

## Assignment 1: Grapheme to phoneme conversion

Due: Sept. 16

The goal of this assignment is to get familiar with some Perl syntax, while creating a program that does something like phonology. Your task is to write a script that converts Italian orthography into phonemic transcription. It should take an input file, read it, and perform whatever substitutions are necessary to produce a first approximation of a (broad) phonemic transcription. When you have completed your script, you can send it to me as an email attachment,

- If you have a particular interest in a language other than Italian, it is possible to take on a different language! (Consult with me first, though, so we can determine whether it looks an appropriate task.)

In addition, you should ponder the following questions, for discussion next week: (1) in some cases, two or more replacement rules must be ordered in a particular way. Is it possible to examine two rules and determine that their order might make a difference? How might you go about diagnosing that an ordering is necessary? Can you predict what sequences would be informative, so a learner could look out for them? (2) How would you write a program that could detect whether an input file was in orthography or transcription?

### 1 Background on Italian

Italian has (roughly) the following segments, which your transcription should make use of:

p, b	t, d		k, g		i	u
m	n	ɲ	ŋ		e	o
f, v	s, z	ʃ			ɛ	ɔ
	ʃ̂, d͡ẑ	t͡ʃ̂, d͡ʒ̂				a
	l, r	ʎ	j, w			

- We'll ignore [ɛ] and [ɔ] (they can't be predicted from spelling, and not all dialects have them anyway)
- We'll also ignore the phoneme [d͡ẑ], and lump it together with [ʃ̂]

When creating a transcription for computational modeling, it is generally convenient to represent segments using symbols that can be typed and transmitted easily. I recommend the following phonemic “alphabet” for Italian (but feel free to modify as you see fit):

p, b	t, d		k, g		i	u
m	n	N	G		e	o
f, v	s, z	ʃ				a
	Z	t͡ʃ̂, d͡ʒ̂				
	l, r	L	j, w			

### 2 Pronunciation rules

Here is the core of what makes Italian orthography not “purely phonetic” (and what you need to undo to get back to phonetic representation). A summary can be found at:

- <http://italian.about.com/library/nosearch/nblfare103a.htm>

Grapheme-to-phoneme correspondence in Italian is largely, but not completely predictable. The following graphemes correspond straightforwardly to the equivalent phonemes: *p, t, b, d, m, n, l, r*. Here are the rules for more complicated cases: (simplified somewhat)

- *gli* is pronounced [ʎ]<sup>1</sup> before another vowel (*miglio* = [meʎo], *famiglia* = [famiʎa]). It is pronounced [ʎi] at the end of a word, or before a consonant (*figli* = [fiʎi])
- *sci* is pronounced as [ʃ] before a vowel (*sciolta* = [ʃolta]), [ʃi] before a consonant or word-finally (*pesce* = [peʃi], *scioppo* = [ʃioppo])
- *sc* is pronounced as [ʃ] before [e], [i] (*pesce* = [peʃe]) (See previous for rules about [i] in particular). It is pronounced as [sk] otherwise (*pesca* = [peska])
- *ci* and *gi* are pronounced as [tʃ] and [dʒ] before another vowel (*cielo* = [tʃelo], *giusto* = [dʒusto]), and as [tʃi] and [dʒi], respectively, before a consonant or word-finally (*cibo* = [tʃibo], *undici* = [unditʃi])
- *c* is pronounced [k] before back vowels ([a, o, u]), [tʃ] before front vowels ([e, i]) (e.g., *cera* = [tʃera])
- *ch* is pronounced as [k] before front vowels (it doesn't occur before back vowels). Hence, *sch* is [sk] (*dischi* = [diski])
- *g* is pronounced [g] before back vowels, [dʒ] before front vowels (e.g., *gelo* = [dʒelo])
- *gh* is pronounced as [g] before front vowels (it doesn't occur before back vowels)
- *gn* is pronounced as [ɲ]
- *n* is pronounced as [ŋ] before velar stops (*funghi* = [fuŋgi], *banca* = [baŋka])
- *ng* is pronounced as [ŋg] before back vowels (*mango* = [maŋgo]), [ndʒ] before front vowels (*mangio* = [mandʒo]) (This is not a separate rule, but illustrates the interaction of two rules above)
- *qu* is pronounced as [kw]
- *h* is silent word-initially (pronounced as nothing)
- *s* is pronounced as [s] or [z]; for present purposes we'll just assume the following rough distribution:
  - [z] intervocalically, and before voiced stops (*rosa* = [roza], *frase* = [frazze], *sbarco* = [zbarko], *sgarbato* = [zgarbato])
  - [s] elsewhere (*sale* = [sale], *pasta* = [pasta], *pensa* = [pensa])
- *z* is pronounced somewhat unpredictable as [ts] or [dz]; we'll cheat and call them all Z
- Italian also has geminates, which you can represent as a sequences ([pp], [ZZ], etc.)
  - ☞ Watch out for sequences like *cci*, which should be CC, not kC (*braccio* = braCCo)
- Vowels are sometimes written with accents to indicate stress (*à*, *è*, etc.) There are also diphthongs (e.g., *ao* = [av]). In addition, [u], and [i] are sometimes realized as glides ([w, j]). Since the goal here is to learn how to implement a set of replacement rules, I recommend focusing on the consonants and ignoring these vowel issues.

A proposed alphabet of symbols to use for phonetic transcription can be found in `ItalianPhones.txt` on the course website. Here are some examples of transcriptions using this system:

<i>cara</i>	kara	<i>cera</i>	Cera	<i>cielo</i>	Celo
<i>gusto</i>	gusto	<i>giusto</i>	Justo	<i>mango</i>	maGgo
<i>mangio</i>	manJo	<i>anche</i>	aGke	<i>braccio</i>	braCCo
<i>sogno</i>	soNo	<i>questo</i>	kwesto	<i>hanno</i>	anno
<i>lascia</i>	laSa	<i>pesce</i>	peSe	<i>pesce</i>	peSi
<i>scioppo</i>	Sioppo	<i>schivo</i>	skivo	<i>pazzo</i>	paZZo
<i>funghi</i>	fuGgi	<i>ghigno</i>	giNo	<i>sbaglio</i>	zbaLo
<i>chiuso</i>	kiuzo	<i>festa</i>	festa	<i>esercizio</i>	ezerCiZio

<sup>1</sup>In point of fact, both [ʎ] and [ʃ] are always long intervocalically: [meʎo], [peʃi]. We will ignore this detail here.