

UTIL: UNIFIED TRANSLITERATION OF INDIC LANGUAGES

:NOTE: Download [PDF](#) of this article. You need good Unicode fonts to read the article. I recommend [Noto font](#) family from Google. You can apt-get it. IPA symbols are within square brackets, like [ʃ] and transliterated symbols are within slashes, like /t/.

UTIL is a romanization scheme for Indic languages. It is designed as pan-Indian transliteration scheme. It covers 20+ languages: Bengali, Dogra, Gujarati, Hindi, Kannada, Kashmiri, Konkani, Lepcha, Limbu, Manipuri (Meitei), Maithili, Malayalam, Marathi, Nepali, Oriya, Punjabi, Sanskrit, Santali, Sindhi, Sinhala, Tamil, Telugu, Urdu and probably many more.

So, why yet another scheme?

1. [IAST](#) is insufficient. It serves Sanskrit and Pali but is incomplete for pretty much everything else (e.g. Bengali, Gujarati, etc.).
2. [ISO-15919](#) is also insufficient. It ignores Kashmiri and Sindhi, which are integral Indian languages. Plus, it lacks symbols for newly-assigned Unicode codepoints (e.g. [𑆚](#) or [𑆛](#)). Also [𑆚](#)/Kṛṣṇa/ is typographically more consistent than /Kṛṣṇa/.
3. [ALA-LC](#) is designed as a single-language model ignoring the inherent similarity of Brahmic scripts. This leads to inconsistencies. For example, Tamil [𑤢](#) /l/ and Kannada [𑤢](#) /l/ correspond to same character and sound (“retorflex approximant”) and yet have different representations. Conversely, the same symbol /ʃ/ represents Hindi [𑆚](#) [ʃ] and Urdu [ص](#) [sʰ] even though they’re completely different sounds.
4. Other schemes like Hunterian or Gretil are as bad as the above or even worse sometimes.

So, how is UTIL better?

- Covers the entire character set of ISO15919 plus more (Kashmiri, Sindhi)
- Long vowels always have macron above (ā, ē, ...)
- Aspirated consonants always have ‘h’ as second letter (kh, gh, ...)
- Minimum number of diacritical marks:
 - Only four diacritics are used: “dot above”, “dot below”, “macro above”, “macro below” (or their combination).
 - Only three diacritics are needed for Sanskrit, instead of IAST’s five.

- Prefers precomposed characters in Unicode repertoire, but not required.

VOWELS

Primary vowels and diphthongs:

अ	आ	इ	ई	उ	ऊ	ऋ	ॠ	ऌ	ॡ	ऐ	ए	ऐ	ओ	औ	औ
a	ā	i	ī	u	ū	r̥	r̄	l̥	l̄	e	ē	ai	o	ō	au

Additional ones, all have a dot below:

अं	आं	एँ	अँ	औँ	अु	अु
ạ	ạ̄	ẹ	ọ	ọ̄	ụ	ụ̄

CONSONANTS

Consonants with their Sanskrit names:

	Plosives				Nasal	Implosives	Fricatives		Vibrants		Approximants	
	स्पर्श				नासिक	विस्पर्श	ऊष्मन्		द्रव		अन्तस्थ	
कण्ठ्य Velar	क	ख	ग	घ	ङ	ग	ख	ग			ह	
	k	kh	g	gh	ṅ	ḡ	kh	g̃			h	
तालव्य Palatal	च	छ	ज	झ	ञ	ञ	श	झ			य	य
	c	ch	j	jh	ñ	j̣	ś	zh			y	ỵ
मूर्धन्य Retroflex	ट	ठ	ड	ढ	ण		ष		ड	ढ	ळ	ळ
	ṭ	ṭh	ḍ	ḍh	ṇ		ṣ		r̥	r̄h	ḷ	ḷ̣
दन्त्य Dental	त	थ	द	ध	न		स				ल	
	t	th	d	dh	n		s				l	
वर्त्स्य Alveolar	च	छ			न	ड	स	ज	र	र		
	ç	çh			ṅ	ḍ	s	z	r	r̄		
ओष्ठ्य Labial	प	फ	ब	भ	म	ब	फ				व	व
	p	ph	b	bh	m	ḅ	f				v	w

Affricate glide ष ('JJYA') is transcribed /j/.

OTHER SYMBOLS

Anusvāra: ṁ	Anunāsika: ◌̣	Avagraha: '̣
Visarga: ḥ	Jihvāmūliya: ḥ̣	Upadhmaniya: ḥ̣̣

Vedic Udātta: ´	Svarita (independent): `	Anudātta: _
Arabic hamza ء: ˀ	Arabic ain ع: ˁ	
Rising tone: ˊ	Falling tone: ˋ	

Udātta and svarita use *combining* grave and acute accent respectively. Whereas hamza and ain use *non-combining* modifier letters U+02BC and U+02BD respectively. Tone modifiers are used in Maithili, Dogra and other Pahari languages.

GENERAL NOTES

- Anunāsika is denoted by a combining candrabindu. Note the difference between हंस (swan) /haṃsa/ and हँस (laugh) /hāṃs/. Diacritic only on second letter in a digraph. Example: हैँ /haiṅ/
- A colon is used to denote vowel hiatus or resolve ambiguity. Example: बई /ba:i/ (not /bai/)

SCRIPT NOTES

- ऋ /ṛ/, ॠ /ṝ/ and ॡ /ṝ̄/ are used only in Sanskrit.
- ऐ /e/ = short ए in Southern scripts (எ, എ, ಎ, ೆ)
- ओ /o/ = short ओ in Southern scripts (ஒ, ഒ, ೆ, ೆ)
- ે /e/ = Gujarati એ, Sinhala ે, pronounced [æ] as in “bat”
- ો /ō/ = Gujarati ઓ, pronounced [ɔ:] as in “ball”
- ৳ /ra/ = Bengali ৳, Punjabi ੳ, Oriya ୳ (RRA, “retroflex flap”)
- ৳̄ /r̄ha/ = Bengali ৳̄, Oriya ୳̄ (RHA, “aspirated retroflex flap”)
- 𑌗 /la/ = used in Marathi, Tamil ள, Malayalam ല, Kannada ಳ, Telugu ళ (LLA, “retroflex lateral approximant”)
- 𑌗̄ /l̄a/ = Tamil ள̄, Malayalam ല̄, Kannada ಳ̄, Telugu ళ̄ (LLLA, “retroflex approximant” = zha)
- 𑌎 /na/ = Tamil ன, Kannada ನ, Malayalam ണ (NNNA, “alveolar n”)
- 𑌖 /ra/ = Tamil ற, Malayalam റ, Kannada ರ, Telugu ర (RRA, “alveolar r”)
- 𑌖̄ /r̄/ = repha in Marathi
- 𑌙 /ya/ = য in Bengali and Oriya, while 𑌙 /y/ = য
- 𑌚 /wa/ = Urdu و, Assamese ব, Oriya ୱ
- 𑌘 /k̄ha/ = used in Urdu, Punjabi 𑌘
- 𑌙 /ga/ = used in Urdu, Punjabi 𑌙
- 𑌚 /ca/ = used in Kashmiri, Telugu డ
- 𑌚 /za/ = Urdu ز, Gurmukhi ਜ, Bengali জ, Kannada ಜ, Telugu జ [d̪z]
- 𑌚, 𑌚 /zh/ = Urdu ذ, Gujarati જ, Avestan uses 𑌚
- Kashmiri vowels 𑌚, 𑌚̄, 𑌚̄̄, 𑌚̄̄̄, 𑌚̄̄̄̄ are pronounced [ə], [ə:], [ɔ], [ɔ:], [i], [i:].

Another vowel form, औ, is sometimes used for [ɔ] (and [ɔ:] is skipped altogether). These symbols are taken from ALA-LC as it is and follows Wikipedia. Kashmiri consonants: च [tʃ], छ [tʃʰ] and ज [z]

- Sindhi implosives: ڱ /g/, ڃ /j/, ڌ /d/ and ڙ /b/
- Sinhalese nasals: ඞ /ŋga/, ඞ් /ŋja/, ඞ් /nda/, ඞ් /nda/ and ඞ් /mba/
- Sinhala long vowel ඌ and Devanagari vowel sign candra long E (U+0955), used in Avestan, are transliterated /ē/
- The six Malayalam chillu characters represent dead consonants (without implicit vowel). As such, they are simply transliterated without adding an a next to the consonant. Hence, ഛ, ണ, ശ, റ, റ and റ are respectively transliterated as /k-/ , /l-/ , /ll-/ , /n-/ , /nn-/ and /rr-/.

URDU

Perso-Arabic characters are chosen in a non-conflicting way with the Brahmic scripts. Urdu introduces six sounds [f, z, ʒ, q, x, ɣ] on top of Hindi (see [Hindustani phonology](#)). Note that [f, z, x, ɣ] are fricatives, just like ष [ʃ], स [s], ह [h]. Excluding these IPA signs, the ones in the below table are indicative only.

Urdu	ق	ح	خ	ء	ع	غ	ط	ظ	ز	ذ	ض	ص	ث	ش	ژ	ف	و
UTIL	q	h	kh	'	'	g	t	z	z	z	z	s	s	s	zh	f	w
IPA	[q]	[ħ]	[x]	[ʔ]	[ʔ]	[ɣ]	[tʃ]	[zʃ]	[z]	[ð]	[dʃ]	[sʃ]	[θ]	[ʃ]	[ʒ]	[f]	[w]
Devanagari	क	ह	ख	?	?	ग	त	ज	ज	ज	ज	स	स	श	झ	फ	व

INPUT METHODS FOR IME:S

Of course, a transliteration scheme is not fruitful if it cannot be entered into a computer, for which Input Method Editors (IMEs) are used. This can be thought of as an ASCII transliteration of UTIL.