

---

# yelp-clog Documentation

*Release 2.16.0*

**Yelp Infra Team**

**Mar 28, 2018**



---

## Contents

---

<b>1</b>	<b>Python Logging Handler</b>	<b>3</b>
<b>2</b>	<b>Logging Operational and Mission Critical Data</b>	<b>5</b>
<b>3</b>	<b>Reading Scribe Logs</b>	<b>7</b>
<b>4</b>	<b>Contents</b>	<b>9</b>
4.1	Release Notes . . . . .	9
4.2	Configuration . . . . .	12
4.3	Global Logger . . . . .	12
4.4	Handlers . . . . .	13
4.5	Loggers . . . . .	14
4.6	Scribe Readers . . . . .	15
<b>5</b>	<b>Indices and tables</b>	<b>19</b>
	<b>Python Module Index</b>	<b>21</b>



*clog* is a package for handling log data. It can be used for the following:



# CHAPTER 1

---

## Python Logging Handler

---

*clog.handlers.ScribeHandler* can be used to send standard python logging to a scribe stream.



# CHAPTER 2

---

## Logging Operational and Mission Critical Data

---

`clog.loggers.ScribeLogger.log_line()` can be used to log mission critical, machine readable, and operational data to scribe. There is also a global `clog.log_line()` which has the same purpose but requires global configuration (see `clog.config`). Use of the global is discouraged.



# CHAPTER 3

---

## Reading Scribe Logs

---

*clog.readers* provides classes for reading scribe logs locally or from a server.



# CHAPTER 4

---

## Contents

---

### 4.1 Release Notes

#### 4.1.1 2.16.0

- Remove backend\_map logic from ScribeLogger and MonkLogger.
- When Monk is enabled, use ScribeMonkLogger instead of using MonkLogger directly.
- Implement preferred\_backend logic in ScribeMonkLogger. See preferred\_backend and preferred\_backend\_map in config.py.
- Upgrade Monk client library to 0.4.1.
- Implement memory buffering for MonkLogger.

#### 4.1.2 2.15.0

- Use scribify function to convert MonkLogger stream names, for backward compatibility with ScribeLogger.

#### 4.1.3 2.14.0

- Add close() method to MonkLogger

#### 4.1.4 2.13.0

- When creating a ScribeTailer instance, fall back to find\_tail\_host() if no tail service host is specified and the configuration is not set.
- Implement size limit for MonkLogger (5MB)

#### **4.1.5 2.12.1**

- Fix a compatibility issue where gzipped logs weren't decompressed correctly in python 3.

#### **4.1.6 2.12.0**

- Change reader configuration file path (from /etc/yelp\_clog.json to /nail/srv/configs/yelp\_clog.json)

#### **4.1.7 2.11.0**

- Prevent exceptions if Monk is enabled but not installed

#### **4.1.8 2.6.2**

- Remove references to yelp\_lib

#### **4.1.9 2.6.1**

- Add type checking to MockLogger log\_line function

#### **4.1.10 2.6.0**

- Make StreamTailer connection message more flexible
- Drop py33 and add py35

#### **4.1.11 2.5.2**

- Fix FileLogger when using python3

#### **4.1.12 2.5.1**

- Fix use of unicode

#### **4.1.13 2.5.0**

- Use six instead of future

#### **4.1.14 2.4.4**

- Fix ScribeHandler under Python 3.

#### **4.1.15 2.4.3**

- Fix tail/nc process leaks during testing.

## 4.1.16 2.4.2

- Improvements to testing.
- Fix `scribify` function under Python 3.

## 4.1.17 2.4.1

- Remove the `thriftpy` hard-pinning. Compatible now with `thriftpy` 0.1+ and 0.2+.

## 4.1.18 2.4.0

- Add `clog_enable_stdout_logging` config option to dump log lines to stdout. Defaults to false.

## 4.1.19 1.4.0

- Switched to `thriftpy`, a third-party Thrift implementation.
- Now compatible with Python 3.3+ and PyPy2.

## 4.1.20 1.3.0

- Add an enforcement on the size of scribe log lines. Log lines are expected to be less than 5MB; they are processed as normal. For log lines with size between 5MB and 50MB, warnings are issued by logging them additionally to the scribe log “`tmp_who_clog_large_line`”. Any line over 50 MB is treated as an error, and the exception `LogLineIsTooLongError` is raised.

## 4.1.21 1.2.0

- The `logging_timeout` argument of `loggers.ScribeLogger` is added. This allows scribe logging to timeout instead of being blocked if the scribe server is non-responding promptly; this benefits the applications prioritizing user experience over logging.

## 4.1.22 1.1.4

- The `_recv_compressed` method of `StreamTailer` and the `python-snappy` dependency have been removed. Now it's no more possible to automatically decompress a snappy-compressed stream.

## 4.1.23 1.0.0

This is a major release as it breaks backwards compatibility.

- Many imports were removed from the top level `clog` namespace. They are still available from the full module name (ex: `clog.loggers`, `clog.config`, etc)
- Configuration defaults have changed. The default is now to log to a local file in `/tmp`, instead of raising a `ValueError` if scribe is not configured. Note this only applies to `clog.log_line()`

## 4.2 Configuration

Configuration for `clog.global_state`. The following configuration settings are supported:

- scribe\_host** (string) hostname of a scribe service used by `clog.log_line()` to write logs to scribe
- scribe\_port** (int) port of the abovementioned scribe service
- scribe\_retry\_interval** number of seconds to wait before retrying a connection to scribe. Used by `clog.log_line()` (default 10)
- log\_dir** directory used to store files when `clog_enable_file_logging` is enabled. Defaults to the value of `$TMPDIR`, or `/tmp` if unset
- scribe\_disable** disable writing any logs to scribe (default True)
- scribe\_errors\_to\_syslog** flag to enable sending errors to syslog, otherwise to stderr (default False)
- scribe\_logging\_timeout** number of milliseconds to wait on socket connection to scribe server. This prevents from being blocked forever on, e.g., writing to scribe server. If a write times out, the delivery of the log line is not guaranteed and the status of the delivery is unknown; it can be either succeeded or failed. For backward compatibility, the old default blocking behavior is on when this parameter is left unset or set to 0; both the values `None` and `0` are treated as “infinite”.
- clog\_enable\_file\_logging** flag to enable logging to local files. (Default False)
- clog\_enable\_stdout\_logging** flag to enable logging to stdout. Each log line is prefixed with the stream name. (Default False)
- localS3** If True, will fetch s3 files directly rather than talking to a service.

`clog.config.configure(scribe_host, scribe_port, **kwargs)`  
Configure the `clog` package from arguments.

### Parameters

- **scribe\_host** – the scribe service hostname
- **scribe\_port** – the scribe service port
- **kwargs** – other configuration parameters

`clog.config.configure_from_dict(config_dict)`  
Configure the `clog` package from a dictionary.

### Parameters `config_dict` – a dict of config data

`clog.config.configure_from_object(config_obj)`  
Configure the `clog` package from an object (or module).

### Parameters `config_obj` – an object or module with config attributes

## 4.3 Global Logger

Log lines to scribe using the default global logger.

`exception clog.global_state.LoggingNotConfiguredError`

`clog.global_state.check_create_default_loggers()`  
Set up global loggers, if necessary.

```
clog.global_state.create_preferred_backend_map()
```

PyStaticConfig doesn't support having a map in the configuration, so we represent a map as a list, and we use this function to generate an actual python dictionary from it.

```
clog.global_state.log_line(stream, line)
```

Log a single line to the global logger(s). If the line contains any newline characters each line will be logged as a separate message. If this is a problem for your log you should encode your log messages.

#### Parameters

- **stream** – name of the scribe stream to send this log
- **line** – contents of the log message

```
clog.global_state.reset_default_loggers()
```

Destroy the global `clog` loggers. This must be done when forking to ensure that children do not share a desynchronized connection to Scribe

Any writes *after* this call will cause the loggers to be rebuilt, so this must be the last thing done before the fork or, better yet, the first thing after the fork.

## 4.4 Handlers

`logging.Handler` objects which can be used to send standard python logging to a scribe stream.

```
class clog.handlers.ScribeHandler(host, port, stream, retry_interval=0)
```

Handler for sending python standard logging messages to a scribe stream.

```
import clog.handlers, logging
log = logging.getLogger(name)
log.addHandler(clog.handlers.ScribeHandler('localhost', 3600, 'stream', retry_
interval=3))
```

#### Parameters

- **host** – hostname of scribe server
- **port** – port number of scribe server
- **stream** – name of the scribe stream logs will be sent to
- **retry\_interval** – default 0, number of seconds to wait between retries

```
emit(record)
```

Do whatever it takes to actually log the specified logging record.

This version is intended to be implemented by subclasses and so raises a `NotImplementedError`.

```
class clog.handlers.CLogHandler(stream, logger=None)
```

Deprecated since version 0.1.6.

**Warning:** Use `ScribeHandler` if you want to log to scribe, or a `logging.handlers.FileHandler` to log to a local file.

Handler for the standard logging library that logs to clog.

**emit** (*record*)

Do whatever it takes to actually log the specified logging record.

This version is intended to be implemented by subclasses and so raises a `NotImplementedError`.

```
clog.handlers.add_logger_to_scribe(logger, log_level=20, fmt='%(process)s\t%(asctime)s\t%(name)-12s\t%(levelname)-8s:\t\t\t%(message)s', clogger_object=None)
```

Sets up a logger to log to scribe.

By default, messages at the INFO level and higher will go to scribe.

Deprecated since version 0.1.6.

**Warning:** This function is deprecated in favor of using `clog.log_line()` or `ScribeHandler` directly.

**Parameters**

- **logger** – A `logging.Logger` instance
- **log\_level** – The level to log at
- **clogger\_object** – for use in testing

```
clog.handlers.get_scribed_logger(log_name, *args, **kwargs)
```

Get/create a logger and adds it to scribe.

Deprecated since version 0.1.6.

**Warning:** This function is deprecated in favor of using `clog.log_line()` directly.

**Parameters**

- **log\_name** – name of log to write to using `logging.getLogger`
- **kwargs (args,)** – passed to `add_logger_to_scribe`

**Returns** a `logging.Logger`

## 4.5 Loggers

Loggers which implement `log_line(stream, data)` which are used to send log data to a stream.

```
class clog.loggers.FileLogger
```

Implementation that logs to local files under a directory

```
class clog.loggers.GZipFileLogger(day=None)
```

Implementation of a logger that logs to local gzipped files.

```
exception clog.loggers.LogLineIsTooLongError
```

```
class clog.loggers.MockLogger
```

Mock implementation for testing

```
class clog.loggers.MonkLogger(client_id, host=None, port=None)
```

Wrapper around MonkProducer

```
exception clog.loggers.ScribeIsNotForkSafeError
class clog.loggers.ScribeLogger(host, port, retry_interval, report_status=None, logging_timeout=None)
Implementation that logs to a scribe server. If errors are encountered, drop lines and retry occasionally.
```

**Parameters**

- **host** – hostname of the scribe server
- **port** – port number of the scribe server
- **retry\_interval** – number of seconds to wait between retries
- **report\_status** – a function *report\_status(is\_error, msg)* which is called to print out errors and status messages. The first argument indicates whether what is being printed is an error or not, and the second argument is the actual message.
- **logging\_timeout** – milliseconds to time out scribe logging; “0” means blocking (no timeout)

**log\_line**(*stream*, *line*)

Log a single line. It should not include any newline characters. If the line size is over 50 MB, an exception raises and the line will be dropped. If the line size is over 5 MB, a message consisting origin stream information will be recorded at WHO\_CLOG\_LARGE\_LINE\_STREAM (in json format).

```
class clog.loggers.ScribeMonkLogger(config, scribe_logger, monk_logger, preferred_backend_map={})
```

The ScribeMonkLogger is a wrapper around both the ScribeLogger and the MonkLogger. The actual logger being used will depend on the preferred\_backend and preferred\_backend\_map.

**class** clog.loggers.StdoutLogger

Implementation that logs to stdout with stream name as a prefix.

**clog.loggers.get\_default\_reporter**(*use\_syslog=None*)

Returns the default reporter based on the value of the argument

**Parameters** **report\_to\_syslog** – Whether to use syslog or stderr. Defaults to the value of *config.scribe\_errors\_to\_syslog*

**clog.loggers.report\_to\_syslog**(*is\_error*, *msg*)

Report errors into Syslog.

Use Syslog with UDP in order not to be forced to write to a unix socket, useful if yelp-clog is run inside docker containers with Syslog running on the actual host.

Syslogs prepends ‘ident: ‘ to messages with an identifier. This function uses *clog* as identifier.

## 4.6 Scribe Readers

Classes which read log data from scribe.

```
class clog.readers.CLogStreamIterator(stream_reader, line_num=-1, current_chunk=None, chunk_line_num=-1)
```

Iterator used by ClogStreamReader for iterating over lines of chunks of a stream.

```
class clog.readers.CLogStreamReader(stream_name, stream_dir, date, fail_on_missing=False)
```

Make a stream reader for a day of *clog* entries

*clog* entries are stored by stream name and date and broken into separate chunks which may or may not be compressed with gzip or bzip or be plaintext.

For instance, the entries for a stream called ‘foo’ on New Years Day 2009 will be laid out in the file system like

```
STREAM_DIR/foo/foo-2009-01-01_00000.gz  
STREAM_DIR/foo/foo-2009-01-01_00001.gz  
...
```

Example usage:

```
reader = CLogStreamReader('stream_name', '/path/to/logs', date.today())  
for line in reader:  
    print line
```

#### Parameters

- **stream\_dir** – the stream directory like */storage/coraid5/scribe\_logs*
- **stream\_name** – the stream name like *biz\_views*
- **date** – the date of the logs
- **fail\_on\_missing** – Fail if there are no log files for the specified stream and date

#### chunk\_filenames()

Get an iterator for all the chunk filenames

```
class clog.readers.NetCLogStreamReader(bufsize=1024, host=None, port=None, au-  
tomagic_recovery=False, localS3=False)
```

Read logs from a scribe server

---

**Note:** This reader will stream logs from the source, it is not recommended for large logs. Use a mrjob instead.

---

Example usage:

```
stream_reader = NetCLogStreamReader()  
with stream_reader.read_date_range(  
    'ranger',  
    date(2010, 1, 1),  
    date(2010, 12, 31)  
) as reader:  
    for line in reader:  
        print line
```

Configuration:

This class can be configured either by passing a *host* and *port* to the constructor, or by using staticconf to with the following settings

**scribe\_net\_reader.host** hostname of the scribe server used to stream scribe logs

**scribe\_net\_reader.port** port of the scribe server used to stream scribe logs

#### Parameters

- **bufsize** – How many bytes to buffer internally
- **host** – The host to connect to (defaults to scribe\_net\_reader.host)

- **port** – The port to connect to (defaults to scribe\_net\_reader.port)
- **automagic\_recovery** – Whether to tail the stream, continuously retrying the connection (defaults to False)

```
list_streams()
    Get a context manager to use for reading list names

read_date_range(stream_name, start_date, end_date)
    Get a context manager to use for reading a stream for a date range
```

```
exception clog.readers.NoLogDataError

class clog.readers.StreamTailer(stream, host=None, port=None, bufsize=4096,
                                   automagic_recovery=True, add_newlines=True,
                                   raise_on_start=True, timeout=None, reconnect_callback=None,
                                   use_kafka=False, lines=None, protocol_opts=None)
```

Tail a Scribe stream from a tailing server

Example Usage:

```
tailer = StreamTailer('stream_name', host, port)
for line in tailer:
    print line
```

Configuration:

This class can be configured by passing a host and port to the constructor or by using staticconf with the following setting:

**scribe\_tail\_services** (list of dicts {‘host’: host, ‘port’: port}) list of host and port addresses of scribe endpoints for tailing logs in real time.

### Parameters

- **stream** (*string*) – the name of the string like ‘ranger’
- **host** (*string*) – the host name
- **port** – the port to connect to
- **bufsize** – the number of bytes to buffer
- **automagic\_recovery** (*bool*) – continue to retry connection forever
- **add\_newline** – add newlines to the items yielded in the iter
- **raise\_on\_start** (*bool*) – raise an error if you get a disconnect immediately after starting (otherwise, returns silently), Default True
- **timeout** (*int*) – connection timeout
- **reconnect\_callback** (*function*) – callback called when reconnecting
- **protocol\_opts** (*dict*) – optional protocol parameters

```
list_streams()
    Get a context manager to use for reading list names

exception clog.readers.StreamTailerSetupError(host, port, message)

clog.readers.construct_conn_msg(stream, lines=None, protocol_opts=None)
    Return a connection message
```

## Parameters

- **stream** – stream name
- **lines** – number of messages to consume
- **protocol\_opts** – optional arguments

`clog.readers.get_s3_info(hostname, stream_name=None)`

Returns (s3\_host, s3\_bucket(s), s3\_prefix)

If no stream name is provided (i.e. None), both normal and tmp buckets are returned as a dict.

# CHAPTER 5

---

## Indices and tables

---

- genindex
- modindex
- search



---

## Python Module Index

---

### C

`clog`, 3  
`clog.config`, 12  
`clog.global_state`, 12  
`clog.handlers`, 13  
`clog.loggers`, 14  
`clog.readers`, 15



---

## Index

---

### A

add\_logger\_to\_scribe() (in module clog.handlers), 14

### C

check\_create\_default\_loggers() (in module clog.global\_state), 12  
chunk\_filenames() (clog.readers.CLogStreamReader method), 16  
clog (module), 1  
clog.config (module), 12  
clog.global\_state (module), 12  
clog.handlers (module), 13  
clog.loggers (module), 14  
clog.readers (module), 15  
CLogHandler (class in clog.handlers), 13  
CLogStreamIterator (class in clog.readers), 15  
CLogStreamReader (class in clog.readers), 15  
configure() (in module clog.config), 12  
configure\_from\_dict() (in module clog.config), 12  
configure\_from\_object() (in module clog.config), 12  
construct\_conn\_msg() (in module clog.readers), 17  
create\_preferred\_backend\_map() (in module clog.global\_state), 12

### E

emit() (clog.handlers.CLogHandler method), 13  
emit() (clog.handlers.ScribeHandler method), 13

### F

FileLogger (class in clog.loggers), 14

### G

get\_default\_reporter() (in module clog.loggers), 15  
get\_s3\_info() (in module clog.readers), 18  
get\_scribed\_logger() (in module clog.handlers), 14  
GZipFileLogger (class in clog.loggers), 14

### L

list\_streams() (clog.readers.NetCLogStreamReader method), 17

list\_streams() (clog.readers.StreamTailer method), 17  
log\_line() (clog.loggers.ScribeLogger method), 15  
log\_line() (in module clog.global\_state), 13  
LoggingNotConfiguredError, 12  
LogLineIsTooLongError, 14

### M

MockLogger (class in clog.loggers), 14  
MonkLogger (class in clog.loggers), 14

### N

NetCLogStreamReader (class in clog.readers), 16  
NoLogDataError, 17

### R

read\_date\_range() (clog.readers.NetCLogStreamReader method), 17  
report\_to\_syslog() (in module clog.loggers), 15  
reset\_default\_loggers() (in module clog.global\_state), 13

### S

ScribeHandler (class in clog.handlers), 13  
ScribeIsNotForkSafeError, 14  
ScribeLogger (class in clog.loggers), 15  
ScribeMonkLogger (class in clog.loggers), 15  
StdoutLogger (class in clog.loggers), 15  
StreamTailer (class in clog.readers), 17  
StreamTailerSetupError, 17