

# Package ‘maestro’

September 23, 2024

**Type** Package

**Title** Orchestration of Data Pipelines

**Version** 0.3.0

**Maintainer** Will Hipson <will.e.hipson@gmail.com>

**Description** Framework for creating and orchestrating data pipelines. Organize, orchestrate, and monitor multiple pipelines in a single project. Use tags to decorate functions with scheduling parameters and configuration.

**License** MIT + file LICENSE

**URL** <https://github.com/whipson/maestro>,  
<https://whipson.github.io/maestro/>

**BugReports** <https://github.com/whipson/maestro/issues>

**Encoding** UTF-8

**LazyData** true

**Imports** cli (>= 3.3.0), dplyr (>= 1.1.0), glue, lifecycle, logger, lubridate (>= 1.9.1), purrr (>= 1.0.0), R.utils, R6, rlang (>= 1.0.0), roxygen2, tictoc, timechange, utils

**RoxygenNote** 7.3.2

**Depends** R (>= 4.1.0)

**Suggests** furrr, future, knitr, rmarkdown, rstudioapi, testthat (>= 3.0.0), withr

**Config/testthat.edition** 3

**VignetteBuilder** knitr

**NeedsCompilation** no

**Author** Will Hipson [cre, aut, cph] (<<https://orcid.org/0000-0002-3931-2189>>), Ryan Garnett [aut, ctb, cph]

**Repository** CRAN

**Date/Publication** 2024-09-23 15:10:06 UTC

## Contents

build_schedule . . . . .	2
create_maestro . . . . .	3
create_orchestrator . . . . .	4
create_pipeline . . . . .	4
get_artifacts . . . . .	6
get_schedule . . . . .	7
get_status . . . . .	7
invoke . . . . .	8
last_build_errors . . . . .	9
last_run_errors . . . . .	10
last_run_messages . . . . .	10
last_run_warnings . . . . .	11
MaestroSchedule . . . . .	11
maestro_tags . . . . .	13
run_schedule . . . . .	14
suggest_orch_frequency . . . . .	16

<b>Index</b>	<b>17</b>
--------------	-----------

---

<b>build_schedule</b>	<i>Build a schedule table</i>
-----------------------	-------------------------------

---

### Description

Builds a schedule data.frame for scheduling pipelines in `run_schedule()`.

### Usage

```
build_schedule(pipeline_dir = "./pipelines", quiet = FALSE)
```

### Arguments

<code>pipeline_dir</code>	path to directory containing the pipeline scripts
<code>quiet</code>	silence metrics to the console (default = FALSE)

### Details

This function parses the maestro tags of functions located in `pipeline_dir` which is conventionally called 'pipelines'. An orchestrator requires a schedule table to determine which pipelines are to run and when. Each row in a schedule table is a pipeline name and its scheduling parameters such as its frequency.

The schedule table is mostly intended to be used by `run_schedule()` immediately. In other words, it is not recommended to make changes to it.

### Value

`MaestroSchedule`

## Examples

```
# Creating a temporary directory for demo purposes! In practice, just
# create a 'pipelines' directory at the project level.

if (interactive()) {
  pipeline_dir <- tempdir()
  create_pipeline("my_new_pipeline", pipeline_dir, open = FALSE)
  build_schedule(pipeline_dir = pipeline_dir)
}
```

---

create_maestro	<i>Creates a new maestro project</i>
----------------	--------------------------------------

---

## Description

Creates a new maestro project

## Usage

```
create_maestro(path, type = "R", overwrite = FALSE, quiet = FALSE, ...)
```

## Arguments

path	file path for the orchestrator script
type	file type for the orchestrator (supports R, Quarto, and RMarkdown)
overwrite	whether to overwrite an existing orchestrator or maestro project
quiet	whether to silence messages in the console (default = FALSE)
...	unused

## Value

invisible

## Examples

```
# Creates a new maestro project with an R orchestrator
if (interactive()) {
  new_proj_dir <- tempdir()
  create_maestro(new_proj_dir, type = "R", overwrite = TRUE)

  create_maestro(new_proj_dir, type = "Quarto", overwrite = TRUE)
}
```

`create_orchestrator`     *Create a new orchestrator*

### Description

Create a new orchestrator

### Usage

```
create_orchestrator(
  path,
  type = c("R", "Quarto", "RMarkdown"),
  open = interactive(),
  quiet = FALSE,
  overwrite = FALSE
)
```

### Arguments

<code>path</code>	file path for the orchestrator script
<code>type</code>	file type for the orchestrator (supports R, Quarto, and RMarkdown)
<code>open</code>	whether or not to open the script upon creation
<code>quiet</code>	whether to silence messages in the console (default = FALSE)
<code>overwrite</code>	whether to overwrite an existing orchestrator or maestro project

### Value

`invisible`

`create_pipeline`     *Create a new pipeline in a pipelines directory*

### Description

Allows the creation of new pipelines (R scripts) and fills in the maestro tags as specified.

### Usage

```
create_pipeline(
  pipe_name,
  pipeline_dir = "pipelines",
  frequency = "1 day",
  start_time = Sys.Date(),
  tz = "UTC",
```

```
    log_level = "INFO",
    quiet = FALSE,
    open = interactive(),
    overwrite = FALSE
)
```

## Arguments

pipe_name	name of the pipeline and function
pipeline_dir	directory containing the pipeline scripts
frequency	how often the pipeline should run (e.g., 1 day, daily, 3 hours, 4 months). Fills in maestroFrequency tag
start_time	start time of the pipeline schedule. Fills in maestroStartTime tag
tz	timezone that pipeline will be scheduled in. Fills in maestroTz tag
log_level	log level for the pipeline (e.g., INFO, WARN, ERROR). Fills in maestroLogLevel tag
quiet	whether to silence messages in the console (default = FALSE)
open	whether or not to open the script upon creation
overwrite	whether or not to overwrite an existing pipeline of the same name and location.

## Value

invisible

## Examples

```
if (interactive()) {
  pipeline_dir <- tempdir()
  create_pipeline(
    "extract_data",
    pipeline_dir = pipeline_dir,
    frequency = "1 hour",
    open = FALSE,
    quiet = TRUE,
    overwrite = TRUE
  )

  create_pipeline(
    "new_job",
    pipeline_dir = pipeline_dir,
    frequency = "20 minutes",
    start_time = as.POSIXct("2024-06-21 12:20:00"),
    log_level = "ERROR",
    open = FALSE,
    quiet = TRUE,
    overwrite = TRUE
  )
}
```

<code>get_artifacts</code>	<i>Get the artifacts (return values) of the pipelines in a MaestroSchedule object.</i>
----------------------------	--

## Description

Artifacts are return values from pipelines. They are accessible as a named list where the names correspond to the names of the pipeline.

## Usage

```
get_artifacts(schedule)
```

## Arguments

schedule	object of type MaestroSchedule created using build_schedule()
----------	---

## Value

named list

## Examples

```
if (interactive()) {
  pipeline_dir <- tempdir()
  create_pipeline("my_new_pipeline", pipeline_dir, open = FALSE)
  schedule <- build_schedule(pipeline_dir = pipeline_dir)

  schedule <- run_schedule(
    schedule,
    orch_frequency = "1 day",
    quiet = TRUE
  )

  get_artifacts(schedule)

  # Alternatively, use the underlying R6 method
  schedule$get_artifacts()
}
```

---

get_schedule	<i>Get the schedule from a MaestroSchedule object</i>
--------------	---

---

## Description

A schedule is represented as a table where each row is a pipeline and the columns contain scheduling parameters such as the frequency and start time.

## Usage

```
get_schedule(schedule)
```

## Arguments

schedule      object of type MaestroSchedule created using build\_schedule()

## Details

The schedule table is used internally in a MaestroSchedule object but can be accessed using this function or accessing the R6 method of the MaestroSchedule object.

## Value

data.frame

## Examples

```
if (interactive()) {  
  pipeline_dir <- tempdir()  
  create_pipeline("my_new_pipeline", pipeline_dir, open = FALSE)  
  schedule <- build_schedule(pipeline_dir = pipeline_dir)  
  
  get_schedule(schedule)  
  
  # Alternatively, use the underlying R6 method  
  schedule$get_schedule()  
}
```

---

get_status	<i>Get the statuses of the pipelines in a MaestroSchedule object</i>
------------	--

---

## Description

A status data.frame contains the names and locations of the pipelines as well as information around whether they were invoked, the status (error, warning, etc.), and the run time.

**Usage**

```
get_status(schedule)
```

**Arguments**

schedule      object of type MaestroSchedule created using build\_schedule()

**Value**

data.frame

**Examples**

```
if (interactive()) {
  pipeline_dir <- tempdir()
  create_pipeline("my_new_pipeline", pipeline_dir, open = FALSE)
  schedule <- build_schedule(pipeline_dir = pipeline_dir)

  schedule <- run_schedule(
    schedule,
    orch_frequency = "1 day",
    quiet = TRUE
  )

  get_status(schedule)

  # Alternatively, use the underlying R6 method
  schedule$get_status()
}
```

**invoke**

*Manually run a pipeline regardless of schedule*

**Description**

Instantly run a single pipeline from the schedule. This is useful for testing purposes or if you want to just run something one-off.

**Usage**

```
invoke(schedule, pipe_name, resources = list(), ...)
```

**Arguments**

schedule	object of type MaestroSchedule created using build_schedule()
pipe_name	name of a single pipe name from the schedule
resources	named list of shared resources made available to pipelines as needed
...	other arguments passed to run_schedule()

**Details**

Scheduling parameters such as the frequency, start time, and specifiers are ignored. The pipeline will be run even if maestroSkip is present.

**Value**

invisible

**Examples**

```
if (interactive()) {  
    pipeline_dir <- tempdir()  
    create_pipeline("my_new_pipeline", pipeline_dir, open = FALSE)  
    schedule <- build_schedule(pipeline_dir = pipeline_dir)  
  
    invoke(schedule, "my_new_pipeline")  
}
```

---

**last\_build\_errors**

*Retrieve latest maestro build errors*

---

**Description**

Gets the latest schedule build errors following use of build\_schedule(). If the build succeeded or build\_schedule() has not been run it will be NULL.

**Usage**

```
last_build_errors()
```

**Value**

error messages

**Examples**

```
last_build_errors()
```

---

<code>last_run_errors</code>	<i>Retrieve latest maestro pipeline errors</i>
------------------------------	--

---

### Description

Gets the latest pipeline errors following use of `run_schedule()`. If the all runs succeeded or `run_schedule()` has not been run it will be NULL.

### Usage

```
last_run_errors()
```

### Value

error messages

### Examples

```
last_run_errors()
```

---

<code>last_run_messages</code>	<i>Retrieve latest maestro pipeline messages</i>
--------------------------------	--

---

### Description

Gets the latest pipeline messages following use of `run_schedule()`. If there are no messages or `run_schedule()` has not been run it will be NULL.

### Usage

```
last_run_messages()
```

### Details

Note that setting `maestroLogLevel` to something greater than `INFO` will ignore messages.

### Value

messages

### Examples

```
last_run_messages()
```

---

last_run_warnings	<i>Retrieve latest maestro pipeline warnings</i>
-------------------	--

---

## Description

Gets the latest pipeline warnings following use of `run_schedule()`. If there are no warnings or `run_schedule()` has not been run it will be NULL.

## Usage

```
last_run_warnings()
```

## Details

Note that setting `maestroLogLevel` to something greater than `WARN` will ignore warnings.

## Value

warning messages

## Examples

```
last_run_warnings()
```

---

MaestroSchedule	<i>Class for a schedule of pipelines</i>
-----------------	--

---

## Description

Class for a schedule of pipelines

Class for a schedule of pipelines

## Public fields

`PipelineList` object of type `MaestroPipelineList`

## Active bindings

status deprecation

artifacts deprecation

## Methods

### Public methods:

- [MaestroSchedule\\$new\(\)](#)
- [MaestroSchedule\\$print\(\)](#)
- [MaestroSchedule\\$run\(\)](#)
- [MaestroSchedule\\$get\\_schedule\(\)](#)
- [MaestroSchedule\\$get\\_status\(\)](#)
- [MaestroSchedule\\$get\\_artifacts\(\)](#)
- [MaestroSchedule\\$clone\(\)](#)

**Method new():** Create a MaestroSchedule object

*Usage:*

```
MaestroSchedule$new(Pipelines = NULL)
```

*Arguments:*

Pipelines list of MaestroPipelines

*Returns:* MaestroSchedule

**Method print():** Print the schedule object

*Usage:*

```
MaestroSchedule$print()
```

*Returns:* print

**Method run():** Run a MaestroSchedule

*Usage:*

```
MaestroSchedule$run(..., quiet = FALSE, run_all = FALSE, n_show_next = 5)
```

*Arguments:*

... arguments passed to MaestroPipelineList\$run

quiet whether or not to silence console messages

run\_all run all pipelines regardless of the schedule (default is FALSE) - useful for testing.

n\_show\_next show the next n scheduled pipes

*Returns:* invisible

**Method get\_schedule():** Get the schedule as a data.frame

*Usage:*

```
MaestroSchedule$get_schedule()
```

*Returns:* data.frame

**Method get\_status():** Get status of the pipelines as a data.frame

*Usage:*

```
MaestroSchedule$get_status()
```

*Returns:* data.frame

**Method** `get_artifacts()`: Get artifacts (return values) from the pipelines

*Usage:*

```
MaestroSchedule$get_artifacts()
```

*Returns:* list

**Method** `clone()`: The objects of this class are cloneable with this method.

*Usage:*

```
MaestroSchedule$clone(deep = FALSE)
```

*Arguments:*

`deep` Whether to make a deep clone.

## Examples

```
if (interactive()) {
  pipeline_dir <- tempdir()
  create_pipeline("my_new_pipeline", pipeline_dir, open = FALSE)
  schedule <- build_schedule(pipeline_dir = pipeline_dir)
}
```

maestro\_tags

*Maestro Tags*

## Description

maestro tags are roxygen2 comments for configuring the scheduling and execution of pipelines.

## Details

maestro tags follow the format #' @maestroTagName

### Tag List:

tagName	description	value	examples (continued)
maestroFrequency	Time unit for scheduling	string	1 hour, daily, 3
maestroLogLevel	Threshold for logging when logging is requested	string	INFO, WARN
maestroSkip	Skips the pipeline when running (presence of tag indicates to skip)	n/a	
maestroStartTime	Start time of the pipeline; sets the point in time for recurrence	date or timestamp	1970-01-01 00:00:00
maestroTz	Timezone of the start time	string	UTC, America/New_York
maestroHours	Hours of day to run pipeline	ints	0 12 23
maestroDays	Days of week or days of month to run pipeline	ints or strings	1 14 30, Mon
maestroMonths	Months of year to run pipeline	ints	1 3 9 12

---

run_schedule	<i>Run a schedule</i>
--------------	-----------------------

---

## Description

Given a schedule in a `maestro` project, runs the pipelines that are scheduled to execute based on the current time.

## Usage

```
run_schedule(
  schedule,
  orch_frequency = "1 day",
  check_datetime = lubridate::now(tzone = "UTC"),
  resources = list(),
  run_all = FALSE,
  n_show_next = 5,
  cores = 1,
  logging = FALSE,
  log_file = NULL,
  log_file_max_bytes = 1e+06,
  quiet = FALSE
)
```

## Arguments

<code>schedule</code>	object of type <code>MaestroSchedule</code> created using <code>build_schedule()</code>
<code>orch_frequency</code>	of the orchestrator, a single string formatted like "1 day", "2 weeks", "hourly", etc.
<code>check_datetime</code>	datetime against which to check the running of pipelines (default is current system time in UTC)
<code>resources</code>	named list of shared resources made available to pipelines as needed
<code>run_all</code>	run all pipelines regardless of the schedule (default is FALSE) - useful for testing. Does not apply to pipes with a <code>maestroSkip</code> tag.
<code>n_show_next</code>	show the next n scheduled pipes
<code>cores</code>	number of cpu cores to run if running in parallel. If > 1, <code>furrr</code> is used and a multisession plan must be executed in the orchestrator (see details)
<code>logging</code>	whether or not to write the logs to a file (default = FALSE)
<code>log_file</code>	path to the log file (ignored if <code>logging == FALSE</code> )
<code>log_file_max_bytes</code>	numeric specifying the maximum number of bytes allowed in the log file before purging the log (within a margin of error)
<code>quiet</code>	silence metrics to the console (default = FALSE)

## Details

### Pipeline schedule logic:

The function `run_schedule()` examines each pipeline in the schedule table and determines whether it is scheduled to run at the current time using some simple time arithmetic. We assume `run_schedule(schedule, check_datetime = Sys.time())`, but this need not be the case.

### Output:

`run_schedule()` returns the same MaestroSchedule object with modified attributes. Use `get_status()` to examine the status of each pipeline and use `get_artifacts()` to get any return values from the pipelines as a list.

### Pipelines with arguments (resources):

If a pipeline takes an argument that doesn't include a default value, these can be supplied in the orchestrator via `run_schedule(resources = list(arg1 = val))`. The name of the argument used by the pipeline must match the name of the argument in the list. Currently, each named resource must refer to a single object. In other words, you can't have two pipes using the same argument but requiring different values.

### Running in parallel:

Pipelines can be run in parallel using the `cores` argument. First, you must run `future::plan(future::multisession)` in the orchestrator. Then, supply the desired number of cores to the `cores` argument. Note that console output appears different in multicore mode.

## Value

MaestroSchedule object

## Examples

```
if (interactive()) {  
  pipeline_dir <- tempdir()  
  create_pipeline("my_new_pipeline", pipeline_dir, open = FALSE)  
  schedule <- build_schedule(pipeline_dir = pipeline_dir)  
  
  # Runs the schedule every 1 day  
  run_schedule(  
    schedule,  
    orch_frequency = "1 day",  
    quiet = TRUE  
  )  
  
  # Runs the schedule every 15 minutes  
  run_schedule(  
    schedule,  
    orch_frequency = "15 minutes",  
    quiet = TRUE  
  )  
}
```

---

**suggest\_orch\_frequency**

*Suggest orchestrator frequency based on a schedule*

---

**Description**

Suggests a frequency to run the orchestrator based on the frequencies of the pipelines in a schedule.

**Usage**

```
suggest_orch_frequency(  
  schedule,  
  check_datetime = lubridate::now(tzone = "UTC")  
)
```

**Arguments**

schedule	MaestroSchedule object created by build_schedule()
check_datetime	datetime against which to check the running of pipelines (default is current system time in UTC)

**Details**

This function attempts to find the smallest interval of time between all pipelines. If the smallest interval is less than 15 minutes, it just uses the smallest interval.

Note this function is intended to be used interactively when deciding how often to schedule the orchestrator. Programmatic use is not recommended.

**Value**

frequency string

**Examples**

```
if (interactive()) {  
  pipeline_dir <- tempdir()  
  create_pipeline("my_new_pipeline", pipeline_dir, open = FALSE)  
  schedule <- build_schedule(pipeline_dir = pipeline_dir)  
  suggest_orch_frequency(schedule)  
}
```

# Index

build\_schedule, 2  
create\_maestro, 3  
create\_orchestrator, 4  
create\_pipeline, 4  
get\_artifacts, 6  
get\_schedule, 7  
get\_status, 7  
invoke, 8  
last\_build\_errors, 9  
last\_run\_errors, 10  
last\_run\_messages, 10  
last\_run\_warnings, 11  
maestro\_tags, 13  
MaestroSchedule, 11  
run\_schedule, 14  
suggest\_orch\_frequency, 16