

Machina

None

Table of contents

1. Welcome to the Machina documentation	3
2. Installation	4
2.1 System	4
2.2 CLI	8
3. Administration	11
3.1 Start	11
3.2 Scale	11
3.3 Stop	11
3.4 System Stats	11
3.5 Services	11
4. Usage	12
4.1 Submission	12
4.2 Examples	13
5. Workers	14
6. Development	15
6.1 Worker development	15
6.2 Ghidra Worker development	18
6.3 Periodic Worker development	20
6.4 API	22

1. Welcome to the Machina documentation

Machina is a scalable and modular analysis framework. Machina enables the rapid integration of both existing open source and novel analysis as worker modules.

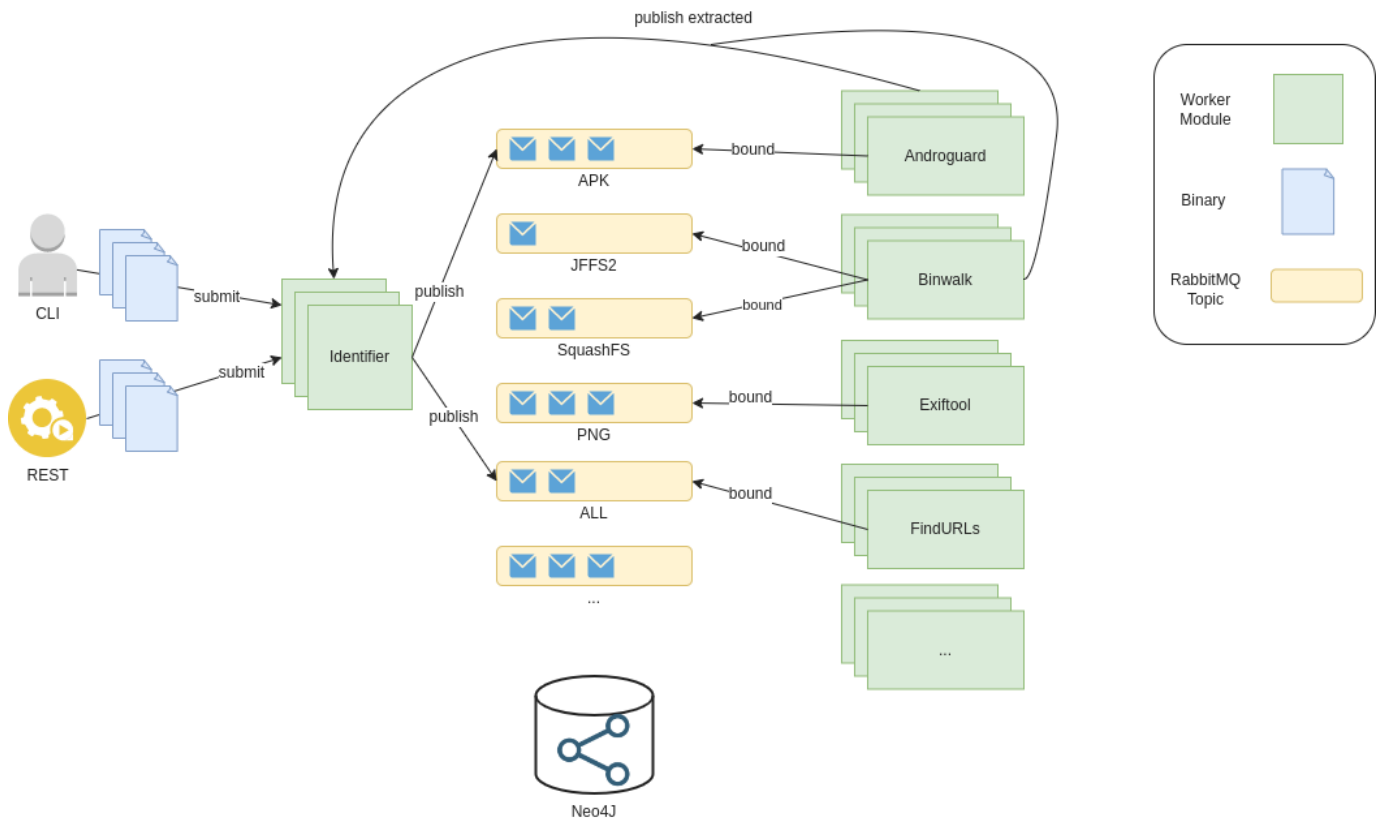
Machina's key features

- **automated identification of input data** - Machina's Identifier worker module classifies input data. Optionally, a subject matter expert can provide a data type with their data if they wish
- **modular worker development** - A Worker base class provides boilerplate functionality, enabling developers to focus only on their analysis implementation
- **recursive analysis** - Worker modules may publish data back to the entrypoint of the system (Identifier). This allows for recursive extraction (firmware, archives, compressed data) of data, or re-typing of data after initial triage.
- **graph storage** - Machina uses a graph storage to store input and discovered analysis objects (as nodes), and their relations with other nodes as edges. For example: edges between nodes can represent origin ("extracted from") or similarity

Machina uses RabbitMQ for message passing into and within the system.

Each Machina Worker Module gets its own queue for direct publishing of data. More importantly, each Worker module also subscribes to topic(s) which correspond to abstract file types (e.g. an apk, pe, elf, etc..). This allows for a Worker Module to receive and handle any or all supported abstract file types.

When data is submitted to the pipeline (directly to the Identifier queue), the Identifier determines its type and publishes a message with the classified type as the routing key. This message is consumed by all Worker Modules that are configured to support that type. The Identifier prioritizes "detailed type" data over "mimetype" for classification.



This documentation is also available in [PDF](#)

2. Installation

2.1 System

2.1.1 Dependencies

- Tested on Ubuntu 20.04+
- git
- [docker](#)
- See [Post-Installation steps](#) to ensure Docker can be run without 'sudo'
- Ensure to install the 'docker-compose-plugin' that's described in the installation guide above

2.1.2 Production

Clone

```
1 git clone https://github.com/ehrenb/machina.git
```

Update/Pull

```
1 docker compose pull
```

 **Note**

By default the 'latest' stable image versions will be used to pull. To change to a specific version, modify the 'latest' tag in the docker-compose.yml file cloned

2.1.3 Development

Clone

```
1 git clone --recurse-submodules https://github.com/ehrenb/machina.git &&\
2 cd machina &&\
3 git submodule foreach git checkout main &&\
4 git submodule foreach git pull
```

Build

Build in the proper order to ensure that parent base images are built first:

```
1 docker compose build base-alpine base-ubuntu &&\
2 docker compose build base-ghidra &&\
3 docker compose build
```

2.1.4 Validate

Validate built or pulled images:

```
1 docker images | grep machina
```

Output:

```

1 behren/machina-docs          latest          7b6f5d4fb34f  6 days ago    601MB
2 behren/machina-androguard    latest         de457208b6c2  7 days ago    766MB
3 behren/machina-binwalk       latest         829efd9691a3  7 days ago    1.8GB
4 behren/machina-findurls      latest         b09157ce035a  7 days ago    524MB
5 behren/machina-similarity     latest         287bc579e7c0  7 days ago    657MB
6 behren/machina-ssdeep        latest         226791c470f0  7 days ago    657MB
7 behren/machina-identifier     latest         56d18aaeb1d4  7 days ago    446MB
8 behren/machina-exif          latest         d737587a3d89  7 days ago    446MB
9 behren/machina-ghidra-project-creator latest         9104987169d1  7 days ago    2.72GB
10 behren/machina-bz2           latest         56262d5b591f  7 days ago    446MB
11 behren/machina-zip            latest         65c8cd3611db  7 days ago    446MB
12 behren/machina-tar           latest         9bb581d6779c  7 days ago    446MB
13 behren/machina-jar            latest         3a05ae24f59c  7 days ago    446MB
14 behren/machina-gzip           latest         2a0669038534  7 days ago    446MB
15 behren/machina-initializer    latest         9ae05328c766  7 days ago    446MB
16 behren/machina-base-ghidra    latest         8312415a7460  7 days ago    2.72GB
17 behren/machina-base-ubuntu    latest         679df4f10abf  7 days ago    648MB
18 behren/machina-base-alpine    latest         a97f8b394af8  7 days ago    446MB

```

2.2 CLI

2.2.1 Dependencies

It is recommended to set up and use a Python 3 virtual environment to isolate dependency installation. Install [virtualenv](#).

```
1 virtualenv -p python3 machina-cli
2 source machina-cli/bin/activate
```


2.2.2 CLI

Clone

```
1 git clone https://github.com/ehrenb/machina.git
```

Install

Within your virtualenv:

```
1 cd machina/cli/  
2 python3 setup.py install
```

2.2.3 Validate

```
1 python3 bin/machinacli.py --help
```

Output:

```
1 usage: machinacli.py [-h] [--verbose] {submit} ...
2
3 positional arguments:
4 {submit}
5     submit            submit a file
6
7 options:
8 -h, --help          show this help message and exit
9 --verbose, -v
```

3. Administration

3.1 Start

Start the system in the background:

```
1 docker compose up -d
```

3.2 Scale

Scale up worker modules to support parallel analyses

```
1 docker compose up -d --scale identifier=2 androguardanalysis=5
```

3.3 Stop

```
1 docker compose down
```

3.4 System Stats

```
1 docker stats $(docker compose ps | awk 'NR>2 {print $1}')
```

3.5 Services

- Neo4j GUI
- <http://127.0.0.1:7474>
- (default) username: neo4j
- (default) password: tXOCq81bn7QfGTMJMrkQqP4J1
- RabbitMQ Management GUI
- <http://127.0.0.1:15672>
- (default) username: rabbitmq
- (default) password: rabbitmq

4. Usage

4.1 Submission

4.1.1 CLI

`machinacli.py` can be used to submit either a single file or a batch of individual files:

```
1 python3 bin/machinacli.py submit /usr/bin/ls /usr/bin/ps
```

Alternatively, if there are too many files to conveniently list as command line options, files can be compressed or archived and submitted. Machina's decompression/unarchiving modules would then handle the unpacking and resubmission automatically to the system:

```
1 zip test.zip /usr/bin/ls /usr/bin/ps
2 python3 bin/machinacli.py submit test.zip
```

As analyses are completed, they are viewable in the Neo4J dashboard described within the 'Administration' section.

4.1.2 JSON

Optionally, you can publish your own message to the RabbitMQ Server, with the routing key set to 'Identifier' and a JSON body containing a base-64 encoded payload like the following:

```
1 {"data": "<b64encoded_data>"}
```

Or, to assert a type (must be available in the 'available_types' configuration within `machina/configs/types.json`), provide it in the 'type' key-value pair. Providing this key forces the Identifier to skip type resolution and accept your own:

```
1 {
2   "data": "<b64encoded_data>",
3   "type": "apk"
4 }
```

4.2 Examples

4.2.1 Firmware (squashfs filesystem)

Obtain a sample of a squashfs file system:

```
1 wget https://github.com/ehrenb/machina-test/raw/main/data/squashfs/firmware.squashfs
```

Submit the file:

```
1 python3 bin/machinacli.py submit firmware.squashfs
```

After several minutes, all files in the submitted squashfs file are extracted and stored in Neo4J:

The screenshot shows the Neo4j web interface. On the left, the 'Database Information' sidebar is visible, showing the database 'machina' selected. Under 'Node Labels', there are 1769 nodes with various file types like APK, Artifact, BZ2, CPIO, Dex, Elf, Emf, Excel, Gzip, HTML, JFFS2, JPEG, Jar, LZMA, Macho, MemoryDump, Msg, PDF, PE, PNG, Powerpoint, RTF, SquashFS, TIFF, Tar, URL, Word, and Zip. Under 'Relationship Types', there are 1,872 relationships of types EXTRACTS and SIMILAR. The main area shows a graph visualization of the data, with a large number of nodes (1,605) and relationships (0). The graph is a dense, circular cluster of nodes, with a color gradient from blue to yellow. The 'Overview' panel on the right shows the following statistics:

- Node labels: (1605) SquashFS (1), Artifact (1275), HTML (67), Elf (248), PNG (12), JPEG (2)
- Relationship Types: (1684) EXTRACTS (1604), SIMILAR (80)
- Displaying 1,605 nodes, 0 relationships.

5. Workers

Name	Description	Input Types	Repository
androguard	Extract metadata from an Android Application (APK). Resubmit extracted files	apk	https://github.com/ehrenb/machina-androguard
binwalk	Unpack known firmware formats and resubmit extracted files	jffs2, squashfs, lzma, cpio	https://github.com/ehrenb/machina-binwalk
bz2	Decompress bzip2 data and resubmit extracted files	bz2	https://github.com/ehrenb/machina-bz2
exif	Extract exifdata from image files	png, jpeg, tiff	https://github.com/ehrenb/machina-exif
findurls	Extract URLs from any kind of data	all except url	https://github.com/ehrenb/machina-findurls
ghidra-project-creator	Perform default Ghidra analysis on executable files	elf, pe	https://github.com/ehrenb/machina-ghidra-project-creator
gzip	Decompress gzip data and resubmit extracted files	gz	https://github.com/ehrenb/machina-gzip
identifier	Classify input data and apply Machina type data	all	https://github.com/ehrenb/machina-identifier
jar	Detect APK from JAR format and resubmit data as APK	jar	https://github.com/ehrenb/machina-jar
similarity	Analyze previously determined SSDeep hash with a configured threshold to create a 'similarity' edge between nodes	None, invoked at interval	https://github.com/ehrenb/machina-similarity
ssdeep	Calculate SSDeep hash of file	all	https://github.com/ehrenb/machina-ssdeep
tar	Unarchive a tar file, resubmit all extracted files for analysis	tar	https://github.com/ehrenb/machina-tar
zip	Unzip a zip file, resubmit all extracted files for analysis	zip	https://github.com/ehrenb/machina-zip

6. Development

6.1 Worker development

This type of Machina worker triggers an analysis when new files are ingested into the system and tagged with a worker-compatible Machina type.

6.1.1 Dockerfile

Install any system dependencies required within your worker's Dockerfile. There are two base image options provided:

Base image options:

- [behren/machina-base-alpine](#)
- [behren/machina-base-ubuntu](#)

Dockerfile example:

```
1 FROM behren/machina-base-ubuntu:latest
2 ...
3 RUN apt update && apt install libz-dev
4 ...
```

6.1.2 requirements.txt

Put any Python 3 requirements required by your worker module into requirements.txt and ensure that you copy requirements.txt into the image and 'pip3 install -r requirements.txt' to install the dependencies.

6.1.3 youranalysismodule.py

This file contains the implementation of your worker module. subclass the [Worker](#) class, this will ensure your worker module has boilerplate connectivity to the database, RabbitMQ, and configurations. Choose any Machina types (see 'machina/configs/types.json') your worker module supports, or specify '*' for all. The [callback](#) function provides your analysis implementation with data that your module is configured to support. This callback function fires whenever the system identifies a compatible sample.

```
1 class YourAnalysisModule(Worker):
2     types = ["zip"]
3
4     def __init__(self, *args, **kwargs):
5         super(YourAnalysisModule, self).__init__(*args, **kwargs)
6         ...
7
8     def callback(self, data, properties):
9         data = json.loads(data)
```

```
1 class YourAnalysisModule(Worker):
2     types = ["*"]
3
4     def __init__(self, *args, **kwargs):
5         super(YourAnalysisModule, self).__init__(*args, **kwargs)
6         ...
7
8     def callback(self, data, properties):
9         data = json.loads(data)
```

```
1 class YourAnalysisModule(Worker):
2     types_blacklist = ["zip"]
3
4     def __init__(self, *args, **kwargs):
5         super(YourAnalysisModule, self).__init__(*args, **kwargs)
6         ...
7
8     def callback(self, data, properties):
9         data = json.loads(data)
```



If 'behren/machina-base-ghidra' was selected as your base, and Pythonic access to Ghidra is desired, see the Ghidra Worker Development documentation

6.1.4 YourAnalysisModule.json (schema)

This schema file provides validation constraints that are applied to data incoming to your worker module before it handles the data. The Schema name must match the class name that it belongs to (e.g. for the worker module 'AndroguardAnalysis'). This file belongs at the top level of your worker's directory, and must be copied within the Dockerfile.

Typically, since workers are handling data published by the Identifier, they inherit from the 'binary.json' schema. Additional input requirements can be specified in "properties"

```
1 {
2   "allof": [{ "$ref": "binary.json" }],
3   "properties": {}
4 }
```

6.1.5 YourAnalysisModule.json (configuration)

This top-level configuration file belongs in machina/configs/workers/youranalysismodule.json. This file allows for reconfiguration without rebuilding of images or code. This file must be named after the worker class name that it corresponds to. Configuration data set in this file is made available through the worker module's 'self.config["worker"]' attribute. Log level is handled by the [Worker](#) base class to automatically adjust the subclass logging level if it is overridden in the configuration.

```
1 {
2   "log_level": "debug",
3   "hash_algorithms": ["md5", "sha256"]
4 }
```

```
1 class YourAnalysisModule(Worker):
2     types = ["zip"]
3     ...
4     def callback(self, data, properties):
5         self.logger.info(self.config['worker']['hash_algorithms'])
```

6.1.6 Other notes

Republishing

Worker modules are not intended to create new nodes (e.g. files, binary data) in the database directly, only update node attributes or create edges (relationships). They should publish any extracted data of interest to the Identifier queue so that it re-enters the pipeline, e.g.:

```
1 class YourAnalysisModule(Worker):
2     next_queues = ['Identifier']
3     ...
4
5     def callback(self, data, properties):
6         ...
7         self.publish_next(json.dumps(data)) # publish to queues configured in 'next_queues'
```

OR

```
1 class MyWorker(Worker):
2     ...
3     def callback(self, data, properties):
4         ...
5         self.publish(json.dumps(data), queues=['Identifier']) # publish to 'Identifier'
```


Retyping

File typing through mimetypes or file magic is not always granular enough to accurately determine a type. Sometimes it requires a bit of context, e.g. an Android APK is technically a zip file, and can only really be identified by peering into the zip and searching for common APK files. Only then can we retype the file properly as an APK. This burden should be on the Zip module to discover, not the Identifier.

The snippet below is an example of when the Zip analysis module detects that it is actually working on an APK. The Zip module resubmits most of the same data that consumed from the queue, except it manually specifies the 'type' to 'apk', which the Identifier will take at face value.

```
1 def callback(self, data, properties):
2     ...
3     body = {
4         "data": data_encoded,
5         "origin": {
6             "ts": data['ts'],
7             "md5": data['hashes']['md5'],
8             "uid": data['uid'],
9             "type": data['type']
10        },
11        'type': 'apk'
12    }
13
14    self.publish(json.dumps(data), queues=['Identifier']) # publish to 'Identifier'
```

6.2 Ghidra Worker development

6.2.1 Dockerfile

Install any additional system dependencies required within your worker's Dockerfile.

Base image options:

- [behren/machina-base-ghidra](#)

Dockerfile example:

```
1 FROM behren/machina-base-ghidra:latest
2 ...
3 RUN apt update && apt install libz-dev
4 ...
```

6.2.2 requirements.txt

Put any Python 3 requirements required by your worker module into requirements.txt and ensure that you copy requirements.txt into the image and 'pip3 install -r requirements.txt' to install the dependencies.

6.2.3 youranalysismodule.py

This file contains the implementation of your worker module. subclass the [GhidraWorker](#) class, this will ensure your worker module has boilerplate connectivity to the database, RabbitMQ, and configurations. Choose any Machina types (see 'machina/configs/types.json') your worker module supports, or specify '*' for all.

The [callback](#) function provides your analysis implementation with data that your module is configured to support. This callback function fires whenever the system identifies a compatible sample.

```
1 class YourAnalysisModule(GhidraWorker):
2     types = ["elf"]
3
4     def __init__(self, *args, **kwargs):
5         super(YourAnalysisModule, self).__init__(*args, **kwargs)
6         ...
7
8     def callback(self, data, properties):
9
10        # resolve path
11        target = self.get_binary_path(data['ts'], data['hashes']['md5'])
12        self.logger.info(f"resolved path: {target}")
13
14        self.analyze_headless(
15            str(Path(target).parent),
16            f'proj-{data["hashes"]["md5"]}-{self.cls_name}',
17            import_files=[target]
18        )
```

6.2.4 YourAnalysisModule.json (schema)

This schema file provides validation constraints that are applied to data incoming to your worker module before it handles the data. The Schema name must match the class name that it belongs to (e.g. for the worker module 'AndroguardAnalysis'). This file belongs at the top level of your worker's directory, and must be copied within the Dockerfile.

Typically, since workers are handling data published by the Identifier, they inherit from the 'binary.json' schema. Additional input requirements can be specified in "properties"

```
1 {
2     "allof": [{"$ref": "binary.json"}],
3     "properties": {}
4 }
```

6.2.5 YourAnalysisModule.json (configuration)

This top-level configuration file belongs in `machina/configs/workers/youranalysismodule.json`. This file allows for reconfiguration without rebuilding of images or code. This file must be named after the worker class name that it corresponds to. Configuration data set in this file is made available through the worker module's `'self.config["worker"]'` attribute. Log level is handled by the `Worker` base class to automatically adjust the subclass logging level if it is overridden in the configuration.

```
1 {
2   "log_level": "debug",
3   "analysis_timeout_per_file": 600
4 }
```

```
1 class YourAnalysisModule(Worker):
2     types = ["elf"]
3     ...
4     def callback(self, data, properties):
5         self.logger.info(self.config['worker']['analysis_timeout_per_file'])
6
7         # resolve path
8         target = self.get_binary_path(data['ts'], data['hashes']['md5'])
9         self.logger.info(f"resolved path: {target}")
10
11        self.analyze_headless(
12            str(Path(target).parent),
13            f'proj-{data["hashes"]["md5"]}-{self.cls_name}',
14            import_files=[target],
15            analysis_timeout_per_file=self.config['worker']['analysis_timeout_per_file']
16        )
```

6.3 Periodic Worker development

This type of Machina worker triggers an analysis periodically. This is useful for longer-running analysis work that isn't feasible to be run on-demand for every new ingested sample. One Periodic Worker example is `SimilarityAnalysis`. `SimilarityAnalysis` runs every hour, over all data in the graph, and updates/creates similarity relationships.

6.3.1 Dockerfile

Install any system dependencies required within your worker's Dockerfile. There are two base image options provided:

Base image options:

- [behren/machina-base-alpine](#)
- [behren/machina-base-ubuntu](#)

```
1 FROM behren/machina-base-ubuntu:latest
2 ...
3 RUN apt update && apt install libz-dev
4 ...
```

6.3.2 requirements.txt

Put any Python 3 requirements required by your worker module into `requirements.txt` and ensure that you copy `requirements.txt` into the image and `'pip3 install -r requirements.txt'` to install the dependencies

6.3.3 youranalysismodule.py

This file contains the implementation of your worker module. Subclass the `PeriodicWorker` class, this will ensure your worker module has boilerplate connectivity to the database, and configurations. The `callback` function fires at the interval your analysis module is configured to.

```
1 class YourAnalysisModule(PeriodicWorker):
2
3     def __init__(self, *args, **kwargs):
4         super(YourAnalysisModule, self).__init__(*args, **kwargs)
5         ...
6
7     def callback(self):
8         self.logger.info("I'm firing!")
```

The `PeriodicWorker` provides some common triggers that can be used to further constrain execution at the configured interval.

For example, `'n_nodes_added_since'` can be used to fire an analysis only if 1,000 new nodes of class type `Elf` have been added within the last 1 hour. These available triggers are documented within the machina core API.

```
1 from datetime import timedelta
2 from machina.core.models import Elf
3
4 class YourAnalysisModule(PeriodicWorker):
5
6     def __init__(self, *args, **kwargs):
7         super(YourAnalysisModule, self).__init__(*args, **kwargs)
8         ...
9
10    def callback(self):
11        if self.n_nodes_added_since(
12            1000,
13            Elf,
14            timedelta(seconds=60)
15        ):
16            self.logger.info("I'm firing under a special constraint!")
```

6.3.4 YourAnalysisModule.json (configuration)

This top-level configuration file belongs in `machina/configs/workers/youranalysismodule.json`. This file allows for reconfiguration without rebuilding of images or code. This file must be named after the worker class name that it corresponds to.

Configuration data set in this file is made available through the worker module's `'self.config["worker"]'` attribute. Log level is handled by the `PeriodicWorker` base class to automatically adjust the subclass logging level if it is overridden in the configuration. The interval for invoking `callback` is also handled by the `PeriodicWorker` base class, and can be overridden in the configuration.

```
1 {
2   "log_level": "debug",
3   "interval": "minutely",
4   "new_value": "I'm a new configuration!"
5 }
```

Note

Complex intervals can be set up. Under the hood, the `PeriodicWorker` uses `rocketry` for scheduling. Rocketry provides a verbose syntax for describing intervals, outlined [here](#). This syntax can be used within the `'interval'` configuration value to specify complex intervals.

```
1 class YourAnalysisModule(Worker):
2     ...
3     def callback(self):
4         self.logger.info(self.config['worker']['interval'])
5         self.logger.info(self.config['worker']['new_value'])
```

6.4 API

6.4.1 Machina Worker API

Analysis worker base class. Workers inheriting from this class receive data to analyze based on the chosen data type bindings

Source code in `machina/core/worker.py` 

16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115

116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315

```

316 class Worker():
317     """Analysis worker base class. Workers inheriting from this class receive data to analyze based on the chosen data type bindings"""
318
319     next_queues = [] # passes data to another queue in the sequence, this should allow for chaining
320     types = [] # indicates what type data to bind to, this should be completed in the subclass
321     types_blacklist = [] # indicates what type data NOT to bind to, this should be completed in the subclass. implies binding to all types not in this
322     list. cannot be combined with types
323
324     def __init__(self):
325         self.cls_name = self.__class__.__name__
326         self.config = self._load_configs()
327         self.schema = self._load_schema()
328
329         # Logging
330         level = logging.getLevelName(self.config['worker']['log_level'])
331         logging.basicConfig(level=level, format='[*] %(message)s')
332         self.logger = logging.getLogger(__name__)
333
334         if self.types and self.types_blacklist:
335             self.logger.error("both types and types_blacklist cannot be set at the same time")
336             raise Exception
337
338         # validate whitelist types
339         # if *, then set types to all available types
340         if self.types:
341             self.logger.info(f"Validating types: {pformat(self.types)}")
342             if '*' in self.types:
343                 self.types = self.config['types']['available_types']
344             else:
345                 types_valid, t = self._types_valid()
346                 if not types_valid:
347                     self.logger.error(f"{t} is not configured as a type in types.json")
348                     raise Exception
349
350         # validate black list types
351         if self.types_blacklist:
352             self.logger.info(f"Validating types blacklist: {pformat(self.types_blacklist)}")
353             types_blacklist_valid, t = self._types_blacklist_valid()
354             if not types_blacklist_valid:
355                 self.logger.error(f"{t} is not configured as a type in types.json, so cannot be blacklisted")
356                 raise Exception
357
358         # if valid, set types to all except the ones in valid blacklist, and '*'
359         self.types = [t for t in self.config['types']['available_types'] if t not in self.types_blacklist]
360         # self.types.remove('*')
361
362         # neo4j set connection
363         _cfg = self.config['neo4j']
364         config.DATABASE_URL = f"bolt://{_cfg['user']}:{_cfg['pass']}@{_cfg['host']}:{_cfg['port']}/{_cfg['db_name']}"
365
366         # Initializer does no queue consumption, so
367         # dont create a connection or queue for it
368         if self.cls_name != 'Initializer':
369
370             # RabbitMQ Connection info
371             # note this is not thread-safe
372             self.rmq_conn = self.get_rmq_conn()
373             self.rmq_recv_channel = self.rmq_conn.channel()
374
375             # reduce Pika logging level
376             logging.getLogger('pika').setLevel(logging.ERROR)
377
378             # The queue to bind to is the name of the class
379             bind_queue = self.cls_name
380
381             # Initialize an exchange
382             self.rmq_recv_channel.exchange_declare('machina')
383
384             # Initialize direct queue w/ subclass name
385             self.logger.info(f"Binding to direct queue: {bind_queue}")
386             self.rmq_recv_channel.queue_declare(self.cls_name, durable=True)
387
388             # Ensure that the worker's queue is bound to the exchange
389             self.rmq_recv_channel.queue_bind(exchange='machina',
390                                             queue=self.cls_name)
391
392             # multiple-bindings approach:
393             # https://www.rabbitmq.com/tutorials/tutorial-four-python.html
394             # Bind using type strings as routing_keys
395             # Publish using only a routing_key should go to all queues
396             # that are bound to that routing_key
397             for t in self.types:
398                 self.logger.info(f'binding to type: {t}')
399                 self.rmq_recv_channel.queue_bind(exchange='machina',
400                                                 queue=self.cls_name,
401                                                 routing_key=t)
402
403             self.rmq_recv_channel.basic_qos(prefetch_count=1)
404             self.rmq_recv_channel.basic_consume(self.cls_name,
405                                                on_message_callback=self._callback)
406
407     #####
408     # Privates
409
410     def _load_configs(self) -> dict:
411         """load configuration files from expected path, return as dictionary
412
413         :return: the configuration dictionary
414         :rtype: dict
415         """

```

```

fdir = '/configs'

paths_cfg_fp = Path(fdir, 'paths.json')
with open(paths_cfg_fp, 'r') as f:
    paths_cfg = json.load(f)

rabbitmq_cfg_fp = Path(fdir, 'rabbitmq.json')
with open(rabbitmq_cfg_fp, 'r') as f:
    rabbitmq_cfg = json.load(f)

neo4j_cfg_fp = Path(fdir, 'neo4j.json')
with open(neo4j_cfg_fp, 'r') as f:
    neo4j_cfg = json.load(f)

types_fp = Path(fdir, 'types.json')
with open(types_fp, 'r') as f:
    types_cfg = json.load(f)

# Base-worker configurations, will be overridden by worker-specific
# configurations if there is overlap
base_worker_cfg_fp = Path(fdir, 'workers', 'Worker.json')
with open(base_worker_cfg_fp, 'r') as f:
    worker_cfg = json.load(f)

# Worker-specific configuration
worker_cfg_fp = Path(fdir, 'workers', self.cls_name+'.json')
with open(worker_cfg_fp, 'r') as f:
    worker_cfg.update(json.load(f))

return dict(paths=paths_cfg,
            rabbitmq=rabbitmq_cfg,
            neo4j=neo4j_cfg,
            types=types_cfg,
            worker=worker_cfg)

def _load_schema(self) -> dict:
    """automatically resolve schema name based on class name

    :return: the schema dictionary
    :rtype: dict
    """
    class_schema = Path(self.config['paths']['schemas'], self.cls_name+'.json')
    with open(class_schema, 'r') as f:
        schema_data = json.load(f)
    return schema_data

def _callback(
    self,
    ch: pika.channel.Channel,
    method: pika.spec.Basic.Deliver,
    properties: pika.spec.BasicProperties,
    body: bytes):
    """do last-second validation before handling the callback

    :param ch: pika channel
    :type ch: pika.channel.Channel
    :param method: pika method
    :type method: pika.spec.Basic.Deliver
    :param properties: pika properties
    :type properties: pika.spec.BasicProperties
    :param body: message body
    :type body: bytes
    """
    self._validate_body(body)

    self.logger.info("entering callback")
    thread = threading.Thread(target=self.callback, args=(body, properties))
    thread.start()
    while thread.is_alive():
        self.rmq_recv_channel._connection.sleep(1.0)
    self.logger.info("exiting callback")

    self.rmq_recv_channel.basic_ack(delivery_tag=method.delivery_tag)

def _validate_body(self, body: bytes):
    """apply subclass worker schema and validate

    :param body: message body
    :type body: bytes
    """
    self.logger.info("validating schema")
    data = json.loads(body)
    # fixed resolver to ensure base schema uri is resolved
    # e.g. https://stackoverflow.com/questions/53968770/how-to-set-up-local-file-references-in-python-jjsonschema-document
    # resolver = jsonschema.RefResolver('file://{}'.format(os.path.join(self.config['paths']['schemas'], 'binary.json')), self.schema)
    resolver = jsonschema.RefResolver(f"file:{Path(self.config['paths']['schemas'], self.cls_name+'.json')}", self.schema)

    jsonschema.validate(
        instance=data,
        schema=self.schema,
        resolver=resolver)

def _types_valid(self) -> tuple:
    """ensure that the type to bind to is configured in types.json

    :return: tuple where first element is True if all requested type bindings are valid, or False if not. If invalid, set the second element to the first
    discovered invalid type
    :rtype: tuple
    """

```

```

for t in self.types:
    if t not in self.config['types']['available_types']:
        return False, t
    return True, None

def _types_blacklist_valid(self) -> tuple:
    """ensure that the type to bind to is configured in types.json

    :return: tuple where first element is True if all requested type bindings are valid, or False if not. If invalid, set the second element to the first
    discovered invalid type
    :rtype: tuple
    """
    for t in self.types_blacklist:
        if t not in self.config['types']['available_types']:
            return False, t
        return True, None
#####
#####
# RabbitMQ Helpers
def get_rmq_conn(
    self,
    max_attempts:int=10,
    delay_seconds:int=1) -> pika.BlockingConnection:
    """get RabbitMQ connection instance

    :param max_attempts: max number of attempts to try to get the connection, defaults to 10
    :type max_attempts: int, optional
    :param delay_seconds: the delay between attempts to get the connection, defaults to 1
    :type delay_seconds: int, optional
    :return: the connection instance
    :rtype: pika.BlockingConnection
    """

    rabbitmq_user = self.config['rabbitmq']['rabbitmq_user']
    rabbitmq_password = self.config['rabbitmq']['rabbitmq_password']
    rabbitmq_host = self.config['rabbitmq']['rabbitmq_host']
    rabbitmq_port = self.config['rabbitmq']['rabbitmq_port']
    rabbitmq_heartbeat = self.config['rabbitmq']['rabbitmq_heartbeat']

    connection = None
    credentials = pika.PlainCredentials(rabbitmq_user, rabbitmq_password)
    parameters = pika.ConnectionParameters(rabbitmq_host,
        int(rabbitmq_port),
        '/',
        credentials,
        heartbeat=int(rabbitmq_heartbeat),
        socket_timeout=2)

    attempt = 0
    while attempt < max_attempts:
        try:
            connection = pika.BlockingConnection(parameters)
            break
        except pika.exceptions.AMQPConnectionError as e:
            self.logger.info(f"Attempt {attempt}/{max_attempts} to connect to RabbitMQ at {rabbitmq_host}:{rabbitmq_port}")
            self.logger.warn("Error connecting to RabbitMQ")

            attempt += 1
            time.sleep(delay_seconds)

    if not connection:
        self.logger.error('max attempts exceeded')
        sys.exit()

    return connection

def start_consuming(self):
    """start consuming"""
    self.logger.info(f'{self.cls_name} worker started')
    try:
        self.rmqs_channel.start_consuming()
    except Exception as e:
        self.logger.error(e, exc_info=True)
        self.rmqs_channel.stop_consuming()
    self.connection.close()

def callback(self, data: bytes, properties: pika.spec.BasicProperties):
    """callback for worker, implement in subclass

    :param data: incoming string payload
    :type data: bytes
    :param properties: message properties
    :type properties: pika.spec.BasicProperties
    :raises NotImplementedError:
    """
    raise NotImplementedError

def publish_next(self, data: bytes):
    """publish to configured next_queues

    :param data: the data to publish
    :type data: bytes
    """
    if not self.next_queues:
        self.logger.warn('attempting to publish to next queue, but no next_queues defined in worker class')
    self.publish(data, queues=self.next_queues)

def publish(self, data: bytes, queues: list):
    """publish directly to a list of arbitrary queues

```

```

:param data: the data to publish
:type data: bytes
:param queues: the list of queue names (as strings) to publish data to
:type queues: list
"""
rmq_conn = self.get_rmq_conn()
for q in queues:
    self.logger.info(f"publishing directly to {q}")
    rmq_channel = rmq_conn.channel()
    rmq_channel.basic_publish(
        exchange='machina',
        routing_key=q,
        body=data)
rmq_conn.close()

#####

#####
# Misc
def get_binary_path(self, ts:str, md5:str) -> str:
    """get path to a binary on disk given a timestamp and its md5

    :param ts: the timestamp of the binary
    :type ts: str
    :param md5: the md5 of the binary
    :type md5: str
    :return: the path to the binary
    :rtype: str
    """
    binary_path = Path(self.config['paths']['binaries'], ts, md5)
    return str(binary_path)

```

`callback(data, properties)`

callback for worker, implement in subclass

Parameters:

Name	Type	Description	Default
<code>data</code>	<code>bytes</code>	incoming string payload	<i>required</i>
<code>properties</code>	<code>pika.spec.BasicProperties</code>	message properties	<i>required</i>

Raises:

Type	Description
<code>NotImplementedError</code>	

Source code in `machina/core/worker.py`

```

290 def callback(self, data: bytes, properties: pika.spec.BasicProperties):
291     """callback for worker, implement in subclass
292
293     :param data: incoming string payload
294     :type data: bytes
295     :param properties: message properties
296     :type properties: pika.spec.BasicProperties
297     :raises NotImplementedError:
298     """
299     raise NotImplementedError

```

`get_binary_path(ts, md5)`

get path to a binary on disk given a timestamp and its md5

Parameters:

Name	Type	Description	Default
<code>ts</code>	<code>str</code>	the timestamp of the binary	<i>required</i>
<code>md5</code>	<code>str</code>	the md5 of the binary	<i>required</i>

Returns:

Type	Description
<code>str</code>	the path to the binary

Source code in `machina/core/worker.py`

```

333 def get_binary_path(self, ts: str, md5: str) -> str:
334     """get path to a binary on disk given a timestamp and its md5
335
336     :param ts: the timestamp of the binary
337     :type ts: str
338     :param md5: the md5 of the binary
339     :type md5: str
340     :return: the path to the binary
341     :rtype: str
342     """
343     binary_path = Path(self.config['paths']['binaries'], ts, md5)
344     return str(binary_path)

```

```
get_rmq_conn(max_attempts=10, delay_seconds=1)
```

get RabbitMQ connection instance

Parameters:

Name	Type	Description	Default
max_attempts	int	max number of attempts to try to get the connection, defaults to 10	10
delay_seconds	int	the delay between attempts to get the connection, defaults to 1	1

Returns:

Type	Description
pika.BlockingConnection	the connection instance

Source code in machina/core/worker.py

```

233 def get_rmq_conn(
234     self,
235     max_attempts:int=10,
236     delay_seconds:int=1) -> pika.BlockingConnection:
237     """get RabbitMQ connection instance
238
239     :param max_attempts: max number of attempts to try to get the connection, defaults to 10
240     :type max_attempts: int, optional
241     :param delay_seconds: the delay between attempts to get the connection, defaults to 1
242     :type delay_seconds: int, optional
243     :return: the connection instance
244     :rtype: pika.BlockingConnection
245     """
246
247     rabbitmq_user = self.config['rabbitmq']['rabbitmq_user']
248     rabbitmq_password = self.config['rabbitmq']['rabbitmq_password']
249     rabbitmq_host = self.config['rabbitmq']['rabbitmq_host']
250     rabbitmq_port = self.config['rabbitmq']['rabbitmq_port']
251     rabbitmq_heartbeat = self.config['rabbitmq']['rabbitmq_heartbeat']
252
253     connection = None
254     credentials = pika.PlainCredentials(rabbitmq_user, rabbitmq_password)
255     parameters = pika.ConnectionParameters(rabbitmq_host,
256     int(rabbitmq_port),
257     '/',
258     credentials,
259     heartbeat=int(rabbitmq_heartbeat),
260     socket_timeout=2)
261
262     attempt = 0
263     while attempt < max_attempts:
264         try:
265             connection = pika.BlockingConnection(parameters)
266             break
267         except pika.exceptions.AMQPConnectionError as e:
268             self.logger.info(f"Attempt {attempt}/{max_attempts} to connect to RabbitMQ at {rabbitmq_host}:{rabbitmq_port}")
269             self.logger.warn("Error connecting to RabbitMQ")
270
271         attempt += 1
272         time.sleep(delay_seconds)
273
274     if not connection:
275         self.logger.error("max attempts exceeded")
276         sys.exit()
277
278     return connection

```

```
publish(data, queues)
```

publish directly to a list of arbitrary queues

Parameters:

Name	Type	Description	Default
data	bytes	the data to publish	<i>required</i>
queues	list	the list of queue names (as strings) to publish data to	<i>required</i>

Source code in machina/core/worker.py 

```

311 def publish(self, data: bytes, queues: list):
312     """publish directly to a list of arbitrary queues
313
314     :param data: the data to publish
315     :type data: bytes
316     :param queues: the list of queue names (as strings) to publish data to
317     :type queues: list
318     """
319     rmq_conn = self.get_rmq_conn()
320     for q in queues:
321         self.logger.info(f"publishing directly to {q}")
322         rmq_channel = rmq_conn.channel()
323         rmq_channel.basic_publish(
324             exchange='machina',
325             routing_key=q,
326             body=data)
327     rmq_conn.close()

```

publish_next(data)

publish to configured next_queues

Parameters:

Name	Type	Description	Default
data	bytes	the data to publish	<i>required</i>

Source code in machina/core/worker.py 

```

301 def publish_next(self, data: bytes):
302     """publish to configured next_queues
303
304     :param data: the data to publish
305     :type data: bytes
306     """
307     if not self.next_queues:
308         self.logger.warn('attempting to publish to next queue, but no next_queues defined in worker class')
309     self.publish(data, queues=self.next_queues)

```

start_consuming()

start consuming

Source code in machina/core/worker.py 

```

280 def start_consuming(self):
281     """start consuming"""
282     self.logger.info(f'{self.cls_name} worker started')
283     try:
284         self.rmq_recv_channel.start_consuming()
285     except Exception as e:
286         self.logger.error(e, exc_info=True)
287         self.rmq_recv_channel.stop_consuming()
288     self.connection.close()

```

6.4.2 Machina Ghidra Worker API

Bases: [Worker](#)

Ghidra Analysis worker base class

Source code in `machina/core/ghidra_worker.py` 

13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112

```

113 class GhidraWorker(Worker):
114     """Ghidra Analysis worker base class"""
115
116     def __init__(self):
117         super(GhidraWorker, self).__init__()
118
119         self.logger.info(f"GhidraWorker subclass name: {self.cls_name}")
120
121         # update worker config with GhidraWorker.json base default configurations,
122         # override with child class if specified.
123         self.config['worker'].update(self._load_ghidra_configs()['worker'])
124
125         # set GHIDRA_MAXMEM env var which is used in the patched analyzeHeadless script to control max memory
126         # os.environ['GHIDRA_MAXMEM'] = self.config['worker']['maxmem']
127
128         # update base configs
129         self.ghidra_analyzeheadless_path = Path(os.environ['GHIDRA_HOME'], 'support', 'analyzeHeadless')
130
131
132     def _load_ghidra_configs(self) -> dict:
133         """reload config again using GhidraWorker.json as default base
134
135         :return: the configuration dictionary
136         :rtype: dict
137         """
138
139         fdir = '/configs'
140
141         # Base-worker configurations, will be overridden by worker-specific
142         # configurations if there is overlap
143         base_worker_cfg_fp = Path(fdir, 'workers', 'GhidraWorker.json')
144         with open(base_worker_cfg_fp, 'r') as f:
145             worker_cfg = json.load(f)
146
147         # Worker-specific configuration
148         worker_cfg_fp = Path(fdir, 'workers', self.cls_name+'.json')
149         with open(worker_cfg_fp, 'r') as f:
150             worker_cfg.update(json.load(f))
151
152         return dict(
153             worker=worker_cfg
154         )
155
156
157     #####
158     # Ghidra API wrappers
159     def analyze_headless(
160         self,
161         project_location:str,
162         project_name:str,
163         import_files: List[str]=[],
164         process:List[str]=[],
165         pre_script:List[Tuple[str,List[str]]]=[],
166         post_script:List[Tuple[str,List[str]]]=[],
167         overwrite: bool=False,
168         recursive: bool=False,
169         read_only: bool=False,
170         delete_project: bool=False,
171         no_analysis: bool=False,
172         analysis_timeout_per_file: int=300
173     ):
174         """run Ghidra's analyzeHeadless
175
176         :param project_location: root directory of an existing project to import (if 'import' set or 'process' set), or where a new project will be created
177         (if 'import' set)
178         :type project_location: str
179         :param project_name: the name of an existing project to import (if 'import' set or 'process' set), or a new project (if 'import' set). if the
180         project name includes a folder path, imports will be rooted under that folder
181         :type project_name: str
182         :param import_files: list of executables and/or directories containing executables to import into a project. cannot be set when 'process' is set
183         :type import_files: List[str], optional
184         :param process: perform processing (pre/post scripts and/or analysis) on files in the list. if not set, all files in project_name will be analyzed.
185         cannot be set with 'import_files'
186         :type process: List[str], optional
187         :param pre_script: a list of tuples where the first element in a tuple is the pre script name (with extension), and the second element is a list of
188         strings containing arguments (if any) for the script. example: [ ('prescript1.java', ['-my-arg1','-my1-arg2']), ('prescript2.java', ['-my2-arg1','-my2-
189         arg2']) ] defaults to []
190         :type pre_script: List[Tuple[str,List[str]]], optional
191         :param post_script: a list of tuples where the first element in a tuple is the post script name (with extension), and the second element is a list
192         of strings containing arguments (if any) for the script. example: [ ('postscript1.java', ['-my-arg1','-my1-arg2']), ('postscript2.java', ['-my2-arg1','-my2-
193         arg2']) ] defaults to []
194         :type post_script: List[Tuple[str,List[str]]], optional
195         :param overwrite: overwrite any existing project files that conflict with the ones being imported. applies only if 'import' is set, and is ignored
196         if 'read_only' is set. defaults to False
197         :type overwrite: bool, optional
198         :param recursive: enables recursive descent into directories and project sub-folders when a directory/folder has been specified in 'import' or
199         'process'. defaults to False
200         :type recursive: bool, optional
201         :param read_only: ensures imported files will not be written if 'import' is set. also ensures any changes made when 'process' is set are discarded
202         :param delete_project: delete the project once all scripts have completed. only works on new projects created when 'import' is set. defaults to
203         False
204         :type delete_project: bool, optional
205         :param no_analysis: disables auto-analysis from running, defaults to False
206         :type no_analysis: bool, optional
207         :param analysis_timeout_per_file: timeout value in seconds per file analysis. defaults to 300
208         :type analysis_timeout_per_file: int
209
210         """
211
212         if not (import_files or process) or (import_files and process):

```

```

        self.logger.error('import_files OR process must be set')
        return

    cmd = f'{self.ghidra_analyzeheadless_path} {project_location} {project_name} '

    if import_files:
        cmd += '-import ' + ' '.join(import_files) + ' '
    if process:
        cmd += '-process ' + ' '.join(process) + ' '

    for script_cfg in pre_script:
        script_name = script_cfg[0]
        cmd += f'-preScript {script_name} '
        if len(script_cfg) > 1:
            script_args = script_cfg[1]
            for script_arg in script_args:
                cmd += f'{script_arg} '

    for script_cfg in post_script:
        script_name = script_cfg[0]
        cmd += f'-postScript {script_name} '
        if len(script_cfg) > 1:
            script_args = script_cfg[1]
            for script_arg in script_args:
                cmd += f'{script_arg} '

    if overwrite:
        cmd += '-overwrite '
        # https://static.grumpycoder.net/pixel/support/analyzeHeadlessREADME.html#overwrite
        if not import_files:
            self.logger.warn('overwrite only applies when import_files specified')

    if recursive:
        cmd += '-recursive '

    if read_only:
        cmd += '-readOnly '

        # https://static.grumpycoder.net/pixel/support/analyzeHeadlessREADME.html#readOnly
        if import_files and overwrite:
            self.logger.warn('overwrite option ignored')

    if delete_project:
        cmd += '-deleteProject '

    if no_analysis:
        cmd += '-noanalysis '

    cmd += f'-analysisTimeoutPerFile {analysis_timeout_per_file} '

    cmd += f'-max-cpu {self.config["worker"]["max_cpu"]}'

    self.logger.debug(f"running command: {cmd}")
    p = subprocess.run(
        shlex.split(cmd),
        env={
            'GHIDRA_MAXMEM': self.config['worker']['maxmem']
        })

```

```
analyze_headless(project_location, project_name, import_files=[], process=[], pre_script=[], post_script=[], overwrite=False,  
recursive=False, read_only=False, delete_project=False, no_analysis=False, analysis_timeout_per_file=300)
```

run Ghidra's analyzeHeadless

Parameters:

Name	Type	Description	Default
<code>project_location</code>	<code>str</code>	root directory of an existing project to import (if 'import' set or 'process' set), or where a new project will be created (if 'import' set)	<i>required</i>
<code>project_name</code>	<code>str</code>	the name of an existing project to import (if 'import' set or 'process' set), or a new project (if 'import' set). if the project name includes a folder path, imports will be rooted under that folder	<i>required</i>
<code>import_files</code>	<code>List[str]</code>	list of executables and/or directories containing executables to import into a project. cannot be set when 'process' is set	<code>[]</code>
<code>process</code>	<code>List[str]</code>	perform processing (pre/post scripts and/or analysis) on files in the list. if not set, all files in <code>project_name</code> will be analyzed. cannot be set with 'import_files'	<code>[]</code>
<code>pre_script</code>	<code>List[Tuple[str, List[str]]]</code>	a list of tuples where the first element in a tuple is the pre script name (with extension), and the second element is a list of strings containing arguments (if any) for the script. example: [('prescript1.java', ['-my-arg1', '-my1-arg2']), ('prescript2.java', ['-my2-arg1', '-my2-arg2'])] defaults to []	<code>[]</code>
<code>post_script</code>	<code>List[Tuple[str, List[str]]]</code>	a list of tuples where the first element in a tuple is the post script name (with extension), and the second element is a list of strings containing arguments (if any) for the script. example: [('postscript1.java', ['-my-arg1', '-my1-arg2']), ('postscript2.java', ['-my2-arg1', '-my2-arg2'])] defaults to []	<code>[]</code>
<code>overwrite</code>	<code>bool</code>	overwrite any existing project files that conflict with the ones being imported. applies only if 'import' is set, and is ignored if 'read_only' is set. defaults to False	<code>False</code>
<code>recursive</code>	<code>bool</code>	enables decursive descent into directories and project sub-folders when a directory/folder has been specified in 'import' or 'process'. defaults to False	<code>False</code>
<code>read_only</code>	<code>bool</code>	ensures imported files will not be written if 'import' is set. also ensures any changes made when 'process' is set are discarded	<code>False</code>
<code>delete_project</code>	<code>bool</code>	delete the project once all scripts have completed. only works on new projects created when 'import' is set. defaults to False	<code>False</code>
<code>no_analysis</code>	<code>bool</code>	disables auto-analysis from running, defaults to False	<code>False</code>
<code>analysis_timeout_per_file</code>	<code>int</code>	timeout value in seconds per file analysis. defaults to 300	<code>300</code>

Source code in `machina/core/ghidra_worker.py` 

```

59 def analyze_headless(
60     self,
61     project_location:str,
62     project_name:str,
63     import_files: List[str]=[],
64     process:List[str]=[],
65     pre_script:List[Tuple[str,List[str]]]=[],
66     post_script:List[Tuple[str,List[str]]]=[],
67     overwrite: bool=False,
68     recursive: bool=False,
69     read_only: bool=False,
70     delete_project: bool=False,
71     no_analysis: bool=False,
72     analysis_timeout_per_file: int=300
73 ):
74     """run Ghidra's analyzeHeadless
75
76     :param project_location: root directory of an existing project to import (if 'import' set or 'process' set), or where a new project will be created (if
77 'import' set)
78     :type project_location: str
79     :param project_name: the name of an existing project to import (if 'import' set or 'process' set), or a new project (if 'import' set). if the project
80 name includes a folder path, imports will be rooted under that folder
81     :type project_name: str
82     :param import_files: list of executables and/or directories containing executables to import into a project. cannot be set when 'process' is set
83     :type import_files: List[str], optional
84     :param process: perform processing (pre/post scripts and/or analysis) on files in the list. if not set, all files in project_name will be analyzed.
85 cannot be set with 'import_files'
86     :type process: List[str], optional
87     :param pre_script: a list of tuples where the first element in a tuple is the pre script name (with extension), and the second element is a list of
88 strings containing arguments (if any) for the script. example: [ ('prescript1.java', ['-my-arg1','-my1-arg2']), ('prescript2.java', ['-my2-arg1','-my2-
89 arg2']) ] defaults to []
90     :type pre_script: List[Tuple[str,List[str]]], optional
91     :param post_script: a list of tuples where the first element in a tuple is the post script name (with extension), and the second element is a list of
92 strings containing arguments (if any) for the script. example: [ ('postscript1.java', ['-my-arg1','-my1-arg2']), ('postscript2.java', ['-my2-arg1','-my2-
93 arg2']) ] defaults to []
94     :type post_script: List[Tuple[str,List[str]]], optional
95     :param overwrite: overwrite any existing project files that conflict with the ones being imported. applies only if 'import' is set, and is ignored if
96 'read_only' is set. defaults to False
97     :type overwrite: bool, optional
98     :param recursive: enables recursive descent into directories and project sub-folders when a directory/folder has been specified in 'import' or
99 'process'. defaults to False
100     :type recursive: bool, optional
101     :param read_only: ensures imported files will not be written if 'import' is set. also ensures any changes made when 'process' is set are discarded
102     :param delete_project: delete the project once all scripts have completed. only works on new projects created when 'import' is set. defaults to False
103     :type delete_project: bool, optional
104     :param no_analysis: disables auto-analysis from running, defaults to False
105     :type no_analysis: bool, optional
106     :param analysis_timeout_per_file: timeout value in seconds per file analysis. defaults to 300
107     :type analysis_timeout_per_file: int
108
109     """
110
111     if not (import_files or process) or (import_files and process):
112         self.logger.error("import_files OR process must be set")
113         return
114
115     cmd = f'{self.ghidra_analyzeheadless_path} {project_location} {project_name} '
116
117     if import_files:
118         cmd += '-import ' + ' '.join(import_files) + ' '
119     if process:
120         cmd += '-process ' + ' '.join(process) + ' '
121
122     for script_cfg in pre_script:
123         script_name = script_cfg[0]
124         cmd += f'-preScript {script_name} '
125         if len(script_cfg) > 1:
126             script_args = script_cfg[1]
127             for script_arg in script_args:
128                 cmd += f'{script_arg} '
129
130     for script_cfg in post_script:
131         script_name = script_cfg[0]
132         cmd += f'-postScript {script_name} '
133         if len(script_cfg) > 1:
134             script_args = script_cfg[1]
135             for script_arg in script_args:
136                 cmd += f'{script_arg} '
137
138     if overwrite:
139         cmd += '-overwrite '
140         # https://static.grumpycoder.net/pixel/support/analyzeHeadlessREADME.html#overwrite
141         if not import_files:
142             self.logger.warn('overwrite only applies when import_files specified')
143
144     if recursive:
145         cmd += '-recursive '
146
147     if read_only:
148         cmd += '-readOnly '
149         # https://static.grumpycoder.net/pixel/support/analyzeHeadlessREADME.html#readOnly
150         if import_files and overwrite:
151             self.logger.warn('overwrite option ignored')
152
153     if delete_project:
154         cmd += '-deleteProject '
155
156     if no_analysis:
157         cmd += '-noanalysis '

```

```
cmd += f'-analysisTimeoutPerFile {analysis_timeout_per_file} '  
cmd += f'-max-cpu {self.config["worker"]["max_cpu"]}'  
  
self.logger.debug(f"running command: {cmd}")  
p = subprocess.run(  
    shlex.split(cmd),  
    env={'GHIDRA_MAXMEM':self.config['worker']['maxmem']})
```

6.4.3 Machina Periodic Worker API

Analysis worker base class. Workers inheriting from this class analyze based on a schedule

Source code in `machina/core/periodic_worker.py` 

```

12 class PeriodicWorker():
13     """Analysis worker base class. Workers inheriting from this class analyze based on a schedule"""
14
15     def __init__(self):
16         self.cls_name = self.__class__.__name__
17         self.config = self._load_configs()
18         self.app = Rocketry()
19
20         # Logging
21         level = logging.getLevelName(self.config['worker']['log_level'])
22         logging.basicConfig(level=level, format='[*] %(message)s')
23         self.logger = logging.getLogger(__name__)
24
25         # reduce Rocketry logging level
26         logging.getLogger('rocketry.scheduler').setLevel(logging.ERROR)
27
28         # neo4j set connection
29         _cfg = self.config['neo4j']
30         config.DATABASE_URL = f'bolt://{_cfg['user']};{_cfg['pass']}@{_cfg['host']}:{_cfg['port']}/{_cfg['db_name']}'
31
32     #####
33     # Privates
34     def _load_configs(self) -> dict:
35         """load configuration files from expected path, return as dictionary
36
37         :return: the configuration dictionary
38         :rtype: dict
39         """
40
41         fdir = '/configs'
42
43         paths_cfg_fp = Path(fdir, 'paths.json')
44         with open(paths_cfg_fp, 'r') as f:
45             paths_cfg = json.load(f)
46
47         neo4j_cfg_fp = Path(fdir, 'neo4j.json')
48         with open(neo4j_cfg_fp, 'r') as f:
49             neo4j_cfg = json.load(f)
50
51         types_fp = Path(fdir, 'types.json')
52         with open(types_fp, 'r') as f:
53             types_cfg = json.load(f)
54
55         # Base-worker configurations, will be overridden by worker-specific
56         # configurations if there is overlap
57         base_worker_cfg_fp = Path(fdir, 'workers', 'PeriodicWorker.json')
58         with open(base_worker_cfg_fp, 'r') as f:
59             worker_cfg = json.load(f)
60
61         # Worker-specific configuration
62         worker_cfg_fp = Path(fdir, 'workers', self.cls_name+'.json')
63         with open(worker_cfg_fp, 'r') as f:
64             worker_cfg.update(json.load(f))
65
66         return dict(paths=paths_cfg,
67                   neo4j=neo4j_cfg,
68                   types=types_cfg,
69                   worker=worker_cfg)
70     #####
71
72     def start(self):
73         """start running callback at interval"""
74         self.logger.info(f"starting with interval: {self.config['worker']['interval']}")
75         self.app.task(self.config['worker']['interval'], func=self.callback)
76         self.app.run()
77
78     def callback(self):
79         """implement in subclass"""
80         raise NotImplemented
81
82     #####
83     # Triggers
84     def n_nodes_added_since(
85         n: int,
86         node_cls: Type[Base],
87         duration: timedelta):
88         """return True if 'n' nodes of 'node_cls' type have been added within a duration
89         of time, return False if either (or both) condition was not met.
90
91         :param n: the threshold number of nodes to consider before returning True
92         :type n: int
93         :param node_cls: the neomodel OGM class to use for counting node instances
94         :type node_cls: type[Base]
95         :param duration: the datetime.timedelta object specifying the threshold duration of time
96         :type duration: timedelta
97         :return: True if both conditions were met, False of neither (or both) not met
98         :rtype: bool
99         """
100
101         # get adjusted timestamp for provided
102         # threshold duration
103         now = datetime.now(timezone.utc)
104         duration_ts = (now - duration)
105
106         # filter nodes within the duration window
107         nodes = node_cls.nodes.filter(ts__gte=duration_ts).order_by('ts')
108         if len(nodes) >= n:
109             return True
110         return False

```

callback()

implement in subclass

Source code in `machina/core/periodic_worker.py`

```
78 def callback(self):
79     """implement in subclass"""
80     raise NotImplementedError
```

n_nodes_added_since(n, node_cls, duration)

return True if 'n' nodes of 'node_cls' type have been added within a duration of time, return False if either (or both) condition was not met.

Parameters:

Name	Type	Description	Default
n	int	the threshold number of nodes to consider before returning True	required
node_cls	Type[Base]	the neomodel OGM class to use for counting node instances	required
duration	timedelta	the datetime.timedelta object specifying the threshold duration of time	required

Returns:

Type	Description
bool	True if both conditions were met, False of neither (or both) not met

Source code in `machina/core/periodic_worker.py`

```
84 def n_nodes_added_since(
85     n: int,
86     node_cls: Type[Base],
87     duration: timedelta):
88     """return True if 'n' nodes of 'node_cls' type have been added within a duration
89     of time, return False if either (or both) condition was not met.
90
91     :param n: the threshold number of nodes to consider before returning True
92     :type n: int
93     :param node_cls: the neomodel OGM class to use for counting node instances
94     :type node_cls: type[Base]
95     :param duration: the datetime.timedelta object specifying the threshold duration of time
96     :type duration: timedelta
97     :return: True if both conditions were met, False of neither (or both) not met
98     :rtype: bool
99     """
100
101     # get adjusted timestamp for provided
102     # threshold duration
103     now = datetime.now(timezone.utc)
104     duration_ts = (now - duration)
105
106     # filter nodes within the duration window
107     nodes = node_cls.nodes.filter(ts_gte=duration_ts).order_by('ts')
108     if len(nodes) >= n:
109         return True
110     return False
```

start()

start running callback at interval

Source code in `machina/core/periodic_worker.py` 

```
72 def start(self):
73     """start running callback at interval"""
74     self.logger.info(f"starting with interval: {self.config['worker']['interval']}")
75     self.app.task(self.config['worker']['interval'], func=self.callback)
76     self.app.run()
```


6.4.4 Models API

Nodes

Bases: [StructuredNode](#)

Base node type

Source code in [machina/core/models/nodes/base.py](#) ▾

```

9 class Base(StructuredNode):
10     """Base node type"""
11
12     __abstract_node__ = True
13
14     # Common attributes
15     uid = UniqueIdProperty()
16
17     md5 = StringProperty(required=True)
18     sha256 = StringProperty(required=True)
19     size = IntegerProperty(required=True)
20     ts = DateTimeProperty(required=True)
21     type = StringProperty(required=True)
22
23     ssdeep = StringProperty(default=None) # set later
24
25     extracts = RelationshipTo('Base', 'EXTRACTS', model=Extracts)
26     similar = Relationship('Base', 'SIMILAR', model=Similar)
27     retyped = RelationshipTo('Base', 'RETYPED', model=Retyped)

```

Bases: [Base](#)

A generic artifact for unknown/untyped data

Source code in [machina/core/models/nodes/artifact.py](#) ▾

```

3 class Artifact(Base):
4     """A generic artifact for unknown/untyped data"""
5     pass

```

Bases: [Base](#)

Source code in [machina/core/models/nodes/apk.py](#) ▾

```

6 class APK(Base):
7
8     # APK Attribute
9     package = StringProperty()
10    name = StringProperty()
11    androidversion_code = StringProperty()
12    androidversion_name = StringProperty()
13    permissions = ArrayProperty(StringProperty())
14    activities = ArrayProperty(StringProperty())
15    providers = ArrayProperty(StringProperty())
16    receivers = ArrayProperty(StringProperty())
17    services = ArrayProperty(StringProperty())
18    min_sdk_version = StringProperty()
19    max_sdk_version = StringProperty()
20    max_sdk_version = StringProperty()
21    effective_target_sdk_version = StringProperty()
22    libraries = ArrayProperty(StringProperty())
23    main_activity = StringProperty()
24    content_provider_uris = ArrayProperty(StringProperty())
25
26    classes = ArrayProperty(JSONProperty())

```

Bases: [Base](#)

bz2 compressed file

Source code in `machina/core/models/nodes/bz2.py` ▾

```
3 class BZ2(Base):
4     """bz2 compressed file"""
5
6     pass
```

Bases: `Base`

CPIO firmware

Source code in `machina/core/models/nodes/cpio.py` ▾

```
3 class CPIO(Base):
4     """CPIO firmware"""
5
6     pass
```

Bases: `Base`

Android DEX file

Source code in `machina/core/models/nodes/dex.py` ▾

```
3 class Dex(Base):
4     """Android DEX file"""
5
6     pass
```

Bases: `Base`

Linux ELF file

Source code in `machina/core/models/nodes/elf.py` ▾

```
3 class Elf(Base):
4     """Linux ELF file"""
5
6     pass
```

Bases: `Base`

EML

Source code in `machina/core/models/nodes/eml.py` ▾

```
3 class Eml(Base):
4     """EML"""
5
6     pass
```

Bases: `Base`

Excel file

Source code in `machina/core/models/nodes/excel.py` ▾

```
3 class Excel(Base):
4     """Excel file"""
5     pass
```

Bases: [Base](#)

GZIP compressed file

Source code in `machina/core/models/nodes/gzip.py` ▾

```
3 class Gzip(Base):
4     """GZIP compressed file"""
5     pass
```

Bases: [Base](#)

HTML file

Source code in `machina/core/models/nodes/html.py` ▾

```
3 class HTML(Base):
4     """HTML file"""
5     pass
```

Bases: [Base](#)

JAR file

Source code in `machina/core/models/nodes/jar.py` ▾

```
3 class Jar(Base):
4     """JAR file"""
5     pass
```

Bases: [Base](#)

JFFS2 firmware

Source code in `machina/core/models/nodes/jffs2.py` ▾

```
3 class JFFS2(Base):
4     """JFFS2 firmware"""
5     pass
```

Bases: [Base](#)

JPEF Image file

Source code in `machina/core/models/nodes/jpeg.py` 

```

5 class JPEG(Base):
6     """JPEF Image file"""
7
8     # PNG attributes
9     exif = JSONProperty()

```

Bases: [Base](#)

LZMA compressed file

Source code in `machina/core/models/nodes/lzma.py` 

```

3 class LZMA(Base):
4     """LZMA compressed file"""
5     pass

```

Bases: [Base](#)

MachO file

Source code in `machina/core/models/nodes/macho.py` 

```

3 class Macho(Base):
4     """MachO file"""
5     pass

```

Bases: [Base](#)

memory dump file

Source code in `machina/core/models/nodes/memory_dump.py` 

```

3 class MemoryDump(Base):
4     """memory dump file"""
5     pass

```

Bases: [Base](#)

MSG file

Source code in `machina/core/models/nodes/msg.py` 

```

3 class Msg(Base):
4     """MSG file"""
5     pass

```

Bases: [Base](#)

PDF file

Source code in `machina/core/models/nodes/pdf.py` 

```

3 class PDF(Base):
4     """PDF file"""
5     pass

```

Bases: [Base](#)

PE file

Source code in `machina/core/models/nodes/pe.py` 

```

3 class PE(Base):
4     """PE file"""
5     pass

```

Bases: [Base](#)**Source code in** `machina/core/models/nodes/png.py` 

```

5 class PNG(Base):
6
7     # PNG attributes
8     exif = JSONProperty()

```

Bases: [Base](#)

Powerpoint file

Source code in `machina/core/models/nodes/powerpoint.py` 

```

3 class Powerpoint(Base):
4     """Powerpoint file"""
5     pass

```

Bases: [Base](#)

RTF file

Source code in `machina/core/models/nodes/rtf.py` 


```

3 class RTF(Base):
4     """RTF file"""
5     pass

```

Bases: [Base](#)

SquashFS firmware

Source code in `machina/core/models/nodes/squashfs.py` 

```

3 class SquashFS(Base):
4     """SquashFS firmware"""
5     pass

```

Bases: [Base](#)

TAR archive file

Source code in `machina/core/models/nodes/tar.py` 

```

3 class Tar(Base):
4     """TAR archive file"""
5     pass

```

Bases: [Base](#)

extracted URL

Source code in `machina/core/models/nodes/url.py` 

```

5 class URL(Base):
6     """extracted URL"""
7
8     # URL Attribute
9     url = StringProperty()

```

Bases: [Base](#)

TIFF image file

Source code in `machina/core/models/nodes/tiff.py` 

```

5 class TIFF(Base):
6     """TIFF image file"""
7
8     # attributes
9     exif = JSONProperty()

```

Bases: [Base](#)

Word file

Source code in `machina/core/models/nodes/word.py` 

```

3 class Word(Base):
4     """Word file"""
5     pass

```

Bases: [Base](#)

Zip compressed file

Source code in `machina/core/models/nodes/zip.py` ▾

```
3 class Zip(Base):
4     """Zip compressed file"""
5     pass
```

Relationships

Bases: `StructuredRel`

Base relationship

Source code in `machina/core/models/relationships/base.py` ▾

```
5 class BaseRelationship(StructuredRel):
6     """Base relationship"""
7
8     ts = DateTimeProperty(default=lambda: datetime.now())
```

Bases: `BaseRelationship`

Establish a node (some binary data) as being extracted from another node

Source code in `machina/core/models/relationships/extracts.py` ▾

```
5 class Extracts(BaseRelationship):
6     """Establish a node (some binary data) as being
7     extracted from another node"""
8     label = 'extracts'
9
10    # E.g. 'dynamic', 'static'
11    method = StringProperty()
```

Bases: `BaseRelationship`

Establish a node (some binary data) as being extracted from another node

Source code in `machina/core/models/relationships/retyped.py` ▾

```
5 class Retyped(BaseRelationship):
6     """Establish a node (some binary data) as being
7     extracted from another node"""
8     pass
```

Utils

`db_ts_to_fs_fmt(ts)`

convert database timestamp to file-system formatted timestamp string

Source code in `machina/core/models/utils.py` 

```
20 def db_ts_to_fs_fmt(ts:datetime) -> str:
21     """convert database timestamp to file-system formatted timestamp string"""
22     return ts.strftime("%Y%m%d%H%M%S%f")
```

`resolve_db_node_cls(resolved_type)`

resolve a OGM subclass given a resolved machina type (e.g. in types.json) if not resolved, we expect unresolved to be stored as a generic Artifact, so return that cls

Returns:

Type	Description
str	the type string to resolve to a class

Source code in `machina/core/models/utils.py` 

```
6 def resolve_db_node_cls(resolved_type: str) -> Type[Base]:
7     """resolve a OGM subclass given a resolved machina type (e.g. in types.json)
8     if not resolved, we expect unresolved to be stored as a generic Artifact, so return that cls
9
10    :return: the type string to resolve to a class
11    :rtype: str
12    """
13    all_models = Base.__subclasses__()
14    for c in all_models:
15        # if c.element_type.lower() == resolved_type.lower():
16        if c.__name__.lower() == resolved_type.lower():
17            return c
18    return Artifact
```