

FINAL
DRAFT

INTERNATIONAL
STANDARD

ISO/IEC
FDIS
15938-16

ISO/IEC JTC 1/SC 29

Secretariat: JISC

Voting begins on:
2021-03-09

Voting terminates on:
2021-05-04

Information technology — Multimedia content description interface —

Part 16: Conformance and reference software for compact descriptors for video analysis

iTeh STANDARD PREVIEW
(standards.iteh.ai)

[ISO/IEC FDIS 15938-16](https://standards.iteh.ai/catalog/standards/sist/fc65c6d0-f3b7-4be4-8440-ab83f84ac515/iso-iec-fdis-15938-16)

<https://standards.iteh.ai/catalog/standards/sist/fc65c6d0-f3b7-4be4-8440-ab83f84ac515/iso-iec-fdis-15938-16>

RECIPIENTS OF THIS DRAFT ARE INVITED TO SUBMIT, WITH THEIR COMMENTS, NOTIFICATION OF ANY RELEVANT PATENT RIGHTS OF WHICH THEY ARE AWARE AND TO PROVIDE SUPPORTING DOCUMENTATION.

IN ADDITION TO THEIR EVALUATION AS BEING ACCEPTABLE FOR INDUSTRIAL, TECHNOLOGICAL, COMMERCIAL AND USER PURPOSES, DRAFT INTERNATIONAL STANDARDS MAY ON OCCASION HAVE TO BE CONSIDERED IN THE LIGHT OF THEIR POTENTIAL TO BECOME STANDARDS TO WHICH REFERENCE MAY BE MADE IN NATIONAL REGULATIONS.



Reference number
ISO/IEC FDIS 15938-16:2021(E)

© ISO/IEC 2021

iTeh STANDARD PREVIEW (standards.iteh.ai)

ISO/IEC FDIS 15938-16

<https://standards.iteh.ai/catalog/standards/sist/fc65c6d0-f3b7-4be4-8440-ab83f84ac515/iso-iec-fdis-15938-16>



COPYRIGHT PROTECTED DOCUMENT

© ISO/IEC 2021

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 | iv |
| Introduction | v |
| 1 Scope | 1 |
| 2 Normative references | 1 |
| 3 Terms and definitions | 1 |
| 4 Abbreviated terms, relational operators, bitwise operators and functions | 1 |
| 5 Conformance testing | 1 |
| 5.1 Conformance video dataset and reference bitstreams | 1 |
| 5.2 Conformance test conditions for strict conformance | 2 |
| 5.2.1 Test setup | 2 |
| 5.2.2 Test conditions | 2 |
| 5.2.3 Global descriptor test | 2 |
| 5.2.4 Local descriptor test | 2 |
| 5.2.5 Deep feature descriptor test with default NN | 2 |
| 5.2.6 Combined global and local descriptor test | 2 |
| 5.2.7 Combined global and deep feature descriptor test with default NN | 3 |
| 5.3 Conformance test conditions using custom NN | 3 |
| 6 Reference software | 3 |
| 6.1 Reference software location | 3 |
| 6.2 Reference software licence | 3 |
| 6.3 Reference software documentation | 3 |
| 6.4 Reference software installation and compilation | 3 |
| 6.5 Reference software architecture | 3 |
| 6.5.1 General | 3 |
| 6.5.2 Extraction | 4 |
| 6.5.3 Pairwise matching | 6 |
| 6.5.4 Retrieval | 7 |
| 6.5.5 Parameters | 8 |
| 6.5.6 Implementation details for CDVA reference software | 9 |
| 6.5.7 Usage | 9 |
| Annex A (informative) Deployment of the reference software using Docker containers | 13 |

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).

Attention is drawn to the possibility that some of the elements of this document may be the subject of patent rights. ISO and IEC shall not be held responsible for identifying any or all such patent rights. Details of any patent rights identified during the development of the document will be in the Introduction and/or on the ISO list of patent declarations received (see www.iso.org/patents) or the IEC list of patent declarations received (see patents.iec.ch).

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 15938 series can be found on the ISO website.

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

ISO/IEC 15938 (all parts) provides a standardized set of technologies for describing multimedia content. It addresses a broad spectrum of multimedia applications and requirements by providing a metadata system for describing the features of multimedia content.

The following are specified in ISO/IEC 15938 (all parts):

Description schemes (DS) describe entities or relationships pertaining to multimedia content. Description schemes specify the structure and semantics of their components, which can be description schemes, descriptors or datatypes.

Descriptors (D) describe features, attributes or groups of attributes of multimedia content.

Datatypes are the basic reusable datatypes employed by description schemes and descriptors.

Description definition language (DDL) defines description schemes, descriptors and datatypes by specifying their syntax, and allows their extension.

Systems tools support delivery of descriptions, multiplexing of descriptions with multimedia content, synchronization, file format, etc.

ISO 15938 (all parts) is subdivided into 16 published parts with further parts in development:

- **Part 1: Systems:** specifies the tools for preparing descriptions for efficient transport and storage, compressing descriptions, and allowing synchronization between content and descriptions.
- **Part 2: Description definition language:** specifies the language for defining the series set of description tools (DSs, Ds and datatypes) and for defining new description tools.
- **Part 3: Visual:** specifies the description tools pertaining to visual content.
- **Part 4: Audio:** specifies the description tools pertaining to audio content.
- **Part 5: Multimedia description schemes:** specifies the generic description tools pertaining to multimedia including audio and visual content.
- **Part 6: Reference software:** provides a software implementation of the series.
- **Part 7: Conformance testing:** specifies the guidelines and procedures for testing conformance of implementations of the series.
- **Part 8: Extraction and use of MPEG-7 descriptions:** provides guidelines and examples of the extraction and use of descriptions.
- **Part 9: Profiles and levels:** provides guidelines and standard profiles.
- **Part 10: Schema definition:** specifies the schema using description definition language.
- **Part 11: MPEG-7 profile schemas:** listing of profile schemas using description definition language.
- **Part 12: Query format:** contains the tools of the MPEG query format (MPQF).
- **Part 13: Compact descriptors for visual search:** specifies an image description tool for visual search applications.
- **Part 14: Reference software, conformance and usage guidelines for compact descriptors for visual search:** provides the reference software and guidelines, specifies the conformance testing.
- **Part 15: Compact descriptors for video analysis:** specifies a video description tool designed to enable efficient and interoperable video analysis applications, allowing visual content matching in videos.

- **Part 16** (this document): **Conformance and reference software for compact descriptors for video analysis**: describes conformance testing of descriptors specified in ISO/IEC 15938-15, and provides the reference software for extracting and matching such descriptors.

The structure of this document is as follows:

- [Clause 5](#) specifies conformance testing of descriptors.
- [Clause 6](#) describes the reference software for extracting and matching descriptors. The deployment of the reference software using Docker containers is described in [Annex A](#).

iTeh STANDARD PREVIEW (standards.iteh.ai)

ISO/IEC FDIS 15938-16

<https://standards.iteh.ai/catalog/standards/sist/fc65c6d0-f3b7-4be4-8440-ab83f84ac515/iso-iec-fdis-15938-16>

Information technology — Multimedia content description interface —

Part 16: Conformance and reference software for compact descriptors for video analysis

1 Scope

This document specifies the assessment of conformance to ISO/IEC 15938-15 as well as the reference software.

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 document. For dated references, only the edition cited applies. For undated references, the latest edition of the referenced document (including any amendments) applies.

ISO/IEC 15938-14, *Information technology — Multimedia content description interface — Part 14: Reference software, conformance and usage guidelines for compact descriptors for visual search*

ISO/IEC 15938-15, *Information technology — Multimedia content description interface — Part 15: Compact descriptors for video analysis*

<https://standards.iteh.ai/catalog/standards/sist/fc65c6d0-f3b7-4be4-8440-ab83f84ac515/iso-iec-fdis-15938-16>

3 Terms and definitions

For the purposes of this document, the terms and definitions in ISO/IEC 15938-15 apply.

ISO and IEC maintain terminological databases for use in standardization at the following addresses:

- ISO Online browsing platform: available at <https://www.iso.org/obp>
- IEC Electropedia: available at <http://www.electropedia.org/>

4 Abbreviated terms, relational operators, bitwise operators and functions

For the purposes of this document, the abbreviated terms, relational operators, bitwise operators and functions in ISO/IEC 15938-15 apply.

5 Conformance testing

5.1 Conformance video dataset and reference bitstreams

The CDVA conformance data set, containing 3 318 items, is used for conformance testing.

The data set is available at <https://standards.iso.org/iso-iec/15938/-16/ed-1/en>.

For deep feature descriptors, two levels of conformance are considered:

- Conformance using default neural network: This is a stricter level of conformance, which ensures that the bitstreams match those generated with the reference software as-is.

- Conformance using custom neural network: This is a less strict level of conformance, that allows using a custom neural network. The deep feature descriptors generated with implementations using different neural networks are not interoperable.

5.2 Conformance test conditions for strict conformance

5.2.1 Test setup

- Bitstreams extracted with the reference software on Ubuntu 16.04 are provided.
- Implementation under test is used to extract bitstreams from the conformance test set, including all three descriptor components, and varying the parameters listed in subclause [6.5.4](#).
- Use the reference implementation to perform pairwise matching between the bitstreams of corresponding videos.
- Check that the resulting similarity is above the thresholds listed in the following subclauses.

5.2.2 Test conditions

Vary *minLocalDiff* for the local feature descriptor in the conformance tests. At least the following values for *minLocalDiff* should be tested: 1; 31.

5.2.3 Global descriptor test

Perform pairwise matching of global descriptors.

An implementation is considered conforming to the specification, if at least 99% of the descriptors of the conformance dataset yield a score greater or equal to 0,9s, when matching bitstreams extracted with the implementation under test with the reference bitstreams for the same video. *s* is the matching score for the same pair of videos when matching the reference bitstreams to itself using the reference implementation.

5.2.4 Local descriptor test

Perform pairwise matching of local descriptors.

An implementation is considered conforming to the specification if at least 99% of the descriptors of the conformance video dataset yield a score greater or equal to 0,9s, when matching bitstreams extracted with the implementation under test with the reference bitstreams for the same video. *s* is the matching score for the same pair of videos when matching the reference bitstreams to itself using the reference implementation.

5.2.5 Deep feature descriptor test with default NN

Perform pairwise matching of deep feature descriptors, using the default NN.

An implementation is considered conforming to the specification if at least 99% of the descriptors of the conformance video dataset yield a score greater or equal to 0,9s, when matching bitstreams extracted with the implementation under test with the reference bitstreams for the same video. *s* is the matching score for the same pair of videos when matching the reference bitstreams to itself using the reference implementation.

5.2.6 Combined global and local descriptor test

Perform pairwise matching combining local and global descriptors.

An implementation is considered conforming to the specification if at least 99% of the descriptors of the conformance video dataset yield a score greater or equal to 0,9s, when matching bitstreams extracted

with the implementation under test with the reference bitstreams for the same video. s is the matching score for the same pair of videos when matching the reference bitstreams to itself using the reference implementation.

5.2.7 Combined global and deep feature descriptor test with default NN

Perform pairwise matching combining local and global descriptors, using the default NN.

An implementation is considered conforming to the specification if at least 99% of the descriptors of the conformance video dataset yield a score greater or equal to $0,9s$, when matching bitstreams extracted with the implementation under test with the reference bitstreams for the same video. s is the matching score for the same pair of videos when matching the reference bitstreams to itself using the reference implementation.

5.3 Conformance test conditions using custom NN

As a prerequisite for this test, at least the conformance tests for global descriptors shall be run.

No reference data is provided for this test.

Only syntactic checking can be performed by using the reference software to parse and match bitstreams generated with the implementation under test. This shall complete without errors.

6 Reference software

6.1 Reference software location

The reference software is available at <https://standards.iso.org/iso-iec/15938/-16/ed-1/en>.
ISO/IEC FDIS 15938-16

6.2 Reference software licence

The licence of the software is described in the file `CDVA_evaluation_framework/COPYING` included with the reference software.

6.3 Reference software documentation

The documentation of the software is contained in the folder `CDVA_evaluation_framework/docs` included with the reference software. An overview is provided in the file `CDVA_evaluation_framework/README`.

6.4 Reference software installation and compilation

The installation of the software is described in file `CDVA_evaluation_framework/INSTALLING` included with the reference software.

The CDVA reference software is entirely written in C and C++. On Windows, it has been compiled and tested on Windows 10 Enterprise 64-bit using Visual C++ 2015 and 2017 (64 bit). On Linux, it has been compiled and tested on Ubuntu 16.04 LTS (64-bit).

The CDVA reference software has a number of dependencies, in particular due to requiring feature extraction with CNNs. The build process including all dependencies is thus not trivial.

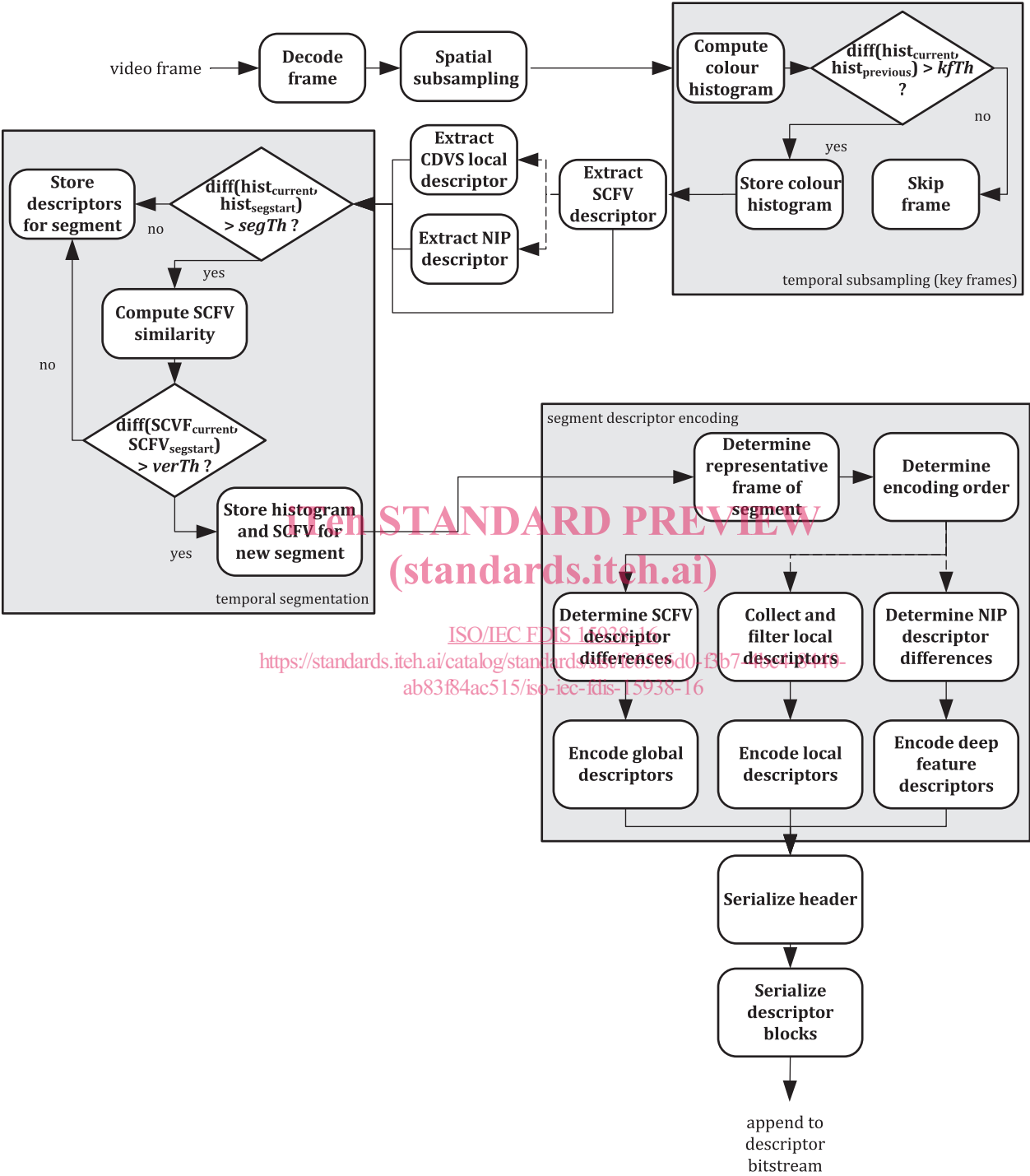
6.5 Reference software architecture

6.5.1 General

The software implementation uses the CDVS reference software (as published in ISO/IEC 15938-14) for extraction of the global and local descriptor components. It includes modules for temporal sampling

and segmentation, extraction of deep features and temporal encoding of global and local CDVS features and deep features.

6.5.2 Extraction



NOTE Dashed lines indicate optional steps in the process.

Figure 1 — CDVA descriptor extractor

Figure 1 illustrates how the CDVA reference software produces a compact descriptor of a video segment in a series of processing steps, when a video frame is given in input to the system. The process is repeated for all frames in the input video segment. The resulting descriptors are grouped by segments.

Segments are represented with their first key frame being identified as segment boundary. All following frames until the next segment boundary or the end of the file belong to the same segment. The output descriptor is updated by appending the output of the CDVA extraction process for a segment to a single CDVA descriptor.

1. **Frame subsampling:** Performs temporal frame subsampling (typically by a factor of 2–10).
2. **Decode frame:** Decode a frame present in the video.
3. **Compute colour histogram:** A histogram of the R,G,B planes is computed, using 32 bins for each plane.
4. **Check the difference between current and previous colour histograms:** If the difference is greater than a given threshold ($kfTh$), the frame is selected as *keyframe* and further processed. If not, the frame is dropped.
5. **Frame drop module:** If, according to step 4, the current frame is similar to the previously encoded one, the current frame is dropped.
6. **Store colour histogram:** The colour histogram is stored in memory, to be used as “previous histogram” in the next iteration.
7. **Extract SCFV descriptor:** Extracts the SCFV descriptors from individual frames, using mode 0 of ISO/IEC 15938-14.
8. **Extract CDVS local descriptor:** Extracts the CDVS local descriptors from individual frames, using mode 0 of ISO/IEC 15938-14.
9. **Extract deep feature descriptor:** Extract the deep feature descriptor using a pretrained VGG16 network, and optionally binarize the resulting descriptor.
10. **Check the difference of colour histograms between current frame and segment start:** If the difference is greater than a given threshold ($segTh$), the frame is selected as candidate of a segment boundary (first frame of a new segment).
11. **Compute global SCFV similarity between current descriptor and descriptor of first frame of segment:** If the similarity is below a given threshold ($verTh$), the frame is confirmed to be a segment boundary.
12. **Store colour histogram and SCFV descriptor:** For frames identified as segment boundaries, store the descriptors for comparisons with the subsequent frames.
13. **Store descriptors for segment:** Store all extracted descriptors for the current segment together with their time index.
14. **Determine representative frame:** Select a representative frame for the segment, selected as the frame with the medoid SCFV descriptor. This frame is used as the reference for encoding the global and local descriptor of the segment.
15. **Determine encoding order:** The encoding order is determined from the SCFV descriptors of the key frames in the segment. The first is the representative frame, followed by frames with decreasing distance to any of the frames encoded so far.
16. **Determine global descriptor differences:** For each key frame other than the representative frame, determine a difference descriptor as XOR between its SCFV descriptor and the SCFV descriptor of the reference frame.
17. **Determine deep feature descriptor differences:** For each key frame other than the representative frame, determine a difference descriptor as XOR between its NIP descriptor and the NIP descriptor of the reference frame.