

## IEEE Standards Interpretations for IEEE Std 1003.2™-1992 IEEE Standard for Information Technology--Portable Operating System Interfaces (POSIX®)-- Part 2: Shell and Utilities

Copyright © 1996 by the Institute of Electrical and Electronics Engineers, Inc. 3 Park Avenue New York, New York 10016-5997 USA All Rights Reserved.

Interpretations are issued to explain and clarify the intent of a standard and **do not** constitute an alteration to the original standard. In addition, interpretations are not intended to supply consulting information. Permission is hereby granted to download and print one copy of this document. Individuals seeking permission to reproduce and/or distribute this document in its entirety or portions of this document must contact the IEEE Standards Department for the appropriate license. Use of the information contained in this document is at your own risk.

IEEE Standards Department, Copyrights and Permissions, 445 Hoes Lane, Piscataway, New Jersey 08855-1331, USA

### Interpretation Request #40

**Topic:** I18N issues **Relevant Clauses:** 2.5.2.2

A POSIX defined a mechanism for talking about multi-character sequences as a single unit, namely as collating elements (CEs). Although CEs are motivated by sorting issues, they appear in REs. This obviously leads to the question of how to parse text into CEs?

There are many possible answers, and furthermore, the parsing might be affected by context. For example, given the usual alphabet augmented by the collating element `<ij>` defined as `<i><j>`, can the string `ij` ever be parsed as two collating elements? [1] 2.5.2.2 says in the context of sorting, "strings are first broken up into a series of collating elements" (line 1668). Does this apply to pattern matching? And if so, how exactly is this done (for sorting or pattern matching)?

Proposed suggestion: Add the following text somewhere; this text should be referred to by line 1668 and by the general RE introduction (2.8.2). "When a string is interpreted as a sequence of CEs, the sequence shall be as found by the following process: starting at the first character of the string, determine the longest prefix of the string that matches a CE, add that CE to the sequence and continue this process with the character after that prefix until the string is exhausted."

Note that this applies even if a sort key indicates that a piece of the text is processed in backwards (right- to-left) order; that is, the right-to-left processing applies to the CEs found by a left-to-right lexical scan. This is the greedy algorithm normally done in lexical analysis. Any other choice would require backtracking with potentially exponential runtime. It implies that, when `<i><j>` is a collating element, under no circumstances can a bracket expression match the `i` alone in the string `ij`. In particular, neither `[[.i.][.ij.]]j` nor

`[[.i.]]j` matches `ij`. By contrast, `i[[.j.]]` does not match `ij`, because in this regular expression `i` denotes a character and is unaffected by concerns about collating elements.

### **Interpretation Response**

The standard is unclear on this issue, and no conformance distinction can be made between alternative implementations based on this. This is being referred to the sponsor.

### **Rationale for Interpretation**

None.