

The `rtkinenc` package

Lars Hellström*

2000/01/24

Abstract

The `rtkinenc` package is functionally similar to the standard \LaTeX package `inputenc`—both set up active characters so that input character outside 7-bit ASCII are converted to the corresponding \LaTeX commands. Names of commands in `rtkinenc` have been selected so that it can read `inputenc` encoding definition files, and the aim is that `rtkinenc` should be backwards compatible with `inputenc`. `rtkinenc` is not a new version of `inputenc` though, nor is it part of standard \LaTeX .

The main difference between the two packages lies in the view on the input. With `inputenc`, the non-ASCII characters in the input are considered as shorthand representations of the “true” document contents (usually one or several commands), and the conversion is therefore irreversible. With `rtkinenc` the input file itself is considered the true document, and the conversion of non-ASCII characters to \LaTeX commands is merely done because it is the first step on the preferred route to typeset output. If the command is for an unavailable text symbol, then it is possible to return to the raw input and try some fallback method of typesetting the character.

The `inputenc` approach is natural for normal \LaTeX documents, but the `rtkinenc` approach is advantageous for program source code, where the *true* meaning of a file is not defined by \TeX , but by the compiler, interpreter, or whatever.

1 Implementation

```
1 <*pkg>
2 \NeedsTeXFormat{LaTeX2e}[1995/12/01]
3 \ProvidesPackage{rtkinenc}
4   [2000/01/24 v1.0 Rethinking input encoding package]
```

1.1 Basic mechanisms

`\RIE@last@char` `\RIE@last@char` is a `\count` register for storing the code of the raw character currently being typeset. It is minus one if no raw character is being typeset. This register should always be set globally.

*E-mail: `Lars.Hellstrom@math.umu.se`

```

5 \newcount\RIE@last@char
6 \global\RIE@last@char=\m@ne

```

```

\RIeC The syntax for the \RIeC command is
\RIE@char
\RIE@text@char \RIeC{<code>}{<definition>}
\RIE@code@char

```

Here *<code>* is assumed to be a sequence of digits giving some raw character code, and *<definition>* is assumed to be robust L^AT_EX code for typesetting (some representation of) that raw character. `\RIeC` is used as `inputenc`'s `\IeC` command—an `\RIeC` with arguments form the definition of an active character—but it carries extra information in the *<code>* argument.

Depending on what the control sequences `\protect` and `\RIE@char` are there are three different things this can do.

- It can execute the *<definition>* straight off. This happens when `\protect` is `\@typeset@protect` and `\RIE@char` is `\RIE@text@char`.
- It can set `\RIE@last@char` to *<code>*, then execute the *<definition>*, and finally set `\RIE@last@char` to minus one. This happens when `\protect` is `\@typeset@protect` and `\RIE@char` is `\RIE@code@char`.

Setting `\RIE@last@char` like that has the side-effect of prohibiting kerns and ligatures between the *<defintion>* and what surrounds it. Therefore it is inappropriate to have `\RIE@char` equal to `\RIE@code@char` in many types of text, and by default it will not be.

- It can expand to itself.¹ This happens when `\protect` is not `\@typeset@protect`, e.g. when writing to a file.

`\RIeC` cannot be defined using `\DeclareRobustCommand`, since that would insert `\protect`s that would prohibit normal kerning and ligaturing. Therefore the command robustness is maintained through an ad hoc definition (the `\def`). The reason for starting with an `\@ifundefined` is that the user shouldn't get less info about a redefinition than he/she would with a `\DeclareRobustCommand`.

```

7 \@ifundefined{RIeC}{-}{%
8   \PackageError{rtnkinenc}{Redefining \protect\RIeC}{%
9     I had expected \protect\RIeC\space to be undefined.\MessageBreak
10    Since it wasn't, there's a chance I have\MessageBreak
11    broken something.\MessageBreak\@ehc
12  }
13 }
14 \def\RIeC{%
15   \ifx \protect\@typeset@protect
16     \expandafter\RIE@char
17   \else
18     \noexpand\RIeC

```

¹Unless some command in the *<definition>* was defined using `\DeclareRobustCommand`, in which case it is the one level expansion of that command that will expand to itself. It all works out right in the end anyway.

```

19  \fi
20 }

21 \let\RIE@text@char=\@secondoftwo
22 \let\RIE@char=\RIE@text@char

23 \def\RIE@code@char#1#2{%
24  \global\RIE@last@char=#1
25  #2%
26  \global\RIE@last@char=\m@ne
27 }

```

RieBC The syntax for the `\RieBC` command is
RIE@both@char

```
\RieBC{<code>}{<text definition>}{<math definition>}
```

Here *<code>* is assumed to be a sequence of digits giving some raw character code, whereas *<text definition>* and *<math definition>* are assumed to be robust L^AT_EX code for typesetting (some representation of) that raw character in text and math mode, respectively. `\RieBC` is used like `\RieC`, but offers the possibility of alternative definitions for increased typesetting quality.

```

28 \@ifundefined{RieBC}{-}{%
29  \PackageError{rtnkinenc}{Redefining \protect\RieBC}{%
30   I had expected \protect\RieBC\space to be undefined.\MessageBreak
31   Since it wasn't, there's a chance I have\MessageBreak
32   broken something.\MessageBreak\@ehc
33  }
34 }

35 \def\RieBC{%
36  \ifx \protect\@typeset@protect
37   \expandafter\RIE@both@char
38  \else
39   \noexpand\RieBC
40  \fi
41 }

42 \def\RIE@both@char#1#2#3{%
43  \ifx \RIE@char\RIE@code@char
44   \global\RIE@last@char=#1
45  \fi
46  \ifmmode #3\else #2\fi
47  \ifx \RIE@char\RIE@code@char
48   \global\RIE@last@char=\m@ne
49  \fi
50 }

```

\TextSymbolUnavailable
@@TextSymbolUnavailable
\SetUnavailableAction
\RIE@symbol@unavailable

The only way to know whether a particular text command could be rendered as intended or not is to seize control of the standard L^AT_EX command `\TextSymbolUnavailable`. This command will then be given a new definition which selects whether some raw character fallback macro or the standard L^AT_EX error message should be given. The raw character fallback macro can be set using the `\SetUnavailableAction` command.

Before seizing control of `\TextSymbolUnavailable`, one must make sure that it does not have its autoload definition. Then the L^AT_EX definition is saved away in `\@@TextSymbolUnavailable`.

```
51 \def\@tempa{\@autoerr\TextSymbolUnavailable}
52 \ifx \@tempa\TextSymbolUnavailable
53   \@autoerr\relax
54 \fi
55 \let\@@TextSymbolUnavailable=\TextSymbolUnavailable
```

Then the new definition is given. It is pretty straightforward.

```
56 \def\TextSymbolUnavailable{%
57   \ifnum \m@ne<\RIE@last@char
58     \expandafter\RIE@symbol@unavailable \expandafter\RIE@last@char
59   \else
60     \expandafter\@@TextSymbolUnavailable
61   \fi
62 }
63 \PackageInfo{rtkinenc}{Redefining \protect\TextSymbolUnavailable}
```

The `\SetUnavailableAction` command locally defines the `\RIE@symbol@unavailable` macro, which is executed instead of the standard L^AT_EX `\TextSymbolUnavailable` when the text symbol in case was the representation of a raw input character. `\SetUnavailableAction` is used as

```
\SetUnavailableAction{<definition>}
```

where `<definition>` is like the last argument of `\newcommand`. The `<definition>` may contain `#1` and `#2`, where `#1` will contain the current raw character number and `#2` will contain the text command for which no definition could be found.

```
64 \newcommand\SetUnavailableAction{\def\RIE@symbol@unavailable##1##2}
```

The default fallback action is to call `\@@TextSymbolUnavailable`, but most of the definition deals with recognizing and handling the case that the input character hasn't been declared, rather than that its definition is not available. Normally, this shouldn't be used at all; instead the user should have installed another fallback.

```
65 \SetUnavailableAction{%
66   \ifx #2\relax
67     \begingroup
68       \let\RIE@char=\RIE@text@char
69       \RIE@undefined{#1}%
70     \endgroup
71   \else
72     \@@TextSymbolUnavailable{#2}%
73   \fi
74 }
```

`\RIE@undefined` The `\RIE@undefined` macro is used in the definition of input characters which are not defined in the current input encoding. It takes one argument, namely the code for the character in question. In text mode, this is an error. In code mode, it is passed on to the unavailable-action macro `\RIE@symbol@unavailable`.

The default value of `\RIE@symbol@unavailable` recognizes the `\relax` passed as second argument below as a flag that in reality it's the input character that hasn't been defined.

```

75 \def\RIE@undefined#1{%
76   \ifx \RIE@char\RIE@text@char
77     \PackageError{rtkinenc}{%
78       Input character #1 is undefined\MessageBreak
79       in inputencoding \RIE@encoding}\@eha
80   \else
81     \RIE@symbol@unavailable{#1}\relax
82   \fi
83 }
```

`\InputModeCode` The `\InputModeCode` and `\InputModeText` commands switch to the ‘code’ and ‘text’ respectively modes for the `rtkinenc` package. They both act locally, since it is often convenient to have the previous mode restored at the end of an environment.

The `\IfInputModeCode` command can be used to test which mode the `rtkinenc` package currently is in. Is is used as

```
\IfInputModeCode{<code>}{<text>}
```

This will expand to `<code>` or `<text>` when the current mode is code mode or text mode, respectively.

```

84 \newcommand\InputModeCode{\let\RIE@char=\RIE@code@char}
85 \newcommand\InputModeText{\let\RIE@char=\RIE@text@char}
86 \newcommand\IfInputModeCode{%
87   \ifx \RIE@char\RIE@code@char
88     \expandafter\@firstoftwo
89   \else
90     \expandafter\@secondoftwo
91   \fi
92 }
```

1.2 Setting the input encoding

The first two commands are identical; the duplication is only for being compabile with `inputenc`. The reason that there are two different commands in `inputenc` is that `\DeclareInputMath` saves a little memory by not taking special measures to see to that a space (if there is any) that follows the input character is not gobbled in case it is written to a temporary file and subsequently read back. Saving that small amount of memory is not the aim for `rtkinenc`, which is instead using up even more memory by including the character code in the definition.

`\DeclareInputText` These two commands are used as

```

\DeclareInputText{<slot>}{<definition>}
\DeclareInputMath{<slot>}{<definition>}
```

This makes the active character whose character code is `<slot>` a parameterless macro whose expansion is

`\RieC{⟨slot (sanitized)⟩}{⟨definition⟩}`

`⟨slot (sanitized)⟩` has the same numerical value as `⟨slot⟩`, but consists only of decimal digits. `⟨definition⟩` is the same in both cases.

```

93 \expandafter\ifx \csname DeclareInputText\endcsname\relax
94   \begingroup
95     \catcode\z@=13
96     \gdef\DeclareInputText#1#2{%
97       \@inpec@test
98       \begingroup
99         \uccode\z@=#1%
100        \uppercase{%
101          \endgroup
102          \expandafter\def \expandafter^^@%
103          }\expandafter{\expandafter\RieC \expandafter{\number#1}{#2}}%
104        }%
105      \endgroup
106 \else
107   \PackageError{rtkinenc}{\protect\DeclareInputText\space
108     already defined}{\@ehd\MessageBreak
109     Likely cause: you are already using the inputenc package.}
110 \fi

111 \@ifundefined{DeclareInputMath}{%
112   \let\DeclareInputMath=\DeclareInputText
113 }{%
114   \PackageError{rtkinenc}{\protect\DeclareInputMath\space
115     already defined}{\@ehd\MessageBreak
116     Likely cause: you are already using the inputenc package.}
117 }

```

`\DeclareInputBoth` The `\DeclareInputBoth` command is similar to `\DeclareInputText` and `\DeclareInputMath` commands, but it offers an extra feature—the text and math definitions of a character can be different. `\DeclareInputBoth` is used as

`\DeclareInputBoth{⟨slot⟩}{⟨text⟩}{⟨math⟩}`

where `⟨text⟩` and `⟨math⟩` are the text and math mode definitions respectively.

```

118 \expandafter\ifx \csname DeclareInputBoth\endcsname\relax
119   \begingroup
120     \catcode\z@=13
121     \gdef\DeclareInputBoth#1#2#3{%
122       \@inpec@test
123       \begingroup
124         \uccode\z@=#1%
125         \uppercase{%
126           \endgroup
127           \expandafter\def \expandafter^^@%
128           }\expandafter{\expandafter\RieBC \expandafter{\number#1}%
129             {#2}{#3}}%

```

```

130     }%
131   }%
132 \endgroup
133 \else
134   \PackageError{rtkinenc}{\protect\DeclareInputBoth\space
135     already defined}\@ehd
136 \fi

```

`\inputencoding` The `\inputencoding` command sets the current input encoding to be the one specified in its only argument. First all characters are set to be active and defined as `\RIE@undefined{<slot>}`, except for Null (`^^@`), tab (`^^I`), line feed (`^^J`), form feed (`^^L`), carriage return (`^^M`), and space— which are left as they were. Then `#1.def` is inputted; this file is expected to contain all the `\DeclareInput...` commands that are needed for the wanted input encoding.

Besides that, `\inputencoding` also does a check to see that the file `#1.def` actually did execute some `\DeclareInput...` command. Since it wouldn't at all surprise me if someone likes to tinker with this test, it is done exactly as in `inputenc`.

`\inputencoding` should not be used in horizontal mode since space tokens within the file inputted will produce unwanted spaces in the output.

`\RIE@encoding` stores the name of the current input encoding. It is used in an error message by `\RIE@undefined`.

```

137 \def\inputencoding#1{%
138   \gdef\@inpenc@test{\global\let\@inpenc@test\relax}%
139   \protected@edef\RIE@encoding{#1}%
140   \ifvmode
141     \RIE@loop^^A^^H%
142     \RIE@loop^^K^^K%
143     \RIE@loop^^N^^_%
144     \RIE@loop^^?\^^ff%
145     \input{#1.def}%
146   \fi
147   \ifx \@inpenc@test\relax \else
148     \PackageWarning{rtkinenc}%
149       {No characters defined\MessageBreak
150        by input encoding change to '#1'}%
151   \fi
152 }

```

`\RIE@loop` The `\RIE@loop` command makes characters `#1` to `#2` inclusive active and undefined.

```

153 \begingroup
154   \catcode\z@=\active
155   \gdef\RIE@loop#1#2{%
156     \@tempcnta='#1\relax
157     \count@=\uccode\z@
158     \loop
159       \catcode\@tempcnta\active
160       \uccode\z@=\@tempcnta

```

```

161     \uppercase{%
162     \expandafter\def \expandafter^^@\expandafter{%
163     \expandafter\RIE@undefined\expandafter{\the\@tempcnta}%
164     }%
165     }%
166     \ifnum '#2>\@tempcnta
167     \advance \@tempcnta \@ne
168     \repeat
169     \uccode\z@=\count@
170   }
171 \endgroup

```

1.3 Miscellanea

`\TypesetHexNumber`
`\TypesetOctalNumber`

In many computer languages, special character escape sequences based on character codes must be written with the character code in hexadecimal or octal notation. These commands take a T_EX number in the interval 0–255 as argument and typesets it using the hexadecimal (figures 0–9 and a–f) or octal (figures 0–7) notation. No font changes or T_EX mode changes are made. `\TypesetHexNumber` always typesets two characters, `\TypesetOctalNumber` always typesets three characters.

Care has been taken to see to that every count register can be used as the argument of these two macros, and they only make local assignments. They do not check that the argument is in range, though.

```

172 \newcommand\TypesetHexNumber[1]{%
173   \begingroup
174     \count@=#1\relax
175     \chardef\@tempa=\count@
176     \divide \count@ \sixt@@n
177     \ifcase\count@
178     0\or 1\or 2\or 3\or 4\or 5\or 6\or 7\or 8\or 9\or a\or b\or
179     c\or d\or e\else f%
180   \fi
181   \multiply \count@ -\sixt@@n
182   \advance \count@ \@tempa
183   \ifcase\count@
184   0\or 1\or 2\or 3\or 4\or 5\or 6\or 7\or 8\or 9\or a\or b\or
185   c\or d\or e\else f%
186   \fi
187 \endgroup
188 }

189 \newcommand\TypesetOctalNumber[1]{%
190   \begingroup
191     \count@=#1\relax
192     \chardef\@tempa=\count@
193     \divide \count@ 64\relax
194     \the\count@
195     \multiply \count@ -64%
196     \advance \count@ \@tempa

```



```

197     \chardef\@tempa=\count@
198     \divide \count@ 8\relax
199     \the\count@
200     \multiply \count@ -8%
201     \advance \count@ \@tempa
202     \the\count@
203 \endgroup
204 }

```

`\verifycharcode` The `\verifycharcode` command is used as

```
\verifycharcode{⟨character⟩}{⟨code⟩}
```

Here $\langle character \rangle$ can be any $\langle character token \rangle$ (as defined in *The T_EXbook*); e.g. a control sequence whose name consists of one character. $\langle code \rangle$ can be any valid $\langle number \rangle$. If $\langle code \rangle$ is not the character code of the $\langle character \rangle$, then `\verifycharcode` makes a warning about this.

The purpose of this command is to detect when the code of some character used in a document is changed. Today these things happen mainly when transferring a document between two systems which use different encodings, and it is usually the right thing to do. Some computer programs (and now I mean the source) do however rely on the exact character codes used in them, and documents containing such programs may use the `\verifycharcode` to test that none of the critical characters have been altered.

```

205 \newcommand\verifycharcode[2]{%
206   \ifnum '#1=#2 \else
207     \PackageWarning{rtninc}{%
208       Input character with code \number#2\MessageBreak
209       should be the character with code \number'#1}%
210   \fi
211 }

```

1.4 Option processing

Each option is interpreted as the name of an input encoding. That input encoding definition file is inputed.

```

212 \DeclareOption*{\inputencoding{\CurrentOption}}
213 \ProcessOptions
214 \</pkg>

```