

4 The BAM Format Specification

4.1 The BGZF compression format

BGZF is block compression implemented on top of the standard gzip file format.²⁶ The goal of BGZF is to provide good compression while allowing efficient random access to the BAM file for indexed queries. The BGZF format is ‘gunzip compatible’, in the sense that a compliant gunzip utility can decompress a BGZF compressed file.²⁷

A BGZF file is a series of concatenated BGZF blocks, each no larger than ~~64Kb-before or 64KB before and~~ after compression. Each BGZF block is itself a spec-compliant gzip archive which contains an “extra field” in the format described in RFC1952. The gzip file format allows the inclusion of application-specific extra fields and these are ignored by compliant decompression implementation. The gzip specification also allows gzip files to be concatenated. The result of decompressing concatenated gzip files is the concatenation of the uncompressed data.

Each BGZF block contains a standard gzip file header with the following standard-compliant extensions:

1. The `F.EXTRA` bit in the header is set to indicate that extra fields are present.
2. The extra field used by BGZF uses the two subfield ID values 66 and 67 (ASCII ‘BC’).
3. The length of the BGZF extra field payload (field `LEN` in the gzip specification) is 2 (two bytes of payload).
4. The payload of the BGZF extra field is a 16-bit unsigned integer in little endian format. This integer gives the size of the containing BGZF block minus one.

On disk, a complete BGZF file is a series of blocks as shown in the following table. (All integers are little endian as is required by RFC1952.)

Field	Description	Type	Value
<i>List of compression blocks (until the end of the file)</i>			
<code>ID1</code>	gzip IDentifier1	<code>uint8_t</code>	31
<code>ID2</code>	gzip IDentifier2	<code>uint8_t</code>	139
<code>CM</code>	gzip Compression Method	<code>uint8_t</code>	8
<code>FLG</code>	gzip FLaGs	<code>uint8_t</code>	4
<code>MTIME</code>	gzip Modification TIME	<code>uint32_t</code>	
<code>XFL</code>	gzip eXtra FLags	<code>uint8_t</code>	
<code>OS</code>	gzip Operating System	<code>uint8_t</code>	
<code>XLEN</code>	gzip eXtra LENgth	<code>uint16_t</code>	
<i>Extra subfield(s) (total size=XLEN)</i>			
<i>Additional RFC1952 extra subfields if present</i>			
<code>SI1</code>	Subfield Identifier1	<code>uint8_t</code>	66
<code>SI2</code>	Subfield Identifier2	<code>uint8_t</code>	67
<code>SLEN</code>	Subfield LENgth	<code>uint16_t</code>	2
<code>BSIZE</code>	total Block SIZE minus 1	<code>uint16_t</code>	
<i>Additional RFC1952 extra subfields if present</i>			
<code>CDATA</code>	Compressed DATA by <code>zlib::deflate()</code>	<code>uint8_t[BSIZE-XLEN-19]</code>	
<code>CRC32</code>	CRC-32	<code>uint32_t</code>	
<code>ISIZE</code>	Input SIZE (length of uncompressed data)	<code>uint32_t</code>	

The random access method to be described next limits the uncompressed contents of each BGZF block to a maximum of 2^{16} bytes of data. Thus while `ISIZE` is stored as a `uint32_t` as per the gzip format, in BGZF it is limited to the range [0, 65536]. `BSIZE` can represent BGZF block sizes in the range [1, 65536], though typically `BSIZE` will be rather less than `ISIZE` due to compression.

²⁶L. Peter Deutsch, **GZIP file format specification version 4.3**, RFC 1952.

²⁷It is worth noting that there is a known bug in the Java `GZIPInputStream` class that concatenated gzip archives cannot be successfully decompressed by this class. BGZF files can be created and manipulated using the built-in Java `util.zip` package, but naive use of `GZIPInputStream` on a BGZF file will not work due to this bug.