

**INTERNATIONAL
STANDARD**

**ISO/IEC
23092-6**

First edition
2023-11

**Information technology — Genomic
information representation —**

**Part 6:
Coding of genomic annotations**



Reference number
ISO/IEC 23092-6:2023(E)

© ISO/IEC 2023



COPYRIGHT PROTECTED DOCUMENT

© ISO/IEC 2023

All rights reserved. Unless otherwise specified, or required in the context of its implementation, no part of this publication may be reproduced or utilized otherwise in any form or by any means, electronic or mechanical, including photocopying, or posting on the internet or an intranet, without prior written permission. Permission can be requested from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office
CP 401 • Ch. de Blandonnet 8
CH-1214 Vernier, Geneva
Phone: +41 22 749 01 11
Email: copyright@iso.org
Website: www.iso.org

Published in Switzerland

Contents

Page

Foreword.....	vi
Introduction.....	vii
1 Scope.....	1
2 Normative references.....	1
3 Terms and definitions.....	1
4 Abbreviated terms.....	8
5 Conventions.....	8
5.1 General.....	8
5.2 Logical operators.....	8
5.3 Arithmetic operators.....	9
5.4 Relational operators.....	9
5.5 Bit-wise operators.....	9
5.6 Assignment operators.....	10
5.7 Range notation.....	10
5.8 Mathematical functions.....	10
5.9 Array and strings operation functions.....	11
5.10 Order of operation precedence.....	11
5.11 Variables, syntax elements and tables.....	12
5.12 Text description of logical operators.....	13
5.13 Processes.....	15
5.13.1 General.....	15
5.13.2 Process output operators.....	15
5.14 Method of specifying syntax in tabular form.....	16
5.15 Bit ordering.....	17
5.16 Specification of syntax functions and data types.....	17
5.17 Semantics.....	17
6 Data Structures.....	18
6.1 General.....	18
6.2 Data unit.....	18
6.3 Annotation parameter set.....	19
6.3.1 General.....	19
6.3.2 Tile configuration.....	20
6.3.3 Annotation encoding parameters.....	23
6.3.4 Descriptor configuration.....	24
6.3.5 Compressor parameter set.....	31
6.3.6 Attribute parameter set.....	32
6.4 Annotation access unit.....	34
6.4.1 General.....	34
6.4.2 Annotation access unit header.....	35
6.4.3 Annotation access unit types.....	36
6.4.4 Block.....	37
7 Descriptors and attributes semantics.....	46
7.1 General.....	46
7.2 Descriptors.....	48
7.2.1 General.....	48
7.2.2 Genomic intervals.....	48
7.2.3 Genomic variants.....	48
7.2.4 Functional annotations.....	48
7.2.5 Contact matrices.....	48
7.3 Attributes.....	48
7.4 Data types.....	48
7.4.1 General.....	48

	7.4.2	Typed data.....	49
8		Decompression codecs.....	50
	8.1	General.....	50
	8.2	Inverse transformation algorithms.....	52
	8.2.1	General.....	52
	8.2.2	Lempel-Ziv-Welch transform.....	52
	8.2.3	Binarization transform.....	53
	8.2.4	Sparse transform.....	54
	8.2.5	Delta transform.....	55
	8.2.6	Run-Length Encoding transform.....	57
	8.2.7	Serialization transform.....	57
	8.3	Decompression algorithms.....	58
	8.3.1	General.....	58
	8.3.2	Context-Adaptive Binary Arithmetic Coding.....	59
	8.3.3	Lempel-Ziv-Markov Chain Algorithm.....	59
	8.3.4	Zstandard.....	59
	8.3.5	JBIG.....	59
	8.3.6	Block Sorting Coder.....	60
9		Decoding process.....	60
	9.1	General.....	60
	9.2	Access Units decoding process.....	60
	9.2.1	General.....	60
	9.2.2	Genomic variant access units.....	62
	9.2.3	Functional annotation Access Units.....	64
	9.2.4	Gene expression Access Units.....	65
	9.2.5	Position-to-position contact intensity Access Units.....	66
	9.2.6	Genome browser track Access Units.....	66
	9.3	Descriptors decoding process.....	67
	9.3.1	General.....	67
	9.3.2	Common descriptors.....	68
	9.3.3	Variant site information descriptors.....	70
	9.3.4	Functional annotation descriptors.....	73
	9.3.5	Genotype descriptor.....	75
	9.3.6	Likelihood descriptor.....	84
	9.3.7	Contact matrix descriptor.....	87
	9.4	Attributes decoding process.....	102
	9.5	Generic block payload decoding process.....	103
	9.5.1	Descriptor payload decoding process.....	103
	9.5.2	Attribute payload decoding process.....	104
10		Output format.....	107
	10.1	Variant site record.....	107
	10.1.1	General.....	107
	10.1.2	Semantics.....	108
	10.1.3	Initialization.....	110
	10.2	Variant genotype record.....	111
	10.2.1	General.....	111
	10.2.2	Semantics.....	112
	10.2.3	Initialization.....	113
	10.3	Sample record.....	114
	10.3.1	General.....	114
	10.3.2	Semantics.....	114
	10.3.3	Initialization.....	115
	10.4	Functional annotation record.....	115
	10.4.1	General.....	115
	10.4.2	Semantics.....	116
	10.4.3	Initialization.....	117
	10.5	Track property record.....	118

10.5.1	General	118
10.5.2	Semantics	119
10.5.3	Initialization	119
10.6	Track data record	120
10.6.1	General	120
10.6.2	Semantics	121
10.6.3	Initialization	121
10.7	Expression record	122
10.7.1	General	122
10.7.2	Semantics	123
10.7.3	Initialization	123
10.8	Feature record	124
10.8.1	General	124
10.8.2	Semantics	125
10.8.3	Initialization	125
10.9	Contact matrix record	126
10.9.1	General	126
10.9.2	Semantics	127
10.9.3	Initialization	128
Bibliography		129

Foreword

ISO (the International Organization for Standardization) and IEC (the International Electrotechnical Commission) form the specialized system for worldwide standardization. National bodies that are members of ISO or IEC participate in the development of International Standards through technical committees established by the respective organization to deal with particular fields of technical activity. ISO and IEC technical committees collaborate in fields of mutual interest. Other international organizations, governmental and non-governmental, in liaison with ISO and IEC, also take part in the work.

The procedures used to develop this document and those intended for its further maintenance are described in the ISO/IEC Directives, Part 1. In particular, the different approval criteria needed for the different types of document should be noted. This document was drafted in accordance with the editorial rules of the ISO/IEC Directives, Part 2 (see www.iso.org/directives or www.iec.ch/members_experts/refdocs).

ISO and IEC draw attention to the possibility that the implementation of this document may involve the use of (a) patent(s). ISO and IEC take no position concerning the evidence, validity or applicability of any claimed patent rights in respect thereof. As of the date of publication of this document, ISO and IEC had not received notice of (a) patent(s) which may be required to implement this document. However, implementers are cautioned that this may not represent the latest information, which may be obtained from the patent database available at www.iso.org/patents and <https://patents.iec.ch>. ISO and IEC shall not be held responsible for identifying any or all such patent rights.

Any trade name used in this document is information given for the convenience of users and does not constitute an endorsement.

For an explanation of the voluntary nature of standards, the meaning of ISO specific terms and expressions related to conformity assessment, as well as information about ISO's adherence to the World Trade Organization (WTO) principles in the Technical Barriers to Trade (TBT) see www.iso.org/iso/foreword.html. In the IEC, see www.iec.ch/understanding-standards.

This document was prepared by Joint Technical Committee ISO/IEC JTC 1, *Information technology*, Subcommittee SC 29, *Coding of audio, picture, multimedia and hypermedia information*.

A list of all parts in the ISO/IEC 23092 series can be found on the ISO and IEC websites.

Any feedback or questions on this document should be directed to the user's national standards body. A complete listing of these bodies can be found at www.iso.org/members.html and www.iec.ch/national-committees.

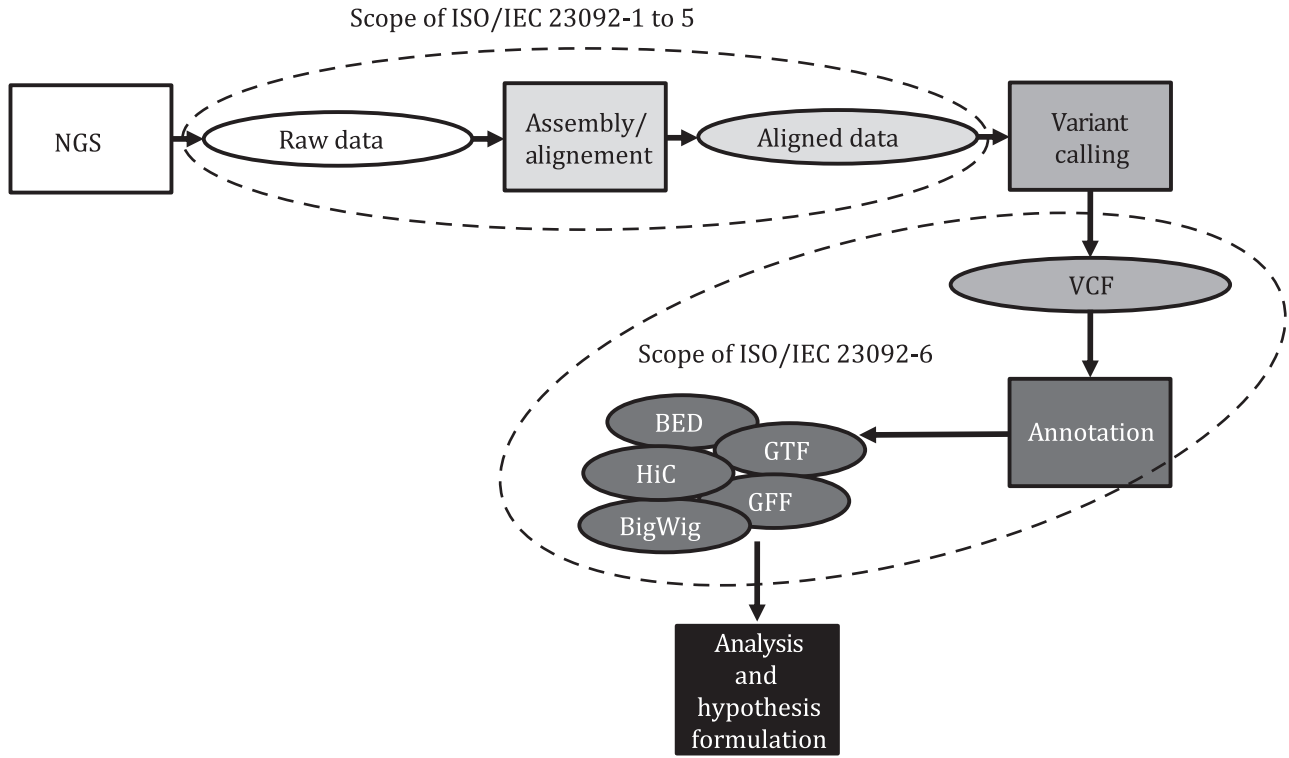
Introduction

While ISO/IEC 23092-1 to ISO/IEC 23092-5 (MPEG-G) deal with the representation of genomic information derived from the primary analysis of high-throughput sequencing (HTS) data – sequencing reads and qualities, and their alignment to a reference genome – which is only the first step in a long series. In particular, the results of primary analysis are usually processed further in order to obtain higher-level information. Such a process of aggregating information deduced from single reads and their alignments to the genome into more complex results is generally known as secondary analysis. In most HTS-based biological studies, the output of secondary analysis is usually represented as different types of annotations associated to one or more genomic intervals on the reference sequences.

Biological studies typically produce genomic annotation data such as mapping statistics, quantitative browser tracks, variants, genome functional annotations, gene expression data and Hi-C contact matrices. These diverse types of downstream genomic data are currently represented in different formats such as VCF, BED, WIG, etc., with loosely defined semantics, leading to issues with interoperability, the need for frequent conversions between formats, difficulty in the visualization of multi-modal data and complicated information exchange. [Figure 1](#) depicts a typical pipeline for the primary and secondary analyses of HTS data, the file formats involved and the scopes of different parts of the ISO/IEC 23092 series.

Furthermore, the lack of a single format has stifled the work on compression algorithms and has led to the widespread use of general compression algorithms with suboptimum performance. These algorithms do not exploit the fact the annotation data typically comprises of multiple fields (attributes) with different statistical characteristics and instead compress them together. Therefore, while these algorithms support efficient random access with respect to genomic position, they do not allow extraction of specific fields without decompressing all the whole file.

In response to the aforementioned challenges, this document details a unified data format for the efficient representation and compression of diverse genomic annotation data for file storage or data transport. The benefits are manifold: reducing the cost of data storage, improving the speed of random data access and processing, providing support for data security and privacy in selective genomic regions, and creating linkages across different types of genomic annotation and sequencing data. The ultimate goal is to enable the secured and seamless sharing, processing and analysis of multi-modal genomic data in order to reduce the burden of data manipulation and management, so scientists can focus on biological interpretation and discovery.



- Key**
- sequencing generates raw reads
 - read alignment
 - variant calling
 - variants annotations
 - analysis

Figure 1 — Typical pipeline for the primary and secondary analyses of HTS data

Information technology — Genomic information representation —

Part 6: Coding of genomic annotations

1 Scope

This document provides specifications for the normative representation of the following types of genomic information:

- variants with genotyping information
- functional annotations
- tracks
- expression matrices
- contact matrices (from Hi-C experiments or similar).

2 Normative references

The following documents are referred to in the text in such a way that some or all of their content constitutes requirements of this specification. For dated references, only the edition cited applies. For undated references, the latest edition of the referenced document (including any amendments) applies.

ISO/IEC 10646, *Information technology — Universal coded character set (UCS)*

ISO/IEC 11544, *Information technology — Coded representation of picture and audio information — Progressive bi-level image compression*

ISO/IEC 23092-1, *Information technology — Genomic information representation — Part 1: Transport and storage of genomic information*

ISO/IEC 23092-2, *Information technology — Genomic Information Representation — Part 2: Coding of Genomic Information*