
**Information technology — International
string ordering and comparison —
Method for comparing character strings
and description of the common template
tailorable ordering**

*Technologies de l'information — Classement international et
comparaison de chaînes de caractères — Méthode de comparaison de
chaînes de caractères et description du modèle commun et adaptable
d'ordre de classement*

IECNORM.COM : Click to view the full text of ISO/IEC 14651:2011

IECNORM.COM : Click to view the full PDF of ISO/IEC 14651:2011



COPYRIGHT PROTECTED DOCUMENT

© ISO/IEC 2011

All rights reserved. Unless otherwise specified, no part of this publication may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying and microfilm, without permission in writing 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	v
Introduction.....	vi
1 Scope	1
2 Conformance	2
3 Normative references	2
4 Terms and definitions	2
5 Symbols and abbreviations	4
6 String comparison	4
6.1 Preparation of character strings prior to comparison	4
6.2 Key building and comparison	4
6.2.1 Preliminary considerations	4
6.2.2 Reference ordering key formation	5
6.2.3 Reference comparison method for ordering character strings	7
6.3 Common Template Table: formation and interpretation	8
6.3.1 BNF syntax rules for the Common Template Table in Annex A	8
6.3.2 Well-formedness conditions	11
6.3.3 Interpretation of tailored tables	12
6.3.4 Evaluation of weight tables	13
6.3.5 Conditions for considering specific table equivalences	14
6.3.6 Conditions for results to be considered equivalent	14
6.4 Declaration of a delta	14
6.5 Name of the Common Template Table and name declaration	16
Annex A (normative) Common Template Table	17
Annex B (informative) Example tailoring deltas	18
B.1 Example 1 — Minimal tailoring	18
B.2 Example 2 — Reversing the order of lowercase and uppercase letters	18
B.3 Example 3 — Canadian delta and benchmark	18
B.4 Example 4 — Danish delta and benchmark	21
B.5 Example 5 — A tailoring for Khmer	24
Annex C (informative) Preparation	26
C.1 General considerations	26
C.2 Thai string ordering	26
C.2.1 Thai ordering principles	27
C.2.2 Vowel/consonant rearrangement	27
C.2.3 Example ordered strings	29
C.3 Handling of numeral substrings in collation	29
C.3.1 Handling of ‘ordinary’ numerals for natural numbers	30
C.3.2 Handling of positional numerals in other scripts	32
C.3.3 Handling of other non-pure positional system numerals or non-positional system numerals (e.g. Roman numerals)	32
C.3.4 Handling of numerals for whole numbers	33
C.3.5 Handling of positive positional numerals with fractional parts	34
C.3.6 Handling of positive positional numerals with fraction parts and exponent parts	35
C.3.7 Handling of date and time of day indications	36
C.3.8 Making numbers less significant than letters	37
C.3.9 Maintaining determinacy	37

C.4	Preprocessing Hangul.....	38
C.4.1	Step 1	38
C.4.2	Step 2	38
C.4.3	Step 3	39
C.4.4	Step 4	39
Annex D (informative)	Tutorial on solutions brought by this International Standard to problems of lexical ordering	40
D.1	Problems.....	40
D.2	Solution.....	41
D.3	Tailoring.....	42
Bibliography	44

IECNORM.COM : Click to view the full PDF of ISO/IEC 14651:2011

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 14651 was prepared by Joint Technical Committee ISO/IEC JTC 1, *Information technology*, Subcommittee SC 2, *Coded character sets*.

This third edition cancels and replaces the second edition (ISO/IEC 14651:2007), which has been technically revised. It also incorporates the Amendment ISO/IEC 14651:2007/Amd.1:2009.

Introduction

This International Standard provides a method, applicable around the world, for ordering text data, and provides a Common Template Table which, when tailored, can meet a given language's ordering requirements while retaining reasonable ordering for other scripts.

The Common Template Table requires some tailoring in different local environments. Conformance to this International Standard requires that all deviations from the Template, called "deltas", be declared to document resultant discrepancies.

This International Standard describes a method to order text data independently of context.

ISO/IEC TR 14652 has specifications for ordering that informatively complement the specifications in this International Standard, and indicates where additional information can be sought on ordering keywords defined in this International Standard.

IECNORM.COM : Click to view the full PDF of ISO/IEC 14651:2011

Information technology — International string ordering and comparison — Method for comparing character strings and description of the common template tailorable ordering

1 Scope

This International Standard defines the following.

- A reference comparison method. This method is applicable to two character strings to determine their collating order in a sorted list. The method can be applied to strings containing characters from the full repertoire of ISO/IEC 10646. This method is also applicable to subsets of that repertoire, such as those of the different ISO/IEC 8-bit standard character sets, or any other character set, standardized or not, to produce ordering results valid (after tailoring) for a given set of languages for each script. This method uses collation tables derived either from the Common Template Table defined in this International Standard or from one of its tailorings. This method provides a reference format. The format is described using the Backus-Naur Form (BNF). This format is used to describe the Common Template Table. The format is used normatively *within* this International Standard.
- A Common Template Table. A given tailoring of the Common Template Table is used by the reference comparison method. The Common Template Table describes an order for all characters encoded in ISO/IEC 10646:2011. It allows for a specification of a fully deterministic ordering. This table enables the specification of a string ordering adapted to local ordering rules, without requiring an implementer to have knowledge of all the different scripts already encoded in the Universal Coded Character Set (UCS).

NOTE 1 This Common Template Table is to be modified to suit the needs of a local environment. The main worldwide benefit is that, for other scripts, often no modification is required and the order will remain as consistent as possible and predictable from an international point of view.

NOTE 2 The character repertoire used in this International Standard is equivalent to that of the Unicode Standard version 6.0.

- A reference name. The reference name refers to this particular version of the Common Template Table, for use as a reference when tailoring. In particular, this name implies that the table is linked to a particular stage of development of the ISO/IEC 10646 Universal multiple-octet coded character set.
- Requirements for a declaration of the differences (delta) between the collation table and the Common Template Table.

This International Standard does *not* mandate the following.

- A specific comparison method; any equivalent method giving the same results is acceptable.
- A specific format for describing or tailoring tables in a given implementation.
- Specific symbols to be used by implementations, except for the name of the Common Template Table.
- Any specific user interface for choosing options.
- Any specific internal format for intermediate keys used when comparing, nor for the table used. The use of numeric keys is not mandated either.

- A context-dependent ordering.
- Any particular preparation of character strings prior to comparison.

NOTE 1 It is normally necessary to do preparation of character strings prior to comparison even if it is not prescribed by this International Standard (see Annex C).

NOTE 2 Although no user interface is required to choose options or to specify tailoring of the Common Template Table, conformance requires always declaring the applicable delta, a declaration of differences with this table. It is recommended that processes present available tailoring options to users.

2 Conformance

A process is conformant to this International Standard if it produces results identical to those that result from the application of the specifications given in 6.2 to 6.5.

A declaration of conformity to this International Standard shall be accompanied by a statement, either directly or by reference, of the following.

- The number of levels that the process supports; this number shall be at least three.
- Whether the process supports the `forward, position` processing parameter.
- Whether the process supports the `backward` processing parameter and at which level.
- The tailoring *delta* described in 6.4 and how many levels are defined in the delta.
- If a preparation process is used, the method used shall be declared.

It is the responsibility of implementers to show how their delta declaration is related to the table syntax described in 6.3, and how the comparison method they use, if different from the one mentioned in Clause 6, can be considered as giving the same results as those prescribed by the method specified in Clause 6. The use of a preparation process is optional and its details are not specified in this International Standard.

3 Normative references

The following referenced documents are indispensable for the application 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 10646:2011, *Information technology — Universal Coded Character Set (UCS)*

4 Terms and definitions

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

4.1

character string

sequence of characters considered as a single object

4.2

collating symbol

symbol used to specify weights assigned to a collating element

4.3**collation (weighting) table**

mapping from collating elements to weighting elements

4.4**collating element**

sequence of one or more characters that are considered a single entity for ordering

4.5**delta**

list of the differences between a given collation table and another one

NOTE 1 The given collation table, together with a given delta, forms a new collation table.

NOTE 2 Unless otherwise specified in this International Standard, the term “delta” always refers to differences from the Common Template Table as defined in this International Standard.

4.6**(collation) level**

sequence number for a subkey in the series of subkeys forming a key

4.7**ordering****collation**

process by which, given two strings, it is determined whether the first one is less than, equal to, or greater than the second one

4.8**ordering key**

sequence of subkeys used to determine an order

4.9**(collation) preparation**

process in which given character strings are mapped to (other) character strings before the calculation of the ordering key for each of the strings

4.10**reference comparison method**

method for establishing an order between two ordering keys

NOTE 1 See Clause 6.

4.11**subkey**

sequence of weights computed for a character string

4.12**symbol**

collating element

4.13**(collation) weight**

positive integer value, used in subkeys, reflecting the relative order of collating elements

4.14**weighting element**

list of a given number of weights sequentially ordered by level

5 Symbols and abbreviations

Following ISO/IEC 10646, characters are referenced as UX where X stands for a series of one to eight hexadecimal digits (where all the letters in the hexadecimal string are in upper case) and refers to the value of that character in ISO/IEC 10646. This convention is used throughout this International Standard.

In the Common Template Table arbitrary symbols representing weights are used according to the BNF notation description in 6.3.1.

6 String comparison

6.1 Preparation of character strings prior to comparison

It may be necessary to transform character strings before the reference comparison method is applied to them (see Annex C for an example of such preparation). Although not part of the scope of this International Standard, preparation may be an important part of the ordering process. See Annex C for some examples of preparation.

Where applicable, it can be an important part of the preparation phase to map characters from a non-UCS encoding scheme to the UCS for input to the comparison method. This task can amongst other things encompass the correct handling of escape sequences in the originating encoding scheme, the mapping of characters without an allocated UCS codepoint to an application-defined codepoint in the private zone area and change the sequence of characters in strings that are not stored in logical order. For example, for visual order Arabic code sets, input strings must be put into logical order, and for some bibliographic code sets, strings with combining accents stored before their respective base character require that the combining accents be put after their base character. The resulting string sequence may then have to be remapped into its original encoding scheme.

NOTE 1 The Common Template Table is designed so that combining sequences and corresponding single characters (precomposed) will have precisely the same ordering. To avoid inadvertently breaking this invariant (and in the process breaking Unicode conformance), tailoring should reorder combining sequences when corresponding precomposed characters are reordered. For example, if Å is reordered after Z, then the sequence <A>+<combining diaeresis> should also be reordered. To avoid exposing encoding differences that may be invisible to the end-user, it is recommended that strings be normalized according to Unicode normalization form NFD to achieve this equivalence – see Bibliography, Unicode Technical Report no. 15.

NOTE 2 Escape sequences and control characters constitute very sensitive data to interpret, and it is highly recommended that preparation should filter out or transform these sequences.

NOTE 3 Since the reference method is a logical statement for the mechanism for string comparison, it does not preclude an implementation from using a non-UCS character encoding only, as long as it produces results as if it were using the reference comparison method.

6.2 Key building and comparison

6.2.1 Preliminary considerations

6.2.1.1 Assumptions

The collation table is a mapping from collating elements to weighting elements. In each weighting element, four levels are described in the Common Template Table. This number of levels can be extended or reduced, but can not be less than 3, in tailoring.

NOTE In the Common Template Table, levels generally have the following characteristics, although this purpose is not absolute:

Level 1: This level generally corresponds to the set of common letters of the alphabets for that script, if the script is alphabetic, and to the set of common characters of the script if the script is ideographic or syllabic.

Level 2: This level generally corresponds to diacritical marks affecting each basic character of the script. For some languages, letters with diacritics are always considered an integral part of the basic letters of the alphabet, and are not considered at this second level, but rather at the first. For example, in Spanish, N TILDE is considered a basic letter of the Latin script. Therefore, tailoring for Spanish will change the definition of N TILDE from "the weight of an N in the first level and the weight of a TILDE in the second level" to "the weight of an N TILDE (placed after N and before O) in the first level, and indication of the absence of a diacritic in the second level". For some characters, variant letter shapes are also dealt with on level 2. An example of this is ß, the LATIN SMALL LETTER SHARP S, which is treated as equivalent to ss on level 1, but traditionally distinguished from it on level 2.

Level 3: This level generally corresponds to case distinctions (upper and lower case) or to distinctions based on variant letter shapes (like the distinction between Hiragana and Katakana).

Level 4: This level generally corresponds to weighting differences that are less significant than those at the other levels. Often the last level (level 4 in the Common Template Table) is intended to specify additional weighting for "special" characters, i.e., characters normally not part of the spelling of words of a language (such as dingbats, punctuation, etc.), sometimes called "ignorable" characters in the context of computerized ordering.

6.2.1.2 Processing properties

A given tailored table has specific scanning and ordering properties. These properties may have been changed by the tailoring.

A scanning direction (forward or backward) for each level is used to indicate how to process the string. The scanning direction is a global property of each level defined in the tailored table.

If the last level is greater than three, there is an optional property of this level of comparison, called "position" option: when active, a comparison on the numeric position of each "ignorable" character in the two strings is effected, before comparing their weights. In other words, for two strings equivalent at all levels except the last one, the string having an ignorable character in the lowest position comes before the other one. In case corresponding ignorable characters are at the same position, then their weights are considered, until a difference is found. *Support* for this kind of processing is *optional* and is not necessary to claim conformance to this International Standard.

NOTE The scanning direction (forward or backward) is not normally related to the natural writing direction of scripts. The scanning direction applies to the logical sequence of the coded character string.

According to ISO/IEC 10646, for scripts written right to left, such as Arabic, the first characters in the logical sequence correspond to the rightmost characters in their natural presentation sequence. Conversely, for the Latin script, written left to right, the first characters in the logical sequence correspond to the leftmost characters in their natural presentation sequence.

Scanning forward starts with the lowest position in the logical sequence, while scanning backward starts from the highest position, independently of the presentation sequence. The scanning direction for ordering purposes is a global property of each level described in the table.

In ISO/IEC 10646, the Arabic script is artificially separated into two pseudo-scripts: 1) the logical, intrinsic Arabic, coded independently of contextual shapes, and 2) the Arabic presentation forms. Both allow the complete coding of Arabic, but intrinsic Arabic is normally preferred for better processing, while presentation-form Arabic is preferred by some presentation-oriented applications. ISO/IEC 10646 does not prescribe that the presentation forms be stored in any specific order, and in some implementations, the storage order for the latter is the reverse of the storage order used for intrinsic Arabic. It is therefore advisable that the preparation phase be used to make sure that Arabic presentation forms and other Arabic characters be fed to the comparison method in logical order.

A tailored sort table may be separated into sections for ease of tailoring. Each section is then assigned a name consistent with the specification in 6.3.1. One of the tailoring possibilities is to assign a given order to each section and to change the relative order of an entire section relative to other sections.

6.2.2 Reference ordering key formation

When two strings are to be compared to determine their relative order, the two strings are first parsed into a sequence of collating elements taking into account the multi-character "collating-element" statements declared and used in a tailored table (*if* the syntax of 6.3.1 is used). For the syntax used for expressing the Common Template Table, the name of a collating element consisting of a single character, is formed by the

UCS value of the character, expressed as a hexadecimal string, prefixed with “U”. For multi-character collating elements, the name and association to characters can be found via the collating elements declarations.

NOTE Collating elements with more characters have preference over shorter ones. As an example, if a multicharacter collating element is defined for “abc” and another one is defined for “ab” or for “bc”, then if “abc” is encountered, the collating element for “abc” will apply and not the one for “ab” or “bc”.

Then a sequence of *m* intermediary subkeys is formed out of a character string, where *m* is the number of levels described in a tailored collation weighting table.

Each ordering key is a sequence of subkeys. Each subkey is a list of numeric weights. A subkey is formed by successively appending the list of the weights assigned, at the level of the subkey, to each collating element of the string. The keyword “IGNORE” in the Common Template Table at the place of a sequence of collating symbols at a level, indicates that the sequence of weights at that level for that collating element is an empty sequence of weights.

There are three ways of forming subkeys: subkeys formed using the “forward” processing parameter; subkeys formed using the “backward” processing parameter; and subkeys formed using the “forward,position” processing parameter. Subkeys that use the “position” option can only occur at the last level, and only if that level is greater than three. Support of the “position” option is not required for conformance. If the processing parameter “forward,position” is not supported, “forward,position” shall be interpreted as if the processing parameter had been “forward”.

If there is no entry in the tailored table for a character of the input string, then the character’s weights are undefined. Characters with undefined weights should be ordered, with respect to characters that have defined weights, as if the undefined ones were given the weight named “UNDEFINED” at the first level. If there is no weight assignment to the symbol “UNDEFINED” before the symbol <SFFFF>’s weight assignment in a given tailored table, then the table shall be interpreted as if “UNDEFINED” was weighted just before <SFFFF>. The ordering of characters with undefined weights with respect to other characters with undefined weights is not specified in this International Standard.

NOTE 1 A possible way to order characters with undefined weights is *as if* there were tailoring lines like this one added to the table, in UCS code point order (call the maximal level 4 weight <PLAIN> here):

<UXXXX> “<UNDEFINED><UXXXX>”;<BASE>;<MIN>;<PLAIN>

NOTE 2 <SFFFF> is the maximal level 1 weight in the Common Template Table.

6.2.2.1 Formation of a subkey with the “forward” level processing parameter

Subkeys, at a particular level, formed with the “forward” level processing parameter, are built in the following way:

During forward scanning of each collating element of the input character string, one or more weights are obtained. These weights are obtained by matching the collating element in the given tailored collation weighting table, obtaining the list of weights assigned to the collating element at the particular level. The obtained weight list is appended to the end of the subkey.

6.2.2.2 Formation of a subkey with the “backward” level processing parameter

Subkeys, at a particular level, formed with the “backward” level processing parameter are built by forming a subkey as with the “forward” parameter, then reversing that subkey weight by weight.

6.2.2.3 Formation of a subkey with the “forward,position” level processing parameter

Subkeys, at the last level, formed with the “forward,position” level processing parameter are formed by forming a subkey as with the “forward” parameter, but for collating elements that are not ignored at all levels but the last one, their last level weighting (list of weights) is replaced by a single weight (call it <PLAIN> here) that is larger than all other weights at the last level in the given tailored table. Collating elements that are ignored at all levels but the last one, retain their weighting according to the given tailored table. Finally, any

trailing sequence of the maximal weight (<PLAIN>) is removed from the subkey, effectively replacing each trailing maximal weight with a zero weight.

NOTE For any level, implementations are allowed to apply an order-preserving reduction of all subkeys at that level. Such a reduction is useful for levels 2, 3, and 4. Level 2 often has long stretches of the weight named <BASE> in Annex A. Level 3 often has long stretches of the weight named <MIN> in Annex A. Level 4 often has long stretches of the weight named <PLAIN> here. One such ordering preserving subkey reduction technique effectively encodes, in the last level subkey, the (relative) position, as single number each, of each otherwise ignored character; hence the name of the “*position*” option.

6.2.3 Reference comparison method for ordering character strings

The reference comparison method for ordering two given character strings (*after* collation preparation, which is not part of the reference comparison method itself) is to compare ordering keys generated by the reference key formation method described in 6.2.2 of this International Standard:

- Begin by building an ordering key, using a given tailored collation weighting table, for each of the two given character strings being compared.
- Then compare the resulting keys according to the key ordering definition below in this subclause. Keys can be compared either up to a given level, or up to the last level of the given tailored collation weighting table.

NOTE 1 The comparison may be made *while* generating the ordering keys for two strings to be compared, stopping the key generation when the order of the strings can be determined. Such a technique is sometimes termed lazy evaluation, and some systems support it by default. This avoids generating the full ordering key, when an ordering difference may be found early in the keys. When a bigger set of strings are to be ordered, it may be advisable to generate the ordering keys, and store each key or an initial segment of each, before comparing the keys.

Weights for different levels should not be compared, which implies that subkeys at different levels should not be compared. Nor should keys generated from different tailored tables be compared.

NOTE 2 This allows implementations to assign weightings at each level independently of the other levels, and independently of other tailorings.

m is the maximal level of a given tailored table. Recall that a key is a list, of length m , of subkeys; a subkey is a list of weights; and a weight is a positive integer. Other notations used below are:

- L_z is the length of the subkey z , i.e., the number of weights in the subkey.
- $z_{wt(a)}$, where $1 \leq a \leq L_z$, is the weight at index position a (an integer > 0) of the subkey z .
- $u_{sk(b)}$, where $1 \leq b \leq m$, is the subkey at level b (an integer > 0) of the key u .

The orderings of weights, subkeys, and ordering keys (up to a given level, or up to the last level) are total order relations, defined for a given tailored collation table as follows:

1. Weights are positive integers (in the reference method), and are compared as such for the purposes of collation.
2. A subkey v is *less than* a subkey w (written $v < w$) **if and only if** there exists an integer i , where $1 \leq i \leq L_v + 1$ and $i \leq L_w$, such that
 - $i = 1$ and $v_{wt(i)} < w_{wt(i)}$, or
 - for all integers j , where $1 \leq j < i$, the equality $v_{wt(j)} = w_{wt(j)}$ holds, and either
 - $i \leq L_v$ and $v_{wt(i)} < w_{wt(i)}$, or
 - $i = L_v + 1$ and $0 < w_{wt(i)}$

A subkey v is *greater than* a subkey w (written $v > w$) **if and only if** w is less than v . A subkey v is *equal to* a subkey w (written $v = w$) **if and only if** neither v is less than w , nor w is less than v .

3. An ordering key x is *less than* an ordering key y at level s (written $x <_s y$) **if and only if** there exists an integer i , where $1 \leq i \leq s$ and $i \leq m$, such that

- $i = 1$ and $x_{sk(i)} < y_{sk(i)}$, or
- for all integers j , where $1 \leq j < i$, the equality $x_{sk(j)} = y_{sk(j)}$ holds, and $x_{sk(i)} < y_{sk(i)}$.

An ordering key x is *greater than* an ordering key y at level s (written $x >_s y$) **if and only if** y is less than x at level s . An ordering key x is *equal to* an ordering key y at level s (written $x =_s y$) **if and only if** neither x is less than y at level s , nor y is less than x at level s .

4. For ordering keys, $<$, $>$, and $=$ are defined as $<_m$, $>_m$, and $=_m$ respectively.

NOTE 3 For ordering keys, $x <_t y$ implies $x <_{t+1} y$, $x >_t y$ implies $x >_{t+1} y$, $x =_t y$ implies $x =_{t+1} y$, $x <_0 y$ is false, $x >_0 y$ is false, and $x =_0 y$ is true. Above level m , for a given tailored table, there are no further ordering distinctions. Note that this definition implies that if two ordering keys are in the 'less than' relationship at level 1, they will also be in the 'less than' relationship at levels 2, 3, 4, etc. In general, whenever two ordering keys are less than at a given level, they will also automatically be less than at all subsequent, higher levels. Conversely, if two ordering keys are equal at a given level, they will also automatically be equal at all preceding, lower levels.

6.3 Common Template Table: formation and interpretation

This subclause specifies:

- The syntax used to form the Common Template Table in Annex A of this International Standard or a tailored table based upon the Common Template Table as expressed in Annex A.
- Conditions of well-formedness of a table using this syntax.
- Interpretation of tailoring statements in deltas for tables formed using this syntax.
- Evaluation from symbols to weights of tailored tables formed using this syntax.
- Conditions for considering two tables as equivalent.
- Conditions for considering comparison results as equivalent.

6.3.1 BNF syntax rules for the Common Template Table in Annex A

Definitions between \langle angle brackets \rangle make use of terms not defined in this BNF syntax, and assume general English usage.

Other conventions:

- * indicates 0 or more repetitions of a token or a group of tokens;
- + indicates 1 or more repetitions of a token or a group of tokens;
- ? indicates optional occurrence of a token or a group of tokens (0 or 1 occurrences);
- parentheses are used to group tokens;
- production rules are terminated by a semicolon;

Define collation tables as sequences of lines:

```
weight_table = common_template_table | tailored_table ;
common_template_table =
    simple_line+ ;
tailored_table = table_line+ ;
```

Define the line types:

```

simple_line = (symbol_definition | collating_element |
              weight_assignment | order_end)? line_completion ;

table_line = simple_line | tailoring_line ;

tailoring_line = (reorder_after | order_start | reorder_end |
                 section_definition | reorder_section_after)
                 line_completion ;

```

Define the basic syntax for collation weighting:

```

symbol_definition =
    'collating-symbol' space+ symbol_element ;

symbol_element = symbol | symbol_range ;

symbol_range = symbol '..' symbol ;

symbol = simple_symbol | ucs_symbol ;

ucs_symbol = ('<U' one_to_eight_digit_hex_string '>') |
             ('<U-' one_to_eight_digit_hex_string '>') ;

simple_symbol = '<' identifier '>' ;

collating_element =
    'collating-element' space+ symbol space+
    'from' space+ quoted_symbol_sequence ;

quoted_symbol_sequence =
    '"' simple_weight+ '"' ;

weight_assignment =
    simple_weight | symbol_weight ;

simple_weight = symbol_element | 'UNDEFINED' ;

symbol_weight = symbol_element space+ weight_list ;

weight_list = level_token (semicolon level_token)* ;

level_token = symbol_group | 'IGNORE' ;

symbol_group = symbol_element | quoted_symbol_sequence ;

order_end = 'order_end' ;

```

Define the tailoring syntax:

```

reorder_after = 'reorder-after' space+ target_symbol ;

target_symbol = symbol ;

order_start = 'order_start' space+ multiple_level_direction ;

multiple_level_direction =
    (direction semicolon)* direction (' ,position')? ;

direction = 'forward' | 'backward' ;

reorder_end = 'reorder-end' ;

```

```

section_definition =
    section_definition_simple |
    section_definition_list ;

section_definition_simple =
    'section' space+ section_identifier ;

section_identifier =
    identifier ;

section_definition_list =
    'section' space+ section_identifier space+ symbol_list ;

symbol_list =
    symbol_element (semicolon symbol_element)* ;

reorder_section_after =
    'reorder-section-after' space+ section_identifier space+
    target_symbol ;

```

Define low-level tokens used by the rest of the syntax:

```

identifier =
    (letter | digit) id_part* ;

id_part =
    letter | digit | '-' | '_' ;

line_completion =
    space* comment? EOL ;

comment =
    comment_char character* ;

one_to_eight_digit_hex_string =
    hex_upper | hex_upper hex_upper |
    hex_upper hex_upper hex_upper |
    hex_upper hex_upper hex_upper hex_upper |
    hex_upper hex_upper hex_upper hex_upper hex_upper hex_upper |
    hex_upper hex_upper hex_upper hex_upper hex_upper hex_upper ;

hex_numeric_string =
    hex_upper+ ;

space =
    ' ' | <TAB> ;

semicolon =
    ';' ;

comment_char =
    '%' ;

digit =
    '0' | '1' | '2' | '3' | '4' |
    '5' | '6' | '7' | '8' | '9' ;

hex_upper =
    digit | 'A' | 'B' | 'C' | 'D' | 'E' | 'F' ;

letter =
    'a' | 'b' | 'c' | 'd' | 'e' | 'f' | 'g' |
    'h' | 'i' | 'j' | 'k' | 'l' | 'm' | 'n' |
    'o' | 'p' | 'q' | 'r' | 's' | 't' | 'u' |
    'v' | 'w' | 'x' | 'y' | 'z' |
    'A' | 'B' | 'C' | 'D' | 'E' | 'F' | 'G' |
    'H' | 'I' | 'J' | 'K' | 'L' | 'M' | 'N' |
    'O' | 'P' | 'Q' | 'R' | 'S' | 'T' | 'U' |
    'V' | 'W' | 'X' | 'Y' | 'Z' ;

```

EOL = <end-of-line in the text conventions in use> ;

character = <any member of the repertoire of the encoded character set in use, not including any characters used to delimit the end of lines> ;

6.3.1.1 Keyword usage

The usage of the following locale syntax keywords is as follows:

collating-symbol	Defines a collating symbol representing a weight.
collating-element	Defines a collating element symbol representing a multi-character collating element.
order_start	Defines collation rules. This statement is (after reordering is done) followed by one or more collation order statements, assigning multi-level collation weightings to collating elements. A collating element is either a character or a defined substring.
order_end	Specifies the end of the collating order statements.
reorder-after	Redefines collating rules. Specifies after which collating symbol weighting the lines between the "reorder-after" and the next (or "reorder-end") are to be moved. This statement is followed by one or more collation table lines. The result is to reassign collating symbol's values or collating element's weightings.
reorder-end	Specifies the end of the "reorder-after" collating order statements.
section	Defines a section of the table. A section can be moved as a whole by "reorder-section-after".
reorder-section-after	Redefines the order of sections. This statement is followed by a section symbol and a target symbol. The result is to reassign collating symbol's values or collating element's weightings.

6.3.2 Well-formedness conditions

WF 1. Any *simple_symbol* occurring in a *weight_list* shall also occur in the initial *symbol_element* of a *symbol_weight*, or *symbol_definition* that occurs no later in the sequence of lines that constitute a *weight_table*.

NOTE All *simple_symbols* must be "defined" before they are "used".

WF 2. No *symbol* that occurs in a *symbol_definition* in a *weight_table* that contains no *tailoring_lines* may occur in another *symbol_definition* in the same *weight_table*.

NOTE Duplication of collation weighting symbols is prohibited. This is true for the Common Template Table itself. It must remain true for a *tailored_table* after all reordering of lines has been applied.

WF 3. All *weight_lists* in a *tailored_table* shall contain the same number of *level_tokens*. An empty *level_token* shall be interpreted as the *collating_element* itself.

NOTE A tailorable table must be consistent in its use of levels throughout.

WF 4. A *tailored_table* shall contain exactly one *order_start* statement followed later by exactly one *order_end* statement. The *order_start* statement shall appear after the *symbol_definition* entries and before the *symbol_weight* entries after all reordering of lines has been applied.

WF 5. A *multiple_level_direction* in a *tailored_table* shall contain the same number of *directions* as the number of *level_tokens* of any *weight_list* in that *tailored_table*.

NOTE Any *order_start* must have the same number of levels as is generally used in the table.

WF 6. If a *level_token* in a *weight_list* consists of a *symbol_group*, all successive *level_tokens* in that *weight_list* shall also consist of a *symbol_group*.

NOTE *IGNORE* must not be used at a level after an explicit symbol for a weighting.

WF 7. Any *section_identifier* occurring in a *reorder_section_after* shall occur in a *section_definition* that occurs earlier in the sequence of *table_lines* that constitutes a *tailored_table*.

NOTE All *section_identifiers* must be “defined” before they are “used”.

WF 8. No two *section_definitions* in a *tailored_table* shall contain the same values in their *section_identifiers*.

NOTE Multiple definitions of *sections* are prohibited; *section_identifiers* must be unique.

WF 9. Each *reorder_after* in a *tailored_table* shall be followed at a later point in that tailored table by a *reorder_end* or another *reorder_after*.

WF 10. If the *section_identifier* of a statement *reorder_section_after* corresponds to the *section_identifier* of a *section_definition_list*, then the *target_symbol* of that statement *reorder_section_after* shall not have the same value as any of the *symbol*'s of the *section_definition_list*.

NOTE A *section* must not be reordered after a line which the *section* itself contains; recursive relocation of lines is prohibited.

WF 11. Any *symbol_range* shall contain two *symbols* which meet the following conditions: 1) Each of the two *symbols* shall contain a common prefix of length one. The prefix can be any letter except 'U'. 2) The portions of *identifier* of each of the two *symbols* following the common prefix shall be a *hex_numeric_string* containing the same number of hexadecimal digits. 3) When interpreted as numeric values, the *hex_numeric_string* of the first *symbol* shall be less than the *hex_numeric_string* of the second *symbol*. One plus the positive integral difference between the *hex_numeric_string* of the second *symbol*, interpreted as a numeric value, and the *hex_numeric_string* of the first *symbol*, interpreted as a numeric value, constitutes the number of *values* of the *symbol_range*.

NOTE A well-formed *symbol_range* is of a form such as <S4E00>..<>S9FA5>, where the common prefix is “S”, and the rest of the *identifier* portion of each *symbol* is a *hex_numeric_string*.

WF 12. Any *symbol_weight* that contains more than one *symbol_range* shall contain only *symbol_ranges* that meet the following requirement: Each *symbol_range* following the first *symbol_range* shall have the same number of *values* in its *range* as that of the first *symbol_range*.

NOTE This condition guarantees that all expanded ranges will be well-formed, since for any one *symbol_weight*, all of the range expansions will have the same number of *values*.

6.3.3 Interpretation of tailored tables

I 1. A *section* consists either 1) of the list of *simple_lines* which contain a *symbol_definition* whose value is equal to any *symbol* contained in the *symbol_list* in a *section_definition_list*, or 2) of the list of *simple_lines* following a *section_definition_simple* in a *tailored_table*.

NOTE A *section* is defined 1) by a specific *symbol_list*, or 2) by taking all the lines following the *section_definition* until another tailoring line such as an *order_start*, a *reorder_section_after*, another *section_definition*, or the end of the table is encountered.

I 2. A *simple_line* consisting of a *symbol_definition* containing a *symbol_range* is equivalent to a sequence of *simple_lines*, where each of those lines contain a *symbol* in place of the *symbol_range*. The *symbol* for each successive *simple_line* is generated by concatenating a *hex_numeric_string* to the common prefix of the *symbol_range*, in numeric order, starting with the hex value associated with the *hex_numeric_string* of the first *symbol* the range, and ending with the hex value associated with the *hex_numeric_string* of the second *symbol*. The *hex_numeric_string* concatenated to the common prefixes must contain the same number of digits as the *hex_numeric_string* of the first *symbol*. The number of *simple_lines* thus generated is equal to the number of *symbols* in the *symbol_range*.

NOTE A *symbol_definition* of the form “collating-symbol <S0301>..*<S0303>*” is equivalent to the three lines:

collating-symbol <S0301>
collating-symbol <S0302>
collating-symbol <S0303>

I 3. A *simple_line* consisting of a *symbol_weight* containing one or more *symbol_ranges* is equivalent to a sequence of *simple_lines*, where each *symbol_range* has been expanded into a sequence of *symbols*, as described in I 2 for *symbol_definition*

NOTE A *symbol_weight* of the form “<U2000>..*<U2002>* <S0301>..*<S0303>*; <BASE>; <MIN>; <U2000>..*<U2002>*” is equivalent to the three lines

<U2000> <S0301>; <BASE>; <MIN>; <U2000>
<U2001> <S0302>; <BASE >; <MIN>; <U2001>
<U2002> <S0303>; <BASE >; <MIN>; <U2002>

I 4a. A *tailored_table* containing a *reorder_after* is equivalent to the *tailored_table* where:

1. all *table_lines* that were ahead of the *reorder_after* and that contained *weight_assignments* whose initial *symbol* matches the initial *symbol* of any *weight_assignment* in the *table_lines* between the *reorder_after* and *reorder_end* have been removed,
2. the *table_lines* between that *reorder_after* and the first subsequent *reorder_end* to immediately follow the first *table_line* in the *tailored_table* containing a *weight_assignment* whose initial *symbol* is the same as the *target_symbol* in the *reorder_after* have been reordered, and
3. that *reorder_after* and that *reorder_end* have been removed.

NOTE Move the block of lines between the *reorder_after* and the *reorder_end* to follow the *target_symbol*, delete any prior lines that duplicate the *symbol_definitions* of the reordered lines, and remove the *reorder_after* and *reorder_end* themselves.

I 4b. When a *tailored_table* contains multiple groups of lines to be reordered, the table is interpreted by processing each *reorder_after* sequentially, starting from the first line of the table.

NOTE Subsequent line reorderings may impact lines that themselves were reordered by prior reorderings.

I 5. A *tailored_table* containing a *reorder_section_after* is equivalent to the *tailored_table* with the *section* associated with that *reorder_section_after* reordered (in the same relative order as the *table_lines* have in that *section*) to immediately follow the last *table_line* in the *tailored_table* containing a *symbol_definition* whose *symbol* is the same as the *target_symbol* in the *reorder_section_after*, and with that *reorder_section_after* removed.

I 6. A *weight_table* is said to be in normal form when it contains no *reorder_afters* nor *reorder_section_afters*.

NOTE A *tailored_table* can be put into normal form by the operations implied by I 4 and I 5.

6.3.4 Evaluation of weight tables

E 1. A *weight_table* in normal form is said to be evaluated when each *weight_assignment* in the *weight_table* is mapped to a positive integer value (a weight) such that those values increase monotonically by the order in which the *weight_assignments* occur in the *weight_table*.

NOTE 1 The *table_lines* of the *weight_table* can first be mapped to the set of positive integers, by sequential order in the table. This mapping defines an ordered set of line numbers. The *weight_assignments* are then mapped to a set of positive integers (weights) that varies monotonically with the set of line numbers.

NOTE 2 This does not restrict the starting number for the weight of the first *weight_assignment* (other than it must be positive) nor does it require that the numbers for these weights be immediately consecutive.

E 2. An evaluated *weight_table* is said to be collating-element-weighted when each *simple_symbol* occurring in each *weight_list* in that evaluated *weight_table* has been mapped to the weight that corresponds to the *weight_assignment* that contains the same *simple_symbol*.

NOTE 3 Each *weight_list* can be interpreted as containing either *symbol*'s mapped to integral weight values or as instances of the string 'IGNORE', which denotes the empty sequence of weights. At this point, the mathematical injection of strings can be defined using the *weight_table*.

NOTE 4 In a tailored table the value of any *hex_numeric_string* associated with a *symbol* generally does not reflect the numeric weighting of the *symbol*.

6.3.5 Conditions for considering specific table equivalences

A *weight_table* TBL1 and a *weight_table* TBL2 are said to be equivalent at a particular level if any comparison of strings using those tables up to that level results in the same ordering.

NOTE If two tables were claimed to be equivalent, one should always get the same results in forming keys for two strings based on TBL1 and comparing them than by forming keys based on TBL2 and comparing them.

6.3.6 Conditions for results to be considered equivalent

An implementation of international string ordering is conformant with this International Standard if, for any set of strings *S* defined on a repertoire *R*, the implementation can duplicate the same comparisons as those resulting from comparison of the numbers from an injection constructed according to the rules of 6.2.3 of this International Standard.

6.4 Declaration of a delta

Tailoring shall be based upon the Common Template Table described in Annex A. Tailoring may be accomplished using any syntax that is equivalent to the one described in this International Standard.

NOTE 1 For example, ISO/IEC TR 14652, uses a compatible extension of the syntax used in this International Standard for tailoring. A tailoring delta can also be expressed using the syntax of the Unicode collation algorithm (see Bibliography - Unicode Technical Report no. 10). It has also been demonstrated that a tailoring delta can also be expressed using an XML-conformant mark-up scheme.

Any declaration of conformance to this International Standard shall be accompanied with a declaration of the differences between the collation weighting table and the Common Template Table. A delta shall contain the equivalent of:

1. At least one valid *order_start* entry described in 6.3.1; an unlimited number of sections containing an *order_start* entry and an *order_end* entry may be declared.
2. The number of levels used for comparison.
3. The list of *symbol_definition* weights (as defined in 6.3.1) added and the position of the *symbol_definition* entry after which each insertion is made.
4. The list of *simple_line* entries (as defined in 6.3.1) deleted or inserted, referencing the position of the *simple_line* entry in the Common Template Table after which the insertions are made

NOTE 2 It is recommended that a delta should not be bigger than necessary.

In cases where a process allows the end-user to tailor the table himself or herself, a statement of conformance shall indicate which of the 4 elements of the previous list are tailorable and which of those 4 elements are not tailorable. For those which are not tailorable, the delta between fixed elements and the Common Template Table shall be declared.

NOTE 3 The declaration may use a different syntax from the one specified in 6.3 provided that the relationship with this syntax can be reasonably established. For example, the following declarations are valid:

"Collate U00E5 after U00FE at the primary level.
Collate U00E4 after U00E5 at the primary level. "

or

"The primary alphabet order is modified so that in all cases $z < \text{p} < \text{å} < \text{ä}$ ".

Note that the letters á and ä are sorted after Icelandic letter thorn (þ), itself already coming after all the variants of the letter z, i.e. á and ä have a weight value higher at level 1 than the one for þ, which itself comes after the ones for all variants of z.

The above two informal expressions can reasonably be considered to be equivalent to the following more precise expression (which also gives weights at levels 2 and 3 and explicitly takes care of accented á's, accented ä's and the Ångström sign):

reorder-after <S00FE> % Weight for THORN (after z; and unlike z, þ has no variants).

% Declare new collation symbols (weight names):

collating-symbol <S00E5> % for á

collating-symbol <S00E4> % for ä

% Declare new collating elements for the decompositions (substring names):

collating-element <U0061_030A> **from** "<U0061><U030A>" % decomposition of á

collating-element <U0041_030A> **from** "<U0041><U030A>" % decomposition of Á

collating-element <U0061_0308> **from** "<U0061><U0308>" % decomposition of ä

collating-element <U0041_0308> **from** "<U0041><U0308>" % decomposition of Ä

% Assign weights to the new collation symbols (after THORN):

<S00E5> % **for á**

<S00E4> % **for ä**

reorder-end

reorder-after <SFFFF> % The only place where we can put the order_start line.

order_start forward;forward;forward;forward

% Use the new weighted collation symbols and collating elements to tailor the collation rules:

% The letter Å:

<U00E5> <S00E5>;<BASE>;<MIN>;<U00E5> % LATIN SMALL LETTER A WITH RING ABOVE

<U0061_030A> <S00E5>;<BASE>;<MIN>;"<U0061_030A>" % decomposition of á

<U00C5> <S00E5>;<BASE>;<CAP>;<U00C5> % LATIN CAPITAL LETTER A WITH RING ABOVE

<U0041_030A> <S00E5>;<BASE>;<CAP>;"<U0041_030A>" % decomposition of Á

<U212B> <S00E5>;<BASE>;<CAP>;<U212B> % ÅNGSTROM SIGN (**the letter Å really**)

<U01FB> <S00E5>;"<BASE><AIGUT>";"<MIN><MIN>";<U01FB> % LATIN SMALL LETTER A WITH RING ABOVE AND ACUTE

<U01FA> <S00E5>;"<BASE><AIGUT>";"<CAP><MIN>";<U01FA> % LATIN CAPITAL LETTER A WITH RING ABOVE AND ACUTE

% The letter Ä:

<U00E4> <S00E4>;<BASE>;<MIN>;<U00E4> % LATIN SMALL LETTER A WITH DIAERESIS

<U0061_0308> <S00E4>;<BASE>;<MIN>;"<U0061_0308>" % decomposition of ä

<U00C4> <S00E4>;<BASE>;<CAP>;<U00C4> % LATIN CAPITAL LETTER A WITH DIAERESIS

<U0041_0308> <S00E4>;<BASE>;<CAP>;"<U0041_0308>" % decomposition of Ä

<U01DF> <S00E4>;"<BASE><MACRO>";"<MIN><MIN>";<U01DF> % LATIN SMALL LETTER A WITH DIAERESIS AND MACRON

<U01DE> <S00E4>;"<BASE><MACRO>";"<CAP><MIN>";<U01DE> % LATIN CAPITAL LETTER A WITH DIAERESIS AND MACRON

reorder-end

6.5 Name of the Common Template Table and name declaration

The name ISO14651_2010_TABLE1 shall be used whenever the Common Template Table is referred to externally as a base point in a given context, whether in a process, contract, or procurement requirement. If another name is used due to practical constraints, a declaration of conformance shall indicate the correspondence between this other name and the name ISO14651_2010_TABLE1.

The use of a defined name is necessary to manage the different stages of development of this table. This follows from the nature of the reference character repertoire, for which development will be ongoing for a number of years or even decades.

IECNORM.COM : Click to view the full PDF of ISO/IEC 14651:2011

Annex A (normative)

Common Template Table

In order to minimize formatting problems and the risk of errors in reproduction, the common template table is provided separately in a machine-readable file as a normative component of this International Standard. The file name for this language version is different from the normative reference name specified in 6.5 of this International Standard due to the existence of file versions commented in other natural languages. The file for this language version can also be retrieved on the ITTF web site at the following URL:

http://www.iso.org/ittf/ISO14651_2010_TABLE1_en.txt

There is an official French version of the file which only differs in its comments (its technical content is identical), and its name is: ISO14651_2010_TABLE1_fr.txt

NOTE 1 This International standard deprecates, but does not preclude specific reference to, the previous tables, which contained and still contain ordering information applicable to the repertoires of previous versions of ISO/IEC 10646 and their amendments. The previous tables can be found at the following URLs:

[ordering information on the repertoire of characters as defined in ISO/IEC 10646-1:1993 including Amendments 1-9] http://www.iso.org/ittf/ISO14651_2000_TABLE1.htm

[ordering information on the combined repertoire of characters of ISO/IEC 10646-1:2000 and ISO/IEC 10646-2:2001] http://www.iso.org/ittf/ISO14651_2002_TABLE1_en.txt

[ordering information on the repertoire of characters as defined in ISO/IEC 10646:2003] http://www.iso.org/ittf/ISO14651_2003_TABLE1_en.txt

[ordering information on the repertoire of characters as defined in ISO/IEC 10646-1:2003 including Amendments 1-2] http://www.iso.org/ittf/ISO14651_2006_TABLE1_en.txt

[ordering information on the repertoire of characters as defined in ISO/IEC 10646-1:2003 including Amendments 3-4] http://www.iso.org/ittf/ISO14651_2008_TABLE1_en.txt

The current Common Template Table reflects the repertoire of characters as defined in ISO/IEC 10646:2011, as indicated in the scope.

NOTE 2 The repertoire targeted by this International standard is equivalent to the repertoire of *The Unicode Standard Version 6.0*, published by *The Unicode Consortium*.

When ordering data applicable to other amendments of ISO/IEC 10646 becomes available, this International Standard and specifically its Common Template Table will be amended accordingly to cover the ordering of the additional characters and scripts. To meet cultural requirements of specific communities, delta declarations will have to be applied to the amended table as defined in this International Standard.

ISO_14651_2010_TABLE1 is the name that is used for referring to this table in this version of this International Standard.

Annex B (informative)

Example tailoring deltas

B.1 Example 1 — Minimal tailoring

The following is a minimal tailoring of the Common Template Table:

```
reorder_after <SFFFF>
order_start forward;forward;forward;forward
reorder-end
```

B.2 Example 2 — Reversing the order of lowercase and uppercase letters

The following is a simple tailoring example to show how to reverse the order of uppercase versus lowercase from the order specified in the Common Template Table.

```
% Make uppercase letters sort before lowercase
% and scanning of accents done forward at level 2.

% The entire range of tertiary weight symbols
% <MIN>..

```

B.3 Example 3 — Canadian delta and benchmark

This annex describes benchmark 1, based on Canadian standard CAN/CSA Z243.4.1-1998 (and 1992). The delta that precedes the benchmark *has been simplified* for illustration here; a larger delta is required, mainly for special characters, for full conformance to this Canadian standard, and is given here as an example only, limited to what is required for the benchmark. For complete information, the Canadian standard CAN/CSA Z243.4.1 should be consulted. The example tailoring is to be applied to the Common Template Table of Annex A, with the following delta:

1. Level processing properties:

```
forward; backward; forward; forward,position
```

2. Number of levels: 4 (unchanged).

3. No symbol changes.

4. The following ordering changes are done:

- æ sorted as if it were separate letters "ae" at level 1. The letters "ae" are distinguished at level 2 from the character "æ" and is ordered before it.
- ð sorted as if it were the letter "d" at level 1. The letter "ð" is distinguished at level 2 from the letter "d" and is ordered after it.
- þ sorted as if it were separate letters "th" at level 1. The letters "th" is distinguished at level 2 from the letter "þ" and is ordered before it.

A Canadian tailoring expressed in the tailoring syntax for this International Standard (normative only for Annex A) can be:

```
% copy ISO14651_2002_TABLE1

reorder-after <SFFFF>
order_start forward;backward;forward;forward,position

<U00E6> "<S0061><S0065>";"<BASE><VRNT1><BASE>";"<MIN><COMPAT><MIN>";<U00E6> % æ
<U00C6> "<S0061><S0065>";"<BASE><VRNT1><BASE>";"<CAP><COMPAT><CAP>";<U00C6> % Æ

<U01E3> "<S0061><S0065>";"<BASE><VRNT1><BASE><MACRO>";"<MIN><COMPAT><MIN><MIN>";<U01E3>
% æ WITH MACRON
<U01E2> "<S0061><S0065>";"<BASE><VRNT1><BASE><MACRO>";"<CAP><COMPAT><CAP><MIN>";<U01E2>
% Æ WITH MACRON

<U01FD> "<S0061><S0065>";"<BASE><VRNT1><BASE><AIGUT>";"<MIN><COMPAT><MIN><MIN>";<U01FD> % æ̆
<U01FC> "<S0061><S0065>";"<BASE><VRNT1><BASE><AIGUT>";"<CAP><COMPAT><CAP><MIN>";<U01FC> % Æ̆

<U00F0> <S0064>;<VRNT1>;<MIN>;<U00F0> % ð
<U00D0> <S0064>;<VRNT1>;<CAP>;<U00D0> % Ð

<U00FE> "<S0074><S0068>";"<BASE><VRNT1><BASE>";"<MIN><COMPAT><MIN>";<U00FE> % þ̆
<U00DE> "<S0074><S0068>";"<BASE><VRNT1><BASE>";"<CAP><COMPAT><CAP>";<U00DE> % Þ̆

reorder-end
```

Unordered list (required test case as per Canadian standard CAN/CSA Z243.4.1-1998, plus additions)			
ou	pêcher	Grossist	août
lésé	les	vice-presidents' offices	NOËL
péché	CÔTÉ	Copenhagen	@@@@@
vice-président	résumé	côte	L'Hay-les-Roses
9999	Ålborg	McArthur	CÔTE
OÙ	cañon	Mc Mahon	COTE
haïe	du	Aalborg	côté
coop	haie	Größe	coté
caennais	pêcher	vice-president's offices	aide
lèse	Mc Arthur	cølibat	air
dù	cote	PÉCHÉ	vice-president
air@@@	colon	COOP	modélé
côlon	l'âme	@@@air	Thorvardur
bohème	resume	VICE-VERSA	MODÈLE
gêne	élève	gêne	maçon
meðal	Þorvarður	CO-OP	MÂCON
lamé	Canon	révélé	pêche
pêche	lame	révèle	péché
LÈS	Bohême	çà et là	medal
vice versa	0000	MacArthur	ovoïde
C.A.F.	relève	Noël	pechère
Þorsmörk	gène	île	ode
cæsiium	casanier	aïeul	péchère
resumé	élevé	Île d'Orléans	œil
Bohémien	COTÉ	nôtre	
co-op	relevé	notre	

List with required results as per Canadian standard CAN/CSA Z243.4.1-1998			
@@@@@	COOP	lamé	pêche
0000	CO-OP	les	péché
9999	Copenhagen	LÈS	PÉCHÉ
Aalborg	cote	lèse	pêché
aide	COTE	lésé	pêcher
aïeul	côte	L'Hay-les-Roses	pêcher
air	CÔTE	MacArthur	pechère
@@@air	coté	MÂCON	péchère
air@@@	COTÉ	maçon	relève
Ålborg	côté	McArthur	relevé
août	CÔTÉ	Mc Arthur	resume
bohème	du	Mc Mahon	resumé
Bohême	dù	medal	résumé
Bohémien	élève	meðal	révèle
caennais	élevé	MODÈLE	révélé
cæsiium	gène	Modelé	Þorsmörk
çà et là	gène	Noël	Thorvardur
C.A.F.	géné	NOËL	Þorvarður
Canon	Größe	notre	vice-president
cañon	Grossist	nôtre	vice-président
casanier	haie	ode	vice-president's offices
cølibat	haïe	œil	vice-presidents' offices
colon	île	ou	vice versa
côlon	Île d'Orléans	OÙ	VICE-VERSA
coop	lame	ovoïde	
co-op	l'âme	pêche	

B.4 Example 4 — Danish delta and benchmark

The following is a Danish example tailoring delta. This formal specification corresponds to Danish standard DS 377 and to "Retskrivningsordbogen", the Danish orthography specification.

```
% This tailoring is in accordance with Danish Standard DS 377 (1980)
% and the Danish Orthography Dictionary (Retskrivningsordbogen par 4, 1986).
% It is also in accordance with Greenlandic orthography.
```

```
collating-symbol <LIGHT> % Symbolic weight for lighter than <BASE>
```

```
% Define collating elements for <AA> - LETTER A WITH RING ABOVE
% and combinations with combining accents
```

```
collating-symbol <A-A> % symbolic weight for <AA>
collating-element <A-plus-combining-ring> from "<U0041><U030A>"
collating-element <a-plus-combining-ring> from "<U0061><U030A>"
```

```
% Define collating elements for sequences of a-plus-a
```

```
collating-element <A-plus-A> from "<U0041><U0041>"
collating-element <A-plus-a> from "<U0041><U0061>"
collating-element <a-plus-A> from "<U0061><U0041>"
collating-element <a-plus-a> from "<U0061><U0061>"
```

```
% Define collating elements for combinations with combining accents
```

```
collating-element <U-plus-combining-diaeresis> from "<U0055><U0308>"
collating-element <u-plus-combining-diaeresis> from "<U0075><U0308>"
collating-element <U-plus-combining-doubleacute> from "<U0055><U030B>"
collating-element <u-plus-combining-doubleacute> from "<U0075><U030B>"
collating-element <O-plus-combining-diaeresis> from "<U004F><U0308>"
collating-element <o-plus-combining-diaeresis> from "<U006F><U0308>"
collating-element <O-plus-combining-doubleacute> from "<U004F><U030B>"
collating-element <o-plus-combining-doubleacute> from "<U006F><U030B>"
```

```
% Add the obligatory order_start line.
```

```
reorder-after <SFFFF>
order_start forward;backward;forward;forward,position
```

```
% copy ISO14651_2002_TABLE1
```

```
% Make capital letters sort before lowercase.
```

```
% Cf. example 2 for more explanation.
```

```
reorder-after <CIRCLECAP>
```

```
<CAP>
<WIDECAP>
<COMPATCAP>
<FONTCAP>
<CIRCLECAP>
<MIN>
<WIDE>
<COMPAT>
<FONT>
<CIRCLE>
```

```
% Introduce a weight that is lighter than <BASE>
```

```
reorder-after <BASE>
```

```
<LIGHT>
<BASE>
```

```
% A list of reweighting statements to deal with specific collation
% behaviour for Danish. All of these define or redefine weight_list's,
% and so the entire block could simply be reordered after the
% order-start entry in the table. However, for clarity here and for
% stability, each separate set of weightings is reordered locally in
% the table around the first entry for that set of weightings.
```

```
% Actually a number of other reweighting statements should be specified
% with respect to the ISO/IEC 14651 table so that all accents be
% ignored on the first level, while the 14651 table distinguish
% for example between different accented versions of <l> and a number
% of other latin letters. This is considered less important
% and too elaborate for this example.
```

% Also this example delta does not include specifications for reordering
 % of special characters like the currency denominators DOLLAR SIGN,
 % CENT SIGN, and POUND SIGN, and unit denominators like the
 % ANGSTROM SIGN, which have first-level weights in the CTT, while
 % rules in DS 377 prescribe that special characters are ignored.

% Reorder Danish letters at the end of the alphabet, after z

```
reorder-after <S007A> % z
<S00E6> % <AE> - LETTER AE
<S00F8> % <O/> - LETTER O WITH STROKE
<A-A> % <AA> - LETTER A WITH RING ABOVE
```

reorder-after <U007A> % z - this *line* only given for stability

```
% The letter ae is a separate letter in Danish
<U00C6> <S00E6>;<BASE>;<CAP>;<U00C6> % AE
<U00E6> <S00E6>;<BASE>;<MIN>;<U00E6> % ae
<U01FC> <S00E6>;"<BASE><AIGUT>";"<CAP><MIN>";<U01FC> % AE WITH ACUTE
<U01FD> <S00E6>;"<BASE><AIGUT>";"<MIN><MIN>";<U01FD> % ae WITH ACUTE
```

% The letter <a:> is given the same primary

```
% weight as <ae>, with unique variant weights at the secondary level.
<U00D6> <S00E6>;"<BASE><VRNT1>";"<CAP><MIN>";<U00D6> % A WITH DIAERESIS
<U00F6> <S00E6>;"<BASE><VRNT1>";"<MIN><MIN>";<U00F6> % a WITH DIAERESIS
```

% And replicate the weighting for the collating-element's formed with combining accents

```
<A-plus-combining-diaeresis> <S00E6>;"<BASE><VRNT1>";"<CAP><MIN>";"<U0041><U0308>"
<a-plus-combining-diaeresis> <S00E6>;"<BASE><VRNT1>";"<MIN><MIN>";"<U0061><U0308>"
```

% The letter <o/> - O WITH STROKE - is a separate letter in Danish

```
<U00D8> <S00F8>;<BASE>;<CAP>;<U00D8> % <O/>
<U00F8> <S00F8>;<BASE>;<MIN>;<U00F8> % <o/>
<U01FE> <S00F8>;"<BASE><AIGUT>";"<CAP><MIN>";<U01FE> % <O/> WITH ACUTE
<U01FF> <S00F8>;"<BASE><AIGUT>";"<MIN><MIN>";<U01FF> % <o/> WITH ACUTE
```

% The letters <o:> and <o"&> are given the same primary

```
% weight as <o/>, with unique variant weights at the secondary level.
<U00D6> <S00F8>;"<BASE><VRNT1>";"<CAP><MIN>";<U00D6> % O WITH DIAERESIS
<U00F6> <S00F8>;"<BASE><VRNT1>";"<MIN><MIN>";<U00F6> % o WITH DIAERESIS
<U0150> <S00F8>;"<BASE><VRNT2>";"<CAP><MIN>";<U0150> % O WITH DOUBLE ACUTE
<U0151> <S00F8>;"<BASE><VRNT2>";"<MIN><MIN>";<U0151> % o WITH DOUBLE ACUTE
```

% Replicate the weighting for the collating-element's formed with combining accents

```
<O-plus-combining-diaeresis> <S00F8>;"<BASE><VRNT1>";"<CAP><MIN>";"<U004F><U0308>"
<o-plus-combining-diaeresis> <S00F8>;"<BASE><VRNT1>";"<MIN><MIN>";"<U006F><U0308>"
<O-plus-combining-doubleacute> <S00F8>;"<BASE><VRNT2>";"<CAP><MIN>";"<U004F><U030B>"
<o-plus-combining-doubleacute> <S00F8>;"<BASE><VRNT2>";"<MIN><MIN>";"<U006F><U030B>"
```

% The letter <aa> - A WITH RING ABOVE - is weighted following the letter <o/> (see above)

```
<U00C5> <A-A>;<BASE>;<CAP>;<U00C5> % <AA>
<U00E5> <A-A>;<BASE>;<MIN>;<U00E5> % <aa>
<U01FA> <A-A>;"<BASE><AIGUT>";"<CAP><MIN>";<U01FA> % <AA> WITH ACUTE
<U01FB> <A-A>;"<BASE><AIGUT>";"<MIN><MIN>";<U01FB> % <aa> WITH ACUTE
```

% And replicate the weighting for the collating-element's formed with combining accents

```
<A-plus-combining-ring> <A-A>;<BASE>;<CAP>;<U00C5>
<a-plus-combining-ring> <A-A>;<BASE>;<MIN>;<U00E5>
```

% The sequences of letters a-plus-a are weighted as secondary variants of <AA>

```
<A-plus-A> <A-A>;"<BASE><VRNT1>";"<CAP><CAP>";"<U0041><U0041>" % AA
<A-plus-a> <A-A>;"<BASE><VRNT1>";"<CAP><MIN>";"<U0041><U0061>" % Aa
<a-plus-A> <A-A>;"<BASE><VRNT1>";"<MIN><CAP>";"<U0061><U0041>" % aA
<a-plus-a> <A-A>;"<BASE><VRNT1>";"<MIN><MIN>";"<U0061><U0061>" % aa
```

% The letters u with diaeresis and u with double-acute are given the same primary

% weight as y, with unique variant weights at the secondary level.

```
reorder-after <U00DC> % this *line* only given for stability
<U00DC> <S0079>;"<BASE><VRNT1>";"<CAP><MIN>";<U00DC> % U WITH DIAERESIS
<U00FC> <S0079>;"<BASE><VRNT1>";"<MIN><MIN>";<U00FC> % u WITH DIAERESIS
<U0170> <S0079>;"<BASE><VRNT2>";"<CAP><MIN>";<U0170> % U WITH DOUBLE ACUTE
<U0171> <S0079>;"<BASE><VRNT2>";"<MIN><MIN>";<U0171> % u WITH DOUBLE ACUTE
```

% And replicate the weighting for the collating-element's formed with combining accents

```
<U-plus-combining-diaeresis> <S0079>;"<BASE><VRNT1>";"<CAP><MIN>";"<U0055><U0308>"
<u-plus-combining-diaeresis> <S0079>;"<BASE><VRNT1>";"<MIN><MIN>";"<U0075><U0308>"
<U-plus-combining-doubleacute> <S0079>;"<BASE><VRNT2>";"<CAP><MIN>";"<U0055><U030B>"
<u-plus-combining-doubleacute> <S0079>;"<BASE><VRNT2>";"<MIN><MIN>";"<U0075><U030B>"
```

```

% The letter eth is equated to d
% with a secondary difference to distinguish it from d
reorder-after <U0064> % this *line* only given for stability
<U00D0> <S0064>; "<BASE><VRNT1>"; "<CAP><MIN>"; <U00D0> % ETH
<U00F0> <S0064>; "<BASE><VRNT1>"; "<MIN><MIN>"; <U00F0> % eth

% The letter thorn is treated as a sequence of t + h, with a variant weight
% at the secondary level
reorder-after <U00DE> % this *line* only given for stability
<U00DE> "<S0074><S0068>"; "<BASE><VRNT1><BASE>"; "<CAP><MIN><MIN>"; <U00DE> % THORN
<U00FE> "<S0074><S0068>"; "<BASE><VRNT1><BASE>"; "<MIN><MIN><MIN>"; <U00FE> % thorn

% The letter oe is treated as a sequence of o + e, with a special weight
% at the secondary level
reorder-after <U006F> % this *line* only given for stability
<U01D2> "<S006F><S0065>"; "<BASE><LIGHT>"; "<CAP><MIN>"; <U01D2> % OE
<U01D3> "<S006F><S0065>"; "<BASE><LIGHT>"; "<MIN><MIN>"; <U01D3> % oe

% Space, hyphen-minus, hyphen, and solidus are given a primary weight
% before any letter or digit, with hyphen-minus and solidus
% given a secondary difference from the weight for space.

reorder-after <U0020> % this *line* only given for stability
<U0020> <S0020>; "<BASE>"; "<MIN>"; <U0020> % SPACE
<U002D> <S0020>; "<BASE><VRNT1>"; "<U002D><MIN>"; <U002D> % HYPHEN-MINUS
<U2010> <S0020>; "<BASE><VRNT1>"; "<U0010><MIN>"; <U2010> % HYPHEN
<U002F> <S0020>; "<BASE><VRNT2>"; "<U002F><MIN>"; <U002F> % SOLIDUS

% The letter kra (for Greenlandic) is equated to a lowercase q
% with a secondary difference to distinguish it from q itself.
reorder-after <U0071> % q - this *line* only given for stability
<U0138> <S0071>; "<BASE><VRNT1>"; "<MIN><MIN>"; <U0138> % kra

% The letter sharp s is treated as a sequence of s + s
% with a special weight at the secondary level to make it come
% before s-plus-s - shorter precedes longer.
reorder-after <U0073> % s - this *line* only given for stability
<U00DF> "<S0073><S0073>"; "<BASE><LIGHT>"; "<MIN><MIN>"; <U00DF> % SHARP S

% To facilitate deterministic ordering, all controls have a unique
% weight at the 4th level.
reorder-after <U0000>
<U0000>..<U001F> IGNORE; IGNORE; IGNORE; <U0000>..<U001F>
<U007F>..<U009F> IGNORE; IGNORE; IGNORE; <U007F>..<U009F>

reorder-end

% End of the example tailoring for Danish

```

Benchmark for Danish (sorted order)			
A/S	D.S.B.	RÉE, A	STORM PETERSEN
ANDRE	DSC	REE, B	STORMLY
ANDRÉ	EKSTRA-ARBEJDE	RÉE, L	THORVALD
ANDREAS	EKSTRABUD	REE, V	THORVARDUR
AS	EKSTRAARBEJDE	SCHYTT, B	ÞORVARÐUR
CA	HØST	SCHYTT, H	THYGESEN
ÇA	HAAG	SCHÜTT, H	VESTERGÅRD, A
CB	HÅNDBOG	SCHYTT, L	VESTERGAARD, A
ÇC	HAANDVÆRKS BANKEN	SCHÜTT, M	VESTERGÅRD, B
DA	Karl	ß	ÆBLE
ÐA	karl	SS	ÅBLE
DB	NIELS JØRGEN	SSA	ØBERG
ÐC	NIELS-JØRGEN	STORE VILDMOSE	ÖBERG
DSB	NIELSEN	STOREKÆR	Århus

B.5 Example 5 — A tailoring for Khmer

The Khmer script is mainly used in Cambodia. The tailoring given below is not included in the CTT (see Annex A) itself in order to keep the CTT simple, especially for rare letterforms. E.g. the Khmer ROBAT for which the tailoring below may not be desirable for efficiency reasons, since this letter occurs very rarely, but the tailoring for handling it correctly may affect the efficiency of collation also for texts that do not contain any ROBAT.

```
reorder-after <MAX>

% Khmer:
collating-symbol <S1794_S17C9> % KHMER LETTER BA, KHMER SIGN MUUSIKATOAN
collating-symbol <S1794_S17CA> % KHMER LETTER BA, KHMER SIGN TRIISAP
collating-symbol <S17BB_S17C6> % KHMER VOWEL SIGN U, KHMER SIGN NIKAHIT
collating-symbol <S17B6_S17C6> % KHMER VOWEL SIGN AA, KHMER SIGN NIKAHIT
collating-symbol <C1780>..<C179C>

    % Declaration of Khmer contractions
collating-element <U1794_17C9> from "<U1794><U17C9>" % KHMER LETTER BA, KHMER SIGN MUUSIKATOAN
collating-element <U1794_17CA> from "<U1794><U17CA>" % KHMER LETTER BA, KHMER SIGN TRIISAP
collating-element <SW_17CC_1780>..<SW_17CC_17A2> from "<U1780>..<U17A2><U17CC>"
% KHMER LETTER KA, KHMER SIGN ROBAT..KHMER LETTER QA, KHMER SIGN ROBAT
collating-element <SW_17CC_17A5>..<SW_17CC_17B3> from "<U17A5>..<U17B3><U17CC>"
% KHMER INDEPENDENT VOWEL QI, KHMER SIGN ROBAT..KHMER INDEPENDENT VOWEL QAU, KHMER SIGN ROBAT
collating-element <U17C6_17BB> from "<U17BB><U17C6>" % KHMER VOWEL SIGN U, KHMER SIGN NIKAHIT
(OM properly spelled)
collating-element <U17BB_17C6> from "<U17C6><U17BB>" % KHMER SIGN NIKAHIT, KHMER VOWEL SIGN U
(OM with the wrong sequence of the characters)
collating-element <U17C6_17B6> from "<U17B6><U17C6>" % KHMER VOWEL SIGN AA, KHMER SIGN NIKAHIT
(AM properly spelled)
collating-element <U17B6_17C6> from "<U17C6><U17B6>" % KHMER SIGN NIKAHIT, KHMER VOWEL SIGN AA
(AM with the wrong sequence of the characters)
collating-element <U17D2_1780>..<U17D2_179C> from "<U17D2><U1780>..<U179C>"
% COENG, KHMER LETTER KA..COENG, KHMER LETTER QA
collating-element <U17D2_17A5>..<U17D2_17B3> from "<U17D2><U17A5>..<U17B3>"
% COENG, KHMER INDEPENDENT VOWEL QI..COENG, KHMER INDEPENDENT VOWEL QAU

reorder-after <S1794> % KHMER LETTER BA
<S1794_17C9> % KHMER LETTER BA, KHMER SIGN MUUSIKATOAN
<S1794_17CA> % KHMER LETTER BA, KHMER SIGN TRIISAP

reorder-after <S17C5> KHMER VOWEL SIGN AU
<S17BB_17C6> % KHMER VOWEL SIGN U, KHMER SIGN NIKAHIT

reorder-after <S17C6> KHMER SIGN NIKAHIT
<S17B6_17C6> % KHMER VOWEL SIGN AA, KHMER SIGN NIKAHIT

reorder-after <S17D2>
<C1780>..<C1794> % COENG, KHMER LETTER KA..COENG, KHMER LETTER BA
<C1795>..<C179A> % COENG, KHMER LETTER PHA..COENG, KHMER LETTER RO
<C17AB> % COENG, KHMER INDEPENDENT VOWEL RY
<C17AC> % COENG, KHMER INDEPENDENT VOWEL RYY
<C179B> % COENG, KHMER LETTER LO
<C17AD> % COENG, KHMER INDEPENDENT VOWEL LY
<C17AE> % COENG, KHMER INDEPENDENT VOWEL LYY
<C179C>..<C17A2> % COENG, KHMER LETTER VO..COENG, KHMER LETTER QA

reorder-after <SFFFF>
order_start forward;forward;forward;forward

<U1794_17C9> <S1794_17C9>;<BASE>;<MIN>;<U1794_17C9> % KHMER LETTER BA, KHMER SIGN MUUSIKATOAN
<U1794_17CA> <S1794_17CA>;<BASE>;<MIN>;<U1794_17CA> % KHMER LETTER BA, KHMER SIGN TRIISAP

%% The ROBAT contractions should be used only in an "advanced" tailoring for
%% Khmer, since ROBAT is rather rarely used, and these contractions
%% may impact on the efficiency of the key computation even if ROBAT does not
%% occur, since these contractions begin with commonly used letters.

<SW_17CC_1780>..<SW_17CC_17A2> "<S179A><S17D2><S1780>..<S17A2>";
    "<BASE><VRNT1><BASE><BASE>";<MIN><MIN><MIN><MIN>"; <SW_17CC_1780>..<SW_17CC_17A2>
```

```

% KHMER LETTER KA, KHMER SIGN ROBAT..KHMER LETTER QA, KHMER SIGN ROBAT
<SW_17CC_17A5>..

```

Annex C (informative)

Preparation

C.1 General considerations

Preparation is necessary only for modification and/or duplication of original strings to render them context-independent prior to the comparison phase. A non-duplicating preparation maps a given string to one string. A non-duplicating preparation can be composed with the key generation and comparison, as e.g. is needed for Lao and Thai (see Annex C.2), or proper ordering of numerals (see Annex C.3). A duplicating preparation can map a given string to several strings (to be sorted).

Examples of non-duplicating preparations are:

- Vowel-consonant rearrangement, as is needed for Thai (see C.2) and Lao.
- Transformation of numbers so that the result will be ordered in numerical order, as opposed to positional order (see C.3). Numeric ordering is particularly delicate and requires special consideration in many cases.
- Removal or rotation of characters that are a nuisance for special requirements of ordering; for example, removing articles (language dependent) in sorting book names as in:

Tale of two cities, A

- Transformation of abbreviated data into a fuller form. For example: transformation of "McArthur" to give "MacArthur".

Some examples of a duplicating preparation are:

- Duplicating a string into several 'rotations', like when producing a keyword-in-context index:

International string	International string ordering
International	ordering
	string ordering

- Duplication of a string such as "41" for as it is spelled out in different languages (Irish Gaelic, German, English, and French):

daichead a haon
einundvierzig
forty-one
quarante-et-un

C.2 Thai string ordering

This annex explains some of the principles behind the tailoring of the CTT given in Annex B.5 above, as well as the CTT ordering for Thai (and to some extent Lao).

by means of collating-element formation – the approach taken both by the CTT of this International Standard and by the collation weighting table of the Unicode Collation Algorithm (UTS #10). Every possible pair of leading vowel and consonant is defined as a collating-element, whose weight is equal to that of the rearranged substring. In addition, since two ๒ in sequence look just like a ๓, two ๒ in sequence should be handled just like a ๓.

Note that the rearrangement of each leading vowel is simply performed with its immediate succeeding consonant. No consonant cluster analysis is needed. Indeed, doing so would result in ambiguities or yield a different order than that specified in the Royal Institute Dictionary. For example:

1. Ambiguities: The problem with ambiguity is illustrated by the word “เพลลา”. It has two potential pronunciations: either as a two-syllable word, “phe-la” (meaning “time”), or as a one-syllable word, “phlao” (meaning “axle” or “abate”). A rearrangement algorithm which follows the distinct pronunciation of the potential cluster ‘พล’ in this string would result in distinct keys, “เพลลา” and “พลลา”, and therefore different weights, which are equally legal. Both words need to have the same weight to be sortable, however.
2. Non-conforming ordering: To illustrate the difference in ordering caused by the treatment of consonant clusters, consider these words, shown in conforming order: “เพล, เพลง, เพศ”. The correct rearrangement ignores any clusters and results in the following: “เพล, เพลง, เพศ”, which sorts in the order shown. If, however, pairs of consonants that form legal clusters were grouped as single collation elements (regardless of actual pronunciation where the potential pronunciation is ambiguous), then the results of rearrangement would be “<พล>๒, <พล>๓, เพศ”, which would yield the (non-conforming) ordering “เพศ, เพล, เพลง”. Again, if actual clusters were grouped as single collation elements (with some disambiguation effort), then the results of rearrangement would be “เพล, <พล>๓, เพศ”, which would yield the (non-conforming) ordering “เพล, เพศ, เพลง”.

IECNORM.COM : Click to view the full PDF of ISO/IEC 14651:2011

C.2.3 Example ordered strings

Here is an ordering example: Example for Thai (sorted order)			
กก	โกน	แข่งขัน	คัด
กรรม	โกรน	แขน	ฯพณฯ
กรรม	ใกล้	ครรภ-	พณิชนัย
-กระแย่ง	ไก่อ	ครรภ	ย่อง
กราบ	ไกล	จุมพล	รอง
กะเกณฑ์	ขัน	จูปล	ฤทธิ
กัก	ขนาน	ชาย	ฤณี
ก้าว	ข้าง	เดมา	ฤณี
กำ	ข้างๆ	เณร	ลลิตา
กิน	ข้างกระดาน	ตลาด	ภาชา
กึ่ง	ข้างขึ้น	ทูลเกล้า	วก
กั้น	ข้างควาย	ทูลเกล้าฯ	ศาล
กุน	ข้างๆ คู่ๆ	ทูลเกล้าทูลกระหม่อม	หริภุชชัย
กุด	ข้างเงิน	น้ำ	หฤทัย
แก้ง	ข้างออก	น้ำ	หลง
เกล้า	เข็ด	นี้	แห่ง
เกลียว	เขน	บุญหลง	แห่ง
เก้า	เขิน	บุญ-หลง	แหนม
เกาะ	เข่น	ป่า	แหนหวง
เกี้ยว	แข็ง	ป่า	แหบ
เกี้ยว	แข็ง	ป่า	แหม
เกือก	แข็ง	ป่า	อาน
แกง	แข็งขวา	ป่า	ฮ่า
แกะ	แข็งขึ้น	ป่าน	

C.3 Handling of numeral substrings in collation

A numeral is a string representing a number. The examples here deal with numerals that represent values in R , the real numbers, or, really, subsets of R , as these have a predetermined order. Only decimal numerals are dealt with in the examples given here.

The presentation below will first give positional system decimal numerals for natural numbers using the digits 0-9. It will progress to numerals for whole numbers, numerals with a fraction part, a fraction part and an exponent. There is also a brief discussion on numerals with digits from other scripts, scripts that sometimes uses another syntax with digits for numerals (such as Hân numerals), and Roman numerals. There are circumstances where digits do not represent numerical values, such as in part numbers and telephone numbers. The preparations described below have undesirable consequences in cases where *apparent* numerals do not represent numerical values, such as when the ordering telephone numbers or part 'numbers', and should be avoided in those cases.

C.3.1 Handling of 'ordinary' numerals for natural numbers

The Common Template Table has no means of ordering strings with numbers in such a way that the resulting order reflects the number values represented by the numerals. For example, given the following randomly-arranged strings:

Release 1
Release 20
Release 12
Release 2
Release 9

the method described in the this International Standard yields the following order for these strings:

Release 1
Release 12
Release 2
Release 20
Release 9

(It is sufficient simply to look positionally at just the first digit in each numeral to see why this ordering results.)
A more acceptable ordering is:

Release 1
Release 2
Release 9
Release 12
Release 20

The Common Template Table defined in this International Standard cannot be tailored to give this result. However, preparation can be done prior to the basic collation step to achieve the desired results when numeric value order is desired. The prepared strings are normally not presented to the user; only the original strings are. The prepared strings are normally only used for the collation key construction. A simple, but not very general, way of preparing numerals for natural ordering is to pad them with zeroes to a given number of digits. If one pads the numerals in our original example strings up to three digits, the following will result:

Release 001
Release 020
Release 012
Release 002
Release 009

Using the Common Template Table defined in this International Standard one then obtains the strings in a better order (here showing the strings as they are after preparation, which are normally not shown in the result):

Release 001
Release 002
Release 009
Release 012
Release 020

However, there are two problems with this approach:

- One must determine beforehand a (usually small) number of digits to pad up to. If the number of digits to pad up to is too large, the strings after preparation can become rather long, especially if there are several numerals in each string. If the number of digits to pad up to is too small, however, the risk is greater that there are actually occurring numerals with more digits than one has padded up to, which results in partially getting back to the original situation, where the numerals' values are not taken entirely into account.

- Determinacy is lost, if some of the original numerals were already partially zero-padded. For example, if the original strings were:

Release 01
Release 1

the strings after preparation are identical, and the end result (as the user would normally see it) could be either

Release 01
Release 1

or

Release 1
Release 01

and the relative order may come out differently for different occurrences of numerals, or different runs of the collation process applying the same rules. This kind of indeterminacy is undesirable.

There are many ways to deal with these problems. The following is one such way.

To each maximal digit subsequence prepend a fixed-number-of-digits numeral which represents the original number of digits in the numeral. For most cases a two-digit count would suffice (allowing up to 99 digits in the original integer numerals). For example, given the original strings:

Release 1
Release 01
Release 20
Release 12
Release 2
Release 09
Release 9

one obtains after this preparation the following strings:

Release 011
Release 0201
Release 0220
Release 0212
Release 012
Release 0209
Release 019

which would be ordered by the basic mechanism of this International Standard to:

Release 011
Release 012
Release 019
Release 0201
Release 0209
Release 0212
Release 0220

and are normally presented to the user as:

Release 1
Release 2
Release 9
Release 01
Release 09
Release 12
Release 20

This particular method puts numerals with a like original number of digits close to each other, even if the actual value represented is smaller due to the original zero-padding. If the represented *values* should be kept close together, one should instead duplicate the numeral: first a count of digits for the leading-zero-stripped numeral, the leading-zero-stripped numeral itself, followed by the original numeral. The duplication is needed to get determinacy relative to the original strings. For example, using the same original strings as above:

Release 011 1
Release 011 01
Release 0220 20
Release 0212 12
Release 012 2
Release 019 09
Release 019 9

which would be ordered by the basic mechanism of this International Standard to:

Release 011 01
Release 011 1
Release 012 2
Release 019 09
Release 019 9
Release 0212 12
Release 0220 20

and normally be presented to the user as:

Release 01
Release 1
Release 2
Release 09
Release 9
Release 12
Release 20

The originally zero-padded numerals consistently come before the numeral without (or with less) original zero-padding. The preparation processing could move the original numerals (in order of occurrence) to the very end of each string, if one wants to give the original zero-padding lesser significance than the text following the numerals.

The presence of several natural numerals in each string causes no additional problem.

Taking care of the natural number numerals is in most cases sufficient, and it is recommended that it be included as part of the usual preparation of strings to be collated. However, such preparation is not required by this International Standard.

C.3.2 Handling of positional numerals in other scripts

ISO/IEC 10646 encodes decimal digits for a number of scripts. In most cases these are used in a positional system, just like 0-9 usually are. However, one should not regard a sequence of numerals mixed from different scripts as a single numeral; rather, one should consider each maximal substring of digits of the same script to be a numeral.

C.3.3 Handling of other non-pure positional system numerals or non-positional system numerals (e.g. Roman numerals)

Chinese and some other languages can use decimal digits (in the Hànn script, for instance) interspersed with ideographs for “one thousand”, “ten”, etc. If such numerals are to be collated according to the value they represent, one can proceed as above, adding a step just after the initial duplication: convert the copy to the corresponding positional system numeral in the syntax used here for whole numerals.

Roman numerals, if handled, can be handled in a similar fashion to that described above. Duplicate, and replace the first copy with the same natural number expressed in the decimal positional system. E.g. "Louis V", where the V is determined to be a Roman numeral, can be modified to "Louis 5 V".

Caveat: In this case human interactive intervention or an expert system may be required, as in the following example involving the French language: CHAPITRE DIX might mean CHAPTER 10 or CHAPTER 509 ("dix" is the French word for 10, it is also the Roman numeral for 509).

C.3.4 Handling of numerals for whole numbers

If numerals for negative whole numbers are also to be ordered according to their value, there are a number of issues to be considered. Most frequently, negative whole values are given numerals that begin with a negation sign. The negation sign may be HYPHEN-MINUS U002D (but this character may often represent a true hyphen, rather than a negation), or MINUS SIGN U2212. However, there are other conventions also, like using a SOLIDUS U002F or a PERCENT SIGN U0025 to indicate negativeness; or the negation indicator can come *after* the digits rather than before; or negativeness can be indicated by putting the digits between parenthesis, and/or putting the digits in a contrasting colour (often red). In the examples here, only the case that negativeness is indicated by an immediately prepended MINUS SIGN is dealt with. Positiveness is indicated by either the absence of a MINUS SIGN, or the presence of a PLUS SIGN U002B.

Example strings:

Temperature: -9 °C
 Temperature: 0 °C
 Temperature: -14 °C
 Temperature: 05 °C
 Temperature: +5 °C
 Temperature: -0 °C
 Temperature: -09 °C
 Temperature: 105 °C
 Temperature: +05 °C
 Temperature: 5 °C

One preparation to get an acceptable and determinate order for numerals (in this syntax) for whole numbers is as follows (actual implementations should do something equivalent, but more efficient):

1. Duplicate the numerals in the string (including sign indications), putting the 'original' ones (not to be touched by the following steps) in order of original occurrence at the end of the string, leaving the copies to be modified at the original positions. This step ensures determinacy.
2. Ensure that all of the copies have an explicit initial sign indicator.
3. Remove leading zeroes in the copies of the numerals (systematically either leaving one zero digit for zero or representing 0 by the empty string of digits); alternatively, let all numeral copies have exactly one leading zero.
4. Between the sign indicator and the digits in the copies of the numerals, insert a (two-digit) count of how many digits there were (after removing the leading zeroes).
5. Do 9's complement on each digit in each copy of a negated numeral. 9's complement of a digit that individually represents the value x , is $9-x$. That is, 9's complement of 0 is 9, of 9 is 0, of 5 is 4, etc.

For the basic collation step, use a tailoring of the template given in this International Standard, namely, a tailoring where the PLUS SIGN and the MINUS SIGN are significant at the same level as the digits, and where the MINUS SIGN has less weight than the PLUS SIGN. (In the example below, it is assumed that the weight of PLUS SIGN is less than the weight of 0, but this is not a prerequisite for getting an acceptable ordering.)

Our example strings after this preparation:

Temperature: -980 °C -9
Temperature: +00 °C 0
Temperature: -9785 °C -14
Temperature: +015 °C 05
Temperature: +015 °C +5
Temperature: -99 °C -0
Temperature: -980 °C -09
Temperature: +03105 °C 105
Temperature: +015 °C +05
Temperature: +015 °C 5

The ordering for these, using the basic mechanism of this International Standard, is:

Temperature: -9785 °C -14
Temperature: -980 °C -09
Temperature: -980 °C -9
Temperature: -99 °C -0
Temperature: +00 °C 0
Temperature: +015 °C +05
Temperature: +015 °C +5
Temperature: +015 °C 05
Temperature: +015 °C 5
Temperature: +03105 °C 105

and normally presented to the user as:

Temperature: -14 °C
Temperature: -09 °C
Temperature: -9 °C
Temperature: -0 °C
Temperature: 0 °C
Temperature: +05 °C
Temperature: +5 °C
Temperature: 05 °C
Temperature: 5 °C
Temperature: 105 °C

This preparation results in a determinate ordering of strings that have zero or more numerals for whole numbers in them, such that the numerals are ordered according to the integer value they represent.

The process for other syntaxes for whole numbers can be similar. Just add a step to convert the copies to the syntax used here for whole numbers.

This technique for handling negative numerals can be used also for numerals with a fractional part, and so on (see below).

C.3.5 Handling of positive positional numerals with fractional parts

The method presented above can easily be adapted to the case where fraction parts may occur and are to be taken into account. A problem is, however, that the characters often used to delimit the integer part from the fraction part are also used for other purposes. The separator character is generally either FULL STOP U002E, or COMMA U002C. These characters also have other uses, also in conjunction with digits.

For the examples here, assume that FULL STOP is used (only) as a fraction part delimiter.

Do as above, but count only the digits in the integer part of the numeral for the count of digits to be prepended. The fraction part delimiter character (here: FULL STOP) can be removed.

For example:

–12.34
12.34
3.1415
3.14

After preparation:

–978765 –12.34
+021234 12.34
+013.1415 3.1415
+01314 3.14

Ordered:

–978765 –12.34
+01314 3.14
+0131415 3.1415
+021234 12.34

As presented to the user:

–12.34
3.14
3.1415
12.34

C.3.6 Handling of positive positional numerals with fraction parts and exponent parts

For very large, or very small, values, one often uses formats like 2.5×10^7 (to illustrate just one possible way of writing these for the purposes of the examples here). Here there is already an exponent, which must be combined with the “number of integer part digits” (here: digits before the decimal point), by adding those two numbers to get a resulting fixed-number-of-digits exponent to prepend just before the first digit. For this example, with a three-digit exponent: we get +00825. One problem here is that the resulting exponent may be negative. To handle this, use an exponent bias. For a three-digit exponent a bias of 500 may be suitable, which gives us for this example numeral: +50825, and for the numeral 2.5×10^{-7} we get +49425. Negative values are handled as before, with 9’s complement. -2.5×10^7 gives –49174, and -2.5×10^{-7} gives –50574. This method should be familiar to anyone with knowledge about (radix 10) floating point arithmetic.

Thus:

2.5×10^{-7}
 -2.5×10^7
 2.5×10^7
 -2.5×10^{-7}

After preparation (including a duplicate of the original, for determinacy):

+49425 2.5×10^{-7}
–49174 -2.5×10^7
+50825 2.5×10^7
–50574 -2.5×10^{-7}

Ordered:

–49174 -2.5×10^7
–50574 -2.5×10^{-7}
+49425 2.5×10^{-7}
+50825 2.5×10^7