflow.

Let us say we want to select only the prime numbers in a user defined range. One possible brute force approach is to generate the list of integers between the range extremes and check if it is prime or not. If yes, then add it to the list of prime numbers, if not continue with the next element.

This is a perfect application for a loop execution workflow.

Processor parameters

- **start_from:** From which number to start the search (default: 50)
- **stop_at:** At which number to stop the search (default: 100)

Processor parameters:

Parameters

- **start_from** (*int*) – First element of the range under investigation.
- **stop_at** (*int*) – Last element of the range under investigation.

finish() → None

Overload of the finish method.

Remember: to call the super method when you overload the finish method.

In this case, we just print out some information about the prime number found in the range.

format_progress_message() → None

Customizes the progress message with information about the current item.

The user can overload this method in order to modify the message being displayed during the process loop with information about the current item.

The user can access the current value, its position in the looping cycle and the total number of items using *Processor.item*, *Processor.i_item* and *Processor.n_item*.

get_items() → Collection[Any]

Overload of the get_items method.

This method must be overloaded when you select a for loop workflow.

Here we generate the list of odd numbers between the start and stop that we need to check. We also check that the stop is actually larger than the start, otherwise we print an error message, and we return an empty list of items.

Returns

A list of odd integer numbers between start_from and stop_at.

Return type

list[int]

process() → None

The process method.

In this case, it is very simple. We check if *Processor.item* is a prime number, if so we added to the list, otherwise we let the loop continue.

start() → None

Overload of the start method.

Remember: to call the super method when you overload the start.

In this specific case, we just make sure that the list of found prime numbers is empty.

prime_num_found: list[int]

The list with the found prime numbers

class mafw.examples.loop_modifier.ModifyLoopProcessor(*args: Any, **kwargs: Any)

Bases: *Processor*

Example processor demonstrating how it is possible to change the looping structure.

It is a looping processor where some events will be skipped, and at some point one event will trigger an abort.

Processor parameters

- **item_to_abort:** Item to abort (default: 65)
- **items_to_skip:** List of items to be skipped. (default: [12, 16, 25])
- **total_item:** Total item in the loop. (default: 100)

Processor Parameters:

Parameters

- **total_item** (*int*) – The total number of items
- **items_to_skip** (*list[int]*) – A list of items to skip.

- **item_to_abort** (*int*) – The item where to trigger an abort.

get_items() → list[int]

Returns the list of items, the range from 0 to total_item.

process()

Processes the item

skip_item()

Add skipped item to the skipped item list.

start()

Resets the skipped item container.

skipped_items: [list[int]]

A list with the skipped items.

`maf.w.examples.loop_modifier.is_prime(n: int)` → bool

Check if *n* is a prime number.

Parameters

n (*int*) – The integer number to be checked.

Returns

True if *n* is a prime number. False, otherwise.

Return type

bool

maf.w.examples.multi_primary

Module demonstrates how to use multicolumn primary keys and foreign keys.

maf.w.examples.processor_list

The module provides some examples on how to use ProcessorList to combine several processors.

Functions

<code>run_processor_list_with_loop_modifier()</code>	Example on deal with processors inside a processor list changing the loop structure.
<code>run_simple_processor_list()</code>	Simplest way to run several processors in a go.

`maf.w.examples.processor_list.run_processor_list_with_loop_modifier()`

Example on deal with processors inside a processor list changing the loop structure.

In this example there are two processors, one that will run until the end and the other that will set the looping status to abort half way. The user can see what happens when the *ProcessorList* is executed.

`maf.w.examples.processor_list.run_simple_processor_list()`

Simplest way to run several processors in a go.

maf.w.examples.sum_processor

The module provides some examples for the user to develop their own processors.

Those implemented here are mainly used in the test suit.

Classes

<code>AccumulatorProcessor(*args, **kwargs)</code>	A processor to calculate the sum of the first n values via a looping approach.
<code>GaussAdder(*args, **kwargs)</code>	A processor to calculate the sum of the first n values via the so called <i>Gauss formula</i> .

class mafw.examples.sum_processor.**AccumulatorProcessor**(*args: Any, **kwargs: Any)

Bases: *Processor*

A processor to calculate the sum of the first n values via a looping approach.

In mathematical terms, this processor solves this easy equation:

$$N = \sum_{i=0}^n i$$

by looping. It is a terribly inefficient approach, but it works as a demonstration of the looping structure.

The user can get the results by retrieving the *accumulated_value* parameter at the end of the processor execution.

Processor parameters

- **last_value**: Last value of the series (default: 100)

Constructor parameters:

Parameters

- **last_value** (*int*) – The *n* in the equation above. Defaults to 100
- **accumulated_value** (*int*) – The *N* in the equation above at the end of the process.

get_items() → list[int]

Returns the list of the first *last_value* integers.

process()

Increase the accumulated value by the current item.

start()

Resets the accumulated value to 0 before starting.

class mafw.examples.sum_processor.**GaussAdder**(*args: Any, **kwargs: Any)

Bases: *Processor*

A processor to calculate the sum of the first n values via the so called *Gauss formula*.

In mathematical terms, this processor solves this easy equation:

$$N = \frac{n * (n - 1)}{2}$$

without any looping

The user can get the results by retrieving the *sum_value* parameter at the end of the processor execution.

Processor parameters

- **last_value**: Last value of the series. (default: 100)

Constructor parameters:

Parameters

- **last_value** (*int*) – The *n* in the equation above. Defaults to 100
- **sum_value** (*int*) – The *N* in the equation above.

process()

Compute the sum using the Gauss formula.

start()

Sets the sum value to 0.

11.1.6 mafw.hookspecs

Defines the hook specification decorator bound the MAFw library.

Functions

<i>register_processors()</i>	Register multiple processor classes
<i>register_standard_tables()</i>	Register standard tables
<i>register_user_interfaces()</i>	Register multiple user interfaces

`maf.w.hookspecs.register_processors()` → list[*Processor*]

Register multiple processor classes

`maf.w.hookspecs.register_standard_tables()` → list[*StandardTable*]

Register standard tables

`maf.w.hookspecs.register_user_interfaces()` → list[*UserInterfaceBase*]

Register multiple user interfaces

11.1.7 mafw.mafw_errors

Module defines MAFw exceptions

Exceptions

<i>AbortProcessorException</i>	Exception raised during the execution of a processor requiring immediate exit.
<i>InvalidSteeringFile</i>	Exception raised when validating an invalid steering file
<i>MAFwException</i>	Base class for MAFwException
<i>MissingAttribute</i>	Exception raised when an attempt is made to execute a statement without a required parameter/attributes
<i>MissingDatabase</i>	Exception raised when a processor requiring a database connection is being operated without a database
<i>MissingOptionalDependency</i>	UserWarning raised when an optional dependency is required
<i>MissingOverloadedMethod</i>	Warning issued when the user did not overload a required method.
<i>MissingSQLStatement</i>	Exception raised when a Trigger is created without any SQL statements.
<i>MissingSuperCall</i>	Warning issued when the user did not invoke the super method for some specific processor methods.
<i>ModelError</i>	Exception raised when an error in a DB Model class occurs
<i>ParserConfigurationError</i>	Exception raised when an error occurred during the configuration of a filename parser
<i>ParsingError</i>	Exception raised when a regular expression parsing failed

continues on next page

Table 11.34 – continued from previous page

<i>PlotterMixinNotInitialized</i>	Exception raised when a plotter mixin has not properly initialized
<i>ProcessorParameterError</i>	Error with a processor parameter
<i>RunnerNotInitialized</i>	Exception raised when attempting to run a not initialized Runner.
<i>UnknownDBEngine</i>	Exception raised when the user provided an unknown db engine
<i>UnknownProcessor</i>	Exception raised when an attempt is made to create an unknown processor
<i>UnsupportedDatabaseError</i>	Error raised when a feature is not supported by the database.

exception mafw.mafw_errors.**AbortProcessorException**Bases: *MAFwException*

Exception raised during the execution of a processor requiring immediate exit.

exception mafw.mafw_errors.**InvalidSteeringFile**Bases: *MAFwException*

Exception raised when validating an invalid steering file

exception mafw.mafw_errors.**MAFwException**

Bases: Exception

Base class for MAFwException

exception mafw.mafw_errors.**MissingAttribute**Bases: *MAFwException*

Exception raised when an attempt is made to execute a statement without a required parameter/attributes

exception mafw.mafw_errors.**MissingDatabase**Bases: *MAFwException*

Exception raised when a processor requiring a database connection is being operated without a database

exception mafw.mafw_errors.**MissingOptionalDependency**

Bases: UserWarning

UserWarning raised when an optional dependency is required

exception mafw.mafw_errors.**MissingOverloadedMethod**

Bases: UserWarning

Warning issued when the user did not overload a required method.

It is a warning and not an error because the execution framework might still work, but the results might be different from what is expected.

exception mafw.mafw_errors.**MissingSQLStatement**Bases: *MAFwException*

Exception raised when a Trigger is created without any SQL statements.

exception mafw.mafw_errors.**MissingSuperCall**

Bases: UserWarning

Warning issued when the user did not invoke the super method for some specific processor methods.

Those methods (like *start()* and *finish()*) have not empty implementation also in the base class, meaning that if the user forgets to call *super* in their overloads, then the basic implementation will be gone.

It is a warning and not an error because the execution framework might still work, but the results might be different from what is expected.

exception mafw.mafw_errors.**ModelError**Bases: *MAFWException*

Exception raised when an error in a DB Model class occurs

exception mafw.mafw_errors.**ParserConfigurationError**Bases: *MAFWException*

Exception raised when an error occurred during the configuration of a filename parser

exception mafw.mafw_errors.**ParsingError**Bases: *MAFWException*

Exception raised when a regular expression parsing failed

exception mafw.mafw_errors.**PlotterMixinNotInitialized**Bases: *MAFWException*

Exception raised when a plotter mixin has not properly initialized

exception mafw.mafw_errors.**ProcessorParameterError**Bases: *MAFWException*

Error with a processor parameter

exception mafw.mafw_errors.**RunnerNotInitialized**Bases: *MAFWException*

Exception raised when attempting to run a not initialized Runner.

exception mafw.mafw_errors.**UnknownDBEngine**Bases: *MAFWException*

Exception raised when the user provided an unknown db engine

exception mafw.mafw_errors.**UnknownProcessor**Bases: *MAFWException*

Exception raised when an attempt is made to create an unknown processor

exception mafw.mafw_errors.**UnsupportedDatabaseError**Bases: *Exception*

Error raised when a feature is not supported by the database.

11.1.8 mafw.plugin_manager

Provides utilities to retrieve internal and external plugins.

Functions

<code>get_plugin_manager([force_recreate])</code>	Create a new or return an existing plugin manager for a given project
---------------------------------------------------	-----------------------------------------------------------------------

`mafw.plugin_manager.get_plugin_manager(force_recreate: bool = False) → PluginManager`

Create a new or return an existing plugin manager for a given project

Parameters**force_recreate** (*bool*, *Optional*) – Flag to force the creation of a new plugin manager. Defaults to False**Returns**

The plugin manager

Return type
pluggy.PluginManager

11.1.9 mafw.plugins

Exports Processor classes to the execution script.

Functions

<code>register_processors()</code>	Returns a list of processors to be registered
<code>register_standard_tables()</code>	Returns a list of Models to be used as Standard Tables in MAFw
<code>register_user_interfaces()</code>	Returns a list of user interfaces that can be used

`mafw.plugins.register_processors()` → list[type[*Processor*]]

Returns a list of processors to be registered

`mafw.plugins.register_standard_tables()` → list[type[*StandardTable*]]

Returns a list of Models to be used as Standard Tables in MAFw

`mafw.plugins.register_user_interfaces()` → list[type[*UserInterfaceBase*]]

Returns a list of user interfaces that can be used

11.1.10 mafw.processor

Module implements the basic Processor class, the ProcessorList and all helper classes to achieve the core functionality of the MAFw.

Module Attributes

<i>ParameterType</i>	Generic variable type for the <i>ActiveParameter</i> and <i>PassiveParameter</i> .
<i>F</i>	Type variable for generic callable with any return value.

Functions

<code>ensure_parameter_registration(func)</code>	Decorator to ensure that before calling <i>func</i> the processor parameters have been registered.
<code>validate_database_conf([database_conf])</code>	Validates the database configuration.

Classes

<code>ActiveParameter(name[, value, default, help_doc])</code>	The public interface to the processor parameter.
<code>PassiveParameter(name[, value, default, ...])</code>	A processor parameter that can be registered and configured.
<code>Processor(*args, **kwargs)</code>	The basic processor.
<code>ProcessorList(*args[, name, description, ...])</code>	A list like collection of processors.
<code>ProcessorMeta</code>	A metaclass to implement the post-init method.

```
class mafw.processor.ActiveParameter(name: str, value: ParameterType | None = None, default:
    ParameterType | None = None, help_doc: str = "")
```

Bases: `Generic[ParameterType]`

The public interface to the processor parameter.

The behaviour of a *Processor* can be customized by using processor parameters. The value of these parameters can be either set via a configuration file or directly when creating the class.

If the user wants to benefit from this facility, they have to add in the instance of the Processor subclass an *ActiveParameter* instance in this way:

```
class MyProcessor(Processor):
    # this is the input folder
    input_folder = ActiveParameter('input_folder', Path(r'C:\'), help_doc='This_
→is where to look for input files')

    def __init__(self, *args, **kwargs):
        super().__init__(*args, **kwargs)

    # change the input folder to something else
    self.input_folder = Path(r'D:\data')

    # get the value of the parameter
    print(self.input_folder)
```

The *ActiveParameter* is a *descriptor*, it means that when you create one of them, a lot of work is done behind the scene.

In simple words, a processor parameter is made by two objects: a public interface where the user can easily access the value of the parameter and a private interface where all other information (default, documentation...) is also stored.

The user does not have to take care of all of this. When a new *ActiveParameter* instance is added to the class as in the code snippet above, the private interface is automatically created and will stay in the class instance until the end of the class lifetime.

To access the private interface, the user can use the *Processor.get_parameter()* method using the parameter name as a key.

➔ See also

The private counter part in the *PassiveParameter*.

Constructor parameters:

Parameters

- **name** (*str*) – The name of the parameter.
- **value** (*ParameterType*, *Optional*) – The initial value of the parameter. Defaults to `None`.
- **default** (*ParameterType*, *Optional*) – The default value of the parameter, to be used when `value` is not set., Defaults to `None`.
- **help_doc** (*str*, *Optional*) – An explanatory text describing the parameter.

```
class mafw.processor.F
```

Type variable for generic callable with any return value.

alias of `TypeVar('F', bound=Callable[[], Any])`

class mafw.processor.**ParameterType**

Generic variable type for the *ActiveParameter* and *PassiveParameter*.

alias of `TypeVar('ParameterType')`

class mafw.processor.**PassiveParameter**(*name: str, value: ParameterType | None = None, default: ParameterType | None = None, help_doc: str = ""*)

Bases: `Generic[ParameterType]`

A processor parameter that can be registered and configured.

For Processors to perform their analytical task, it may be necessary to have some configurable parameters, like a DB input table or the output folder or a numeric parameter.

The name of the parameter must be unique within the Processor scope and a valid python identifier.

When defined as a *PassiveParameter*, an instance variable can have a default value if the user did not provide one, and it is much easier to configure via a configuration file.

A parameter for which only a default value is provided is automatically considered *optional*.

If both the value and the default value are not provided, an exception is raised.

This class is working behind the scene, that is why it is named **passive**. The user will very likely add class instances of *ActiveParameter*, that are publicly exposed in the processor class namespace, and an *PassiveParameter* will automatically added to the class. To access this passive parameter the user can use the *Processor.get_parameter()* using the name as key. The *value* of the passive parameter is always accessible using calling the corresponding *ActiveParameter*.

See also

An explanation on how processor parameters work and should be used is given in *Understanding processor parameters*

Constructor parameters:

Parameters

- **name** (*str*) – The name of the parameter. It must be a valid python identifier.
- **value** (*ParameterType, Optional*) – The set value of the parameter. If *None*, then the default value will be used. Defaults to *None*.
- **default** (*ParameterType, Optional*) – The default value for the parameter. It is used if the *value* is not provided. Defaults to *None*.
- **help_doc** (*str, Optional*) – A brief explanation of the parameter.

Raises

ProcessorParameterError – if both *value* and *default* are not provided or if *name* is not a valid identifier.

property is_optional: bool

Property to check if the parameter is optional.

Returns

True if the parameter is optional

Return type

bool

property is_set: bool

Property to check if the value has been set.

It is useful for optional parameter to see if the current value is the default one, or if the user set it.

property value: *ParameterType*

Gets the parameter value.

Returns

The parameter value.

Return type

ParameterType

Raises

ProcessorParameterError – if both value and default were not defined.

class mafw.processor.Processor(*args: Any, **kwargs: Any)

Bases: object

The basic processor.

A very comprehensive description of what a Processor does and how it works is available at *Processor: The core of MAFw*.

Constructor parameters

Parameters

- **name** (*str*, *Optional*) – The name of the processor. If None is provided, the class name is used instead. Defaults to None.
- **description** (*str*, *Optional*) – A short description of the processor task. Defaults to the processor name.
- **config** (*dict*, *Optional*) – A configuration dictionary for this processor. Defaults to None.
- **looper** (*LoopType*, *Optional*) – Enumerator to define the looping type. Defaults to *LoopType.ForLoop*
- **user_interface** (*UserInterfaceBase*, *Optional*) – A user interface instance to be used by the processor to interact with the user.
- **timer** (*Timer*, *Optional*) – A timer object to measure process duration.
- **timer_params** (*dict*, *Optional*) – Parameters for the timer object.
- **database** (*Database*, *Optional*) – A database instance. Defaults to None.
- **database_conf** (*dict*, *Optional*) – Configuration for the database. Default to None.
- **remove_orphan_files** (*bool*, *Optional*) – Boolean flag to remove files on disc without a reference to the database. See *Standard tables* and *_remove_orphan_files()*. Defaults to True
- **kwargs** – Keyword arguments that can be used to set processor parameters.

_check_method_overload() → None

Check if the user overloaded the required methods.

Depending on the loop type, the user must overload different methods. This method is doing the check and if the required methods are not overloaded a warning is emitted.

_check_method_super() → None

Check if some specific methods are calling their super.

For some specific methods (for example: start and finish), the user should always call their super method. This method verifies that the user implementation of these methods is including a super call, otherwise a warning is emitted to inform the user about the problem and possible misbehaviour of the processor.

The list of methods to be verified is stored in a private class attribute `_methods_to_be_checked_for_super` as a list of tuples, made by the name of the methods to be verified and the base class for comparison. The base class is required because Processor subclasses may be extending this list with methods that are not present in the base Processor. See, for example, the `patch_data_frame()` that is required to have a super call, but it is not present in the base Processor.

`_execute_for_loop()` → None

Executes the processor within a for loop.

Private method. Do not overload nor invoke it directly. The `execute()` method will call the appropriate implementation depending on the processor `LoopType`.

`_execute_single()` → None

Execute the processor in single mode.

Private method. Do not overload nor invoke it directly. The `execute()` method will call the appropriate implementation depending on the processor `LoopType`.

`_execute_while_loop()` → None

Executes the processor within a while loop.

Private method. Do not overload nor invoke it directly. The `execute()` method will call the appropriate implementation depending on the processor `LoopType`.

`_remove_orphan_files()` → None

Remove orphan files.

If a connection to the database is available, then the `OrphanFile` standard table is queried for all its entries, and all the files are then removed.

The user can turn off this behaviour by switching the `remove_orphan_files` to `False`.

`accept_item()` → None

Does post process actions on a successfully processed item.

Within the `process()`, the user left the looping status to `Continue`, so it means that everything looks good and this is the right place to perform database updates or file savings.

`acquire_resources()` → None

Acquires resources and add them to the resource stack.

The whole body of the `execute()` method is within a context structure. The idea is that if any part of the code inside should throw an exception that breaking the execution, we want to be sure that all stateful resources are properly closed.

Since the number of resources may vary, the variable number of nested `with` statements has been replaced by an `ExitStack`. Resources, like open files, timers, db connections, need to be added to the resource stacks in this method.

In the case a processor is being executed within a `ProcessorList`, then some resources might be shared, and for this reason they are not added to the stack. This selection can be done via the private `local_resource_acquisition`. This is normally `True`, meaning that the processor will handle its resources independently, but when the processor is executed from a `ProcessorList`, this flag is automatically turned to `False`.

If the user wants to add additional resources, he has to overload this method calling the super to preserve the original resources. If he wants to have shared resources among different processors executed from inside a processor list, he has to overload the `ProcessorList` class as well.

`delete_parameter(name: str)` → None

Deletes a processor parameter.

Parameters

name (`str`) – The name of the parameter to be deleted.

Raises

ProcessorParameterError – If a parameter with *name* is not registered.

dump_parameter_configuration(*option: int = 1*) → dict[str, Any]

Dumps the processor parameter values in a dictionary.

The snippet below explains the meaning of *option*.

```
# option 1
conf_dict1 = {
    'Processor': {'param1': 5, 'input_table': 'my_table'}
}

# option 2
conf_dict2 = {'param1': 5, 'input_table': 'my_table'}
```

Parameters

option (*int, Optional*) – Select the dictionary style. Defaults to 1.

Returns

A parameter configuration dictionary.

Return type

dict

execute() → None

Execute the processor tasks.

This method works as a dispatcher, reassigning the call to a more specific execution implementation depending on the *loop_type*.

finish() → None

Concludes the execution.

The user can reimplement this method if there are some conclusive tasks that must be achieved. Always include a call to `super()`.

format_progress_message() → None

Customizes the progress message with information about the current item.

The user can overload this method in order to modify the message being displayed during the process loop with information about the current item.

The user can access the current value, its position in the looping cycle and the total number of items using *Processor.item*, *Processor.i_item* and *Processor.n_item*.

get_filter(*model_name: str*) → *Filter*

Returns a registered *Filter* via the model name.

If a filter for the provided *model_name* does not exist, a `KeyError` is raised.

Parameters

model_name (*str*) – The model name for which the filter will be returned.

Returns

The registered filter

Return type

maf.db.db_filter.Filter

Raises

`KeyError` is a filter with the give name is not found.

get_items() → Collection[Any]

Returns the item collections for the processor loop.

This method must be overloaded for the processor to work. Generally, this is getting a list of rows from the database, or a list of files from the disk to be processed.

Returns

A collection of items for the loop

Return type

Collection[Any]

get_parameter(name: str) → *PassiveParameter[ParameterType]*

Gets the processor parameter named name.

Parameters

name (str) – The name of the parameter.

Returns

The processor parameter

Return type

PassiveParameter

Raises

ProcessorParameterError – If a parameter with *name* is not registered.

get_parameters() → dict[str, *PassiveParameter[ParameterType]*]

Returns the full dictionary of registered parameters for this processor.

Useful when dumping the parameter specification in a configuration file, for example.

Returns

The dictionary with the registered parameters.

Return type

dict[str, *PassiveParameter[ParameterType]*]

on_looping_status_set(status: LoopingStatus) → None

Call back invoked when the looping status is set.

The user can overload this method according to the needs.

Parameters

status (*LoopingStatus*) – The set looping status.

on_processor_status_change(old_status: ProcessorStatus, new_status: ProcessorStatus) → None

Callback invoked when the processor status is changed.

Parameters

- **old_status** (*ProcessorStatus*) – The old processor status.
- **new_status** (*ProcessorStatus*) – The new processor status.

print_process_statistics() → None

Print the process statistics.

A utility method to display the fastest, the slowest and the average timing required to process on a single item. This is particularly useful when the looping processor is part of a ProcessorList.

process() → None

Processes the current item.

This is the core of the Processor, where the user has to define the calculations required.

set_parameter_value(*name*: str, *value*: ParameterType) → None

Sets the value of a processor parameter.

Parameters

- **name** (str) – The name of the parameter to be deleted.
- **value** (ParameterType) – The value to be assigned to the parameter.

Raises

ProcessorParameterError – If a parameter with *name* is not registered.

skip_item() → None

Does post process actions on a *NOT* successfully processed item.

Within the `process()`, the user set the looping status to Skip, so it means that something went wrong and here corrective actions can be taken if needed.

start() → None

Start method.

The user can overload this method, including all steps that should be performed at the beginning of the operation.

If the user decides to overload it, it should include a call to the super method.

while_condition() → bool

Return the while condition

Returns

True if the while loop has to continue, false otherwise.

Return type

bool

_ids = count(0)

A counter for all processor instances

_methods_to_be_checked_for_super

List of methods to be checked for super inclusion.

It is a list of tuple, with the first element the name of the method to be checked and the second the base class to be compared.

property database: Database

Returns the database instance

Returns

A database object.

Raises

MissingDatabase – If the database connection has not been established.

description

A short description of the processor task.

filter_register: FilterRegister

The DB filter register of the Processor.

property i_item: int

The enumeration of the current item being processed.

item: Any

The current item of the loop.

property local_resource_acquisition: bool

Checks if resources should be acquired locally.

When the processor is executed in stand-alone mode, it is responsible to acquire and release its own external resources, but when it is executed from a ProcessorList, then is a good practice to share and distribute resources among the whole processor list. In this case, resources should not be acquired locally by the single processor, but from the parent execution context.

Returns

True if resources are to be acquired locally by the processor. False, otherwise.

Return type

bool

loop_type: *LoopType*

The loop type.

The value of this parameter can also be changed by the *execution_workflow()* decorator factory.

See *LoopType* for more details.

looping_status

Looping modifier

property n_item: int | None

The total number of items to be processed or None for an undefined loop

name

The name of the processor.

processor_exit_status

Processor exit status

processor_status

Processor execution status

progress_message: str = 'Processor is working'

Message displayed to show the progress.

It can be customized with information about the current item in the loop by overloading the *format_progress_message()*.

remove_orphan_files: bool

The flag to remove or protect the orphan files. Defaults to True

unique_id

A unique identifier representing how many instances of Processor has been created.

property unique_name: str

Returns the unique name for the processor.

```
class mafw.processor.ProcessorList(*args: Processor | ProcessorList, name: str | None = None,
                                  description: str | None = None, timer: Timer | None = None,
                                  timer_params: dict[str, Any] | None = None, user_interface:
                                  UserInterfaceBase | None = None, database: Database | None =
                                  None, database_conf: dict[str, Any] | None = None)
```

Bases: list[Processor | ProcessorList]

A list like collection of processors.

ProcessorList is a subclass of list containing only Processor subclasses or other ProcessorList.

An attempt to add an element that is not a Processor or a ProcessorList will raise a TypeError.

Along with an iterable of processors, a new processor list can be built using the following parameters.

Constructor parameters:

Parameters

- **name** (*str*, *Optional*) – The name of the processor list. Defaults to ProcessorList.
- **description** (*str*, *Optional*) – An optional short description. Default to ProcessorList.
- **timer** (*Timer*, *Optional*) – The timer object. If None is provided, a new one will be created. Defaults to None.
- **timer_params** (*dict*, *Optional*) – A dictionary of parameter to build the timer object. Defaults to None.
- **user_interface** (*UserInterfaceBase*, *Optional*) – A user interface. Defaults to None
- **database** (*Database*, *Optional*) – A database instance. Defaults to None.
- **database_conf** (*dict*, *Optional*) – Configuration for the database. Default to None.

static validate_item(*item*: Processor | ProcessorList) → Processor | ProcessorList

Validates the item being added.

static validate_items(*items*: tuple[Processor | ProcessorList, ...] = ()) → tuple[Processor | ProcessorList, ...]

Validates a tuple of items being added.

acquire_resources() → None

Acquires external resources.

append(*_ProcessorList__object*: Processor | ProcessorList) → None

Appends a new processor at the end of the list.

distribute_resources(*processor*: Processor | Self) → None

Distributes the external resources to the items in the list.

execute() → ProcessorExitStatus

Execute the list of processors.

Similarly to the *Processor*, ProcessorList can be executed. In simple words, the execute method of each processor in the list is called exactly in the same sequence as they were added.

extend(*_ProcessorList__iterable*: Iterable[Processor | ProcessorList]) → None

Extends the processor list with a list of processors.

insert(*_ProcessorList__index*: SupportsIndex, *_ProcessorList__object*: Processor | ProcessorList) → None

Adds a new processor at the specified index.

property database: Database

Returns the database instance

Returns

A database instance

Raises

MissingDatabase – if a database connection is missing.

property name: str

The name of the processor list

Returns

The name of the processor list

Return type

str

property processor_exit_status: *ProcessorExitStatus*

The processor exit status.

It refers to the whole processor list execution.

class mafw.processor.ProcessorMeta

Bases: type

A metaclass to implement the post-init method.

mafw.processor.ensure_parameter_registration(*func*: F) → F

Decorator to ensure that before calling *func* the processor parameters have been registered.

mafw.processor.validate_database_conf(*database_conf*: dict[str, Any] | None = None) → dict[str, Any] | None

Validates the database configuration.

Parameters

database_conf (*dict*, *Optional*) – The input database configuration. Defaults to None.

Returns

Either the validated database configuration or None if it is invalid.

Return type

dict, None

11.1.11 mafw.processor_library

Library of generic processors.

Subpackage containing a library of general purpose processors that the user can either use directly or overload them to implement their analytical tasks.

Modules

<i>abstract_plotter</i>	Module implements the abstract base interface to a processor to generate plots.
<i>importer</i>	Provides a basic element importer.
<i>sns_plotter</i>	Module implements a Seaborn plotter processor with a mixin structure to generate seaborn plots.

mafw.processor_library.abstract_plotter

Module implements the abstract base interface to a processor to generate plots.

This abstract interface is needed because MAFw does not force the user to select a specific plot and data manipulation library.

The basic idea is to have a *basic processor class* featuring a modified *process()* method where a skeleton of the standard operations required to generate a graphical representation of a dataset is provided.

The user has the possibility to compose the *GenericPlotter* by mixing it with one *DataRetriever* and a *FigurePlotter*.

For a specific implementation based on *seaborn*, please refer to *sns_plotter*.

Classes

<i>DataRetriever</i> (*args, **kwargs)	Base mixin class to retrieve a data frame from an external source
----------------------------------------	-------------------------------------------------------------------

continues on next page

Table 11.41 – continued from previous page

<i>FigurePlotter</i> ()	
<i>GenericPlotter</i> (*args, **kwargs)	The Generic Plotter processor.
<i>PlotterMeta</i> (name, bases, namespace, /, **kwargs)	Metaclass for the plotter mixed classes

class mafw.processor_library.abstract_plotter.**DataRetriever**(*args: Any, **kwargs: Any)

Bases: ABC

Base mixin class to retrieve a data frame from an external source

The dataframe instance. It will be filled for the main class

abstractmethod `get_data_frame()` → None

The mixin implementation of the shared method with the base class

abstractmethod `patch_data_frame()` → None

The mixin implementation of the shared method with the base class

class mafw.processor_library.abstract_plotter.**FigurePlotter**

Bases: ABC

class mafw.processor_library.abstract_plotter.**GenericPlotter**(*args: Any, **kwargs: Any)

Bases: *Processor*

The Generic Plotter processor.

This is a subclass of a Processor with advanced functionality to fetch data in the form of a dataframe and to produce plots. When mentioning dataframe in the context of the generic plotter, we do not have in mind any specific dataframe implementation.

The GenericPlotter is actually a kind of abstract class: since MAFw is not forcing you to use any specific plotting and data manipulation library, you need to subclass the GenericPlotter in your code, be sure that the required dependencies are available for import and use it as a normal processor.

If you are ok with using *seaborn* (with *matplotlib* as a graphical backend and *pandas* for data storage and manipulation), then be sure to install mafw with the optional feature *seaborn* (`pip install mafw[seaborn]`) and have a look at the *sns_plotter* for an already prepared implementation of a Plotter.

The key difference with respect to a normal processor is its `process()` method that has been already implemented as follows:

```
def process(self) -> None:
    """
    Process method overload.

    In the case of a plotter subclass, the process method is already implemented,
    and the user should not overload it. On the contrary, the user must overload the other implementation methods,
    described in the general
    :class:`class description <.SNSPlotter>`.
    """
    if self.filter_register.new_only:
        if self.is_output_existing():
            return

    self.in_loop_customization()
    self.get_data_frame()
    self.patch_data_frame()
    self.slice_data_frame()
    self.group_and_aggregate_data_frame()
```

(continues on next page)

```
if not self.is_data_frame_empty():
    self.plot()
    self.customize_plot()
    self.save()
    self.update_db()
```

This actually means that when you are subclassing a `GenericPlotter` you do not have to implement the `process` method as you would do for a normal `Processor`, but you will have to implement the following methods:

- `in_loop_customization()`.

The processor execution workflow (`LoopType`) can be any of the available, so actually the `process` method might be invoked only once, or multiple times inside a loop structure (for or while). If the execution is cyclic, then you may want to have the possibility to do some customisation for each iteration, for example, changing the plot title, or modifying the data selection, or the filename where the plots will be saved.

You can use this method also in case of a single loop processor, in this case you will not have access to the loop parameters.

- `get_data_frame()`.

This method has the task to get the data to be plotted. Since it is an almost abstract class, you need to

- `patch_data_frame()`.

A convenient method to apply data frame manipulation to the data just retrieved. A typical use case is for conversion of unit of measurement. Imagine you saved the data in the S.I. units, but for the visualization you prefer to use practical units, so you can subclass this method to add a new column containing the same converted values of the original one.

- `slice_data_frame()`.

Slicing a dataframe is similar as applying a where clause in a SQL query. Implement this method to select which row should be used in the generation of your plot.

- `group_and_aggregate_data_frame()`.

In this method, you can manipulate your data frame to perform row grouping and aggregation.

- `is_data_frame_empty()`.

A simple method to test if the dataframe contains any data to be plotted. In fact, after the slicing, grouping and aggregation operations, it is possible that the dataframe is now left without any row. In this case, it makes no sense to waste time in plotting an empty graph.

- `plot()`.

This method is where the actual plotting occurs.

- `customize_plot()`.

This method can be optionally used to customize the appearance of the facet grid produced by the `plot()` method. It is particularly useful when the user is mixing this class with one of the `FigurePlotter` mixin, thus not having direct access to the plot method.

- `save()`.

This method is where the produced plot is saved in a file. Remember to append the output file name to the *list of produced outputs* so that the `_update_plotter_db()` method will automatically store this file in the database during the `finish()` execution.

- `update_db()`.

If the user wants to update a specific table in the database, they can use this method.

It is worth reminding that all plotters are saving all generated files in the standard table `PlotterOutput`. This is automatically done by the `_update_plotter_db()` method that is called in the `finish()` method.

Constructor parameters

Parameters

- **name** (*str*, *Optional*) – The name of the processor. If `None` is provided, the class name is used instead. Defaults to `None`.
- **description** (*str*, *Optional*) – A short description of the processor task. Defaults to the processor name.
- **config** (*dict*, *Optional*) – A configuration dictionary for this processor. Defaults to `None`.
- **looper** (`LoopType`, *Optional*) – Enumerator to define the looping type. Defaults to `LoopType.ForLoop`
- **user_interface** (`UserInterfaceBase`, *Optional*) – A user interface instance to be used by the processor to interact with the user.
- **timer** (`Timer`, *Optional*) – A timer object to measure process duration.
- **timer_params** (*dict*, *Optional*) – Parameters for the timer object.
- **database** (`Database`, *Optional*) – A database instance. Defaults to `None`.
- **database_conf** (*dict*, *Optional*) – Configuration for the database. Default to `None`.
- **remove_orphan_files** (*bool*, *Optional*) – Boolean flag to remove files on disc without a reference to the database. See *Standard tables* and `_remove_orphan_files()`. Defaults to `True`
- **kwargs** – Keyword arguments that can be used to set processor parameters.

_update_plotter_db() → `None`

Updates the Plotter DB.

A plotter subclass primarily generates plots as output in most cases, which means that no additional information needs to be stored in the database. This is sufficient to prevent unnecessary execution of the processor when it is not required.

This method is actually protected against execution without a valid database instance.

customize_plot() → `None`

The customize plot method.

The user can overload this method to customize the output produced by the `plot()` method, like, for example, adding meaningful axis titles, changing format, and so on.

As usual, it is possible to use the `item`, `i_item` and `n_item` to access the loop parameters.

finish() → `None`

Concludes the execution.

The user can reimplement this method if there are some conclusive tasks that must be achieved. Always include a call to `super()`.

format_progress_message() → `None`

Customizes the progress message with information about the current item.

The user can overload this method in order to modify the message being displayed during the process loop with information about the current item.

The user can access the current value, its position in the looping cycle and the total number of items using `Processor.item`, `Processor.i_item` and `Processor.n_item`.

get_data_frame() → None

Get the data frame with the data to be plotted.

This method can be either implemented in the SNSPlotter subclass or via a *DataRetriever* mixin class.

in_loop_customization() → None

Customize the parameters for the output or input data for each execution iteration.

is_data_frame_empty() → bool

Check if the data frame is empty

is_output_existing() → bool

Check for plotter output existence.

Generally, plotter subclasses do not have a real output that can be saved to a database. This class is meant to generate one or more graphical output files.

One of the biggest advantages of having the output of a processor stored in the database is the ability to conditionally execute the processor if, and only if, the output is missing or changed.

In order to allow also plotter processor to benefit from this feature, a *dedicated table* is available among the *standard tables*.

If a connection to the database is provided, then this method is invoked at the beginning of the *process()* and a select query over the *PlotterOutput* model is executed filtering by processor name. All files in the filename lists are checked for existence and also the checksum is verified.

This method will return True, if the output of the processor is already existing and valid, False, otherwise.

Returns

True if the processor output exists and it is valid.

Return type

bool

patch_data_frame() → None

Modify the data frame

This method can be used to perform operation on the data frame, like adding new columns. It can be either implemented in the plotter processor subclasses or via a mixin class.

plot() → None

The plot method.

This is where the user has to implement the real plot generation

process() → None

Process method overload.

In the case of a plotter subclass, the process method is already implemented and the user should not overload it. On the contrary, the user must overload the other implementation methods described in the general *class description*.

save() → None

The save method.

This is where the user has to implement the saving of the plot on disc.

update_db() → None

The update database method.

This is where the user has to implement the optional update of the database.

```
class mafw.processor_library.abstract_plotter.PlotterMeta(name, bases, namespace, /,
                                                         **kwargs)
```

Bases: `_ProtocolMeta`, `ProcessorMeta`

Metaclass for the plotter mixed classes

mafw.processor_library.importer

Provides a basic element importer.

The first step in the setting up of the analytical framework of a data analysis procedure is to add new elements to the input set. These elements can encompass a wide range of data, including results from experiments or simulations, as well as information gathered through from webscraping or other data sources.

Independently of where the data are coming from, one common task is to add those data to your collection inside the DB, so that the following analytical steps know where the data are and what they are.

This module provides a generic processor that the user can subclass and customize to their needs to import input files. Thanks to a smart filename parsing, other information can be extracted from the filename itself and used to populate additional columns in the dedicated database table.

Classes

<code>FilenameElement</code> (name, regex, value_type, ...)	Helper class for the definition of filename element.
<code>FilenameParser</code> (configuration_file[, filename])	Helper class to interpret all elements in a filename.
<code>Importer</code> (*args, **kwargs)	Importer is the base class for importing elements in the Database structure.

```
class mafw.processor_library.importer.FilenameElement(name: str, regex: str | ~re.Pattern[str],
                                                       value_type: type = <class 'str'>,
                                                       default_value: str | int | float | None =
                                                       None)
```

Bases: `object`

Helper class for the definition of filename element.

While importing an element to the DB, several parameters can be retrieved directly from the filename. The role of this helper class is to provide an easy way to define patterns in the filename representing a specific piece of information that has to be transferred to the DB.

The element is characterized by a name, a regular expression, the expected python type for the parsed value and an optional default value. The regular expression should contain a named group in the form `?P<name>` where name is matching the `FilenameElement` name.

To make a filename element optional, it is enough to provide a default value different from `None`. In this case, if the parsing is failing, then the default value will be returned.

Constructor parameters:

Parameters

- **name** (`str`) – The name of the filename element
- **regex** (`str | re.Pattern[str]`) – The regular expression associated to this filename element. It must contain a named group in the form `?P<name>`.
- **value_type** (`type, Optional`) – The type the output value should be converted into. It defaults to `str`.
- **default_value** (`Any, Optional`) – The default value to assign to the filename element if the pattern is not found in the filename. It defaults to `None`

classmethod `_get_value_type`(*type_as_string: str*) → *type*

Returns the value type.

This method is used by the class method constructor to check if the user provided type in the form of a string is a valid one.

If so, then the corresponding python type is returned, otherwise a `ValueError` exception is raised.

Parameters

type_as_string (*str*) – The type of the value as a string.

Returns

The corresponding python type.

Return type

type

Raises

ValueError – if *type_as_string* is not any of the acceptable type for the value.

classmethod `from_dict`(*name: str, info_dict: dict[str, str | int | float]*) → *FilenameElement*

Generates a `FilenameElement` starting from external information stored in a dictionary.

info_dict should contain the following three keys:

- `regexp`: the Regular expression for the element search.
- `type`: a string with the python type name (`int`, `float`, `str`) for the element conversion.
- `default` (*optional*): a default value.

Parameters

- **name** (*str*) – The name of the element.
- **info_dict** (*dict*) – The dictionary with the required parameters for the class constructor.

Returns

An instance of `FilenameElement`.

Return type

FilenameElement

classmethod `_validate_default_type`() → *None*

Checks that the default has a type matching the value type. The check is actually performed if and only if a default value is provided. If `None`, then the validation is skipped.

Raises

TypeError – if the default value type does not match the declared value type.

classmethod `_validate_regexp`() → *None*

Checks if the regular expression contains a named group named after the element itself.

Raises

ValueError – if the regular expression is not valid.

method `reset`() → *None*

Resets the value to the default value.

Remember: that the default value is `None` for compulsory elements.

method `search`(*string: str | Path*) → *None*

Searches the string for the regular expression.

If the pattern is found in the string, then the matched value is transferred to the `FilenameElement` value.

Note

This method is not returning the match value. It is only searching the input string for the registered pattern. If the pattern is found, then the user can retrieve the matched value by invoking the `value()` method. If the pattern is not found, the `value()` will return either `None`, for a compulsory element, or the default value for an optional one.

Parameters

string (*str* | *Path*) – The string to be parsed. In most of the case, this is a filename, that is why the method is accepting also a *Path* type.

property is_found: **bool**

Returns if the file element is found

property is_optional: **bool**

Returns if the element is optional

property name: **str**

Returns the class name

property pattern: **str** | **bytes**

Returns the regular expression pattern

type_lut: **dict**[**str**, **type**[**str**] | **type**[**int**] | **type**[**float**]] = {'float': <class 'float'>, 'int': <class 'int'>, 'str': <class 'str'>}

A lookup table for converting type definition as string into python types

property value: **str** | **int** | **float** | **None**

Returns the class value

class `maf.processor_library.importer.FilenameParser`(*configuration_file*: *str* | *Path*, *filename*: *str* | *Path* | *None* = *None*)

Bases: `object`

Helper class to interpret all elements in a filename.

Inside a filename, there might be many elements containing information about the item that must be stored in the DB. This class will parse the filename, and after a successful identification of them all, it will make them available for the importer class to fill in the fields in the database.

The `FilenameParser` needs to be configured to be able to recognise each element in the filename. Such configuration is saved in a *toml* file. An example of such a configuration is provided here.

Each element must start with its name and a valid regular expression and a python type (in string). If an element is optional, then a default value must be provided as well.

After the configuration, the filename can be interpreted invoking the `interpret()` method. This will perform the actual parsing of the filename. If an error occurs during the parsing process, meaning that a compulsory element is not found, then the `ParsingError` exception will be raised. So remember to protect the interpretation with a try/except block.

The value of each file element is available upon request. The user has simply to invoke the `get_element_value()` providing the element name.

Constructor parameters:

Parameters

- **filename** (*str* | *Path*) – The filename to be interpreted.
- **configuration_file** (*str* | *Path*) – The configuration file for the interpreter.

Raises

`ParserConfigurationException` – If the configuration file is invalid.

_parser_configuration() → None

Loads the parser configuration, generates the required `FilenameElement` and adds them element dictionary.

The configuration file is stored in a TOML file.

This private method is automatically invoked by the class constructor.

Raises

`ParserConfigurationError` – if the provided configuration file is invalid.

get_element(*element_name: str*) → `FilenameElement` | None

Gets the `FilenameElement` named `element_name`

get_element_value(*element_name: str*) → str | int | float | None

Gets the value of the `FilenameElement` named `element_name`.

It is equivalent to call `self.get_element('element_name').value`

interpret(*filename: str | Path | None = None*) → None

Performs the interpretation of the filename.

The filename can be provided either as constructor argument or here as an argument. If both, then the local one will have the precedence.

Raises

- **`ParsingError`** – if a compulsory element is not found in the filename
- **`MissingAttribute`** – if no filename has been specified.

reset() → None

Resets all filename elements

_configuration_file

The configuration file for the interpreter.

_element_dict: dict[str, `FilenameElement`]

A dictionary with all the `FilenameElement`

_filename

The filename for this interpreter. If None, it should be specified before interpretation.

property elements: dict[str, `FilenameElement`]

Returns the filename element dictionary

class `mafw.processor_library.importer.Importer`(*args: Any, **kwargs: Any)

Bases: `Processor`

`Importer` is the base class for importing elements in the Database structure.

It provides an easy skeleton to be subclassed by a more specific importer related to a certain project.

It can be customised with three processor parameters:

- The `parser_configuration`: the path to the configuration file for the `FilenameParser`.
- The `input_folder`: the path where the input files to be imported are.
- The recursive flag: to specify if all subfolders should be also scanned.

For a concrete implementation, have a look at the `ImporterExample` from the example library.

Processor parameters

- **input_folder**: The input folder from where the images have to be imported. (default: `‘/tmp/mafwdocs-_mot8xgc/v1.2.0_pdf’`)

- **parser_configuration**: The path to the TOML file with the filename parser configuration (default: 'parser_configuration.toml')
- **recursive**: Extend the search to sub-folder (default: True)

Constructor parameters

Parameters

- **name** (*str*, *Optional*) – The name of the processor. If None is provided, the class name is used instead. Defaults to None.
- **description** (*str*, *Optional*) – A short description of the processor task. Defaults to the processor name.
- **config** (*dict*, *Optional*) – A configuration dictionary for this processor. Defaults to None.
- **looper** (*LoopType*, *Optional*) – Enumerator to define the looping type. Defaults to `LoopType.ForLoop`
- **user_interface** (*UserInterfaceBase*, *Optional*) – A user interface instance to be used by the processor to interact with the user.
- **timer** (*Timer*, *Optional*) – A timer object to measure process duration.
- **timer_params** (*dict*, *Optional*) – Parameters for the timer object.
- **database** (*Database*, *Optional*) – A database instance. Defaults to None.
- **database_conf** (*dict*, *Optional*) – Configuration for the database. Default to None.
- **remove_orphan_files** (*bool*, *Optional*) – Boolean flag to remove files on disc without a reference to the database. See *Standard tables* and *_remove_orphan_files()*. Defaults to True
- **kwargs** – Keyword arguments that can be used to set processor parameters.

format_progress_message() → None

Customizes the progress message with information about the current item.

The user can overload this method in order to modify the message being displayed during the process loop with information about the current item.

The user can access the current value, its position in the looping cycle and the total number of items using *Processor.item*, *Processor.i_item* and *Processor.n_item*.

start() → None

The start method.

The filename parser is created using the provided configuration file.

Raises

ParserConfigurationError – If the configuration file is not valid.

_filename_parser: *FilenameParser*

The filename parser instance

maf.processor_library.sns_plotter

Module implements a Seaborn plotter processor with a mixin structure to generate seaborn plots.

This module implements the *abstract_plotter* functionalities using *seaborn* and *pandas*.

These two packages are not installed in the default installation of MAFw, unless the user decided to include the optional feature *seaborn*.

Along with the *SNSPlotter*, it includes a set of standard data retriever specific for pandas data frames.

Classes

<code>CatPlot</code> ([x, y, hue, row, col, palette, ...])	The categorical plot mixin.
<code>DisPlot</code> ([x, y, hue, row, col, palette, ...])	The distribution plot mixin.
<code>FromDatasetDataRetriever</code> ([dataset_name])	A data retriever to get a dataframe from a seaborn dataset
<code>HDFPdDataRetriever</code> ([hdf_filename, key])	Retrieve a data frame from a HDF file
<code>PdDataRetriever</code> (*args, **kwargs)	The dataframe instance.
<code>RelPlot</code> ([x, y, hue, row, col, palette, ...])	The relational plot mixin.
<code>SNSFigurePlotter</code> (*args, **kwargs)	Base mixin class to generate a seaborn Figure level plot
<code>SNSPlotter</code> (*args, **kwargs)	The Generic Plotter processor.
<code>SQLPdDataRetriever</code> ([table_name, ...])	A specialized data retriever to get a data frame from a database table.

```
class mafw.processor_library.sns_plotter.CatPlot(x: str | bytes | date | datetime | timedelta | bool |
complex | Timestamp | Timedelta | Iterable[float
| complex | int] | None = None, y: str | bytes |
date | datetime | timedelta | bool | complex |
Timestamp | Timedelta | Iterable[float | complex
| int] | None = None, hue: str | bytes | date |
datetime | timedelta | bool | complex |
Timestamp | Timedelta | Iterable[float | complex
| int] | None = None, row: str | bytes | date |
datetime | timedelta | bool | complex |
Timestamp | Timedelta | Iterable[float | complex
| int] | None = None, col: str | bytes | date |
datetime | timedelta | bool | complex |
Timestamp | Timedelta | Iterable[float | complex
| int] | None = None, palette: str |
Sequence[tuple[float, float, float]] | str |
tuple[float, float, float, float] | tuple[tuple[float,
float, float] | str, float] | tuple[tuple[float, float,
float, float], float]] | Mapping[Any, tuple[float,
float, float] | str | tuple[float, float, float, float] |
tuple[tuple[float, float, float] | str, float] |
tuple[tuple[float, float, float, float], float]] |
None = None, kind: Literal['strip', 'swarm',
'box', 'violin', 'boxen', 'point', 'bar', 'count'] =
'strip', legend: Literal['auto', 'brief', 'full'] | bool
= 'auto', native_scale: bool = False, plot_kws:
Mapping[str, Any] | None = None, facet_kws:
dict[str, Any] | None = None, *args: Any,
**kwargs: Any)
```

Bases: `SNSFigurePlotter`

The categorical plot mixin.

This mixin will produce a figure level categorical plot as described [here](#).

Constructor parameters:

Parameters

- **x** (`str | Iterable, Optional`) – The name of the x variable or an iterable containing the x values.
- **y** (`str | Iterable, Optional`) – The name of the y variable or an iterable containing the y values.

- **hue** (*str* | *Iterable*, *Optional*) – The name of the hue variable or an iterable containing the hue values.
- **row** (*str* | *Iterable*, *Optional*) – The name of the row category or an iterable containing the row values.
- **col** (*str* | *Iterable*, *Optional*) – The name of the column category or an iterable containing the column values.
- **palette** (*str*, *Optional*) – The colour palette to be used.
- **kind** (*str*, *Optional*) – The type of relational plot (scatter or line). Defaults to scatter.
- **legend** (*str* | *bool*, *Optional*) – How to draw the legend. If “brief”, numeric hue and size variables will be represented with a sample of evenly spaced values. If “full”, every group will get an entry in the legend. If “auto”, choose between brief or full representation based on number of levels. If False, no legend data is added and no legend is drawn. Defaults to auto.
- **native_scale** (*bool*, *Optional*) – When True, numeric or datetime values on the categorical axis will maintain their original scaling rather than being converted to fixed indices. Defaults to False.
- **plot_kws** (*dict*[*str*, *Any*], *Optional*) – A dictionary like list of keywords passed to the underlying `seaborn.catplot`.
- **facet_kws** (*dict*[*str*, *Any*], *Optional*) – A dictionary like list of keywords passed to the underlying `seaborn.FacetGrid`

`plot()` → None

Implements the plot method of a figure-level categorical graph.

```
class mafw.processor_library.sns_plotter.DisPlot(x: str | bytes | date | datetime | timedelta | bool |
complex | Timestamp | Timedelta | Iterable[float
| complex | int] | None = None, y: str | bytes |
date | datetime | timedelta | bool | complex |
Timestamp | Timedelta | Iterable[float | complex
| int] | None = None, hue: str | bytes | date |
datetime | timedelta | bool | complex |
Timestamp | Timedelta | Iterable[float | complex
| int] | None = None, row: str | bytes | date |
datetime | timedelta | bool | complex |
Timestamp | Timedelta | Iterable[float | complex
| int] | None = None, col: str | bytes | date |
datetime | timedelta | bool | complex |
Timestamp | Timedelta | Iterable[float | complex
| int] | None = None, palette: str |
Sequence[tuple[float, float, float]] | str |
tuple[float, float, float, float] | tuple[tuple[float,
float, float] | str, float] | tuple[tuple[float, float,
float, float], float]] | Mapping[Any, tuple[float,
float, float] | str | tuple[float, float, float, float] |
tuple[tuple[float, float, float] | str, float] |
tuple[tuple[float, float, float, float], float]] |
Colormap | None = None, kind: Literal['hist',
'kde', 'ecdf'] = 'hist', legend: bool = True, rug:
bool = False, rug_kws: dict[str, Any] | None =
None, plot_kws: Mapping[str, Any] | None =
None, facet_kws: dict[str, Any] | None = None,
*args: Any, **kwargs: Any)
```

Bases: `SNSFigurePlotter`

The distribution plot mixin.

This mixin is the MAFw implementation of the [seaborn displot](#) and will produce one of the following figure level plots:

- **histplot**: a simple [histogram plot](#)
- **kdeplot**: a [kernel density estimate plot](#)
- **ecdfplot**: an [empirical cumulative distribution functions plot](#)
- **rugplot**: a plot of the [marginal distributions](#) as ticks.

Constructor parameters:

Parameters

- **x** (*str* | *Iterable*, *Optional*) – The name of the x variable or an iterable containing the x values.
- **y** (*str* | *Iterable*, *Optional*) – The name of the y variable or an iterable containing the y values.
- **hue** (*str* | *Iterable*, *Optional*) – The name of the hue variable or an iterable containing the hue values.
- **row** (*str* | *Iterable*, *Optional*) – The name of the row category or an iterable containing the row values.
- **col** (*str* | *Iterable*, *Optional*) – The name of the column category or an iterable containing the column values.
- **palette** (*str* | *Colormap*, *Optional*) – The colour palette to be used.
- **kind** (*str*, *Optional*) – The type of distribution plot (hist, kde or ecdf). Defaults to hist.
- **legend** (*bool*, *Optional*) – If false, suppress the legend for the semantic variables. Defaults to True.
- **rug** (*bool*, *Optional*) – If true, show each observation with marginal ticks. Defaults to False.
- **rug_kws** (*Mapping*[*str*, *Any*], *Optional*) – Parameters to control the appearance of the rug plot.
- **plot_kws** (*Mapping*[*str*, *Any*], *Optional*) – Parameters passed to the underlying plotting object.
- **facet_kws** (*Mapping*[*str*, *Any*], *Optional*) – Parameters passed to the facet grid object.

plot() → None

Implements the plot method for a figure-level distribution graph

```
class mafw.processor_library.sns_plotter.FromDatasetDataRetriever(dataset_name: str | None
                                                                = None, *args: Any,
                                                                **kwargs: Any)
```

Bases: [PdDataRetriever](#)

A data retriever to get a dataframe from a seaborn dataset

The dataframe instance. It will be filled for the main class

_attributes_valid() → bool

Checks if the attributes of the mixin are all valid

get_data_frame() → None

Gets the data frame from the standard seaborn datasets

```
class mafw.processor_library.sns_plotter.HDFPdDataRetriever(hdf_filename: str | Path | None = None, key: str | None = None, *args: Any, **kwargs: Any)
```

Bases: *DataRetriever*

Retrieve a data frame from a HDF file

This data retriever is getting a dataframe from a HDF file provided the filename and the object key.

Constructor parameters:

Parameters

- **hdf_filename** (*str | Path, Optional*) – The filename of the HDF file
- **key** (*str, Optional*) – The key of the HDF store with the dataframe

```
get_data_frame() → None
```

Retrieve the dataframe from a HDF file

Raises

PlotterMixinNotInitialized – if some of the required attributes are not initialised or invalid.

```
patch_data_frame() → None
```

The mixin implementation of the shared method with the base class

```
class mafw.processor_library.sns_plotter.PdDataRetriever(*args: Any, **kwargs: Any)
```

Bases: *DataRetriever*

The dataframe instance. It will be filled for the main class

```
get_data_frame() → None
```

The mixin implementation of the shared method with the base class

```
patch_data_frame() → None
```

The mixin implementation of the shared method with the base class

```
class mafw.processor_library.sns_plotter.RelPlot(x: str | bytes | date | datetime | timedelta | bool | complex | Timestamp | Timedelta | Iterable[float | complex | int] | None = None, y: str | bytes | date | datetime | timedelta | bool | complex | Timestamp | Timedelta | Iterable[float | complex | int] | None = None, hue: str | bytes | date | datetime | timedelta | bool | complex | Timestamp | Timedelta | Iterable[float | complex | int] | None = None, row: str | bytes | date | datetime | timedelta | bool | complex | Timestamp | Timedelta | Iterable[float | complex | int] | None = None, col: str | bytes | date | datetime | timedelta | bool | complex | Timestamp | Timedelta | Iterable[float | complex | int] | None = None, palette: str | Sequence[tuple[float, float, float]] | str | tuple[float, float, float] | tuple[tuple[float, float, float]] | str, float] | tuple[tuple[float, float, float], float]] | Mapping[Any, tuple[float, float, float]] | str | tuple[float, float, float, float] | tuple[tuple[float, float, float]] | str, float] | tuple[tuple[float, float, float, float], float]] | Colormap | None = None, kind: Literal['scatter', 'line'] = 'scatter', legend: Literal['auto', 'brief', 'full'] | bool = 'auto', plot_kws: Mapping[str, Any] | None = None, facet_kws: dict[str, Any] | None = None, *args: Any, **kwargs: Any)
```

Bases: *SNSFigurePlotter*

The relational plot mixin.

This mixin will produce either a scatter or a line figure level plot.

The full documentation of the relplot object can be read at [this link](#).

Constructor parameters:

Parameters

- **x** (*str* | *Iterable*, *Optional*) – The name of the x variable or an iterable containing the x values.
- **y** (*str* | *Iterable*, *Optional*) – The name of the y variable or an iterable containing the y values.
- **hue** (*str* | *Iterable*, *Optional*) – The name of the hue variable or an iterable containing the hue values.
- **row** (*str* | *Iterable*, *Optional*) – The name of the row category or an iterable containing the row values.
- **col** (*str* | *Iterable*, *Optional*) – The name of the column category or an iterable containing the column values.
- **palette** (*str* | *Colormap*, *Optional*) – The colour palette to be used.
- **kind** (*str*, *Optional*) – The type of relational plot (scatter or line). Defaults to scatter.
- **legend** (*str* | *bool*, *Optional*) – How to draw the legend. If “brief”, numeric hue and size variables will be represented with a sample of evenly spaced values. If “full”, every group will get an entry in the legend. If “auto”, choose between brief or full representation based on number of levels. If False, no legend data is added and no legend is drawn. Defaults to auto.
- **plot_kws** (*dict*[*str*, *Any*], *Optional*) – A dictionary like list of keywords passed to the underlying *seaborn.relplot*.
- **facet_kws** (*dict*[*str*, *Any*], *Optional*) – A dictionary like list of keywords passed to the underlying *seaborn.FacetGrid*

plot() → None

Implements the plot method of a figure-level relational graph.

```
class mafw.processor_library.sns_plotter.SNSFigurePlotter(*args: Any, **kwargs: Any)
```

Bases: *FigurePlotter*

Base mixin class to generate a seaborn Figure level plot

data_frame: **DataFrame**

The dataframe instance shared with the main class

facet_grid: **FacetGrid**

The facet grid instance shared with the main class

```
class mafw.processor_library.sns_plotter.SNSPlotter(*args: Any, **kwargs: Any)
```

Bases: *GenericPlotter*

The Generic Plotter processor.

This is a subclass of a Processor with advanced functionality to fetch data in the form of a dataframe and to produce plots.

The key difference with respect to a normal processor is it *process()* method that has been already implemented as follows:

```

def process(self) -> None:
    """
    Specific implementation of the process method for the Seaborn plotter.

    It is almost the same as the GenericProcessor, with the addition that all
    ↪open pyplot figures are closed
    after the process is finished.

    This part cannot be moved upward to the :class:`~.GenericPlotter` because the
    ↪user might have selected
    another plotting library different from :link:`matplotlib`.
    """
    super().process()
    if not self.is_data_frame_empty():
        plt.close('all')

```

This actually means that when you are subclassing a SNSPlotter you do not have to implement the process method as you would do for a normal Processor, but you will have to implement the following methods:

- `in_loop_customization()`.

The processor execution workflow (LoopType) can be any of the available, so actually the process method might be invoked only once, or multiple times inside a loop structure (for or while). If the execution is cyclic, then you may want to have the possibility to do some customisation for each iteration, for example, changing the plot title, or modifying the data selection, or the filename where the plots will be saved.

You can use this method also in case of a single loop processor, in this case you will not have access to the loop parameters.

- `get_data_frame()`.

This method has the task to get the data to be plotted in the form of a pandas DataFrame. The processor has the `data_frame` attribute where the data will be stored to make them accessible from all other methods.

- `patch_data_frame()`.

A convenient method to apply data frame manipulation to the data just retrieved.

- `plot()`.

This method is where the actual plotting occurs. Use the `data_frame` to plot the quantities you want.

- `customize_plot()`.

This method can be optionally used to customize the appearance of the facet grid produced by the `plot()` method. It is particularly useful when the user is mixing this class with one of the `FigurePlotter` mixin, thus not having direct access to the plot method.

- `save()`.

This method is where the produced plot is saved in a file. Remember to append the output file name to the `list of produced outputs` so that the `_update_plotter_db()` method will automatically store this file in the database during the `finish()` execution.

- `update_db()`.

If the user wants to update a specific table in the database, they can use this method.

It is worth reminding that all plotters are saving all generated files in the standard table `PlotterOutput`. This is automatically done by the `_update_plotter_db()` method that is called in the `finish()` method.

You do not need to overload the `slice_data_frame()` nor the `group_and_aggregate_data_frame()` methods, but you can simply use them by setting the `slicing_dict` and the `grouping_columns` and the `aggregation_functions`.

Constructor parameters:

Parameters

- **slicing_dict** (`dict[str, Any]`, *Optional*) – A dictionary with key, value pairs to slice the input data frame before the plotting occurs.
- **grouping_columns** (`list[str]`, *Optional*) – A list of columns for the groupby operation on the data frame.
- **aggregation_functions** (`list[str | Callable[[Any], Any]`, *Optional*) – A list of functions for the aggregation on the grouped data frame.
- **matplotlib_backend** (`str`, *Optional*) – The name of the matplotlib backend to be used. Defaults to 'Agg'

get_data_frame() → None

Specific implementation of the get data frame for the Seaborn plotter.

It must be overloaded.

The method is **NOT** returning the `data_frame`, but in your implementation you need to assign the data frame to the class `data_frame` attribute.

group_and_aggregate_data_frame() → None

Performs groupby and aggregation of the data frame.

If the user provided some `grouping_columns` and `aggregation_functions` then the `group_and_aggregate_data_frame()` is invoked accordingly.

The user can update the values of those attributes during each cycle iteration within the implementation of the `in_loop_customization()`.

➔ See also

This method is simply invoking the `group_and_aggregate_data_frame()` function from the `pandas_tools`.

is_data_frame_empty() → bool

Check if the data frame is empty

process() → None

Specific implementation of the process method for the Seaborn plotter.

It is almost the same as the `GenericProcessor`, with the addition that all open pyplot figures are closed after the process is finished.

This part cannot be moved upward to the `GenericPlotter` because the user might have selected another plotting library different from `matplotlib`.

slice_data_frame() → None

Perform data frame slicing

The user can set some slicing criteria in the `slicing_dict` to select some specific data subset. The values of the slicing dict can be changed during each iteration within the implementation of the `in_loop_customization()`.

➔ See also

This method is simply invoking the `slice_data_frame()` function from the `pandas_tools`.

start() → None

Overload of the start method.

The *SNSPlotter* is overloading the *start()* in order to allow the user to change the matplotlib backend.

The user can selected which backend to use either directly in the class constructor or assign it to the class attribute *matplotlib_backend*.

aggregation_functions: `Iterable[str | Callable[[Any], Any]] | None`

The list of aggregation functions to be applied to the grouped dataframe

data_frame: `pd.DataFrame`

The pandas DataFrame containing the data to be plotted.

facet_grid: `sns.FacetGrid | None`

The reference to the facet grid.

grouping_columns: `Iterable[str] | None`

The list of columns for grouping the data frame

matplotlib_backend: `str`

The backend to be used for matplotlib.

output_filename_list: `list[Path]`

The list of produced filenames.

slicing_dict: `MutableMapping[str, Any] | None`

The dictionary for slicing the input data frame

```
class mafw.processor_library.sns_plotter.SQLPdDataRetriever(table_name: str | None = None,
                                                         required_cols: Iterable[str] | str |
                                                         None = None, where_clause: str |
                                                         None = None, *args: Any,
                                                         **kwargs: Any)
```

Bases: *PdDataRetriever*

A specialized data retriever to get a data frame from a database table.

The idea is to implement an interface to the pandas *read_sql*. The user has to provide the *table name*, the *the list of required columns* and an optional *where clause*.

Constructor parameters:

Parameters

- **table_name** (*str, Optional*) – The name of the table from where to get the data
- **required_cols** (*Iterable[str] | str | None, Optional*) – A list of columns to be selected from the table and transferred as column in the dataframe.
- **where_clause** (*str, Optional*) – The where clause used in the select SQL statement. If None is provided, then all rows will be selected.

_attributes_valid() → bool

Check if all required parameters are provided and valid.

get_data_frame() → None

Retrieve the dataframe from a database table.

Raises

PlotterMixinNotInitialized – If some of the required attributes are missing.

database: `Database`

The database instance. It comes from the main class

required_columns: `Iterable[str]`

The iterable of columns.

Those are the column names to be selected from the `table_name` and included in the dataframe.

table_name: `str`

The table from where the data should be taken.

where_clause: `str`

The where clause of the SQL statement

11.1.12 mafw.runner

Provides a container to run configurable and modular analytical tasks.

Classes

<code>MAFwApplication</code> ([steering_file, ...])	The MAFw Application.
-----------------------------------------------------	-----------------------

```
class mafw.runner.MAFwApplication(steering_file: Path | str | None = None, user_interface:
    UserInterfaceBase | type[UserInterfaceBase] | str | None = None,
    plugin_manager: PluginManager | None = None)
```

Bases: object

The MAFw Application.

This class takes care of reading a steering file and from the information retrieved construct a `ProcessorList` from the processor listed there and execute it.

It is very practical because any combination of processors can be run without having to write dedicated scripts but simply modifying a steering file.

The application will search for processors not only among the ones available in the MAFw library, but also in all other packages exposing processors via the *plugin mechanism*.

All parameters in the constructor are optional.

An instance can be created also without the `steering_file`, but such an instance cannot be executed. The steering file can be provided in a later stage via the `init()` method or directly to the `run()` method.

The user interface can be either provided directly in the constructor, or it will be taken from the steering file. In the worst case, the fallback `ConsoleInterface` will be used.

The plugin manager, if not provided, the global plugin manager will be retrieved from the `get_plugin_manager()`.

A simple example is provided here below:

Listing 11.1: Creation and execution of a MAFwApplication

```
import logging
from pathlib import Path

from mafw.runner import MAFwApplication

log = logging.getLogger(__name__)

# put here your steering file
steering_file = Path('path_to_my_steering_file.toml')

try:
    # create the app
```

(continues on next page)

(continued from previous page)

```

app = MAFwApplication(steering_file)

# run it!
app.run()

except Exception as e:
    log.error('An error occurred!')
    log.exception(e)

```

Constructor parameters:

Parameters

- **steering_file** (*Path | str, Optional*) – The path to the steering file.
- **user_interface** (*UserInterfaceBase | type[UserInterfaceBase] | str, Optional*) – The user interface to be used by the application.
- **plugin_manager** (*PluginManager, Optional*) – The plugin manager.

get_user_interface(*user_interface: str*) → *UserInterfaceBase*

Retrieves the user interface from the plugin managers.

User interfaces are exposed via the plugin manager.

If the requested *user_interface* is not available, then the fallback console interface is used.

Parameters

user_interface (*str*) – The name of the user interface to be used. Normally rich or console.

init(*steering_file: Path | str*) → *None*

Initializes the application.

This method is normally automatically invoked by the class constructor. It can be called in a later moment to force the parsing of the provided steering file.

Parameters

steering_file (*Path | str*) – The path to the steering file.

run(*steering_file: Path | str | None = None*) → *ProcessorExitStatus*

Runs the application.

This method builds the *ProcessorList* with the processors listed in the steering file and launches its execution.

A steering file can be provided at this stage if it was not done before.

Parameters

steering_file (*Path | str, Optional*) – The steering file. Defaults to None.

Raises

- **RunnerNotInitialized** – if the application has not been initialized. Very likely a steering file was never provided.
- **UnknownProcessor** – if a processor listed in the steering file is not available in the plugin library.

exit_status

the exit status of the application

name

the name of the application instance

plugin_manager

the plugin manager of the application instance

11.1.13 mafw.scripts

Executables.

This package includes the executables necessary to run MAFw.

Modules

<i>mafw_exe</i>	The execution framework.
<i>update_changelog</i>	Module provides a simplified script to perform MAFw changelog update using auto-changelog.

mafw.scripts.mafw_exe

The execution framework.

This module provides the run functionality to the whole library.

It is heavily relying on `click` for the generation of commands, options, and arguments.

mafw

The Modular Analysis Framework execution.

This is the command line interface where you can configure and launch your analysis tasks.

More information on our documentation page.

param ctx

The click context.

type ctx

`click.core.Context`

param log_level

The logging level as a string. Choice from debug, info, warning, error and critical.

type log_level

`str`

param ui

The user interface as a string. Choice from console and rich.

type ui

`str`

param debug

Flag to show debug information about exception.

type debug

`bool`

Usage

```
mafw [OPTIONS] COMMAND [ARGS] ...
```

Options

--log-level <log_level>

Log level

Default

'info'

Options

debug | info | warning | error | critical

--ui <ui>

The user interface

Default

'rich'

Options

console | rich

-D, --debug

Show debug information about errors

-v, --version

Show the version and exit.

db

Advanced database commands.

The db group of commands offers a set of useful database operations. Invoke the help option of each command for more details.

param ctx

The click context.

type ctx

click.core.Context

Usage

```
mafw db [OPTIONS] COMMAND [ARGS]...
```

wizard

Reflect an existing DB into a python module.

```
mafw db wizard [Options] Database
```

Database Name of the Database to be reflected.

About connection options (user / host / port):

That information will be used only in case you are trying to access a network database (MySQL or PostgreSQL). In case of Sqlite, the parameters will be discarded.

About passwords:

If you need to specify a password to connect to the DB server, just add `-password` in the command line without typing your password as clear text. You will be prompted to insert the password with hidden characters at the start of the processor.

About engines:

The full list of supported engines is provided in the option below. If you do not specify any engine and the database is actually an existing filename, then engine is set to Sqlite, otherwise to postgresql.

param database

The name of the database.

type database

str

param schema

The database schema to be reflected.

type schema

str

param engine

The database engine. A selection of possible values is provided in the script help.

type engine

str

param password

The password for the DB connection. Not used in case of Sqlite.

type password

str

param user

The username for the DB connection. Not used in case of Sqlite.

type user

str

param port

The port number of the database server. Not used in case of Sqlite.

type port

int

param host

The database hostname. Not used in case of Sqlite.

type host

str

param output_file

The filename for the output python module.

type output_file

click.Path | pathlib.Path | str

param snake_case

Flag to select snake_case convention for table and field names, or all small letter formatting.

type snake_case

bool

param ignore_unknown

Flag to ignore unknown fields. If False, an unknown field will be labelled with UnknownField.

type ignore_unknown

bool

param with_views

Flag to include views in the reflected elements.

type with_views

bool

param preserve_order

Flag to select if table fields should be reflected in the original order (True) or in alphabetical order (False)

type preserve_order

bool

param tables

A tuple containing a selection of table names to be reflected.

type tables

tuple[str, ...]

param overwrite

Flag to overwrite the output file if exists. If False and the output file already exists, the user can decide what to do.

type overwrite

bool

param ctx

The click context, that includes the original object with global options.

type ctx

click.core.Context

return

The script return value

Usage

```
mafw db wizard [OPTIONS] DATABASE
```

Options

- o, --output-file** <output_file>
The name of the output file with the reflected model.
- s, --schema** <schema>
The name of the DB schema
- t, --tables** <tables>
Generate model for selected tables. Multiple option possible.
- overwrite, --no-overwrite**
Overwrite output file if already exists.
- preserve-order, --no-preserve-order**
Preserve column order.
- with-views, --without-views**
Include also database views.
- ignore-unknown, --no-ignore-unknown**
Ignore unknown fields.
- snake-case, --no-snake-case**
Use snake case for table and field names.
- host** <host>
Hostname for the DB server.
- p, --port** <port>
Port number for the DB server.
- u, --user, --username** <user>
Username for the connection to the DB server.

--password <password>

Insert password when prompted

-e, --engine <engine>

The DB engine

Options

cockroach | cockroachdb | crdb | mysql | mysqldb | postgres | postgresql | sqlite | sqlite3

Arguments

DATABASE

Required argument

list

Display the list of available processors.

This command will retrieve all available processors via the plugin manager. Both internal and external processors will be listed if the ext-plugins option is passed.

param ext_plugins

Flag to enable searching for processor external libraries. Defaults to True

type ext_plugins

bool, Optional

Usage

```
mafw list [OPTIONS]
```

Options

--ext-plugins, --no-ext-plugin

Load external plugins

run

Runs a steering file.

STEERING_FILE A path to the steering file to execute.

param obj

The context object being passed from the main command.

type obj

dict

param steering_file

The path to the output steering file.

type steering_file

Path

Usage

```
mafw run [OPTIONS] STEERING_FILE
```

Arguments

STEERING_FILE

Required argument

steering

Generates a steering file with the default parameters of all available processors.

STEERING_FILE A path to the steering file to execute.

The user must modify the generated steering file to ensure it can be executed using the run command.

param obj

The context object being passed from the main command.

type obj

dict

param show

Display the steering file in the console after the generation. Defaults to False.

type show

bool

param ext_plugins

Extend the search for processor to external libraries.

type ext_plugins

bool

param open_editor

Open a text editor after the generation to allow direct editing.

type open_editor

bool

param steering_file

The steering file path.

type steering_file

Path

param db_engine

The name of the db engine.

type db_engine

str

param db_url

The URL of the database.

type db_url

str

Usage

```
mafws steering [OPTIONS] STEERING_FILE
```

Options

--show, --no-show

Display the generated steering file on console

--ext-plugins, --no-ext-plugin

Load external plugins

--open-editor, --no-open-editor

Open the file in your editor.

--db-engine <db_engine>

Select a DB engine

Options

sqlite | mysql | postgresql

--db-url <db_url>

URL to the DB

Arguments

STEERING_FILE

Required argument

Functions

<code>custom_formatwarning</code> (message, category, ...)	Return the pure message of the warning.
<code>logger_setup</code> (level, ui, tracebacks)	Set up the logger.

Classes

<code>MAFwGroup</code> ([name, commands, ...])	Custom Click Group for MAFw runner.
<code>ReturnValue</code> (*values)	Enumerator to handle the script return value.

```
class mafw.scripts.mafw_exe.MAFwGroup(name: str | None = None, commands: MutableMapping[str,
    Command] | Sequence[Command] | None = None,
    invoke_without_command: bool = False, no_args_is_help:
    bool | None = None, subcommand_metavar: str | None = None,
    chain: bool = False, result_callback: Callable[[...], Any] |
    None = None, **kwargs: Any)
```

Bases: Group

Custom Click Group for MAFw runner.

It implements two main features:

1. Support commands abbreviation. Instead of providing the whole command, the user can use whatever abbreviation instead as long as it is unique. So for example, instead of `mafw list`, the user can provide `mafw l` and the result will be the same.
2. Implements the cascading of return values among different command levels.

get_command(ctx: Context, cmd_name: str) → Command | None

Return a command.

Given a context and a command name as passed from the CLI, the click.Command is returned. This method overloads the basic one allowing to use command abbreviations.

If more than one match is found, then an error is raised.

If no matches are found, then click will handle this case as in the standard situation.

Parameters

- **ctx** – The click context
- **cmd_name** – The command name as provided from the CLI

Returns

The corresponding command or None if no command is found.

invoke(*ctx: Context*) → Any

Invoke the command.

This override method is just wrapping the base invoke call in a try / except block.

In the case of a ClickException, then this is shown and its exit code is used passed to the sys.exit call. In case of a SystemExit or click.exceptions.Exit, then this is simply re-raised, so that Click can handle it as in normal circumstances. In all other cases, the exception is caught and the sys.exit is called with the *ReturnValue.Error*.

Parameters

ctx – The click context

Returns

The return value of the invoked command

main(*args: Any, **kwargs: Any) → None

Override main to handle return values properly.

class mafw.scripts.mafw_exe.**ReturnValue**(*values)

Bases: IntEnum

Enumerator to handle the script return value.

Error = 1

Generic error

OK = 0

No error

mafw.scripts.mafw_exe.**custom_formatwarning**(*message: Warning | str, category: type[Warning], filename: str, lineno: int, line: str | None = None*) → str

Return the pure message of the warning.

mafw.scripts.mafw_exe.**logger_setup**(*level: str, ui: str, tracebacks: bool*) → None

Set up the logger.

This function is actually configuring the root logger level from the command line options and it attaches either a RichHandler or a StreamHandler depending on the user interface type.

The *tracebacks* flag is used only by the RichHandler. Printing the tracebacks is rather useful when debugging the code, but it could be detrimental for final users. In normal circumstances, tracebacks is set to False, and is turned on when the debug flag is activated.

Parameters

- **level** (*str*) – Logging level as a string.
- **ui** (*str*) – User interface as a string ('rich' or 'console').
- **tracebacks** – Enable/disable the logging of exception tracebacks.

mafw.scripts.update_changelog

Module provides a simplified script to perform MAFw changelog update using auto-changelog. It can be used as pre-commit entry point and also in CI.

The basic idea is that this command is invoking the auto-changelog tool to generate a temporary changelog. The checksum of the temporary changelog is compared with the existing one. If the two checksums differs, the current changelog is replaced with the newly created version.

When committing the changelog update please use mute as commit type, to avoid having a new changelog generated containing the changelog update commit.

Functions

<code>commit_changelog_changes()</code>	Commit the changes to CHANGELOG.md.
<code>get_last_commit_message()</code>	Get the message of the last commit.
<code>get_latest_tag()</code>	Get the latest git tag.
<code>main()</code>	Script entry point

`mafwscripsts.update_changelog.commit_changelog_changes()` → None
Commit the changes to CHANGELOG.md.

`mafwscripsts.update_changelog.get_last_commit_message()` → str
Get the message of the last commit.

Returns

The last commit message

Return type

str

`mafwscripsts.update_changelog.get_latest_tag()` → str
Get the latest git tag.

Returns

The last git tag.

Return type

str

`mafwscripsts.update_changelog.main()` → None
Script entry point

11.1.14 mafw.timer

Module implements a simple timer to measure the execution duration.

Basic usage:

```
from mafw import timer

with Timer() as timer:
    do_long_lasting_operation()
```

When exiting from the context manager, a message with the duration of the process is printed.

Functions

<code>pretty_format_duration(duration_s[, n_digits])</code>	Return a formatted version of the duration with increased human readability.
<code>rreplace(inp_string, old_string, new_string, ...)</code>	Utility function to replace a substring in a given string a certain number of times starting from the right-most one.

Classes

<code>Timer([suppress_message])</code>	The timer class.
----------------------------------------	------------------

class mafw.timer.Timer(*suppress_message: bool = False*)

Bases: object

The timer class.

Constructor parameter:

Parameters

suppress_message (*bool*) – A boolean flag to mute the timer

format_duration() → str

Nicely format the timer duration.

Returns

A string with the timer duration in a human-readable formatted string

Return type

str

property duration: float

The elapsed time of the timer.

Returns

Elapsed time in seconds.

mafw.timer.pretty_format_duration(*duration_s: float, n_digits: int = 1*) → str

Return a formatted version of the duration with increased human readability.

Parameters

- **duration_s** (*float*) – The duration to be printed in seconds. If negative, a ValueError exception is raised.
- **n_digits** (*int, Optional*) – The number of decimal digits to show. Defaults to 1. If negative, a ValueError exception is raised.

Returns

The formatted string.

Return type

str

Raises

ValueError – if a negative duration or a negative number of digits is provided

mafw.timer.rreplace(*inp_string: str, old_string: str, new_string: str, counts: int*) → str

Utility function to replace a substring in a given string a certain number of times starting from the right-most one.

This function is mimicking the behavior of the string.replace method, but instead of replacing from the left, it is replacing from the right.

Parameters

- **inp_string** (*str*) – The input string
- **old_string** (*str*) – The old substring to be matched. If empty, a ValueError is raised.
- **new_string** (*str*) – The new substring to be replaced
- **counts** (*int*) – The number of times the old substring has to be replaced.

Returns

The modified string

Return type

str

Raises

ValueError – if old_string is empty.

11.1.15 mafw.tools

The package provides a set of tools for different range of applications.

Modules

<code>file_tools</code>	The module provides utilities for handling file, filename, hashing and so on.
<code>pandas_tools</code>	A collection of useful convenience functions for common pandas operations
<code>regex</code>	Module implements some basic functions involving regular expressions.
<code>toml_tools</code>	The module provides tools to read / write / modify specific TOML files.

mafw.tools.file_tools

The module provides utilities for handling file, filename, hashing and so on.

Functions

<code>file_checksum</code> (filenames[, buf_size])	Generates the hexadecimal digest of a file or a list of files.
<code>remove_widow_db_rows</code> (models)	Removes widow rows from a database table.
<code>verify_checksum</code> (models)	Verifies the goodness of FileChecksumField.

```
mafw.tools.file_tools.file_checksum(filenames: str | Path | Sequence[str | Path], buf_size: int = 65536) → str
```

Generates the hexadecimal digest of a file or a list of files.

The digest is calculated using the sha256 algorithm.

Parameters

- **filenames** (*str*, *Path*, *list*) – The filename or the list of filenames for digest calculations.
- **buf_size** (*int*, *Optional*) – The buffer size in bytes for reading the input files. Defaults to 64kB.

Returns

The hexadecimal digest.

Return type

str

```
mafw.tools.file_tools.remove_widow_db_rows(models: list[Model | type[Model]] | Model | type[Model]) → None
```

Removes widow rows from a database table.

According to MAFw architecture, the Database is mainly providing I/O support to the various processors.

This means that the processor retrieves a list of items from a database table for processing and subsequently updates a result table with the newly generated outputs.

Very often the input and output data are not stored directly in the database, but rather in files saved on the disc. In this case, the database is just providing a valid path where the input (or output) data can be found.

From this point of view, a **widow row** is a database entry in which the file referenced by the FilenameField has been deleted. A typical example is the following: the user wants a certain processor to regenerate a given result stored inside an output file. Instead of setting up a complex filter so that the processor receives only

this element to process, the user can delete the actual output file and ask the processor to process all new items.

The provided `models` can be either a list or a single element, representing either an instance of a DB model or a model class. If a model class is provided, then a select over all its entries is performed.

The function will look at all fields of `FileNameField` and `FileNameListField` and check if it corresponds to an existing path or list of paths. If not, then the corresponding row is removed from the DB table.

Parameters

`models` (`list[Model | type(Model)] | Model | type(Model)`) – A list or a single Model instance or Model class for widow rows removal.

Raises

TypeError – if `models` is not of the right type.

`maf.tools.file_tools.verify_checksum(models: list[Model | type[Model]] | Model | type[Model]) → None`

Verifies the goodness of `FileChecksumField`.

If in a model there is a `FileChecksumField`, this must be connected to a `FileNameField` or a `FileNameListField` in the same model. The goal of this function is to recalculate the checksum of the `FileNameField` / `FileNameListField` and compare it with the actual stored value. If the newly calculated value differs from the stored one, the corresponding row in the model will be removed, as it is no longer valid.

If a file is missing, then the checksum check is not performed, but the row is removed right away.

This function can be CPU and I/O intensive and last a lot, so use it with care, especially when dealing with long tables and large files.

Parameters

`models` (`list[Model | type(Model)] | Model | type(Model)`) – A list or a single Model instance or Model class for checksum verification.

Raises

- **TypeError** – if `models` is not of the right type.
- `maf.mafw_errors.ModelError` – if the `FileChecksumField` is referring to a `FileNameField` that does not exist.

maf.tools.pandas_tools

A collection of useful convenience functions for common pandas operations

Functions

<code>group_and_aggregate_data_frame(data_frame, ...)</code>	Utility function to perform dataframe groupby and aggregation.
<code>slice_data_frame(input_data_frame[, ...])</code>	Slice a data frame according to <code>slicing_dict</code> .

`maf.tools.pandas_tools.group_and_aggregate_data_frame(data_frame: DataFrame, grouping_columns: Iterable[str], aggregation_functions: Iterable[str | Callable[[Any], Any]]) → DataFrame`

Utility function to perform dataframe groupby and aggregation.

This function is a simple wrapper to perform group by and aggregation operations on a dataframe. The user must provide a list of columns to perform the group by on and a list of functions for the aggregation of the other columns.

The output dataframe will have the aggregated columns renamed as `originalname_aggregationfunction`.

Note

Only numeric columns (and columns that can be aggregated) will be included in the aggregation. String columns that are not used for grouping will be automatically excluded from aggregation.

Parameters

- **data_frame** (*pandas.DataFrame*) – The input data frame
- **grouping_columns** (*Iterable[str]*) – The list of columns to group by on.
- **aggregation_functions** (*Iterable[str | Callable[[Any], Any]*) – The list of functions to be used for the aggregation of the not grouped columns.

Returns

The aggregated dataframe after the groupby operation.

Return type

pandas.DataFrame

`maf.tools.pandas_tools.slice_data_frame`(*input_data_frame: DataFrame, slicing_dict: MutableMapping[str, Any] | None = None, **kwargs: Any*) → *DataFrame*

Slice a data frame according to *slicing_dict*.

The input data frame will be sliced using the items of the *slicing_dict* applying the loc operator in this way: `sliced = input_data_frame[(input_data_frame[key]==value)]`.

If the *slicing_dict* is empty, then the full *input_data_frame* is returned.

Instead of the *slicing_dict*, the user can also provide key and value pairs as keyword arguments.

```
slice_data_frame(data_frame, {'A':14})
```

is equivalent to

```
slice_data_frame(data_frame, A=14).
```

If the user provides a keyword argument that also exists in the *slicing_dict*, then the keyword argument will update the *slicing_dict*.

No checks on the column name is done, should a label be missing, the loc method will raise a *KeyError*.

Parameters

- **input_data_frame** (*pd.DataFrame*) – The data frame to be sliced.
- **slicing_dict** (*dict, Optional*) – A dictionary with columns and values for the slicing. Defaults to *None*
- **kwargs** – Keyword arguments to be used instead of the slicing dictionary.

Returns

The sliced dataframe

Return type

pd.DataFrame

maf.tools.regex

Module implements some basic functions involving regular expressions.

Functions

<code>extract_protocol(url)</code>	Extract the protocol portion from a database connection URL.
<code>normalize_sql_spaces(sql_string)</code>	Normalize multiple consecutive spaces in SQL string to single spaces.

`maf.tools.regex.extract_protocol(url: str) → str | None`

Extract the protocol portion from a database connection URL.

The `extract_protocol` function takes a database connection URL string as input and extracts the protocol portion (the part before “://”). This function is useful for identifying the database type from connection strings.

Parameters

url (*str*) – The url from which the protocol will be extracted.

Returns

The protocol or None, if the extraction failed

Return type

`str | None`

`maf.tools.regex.normalize_sql_spaces(sql_string: str) → str`

Normalize multiple consecutive spaces in SQL string to single spaces. Only handles spaces, preserves other whitespace characters.

Parameters

sql_string (*str*) – The SQL string for space normalization.

Returns

The normalized SQL command.

Return type

`str`

maf.tools.toml_tools

The module provides tools to read / write / modify specific TOML files.

Functions

<code>dump_processor_parameters_to_toml(...)</code>	Dumps a toml file with processor parameters.
<code>generate_steering_file(output_file, processors)</code>	Generates a steering file.
<code>load_steering_file(steering_file[, validate])</code>	Load a steering file for the execution framework.
<code>path_encoder(obj)</code>	Encoder for PathItem.
<code>processor_validator(processors)</code>	Validates that all items in the list are valid processor instances or classes.

Classes

<code>PathItem(t, value, original, trivia)</code>	TOML item representing a Path
---------------------------------------------------	-------------------------------

`class maf.tools.toml_tools.PathItem(t, value, original, trivia)`

Bases: `String`

TOML item representing a Path

`unwrap()` → Path

Returns as pure python object (ppo)

`mafw.tools.toml_tools._add_db_configuration(database_conf: dict[str, Any] | None = None, db_engine: str = 'sqlite', doc: TOMLDocument | None = None) → TOMLDocument`

Add the DB configuration to the TOML document

The expected structure of the `database_conf` dictionary is one of these two:

```
option1 = {
    'DBConfiguration': {
        'URL': 'sqlite:///memory:',
        'pragmas': {
            'journal_mode': 'wal',
            'cache_size': -64000,
            'foreign_keys': 1,
            'synchronous': 0,
        },
    },
}

option2 = {
    'URL': 'sqlite:///memory:',
    'pragmas': {
        'journal_mode': 'wal',
        'cache_size': -64000,
        'foreign_keys': 1,
        'synchronous': 0,
    },
}
```

We will always convert the `option1` in `option2`.

Parameters

- **database_conf** (*dict*) – A dictionary with the database configuration. See comments above. If `None`, then the default is used.
- **db_engine** (*str, Optional*) – The database engine. It is used only in case the provided database configuration is invalid to retrieve the default configuration. Defaults to `sqlite`.
- **doc** (*TOMLDocument, Optional*) – The TOML document to add the DB configuration. If `None`, one will be created.

Returns

The modified document.

Return type

TOMLDocument

Raises

UnknownDBEngine – if the `database_conf` is invalid and the `db_engine` is not yet implemented.

`mafw.tools.toml_tools.dump_processor_parameters_to_toml(processors: list[type[Processor] | Processor] | type[Processor] | Processor, output_file: Path | str) → None`

Dumps a toml file with processor parameters.

This helper function can be used when the parameters of one or many processors have to be dumped to a TOML file. For each Processor in the *processors* a table in the TOML file will be added with their parameters is the shape of parameter name = value.

It must be noted that *processors* can be:

- a list of processor classes (list[type[Processor]])
- a list of processor instances (list[Processor])
- one single processor class (type[Processor])
- one single processor instance (Processor)

What value of the parameters will be dumped?

Good question, have a look at this [explanation](#).

param processors

One or more processors for which the parameters should be dumped.

type processors

list[type[Processor | Processor]] | type[Processor] | Processor

param output_file

The name of the output file for the dump.

type output_file

Path | str

raise KeyAlreadyPresent

if an attempt to add twice, the same processor is made.

raise TypeError

if the list contains items different from Processor classes and instances.

```
mafw.tools.toml_tools.generate_steering_file(output_file: Path | str, processors:
                                             list[type[Processor] | Processor] | type[Processor] |
                                             Processor, database_conf: dict[str, Any] | None =
                                             None, db_engine: str = 'sqlite') → None
```

Generates a steering file.

Parameters

- **output_file** (*Path | str*) – The output filename where the steering file will be save.
- **processors** (*list[type[Processor] | Processor], type[Processor], Processor*) – The processors list for which the steering file will be generated.
- **database_conf** (*dict, Optional*) – The database configuration dictionary
- **db_engine** – A string representing the DB engine to be used. Possible values are: *sqlite*, *postgresql* and *mysql*.

Type

str

```
mafw.tools.toml_tools.load_steering_file(steering_file: Path | str, validate: bool = True) → dict[str,
                                             Any]
```

Load a steering file for the execution framework.

Parameters

- **steering_file** (*Path, str*) – The path to the steering file.
- **validate** (*bool, Optional*) – A flag to validate the content. Defaults to True.

Returns

The configuration dictionary.

Return type

dict

Raises**FileNotFound** – if `steering_file` does not exist.`maf.tools.toml_tools.path_encoder(obj: Any) → Item`

Encoder for PathItem.

`maf.tools.toml_tools.processor_validator(processors: list[type[Processor] | Processor]) → bool`

Validates that all items in the list are valid processor instances or classes.

Parameters**processors** (`list[type[Processor] | Processor]`) – The list of items to be validated.**Returns**

True if all items are valid.

Return type

bool

11.1.16 mafw.ui

User interface modules.

In this subpackage, there are some implementations of possible user interfaces to connect the library to the execution framework.

Modules

<code>abstract_user_interface</code>	An abstract generic user interface.
<code>console_user_interface</code>	The console user interface.
<code>rich_user_interface</code>	The rich user interface.

mafw.ui.abstract_user_interface

An abstract generic user interface.

The module provides a generic user interface that can be implemented to allow MAFw to communicate with different user interfaces.

MAFw is designed to operate seamlessly without a user interface; however, users often appreciate the added benefit of communication between the process execution and themselves.

There are several different interfaces and different interface types (Command Line, Textual, Graphical...) and everyone has its own preferences. In order to be as generic as possible, MAFw is allowing for an abstract interface layer so that the user can either decide to use one of the few coming with MAFw or to implement the interface to their favorite interface.

Classes

<code>UserInterfaceBase()</code>	The abstract base user interface class.
<code>UserInterfaceMeta</code>	A metaclass used for the creation of user interface

class `mafw.ui.abstract_user_interface.UserInterfaceBase`Bases: `object`

The abstract base user interface class.

create_task(*task_name: str, task_description: str = "", completed: int = 0, increment: int | None = None, total: int | None = None, **kwargs: Any*) → None

Create a new task.

Parameters

- **task_name** (*str*) – A unique identifier for the task. You cannot have more than 1 task with the same name in the whole execution. If you want to use the processor name, it is recommended to use the *unique_name*.
- **task_description** (*str, Optional*) – A short description for the task. Defaults to None.
- **completed** (*int, Optional*) – The amount of task already completed. Defaults to None.
- **increment** (*int, Optional*) – How much of the task has been done since last update. Defaults to None.
- **total** (*int, Optional*) – The total amount of task. Defaults to None.

display_progress_message(*message: str, i_item: int, n_item: int | None, frequency: float*) → None

Display a message during the process execution.

Parameters

- **message** (*str*) – The message to be displayed.
- **i_item** (*int*) – The current item enumerator.
- **n_item** (*int | None*) – The total number of items or None for an indeterminate progress (while loop).
- **frequency** (*float*) – How often (in percentage of *n_item*) to display the message.

update_task(*task_name: str, completed: int = 0, increment: int | None = None, total: int | None = None, **kwargs: Any*) → None

Update an existing task.

Parameters

- **task_name** (*str*) – A unique identifier for the task. You cannot have more than 1 task with the same name in the whole execution. If you want to use the processor name, it is recommended to use the *unique_name*.
- **completed** (*int, Optional*) – The amount of task already completed. Defaults to None.
- **increment** (*int, Optional*) – How much of the task has been done since last update. Defaults to None.
- **total** (*int, Optional*) – The total amount of task. Defaults to None.

always_display_progress_message = 10

Threshold for displaying progress messages.

If the total number of events is below this value, then the progress message is always displayed, otherwise follow the standard update frequency.

name = 'base'

The name of the interface

class mafw.ui.abstract_user_interface.UserInterfaceMeta

Bases: type

A metaclass used for the creation of user interface

mafwi.console_user_interface

The console user interface.

The module provides a simple, still efficient user interface ideal for code execution of a headless system where it is not possible to observe the output in real-time. Nevertheless, important messages are logged via the logging library and thus it is also possible to save them to a file, if a proper logging handler is set up.

Classes

<i>ConsoleInterface()</i>	A console user interface.
---------------------------	---------------------------

class mafwi.console_user_interface.ConsoleInterface

Bases: *UserInterfaceBase*

A console user interface.

Ideal for execution in a headless environment.

Messages are sent via the logging system, so they can also be saved to a file if a logging handler is properly set up in the execution framework.

create_task(*task_name: str, task_description: str = "", completed: int = 0, increment: int | None = None, total: int | None = None, **kwargs: Any*) → None

Create a new task.

Parameters

- **task_name** (*str*) – A unique identifier for the task. You cannot have more than 1 task with the same name in the whole execution. If you want to use the processor name, it is recommended to use the *unique_name*.
- **task_description** (*str, Optional*) – A short description for the task. Defaults to None.
- **completed** (*int, Optional*) – The amount of task already completed. Defaults to None.
- **increment** (*int, Optional*) – How much of the task has been done since last update. Defaults to None.
- **total** (*int, Optional*) – The total amount of task. Defaults to None.

display_progress_message(*message: str, i_item: int, n_item: int | None, frequency: float*) → None

Display a message during the process execution.

Parameters

- **message** (*str*) – The message to be displayed.
- **i_item** (*int*) – The current item enumerator.
- **n_item** (*int | None*) – The total number of items or None for an indeterminate progress (while loop).
- **frequency** (*float*) – How often (in percentage of *n_item*) to display the message.

update_task(*task_name: str, completed: int | None = None, increment: int | None = None, total: int | None = None, **kwargs: Any*) → None

Update an existing task.

Parameters

- **task_name** (*str*) – A unique identifier for the task. You cannot have more than one task with the same name in the whole execution. If you want to use the processor name, it is recommended to use the *unique_name*.

- **completed** (*int*, *Optional*) – The amount of task already completed. Defaults to `None`.
- **increment** (*int*, *Optional*) – How much of the task has been done since last update. Defaults to `None`.
- **total** (*int*, *Optional*) – The total amount of task. Defaults to `None`.

`name = 'console'`

The name of the interface

mafwi.rich_user_interface

The rich user interface.

The module provides an implementation of the abstract user interface that takes advantage from the *rich* library. Progress bars and spinners are shown during the processor execution along with log messages including markup language. In order for this logging message to appear properly rendered, the logger should be connected to a `RichHandler`.

Classes

<code>RichInterface([progress_kws])</code>	Implementation of the interface for rich.
--------------------------------------------	-------------------------------------------

class `mafwi.rich_user_interface.RichInterface`(*progress_kws: dict[str, Any] | None = None*)

Bases: `UserInterfaceBase`

Implementation of the interface for rich.

Parameters

progress_kws (*dict*, *Optional*) – A dictionary of keywords passed to the *rich.Progress*. Defaults to `None`

create_task(*task_name: str*, *task_description: str = ""*, *completed: int = 0*, *increment: int | None = None*, *total: int | None = None*, ***kwargs: Any*) → `None`

Create a new task.

Parameters

- **task_name** (*str*) – A unique identifier for the task. You cannot have more than 1 task with the same name in the whole execution. If you want to use the processor name, it is recommended to use the *unique_name*.
- **task_description** (*str*, *Optional*) – A short description for the task. Defaults to `''`.
- **completed** (*int*, *Optional*) – The amount of task already completed. Defaults to `0`.
- **increment** (*int*, *Optional*) – How much of the task has been done since last update. Defaults to `None`.
- **total** (*int*, *Optional*) – The total amount of task. Defaults to `None`.

display_progress_message(*message: str*, *i_item: int*, *n_item: int | None*, *frequency: float*) → `None`

Display a message during the process execution.

Parameters

- **message** (*str*) – The message to be displayed.
- **i_item** (*int*) – The current item enumerator.
- **n_item** (*int | None*) – The total number of items or `None` for an indeterminate progress (while loop).

- **frequency** (*float*) – How often (in percentage of *n_item*) to display the message.

update_task(*task_name: str, completed: int = 0, increment: int | None = None, total: int | None = None, **kwargs: Any*) → None

Update an existing task.

Parameters

- **task_name** (*str*) – A unique identifier for the task. You cannot have more than one task with the same name in the whole execution. If you want to use the processor name, it is recommended to use the *unique_name*.
- **completed** (*int, Optional*) – The amount of task already completed. Defaults to 0.
- **increment** (*int, Optional*) – How much of the task has been done since last update. Defaults to None.
- **total** (*int, Optional*) – The total amount of task. Defaults to None.

name = 'rich'

The name of the interface

a

mafww.active, 106

d

mafww.db, 107
mafww.db.db_configurations, 107
mafww.db.db_filter, 108
mafww.db.db_model, 111
mafww.db.db_types, 113
mafww.db.db_wizard, 114
mafww.db.fields, 115
mafww.db.std_tables, 117
mafww.db.trigger, 120
mafww.decorators, 125

e

mafww.enumerators, 129
mafww.examples, 131
mafww.examples.db_processors, 131
mafww.examples.importer_example, 133
mafww.examples.loop_modifier, 135
mafww.examples.multi_primary, 139
mafww.examples.processor_list, 139
mafww.examples.sum_processor, 139

h

mafww.hookspeccs, 141

m

mafww, 105
mafww.mafww_errors, 141

p

mafww.plugin_manager, 143
mafww.plugins, 144
mafww.processor, 144
mafww.processor_library, 154
mafww.processor_library.abstract_plotter,
154
mafww.processor_library.importer, 159
mafww.processor_library.sns_plotter, 163

r

mafww.runner, 172

s

mafww.scripts, 174
mafww.scripts.mafww_exe, 174
mafww.scripts.update_changelog, 181

t

mafww.timer, 182
mafww.tools, 184
mafww.tools.file_tools, 184
mafww.tools.pandas_tools, 185
mafww.tools.regex, 186
mafww.tools.toml_tools, 187

u

mafww.ui, 190
mafww.ui.abstract_user_interface, 190
mafww.ui.console_user_interface, 192
mafww.ui.rich_user_interface, 193

Symbols

- `_add_db_configuration()` (in module `mafw.tools.toml_tools`), 188
- `_attributes_valid()` (`mafw.processor_library.sns_plotter.FromDatasetDataRetriever` method), 166
- `_attributes_valid()` (`mafw.processor_library.sns_plotter.SQLPdDataRetriever` method), 171
- `_check_method_overload()` (`mafw.processor.Processor` method), 147
- `_check_method_super()` (`mafw.processor.Processor` method), 147
- `_clean_sql()` (`mafw.db.trigger.PostgreSQLDialect` method), 121
- `_configuration_file` (`mafw.processor_library.importer.FilenameParser` attribute), 162
- `_element_dict` (`mafw.processor_library.importer.FilenameParser` attribute), 162
- `_execute_for_loop()` (`mafw.processor.Processor` method), 148
- `_execute_single()` (`mafw.processor.Processor` method), 148
- `_execute_while_loop()` (`mafw.processor.Processor` method), 148
- `_filename` (`mafw.processor_library.importer.FilenameParser` attribute), 162
- `_filename_parser` (`mafw.processor_library.importer.Importer` attribute), 163
- `_generate_next_value_()` (`mafw.db.trigger.TriggerAction` static method), 124
- `_generate_next_value_()` (`mafw.db.trigger.TriggerWhen` static method), 125
- `_generate_next_value_()` (`mafw.enumerators.LoopType` static method), 129
- `_generate_next_value_()` (`mafw.enumerators.ProcessorStatus` static method), 130
- `_get_dialect()` (`mafw.db.trigger.Trigger` method), 122
- `_get_value_type()` (`mafw.processor_library.importer.FilenameElement` class method), 159
- `_ids` (`mafw.processor.Processor` attribute), 151
- `_methods_to_be_checked_for_super` (`mafw.processor.Processor` attribute), 151
- `_parser_configuration()` (`mafw.processor_library.importer.FilenameParser` method), 162
- `_remove_orphan_files()` (`mafw.processor.Processor` method), 148
- `_update_plotter_db()` (`mafw.processor_library.abstract_plotter.GenericPlotter` method), 157
- `_validate_default_type()` (`mafw.processor_library.importer.FilenameElement` method), 160
- `_validate_regexp()` (`mafw.processor_library.importer.FilenameElement` method), 160
- `-D` mafw command line option, 175
- `--db-engine` mafw-steering command line option, 180
- `--db-url` mafw-steering command line option, 180
- `--debug` mafw command line option, 175
- `--engine` mafw-db-wizard command line option, 178
- `--ext-plugins` mafw-list command line option, 178
- `mafw-list` mafw-steering command line option, 179
- `--host` mafw-db-wizard command line option, 177
- `--ignore-unknown` mafw-db-wizard command line option, 177
- `--log-level` mafw command line option, 175
- `--no-ext-plugin` mafw-list command line option, 178
- `mafw-steering` mafw-steering command line option, 179
- `--no-ignore-unknown` mafw-db-wizard command line option, 177
- `--no-open-editor`

mafw-steering command line option, 179
 --no-overwrite
 mafw-db-wizard command line option, 177
 --no-preserve-order
 mafw-db-wizard command line option, 177
 --no-show
 mafw-steering command line option, 179
 --no-snake-case
 mafw-db-wizard command line option, 177
 --open-editor
 mafw-steering command line option, 179
 --output-file
 mafw-db-wizard command line option, 177
 --overwrite
 mafw-db-wizard command line option, 177
 --password
 mafw-db-wizard command line option, 177
 --port
 mafw-db-wizard command line option, 177
 --preserve-order
 mafw-db-wizard command line option, 177
 --schema
 mafw-db-wizard command line option, 177
 --show
 mafw-steering command line option, 179
 --snake-case
 mafw-db-wizard command line option, 177
 --tables
 mafw-db-wizard command line option, 177
 --ui
 mafw command line option, 175
 --user
 mafw-db-wizard command line option, 177
 --username
 mafw-db-wizard command line option, 177
 --version
 mafw command line option, 175
 --with-views
 mafw-db-wizard command line option, 177
 --without-views
 mafw-db-wizard command line option, 177
 -e
 mafw-db-wizard command line option, 178
 -o
 mafw-db-wizard command line option, 177
 -p
 mafw-db-wizard command line option, 177
 -s
 mafw-db-wizard command line option, 177
 -t
 mafw-db-wizard command line option, 177
 -u
 mafw-db-wizard command line option, 177
 -v
 mafw command line option, 175

A

Abort (*mafw.enumerators.LoopingStatus* attribute),

130

Aborted (*mafw.enumerators.ProcessorExitStatus* attribute), 130

AbortProcessorException, 142

accept_item() (*mafw.processor.Processor* method), 148

accessor_class (*mafw.db.fields.FileChecksumField* attribute), 115

accessor_class (*mafw.db.fields.FileNameField* attribute), 116

AccumulatorProcessor (class in *mafw.examples.sum_processor*), 140

acquire_resources() (*mafw.processor.Processor* method), 148

acquire_resources() (*mafw.processor.ProcessorList* method), 153

ActivableType (class in *mafw.active*), 106

Active (class in *mafw.active*), 106

ActiveParameter (class in *mafw.processor*), 144

ActiveType (class in *mafw.active*), 107

add_sql() (*mafw.db.trigger.Trigger* method), 122

add_when() (*mafw.db.trigger.Trigger* method), 123

aggregation_functions (*mafw.processor_library.sns_plotter.SNSPlotter* attribute), 171

always_display_progress_message (*mafw.ui.abstract_user_interface.UserInterfaceBase* attribute), 191

and_() (in module *mafw.db.trigger*), 125

append() (*mafw.processor.ProcessorList* method), 153

B

bind() (*mafw.db.db_filter.Filter* method), 110

bind_all() (*mafw.db.db_filter.FilterRegister* method), 111

C

CatPlot (class in *mafw.processor_library.sns_plotter*), 164

class_depends_on_optional() (in module *mafw.decorators*), 126

commit_changelog_changes() (in module *mafw.scripts.update_changelog*), 182

ConsoleInterface (class in *mafw.ui.console_user_interface*), 192

Continue (*mafw.enumerators.LoopingStatus* attribute), 130

CountStandardTables (class in *mafw.examples.db_processors*), 131

create() (*mafw.db.trigger.Trigger* method), 123

create_table() (*mafw.db.db_model.MAFwBaseModel* class method), 112

create_task() (*mafw.ui.abstract_user_interface.UserInterfaceBase* method), 190

create_task() (*mafw.ui.console_user_interface.ConsoleInterface* method), 192

- create_task() (*mafwi.ui.rich_user_interface.RichInterface* method), 193
 create_trigger_sql() (*mafwi.db.trigger.MySQLDialect* method), 121
 create_trigger_sql() (*mafwi.db.trigger.PostgreSQLDialect* method), 121
 create_trigger_sql() (*mafwi.db.trigger.SQLiteDialect* method), 122
 create_trigger_sql() (*mafwi.db.trigger.TriggerDialect* method), 124
 custom_formatwarning() (in module *mafwi.scripts.mafwi_exe*), 181
 customize_plot() (*mafwi.processor_library.abstract_plotter.GeneratorPlotter* method), 157
- ## D
- data_frame (*mafwi.processor_library.sns_plotter.SNSFigurePlotter* attribute), 168
 data_frame (*mafwi.processor_library.sns_plotter.SNSPlotter* attribute), 171
 DATABASE
 mafwi-db-wizard command line option, 178
 database (*mafwi.processor.Processor* property), 151
 database (*mafwi.processor.ProcessorList* property), 153
 database (*mafwi.processor_library.sns_plotter.SQLPdDataRetriever* attribute), 171
 database_proxy (in module *mafwi.db.db_model*), 113
 database_required() (in module *mafwi.decorators*), 127
 DataRetriever (class in *mafwi.processor_library.abstract_plotter*), 155
 db_scheme (in module *mafwi.db.db_configurations*), 107
 db_value() (*mafwi.db.fields.FileChecksumField* method), 115
 db_value() (*mafwi.db.fields.FileNameField* method), 116
 db_value() (*mafwi.db.fields.FileNameListField* method), 117
 default_conf (in module *mafwi.db.db_configurations*), 108
 delete_parameter() (*mafwi.processor.Processor* method), 148
 depends_on_optional() (in module *mafwi.decorators*), 127
 description (*mafwi.processor.Processor* attribute), 151
 disable() (*mafwi.db.std_tables.TriggerDisabler* method), 120
 display_progress_message() (*mafwi.ui.abstract_user_interface.UserInterfaceBase* method), 191
 display_progress_message() (*mafwi.ui.console_user_interface.ConsoleInterface* method), 192
 display_progress_message() (*mafwi.ui.rich_user_interface.RichInterface* method), 193
 DisPlot (class in *mafwi.processor_library.sns_plotter*), 165
 distribute_resources() (*mafwi.processor.ProcessorList* method), 153
 DoesNotExist (*mafwi.db.db_model.MAFwBaseModel* attribute), 112
 DoesNotExist (*mafwi.db.std_tables.OrphanFile* attribute), 118
 DoesNotExist (*mafwi.db.std_tables.PlotterOutput* attribute), 118
 DoesNotExist (*mafwi.db.std_tables.StandardTable* attribute), 119
 DoesNotExist (*mafwi.db.std_tables.TriggerStatus* attribute), 120
 DoesNotExist (*mafwi.examples.db_processors.File* attribute), 132
 DoesNotExist (*mafwi.examples.importer_example.InputElement* attribute), 135
 drop() (*mafwi.db.trigger.Trigger* method), 123
 drop_trigger_sql() (*mafwi.db.trigger.MySQLDialect* method), 121
 drop_trigger_sql() (*mafwi.db.trigger.PostgreSQLDialect* method), 121
 drop_trigger_sql() (*mafwi.db.trigger.SQLiteDialect* method), 122
 drop_trigger_sql() (*mafwi.db.trigger.TriggerDialect* method), 124
 dump_models() (in module *mafwi.db.db_wizard*), 114
 dump_parameter_configuration() (*mafwi.processor.Processor* method), 149
 dump_processor_parameters_to_toml() (in module *mafwi.tools.toml_tools*), 188
 duration (*mafwi.timer.Timer* property), 183
- ## E
- elements (*mafwi.processor_library.importer.FileNameParser* property), 162
 enable() (*mafwi.db.std_tables.TriggerDisabler* method), 120
 ensure_parameter_registration() (in module *mafwi.processor*), 154
 Error (*mafwi.scripts.mafwi_exe.ReturnValue* attribute), 181
 execute() (*mafwi.processor.Processor* method), 149
 execute() (*mafwi.processor.ProcessorList* method), 153

- execution_workflow() (in module *mafw.decorators*), 127
 exit_status (*mafw.runner.MAFwApplication* attribute), 173
 extend() (*mafw.processor.ProcessorList* method), 153
 extract_protocol() (in module *mafw.tools.regexp*), 187
- ## F
- F (class in *mafw.decorators*), 126
 F (class in *mafw.processor*), 145
 facet_grid(*mafw.processor_library.sns_plotter.SNSFigure* attribute), 168
 facet_grid(*mafw.processor_library.sns_plotter.SNSPlotter* attribute), 171
 Failed (*mafw.enumerators.ProcessorExitStatus* attribute), 130
 field_names (*mafw.db.db_filter.Filter* property), 110
 FigurePlotter (class in *mafw.processor_library.abstract_plotter*), 155
 File (class in *mafw.examples.db_processors*), 132
 file_checksum() (in module *mafw.tools.file_tools*), 184
 FileChecksumField (class in *mafw.db.fields*), 115
 FileChecksumFieldAccessor (class in *mafw.db.fields*), 115
 FileDoesNotExist, 131
 FilenameElement (class in *mafw.processor_library.importer*), 159
 FileNameField (class in *mafw.db.fields*), 116
 FileNameFieldAccessor (class in *mafw.db.fields*), 116
 FileNameListField (class in *mafw.db.fields*), 117
 FilenameParser (class in *mafw.processor_library.importer*), 161
 FillFileTableProcessor (class in *mafw.examples.db_processors*), 132
 Filter (class in *mafw.db.db_filter*), 108
 filter() (*mafw.db.db_filter.Filter* method), 110
 filter_all() (*mafw.db.db_filter.FilterRegister* method), 111
 filter_register (*mafw.processor.Processor* attribute), 151
 FilterRegister (class in *mafw.db.db_filter*), 110
 FindNPrimeNumber (class in *mafw.examples.loop_modifier*), 136
 FindPrimeNumberInRange (class in *mafw.examples.loop_modifier*), 137
 Finish (*mafw.enumerators.ProcessorStatus* attribute), 130
 finish() (*mafw.examples.db_processors.FillFileTableProcessor* method), 133
 finish() (*mafw.examples.importer_example.ImporterExample* method), 134
 finish() (*mafw.examples.loop_modifier.FindNPrimeNumber* method), 136
 finish() (*mafw.examples.loop_modifier.FindPrimeNumberInRange* method), 137
 finish() (*mafw.processor.Processor* method), 149
 finish() (*mafw.processor_library.abstract_plotter.GenericPlotter* method), 157
 for_loop() (in module *mafw.decorators*), 127
 ForLoop (*mafw.enumerators.LoopType* attribute), 129
 format_duration() (*mafw.timer.Timer* method), 183
 format_progress_message() (*mafw.examples.db_processors.FillFileTableProcessor* method), 133
 format_progress_message() (*mafw.examples.importer_example.ImporterExample* method), 135
 format_progress_message() (*mafw.examples.loop_modifier.FindNPrimeNumber* method), 136
 format_progress_message() (*mafw.examples.loop_modifier.FindPrimeNumberInRange* method), 138
 format_progress_message() (*mafw.processor.Processor* method), 149
 format_progress_message() (*mafw.processor_library.abstract_plotter.GenericPlotter* method), 157
 format_progress_message() (*mafw.processor_library.importer.Importer* method), 163
 from_conf() (*mafw.db.db_filter.Filter* class method), 109
 from_dict() (*mafw.processor_library.importer.FilenameElement* class method), 160
 FromDatasetDataRetriever (class in *mafw.processor_library.sns_plotter*), 166
- ## G
- GaussAdder (class in *mafw.examples.sum_processor*), 140
 generate_steering_file() (in module *mafw.tools.toml_tools*), 189
 GenericPlotter (class in *mafw.processor_library.abstract_plotter*), 155
 get_command() (*mafw.scripts.mafw_exe.MAFwGroup* method), 180
 get_data_frame() (*mafw.processor_library.abstract_plotter.DataRetriever* method), 155
 get_data_frame() (*mafw.processor_library.abstract_plotter.GenericPlotter* method), 157
 get_data_frame() (*mafw.processor_library.sns_plotter.FromDatasetDataRetriever* method), 166
 get_data_frame() (*mafw.processor_library.sns_plotter.HDFPdDataRetriever* method), 167
 get_data_frame() (*mafw.processor_library.sns_plotter.PdDataRetriever* method), 167
 get_data_frame() (*mafw.processor_library.sns_plotter.SNSPlotter* method), 170

- get_data_frame() (*mafww.processor_library.sns_plotter.SNSPlotter* method), 171
 get_element() (*mafww.processor_library.importer.FilenameParser* method), 162
 get_element_value() (*mafww.processor_library.importer.FilenameParser* method), 162
 get_field() (*mafww.db.db_filter.Filter* method), 110
 get_filter() (*mafww.processor.Processor* method), 149
 get_items() (*mafww.examples.db_processors.FillFileTableProcessor* method), 133
 get_items() (*mafww.examples.importer_example.ImporterExample* method), 135
 get_items() (*mafww.examples.loop_modifier.FindPrimeNumbers* method), 138
 get_items() (*mafww.examples.loop_modifier.ModifyLoopProcessor* method), 139
 get_items() (*mafww.examples.sum_processor.AccumulatorProcessor* method), 140
 get_items() (*mafww.processor.Processor* method), 149
 get_last_commit_message() (in module *mafww.scripts.update_changelog*), 182
 get_latest_tag() (in module *mafww.scripts.update_changelog*), 182
 get_parameter() (*mafww.processor.Processor* method), 150
 get_parameters() (*mafww.processor.Processor* method), 150
 get_plugin_manager() (in module *mafww.plugin_manager*), 143
 get_user_interface() (*mafww.runner.MAFwApplication* method), 173
 group_and_aggregate_data_frame() (in module *mafww.tools.pandas_tools*), 185
 group_and_aggregate_data_frame() (*mafww.processor_library.sns_plotter.SNSPlotter* method), 170
 grouping_columns (*mafww.processor_library.sns_plotter.SNSPlotter* attribute), 171
- ## H
- HDFPdDataRetriever (class in *mafww.processor_library.sns_plotter*), 166
- ## I
- i_item (*mafww.processor.Processor* property), 151
 Importer (class in *mafww.processor_library.importer*), 162
 ImporterExample (class in *mafww.examples.importer_example*), 133
 in_loop_customization() (*mafww.processor_library.abstract_plotter.GenericPlotter* method), 158
 Init (*mafww.enumerators.ProcessorStatus* attribute), 130
- init() (*mafww.runner.MAFwApplication* method), 173
 InputElement (class in *mafww.examples.importer_example*), 135
 InputElementDoesNotExist, 133
 insert() (*mafww.processor.ProcessorList* method), 153
 interpret() (*mafww.processor_library.importer.FilenameParser* method), 162
 InvalidSteeringFile, 142
 invoke() (*mafww.scripts.mafww_exe.MAFwGroup* method), 181
 is_bound() (*mafww.db.db_filter.Filter* property), 110
 is_data_frame_empty() (*mafww.processor_library.abstract_plotter.GenericPlotter* method), 158
 is_data_frame_empty() (*mafww.processor_library.sns_plotter.SNSPlotter* method), 170
 is_found (*mafww.processor_library.importer.FilenameElement* property), 161
 is_optional (*mafww.processor.PassiveParameter* property), 146
 is_optional (*mafww.processor_library.importer.FilenameElement* property), 161
 is_output_existing() (*mafww.processor_library.abstract_plotter.GenericPlotter* method), 158
 is_prime() (in module *mafww.examples.loop_modifier*), 139
 is_set (*mafww.processor.PassiveParameter* property), 146
 item (*mafww.processor.Processor* attribute), 151
- ## L
- load_steering_file() (in module *mafww.tools.toml_tools*), 189
 local_resource_acquisition (*mafww.processor.Processor* property), 151
 logger_setup() (in module *mafww.scripts.mafww_exe*), 181
 loop_type (*mafww.processor.Processor* attribute), 152
 looping_status (*mafww.processor.Processor* attribute), 152
 LoopingStatus (class in *mafww.enumerators*), 129
 LoopType (class in *mafww.enumerators*), 129
- ## M
- mafw module, 105
 mafw command line option
 -D, 175
 -debug, 175
 --log-level, 175
 --ui, 175
 --version, 175

- v, 175
- mafwb.active
 - module, 106
- mafwb.db
 - module, 107
- mafwb.db.db_configurations
 - module, 107
- mafwb.db.db_filter
 - module, 108
- mafwb.db.db_model
 - module, 111
- mafwb.db.db_types
 - module, 113
- mafwb.db.db_wizard
 - module, 114
- mafwb.db.fields
 - module, 115
- mafwb.db.std_tables
 - module, 117
- mafwb.db.trigger
 - module, 120
- mafwb.decorators
 - module, 125
- mafwb.enumerators
 - module, 129
- mafwb.examples
 - module, 131
- mafwb.examples.db_processors
 - module, 131
- mafwb.examples.importer_example
 - module, 133
- mafwb.examples.loop_modifier
 - module, 135
- mafwb.examples.multi_primary
 - module, 139
- mafwb.examples.processor_list
 - module, 139
- mafwb.examples.sum_processor
 - module, 139
- mafwb.hookspecs
 - module, 141
- mafwb.mafwb_errors
 - module, 141
- mafwb.plugin_manager
 - module, 143
- mafwb.plugins
 - module, 144
- mafwb.processor
 - module, 144
- mafwb.processor_library
 - module, 154
- mafwb.processor_library.abstract_plotter
 - module, 154
- mafwb.processor_library.importer
 - module, 159
- mafwb.processor_library.sns_plotter
 - module, 163
- mafwb.runner
 - module, 172
- mafwb.scripts
 - module, 174
- mafwb.scripts.mafwb_exe
 - module, 174
- mafwb.scripts.update_changelog
 - module, 181
- mafwb.timer
 - module, 182
- mafwb.tools
 - module, 184
- mafwb.tools.file_tools
 - module, 184
- mafwb.tools.pandas_tools
 - module, 185
- mafwb.tools.regex
 - module, 186
- mafwb.tools.toml_tools
 - module, 187
- mafwb.ui
 - module, 190
- mafwb.ui.abstract_user_interface
 - module, 190
- mafwb.ui.console_user_interface
 - module, 192
- mafwb.ui.rich_user_interface
 - module, 193
- mafwb_hookimpl (*in module mafwb*), 105
- mafwb-db-wizard command line option
 - engine, 178
 - host, 177
 - ignore-unknown, 177
 - no-ignore-unknown, 177
 - no-overwrite, 177
 - no-preserve-order, 177
 - no-snake-case, 177
 - output-file, 177
 - overwrite, 177
 - password, 177
 - port, 177
 - preserve-order, 177
 - schema, 177
 - snake-case, 177
 - tables, 177
 - user, 177
 - username, 177
 - with-views, 177
 - without-views, 177
 - e, 178
 - o, 177
 - p, 177
 - s, 177
 - t, 177
 - u, 177
- DATABASE, 178
- mafwb-list command line option
 - ext-plugins, 178
 - no-ext-plugin, 178

- mafwr-run command line option
 - STEERING_FILE, 179
 - mafwr-steering command line option
 - db-engine, 180
 - db-url, 180
 - ext-plugins, 179
 - no-ext-plugin, 179
 - no-open-editor, 179
 - no-show, 179
 - open-editor, 179
 - show, 179
 - STEERING_FILE, 180
 - MAFwApplication (class in mafwr.runner), 172
 - MAFwBaseModel (class in mafwr.db.db_model), 112
 - MAFwBaseModelDoesNotExist, 111
 - MAFwException, 142
 - MAFwGroup (class in mafwr.scripts.mafwr_exe), 180
 - main() (in module mafwr.scripts.update_changelog), 182
 - main() (mafwr.scripts.mafwr_exe.MAFwGroup method), 181
 - matplotlib_backend (mafwr.processor_library.sns_plotter.SNSPlotter attribute), 171
 - MissingAttribute, 142
 - MissingDatabase, 142
 - MissingOptionalDependency, 142
 - MissingOverloadedMethod, 142
 - MissingSQLStatement, 142
 - MissingSuperCall, 142
 - ModelError, 143
 - ModifyLoopProcessor (class mafwr.examples.loop_modifier), 138
 - module
 - mafwr, 105
 - mafwr.active, 106
 - mafwr.db, 107
 - mafwr.db.db_configurations, 107
 - mafwr.db.db_filter, 108
 - mafwr.db.db_model, 111
 - mafwr.db.db_types, 113
 - mafwr.db.db_wizard, 114
 - mafwr.db.fields, 115
 - mafwr.db.std_tables, 117
 - mafwr.db.trigger, 120
 - mafwr.decorators, 125
 - mafwr.enumerators, 129
 - mafwr.examples, 131
 - mafwr.examples.db_processors, 131
 - mafwr.examples.importer_example, 133
 - mafwr.examples.loop_modifier, 135
 - mafwr.examples.multi_primary, 139
 - mafwr.examples.processor_list, 139
 - mafwr.examples.sum_processor, 139
 - mafwr.hookspeccs, 141
 - mafwr.mafwr_errors, 141
 - mafwr.plugin_manager, 143
 - mafwr.plugins, 144
 - mafwr.processor, 144
 - mafwr.processor_library, 154
 - mafwr.processor_library.abstract_plotter, 154
 - mafwr.processor_library.importer, 159
 - mafwr.processor_library.sns_plotter, 163
 - mafwr.runner, 172
 - mafwr.scripts, 174
 - mafwr.scripts.mafwr_exe, 174
 - mafwr.scripts.update_changelog, 181
 - mafwr.timer, 182
 - mafwr.tools, 184
 - mafwr.tools.file_tools, 184
 - mafwr.tools.pandas_tools, 185
 - mafwr.tools.regex, 186
 - mafwr.tools.toml_tools, 187
 - mafwr.ui, 190
 - mafwr.ui.abstract_user_interface, 190
 - mafwr.ui.console_user_interface, 192
 - mafwr.ui.rich_user_interface, 193
 - MySQLDialect (class in mafwr.db.trigger), 121
- ## N
- n_item (mafwr.processor.Processor property), 152
 - name (mafwr.processor.Processor attribute), 152
 - name (mafwr.processor.ProcessorList property), 153
 - name (mafwr.processor_library.importer.FilenameElement property), 161
 - name (mafwr.runner.MAFwApplication attribute), 173
 - name (mafwr.ui.abstract_user_interface.UserInterfaceBase attribute), 191
 - in name (mafwr.ui.console_user_interface.ConsoleInterface attribute), 193
 - name (mafwr.ui.rich_user_interface.RichInterface attribute), 194
 - new_only (mafwr.db.db_filter.FilterRegister property), 111
 - normalize_sql_spaces() (in module mafwr.tools.regex), 187
- ## O
- OK (mafwr.scripts.mafwr_exe.ReturnValue attribute), 181
 - on_looping_status_set() (mafwr.processor.Processor method), 150
 - on_processor_status_change() (mafwr.processor.Processor method), 150
 - or_() (in module mafwr.db.trigger), 125
 - orphan_protector() (in module mafwr.decorators), 127
 - OrphanFile (class in mafwr.db.std_tables), 118
 - OrphanFileDoesNotExist, 118
 - output_filename_list (mafwr.processor_library.sns_plotter.SNSPlotter attribute), 171
- ## P
- P (class in mafwr.decorators), 126
 - ParameterType (class in mafwr.processor), 145

- ParserConfigurationError, 143
- ParsingError, 143
- PassiveParameter (class in mafw.processor), 146
- patch_data_frame() (mafw.processor_library.abstract_plotter.DataRetriever method), 155
- patch_data_frame() (mafw.processor_library.abstract_plotter.GenericPlotter method), 158
- patch_data_frame() (mafw.processor_library.sns_plotter.HDFPdDataRetriever method), 167
- patch_data_frame() (mafw.processor_library.sns_plotter.PdDataRetriever method), 167
- path_encoder() (in module mafw.tools.toml_tools), 190
- PathItem (class in mafw.tools.toml_tools), 187
- pattern (mafw.processor_library.importer.FileNameElement property), 161
- PdDataRetriever (class in mafw.processor_library.sns_plotter), 167
- PeeweeModelWithMeta (class in mafw.db.db_types), 113
- plot() (mafw.processor_library.abstract_plotter.GenericPlotter method), 158
- plot() (mafw.processor_library.sns_plotter.CatPlot method), 165
- plot() (mafw.processor_library.sns_plotter.DisPlot method), 166
- plot() (mafw.processor_library.sns_plotter.RelPlot method), 168
- PlotterMeta (class in mafw.processor_library.abstract_plotter), 158
- PlotterMixinNotInitialized, 143
- PlotterOutput (class in mafw.db.std_tables), 118
- PlotterOutputDoesNotExist, 118
- plugin_manager (mafw.runner.MAFwApplication attribute), 173
- PostgreSQLDialect (class in mafw.db.trigger), 121
- pretty_format_duration() (in module mafw.timer), 183
- prime_num_found (mafw.examples.loop_modifier.FindNPrimeNumber attribute), 137
- prime_num_found (mafw.examples.loop_modifier.FindPrimeNumberInRange attribute), 138
- print_process_statistics() (mafw.processor.Processor method), 150
- process() (mafw.examples.db_processors.CountStandardTables method), 132
- process() (mafw.examples.db_processors.FillFileTableProcessor method), 133
- process() (mafw.examples.importer_example.ImporterExample method), 135
- process() (mafw.examples.loop_modifier.FindNPrimeNumber method), 136
- process() (mafw.examples.loop_modifier.FindPrimeNumberInRange method), 138
- process() (mafw.examples.loop_modifier.ModifyLoopProcessor method), 139
- process() (mafw.examples.sum_processor.AccumulatorProcessor method), 140
- process() (mafw.examples.sum_processor.GaussAdder method), 140
- process() (mafw.processor.Processor method), 150
- process() (mafw.processor_library.abstract_plotter.GenericPlotter method), 158
- process() (mafw.processor_library.sns_plotter.SNSPlotter method), 170
- Processor (class in mafw.processor), 147
- processor_depends_on_optional() (in module mafw.decorators), 127
- processor_exit_status (mafw.processor.Processor attribute), 152
- processor_exit_status (mafw.processor.ProcessorList property), 153
- processor_status (mafw.processor.Processor attribute), 152
- processor_validator() (in module mafw.tools.toml_tools), 190
- ProcessorExitStatus (class in mafw.enumerators), 130
- ProcessorList (class in mafw.processor), 152
- ProcessorMeta (class in mafw.processor), 154
- ProcessorParameterError, 143
- ProcessorStatus (class in mafw.enumerators), 130
- progress_message (mafw.processor.Processor attribute), 152
- python_value() (mafw.db.fields.FileNameField method), 116
- python_value() (mafw.db.fields.FileNameListField method), 117
- ## Q
- Quit (mafw.enumerators.LoopingStatus attribute), 130
- ## R
- register_processors() (in module mafw.hookspecs), 141
- register_processors() (in module mafw.plugins), 144
- register_standard_tables() (in module mafw.hookspecs), 141
- register_standard_tables() (in module mafw.plugins), 144
- register_user_interfaces() (in module mafw.hookspecs), 141
- register_user_interfaces() (in module mafw.plugins), 144
- RelPlot (class in mafw.processor_library.sns_plotter), 167
- remove_orphan_files (mafw.processor.Processor attribute), 152

remove_widow_db_rows() (in module StandardTableDoesNotExist, 118
mafww.tools.file_tools), 184
 start (*mafww.enumerators.ProcessorStatus* attribute),
 required_columns (*mafww.processor_library.sns_plotter.SQLPdDataRetriever*
 attribute), 171
 start () (*mafww.examples.db_processors.CountStandardTables*
 method), 132
 reset () (*mafww.processor_library.importer FilenameElement*
 method), 160
 start () (*mafww.examples.db_processors.FillFileTableProcessor*
 method), 133
 reset () (*mafww.processor_library.importer FilenameParser*
 method), 162
 start () (*mafww.examples.importer_example.ImporterExample*
 method), 135
 ReturnValue (class in *mafww.scripts.mafww_exe*), 181
 start () (*mafww.examples.loop_modifier.FindNPrimeNumber*
 method), 136
 RichInterface (class in *mafww.ui.rich_user_interface*), 193
 start () (*mafww.examples.loop_modifier.FindPrimeNumberInRange*
 method), 138
 rreplace () (in module *mafww.timer*), 183
 start () (*mafww.examples.loop_modifier.ModifyLoopProcessor*
 method), 139
 Run (*mafww.enumerators.ProcessorStatus* attribute), 131
 start () (*mafww.examples.sum_processor.AccumulatorProcessor*
 method), 140
 run () (*mafww.runner.MAFwApplication* method), 173
 start () (*mafww.examples.sum_processor.GaussAdder*
 method), 141
 run_processor_list_with_loop_modifier () (in
 module *mafww.examples.processor_list*), 139
 start () (*mafww.processor.Processor* method), 151
 run_simple_processor_list () (in module
mafww.examples.processor_list), 139
 start () (*mafww.processor_library.importer.Importer*
 method), 163
 RunnerNotInitialized, 143
 start () (*mafww.processor_library.sns_plotter.SNSPlotter*
 method), 171
 std_upsert () (*mafww.db.db_model.MAFwBaseModel*
 class method), 112
 std_upsert_many () (*mafww.db.db_model.MAFwBaseModel*
 class method), 113
 STEERING_FILE
 mafww-run command line option, 179
 mafww-steering command line option, 180
 Successful (*mafww.enumerators.ProcessorExitStatus*
 attribute), 130
 supports_safe_create ()
 (*mafww.db.trigger.MySQLDialect* method),
 121
 supports_safe_create ()
 (*mafww.db.trigger.PostgreSQLDialect*
 method), 121
 supports_safe_create ()
 (*mafww.db.trigger.SQLiteDialect* method),
 122
 supports_safe_create ()
 (*mafww.db.trigger.TriggerDialect* method),
 124
 supports_trigger_type ()
 (*mafww.db.trigger.MySQLDialect* method),
 121
 supports_trigger_type ()
 (*mafww.db.trigger.PostgreSQLDialect*
 method), 121
 supports_trigger_type ()
 (*mafww.db.trigger.SQLiteDialect* method),
 122
 supports_trigger_type ()
 (*mafww.db.trigger.TriggerDialect* method),
 124
 supports_update_of_columns ()

- (*mafw.db.trigger.MySQLDialect* method), 121
- supports_update_of_columns()* (*mafw.db.trigger.PostgreSQLDialect* method), 121
- supports_update_of_columns()* (*mafw.db.trigger.SQLiteDialect* method), 122
- supports_update_of_columns()* (*mafw.db.trigger.TriggerDialect* method), 124
- supports_when_clause()* (*mafw.db.trigger.MySQLDialect* method), 121
- supports_when_clause()* (*mafw.db.trigger.PostgreSQLDialect* method), 121
- supports_when_clause()* (*mafw.db.trigger.SQLiteDialect* method), 122
- supports_when_clause()* (*mafw.db.trigger.TriggerDialect* method), 124
- suppress_warnings()* (in module *mafw.decorators*), 128
- ## T
- table_name* (*mafw.processor_library.sns_plotter.SQLPdDataRetriever* attribute), 172
- Timer* (class in *mafw.timer*), 182
- Trigger* (class in *mafw.db.trigger*), 122
- TriggerAction* (class in *mafw.db.trigger*), 124
- TriggerDialect* (class in *mafw.db.trigger*), 124
- TriggerDisabler* (class in *mafw.db.std_tables*), 119
- triggers()* (*mafw.db.db_model.MAFwBaseModel* class method), 113
- triggers()* (*mafw.db.std_tables.PlotterOutput* class method), 119
- TriggerStatus* (class in *mafw.db.std_tables*), 120
- TriggerStatusDoesNotExist*, 118
- TriggerWhen* (class in *mafw.db.trigger*), 125
- type_lut* (*mafw.processor_library.importer.FileNameElement* attribute), 161
- ## U
- unique_id* (*mafw.processor.Processor* attribute), 152
- unique_name* (*mafw.processor.Processor* property), 152
- Unknown* (*mafw.enumerators.ProcessorStatus* attribute), 131
- UnknownDBEngine*, 143
- UnknownField* (class in *mafw.db.db_wizard*), 114
- UnknownProcessor*, 143
- UnsupportedDatabaseError*, 143
- unwrap()* (*mafw.tools.toml_tools.PathItem* method), 187
- update_db()* (*mafw.processor_library.abstract_plotter.GenericPlotter* method), 158
- update_task()* (*mafw.ui.abstract_user_interface.UserInterfaceBase* method), 191
- update_task()* (*mafw.ui.console_user_interface.ConsoleInterface* method), 192
- update_task()* (*mafw.ui.rich_user_interface.RichInterface* method), 194
- UserInterfaceBase* (class in *mafw.ui.abstract_user_interface*), 190
- UserInterfaceMeta* (class in *mafw.ui.abstract_user_interface*), 191
- ## V
- validate_database_conf()* (in module *mafw.processor*), 154
- validate_item()* (*mafw.processor.ProcessorList* static method), 153
- validate_items()* (*mafw.processor.ProcessorList* static method), 153
- value* (*mafw.processor.PassiveParameter* property), 146
- value* (*mafw.processor_library.importer.FileNameElement* property), 161
- verify_checksum()* (in module *mafw.tools.file_tools*), 185
- ## W
- where_clause* (*mafw.processor_library.sns_plotter.SQLPdDataRetriever* attribute), 172
- while_condition()* (*mafw.examples.loop_modifier.FindNPrimeNumber* method), 137
- while_condition()* (*mafw.processor.Processor* method), 151
- while_loop()* (in module *mafw.decorators*), 129
- WhileLoop* (*mafw.enumerators.LoopType* attribute), 129