

---

# **NetSA Python Documentation**

***Release 1.4***

**Carnegie Mellon University**

September 30, 2011



# CONTENTS

<b>1</b>	<b>net<code>sa</code>.script — The NetSA Scripting Framework</b>	<b>1</b>
1.1	Overview . . . . .	1
1.2	Exceptions . . . . .	3
1.3	Metadata Functions . . . . .	4
1.4	Script Parameters . . . . .	4
1.5	Verbose Output . . . . .	7
1.6	Flow Data Parameters . . . . .	7
1.7	Producing Output . . . . .	10
1.8	Script Execution . . . . .	12
<b>2</b>	<b>net<code>sa</code>.script.golem — Golem Script Automation</b>	<b>13</b>
2.1	Overview . . . . .	13
2.2	Command Line Usage . . . . .	13
2.3	Metadata Functions . . . . .	15
2.4	Configuration Functions . . . . .	16
2.5	Parameter Functions . . . . .	21
2.6	Processing and Status Functions . . . . .	22
2.7	Utility Functions . . . . .	22
2.8	Additional Functions . . . . .	23
2.9	Usage of Tags, Loops, and Templates . . . . .	23
2.10	Intervals and Spans Explained . . . . .	26
2.11	Examples . . . . .	27
2.12	Classes . . . . .	41
<b>3</b>	<b>net<code>sa</code>.sql — SQL Database Access</b>	<b>47</b>
3.1	Overview . . . . .	47
3.2	Exceptions . . . . .	47
3.3	Connecting . . . . .	48
3.4	Connections and Result Sets . . . . .	48
3.5	Compiled Queries . . . . .	49
3.6	Implementing a New Driver . . . . .	50
3.7	Experimental Connection Pooling . . . . .	51
3.8	Why Not DB API 2.0? . . . . .	51
<b>4</b>	<b>net<code>sa</code>.util.shell — Robust Shell Pipelines</b>	<b>53</b>
4.1	Overview . . . . .	53
4.2	Exceptions . . . . .	55
4.3	Building Commands and Pipelines . . . . .	55
4.4	Running Pipelines . . . . .	58

<b>5</b>	<b>net<code>sa</code>_sil<code>k</code> — NetSA Python/PySiLK Support</b>	<b>61</b>
5.1	IPv6 Support . . . . .	61
5.2	IP Addresses . . . . .	61
5.3	IP Sets . . . . .	63
5.4	IP Wildcards . . . . .	65
5.5	TCP Flags . . . . .	66
5.6	Support for SiLK versions before 3.0 . . . . .	67
<b>6</b>	<b>Data Manipulation</b>	<b>69</b>
6.1	<code>net<code>sa</code>.data.countries</code> — Country and Region Codes . . . . .	69
6.2	<code>net<code>sa</code>.data.format</code> — Formatting Data for Output . . . . .	70
6.3	<code>net<code>sa</code>.data.nice</code> — “Nice” Numbers for Chart Bounds . . . . .	74
6.4	<code>net<code>sa</code>.data.times</code> — Time and Date Manipulation . . . . .	75
<b>7</b>	<b>Miscellaneous Facilities</b>	<b>79</b>
7.1	<code>net<code>sa</code>.files</code> — File and Path Manipulation . . . . .	79
7.2	<code>net<code>sa</code>.json</code> — JSON Wrapper Module . . . . .	80
7.3	<code>net<code>sa</code>.util.clitest</code> — Utility for testing CLI tools . . . . .	80
7.4	<code>net<code>sa</code>.util.compat</code> — Python version compatibility code . . . . .	82
<b>8</b>	<b>Internal Facilities</b>	<b>83</b>
8.1	<code>net<code>sa</code>.dist</code> — Common Installation Procedures . . . . .	83
<b>9</b>	<b>Deprecated Features</b>	<b>89</b>
9.1	Deprecated functions from <code>net<code>sa</code>.files</code> . . . . .	89
9.2	Deprecated module <code>net<code>sa</code>.files.datefiles</code> . . . . .	90
9.3	Deprecated functions from <code>net<code>sa</code>.script</code> . . . . .	92
9.4	Deprecated module <code>net<code>sa</code>.tools.service</code> . . . . .	93
<b>10</b>	<b>Changes</b>	<b>95</b>
10.1	Version 1.4 - 2011-09-30 . . . . .	95
10.2	Version 1.3 - 2011-03-28 . . . . .	95
10.3	Version 1.2 - 2011-01-12 . . . . .	96
10.4	Version 1.1 - 2010-10-04 . . . . .	96
10.5	Version 1.0 - 2010-09-14 . . . . .	96
10.6	Version 0.9 - 2010-01-19 . . . . .	96
<b>11</b>	<b>Licenses</b>	<b>97</b>
11.1	License for <code>net<code>sa</code>-python</code> . . . . .	97
11.2	License for <code>simplejson</code> . . . . .	98
	<b>Index</b>	<b>99</b>

---

# NETSA.SCRIP — THE NETSA SCRIPTING FRAMEWORK

## 1.1 Overview

The `netsa.script` module provides a common framework for building SiLK-based analysis scripts. This framework is intended to make scripts re-usable and automatable without much extra work on the part of script authors. The primary concerns of the scripting framework are providing metadata for cataloging available scripts, standardizing handling of command-line arguments (particularly for flow data input), and locating output files.

Here's an example of a simple Python script using the `netsa.script` framework.

First is a version without extensive comments, for reading clarity. Then the script is repeated with comments explaining each section.

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

# Import the script framework under the name "script".
from netsa import script

# Set up the metadata for the script, including the title, what it
# does, who wrote it, who to ask questions about it, etc.
script.set_title("Sample Framework Script")
script.set_description("""
    An example script to demonstrate the basic features of the
    netsa.script scripting framework.  This script counts the
    number of frobnitzim observed in each hour (up to a maximum
    volume of frobs per hour.)
""")
script.set_version("0.1")
script.set_contact("H. Bovik <hbovik@example.org>")
script.set_authors(["H. Bovik <hbovik@example.org>"])

script.add_int_param("frob-limit",
    "Maximum volume of frobs per hour to observe.",
    default=10)

script.add_float_param("frobnitz-sensitivity",
    "Sensitivity (between 0.0 and 1.0) of frobnitz categorizer.",
    default=0.61, expert=True, minimum=0.0, maximum=1.0)

script.add_flow_params(require_pull=True)
```

```
script.add_output_file_param("output-path",
    "Number of frobnitzim observed in each hour of the flow data.",
    mime_type="text/csv")

# See the text for discussion of the next two functions.

def process_hourly_data(out_file, flow_params, frob_limit, frob_sense):
    ...

def main():
    frob_limit = script.get_param("frob-limit")
    frobnitz_sensitivity = script.get_param("frobnitz-sensitivity")
    out_file = script.get_output_file("output-path")
    for hour_params in script.get_flow_params().by_hour():
        process_hourly_data(out_file, hour_params, frob_limit,
                            frobnitz_sensitivity)

script.execute(main)
```

Let's break things down by section:

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

from netsa import script
```

This is basic Python boilerplate. Any other libraries we use would also be imported at this time.

```
script.set_title("Sample Framework Script")
script.set_description("""
    An example script to demonstrate the basic features of the
    netsa.script scripting framework. This script counts the
    number of frobnitzim observed in each hour (up to a maximum
    volume of frobs per hour.)
""")
script.set_version("0.1")
script.set_contact("H. Bovik <hbovik@example.org>")
script.set_authors(["H. Bovik <hbovik@example.org>"])
```

Script metadata allows users to more easily find out information about a script, and browse available scripts stored in a central repository. The above calls define all of the metadata that the *netsa.script* framework currently supports. It is possible that a future version will include additional metadata fields.

```
script.add_int_param("frob-limit",
    "Maximum volume of frobs per hour to observe.",
    default=10)

script.add_float_param("frobnitz-sensitivity",
    "Sensitivity (between 0.0 and 1.0) of frobnitz categorizer.",
    default=0.61, expert=True, minimum=0.0, maximum=1.0)
```

Script parameters are defined by calling `netsa.script.add_X_param` (where *X* is a type) for each parameter. Depending on the type of the parameter, there may be additional configuration options (like *minimum* and *maximum* for the float parameter above) available. See the documentation for each function later in this document.

Expert parameters are tuning parameters that are intended for expert use only. An expert parameter is created by setting *expert* to `True` when creating a new parameter. This parameter will then be displayed only if the user asks for `--help-expert`, and the normal help will indicate that expert options are available.

```
script.add_flow_params(require_pull=True)
```

Parameters involving flow data are handled separately, in order to ensure that flows are handled consistently across all of our scripts. The `netса.script.add_flow_params` function is used to add all of the flow related command-line arguments at once. There is currently only one option. If the `require_pull` option is set, the flow data must come from an `rwfilter` data pull (including switches like `--start-date`, `--end-date`, `--class`, etc.) If `require_pull` is not set, then it is also possible for input files or pipes to be given on the command-line.

```
script.add_output_file_param("output-path",
    "Number of frobnitzim observed in each hour of the flow data.",
    mime_type="text/csv")
```

Every output file (not temporary working file) that the script produces must also be defined using calls to the framework—this ensures that when an automated tool is used to run the script, it can find all of the relevant output files. It's preferable, but not required, for a MIME content-type (like `"text/csv"`) and a short description of the contents of the file be included.

```
def process_hourly_data(out_file, flow_params, frob_limit, frob_sense):
    ...
```

In this example, the `process_hourly_data` function would be expected to use the functions in `netса.util.shell` to acquire and process flow data for each hour (based on the `flow_params` argument). The details have been elided for simplicity in this example.

```
def main():
    frob_limit = script.get_param("frob-limit")
    frobnitz_sensitivity = script.get_param("frobnitz-sensitivity")
    out_file = script.get_output_file("output-path")
    for hour_data in script.get_flow_params().by_hour():
        process_hourly_data(out_file, hour_params, frob_limit,
                           frobnitz_sensitivity)
```

It is important that no work is done outside the `main` function (which can be given any name you wish). If instead you do work in the body of the file outside of a function, that work will be done whether or not the script has actually been asked to do work. (For example, if the script is given `--help`, it will not normally call your `main` function.) So make sure everything is in here.

```
script.execute(main)
```

The final statement in the script should be a call to `netса.script.execute`, as shown above. This allows the framework to process any command-line arguments (including producing help output, etc.), then call your `main` function, and finally do clean-up work after the completion of your script.

See the documentation for functions in this module for more details on individual features, including further examples.

## 1.2 Exceptions

**exception** `netса.script.ParamError` (*param, value, message*)

This exception represents an error in the arguments provided to a script at the command-line. For example, `ParamError('foo', '2x5', 'not a valid integer')` is the exception generated when the value given for an integer param is not parsable, and will produce the following error output when thrown from a script's `main` function:

```
<script-name>: Invalid foo '2x5': not a valid integer
```

**exception** `netsa.script.UserError` (*message*)

This exception represents an error reported by the script that should be presented in a standard way. For example, `UserError('your message here')` will produce the following error output when thrown from a script's main function:

```
<script-name>: your message here
```

**exception** `netsa.script.ScriptError` (*message*)

This exception represents an error in script definition or an error in processing script data. This is thrown by some `netsa.script` calls.

## 1.3 Metadata Functions

The following functions define “metadata” for the script—they provide information about the name of the script, what the script is for, who to contact with problems, and so on. Automated tools can use this information to allow users to browse a list of available scripts.

`netsa.script.set_title` (*script\_title* : *str*)

Set the title for this script. This should be the human-readable name of the script, and denote its purpose.

`netsa.script.set_description` (*script\_description* : *str*)

Set the description for this script. This should be a longer human-readable description of the script's purpose, including simple details of its behavior and required inputs.

`netsa.script.set_version` (*script\_version* : *str*)

Set the version number of this script. This can take any form, but the standard *major . minor ( . patch )* format is recommended.

`netsa.script.set_package_name` (*script\_package\_name* : *str*)

Set the package name for this script. This should be the human-readable name of a collection of scripts.

`netsa.script.set_contact` (*script\_contact* : *str*)

Set the point of contact email for support of this script, which must be a single string. The form should be suitable for treatment as an email address. The recommended form is a string containing:

```
Full Name <full.name@contact.email.org>
```

`netsa.script.set_authors` (*script\_authors* : *str list*)

Set the list of authors for this script, which must be a list of strings. It is recommended that each author be listed in the form described for `set_contact`.

`netsa.script.add_author` (*script\_author* : *str*)

Add another author to the list of authors for this script, which must be a single string. See `set_authors` for notes on the content of this string.

## 1.4 Script Parameters

These calls are used to add parameters to a script. When the script is called from the command-line, these are command-line arguments. When a GUI is used to invoke the script, the params might be presented in a variety of ways. This need to support both command-line and GUI access to script parameters is the reason that they've been standardized here. It's also the reason that you'll find no “add an argument with this arbitrary handler function” here.

If you do absolutely need deeper capabilities than are provided here, you can use one of the basic param types and then do additional checking in the `main` function. Note, however, that a GUI will not aid users in choosing acceptable



values for params defined in this way. Also, make sure to raise `ParamError` with appropriate information when you reject a value, so that the error can be most effectively communicated back to the user.

```
netsa.script.add_text_param(name : str, help : str[, required=False, default : str, default_help :
                             str, expert=False, regex : str])
```

Add a text parameter to this script. This parameter can later be fetched as a `str` by the script using `netsa.script.get_param`. The `required`, `default`, `default_help`, and `expert` arguments are used by all `add_X_param` calls, but each kind of parameter also has additional features that may be used. See below for a list of these features for text params.

**Example:** Add a new parameter which is required for the script to run.

```
add_text_param("graph-title",
               "Display this title on the output graph.",
               required=True)
```

It is an error if this parameter is not set, and the script will exit with a usage message when it is run at the command-line.

**Example:** Add a new parameter with a default value of "" (the empty string):

```
add_text_param("graph-comment",
               "Display this comment on the output graph.",
               default="")
```

If the parameter is not provided, the default value will be used.

**Example:** Display something different in the help text than the actual default value:

```
add_text_param("graph-date",
               "Display data for the given date.",
               default=date_for_today(), default_help="today")
```

Sometimes a default value should be computed but not displayed as the default to the user when they ask for help at the command-line. In this case, a default value should be provided (which will be displayed to users of a GUI), while a value for `default_help` will be presented in the `-help` output. In addition, GUIs will also display the value of `default_help` in some way next to the entry field for this parameter.

It is perfectly legal to provide a value for `default_help` and not provide a value for `default`. This makes sense when the only way to compute the default value for the field is at actual execution time. (For example, if the end-date defaults to be the same as the provided start-date.)

**Example:** Add a new “expert” parameter:

```
add_text_param("gnuplot-extra-commands",
               "Give these extra command to gnuplot when writing output.",
               expert=True)
```

Expert parameters are not listed for users unless they explicitly ask for them. (For example, by using `--help-expert` at the command line.)

Other keyword arguments meaningful for text params:

***regex*** Require strings to match this regular expression.

**Example:** Add a new text parameter that is required to match a specific pattern for phone numbers:

```
add_text_param("phone-number",
               "Send reports to this telephone number.",
               regex=r"[0-9]{3}-[0-9]{3}-[0-9]{4}")
```

```
netlsa.script.add_int_param(name : str, help : str[, required=False, default : int, default_help : str,
                             expert=False, minimum : int, maximum : int])
```

Add an integer parameter to this script. This parameter can later be fetched as an `int` by the script using `netlsa.script.get_param`. The *required*, *default*, *default\_help*, and *expert* arguments are described in the help for `netlsa.script.add_text_param`.

Other keyword arguments meaningful for integer parameters:

***minimum*** Only values greater than or equal to this value are allowed for this param.

***maximum*** Only values less than or equal to this value are allowed for this param.

**Example:** Add a new int parameter which is required to be in the range  $0 \leq x \leq 65535$ .

```
add_int_param("targeted-port",
              "Search for attacks targeting this port number.",
              required=True, minimum=0, maximum=65535)
```

```
netlsa.script.add_float_param(name : str, help : str[, required=False, default : float, default_help
                                     : str, expert=False, minimum : float, maximum : float])
```

Add a floating-point parameter to this script. This parameter can later be fetched as a `float` by the script using `netlsa.script.get_param`. The *required*, *default*, *default\_help* and *expert* arguments are described in the help for `netlsa.script.add_text_param`.

Other keyword arguments meaningful for floating-point parameters:

***minimum*** Only values greater than or equal to this value are allowed for this param.

***maximum*** Only values less than or equal to this value are allowed for this param.

```
netlsa.script.add_date_param(name : str, help : str[, required=False, default : datetime, de-
                                     fault_help : str, expert=False])
```

Add a date parameter to this script. This parameter can later be fetched by the script as a `datetime.datetime` object using `netlsa.script.get_param`. The *required*, *default*, *default\_help*, and *expert* arguments are described in the help for `netlsa.script.add_text_param`.

```
netlsa.script.add_label_param(name : str, help : str[, required=False, default : str, default_help :
                                     str, expert=False, regex : str])
```

Add a label parameter to this script. This parameter can later be fetched by the script as a Python `str` using `netlsa.script.get_param`. The *required*, *default*, *default\_help*, and *expert* arguments are described in the help for `netlsa.script.add_text_param`.

Other keyword arguments meaningful for label params:

***regex*** Require strings to match this regular expression, instead of the default `r"[\^\\S, ]+"` (no white space or commas).

**Example:** Add a new label parameter that is required to match a specific pattern for phone numbers:

```
add_label_param("output-label",
                "Store output to the destination with this label.",
                regex=r"[0-9]{3}-[0-9]{3}-[0-9]{4}")
```

```
netlsa.script.add_file_param(name : str, help : str[, required=False, default : str, default_help :
                                     str, expert=False, mime_type : str])
```

Add a file parameter to this script. This parameter can later be fetched by the script as a Python `str` filename using `netlsa.script.get_param`. The *required*, *default*, *default\_help*, and *expert* arguments are described in the help for `netlsa.script.add_text_param`.

When the script is run at the command-line, an error will be reported to the user if they specify a file that does not exist, or the path of a directory.

Other keyword arguments meaningful for file params:

***mime\_type*** The expected MIME Content-Type of the file, if any.

```
netlsa.script.add_dir_param(name : str, help : str[, required=False, default : str, default_help : str,
                                expert=False])
```

Add a directory parameter to this script. This parameter can later be fetched by the script as a Python `str` filename using `netlsa.script.get_param`. The *required*, *default*, *default\_help*, and *expert* arguments are described in the help for `netlsa.script.add_text_param`.

When the script is run at the command-line, an error will be reported to the user if they specify a directory that does not exist, or the path of a file.

```
netlsa.script.add_path_param(name : str, help : str[, required=False, default : str, default_help :
                                str, expert=False])
```

Add a path parameter to this script. This parameter can later be fetched by the script as a Python `str` using `netlsa.script.get_param`. The *required*, *default*, *default\_help*, and *expert* arguments are described in the help for `netlsa.script.add_text_param`.

```
netlsa.script.add_path_param(name : str, help : str[, required=False, default : str, default_help :
                                str, expert=False])
```

Add a path parameter to this script. This parameter can later be fetched by the script as a Python `str` using `netlsa.script.get_param`. The *required*, *default*, *default\_help*, and *expert* arguments are described in the help for `netlsa.script.add_text_param`.

```
netlsa.script.add_flag_param(name : str, help : str[, default=False, default_help : str, ex-
                                pert=False])
```

Add a flag parameter to this script. This parameter can later be fetched by the script as a `bool` using `netlsa.script.get_param`. The *default*, *default\_help*, and *expert* arguments are described in the help for `netlsa.script.add_text_param`.

```
netlsa.script.get_param(name : str) → value
```

Returns the value of the parameter given by the `str` argument *name*. This parameter will be in the type specified for the param when it was added (for example, date parameters will return a `datetime.datetime` object.) Note that a parameter with no default that is not required may return `None`.

## 1.5 Verbose Output

```
netlsa.script.get_verbosity() → int
```

Returns the current verbosity level (default 0) for the script invocation. The `message` function may be used to automatically emit messages based on the verbosity level set for the script. Verbosity is set from the command-line via the `--verbose` or `-v` flags.

```
netlsa.script.display_message(text[, min_verbosity=1])
```

Writes the string *text* to `stderr`, as long as the script's verbosity is greater than or equal to *min\_verbosity*. Verbosity is set from the command-line via the `--verbose` or `-v` flags. The current verbosity level may be retrieved by using the `get_verbosity` function.

Use this function to write debugging or informational messages from your script for command-line use. For example, writing out which file you are processing, or what stage of processing is in progress.

Do not use it to write out important information such as error messages or actual output. (See `UserError` or `add_output_file_param` and `add_output_dir_param` for error messages and output.)

## 1.6 Flow Data Parameters

In order to standardize the large number of scripts that work with network flow data using the SiLK tool suite, the following calls can be used to work with flow data input.

```
netlsa.script.add_flow_annotation (script_annotation : str)
```

Add a note that will automatically be included in SiLK data pulls generated by this script. This will be included only by `rwfilter` pulls created by this script using `Flow_params`.

```
netlsa.script.add_flow_params ([require_pull=False, without_params : str list])
```

Add standard flow parameters to this script. The following params are added by default, but individual params may be disabled by including their names in the `without_params` argument. You might wish to disable the `--type` param, for example, if your script will run the same pull multiple times, once with `--type=in`, `inweb`, then again with `--type=out`, `outweb`. (Of course, you might then also want to add `in-type` and `out-type` params to the script.)

**--class** Req Arg. Class of data to process

**--type** Req Arg. Type(s) of data to process within the specified class. The type names and default type(s) vary by class. Use `all` to process every type for the specified class. Use `rwfilter -help` for details on valid class/type pairs.

**--flowtypes** Req Arg. Comma separated list of class/type pairs to process. May use `all` for class and/or type. This is alternate way to specify class/type; switch cannot be used with `--class` and `--type`

**--sensors** Req Arg. Comma separated list of sensor names, sensor IDs, and ranges of sensor IDs. Valid sensors vary by class. Use `mapsid` to see a mapping of sensor names to IDs and classes.

**--start-date** Req Arg. First hour of data to process. Specify date in YYYY/MM/DD[:HH] format: time is in UTC. When no hour is specified, the entire date is processed. Def. Start of today

**--end-date** Req Arg. Final hour of data to process specified as YYYY/MM/DD[:HH]. When no hour specified, end of day is used unless `start-date` includes an hour. When switch not specified, defaults to value in `start-date`.

If the `require_pull` argument to `netlsa.script.add_flow_params` is not `True`, input filenames may be specified bare on the command-line, and the following additional options are recognized:

**--input-pipe** Req Arg. Read SiLK flow records from a pipe: `stdin` or path to named pipe.  
No default

**--xargs (expert)** Req Arg. Read list of input file names from a file or pipe pathname or `stdin`.  
No default

The values of these parameters can later be retrieved as a `netlsa.script.Flow_params` object using `netlsa.script.get_flow_params`.

```
netlsa.script.get_flow_params () → Flow_params
```

Returns a `Flow_params` object encapsulating the `rwfilter` flow selection parameters the script was invoked with. This object is filled in based on the command-line arguments described in `add_flow_params`.

```
class netlsa.script.Flow_params ([flow_class : str, flow_type : str, flowtypes : str list, sensors : str list,
                                start_date : datetime, end_date : datetime, input_pipe : str, xargs :
                                str, filenames : str list])
```

This object represents the flow selection arguments to an `rwfilter` data pull. In typical use it is built automatically from command-line arguments by the `netlsa.script.get_flow_params` call. Afterwards, methods such as `by_hour` are used to modify the scope of the data pull, and then the parameters are included in a call to `rwfilter` using the functions in `netlsa.util.shell`.

**Example:** Process SMTP data from the user's requested flow data:

```
netlsa.util.shell.run_parallel(
    ["rwfilter %(flow_params)s --protocol=6 --aport=25 --pass=stdout",
     "rwuniq --fields=sip",
```

```
>>>output_file.txt"],
vars={'flow_params': script.get_flow_params()})
```

**Example:** Separately process each hour's SMTP data from the user's request flow data:

```
flow_params = script.get_flow_params()
# Iterate over each hour individually
for hourly_params in flow_params.by_hour():
    # Format ISO-style datetime for use in a filename
    sdate = iso_datetime(hourly_params.get_start_date())
    netsa.util.shell.run_parallel(
        ["rwfilter %(flow_params)s --protocol=6 --pass=stdout",
         "rwuniq --fields=dport",
         ">>>output_file_%(sdate)s.txt"],
        vars={'flow_params': hourly_params,
              'sdate': sdate})
```

**by\_day()** → Flow\_params iter

Given a `Flow_params` object including a start-date and an end-date, returns an iterator yielding a `Flow_params` for each individual day in the time span.

If the original `Flow_params` starts or ends on an hour that is not midnight, the first or last yielded pulls will not be for full days. All of the other pulls will be full days stretching from midnight to midnight.

See also `by_hour` which iterates over the time span of the `Flow_params` by hours instead of days.

Raises a `ScriptError` if the `Flow_params` has no date information (for example, the script user specified input files rather than a data pull.) This can be prevented by using *require\_pull* in your call to `script.add_flow_params`.

**by\_hour()** → Flow\_params iter

Given a `Flow_params` object including a start-date and an end-date, returns an iterator yielding new `Flow_params` object identical to this one specialized for each hour in the time period.

Example (strings are schematic of the `Flow_params` involved):

```
>>> # Note: Flow_params cannot actually take a str argument like this.
>>> some_flows = Flow_params('--type in,inweb --start-date 2009/01/01T00 '
>>>                          '--end-date 2009/01/01T02')
>>> list(some_flows.by_hour())
[netsa.script.Flow_params('--type in,inweb --start-date 2009/01/01T00 '
                          '--end-date 2009/01/01T00'),
 netsa.script.Flow_params('--type in,inweb --start-date 2009/01/01T01 '
                          '--end-date 2009/01/01T01'),
 netsa.script.Flow_params('--type in,inweb --start-date 2009/01/01T02 '
                          '--end-date 2009/01/01T02')]
```

See also `by_day` which iterates over the time span of the `Flow_params` by days instead of hours.

Raises a `ScriptError` if the `Flow_params` has no date information (for example, the script user specified input files rather than a data pull.) This can be prevented by using *require\_pull* in your call to `script.add_flow_params`.

**by\_sensor()** → Flow\_params iter

Given a `Flow_params` object including a data pull, returns an iterator yielding a `Flow_params` for each individual sensor defined in the system.

**get\_argument\_list()** → str list or None

Returns the bundle of flow selection parameters as a list of strings suitable for use as command-line arguments in an `rwfilter` call. This is automatically called by the `netsa.util.shell` routines when a `Flow_params` object is used as part of a command.

**get\_class()** → str or None  
Returns the `rwfilter pull --class` argument as a str.

**get\_end\_date()** → datetime or None  
Returns the `rwfilter pull --end-date` argument as a `datetime.datetime` object.

**get\_filenames()** → str list or None  
Returns any files given on the command-line for an `rwfilter pull` as a str.

**get\_flowtypes()** → str list or None  
Returns the `rwfilter pull --flowtypes` argument as a str.

**get\_input\_pipe()** → str or None  
Returns the `rwfilter pull --input-pipe` argument as a str.

**get\_sensors()** → str list or None  
Returns the `rwfilter pull --sensors` argument as a list of str.

**get\_start\_date()** → datetime or None  
Returns the `rwfilter pull --start-date` argument as a `datetime.datetime` object.

**get\_type()** → str or None  
Returns the `rwfilter pull --type` argument as a str.

**get\_xargs()** → str or None  
Returns the `rwfilter pull --xargs` argument as a str.

**is\_files()** → bool  
Returns True if this `Flow_params` object represents processing of already retrieved files.

**is\_pull()** → bool  
Returns True if this `Flow_params` object represents a data pull from the repository. (i.e. it contains selection switches.)

**using**(`[flow_class : str, flow_type : str, flowtypes : str list, sensors : str list, start_date : datetime, end_date : datetime, input_pipe : str, xargs : str, filenames : str list]`) → `Flow_params`  
Returns a new `Flow_params` object in which the arguments in this call have replaced the parameters in *self*, but all other parameters are the same.  
  
Raises a `ScriptError` if the new parameters are inconsistent or incorrectly typed.

## 1.7 Producing Output

Every output file that a script produces needs to be registered with the system, so that automated tools can be sure to collect everything. Some scripts produce one or more set outputs. For example “the report”, or “the HTML version of the report”. Others produce a number of outputs based on the content of the data they process. For example “one image for each host we identify as suspicious.”

`netsa.script.add_output_file_param`(`name : str, help: str[, required=True, expert=False, description : str, mime_type='application/octet-stream']`)

Add an output file parameter to this script. This parameter can later be fetched by the script as a Python str filename or a Python file object using `netsa.script.get_output_file_name` or `netsa.script.get_output_file`. Note that if you ask for the file name, you may wish to handle the filenames `stdout`, `stderr`, and `-` specially to be consistent with other tools. (See the documentation of `netsa.script.get_output_file_name` for details.) Unlike most parameters, output file parameters never have default values, and are required by default. If an output file parameter is not required, the implication is that if the user does not specify this argument, then this output is not produced.

In keeping with the behavior of the SiLK tools, it is an error for the user to specify an output file that already exists. If the environment variable `SILK_CLOBBER` is set, this restriction is relaxed and existing output files may be overwritten.

The `mime_type` argument is advisory., but it should be set to an appropriate MIME content type for the output file. The framework will not report erroneous types, nor will it automatically convert from one type to another. Examples:

**text/plain** Human readable text file.

**text/csv** Comma-separated-value file.

**application/x-silk-flows** SiLK flow data

**application/x-silk-ipset** SiLK ipset data

**application/x-silk-bag** SiLK bag data

**application/x-silk-pmap** SiLK prefix map data

**image/png** etc. Various standard formats, many of which are listed [on IANA's website](#).

It is by no means necessary to provide a useful MIME type, but it is helpful to automated systems that wish to interpret or display the output of your script.

The `description` argument may also be provided, with a long-form text description of the contents of this output file. Note that `description` describes the contents of the file, while `help` describes the meaning of the command-line argument.

`netlsa.script.get_output_file_name(name : str) → str`

Returns the filename for the output parameter `name`. Note that many SiLK tools treat the names `stdout`, `stderr`, and `-` as meaning something special. `stdout` and `-` imply the output should be written to standard out, and `stderr` implies the output should be written to standard error. It is not required that you handle these special names, but it helps with interoperability. Note that you may need to take care when passing these filenames to SiLK command-line tools for output or input locations, for the same reason.

If you use `netlsa.script.get_output_file`, it will automatically handle these special filenames.

If this output file is optional, and the user has not specified a location for it, this function will return `None`.

`netlsa.script.get_output_file(name : str) → file`

Returns an open file object for the output parameter `name`. The special names `stdout`, `-` are both translated to standard output, and `stderr` is translated to standard error.

If you need the output file name, use `netlsa.script.get_output_file_name` instead.

If `append` is `True`, then the file is opened for append. Otherwise it is opened for write.

`netlsa.script.add_output_dir_param(name : str, help : str[, required=True, expert=False, description : str, mime_type : str])`

Add an output directory parameter to this script. This parameter can later be used to construct a `str` filename or a Python file object using `netlsa.script.get_output_dir_file_name` or `netlsa.script.get_output_dir_file`. Unlike most parameters, output directory parameters never have default values, and are required by default. If an output directory parameter is not required, the implication is that if the user does not specify this argument, then this output is not produced.

See `add_output_file_param` for the meanings of the `description` and `mime_type` arguments. In this context, these arguments provide default values for files created in this output directory. Each individual file can be given its own `mime_type` and `description` when using the `netlsa.script.get_output_dir_file_name` and `netlsa.script.get_output_dir_file` functions.



```
netса.script.get_output_dir_file_name(dir_name : str, file_name : str[, description : str,  
                                     mime_type : str]) → str
```

Returns the path for the file named *file\_name* in the output directory specified by the parameter *dir\_name*. Also lets the `netса.script` system know that this output file is about to be used. If provided, the *description* and *mime\_type* arguments have meanings as described in `add_output_file_param`. If these arguments are not provided, the defaults from the call where *dir\_name* was defined in `add_output_dir_param` are used.

If the output directory parameter is optional, and the user has not specified a location for it, this function will return `None`.

```
netса.script.get_output_dir_file(dir_name : str, file_name : str[, description : str, mime_type  
                               : str]) → file
```

Returns the an open `file` object for the file named *file\_name* in the output directory specified by the parameter *dir\_name*. Also lets the `netса.script` system know that this output file is about to be used. If provided, the *description* and *mime\_type* arguments have meanings as described in `add_output_file_param`. If these arguments are not provided, the defaults from the call where *dir\_name* was defined in `add_output_dir_param` are used.

If the output dir param is optional, and the user has not specified a location for it, this function will return `None`.

If *append* is `True`, the file is opened for append. Otherwise, the file is opened for write.

## 1.8 Script Execution

```
netса.script.execute(func : callable)
```

Executes the main function of a script. This should be called as the last line of any script, with the script's main function (whatever it might be named) as its only argument.

It is important that all work in the script is done within this function. The script may be loaded in such a way that it is not executed, but only queried for metadata information. If the script does work outside of the `main` function, this will cause metadata queries to be very inefficient.



---

# NETSA . SCRIPT . GOLEM — GOLEM SCRIPT AUTOMATION

## 2.1 Overview

The *Golem Script Automation* framework is a specialized extension of the *NetSA Scripting Framework* for constructing automated analysis scripts. Such scripts might be launched periodically from a cron job in order to produce regular sets of results in a data repository. In addition to the golem-specific extensions, `netsa.script.golem` offers the same functionality as `netsa.script`.

At its heart, Golem is a template engine combined with some synchronization logic across analytic time bins. The templates define the output paths for result data and the command line sequences and pipelines used to generate these results. Template variables are known as *tags*. Golem provides certain tags by default and others can be added or modified via the golem *configuration functions*.

Golem Automation is designed to allow developers to build scripts that easily consume the output results of other golem scripts without the need for detailed knowledge of the implementation details of the external script (e.g. how often it runs or what pathnames it uses for populating its data repository). When provided the path to the external script in its configuration, a golem script will interrogate the external script for information regarding its output results, automatically synchronize across processing windows, and make the result paths available for use in command templates as input paths.

In addition to analysis automation, Golem scripts also offer a query mode so that results in the data repository can be easily examined or pulled to local directories.

Golem offers a number of shortcuts and convenience functions specific to the *SiLK Flow Analysis suite*, but is not limited to using SiLK for analysis.

See the *examples* at the end of this chapter to learn how to write a golem script. See the description of *template and tag usage* for details on tags provided for use within templates. See the *API reference* for more thorough documentation of the features and interface. See the *CLI reference* for the standard command line parameters that golem enables.

## 2.2 Command Line Usage

Golem-enabled scripts offer standard command line parameters grouped into three categories: *basic*, *repository-related*, and *query-related*. Parameters must be enabled by the script author. Parameters can be enabled individually or by category. The `add_golem_params` function will enable all parameters.

Golem scripts also have the standard `netsa.script` command line options `--help`, `--help-expert`, and `--verbose`.

### 2.2.1 Basic Parameters

The following `golem` command line parameters control the filtering of processing windows, as well as some input and output behavior. These options affect both repository and query operations.

**-last-date** <date>

Date specifying the last time bin of interest. The provided date will be rounded down to the first date of the processing interval in which it resides. Default: most recent

**-first-date** <date>

Date specifying the first time bin of interest. The provided date will be rounded down to the first date of the processing interval in which it resides. Default: value of `--last-date`

**-intervals** <count>

Process or query the last *count* intervals (as defined by `set_interval` within the script header) of data from the current date. This will override any values provided by `--first-date` or `--last-date`.

**-skip-incomplete**

Skip processing intervals that have incomplete input requirements (i.e. ignore the date if any source dependencies have incomplete results).

**-overwrite**

Overwrite output results if they already exist.

**-<loop-select>**

Each template loop that has been defined in the `golem` script (via the `add_loop` function) has an associated parameter that allows a comma-separated selection of a subset of that loop's values. For example, if a sensor loop was defined under the tag `sensor`, the parameter would be `--sensor` by default. If grouping was enabled it would be `--sensor-group` instead.

### 2.2.2 Repository Parameters

Repository parameters allow maintenance of the result repository, typically via a cron job. They are categorized as 'expert' options, therefore showing up in `--help-expert` rather than `--help`.

**-data-load**

Generate and store incomplete or missing analysis results in the data repository, skipping those that are complete (unless `--overwrite` is given).

**-data-status**

Show the status of repository processing bins, for the provided options. No processing is performed.

**-data-queue**

List the dates of all pending repository results for the provided options. No processing is performed.

**-data-complete**

List the dates of all completed repository results, for the provided options. No processing is performed.

**-data-inputs**

Show the status of input dependencies from other `golem` scripts or defined templates, for the provided options. Use `-v` or `-vv` for less abbreviated paths. No processing is performed.

**-data-outputs**

Show the status of repository output results for the provided options. Use `-v` or `-vv` for less abbreviated paths. No processing is performed.

### 2.2.3 Query Parameters

Query parameters control how local copies of results are stored, both when copying from the repository or when performing a fresh analysis.

**-output-path** <path>

Generate a single query result in the specified output file for the given parameters. Also accepts '-' and 'stdout'.

**-output-select** <arg1[, arg2[, ...]]>

For golem scripts that have more than one output defined, limit output results to the comma-separated names provided.

**-output-dir** <path>

Copy query result files in the specified output directory for the given parameters. The files will follow the same naming scheme as specified for the repository. These names can be previewed via the `--show-outputs` option.

**-show-inputs**

Show the status of input dependencies from other golem scripts or defined templates for the provided options. Use `-v` or `-vv` for less abbreviated paths. No processing is performed.

**-show-outputs**

Show relative output paths that would be generated within `--output-dir`, given the options provided. Use `-v` or `-vv` for less abbreviated paths. No processing is performed.

## 2.3 Metadata Functions

Golem scripts are an extension of `netsa.script`. As such, golem scripts offer the same functions as `netsa.script`. Script authors are encouraged to use the following *metadata functions*:

`netsa.script.golem.set_title(script_title : str)`

Set the title for this script. This should be the human-readable name of the script, and denote its purpose.

`netsa.script.golem.set_description(script_description : str)`

Set the description for this script. This should be a longer human-readable description of the script's purpose, including simple details of its behavior and required inputs.

`netsa.script.golem.set_version(script_version : str)`

Set the version number of this script. This can take any form, but the standard *major . minor ( . patch )* format is recommended.

`netsa.script.golem.set_package_name(script_package_name : str)`

Set the package name for this script. This should be the human-readable name of a collection of scripts.

`netsa.script.golem.set_contact(script_contact : str)`

Set the point of contact email for support of this script, which must be a single string. The form should be suitable for treatment as an email address. The recommended form is a string containing:

Full Name <full.name@contact.email.org>

`netsa.script.golem.set_authors(script_authors : str list)`

Set the list of authors for this script, which must be a list of strings. It is recommended that each author be listed in the form described for `set_contact`.

`netsa.script.golem.add_author(script_author : str)`

Add another author to the list of authors for this script, which must be a single string. See `set_authors` for notes on the content of this string.

Please see `netsa.script` for additional functions, such as those used to *add custom command line parameters*.

## 2.4 Configuration Functions

Golem scripts are configured by calling functions from within a particular imported module. You can either import `netsa.script.golem` directly:

```
from netsa.script import golem
```

Or import as `netsa.script.golem.script`:

```
from netsa.script.golem import script
```

In either case, the imported module will provide identical functionality. All functions available within `netsa.script` are also available from within `golem`, with the addition of the golem-specific functions and classes. Please consult the `netsa.script` documentation for details on functions offered by that module.

The following functions configure the behavior of golem scripts.

`netsa.script.golem.set_default_home(path : str[, path : str, ...])`

Sets the default base path for this golem script in cases where the `GOLEM_HOME` environment variable is not set. Multiple arguments will be joined together as with `os.path.join`. If the provided path is relative, it is assumed to be relative to the directory in which script resides.

The actual home path will be decided by the first available source in the following order:

- 1.The `GOLEM_HOME` environment variable
- 2.The default home if set by this function
- 3.The directory in which the script resides

Subsequent path settings (e.g. `set_repository`) will be relative to the script home if they are not absolute paths.

`netsa.script.golem.set_repository(path : str[, path : str, ...])`

Sets the default path for the output results data repository. Multiple arguments will be joined together as with `os.path.join`. Output results will be stored in this directory or in a subdirectory beneath, depending on how each output template is specified. Relative paths are considered to be relative to the home path.

`netsa.script.golem.set_suite_name(name : str)`

Set the short suite name for this script, if it belongs to a suite of related scripts. This should be a simple label suitable for use as a component in paths or filenames. Defaults to None.

`netsa.script.golem.set_name(name : str)`

Set the short name for this script. The name should be a simple label suitable for use as a component in paths or filenames. This defaults to the basename of the script file itself, minus the `.py` extension, if present.

`netsa.script.golem.set_interval([days : int, minutes : int, hours : int, weeks : int])`

Set how often this golem script is expected to generate results.

The interval roughly corresponds to how often the script should be run (such as from a cron job). Golem scripts will only process data for incomplete intervals over a provided date range, unless told otherwise via the `--overwrite` option.

`netsa.script.golem.set_span([days : int, minutes : int, hours : int, weeks : int])`

Set the span over which this golem script will expect input data for each processing interval. Defaults to one day.

The span will manifest as how much data is being pulled from a SiLK repository or possibly how many outputs from another golem script are being consumed. For example, a script having a 4 week span might run once a week, pulling 4 weeks worth of data each time it runs.

See *Intervals and Spans Explained* for more details.

```
netlsa.script.golem.set_lag([days : int, minutes : int, hours : int, weeks : int])
```

Set the lag for this golem script relative to the current date and time.

The default lag is 3 hours, a typical value for data to finish accumulating in a given hour within a SiLK repository. The lag effectively shifts the script's concept of the current time this far into the past.

```
netlsa.script.golem.set_realtime(enable=True)
```

Set whether or not this golem script will report output results in real time or not. Defaults to `False`.

Normally a golem script will wait until a processing interval has completely passed before performing any processing or reporting any output results. With `set_realtime` enabled, the golem script will consider the current processing bin to be the most recent, even if it extends to a future date.

Enabling realtime has a side effect of setting `lag` to zero.

```
netlsa.script.golem.set_tty_safe(enable=True)
```

Set whether or not query results are safe to send to the terminal. Defaults to `False`.

```
netlsa.script.golem.set_passive_mode(enable=False)
```

Controls whether or not an output option is required before running the main loop of the script. Defaults to `False`. Normally at least one repository-related or query-related option must be present or the script aborts. This is useful for scripts that require query behavior by default or are maintaining the repository in a custom fashion. If enabled, script authors should explicitly check whether or not repository updates were requested via the command line prior to updating the repository.

```
netlsa.script.golem.add_tag(name : str, value : str or func)
```

Set a command template tag with key `name`. The provided value can be callable, in which case it is resolved once the main function is invoked. Tags can reference other tags. See [template and tag usage](#) for more information.

```
netlsa.script.golem.add_loop(name : str, value : str list or func[, group_by : str or func,
                                                                    group_name : str, sep=', '])
```

Add a template tag under key `name` whose values cycle through those provided by `values`, either as an iterable or callable. In the latter case, the values are not resolved until the main function is invoked.

Optional keyword arguments:

**group\_by** Specifies how to group entries from `values` into a single loop entry. If the provided value is callable, the function should accept a single entry from `values` and return the group label for that entry or the original string, e.g. 'LAB2' might return 'LAB'. If the value is a dictionary, it will resolve to the mapped value if present. If it is an iterable of prefix matches, they will be converted into a regular expression to be applied to the beginning of each entry. In all cases, if there is no match or result, the original entry becomes its own group label.

**group\_name** This is the name of the template tag under which the group label appears, defaulting to the value of `name` appended with `'_group'`. In the above example, if `name` is 'sensor', then `%(sensor)s` might resolve to 'LAB1,LAB2,LAB3' whereas `%(sensor_group)s` would merely resolve to 'LAB'.

**sep** Depending on how these loops are visited, the template value under `name` might contain multiple values from `values`. Under these circumstances, the resulting string is joined by the value of 'sep' (default: ',')

Note that adding a loop will automatically add an additional query command line parameter named after the `name` for limiting which values to process within the loop. If grouping was requested, an additional parameter named after the `group_name` is also added. If these parameters are not desired or need to be modified, use the `modify_golem_param` function.

See [template usage](#) for more information on templates and loop values.

```
netlsa.script.golem.add_sensor_loop([name='sensor' : str, sensors : str list or func, group_by :
                                                                    str or func, group_name : str, auto_group=False])
```

This is a convenience function for adding a template loop tag under key `name` whose values are based on

sensors defined in a SiLK repository. Special note is taken that these loop values represent SiLK sensors, so any `netsa.script.Flow_params` tags that are defined (see `add_flow_tag`) will have their *sensors* parameter automatically bound (if not explicitly bound to something else) to the last sensor loop defined by this function.

Optional keyword arguments:

**name** The tag name within the template. Defaults to ‘sensor’, accessible from within templates as the tag `%(sensor)s`

**sensors** The source of sensor names, specified either as an iterable or callable. Defaults to the `get_sensors` function, which will interrogate the local SiLK repository once the main function is invoked.

**group\_by** Same as with `add_loop`

**group\_name** Same as with `add_loop`

**auto\_group** Causes *group\_by* to be set to the `get_sensor_group` function, a convenience function that strips numbers, possibly preceded by an underscore, from the end of sensor names. This can provide serviceable sensor grouping for descriptive sensor names (e.g. ‘LAB0’, ‘LAB1’, ‘LAB2’) but will not be of much use if they are generically named (e.g. ‘S0’, ‘S1’, etc).

```
netsa.script.golem.add_flow_tag(name : str[, flow_class : str, flow_type : str, flowtypes : str,
                                     sensors : str, start_date : str, end_date : str, input_pipe : str,
                                     xargs : str, filenames : str])
```

Add a `netsa.script.Flow_params` object as a template tag under key *name*. The rest of the keyword arguments correspond to the same parameters accepted by the `netsa.script.Flow_params` constructor and serve to map these fields to specific template tags in the golem script. If not otherwise specified, the *start\_date* and *end\_date* attributes are bound to the values of `%(golem_start_date)s` and `%(golem_end_date)s`, respectively, for each loop iteration. Additionally, if any sensor-specific loops were specified via `add_sensor_loop`, the *sensors* parameter defaults to the tag associated with the last defined sensor loop (typically `%(sensor)s`). The resulting `netsa.script.Flow_params` object is associated with a tag entry specified by *name*. The values of tags associated with `netsa.script.Flow_params` attributes in this way are still accessible under their original tag names.

The following optional keyword arguments are available to map template tag values to their corresponding attributes in the flow params object: *flow\_class*, *flow\_type*, *flowtypes*, *sensors*, *start\_date*, *end\_date*, *input\_pipe*, *xargs*, and *filenames*.

See *template and tag usage* for more information on templates.

```
netsa.script.golem.add_output_template(name : str, template : str[, scope : int, mime_type :
                                                                    str, description : str])
```

Define an output template tag key *name* with the provided *template*. The provided template can avail itself of the same set of tags available to command line templates. Absolute paths are not allowed.

The following optional keyword arguments are available. Pass any of them a value of `None` to disable entirely:

**scope** Defines how many intervals of this output are required to represent a complete analysis result in cases where the output from a single interval represents a partial result. For example, a golem script might have an interval of 1 day whereas a “complete” set of results is 7 days worth relative to any particular day.

**mime\_type** The expected MIME Content-Type of the output file, if any.

**description** A long-form text description, if any, of the contents of this output file.

See *template and tag usage* for more information on templates.

Templates are not required to reference all distinguishing tags and can therefore be ‘lossy’ across loop iterations if such a thing is desired.

```
netsa.script.golem.add_input_template(name : str, template : str[, required=True, mime_type : str, descriptions : str])
```

Define an input template tag under key *name*. This is useful for defining inputs not produced by other golem scripts. The template can be specified as a callable function, in which case the function is given a dictionary of currently defined template tags for each iteration and should return a template string if one is available for the current iteration.

The provided template can use the same set of tags available to output and command templates. The resolved string is available to command templates under key *name*.

Optional keyword arguments:

***required*** When `False` will ignore missing inputs once the template is resolved. Defaults to `True`.

***mime\_type*** The expected MIME Content-Type of the input.

***description*** A long-form text description of the expected contents of this input.

See [template and tag usage](#) for more information on templates.

```
netsa.script.golem.add_golem_input(golem_script : str, name : str[, output_name : str, count : int, cover=False, offset : int, span : timedelta, join_on : str or str list, join : dict, required=True])
```

Specify a tag, under key *name*, that represents the path (or paths) to an output of an external *golem script*. Only a single output can be associated at a time – if the external script has multiple outputs defined, then additional calls to this function are necessary for each output of interest.

For each output template, efforts are made to synchronize across time intervals and loop tags as appropriate. By default, this is the output of the most recent interval of the other golem script that corresponds to the interval currently under consideration within the local script.

Optional keyword arguments:

***output\_name*** The tag name used for this output in the external script if it differs from the value of *name* used locally for this input. By default the names are assumed to be identical.

***count*** Specifies how many intervals (as defined by the external script) of output data are to be used as input. By default, the most recent interval of the other golem script that corresponds to the local interval currently under consideration is provided.

For example, if the other golem script has an *interval* of one week, a *count* of 4 will provide the last 4 weeks of output from that script.

***offset*** Specify how far back (in units of the other script's interval) to reference for this input. Defaults to the most recent corresponding interval. If a *count* has been specified, the offset shifts the entire count of intervals.

For example, if the other script's interval is one week, an offset of -1 will reference the output from the week prior to the most recent week. Negative and positive offsets are equivalent for these purposes, they always reach backwards through time.

***cover*** If `True`, a *count* is calculated that will fully cover the local interval under consideration. This option cannot be used with the *count* or *offset* options. Defaults to `False`.

For example, if the local interval is one week, and the other interval is one day, then this is equivalent to specifying a *count* of 7 (days).

***span*** A `datetime` object that represents a span of time that covers the intervals of interest. Based on the other script's interval, a *count* is calculated that will cover the provided span. Cannot be used simultaneously with *count*, *offset*, or *cover*.

***join\_on*** A string or an iterable that specifies equivalent loop tags shared between the other golem input and the local script. Loops joined in this way will be synchronized in their iterations. By



using this parameter, the loop tags in both scripts are assumed to share the same name. If this is not the case, use the *join* parameter instead (described below).

**join** A dictionary or iterable of tuples that provides an equivalence mapping between template tags defined in the other golem script and locally defined tags.

For example, if the other script defines a template loop on the `% (my_sensor) s` tag, and the local script defines a loop on `% (sensor) s`, a mapping from ‘my\_sensor’ to ‘sensor’ will ensure that for each iteration over the values of `% (sensor) s` the input tag value is also sensor-specific based on its `% (my_sensor) s` loop. Without this association, the input tag would resolve to all outputs across all sensors, regardless of which sensor or sensor group was currently under local consideration. Iterations with no valid mapping are ignored, as opposed to when a valid association exists but the other output is missing.

**required** If `False` or `0` and the expected output from the other golem script is missing, continue processing rather than raising an exception. If specified as a positive integer, it means *at least* that many of the other inputs should exist, otherwise an exception is thrown. By default, at least one input is required. It is up to the developer to handle missing inputs (i.e. empty template tags) appropriately in these cases.

```
netlsa.script.golem.add_query_handler(name: str, query_handler: func)
```

Define a query handler under tag key *name* to be processed by the callable *query\_handler*. The output will only be generated dynamically when specifically requested via *query-related parameters*.

The provided callable will be passed the name of this query and a ‘tags’ dictionary for use with templates. And additional tag `% (golem_query_tgt) s` will be provided in the standard dictionary that contains the output path for this query. The function is responsible for creating this output.

```
netlsa.script.golem.add_self_input(name: str, output_name: str[, count: int, offset: int, span:
                                : timedelta])
```

Specify an input template tag, under key *name*, associated with this script’s own output from prior intervals. Since output tag names will necessarily collide, *name* is the new tag name for the input and *output\_name* is the name of the output template.

Optional keyword arguments:

**count** Same as with `add_golem_input`

**offset** Same as with `add_golem_input`, except defaults to `-1`. If self-referencing for purposes of delta-encoding, `0` should probably be specified.

**span** Same as with `add_golem_input`

Assuming a local loop over the template tag `% (sensor) s`, the following example:

```
>>> script.add_self_input('result', 'prior_result')
```

...is equivalent to this:

```
>>> script.add_golem_input(script.get_script_path(), 'prior_result',
>>>     output_name='result',
>>>     offset=-1,
>>>     join_on=['sensor'],
>>>     required=False)
```

---

**Note:** If this output happens to have a *scope* the developer should ensure that this self-reference means what is intended for ‘most recent’ output.

For example, if this golem script has an *interval* of 1 day, but the output has a *scope* of 28 (days), the example above would capture the prior 28 days of output due to the offset of `-1`.



If the script is delta-encoding its current result with its own prior results, however, probably what is desired would be the prior 27 days, in which case offset should be specified as 0. The 28th day in this scope scenario is the very result currently being generated, which does not exist yet, and therefore will not appear in the template as a ‘prior’ result. After processing is complete, however, it *will* appear in the collected outputs if a different golem script references this interval as input.

---

```
netsa.script.golem.add_golem_source(path : str)
```

Adds a directory within which to search for other golem scripts. Directories are searched in the following order: 1) paths within the colon-separated list of directories in the `GOLEM_SOURCES` environment variable; 2) any paths added by this function (multiple invocations allowed); 3) this script’s own directory as reported by the `get_script_dir` function, and finally 4) this script’s home directory as reported by the `get_home` function, if different than the script directory. These directories are only searched if the external script is specified as a relative path.

```
netsa.script.golem.get_script_path() → str
```

Returns the normalized absolute path to this script.

```
netsa.script.golem.get_script_dir() → str
```

Returns the normalized absolute path to the directory in which this script resides.

```
netsa.script.golem.get_home() → str
```

Returns the current value of this script’s home path if it has been set. Otherwise, defaults to the contents of the `GOLEM_HOME` environment variable (if set), or finally, the directory in which this script resides.

```
netsa.script.golem.get_repository() → str
```

Returns the current path for this script’s data output repository, or None if not set.

## 2.5 Parameter Functions

The following functions are used to enable and modify the standard command-line parameters available for golem scripts. No command line parameters are enabled by default, so at least one of these should be invoked in a typical golem script:

```
netsa.script.golem.add_golem_params([without_params : str list])
```

Enables all golem command line parameters. Equivalent to individually invoking each of `add_golem_basic_params`, `add_golem_repository_params`, and `add_golem_query_params`. Optionally accepts a list of parameters to exclude.

```
netsa.script.golem.add_golem_basic_params([without_params : str list])
```

Enables *basic* golem command line parameters. Optionally accepts a list of parameters to exclude.

```
netsa.script.golem.add_golem_query_params([without_params : str list])
```

Enables *query-related* golem command line parameters. Optionally accepts a list of parameters to exclude.

```
netsa.script.golem.add_golem_repository_params([without_params : str list])
```

Adds *repository-related* golem command line parameters. Optionally accepts a list of parameters to exclude.

```
netsa.script.golem.add_golem_param(name : str[, alias : str])
```

Enables a particular golem command line parameter. Accepts the same optional keyword parameters as the `modify_golem_param` function with the exception of `enable`, which is implied. An example is the `alias` parameter, which which can be provided to change the default parameter string (for example, aliasing `--last-date` to `--date` in cases where date ranges are not desired).

```
netsa.script.golem.modify_golem_param(name : str[, enabled : bool, alias : str, help : str, ...])
```

Modifies the settings for the given golem script parameter. Accepts new values for the golem-specific keywords

*enabled* and *alias*, along with the usual `netса.script` parameter keywords (e.g. *help*, *required*, *default*, *default\_help*, *description*, *mime\_type*)

## 2.6 Processing and Status Functions

The following functions are intended for use within the `main` function during processing or examination of status.

`netса.script.golem.execute(func)`

Executes the main function of a golem script. This should be called as the last line of any golem script, with the script's main function (whatever it might be named) as its only argument.

**Warning:** It is important that most, if not all, actual work the script does is done within this function. Golem scripts (as with all *NetSA Scripting Framework* scripts) may be loaded in such a way that they are not executed, but merely queried for metadata instead. If the golem script performs significant work outside of the `main` function, metadata queries will no longer be efficient. Golem scripts must use this `execute` function rather than `netса.script.execute`.

`netса.script.golem.process([golem_view : GolemView])`

Returns a `GolemProcess` wrapper around the given golem view, which defaults to the main script view. The result behaves much like a `GolemTags` object. Iterating over it returns a dictionary of resolved template tags while performing system level interactions (such as checking for input existence and creating output directories and/or files) in preparation for whatever processing the developer specifies in the processing loop. Views that have already completed processing are ignored.

`netса.script.golem.loop([golem_view : GolemView])`

Returns a `GolemTags` view of the given golem view, which defaults to the main script view. No system level processing happens while iterating over or interacting with this view.

`netса.script.golem.inputs([golem_view : GolemView])`

Returns a `GolemInputs` view of the given golem view, which defaults to the main script view. No system level processing happens while iterating over or interacting with this view.

`netса.script.golem.outputs([golem_view : GolemView])`

Returns a `GolemOutputs` view of the given golem view, which defaults to the main script view.. No system level processing happens while iterating over or interacting with this view.

`netса.script.golem.is_complete([golem_view : GolemView])`

Examines the outputs of the optionally provided `GolemView` object, which defaults to the main script view, and examines the status of the outputs for each processing interval. If all appear to be complete, returns `True`, otherwise `False`.

`netса.script.golem.script_view()`

Returns the currently defined global `GolemView` object.

`netса.script.golem.current_view([golem_view : GolemView])`

Returns a version of the given `GolemView` object, which defaults to the main script view, for the most recent interval available.

## 2.7 Utility Functions

The following functions provide some SiLK-specific tools and other potentially useful features for script authors.

`netса.script.golem.get_sensors() → str list`

Retrieves a list of sensors as defined in the local SiLK repository configuration.

`netса.script.golem.get_sensor_group(sensor : str) → str`

Convenience function, such as can be passed as a value for the *group\_by* parameter in `add_sensor_loop`, for extracting a sensor ‘group’ out of a sensor name. Groups are determined by extracting prefixes made of ‘word’ characters excluding ‘\_’. For example, two sensors called ‘LAB0’ and ‘LAB1’ would be grouped under ‘LAB’.

`netса.script.golem.get_sensors_by_group([grouper : func, sensors : str list])`

Convenience function that uses the callable *grouper* to construct a tuple of named pairs of the form (group\_name, members) suitable for use in constructing a dictionary. Defaults to the `get_sensor_group` function.

`netса.script.golem.get_args() → GolemArgs`

Returns a `GolemArgs` object containing any non-parameter arguments that were provided to the script on the command line.

## 2.8 Additional Functions

Please see the following sections in the `netса.script` documentation for details regarding other functions available within `netса.script.golem`:

- *Adding Script Parameters*
- *Verbose Output*
- *Flow Data Parameters*
- *Adding Additional Output*

## 2.9 Usage of Tags, Loops, and Templates

Golem is a templating engine. Based on a script’s configuration, the processing loop will iterate over time bins (processing intervals) and the values of any other loops defined for tags in the script. For each iteration, a dictionary of template tags is produced for use with such utilities as `netса.util.shell`. A typical golem script looks something like this:

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

from netса.script import golem
from netса.util import shell

# golem configuration and command line parameter
# configuration up here

def main():

    # All 'real' work should happen in this function.
    # Before invoking the processing loop, perhaps do some
    # configuration and prep work

    for tags in golem.process():

        # Set up per-iteration prep work, such as perhaps
        # some temp files

        # maybe modify contents of tags
        tags['temp_file'] = ...
```

```

# set up command templates

cmd1 = ...
cmd2 = ...
...

# run the commands
shell.run_parallel(cmd1, cmd2, vars=tags)

# pass main function to module for invocation and handling
golem.execute(main)

```

The following template tags are automatically available to each iteration over a golem processing loop:

Tag Name	Contents
golem_name	golem script name
golem_suite	suite name, if any
golem_span	timedelta obj for span
golem_interval	timedelta obj for interval
golem_span_iso	iso string repr of golem_span
golem_interval_iso	iso string repr of golem_interval
golem_bin_date	start datetime for this interval
golem_bin_year	interval datetime component 'year'
golem_bin_month	interval datetime component 'month'
golem_bin_day	interval datetime component 'day'
golem_bin_hour	interval datetime component 'hour'
golem_bin_second	interval datetime component 'second'
golem_bin_microsecond	interval datetime component 'microsecond'
golem_bin_iso	iso string repr for golem_bin_date
golem_bin_basic	iso basic string repr for golem_bin_date
golem_bin_silk	silk string repr for golem_bin_date
golem_start_date	start datetime for this span
golem_start_year	start datetime component 'year'
golem_start_month	start datetime component 'month'
golem_start_day	start datetime component 'day'
golem_start_hour	start datetime component 'hour'
golem_start_second	start datetime component 'second'
golem_start_microsecond	start datetime component 'microsecond'
golem_start_iso	iso string repr for golem_start_date
golem_start_basic	iso basic string repr for golem_start_date
golem_start_silk	silk string repr for golem_start_date
golem_end_date	end datetime for this span
golem_end_year	end datetime component 'year'
golem_end_month	end datetime component 'month'
golem_end_day	end datetime component 'day'
golem_end_hour	end datetime component 'hour'
golem_end_second	end datetime component 'second'
golem_end_microsecond	end datetime component 'microsecond'
golem_end_iso	iso string repr for golem_end_date
golem_end_basic	iso basic string repr for golem_end_date
golem_end_silk	silk string repr for golem_end_date

Continued on next page

Table 2.1 – continued from previous page

golem_view	the GolemView object which produced these tags
------------	--

For both intervals and spans, time bins are represented by the first date within that bin. For details on how intervals and spans relate to one another and format precisions, see *Intervals and Spans Explained*.

In addition to the standard golem tags defined above, all other tags, loops, inputs, and outputs defined in the initial script configuration are available for use in templates. For example, assume a script makes the following declarations:

```
script.add_tag('in_types', 'in,inweb')
script.add_tag('out_types', 'out,outweb')
script.add_tag('month', "%(golem_month)02d/%(golem_year)d")
script.add_sensor_loop()
script.add_flow_tag('in_flow', flow_type='in_types')
script.add_flow_tag('out_flow', flow_type='out_types')
script.add_output_template('juicy_set', ext='.set',
    description="Target 'juicy' set to generate.",
    mime_type='application/x-silk-ipset')
```

During each iteration of the processing loop, the tags dictionary will now include the following additional template entries:

#### in\_types

```
"in,inweb"
```

#### out\_types

```
"out,outweb"
```

#### month

```
"%(golem_month)02d/%(golem_year)d" \
    % (tags['golem_month'], tags['golem_year'])
```

#### sensor

```
the current iteration value of script.get_sensors()
```

#### in\_flow

```
Flow_params(
    start_date = tags['golem_start_date'],
    end_date   = tags['golem_end_date'],
    sensors    = tags['sensor'],
    flow_type  = tags['in_types'])
```

#### out\_flow

```
Flow_params(
    start_date = tags['golem_start_date'],
    end_date   = tags['golem_end_date'],
    sensors    = tags['sensor'],
    flow_type  = tags['out_types'])
```

#### juicy\_set

```
"%(golem_name)s/%(sensor)s/" \
    "%(golem_name)s.%(sensor)s.%(golem_start_date_iso)s.set" \
```

```
% (tags['golem_name'], tags['sensor'], \
   tags['golem_start_date_iso'])
```

Sometimes, depending on how loops and dependencies are arranged and how views are being manipulated, an input or output tag for the current view might contain multiple values (e.g. multiple filenames). In these cases, the resolved values are bundled into a `GolemArgs` object, which in turn resolves as a string of paths separated by spaces in the final command template.

## 2.10 Intervals and Spans Explained

Intervals and spans represent two different concepts.

An interval is a *processing interval* which represents how frequently the script is intended to produce results. This is roughly analogous to how frequently the script might be invoked via a cron job, except that golem scripts will back-fill missing results upon request and ignore intervals that appear to already have results present. An interval is always represented by the *first* timestamp contained within the interval.

A span, on the other hand, is a *data window* which represents how much input data the script is expected to consume, whether it be from a SiLK repository, results from other golem scripts, or other sources.

By default, the span of a golem script is the same as its interval. Not much surprising happens when the values are equal. They can be different, however. For example, a script might have a weekly interval yet consume 4 weeks worth of data for each of those weeks. Alternatively, a script might run every 4 weeks yet consume only a day's worth of data, akin to a monthly snapshot.

Intervals are *anchored* relative to a particular epoch. Intervals are always relative to midnight of the first Monday after January 1st, 1970 which was January 5th. Weeks therefore begin with Monday and multiples of weeks are always relative to that particular Monday.

Spans are always anchored relative to the *end* of the processing interval.

In the tags dictionary provided for each processing loop, the interval is represented by the `golem_bin_date` entry and the span is represented by `golem_start_date` and `golem_end_date`. Given a 3 week interval and a 4 week span, for example, these values are aligned like so:

```
interval:                bin_date                next_bin_date
                        |                         |
                        |-----|-----|-----|-----|
                        |                         |
span:      start_date                end_date
```

Note that `end_date` is not inclusive—its actual value is the value of `next_bin_date` minus one millisecond.

Each of these entries are represented by `datetime.datetime` objects along with an assortment of formatted string representations. If both the interval and span have a magnitude of at least a day or more, the formatted string variations look like so:

Variation	Format
iso	YYYY-MM-DD
basic	YYYYMMDD
silk	YYYY/MM/DD

If either the interval or span is less than a day, hours are included:

Variation	Format
iso	YYYY-MM-DDTHH
basic	YYYYMMDDTHH
silk	YYYY/MM/DDTHH

In all of the examples covered in this documentation, result templates are all based on the `golem_bin_date` values, i.e. the processing interval. As the example diagram above illustrates, this may not be intuitive as to what data is represented in the results. It is up to script authors to decide how to name their results, but they should choose a convention, stay consistent with it, and document the decision.

## 2.11 Examples

### 2.11.1 Trivial Example

The following is a simple example using the Golem API that demonstrates basic templating. The script monitors an “incoming” directory for daily text files containing IP addresses and converts them into rwsset files. A line by line explanation follows after the script:

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

from netsa.script.golem import script
from netsa.util import shell

script.set_title('Daily IP Sets')
script.set_description("""
    Convert daily IP lists into rwsset files.
""")
script.set_version("0.1")
script.set_contact("H.P. Fnord <fnord@example.com>")
script.set_authors(["H.P. Fnord <fnord@example.com>"])

script.set_name('daily_set')
script.set_interval(days=1)
script.set_span(days=1)

script.add_golem_params()

script.set_repository('dat')

script.add_input_template('daily_txt',
    "/data/incoming/daily.%(golem_bin_iso)s.txt",
    description="Daily IP text files",
    mime_type='text/plain')

script.add_output_template('daily_set',
    "daily/daily.%(golem_bin_iso)s.set",
    description="Daily IP Sets",
    mime_type='application/x-silk-ipset')

def main():
    cmd = "rwssetbuild %(daily_txt)s %(daily_set)s"
    for tags in script.process():
        shell.run_parallel(cmd, vars=tags)

script.execute(main)
```

Here is the breakdown, line by line:

```
from netsa.script.golem import script
from netsa.util import shell
```

The first two lines import `golem` itself as well as `netsa.util.shell`, which assists in constructing command line templates and executing the resulting system commands and pipelines.

Next are some lines for configuring meta information about the script:

```
script.set_title('Daily IP Sets')
script.set_description("""
    Convert daily IP lists into rwset files.
""")
script.set_version("0.1")
script.set_contact("H.P. Fnord <fnord@example.com>")
script.set_authors(["H.P. Fnord <fnord@example.com>"])
```

Setting the title, description, and other meta-data of the script is the same as with the regular `netsa.script meta-data functions`.

Now for the `golem`-specific configuration:

```
script.set_name('daily_set')
```

Though optional, every `golem` script should have a short name, suitable for inclusion withing directory paths and filenames. It will be made available for use in templates as the `%(golem_name)s` tag. For groups of related scripts, the `set_suite_name` function is also available.

Next, the script must be told the size of its processing intervals and the size of its data window:

```
script.set_interval(days=1)
script.set_span(days=1)
```

These two parameters, *interval* and *span*, are the core configuration parameters for any `golem` script. The *interval* represents ‘how often’ this script is expected to generate results. Typically this would correspond to the schedule by which the script is invoked via a cron job. The *span* represents how far back the script will look for input data. The interval and span do not have to match as they do here—for example, a script might have a ‘daily’ interval which processes one week of data for each of those days.

Next, the script author will almost always want to enable the standard `golem` command line parameters:

```
script.add_golem_params()
```

There are three general categories of parameters (*basic*, *repository-related*, and *query-related*) which can be separately enabled; the line above enables all of these.

Next, the script can be told where its results will live:

```
script.set_repository('dat')
```

This line defines the location of the scripts output data repository. Note that some scripts can be designed for query purposes only and will therefore not need to define a repository location.

If the path provided is a relative path, it is assumed to be relative to the script’s *home* path. See the `set_default_home` function for details on how the home path is configured or determined.

The script then defines a template for its input data:

```
script.add_input_template('daily_txt',
    "/data/incoming/daily.%(golem_bin_iso)s.txt",
    description="Daily IP text files",
    mime_type='text/plain')
```

This template assumes that the incoming files will correspond to the standard ISO-formatted datetime (‘YYYY-MM-DD’).



Next, an output template is defined:

```
script.add_output_template('daily_set',
    "daily/daily.%(golem_bin_iso)s.set",
    description="Daily IP Sets",
    mime_type='application/x-silk-ipset')
```

For each day of processing, a single rwsset file will be generated. Once again, the standard ISO-formatted date is chosen for the template.

In both the input and output templates, the script uses the tag `%(golem_bin_iso)s`. This tag is an implicit template tag, automatically available for use within golem scripts. Other timestamps are also available, including portions of each timestamp (such as year, month, and day) for constructing more elaborate templates. For more details about the rest of these ‘implicit’ template tags, see *Usage of Tags, Loops, and Templates*.

Now for the actual processing loop:

```
def main():
    cmd = "rwsetbuild %(daily_txt)s %(daily_set)s"
    for tags in script.process():
        shell.run_parallel(cmd, vars=tags)

script.execute(main)
```

All ‘real work’ in a golem script should take place in a main function, which is subsequently passed to the `execute` function. In order for golem scripts to work properly, this must always be the case. Using the `netlsa.script.execute` function instead will not work for a golem script.

The main entry point for looping across template values is the `process` function. This construct does a number of things, including creating output paths and checking for the existence of required inputs. On each iteration, a dictionary of template tags with resolved values is provided.

Every golem script has an intrinsic loop over processing intervals. In this example, our processing interval is once a day. If, via the command line parameters `--first-date` and `--last-date`, a window of 1 week had been specified, it would result in seven main iterations with the tag `%(golem_bin_iso)s` corresponding to a string representation of the beginning timestamp of each daily interval covered in the requested range. Unless told otherwise, a golem script will skip iterations which have already generated results.

## 2.11.2 Basic Example

The *Trivial Example* defined an input template that described daily text files without any explanation about where or how those files were produced. The following example assumes that the resulting set of addresses represent “observed internal hosts” and illustrates how the daily set files might be produced directly from queries to a SiLK repository. In order to do so, the script relies on a couple of SiLK command line tools:

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

from netlsa.script.golem import script
from netlsa.util import shell
from netlsa.data.format import datetime_silk_day

script.set_title('SiLK Daily Active Internal Hosts')
script.set_description("""
    Daily inventory of observed internal host activity.
""")
script.set_version("0.1")
script.set_contact("H.P. Fnord <fnord@example.com>")
script.set_authors(["H.P. Fnord <fnord@example.com>"])
```

```
script.set_name('daily_set')
script.set_interval(days=1)
script.set_span(days=1)

script.add_golem_params()

script.set_repository('dat')

script.add_tag('in_types', 'in,inweb')
script.add_tag('out_types', 'out,outweb')

script.add_output_template('daily_set',
    "internal/daily/daily.%(golem_bin_iso)s.set",
    description="Daily Internal host activity",
    mime_type='application/x-silk-ipset')

def main():
    for tags in script.process():
        tags['silk_start'] = datetime_silk_day(tags['golem_start_date'])
        tags['silk_end'] = datetime_silk_day(tags['golem_end_date'])
        tags['out_fifo'] = script.get_temp_dir_pipe_name()
        tags['in_fifo'] = script.get_temp_dir_pipe_name()
        cmd1 = [
            "rwfilter --start-date=%(silk_start)s"
            " --end-date=%(silk_end)s"
            " --type=%(in_types)s"
            " --proto=0-255 --pass=stdout",
            "rwset --sip=%(out_fifo)s"]
        cmd2 = [
            "rwfilter --start-date=%(silk_start)s"
            " --end-date=%(silk_end)s"
            " --type=%(out_types)s"
            " --proto=0-255 --pass=stdout",
            "rwset --dip=%(in_fifo)s"]
        cmd3 = [
            "rwsettool --union --output-path=%(internal_set)s"
            " %(out_fifo)s %(in_fifo)s"]
        shell.run_parallel(cmd1, cmd2, cmd3, vars=tags)

script.execute(main)
```

There are a couple of new techniques to note with this script. Below the standard meta-configuration are the following lines:

```
script.add_tag('in_types', 'in,inweb')
script.add_tag('out_types', 'out,outweb')
```

These two statements add a couple of simple template tags. All templates will now have access to the tags `%(in_types)s` and `%(out_types)s`, which will resolve to the strings `'in,inweb'` and `'out,outweb'`, respectively. This is equivalent to manually adding these entries to the `tags` dictionary down in the processing loop; predefining them here is a matter of style preference.

Next comes the main processing loop, which illustrates some more advanced usage of the `netlsa.util.shell` module:

```
def main():
    for tags in script.process():
        tags['silk_start'] = datetime_silk_day(tags['golem_start_date'])
```

```
tags['silk_end']    = datetime_silk_day(tags['golem_end_date'])
tags['out_fifo']   = script.get_temp_dir_pipe_name()
tags['in_fifo']    = script.get_temp_dir_pipe_name()
```

As mentioned earlier, a `tags` dictionary is provided for each processing interval and sensor. In the first four lines within the processing loop, some additional tags are added to the dictionary. The first two are reformat-  
ted date parameters destined to be used in the `rwfilter` commands. The `%(golem_start_date)s` and `%(golem_end_date)s` are template tags automatically provided by `golem`. The second two lines illustrate how the `netsa.script` module can be used to create temporary named pipes so that data can be fed from one command to another.

These new template additions are then used in the construction of some command templates used to pull data from the SiLK repository:

```
cmd1 = [
    "rwfilter --start-date=%(silk_start)s"
    " --end-date=%(silk_end)s"
    " --type=%(in_types)s"
    " --proto=0-255 --pass=stdout",
    "rwset --sip=%(out_fifo)s"]
cmd2 = [
    "rwfilter --start-date=%(silk_start)s"
    " --end-date=%(silk_end)s"
    " --type=%(out_types)s"
    " --proto=0-255 --pass=stdout",
    "rwset --dip=%(in_fifo)s"]
cmd3 = [
    "rwsettool --union --output-path=%(internal_set)s"
    " %(out_fifo)s %(in_fifo)s"]
shell.run_parallel(cmd1, cmd2, cmd3, vars=tags)
```

The first two command templates utilize the template definitions defined earlier, `%(in_types)s` and `%(out_types)s`, along with the date ranges associated with each processing loop. Each of these commands sends its results into its respective named pipe. Finally, the third command uses `rwsettool` to create a union from the output of these named pipes and creates the `rwset` file defined by the output template. All three commands are run in parallel using the facilities of the `netsa.util.shell` module.

### 2.11.3 Basic Golem Dependency Example

The *Basic Example* provides a `golem` script that produces daily `rwset` files produced from queries to a SiLK repository. What if a weekly, rather than daily, summary of IP addresses is desired? One option would be to adjust the processing interval and span of the script, thereby pulling an entire week's worth of data from SiLK in the calls to `rwfilter`. An alternative is to utilize the daily sets from the original script as inputs and construct a weekly summary via the union of the daily sets for the week in question.

One of the core features of `golem` scripts is that they can be assigned as inputs to one other. Details such as how often the inputs are produced, the naming scheme, and synchronization across time bins is sorted out automatically by the `golem` scripts involved. Assume that the script in the *Basic Example* is called `daily_set.py`. The following example illustrates how to configure the dependency on this external script:

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

from netsa.script.golem import script
from netsa.util import shell

script.set_title('Weekly Active Internal Host Set')
script.set_description("""
```

```

    Aggregate daily internal activity sets over the last week.
    """
    script.set_version("0.1")
    script.set_contact("H.P. Fnord <fnord@example.com>")
    script.set_authors(["H.P. Fnord <fnord@example.com>"])

    script.set_name('weekly_set')
    script.set_interval(weeks=1)
    script.set_span(weeks=1)

    script.add_golem_params()

    script.set_repository('dat')

    script.add_golem_input('daily_set.py', 'daily_set', cover=True)

    script.add_output_template('weekly_set',
                              "weekly/weekly.%(golem_bin_iso)s.set",
                              description="Aggregated weekly sets.",
                              mime_type='application/x-silk-ipset')

def main():
    cmd = "rwsettool --union --output-path=%(weekly_set)s %(daily_set)s"
    for tags in script.process():
        shell.run_parallel(cmd, vars=tags)

script.execute(main)

```

The first thing to note is that this script has a different interval and span:

```

script.set_interval(weeks=1)
script.set_span(weeks=1)

```

The script will produce weekly results and will expect to consume a week of data while doing so.

The input template is now defined as a dependency on the external script like so:

```

script.add_golem_input('daily_set.py', 'daily_set', cover=True)

```

The first argument is the name of the external script. For details on how relative paths to scripts are resolved, see the [add\\_golem\\_source](#) function.

The second argument is the output as defined within that external script. Golem scripts can have multiple outputs, so the specific output desired must be explicitly defined.

The third argument, the *cover* parameter, controls how this external output is synchronized across local processing intervals—in this case the processing interval of 1 week will be ‘covered’ by 7 days worth of inputs.

After the configuration of the weekly output template comes the main processing loop:

```

def main():
    cmd = "rwsettool --union --output-path=%(weekly_set)s %(daily_set)s"
    for tags in script.process():
        shell.run_parallel(cmd, vars=tags)

script.execute(main)

```

The output tag `%(weekly_set)s` is based on the `%(golem_bin_iso)s` timestamp, which in this case ends up being the date of the first Monday of each week in question. The `%(daily_set)s` tags represents 7 days of results—this will resolve to 7 individual filenames separated by whitespace in the eventual call to `rwsettool`.

### 2.11.4 Loop, Interval and Span Example

Golem scripts can define additional loops addition to the intrinsic loop over processing intervals. The following script is a modification of the script in the *Basic Example* which builds daily inventories by directly querying the SiLK repository. Rather than construct a monolithic inventory across all sensors, this version will construct inventories on a per-sensor basis by defining a template loop over a list of sensor names. Finally, it will illustrate the difference between intervals and spans by using a less frequent interval and a larger data window:

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

from netsa.script.golem import script
from netsa.util import shell
from netsa.data.format import datetime_silk_day

script.set_title('Active Internal Hosts')
script.set_description("""
    Per-sensor inventory of observed internal host activity over a
    four week window of observation, generated every three weeks.
""")
script.set_version("0.1")
script.set_contact("H.P. Fnord <fnord@example.com>")
script.set_authors(["H.P. Fnord <fnord@example.com>"])

script.set_name('internal')
script.set_interval(weeks=3)
script.set_span(weeks=4)

script.add_golem_params()

script.set_repository('dat')

script.add_tag('in_types', 'in,inweb')
script.add_tag('out_types', 'out,outweb')

script.add_loop('sensor', ["S0", "S1", "S2", "S3"])

script.add_output_template('internal_set',
    "internal/internal.%(sensor)s.%(golem_bin_iso)s.set",
    description="Internal host activity",
    mime_type='application/x-silk-ipset')

def main():
    for tags in script.process():
        tags['silk_start'] = datetime_silk_day(tags['golem_start_date'])
        tags['silk_end'] = datetime_silk_day(tags['golem_end_date'])
        tags['out_fifo'] = script.get_temp_dir_pipe_name()
        tags['in_fifo'] = script.get_temp_dir_pipe_name()
        cmd1 = [
            "rwfilter --start-date=%(silk_start)s"
            " --end-date=%(silk_end)s"
            " --type=%(in_types)s"
            " --sensors=%(sensor)s"
            " --proto=0-255 --pass=stdout",
            "rwset --sip=%(out_fifo)s"]
        cmd2 = [
            "rwfilter --start-date=%(silk_start)s"
            " --end-date=%(silk_end)s"
            " --type=%(out_types)s"]
```

```

        " --sensors=%(sensor)s"
        " --proto=0-255 --pass=stdout",
        "rwset --dip=%(in_fifo)s"]
cmd3 = [
    "rwsettool --union --output-path=%(internal_set)s"
    " %(out_fifo)s %(in_fifo)s"]
shell.run_parallel(cmd1, cmd2, cmd3, vars=tags)

script.execute(main)

```

The first thing to note is the new interval and span definitions:

```

script.set_interval(weeks=3)
script.set_span(weeks=4)

```

The script will produce results every 3 weeks and will expect to consume 4 weeks of data while doing so. This is the first example in which the interval and span are not equal. For more detail on the implications of this see *Intervals and Spans Explained*.

Further down in the script is the new loop definition:

```

script.add_loop('sensor', ["S0", "S1", "S2", "S3"])

```

With the addition of this line, for each 3-week processing interval, the script will return a separate `tags` dictionary for each sensor, setting the value of the `sensor` entry accordingly. Logically speaking this is equivalent to having two embedded ‘for’ loops, one for intervals and one for sensors.

This newly defined `%(sensor)s` tag is then used in the modified definition of the output template:

```

script.add_output_template('internal_set',
    "internal/internal.%(sensor)s.%(golem_bin_iso)s.set",
    description="Internal host activity",
    mime_type='application/x-silk-ipset')

```

Next comes the main processing loop. Note that it is *identical* to the processing loop in the earlier incarnation of the script. The interval and span were changed, an extra loop was introduced, and the output template was modified, but the essential processing logic remains unchanged.

### 2.11.5 SiLK Integration Example

The Golem API and *NetSA Scripting Framework* include a number of convenience functions and classes for interacting with a SiLK repository. The *Loop, Interval and Span Example* can be simplified using a few of these features as illustrated below:

```

#!/usr/bin/env python

from netsa.script.golem import script
from netsa.util import shell
from netsa.data.format import datetime_silk_day

script.set_title('Active Internal Hosts')
script.set_description("""
    Per-sensor inventory of observed internal host activity over a
    four week window of observation, generated every three weeks.
""")
script.set_version("0.2")
script.set_contact("H.P. Fnord <fnord@example.com>")
script.set_authors(["H.P. Fnord <fnord@example.com>"])

```

```

script.set_name('internal')
script.set_interval(weeks=3)
script.set_span(weeks=4)

script.add_golem_params()

script.set_repository('dat')

script.add_tag('in_types', 'in,inweb')
script.add_tag('out_types', 'out,outweb')

script.add_sensor_loop()

script.add_flow_tag('in_flow', flow_type='in_types')
script.add_flow_tag('out_flow', flow_type='out_types')

script.add_output_template('internal_set',
    "internal/internal.%(sensor)s.%(golem_bin_iso)s.set",
    description="Internal host activity",
    mime_type='application/x-silk-ipset')

def main():
    for tags in script.process():
        tags['out_fifo'] = script.get_temp_dir_pipe_name()
        tags['in_fifo'] = script.get_temp_dir_pipe_name()
        cmd1 = [
            "rwfilter %(out_flow)s --proto=0-255 --pass=stdout",
            "rwset --sip=%(out_fifo)s"]
        cmd2 = [
            "rwfilter %(in_flow)s --proto=0-255 --pass=stdout",
            "rwset --dip=%(in_fifo)s"]
        cmd3 = [
            "rwsettool --union --output-path=%(internal_set)s"
            " %(out_fifo)s %(in_fifo)s"]
        shell.run_parallel(cmd1, cmd2, cmd3, vars=tags)

script.execute(main)

```

The first difference to note is that rather than manually defining a loop over sensors, the following shorthand is used:

```
script.add_sensor_loop()
```

This line sets up a loop on the template tag `sensor` as before, but the list of sensors is automatically determined from the SiLK repository itself (see the `mapsid` command). The script also remembers that this particular loop involves sensors.

The next modification to note is the definition of two special SiLK-related compound tags:

```

script.add_flow_tag('in_flow', flow_type='in_types')
script.add_flow_tag('out_flow', flow_type='out_types')

```

These statements create template entries bound to `netsa.script.Flow_params` objects which serve to simplify the construction of `rwfilter` command line templates.

Each call to `add_flow_tag` implicitly binds the `start_date` and `end_date` object attributes to the value of the template tags `golem_start_date` and `golem_end_date`. Given that a sensor-specific loop was declared earlier, the function calls will also bind the `sensors` attribute to the value of the `sensor` tag for each loop.

Additional tags can be bound to `netsa.script.Flow_params` attributes using keyword arguments. In this

example, the `in_types` and `out_types` tags defined earlier in the script are bound to the `flow_type` attribute of each object.

The rest of the script proceeds as before, except that in the processing loop the `rwfilter` command templates are far more compact:

```
cmd1 = [
    "rwfilter %(out_flow)s --proto=0-255 --pass=stdout",
    "rwset --sip=%(out_fifo)s"]
cmd2 = [
    "rwfilter %(in_flow)s --proto=0-255 --pass=stdout",
    "rwset --dip=%(in_fifo)s"]
```

The `%(out_flow)s` and `%(in_flow)s` tags will each expand into four parameters in the eventual command string.

### 2.11.6 Synchronization Example

The following example will build a daily inventory of internal addresses that exhibit activity on source port 25. In order to limit the pool of addresses under consideration, it will utilize the internal inventory results generated by the *SiLK Integration Example*. Furthermore, it will utilize some additional SiLK-related tools in order to organize results into ‘sensor groups’ rather than under individual sensors. The following assumes that the prior inventory script is called `internal.py`:

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

from netsa.script.golem import script
from netsa.util import shell

script.set_title('Daily Internal Port 25 Activity')
script.set_description("""
    Daily per-sensor-group inventory of observed internal host
    activity on port 25.
""")
script.set_version("0.1")
script.set_contact("H.P. Fnord <fnord@example.com>")
script.set_authors(["H.P. Fnord <fnord@example.com>"])

script.set_name('p25')
script.set_interval(days=1)
script.set_span(days=1)

script.add_golem_params()

script.set_repository('dat')

script.add_sensor_loop(auto_group=True)

script.add_tag('out_type', 'out,outweb')
script.add_flow_tag('out_flow', flow_type='out_type')

script.add_golem_input('internal.py', 'internal_set', join_on='sensor')

script.add_output_template('p25_set',
    "p25/p25.%(sensor_group)s.%(golem_bin_iso)s.set",
    description="Daily IP set for internal port 25 activity",
    mime_type='application/x-silk-ipset')
```



```
def main():
    for tags in script.process():
        tags['in_fifo'] = script.get_temp_dir_pipe_name()
        cmd1 = [
            "rwsettool --union"
            " --output-path=%(in_fifo)s"
            " %(internal_set)s"]
        cmd2 = [
            "rwfilter %(out_flow)s"
            " --proto=6"
            " --sport=25"
            " --packets=2-"
            " --sipset=%(in_fifo)s"
            " --pass=stdout",
            "rwset --sip-set=%(p25_set)s"]
        shell.run_parallel(cmd1, cmd2, vars=tags)

script.execute(main)
```

First, the script is configured to generate once per day using a span of one day:

```
script.set_interval(days=1)
script.set_span(days=1)
```

Next, the sensor loop is configured:

```
script.add_sensor_loop(auto_group=True)
```

This invocation of `add_sensor_loop` uses a new named parameter, *auto\_group*, which loops over *groups* of related sensors rather than individual sensors. Normally, a single template tag `sensor` is added. When grouping is enabled for a sensor loop another tag `sensor_group` is added in addition to the `sensor` tag. So, for example, if there are three sensors in a group labeled 'LAB0', 'LAB1', and 'LAB2', these two template tags would expand into strings like so:

Tag	Value
<code>%(sensor)s</code>	LAB0, LAB1, LAB2
<code>%(sensor_group)s</code>	LAB

See the `add_loop` function for the details of how the above features work for generic, non-sensor-related, loops.

Next the script sets up the input dependency from the script in the *Basic Example* called `internal.py`:

```
script.add_golem_input('internal.py', 'internal_set', join_on='sensor')
```

When golem scripts use other golem script results as inputs, they are automatically synchronized across processing intervals. The basic rule is to synchronize on the latest external interval containing an end-point less than or equal to the end-point of the local interval under consideration.

The synchronization of any loops other than intervals must be explicitly configured. In this case, the *join\_on* parameter is used to indicate that the external sensor loop and local sensor loop should align on each value of the `sensor` tag. This synchronization happens per-sensor and does not affect the eventual sensor grouping behavior.

Next, the output template is defined. Note the use of the `sensor_group` tag rather than `sensor`:

```
script.add_output_template('p25_set',
    "p25/p25.%(sensor_group)s.%(golem_bin_iso)s.set",
    description="Daily IP set for internal port 25 activity",
    mime_type='application/x-silk-ipset')
```

Followed by the processing loop:

```
def main():
    for tags in script.process():
        tags['in_fifo'] = script.get_temp_dir_pipe_name()
        cmd1 = [
            "rwsettool --union"
            " --output-path=%(in_fifo)s"
            " %(internal_set)s"]
        cmd2 = [
            "rwfilter %(out_flow)s"
            " --proto=6"
            " --sport=25"
            " --packets=2-"
            " --sipset=%(in_fifo)s"
            " --pass=stdout",
            "rwset --sip-set=%(p25_set)s"]
        shell.run_parallel(cmd1, cmd2, vars=tags)

script.execute(main)
```

Since sensors are being grouped, the `%(internal_set)s` tag for each loop potentially represents multiple input files, one for each individual sensor. The first command defines a template for `rwsettool` that sends a union of these per-sensor sets into the named pipe. The second command pipeline uses this merged set to filter the initial flows being examined by the `rwfilter` query.

When invoked on a regular basis, this script will produce a daily subset of the most recent per-sensor-group inventory for those internal IP addresses that have exhibited activity on source port 25.

### 2.11.7 Self Dependency Example

The *Synchronization Example* demonstrates how to configure an input dependency on the results of another golem script. It is also possible to configure dependencies on a golem script's *own* past results.

Recall that the *SiLK Integration Example* is configured with a 3-week interval and 4-week span. The 3-week interval was chosen due to the resource-intensive query across 4 weeks of data. Whereas this does produce internal inventories, the information is potentially less accurate over time (particularly during the final few days of the 3-week processing interval).

The inventory script can be modified to consume its own outputs and produce *delta encoded* results on a daily basis:

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

from netsa.script.golem import script
from netsa.util import shell
from netsa.data.format import datetime_silk_day

script.set_title('Active Internal Hosts')
script.set_description("""
    Daily per-sensor inventory of observed internal host activity,
    delta-encoded using the prior four weeks of results.
""")
script.set_version("0.3")
script.set_contact("H.P. Fnord <fnord@example.com>")
script.set_authors(["H.P. Fnord <fnord@example.com>"])

script.set_name('internal')
script.set_interval(days=1)
script.set_span(days=1)
```

```

script.add_golem_params()

script.set_repository('dat')

script.add_tag('in_types', 'in,inweb')
script.add_tag('out_types', 'out,outweb')

script.add_sensor_loop()

script.add_flow_tag('in_flow', flow_type='in_types')
script.add_flow_tag('out_flow', flow_type='out_types')

script.add_output_template('internal_set',
    "internal/internal.delta.%(sensor)s.%(golem_bin_iso)s.set",
    description="Delta set of internal host activity",
    mime_type='application/x-silk-ipset',
    scope=28)

script.add_self_input('prior_set', 'internal_set', offset=0)

def main():
    for tags in script.process():
        tags['out_fifo'] = script.get_temp_dir_pipe_name()
        tags['in_fifo'] = script.get_temp_dir_pipe_name()
        cmds = []
        cmds.append([
            "rwfilter %(out_flow)s --proto=0-255 --pass=stdout",
            "rwset --sip=%(out_fifo)s"])
        cmds.append([
            "rwfilter %(in_flow)s --proto=0-255 --pass=stdout",
            "rwset --dip=%(in_fifo)s"])
        cmds.append([
            "rwsettool --union --output-path=%(current_set)s"
            " %(out_fifo)s %(in_fifo)s"])
        if tags['prior_set']:
            tags['current_set'] = script.get_temp_dir_pipe_name()
            cmds.append([
                "rwsettool --difference"
                " --output-path=%(internal_set)s"
                " %(current_out)s %(prior_set)s"])
        else:
            tags['current_set'] = tags['internal_set']

        shell.run_parallel(vars=tags, *cmds)

script.execute(main)

```

The goal is to generate a viable internal inventory on a daily basis with minimal overhead. The naive approach would be to define an interval of 1 day and leave the span as 4 weeks. This would pull 4 weeks of data *every single day* and construct a full inventory for that day. This is inefficient in terms of processing and storage. Instead, this script introduces a new concept called *scope*. Scope is used to indicate situations where a single interval of processing does not represent a complete analysis result.

First, the basics are configured:

```

script.set_interval(days=1)
script.set_span(days=1)

```

The script produces a daily result and expects to consume a single day's worth of 'regular' data while doing so. Next,

the script must define its daily output template:

```
script.add_output_template('internal_set',
    "internal/internal.delta.%(sensor)s.%(golem_bin_iso)s.set",
    description="Delta set of internal host activity",
    mime_type='application/x-silk-ipset',
    scope=28)
```

This declaration shows the use of the new `scope` parameter. The scope indicates the number of processing interval outputs required to represent a *complete* result. Here, the scope is defined as 28 intervals (days in this case).

Now when other golem scripts use this script output as an input dependency, they will see 4 weeks of files relative to each day of interest. This also applies in cases where a golem script asks *itself* for prior results. An example of this is shown next:

```
script.add_self_input('prior_set', 'internal_set', offset=0)
```

This self-referential input dependency maps `internal_set` to a new template tag called `prior_set`.

By default, self-referential inputs have an offset of `-1` which excludes the results for the current processing interval. In cases such as this, where the goal is delta-encoding, the offset should be `0`. (The daily result being generated for the current day represents addresses not present in the last 27 days).

Next is the main processing loop:

```
def main():
    for tags in script.process():
        tags['out_fifo'] = script.get_temp_dir_pipe_name()
        tags['in_fifo'] = script.get_temp_dir_pipe_name()
        cmds = []
        cmds.append([
            "rwfilter %(out_flow)s --proto=0-255 --pass=stdout",
            "rwset --sip=%(out_fifo)s"])
        cmds.append([
            "rwfilter %(in_flow)s --proto=0-255 --pass=stdout",
            "rwset --dip=%(in_fifo)s"])
        cmds.append([
            "rwsettool --union --output-path=%(current_set)s"
            " %(out_fifo)s %(in_fifo)s"])
        if tags['prior_set']:
            tags['current_set'] = script.get_temp_dir_pipe_name()
            cmds.append([
                "rwsettool --difference"
                " --output-path=%(internal_set)s"
                " %(current_out)s %(prior_set)s"])
        else:
            tags['current_set'] = tags['internal_set']

        shell.run_parallel(vars=tags, *cmds)

script.execute(main)
```

The core logic is similar to the earlier version. A new template tag, `current_set` is added to the `tags` dictionary for each iteration. Depending on circumstances, the value of this tag is set to one of two things: If no prior results are available, a regular `rwset` is constructed just as it was before. If prior results are available, however, the difference is taken between the current day's results and the union of up to 27 days of prior results.

This technique allows the reconstruction of an accurate 4-week internal inventory, for any particular day, by taking the union over the 28 days ending on that day.

Having made these changes, what now needs to be changed in the script from the *Synchronization Example* which depends on these internal sets as input?

Not a single thing.

The script in the *Synchronization Example* is already performing a union with `rwsettool` on the tag `%(internal_set)s` in order to merge data across sensors into sensor groups. Due to the `scope` declaration, the `%(internal_set)s` tag will now also include paths to the files for each of the 28 days required to reconstruct results.

## 2.12 Classes

### 2.12.1 GolemView

**class** `netsa.script.golem.GolemView` (`golem : Golem`, `first_date : datetime`, `last_date : datetime` )

A `GolemView` object encapsulates a golem script model and is used to view and manipulate it in various ways. These different views are primarily accessed through the `loop`, `outputs`, and `inputs` methods.

Optional keyword arguments:

**`last_date`** The interval containing this `datetime` object is the last to be considered for processing.  
(default: most recent)

**`first_date`** The interval containing this `datetime` object is the first to be considered for processing.  
(default: `last_date`)

**`golem`**

The golem script model which this view manipulates.

**`first_bin`**

A `datetime` object representing the first processing interval for this view, as determined by the `first_date` and `last_date` parameters during construction. Defaults to `last_bin`.

**`last_bin`**

A `datetime` object representing the last processing interval for this view, as determined by the `last_date` and `first_date` parameters during construction. Defaults to the ‘most recent’ interval that does not overlap into the future, taking into account `lag`.

**`start_date`**

A `datetime` object representing the beginning of the first data span covered by this view. Spans can be larger (or smaller) than the defined interval, so this value is not necessarily equal to `first_bin`.

**`end_date`**

A `datetime` object representing the end of the last data span represented by this view. This is determined based the interval defined by `last_bin`: `last_bin + interval - 1 millisecond`.

**`using`** (`[golem : Golem, first_date : datetime, last_date : datetime]`)

Return a copy of this `GolemView` object, optionally using new values for the following keyword arguments:

**`golem`** Use a different golem script model.

**`first_date`** Select a different starting time bin based on the provided `datetime` object.

**`last_date`** Select a different ending time bin based on the provided `datetime` object.

**`bin_dates`** () → `datetime` iter

Provide an iterator over `datetime` objects representing all processing intervals represented by this view.

**by\_bin\_date** () → GolemView iter

Provide an iterator over `GolemView` objects for each interval represented by this view.

**group\_by** (key : str[, ...]) → (str tuple, GolemView) iter

Returns an iterator that yields a tuple with a primary key and `GolemView` object grouped by the provided keys. Each primary key is a tuple containing the current values of the keys provided to `group_by`. Iterating over the provided view objects will resolve any remaining loops if any remain that were not used for the provided key.

**by\_key** (key : str) → (str, GolemView) iter

Similar to `group_by` but takes only a single key as an argument. Returns an iterator that yields view objects for each value of the key; iterating over the provided view objects will resolve any remaining loops, if present.

**product** () → GolemView iter

Fully resolve the loops defined within this view. The ‘outer’ loop is always over intervals, followed by any other loops in the order in which they were defined. Each view thus provided is therefore fully resolved, with no loops remaining.

**bin\_count** () → int

Return the number of intervals represented by this view, as defined by `first_bin` and `last_bin`.

**loop\_count** () → int

Return the number of non-interval iterations represented by this view that are produced by resolving any defined loops.

**sync\_to** (other : GolemView [, count : int, offset : int, cover=False) → GolemView

Given another `GolemView` object, return a version of *self* that has been synchronized to the given view object.

Optional keyword arguments:

**count** Synchronize to this many intervals of the given object (default: 1)

**offset** Synchronize to this many interval offsets behind the given object (default: 0)

**cover** Calculate a *count* necessary to cover all intervals represented by the given object (overrides *count* and *offset*)

**loop** () → GolemTags

Return a `GolemTags` object representing this view.

**outputs** () → GolemOutputs

Return a `GolemOutputs` object representing this view.

**inputs** () → GolemInputs

Return a `GolemInputs` object representing this view.

**\_\_len\_\_** () → int

Return the number of fully-resolved iterations represented by this view, over intervals as well as any defined loops.

**\_\_iter\_\_** () → GolemView iter

Iterates over the views returned by the `product` method.

## 2.12.2 GolemTags

**class** `netса.script.golem.GolemTags` (golem : Golem[, first\_date : datetime, last\_date : datetime])

Bases: `netса.script.golem.GolemView`

A `GolemTags` object is used to examine resolved template tags produced by looping over intervals and other defined loops.

As well as the methods and attributes of `GolemView`, the following additional and overridden methods are available:

**tags** () → dict

Return a dictionary of resolved template tags for the current view (flattens tags across the loops that would result from invoking the `product` method).

**\_\_iter\_\_** () → dict iter

### 2.12.3 GolemOutputs

**class** `netса.script.golem.GolemOutputs` (*golem* : *Golem*[, *first\_date* : *datetime*, *last\_date* : *datetime*])

Bases: `netса.script.golem.GolemView`

A `GolemOutputs` object is used to examine resolved output templates, either for a specific iteration or aggregated across multiple iterations.

As well as the methods and attributes of `GolemView`, the following additional and overridden methods are available:

**expand** () → `GolemArgs`

Returns a `GolemArgs` object representing all resolved output templates for the current view.

**\_\_len\_\_** () → int

Return the number of resolved output templates for the current view.

**\_\_iter\_\_** () → str iter

Iterate over each resolved output template for the current view.

### 2.12.4 GolemInputs

**class** `netса.script.golem.GolemInputs` (*golem* : *Golem*[, *first\_date* : *datetime*, *last\_date* : *datetime*])

Bases: `netса.script.golem.GolemView`

A `GolemInputs` object is used to examine resolved input templates, either for a specific iteration or aggregated across multiple iterations.

As well as the methods and attributes of `GolemView`, the following additional and overridden methods are available:

**expand** () → `GolemArgs`

Returns a `GolemArgs` object representing all resolved input templates for the current view.

**members** () → `GolemOutputs` iter

Iterate over each golem script that provides inputs for the for this golem script, returning each as a synchronized `GolemOutputs` object.

**\_\_len\_\_** () → int

Return the number of resolved input templates for the current view.

**\_\_iter\_\_** () → str iter

Iterate over each resolved input template for the current view.

## 2.12.5 GolemArgs

**class** `netsa.script.golem.GolemArgs` (*item* : *str* or *str iter*[, ...])

A `GolemArgs` object encapsulates a list of resolved input or output templates destined to be used as a parameter in a tags dictionary. The constructor takes any number of strings, or string iterators, and flattens them into a unique list in the order they were seen.

It will resolve to a string of space-separated values and will properly resolve when passed to the `netsa.util.shell` module for command and pipeline execution.

Note that some file-related python functions (such as `open`) will complain if passed a single-value `GolemArgs` object (representing a single file name) without having first explicitly converted it to a string via `str`.

The length of a `GolemArgs` object represents the number of items it contains. These can be accessed via an index like a list. Two objects can be added and subtracted from one another, as with sets.

## 2.12.6 GolemProcess

**class** `netsa.script.golem.GolemProcess` (*gview* : *GolemView*[, *overwrite\_outputs*=*False*,  
*skip\_complete*=*True*, *keep\_empty\_outputs*=*False*,  
*skip\_missing\_inputs*=*False*, *optional\_inputs* : *dict*])

A utility class for performing system-level interactions (such as checking for required inputs, pre-existing outputs, creating output paths, etc) while iterating over the provided view.

***overwrite\_outputs*** Delete existing outputs prior to processing. (default: *False*)

***keep\_empty\_outputs*** Consider zero-byte output results to be valid, otherwise they will be ignored or deleted when encountered, regardless of the value of *overwrite\_results*. (default: *False*)

***skip\_complete*** Ignore processing bins whose results appear to be completed. (default: *True*)

***skip\_missing\_inputs*** If insufficient inputs are present, determines whether this iteration should be skipped, as opposed to raising an exception. (default: *False*)

Most methods and attributes available from the `GolemTags` class are available through this class as well, with some behavioral changes as noted below. The following methods are in addition to those available from `GolemTags`:

***is\_complete*** () → *bool*

Returns a boolean value indicating whether processing has been completed for the intervals represented by this view.

***status*** (*label* : *str*) → (*str*, *bool*) *iter*

Iterate over items within the given tag, returning a tuple containing the item string and its current status. Status is typically the size in bytes of each input or output, or *None* if it does not exist.

The following methods have slightly different behavior than that of `GolemTags`:

***using*** ([*gview* : *GolemView*, *overwrite\_outputs*=*False*, *skip\_complete*=*True*,  
*keep\_empty\_outputs*=*False*, *skip\_missing\_inputs*=*False*, *create\_slots*=*True*, *optional\_inputs* :  
*dict*])

Return a copy of this `GolemProcess` object, possibly replacing certain attributes corresponding to the keyword arguments in the constructor.

***product*** () → *GolemProcess iter*

Return a `GolemProcess` object for each iteration over the processing intervals and loops defined for this process view, possibly performing system level tasks along the way (such as creating output paths and performing input checks). Iterations where processing is complete will be skipped, unless *overwrite\_outputs* has been enabled for this object.



**by\_bin\_date**() → GolemProcess iter

Provide an iterator over `GolemProcess` objects for each processing interval represented by this view, possibly performing system level tasks along the way. Iterations for which processing is complete will be skipped, unless `overwrite_outputs` has been enabled for this object.

**group\_by**(key : str[, ...]) → (str tuple, GolemProcess) iter

Returns an iterator that yields a tuple with a primary key and a `GolemProcess` object, grouped by the provided keys, possibly performing system level tasks along the way. Each primary key is a tuple containing the current values of the keys provided. Iterating over the resulting process objects process objects will resolve any remaining loops remaining in that view, if any. Views for which processing is complete will be skipped, unless `overwrite_outputs` has been enabled for this object.

**by\_key**(str) → GolemProcess iter

Similar to `group_by` but takes a single key as an argument. Returns and iterator that yields `GolemProcess` objects for each value of the key, possibly performing system level tasks along the way. Views for which processing is complete will be skipped, unless `overwrite_outputs` has been enabled for this object.

**\_\_iter\_\_**() → dict iter

Iterate over the views produced by the `product` method, yielding a dictionary of fully resolved template tags. Iterations for which processing is complete will be skipped, unless `overwrite_outputs` has been enabled for this object.



# NETSA.SQL — SQL DATABASE ACCESS

## 3.1 Overview

The normal flow of code that works with databases using the `netsa.sql` API looks like this:

```
from netsa.sql import *

select_stuff = db_query("""
    select a, b, c
    from test_table
    where a + b <= :threshold
    limit 10
""")

conn = db_connect("nsql-sqlite:/var/tmp/test_db.sqlite")

for (a, b, c) in conn.execute(select_stuff, threshold=5):
    print ("a: %d, b: %d, c: %d, a + b: %d" % (a, b, c, a+b))

# Alternatively:
for (a, b, c) in select_stuff(conn, threshold=5):
    print ("a: %d, b: %d, c: %d, a + b: %d" % (a, b, c, a+b))
```

First, the required queries are created as instances of the `db_query` class. Some developers prefer to have a separate module containing all of the queries grouped together. Others prefer to keep the queries close to where they are used.

When the database is to be used, a connection is opened using `db_connect`. The query is executed using `db_connection.execute`, or by calling the query directly. The result of that call is then iterated over and the data processed.

Connections and result sets are automatically closed when garbage collected. If you need to make sure that they are collected as early as possible, make sure the values are not kept around in the environment (for example, by assigning `None` to the variable containing them when your work is complete, if the variable won't be leaving scope for a while.)

## 3.2 Exceptions

**exception** `netsa.sql.sql_exception` (*message* : *str*)  
Specific exceptions generated by `netsa.sql` derive from this.

**exception** `net.sa.sql.sql_no_driver_exception(message : str)`

This exception is raised when no driver is installed that can handle a URL opened via `db_connect`.

**exception** `net.sa.sql.sql_invalid_uri_exception(message : str)`

This exception is raised when the URI passed to `db_connect` cannot be parsed.

### 3.3 Connecting

`net.sa.sql.db_connect(uri[, user : str, password : str]) → db_connection`

Given a database URI and an optional *user* and *password*, attempts to connect to the specified database and return a `db_connection` subclass instance.

If a user and password are given in this call as well as in the URI, the values given in this call override the values given in the URI.

Database URIs have the form:

```
<scheme>://<user>:<password>@hostname:port/<path>;<param>=<value>;...?<query>#<fragment>
```

Various pieces can be left out in various ways. Typically, the following form is used for databases with network addresses:

```
<scheme>://[user[:password]@]hostname[:port]/<dbname>[;<parameters>]
```

While the following form is used for databases without network addresses, or sometimes for connections to databases on the local host:

```
<scheme>:<dbname>[;user=<user>][;password=<password>][;<params>]
```

The user and password may always be given either in the network location or in the params. Values given in the `db_connect` call override either of those, and values given in the network location take priority over those given in the params.

Refer to a specific database driver for details on what URI scheme to use, and what other params or URI pieces may be meaningful.

### 3.4 Connections and Result Sets

**class** `net.sa.sql.db_connection(driver : db_driver, variants : str list)`

An open database connection, returned by `db_connect`.

**get\_driver()** → `db_driver`

Returns the `db_driver` used to open this connection.

**clone()** → `db_connection`

Returns a fresh open `db_connection` open to the same database with the same options as this connection.

**execute(query\_or\_sql : db\_query or str[, <param\_name>=<param\_value>, ...])** → `db_result`

Executes the given SQL query (either a SQL string or a query compiled with `db_query`) with the provided variable bindings for side effects. Returns a `db_result` result set if the query returns a result set, an `int` with the number of rows affected if available, or `None` otherwise.

**commit()**

Commits the current database transaction in progress. Note that if a `db_connection` closes without `commit` being called, the transaction will automatically be rolled back.

**rollback()**

Rolls back the current database transaction in progress. Note that if a `db_connection` closes without `commit` being called, the transaction will automatically be rolled back.

**get\_variants()** → str seq

Returns which variant tags are associated with this connection.

**class** `netsa.sql.db_result` (*connection : db\_connection, query : db\_query, params : dict*)

A database result set, which may be iterated over.

**get\_connection()** → db\_connection

Returns the `db_connection` which produced this result set.

**get\_query()** → db\_query

Returns the `db_query` which was executed to produce this result set. (Note that if a string query is given to `db_connection.execute`, it will automatically be wrapped in a `db_query`, so this is always a `db_query`.)

**get\_params()** → dict

Returns the `dict` of params which was given when this query was executed.

**\_\_iter\_\_()** → iter

Returns an iterator over the rows of this result set. Each row returned is a `tuple` with one item for each column. If there is only one column in the result set, a tuple of one column is returned. (e.g. `(5,)`, not just `5` if there is a single column with the value five in it.)

It is an error to attempt to iterate over a result set more than once, or multiple times at once.

## 3.5 Compiled Queries

**class** `netsa.sql.db_query` (*sql : str[, <variant> : str, ...]*)

A `db_query` represents a “compiled” database query, which will be used one or more times to make requests.

Whenever a query is executed using the `db_connection.execute` method, it may be provided as either a string or as a `db_query` object. If an object is used, it can represent a larger variety of possible behaviors. For example, it might give both a “default” SQL to run for the query, but also several specific versions meant to work with or around features of specific RDBMS products. For example:

```
test_query = db_query(
    """
        select * from blah
    """,
    postgres="""
        select * from pg_blah
    """,
    oracle="""
        select rownum, * from ora_blah
    """)
```

A `db_query` object is a callable object. If called on a connection, it will execute itself on that connection. Specifically:

```
test_query(conn, ...)
```

has the same effect as:

```
conn.execute(test_query, ...)
```

`__call__` (*self*, *\_conn* : *db\_connection* [, *<param\_name>=<param\_value>*, ... ]) → *db\_result*  
Execute this *db\_query* on the given *db\_connection* with parameters.

Note that the following methods are primarily of interest to driver implementors.

**`get_variant_sql`** (*accepted\_variants* : *str seq*) → *str*

Given a list of accepted variant tags, returns the most appropriate SQL for this query. Specifically, this returns the first variant SQL given in the query which is acceptable, or the default SQL if none is acceptable.

**`get_variant_qmark_params`** (*accepted\_variants* : *str seq*, *params* : *dict*) → *str*, *seq*

Like `get_variant_format_params`, but for the DB API 2.0 ‘format’ paramstyle (i.e. %s placeholders). This also escapes any percent signs originally present in the query.

**`get_variant_numeric_params`** (*accepted\_variants* : *str seq*, *params* : *dict*) → *str*, *seq*

Like `get_variant_format_params`, but for the DB API 2.0 ‘numeric’ paramstyle (i.e. :<n> placeholders).

**`get_variant_named_params`** (*accepted\_variants* : *str seq*, *params* : *dict*) → *str*, *dict*

Like `get_variant_format_params`, but for the DB API 2.0 ‘named’ paramstyle (i.e. :<name> placeholders). Note that this paramstyle is the native style required by the `net.sa.sql` API.

**`get_variant_format_params`** (*accepted\_variants* : *str seq*, *params* : *dict*) → *str*, *seq*

Converts the SQL and params of this query to a form appropriate for databases that use the DB API 2.0 ‘format’ paramstyle (i.e. %s placeholders). Given a list of accepted variants and a dict of params, this returns the appropriate SQL with param placeholders converted to ‘format’ style, and a list of params suitable for filling those placeholders.

**`get_variant_pyformat_params`** (*accepted\_variants* : *str seq*, *params* : *dict*) → *str*, *dict*

Like `get_variant_format_params`, but for the DB API 2.0 ‘pyformat’ paramstyle (i.e. %(<name>)s placeholders). This also escapes any percent signs originally present in the query.

## 3.6 Implementing a New Driver

In order to implement a new database driver, you should create a new module that implements a subclass of *db\_driver*, then calls `register_driver` with an instance of that subclass in order to register the new driver.

Your *db\_driver* subclass will, of course, return subclasses of *db\_connection* and *db\_result* specific to your database as well. It should never be necessary to subclass *db\_query*—that class is meant to be a database-neutral representation of a “compiled” query.

For most drivers, one of the `get_variant_...` methods of *db\_query* should provide the query in a form that the underlying database can easily digest.

**`class net.sa.sql.db_driver`**

A database driver, which holds the responsibility of deciding which database URLs it will attempt to open, and returning *db\_connection* objects when a connection is successfully opened.

**`can_handle`** (*uri\_scheme* : *str*) → *bool*

Returns True if this *db\_driver* believes it can handle this database URI scheme.

**`connect`** (*uri* : *str*, *user* : *str* or *None*, *password* : *str* or *None*) → *db\_connection*

Returns None if this *db\_driver* cannot handle this database URI, or a *db\_connection* subclass instance connected to the database if it can. The *user* and *password* parameters passed in via this call override any values from the URI.

**`net.sa.sql.register_driver`** (*driver* : *db\_driver*)

Registers a *db\_driver* database driver object with the `net.sa.sql` module. Driver modules generally register themselves, and this function is only of interest to driver writers.

`net.sa.sql.unregister_driver(driver : db_driver)`

Removes a `db_driver` database driver object from the set of drivers registered with the `net.sa.sql` module.

## 3.7 Experimental Connection Pooling

This version of `net.sa.sql` contains experimental support for connection pooling. Connections in a pool will be created before they're needed and kept available for re-use. Note that since this API is still in the early stages of development, it is very likely to change between versions of *net.sa-python*.

`net.sa.sql.db_create_pool(uri[, user : str, password : str], ...) → db_pool`

Given a database URI, an optional `user` and `password`, and additional parameters, creates a driver-specific connection pool. Returns a `db_pool` from which connections can be obtained.

If a user and password (or other parameter) is given in this call as well as in the URI, the values given in this call override the values given in the URI.

See `db_connect` for details on database URIs.

**class** `net.sa.sql.db_pool`

A pool of database connections for a single specific connection specification and pool configuration. See `db_create_pool`.

`get_driver()` → `db_driver`

Returns the `db_driver` used to open this connection.

`connect()` → `db_connection`

Returns a `db_connection` subclass instance from the pool, open on the database specified when the pool was created.

**class** `net.sa.sql.db_driver`

`create_pool(uri, user : str or None, password : str or None, ...) → db_pool`

Returns `None` if this `db_driver` does not support pooled connections or cannot handle this database URI, or a `db_pool` subclass instance which can be used to obtain connections from a pool. The `user` and `password` parameters and any other parameters passed in via this call override any values from the URI.

## 3.8 Why Not DB API 2.0?

If you have experience with Python database APIs, you may be wondering why we have chosen to implement a new API rather than simply using the standard [DB API 2.0](#).

In short, the problem is that the standard database API isn't really an API, but more a set of guidelines. For example, each database driver may use a different mechanism for providing query parameters. As another example, each API may also have different behaviors in the presence of threads.

Specifically, the `sqlite` module uses the 'pyformat' param style, which allows named parameters to queries which are passed as a dict, using Python-style formats. The `sqlite3` module, on the other hand, uses the 'qmark' param style, where `?` is used as a place-holder in queries, and the parameters are positional and passed in as a sequence.

We've done work to make sure that it's simple to implement `net.sa.sql`-style drivers over the top of [DB API 2.0](#)-style drivers. In fact, all of the currently deployed drivers are of this variety. The only work that has to be done for such a driver is to start with one of the existing drivers, determine which paramstyle is being used, do any protection against threading issues that might be necessary, and turn the connection URI into a form that the driver you're using can handle.

Once that's done, you still have the issue that different databases may require different SQL to operate—but that's a lot easier to handle than “some databases use named parameters and some use positional”. And, the variant system makes it easy to put different compatibility versions of the same query together.



---

# NETSA.UTIL.SHELL — ROBUST SHELL PIPELINES

## 4.1 Overview

The `netsa.util.shell` module provides a facility for securely and efficiently running UNIX command pipelines from Python. To avoid text substitution attacks, it does not actually use the UNIX shell to process commands. In addition, it runs commands directly in a way that allows easier clean-up in the case of errors.

The following standard Python library functions provide similar capabilities, but without either sufficient text substitution protections or sufficient error-checking and recovery mechanisms:

- The `os.system` function
- The `subprocess` module
- The `popen2` module

Here are some examples, in increasing complexity, of the use of the `run_parallel` and `run_collect` functions:

Run a single process and wait for it to complete:

```
# Shell: rm -rf /tmp/test
run_parallel("rm -rf /tmp/test")
```

Start two processes and wait for both to complete:

```
# Shell: rm -rf /tmp/test1 & rm -rf /tmp/test2 & wait
run_parallel("rm -rf /tmp/test_dir_1",
             "rm -rf /tmp/test_dir_2")
```

Store the output of a command into a file:

```
# Shell: echo test > /tmp/testout
run_parallel(["echo test", ">/tmp/testout"])
```

Read the input of a command from a file (and put the output into another file):

```
# Shell: cat < /tmp/test > /tmp/testout
run_parallel(["</tmp/test", "cat", ">/tmp/testout"])
```

Append the output of a command to a file:

```
# Shell: echo test >> /tmp/testout
run_parallel(["echo test", ">>/tmp/testout"])
```

Pipe the output of one command into another command (and put the output into a file):

```
# Shell: echo test | sed 's/e/f/' > /tmp/testout
run_parallel(["echo test", "sed 's/e/f/'", ">/tmp/testout"])
```

Run two pipelines in parallel and wait for both to complete:

```
# Shell:
#   echo test | sed 's/e/f/' > /tmp/testout &
#   cat /etc/passwd | cut -f1 -d'|' > /tmp/testout2 &
#   wait
run_parallel(["echo test", "sed 's/e/f/'", ">/tmp/testout"],
             ["cat /etc/passwd", "cut -f1 -d'|'", ">/tmp/testout2"])
```

Run a single pipeline and collect the output and error output in the variables *out* and *err*:

```
# Shell: foo='cat /etc/passwd | cut -f1 -d'|'
(foo, foo_err) = run_collect("cat /etc/passwd", "cut -f1 -d'|'")
```

The following examples are more complicated, and require the use of the long forms of `command` and `pipeline` specifications. (All of the examples above have used the short-hand forms.) You should read the documentation for `command` and `pipeline` to see how the long forms and short-hand forms are related.

Run a pipeline, collect standard output of the pipeline to one file, and append standard error from all of the commands to another file:

```
# Shell: ( gen-data | cut -f1 -d'|' > /tmp/testout ) 2>> /tmp/testlog
run_parallel(pipeline("gen-data", "cut -f1 -d'|'", ">/tmp/testout",
                      stderr="/tmp/testlog", stderr_append=True))
```

Run a pipeline, collect standard output of the pipeline to one file, and collect standard error from one command to another file:

```
# Shell: ( gen-data 2> /tmp/testlog ) | cut -f1 -d'|' > /tmp/testout
run_parallel([command("gen-data", stderr="/tmp/testlog"),
              "cut -f1 -d'|'", ">/tmp/testout"])
```

Run a pipeline, collect standard output of the pipeline to a file, and ignore the potentially non-zero exit status of the `gen-data` command:

```
# Shell: (gen-data | cut -f1 -d'|' > /tmp/testout) || true
run_parallel([command("gen-data", ignore_exit_status=True),
              "cut -f1 -d'|'", ">/tmp/testout"])
```

Use long pipelines to process data using multiple named pipes:

```
# Shell:
#   mkfifo /tmp/fifo1
#   mkfifo /tmp/fifo2
#   tee /tmp/fifo1 < /etc/passwd | cut -f1 -d'|' | sort > /tmp/out1 &
#   tee /tmp/fifo2 < /tmp/fifo1 | cut -f2 -d'|' | sort > /tmp/out2 &
#   cut -f3 -d'|' < /tmp/fifo2 | sort | uniq -c > /tmp/out3 &
#   wait
run_parallel("mkfifo /tmp/fifo1",
             "mkfifo /tmp/fifo2")
run_parallel(
    ["</etc/passwd", "tee /tmp/fifo1", "cut -f1 -d'|'", ">/tmp/out1"],
    ["</tmp/fifo1", "tee /tmp/fifo2", "cut -f2 -d'|'", ">/tmp/out2"],
    ["</tmp/fifo2", "cut -f3 -d'|'", "sort", "uniq -c", ">/tmp/out3"])
```

## 4.2 Exceptions

**exception** `net.sa.util.shell.PipelineException` (*message*, *exit\_statuses*)

This exception represents a failure to process a pipeline in either `run_parallel` or `run_collect`. It can be triggered by any of the commands being run by the function failing (either because the file was not found or because the command's exit status was unacceptable.) The message contains a summary of the status of all of the sub-commands at the time the problem was discovered, including `stderr` output for each sub-command if available.

## 4.3 Building Commands and Pipelines

`net.sa.util.shell.command` (<command spec>[, *stderr* : *str* or *file*, *stderr\_append*=*False*, *ignore\_exit\_status*=*False*, *ignore\_exit\_statuses* : *int seq*]) → *command*

Interprets the arguments as a “command specification”, and returns that specification as a value.

If there is only a single argument and it is a `command`, then a new command is returned with the options provided by this call. For example:

```
new_command = command(old_command, ignore_exit_status=True)
```

If there is only a single argument and it is a `str`, the string is parsed as if it were a simple shell command. (i.e. respecting single and double quotation marks, backslashes, etc.) For example:

```
new_command = command("ls /etc")
```

If there is only a single argument and it is a `list` or a `tuple`, interpret it as being the argument vector for the command (with the first argument being the command to be executed.) For example:

```
new_command = command(["ls", "/etc"])
```

If there are multiple arguments, each argument is taken as being one element of the argument vector, with the first bring the command to be executed. For example:

```
new_command = command("ls", "/etc")
```

The following keyword arguments may be given as options to a command specification:

***stderr*** Filename (*str*) or open *file* object of destination for *stderr*.

***stderr\_append*** *True* if *stderr* should be opened for append. Does nothing if *stderr* is already an open file.

***ignore\_exit\_status*** If *True*, then the exit status for this command is completely ignored.

***ignore\_exit\_statuses*** A list of numeric exit statuses that should not be considered errors when they are encountered.

In addition, these options may be “handed down” from the `pipeline` call, or from `run_parallel` or `run_collect`. If so, then options given locally to the command take precedence.

Example: Define a command spec using a single string:

```
c = command("ls -lR /tmp/foo")
```

Example: Define a command as the same as an old command with different options:

```
d = command(c, ignore_exit_status=True)
```

Example: Define a command using a list of strings:

```
e = command(["ls", "-lR", "/tmp/foo"])
```

Example: Define a command using individual string arguments:

```
f = command("ls", "-lR", "/tmp/foo")
```

### Short-hand Form:

In the `pipeline`, `run_parallel`, and `run_collect` functions, commands may be given in a short-hand form where convenient. The short-hand form of a command is a single string. Here are some examples:

```
"ls -lR"                =>  command(["ls", "-lR"])
"echo test test a b"    =>  command(["echo", "test", "test", "a", "b"])
"echo 'test test' a"    =>  command(["echo", "test test", "a"])
"'weird program'"       =>  command(["weird program"])
```

There is no way to associate options with a short-hand `command`. If you wish to redirect error output or ignore exit statuses, you will need to use the long form.

### Variable Expansion:

When commands are executed, variable expansion is performed. The expansions are provided by the argument `vars` to `run_parallel` or `run_collect`. Note that commands are split into arguments *before* this expansion occurs, which is a security measure. This means that no matter what whitespace or punctuation is in an expansion, it can't change the sense of the command. The down side of that is that on occasions when you would like to add multiple arguments to a command, you must construct the `command` using the list syntax.

Expansion variable references are placed using the [Python String formatting operations](#).

Here is an example substitution, showing how `%(target)s` becomes a single argument before the substitution occurs.

```
("ls -lR %(target)s", vars={'target': 'bl ah'}) =>
("ls", "-lR", "%(target)s", vars={'target': 'bl ah'}) =>
("ls", "-lR", 'bl ah')
```

If the value to be substituted implements the method `get_argument_list`, which takes no arguments and returns a list of strings, then those strings are included as multiple separate arguments. This is an expert technique for extending commands at call-time for use internal to APIs.

```
("ls -lR %(targets)s", vars={'targets': special_container}) =>
("ls", "-lR", "target1", "target2", ...)
```

### Functions as Commands:

In addition to executable programs, Python functions may also be used as commands. This is useful if you wish to do processing of data in a sub-process as part of a pipeline without needing to have auxilliary Python script files. However, this is an advanced technique and you should fully understand the subtleties before making use of it.

When a Python function is used as a command, the process will *fork* as normal in preparation for executing a new command. However, instead of *exec*-ing a new executable, the Python function is called. When the Python function completes (either successfully or unsuccessfully), the child process exits immediately.

If you intend to use this feature, be sure that you know how the lifecycles of various objects will behave when the Python interpreter is forked and two copies are running at once.

The `command` function is called with `vars` (as given to `run_parallel` or `run_collect`) as its first argument, and the remainder of `argv` from calling `command` as its remaining arguments.

`net.sa.util.shell.pipeline(<pipeline spec>[, stdin : str or file, stdout : str or file, stdout_append=False, ...]) → pipeline`

Interprets the arguments as a “pipeline specification”, and returns that specification as a value.

If there is only a single argument and it is a `pipeline`, then a new pipeline is returned with the options provided by this call. For example:

```
new_pipeline = pipeline(old_pipeline, stdout="/tmp/newfile")
```

If there is only a single argument and it is a `list` or a `tuple`, interpret it as being a list of commands and I/O redirection short-hands to run in the pipeline. For example:

```
new_pipeline = pipeline(["ls /etc", "tac"])
```

If there are multiple arguments, these arguments are treated as a list of commands and I/O redirection short-hands (as if they were passed as a single list.) For example:

```
new_pipeline = pipeline("ls /etc", "tac")
```

The following keyword arguments may be given as options to a pipeline specification:

**stdin** Filename (`str`) or open file object of source for stdin.

**stdout** Filename (`str`) or open file object of destination for stdout.

**stdout\_append** True if `stdout` should be opened for append. Does nothing if `stdout` is already an open file.

Because these options are so common, they may also be given in short-hand form. If the first command in the pipeline is a string starting with `<`, the remainder of the string is interpreted as a filename for stdin. If the last command in the pipeline is a string starting with `>` or `>>`, the remainder of the string is interpreted as a filename for stdout (and if `>>` was used, it is opened for append.)

In addition, any unrecognized keyword arguments will be provided as defaults for any `command` specifications used in this pipeline. (So, for example, if you give the `ignore_exit_status` option to `pipeline`, all of the commands in that pipeline will use the same value of `ignore_exit_status` unless they have their own overriding setting.)

Example: Define a pipeline using a list of commands:

```
a = pipeline(command("ls -lR /tmp/foo"),
              command("sort"),
              stdout="/tmp/testout")
```

Example: Define the same pipeline using the short-hand form of commands, and the shorthand method of setting stdout:

```
b = pipeline("ls -lR /tmp/foo",
              "sort",
              ">/tmp/testout")
```

Example: Define the same pipeline using a list instead of multiple arguments:

```
c = pipeline(["ls -lR /tmp/foo",
              "sort",
              ">/tmp/testout"])
```

Example: Define a new pipeline which is the same as an old pipeline but with different options:

```
d = pipeline(c, stdout="/tmp/newout")
```

**Short-hand Form:**

In the `run_parallel` command, pipelines may be given in a short-hand form where convenient. The short-hand form of a pipeline is a list of commands and I/O redirection short-hands. Here are some examples:

```
["ls /tmp/die", "xargs rm"] => pipeline(["ls /tmp/die", "xargs rm"])
["</tmp/testin", "sort", ">/tmp/testsort"] =>
    pipeline(["sort"], stdin="/tmp/testin", stdout="/tmp/testsort")
```

Note that although you can set `stdin`, `stdout`, and `stdout_append` using the short-hand form (by using the I/O redirection strings at the start and end of the list), you cannot set these options to open file objects, only to filenames. You also set other options to be passed down to the individual commands.

**Variable Expansion:**

As in `command`, pipelines have variable expansion. Most variable expansion happens inside the actual commands in the pipeline. However, variable expansion also occurs in filenames provided for the `stdin` and `stdout` options. For example:

```
pipeline("ls -lR", ">% (output_file)s")
pipeline("ls -lR", stdout="% (output_file)s")
```

## 4.4 Running Pipelines

`netsa.util.shell.run_parallel (<pipeline spec>, ...[, vars : dict, ...])`

Runs a series of commands (as specified by the arguments provided) by forking and establishing pipes between commands. Raises `PipelineException` and kills off all remaining subprocesses if any one command fails.

Each argument is passed to the `pipeline` function to create a new pipeline, which allows the short-hand form of pipelines (as `list` short-hands) to be used.

The following keyword arguments may be given as *options* to `run_parallel`:

**vars** A dictionary of variable substitutions to make in the `command` and `pipeline` specifications in this `run_parallel` call.

Additional keyword arguments will be passed down as default values to the `pipeline` and `command` specifications making up this `run_parallel` call.

The `run_parallel` function returns the list of exit codes of the processes in each pipeline as a list of lists. Each list corresponds to a pipeline, in the order in which they were passed into the function. Each element represents a process in the pipeline, in the order they were defined in the pipeline. If a process is not run (e.g., because a process preceding it in the pipeline fails), the exit status will be `None`.

Example: Run three `mkdirs` in parallel and fail if any of them fails:

```
run_parallel("mkdir a", "mkdir b", "mkdir c")
```

Example: Make a `fifo`, then afterwards, use it to do some work. (Try making a typo in here and watch it kill everything off instead of hanging forever.)

```
run_parallel("mkfifo test.fifo")
run_parallel(["cat /etc/passwd", "sort -r", "cut -f1 -d:", ">%(f)s"],
             ["cat %(f)s", "sed -e 's/a/b/g'", ">%(f2)s"],
             vars={'f': 'test.fifo', 'f2': 'test.'})
```

Example: run two pipelines in parallel, then investigate their processes' exit statuses:

```
exits = run_parallel(["ls -l", "grep ^d",
                    ["cat /etc/passwd", "sort -r", "cut -f1 -d:"]])
# If all complete successfully, exits will be:
# [[0, 0], [0, 0, 0]]
```

`netsa.util.shell.run_collect(<command spec>, ...[, vars : dict, ...]) → str, str`

Runs a series of commands specifying a single pipeline by forking and establishing pipes between commands. The output of the final command is collected and returned in the result. `stderr` across all commands is returned in the result. The final result is a tuple (`stdout`, `stderr`)

Raises `PipelineException` and kills off all remaining subprocesses if any one command fails.

The arguments are passed as arguments to a single call of the `pipeline` function to create a pipeline specification. That is: each argument is a `command` specification. Note that this is not the same as `run_parallel`, which interprets its arguments as multiple `pipeline` specifications.

You can also redirect `stderr` independently for each command if needed, allowing you to send some `stderr` to `/dev/null` or another destination instead of collecting it.

Example: Reverse sort the output of `ls -l` and store the output and error in the variables `a_stdout` and `a_stderr`:

```
# Reverse sort the output of ls -l
(a_stdout, a_stderr) = run_collect("ls -l", "sort -r")
```

Example: Do the same as the above, but run `ls -l` on a named directory instead of the current working directory:

```
# The same with a named directory
(b_stdout, b_stderr) = run_collect("ls -l %(dir)s", "sort -r",
                                  vars={'dir': 'some_directory'})
```

Example: The following *does not collect output*, but instead writes it to a file. If there were any error output, it would be returned in the variable `c_stderr`:

```
(empty_stdout, c_stderr) = run_collect("ls -l", "sort -r", ">test.out")
```

`netsa.util.shell.run_collect_files(<command spec>, ...[, vars : dict, ...]) → file, file`

Runs a series of commands like `run_collect`, but returns open file objects for `stdout` and `stderr` instead of strings.

Example: Iterate over the lines of `ls -l | sort -r` and print them out with line numbers:

```
(f_stdout, f_stderr) = run_collect_files("ls -l", "sort -r")
for (line_no, line) in enumerate(f_stdout):
    print ("%3d %s" % (line_no, line[:-1]))
```





---

# NETSA\_SILK — NETSA PYTHON/PYSILK SUPPORT

The `netsa_silk` module contains a shared API for working with common Internet data in both `netsa-python` and `PySiLK`. If `netsa-python` is installed but `PySiLK` is not, the less efficient but more portable pure-Python version of this functionality that is included in `netsa-python` is used. If `PySiLK` is installed, then the high-performance C version of this functionality that is included in `PySiLK` is used.

## 5.1 IPv6 Support

Depending on which version of the `netsa_silk` functionality is in use, IPv6 support may not be present or may be limited. The following functions allow determining what variety of IPv6 support is available.

`netsa_silk.has_IPv6Addr()` → bool

Returns True if the most basic form of IPv6 support—support for IPv6 addresses—is available. If it is not available, then this function returns False, `IPAddr` will raise `ValueError` when given an IPv6 address, and any call to `IPv6Addr` will raise `NotImplementedError`.

See also `ip_set.supports_ipv6`.

## 5.2 IP Addresses

An IP address is represented by either an `IPv4Addr` or an `IPv6Addr`. Both of these are subclasses of the generic `IPAddr` class.

**class** `netsa_silk.IPAddr` (*address* : str or `IPAddr`) → `IPv4Addr` or `IPv6Addr`

Converts the input into an IP address, either IPv4 or IPv6. Returns either an `IPv4Addr` or `IPv6Addr` object, depending on whether the given input is parsed as an IPv4 or an IPv6 address.

If IPv6 address support is not available (`has_IPv6Addr` returns False), then attempting to parse an IPv6 address will raise `ValueError`.

Examples:

```
>>> addr1 = IPAddr('192.160.1.1')
>>> addr2 = IPAddr('2001:db8::1428:57ab')
>>> addr3 = IPAddr('::ffff:12.34.56.78')
>>> addr4 = IPAddr(addr1)
>>> addr5 = IPAddr(addr2)
```

**class** `netsa_silk.IPv4Addr` (*address : int or str or IPAddr*)

Converts the input into a new IPv4 address. If the integer input is too large a value, if the string input is unparseable as an IPv4 address, or if the `IPAddr` input is not convertible to an IPv4 address, raises a `ValueError`.

`IPv4Addr` is a subclass of `IPAddr`.

Examples:

```
>>> addr1 = IPv4Addr('192.160.1.1')
>>> addr2 = IPv4Addr(IPAddr('::ffff.12.34.56.78'))
>>> addr3 = IPv4Addr(addr1)
>>> addr4 = IPv4Addr(0x10000000)
```

**class** `netsa_silk.IPv6Addr` (*address : int or str or IPAddr*)

Converts the input into a new IPv6 address. If the integer input is too large a value, or if the string input is unparseable as an IPv6 address, raises a `ValueError`. If the input is an `IPv4Addr`, the address is converted to IPv6 via IPv4-mapped address embedding. (“1.2.3.4” becomes “::ffff:1.2.3.4”).

If IPv6 address support is not available (`has_IPv6Addr` returns `False`), then calling `IPv6Addr` will raise `NotImplementedError`.

`IPv6Addr` is a subclass of `IPAddr`.

Examples:

```
>>> addr1 = IPv6Addr('2001:db8::1428:57ab')
>>> addr2 = IPv6Addr(IPAddr('192.168.1.1'))
>>> addr3 = IPv6Addr(addr1)
>>> addr4 = IPv6Addr(0x10000000000000000000000000000000)
```

## 5.2.1 Comparisons

Whenever an IPv4 address is compared to an IPv6 address, the IPv4 address is converted to IPv6 using IPv4-mapped address embedding. This means that `IPAddr('0.0.0.0')` equals `IPAddr('::ffff:0.0.0.0')`. You can distinguish IPv4 addresses from IPv6 address by using the `is_ipv6()` method.

Operation	Result
<code>'a == b'</code>	if <i>a</i> is equal to <i>b</i> , then <code>True</code> , else <code>False</code>
<code>'a != b'</code>	if <i>a</i> is equal to <i>b</i> , then <code>False</code> , else <code>True</code>
<code>'a &lt; b'</code>	if <i>a</i> 's integer representation is less than <i>b</i> 's, then <code>True</code> , else <code>False</code>
<code>'a &lt;= b'</code>	if <i>a</i> 's integer representation is less than or equal to <i>b</i> 's, then <code>True</code> , else <code>False</code>
<code>'a &gt;= b'</code>	if <i>a</i> 's integer representation is greater than or equal to <i>b</i> 's, then <code>True</code> , else <code>False</code>
<code>'a &gt; b'</code>	if <i>a</i> 's integer representation is greater than <i>b</i> 's, then <code>True</code> , else <code>False</code>

## 5.2.2 Conversions

The following operations and methods may be used to convert between IPv4 and IPv6 addresses and between IP addresses and other types.

Operation	Result	Notes
<code>'addr.is_ipv6()'</code>	if <i>addr</i> is an IPv6 address, then <code>True</code> , else <code>False</code>	
<code>'addr.to_ipv4()'</code>	the IPv4 equivalent of <i>addr</i> , or <code>None</code> if no such equivalent exists	(1)
<code>'addr.to_ipv6()'</code>	the IPv6 equivalent of <i>addr</i>	(2)
<code>'int(addr)'</code>	the integer representation of <i>addr</i>	(3)
<code>'str(addr)'</code>	the human-readable string representation of <i>addr</i>	(4)
<code>'addr.padded()'</code>	a zero-padded human-readable string representation of <i>addr</i>	(5)
<code>'addr.octets()'</code>	a tuple containing each octet of <i>addr</i> in network byte order as an unsigned integer	

Notes:

1. If the address is already an IPv4 address, does nothing. If the address is an IPv6 address using IPv4-mapped address embedding (e.g. `::ffff:1.2.3.4`), returns the equivalent IPv4 address. Otherwise, returns `None`.
2. If the address is already an IPv6 address, does nothing. If the address is an IPv4 address, returns the address converted to an IPv6 address using IPv4-mapped address embedding.
3. If the address is an IPv4 address, returns an unsigned 32-bit integer value. If the address is an IPv6 address, returns an unsigned 128-bit integer value.
4. The address is returned in its canonical form.
5. If the address is an IPv4 address, returns a string of the form `"xxx.xxx.xxx.xxx"`, where each field is one octet of the address as a zero-padded base-10 integer. If the address is an IPv6 address, returns a string of the form `"xxxx:xxxx:xxxx:xxxx:xxxx:xxxx:xxxx:xxxx"`, where each field is two octets of the address as a zero-padded base-16 integer.

## 5.2.3 Masking

Masking operations return a copy of the address with all non-masked bits set to zero.

Operation	Result	Notes
<code>'addr.mask(mask)'</code>	the address <i>addr</i> masked by the bits of the address <i>mask</i>	(1)
<code>'addr.mask_prefix(len)'</code>	the address <i>addr</i> masked to a length of <i>len</i> prefix bits	

Notes:

1. If *addr* is an IPv6 address but *mask* is an IPv4 address, *mask* is converted to IPv6 and then the mask is applied. If *addr* was not an IPv4-mapped embedded IPv6 address, the result may not be what was expected. Prefer `'addr.mask_prefix(len)'` for IPv6 addresses when possible.

## 5.3 IP Sets

While there are multiple different IP set implementations with different qualities, `netsa_silk` provides a standard API for these sets, and a standard mechanism for acquiring a value of some set when you don't need a specific implementation.

**class** `netsa_silk.ip_set([iterable])`

Returns a new IP set object from an unspecified implementation that matches the following API. The elements of the set must be IP addresses. The values in *iterable* may be `IPAddr` objects, strings parsable by `IPAddr`, `IPWildcard` objects, or strings parsable by `IPWildcard`.

In all of the following descriptions, *s*, *s1*, and *s2* must be `ip_set` objects, *addr* may be an `IPAddr` object or a string that is parsable as an `IPAddr`. *iterable* may be any iterable containing `IPAddr` objects, `IPWildcard` objects, and strings, as described above.

### 5.3.1 Cardinality

Since IP sets can grow very large, an additional method for querying the cardinality which supports large values is available.

Operation	Result	Notes
<code>'s.cardinality()'</code> <code>'len(s)'</code>	the cardinality of <i>s</i> the cardinality of <i>s</i>	(1)

Notes:

1. If the cardinality of *s* is exceptionally large, this may raise `OverflowError` due to limitations in Python. Using `'s.cardinality()'` is highly preferred for IP sets.

### 5.3.2 Membership

Operation	Result
<code>'addr in s'</code>	if <i>addr</i> is a member of <i>s</i> , then <code>True</code> , else <code>False</code>
<code>'addr not in s'</code>	if <i>addr</i> is a member of <i>s</i> , then <code>False</code> , else <code>True</code>

### 5.3.3 Comparison

Operation	Result
<code>'s1 == s2'</code>	if <i>s</i> has exactly the same elements as <i>s2</i> , then <code>True</code> , else <code>False</code>
<code>'s1 != s2'</code>	if <i>s1</i> does not have exactly the same elements as <i>s2</i> , then <code>True</code> , else <code>False</code>
<code>'s.isdisjoint(iterable)'</code>	if <i>s</i> has no elements in common with <i>iterable</i> , then <code>True</code> , else <code>False</code>
<code>'s.issubset(iterable)'</code>	if every element of <i>s</i> is also in <i>iterable</i> , then <code>True</code> , else <code>False</code>
<code>'s1 &lt;= s2'</code>	if every element of <i>s1</i> is also in <i>s2</i> , then <code>True</code> , else <code>False</code>
<code>'s1 &lt; s2'</code>	if <code>'s1 &lt; s2'</code> and <code>'s1 != s2'</code> , then <code>True</code> , else <code>False</code>
<code>'s.issuperset(iterable)'</code>	if every element of <i>iterable</i> is also in <i>s</i> , then <code>True</code> , else <code>False</code>
<code>'s1 &gt;= s2'</code>	if every element of <i>s2</i> is also in <i>s1</i> , then <code>True</code> , else <code>False</code>
<code>'s1 &gt; s2'</code>	if <code>'s1 &gt; s2'</code> and <code>'s1 != s2'</code> , then <code>True</code> , else <code>False</code>

### 5.3.4 Manipulation

The following operations return a new IP set with the desired changes, while leaving the original IP set unmodified.

Operation	Result
<code>'s.union(iterable, ...)'</code>	a set with all elements that are in <i>s</i> or any <i>iterable</i>
<code>'s1   s2'</code>	a set with all elements that are in <i>s1</i> or <i>s2</i>
<code>'s.intersection(iterable, ...)'</code>	a set with only elements that are in <i>s</i> and every <i>iterable</i>
<code>'s1 &amp; s2'</code>	a set with only elements that are in both <i>s1</i> and <i>s2</i>
<code>'s.difference(iterable, ...)'</code>	a set with all elements that are in <i>s</i> but in no <i>iterable</i>
<code>'s1 - s2'</code>	a set with all elements that are in <i>s1</i> but not in <i>s2</i>
<code>'s.symmetric_difference(iterable)'</code>	a set with all elements that are in <i>s</i> or <i>iterable</i> but not both
<code>'s1 ^ s2'</code>	a set with all elements that are in <i>s1</i> or <i>s2</i> but not both
<code>'s.copy()'</code>	a shallow copy of <i>s</i>

### 5.3.5 Modification

The following operations modify the target IP set in place.

Operation	Result	Notes
<code>'s.update(iterable, ...)</code>	updates <i>s</i> by adding all elements from each <i>iterable</i>	
<code>'s1  = s2'</code>	updates <i>s1</i> by adding all elements from <i>s2</i>	
<code>'s.intersection_update(iterable, ...)</code>	updates <i>s</i> by removing all elements that do not appear in every <i>iterable</i>	
<code>'s1 &amp;= s2'</code>	updates <i>s1</i> by removing all elements that do not appear in <i>s2</i>	
<code>'s.difference_update(iterable, ...)</code>	updates <i>s</i> by removing all elements that appear in any <i>iterable</i>	
<code>'s1 -= s2'</code>	updates <i>s1</i> by removing all elements that appear in <i>s2</i>	
<code>'s.symmetric_difference_update(iterable, ...)</code>	updates <i>s</i> , keeping only elements found in <i>s</i> or <i>iterable</i> but not both	
<code>'s1 ^= s2'</code>	updates <i>s1</i> , keeping only elements found in <i>s1</i> or <i>s2</i> but not both	
<code>'s.add(addr)'</code>	adds <i>addr</i> to <i>s</i>	
<code>'s.remove(addr)'</code>	removes <i>addr</i> from <i>s</i>	(1)
<code>'s.discard(addr)'</code>	removes <i>addr</i> from <i>s</i>	
<code>'s.pop()'</code>	removes and returns an arbitrary element from <i>s</i>	(2)
<code>'s.clear()'</code>	removes all elements from <i>s</i>	

Notes:

1. Raises `KeyError` if *addr* is not in *s*.
2. Raises `KeyError` if *s* was empty.

### 5.3.6 CIDR Block Iteration

`ip_set.cidr_iter()` → (IPAddr, int) iter

Returns an iterator over the CIDR blocks covered by this IP set. Each value in the iterator is a pair '(addr, prefix\_len)' where *addr* is the first IP address in the block, and *prefix\_len* is the prefix length of the block.

### 5.3.7 IPv6 Support

Some `ip_set` implementations do not provide IPv6 support. Such an implementation will raise an `exception.TypeError` on any attempt to add an `IPv6Addr` to the set. The following class method can be used to determine if a given implementation has IPv6 support:

**classmethod** `ip_set.supports_ipv6()` → bool

Returns True if this IP set implementation provides support for IPv6 addresses, or False otherwise.

## 5.4 IP Wildcards

An `IPWildcard` object represents the specification of a set of IP addresses using SiLK IP wildcard syntax. Not all sets of IP addresses can be represented by a single IP wildcard.

**class** `netesa_silk.IPWildcard(wildcard : str or IPWildcard)`

Returns a new `IPWildcard` object constructed from *wildcard*. The string *wildcard* may contain an IP address, an IP address with a CIDR prefix designation, an integer, an integer with a CIDR prefix designation, or a SiLK wildcard expression. In SiLK wildcard notation, a wildcard is represented as an IP address in canonical form with each octet (for IPv4 addresses) or octet pair (IPv6) holding a single value, a range of values, a comma-separated list of values and ranges, or the character 'x' to accept all values.

Examples:

```

>>> wild1 = IPWildcard('1.2.3.0/24')
>>> wild2 = IPWildcard('ff80::/16')
>>> wild3 = IPWildcard('1.2.3.4')
>>> wild4 = IPWildcard('::FFFF:0102:0304')
>>> wild5 = IPWildcard('16909056')
>>> wild6 = IPWildcard('16909056/24')
>>> wild7 = IPWildcard('1.2.3.x')
>>> wild8 = IPWildcard('1:2:3:4:5:6:7:x')
>>> wild9 = IPWildcard('1.2,3.4,5.6,7')
>>> wild10 = IPWildcard('1.2.3.0-255')
>>> wild11 = IPWildcard('::2-4')
>>> wild12 = IPWildcard('1-2:3-4:5-6:7-8:9-a:b-c:d-e:0-ffff')

```

### 5.4.1 Membership

The primary operation on `IPWildcard` objects is testing whether an address is contained in the set covered by the wildcard. Both `IPAddr` and `str` values may be tested for membership.

Operation	Results	Notes
'addr in wildcard'	<i>addr</i> matches <i>wildcard</i> , then True, else False	(1)
'addr not in wildcard'	<i>addr</i> matches <i>wildcard</i> , then False, else True	(1)

Notes:

1. *addr* may be an `IPAddr` or a string. Strings are automatically converted as with `'IPAddr(addr)'`

### 5.4.2 Other

The following additional operations are available on `IPWildcard` objects:

Operation	Result
'str(wildcard)'	the string that was used to construct <i>wildcard</i>
'wildcard.is_ipv6()'	if <i>wildcard</i> contains IPv6 addresses then True, else False

## 5.5 TCP Flags

A `TCPFlags` object represents the eight bits of flags from a TCP session.

**class** `netsa_silk.TCPFlags` (*value* : *int* or *str* or *TCPFlags*)

Returns a new `TCPFlags` object with the given flags set. If *value* is an integer, it is interpreted as the bitwise integer representation of the flags. If *value* is a string, it is interpreted as a case-insensitive sequence of letters indicating individual flags, and optional white space. The mapping is described below.

Each supported flag has an assigned letter in string representations, is available as an attribute on `TCPFlags` values, and is available as a `TCPFlags` constant in `netsa_silk`:

Flag	Meaning	Letter	<code>TCPFlags</code> attribute	<code>netsa_silk</code> constant
FIN	No more data from sender	F	'flags.fin'	TCP_FIN
SYN	Synchronize sequence numbers	S	'flags.syn'	TCP_SYN
RST	Reset the connection	R	'flags.rst'	TCP_RST
PSH	Push Function	P	'flags.psh'	TCP_PSH
ACK	Acknowledgment field significant	A	'flags.ack'	TCP_ACK
URG	Urgent Pointer field significant	U	'flags.urg'	TCP_URG
ECE	ECN-echo (RFC 3168)	E	'flags.ece'	TCP_ECE
CWR	Congestion window reduced (RFC 3168)	C	'flags.cwr'	TCP_CWR

### 5.5.1 Bit-Manipulation

The following bit-manipulation operations are available on `TCPFlags` objects:

Operation	Result
<code>'~flags'</code>	the bitwise inversion (not) of <i>flags</i>
<code>'flags1 &amp; flags2'</code>	the bitwise intersection (and) of <i>flags1</i> and <i>flags2</i>
<code>'flags1   flags2'</code>	the bitwise union (or) of <i>flags1</i> and <i>flags2</i>
<code>'flags1 ^ flags2'</code>	the bitwise exclusive disjunction (xor) of <i>flags1</i> and <i>flags2</i>

### 5.5.2 Conversions

The following operations and methods may be used to convert `TCPFlags` objects into other types.

Operation	Result
<code>'int(flags)'</code>	the integer value of the flags set in <i>flags</i>
<code>'str(flags)'</code>	a string representation of the flags set in <i>flags</i>
<code>'flags.padded()'</code>	a space-padded column aligned string representation of the flags set in <i>flags</i>
<code>'bool(flags)'</code>	if any flag is set in <i>flags</i> , then <code>True</code> , else <code>False</code>

### 5.5.3 Matching

`TCPFlags.matches(flagmask: str) → bool`

The `TCPFlags.matches` method may be used to determine if a `TCPFlags` value matches a given flag/mask specification. The specification is given as a string containing a set of flags that must be set, optionally followed by a slash and a set of flags that must be checked. (i.e. if “A” is not in the flag list but is in the mask, it must be false. If “U” is not in either, it may have any value.) For example, `'flags.matches('S/SA')` would return `True` if SYN was set and ACK was not set in *flags*.

Examples:

```
>>> flags = TCPFlags('SAU')
>>> flags.matches('S')
True
>>> flags.matches('SA/SA')
True
>>> flags.matches('S/SP')
True
>>> flags.matches('S/SA')
False
>>> flags.matches('SP/SP')
False
>>> flags.matches('A/SA')
False
```

## 5.6 Support for SiLK versions before 3.0

Although `netsa_silk` is only fully supported by SiLK as of version 3.0, some legacy support for older versions of PySiLK is available. Some specific things to watch for if you need to work with older SiLK versions:

1. In `IPv6Addr` the `octets` method is not available. Other conversion operations still work, however.
2. For `TCPFlags`, the lower-case flag name attributes (e.g. `'flags.syn'`) on the object are not available. To work around this, use the `TCPFlags.matches` method, or perform bitwise operations on the constants in

the module. (For example, instead of `if flags.fin: stuff...`, use `if flags & TCP_FIN: stuff ...`.)



---

# DATA MANIPULATION

## 6.1 `netsa.data.countries` — Country and Region Codes

Definitions of country and region names and codes as defined by ISO 3166-1 and the UN Statistics Division. The information in this module is current as of January 2010.

`netsa.data.countries.get_area_numeric` (*code : int or str*) → int

Given a country or region code as one of the following:

- String containing ISO 3166-1 alpha-2 code
- String containing ISO 3166-1 alpha-3 code
- String or integer containing ISO 3166-1 numeric code
- String containing DNS top-level domain alpha-2 code
- String or integer containing UN Statistics Division numeric region code

Returns the appropriate ISO 3166-1 or UN Statistics Division numeric code as an integer.

Note that some regions and other special items that are not defined by ISO 3166-1 or the UN Statistics Division are encoded as ISO 3166-1 user-assigned code elements.

Raises `KeyError` if the code is unrecognized.

`netsa.data.countries.get_area_name` (*code : int or str*) → str

Given a country or region code as a string or integer, returns the name for the country or region.

Raises `KeyError` if the country or region code is unrecognized.

`netsa.data.countries.get_area_tlds` (*code : int or str*) → str list

Given a country or region code as a string or integer, returns a list of zero or more DNS top-level domains for that country or region.

Raises `KeyError` if the country or region code is unrecognized.

`netsa.data.countries.get_country_numeric` (*code : int or str*) → int

Given a country code as a string or integer, returns the ISO 3166-1 numeric code for the country.

Raises `KeyError` if the country code is unrecognized.

`netsa.data.countries.get_country_name` (*code : int or str*) → str

Given a country code as a string or integer, returns the name for the country.

Raises `KeyError` if the country code is unrecognized.

`netlsa.data.countries.get_country_alpha2 (code : int or str) → str`

Given a country code as a string or integer, returns the ISO 3166-1 alpha-2 code for the country, or `None` if that is not possible.

Raises `KeyError` if the country code is unrecognized.

`netlsa.data.countries.get_country_alpha3 (code : int or str) → str`

Given a country code as a string or integer, returns the ISO 3166-1 alpha-3 code for the country, or `None` if that is not possible.

Raises `KeyError` if the country code is unrecognized.

`netlsa.data.countries.get_country_tlds (code : int or str) → str list`

Given a country code as a string or integer, returns a list of zero or more DNS top-level domains for that country.

Raises `KeyError` if the country code is unrecognized.

`netlsa.data.countries.iter_countries () → int iter`

Returns an iterator which yields all known ISO 3166-1 numeric country codes as integers, including user-assigned code elements in use.

`netlsa.data.countries.get_region_numeric (code : int or str) → int`

Given a UN Statistics Division region code as a string or integer, returns the code as an integer.

Raises `KeyError` if the region code is unrecognized.

`netlsa.data.countries.get_region_name (code : int or str) → str`

Given a region code as a string or integer, returns the name for the region.

Raises `KeyError` if the region code is unrecognized.

`netlsa.data.countries.get_region_tlds (code : int or str) → str list`

Given a region code as a string or integer, returns a list of zero or more DNS top-level domains for that region.

Raises `KeyError` if the region code is unrecognized.

`netlsa.data.countries.iter_regions () → int iter`

Returns an iterator which yields all top-level UN Statistics Division numeric region codes as integers. This includes Africa, the Americas, Asia, Europe, Oceania, and Other.

`netlsa.data.countries.iter_region_subregions (code : int or str) → int iter`

Given the code for a containing region, returns an iterator which yields all second-level UN Statistics Division numeric region codes as integers.

Raises `KeyError` if the region code is unrecognized.

`netlsa.data.countries.iter_region_countries (code : int or str) → int iter`

Given the code for a containing region, returns an iterator which yields as integers all ISO 3166-1 numeric country codes that are part of that region.

## 6.2 netlsa.data.format — Formatting Data for Output

The `netlsa.data.format` module contains functions useful for formatting data to be displayed in human-readable output.

### 6.2.1 Numbers

`netlsa.data.format.num_fixed (value : num[, units : str, dec_fig=2, thousands_sep : str]) → str`

Format *value* using a fixed number of figures after the decimal point. (e.g. “1234” is formatted as “1234.00”)

If *units* is provided, this unit of measurement is included in the output. *dec\_fig* specifies the number of figures after the decimal point.

If *thousands\_sep* is given, it is used to separate each group of three digits to the left of the decimal point.

Examples:

```
>>> num_fixed(1234, 'm')
'1234.00m'
>>> num_fixed(1234, 'm', dec_fig=4)
'1234.0000m'
>>> num_fixed(1234.5678, 'm', dec_fig=0)
'1235m'
>>> num_fixed(123456789, dec_fig=3, thousands_sep=",")
'123,456,789.000'
```

`netса.data.format.num_exponent (value : num[, units : str, sig_fig=3]) → str`

Format *value* using exponential notation. (i.e. “1234” becomes “1.23e+3” for three significant digits, or “1.234e+4” for four significant digits.) If *units* is provided, this unit of measurement is included in the output. *sig\_fig* is the number of significant figures to display in the formatted result.

Examples:

```
>>> num_exponent(1234, 'm')
'1.23e+3m'
>>> num_exponent(1234, 'm', sig_fig=4)
'1.234e+3m'
>>> num_exponent(1234.5678, 'm', sig_fig=6)
'1.23457e+3m'
>>> num_exponent(123456789, sig_fig=2)
'1.2e+8'
>>> num_exponent(123456, sig_fig=6)
'1.23456e+5'
```

`netса.data.format.num_prefix (value : num[, units : str, sig_fig=3, use_binary=False, thousands_sep : str]) → str`

Format *value* using SI prefix notation. (e.g. 1k is 1000) If *units* is provided, this unit of measurement is included in the output. *sig\_fig* is the number of significant figures to display in the formatted result.

If *use\_binary* is True, then SI binary prefixes are used (e.g. 1Ki is 1024). Note that there are no binary prefixes for negative exponents, so standard prefixes are always used for such cases.

For very large or very small values, exponential notation (e.g. “1e-30”) is used.

If *thousands\_sep* is given, it is used to separate each group of three digits to the left of the decimal point.

Examples:

```
>>> num_prefix(1024, 'b')
'1.02kb'
>>> num_prefix(1024, 'b', use_binary=True)
'1.00Kib'
>>> num_prefix(12345, 'b', sig_fig=2)
'12kb'
>>> num_prefix(12345, 'b', sig_fig=7)
'12345.00b'
>>> num_prefix(12345678901234567890, 'b')
'12.3Eb'
>>> num_prefix(12345678901234567890, 'b', sig_fig=7)
'12345.68Pb'
>>> num_prefix(1234567890123456789012345, 's')
'1.23e+24s'
```

```
>>> num_prefix(0.001, 's')
'1.00ms'
>>> num_prefix(0.001, 's', use_binary=True)
'1.00ms'
```

## 6.2.2 Dates and Times

Dates and times may be formatted to a variety of precisions. The formatting functions support the following precisions, except where otherwise noted: `DATETIME_YEAR`, `DATETIME_MONTH`, `DATETIME_DAY`, `DATETIME_HOUR`, `DATETIME_MINUTE`, `DATETIME_SECOND`, `DATETIME_MSEC`, and `DATETIME_USEC`.

`netlsa.data.format.datetime_silk(value : datetime[, precision=DATETIME_SECOND]) → str`  
Format *value* as a SiLK format date and time (YYYY/MM/DDTHH:MM:SS.SSS). Implicitly coerces the time into UTC.

*precision* is the amount of precision that should be included in the output.

For a more general way to round times, see `netlsa.data.times.bin_datetime`. See also `datetime_silk_hour` and `datetime_silk_day` for the most common ways to format incomplete dates in SiLK format.

Examples:

```
>>> t = netlsa.data.times.make_datetime("2010-02-03T04:05:06.007")
>>> datetime_silk(t)
'2010/02/03T04:05:06'
>>> datetime_silk(t, precision=DATETIME_YEAR)
'2010'
>>> datetime_silk(t, precision=DATETIME_MONTH)
'2010/02'
>>> datetime_silk(t, precision=DATETIME_DAY)
'2010/02/03'
>>> datetime_silk(t, precision=DATETIME_HOUR)
'2010/02/03T04'
>>> datetime_silk(t, precision=DATETIME_MINUTE)
'2010/02/03T04:05'
>>> datetime_silk(t, precision=DATETIME_SECOND)
'2010/02/03T04:05:06'
>>> datetime_silk(t, precision=DATETIME_MSEC)
'2010/02/03T04:05:06.007'
>>> datetime_silk(t, precision=DATETIME_USEC)
'2010/02/03T04:05:06.007000'
```

`netlsa.data.format.datetime_silk_hour(value : datetime) → str`  
Format *value* as a SiLK format datetime to the precision of an hour (YYYY/MM/DDTHH). Implicitly coerces time into UTC. This is shorthand for `datetime_silk(value, precision=DATETIME_HOUR)`.

Example:

```
>>> t = netlsa.data.times.make_datetime("2010-02-03T04:05:06.007")
>>> datetime_silk_hour(t)
'2010/02/03T04'
```

`netlsa.data.format.datetime_silk_day(v : datetime) → str`  
Format *value* as a SiLK format datetime to the precision of a day (YYYY/MM/DD). Implicitly coerces time into UTC. This is shorthand for `datetime_silk(value, precision=DATETIME_DAY)`.

Example:

```
>>> t = netsa.data.times.make_datetime("2010-02-03T04:05:06.007")
>>> datetime_silk_day(t)
'2010/02/03'
```

`netsa.data.format.datetime_iso(value : datetime[, precision=DATETIME_SECOND]) → str`  
 Format *value* as an ISO 8601 extended format date and time (YYYY/MM/DDTHH:MM:SS.SSSSSS[TZ]). Includes timezone offset unless the value has no timezone or the value's timezone is UTC.

*precision* is the amount of precision that should be included in the output.

For a more general way to round times, see `netsa.data.times.bin_datetime`. See also `datetime_silk_hour` and `datetime_silk_day` for the most common ways to format incomplete dates in SILK format.

Examples:

```
>>> t = netsa.data.times.make_datetime("2010-02-03T04:05:06.007008")
>>> datetime_iso(t)
'2010-02-03T04:05:06'
>>> datetime_iso(t, precision=DATETIME_YEAR)
'2010'
>>> datetime_iso(t, precision=DATETIME_MONTH)
'2010-02'
>>> datetime_iso(t, precision=DATETIME_DAY)
'2010-02-03'
>>> datetime_iso(t, precision=DATETIME_HOUR)
'2010-02-03T04'
>>> datetime_iso(t, precision=DATETIME_MINUTE)
'2010-02-03T04:05'
>>> datetime_iso(t, precision=DATETIME_SECOND)
'2010-02-03T04:05:06'
>>> datetime_iso(t, precision=DATETIME_MSEC)
'2010-02-03T04:05:06.007'
>>> datetime_iso(t, precision=DATETIME_USEC)
'2010-02-03T04:05:06.007008'
>>> t = netsa.data.times.make_datetime("2010-02-03T04:05:06.007008+09:10", utc_only=False)
>>> datetime_iso(t)
'2010-02-03T04:05:06+09:10'
```

`netsa.data.format.datetime_iso_day(value : datetime) → str`  
 Format *value* as an ISO 8601 extended format date to the precision of a day (YYYY-MM-DD[TZ]). Includes timezone offset unless the value has no timezone or the value's timezone is UTC. This is shorthand for `datetime_iso(value, precision=DATETIME_DAY)`.

Example:

```
>>> t = netsa.data.times.make_datetime("2010-02-03T04:05:06.007")
>>> datetime_iso_day(t)
'2010-02-03'
>>> t = netsa.data.times.make_datetime("2010-02-03T04:05:06.007+03:00", utc_only=False)
>>> datetime_iso_day(t)
'2010-02-03+03:00'
```

`netsa.data.format.datetime_iso_basic(value : datetime[, precision=DATETIME_SECOND]) → str`  
 Format *value* as an ISO 8601 basic (compact) format date and time (YYYYMMDDTHHMMSS.SSSSSS[TZ]). Includes timezone offset unless the value has no timezone or the value's timezone is UTC.

*precision* is the amount of precision that should be included in the output. Note that in accordance with the ISO 8601 specification, this format does not support the `DATETIME_MONTH` precision, because `YYYYMM` and

YYMMDD would be potentially ambiguous.

For a more general way to round times, see `netsa.data.times.bin_datetime`. See also `datetime_silk_hour` and `datetime_silk_day` for the most common ways to format incomplete dates in SiLK format.

Examples:

```
>>> t = netsa.data.times.make_datetime("2010-02-03T04:05:06.007008")
>>> datetime_iso_basic(t)
'20100203T040506'
>>> datetime_iso_basic(t, precision=DATETIME_YEAR)
'2010'
>>> datetime_iso_basic(t, precision=DATETIME_DAY)
'20100203'
>>> datetime_iso_basic(t, precision=DATETIME_HOUR)
'20100203T04'
>>> datetime_iso_basic(t, precision=DATETIME_MINUTE)
'20100203T0405'
>>> datetime_iso_basic(t, precision=DATETIME_SECOND)
'20100203T040506'
>>> datetime_iso_basic(t, precision=DATETIME_MSEC)
'20100203T040506.007'
>>> datetime_iso_basic(t, precision=DATETIME_USEC)
'20100203T040506.007008'
>>> t = netsa.data.times.make_datetime("2010-02-03T04:05:06.007008+09:10", utc_only=False)
>>> datetime_iso_basic(t)
'20100203T040506+0910'
```

`netsa.data.format.timedelta_iso` (*value : timedelta*)  $\rightarrow$  str

Format a `datetime.timedelta` object as a str in ISO 8601 duration format, minus ‘year’ and ‘month’ designators (P[n]DT[n]H[n]M[n]S). Fractional seconds will be represented using decimal notation in the seconds field.

Note that conversions between units are precise and do not take into account any calendrical context. In particular, a day is exactly 24\*3600 seconds, just like `datetime.timedelta` uses.

If you apply the resulting `timedelta` to a `datetime` and the interval happens to include something like leap seconds adjust your expectations accordingly.

Since `datetime.timedelta` has no internal representation of months or years, these units are never included in the result.

Examples:

```
>>> t1 = netsa.data.times.make_datetime("2010-02-03T04:05:06.007008")
>>> t2 = netsa.data.times.make_datetime("2010-02-04T05:06:07.008009")
>>> t3 = netsa.data.times.make_datetime("2010-02-03T04:06:06.008009")
>>> d1 = t2 - t1
>>> d2 = t3 - t1
>>> timedelta_iso(d1)
'P1DT1H1M1.001001S'
>>> timedelta_iso(d2)
'PT1M0.001001S'
```

## 6.3 `netsa.data.nice` — “Nice” Numbers for Chart Bounds

A set of functions to produce ranges of aesthetically-pleasing numbers that have the specified length and include the specified range. Functions are provided for producing nice numeric and time-based ranges.

```
net.sa.data.nice.nice_ticks(lo : num, hi : num[, ticks=5, inside=False]) → num, num, num iter
```

Find ‘nice’ places to put *ticks* tick marks for numeric data spanning from *lo* to *hi*. If *inside* is `True`, then the nice range will be contained within the input range. If *inside* is `False`, then the nice range will contain the input range. To find nice numbers for time data, use `nice_time_ticks`.

The result is a tuple containing the minimum value of the nice range, the maximum value of the nice range, and an iterator over the tick marks.

See also `nice_ticks_seq`.

```
net.sa.data.nice.nice_ticks_seq(lo : num, hi : num[, ticks=5, inside=False]) → num seq
```

A convenience wrapper of `nice_ticks` to return the nice range as a sequence.

```
net.sa.data.nice.nice_time_ticks(lo : datetime, hi : datetime[, ticks=5, inside=False,
                                                                    as_datetime=True]) → datetime/int, datetime/int, datetime/int
                                                                    iter
```

Find ‘nice’ places to put *ticks* tick marks for time data spanning from *lo* to *hi*. If *inside* is `True`, then the nice range will be contained within the input range. If *inside* is `False`, then the nice range will contain the input range. To find nice numbers for numerical data, use `nice_ticks`.

The result is a tuple containing the minimum value of the nice range, the maximum value of the nice range, and an iterator over the ticks marks. If *as\_datetime* is `True`, the result values will be `datetime.datetime` objects. Otherwise, the result values will be numbers of seconds since UNIX epoch. Regardless, the return value is expressed in UTC.

See also `nice_time_ticks_seq`.

```
net.sa.data.nice.nice_time_ticks_seq(lo : datetime, hi : datetime[, ticks=5, inside=False,
                                                                    as_datetime=True]) → datetime/int seq
```

A convenience wrapper of `nice_time_ticks` to return the nice range as a sequence.

## 6.4 net.sa.data.times — Time and Date Manipulation

```
net.sa.data.times.make_datetime(v : num or str or datetime or mxDateTime[, utc_only=True]) →
                                                                    datetime
```

Produces a `datetime.datetime` object from a number (seconds from UNIX epoch), a string (in ISO format, SiLK format, or old SiLK format), or a `datetime.datetime` object. If *utc\_only* is `True`, coerces the result to be in the UTC time zone.

If the `mxDateTime` library is installed, this function also accepts `mxDateTime` objects.

```
net.sa.data.times.bin_datetime(dt : timedelta, t : datetime[, z=UNIX_EPOCH : datetime]) →
                                                                    datetime
```

Returns a new `datetime.datetime` object which is the floor of the `datetime.datetime` *t* in a *dt*-sized bin. For example:

```
bin_datetime(timedelta(minutes=5), t)
```

will return the beginning of a five-minute bin containing the time *t*. If you have very specific requirements, you can replace the origin point for binning (*z*) with a time of your choice. By default, the UNIX epoch is used, which is appropriate for most uses.

```
net.sa.data.times.make_timedelta(v : timedelta or str) → timedelta
```

Produces a `datetime.timedelta` object from a string (in ISO 8601 duration format) or a `datetime.timedelta` object.

Since `datetime.timedelta` objects do not internally support units larger than ‘days’, ISO 8601 strings containing month or year designations are discouraged. If these units are encountered in the string, however, they converted to days using a precise formula. This is an exact conversion that does not take into account any

calendrical context. If you apply the result to a datetime, and the interval happens to include leapseconds, or if you expect to land on the same day of the month while adding ‘months’ or ‘years’, adjust your expectations accordingly.

Example of ISO 8601: ‘P1DT1H1M1S’

This translates as a period of ‘1 day’ with time offset of ‘1 hour, 1 minute, and 1 second’. Fields are optional, the ‘P’ is required, as is the ‘T’ if using any units smaller than a day. A zero-valued timedelta can be represented as ‘P0D’.

`netسا.ااتا.اااس.اااamoa_aamoa_aa (n : aamoa_aa, d : aamoa_aa) → int, aamoa_aa`

Given two `اااااا.اamoa_aa` objects, return the number of times the second one (denominator) fits into the first one (numerator), along with any remainder expressed as another `اamoa_aa`.

## 6.4.1 Date Snappers

**class** `netسا.ااتا.اااس.اااSnapper` (*size : aamoa\_aa* [, *epoch=UNIX\_EPOCH : اااااا* ])

Class for date bin manipulations

**اamoa\_aa\_aa** (*اamoa*) → bool

Tests whether or not the provided date is the beginning `اااااا.اamoa_aa` for the containing time bin.

See `اamoa_aamoa_aa` for more detail on acceptable formats for date descriptors.

**اamoa\_aa** (*اamoa*) → اااااا

Returns a `اااااا.اamoa_aa` object representing the beginning of the date bin containing the provided date (‘snapping’ the date into place)

See `اamoa_aamoa_aa` for more detail on acceptable formats for date descriptors.

**اamoa\_aa\_aa** (*اamoa*) → اااااا

Returns a `اااااا.اamoa_aa` object representing the last date of the date bin which contains the provided date.

See `اamoa_aamoa_aa` for more detail on acceptable formats for date descriptors.

**اamoa\_aa\_aa** (*اamoa : اamoa\_aa*) → seq

Given a list of datetimes, returns an iterator which produces tuples containing two datetime objects for each provided datetime. The first value of the tuple is the beginning of the date bin containing the datetime in question and the second value is the original datetime.

See `اamoa_aamoa_aa` for more detail on acceptable formats for datetime descriptors.

**اamoa\_aamoa\_aa** (*اamoa\_aamoa : seq*) → اamoa\_aa\_aa

Given a list of date ranges, return a list of date bins that intersect the union of the given date ranges. Each date range in the provided list can be a single datetime descriptor or a tuple representing a beginning and end datetime for the range.

See `اamoa_aamoa_aa` for more detail on acceptable formats for date descriptors.

**اamoa\_aa\_aa\_aa** (*اamoa\_aa : اamoa\_aa*) → seq

Given a list of datetimes, returns an iterator which produces tuples containing two datetime objects. The first value of the tuple is the beginning of the date bin and the second value is the original datetime. Both bins and datetimes will be repeated where necessary to fill in gaps not present in the original list of datetimes.

If, for example, the span between each successive datetime in the provided list is smaller than the defined bin size, the same bin will be returned for each datetime residing in that bin. (An example of this would be a bin size of 7 days and a list of daily dates – the same bin would be returned for each week of dates within that bin).



If, on the other hand, the span between successive datetimes in the provided list is larger than the defined bin size, each provided date will be repeatedly returned with each bin that exists between the provided datetimes. (An example of this would be a bin size of 7 days and a list of monthly dates – each monthly date would be successively returned with the 4 (or so) bins touched by that month)

See `make_datetime` for more detail on acceptable formats for date descriptors.

**next\_date\_bin**(*date*) → datetime

Returns a `datetime.datetime` object representing the beginning of the date bin following the date bin in which the given date resides.

See `make_datetime` for more detail on acceptable formats for date descriptors.

**prior\_date\_bin**(*date*) → datetime

Returns a `datetime.datetime` object representing the beginning of the date bin prior to the date bin in which the given date resides.

See `make_datetime` for more detail on acceptable formats for date descriptors.

**today\_bin**() → datetime

Returns a `datetime.datetime` object representing the beginning of the date bin containing the current date.

`netsa.data.times.dow_day_snapper`(*size* : int[, *dow*=0]) → DateSnapper

Given an integer size in days and an integer day-of-the-week, returns a `:class:DateSnapper` object anchored on the first occurring instance of that DOW after the given epoch, which defaults to the UNIX epoch. Monday is the 0th DOW. DOW values are modulo 7, so the 7th DOW would also represent Monday.



# MISCELLANEOUS FACILITIES

## 7.1 `netsa.files` — File and Path Manipulation

The routines in `netsa.files` are intended to help with manipulation of files in the filesystem as well as pathnames.

### 7.1.1 Paths

`netsa.files.relpath` (*p* : *str*, *base* : *str*) → *str*

Given a target path along with a reference path, return the relative path from the target to the reference.

This is a logical operation that does not consult the physical filesystem.

`os.path.relpath` in Python 2.6 adds something similar to this.

`netsa.files.is_relpath` (*p* : *str*, *base* : *str*) → *bool*

Given a target path along with a base reference path, return whether or not the base path subsumes the target path.

This is a logical operation that does not consult the physical filesystem.

### 7.1.2 Process ID Locks

This form of lock is generally used by services that wish to ensure that only one copy of the service is running at a time.

`netsa.files.acquire_pidfile_lock` (*path* : *str*) → *bool*

Attempts to acquire a locking PID file at the requested pathname *path*. If the file does not exist, creates it with the current process ID. If the file does exist but refers to a no-longer-existing process, attempts to replace it. If the file does exist and refers to a running process, does nothing.

Registers the lock to be released (as with `release_pidfile_lock`) when the currently running process exits, if it has not already been released.

Returns *True* if this process holds the lock, or *False* if another running process holds the lock.

`netsa.files.examine_pidfile_lock` (*path* : *str*) → (*int*, *bool*) or *None*

Examines the state of a locking PID file at *path* and returns it. If the file does not exist or is not in the proper format, returns *None*. If the file does exist and contains a PID, returns ‘(*pid*, *state*)’ where *pid* is the process ID holding the lock, and *state* is *True* if a running process has that process ID, or *False* otherwise.

`netsa.files.release_pidfile_lock` (*path* : *str*)

Attempts to release a locking PID file at the requested pathname *path*. If the file does not exist or contains a lock for a different PID, does nothing. Otherwise, unlinks the file.

### 7.1.3 Temporary Files

`netsa.files.get_temp_file_name([file_name : str]) → str`

Return the path to a file named *file\_name* in a temporary directory that will be cleaned up when the process exits. If *file\_name* is `None` then a new file name is created that has not been used before.

A temporary file at the named location will be cleaned up as long as the Python interpreter exits normally.

`netsa.files.get_temp_file([file_name : str, mode='r']) → file`

Returns an open `file` object for the file named *file\_name* in the temporary working directory, with the given *mode* (as described in `open`). If *file\_name* is `None` then a new file name is used that has not been used before.

The resulting temporary file will be cleaned up as long as the Python interpreter exits normally.

`netsa.files.get_temp_pipe_name([pipe_name : str]) → str`

Returns the path to a named pipe *file\_name* that has been created in a temporary directory that will be cleaned up when the process exits. If *file\_name* is `None` then a new file name is created that has not been used before.

## 7.2 `netsa.json` — JSON Wrapper Module

The `netsa.json` module provides a wrapper module for either the Python standard library `json` module, if it is available, or an included copy of the `simplejson` module, otherwise. Please see the standard library documentation for details.

## 7.3 `netsa.util.clitest` — Utility for testing CLI tools

### 7.3.1 Overview

The `netsa.util.clitest` module provides an API for driving automated tests of command-line applications. It doesn't do the work of a test framework; for that, use a framework library such as `unittest` or `functest`.

Enough of `netsa.util.clitest` has been implemented to fulfill a minimal set of requirements. Additional features will be added as necessary to support more complex testing.

This module is influenced by <http://pythonpaste.org/scripttest/>.

A usage example:

```
from clitest import *
env = Environment("./test-output")
# Run the command
result = env.run("ryscatterplot --help")
assert(result.success())
assert(result.stdout() == "whatever the help output is")
assert(result.stderr() == "")
# Clean up whatever detritus the command left
env.cleanup()
```

### 7.3.2 Exceptions

**exception** `netsa.util.clitest.TestingException`

Class of exceptions raised by the `clitest` module.

### 7.3.3 Classes

```
class netsa.util.clitest.Environment ([work_dir : str][, save_work_dir : bool][, debug : bool][, <env_name>=<env_val>, ... ])
```

An environment for running commands, including a set of environment variables and a working directory.

The *work\_dir* argument is the working directory in which the commands are run. If *work\_dir* is `None`, a directory will be made using `tempfile.mkdtemp` with default values.

*work\_dir* must not already exist or `run` will raise a `TestingException`. If *save\_work\_dir* is `False`, `cleanup` will remove this directory when it is called.

If *debug* is `True`, several debug messages will be emitted on `stderr`.

Any additional keyword arguments are used as environment variables.

```
get_env (env_name : str) → str
```

Returns the value of *env\_name* in the environment. If *env\_name* does not exist in the environment, this method returns `None`.

```
set_env (env_name : str, env_val : str)
```

Sets the value of *env\_name* in the environment to *env\_val*. *env\_val* must be a string.

```
del_env (env_name : str)
```

Removes *env\_name* from the environment. If *env\_name* doesn't exist, this method has no effect.

```
get_work_dir () → str
```

Returns the working directory in which the commands are run.

```
run (command : str[, <keyword>=<value>, ... ]) → Result
```

Runs a single command, capturing and returning result information. Keyword arguments are passed to `netsa.util.shell.run_parallel`. See the documentation of that function for an explanation of how such arguments are interpreted.

Returns a `Result` object representing the outcome.

```
cleanup ()
```

Cleans up resources left behind by the test process.

```
class netsa.util.clitest.Result (command, envvars, exit_codes, stdout, stderr, debug=False)
```

Contains information on a command's exit status and output.

```
success () → bool
```

Returns `True` if the exit code of the process was 0. This usually, but not always, indicates that the process ran successfully. Know Your Tool before relying on this function.

```
exited ([code]) → bool
```

Returns `True` if the process exited with the specified exit code. If the exit code is `None` or unsupplied, returns `True` if the process terminated normally (e.g., not on a signal).

```
exit_status () → int or None
```

Returns the exit status of the process, if the process exited normally (e.g., was not terminated on a signal). Otherwise, this function returns `None`.

```
signal () → int or None
```

Returns the signal on which the process terminated, if the process terminated on a signal. Otherwise, this function returns `None`.

```
signaled () → bool
```

Returns `True` if the process terminated on the specified signal. If the signal is `None` or unsupplied, returns `True` if the process terminated on a signal.

**format\_status()** → str

Returns a human-readable representation of how the process exited.

**get\_status()** → int

Returns the raw exit status of the process, as an integer formatted in the style of `os.wait`.

**get\_stdout()** → str

Returns the standard output of the process as a string.

**get\_stderr()** → str

Returns the standard error of the process as a string.

**get\_info()** → str

Returns the information contained in the result as a human-readable string.

## 7.4 `netsa.util.compat` — Python version compatibility code

The `netsa.util.compat` module provides some additional functionality introduced between Python 2.4 and the latest versions of Python. Obviously new syntax features cannot be supported, but certain utility functions in modules or built-in functions can be added on for the sake of sanity.

The list of provided features is currently small, but is likely to grow over time.

To use the compatibility features, simply import this module:

```
import netsa.util.compat
```

There is no need to import any specific symbols from the module—it will add the symbols directly where needed so that they may be imported as normal. Built-ins will also work wherever used.

The additional functions currently provided by this module are:

- `all`
- `any`
- `heapq.merge`
- `itertools.product`
- `os.path.relpath`

# INTERNAL FACILITIES

The following features are meant primarily for internal NetSA use. If you do make use of them, do so with the understanding that their interfaces are much more likely to change than other APIs in this document.

## 8.1 `net.sa.dist` — Common Installation Procedures

The `net.sa.dist` module is intended primarily for NetSA development team, to provide a common set of practices for generating documentation, running tests, and distributing and installing our software. If you are not a member of the NetSA dev team, you're likely to be better served by using the standard `distutils` module or the more powerful `setuptools` package.

### 8.1.1 Overview

`net.sa.dist` provides a set of extensions to `distutils`, along with an alternative API for specifying the contents of the distribution. The extensions provide for automatic generation of documentation using `Sphinx` (including PDF, HTML, and man pages), automatic generation of Python-readable version information from configuration metadata, and targets for running automated tests. The alternative API allows for specification of project metadata in a similar style to the metadata of `net.sa.script`.

### 8.1.2 New `setup.py` Commands

When running `setup.py` for a project that uses `net.sa.dist`, all of the normal commands (`build`, `install`, `sdist`, etc.) are available, along with the following:

**check** Run all automated tests for the project.

**check\_unit** Run automated unit tests for the project.

**check\_other** Run all other automated tests for the project.

**gen\_version** Generates any “version” files required by the project. See `add_version_file`.

**gen\_doc\_html** Generates an HTML manual for this project, placing it in `doc/html`. This manual is in the normal style for Python documentation. It is never automatically generated.

**gen\_doc\_tools\_web** Generates an HTML manual for this project and create a tarball out of it, placing the results in `dist/<name>-<version>-doc-web.tar.gz`. This manual is designed for use on the [NetSA Tools](#) website, and this command is used only to generate documentation to be deployed at that site.

**gen\_doc\_man** Creates generated man pages for the project, placing them under `doc/man`. This is automatically called when generating a source distribution, and the results will be included in the source tarball. When installing, `netса.dist` will attempt to run this command, but if it fails it will use a pre-generated copy if available.

**gen\_doc\_pdf** Generates a PDF manual for this project, placing it in the top level directory. This is automatically called when generating a source distribution, and the resulting manual will be included in the source tarball. This manual is not installed, however, only included with the distribution.

**netса.dist** Generates the standard items for a NetSA distribution. Namely, a source code release tarball in `dist/<name>-<version>.tar.gz` and a documentation tarball in `dist/<name>-<version>-doc-web.tar.gz` for deployment to the [NetSA Tools](#) website.

**netса\_src\_license** Modifies the source files in place to update their license section. The license in `LICENSE-xxx.txt` is used for a section that begins with `@xxx_HEADER_START@`. Backup files are created, just in case.

### 8.1.3 Project Layout

The files of the project are expected to be arranged mostly as follows:

Directory	Purpose
bin	Contains Python scripts to be installed as executables.
doc	Contains Sphinx documentation sources, including <code>conf.py</code> . See <a href="#">disable_documentation</a> , and <a href="#">Documentation Configuration</a> .
src	Contains Python source code and data files to be installed in Python packages. See <a href="#">add_package</a> , <a href="#">add_package_data</a> , and <a href="#">add_module_py</a> .

The following directories are used to contain outputs, and any extra files in them may be automatically destroyed by the `clean` command:

Directory	Purpose
build	A variety of intermediate products are stored here while building the project.
dist	Final products (tarballs) are stored here.
doc/html	HTML documentation generated with <code>gen_doc_html</code> will go here.
doc/man	manpage documentation generated with <code>gen_doc_man</code> will go here.

### 8.1.4 Documentation Configuration

To support simpler common configuration for documentation output, a convenience module is create during documentation generation. Under most circumstances, you should be able to use the following `conf.py` without changes:

```
from netса_sphinx_config import *

add_static_path("static_html")
```

Importing all symbols from `netса_sphinx_config` sets all of the settings to their normal values for a NetSA project, including producing output appropriate for use on the [NetSA Tools](#) website automatically.

If you need to make modifications, just replace or modify the values of standard [Sphinx build options](#).

The following function is provided to allow automatic generation of man pages. Any man page generated from the documentation will be automatically generated and included in source distributions, and automatically installed in the appropriate location.



```
netса_sphinx_config.add_man_page (source_file : str, man_page : str, description : str[, section = 1])
```

Add a new man page to be generated from the given *source\_file* (without extension). The name of the resulting file is *man\_page.\*section\**. Note that when Sphinx generates man pages, the top-level heading from the input file is ignored, and the title used is “*man\_page - description*” instead. This way, you can use the same input file to produce installed man pages and to produce man pages for display in the HTML output.

## 8.1.5 Project Configuration

The following functions are used to set metadata for the project:

```
netса.dist.set_name (project_name : str)
```

Sets the name of the project. This name is used as part of the name of produced tarballs and documentation files.

```
netса.dist.set_title (project_title : str)
```

Sets the title for this project. This should be the human-readable name of the project. It is displayed in most places as the project name.

```
netса.dist.set_description (project_description : str)
```

Sets the long-form description for this project. This should be a detailed explanation of the project’s purpose.

```
netса.dist.set_version (project_version : str)
```

Sets the version number for this project. The version number is used as part of the filename of distribution files, is included in the documentation, and may be written out as version files (see [add\\_version\\_file](#) and [netса.find\\_version](#)).

```
netса.dist.set_copyright (project_copyright : str)
```

Sets the copyright date for this project. This is used in documentation generation, and in the project metadata. For example:

```
dist.set_copyright ("2008-2011, Carnegie Mellon University")
```

```
netса.dist.set_license (project_license : str)
```

Sets the license type for this project, which defaults to ‘GPL’. This is used for project distribution metadata.

```
netса.dist.set_maintainer (project_maintainer : str)
```

Given a name and email address (i.e. ‘Harry Q. Bovik <bovik@sample.samp>’) sets the maintainer name and email address metadata for the project.

```
netса.dist.set_author (project_author : str)
```

Given a name and email address (i.e. ‘Harry Q. Bovik <bovik@sample.samp>’) sets the author name and email address metadata for the project.

```
netса.dist.set_url (project_url : str)
```

Sets the home page URL metadata for this project.

```
netса.dist.set_download_url (project_download_url : str)
```

Sets the download page URL metadata for this project.

Choosing which files should be installed where is accomplished with the following functions:

```
netса.dist.add_package (package_name)
```

Adds a Python package to be installed, by package name. For example:

```
dist.add_package ("netса.data")
```

The files for this Python package would be found under `src/netса/dist`. Remember that the package directory (and every directory leading up to it) must include an `__init__.py` file to be accepted as a Python package.

```
netlsa.dist.add_package_data(package_name, data_file_glob)
```

Adds one or more data files to be installed within a package. Each file or directory that *data\_file\_glob* expands to is included. The files and directories should be stored under *src/<package\_name>*, just like the Python source files for the package. For a method of installing files in different places, see [add\\_install\\_data](#).

```
netlsa.dist.add_module_py(module_name)
```

Adds a single module by module name. For example:

```
dist.add_module_py("netlsa.util.shell")
```

This file for this module would be found at *src/netlsa/util/shell.py*. Remember that the package directory (and every directory leading up to it) must include an *\_\_init\_\_.py* file, which will also be installed.

```
netlsa.dist.add_module_ext(module_name, module_sources, **kwargs)
```

Adds a single C extension module, given a module name, a list of sources, and optional keyword arguments as accepted by *distutils.core.Extension*. For example:

```
dist.add_module_ext('foo', ['foo.c', 'bar.c'])
```

```
netlsa.dist.add_script(script_name)
```

Adds a single script by script name. For example:

```
dist.add_script("helloworld")
```

The file for this script would be found at *bin/helloworld*. When installed, if the script has a *#!* line and contains *python*, it will automatically be modified to point to the version of Python being used to install this project.

```
netlsa.dist.add_install_data(install_path, data_file_name)
```

Adds an extra data file that should be installed when the project is installed. The *data\_file\_name* should be the path to the file from the top level of the project. *install\_path* should be the path to the installation directory from the install prefix. For example:

```
dist.add_install_data("share/doc/helloworld", "samples/helloworld.ini")
```

This would install the file found at *samples/helloworld.ini* as *.../share/doc/helloworld/helloworld.ini* under the installation prefix.

```
netlsa.dist.add_extra_files(extra_glob)
```

Given a glob string, adds files which match that glob to the distribution. This is used to add any extra files (README, etc.) that should be included in a source distribution but are not to be installed. If you do include in this list a file that's already to be installed, it will still be installed, and it will still be included in the distribution.

In order to avoid recording the version number in both the *setup.py* file and the source code, you can use the following functions to automatically generate a file with the version number in it, and read it back at run time:

```
netlsa.dist.add_version_file(version_file_name : str[, version_file_template = "%s\n"])
```

Adds a “version file” to the project, with the given path and template. By default the template is *"%s\n"*, which simply includes the version number and a newline. The path should be given relative to the base of the project. The version file will be generated automatically before any other processing is done.

See [netlsa.find\\_version](#) for a convenient method for retrieving the version number from this file for your package.

Example:

```
# in setup.py
dist.add_version_file("src/netlsa/VERSION")
```

```
# in netlsa/__init__.py
__version__ = netlsa.find_version(__file__)
```

`netса.find_version(source_file : str[, num_levels = 3]) → str`

Given the path to a Python source file, read in a version number from a file `VERSION` in the same directory, or look for a `setup.py` file in up to `num_levels` directories above the file and attempt to find the version there.

The following functions allow automated tests to be added and run from `setup.py`:

`netса.dist.add_unit_test_module(script_unit_test_module : str)`

Adds a unit test module to be run, by module name. For example:

```
dist.add_unit_test_module("netса.data.test")
```

The provided module is expected to be a `unittest` test module, and the tests will be run in a separate process from the process running `setup.py`. Running tests automatically builds the project, and places the build area in the `PYTHONPATH` for the test process.

`netса.dist.add_other_test_module(script_other_test_module : str)`

Adds a module to be run for testing, by module name. For example:

```
dist.add_other_test_module("crunchy.test")
```

The provided module is called in a subprocess like this:

```
python -m crunchy.test ${source_dir}
```

Where `${source_dir}` is the top level source directory of the project. Running tests automatically builds the project, and places the build area in the `PYTHONPATH` for the test process.

Finally, once the project is fully configured, use this function to handle command-line options and actually running the tasks:

`netса.dist.execute()`

Using the project as so far specified, parse command line options and does what is required to build, install, test, or make a distribution for the project. This should be called as the last thing in `setup.py`.



# DEPRECATED FEATURES

## 9.1 Deprecated functions from `netsa.files`

**exception** `netsa.files.DirLockBlock` (*message*)

*Deprecated* as of netsa-python v1.4. Use `acquire_pidfile_lock` instead.

Raised when an attempt to lock a directory is blocked for too long a time by another process holding the lock.

**class** `netsa.files.DirLocker` (*name* : *str* [, *dir* : *str*, *seize*=*False*, *debug*=*False* ])

*Deprecated* as of netsa-python v1.4. Use `netsa.files.acquire_pidfile_lock` instead.

Provides cheap cross-process locking via `mkdir/rmdir`. *name* is the token identifying the application group. *dir* optionally specifies the directory in which to establish the lock (defaults to `'/var/tmp'` or some sensible temp dir name). If *seize* is `True` and the requested lock appears to have been orphaned, a new lock is established and the old lock debris is removed.

This is not an infallible locking solution. This is advisory locking. It is possible to have an orphaned lock or a ghosted lock.

For simple scenarios such as avoiding long-running cron jobs from trampling over one another, it's probably sufficient.

See also `netsa.tools.service.check_pidfile`.

**exception** `netsa.files.LocalTmpDirError` (*message*)

*Deprecated* as of netsa-python v1.4. Use the *Temporary Files* functions instead.

Raised when an unrecoverable error occurs within `LocalTmpDir`.

**class** `netsa.files.LocalTmpDir` ([*dir* : *str*, *prefix*='tmp', *create*=*True*, *verbose*=*False*, *autodelete*=*True* ])

*Deprecated* as of netsa-python v1.4. Use the *Temporary Files* functions instead.

Provides ephemeral temporary directories, similar to `tempfile.NamedTemporaryFile`. The resulting directory and all of its contents will be unlinked when the object goes out of scope.

The parameter *prefix* is passed as the *prefix* parameter to `tempfile.NamedTemporaryFile` when temporary files are created within the temporary directory.

The parameter *create* controls whether the temporary directory is actually created. If you want to create the directory manually, you can use the method `assert_dir`.

The parameter *verbose* controls whether status messages (such as creation/deletion of files and dirs) are printed to `stderr`.

The parameter *autodelete* controls whether the temporary directory, and all its contents, are deleted once the `LocalTmpDir` object goes out of scope. (mostly useful for debugging)

**assert\_dir()**

*Deprecated* as of netsa-python v1.4. Use the [Temporary Files](#) functions instead.

Checks to see if this temp dir exists and creates it if not. Normally this happens during object creation.

**prefix()**

*Deprecated* as of netsa-python v1.4. Use the [Temporary Files](#) functions instead.

Returns the value of 'prefix' that is passed to `tempfile.NamedTemporaryFile`.

**tmp\_file()** → `NamedTemporaryFile`

*Deprecated* as of netsa-python v1.4. Use the [Temporary Files](#) functions instead.

Returns a `tempfile.NamedTemporaryFile` object within within this temp dir.

**tmp\_filename()** → `str`

*Deprecated* as of netsa-python v1.4. Use the [Temporary Files](#) functions instead.

Returns a new temporary filename within this temp dir.

**tmp\_pipe()** → `str`

*Deprecated* as of netsa-python v1.4. Use the [Temporary Files](#) functions instead.

Returns the filename of a new named pipe within within this temp dir.

## 9.2 Deprecated module `netsa.files.datefiles`

*Deprecated* as of netsa-python v1.4. No general replacement is yet available.

**exception** `netsa.files.datefiles.DateFieldParseError` (*message* : `str`)

*Deprecated* as of netsa-python v1.4.

Raised if a function is unable to parse a date within the provided filename.

`netsa.files.datefiles.date_from_file` (*file* : `str`) → `datetime`

*Deprecated* as of netsa-python v1.4.

Attempt to extract dates from a filename. The filename can be a full pathname or relative path. Dates are presumed to exist somewhere in the pathname. See [split\\_on\\_date](#) for more detail on how dates are parsed from filenames.

`netsa.files.datefiles.split_on_date` (*file* : `str`)

*Deprecated* as of netsa-python v1.4.

Given a string (presumably the pathname to a file) with a date in it, return the directory of the file and the date/non-date components of the file name as an array. This routine is pretty liberal about what constites a valid date format since it expects, by contract, a filename with a date string embedded within it.

For example, the input `"foo/bar-20090120:12:17:21.txt"` parses to:

```
('foo', ['bar-', 2009, None, 1, None, 20, ':', 12, ':', 17, ':', 21, '.txt'])
```

The following are all valid dates:

```
2008
200811
20081105
2008/11
2008/11/05
2008-11-05
2008.11.05
2008-11-05:07
```

```
2008-11-05:07:11
2008-11-05:07:11:00
```

Separators between year/month/day are ‘non digits’. This implies that directories within the path string can contribute to the date along with information in the filename itself. The following is valid:

```
‘/path/to/2008/11/05.txt’
```

The extraction is non-greedy: only the ‘last’ part that looks like a date is extracted. For example:

```
‘/path/to/2008/11/2008-11-05.txt’
```

extracts only ‘2008-11-05’ from the end of the string.

Separators between hour/minute/sec must be ‘:’

`netsa.files.datefiles.date_file_template` (*file* : *str* [, *wildcard*=‘x’]) → *str*  
*Deprecated* as of netsa-python v1.4.

Given a pathname *file*, returns the string with ‘x’ in place of date components. The replacement character can be overridden with the *wildcard* argument.

This is useful for determining what dated naming series are present in a shared directory.

`netsa.files.datefiles.sibling_date_file` (*file* : *str*, *date* : *datetime*) → *str*  
*Deprecated* as of netsa-python v1.4.

Given a filename *file*, along with a date, returns the analagous filename corresponding to that date.

`netsa.files.datefiles.datefile_walker` (*dir* : *str* [, *suffix* : *str*, *silent*=*False*, *snapper* : *DateSnapper*, *descend*=*True*, *reverse*=*False*]) → *iter*  
*Deprecated* as of netsa-python v1.4.

Returns an iterator based on the dated files that exist in directory *dir*. Each value returned by the iterator is a tuple of (date, [file1, file2, ...]), where each file matches the given date. See [split\\_on\\_date](#) for more detail on how dates are parsed from filenames.

If *descend* is *True*, the entire directory tree will be traversed. Otherwise, only the top-level directory is examined.

If *reverse* is *True*, the iterator returns entries for each date in descending order. Otherwise, entries are returns in ascending order.

If a `netsa.data.times.DateSnapper` *snapper* is provided, it will be used to enforce the alignment of dates, throwing a *ValueError* if a misaligned date is encountered.

`netsa.files.datefiles.latest_datefile` (*dir* : *str* [, *suffix* : *str*, *silent*=*False*, *snapper* : *DateSnapper*, *descend*=*True*]) → *tuple* or *None*  
*Deprecated* as of netsa-python v1.4.

Traverses the given directory and returns a single tuple (date, [file1, file2, ...]) where date is the latest date present and each file in the list contains that date. See [split\\_on\\_date](#) for more detail on how dates are parsed from filenames.

If *descend* is *True*, the entire directory tree is traversed recursively. Otherwise, only the top-level directory is examined.

If a `netsa.data.times.DateSnapper` *snapper* is provided, it will be used to enforce the alignment of dates, throwing a *ValueError* if a misaligned date is encountered.

`netsa.files.datefiles.date_snap_walker` (*dir* : *str*, *snapper* : *DateSnapper* [, *suffix* : *str*, *sparse*=*True*]) → *iter*  
*Deprecated* as of netsa-python v1.4.

Returns an iterator based on traversing the given directory. Each value returned by the iterator is a tuple `(date_bin, ((date, file), (date2, file2), ...))` where each filename contains a date which falls within the given date bin (as defined by *snapper*).

The beginning and ending dates for this sequence are determined by what files are present on the system. If *spare* is `True`, then only date bins which are actually occupied by files in the directory are emitted. Otherwise, a tuple is generated for each date between the smallest and largest dates present in the directory. See `netsa.data.times.DateSnapper` for more information on date bins.

If *suffix* is provided, all files not ending with the provided extension are ignored.

```
netsa.files.datefiles.tandem_datefile_walker(sources : str seq[, suffix : str, silent=True,
                                             snapper : DateSnapper, reverse=False])
                                             → iter
```

*Deprecated* as of netsa-python v1.4.

Returns an iterator based on traversing multiple directories given by *sources*. Each value returned by the iterator is a tuple `(date, (dir, file), (dir2, file2), ...)`, where the given directory contains the given file, which contains this date in its name. See `split_on_date` for more detail on how dates are parsed from filenames.

For example, given: `( '/dir/one', '/dir/two' )`, a returned tuple might look like:

```
(date, ('/dir/one', file_from_dir_one_containing_date_in_its_name),
      ('/dir/two', file_from_dir_two_containing_date_in_its_name),
      ('/dir/two', another_file_from_dir_two_with_date_in_its_name))
```

If *suffix* is provided, all files not ending with the provided extension are ignored.

If *reverse* is `True`, tuples are generated with dates in descending order. Otherwise, dates are generated in ascending order.

If a `netsa.data.times.DateSnapper` *snapper* is provided, it will be used to enforce the alignment of dates, throwing a `ValueError` if a misaligned date is encountered.

## 9.3 Deprecated functions from `netsa.script`

```
netsa.script.get_temp_dir_file_name([file_name : str]) → str
```

*Deprecated* as of netsa-python v1.4. Use `netsa.files.get_temp_file_name` instead.

Return the path to a file named *file\_name* in a temporary directory that will be cleaned up when the process exits. If *file\_name* is `None` then a new file name is created that has not been used before.

```
netsa.script.get_temp_dir_file([file_name : str, append=False]) → file
```

*Deprecated* as of netsa-python v1.4. Use `netsa.files.get_temp_file` instead.

Returns an open `file` object for the file named *file\_name* in the script's temporary working directory. If *append* is `True`, the file is opened for append. Otherwise, the file is opened for write. If *file\_name* is `None` then a new file name is used that has not been used before.

```
netsa.script.get_temp_dir_pipe_name([pipe_name : str]) → str
```

*Deprecated* as of netsa-python v1.4. Use `netsa.files.get_temp_pipe_name` instead.

Returns the path to a named pipe *file\_name* that has been created in a temporary directory that will be cleaned up when the process exits. If *file\_name* is `None` then a new file name is created that has not been used before.



## 9.4 Deprecated module `netsa.tools.service`

`netsa.tools.service.check_pidfile` (*path* : *str* [, *unlink=True* ]) → bool

*Deprecated* as of netsa-python v1.4. Use `netsa.files.acquire_pidfile_lock` instead.

Attempts to create a locking PID file at the requested pathname *path*. If the file does not exist, creates it with the current process ID, and sets up an `atexit` process to remove it. If the file does exist but refers to a no-longer-existing process, replaces it and does the above. If the file does exist and refers to a running process, does nothing.

Returns `True` if the PID file was created or replaced (which means we should continue processing), or `False` if the PID file was left in place (which means someone else is processing and we should exit.)



# CHANGES

## 10.1 Version 1.4 - 2011-09-30

- NOTE: Version 1.4 of NetSA Python is the last version that will support Python 2.4. Future major versions of NetSA Python will require Python 2.6 or greater.
- Added new `netsa_silk` module, which provides a bridge between `netsa-python` and `PySiLK`. When `PySiLK` is available, it uses `PySiLK` to provide a fast C implementation of IP address and related functionality. When `PySiLK` is not available, a pure Python implementation from `netsa-python` is used instead.
- Added new `netsa.script.golem` script automation framework, for building scripts to maintain large time-based data sets (among other things.)
- Added `regex_help` argument to `netsa.script.add_text_param` and `add_label_param`, to support providing a more useful error message when the regex doesn't match the input.
- Added `heapq.merge` and `os.path.relpath` to `netsa.util.compat`.
- Replaced all temporary file code with new functions in the `netsa.files` module, to avoid duplication of effort. Similar functions in other locations have been deprecated and are now implemented using this version.
- Replaced PID locking code with new functions in the `netsa.files` module, to avoid duplication of effort.
- Improved `netsa.data.nice`'s results, particularly for time ticks.
- Added a large number of new tests to improve compatibility testing for different versions of Python in the future.
- Moved the `netsa.script.get_temp_dir...` functions into `netsa.files`, with slight renaming. The old functions are still available but deprecated.
- Deprecated `netsa.files.DirLocker`, `netsa.files.LocalTmpDir`, and `netsa.tools.service`. See the new functions `netsa.files.acquire_pidfile_lock`, `examine_pidfile_lock`, and `release_pidfile_lock`.
- Deprecated `netsa.files.datefiles`.
- Added documentation for the `netsa.dist` module, even though it is primarily for internal use.
- Fixed bug involving `SIGPIPE` handling in `netsa.util.shell`.
- Fixed bug that prevented `netsa.logging` from importing correctly under Python 2.7.

## 10.2 Version 1.3 - 2011-03-28

- Switched to new common installation mechanism (based on `distutils`)
- Improved error handling in `netsa.util.script`

- Added new function `netsa.script.get_temp_dir_pipe_name()`
- Added `timedelta` support to `netsa.data.times`
- Added new `netsa.util.compat` to activity “compatibility features”

## **10.3 Version 1.2 - 2011-01-12**

- Added support for Oracle in `netsa.sql` via `cx_Oracle`
- Added support for multi-paragraph help text in `netsa.script`

## **10.4 Version 1.1 - 2010-10-04**

- Added experimental DB connection pooling to `netsa.sql`
- Made `netsa.script flow_params -help` work when site config file is missing.
- Added `netsa.util.shell.run_collect_files`
- Fixed a bug with `netsa.script.Flow_params.using`
- Fixed a bug involving `netsa.script` missing metadata causing crashes

## **10.5 Version 1.0 - 2010-09-14**

- Added `netsa.util.clitest` module to support CLI tool testing.
- Added PyGreSQL support to `netsa.sql`.
- Fixed a bug in `netsa.sql db_query` parsing code.
- Fixed a bug in `netsa.sql database URI` parsing code.
- Fixed bugs in `netsa.data.nice nice_time_ticks`.

## **10.6 Version 0.9 - 2010-01-19**

- First public release.

# LICENSES

## 11.1 License for netsa-python

Copyright 2008-2010 by Carnegie Mellon University

Use of the Network Situational Awareness Python support library and related source code is subject to the terms of the following licenses:

GNU Public License (GPL) Rights pursuant to Version 2, June 1991

Government Purpose License Rights (GPLR) pursuant to DFARS 252.227.7013

### NO WARRANTY

ANY INFORMATION, MATERIALS, SERVICES, INTELLECTUAL PROPERTY OR OTHER PROPERTY OR RIGHTS GRANTED OR PROVIDED BY CARNEGIE MELLON UNIVERSITY PURSUANT TO THIS LICENSE (HEREINAFTER THE “DELIVERABLES”) ARE ON AN “AS-IS” BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESS OR IMPLIED AS TO ANY MATTER INCLUDING, BUT NOT LIMITED TO, WARRANTY OF FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, INFORMATIONAL CONTENT, NONINFRINGEMENT, OR ERROR-FREE OPERATION. CARNEGIE MELLON UNIVERSITY SHALL NOT BE LIABLE FOR INDIRECT, SPECIAL OR CONSEQUENTIAL DAMAGES, SUCH AS LOSS OF PROFITS OR INABILITY TO USE SAID INTELLECTUAL PROPERTY, UNDER THIS LICENSE, REGARDLESS OF WHETHER SUCH PARTY WAS AWARE OF THE POSSIBILITY OF SUCH DAMAGES. LICENSEE AGREES THAT IT WILL NOT MAKE ANY WARRANTY ON BEHALF OF CARNEGIE MELLON UNIVERSITY, EXPRESS OR IMPLIED, TO ANY PERSON CONCERNING THE APPLICATION OF OR THE RESULTS TO BE OBTAINED WITH THE DELIVERABLES UNDER THIS LICENSE.

Licensee hereby agrees to defend, indemnify, and hold harmless Carnegie Mellon University, its trustees, officers, employees, and agents from all claims or demands made against them (and any related losses, expenses, or attorney’s fees) arising out of, or relating to Licensee’s and/or its sub licensees’ negligent use or willful misuse of or negligent conduct or willful misconduct regarding the Software, facilities, or other rights or assistance granted by Carnegie Mellon University under this License, including, but not limited to, any claims of product liability, personal injury, death, damage to property, or violation of any laws or regulations.

Carnegie Mellon University Software Engineering Institute authored documents are sponsored by the U.S. Department of Defense under Contract FA8721-05-C-0003. Carnegie Mellon University retains copyrights in all material produced under this contract. The U.S. Government retains a non-exclusive, royalty-free license to publish or reproduce these documents, or allow others to do so, for U.S. Government purposes only pursuant to the copyright license under the contract clause at 252.227.7013.

## 11.2 License for simplejson

Copyright (c) 2006 Bob Ippolito

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the “Software”), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED “AS IS”, WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

# INDEX

## Symbols

- (netsa\_silk.ip\_set operator), 64  
-<loop-select>  
    command line option, 14  
-data-complete  
    command line option, 14  
-data-inputs  
    command line option, 14  
-data-load  
    command line option, 14  
-data-outputs  
    command line option, 14  
-data-queue  
    command line option, 14  
-data-status  
    command line option, 14  
-first-date <date>  
    command line option, 14  
-intervals <count>  
    command line option, 14  
-last-date <date>  
    command line option, 14  
-output-dir <path>  
    command line option, 15  
-output-path <path>  
    command line option, 15  
-output-select <arg1[,arg2[,...]]>  
    command line option, 15  
-overwrite  
    command line option, 14  
-show-inputs  
    command line option, 15  
-show-outputs  
    command line option, 15  
-skip-incomplete  
    command line option, 14  
-= (netsa\_silk.ip\_set operator), 64  
== (netsa\_silk.IPAddr operator), 62  
== (netsa\_silk.ip\_set operator), 64  
& (netsa\_silk.TCPFlags operator), 66  
& (netsa\_silk.ip\_set operator), 64

&= (netsa\_silk.ip\_set operator), 64  
\_\_call\_\_() (netsa.sql.db\_query method), 49  
\_\_iter\_\_() (netsa.script.golem.GolemInputs method), 43  
\_\_iter\_\_() (netsa.script.golem.GolemOutputs method), 43  
\_\_iter\_\_() (netsa.script.golem.GolemProcess method), 45  
\_\_iter\_\_() (netsa.script.golem.GolemTags method), 43  
\_\_iter\_\_() (netsa.script.golem.GolemView method), 42  
\_\_iter\_\_() (netsa.sql.db\_result method), 49  
\_\_len\_\_() (netsa.script.golem.GolemInputs method), 43  
\_\_len\_\_() (netsa.script.golem.GolemOutputs method), 43  
\_\_len\_\_() (netsa.script.golem.GolemView method), 42  
^ (netsa\_silk.TCPFlags operator), 66  
^ (netsa\_silk.ip\_set operator), 64  
^= (netsa\_silk.ip\_set operator), 64  
~ (netsa\_silk.TCPFlags operator), 66  
| (netsa\_silk.TCPFlags operator), 66  
| (netsa\_silk.ip\_set operator), 64  
|= (netsa\_silk.ip\_set operator), 64  
> (netsa\_silk.IPAddr operator), 62  
> (netsa\_silk.ip\_set operator), 64  
>= (netsa\_silk.IPAddr operator), 62  
>= (netsa\_silk.ip\_set operator), 64  
< (netsa\_silk.IPAddr operator), 62  
< (netsa\_silk.ip\_set operator), 64  
<= (netsa\_silk.IPAddr operator), 62  
<= (netsa\_silk.ip\_set operator), 64

## A

ack (netsa\_silk.TCPFlags attribute), 66  
acquire\_pidfile\_lock() (in module netsa.files), 79  
add() (netsa\_silk.ip\_set method), 64  
add\_author() (in module netsa.script), 4  
add\_date\_param() (in module netsa.script), 6  
add\_dir\_param() (in module netsa.script), 7  
add\_extra\_files() (in module netsa.dist), 86  
add\_file\_param() (in module netsa.script), 6  
add\_flag\_param() (in module netsa.script), 7  
add\_float\_param() (in module netsa.script), 6  
add\_flow\_annotation() (in module netsa.script), 7  
add\_flow\_params() (in module netsa.script), 8  
add\_flow\_tag() (in module netsa.script.golem), 18

[add\\_golem\\_basic\\_params\(\)](#) (in module `netsa.script.golem`), 21  
[add\\_golem\\_input\(\)](#) (in module `netsa.script.golem`), 19  
[add\\_golem\\_param\(\)](#) (in module `netsa.script.golem`), 21  
[add\\_golem\\_params\(\)](#) (in module `netsa.script.golem`), 21  
[add\\_golem\\_query\\_params\(\)](#) (in module `netsa.script.golem`), 21  
[add\\_golem\\_repository\\_params\(\)](#) (in module `netsa.script.golem`), 21  
[add\\_golem\\_source\(\)](#) (in module `netsa.script.golem`), 21  
[add\\_input\\_template\(\)](#) (in module `netsa.script.golem`), 18  
[add\\_install\\_data\(\)](#) (in module `netsa.dist`), 86  
[add\\_int\\_param\(\)](#) (in module `netsa.script`), 5  
[add\\_label\\_param\(\)](#) (in module `netsa.script`), 6  
[add\\_loop\(\)](#) (in module `netsa.script.golem`), 17  
[add\\_module\\_ext\(\)](#) (in module `netsa.dist`), 86  
[add\\_module\\_py\(\)](#) (in module `netsa.dist`), 86  
[add\\_other\\_test\\_module\(\)](#) (in module `netsa.dist`), 87  
[add\\_output\\_dir\\_param\(\)](#) (in module `netsa.script`), 11  
[add\\_output\\_file\\_param\(\)](#) (in module `netsa.script`), 10  
[add\\_output\\_template\(\)](#) (in module `netsa.script.golem`), 18  
[add\\_package\(\)](#) (in module `netsa.dist`), 85  
[add\\_package\\_data\(\)](#) (in module `netsa.dist`), 85  
[add\\_path\\_param\(\)](#) (in module `netsa.script`), 7  
[add\\_query\\_handler\(\)](#) (in module `netsa.script.golem`), 20  
[add\\_script\(\)](#) (in module `netsa.dist`), 86  
[add\\_self\\_input\(\)](#) (in module `netsa.script.golem`), 20  
[add\\_sensor\\_loop\(\)](#) (in module `netsa.script.golem`), 17  
[add\\_tag\(\)](#) (in module `netsa.script.golem`), 17  
[add\\_text\\_param\(\)](#) (in module `netsa.script`), 5  
[add\\_unit\\_test\\_module\(\)](#) (in module `netsa.dist`), 87  
[add\\_version\\_file\(\)](#) (in module `netsa.dist`), 86  
[assert\\_dir\(\)](#) (`netsa.files.LocalTmpDir` method), 89

## B

[bin\\_count\(\)](#) (`netsa.script.golem.GolemView` method), 42  
[bin\\_dates\(\)](#) (`netsa.script.golem.GolemView` method), 41  
[bin\\_datetime\(\)](#) (in module `netsa.data.times`), 75  
[bool\(\)](#) (`netsa_silk.TCPFlags` operator), 67  
[by\\_bin\\_date\(\)](#) (`netsa.script.golem.GolemProcess` method), 44  
[by\\_bin\\_date\(\)](#) (`netsa.script.golem.GolemView` method), 41  
[by\\_day\(\)](#) (`netsa.script.Flow_params` method), 9  
[by\\_hour\(\)](#) (`netsa.script.Flow_params` method), 9  
[by\\_key\(\)](#) (`netsa.script.golem.GolemProcess` method), 45  
[by\\_key\(\)](#) (`netsa.script.golem.GolemView` method), 42  
[by\\_sensor\(\)](#) (`netsa.script.Flow_params` method), 9

## C

[can\\_handle\(\)](#) (`netsa.sql.db_driver` method), 50  
[cardinality\(\)](#) (`netsa_silk.ip_set` method), 63  
[check\\_pidfile\(\)](#) (in module `netsa.tools.service`), 93  
[cidr\\_iter\(\)](#) (`netsa_silk.ip_set` method), 65

[cleanup\(\)](#) (`netsa.util.clitest.Environment` method), 81  
[clear\(\)](#) (`netsa_silk.ip_set` method), 64  
[clone\(\)](#) (`netsa.sql.db_connection` method), 48  
 command line option
 

- [-<loop-select>](#), 14
- [-data-complete](#), 14
- [-data-inputs](#), 14
- [-data-load](#), 14
- [-data-outputs](#), 14
- [-data-queue](#), 14
- [-data-status](#), 14
- [-first-date <date>](#), 14
- [-intervals <count>](#), 14
- [-last-date <date>](#), 14
- [-output-dir <path>](#), 15
- [-output-path <path>](#), 15
- [-output-select <arg1\[,arg2\[,...\]\]>](#), 15
- [-overwrite](#), 14
- [-show-inputs](#), 15
- [-show-outputs](#), 15
- [-skip-incomplete](#), 14

[command\(\)](#) (in module `netsa.util.shell`), 55  
[commit\(\)](#) (`netsa.sql.db_connection` method), 48  
[connect\(\)](#) (`netsa.sql.db_driver` method), 50  
[connect\(\)](#) (`netsa.sql.db_pool` method), 51  
[copy\(\)](#) (`netsa_silk.ip_set` method), 64  
[create\\_pool\(\)](#) (`netsa.sql.db_driver` method), 51  
[current\\_view\(\)](#) (in module `netsa.script.golem`), 22  
[cwr](#) (`netsa_silk.TCPFlags` attribute), 66

## D

[date\\_aligned\(\)](#) (`netsa.data.times.DateSnapper` method), 76  
[date\\_bin\(\)](#) (`netsa.data.times.DateSnapper` method), 76  
[date\\_bin\\_end\(\)](#) (`netsa.data.times.DateSnapper` method), 76  
[date\\_binner\(\)](#) (`netsa.data.times.DateSnapper` method), 76  
[date\\_clumper\(\)](#) (`netsa.data.times.DateSnapper` method), 76  
[date\\_file\\_template\(\)](#) (in module `netsa.files.datefiles`), 91  
[date\\_from\\_file\(\)](#) (in module `netsa.files.datefiles`), 90  
[date\\_sequencer\(\)](#) (`netsa.data.times.DateSnapper` method), 76  
[date\\_snap\\_walker\(\)](#) (in module `netsa.files.datefiles`), 91  
[datefile\\_walker\(\)](#) (in module `netsa.files.datefiles`), 91  
[DateFileParseError](#), 90  
[DateSnapper](#) (class in `netsa.data.times`), 76  
[datetime\\_iso\(\)](#) (in module `netsa.data.format`), 73  
[datetime\\_iso\\_basic\(\)](#) (in module `netsa.data.format`), 73  
[datetime\\_iso\\_day\(\)](#) (in module `netsa.data.format`), 73  
[datetime\\_silk\(\)](#) (in module `netsa.data.format`), 72  
[datetime\\_silk\\_day\(\)](#) (in module `netsa.data.format`), 72  
[datetime\\_silk\\_hour\(\)](#) (in module `netsa.data.format`), 72  
[db\\_connect\(\)](#) (in module `netsa.sql`), 48



db\_connection (class in netsa.sql), 48  
 db\_create\_pool() (in module netsa.sql), 51  
 db\_driver (class in netsa.sql), 50, 51  
 db\_pool (class in netsa.sql), 51  
 db\_query (class in netsa.sql), 49  
 db\_result (class in netsa.sql), 49  
 del\_env() (netsa.util.clitest.Environment method), 81  
 difference() (netsa\_silk.ip\_set method), 64  
 difference\_update() (netsa\_silk.ip\_set method), 64  
 DirLockBlock, 89  
 DirLocker (class in netsa.files), 89  
 discard() (netsa\_silk.ip\_set method), 64  
 display\_message() (in module netsa.script), 7  
 divmod\_timedelta() (in module netsa.data.times), 76  
 dow\_day\_snapper() (in module netsa.data.times), 77

## E

ece (netsa\_silk.TCPFlags attribute), 66  
 end\_date (netsa.script.golem.GolemView attribute), 41  
 Environment (class in netsa.util.clitest), 81  
 environment variable  
     GOLEM\_HOME, 16, 21  
     GOLEM\_SOURCES, 21  
 examine\_pidfile\_lock() (in module netsa.files), 79  
 execute() (in module netsa.dist), 87  
 execute() (in module netsa.script), 12  
 execute() (in module netsa.script.golem), 22  
 execute() (netsa.sql.db\_connection method), 48  
 exit\_status() (netsa.util.clitest.Result method), 81  
 exited() (netsa.util.clitest.Result method), 81  
 expand() (netsa.script.golem.GolemInputs method), 43  
 expand() (netsa.script.golem.GolemOutputs method), 43

## F

fin (netsa\_silk.TCPFlags attribute), 66  
 find\_version() (in module netsa), 86  
 first\_bin (netsa.script.golem.GolemView attribute), 41  
 Flow\_params (class in netsa.script), 8  
 format\_status() (netsa.util.clitest.Result method), 81

## G

get\_area\_name() (in module netsa.data.countries), 69  
 get\_area\_numeric() (in module netsa.data.countries), 69  
 get\_area\_tlds() (in module netsa.data.countries), 69  
 get\_args() (in module netsa.script.golem), 23  
 get\_argument\_list() (netsa.script.Flow\_params method), 9  
 get\_class() (netsa.script.Flow\_params method), 10  
 get\_connection() (netsa.sql.db\_result method), 49  
 get\_country\_alpha2() (in module netsa.data.countries), 69  
 get\_country\_alpha3() (in module netsa.data.countries), 70  
 get\_country\_name() (in module netsa.data.countries), 69  
 get\_country\_numeric() (in module netsa.data.countries), 69

get\_country\_tlds() (in module netsa.data.countries), 70  
 get\_driver() (netsa.sql.db\_connection method), 48  
 get\_driver() (netsa.sql.db\_pool method), 51  
 get\_end\_date() (netsa.script.Flow\_params method), 10  
 get\_env() (netsa.util.clitest.Environment method), 81  
 get\_filenames() (netsa.script.Flow\_params method), 10  
 get\_flow\_params() (in module netsa.script), 8  
 get\_flowtypes() (netsa.script.Flow\_params method), 10  
 get\_home() (in module netsa.script.golem), 21  
 get\_info() (netsa.util.clitest.Result method), 82  
 get\_input\_pipe() (netsa.script.Flow\_params method), 10  
 get\_output\_dir\_file() (in module netsa.script), 12  
 get\_output\_dir\_file\_name() (in module netsa.script), 11  
 get\_output\_file() (in module netsa.script), 11  
 get\_output\_file\_name() (in module netsa.script), 11  
 get\_param() (in module netsa.script), 7  
 get\_params() (netsa.sql.db\_result method), 49  
 get\_query() (netsa.sql.db\_result method), 49  
 get\_region\_name() (in module netsa.data.countries), 70  
 get\_region\_numeric() (in module netsa.data.countries), 70  
 get\_region\_tlds() (in module netsa.data.countries), 70  
 get\_repository() (in module netsa.script.golem), 21  
 get\_script\_dir() (in module netsa.script.golem), 21  
 get\_script\_path() (in module netsa.script.golem), 21  
 get\_sensor\_group() (in module netsa.script.golem), 22  
 get\_sensors() (in module netsa.script.golem), 22  
 get\_sensors() (netsa.script.Flow\_params method), 10  
 get\_sensors\_by\_group() (in module netsa.script.golem), 23  
 get\_start\_date() (netsa.script.Flow\_params method), 10  
 get\_status() (netsa.util.clitest.Result method), 82  
 get\_stderr() (netsa.util.clitest.Result method), 82  
 get\_stdout() (netsa.util.clitest.Result method), 82  
 get\_temp\_dir\_file() (in module netsa.script), 92  
 get\_temp\_dir\_file\_name() (in module netsa.script), 92  
 get\_temp\_dir\_pipe\_name() (in module netsa.script), 92  
 get\_temp\_file() (in module netsa.files), 80  
 get\_temp\_file\_name() (in module netsa.files), 80  
 get\_temp\_pipe\_name() (in module netsa.files), 80  
 get\_type() (netsa.script.Flow\_params method), 10  
 get\_variant\_format\_params() (netsa.sql.db\_query method), 50  
 get\_variant\_named\_params() (netsa.sql.db\_query method), 50  
 get\_variant\_numeric\_params() (netsa.sql.db\_query method), 50  
 get\_variant\_pyformat\_params() (netsa.sql.db\_query method), 50  
 get\_variant\_qmark\_params() (netsa.sql.db\_query method), 50  
 get\_variant\_sql() (netsa.sql.db\_query method), 50  
 get\_variants() (netsa.sql.db\_connection method), 49  
 get\_verbosity() (in module netsa.script), 7

get\_work\_dir() (netsa.util.clitest.Environment method), 81  
get\_xargs() (netsa.script.Flow\_params method), 10  
golem (netsa.script.golem.GolemView attribute), 41  
GOLEM\_HOME, 16, 21  
GOLEM\_SOURCES, 21  
GolemArgs (class in netsa.script.golem), 44  
GolemInputs (class in netsa.script.golem), 43  
GolemOutputs (class in netsa.script.golem), 43  
GolemProcess (class in netsa.script.golem), 44  
GolemTags (class in netsa.script.golem), 42  
GolemView (class in netsa.script.golem), 41  
group\_by() (netsa.script.golem.GolemProcess method), 45  
group\_by() (netsa.script.golem.GolemView method), 42

## H

has\_IPv6Addr() (in module netsa\_silk), 61

## I

in (netsa\_silk.ip\_set operator), 64  
in (netsa\_silk.IPWildcard operator), 66  
inputs() (in module netsa.script.golem), 22  
inputs() (netsa.script.golem.GolemView method), 42  
int() (netsa\_silk.IPAddr operator), 62  
int() (netsa\_silk.TCPFlags operator), 67  
intersection() (netsa\_silk.ip\_set method), 64  
intersection\_update() (netsa\_silk.ip\_set method), 64  
ip\_set (class in netsa\_silk), 63  
IPAddr (class in netsa\_silk), 61  
IPv4Addr (class in netsa\_silk), 61  
IPv6Addr (class in netsa\_silk), 62  
IPWildcard (class in netsa\_silk), 65  
is\_complete() (in module netsa.script.golem), 22  
is\_complete() (netsa.script.golem.GolemProcess method), 44  
is\_files() (netsa.script.Flow\_params method), 10  
is\_ipv6() (netsa\_silk.IPAddr method), 62  
is\_ipv6() (netsa\_silk.IPWildcard method), 66  
is\_pull() (netsa.script.Flow\_params method), 10  
is\_relpsh() (in module netsa.files), 79  
isdisjoint() (netsa\_silk.ip\_set method), 64  
issubset() (netsa\_silk.ip\_set method), 64  
issuperset() (netsa\_silk.ip\_set method), 64  
iter\_countries() (in module netsa.data.countries), 70  
iter\_region\_countries() (in module netsa.data.countries), 70  
iter\_region\_subregions() (in module netsa.data.countries), 70  
iter\_regions() (in module netsa.data.countries), 70

## L

last\_bin (netsa.script.golem.GolemView attribute), 41  
latest\_datefile() (in module netsa.files.datefiles), 91

len() (netsa\_silk.ip\_set operator), 63  
LocalTmpDir (class in netsa.files), 89  
LocalTmpDirError, 89  
loop() (in module netsa.script.golem), 22  
loop() (netsa.script.golem.GolemView method), 42  
loop\_count() (netsa.script.golem.GolemView method), 42

## M

make\_datetime() (in module netsa.data.times), 75  
make\_timedelta() (in module netsa.data.times), 75  
mask() (netsa\_silk.IPAddr method), 63  
mask\_prefix() (netsa\_silk.IPAddr method), 63  
matches() (netsa\_silk.TCPFlags method), 67  
members() (netsa.script.golem.GolemInputs method), 43  
modify\_golem\_param() (in module netsa.script.golem), 21

## N

netsa.data.countries (module), 69  
netsa.data.format (module), 70  
netsa.data.nice (module), 74  
netsa.data.times (module), 75  
netsa.dist (module), 83  
netsa.files (module), 79  
netsa.files.datefiles (module), 90  
netsa.json (module), 80  
netsa.script (module), 1  
netsa.script.golem (module), 13  
netsa.sql (module), 47  
netsa.tools.service (module), 93  
netsa.util.clitest (module), 80  
netsa.util.compat (module), 82  
netsa.util.shell (module), 53  
netsa\_silk (module), 61  
netsa\_sphinx\_config.add\_man\_page() (in module netsa.dist), 84  
next\_date\_bin() (netsa.data.times.DateSnapper method), 77  
nice\_ticks() (in module netsa.data.nice), 74  
nice\_ticks\_seq() (in module netsa.data.nice), 75  
nice\_time\_ticks() (in module netsa.data.nice), 75  
nice\_time\_ticks\_seq() (in module netsa.data.nice), 75  
not in (netsa\_silk.ip\_set operator), 64  
not in (netsa\_silk.IPWildcard operator), 66  
num\_exponent() (in module netsa.data.format), 71  
num\_fixed() (in module netsa.data.format), 70  
num\_prefix() (in module netsa.data.format), 71

## O

octets() (netsa\_silk.IPAddr method), 62  
outputs() (in module netsa.script.golem), 22  
outputs() (netsa.script.golem.GolemView method), 42

## P

[padded\(\)](#) (netsa\_silk.IPAddr method), 62  
[padded\(\)](#) (netsa\_silk.TCPFlags method), 67  
[ParamError](#), 3  
[pipeline\(\)](#) (in module netsa.util.shell), 56  
[PipelineException](#), 55  
[pop\(\)](#) (netsa\_silk.ip\_set method), 64  
[prefix\(\)](#) (netsa.files.LocalTmpDir method), 90  
[prior\\_date\\_bin\(\)](#) (netsa.data.times.DateSnapper method), 77  
[process\(\)](#) (in module netsa.script.golem), 22  
[product\(\)](#) (netsa.script.golem.GolemProcess method), 44  
[product\(\)](#) (netsa.script.golem.GolemView method), 42  
[psh](#) (netsa\_silk.TCPFlags attribute), 66

## R

[register\\_driver\(\)](#) (in module netsa.sql), 50  
[release\\_pidfile\\_lock\(\)](#) (in module netsa.files), 79  
[relpath\(\)](#) (in module netsa.files), 79  
[remove\(\)](#) (netsa\_silk.ip\_set method), 64  
[Result](#) (class in netsa.util.clitest), 81  
[rollback\(\)](#) (netsa.sql.db\_connection method), 48  
[rst](#) (netsa\_silk.TCPFlags attribute), 66  
[run\(\)](#) (netsa.util.clitest.Environment method), 81  
[run\\_collect\(\)](#) (in module netsa.util.shell), 59  
[run\\_collect\\_files\(\)](#) (in module netsa.util.shell), 59  
[run\\_parallel\(\)](#) (in module netsa.util.shell), 58

## S

[script\\_view\(\)](#) (in module netsa.script.golem), 22  
[ScriptError](#), 4  
[set\\_author\(\)](#) (in module netsa.dist), 85  
[set\\_authors\(\)](#) (in module netsa.script), 4  
[set\\_contact\(\)](#) (in module netsa.script), 4  
[set\\_copyright\(\)](#) (in module netsa.dist), 85  
[set\\_default\\_home\(\)](#) (in module netsa.script.golem), 16  
[set\\_description\(\)](#) (in module netsa.dist), 85  
[set\\_description\(\)](#) (in module netsa.script), 4  
[set\\_download\\_url\(\)](#) (in module netsa.dist), 85  
[set\\_env\(\)](#) (netsa.util.clitest.Environment method), 81  
[set\\_interval\(\)](#) (in module netsa.script.golem), 16  
[set\\_lag\(\)](#) (in module netsa.script.golem), 16  
[set\\_license\(\)](#) (in module netsa.dist), 85  
[set\\_maintainer\(\)](#) (in module netsa.dist), 85  
[set\\_name\(\)](#) (in module netsa.dist), 85  
[set\\_name\(\)](#) (in module netsa.script.golem), 16  
[set\\_package\\_name\(\)](#) (in module netsa.script), 4  
[set\\_passive\\_mode\(\)](#) (in module netsa.script.golem), 17  
[set\\_realtime\(\)](#) (in module netsa.script.golem), 17  
[set\\_repository\(\)](#) (in module netsa.script.golem), 16  
[set\\_span\(\)](#) (in module netsa.script.golem), 16  
[set\\_suite\\_name\(\)](#) (in module netsa.script.golem), 16  
[set\\_title\(\)](#) (in module netsa.dist), 85

[set\\_title\(\)](#) (in module netsa.script), 4  
[set\\_tty\\_safe\(\)](#) (in module netsa.script.golem), 17  
[set\\_url\(\)](#) (in module netsa.dist), 85  
[set\\_version\(\)](#) (in module netsa.dist), 85  
[set\\_version\(\)](#) (in module netsa.script), 4  
[sibling\\_date\\_file\(\)](#) (in module netsa.files.datefiles), 91  
[signal\(\)](#) (netsa.util.clitest.Result method), 81  
[signaled\(\)](#) (netsa.util.clitest.Result method), 81  
[split\\_on\\_date\(\)](#) (in module netsa.files.datefiles), 90  
[sql\\_exception](#), 47  
[sql\\_invalid\\_uri\\_exception](#), 48  
[sql\\_no\\_driver\\_exception](#), 47  
[start\\_date](#) (netsa.script.golem.GolemView attribute), 41  
[status\(\)](#) (netsa.script.golem.GolemProcess method), 44  
[str\(\)](#) (netsa\_silk.IPAddr operator), 62  
[str\(\)](#) (netsa\_silk.IPWildcard operator), 66  
[str\(\)](#) (netsa\_silk.TCPFlags operator), 67  
[success\(\)](#) (netsa.util.clitest.Result method), 81  
[supports\\_ipv6\(\)](#) (netsa\_silk.ip\_set class method), 65  
[symmetric\\_difference\(\)](#) (netsa\_silk.ip\_set method), 64  
[symmetric\\_difference\\_update\(\)](#) (netsa\_silk.ip\_set method), 64  
[syn](#) (netsa\_silk.TCPFlags attribute), 66  
[sync\\_to\(\)](#) (netsa.script.golem.GolemView method), 42

## T

[tags\(\)](#) (netsa.script.golem.GolemTags method), 43  
[tandem\\_datefile\\_walker\(\)](#) (in module netsa.files.datefiles), 92  
[TCP\\_ACK](#) (netsa\_silk constant), 66  
[TCP\\_CWR](#) (netsa\_silk constant), 66  
[TCP\\_ECE](#) (netsa\_silk constant), 66  
[TCP\\_FIN](#) (netsa\_silk constant), 66  
[TCP\\_PSH](#) (netsa\_silk constant), 66  
[TCP\\_RST](#) (netsa\_silk constant), 66  
[TCP\\_SYN](#) (netsa\_silk constant), 66  
[TCP\\_URG](#) (netsa\_silk constant), 66  
[TCPFlags](#) (class in netsa\_silk), 66  
[TestingException](#), 80  
[timedelta\\_iso\(\)](#) (in module netsa.data.format), 74  
[tmp\\_file\(\)](#) (netsa.files.LocalTmpDir method), 90  
[tmp\\_filename\(\)](#) (netsa.files.LocalTmpDir method), 90  
[tmp\\_pipe\(\)](#) (netsa.files.LocalTmpDir method), 90  
[to\\_ipv4\(\)](#) (netsa\_silk.IPAddr method), 62  
[to\\_ipv6\(\)](#) (netsa\_silk.IPAddr method), 62  
[today\\_bin\(\)](#) (netsa.data.times.DateSnapper method), 77

## U

[union\(\)](#) (netsa\_silk.ip\_set method), 64  
[unregister\\_driver\(\)](#) (in module netsa.sql), 50  
[update\(\)](#) (netsa\_silk.ip\_set method), 64  
[urg](#) (netsa\_silk.TCPFlags attribute), 66  
[UserError](#), 3  
[using\(\)](#) (netsa.script.Flow\_params method), 10

`using()` (`netsa.script.golem.GolemProcess` method), [44](#)  
`using()` (`netsa.script.golem.GolemView` method), [41](#)