# 13.4. `lzma` — Compression using the LZMA algorithm

*New in version 3.3.*

**Source code:** Lib/lzma.py

This module provides classes and convenience functions for compressing and decompressing data using the LZMA compression algorithm. Also included is a file interface supporting the `.xz` and legacy `.lzma` file formats used by the **xz** utility, as well as raw compressed streams.

The interface provided by this module is very similar to that of the `bz2` module. However, note that `LZMAFile` is *not* thread-safe, unlike `bz2.BZ2File`, so if you need to use a single `LZMAFile` instance from multiple threads, it is necessary to protect it with a lock.

*exception* `lzma.`**`LZMAError`**

This exception is raised when an error occurs during compression or decompression, or while initializing the compressor/decompressor state.

## 13.4.1. Reading and writing compressed files

`lzma.`**`open`**(*filename*, *mode="rb"*, *\**, *format=None*, *check=-1*, *preset=None*, *filters=None*, *encoding=None*, *errors=None*, *newline=None*)

Open an LZMA-compressed file in binary or text mode, returning a file object.

The *filename* argument can be either an actual file name (given as a `str`, `bytes` or path-like object), in which case the named file is opened, or it can be an existing file object to read from or write to.

The *mode* argument can be any of `"r"`, `"rb"`, `"w"`, `"wb"`, `"x"`, `"xb"`, `"a"` or `"ab"` for binary mode, or `"rt"`, `"wt"`, `"xt"`, or `"at"` for text mode. The default is `"rb"`.

When opening a file for reading, the *format* and *filters* arguments have the same meanings as for `LZMADecompressor`. In this case, the *check* and *preset* arguments should not be used.

When opening a file for writing, the *format*, *check*, *preset* and *filters* arguments have the same meanings as for `LZMACompressor`.

For binary mode, this function is equivalent to the `LZMAFile` constructor: `LZMAFile(filename, mode, ...)`. In this case, the *encoding*, *errors* and *newline* arguments must not be provided.

For text mode, a `LZMAFile` object is created, and wrapped in an `io.TextIOWrapper` instance with the specified encoding, error handling behavior, and line ending(s).

*Changed in version 3.4:* Added support for the "x", "xb" and "xt" modes.

*Changed in version 3.6:* Accepts a path-like object.

*class* `lzma.`**`LZMAFile`**(*filename=None*, *mode="r"*, *\**, *format=None*, *check=-1*, *preset=None*, *filters=None*)

Open an LZMA-compressed file in binary mode.

An `LZMAFile` can wrap an already-open file object, or operate directly on a named file. The *filename* argument specifies either the file object to wrap, or the name of the file to open (as a `str`, `bytes` or path-like object). When wrapping an existing file object, the wrapped file will not be closed when the `LZMAFile` is closed.

The *mode* argument can be either "r" for reading (default), "w" for overwriting, "x" for exclusive creation, or "a" for appending. These can equivalently be given as "rb", "wb", "xb" and "ab" respectively.

If *filename* is a file object (rather than an actual file name), a mode of "w" does not truncate the file, and is instead equivalent to "a".

When opening a file for reading, the input file may be the concatenation of multiple separate compressed streams. These are transparently decoded as a single logical stream.

When opening a file for reading, the *format* and *filters* arguments have the same meanings as for `LZMADecompressor`. In this case, the *check* and *preset* arguments should not be used.

When opening a file for writing, the *format*, *check*, *preset* and *filters* arguments have the same meanings as for `LZMACompressor`.

`LZMAFile` supports all the members specified by `io.BufferedIOBase`, except for `detach()` and `truncate()`. Iteration and the `with` statement are supported.

The following method is also provided:

**`peek`**(*size=-1*)

Return buffered data without advancing the file position. At least one byte of data will be returned, unless EOF has been reached. The exact number of bytes returned is unspecified (the *size* argument is ignored).

> **Note:** While calling `peek()` does not change the file position of the `LZMAFile`, it may change the position of the underlying file object (e.g. if the `LZMAFile` was constructed by passing a file object for *filename*).

*Changed in version 3.4:* Added support for the "x" and "xb" modes.

*Changed in version 3.5:* The `read()` method now accepts an argument of None.

*Changed in version 3.6:* Accepts a path-like object.

# 13.4.2. Compressing and decompressing data in memory

*class* lzma. **LZMACompressor**(*format=FORMAT_XZ, check=-1, preset=None, filters=None*)

Create a compressor object, which can be used to compress data incrementally.

For a more convenient way of compressing a single chunk of data, see `compress()`.

The *format* argument specifies what container format should be used. Possible values are:

- `FORMAT_XZ`: The `.xz` container format.
  This is the default format.

- `FORMAT_ALONE`: The legacy `.lzma` container format.
  This format is more limited than `.xz` – it does not support integrity checks or multiple filters.

- `FORMAT_RAW`: A raw data stream, not using any container format.
  This format specifier does not support integrity checks, and requires that you always specify a custom filter chain (for both compression and decompression). Additionally, data compressed in this manner cannot be decompressed using `FORMAT_AUTO` (see `LZMADecompressor`).

The *check* argument specifies the type of integrity check to include in the compressed data. This check is used when decompressing, to ensure that the data has not been corrupted. Possible values are:

- `CHECK_NONE`: No integrity check. This is the default (and the only acceptable value) for `FORMAT_ALONE` and `FORMAT_RAW`.
- `CHECK_CRC32`: 32-bit Cyclic Redundancy Check.
- `CHECK_CRC64`: 64-bit Cyclic Redundancy Check. This is the default for `FORMAT_XZ`.
- `CHECK_SHA256`: 256-bit Secure Hash Algorithm.

If the specified check is not supported, an `LZMAError` is raised.

The compression settings can be specified either as a preset compression level (with the *preset* argument), or in detail as a custom filter chain (with the *filters* argument).

The *preset* argument (if provided) should be an integer between `0` and `9` (inclusive), optionally OR-ed with the constant `PRESET_EXTREME`. If neither *preset* nor *filters* are given, the default behavior is to use `PRESET_DEFAULT` (preset level `6`). Higher presets produce smaller output, but make the compression process slower.

> **Note:** In addition to being more CPU-intensive, compression with higher presets also requires much more memory (and produces output that needs more memory to decompress). With preset `9` for example, the overhead for an `LZMACompressor` object can be as high as 800 MiB. For this reason, it is generally best to stick with the default preset.

The *filters* argument (if provided) should be a filter chain specifier. See Specifying custom filter chains for details.

#### compress(*data*)

Compress *data* (a `bytes` object), returning a `bytes` object containing compressed data for at least part of the input. Some of *data* may be buffered internally, for use in later calls to `compress()` and `flush()`. The returned data should be concatenated with the output of any previous calls to `compress()`.

#### flush()

Finish the compression process, returning a `bytes` object containing any data stored in the compressor's internal buffers.

The compressor cannot be used after this method has been called.

*class* `lzma.` **LZMADecompressor**(*format=FORMAT_AUTO*, *memlimit=None*, *filters=None*)

Create a decompressor object, which can be used to decompress data incrementally.

For a more convenient way of decompressing an entire compressed stream at once, see `decompress()`.

The *format* argument specifies the container format that should be used. The default is `FORMAT_AUTO`, which can decompress both `.xz` and `.lzma` files. Other possible values are `FORMAT_XZ`, `FORMAT_ALONE`, and `FORMAT_RAW`.

The *memlimit* argument specifies a limit (in bytes) on the amount of memory that the decompressor can use. When this argument is used, decompression will fail with an `LZMAError` if it is not possible to decompress the input within the given memory limit.

The *filters* argument specifies the filter chain that was used to create the stream being decompressed. This argument is required if *format* is `FORMAT_RAW`, but should not be used for other formats. See Specifying custom filter chains for more information about filter chains.

> **Note:** This class does not transparently handle inputs containing multiple compressed streams, unlike `decompress()` and `LZMAFile`. To decompress a multi-stream input with `LZMADecompressor`, you must create a new decompressor for each stream.

**decompress**(*data*, *max_length=-1*)

Decompress *data* (a bytes-like object), returning uncompressed data as bytes. Some of *data* may be buffered internally, for use in later calls to `decompress()`. The returned data should be concatenated with the output of any previous calls to `decompress()`.

If *max_length* is nonnegative, returns at most *max_length* bytes of decompressed data. If this limit is reached and further output can be produced, the `needs_input` attribute will be set to `False`. In this case, the next call to `decompress()` may provide *data* as `b''` to obtain more of the output.

If all of the input data was decompressed and returned (either because this was less than *max_length* bytes, or because *max_length* was negative), the `needs_input` attribute will be set to `True`.

Attempting to decompress data after the end of stream is reached raises an *EOFError*. Any data found after the end of the stream is ignored and saved in the `unused_data` attribute.

*Changed in version 3.5:* Added the *max_length* parameter.

### check

The ID of the integrity check used by the input stream. This may be `CHECK_UNKNOWN` until enough of the input has been decoded to determine what integrity check it uses.

### eof

`True` if the end-of-stream marker has been reached.

### unused_data

Data found after the end of the compressed stream.

Before the end of the stream is reached, this will be `b""`.

### needs_input

`False` if the `decompress()` method can provide more decompressed data before requiring new uncompressed input.

*New in version 3.5.*

lzma. **compress**(*data*, *format=FORMAT_XZ*, *check=-1*, *preset=None*, *filters=None*)

Compress *data* (a `bytes` object), returning the compressed data as a `bytes` object.

See `LZMACompressor` above for a description of the *format*, *check*, *preset* and *filters* arguments.

lzma. **decompress**(*data*, *format=FORMAT_AUTO*, *memlimit=None*, *filters=None*)

Decompress *data* (a `bytes` object), returning the uncompressed data as a `bytes` object.

If *data* is the concatenation of multiple distinct compressed streams, decompress all of these streams, and return the concatenation of the results.

See `LZMADecompressor` above for a description of the *format*, *memlimit* and *filters* arguments.

## 13.4.3. Miscellaneous

lzma.**is_check_supported**(*check*)
>    Returns true if the given integrity check is supported on this system.
>
>    `CHECK_NONE` and `CHECK_CRC32` are always supported. `CHECK_CRC64` and `CHECK_SHA256` may be unavailable if you are using a version of **liblzma** that was compiled with a limited feature set.

## 13.4.4. Specifying custom filter chains

A filter chain specifier is a sequence of dictionaries, where each dictionary contains the ID and options for a single filter. Each dictionary must contain the key `"id"`, and may contain additional keys to specify filter-dependent options. Valid filter IDs are as follows:

- Compression filters:
    - `FILTER_LZMA1` (for use with `FORMAT_ALONE`)
    - `FILTER_LZMA2` (for use with `FORMAT_XZ` and `FORMAT_RAW`)

- Delta filter:
    - `FILTER_DELTA`

- Branch-Call-Jump (BCJ) filters:
    - `FILTER_X86`
    - `FILTER_IA64`
    - `FILTER_ARM`
    - `FILTER_ARMTHUMB`
    - `FILTER_POWERPC`
    - `FILTER_SPARC`

A filter chain can consist of up to 4 filters, and cannot be empty. The last filter in the chain must be a compression filter, and any other filters must be delta or BCJ filters.

Compression filters support the following options (specified as additional entries in the dictionary representing the filter):

- `preset`: A compression preset to use as a source of default values for options that are not specified explicitly.
- `dict_size`: Dictionary size in bytes. This should be between 4 KiB and 1.5 GiB (inclusive).
- `lc`: Number of literal context bits.
- `lp`: Number of literal position bits. The sum `lc + lp` must be at most 4.

- `pb`: Number of position bits; must be at most 4.
- `mode`: `MODE_FAST` or `MODE_NORMAL`.
- `nice_len`: What should be considered a "nice length" for a match. This should be 273 or less.
- `mf`: What match finder to use — `MF_HC3`, `MF_HC4`, `MF_BT2`, `MF_BT3`, or `MF_BT4`.
- `depth`: Maximum search depth used by match finder. 0 (default) means to select automatically based on other filter options.

The delta filter stores the differences between bytes, producing more repetitive input for the compressor in certain circumstances. It supports one option, `dist`. This indicates the distance between bytes to be subtracted. The default is 1, i.e. take the differences between adjacent bytes.

The BCJ filters are intended to be applied to machine code. They convert relative branches, calls and jumps in the code to use absolute addressing, with the aim of increasing the redundancy that can be exploited by the compressor. These filters support one option, `start_offset`. This specifies the address that should be mapped to the beginning of the input data. The default is 0.

# 13.4.5. Examples

Reading in a compressed file:

```python
import lzma
with lzma.open("file.xz") as f:
    file_content = f.read()
```

Creating a compressed file:

```python
import lzma
data = b"Insert Data Here"
with lzma.open("file.xz", "w") as f:
    f.write(data)
```

Compressing data in memory:

```python
import lzma
data_in = b"Insert Data Here"
data_out = lzma.compress(data_in)
```

Incremental compression:

```python
import lzma
lzc = lzma.LZMACompressor()
out1 = lzc.compress(b"Some data\n")
```

```
out2 = lzc.compress(b"Another piece of data\n")
out3 = lzc.compress(b"Even more data\n")
out4 = lzc.flush()
# Concatenate all the partial results:
result = b"".join([out1, out2, out3, out4])
```

Writing compressed data to an already-open file:

```
import lzma
with open("file.xz", "wb") as f:
    f.write(b"This data will not be compressed\n")
    with lzma.open(f, "w") as lzf:
        lzf.write(b"This *will* be compressed\n")
    f.write(b"Not compressed\n")
```

Creating a compressed file using a custom filter chain:

```
import lzma
my_filters = [
    {"id": lzma.FILTER_DELTA, "dist": 5},
    {"id": lzma.FILTER_LZMA2, "preset": 7 | lzma.PRESET_EXTREME},
]
with lzma.open("file.xz", "w", filters=my_filters) as f:
    f.write(b"blah blah blah")
```