

# The BAII C Format Specification

Version 0.1

Wei Tian and Joseph R Ecker

BAII C stands for Binary AII Cytosine format, storing DNA 5mC methylation status particularly from whole genome sequencing. It's conceived as a binary alternative to the AII C format, which is a TAB-delimited text format. BAII C not only reduces file sizes and speeds up read/write operations but also eliminates redundancies prevalent in the AII C format, especially for single cell DNA methylation data, leading to a more streamlined and efficient data representation.

## Background on the AII C Format

Before diving into the details of BAII C, it's essential to understand its predecessor, the AII C format. AII C is widely used in bioinformatics for representing cytosine-specific methylation data. The AII C format is a TAB-delimited text, and the typical fields in an AII C file are:

Index	Field name	Example	Note
1	ref	chr12	
2	pos	18283342	1-based
3	strandness	+	either + or -
4	c_context	CGT	3 bases starting with "C"; can use bases defined by IUPAC (eg. H=A,T,orC, N=A,T,C, or G, etc)
5	mc	18	count of reads supporting methylation
6	cov	21	read coverage
7	methyated	1	indicator of significant methylation (1 if no test is performed)
8	(optional) num_matches	3,2,3	number of matched base calls at context nucleotides
9	(optional) num_mismatches	0,1,0	number of mismatches at context nucleotides

## The BAII C Solution

One of the primary inefficiencies of AII C lies in the repeated storage of strandness and c\_context information, particularly for single cell studies. For a given position in a specified genome assembly, the strandness and cytosine context remain consistent across single cells. Storing this redundant information for each cell is wasteful in terms of both storage and processing time.

BAII C addresses the redundancies of the AII C format by segregating the consistent metadata of every cytosine, like strandness and c\_context, into a separate CMeta (cytosine meta) file. The fields 7-9 of the AII C format, which are rarely used in practice, are omitted from BAII C as well.

Furthermore, BALIC adopt a binary representation instead of a text-based one to reduce the file size further.

## Format Structures

Field	Description	Type	Value
<b>header</b>			
magic	magic string	char[6]	"BALLC\1"
version	major version of the BALIC format	uint8_t	e.g. "0"
version_minor	minor version of the BALIC format	uint8_t	e.g. "1"
sc	flag of single cell BALIC format	uint8_t	"0" or "1"
l_assembly	length of assembly text	uint32_t	
assembly_text	user specified assembly info	char[l_assembly]	e.g. "mm10", "hg38", "donor1"
l_text	length of header text	uint32_t	
header_text	note for this balic file	char[l_text]	
n_refs	number of references (chromosomes)	uint16_t	
refs	list of refs (n=n_refs)		
	<b>ref</b>		
	l_name	length of reference name	uint32_t
	ref_name	reference name	char[l_name]
	ref_len	length of reference	uint32_t
			e.g. "248956422"
<b>mc_records</b> list of mc_records			
	<b>mc_record</b>		
	pos	position of cytosine	uint32_t
	ref_id	reference id, $0 \leq \text{ref\_id} < \text{n\_refs}$	uint16_t
	mc	number of counts supporting methylation call	uint16_t when sc=1; uint8_t when sc=0
	cov	number of counts covering the position	uint16_t when sc=1; uint8_t when sc=0

The BALIC file starts with a header that encapsulates essential metadata:

1. Magic String (magic): A signature for the BALIC format. It is always set to "BALLC\1".
2. Version and Version Minor: These fields denote the major and minor versions of the BALIC format, respectively. For instance, "0" and "1" represent version 0.1.
3. Single Cell Flag (sc): A flag that indicates if the BALIC format is for single-cell data. It can be "0" (not single-cell) or "1" (single-cell).
4. Assembly Text: This field provides user-specified assembly information, such as "mm10", "hg38", or "donor1". The length of the assembly text is determined by the l\_assembly field.
5. Header Text: This is a note for the BALIC file. The length of this text is determined by the l\_text field.
6. Number of References (n\_refs): This field denotes the number of references (chromosomes/contigs) in the genome assembly.
7. References (refs): This is a list containing details of each reference. Each reference has a name (like "chr1"), length of which determined by "l\_name" and a length of the reference (for example, "248956422" for "chr1" of human).

Following the header, the file contains data blocks that encapsulate the actual methylation data:

1. Position (pos): This indicates the position of the cytosine in the reference.
2. Reference ID (ref\_id): This ID links the data block to one of the references specified in the header.
3. Methylation Count (mc): Represents the number of counts supporting the methylation call. Its data type varies based on the sc flag.
4. Coverage (cov): Denotes the number of counts covering the position. Like mc, its data type depends on the sc flag.

### CMeta file

The CMeta file serves as a companion to the BAIIIC format, efficiently storing the consistent metadata of each cytosine. It is essentially a TAB-delimited text file that includes following fields:

Index	Field name	Example	Note
1	ref	chr12	
2	pos	18283342	1-based
3	strandness	+	either + or -
4	c_context	CGT	3 bases starting with "C"; can use bases defined by IUPAC (eg. H=A,T,orC; N=A,T,C, or G, etc)

By segregating this information, which remains static across single cells for a given genomic position, into the CMeta file, BAIIIC effectively minimizes redundancy and optimizes storage. When using BAIIIC in tandem with CMeta, researchers can achieve the same comprehensive view of methylation data yet without the inefficiencies associated with the AIIIC format.

## Working with BAIIIC

TODO...