

Sequence Alignment/Map Format Specification

The SAM/BAM Format Specification Working Group

15 May 2023

The master version of this document can be found at <https://github.com/samtools/hts-specs>.
This printing is version 7cfd789 from that repository, last modified on the date shown above.

1 The SAM Format Specification

SAM stands for Sequence Alignment/Map format. It is a TAB-delimited text format consisting of a header section, which is optional, and an alignment section. If present, the header must be prior to the alignments. Header lines start with '@', while alignment lines do not. Each alignment line has 11 mandatory fields for essential alignment information such as mapping position, and variable number of optional fields for flexible or aligner specific information.

This specification is for version 1.6 of the SAM and BAM formats. Each SAM and BAM file may optionally specify the version being used via the @HD VN tag. For full version history see Appendix B.

~~Unless explicitly specified elsewhere, all fields are encoded using~~ SAM file contents are 7-bit US-ASCII¹ ~~in~~ , except for certain field values as individually specified which may contain other Unicode characters encoded in UTF-8. Alternatively and equivalently, SAM files are encoded in UTF-8 but non-ASCII characters are permitted only within certain field values as explicitly specified in the descriptions of those fields.¹

~~Where it makes a difference, SAM file contents should be read and written using the POSIX / C locale. Regular expressions listed use the POSIX~~ For example, floating-point values in SAM always use '.' for the decimal-point character.

~~The regular expressions in this specification are written using the POSIX / IEEE Std 1003.1 extended syntax.~~

1.1 An example

Suppose we have the following alignment with bases in lowercase clipped from the alignment. Read r001/1 and r001/2 constitute a read pair; r003 is a chimeric read; r004 represents a split alignment.

```
Coor      12345678901234  5678901234567890123456789012345
ref       AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT

+r001/1      TTAGATAAAGGATA*CTG
+r002        aaaAGATAA*GGATA
+r003        gcctaAGCTAA
+r004                ATAGCT.....TCAGC
-r003                ttagctTAGGC
-r001/2                CAGCGGCAT
```

¹~~Charset ANSI_X3.4-1968 as defined in RFC1345.~~

¹Hence in particular SAM files must not begin with a byte order mark (BOM) and lines of text are delimited by ASCII line terminator characters only. In addition to the local platform's text file line termination conventions, implementations may wish to support LF and CR LF for interoperability with other platforms.