



SLOVENSKI STANDARD

SIST ISO 28500:2018

01-september-2018

Nadomešča:
SIST ISO 28500:2009

Informatika in dokumentacija - Datotečna oblika zapisa WARC

Information and documentation -- WARC file format

iTeh STANDARD PREVIEW
Information et documentation -- Format de fichier WARC
(standards.iteh.ai)

Ta slovenski standard je istoveten z ~~SIST ISO 28500~~ ISO 28500:2017

<https://standards.iteh.ai/catalog/standards/sist/c3ac43f1-47a5-40e9-815d-a6fed47e1f56/sist-iso-28500-2018>

ICS:

35.240.30	Uporabniške rešitve IT v informatiki, dokumentiranju in založništvu	IT applications in information, documentation and publishing
-----------	---------------------------------------------------------------------	--------------------------------------------------------------

SIST ISO 28500:2018

en,fr,de

iTeh STANDARD PREVIEW
(standards.iteh.ai)

SIST ISO 28500:2018

<https://standards.iteh.ai/catalog/standards/sist/c3ac43f1-47a5-40e9-815d-a6fed47e1f56/sist-iso-28500-2018>

INTERNATIONAL
STANDARD

ISO
28500

Second edition
2017-08

**Information and documentation —
WARC file format**

Information et documentation — Format de fichier WARC

**iTeh STANDARD PREVIEW
(standards.iteh.ai)**

SIST ISO 28500:2018

<https://standards.iteh.ai/catalog/standards/sist/c3ac43f1-47a5-40e9-815d-a6fed47e1f56/sist-iso-28500-2018>



Reference number
ISO 28500:2017(E)

© ISO 2017

iTeh STANDARD PREVIEW
(standards.iteh.ai)

SIST ISO 28500:2018

<https://standards.iteh.ai/catalog/standards/sist/c3ac43f1-47a5-40e9-815d-a6fed47e1f56/sist-iso-28500-2018>



COPYRIGHT PROTECTED DOCUMENT

© ISO 2017, Published in Switzerland

All rights reserved. Unless otherwise specified, 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
Ch. de Blandonnet 8 • CP 401
CH-1214 Vernier, Geneva, Switzerland
Tel. +41 22 749 01 11
Fax +41 22 749 09 47
copyright@iso.org
www.iso.org

Contents

Page

Foreword	v
Introduction	vi
1 Scope	1
2 Normative references	1
3 Terms, definitions and abbreviated terms	2
4 File and record model	3
5 Named fields	5
5.1 General.....	5
5.2 WARC-Record-ID (mandatory).....	5
5.3 Content-Length (mandatory).....	5
5.4 WARC-Date (mandatory).....	6
5.5 WARC-Type (mandatory).....	6
5.6 Content-Type.....	6
5.7 WARC-Concurrent-To.....	7
5.8 WARC-Block-Digest.....	7
5.9 WARC-Payload-Digest.....	7
5.10 WARC-IP-Address.....	8
5.11 WARC-Refers-To.....	8
5.12 WARC-Refers-To-Target-URI.....	8
5.13 WARC-Refers-To-Date.....	8
5.14 WARC-Target-URI.....	9
5.15 WARC-Truncated.....	9
5.16 WARC-Warcinfo-ID.....	9
5.17 WARC-Filename.....	9
5.18 WARC-Profile.....	10
5.19 WARC-Identified-Payload-Type.....	10
5.20 WARC-Segment-Number.....	10
5.21 WARC-Segment-Origin-ID.....	10
5.22 WARC-Segment-Total-Length.....	10
6 WARC record types	11
6.1 General.....	11
6.2 'warcinfo'.....	11
6.3 'response'.....	11
6.3.1 General.....	11
6.3.2 'http' and 'https' schemes.....	12
6.3.3 Other URI schemes.....	12
6.4 'resource'.....	12
6.4.1 General.....	12
6.4.2 'http' and 'https' schemes.....	12
6.4.3 'ftp' scheme.....	12
6.4.4 'dns' scheme.....	13
6.4.5 Other URI schemes.....	13
6.5 'request'.....	13
6.5.1 General.....	13
6.5.2 'http' and 'https' schemes.....	13
6.5.3 Other URI schemes.....	13
6.6 'metadata'.....	13
6.7 'revisit'.....	14
6.7.1 General.....	14
6.7.2 Profile: Identical Payload Digest.....	14
6.7.3 Profile: Server Not Modified.....	15
6.7.4 Other profiles.....	15

ISO 28500:2017(E)

6.8	'conversion'	15
6.9	'continuation'	16
7	Record segmentation	16
8	WARC file name, size and compression	16
Annex A (informative)	Use cases for writing WARC records	18
Annex B (informative)	Examples of WARC records	21
Annex C (informative)	WARC file size and name recommendations	24
Annex D (informative)	Compression recommendations	25
Bibliography	26

iTeh STANDARD PREVIEW
(standards.iteh.ai)

SIST ISO 28500:2018

<https://standards.iteh.ai/catalog/standards/sist/c3ac43f1-47a5-40e9-815d-a6fed47e1f56/sist-iso-28500-2018>

Foreword

ISO (the International Organization for Standardization) is a worldwide federation of national standards bodies (ISO member bodies). The work of preparing International Standards is normally carried out through ISO technical committees. Each member body interested in a subject for which a technical committee has been established has the right to be represented on that committee. International organizations, governmental and non-governmental, in liaison with ISO, also take part in the work. ISO collaborates closely with the International Electrotechnical Commission (IEC) on all matters of electrotechnical standardization.

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 ISO documents 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).

Attention is drawn to the possibility that some of the elements of this document may be the subject of patent rights. ISO 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).

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

For an explanation on 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 the following URL: www.iso.org/iso/foreword.html.

This document was prepared by Technical Committee ISO/TC 46, *Information and documentation*, Subcommittee 4, *Technical interoperability*.

This second edition cancels and replaces the first edition (ISO 28500:2009), which has been technically revised.

ISO 28500:2017(E)

Introduction

Websites and web pages emerge and disappear from the World Wide Web every day. For the past 10 years, memory storage organizations have tried to find the most appropriate ways to collect and keep track of this vast quantity of important material using web-scale tools such as web crawlers. A web crawler is a program that browses the web in an automated manner according to a set of policies; starting with a list of URLs, it saves each page identified by a URL, finds all the hyperlinks in the page (e.g. links to other pages, images, videos, scripting or style instructions, etc.), and adds them to the list of URLs to visit recursively. Storing and managing the billions of saved web page objects itself presents a challenge.

At the same time, those same organizations have a rising need to archive large numbers of digital files not necessarily captured from the web (e.g. entire series of electronic journals, or data generated by environmental sensing equipment). A general requirement that appears to be emerging is for a container format that permits one file simply and safely to carry a very large number of constituent data objects for the purpose of storage, management, and exchange. Those data objects (or resources) need to be of unrestricted type (including many binary types for audio, CAD, compressed files, etc.), but fortunately the container needs only minimal knowledge of the nature of the objects.

The WARC (Web ARChive) file format offers a convention for concatenating multiple resource records (data objects), each consisting of a set of simple text headers and an arbitrary data block into one long file. The WARC format is an extension of the ARC file format (ARC) that has traditionally been used to store “web crawls” as sequences of content blocks harvested from the World Wide Web. Each capture in an ARC file is preceded by a one-line header that very briefly describes the harvested content and its length. This is directly followed by the retrieval protocol response messages and content. The original ARC format file has been used by the Internet Archive (IA) since 1996 for managing billions of objects, and by several national libraries.

The motivation to extend the ARC format arose from the discussion and experiences of the International Internet Preservation Consortium (IIPC), whose members include the national libraries of Australia, Canada, Denmark, Finland, France, Iceland, Italy, Norway, Sweden, The British Library (UK), The Library of Congress (USA), and the Internet Archive (IA). The California Digital Library and the Los Alamos National Laboratory also provided input on extending and generalizing the format.

The WARC format offers a standard way to structure, manage and store billions of resources collected from the web and elsewhere. It is used to build applications for harvesting, managing, accessing, mining and exchanging content. While it represents the unique standard format for web archives, it has been adopted beyond the web archiving community to store born-digital or digitized materials. The way WARC files will be created and resources stored and rendered will depend on software and applications implementations.

Besides the primary content recorded in ARCs, the extended WARC format accommodates related secondary content, such as assigned metadata, abbreviated duplicate detection events, later-date transformations, and segmentation of large resources. The extension may also be useful for more general applications than web archiving. To aid the development of tools that are backwards compatible, WARC content is clearly distinguishable from pre-revision ARC content.

The WARC file format is made sufficiently different from the legacy ARC format files so that software tools can unambiguously detect and correctly process both WARC and ARC records; given the large amount of existing archival data in the previous ARC format, it is important that access and use of this legacy not be interrupted when transitioning to the WARC format.

Information and documentation — WARC file format

1 Scope

This document specifies the WARC file format:

- to store both the payload content and control information from mainstream Internet application layer protocols, such as the HTTP, DNS, and FTP;
- to store arbitrary metadata linked to other stored data (e.g. subject classifier, discovered language, encoding);
- to support data compression and maintain data record integrity;
- to store all control information from the harvesting protocol (e.g. request headers), not just response information;
- to store the results of data transformations linked to other stored data;
- to store a duplicate detection event linked to other stored data (to reduce storage in the presence of identical or substantially similar resources);
- to be extended without disruption to existing functionality;
- to support handling of overly long records by truncation or segmentation, where desired.

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.

RFC1035¹⁾ Mockapetris, P. *Domain names — Implementation and specification*, STD 13, November 1987

RFC2045²⁾ Freed, N. and Borenstein, N. *Multipurpose Internet Mail Extensions (MIME) Part One: Format of Internet Message Bodies*, November 1996

RFC2540³⁾ Eastlake, D. *Detached Domain Name System (DNS) Information*, March 1999

RFC2616⁴⁾ Fielding, R., Gettys, J., Mogul, J., Frystyk, H., Masinter, L., Leach, P. and Berners-Lee, T. *Hypertext Transfer Protocol — HTTP/1.1*. June 1999 (TXT, PS, PDF, HTML, XML)

RFC3629⁵⁾ Yergeau, F. *UTF-8, a transformation format of ISO 10646*. STD 63, November 2003

RFC3986⁶⁾ Berners-Lee, T., Fielding, R., Masinter, L. *Uniform Resource Identifier (URI): Generic Syntax*. STD 66, January 2005 (TXT, HTML, XML)

1) Available at: <https://www.ietf.org/rfc/rfc1035.txt>.

2) Available at: <https://www.ietf.org/rfc/rfc2045.txt>.

3) Available at: <https://tools.ietf.org/html/rfc2540>.

4) Available at: <https://www.ietf.org/rfc/rfc2616.txt>.

5) Available at: <https://tools.ietf.org/html/rfc3629>.

6) Available at: <https://www.ietf.org/rfc/rfc3986.txt>.

ISO 28500:2017(E)

RFC4027⁷⁾ Josefsson, S. *Domain Name System Media Types*, April 2005

RFC4291⁸⁾ Hinden, R. and Deering, S. *IP Version 6 Addressing Architecture*, February 2006

RFC5322⁹⁾ Resnick, P. (ed.) *Internet Message Format*, October 2008

3 Terms, definitions and abbreviated terms**3.1 Terms and definitions**

For the purposes of this document, the following terms and definitions apply.

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

— ISO Online browsing platform: available at <http://www.iso.org/obp>

— IEC Electropedia: available at <http://www.electropedia.org/>

3.1.1**WARC record**

basic constituent of a WARC file, consisting of a sequence of WARC records

3.1.2**WARC record content block**

part (zero or more octets) of a WARC record that follows the header and that forms the main body of a WARC record

iTeh STANDARD PREVIEW
(standards.iteh.ai)

3.1.3**WARC record payload**

data object referred to, or contained by a WARC record as a meaningful subset of the content block

[SIST ISO 28500:2018](https://standards.iteh.ai/catalog/standards/sist/c3ac43f1-47a5-40e9-815d-a6fed47e1f56/sist-iso-28500-2018)

3.1.4**WARC record header**

beginning of a WARC record, consisting of one first line declaring the record to be in the WARC format with a given version number, followed by lines of named fields up to a blank line

<https://standards.iteh.ai/catalog/standards/sist/c3ac43f1-47a5-40e9-815d-a6fed47e1f56/sist-iso-28500-2018>

3.1.5**WARC named fields**

set of elements consisting of a name, a colon, and a value, with long values continued on indented lines

3.1.6**WARC logical record**

<segmentation> record composed of multiple segments, each represented by a WARC record

3.2 Abbreviated terms

ABNF augmented Backus-Naur form

ARC archive

CRLF carriage return line feed

DNS domain name system

FTP file transfer protocol

7) Available at: <https://tools.ietf.org/html/rfc4027>.

8) Available at: <https://tools.ietf.org/html/rfc4291>.

9) Available at: <https://tools.ietf.org/html/rfc5322>.

HTTP	hypertext transport protocol
IANA	Internet Assigned Numbers Authority
IESG	Internet Engineering Steering Group
RFC	request for comments
UR (I/L/N)	uniform resource (identifier/locator/name)
WARC	web archive

4 File and record model

A WARC format file is the simple concatenation of one or more WARC records. The first record usually describes the records to follow. In general, record content is either the direct result of a retrieval attempt (web pages, inline images, URL redirection information, DNS hostname lookup results, stand-alone files, etc.) or is synthesized material (e.g. metadata, transformed content) that provides additional information about archived content.

A WARC record shall consist of a record header followed by a record content block and two new lines. The WARC record header shall consist of one first line declaring the record to be in the WARC format with a given version number, then a variable number of line-oriented named fields terminated by a blank line. The WARC record header format shall follow the general rules of HTTP/1.1 [RFC2616] and [RFC5322] headers with one major exception: it shall also allow UTF-8 characters, as specified in [RFC3629].

The top-level view of a WARC file can be expressed in an ABNF grammar, reusing the augmented constructs defined in section 2.1 of HTTP/1.1 [RFC2616]. (In particular, note that to avoid the risk of confusion, where any WARC rule has the same name as an [RFC2616] rule, the definition here has been made the same, except in the case of the CHAR rule, which in WARC includes multibyte UTF-8 characters.)

```
warc-file      = 1*warc-record
warc-record   = header CRLF
                block CRLF CRLF
header        = version warc-fields
version       = "WARC/1.1" CRLF
warc-fields   = *named-field CRLF
block        = *OCTET
```

The record version shall appear first in every record and hence shall also begin the WARC file itself.

The WARC record relies heavily on named fields. Each named field consists of a name followed by a colon (":") and the field value. Field names are not case-sensitive. The field value may be preceded by any amount of linear white space (LWS), though a single space is preferred. Header fields can be extended over multiple lines by preceding each extra line with at least one space or tab character.

Named fields may appear in any order and field values may contain any UTF-8 character. Both defined-fields and extension-fields follow the generic named-field format. Extension-fields may be used in extensions of the core format.

```
named-field   = field-name ":" [ field-value ]
field-name    = token
field-value   = *(field-content | LWS)          ; further qualified
                                                    ; by field
                                                    ; definitions
field-content = <the OCTETs making up the field-value
                and consisting of either *TEXT or combinations
                of token, separators, and quoted-string>
OCTET        = <any 8-bit sequence of data>
token        = 1*<any US-ASCII character>
                except CTLs or separators>
separators   = "(" | ")" | "<" | ">" | "@"
```

ISO 28500:2017(E)

```

        | ", " | ";" | ":" | "\" | "<">
        | "/" | "[" | "]" | "?" | "="
        | "{" | "}" | SP | HT
TEXT      = <any OCTET except CTLs,
           but including LWS>
CHAR      = <UTF-8 characters; RFC 3629> ; (0-191, 194-244)
DIGIT     = <any US-ASCII digit "0".."9">
CTL       = <any US-ASCII control character
           (octets 0 - 31) and DEL (127)>
CR        = <ASCII CR, carriage return> ; (13)
LF        = <ASCII LF, linefeed> ; (10)
SP        = <ASCII SP, space> ; (32)
HT        = <ASCII HT, horizontal-tab> ; (9)
CRLF     = CR LF
LWS       = [CRLF] 1*(SP | HT) ; semantics same as
           ; single SP
quoted-string = (<"> *(qdtex | quoted-pair) <"> )
qdtex      = <any TEXT except <">>
quoted-pair = "\" CHAR ; single-character quoting
uri        = <'URI' per RFC3986>

```

Although UTF-8 characters are allowed, the ‘encoded-word’ mechanism of [RFC2047] may also be used when writing WARC fields and shall also be understood by WARC reading software.

NOTE In WARC 1.0 standard (ISO 28500:2009), uri was defined as “<” <'URI' per RFC3986> “>”. This rule has been changed to meet requests from implementers.

The rest of the WARC record grammar concerns defined-field parameters such as record identifier, record type, creation time, content length, and content type.

```

defined-field = WARC-Type
               | WARC-Record-ID
               | WARC-Date
               | Content-Length
               | Content-Type
               | WARC-Concurrent-To
               | WARC-Block-Digest
               | WARC-Payload-Digest
               | WARC-IP-Address
               | WARC-Refers-To
               | WARC-Refers-To-Target-URI
               | WARC-Refers-To-Date
               | WARC-Target-URI
               | WARC-Truncated
               | WARC-Warcinfo-ID
               | WARC-Filename ; warcinfo only
               | WARC-Profile ; revisit only
               | WARC-Identified-Payload-Type
               | WARC-Segment-Origin-ID ; continuation only
               | WARC-Segment-Number
               | WARC-Segment-Total-Length ; continuation only

```

Every WARC record shall have a type, reported in the WARC-Type field. Eight WARC record types are defined in this document as follows:

- ‘warcinfo’;
- ‘response’;
- ‘resource’;
- ‘request’;
- ‘metadata’;
- ‘revisit’;
- ‘conversion’;
- ‘continuation’.