
apachelogs

Release 0.4.0

John Thorvald Wodder II

2019 May 19

CONTENTS

| | | |
|----------|---|-----------|
| 1 | Parsing | 3 |
| 2 | Utilities | 5 |
| 2.1 | Log Format Constants | 5 |
| 3 | Exceptions | 7 |
| 4 | Supported Directives | 9 |
| 4.1 | Supported strftime Directives | 11 |
| 5 | Changelog | 13 |
| 5.1 | v0.4.0 (2019-05-19) | 13 |
| 5.2 | v0.3.0 (2019-05-12) | 13 |
| 5.3 | v0.2.0 (2019-05-09) | 13 |
| 5.4 | v0.1.0 (2019-05-06) | 14 |
| 6 | Installation | 15 |
| 7 | Examples | 17 |
| 8 | Indices and tables | 19 |
| | Python Module Index | 21 |

[GitHub](#) | [PyPI](#) | [Documentation](#) | [Issues](#) | *Changelog*

PARSING

```
class apachelogs.LogParser(format, encoding='iso-8859-1', errors=None)
```

A class for parsing Apache access log entries in a given log format. Instantiate with a log format string, and then use the `parse()` and/or `parse_lines()` methods to parse log entries in that format.

Parameters

- **format** (`str`) – an Apache log format
- **encoding** (`str`) – The encoding to use for decoding certain strings in log entries (see [Supported Directives](#)); defaults to 'iso-8859-1'. Set to 'bytes' to cause the strings to be returned as `bytes` values instead of `str`.
- **errors** (`str`) – the error handling scheme to use when decoding; defaults to 'strict'

Raises

- **InvalidDirectiveError** – if an invalid directive occurs in `format`
- **UnknownDirectiveError** – if an unknown directive occurs in `format`

`parse(entry)`

Parse an access log entry according to the log format and return a `LogEntry` object.

Parameters `entry` (`str`) – an access log entry to parse

Return type `LogEntry`

Raises `InvalidEntryError` – if `entry` does not match the log format

`parse_lines(entries, ignore_invalid=False)`

Parse the elements in an iterable of access log entries (e.g., an open text file handle) and return a generator of `LogEntries`. If `ignore_invalid` is `True`, any entries that do not match the log format will be silently discarded; otherwise, such an entry will cause an `InvalidEntryError` to be raised.

Parameters

- **entries** – an iterable of `str`
- **ignore_invalid** (`bool`) – whether to silently discard entries that do not match the log format

Return type `LogEntry` generator

Raises `InvalidEntryError` – if an element of `entries` does not match the log format and `ignore_invalid` is `False`

```
class apachelogs.LogEntry
```

A parsed Apache access log entry. The value associated with each directive in the log format is stored as an attribute on the `LogEntry` object; for example, if the log format contains a %s directive, the `LogEntry` for a

parsed entry will have a `status` attribute containing the status value from the entry as an `int`. See [Supported Directives](#) for the attribute names & types of each directive supported by this library.

If the log format contains two or more directives that are stored in the same attribute (e.g., `%D` and `%{us}T`), the given attribute will contain the first non-`None` directive value.

The values of date & time directives are stored in a `request_time_fields`: `dict` attribute. If this `dict` contains enough information to assemble a complete (possibly naïve) `datetime.datetime`, then the `LogEntry` will have a `request_time` attribute equal to that `datetime.datetime`.

directives = None

New in version 0.3.0.

A `dict` mapping individual log format directives (e.g., `"%h"` or `"%<s"`) to their corresponding values from the log entry. `%{*}t` directives with multiple subdirectives (e.g., `%{Y-%m-%d}t`) are broken up into one entry per subdirective (For `%{Y-%m-%d}t`, this would become the three keys `"%{Y}t"`, `"%{m}t"`, and `"%{d}t"`). This attribute provides an alternative means of looking up directive values besides using the named attributes.

entry = None

The original logfile entry with trailing newlines removed

format = None

The entry's log format string

`apachelogs.parse(format, entry, encoding='iso-8859-1', errors=None)`

A convenience function for parsing a single logfile entry without having to directly create a `LogParser` object.

`encoding` and `errors` have the same meaning as for `LogParser`.

`apachelogs.parse_lines(format, entries, encoding='iso-8859-1', errors=None, ignore_invalid=False)`

A convenience function for parsing an iterable of logfile entries without having to directly create a `LogParser` object.

`encoding` and `errors` have the same meaning as for `LogParser`. `ignore_invalid` has the same meaning as for `LogParser.parse_lines()`.

CHAPTER TWO

UTILITIES

```
apachelogs.parse_apache_timestamp(s)
```

Parse an Apache timestamp into a `datetime.datetime` object. The month name in the timestamp is expected to be an abbreviated English name regardless of the current locale.

```
>>> parse_apache_timestamp('[01/Nov/2017:07:28:29 +0000]')
datetime.datetime(2017, 11, 1, 7, 28, 29, tzinfo=datetime.timezone.utc)
```

Parameters `s` (`str`) – a string of the form DD/Mon/YYYY:HH:MM:SS +HHMM (optionally enclosed in square brackets)

Returns an aware `datetime.datetime`

Raises `ValueError` – if `s` is not in the expected format

2.1 Log Format Constants

The following standard log formats are available as string constants in this package so that you don't have to keep typing out the full log format strings:

```
apachelogs.COMBINED = '%h %l %u %t "%r" %>s %b "%{Referer}i" "%{User-Agent}i"'
NCSA extended/combined log format
```

```
apachelogs.COMBINED_DEBIAN = '%h %l %u %t "%r" %>s %O "%{Referer}i" "%{User-Agent}i"'
Like COMBINED, but with %O (total bytes sent including headers) in place of %b (size of response excluding headers)
```

```
apachelogs.COMMON = '%h %l %u %t "%r" %>s %b'
Common log format (CLF)
```

```
apachelogs.COMMON_DEBIAN = '%h %l %u %t "%r" %>s %O'
Like COMMON, but with %O (total bytes sent including headers) in place of %b (size of response excluding headers)
```

```
apachelogs.VHOST_COMBINED = '%v:%p %h %l %u %t "%r" %>s %O "%{Referer}i" "%{User-Agent}i"
COMBINED_DEBIAN with virtual host & port prepended
```

```
apachelogs.VHOST_COMMON = '%v %h %l %u %t "%r" %>s %b'
COMMON with virtual host prepended
```


EXCEPTIONS

exception apachelogs.Error
Bases: `Exception`

The base class for all custom exceptions raised by apachelogs

exception apachelogs.InvalidDirectiveError
Bases: `apachelogs.errors.Error, ValueError`

Raised by the `LogParser` constructor when given a log format containing an invalid or malformed directive

format = None
The log format string containing the invalid directive

pos = None
The position in the log format string at which the invalid directive occurs

exception apachelogs.InvalidEntryError
Bases: `apachelogs.errors.Error, ValueError`

Raised when attempting to parse a log entry that does not match the given log format

entry = None
The invalid log entry
format = None
The log format string the entry failed to match against

exception apachelogs.UnknownDirectiveError
Bases: `apachelogs.errors.Error, ValueError`

Raised by the `LogParser` constructor when given a log format containing an unknown or unsupported directive

directive = None
The unknown or unsupported directive

CHAPTER
FOUR

SUPPORTED DIRECTIVES

The following table lists the log format directives supported by this library along with the names & types of the attributes at which their parsed values are stored on a *LogEntry*. The attribute names for the directives are based off of the names used internally by the Apache source code.

A directive with the < modifier (e.g., %<s) will be stored at `entry.original_attribute_name`, and a directive with the > modifier will be stored at `entry.final_attribute_name`

A type of `str` marked with an asterisk (*) means that the directive's values are decoded according to the `encoding` option to *LogParser*.

Any directive may evaluate to `None` when it is modified by a set of status codes (e.g., %400, 501T or %!200T).

See the [Apache documentation](#) for information on the meaning of each directive.

| Directive | <i>LogEntry</i> Attribute | Type |
|-----------|---|--|
| %% | N/A | N/A |
| %a | <code>entry.remote_address</code> | <code>str</code> |
| %{c}a | <code>entry.remote_client_address</code> | <code>str</code> |
| %A | <code>entry.local_address</code> | <code>str</code> |
| %b | <code>entry.bytes_sent</code> | <code>int</code> or <code>None</code> |
| %B | <code>entry.bytes_sent</code> | <code>int</code> |
| %{name}c | <code>entry.cryptography[name]¹</code> | <code>str</code> or <code>None</code> |
| %{name}C | <code>entry.cookies[name]¹</code> | <code>str*</code> or <code>None</code> |
| %D | <code>entry.request_duration_microseconds</code> | <code>int</code> |
| %{name}e | <code>entry.env_vars[name]¹</code> | <code>str*</code> or <code>None</code> |
| %f | <code>entry.request_file</code> | <code>str*</code> or <code>None</code> |
| %h | <code>entry.remote_host</code> | <code>str*</code> |
| %{c}h | <code>entry.remote_underlying_host</code> | <code>str*</code> |
| %H | <code>entry.request_protocol</code> | <code>str*</code> or <code>None</code> |
| %{name}i | <code>entry.headers_in[name]¹</code> | <code>str*</code> or <code>None</code> |
| %I | <code>entry.bytes_in</code> | <code>int</code> |
| %k | <code>entry.requests_on_connection</code> | <code>int</code> |
| %l | <code>entry.remote_logname</code> | <code>str*</code> or <code>None</code> |
| %L | <code>entry.request_log_id</code> | <code>str</code> or <code>None</code> |
| %{c}L | <code>entry.connection_log_id</code> | <code>str</code> or <code>None</code> |
| %m | <code>entry.request_method</code> | <code>str*</code> or <code>None</code> |

Continued on next page

Table 1 – continued from previous page

| Directive | <i>LogEntry</i> Attribute | Type |
|-------------------------------------|---|--|
| <code>%{name}n</code> | <code>entry.notes[name]¹</code> | <code>str*</code> or <code>None</code> |
| <code>%{name}o</code> | <code>entry.headers_out[name]¹</code> | <code>str*</code> or <code>None</code> |
| <code>%O</code> | <code>entry.bytes_out</code> | <code>int</code> |
| <code>%p</code> | <code>entry.server_port</code> | <code>int</code> |
| <code>%{canonical}p</code> | <code>entry.server_port</code> | <code>int</code> |
| <code>%{local}p</code> | <code>entry.local_port</code> | <code>int</code> |
| <code>%{remote}p</code> | <code>entry.remote_port</code> | <code>int</code> |
| <code>%P</code> | <code>entry.pid</code> | <code>int</code> |
| <code>%{hextid}P²</code> | <code>entry.tid</code> | <code>int</code> |
| <code>%{pid}P</code> | <code>entry.pid</code> | <code>int</code> |
| <code>%{tid}P</code> | <code>entry.tid</code> | <code>int</code> |
| <code>%q</code> | <code>entry.request_query</code> | <code>str*</code> |
| <code>%r</code> | <code>entry.request_line</code> | <code>str*</code> or <code>None</code> |
| <code>%R</code> | <code>entry.handler</code> | <code>str*</code> or <code>None</code> |
| <code>%s</code> | <code>entry.status</code> | <code>int</code> or <code>None</code> |
| <code>%S</code> | <code>entry.bytes_combined</code> | <code>int</code> |
| <code>%t</code> | <code>entry.request_time_fields["time"]</code> | <code>aware datetime.datetime</code> <code>stamp</code> |
| <code>%{sec}t</code> | <code>entry.request_time_fields["epoch"]</code> | <code>int</code> |
| <code>%{msec}t</code> | <code>entry.request_time_fields["millisecond"]</code> | <code>int</code> |
| <code>%{usec}t</code> | <code>entry.request_time_fields["microsecond"]</code> | <code>int</code> |
| <code>%{msec_frac}t</code> | <code>entry.request_time_fields["msec_frac"]</code> | <code>int</code> |
| <code>%{usec_frac}t</code> | <code>entry.request_time_fields["usec_frac"]</code> | <code>int</code> |
| <code>%{strftime_format}t</code> | <code>entry.request_time_fields</code> (See below) | (See below) |
| <code>%T</code> | <code>entry.request_duration_seconds</code> | <code>int</code> |
| <code>%{ms}T</code> | <code>entry.request_duration_milliseconds</code> | <code>int</code> |
| <code>%{us}T</code> | <code>entry.request_duration_microseconds</code> | <code>int</code> |
| <code>%{s}T</code> | <code>entry.request_duration_seconds</code> | <code>int</code> |
| <code>%u</code> | <code>entry.remote_user</code> | <code>str*</code> or <code>None</code> |
| <code>%U</code> | <code>entry.request_uri</code> | <code>str*</code> or <code>None</code> |
| <code>%v</code> | <code>entry.virtual_host</code> | <code>str*</code> |
| <code>%V</code> | <code>entry.server_name</code> | <code>str*</code> |
| <code>%{name}x</code> | <code>entry.variables[name]¹</code> | <code>str</code> or <code>None</code> |
| <code>%X</code> | <code>entry.connection_status</code> | <code>str</code> |
| <code>%^FB</code> | <code>entry.ttfb</code> | <code>int</code> or <code>None</code> |
| <code>%{name}^ti</code> | <code>entry.trailers_in[name]¹</code> | <code>str*</code> or <code>None</code> |
| <code>%{name}^to</code> | <code>entry.trailers_out[name]¹</code> | <code>str*</code> or <code>None</code> |

¹ The `cookies`, `cryptography`, `env_vars`, `headers_in`, `headers_out`, `notes`, `trailers_in`, `trailers_out`, and

4.1 Supported strftime Directives

The following table lists the `strftime` directives supported for use in the parameter of a `%{*}t` directive along with the keys & types at which they are stored in the `dict entry.request_time_fields`. See any C documentation for information on the meaning of each directive.

A `%{*}t` directive with the `begin:` modifier (e.g., `%{begin:%Y-%m-%d}t`) will have its subdirectives stored in `entry.begin_request_time_fields` (in turn used to set `entry.begin_request_time`), and likewise for the `end:` modifier.

| Directive | request_time_fields key | Type |
|-----------------|---------------------------------|--|
| <code>%%</code> | N/A | N/A |
| <code>%a</code> | <code>"abbrev_wday"</code> | <code>str</code> |
| <code>%A</code> | <code>"full_wday"</code> | <code>str</code> |
| <code>%b</code> | <code>"abbrev_mon"</code> | <code>str</code> |
| <code>%B</code> | <code>"full_mon"</code> | <code>str</code> |
| <code>%C</code> | <code>"century"</code> | <code>int</code> |
| <code>%d</code> | <code>"mday"</code> | <code>int</code> |
| <code>%D</code> | <code>"date"</code> | <code>datetime.date</code> |
| <code>%e</code> | <code>"mday"</code> | <code>int</code> |
| <code>%F</code> | <code>"date"</code> | <code>datetime.date</code> |
| <code>%g</code> | <code>"abbrev_week_year"</code> | <code>int</code> |
| <code>%G</code> | <code>"week_year"</code> | <code>int</code> |
| <code>%h</code> | <code>"abbrev_mon"</code> | <code>str</code> |
| <code>%H</code> | <code>"hour"</code> | <code>int</code> |
| <code>%I</code> | <code>"hour12"</code> | <code>int</code> |
| <code>%j</code> | <code>"yday"</code> | <code>int</code> |
| <code>%m</code> | <code>"mon"</code> | <code>int</code> |
| <code>%M</code> | <code>"min"</code> | <code>int</code> |
| <code>%n</code> | N/A | N/A |
| <code>%p</code> | <code>"am_pm"</code> | <code>str</code> |
| <code>%R</code> | <code>"hour_min"</code> | <code>datetime.time</code> |
| <code>%s</code> | <code>"epoch"</code> | <code>int</code> |
| <code>%S</code> | <code>"sec"</code> | <code>int</code> |
| <code>%t</code> | N/A | N/A |
| <code>%T</code> | <code>"time"</code> | <code>datetime.time</code> |
| <code>%u</code> | <code>"iso_wday"</code> | <code>int</code> |
| <code>%U</code> | <code>"sunday_weeknum"</code> | <code>int</code> |
| <code>%V</code> | <code>"iso_weeknum"</code> | <code>int</code> |
| <code>%w</code> | <code>"wday"</code> | <code>int</code> |
| <code>%W</code> | <code>"monday_weeknum"</code> | <code>int</code> |
| <code>%y</code> | <code>"abbrev_year"</code> | <code>int</code> |
| <code>%Y</code> | <code>"year"</code> | <code>int</code> |
| <code>%z</code> | <code>"timezone"</code> | <code>datetime.timezone or None</code> |
| <code>%Z</code> | <code>"tzname"</code> | <code>str</code> |

variables attributes are case-insensitive `dicts`.

² Apache renders `%{hexid}P` as either a decimal integer or a hexadecimal integer depending on the APR version available. `apachelogs` expects `%{hexid}P` to always be in hexadecimal; if your Apache produces decimal integers instead, you must instead use `%{tid}P` in the log format passed to `apachelogs`.

CHANGELOG

5.1 v0.4.0 (2019-05-19)

- Support the `%{c}h` log directive
- `%f` and `%R` can now be `None`
- **Bugfix:** `%u` can now match the string " " (two double quotes)
- Support `mod_ssl`'s `%{*}c` and `%{*}x` directives
- Support the `%{hextid}P` directive (as a hexadecimal integer)
- Support the `%L` and `%{c}L` directives
- Parameters to `%{*}p`, `%{*}P`, and `%{*}T` are now treated case-insensitively in order to mirror Apache's behavior
- **Refined some directives to better match only the values emitted by Apache:**
 - `%l` and `%m` no longer accept whitespace
 - `%s` and `%{tid}P` now only match unsigned integers
 - `%{*}C` no longer accepts semicolons or leading or trailing spaces
 - `%q` no longer accepts whitespace or pound/hash signs

5.2 v0.3.0 (2019-05-12)

- Gave `LogEntry` a `directives` attribute for looking up directive values by the corresponding log format directives

5.3 v0.2.0 (2019-05-09)

- Changed the capitalization of “User-agent” in the log format string constants to “User-Agent”
- The `cookies`, `env_vars`, `headers_in`, `headers_out`, `notes`, `trailers_in`, and `trailers_out` attributes of `LogEntry` are now all case-insensitive `dicts`.

5.4 v0.1.0 (2019-05-06)

Initial release

apachelogs parses Apache access log files. Pass it a `log format string` and get back a parser for logfile entries in that format. *apachelogs* even takes care of decoding escape sequences and converting things like timestamps, integers, and bare hyphens to `datetime` values, `ints`, and `Nones`.

**CHAPTER
SIX**

INSTALLATION

apachelogs requires Python 3.5 or higher. Just use `pip` for Python 3 (You have `pip`, right?) to install *apachelogs* and its dependencies:

```
python3 -m pip install apachelogs
```

CHAPTER SEVEN

EXAMPLES

Parse a single log entry:

```
>>> from apachelogs import LogParser
>>> parser = LogParser("%h %l %u %t \"%r\" %>s %b \"%{Referer}i\" \"%{User-Agent}i\"")
>>> # The above log format is also available as the constant `apachelogs.COMBINED`.
>>> entry = parser.parse('209.126.136.4 -- [01/Nov/2017:07:28:29 +0000] "GET / HTTP/
->1.1" 301 521 "-" "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36
->(KHTML, like Gecko) Chrome/57.0.2987.133 Safari/537.36"\n')
>>> entry.remote_host
'209.126.136.4'
>>> entry.request_time
datetime.datetime(2017, 11, 1, 7, 28, 29, tzinfo=datetime.timezone.utc)
>>> entry.request_line
'GET / HTTP/1.1'
>>> entry.final_status
301
>>> entry.bytes_sent
521
>>> entry.headers_in["Referer"] is None
True
>>> entry.headers_in["User-Agent"]
'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko)
->Chrome/57.0.2987.133 Safari/537.36'
>>> # Log entry components can also be looked up by directive:
>>> entry.directives["%r"]
'GET / HTTP/1.1'
>>> entry.directives["%>s"]
301
>>> entry.directives["%t"]
datetime.datetime(2017, 11, 1, 7, 28, 29, tzinfo=datetime.timezone.utc)
```

Parse a file full of log entries:

```
>>> with open('/var/log/apache2/access.log') as fp:
...     for entry in parser.parse_lines(fp):
...         print(str(entry.request_time), entry.request_line)
...
2019-01-01 12:34:56-05:00 GET / HTTP/1.1
2019-01-01 12:34:57-05:00 GET /favicon.ico HTTP/1.1
2019-01-01 12:34:57-05:00 GET /styles.css HTTP/1.1
# etc.
```

**CHAPTER
EIGHT**

INDICES AND TABLES

- genindex
- search

PYTHON MODULE INDEX

a

apachelogs, 1

INDEX

A

apachelogs (*module*), 1

C

COMBINED (*in module apachelogs*), 5
COMBINED_DEBIAN (*in module apachelogs*), 5
COMMON (*in module apachelogs*), 5
COMMON_DEBIAN (*in module apachelogs*), 5

D

directive (*apachelogs.UnknownDirectiveError attribute*), 7
directives (*apachelogs.LogEntry attribute*), 4

E

entry (*apachelogs.InvalidEntryError attribute*), 7
entry (*apachelogs.LogEntry attribute*), 4
Error, 7

F

format (*apachelogs.InvalidDirectiveError attribute*), 7
format (*apachelogs.InvalidEntryError attribute*), 7
format (*apachelogs.LogEntry attribute*), 4

I

InvalidDirectiveError, 7
InvalidEntryError, 7

L

LogEntry (*class in apachelogs*), 3
LogParser (*class in apachelogs*), 3

P

parse () (*apachelogs.LogParser method*), 3
parse () (*in module apachelogs*), 4
parse_apache_timestamp () (*in module apachelogs*), 5
parse_lines () (*apachelogs.LogParser method*), 3
parse_lines () (*in module apachelogs*), 4
pos (*apachelogs.InvalidDirectiveError attribute*), 7

U

UnknownDirectiveError, 7

V

VHOST_COMBINED (*in module apachelogs*), 5
VHOST_COMMON (*in module apachelogs*), 5