

OV0

The One Voice Orthography Romanization Scheme

Lance Pollard

`earth@mount.build`

The One Voice Orthography (OV0) is a romanization scheme to allow reading, writing, and pronouncing most human languages on Earth. The goal of this system is to bridge the gap between the International Phonetic Alphabet (IPA) and the general public's use of the Latin alphabet, making it easy for the general public and language learners to read and pronounce words across languages. It makes a tradeoff in accuracy for ease of use, getting you to a close enough approximation of a pronunciation without being overburdeningly accurate. It does this with only using the Latin alphabet and basic diacritics, and handles consonants, vowels, stress, clicks, and tones across the major world languages.

Background

There are thousands of languages which use the Latin script for writing, and currently over 100 romanization schemes for languages which don't use the Latin script. Each of these Latin script usages varies slightly in how they use the letters and how they pronounce the letters in various contexts. The result is that there are hundreds if not thousands of variations in how the Latin script is used across languages in today's conventions, making it a complex and hard task to learn how to read and pronounce words from foreign languages. The romanization schemes don't work across languages.

IPA is the main linguistics system for representing sounds using letter-like symbols. It identifies key points in a fluid spectrum of mouth positions and articulations which are observed in actual languages on Earth. As such, it is not entirely accurate, as it is approximations of a continuous range of mouth shapes and articulations. It is, however, the best system we have today for concisely writing down the pronunciation of words.

However, from anecdotal evidence and primitive interviews, it appears the general public does not have an interest in learning to read IPA and use it for pronunciation. This seems to be because it is a highly domain-specific language (DSL) used by highly trained linguists or amateur language hobbyists. As such, it appears IPA is too detailed for beginners to learn okay-enough pronunciations of words in foreign languages (though not detailed enough if you are coming from the perspective of programmatic speech synthesis). This is largely the reason why we have created a simplified version of IPA for pronunciation purposes. OVO is almost as good as IPA, but makes a few tradeoffs in accuracy for simplicity in understanding and usage.

Implementation

OVO is a Latin orthography that represents all the basic sounds using the Latin letters, with other basic sounds and sound variations using diacritics. It is not a perfect system for representing sounds, and the diacritics are necessary for expanding the 26 Latin characters into a

larger character set to account for other sounds. Nevertheless, it is advantageous because the Latin alphabet is the most commonly used alphabet in the world.

The way this works is simple. We have identified the most clearly distinguishable sounds and mapped them as best as possible to the Latin alphabet. By "clearly distinguishable sounds", we draw the following analogy. Imagine you are trying to distinguish sounds and write letters for the sounds. You would highly probably identify the following sounds as clearly distinct sounds: a knock on a door, a human scream, the sound of rain, the sound of thunder, the sound of a bird, etc.. However, you might not be able to clearly distinguish between a rapid knock and a slow knock on the same door, or a soft knock and a hard knock. Likewise, you might not be able to distinguish between different bird sounds even from different birds. That is to say, certain sounds everyone can clearly distinguish as distinct, with high probability, while other sounds people might not be able to with high probability. In this way, we have matched the sounds to letters. Sounds that are clearly distinguishable and common across languages are given top-priority and mapped to letters. Other sounds which are clearly distinguishable but uncommon, like clicks, are made into diacritics. Also, sounds that are less distinguishable other than to native speakers or experts are relegated to diacritics, such as a nasal vowel sound or a retroflex consonant.

By establishing the romanization system in such a way, we are able to make it so less-skilled (i.e. beginner) language speakers can get pretty close to the correct pronunciation by just knowing the basic letters without diacritics. Then if they learn the diacritics, they will have a more highly accurate pronunciation. This way the system requires minimal learning up front, and just by learning the basic letters you can start pronouncing a wide variety of Earth's languages with a reasonable quality of sound. For those who choose to learn the complete system (representing most of IPA), they will be able to understand how the diacritics modify the letters and so be able to pronounce words with even higher accuracy.

The sounds are shown in the table in the appendix. To summarize that, let's start with the vowels.

There are 5 vowel letters in the roman alphabet: i, e, a, o, u. These 5 vowels take the most common vowel sounds they typically represent. The i is the English "ee" sound, the e is the English "ey" sound, the "a" is the English "ah" sound, the o is the English "oh" sound, and the u is the English "oo" sound. Then there are 5 more common vowel sounds at least in the English language which are represented by adding a dot below the vowel letter. These are: ì as in the English word "bit", ẹ as in the English word "pet", ą as in the English word "cat", ɒ as in the English word "book", and

u as in the English word "hut". Also of particular interest is the English r sound, which we treat as a vowel, since it sounds like a vowel and most languages use r as a rolled sound. This English r is represented with the ring below, as in e.

There are 21 lowercase consonant letters in the Roman alphabet, though we have identified several more than 21 consonant sounds. While the consonant sounds are also included in the appendix, we will go over some of the key ones here. The letters b, d, f, h, k, l, m, n, p, s, t, v, w, y, and z are all like the English sounding equivalents. The letter g is the sound as in the English word "good". The letter j is like the "zh" sound in the English word "measure". Then there are 3 remaining letters, c, q, and x. These take on a unique meaning. The letter x is used to represent the English "sh" sound like in the word "ship". The c is used for the voiceless "th" sound as in "think". And the q is even more far removed, in that it represents the English "ng" sound as in "king". That does it for the consonant letters without diacritics.

Then we add our first layer of diacritics to achieve some more consonants sounds. We will just outline a few here, the rest being in the appendix. First are those with a dot below. The ḍ, ḷ, ṇ, ṙ, ṭ, and ḡ, are like the retroflex Indic consonants. The ƒ and Ƴ are slight variations on the f and v sounds. The ƙ is a deeper k sound like used in Arabic. The ħ is a harsher h sound, and the double dotted ḥ is the extra harsh H sound like in Hebrew and Arabic. Finally, the ʈ sound is the voiced "th" sound as in the English word "these".

Next there are 3 classes of consonant diacritics to handle some less common but present sounds: the ejective consonants, the implosive consonants, and the clicks. The ejective consonants have a forward-tilting slant on them where appropriate, as in k^{h} . The implosive consonants have a backward-tilting slant as in g^{h} . Lastly, the clicks have a ring on them, as in p^{h} or t^{h} . Now, the clicks don't directly correspond to the underlying Latin consonant letter. The click letters were picked because they were the closest or most representative place in the mouth where the click is produced, so you can have that mental association, while at the same time keeping all clicks using the ring diacritic for consistency.

Finally, there are 2 classes of vowel diacritics which are important: the stress marker, and the tone markers. There is only a higher stress (whereas in IPA they also have the uncommon low-stress marker). It is represented with a right-tilting slant above, as in á. Then the tones work differently from Chinese Pinyin but use it as a source of inspiration. In human language on Earth, tones essentially boil down to 3 levels and motion between the levels. They can be regular tone, low tone, and high tone, and they can shift from one level to the next (high-to-low tonal transition,

low-to-high, high-to-low-to-high, high only, etc.). As such, we mark each vowel with a single tone, and duplicate the vowel and mark it for changes in tone. So for example, the Chinese word for horse, written in Pinyin as "mǎ", has a high-to-low-to-high tone, and so would be written mǎââ.

Other than that, there is the tilde added below the vowel for nasal sounding vowels. And do note, diacritics can be combined like combining stress with tone into a single letter, you just stack them. Also note that there are several IPA constructs such as aspiration, velarization, long consonants and vowels, etc. that can be accomplished as well. These are summarized as follow.

- IPA ː suffix means double the letter pattern, for a long sound, on vowels and even on consonants.
- IPA ~ tilde diacritic means adding a tilde suffix to the letter, for a nasal sound.
- IPA ɹing diacritic below a letter means a voiceless sound, which is represented by adding an "h" after the letter.
- IPA ɹ suffix means add a "y" after the letter pattern. You'll notice some IPA letters inherently include this "y" sound in them as well.
- IPA w suffix means add a "w" after the letter pattern.
- IPA ʰ suffix means add a single "h" after the letter pattern, for an aspirated "h" sound.
- IPA ʰ suffix means add a double "hh" after the letter pattern, for a breathy "h" sound.
- IPA ɣ suffix means adding an "ɣ" after the letter pattern, for the velarized sound.
- Implosive consonants add a backward-tilting slant to the letter.
- IPA ' for ejective consonants means adding a forward tilting slant to the letter.
- IPA ̤ suffix means adding a '̤' after the letter pattern, for the glottalized sound.
- IPA letter bridge like d̤ here means the same thing as without the bridge like dɹ.

The main problem with this pronunciation romanization system is that it is difficult to type on a standard keyboard. So we introduce OVO-ASCII as well as a way of writing the OVO script using ASCII, which can then be programmatically converted into regular OVO. We do this to make it easier to type. The basic relationship between IPA, OVO, and ASCII OVO is given in the appendix. Because file systems often are case-insensitive, you can replace capital letters in OVO-ASCII to lowercase letters followed by an underscore, so R becomes r_, etc.

Conclusion

Much thought has been put into designing a romanization system that is easy to understand and use for those who already know the Latin alphabet. It gives you a good enough approximation of a word's pronunciation, making it easy for language learners to read, write, and pronounce words in foreign languages to a high degree of accuracy, without over-burdening them with obscure nomenclature, symbols and sounds. The result is OVO, the One Voice Orthography, a system for reading, writing, and pronouncing the human languages on Earth using the Latin alphabet. This system will serve as the foundation for a new writing script called Hanákana which we save for a future paper.

Appendix

IPA to OVO Map

i	ᵢ	ɪ	b	ᵇ	b?
ɪ	ᵢ̇	I	B	ᵇᵇ	b!b!
ɪ̥	ᵢ̥̇	i@	d	d	d
e	e	e	ɖ	ᵈ	D
ɛ	ᵉ	E	ʈ	ᵈ̥	d*
ø	ᵉ̥	e@	θ	c	c
ə	ᵉ̥̥	E@	ð	ç	C
ɜ	ᵉ̥̥̥	E@	f	f	f
ɜ̥	ᵉ̥̥̥̥	E@	φ	f̥	F
ɑ	a	a	g	g	g
æ	ᵃ	A	ɠ	g	g
ɒ	ᵃ̥	A	ɡ̊	ḡ	g?
o	o	o	ɣ	ḡ	g?
ə	ᵒ̥	O	ʝ	gy	gy
ʊ	ᵒ̥̥	O	f	ḡy	g?y
ʏ	ᵒ̥̥̥	O	h	h	h
ɔ	ᵒ̥̥̥̥	o@	ħ	h"	h"@
u	u	u	ɦ	hh	hh
ʌ	ᵘ	U	x	ᵝ	H
ə	ᵘ̥	U	χ	ᵝ̥	h@
ʉ	ᵘ̥̥	u@	ç	hy	hy
ʊ	ᵘ̥̥̥	u@	ʒ	j	j
æ	ᵘ̥̥̥̥	U@	ʐ	j̥	J
b	b	b	ʑ	ĵ	j@

Ƶ	jy	jy	Ÿ	ṙ	r@
k	k	k	Ɓ	ṙ	r@
k'	ķ	k!	Ɲ	ṙṙ	r@r@
!	ķ̇	k*	s	s	s
q	ķ	K	t	t	t
q'	ķ	K!	t	ṭ	T
m	m	m	l	ṭ̇	t*
n	n	n	v	v	v
ŋ	ŋ	N	u	ṽ	V
ŋ	q	q	β	ṽ	V
ɴ	q	q	w	w	w
ɲ	ny	ny	м	wh	wh
l	l	l	ш	Ẁ	W
l	l̇	L	ʃ	x	x
l̇	l̇	l@	ʂ	ẁ	X
l̇	lṙ	lr@	ʐ	ẁ	x@
ʌ	ly	ly	j	y	y
ll	l̇	l*	ĵ	y	y
p	p	p	y	y	y
p'	ṑ	p!	ყ	yw	yw
ø	ṑ̇	p*	z	z	z
r	r	r	ʔ	'	'
ṙ	ṙ	R	ʃ	"	"

Romanization Schemes

Arabic	Deutsche Morgenländische Gesellschaft (1936)
Arabic	BS 4280 (1968)
Arabic	SATTS (1970s)
Arabic	UNGEGN (1972)
Arabic	DIN 31635 (1982)
Arabic	ISO 233 (1984)
Arabic	Qalam (1985)
Arabic	ISO 233-2 (1993)
Arabic	Buckwalter transliteration (1990s)
Arabic	ALA-LC (1997)
Persian	DMG (1969)
Persian	ALA-LC (1997)
Persian	BGN/PCGN (1958)
Persian	EI (1960)
Persian	EI (2012)
Persian	UN (1967)
Persian	UN (2012)
Amharic	BGN/PCGN (1967)
Armenian	Hübschmann-Meillet (1913)
Armenian	BGN/PCGN (1981)
Armenian	ISO 9985 (1996)
Armenian	ALA-LC (1997)
Greek	ALA-LC

Greek	Beta Code
Greek	ELOT
Greek	Greeklish
Hebrew	ANSI Z39.25 (1975)
Hebrew	UNGEGN (1977)
Hebrew	ISO 259 (1984)
Hebrew	ISO 259-2 (1994)
Hebrew	ISO/DIS 259-3
Hebrew	ALA-LC
Indic scripts	ISO 15919 (2001)
Indic scripts	The National Library at Kolkata romanization
Indic scripts	Harvard-Kyoto
Indic scripts	ITRANS
Indic scripts	ISCII
Mandarin	ALA-LC
Mandarin	EFEO
Mandarin	Lessing-Othmer
Mandarin	Latinxua Sin Wenz
Mandarin	Postal romanization
Mandarin	Wade-Giles
Mandarin	Yale (1942)
Mandarin	Legge romanization
Mainland China	Hanyu Pinyin (1958)
Mainland China	ISO 7098 (1991)
Taiwan	Gwoyeu Romatzyh
Taiwan	Mandarin Phonetic Symbols II (MPS II 1986-2002)
Taiwan	Tongyong Pinyin (2002-2008)

Taiwan	Hanyu Pinyin (since January 1 2009)
Cantonese	Barnett-Chao
Cantonese	Guangdong (1960)
Cantonese	Hong Kong Government
Cantonese	Jyutping
Cantonese	Meyer-Wempe
Cantonese	Sidney Lau
Cantonese	Yale (1942)
Cantonese	Cantonese Pinyin
Japanese	Revised Hepburn
Japanese	Kunrei-shiki
Japanese	Nihon-shiki
Japanese	Hepburn
Japanese	Kunrei-shiki
Japanese	JSL
Japanese	Wāpuro
Japanese	ALA-LC
Korean	McCune-Reischauer (MR; 1937?)
Korean	Yale (1942)
Korean	Revised Romanization of Korean (RR; 2000)
Korean	ISO/TR 11941 (1996)
Korean	Lukoff
Thai	Royal Thai General System of Transcription
Thai	ISO 11940 1998 Transliteration
Thai	ISO 11940-2 2007 Transcription
Thai	ALA-LC
Russian	BGN/PCGN (1947)

Russian	GOST 16876-71 (1971)
Russian	United Nations
Russian	ISO 9 (1995)
Russian	ALA-LC (1997)
Ukrainian	ALA-LC
Ukrainian	ISO 9
Ukrainian	Ukrainian National transliteration
Tibetan	Tibetan pinyin or ZWPY
Tibetan	THL Phonetic Transcription
Tibetan	Tise
Tibetan	Wylie transliteration

Languages Using Latin Script

Acehnese	Croatian	Hmong
Afar	Cree	Hungarian
Afrikaans	Czech	Icelandic
Albanian	Danish	Ido
Aragonese	Dayak	Igbo
Asturian	Dutch	Ilocano
Aymara	English	Indonesian
Azeri	Esperanto	Interlingua
Banjar	Estonian	Innu-aimun
Basque	Faroese	Irish
Belarusian	Fijian	Italian
Betawi	Finnish	Javanese
Berber / Tamazight	French	Judeo-Spanish
Bislama	Fula	Kabylia Berber
Boholano	Gaelic	Khasi
Bosnian	Galician	Kazakh
Breton	German	Kinyarwanda
Catalan	Gikuyu	Klingon language
Cebuano	Guaraní	Kirundi
Chamorro	Haitian Creole	Kongo
Cherokee	Hausa	Konkani
Cornish	Hawaiian	Kurdish
Corsican	Hiri Motu	Latin

Latvian	Oromo	Tagalog
Laz	Palauan	Tahitian
Leonese	Picard	Tetum
Lingala	Polish	Tok Pisin
Lithuanian	Portuguese	Tongan
Luganda	Quechua	Tsonga
Luxembourgish	Rohingya	Tswana
Maori	Romanian	Tunisian Arabic
Malagasy	Romansh	Turkish
Malay	Samoan	Turkmen
Maltese	Sasak	Turoyo
Manx	Saterland Frisian	Uzbek
Marshallese	Scots	Venda
Mauritian Creole	Serbian	Vietnamese
Minangkabau	Seychellois creole	Vastese
Moldovan	Shona	Volapük
Montenegrin	Slovak	Võro
Nahuatl	Slovene	Walloon
Nauruan	Somali	Welsh
Navajo	Sotho	West Frisian
Nias	Sotho	Wolof
Ndebele	Spanish	Xhosa
Ndebele	Sundanese	Yoruba
North Frisian	Swahili	Zazaki
Norwegian	Swedish	Zulu
Occitan	Swati	