

INTERNATIONAL
STANDARD

ISO/IEC
27038

First edition
2014-03-15

**Information technology — Security
techniques — Specification for digital
redaction**

*Technologies de l'information — Techniques de sécurité —
Spécifications pour la rédaction numérique*

Reference number
ISO/IEC 27038:2014(E)



© ISO/IEC 2014



COPYRIGHT PROTECTED DOCUMENT

© ISO/IEC 2014

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
Case postale 56 • CH-1211 Geneva 20
Tel. + 41 22 749 01 11
Fax + 41 22 749 09 47
E-mail copyright@iso.org
Web www.iso.org

Published in Switzerland

Contents

Page

Foreword	iv
Introduction	v
1 Scope	1
2 Terms and definitions	1
3 Symbols and abbreviated terms	2
4 General principles of digital redaction	2
4.1 Introduction.....	2
4.2 Anonymization.....	2
5 Requirements	2
5.1 Overview.....	2
5.2 Redaction principles.....	3
6 Redaction processes	4
6.1 Introduction.....	4
6.2 Paper intermediaries.....	4
6.3 Digital image intermediaries.....	4
6.4 Simple digital redaction.....	4
6.5 Complex digital redaction.....	5
6.6 Contextual information.....	5
7 Keeping records of redaction work	6
8 Characteristics of software redaction tools	6
9 Requirements for redaction testing	7
Annex A (informative) Redacting of PDF documents	9

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. In the field of information technology, ISO and IEC have established a joint technical committee, ISO/IEC JTC 1.

International Standards are drafted in accordance with the rules given in the ISO/IEC Directives, Part 2.

The main task of the joint technical committee is to prepare International Standards. Draft International Standards adopted by the joint technical committee are circulated to national bodies for voting. Publication as an International Standard requires approval by at least 75 % of the national bodies casting a vote.

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.

ISO/IEC 27038 was prepared by Joint Technical Committee ISO/IEC JTC 1, *Information technology*, Subcommittee SC 27, *IT Security techniques*.

Introduction

Some documents can contain information that must not be disclosed to some communities. Modified documents can be released to these communities after an appropriate processing of the original document. This processing can include the removal of sections, paragraphs or sentences with, where appropriate, the mention that they have been removed. This process is called the “redaction” of the document.

The digital redaction of documents is a relatively new area of document management practice, raising unique issues and potential risks. Where digital documents are redacted, removed information must not be recoverable. Hence, care needs to be taken so that redacted information is permanently removed from the digital document (e.g. it must not be simply hidden within nondisplayable portions of the document).

This International Standard specifies methods for digital redaction of digital documents.

Redaction can also involve the removal of document metadata or the removal of some information (e.g. an image) imported into the document.

It can be possible to identify redacted information in a redacted digital document by context. For example, the length of the redaction replacement text can indicate the length of the redacted information, and thus the information itself. This International Standard introduces two levels of redaction:

- BASIC redaction where context is not taken into consideration;
- ENHANCED redaction where context is taken into consideration.

Redaction techniques can be used for the anonymization of the information in a document, for example by the removal of some names within sentences. It can also involve the removal of numbers within sentences and their replacement by “XXX”.

Information technology — Security techniques — Specification for digital redaction

1 Scope

This International Standard specifies characteristics of techniques for performing digital redaction on digital documents. This International Standard also specifies requirements for software redaction tools and methods of testing that digital redaction has been securely completed.

This International Standard does not include the redaction of information from databases.

2 Terms and definitions

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

2.1

anonymization

process by which personally identifiable information (PII) is irreversibly altered in such a way that a PII principal can no longer be identified directly or indirectly, either by the PII controller alone or in collaboration with any other party

[SOURCE: ISO/IEC 29100:2011, definition 2.2]

2.2

document

recorded information which can be treated as a unit

Note 1 to entry: Documents can contain text, pictures, video and audio content, metadata and other associated content.

2.3

personally identifiable information

PII

any information that (a) can be used to identify the PII principal to whom such information relates, or (b) is or might be directly or indirectly linked to a PII principal

Note 1 to entry: To determine whether a PII principal is identifiable, account should be taken of all the means which can reasonably be used by the privacy stakeholder holding the data, or by any other party, to identify that natural person.

[SOURCE: ISO/IEC 29100:2011, definition 2.9]

2.4

redaction

permanent removal of information within a document

2.5

reviewer

individual(s) who assesses a document for redaction requirements

Note 1 to entry: There could be a series of individuals who assess a particular document.

3 Symbols and abbreviated terms

For the purposes of this document, the following abbreviations apply.

PII Personally Identifiable Information

PDF Portable Document Format

OCR Optical Character Recognition

XML Extensible Markup Language

4 General principles of digital redaction

4.1 Introduction

Redaction is carried out in order to permanently remove particular information from a copy of a document. It should be used when, for example, one or two individual words, a sentence or paragraph, an image, a name, address and/or signature needs to be removed from a document prior to it being disclosed to individuals who are not authorized to view the removed information.

The process of digital redaction is not simply to remove information but also to indicate where necessary that some information has been removed, so that the reader knows that the document has been redacted. For example, there can be a need to know that some words or some paragraphs have been deleted in order to maintain the semantics of the non-redacted information.

4.2 Anonymization

As an example, one of the purposes of redaction is to remove personally identifiable information (PII) from a document (anonymization). Where such a purpose is applicable, then redaction processes shall be so designed such that the identity of the individual about whose information is being redacted is protected.

It can be, for example, that even though a name has been redacted from a document, the identity of the individual is evident from the remaining information. Where anonymization is required, all information that could be used to identify the individual shall be redacted. This shall include all information that could be used in conjunction with other information (which can be obtained from other sources) to identify the individual.

5 Requirements

5.1 Overview

Organizations should have the capability to identify documents that need to be redacted prior to their release to other parts of the organization or to others (such as the public).

Redaction should be performed or overseen by reviewers that are knowledgeable about the documents and can determine what information is to be redacted. If reviewers identifying such information do not carry out redaction themselves, their instructions shall be specific e.g. 'Memo dated ..., paragraph no..., line starting... and ending...' etc.

Redaction shall be carried out on copies of the digital document. The redaction process shall result in the creation of a new digital document where complete and irreversible removal of the redacted information is achieved. This new digital document shall be managed and disposed of in the same manner as the original document.

When identifying information that needs to be redacted prior to release, whole sentences or paragraphs should not be identified if only one or two words in that sentence or paragraph are to be redacted, unless the release would enable the identification of the redacted information by context.

Where necessary, information relating to the effect that a digital document has been redacted shall be linked with the digital document. To identify the fact that a redaction process has been undertaken, redacted information may be replaced by a sentence stating that some information has been redacted.

When redaction is performed on a digital document, any metadata included within the digital document shall be reviewed for redaction requirements and appropriate redactions undertaken.

Where redaction is performed on a digital document that contains images, video and/or voice information, redaction techniques that remove the necessary information shall be used.

5.2 Redaction principles

The redaction of digital documents shall be carried out in accordance with the following principles:

- Retention of digital original document

The original or master version of a digital document shall not be redacted – redaction shall be carried out on a copy of the digital document. Original digital documents (e.g. the un-redacted document) shall be retained and be accessible only to those authorized.

- Complete removal of redacted information

Redaction shall irreversibly remove the required information from the redacted version of the document. The information shall be completely removed from the digital document, not simply from the displayable content.

- Security evaluated redaction

Redaction shall always be carried out using methods approved by the organization.

- Controlled environment

Electronic redaction shall be carried out in an environment that provides access only to those trained and authorized to carry out redaction.

- Intermediary stages

All the redacted documents in intermediary stages of the redaction process should be deleted. Only the original digital document and the appropriately redacted version should be retained. Where a particular digital document is to be redacted in different ways (for example for different audiences), then it may be appropriate to temporarily retain intermediate stages until all redaction processes have been completed.

If there are many digital documents that can have many redactions for different reasons, and if over time some of the reasons for redaction expire, it may be necessary to create and retain over time an intermediary copy that indicates the text or objects that should be redacted and the reasons for the redaction. The information is actually redacted when the redacted copy is produced. This provides the reviewers with the capability to re-review the intermediary copy of redacted documents, and remove redactions whose reason for redaction has expired, without having to re-review the documents for unexpired reasons for redaction.

6 Redaction processes

6.1 Introduction

The redaction of digital documents is a relatively new area of information and records management practice, and raises unique issues and potential risks.

Redaction may be carried out using a number of approaches:

- use of paper intermediaries;
- use of digital image intermediaries;
- simple redaction using plain text format files;
- complex redaction using original “complex” format files.

The following specifies the technical aspects of redacting digital documents.

6.2 Paper intermediaries

For digital documents which can be printed as a hardcopy, redaction techniques utilizing paper intermediaries may be used. There are 2 methods that can be used within this technique:

- The digital document is printed onto paper and redaction carried out on the printed copy.

In this case, the equipment / process used for the redaction shall ensure that the redacted information cannot be retrieved. The use of black marker pens may not be sufficient for this requirement. To ensure that redacted information cannot be retrieved, a photocopy of the redacted paper document shall be taken and it shall be used as the final redacted document.

- The information is redacted from a digital copy of the original digital document, and the copy printed to paper.

Where the redacted copy is required in digital format, this can be created by scanning the redacted paper copy into an appropriate format, ensuring that the redacted information is not reproduced on the digital document.

6.3 Digital image intermediaries

Digital images can be produced from digital documents using printer drivers or other similar techniques. Information in digital images to be redacted (this can be text and/or parts of an image) shall be replaced with areas of the same density as the background, such that the redacted information can no longer be retrieved.

6.4 Simple digital redaction

6.4.1 General

The simplest type of digital document to redact is a plain text file, in which there is a direct correspondence between bytes and displayable characters. Because of this direct correspondence, redacting information stored in this format is simply a matter of redacting the displayed information - once the file is saved, the redacted information will not be able to be recovered.

6.4.2 Character encoding

Care shall be taken with regards to encoding of characters used in the plain text file. Extended encoding mechanisms such as Unicode require appropriate editor is used; otherwise the direct correspondence between bytes and displayable characters is lost.

6.5 Complex digital redaction

6.5.1 General

Documents may be redacted electronically using their original format. This may be carried out either using deletion tools within the creating software, or by using specialized redaction software. This approach shall be treated with extreme caution, due to the possibility that deleted information may still be recoverable, and the potential for information to remain hidden within non-displayable portions of the digital document.

When redacting electronically, great care shall be taken over the choice of format for the redacted copy. It is crucial that no evidence of redacted information is retained in a redacted copy. Some binary formats may allow changes to be rolled back (for example by the use of 'mark-up' processes).

The redacted document may be required to be made available in its original format, for example, to preserve complex formatting. In such cases, the conversion of the document to another format, followed by conversion back to the original format should be used, such that the overall process removes all evidence of the redacted information. Information redaction may be carried out either prior to conversion, or in the intermediary format. This approach requires a thorough understanding of the formats and conversion processes involved, and the mechanisms by which information is transferred during format conversion.

Where facilities for the removal of 'mark-up' type information within a redacted document are available, they should be used prior to the release of the redacted copy. Software tools which 'hide' redacted information rather than remove redacted information shall be avoided.

6.5.2 Complex format documents

The majority of digital documents created using modern office software are stored in proprietary, binary-encoded or rich text formats. Neither the rich text nor the binary formats have the simple and direct correlation of plain text, and can contain significant information which is not displayed to the user, the presence of which can therefore not be apparent.

Digital documents using complex formats may incorporate change histories, audit trails, or embedded metadata. Some of this additional information may provide a means by which deleted information can be recovered or simple redaction processes otherwise circumvented. In addition, cryptographic and semantic analysis techniques can potentially be used to identify redacted information.

While the rich text based formats are subject of international standardization, the binary formats are usually proprietary to the software vendor which developed them. The mechanisms by which information is stored within these formats are often poorly understood.

6.5.3 Non-text information

Where documents contain embedded non-textual information such as images, voice and/or video information, redaction software can be required that can irreversibly remove the embedded information. Where only part of an embedded object is to be redacted, then the object shall be extracted, edited with appropriate software and re-embedded into the redacted document.

Where documents are in an image, voice and/or video format, special redaction software that can access and update the file can be required to ensure that the redacted information is not recoverable.

6.6 Contextual information

6.6.1 Introduction

On receipt of a redacted document, it can be possible to identify (or guess) the original contents of an individual redacted item of information. For example, if the redacted information was the month of a

year, then if the redaction suggests that a three-character word had been removed then it is likely that the redacted information would have been 'May'.

In some instances, the contextual information is unlikely to indicate the content of the redacted information. In this case, BASIC redaction shall be used.

Where contextual information is likely to give concordant clues about the nature or meaning of the redacted information, then ENHANCED redaction shall be used.

6.6.2 Basic redaction

Basic redaction shall result in the permanent removal of the required information. Redactions may be marked in any way – for example by the replacement of the redacted information with 'blacked-out' symbols replacing each redacted character.

6.6.3 Enhanced redaction

Enhanced redaction shall result in the permanent removal of the required information, along with such contextual information that could be used to identify the redacted information. All redacted information within a particular document shall be marked in the same manner. The length and other appropriate attributes of the redaction markings shall be consistent within a particular document.

7 Keeping records of redaction work

Organizations performing redaction should keep records of all redactions performed, especially where the reasons behind such redactions can be challenged. Such records should consist of either a copy of the redacted document or a description of the redactions carried out.

A capability may be required to indicate the reason(s) for redaction either in the redacted area or by a digital stamp near the area redacted.

8 Characteristics of software redaction tools

Where electronic documents are to be redacted, software tools that operate in compliance with this International Standard shall be used. Such tools may be part of the document creation software or may be a separate software tool.

Redaction tools shall operate in such a manner that the user marks the appropriate area of the electronic document and then selects the redaction function. There should be a function that enables comments to be assigned to specific areas of redacted information.

Redaction tools shall permanently remove the selected text from the electronic document. Once the text has been selected and the redaction function selected, there shall be no 'undo' function for that operation. Information shall be removed from the electronic document in such a way that the redacted information cannot be recovered by any software tool, including diagnostic and other investigatory tools.

Where information to be redacted is held in multiple instances within a file (e.g. as text and in an image form), all instances of the information to be redacted shall be removed.

Where redaction tools include the facility to redact parts of or whole embedded images, they shall also enable the permanent removal of parts of or whole embedded images and/or other embedded information.

Redaction tools shall have no effect on text or embedded information not marked for redaction.

Redaction tools shall have the ability to redact selected or all document metadata, document property information (including information about who carried out the redaction and when it was carried out) and other 'secondary' information.

Where a separate redaction software tool is used, the electronic document shall be retained in its original software format once redaction is completed.

Where redaction tools allow bulk redactions, they should redact each document in compliance to this International Standard.

Where ENHANCED redaction is required, the text of a document shall be reflowed with the original locations of all redacted text marked in a consistent manner. Such marking can be used to remove 'contextual' details such as the length of the redacted text information.

To enable the redaction of multiple occurrences of a word or phrase, there should be a link between the search facility of the software and the redaction tool. Such a facility enables all occurrences of a particular word or phrase to be redacted in one operation.

There should be an option to apply different marks to the electronic document to indicate different types of redacted information. This could be in the form of different blocked colours or of characters (such as a row of X's) of various sizes and/or fonts.

For information on the redaction of PDF format files, see [Annex A](#).

9 Requirements for redaction testing

This section sets out a number of tests that can be performed on a redacted document to assess whether redaction has been successfully completed. Tests should be selected based on the availability of suitable software and overall security requirements.

These tests do not determine whether the appropriate information has been redacted. They will, however, confirm whether the redaction process has been completed and that the redaction is irreversible.

In each case, the test is carried out on the redacted document. If any information or metadata that should have been redacted is identified, then the redaction process should be repeated.

Test 1: Metadata redaction

Examine the metadata in the redacted document. Such information may be contained in 'File Properties' fields or in separate metadata fields, depending upon the file format and file creation software used.

Test 2: Print to another file

Where an incomplete redaction has been performed, it may be possible to display the redacted information using the print function. Select an area of text that contains some redacted information and some un-redacted information. Print this selected area to paper or to a file (such as PDF). View the results and determine whether redacted information can be read.

Test 3: Copy and paste into another document

Where an incomplete redaction has been performed, it may be possible to display the redacted information by copying information from the redacted document and then pasting it into a new document (which may be in a different format). Select an area of text that contains some redacted information and some un-redacted information. Copy this information to the computers 'clipboard'. Open a blank document and paste this information. View the results and determine whether redacted information can be read.

Test 4: Run OCR

Where an incomplete redaction has been performed, particularly where the original document is in an image format, it may be possible to display the redacted information using Optical Character Recognition (OCR) software. Select an area of text that contains some redacted information and some un-redacted information. Run the OCR software. View the results and determine whether redacted information can be read.

Test 5: Text-to-speech converter

Where an incomplete redaction has been performed, it may be possible to listen to the redacted information using text-to-speech conversion software. Select an area of text that contains some redacted information and some un-redacted information. Run the text-to-speech software. Listen or record the results and determine whether redacted information can be heard.

Annex A

(informative)

Redacting of PDF documents

In the majority of circumstances, users of PDF documents rely upon software that has limited or no facilities to redact digital documents in this format. Where this is the case, then specialist tools for the redaction of information within the PDF document should be used.

In general, a PDF page consists of four primary types of objects as follows:

1. Text Object - information which uses fonts to represent information on a page;
2. Image Object - PDF object typically used to represent pixel information on the page. The same image object may be reused within document;
3. Inline Image Object – Image data embedded within the page content of a single page. These are typically used by OCR application to embed an image of a single word with a low OCR confidence level;
4. Path Object – Vector drawing operators consisting of lines, curves and rectangles. Text on a page may be represented using Path Objects as opposed to Text Objects.

Within a PDF page these objects can be layered in any sequence and to any depth. Redactions of PDF pages need to meet the following criteria:

- a) Within an area of the page to be redacted; all objects containing information, at any layer, are removed.
- b) The appearance of the unredacted portions of the PDF page is not altered.

