

**Convenções usadas na conversão Grafema-Fonema (G2P)
com o alfabeto SAMPA**

Autores:

Fernando Perdigão, Sara Candeias, Arlindo Veiga

Instituto de Telecomunicações – Polo de Coimbra
Coimbra, abril de 2011

1. Introdução

Este documento descreve as convenções que foram usadas no sistema de conversão entre grafemas e fonemas quanto aos símbolos admitidos, tanto para fonemas quanto para grafemas.

O alinhamento entre grafemas e fonemas fica simplificado se forem usados símbolos com um só carácter para fonemas, como por exemplo /õ/ em vez de /o~/, pois evita a inserção do grafema correspondente ao til. De forma idêntica, os dígrafos correspondem a um só fonema, como é o exemplo de <ss> que dá origem ao fonema /s/. Evita-se neste caso o apagamento de um grafema.

Verificámos que a necessidade da inserção de fonemas é particularmente reduzida. De forma a evitar-se completamente a inserção de fonemas, foram definidos símbolos para mais que um fonema (Tabela 2).

As tabelas seguintes convencionam os símbolos usados do conjunto [ISO-8859-1](#)¹, bem como exemplos de utilização com o alfabeto [SAMPA](#)².

Apresentam-se também algumas notas sobre o problema do alinhamento entre grafemas e fonemas e a tabela de "custos" usada no programa de alinhamento.

Optámos por nos referir a fonemas e não a fones, uma vez que o estudo é fundamentalmente fonológico, sem esmiuçar a sua realização acústica em alofones.

¹ http://pt.wikipedia.org/wiki/ISO_8859-1

² <http://www.phon.ucl.ac.uk/home/sampa/>

Tabela 1 - Alfabeto SAMPA e extensões com um só caráter para alinhamento G2P.

nº	Símbolo SAMPA	Uni-caráter (SAMPAuc)	Grafemas possíveis	Exemplo/Observação
Vogais				
1	ɔ	ɔ ^a	a, e, â	cama, senha, câmara
2	a		a, á, à	pá, pala
3	@		e	de
4	e		e, ê	vê, dedo
5	E		e, é,	pé, pele
6	i		i, í, y, e	vi, aí, henry, reunião
7	o		o, ô, ou	oco, avô, louco
8	O		o, ó	pó, pote
9	u		u, ú, o	tu, tio, to, baú
10	ɔ~	ã	ã, an, em, am, êm,	branco
11	e~	ë	ên, en, em	penete, agência, empate
12	i~	ï	i, in, im, ím, ín, m	sim
13	o~	õ	õ, ôn, on, om	ponte, cônsul
14	u~	ü	u, ún, un, um	atum
Semivogais (não usadas)				
15	j			
16	w			
17	j~	ì		
18	w~	ù		
Plosivas/Oclusivas				
19	b		b	belo
20	d		d	dado
21	g		g, gu	gato
22	p		p	pato
23	t		t	tu
24	k		qu, c, k	casa, aquela,
Fricativas				
25	f		f	fé
26	s		s, cc, cç, ç, x, c+{eiéí}	sol, céu, trouxe, caça
27	S		ch, s, z, x	chave, xá, paz, pás
28	v		v	viu
29	z		z, s, x	casa, zé, exemplo
30	Z		j, g, s, z, x	já, gira, desviar, ex-bar
Líquidas				
31	l		l	lá
32	L		lh	velho
33	r		r	caro
34	R		r, rr	carro, rato
Consoantes nasais				
35	m		m	mão
36	n		n	nada
37	J		nh	senha

Tabela 2 - Vogais em posição tónica e símbolos para múltiplos fonemas

nº	Símbolo SAMPA	UniCarácter (SAMPAuc)	Grafemas possíveis	Exemplo/Observação
Vogais em posição tónica				
38	"ɔ	â		cama → /k"amɔ/ → /kâmɔ/
39	"a	á		casa → /k"azɔ/ → /kázɔ/
40	"e	ê		tema → /t"emɔ/ → /têmɔ/
41	"E	É		sete → /s"Et@/ → /sÉt@/
42	"i	í		tio → /tíu/
43	"o	ô		ovo → /ôvu/
44	"O	Ó		logo → /l"Ogu/ → /lÓgu/
45	"u	ú		uva → /úvɔ/
46	"ɔ~	Ã		campo → /kÃpu/
47	"e~	Ë		centro → /s"e~tru/ → /sËtru/
48	"i~	Ï		cinco → /s"i~ku/ → /sÏku/
49	"o~	Õ		conto → /k"o~tu/ → /kÕtu/
50	"u~	Û		assunto → /ɔs"u~tu/ → /asÛtu/
Símbolos para múltiplos fonemas				
51	ɔi	æ	e+{xj}	extrair → /æStrɔír/
52	"ɔi	Æ	e+{xj}	extra → /ÆStrɔ/;
53	ɔ~i~ɔ~	Ê	{tv}+êm	têm → /t"ɔ~i~ɔ~/ → /tÊi/
54	o~i~	ɶ	põem:	/po~i~ɔ~/ → /pɶãi/
55	ks	K	ficção, fax	
56	ai	Å	caem	/k"aiɔ~i~/ → /kÅãi/
57	Oi	®	constroem	/ko~StrOiɔ~i~/ → /kÖStr®ãi/
Fonemas Alternativos/Adicionais/Opcionais (para uso futuro em modelos acústicos)				
58	i e	y	esófago, eleitor	/iz"Ofɔgu/ vs /ez"Ofɔgu/
59	o O	Ô	olhar	/oL"ar/ vs /OL"ar/
60	e E	Ë	eólica, etnia, hebraica	/E"Olikɔ/ vs /e"Olikɔ/
61	u ø	μ	quente	
62	e ɔ	â		
63	l~	£		
64	ø	_	h, -, ' '	(fonema nulo)
65	-	0	-	(marcador de silêncio)
66	#	#	'sp', ' '	(sp; separador de palavras)

Nota: para as semivogais não são usados os símbolos /j/, /w/, /j~/ e /w~/ . Optámos por notá-las como /i/, /u/, /i~/ e /u~/, respetivamente.

Tabela 3 - Grafemas especiais para dígrafos

grafema	dígrafo/letras	fonema(s)	Exemplo com grafemas especiais
C	cc	s, ks	ficcional → fiCional
Ç	cç	s, ks	ficção → fiÇão
R	rr	R	carro → caRo
§	ss	S	massa → ma§a
L	lh	L	molho → moLo
J	nh	J	unha → uJa
S	ch	S	chave → Save
°	ou	o	dourada → d°rada
Ä	an, am	ã, Ã	canto → cÄto, campo → cÄpo
Ë	en, em	ë, Ë	sente → sËte, sempre → sËpre
Ï	in, im	ï, Ï	limbo → lÏbo
Ö	on, om	õ, Ö	conto → cÖto, dom → dÖ
Ü	un, um	ü, Ü	assunto, um
Â	ân, âm	Ã	pântano → pÂtano, lâmpada → lÂpada
Ê	ên, êm	Ë	ênfase → Êfase, êmbolo → Êbolo
Í	ín, ím	Ï	índio → Ídio, límpido → lÍpido
Ô	ôn, ôm	Ö	cônsul → cÔsul, cômputo → cÔputo
Ú	ún, úm	Ü	denúncia → denÚcia, cúmplice → cÚplice

2. Notas sobre alinhamento entre grafemas e fonemas

Princípios:

1. O alinhamento entre grafemas e fonemas permite fazer uma associação ótima entre símbolos, podendo existir apagamentos e inserções de grafemas e fonemas. Exemplos desta associação podem ser os seguintes:

<e_xtra> <põ_em> <tê__m> <fix_o> <campo> <milho> <guerra/
/6jStr6/ /põĩãĩ/ /tãĩãĩ/ /fiksu/ /kã_pu/ /miL_u/ /g_ER_6/

2. A definição dos símbolos para dígrafos e de símbolos para fonemas múltiplos tem como objetivo simplificar a associação entre grafemas e fonemas de tal forma que o alinhamento seja, tanto quanto possível, unívoco. Desta forma, é possível fazer com que não existam inserções de grafemas. Apenas em alguns casos existe inserção de fonemas (por exemplo no caso das formas nasaladas com n ou m: <an/→/ã_/).

Evita-se a inserção de fonemas (ou apagamento de grafemas) considerando 4 símbolos que agregam mais que um fonema:

K = /ks/; exemplo: <ficção> → /fiÇão/ → /fiKãü/; fixar → /fiKar/
æ = /6i/; exemplo: <e_xtra> → /6iStr6/ vs <extra/→/æStr6/
Æ = /6~i~6~/=/ãĩãĩ/; exemplo: <tê__m> → /tãĩãĩ/ vs <têm/→ /tÆi/
ɹ = /o~i~/=/õĩ/; exemplo: <põ_em> → /põĩãĩ/ vs <põem/→ /pɹãĩ/

3. Os custos de apagamento e de inserção de grafemas e/ou de fonemas, bem como os custos de associação entre grafemas e fonemas são previamente calculados. Dá-se preferência à inserção e ao apagamento de fonemas face a uma substituição incorreta. Por exemplo, nos alinhamentos seguintes, apenas o último é aceitável. Isto consegue-se afetando um custo elevado às associações r/@, e/S, n/d e S/@ no 1ª caso e s/e no 2º caso, baixando, em contrapartida, o custo da inserção dos grafemas r, n, s (ou apagamento dos fonemas correspondentes).

```
<correspondesse> <correspondesse> <correspondesse/  
/kuR@S_põde_s@_/ /kuR_@Spõ_d_es@/ /kuR_@Spõ_des_@/
```

Casos mais complicados são os seguintes:

```
<trans_eunte> <t_ens>  
/trã_ziü__t@/ /tãĩ_S/
```

pois o grafema <u> pode ser inserido (ex: guerra) assim como o <n> (em dígrafos nasais), e o <e> pode corresponder a /ĩ/ (ex: <ões> → /õĩS/). O fonema /ã/ também pode ser facilmente inserido, como se viu antes. A solução passa por baixar o custo das associações: <u>→ /ü/, <e> → /i/ e <n> → /ĩ/.

(ver alignG2P.m)

