

## BCF2 Quick Reference (r198)

In BCF2, each key in the **FILTER**, **INFO** and **FORMAT** fields is required to be defined in the VCF header. For each record, a key is stored as an integer which is the index of its first appearance in the header. ‘PASS’ is always indexed at 0, which is special cased as VCF does not require the presence of this word.

In BCF2, a typed value consists of a typing byte and the actual value with type mandated by the typing byte. In the typing byte, the lowest four bits give the atomic type. If the number represented by the higher 4 bits is smaller than 15, it is the size of the following vector; if the number equals 15, the following typed integer is the array size. The highest 4 bits of a **Flag** type equals 0 and in this case, no assumptions can be made about the lower 4 bits. The table below gives the atomic types and their missing values:

Bit 0-3	C type	Missing value	Description
1	<code>int8_t</code>	0x80	signed 8-bit integer
2	<code>int16_t</code>	0x8000	signed 16-bit integer
3	<code>int32_t</code>	0x80000000	signed 32-bit integer
5	<code>float</code>	0x7F800001	IEEE 32-bit floating pointer number
7	<code>char</code>	‘\0’	character

A genotype (GT) is encoded as an integer vector with each integer describing an allele and its phase w.r.t. the previous allele. The first allele does not carry the phase information. In the vector, each integer is organized as ‘(allele+1)<<1|phased’ where **allele** is set to -1 if the allele in GT is a dot ‘.’ (thus the higher bits are all 0). The vector is padded with missing values if the GT having fewer ploidy.

A BCF2 file is BGZF compressed and all multi-byte value are little endian.

Field	Description	Type	Value
<b>magic</b>	BCF2 magic string	<code>char[5]</code>	BCF\2\1
<b>l_text</b>	Length of the header text, including any NULL padding	<code>uint32_t</code>	
<b>text</b>	NULL-terminated plain VCF header text	<code>char[l_text]</code>	
<i>List of VCF records (until the end of the BGZF section)</i>			
<b>l_shared</b>	Data length from <b>CHROM</b> to the end of <b>INFO</b>	<code>uint32_t</code>	
<b>l_indiv</b>	Data length of <b>FORMAT</b> and individual genotype fields	<code>uint32_t</code>	
<b>CHROM</b>	Reference sequence ID	<code>int32_t</code>	
<b>POS</b>	0-based leftmost coordinate	<code>int32_t</code>	
<b>rlen</b>	Length of reference sequence	<code>int32_t</code>	
<b>QUAL</b>	Variant quality; 0x7F800001 for a missing value	<code>float</code>	
<b>n_allele_info</b>	<code>n_allele&lt;&lt;16 n_info</code>	<code>uint32_t</code>	
<b>n_fmt_sample</b>	<code>n_fmt&lt;&lt;24 n_sample</code>	<code>uint32_t</code>	
<b>ID</b>	Variant identifier	<code>typed str</code>	
<i>List of alleles in the REF and ALT fields (n=n_allele)</i>			
<b>allele</b>	A reference or alternate allele	<code>typed str</code>	
<b>FILTER</b>	List of filters; filters are defined in the dictionary	<code>typed vec</code>	
<i>List of key-value pairs in the INFO field (n=n_info)</i>			
<b>info_key</b>	Info key, defined in the dictionary	<code>typed int</code>	
<b>info_value</b>	Value	<code>typed val</code>	
<i>List of FORMATS and sample information (n=n_fmt)</i>			
<b>fmt_key</b>	Format key, defined in the dictionary	<code>typed int</code>	
<b>fmt_type</b>	Typing byte of each individual value, possibly followed by a typed int for the vector length	<code>uint8_t+</code>	
<b>fmt_value</b>	Array of values. The information of each individual is concatenated in the vector. Every value is of the same <b>fmt_type</b> . Variable-length vectors are padded with missing values; a string is stored as a vector of <b>char</b> .	(by <b>fmt_type</b> )	