

ISO/IEC FDIS 23092-3:2024(en)

ISO/IEC JTC 1/SC 29/AWG 11

Date: 2024-11-15

Secretariat: JISC

Date: 2025-02-13

Information technology — Genomic information representation

Part 3:
Metadata and application programming interfaces (APIs)

Technologie de l'information — Représentation des informations génomiques

Partie 3: Métadonnées et interfaces de programmation d'application (API)

ISO/IEC FDIS 23092-3

https://standards.itech.ai/catalog/standards/iso/c0471a33-f16a-4342-8899-c82386dbf1996/iso-iec-fdis-23092-3

FDIS stage

Edited DIS - MUST BE USED FOR FINAL DRAFT

Formatted: Left

Style Definition

Formatted: zzCover large

Formatted: Left: 1.5 cm, Right: 1.5 cm, Gutter: 0 cm, Section start: New page, Header distance from edge: 1.27 cm, Footer distance from edge: 1.27 cm

Formatted: Regular

Formatted: Cover Title_A2

Formatted

Formatted: French (France)

Formatted: Cover Title_B

Formatted: Left

ISO/IEC FDIS 23092-3:2025(en)

© ISO 2024, /IEC 2025 Published in Switzerland

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 • CP 401
CH-1214 Vernier, Geneva, Switzerland
Tel. Phone: + 41 22 749 01 11

Fax + 41 22 749 09 47

E-mail: copyright@iso.org
Website: www.iso.org

Published in Switzerland www.iso.org

Formatted: Default Paragraph Font

Formatted: Adjust space between Latin and Asian text, Adjust space between Asian text and numbers

Formatted: French (Switzerland)

Formatted: zzCopyright address, Adjust space between Latin and Asian text, Adjust space between Asian text and numbers

Formatted: French (Switzerland)

iTeh Standards
(<https://standards.iteh.ai>)
Document Preview

ISO/IEC FDIS 23092-3

<https://standards.iteh.ai/catalog/standards/iso/c0471a33-f16a-4342-8899-c82386dbf996/iso-iec-fdis-23092-3>

Contents

Foreword	vii
Introduction	ix
1 Scope	1
2 Normative references	1
3 Terms and definitions	2
4 Abbreviated terms	2
5 Conventions	2
6 Information metadata	4
7 Metrics metadata	21
8 Clinical data linkage metadata	24
9 Protection metadata	26
10 Access unit information	50
11 Decoding process for metadata	55
12 Application programming interfaces (APIs)	73
Annex A (normative) XML schemas corresponding to metadata information and protection elements	121
Annex B (informative) Example use cases of annotation table linkages metadata	123
Annex C (informative) XML schemas and XML-based data	125
Annex D (informative) Example of key transport	141
Annex E (informative) SAM interoperability	146
Bibliography	154
Foreword	5
Introduction	7
1 Scope	10
2 Normative references	10
3 Terms and definitions	11
4 Abbreviated terms	11
5 Conventions	11
5.1 Character encoding	11
5.2 Bit Ordering	11
5.3 Syntax functions and data types	12
5.4 Graphic notations	13
6 Information metadata	13
6.1 General	13
6.2 Dataset group metadata	13
6.3 Reference metadata	14
6.4 Dataset metadata	14
6.5 Annotation table metadata	16

Formatted: Adjust space between Latin and Asian text,
Adjust space between Asian text and numbers, Tab stops:
Not at 0.71 cm + 6.72 cm

6.5.1	General	16
6.5.2	Annotation table general metadata	17
6.5.3	Annotation table analytics metadata	20
6.5.4	Annotation table linkages metadata	24
6.5.5	Annotation table history metadata	25
6.6	Metadata protection	26
6.7	Mechanism for extensions of the metadata set	27
6.7.1	General	27
6.7.2	Example for dataset metadata extensions	27
6.7.3	Example for obfuscating labels	27
6.7.4	Example for obfuscating sequences	27
6.8	Metadata profiles	28
6.8.1	General	28
6.8.2	Example of metadata profile — Run	28
6.8.3	Example of metadata profile — Genomic data commons	29
7	Metrics metadata	29
7.1	Syntax	29
7.2	Semantics	30
8	Clinical data linkage metadata	32
8.1	General	32
8.2	CDL Metadata protection	34
9	Protection metadata	34
9.1	General	34
9.2	Encryption of <code>gen_info</code> elements and blocks	35
9.2.1	General	35
9.2.2	EncryptionParameters carried in dataset group protection	35
9.2.3	EncryptionParameters carried in dataset protection	36
9.2.4	EncryptionParameters carried in annotation table protection	41
9.2.5	Key retrieval	46
9.2.6	Decryption	47
9.3	Privacy rules for the use of the genomic information	50
9.3.1	General	50
9.3.2	Example of use of privacy rules	51
9.4	Digital signature of <code>gen_info</code> elements and blocks	52
9.4.1	General	52
9.4.2	Signatures carried in dataset group protection	52
9.4.3	Signatures carried in dataset protection	52
9.4.4	Signatures carried in annotation table protection	54
9.4.5	Signatures carried in descriptor stream protection	57
10	Access unit information	57
10.1	General	57
10.2	<code>genAuxRecord</code>	57
10.3	<code>genAux</code>	59
10.4	<code>genTag</code>	59
11	Decoding process for metadata	61
11.1	General	61
11.2	Initialization of parameters	63
11.2.1	General	63
11.2.2	Properties	63
11.2.3	Parameters	63
11.2.4	Constants	64
11.2.5	Process	65

11.3	Macros	67
11.4	Decoding process	69
12	Application programming interfaces (APIs)	77
12.1	General	77
12.2	Structure of the API	78
12.3	Detailed specification of the API	78
12.3.1	Data types	78
12.3.2	Return codes	79
12.3.3	Metadata fields	80
12.3.4	Output structures	80
12.3.5	Filters	90
12.3.6	Genomic information	99
12.3.7	Metadata	105
12.3.8	Protection	108
12.3.9	Reference	111
12.3.10	Statistics	112
Annex A (normative)	XML schemas corresponding to metadata information and protection elements	117
A.1	Dataset group metadata dgmmd XML schema	117
A.2	Dataset metadata dtmd XML schema	117
A.3	Annotation table metadata atmd XML schemas	117
A.4	Clinical data linkage metadata dged/dted/atmd XML schemas	117
A.5	Dataset group protection gen_info XML schema	117
A.6	Dataset protection gen_info XML schema	117
A.7	Annotation table protection gen_info XML schema	117
A.8	(Annotation) access unit protection gen_info XML schema	118
A.9	Descriptor stream protection gen_info XML schema	118
A.10	Dataset group reference metadata XML schema	118
Annex B (informative)	Example use cases of annotation table linkages metadata	119
B.1	Linkage between genomic variants and their supporting reads	119
B.2	Linkage between annotation tables to facilitate join query and data navigation	119
Annex C (informative)	XML schemas and XML-based data	121
C.1	General	121
C.2	EGA sample extension XML schema	121
C.3	EGA experiment extension XML schema	121
C.4	Genomic data commons extension XML schema	121
C.5	Labels obfuscation XML schema	121
C.6	Sequences obfuscation XML schema	121
C.7	Privacy rules and authorization requests	121
Annex D (informative)	Example of key transport	132
D.1	General	132
D.2	Key derivation example	133

D.3	Symmetric key wrap	134
D.4	Asymmetric key wrap	134
D.5	Chaining of KeyTransports	134
D.6	Multiple KeyTransports for the same key	135
Annex E (informative)	SAM interoperability	136
E.1	General	136
E.2	SAM Header	136
E.3	SAM auxiliary fields mapping	136
E.4	Mapping the SAM record to the MPEG-G record	138
E.5	SAM ambiguities resolution	139
Bibliography	143

iTeh Standards
(<https://standards.iteh.ai>)
Document Preview

<https://standards.iteh.ai/catalog/standards/iso/c0471a33-f16a-4342-8899-c82386dbf996/iso-iec-fdis-23092-3>

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 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*.

This third edition cancels and replaces the second edition (ISO/IEC 23092-3:2022), which has been technically revised.

The main changes are as follows:

- ~~The~~ addition of annotation table metadata (~~subclause 6.5~~)(subclause 6.5) that contains general analytics, linkages and access history information associated with an annotation table;
- ~~The~~ addition of metrics metadata (~~Clause 7~~)(Clause 7) that contains pre-computed sequencing data metrics associated with a dataset or an access unit;
- ~~The~~ addition of clinical data linkage metadata (~~Clause 8~~)(Clause 8) that contains linkage information for enabling clinical data interchange (CDI) with external data sources;
- ~~The~~ addition of annotation table protection metadata, including encryption parameters (~~subclause 9.2.4~~)(subclause 9.2.4) and digital signatures (~~subclause 9.4.4~~)(subclause 9.4.4), and updates to the decryption process (~~subclause 9.2.6~~)(subclause 9.2.6) and privacy rules (~~subclause 9.3~~)(subclause 9.3) for enabling the selective protection of annotation data;

Formatted: Adjust space between Latin and Asian text, Adjust space between Asian text and numbers

Formatted: English (United Kingdom)

Field Code Changed

Formatted: English (United Kingdom)

Field Code Changed

Formatted: Font: Cambria

Formatted: No bullets or numbering

ISO/IEC FDIS 23092-3:2025(en)

~~— The extension of the APIs (Clause 12)~~ **(Clause 12)** for supporting the random access and query of annotation data, the retrieval of pre-computed sequencing data statistics, and the return of only the number of matching records without the actual data.

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.

Formatted: English (United Kingdom)

Field Code Changed

iTeh Standards (<https://standards.iteh.ai>) Document Preview

ISO/IEC FDIS 23092-3

<https://standards.iteh.ai/catalog/standards/iso/c0471a33-f16a-4342-8899-c82386dbf996/iso-iec-fdis-23092-3>

Introduction

The advent of high-throughput sequencing (HTS) technologies has the potential to boost the adoption of genomic information in everyday practice, ranging from biological research to personalized genomic medicine in the clinic. As a consequence, the volume of generated data has increased dramatically during the last few years, and an even more pronounced growth is expected in the near future.

At the moment, genomic information is mostly exchanged through a variety of data formats, such as FASTA/FASTQ for unaligned sequencing reads and SAM/BAM/CRAM for aligned reads. With respect to such formats, the ISO/IEC 23092 series provides a new solution for the representation and compression of genome sequencing information by:

- specifying an abstract representation of the sequencing data rather than a specific format with its direct implementation;
- being designed at a time point when technologies and use cases are more mature. This permits the addressing of one limitation of the textual SAM format, for which incremental ad-hoc addition of features followed along the years, resulting in an overall redundant and suboptimal format which at the same time results not general and unnecessarily complicated;
- normatively separating free-field user-defined information with no clear semantics from the normative genomic data representation. This allows a fully interoperable and automatic exchange of information between different data producers;
- allowing multiplexing of relevant metadata information with the data since data and metadata are partitioned at different conceptual levels;
- following a strict and supervised development process which has proven successful in the last 30 years in the domain of digital media for the transport format, the file format, the compressed representation and the application program interfaces.

This document provides the enabling technology that will allow the community to create an ecosystem of novel, interoperable solutions in the field of genomic information processing. In particular, it offers:

- consistent, general and properly designed format definitions and data structures to store sequencing and alignment information. A robust framework which can be used as a foundation to implement different compression algorithms;
- speed and flexibility in the selective access to coded data, by means of newly designed data clustering and optimized storage methodologies;
- low latency in data transmission and consequent fast availability at remote locations, based on transmission protocols inspired by real-time application domains;
- built-in privacy and protection of sensitive information, thanks to a flexible framework which allows customizable secured access at all layers of the data hierarchy;
- reliability of the technology and interoperability among tools and systems, owing to the provision of a normative procedure to assess conformance to the standard on an exhaustive dataset;
- support to the implementation of a complete ecosystem of compliant devices and applications, through the availability of a normative reference implementation covering the totality of the specification.

The fundamental structure of the ISO/IEC 23092 series data representation is the genomic record. The genomic record is a data structure consisting of either a single sequence read, or a paired sequence read, and

Formatted: Adjust space between Latin and Asian text, Adjust space between Asian text and numbers

Formatted: Default Paragraph Font

Formatted: Default Paragraph Font

Formatted: Default Paragraph Font

Formatted: Adjust space between Latin and Asian text, Adjust space between Asian text and numbers, Tab stops: Not at 0.7 cm + 1.4 cm + 2.1 cm + 2.8 cm + 3.5 cm + 4.2 cm + 4.9 cm + 5.6 cm + 6.3 cm + 7 cm

Formatted: Adjust space between Latin and Asian text, Adjust space between Asian text and numbers

Formatted: Adjust space between Latin and Asian text, Adjust space between Asian text and numbers, Tab stops: Not at 0.7 cm + 1.4 cm + 2.1 cm + 2.8 cm + 3.5 cm + 4.2 cm + 4.9 cm + 5.6 cm + 6.3 cm + 7 cm

Formatted: Adjust space between Latin and Asian text, Adjust space between Asian text and numbers

Formatted: Default Paragraph Font

Formatted: Default Paragraph Font

Formatted: Default Paragraph Font

Formatted: Font: Not Italic

ISO/IEC FDIS 23092-3:2025(en)

its associated sequencing and alignment information; it may contain detailed mapping and alignment data, a single or paired read identifier (read name) and quality values.

Without breaking traditional approaches, the genomic record introduced in the ISO/IEC 23092 series provides a more compact, simpler and manageable data structure grouping all the information related to a single DNA template, from simple sequencing data to sophisticated alignment information.

The genomic record, although it is an appropriate logic data structure for interaction and manipulation of coded information, is not a suitable atomic data structure for compression. To achieve high compression ratios, it is necessary to group genomic records into clusters and to transform the information of the same type into sets of descriptors structured into homogeneous blocks. Furthermore, when dealing with selective data access, the genomic record is a too small unit to allow effective and fast information retrieval.

For these reasons, this document introduces the concept of access unit, which is the fundamental structure for coding and access to information in the compressed domain.

The access unit is the smallest data structure that can be decoded by a decoder ~~compliant with~~ confirming to ISO/IEC 23092-2. An access unit is composed of one block for each descriptor used to represent the information of its genomic records; therefore, a block payload is the coded representation of all the data of the same type (i.e. a descriptor) in a cluster.

In addition to clusters of genomic records compressed into access units, reads are further classified in six data classes: five classes are defined according to the result of their alignment against one or more reference sequences; the sixth class contains either reads that could not be mapped or raw sequencing data. The classification of sequence reads into classes enables to develop powerful selective data access. In fact, access units inherit a specific data characterization (e.g. perfect matches in Class P, substitutions in Class M, indels in Class I, half-mapped reads in Class HM) from the genomic records composing them, and thus constitute a data structure capable of providing powerful filtering capability for the efficient support of many different use cases.

Access units are the fundamental, finest grain data structure in terms of content protection and in terms of metadata association. In other words, each access unit can be protected individually and independently. ~~Figure 1~~Figure 1 shows how access units, blocks and genomic records relate to each other in the ISO/IEC 23092 series data structure.

Formatted: Default Paragraph Font

Formatted: Default Paragraph Font

Formatted: Default Paragraph Font

Formatted: Default Paragraph Font

Formatted: Default Paragraph Font

Formatted: Default Paragraph Font

Formatted: Default Paragraph Font

Formatted: Default Paragraph Font

Formatted: Default Paragraph Font

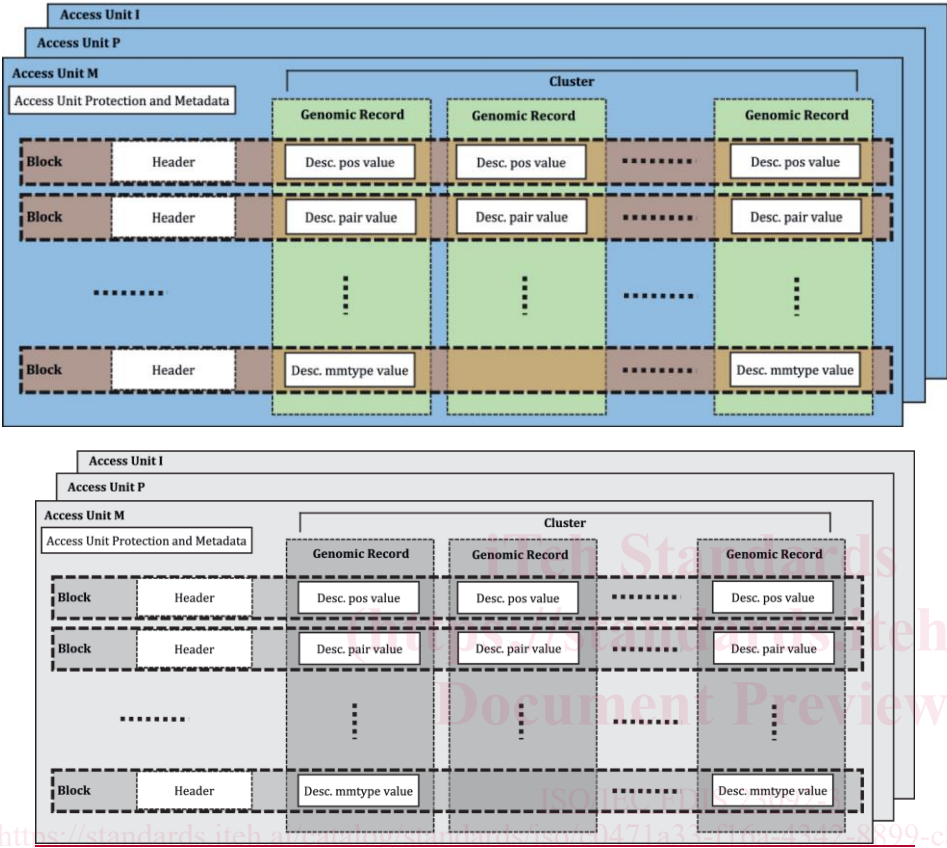


Figure 1.— Access units, blocks and genomic records

Formatted: Space Before: 12 pt, Adjust space between Latin and Asian text, Adjust space between Asian text and numbers

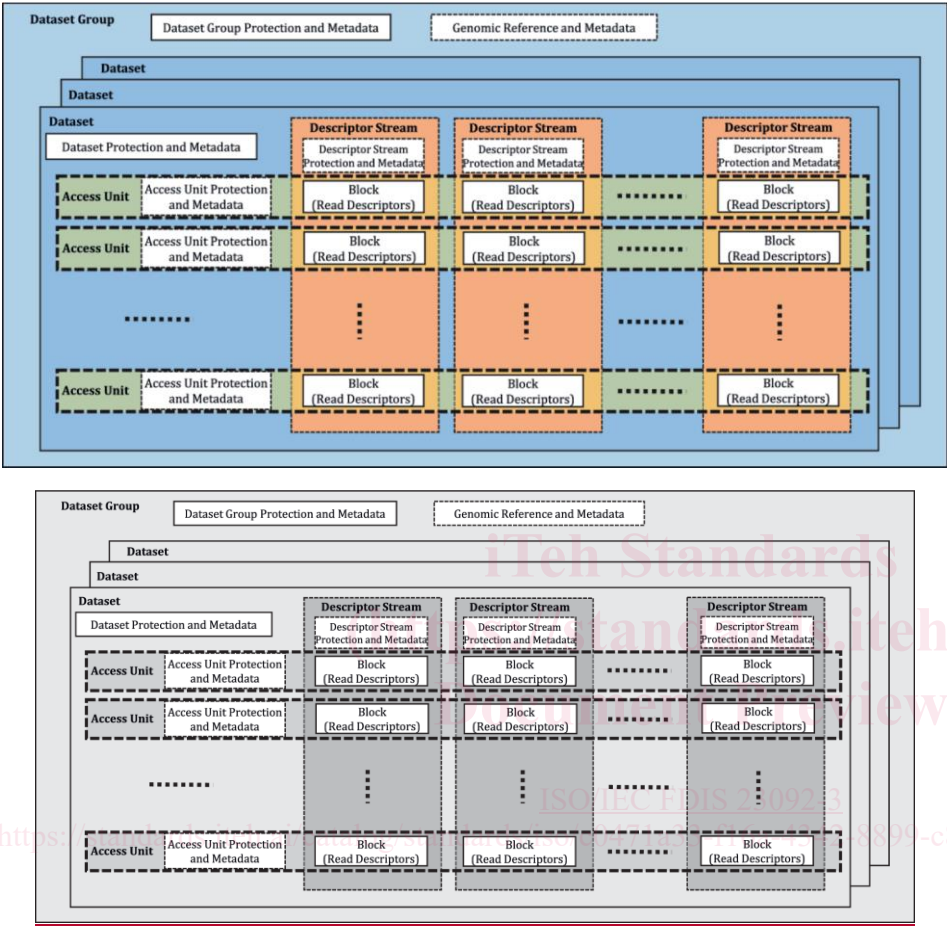


Figure 2.— High-level data structure: datasets and dataset group

Formatted: Adjust space between Latin and Asian text, Adjust space between Asian text and numbers

A dataset is a coded data structure containing headers and one or more access units. Typical datasets ~~could~~can, for example, contain the complete sequencing of an individual, or a portion of it. Other datasets ~~could~~can contain, for example, a reference genome or a subset of its chromosomes. Datasets are grouped in dataset groups, as shown in ~~Figure 2~~Figure 2.

A simplified diagram of the dataset decoding process is shown in ~~Figure 3~~Figure 3.

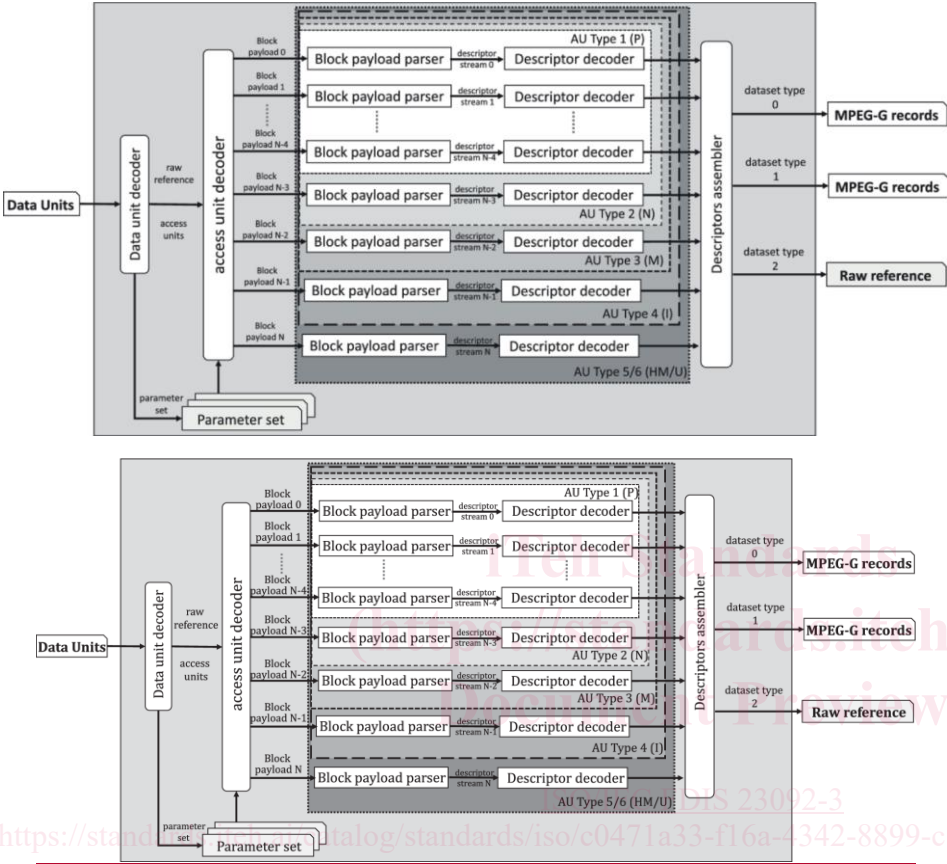


Figure 3.— Decoding process

Formatted: Adjust space between Latin and Asian text,
Adjust space between Asian text and numbers

