

Аббр и проб ставятся перед аббревиатурой и именем собственным

Технический момент: кодировка Вокок

Символы Вофала могут быть записаны в собственной кодировке *Вокок/Wokok* (WOfal KOmplex enKoding) или в Юникоде. В последнем случае систему кодирования будем называть *Вун/Wun* (Wofal UNicode), в ней моды и тюнеры отсутствуют, обратная конвертация в иероглифы невозможна, а в нацик – только если чудесным образом известно в какой, и если одно имаго не записывается в том двумя эмбами. Файлы с Вококом пусть имеют расширение .wkk, с Вуном – .wun

В Вококе патроны простираются в диапазоне 0x00-0x4f (из них цифры – 0x00-0x09, буквы – 0x0a-0x2d), моды – 0x80-0xfe (т.е. 0x80 – это и есть мода ноль). Коды записываются в интелловском, а не в моторолловском формате: младший бит числа расположен ближе к младшим адресам памяти компьютера. Разумеется, “с” – не буква, а шестнадцатеричная цифра 12.

Для уменьшения количества машинных операций при обработке текста, вместе в ряды сгруппированы символы, означающие окончание слова (0x4f-0x7e) или влияющие на отображение последующего символа (аббр и проп влияют, но являются частью слова: вообще все символы 0x0a-0x51 – часть слова). Для уменьшения количества операций преобразования чисел из текстового вида в бинарный и обратно, цифрам присвоены равные им кода. Три раскладки клавиатуры выделены голубым, синим и темно-синим.

	.0	.1	.2	.3	.4	.5	.6	.7	.8	.9	.a	.b	.c	.d	.e	.f
0.	0	1	2	3	4	5	6	7	8	9	a	b	c	d	e	f
1.	g	g	h	h	i	i	h	k	L	m	n	o	o	p	r	s
2.	ş	ş	t	u	v	w	g	z	z	z	'	_	-	hil		
3.				ø	у	ø									>	<
4.																
5.	abbr	prop	.	,	-	:	?	!	/							
6.																
7.																
8.																
9.																
a.																
b.																
c.																
d.																
e.																
f.																NiG

Долготу последней фонемы слова невозможно обозначить двоеточием после буквы (т.к. в этой позиции оно является пунктуационным знаком), а обозначение долготы удвоением буквы мешает заметить одноморфемные слова и бессознательно распознать целое слово как последовательность сгущения и разряжения чернил. Поэтому после буквы ставим подчеркивание

пробела («_»), которое будем называть *дуратором/durator* – чтобы каждый раз не уточнять подчеркивание графемы или пробела. Заметим, в арабском шадда ّ должна конвертироваться раньше огласовки, т.к. первый эмб под ней может быть гласной, а второй – согласной как в сочетаниях ّو (u_w), ّي (i_н).

Введем понятие *байтема/byte me* – последовательность байт, представляющая сигил в той или иной кодировке. Состоит из *патрона/patron* – номера сигила; *моды/mode* – номера, указывающего, какому символу нациков соответствует сигил (например, ٤ и ٥ имеют байтемы t4 и t19 в Вофале. Разумеется, непрограммисты видят только патроны. Spell-checker-ы, орфографические корректоры подставляют невидимые моды по словарю конкретного языка.

Прямая, а затем обратная конвертация дают первоначальный текст. В хорошо урегулированном языке прямая делается с помощью регулярных выражений, в плохо – является подстановкой словарной транскрипции с заменой нескольких транскрипционных знаков снова на буквы, а именно тех, которые звучат по-разному в однокоренных словах или словоформах: например, нельзя в конце слов в русском, немецком, нидерландском заменять звонкие согласные на глухие; в фарси – смыгчать (палатализировать) «k, g» и т.д. А обратная конвертация в любом случае держится только на модах.

1. Обратная конвертация

- 1.1. мода ноль – особая и показывает, что графема должна быть отброшена при обратной конвертации (“t6” из t6s0 конвертируется в «ц», а “s0” отбрасывается)
- 1.2. фонеторы при обратной конвертации отбрасываются, поэтому их неуказанные моды равны нулю
- 1.3. мима удаляется при обратной конвертации
- 1.4. моды, соответствующие графемам разных нациков, должны быть одинаковы, если одинаковы графемы и эмбов, и имаго
- 1.5. если файл начинается с моды ноль, то младшие 15 бит следующих двух байт содержат языковой ключ/*language key* – код языка по ISO 639-3, переведенный из 26-ричной системы счисления (например, “jav” для яванского дает $10 \cdot 26^2 + 1 \cdot 26 + 22 = 6808$); и тогда текст не содержит мод за исключением позиций, где они нужны для разрешения неоднозначности в пределах одного языка (это даже уже, чем в пределах одного нацика); т.о. в каждом языке моды делятся *решающие/solving* и *факультативные/facultative*

2. Некоторые патроны

- 2.1. графемой ʌ обозначим фонему [ʌ] (русскую команду лошади остановиться)
- 2.2. е – [ɛ]: русское «э», украинское «е»
- 2.3. h – [x]; а также [χ] в случае, если в языке есть мфа-шное [ħ], уже обозначенное как «hʒ»: русское “х” [x], чешское и польское “ch” [x] (и устаревшее индонезийское), суахили “kh” [x], арабское ح [χ]
- 2.4. ħ – [χ]; а также [ʁ] в случае, если в языке есть [ʁ], уже обозначенное как «ħʒ»: суахили “gh” [χ], арабское ح [ʁ]
- 2.5. ě – [ɛ̃, h]: чешское и суахили “h”, картавое украинское “г” [h]; арабское ه [h]
- 2.6. g – [g]: общелатинское “g”, русское “г”, украинское “г”
- 2.7. i – [i]: общелатинское “i”, русское “и”
- 2.8. ı – [ɨ, ɪ]: турецкое “ı”, русское “ы”, вьетнамское “ư” [ɨ]; игбо “i” [ɪ], дари [ɪ]
- 2.9. ʏ – [j]: русское “й”, арабское ʔ (например, «е, я, ё» – это русские “е, я, ё” после твердого, мягкого знаков и гласной, “ю” и украинское “ї”: “съел, семья, моя, съёмка, юла” кодируются как «sʔel, sʔemʔa, moʔa, sʔomʔka, ʔula»)

- 2.10. о – [ʏ, o]: при наличии обоих фонем этой буквой обозначается только первая, а вторая [o] – как оʔ
- 2.11. ɔ – [ʌ, ɔ]: при наличии обоих фонем этой буквой обозначается только первая, а вторая [ɔ] – как ɔʔ
- 2.12. r – дрожащие русское “р” и арабское ر [r], недрожащие английское [ɹ] и американское [ɹ] и т.д.
- 2.13. l – общелатинское “l”, русское “л” и т.д.
- 2.14. z – [z]: русское “з”
- 2.15. š – [ʃ]: турецкое ş, сербско-кириллическое и чешское š, сербско-латинянское и русское “ш”, английское sh, немецкое sch, арабское ش
- 2.16. ž – [ʒ]: русское “ж”, сербско-кириллическое и чешское ž, польское ż, английское zh, персидское ژ
- 2.17. ſ – [θ]: (должно быть \$, не стандартизированное Юникодом, ведь различительные возможности неразрывных с буквой нижних видоизменений исчерпаны уже после первого нижнего видоизменения): английское и суахили th, арабское ث, греческое θ
- 2.18. ž – [ð]: (пока Ž не стандартизировано): английское th, суахили dh, арабское ذ, греческое δ
- 2.19. w – [w, ʍ] общелатинское “w”, белорусское “ў”, арабское و
- 2.20. апостраф (') является буквой и обозначает глоттальную смычку, в частности
- 2.20.1. в него конвертируются арабские хамза ء, алиф-хамза إ и إ, вав-хамза و, алифМаксура-хамза ع
3. Дифтонги, трифтонги представляем несколькими буквами, только если удастся подобрать такую их последовательность, прочтение которой носителем языка дает носителю возможность узнать дифтонг, в частности
- 3.1. dz – пуштунское “ځ”, черногорско-кириллическое “з”, черногорско-латинянское “s”, белорусское “дз”
- 3.2. dž – арабское “ج”, английское “j”, тамильское ஜ, сербско-кириллическое “џ”, белорусское “дж”, сербско-латинянское “dž”
- 3.3. ts – пуштунское “څ”, русское и сербско-кириллическое “ц”, чешское и сербско-латинянское “c”
- 3.4. tš – [tʃ] персидское “چ”, английское “ch”, немецкое “tsch”, чешское и сербско-латинянское “č”, сербско-кириллическое “ч”
4. Фонеторы
- 4.1. ʔ – лабиатор/labiator; превращает в лабиализованные
- 4.1.1. eʔ, iʔ – германские, тюркские, угорские “ö” [ø, œ], “ü” [y, y]
- 4.1.2. uʔ – игбо “u” [u], дари [u]
- 4.2. ʕ – эмфатор/emfator; в эмфатические с артикуляцией увулярной или фарингальной: d/dʕ обозначают د/ض, t/tʕ – ط/ت, z/zʕ – ظ/ز, s/sʕ – ص/س; k/kʕ [k/q] – ك (к в персидском, пушту, урду)/ق, к/қ в казахском; h/hʕ [x/ мфа-шное ħ, похожее на громкий шепот] – ح/خ, ħ/hʕ [ħ/ħ] – غ/ع; ɣ/g/ɣʕ – ng/nɣ в серер (дрожащие русское [r]/французское [ʀ] обозначаем не r_/rʕ_, а одной буквой “r”, равно как и недрожащие английское [ɹ]/немецкое [ʁ] не r/rʕ – тоже одной “r”, ибо в рамках одного языка им не противопоставлены парный по эмфатичности или длительности)
- 4.3. ɖ – ретр/retr; в ретрофлексные: в пушто d/dɖ – د/ځ, t/tɖ – ت/ټ, r/rɖ – ر/ړ, n/nɖ – ن/ښ; в урду d/dɖ [d/ɖ] – د/ځ, t/tɖ [t/ɖ] – ت/ټ, r/rɖ [r/ɖ] – ر/ړ; в хауса r/rɖ

- 4.4. >, < – начальный и конечный *сплитеры/spliter*; участки слова от начального сплитера до начала слова (пробела или мима) и от конечного сплитера до конца слова читаются при определенных условиях: например, в арабском когда перед начальным сплитером стоит алиф-васла (إ), оно читается, только если предыдущее слово оканчивается на согласную; во французском когда после конечного сплитера стоит “nt”, оно читается, только если следующее слово начинается с гласной
5. Символ «hil» (hidden letter, скрытая буква), он же «☒» для Вуна
- 5.1. в тексте не отображается, но при обратной конвертации превращается в некоторый эмб: например, в русском пара из “s/z” и “t” не произносится перед “ʃ”; и если “☒5” – это “s”, “☒6” – “z”, “☒7” – “t”, то слова “расчёска, разносчик, извозчик, насчёт” предстают как «ra☒5☒7ʃoska, razno☒5☒7ʃik, izvo☒6☒7ʃik, na☒5☒7ʃot» (но в русском если между “s/z” и “t” есть пробел как во фразе “поднос чистый”, то пара фонем не исчезает)

Арабский и арабопишущие языки

Алиф-васла (إ) в начале слова читается, когда предыдущее слово оканчивается на согласную, и не читается, когда на гласную, поэтому после него ставится сплитер, т.е. он конвертируется в «a8>». По той же причине алиф в составе слитного артикля “ال” – в «a5>». А “ل” артикля может превратиться как в «L», так и в имаго одного из тринадцати других “солнечных” эмбов. В любом случае после такого имаго ставится проп, и весь артикль превращается одну из комбинаций «a5>L△», «a5>ʃ27△», «a5>ʃ5△», «a5>ʒ15△», «a5>s14△», «a5>s14ʒ△», «a5>z7△», «a5>z7ʒ△», «a5>t25△», «a5>d17△», «a5>t25ʒ△», «a5>d17ʒ△», «a5>r7△», «a5>n19△». Предлоги «посредством» (bi/ بـ), «для» (li/ لـ), союз «и» (va/ و), пишущиеся слитно с артиклем как “بالـ, للـ, و الـ” (второй даже без алифа), конвертируем в «biL△», «liL△», «vaL△».

АлифХанджарийя (надстрочный алиф “َ”) преобразуем в «a7_», а фатха-алифМаксура-алифХанджарийя “َى” – в «a6_8». АлифХамза-фатха (إ), вавХамза-фатха (ؤ), алифМаксураХамза-фатха (ئ) – в «’3a6, ’5a6, ’6a6»; алифХамза-кясра (إِ), вавХамза-кясра (ؤِ), алифМаксураХамза-кясра (ئِ) – в «’4i7, ’5i7, ’6i7»; алифХамза-дамма (إُ), вавХамза-дамма (ؤُ), алифМаксураХамза-дамма (ئُ) – в «’3u7, ’5u7, ’6u7». АлифМадда (آ) – сокращенное алифХамза-алиф (إِ), но читаемое уже без глоттальной смычки между двумя гласными – в «a9_». Три отмерших падежных окончания, три танвина «-en َ, -in ِ, -un ُ» и суккун ْ не произносятся, поэтому превращаются в «☒». Алиф и алифМаксура не читаются после танвина «-en», т.е. в сочетаниях « َا , ِا , ُا », поэтому в их составе конвертируются в «☒».

Лям-алиф (لا) отдельной моды не требует, т.к. в Юникоде, из которого и в который делаем прямую и обратную конвертации, на самом деле закодирован как последовательность двух букв (ل ا), замена которых на одну графему прописана в OpenType-функциях файла шрифтов .ttf. Такая же ситуация с лям-алифХамза для даммы (لاُ), для кясры (لاِ) и с лям-алифМадда (لاَ): на машинном уровне они – тоже две буквы (ل ا , ل ا , ل ا). Также в .ttf-файле прописан и выбор начертания в зависимости от расположения пробела до, после, с двух сторон буквы, ни с одной – начальное, конечное, изолированное, срединное.

Факультативные моды не указаны в цветных ячейках, но для полноты картины полные байтемы приведены в серых ячейках под ними. Жухло-зеленым отмечены эмбы и моды, которые исчезнут в следующих языках – персидской группы.

ب	ت	ة	ث	ج	ح	خ	د	ذ	ر	ز	س	ش	ص	ض	ط
b	t4	t19	ʃ3	d9ʒ0	h2ʒu	h4	d4	ʒ3	r5	z3	s3	ʃ6	s4ʒu	d5ʒu	t5ʒu

b3																		
ظ	ع	غ	ف	ق	ك	ل	ل	ل	ل	ل	ل	ل	ل	ل				
z4ۛ	h3ۛ	h2	f	k5ۛ	k3	ۛ10	ş27	ş5	ż15	s14	z7	t25	d17	r7	n19			
			f5															
م	ن	ه	و	و	ي	ى	َ	ِ	ُ	آ	ا							
m	n17	g	w	_5	ۛ	_9	a6	i7	u7	a9_	⊠8							
m4			g5	w6			ۛ3											
أ	َ	ِ	ِ	ِ	ِ	ِ	ى	ا	ي	َ	ِ	ِ	ى	ء	أ	إ	ؤ	ئ
a8	a7	⊠1	⊠2	⊠3	⊠4	⊠9	_3	_4	_2	_6	_7	_8	'2	'3	'4	'5	'6	

В иранском, дари фонемы «g» и «m» записываются двумя эмбами (ح, ه; م, ن), последний из четырех только в сочетаниях «mb», «s» – тремя (ث, ص, س), «z» – четырьмя (ذ, ز, ض, ظ). Для обозначения «k» используется неарабская буква, а артикля «aL» нет вообще, поэтому моды при «L» факультативны. Также появляются четыре новые согласные «p, tʃ, ʒ, g». В иранском «ق, غ» будем конвертировать в разные имаго, в «h, kʁ»; т.к. этот язык, в котором они не различаются, является упрощенным случаем дари и таджикского, в которых различаются.

Эмб “айн” (ع) в начале слова и между гласными не читается никогда; предписывается в конце слова и перед согласными произносить его как гортанную смычку – но никто не произносит (хотя и несколько удлиняют предшествующий гласный). Поскольку в настоящий момент в иранском, дари, таджикском нет противопоставления фонем по долготе, айн во всех случаях кодируется как «⊠» (а не как дуратор «_»). Также не произносится гортанная смычка, обозначаемая хамзой. Начальное и срединное начертание арабской алифМаксура (ى) неопределено, в иранском ее нет вообще, но конечное и изолированное совпадает с иранской «ɪ», поэтому будет говорить, что в иранском пишется алифМаксура (чтобы каждый раз не уточнять «йаАрабское», «йаИранское»).

Уникальные гласные дари [ɪ, ʊ] представим как ɪ, ʊ. Соотносимые фонемы иранского [o, ʊ], дари [ɔ, ʊ], таджикского [ɔ, ʊ] кодируем одинаково, и чтобы не применять фонеторы – как [o, ɔ], т.е. как «o, ɔ» (ближе к иранскому). Также иранские [e, æ], дари [ɛ, a], таджикские [e, a] – как [ɛ, a], т.е. как «e, a» (ближе к дари). Остальные две – [i, u], т.е. «i, u» – одинаковы во всех языках. Дифтонги иранские [ei, ou], дари и таджикские [ai, au] кодируем как в последних двух языках, т.е. как «ai, au». Специфический падеж изафет в иранском маркируется конечной [e], в дари – [a]. Однако хайеХавваз (ه), представляющий не только «g», но и изафет, изредка обозначает [a] и в иранском. Для унификации выберем [a] маркировкой изафета, т.о. хайеХавваз конвертируется в «g» и «a».

Алиф (ا) – не только безгласная подставка (уже не для хамзы, а для огласовок «e, o», т.е. “َ, ُ”), но и самостоятельная фонема «a» (когда без огласовок), поэтому может конвертироваться как в «⊠», так и в «a». АлифМадда (آ) – всегда «ɔ». АлифМаксура (ى) обозначают не только согласную «ɪ», но и «ɪ» в дифтонге «ai» (اي), «i» вне дифтонга и в дифтонге в начале слова (тогда алиф перед ней не читается), изредка «ɔ, a» в конце слова. Вав (و) – не только «v», но и «ɪʔ» в дифтонге «aiʔ» (او), «u»

вне дифтонга и в дифтонге в начале слова (опять тогда алиф не читается), изредка «o» в конце слова, и не произносится в диграфе «h⊠» (خو), если тот в начале слова. Насыщенным и темно зеленым выделены имаги и моды персидской группы, темным – которые исчезнут в пушту.

