


# The akshar package

Vu Van Dung

Version 0.1 — 2020/05/23

 <https://ctan.org/pkg/akshar>

 <https://github.com/jouleev/akshar>

## Abstract

This package provides tools to deal with special characters in a string of South Asian script. Currently supported scripts: Bengali, Gujarati, Gurmukhi, Kannada, Oriya, Malayalam, Sinhala, Tamil and Telugu.

## Contents

### 1 Introduction

When dealing with processing strings in the South Asian scripts, normal  $\LaTeX$  commands usually find some difficulties in distinguishing “normal” characters, like क, and “special” characters, for example ् or ी. Let’s consider this example code:

```
1 \ExplSyntaxOn
2 \tl_set:Nn \l_tmpa_tl { की}
3 \tl_count:N \l_tmpa_tl \c_space_token tokens.
4 \ExplSyntaxOff
```

2 tokens.

The output is 2, but the number of characters in it is only one! The reason is quite simple: the compiler treats ी as a normal character, and it shouldn’t do so.

To tackle that, this package provides `expl3` functions to “convert” a given string, written in South Asian scripts, to a sequence of token lists. each of these token lists is a “true” South Asian language character. You can now do anything you want with this sequence; and this package does provide some front-end macros for some simple actions on the input string.

### 2 User manual

#### 2.1 $\LaTeX 2_{\epsilon}$ macros

---

`\aksharStrLen` `\aksharStrLen {(token list)}`

---

Return the number of South Asian characters in the `(token list)`.

There are 4 characters in नमस्कार.  
`expl3` returns 7, which is wrong.

```
1 There are \aksharStrLen{ नमस्कार} characters in नमस्कार.\par
2 \ExplSyntaxOn
3 \pkg{expl3}~returns~\tl_count:n { नमस्कार},~which~is~wrong.
4 \ExplSyntaxOff
```

---

`\aksharStrHead` `\aksharStrHead {(token list)} {(n)}`

---

Return the first character of the token list.

मं 1 \aksharStrHead { मंलीममड }

---

\aksharStrTail {(token list)} {(n)}

---

Return the last character of the token list.

मं 1 \aksharStrTail { लीममडमं }

---

\aksharStrChar {(token list)} {(n)}

---

Return the  $n$ -th character of the token list.

3rd character of नमस्कार is स्का.  
It is not स.

```
1 3rd character of नमस्कार is \aksharStrChar{ नमस्कार}{3}.\par
2 \ExplSyntaxOn
3 It~is~not~\tl_item:nn { नमस्कार } {3}.
4 \ExplSyntaxOff
```

---

\aksharStrReplace \aksharStrReplace\* \aksharStrReplace {(tl 1)} {(tl 2)} {(tl 3)}

---

Replace all occurrences of <tl 2> in <tl 1> with <tl 3>, and leaves the modified <tl 1> in the input stream.

The starred variant will replace only the first occurrence of <tl 2>, all others are left intact.

expl3 output:  
स्कास्कास्काडडस्कास्कास्काड  
\aksharStrReplace output:  
स्कास्कास्काडडमंलीममड

```
1 \ExplSyntaxOn
2 \pkg{expl3} ~ output:\par
3 \tl_set:Nn \l_tmpa_tl { मममडडमंलीममड }
4 \tl_replace_all:Nnn \l_tmpa_tl { म } { स्का }
5 \tl_use:N \l_tmpa_tl\par
6 \ExplSyntaxOff
7 \cs{aksharStrReplace} output:\par
8 \aksharStrReplace { मममडडमंलीममड } { म } { स्का }
```

expl3 output:  
स्कांममडडमंलीममड  
\aksharStrReplace\* output:  
ममंस्काडडमंलीममड

```
1 \ExplSyntaxOn
2 \pkg{expl3} ~ output:\par
3 \tl_set:Nn \l_tmpa_tl { मममडडमंलीममड }
4 \tl_replace_once:Nnn \l_tmpa_tl { मम } { स्का }
5 \tl_use:N \l_tmpa_tl\par
6 \ExplSyntaxOff
7 \cs{aksharStrReplace*} output:\par
8 \aksharStrReplace* { मममडडमंलीममड } { मम } { स्का }
```

---

\aksharStrRemove \aksharStrRemove\* \aksharStrRemove {(tl 1)} {(tl 2)}

---

Remove all occurrences of <tl 2> in <tl 1>, and leaves the modified <tl 1> in the input stream.

The starred variant will remove only the first occurrence of <tl 2>, all others are left intact.

expl3 output:  
डडंलीड  
\aksharStrRemove output:  
डडमंलीड

```
1 \ExplSyntaxOn
2 \pkg{expl3} ~ output:\par
3 \tl_set:Nn \l_tmpa_tl { मममडडमंलीममड }
4 \tl_remove_all:Nn \l_tmpa_tl { म }
5 \tl_use:N \l_tmpa_tl\par
6 \ExplSyntaxOff
7 \cs{aksharStrRemove} output:\par
8 \aksharStrRemove { मममडडमंलीममड } { म }
```



[illegible]

`\l_akshar_prev_joining_bool` When we get to a normal character, we need to know whether it is joined, i.e. whether the previous character is the joining character. This boolean variable takes care of that.

(End definition for `\l_akshar_prev_joining_bool`.)

`\l_akshar_char_seq` This local sequence stores the output of the converter.

(End definition for \l akshar char seq.)

```
\l akshar tmpa tl Some temporary variables.
```

```

\l__akshar_tmpb_tl
\l__akshar_tmpa_seq 25 \tl_new:N \l__akshar_tmpa_tl
\l__akshar_tmpb_seq 26 \tl_new:N \l__akshar_tmpb_tl
\l__akshar_tmpe_seq 27 \seq_new:N \l__akshar_tmpe_seq
\l__akshar_tmpe_seq 28 \seq_new:N \l__akshar_tmpe_seq
\l__akshar_tmpe_seq 29 \seq_new:N \l__akshar_tmpe_seq
\l__akshar_tmpe_seq 30 \seq_new:N \l__akshar_tmpe_seq
\l__akshar_tmpe_seq 31 \seq_new:N \l__akshar_tmpe_seq
\l__akshar_tmpe_seq 32 \int_new:N \l__akshar_tmpe_int
\l__akshar_tmpe_int 33 \int_new:N \l__akshar_tmpe_int

```

```

\l__akshar_tmpr_seq
27 \seq new:N \l akshar tmpr seq

```

```
27 \seq_new:N \l__akshar_tmppa_seq
28 \seq_new:N \l__akshar_tmppb_seq
```

```
28 \seq_new:N \l__akshar_tmpd_seq
29 \seq_new:N \l__akshar_tmpe_seq
```

```

29 \seq_new:N \t__akshar_tmpe_seq
\l akshar tmpe seq 30 \seq_new:N \l akshar tmpd seq

```

```

30 \seq_new:N \t__akshar_tmpe_seq
31 \seq_new:N \l__akshar_tmpe_seq

```

```

\l __akshar_tmprb_int 31 \seq_new:N \l__akshar_tmpe_seq
\l__akshar_tmprb_int 32 \int_new:N \l__akshar_tmpe_int

```

```

32 \tnc_new:N \t__akshar_tmpr_tnc
33 \int new:N \l_akshar_tmpr_int

```

(End definition for \l akshar tmpa tl and others.)

In `\akshar_convert:Nn` and friends, the argument needs to be a sequence variable. There will be an error if it isn't.

In `\aksharStrChar`, we need to guard against accessing an ‘out-of-bound’ character (like trying to get the 8th character in a 5-character string.)

In `\aksharStrHead` and `\aksharStrTail`, the string must not be blank.

```

58     must ~ not ~ be ~ empty, ~ but ~ the ~ input ~ string ~ is ~ empty.
59     Make ~ sure ~ the ~ string ~ contains ~ something, ~ or ~ proceed ~
60     and ~ I ~ will ~ use ~ \token_to_str:N \scan_stop:.
61 }

```

### 3.3 Utilities

`\tl_if_in:NoTF` When we get to a character which is not the joining one, we need to know if it is a diacritic. The current character is stored in a variable, so an expanded variant is needed. We only need it to expand only once.

```

62 \prg_generate_conditional_variant:Nnn \tl_if_in:Nn { No } { TF }

```

(End definition for `\tl_if_in:NoTF`.)

`\seq_set_split:Nxx` A variant we will need in `\__akshar_var_if_global`.

```

63 \cs_generate_variant:Nn \seq_set_split:Nnn { Nxx }

```

(End definition for `\seq_set_split:Nxx`.)

`\msg_error:nnx` Some variants of `l3msg` functions that we will need when issuing error messages.  
`\msg_error:nnnxx`

```

64 \cs_generate_variant:Nn \msg_error:nnn { nnx }
65 \cs_generate_variant:Nn \msg_error:nnnnn { nnnxx }

```

(End definition for `\msg_error:nnx` and `\msg_error:nnnxx`.)

`\__akshar_var_if_global:NTF` This conditional checks if #1 is a global sequence variable or not. In other words, it returns true iff #1 is a control sequence in the format `\g_⟨name⟩_seq`. If it is not a sequence variable, this function will (TODO) issue an error message.

`\c__akshar_str_g_tl`  
`\c__akshar_str_seq_tl`

```

66 \tl_const:Nx \c__akshar_str_g_tl { \tl_to_str:n {g} }
67 \tl_const:Nx \c__akshar_str_seq_tl { \tl_to_str:n {seq} }
68 \prg_new_conditional:Npnn \__akshar_var_if_global:N #1 { T, F, TF }
69 {
70   \bool_if:nTF
71   { \exp_last_unbraced:Nf \use_iii:nnn { \cs_split_function:N #1 } }
72   {
73     \msg_error:nnx { akshar } { err_not_a_sequence_variable }
74     { \token_to_str:N #1 }
75     \prg_return_false:
76   }
77   {
78     \seq_set_split:Nxx \l__akshar_tmpb_seq { \token_to_str:N _ }
79     { \exp_last_unbraced:Nf \use_i:nnn { \cs_split_function:N #1 } }
80     \seq_get_left:NN \l__akshar_tmpb_seq \l__akshar_tmpa_tl
81     \seq_get_right:NN \l__akshar_tmpb_seq \l__akshar_tmpb_tl
82     \tl_if_eq:NNTF \c__akshar_str_seq_tl \l__akshar_tmpb_tl
83     {
84       \tl_if_eq:NNTF \c__akshar_str_g_tl \l__akshar_tmpa_tl
85       { \prg_return_true: } { \prg_return_false: }
86     }
87     {
88       \msg_error:nnx { akshar } { err_not_a_sequence_variable }
89       { \token_to_str:N #1 }
90       \prg_return_false:
91     }
92   }
93 }

```

(End definition for `\__akshar_var_if_global:NNTF`, `\c__akshar_str_g_tl`, and `\c__akshar_str_seq_tl`.)

`\__akshar_int_append_ordinal:n` Append st, nd, rd or th to interger #1. Will be needed in error messages.

```

94 \cs_new:Npn \__akshar_int_append_ordinal:n #1

```

```

95 {
96   #1
97   \int_case:nnF { #1 }
98   {
99     { 11 } { th }
100    { 12 } { th }
101    { 13 } { th }
102    { -11 } { th }
103    { -12 } { th }
104    { -13 } { th }
105  }
106  {
107    \int_compare:nNnTF { #1 } > { -1 }
108    {
109      \int_case:nnF { #1 - 10 * (#1 / 10) }
110      {
111        { 1 } { st }
112        { 2 } { nd }
113        { 3 } { rd }
114      } { th }
115    }
116    {
117      \int_case:nnF { (- #1) - 10 * ((- #1) / 10) }
118      {
119        { 1 } { st }
120        { 2 } { nd }
121        { 3 } { rd }
122      } { th }
123    }
124  }
125 }

```

(End definition for `\__akshar_int_append_ordinal:n`.)

### 3.4 The `\akshar_convert:Nn` function and its variants

`\akshar_convert:Nn` This converts `#2` to a sequence of true South Asian characters. The sequence is set to `#1`, which should be a sequence variable. The assignment is local.

```

\akshar_convert:cn
\akshar_convert:Nx
\akshar_convert:cx
126 \cs_new:Npn \akshar_convert:Nn #1 #2
127 {

```

Clear anything stored in advance. We don't want different calls of the function to conflict with each other.

```

128   \seq_clear:N \__akshar_char_seq
129   \bool_set_false:N \__akshar_prev_joining_bool

```

Loop through every token of the input.

```

130   \tl_map_variable:NNn {#2} \__akshar_map_tl
131   {
132     \tl_if_in:NoTF \c__akshar_diacritics_tl {\__akshar_map_tl}
133     {

```

It is a diacritic. We append the current diacritic to the last item of the sequence instead of pushing the diacritic to a new sequence item.

```

134       \seq_pop_right:NN \__akshar_char_seq \__akshar_tmpa_tl
135       \seq_put_right:Nx \__akshar_char_seq
136       { \__akshar_tmpa_tl \__akshar_map_tl }
137     }
138   {
139     \tl_if_in:NoTF \c__akshar_joining_tl {\__akshar_map_tl}
140     {

```

In this case, the character is the joining character, ெ. What we do is similar to the above case, but `\__akshar_prev_joining_bool` is set to true so that the next character is also appended to this item.

```

141       \seq_pop_right:NN \__akshar_char_seq \__akshar_tmpa_tl

```

```

142         \seq_put_right:Nx \l__akshar_char_seq
143         { \l__akshar_tmpa_tl \l__akshar_map_tl }
144         \bool_set_true:N \l__akshar_prev_joining_bool
145     }
146     {

```

Now the character is normal. We see if we can push to a new item or not. It depends on the boolean variable.

```

147         \bool_if:NTF \l__akshar_prev_joining_bool
148         {
149             \seq_pop_right:NN \l__akshar_char_seq \l__akshar_tmpa_tl
150             \seq_put_right:Nx \l__akshar_char_seq
151             { \l__akshar_tmpa_tl \l__akshar_map_tl }
152             \bool_set_false:N \l__akshar_prev_joining_bool
153         }
154         {
155             \seq_put_right:Nx
156             \l__akshar_char_seq { \l__akshar_map_tl }
157         }
158     }
159 }
160 }

```

Set #1 to \l\_\_akshar\_char\_seq. The package automatically determines whether the variable is a global one or a local one.

```

161 \__akshar_var_if_global:NTF #1
162 { \seq_gset_eq:NN #1 \l__akshar_char_seq }
163 { \seq_set_eq:NN #1 \l__akshar_char_seq }
164 }

```

Generate variants that might be helpful for some.

```

165 \cs_generate_variant:Nn \akshar_convert:Nn { cn, Nx, cx }

```

(End definition for \akshar\_convert:Nn. This function is documented on page ??.)

### 3.5 Other internal functions

\\_\_akshar\_seq\_push\_seq:NN Append sequence #1 to the end of sequence #2. A simple loop will do.

```

166 \cs_new:Npn \__akshar_seq_push_seq:NN #1 #2
167 { \seq_map_inline:Nn #2 { \seq_put_right:Nn #1 { ##1 } } }

```

(End definition for \\_\_akshar\_seq\_push\_seq:NN.)

\\_\_akshar\_replace:NnnnN If #5 is \c\_false\_bool, this function replaces all occurrences of #3 in #2 by #4 and stores the output sequence to #1. If #5 is \c\_true\_bool, the replacement only happens once.

The algorithm used in this function: We will use \l\_\_akshar\_tmpa\_int to store the “current position” in the sequence of #3. At first it is set to 1.

We will store any subsequence of #2 that may match #3 to a temporary sequence. If it doesn’t match, we push this temporary sequence to the output, but if it matches, #4 is pushed instead.

We loop over #2. For each of these loops, we need to make sure the \l\_\_akshar\_tmpa\_int-th item must indeed appear in #3. So we need to compare that with the length of #3.

- If now \l\_\_akshar\_tmpa\_int is greater than the length of #3, the whole of #3 has been matched somewhere, so we reinitialize the integer to 1 and push #4 to the output.

Note that it is possible that the current character might be the start of another match, so we have to compare it to the first character of #3. If they are not the same, we may now push the current mapping character to the output and proceed; otherwise the current character is pushed to the temporary variable.

- Otherwise, we compare the current loop character of #2 with the `\l__akshar_tmpa_int`-th character of #3.
    - If they are the same, we still have a chance that it will match, so we increase the “iterator” `\l__akshar_tmpa_int` by 1 and push the current mapping character to the temporary sequence.
    - If they are the same, the temporary sequence won’t match. Let’s push that sequence to the output and set the iterator back to 1.
- Note that now the iterator has changed. Who knows whether the current character may start a match? Let’s compare it to the first character of #3, and do as in the case of `\l__akshar_tmpa_int` is greater than the length of #3.

The complexity of this algorithm is  $O(m \max(n, p))$ , where  $m, n, p$  are the lengths of the sequences created from #2, #3 and #4. As #3 and #4 are generally short strings, this is (almost) linear to the length of the original sequence #2.

```

168 \cs_new:Npn \__akshar_replace:NnnnN #1 #2 #3 #4 #5
169 {
170   \akshar_convert:Nn \l__akshar_tmpe_seq {#2}
171   \akshar_convert:Nn \l__akshar_tmpe_seq {#3}
172   \akshar_convert:Nn \l__akshar_tmpe_seq {#4}
173   \seq_clear:N \l__akshar_tmpe_seq
174   \seq_clear:N \l__akshar_tmpe_seq
175   \int_set:Nn \l__akshar_tmpe_int { 1 }
176   \int_set:Nn \l__akshar_tmpe_int { 0 }
177   \seq_map_variable:NNn \l__akshar_tmpe_seq \l__akshar_map_tl
178   {
179     \int_compare:nNnTF { \l__akshar_tmpe_int } > { 0 }
180     { \seq_put_right:NV \l__akshar_tmpe_seq \l__akshar_map_tl }
181     {
182       \int_compare:nNnTF
183       { \l__akshar_tmpe_int } = { 1 + \seq_count:N \l__akshar_tmpe_seq }
184       {
185         \bool_if:NT {#5}
186         { \int_incr:N \l__akshar_tmpe_int }
187         \seq_clear:N \l__akshar_tmpe_seq
188         \__akshar_seq_push_seq:NN
189         \l__akshar_tmpe_seq \l__akshar_tmpe_seq
190         \int_set:Nn \l__akshar_tmpe_int { 1 }
191         \tl_set:Nx \l__akshar_tmpe_tl
192         { \seq_item:Nn \l__akshar_tmpe_seq { 1 } }
193         \tl_if_eq:NNTF \l__akshar_map_tl \l__akshar_tmpe_tl
194         {
195           \int_incr:N \l__akshar_tmpe_int
196           \seq_put_right:NV \l__akshar_tmpe_seq \l__akshar_map_tl
197         }
198         {
199           \seq_put_right:NV \l__akshar_tmpe_seq \l__akshar_map_tl
200         }
201       }
202     }
203     \tl_set:Nx \l__akshar_tmpe_tl
204     {
205       \seq_item:Nn \l__akshar_tmpe_seq { \l__akshar_tmpe_int }
206     }
207     \tl_if_eq:NNTF \l__akshar_map_tl \l__akshar_tmpe_tl
208     {
209       \int_incr:N \l__akshar_tmpe_int
210       \seq_put_right:NV \l__akshar_tmpe_seq \l__akshar_map_tl
211     }
212     {
213       \int_set:Nn \l__akshar_tmpe_int { 1 }
214       \__akshar_seq_push_seq:NN
215       \l__akshar_tmpe_seq \l__akshar_tmpe_seq
216       \seq_clear:N \l__akshar_tmpe_seq
217       \tl_set:Nx \l__akshar_tmpe_tl
218       { \seq_item:Nn \l__akshar_tmpe_seq { 1 } }
219       \tl_if_eq:NNTF \l__akshar_map_tl \l__akshar_tmpe_tl

```



```

220         {
221             \int_incr:N \l__akshar_tmpa_int
222             \seq_put_right:NV
223             \l__akshar_tmpb_seq \l__akshar_map_tl
224         }
225         {
226             \seq_put_right:NV
227             \l__akshar_tmpa_seq \l__akshar_map_tl
228         }
229     }
230 }
231 }
232 }
233 \__akshar_seq_push_seq:NN \l__akshar_tmpa_seq \l__akshar_tmpb_seq
234 \__akshar_var_if_global:NTF #1
235 { \seq_gset_eq:NN #1 \l__akshar_tmpa_seq }
236 { \seq_set_eq:NN #1 \l__akshar_tmpa_seq }
237 }

```

(End definition for `\__akshar_replace:NnnnN`.)

### 3.6 Front-end $\text{\LaTeX}2_{\epsilon}$ macros

`\aksharStrLen` Expands to the length of the string.

```

238 \NewExpandableDocumentCommand \aksharStrLen {m}
239 {
240     \akshar_convert:Nn \l__akshar_tmpa_seq {#1}
241     \seq_count:N \l__akshar_tmpa_seq
242 }

```

(End definition for `\aksharStrLen`. This function is documented on page ??.)

`\aksharStrChar` Returns the  $n$ -th character of the string.

```

243 \NewExpandableDocumentCommand \aksharStrChar {mm}
244 {
245     \akshar_convert:Nn \l__akshar_tmpa_seq {#1}
246     \bool_if:nTF
247     {
248         \int_compare_p:nNn {#2} > {0} &&
249         \int_compare_p:nNn {#2} < {1 + \seq_count:N \l__akshar_tmpa_seq}
250     }
251     { \seq_item:Nn \l__akshar_tmpa_seq { #2 } }
252     {
253         \msg_error:nnnx { akshar } { err_character_out_of_bound }
254         { #1 } { \__akshar_int_append_ordinal:n { #2 } }
255         { \int_eval:n { 1 + \seq_count:N \l__akshar_tmpa_seq } }
256         \scan_stop:
257     }
258 }

```

(End definition for `\aksharStrChar`. This function is documented on page ??.)

`\aksharStrHead` Return the first character of the string.

```

259 \NewExpandableDocumentCommand \aksharStrHead {m}
260 {
261     \akshar_convert:Nn \l__akshar_tmpa_seq {#1}
262     \int_compare:nNnTF { \seq_count:N \l__akshar_tmpa_seq } = {0}
263     {
264         \msg_error:nnn { akshar } { err_character_out_of_bound }
265         { first }
266         \scan_stop:
267     }
268     { \seq_item:Nn \l__akshar_tmpa_seq { 1 } }
269 }

```

(End definition for `\aksharStrHead`. This function is documented on page ??.)

`\aksharStrTail` Return the last character of the string.

```
270 \NewExpandableDocumentCommand \aksharStrTail {m}
271 {
272   \akshar_convert:Nn \l__akshar_tmpa_seq {#1}
273   \int_compare:nNtF { \seq_count:N \l__akshar_tmpa_seq } = {0}
274   {
275     \msg_error:nnn { akshar } { err_character_out_of_bound }
276     { last }
277     \scan_stop:
278   }
279   { \seq_item:Nn \l__akshar_tmpa_seq { \seq_count:N \l__akshar_tmpa_seq } }
280 }
```

(End definition for `\aksharStrTail`. This function is documented on page ??.)

`\aksharStrReplace` Replace occurrences of #3 of a string #2 with another string #4.  
`\aksharStrReplace*`

```
281 \NewExpandableDocumentCommand \aksharStrReplace {smmm}
282 {
283   \IfBooleanTF {#1}
284   {
285     \__akshar_replace:NnnN \l__akshar_tmpa_seq
286     {#2} {#3} {#4} \c_true_bool
287   }
288   {
289     \__akshar_replace:NnnN \l__akshar_tmpa_seq
290     {#2} {#3} {#4} \c_false_bool
291   }
292   \seq_use:Nn \l__akshar_tmpa_seq {}
293 }
```

(End definition for `\aksharStrReplace` and `\aksharStrReplace*`. These functions are documented on page ??.)

`\aksharStrRemove` Remove occurrences of #3 in #2. This is just a special case of `\aksharStrReplace`.  
`\aksharStrRemove*`

```
294 \NewExpandableDocumentCommand \aksharStrRemove {smm}
295 {
296   \IfBooleanTF {#1}
297   {
298     \__akshar_replace:NnnN \l__akshar_tmpa_seq
299     {#2} {#3} {} \c_true_bool
300   }
301   {
302     \__akshar_replace:NnnN \l__akshar_tmpa_seq
303     {#2} {#3} {} \c_false_bool
304   }
305   \seq_use:Nn \l__akshar_tmpa_seq {}
306 }
```

(End definition for `\aksharStrRemove` and `\aksharStrRemove*`. These functions are documented on page ??.)

```
307 </package>
```