

Package ‘flowr’

December 3, 2015

Type Package

Title Streamlining Design and Deployment of Complex Workflows

Description This framework allows you to design and implement complex pipelines, and deploy them on your institution's computing cluster. This has been built keeping in mind the needs of bioinformatics workflows. However, it is easily extendable to any field where a series of steps (shell commands) are to be executed in a (work)flow.

Version 0.9.9.4

Depends R (>= 3.0.2),
methods,
params (>= 0.3),
utils

Imports diagram,
whisker,
tools

Suggests reshape2,
knitr,
ggplot2,
openxlsx,
testthat,
funr

VignetteBuilder knitr

URL <https://github.com/sahilseth/flowr>

BugReports <https://github.com/sahilseth/flowr/issues>

License MIT + file LICENSE

RoxygenNote 5.0.1.9000

R topics documented:

| | |
|------------|---|
| check | 2 |
| check_args | 3 |
| fetch | 4 |
| flow-class | 5 |
| flowopts | 6 |
| get_wds | 8 |

| | |
|------------------------------|-----------|
| job | 8 |
| kill | 9 |
| plot_flow | 10 |
| queue-class | 12 |
| rerun | 13 |
| run | 14 |
| setup | 16 |
| status | 16 |
| submit_flow | 17 |
| submit_job | 18 |
| submit_run | 19 |
| test_queue | 19 |
| to_flow | 20 |
| to_flowdef | 22 |
| to_flowdet | 25 |
| to_flowmat | 26 |
| verbose | 27 |
| whisker_render | 28 |
| write_flow_details | 29 |
| Index | 30 |

| | |
|-------|---|
| check | <i>Check consistency of flowdef and flowmat</i> |
|-------|---|

Description

Check consistency of flowdef and flowmat, using various rules.

Usage

```
check(x, ...)  
  
## S3 method for class 'flowmat'  
check(x, ...)  
  
## S3 method for class 'flowdef'  
check(x, verbose = opts_flow$get("verbose"), ...)
```

Arguments

| | |
|---------|--|
| x | a flowdef or flowmat object |
| ... | Passed onto either check.flowdef OR check.flowmat functions |
| verbose | A numeric value indicating the amount of messages to produce. Values are integers varying from 0, 1, 2, 3, Please refer to the verbose page for more details. opts_flow\$get("verbose") |

Details

A typical output from flowdef with verbose level: 2

```
checking if required columns are present...
checking if resources columns are present...
checking if dependency column has valid names...
checking if submission column has valid names...
checking for missing rows in def...
checking for extra rows in def...
checking submission and dependency types...
jobname prev.sub_type --> dep_type --> sub_type: relationship
1: aln1_a none --> none --> scatter
2: aln2_a scatter --> none --> scatter
3: sampe_a scatter --> serial --> scatter rel: complex one:one
4: fixrg_a scatter --> serial --> scatter rel: complex one:one
5: merge_a scatter --> gather --> serial rel: many:one
6: markdup_a serial --> serial --> serial rel: simple one:one
7: target_a serial --> serial --> serial rel: simple one:one
8: realign_a serial --> burst --> scatter rel: one:many
9: baserecalib_a scatter --> serial --> scatter rel: complex one:one
10: printreads_a scatter --> serial --> scatter rel: complex one:one
```

check_args

Assert none of the arguemnts of a function are null.

Description

Checks all the arguments in the parent function and makes sure that none of them are NULL

Usage

```
check_args(ignore, select)
```

Arguments

| | |
|--------|--|
| ignore | optionally ignore a few variables for checking. |
| select | optionally only check a few variables of the function. |

Details

This function has now been moved to params package.

 fetch

Two generic functions to search for pipelines and configuration files.

Description

These functions help in searching for specific files in the user's space.

`fetch_pipes()`: Fetches pipelines in the following places, in this specific order:

- **user's folder**: `~/flowr/pipelines`
- **current wd**: `./`

NOTE: If same pipeline is available in multiple places; intuitively, one from the later folder would be selected. As such, giving priority to user's home, and current working directories.

`
 fetch_conf()`: Fetches configuration files in ALL of the following places:

- **package**: conf folders in flowr and ngsflows packages.
- **user's folder**: `~/flowr/conf` folder.
- **current wd**: `./`

NOTE: This function would greedily return all matching conf files. One would load all of them in the order returned by this function. If the same variable is repeated in multiple files, value from later files would replace those formerly defined. Thus (as explained above), giving priority to options defined in user's home and current working directories.

By default flowr loads, `flowr.conf` and `ngsflows.conf`. See the details sections, for more explanation on this.

Usage

```
fetch(x, places, urls, verbose = opts_flow$get("verbose"))
```

```
fetch_pipes(x, places, last_only = FALSE,
  urls = opts_flow$get("flowr_pipe_urls"), silent = FALSE,
  verbose = opts_flow$get("verbose"), ask = TRUE)
```

```
fetch_conf(x = "flowr.conf", places, ...)
```

Arguments

| | |
|------------------------|---|
| <code>x</code> | name of the file to search for (without extension). By default fetch_pipes and fetch_conf search for files ending with <code>.R</code> and <code>.conf</code> respectively. |
| <code>places</code> | places (paths) to look for files matching the name. Defaults are already defined in the function. |
| <code>urls</code> | urls to look for, works well for pipelines [not implemented yet] |
| <code>verbose</code> | A numeric value indicating the amount of messages to produce. Values are integers varying from 0, 1, 2, 3, Please refer to the verbose page for more details. <code>opts_flow\$get("verbose")</code> |
| <code>last_only</code> | <code>fetch_pipes()</code> : If multiple pipelines match the pattern, return the last one. [TRUE] |
| <code>silent</code> | <code>fetch_pipes()</code> : logical, be silent even if no such pipeline is available. [FALSE] |
| <code>ask</code> | ask before downloading or copying. [not implemented] |
| <code>...</code> | [not implemented] |

Details

For example flowr has a variable flow_run_path where it puts all the execution logs etc. The default value is picked up from the internal **flowr.conf** file. To redefine this value, one could create a new file called **~/flowr/conf/flowr.conf** and add a line:

flow_run_path TAB my_awesome_path, where TAB is a tab character, since these are tab seperated files.

Also, at any time you can run, [opts_flow\\$load](#); to load custom options.

See Also

[flowopts](#)

Examples

```
## let us find a default conf file
conf = fetch_conf("flowr.conf");conf
## load this
opts_flow$load(conf)

## this returns a list, which prints pretty
pip = fetch_pipes("sleep_pipe")
pip$name
pip$pipe
pip$def
```

| | |
|------------|----------------------------------|
| flow-class | <i>Describing the flow class</i> |
|------------|----------------------------------|

Description

Internal function (used by [to_flow](#)), which aids in creating a flow object.

Usage

```
flow(jobs = list(new("job")), name = "newflow", desc = "my_super_flow",
     mode = c("scheduler", "trigger", "R"),
     flow_run_path = opts_flow$get("flow_run_path"), trigger_path = "",
     flow_path = "", version = "0.0", status = "created",
     module_cmds = opts_flow$get("module_cmds"), execute = "")
```

Arguments

| | |
|---------------|---|
| jobs | list: A list of jobs to be included in this flow |
| name | character: Name of the flow. ['newflow'] |
| desc | character Description of the flow, used to uniquely identify a flow instance. ['my_super_flow'] |
| mode | character Mode of submission of the flow (deprecated). ['scheduler'] |
| flow_run_path | The base path of all the flows you would submit. [~/flows] |

| | |
|--------------|---|
| trigger_path | character [~/flows/trigger]. |
| flow_path | character: A unique path identifying a flow instance, populated by submit_flow . |
| version | version of flowr used to create and execute this flow. |
| status | character: Status of the flow. |
| module_cmds | [advanced use] a character vector of cmds which will be pre-pended to all script of this pipeline. Could be cmds like `module load mytool1;module load mytool2` |
| execute | execution status of flow object. [FALSE] |

Examples

```
cmds = rep("sleep 5", 10)
qobj <- queue(platform='torque')
## run the 10 commands in parallel
jobj1 <- job(q_obj=qobj, cmd = cmds, submission_type = "scatter", name = "job1")

## run the 10 commands sequentially, but WAIT for the previous job to complete
## Many-To-One
jobj2 <- job(q_obj=qobj, cmd = cmds, submission_type = "serial",
  dependency_type = "gather", previous_job = "job1", name = "job2")

## As soon as first job on 'job1' is complete
## One-To-One
jobj3 <- job(q_obj=qobj, cmd = cmds, submission_type = "scatter",
  dependency_type = "serial", previous_job = "job1", name = "job3")

fobj <- flow(jobs = list(jobj1, jobj2, jobj3))

## plot the flow
plot_flow(fobj)
## Not run:
## dry run, only create the structure without submitting jobs
submit_flow(fobj)

## execute the jobs: ONLY works on computing cluster, would fail otherwise
submit_flow(fobj, execute = TRUE)

## End(Not run)
```

flowopts

Default options/params used in flowr and ngsflows

Description

There are three helper functions which attempt to manage parameters used by flowr and ngsflows:

- [get_opts](#) OR `opts_flow$get()`: show all default options
- [set_opts](#) OR `opts_flow$set()`: set default options
- [load_opts](#) OR `opts_flow$load()`: load options specified in a tab separated text file

For more details regarding these funtions refer to [params](#) package.

Usage

```
flowopts
```

```
get_opts(...)
set_opts(...)
load_opts(...)
```

Arguments

- ...
- get: names of options to fetch
 - set: a set of options in a name=value format separated by commas

Format

```
opts_flow
```

Details

By default flowr loads, `~/flowr/conf/flowr.conf` and `~/flowr/conf/ngsflows.conf`

Below is a list of default flowr options, retrieved via

```
opts_flow$get():
```

| name | value | |
|-------------------|--------------------------|--|
| :----- | :----- | |
| default_regex | (.*) | |
| flow_base_path | ~/flowr | |
| flow_conf_path | ~/flowr/conf | |
| flow_parse_lsf | .*(\<[0-9]*\>).* | |
| flow_parse_moab | (.*) | |
| flow_parse_sge | (.*) | |
| flow_parse_slurm | (.*) | |
| flow_parse_torque | (.?)\..* | |
| flow_pipe_paths | ~/flowr/pipelines | |
| flow_pipe_urls | ~/flowr/pipelines | |
| flow_platform | local | |
| flow_run_path | ~/flowr/runs | |
| my_conf_path | ~/flowr/conf | |
| my_dir | path/to/a/folder | |
| my_path | ~/flowr | |
| my_tool_exe | /usr/bin/ls | |
| time_format | %a %b %e %H:%M:%S CDT %Y | |
| verbose | FALSE | |

See Also

[fetch params read_sheet](#)

Examples

```
## Set options: opts_flow$set()
opts = opts_flow$set(flow_run_path = "~/mypath")
## OR if you would like to supply a long list of options:
opts = opts_flow$set(.dots = list(flow_run_path = "~/mypath"))
```

```
## load options from a configuration file: opts_flow$load()
conffile = fetch_conf("flowr.conf")
opts_flow$load(conffile)

## Fetch options: get_opts()
opts_flow$get("flow_run_path")
opts_flow$get()
```

| | |
|---------|---|
| get_wds | <i>Get all the (sub)directories in a folder</i> |
|---------|---|

Description

Get all the (sub)directories in a folder

Usage

```
get_wds(x)
```

Arguments

x path to a folder

| | |
|-----|---|
| job | <i>Describing details of the job object</i> |
|-----|---|

Description

Internal function (used by to_flow), which aids in creating a job object.

Usage

```
job(cmds = "", name = "myjob", q_obj = new("queue"), previous_job = "",
    cpu = 1, memory, walltime, submission_type = c("scatter", "serial"),
    dependency_type = c("none", "gather", "serial", "burst"), ...)
```

Arguments

| | |
|--------------|--|
| cmds | the commands to run |
| name | name of the job |
| q_obj | queue object |
| previous_job | character vector of previous job. If this is the first job, one can leave this empty, NA, NULL, '.', or ''. In future this could specify multiple previous jobs. |
| cpu | no of cpu's reserved |
| memory | The amount of memory reserved. Units depend on the platform used to process jobs |

walltime The amount of time reserved for this job. Format is unique to a platform. Typically it looks like 12:00 (12 hours reserved, say in LSF), in Torque etc. we often see measuring in seconds: 12:00:00

submission_type submission type: A character with values: scatter, serial. Scatter means all the 'cmds' would be run in parallel as separate jobs. Serial, they would be combined into a single job and run one-by-one.

dependency_type dependency type. One of none, gather, serial, burst. If previous_job is specified, then this would not be 'none'. [Required]

... other passed onto object creation. Example: memory, walltime, cpu

Examples

```
qobj <- queue(platform="torque")

## torque job with 1 CPU running command 'sleep 2'
jobj <- job(q_obj=qobj, cmd = "sleep 2", cpu=1)

## multiple commands
cmds = rep("sleep 5", 10)

## run the 10 commands in parallel
jobj1 <- job(q_obj=qobj, cmd = cmds, submission_type = "scatter", name = "job1")

## run the 10 commands sequentially, but WAIT for the previous job to complete
jobj2 <- job(q_obj=qobj, cmd = cmds, submission_type = "serial",
  dependency_type = "gather", previous_job = "job1")

fobj <- flow(jobs = list(jobj1, jobj2))

## plot the flow
plot_flow(fobj)
## Not run:
## dry run, only create the structure without submitting jobs
submit_flow(fobj)

## execute the jobs: ONLY works on computing cluster, would fail otherwise
submit_flow(fobj, execute = TRUE)

## End(Not run)
```

| | |
|------|---|
| kill | <i>Kill all jobs submitted to the computing platform, for one or multiple flows</i> |
|------|---|

Description

NOTE:

This requires files which are created at the end of the [submit_flow](#) command.

Even if you want to kill the flow, its best to let submit_flow do its job, when done simply use kill(flow_wd). If submit_flow is interrupted, files like flow_details.rds etc are not created, thus flowr loses the association of jobs with flow instance and cannot monitor, kill or re-run the flow.

Usage

```
kill(x, ...)

## S3 method for class 'character'
kill(x, force = FALSE, ...)

## S3 method for class 'flow'
kill(x, kill_cmd, verbose = opts_flow$get("verbose"),
     jobid_col = "job_sub_id", ...)
```

Arguments

| | |
|-----------|---|
| x | either path to flow wd or object of class flow |
| ... | not used |
| force | You need to set force=TRUE, to kill multiple flows. This makes sure multiple flows are NOT killed by accident. |
| kill_cmd | The command used to kill. flowr tries to guess this commands, as defined in the <code>detect_kill_cmd()</code> . Supplying it here; fot custom platoforms. |
| verbose | A numeric value indicating the amount of messages to produce. Values are integers varying from 0, 1, 2, 3, Please refer to the verbose page for more details. <code>opts_flow\$get("verbose")</code> |
| jobid_col | Advanced use. The column name in 'flow_details.txt' file used to fetch jobids to kill |

Examples

```
## Not run:

## example for terminal
## flowr kill_flow x=path_to_flow_directory
## In case path matches multiple folders, flowr asks before killing
kill(x='fastq_haplotyper*')
Flowr: streamlining workflows
found multiple wds:
/fastq_haplotyper-MS132-20150825-16-24-04-0Lv1PbpI
/fastq_haplotyper-MS132-20150825-17-47-52-5vFIkrMD
Really kill all of them ? kill again with force=TRUE

## submitting again with force=TRUE will kill them:
kill(x='fastq_haplotyper*', force = TRUE)

## End(Not run)
```

plot_flow

Plot a clean and scalable flowchart describing the (work)flow

Description

Plot a flowchart using a flow object or flowdef

Usage

```

plot_flow(x, ...)

## S3 method for class 'flow'
plot_flow(x, ...)

## S3 method for class 'list'
plot_flow(x, ...)

## S3 method for class 'character'
plot_flow(x, ...)

## S3 method for class 'flowdef'
plot_flow(x, detailed = TRUE, type = c("1", "2"),
  pdf = FALSE, pdffile, ...)

## S3 method for class 'flowdef'
plot(x, detailed = TRUE, type = c("1", "2"),
  pdf = FALSE, pdffile, ...)

## S3 method for class 'flow'
plot(x, ...)

```

Arguments

| | |
|----------|---|
| x | Object of class flow, or a list of flow objects or a flowdef |
| ... | experimental and only for advanced use. |
| detailed | include submission and dependency types in the plot [TRUE] |
| type | 1 is original, and 2 is a ellipse with less details [1] |
| pdf | create a pdf instead of plotting interactively [FALSE] |
| pdffile | output file name for the pdf file. [flow_path/flow_details.pdf] |

Examples

```

qobj = queue(type="lsf")
cmds = rep("sleep 5", 10)
jobj1 <- job(q_obj=qobj, cmd = cmds, submission_type = "scatter", name = "job1")
jobj2 <- job(q_obj=qobj, name = "job2", cmd = cmds, submission_type = "scatter",
  dependency_type = "serial", previous_job = "job1")
fobj <- flow(jobs = list(jobj1, jobj2))
plot_flow(fobj)

### Gather: many to one relationship
jobj1 <- job(q_obj=qobj, cmd = cmds, submission_type = "scatter", name = "job1")
jobj2 <- job(q_obj=qobj, name = "job2", cmd = cmds, submission_type = "scatter",
  dependency_type = "gather", previous_job = "job1")
fobj <- flow(jobs = list(jobj1, jobj2))
plot_flow(fobj)

### Burst: one to many relationship
jobj1 <- job(q_obj=qobj, cmd = cmds, submission_type = "serial", name = "job1")
jobj2 <- job(q_obj=qobj, name = "job2", cmd = cmds, submission_type = "scatter",

```

```

        dependency_type = "burst", previous_job = "job1")
fobj <- flow(jobs = list(jobj1, jobj2))
plot_flow(fobj)

```

queue-class

A queue object defines details regarding how a job is submitted

Description

Internal function (used by [to_flow](#)), to define the format used to submit a job.

Usage

```

queue(object, platform = c("local", "lsf", "torque", "sge", "moab"),
      format = "", queue = "long", walltime, memory, cpu = 1,
      extra_opts = "", submit_exe, nodes = "1", jobname = "name",
      email = Sys.getenv("USER"), dependency = list(), server = "localhost",
      verbose = FALSE, cwd = "", stderr = "", stdout = "", ...)

```

Arguments

| | |
|------------|---|
| object | this is not used currently, ignore. |
| platform | Required and important. Currently supported values are 'lsf' and 'torque'. [Used by class job] |
| format | [advanced use] We have a default format for the final command line string generated for 'lsf' and 'torque'. |
| queue | the type of queue your group usually uses 'bsub' etc. |
| walltime | max walltime of a job. |
| memory | The amount of memory reserved. Units depend on the platform used to process jobs |
| cpu | number of cpus you would like to reserve [Used by class job] |
| extra_opts | [advanced use] Extra options to be supplied while create the job submission string. |
| submit_exe | [advanced use] Already defined by 'platform'. The exact command used to submit jobs to the cluster example 'qsub' |
| nodes | [advanced use] number of nodes you would like to request. Or in case of torque name of the nodes. <i>optional</i> [Used by class job] |
| jobname | [debug use] name of this job in the computing cluster |
| email | [advanced use] Defaults to system user, you may put you own email though may get tons of them. |
| dependency | [debug use] a list of jobs to complete before starting this one |
| server | [not used] This is not implemented currently. This would specify the head node of the computing cluster. At this time submission needs to be done on the head node of the cluster where flow is to be submitted |
| verbose | [logical] TRUE/FALSE |
| cwd | [debug use] Ignore |
| stderr | [debug use] Ignore |
| stdout | [debug use] Ignore |
| ... | other passed onto object creation. Example: memory, walltime, cpu |

Details

Resources: Can be defined ***once*** using a [queue](#) object and recycled to all the jobs in a flow. If resources (like memory, cpu, walltime, queue) are supplied at the job level they overwrite the one supplied in [queue](#) Nodes: can be supplied or extend a job across multiple nodes. This is purely experimental and not supported.

Server: This is a hook which may be implemented in future.

Submission script The 'platform' variable defines the format, and submit_exe; however these two are available for someone to create a custom submission command.

Examples

```
qobj <- queue(platform='lsf')
```

| | |
|-------|--|
| rerun | <i>Re-run a pipeline in case of hardware or software failures.</i> |
|-------|--|

Description

- **hardware** no change required, simple rerun: `rerun(x=flow_wd)`
- **software** either a change to flowmat or flowdef has been made: `rerun(x=flow_wd, mat = new_flowmat, def = r`

NOTE:

flow_wd: flow working directory, same input as used for [status](#)

Usage

```
rerun(x, ...)

## S3 method for class 'character'
rerun(x, ...)

## S3 method for class 'flow'
rerun(x, mat, def, start_from, samplename, execute = TRUE,
      kill = TRUE, select, ignore, verbose = opts_flow$get("verbose"), ...)
```

Arguments

| | |
|------------|---|
| x | flow working directory |
| ... | passed onto to_flow |
| mat | (optional) flowmat fetched from previous submission if missing. For more information regarding the format refer to to_flowmat |
| def | (optional) flowdef fetched from previous submission if missing. For more information regarding the format refer to to_flowdef |
| start_from | (required) which job to start from, this is a job name. |
| samplename | (optional) If flowmat contains multiple samples, provide the samplename, corresponding to the flow working directory provided. |
| execute | [logical] whether to execute or not |

| | |
|---------|---|
| kill | (optional) logical indicating whether to kill the jobs from the previous execution of flow. |
| select | (optional) select a subset of jobs to rerun [character vector] |
| ignore | (optional) ignore a subset of jobs to rerun [character vector] |
| verbose | A numeric value indicating the amount of messages to produce. Values are integers varying from 0, 1, 2, 3, Please refer to the verbose page for more details. <code>opts_flow\$get("verbose")</code> |

Details

This function fetches details regarding the previous execution from the flow working directory (`flow_wd`).

It reads the [flow](#) object from the `flow_details.rds` file, and extracts `flowdef` and `flowmat` from it using [to_flowmat](#) and [to_flowdef](#) functions.

Using new flowmat OR flowdef for re-run:

Optionally, if either of `flowmat` or `flowdef` are supplied; supplied ones are used instead of those extracted from previous submission.

This functions efficiently updates job details of the latest submission into the previous file; thus information regarding previous job ids and their status is not lost.

Examples

```
## Not run:

#
rerun(wd = wd, fobj = fobj, execute = TRUE, kill = TRUE)

## End(Not run)
```

run

Run automated Pipelines

Description

Run complete pipelines, by wrapping several steps into one convinient function.

NOTE: please use flowr version 0.9.8.9010 or higher.

In summary, this function performs the following steps:

- the argument `x` defines the name of the pipeline. Say, for example `sleep_pipe`.
- [fetch_pipes](#): finds the pipeline definition (`sleep_pipe.R`, `sleep_pipe.def` and `sleep_pipe.conf` files)
- `sleep_pipe(...)`: Create all the required commands (`flowmat`)
- [to_flow](#): Use `flowmat` and `sleep_pipe.def` to create a flow object.
- [submit_flow](#): Submit the flow to the cluster.

Usage

```
run(x, platform, def, conf, wd = opts_flow$get("flow_run_path"),
    flow_run_path = wd, rerun_wd, start_from, execute = FALSE, ...)

run_pipe(x, platform, def, conf, wd = opts_flow$get("flow_run_path"),
    flow_run_path = wd, rerun_wd, start_from, execute = FALSE, ...)
```

Arguments

| | |
|---------------|--|
| x | name of the pipeline to run. This is a function called to create a flow_mat. |
| platform | what platform to use, overrides flowdef |
| def | flow definition |
| conf | a tab-delimited configuration file with path to tools and default parameters. See fetch_pipes . |
| wd | an alias to flow_run_path |
| flow_run_path | passed onto to_flow. Default it picked up from flowr.conf. Typically this is ~/flowr/runs |
| rerun_wd | if you need to run, supply the previous working dir |
| start_from | the step to start a rerun from. Intitively, this is ignored in a fresh run and only used in re-running a pipeline. |
| execute | TRUE/FALSE |
| ... | passed onto the pipeline function as specified in x |

Examples

```
## Not run:

## Run a short pipeline (dry run)
run("sleep_pipe")

## Run a short pipeline on the local machine
run("sleep_pipe", platform = "local", execute = TRUE)

## Run a short pipeline on the a torque cluster (qsub)
run("sleep_pipe", platform = "torque", execute = TRUE)

## Run a short pipeline on the a MOAB cluster (msub)
run("sleep_pipe", platform = "moab", execute = TRUE)

## Run a short pipeline on the a IBM (LSF) cluster (bsub)
run("sleep_pipe", platform = "lsf", execute = TRUE)

## Run a short pipeline on the a MOAB cluster (msub)
run("sleep_pipe", platform = "moab", execute = TRUE)

## change parameters of the pipeline
## All extra parameters are passed on to the function function.
run("sleep_pipe", platform = "lsf", execute = TRUE, x = 5)

## End(Not run)
```

| | |
|-------|-----------------------------------|
| setup | <i>Setup and initialize flowr</i> |
|-------|-----------------------------------|

Description

This functions creates a directory structure in user's home directory. Additionally it creates a short-cut to the flowr helper script in ~/bin.

Usage

```
setup(bin = "~/bin", flow_base_path = opts_flow$get("flow_base_path"),
      flow_run_path = opts_flow$get("flow_run_path"),
      flow_conf_path = opts_flow$get("flow_conf_path"),
      flow_pipe_path = opts_flow$get("flow_pipe_paths"))
```

Arguments

| | |
|----------------|--|
| bin | path to bin folder |
| flow_base_path | The base of flowr configuration and execution folders. |
| flow_run_path | base path to be used for execution of this flow. flowr would create a new time-stamped folder in this base path and use it for logs, scripts etc. The default is retrived using <code>opts_flow\$get("flow_run_path")</code> . |
| flow_conf_path | Flowr configuration folder, used by fetch_conf . |
| flow_pipe_path | Folder with all pipelines, used by fetch_pipes . |

Details

Will add more to this, to identify cluster and aid in other things.

| | |
|--------|----------------------------------|
| status | <i>Monitor status of flow(s)</i> |
|--------|----------------------------------|

Description

Summarize status of a flow OR multiple flows OR a high-level summary of all flows in a folder.

Usage

```
status(x, use_cache = FALSE, verbose = opts_flow$get("verbose"),
      out_format = "markdown", ...)

get_status(x, ...)

## S3 method for class 'flow'
get_status(x, verbose, use_cache, out_format, ...)

## S3 method for class 'character'
get_status(x, verbose, use_cache, out_format, ...)

## S3 method for class 'data.frame'
get_status(x, verbose, use_cache, progress = TRUE, ...)
```


Arguments

| | |
|------------|---|
| x | path to the flow root folder or a parent folder to summarize several flows. |
| use_cache | This skips checking status of jobs which have already been completed and assumes no new jobs were submitted in the flow(s) being monitored. [FALSE] |
| verbose | A numeric value indicating the amount of messages to produce. Values are integers varying from 0, 1, 2, 3, Please refer to the verbose page for more details. <code>opts_flow\$get("verbose")</code> |
| out_format | passed onto <code>knitr::kable</code> . supports: markdown, rst, html... [markdown] |
| ... | not used |
| progress | Whether or not to show a progress bar, when fetching/reading files [TRUE] |

Details

`basename(x)` is used in a wild card search.

- Get status of all the flows: (all flows with 'sleep_pipe' in their name are checked and their status is shown)
 `flowr status x=~/flowr/runs/sleep_pipe*`
- Provide a high level summary of ALL flows in a folder:
 `flowr status x=~/flowr/runs`

Use **use_cache=TRUE** to speed up checking the status. This assumes that no new jobs have been submitted and skips (re-)checking status of completed jobs.

Once all the jobs have been submitted to the cluster you may always use `use_cache=TRUE`.

Examples

```
## Not run:
status(x = "~/flowr/runs/sleep_pipe")
## an example for running from terminal
flowr status x=path_to_flow_directory

## End(Not run)
```

submit_flow

Submit a flow to the cluster

Description

Submit a flow to the cluster or perform a dry-run to check and debug issues.

Usage

```
submit_flow(x, verbose = opts_flow$get("verbose"), ...)

## S3 method for class 'list'
submit_flow(x, verbose = opts_flow$get("verbose"), ...)

## S3 method for class 'flow'
submit_flow(x, verbose = opts_flow$get("verbose"),
  execute = FALSE, uuid, plot = TRUE, dump = TRUE, .start_jid = 1, ...)
```

Arguments

| | |
|------------|---|
| x | a object of class flow. |
| verbose | logical. |
| ... | Advanced use. Any additional parameters are passed on to submit_job function. |
| execute | logical whether or not to submit the jobs |
| uuid | character Advanced use. This is the final path used for flow execution. Especially useful in case of re-running a flow. |
| plot | logical whether to make a pdf flow plot (saves it in the flow working directory). |
| dump | dump all the flow details to the flow path |
| .start_jid | Job to start this submission from. Advanced use, should be 1 by default. |

Details

NOTE: Even if you want to kill the flow, its best to let submit_flow do its job, when done simply use kill(flow_wd). If submit_flow is interrupted, files like flow_details.rds etc are not created, thus flowr loses the association of jobs with flow instance and cannot monitor, kill or re-run the flow.

Examples

```
## Not run:
submit_flow(fobj = fobj, ... = ...)
## End(Not run)
```

| | |
|------------|--------------------------------|
| submit_job | <i>Submit a step of a flow</i> |
|------------|--------------------------------|

Description

Internal function (used by submit_flow), which submit a single step of a flow.

Usage

```
submit_job(jobj, fobj, job_id, execute = FALSE, verbose = FALSE, ...)
```

Arguments

| | |
|---------|--|
| jobj | Object of calls job |
| fobj | Object of calls flow |
| job_id | job id |
| execute | A logical vector suggesting whether to submit this job |
| verbose | logical |
| ... | not used |

Examples

```
## Not run:
submit_job(jobj = jobj, fobj = fobj, execute = FALSE,
verbose = TRUE, wd = wd, job_id = job_id)

## End(Not run)
```

| | |
|------------|--|
| submit_run | <i>Submit several flow objects, limit the max running concurrently</i> |
|------------|--|

Description

Submit several flow objects, limit the max running concurrently

Usage

```
submit_run(x, wd, max_processing = 7)
```

Arguments

| | |
|----------------|--|
| x | a list of flow objects |
| wd | a folder to monitor (flow_run_path) |
| max_processing | max number of flow which may be processed concurrently |

| | |
|------------|-------------------|
| test_queue | <i>test_queue</i> |
|------------|-------------------|

Description

This function attempts to test the submission of a job to the queue. We would first submit one single job, then submit another with a dependency to see if configuration works. This would create a folder in home called 'flows'.

[Deprecated]: This function has been superseded by `run("sleep_pipe", platform = "lsf", execute=TRUE)`

Usage

```
test_queue(q_obj, verbose = TRUE, ...)
```

Arguments

| | |
|---------|--|
| q_obj | queue object |
| verbose | toggle |
| ... | These params are passed onto queue. ?queue, for more information |

Examples

```
## Not run:
test_queue(q_obj = q_obj, ... = ...)
## End(Not run)
```

to_flow

*Create flow objects***Description**

Use a set of shell commands (flow mat) and flow definiton to create [flow](#) object.

Usage

```
to_flow(x, ...)
```

```
is.flow(x)
```

```
## S3 method for class 'character'
```

```
to_flow(x, def, grp_col, jobname_col, cmd_col, ...)
```

```
## S3 method for class 'flowmat'
```

```
to_flow(x, def, flowname, grp_col, jobname_col, cmd_col,
        submit = FALSE, execute = FALSE, containerize = TRUE, platform,
        flow_run_path, qobj, verbose = opts_flow$get("verbose"), ...)
```

```
## S3 method for class 'data.frame'
```

```
to_flow(x, ...)
```

```
## S3 method for class 'list'
```

```
to_flow(x, def, flowname, flow_run_path, desc, qobj,
        module_cmds = opts_flow$get("module_cmds"),
        verbose = opts_flow$get("verbose"), ...)
```

Arguments

| | |
|-------------|---|
| x | this can either to a filename, a data.frame or a list. In case it is a file name, it should be a tsv file representing a flow_mat. See to_flowmat for details |
| ... | Supplied to specific functions like to_flow.data.frame |
| def | a flow definition. Basically a table with resource requirements and mapping of the jobs in this flow. See to_flowdef for details on the format. |
| grp_col | name of the grouping column in the supplied flow_mat. See to_flow for details. Default value is [samplename]. |
| jobname_col | name of the job name columnd in flow_mat. Defalt value is [jobname]. |
| cmd_col | name of the command column name in flow_mat. Default value is [cmd]. |
| flowname | name of the flow, this is used as part of the execution foldername. A good simple identifier, which does not support any special characters. Names may use characters (a-z) and numbers (0-9), using underscore (_) as a word seperator. Default value is [flowname]. |
| submit | after creating a flow object, should flowr also use submit_flow to perform a dry-run OR real submission. See below for details. Default value is [FALSE] |
| execute | when calling submit_flow , should flowr execute the flow or perform a dry-run. See below for details. Default value is [FALSE]. |

| | |
|---------------|--|
| containerize | if the flowmat has multiple samples, flowr creates a new date-stamped folder, and includes all flows in this batch inside it. This keeps the logs clean, and containerizes each batch. To disable this behaviour set this to FALSE, default is [TRUE]. |
| platform | a specifying the platform to use, possible values are local, lsf, torque, moab, sge and slurm This over-rides the platform column in the flowdef. (optional) |
| flow_run_path | base path to be used for execution of this flow. flowr would create a new time-stamped folder in this base path and use it for logs, scripts etc. The default is retrieved using <code>opts_flow\$get("flow_run_path")</code> . |
| qobj | Deprecated, modify cluster templates as explained on docs.flowr.space . An object of class queue . |
| verbose | A numeric value indicating the amount of messages to produce. Values are integers varying from 0, 1, 2, 3, Please refer to the verbose page for more details. <code>opts_flow\$get("verbose")</code> |
| desc | Advanced Use. final flow name. |
| module_cmds | A character vector of additional commands, which will be prepended to each script of the flow. Default is retrieved using <code>opts_flow\$get("module_cmds")</code> . |

Details

The parameter `x` can be a path to a `flow_mat`, or a `data.frame` (as read by `read_sheet`). This is a minimum three column table with columns: `samplename`, `jobname` and `cmd`. See [to_flowmat](#) for details.

Value

Returns a flow object. If `execute=TRUE`, `fobj` is rich with information about where and how the flow was executed. It would include details like jobids, path to exact scripts run etc. To use `kill_flow`, to kill all the jobs one would need a rich flow object, with job ids present.

Behaviour: What goes in, and what to expect in return?

- `submit=FALSE & execute=FALSE`: Create and return a flow object
- `submit=TRUE & execute=FALSE`: dry-run, Create a flow object then, create a structured execution folder with all the commands
- `submit=TRUE, execute=TRUE`: Do all of the above and then, submit to cluster

See Also

[to_flowmat](#), [to_flowdef](#), [to_flowdet](#), [flowopts](#) and [submit_flow](#)

Examples

```
## Use this link for a few elaborate examples:
## http://docs.flowr.space/flowr/tutorial.html#define_modules

ex = file.path(system.file(package = "flowr"), "pipelines")
flowmat = as.flowmat(file.path(ex, "sleep_pipe.tsv"))
flowdef = as.flowdef(file.path(ex, "sleep_pipe.def"))
fobj = to_flow(x = flowmat, def = flowdef, flowname = "sleep_pipe", platform = "lsf")

## create a vector of shell commands
```

```

cmds = c("sleep 1", "sleep 2")
## create a named list
lst = list("sleep" = cmds)
## create a flowmat
flowmat = to_flowmat(lst, samplename = "samp")

## Use flowmat to create a skeleton flowdef
flowdef = to_flowdef(flowmat)

## use both (flowmat and flowdef) to create a flow
fobj = to_flow(flowmat, flowdef)

## submit the flow to the cluster (execute=TRUE) or do a dry-run (execute=FALSE)
## Not run:
fobj2 = submit_flow(fobj, execute=FALSE)
fobj3 = submit_flow(fobj, execute=TRUE)

## Get the status or kill all the jobs
status(fobj3)
kill(fobj3)

## End(Not run)

```

to_flowdef

Flow Definition defines how to stitch steps into a (work)flow.

Description

This function enables creation of a skeleton flow definition with several default values, using a flowmat. To customize the flowdef, one may supply parameters such as sub_type and dep_type upfront. As such, these params must be of the same length as number of unique jobs using in the flowmat.

Each row in this table refers to one step of the pipeline. It describes the resources used by the step and also its relationship with other steps, especially, the step immediately prior to it.

Submission types: *This refers to the sub_type column in flow definition.*

Consider an example with three steps A, B and C. A has 10 commands from A1 to A10, similarly B has 10 commands B1 through B10 and C has a single command, C1. Consider another step D (with D1-D3), which comes after C.

step (number of sub-processes) A (10) —> B (10) —> C (1) —> D (3)

- scatter: submit all commands as parallel, independent jobs.
Submit A1 through A10 as independent jobs
- serial: run these commands sequentially one after the other.
- Wrap A1 through A10, into a single job.

Dependency types

This refers to the dep_type column in flow definition.

- none: independent job.

- *Initial step A has no dependency*
- **serial:** *one to one* relationship with previous job.
 - *B1 can start as soon as A1 completes, and B2 starts just after A2 and so on.*
- **gather:** *many to one*, wait for **all** commands in the previous job to finish then start the current step.
 - *All jobs of B (1-10), need to complete before C1 starts*
- **burst:** *one to many* wait for the previous step which has one job and start processing all cmds in the current step.
 - *D1 to D3 are started as soon as C1 finishes.*

Usage

```
to_flowdef(x, ...)

## S3 method for class 'flowmat'
to_flowdef(x, sub_type, dep_type, prev_jobs,
  queue = "short", platform = "torque", memory_reserved = "2000",
  cpu_reserved = "1", nodes = "1", walltime = "1:00", guess = FALSE,
  verbose = opts_flow$get("verbose"), ...)

## S3 method for class 'flow'
to_flowdef(x, ...)

## S3 method for class 'character'
to_flowdef(x, ...)

as.flowdef(x, ...)

is.flowdef(x)
```

Arguments

| | |
|-----------------|---|
| x | can a path to a flowmat, flowmat or flow object. |
| ... | not used |
| sub_type | submission type, one of: scatter, serial. Character, of length one or same as the number of jobnames |
| dep_type | dependency type, one of: gather, serial or burst. Character, of length one or same as the number of jobnames |
| prev_jobs | previous job name |
| queue | Cluster queue to be used |
| platform | platform of the cluster: lsf, sge, moab, torque, slurm etc. |
| memory_reserved | amount of memory required. |
| cpu_reserved | number of cpu's required. [1] |
| nodes | if you tool can use multiple nodes, you may reserve multiple nodes for it. [1] |
| walltime | amount of walltime required |
| guess | should the function, guess submission and dependency types. See details. |
| verbose | A numeric value indicating the amount of messages to produce. Values are integers varying from 0, 1, 2, 3, Please refer to the verbose page for more details. <code>opts_flow\$get("verbose")</code> |

Format

This is a tab separated file, with a minimum of 4 columns:

required columns:

- `jobname`: Name of the step
- `sub_type`: Short for submission type, refers to, how should multiple commands of this step be submitted. Possible values are 'serial' or 'scatter'.
- `prev_jobs`: Short for previous job, this would be the jobname of the previous job. This can be NA/./none if this is a independent/initial step, and no previous step is required for this to start. Additionally, one may use comma(s) to define multiple previous jobs (A,B).
- `dep_type`: Short for dependency type, refers to the relationship of this job with the one defined in 'prev_jobs'. This can take values 'none', 'gather', 'serial' or 'burst'.

resource columns (recommended):

Additionally, one may customize resource requirements used by each step. The format used varies and depends to the computing platform. Thus its best to refer to your institutions guide to specify these.

- `cpu_reserved` integer, specifying number of cores to reserve [1]
- `memory_reserved` Usually in KB [2000]
- `nodes` number of server nodes to reserve, most tools can only use multiple cores on a **single** node [1]
- `walltime` maximum time allowed for a step, usually in a HH:MM or HH:MM:SS format. [1:00]
- `queue` the queue to use for job submission [short]

Details

NOTE: Guessing is an experimental feature, please check the definition carefully. it is provided to help but not replace your best judgement.

Optionally, one may provide the previous jobs and flowr can try guessing the appropriate submission and dependency types. If there are multiple commands, default is submitting them as scatter, else as serial. Further, if previous job has multiple commands and current job has single; its assumed that all of the previous need to complete, suggesting a gather type dependency.

Examples

```
# see ?to_flow for more examples

# read in a tsv; check and confirm format
ex = file.path(system.file(package = "flowr"), "pipelines")

# read in a flowdef from file
flowdef = as.flowdef(file.path(ex, "sleep_pipe.def"))

# check if this a flowdef
is.flowdef(flowdef)

# use a flowmat, to create a sample flowdef
flowmat = as.flowmat(file.path(ex, "sleep_pipe.tsv"))
to_flowdef(flowmat)
```



```

# change the platform
to_flowdef(flowmat, platform = "lsf")

# change the queue name
def = to_flowdef(flowmat,
  platform = "lsf",
  queue = "long")
plot_flow(def)

# guess submission and dependency types
def2 = to_flowdef(flowmat,
  platform = "lsf",
  queue = "long",
  guess = TRUE)
plot_flow(def2)

```

to_flowdet

*Create a flow's submission detail file***Description**

Create a file describing details regarding jobs ids, submission scripts etc.

Usage

```

to_flowdet(x, ...)

## S3 method for class 'rootdir'
to_flowdet(x, ...)

## S3 method for class 'character'
to_flowdet(x, ...)

## S3 method for class 'flow'
to_flowdet(x, ...)

```

Arguments

| | |
|-----|--------------|
| x | this is a wd |
| ... | not used |

Details

The path provided should contain a flow_details.rds file (which is used to extract all the information).

In case a parent folder with multiple flows is provided information regarding jobids is omitted.

if x is char. assumed a path, check if flow object exists in it and read it. If there is no flow object, try using a simpler function

to_flowmat

*Create a flowmat using a list a commands.***Description**

Create a flowmat (data.frame) using a **named** list a commands.

as.flowmat(): reads a file and checks for required columns. If x is data.frame checks for required columns.

Usage

```
to_flowmat(x, ...)

## S3 method for class 'list'
to_flowmat(x, samplename, ...)

## S3 method for class 'data.frame'
to_flowmat(x, ...)

## S3 method for class 'flow'
to_flowmat(x, ...)

as.flowmat(x, grp_col, jobname_col, cmd_col, ...)

is.flowmat(x)
```

Arguments

| | |
|-------------|---|
| x | a named list, where name corresponds to the jobname and value is a vector of commands to run. |
| ... | not used |
| samplename | character of length 1 or that of nrow(x) ['samplename'] |
| grp_col | column used for grouping, default samplename. |
| jobname_col | column specifying jobname, default jobname |
| cmd_col | column specifying commands to run, default cmd |

Examples

```
# Use this link for a few examples:
# http://docs.flowr.space/tutorial.html#define_modules

# create a flow mat, starting with a list of commands.
cmd_sleep = c("sleep 1", "sleep 2")
cmd_echo = c("echo 'hello'", "echo 'hello'")

# create a named list
lst = list("sleep" = cmd_sleep, "echo" = cmd_echo)
flowmat = to_flowmat(lst, samplename = "samp")
```

```

# read in a tsv; check and confirm format
ex = file.path(system.file(package = "flowr"), "pipelines")

flowmat = as.flowmat(file.path(ex, "sleep_pipe.tsv"))

# if your column names are different than defaults, explicitly specify them.
flowmat = as.flowmat(file.path(ex, "sleep_pipe.tsv"), jobname_col = "jobname")

# check if a object is a flowmat
is.flowmat(flowmat)

# create a flowdef, from this flowmat
flowdef = to_flowdef(flowmat)

# create a flow object using flowmat and flowdef
fobj = to_flow(flowmat, flowdef)

# extract a flowmat from a flow (here the samplename also contains the name of the flow)
flowmat2 = to_flowmat(fobj)

## submit the flow to the cluster (execute=TRUE) or do a dry-run (execute=FALSE)
## Not run:
fobj2 = submit_flow(fobj, execute=FALSE)
fobj3 = submit_flow(fobj, execute=TRUE)

## Get the status or kill all the jobs
status(fobj3)
kill(fobj3)

## End(Not run)

```

verbose

Verbose levels, defining verbosity of messages

Description

There are several levels of verbosity one can choose from.

levels:

- level 0 is almost silent, producing only necessary messages
- level 1 is good for most purposes, where as,
- level 2 is good when developing a new pipeline.
- level 3 is good for debugging, especially when getting un-expected results.

One can set the level of verbosity using `opts_flow$set(verbose=2)`, which will be used across `flowr` and `ngsflows` packages. Additionally one may set this value in the configurations files: `~/flowr/conf/flowr.conf` OR `~/flowr/conf/ngsflows.conf`.

Usage

`verbose`

Format

An object of class `NULL` of length 0.

Examples

```
fl = system.file("pipelines/abcd.def", package = "flowr")
def = as.flowdef(fl, verbose = 0)
# def seems to be a file, reading it...
def = as.flowdef(fl, verbose = 1)
# def seems to be a file, reading it...
# checking if required columns are present...
# checking if resources columns are present...
# checking if dependency column has valid names...
# checking if submission column has valid names...
# checking for missing rows in def...
# checking for extra rows in def...
# checking submission and dependency types...
def = as.flowdef(fl, verbose = 2)
# def seems to be a file, reading it...
# checking if required columns are present...
# checking if resources columns are present...
# checking if dependency column has valid names...
# checking if submission column has valid names...
# checking for missing rows in def...
# checking for extra rows in def...
# checking submission and dependency types...
# jobname prev.sub_type --> dep_type --> sub_type: relationship
# 1: A none --> none --> scatter
# 2: B scatter --> serial --> scatter rel: complex one:one
# 3: C scatter --> gather --> serial rel: many:one
# 4: D serial --> burst --> scatter rel: one:many
```

whisker_render

Wrapper around whisker.render with some additional checks

Description

Internal function (used by `submit_job`), which creates a submission script using platform specific templates.

This is a wrapper around [whisker.render](#)

Usage

`whisker_render(template, data)`

Arguments

| | |
|----------|---|
| template | template used |
| data | a list with variables to be used to fill in the template. |

| | |
|--------------------|--|
| write_flow_details | <i>write files desribing this flow</i> |
|--------------------|--|

Description

write files desribing this flow

Usage

```
write_flow_details(x, fobj, summ, flow_det, plot = FALSE)
```

Arguments

| | |
|----------|---------------------------|
| x | path to write to |
| fobj | flow object |
| summ | a status summary. |
| flow_det | a flow details data.frame |
| plot | logical, plot or not |

Index

*Topic **datasets**

flowopts, 6
verbose, 27

*Topic **queue**

queue-class, 12

as.flowdef (to_flowdef), 22
as.flowmat (to_flowmat), 26

check, 2
check_args, 3

definition (to_flowdef), 22

fetch, 4, 7
fetch_conf, 4, 16
fetch_conf (fetch), 4
fetch_pipes, 4, 14–16
fetch_pipes (fetch), 4
flow, 10, 14, 18, 20
flow (flow-class), 5
flow-class, 5
flowdef (to_flowdef), 22
flowopts, 5, 6, 21
flowr (to_flow), 20

get_opts, 6
get_opts (flowopts), 6
get_status (status), 16
get_wds, 8

is.flow (to_flow), 20
is.flowdef (to_flowdef), 22
is.flowmat (to_flowmat), 26

job, 8, 18

kill, 9

load_opts, 6
load_opts (flowopts), 6

opts_flow (flowopts), 6
opts_flow\$get (flowopts), 6
opts_flow\$load, 5

opts_flow\$load (flowopts), 6
opts_flow\$set (flowopts), 6

params, 7
plot.flow (plot_flow), 10
plot.flowdef (plot_flow), 10
plot_flow, 10

queue, 13, 21
queue (queue-class), 12
queue-class, 12

read_sheet, 7
rerun, 13
run, 14
run_flow (run), 14
run_pipe (run), 14

set_opts, 6
set_opts (flowopts), 6
setup, 16
status, 13, 16
submit_flow, 6, 9, 14, 17, 20, 21
submit_job, 18, 18
submit_run, 19

test_queue, 19
to_flow, 5, 12, 14, 20, 20
to_flow.data.frame, 20
to_flowdef, 13, 14, 20, 21, 22
to_flowdet, 21, 25
to_flowmat, 13, 14, 20, 21, 26

verbose, 2, 4, 10, 14, 17, 21, 23, 27

whisker.render, 28
whisker_render, 28
write_flow_details, 29