

# Machine Learning of Phonetic Transcription Rules for Greek

Nikitas N. Karanikolas

Dept. of Informatics and Computer Engineering  
University of West Attica, Athens, Greece

[nnk@teiath.gr](mailto:nnk@teiath.gr)

<http://users.uniwa.gr/nnk/>

# What is phonology

- Phonology is a branch of linguistics concerned with the systematic organization of sounds in languages.
- It has traditionally focused largely on the study of the systems of **phonemes** in particular languages
- It may also cover any linguistic analysis either at a level **beneath the word** (syllable, etc.) or **at all levels of language** where sound is considered to be structured for conveying linguistic meaning.

# Phoneme

- A **phoneme** (/ˈfoʊni:m/) is one of the units of sound (or gesture in the case of sign languages) that distinguish one word from another in a particular language.
- For example, in most dialects of English, the sound patterns **/θʌm/** (*thumb*) and **/dʌm/** (*dumb*) are two separate words distinguished by the substitution of one phoneme, /θ/, for another phoneme, /d/.
- In many other languages these would be interpreted as exactly the same set of phonemes (i.e. /θ/ and /d/ would be considered the same).

# Serbian Cyrillic / Latin alphabet with IPA phonemes

А а	Б б	В в	Г г	Д д	Ђ ђ	Е е	Ж ж	З з	И и
<i>А а</i>	<i>Б б</i>	<i>В в</i>	<i>Г г</i>	<i>Д д</i>	<i>Ђ ђ</i>	<i>Е е</i>	<i>Ж ж</i>	<i>З з</i>	<i>И и</i>
а	бе	ве	ге	де	ђе	е	же	зе	и
a	b	v	g	d	đ	e	ž	z	i
[a]	[b]	[v]	[g]	[d]	[dz]	[e]	[ʒ]	[z]	[i]
Ј ј	К к	Л л	Љ љ	М м	Н н	Њ њ	О о	П п	Р р
<i>Ј ј</i>	<i>К к</i>	<i>Л л</i>	<i>Љ љ</i>	<i>М м</i>	<i>Н н</i>	<i>Њ њ</i>	<i>О о</i>	<i>П п</i>	<i>Р р</i>
је	ка	ле	ље	ме	не	ње	о	пе	ре
j	k	l	lj	m	n	nj	o	p	r
[j]	[k]	[l]	[ʎ]	[m]	[n]	[nj]	[ɔ]	[p]	[r]
С с	Т т	Ђ ђ	У у	Ф ф	Х х	Ц ц	Ч ч	Џ џ	Ш ш
<i>С с</i>	<i>Т т</i>	<i>Ђ ђ</i>	<i>У у</i>	<i>Ф ф</i>	<i>Х х</i>	<i>Ц ц</i>	<i>Ч ч</i>	<i>Џ џ</i>	<i>Ш ш</i>
се	те	ђе	у	фе	ха	це	че	џе	ша
s	t	ć	u	f	h	c	č	dž	š
[s]	[t]	[t͡ɕ]	[u]	[f]	[x/h]	[ts]	[t͡ʃ]	[d͡ʒ]	[ʃ]

# Serbian Language Phoneme Classes

А а    Б б    В в    Г г    Д д    Е е

*А а    Б б    В в    Г г    Д д    Е е*  
 а    бе    ве    ге    де    е  
 a    b    v    g    d    e  
 [a]    [b]    [v]    [g]    [d]    [e]

Ј ј    К к    М м    О о    П п    Р р

*Ј ј    К к    М м    О о    П п    Р р*  
 је    ка    ме    о    пе    ре  
 j    k    m    o    p    r  
 [j]    [k]    [m]    [ɔ]    [p]    [r]

Т т    У у    Ф ф    Х х

*Т т    У у    Ф ф    Х х*  
 те    у    фе    ха  
 t    u    f    h  
 [t]    [u]    [f]    [x/h]

Ђ ђ    Џ џ

*Ђ ђ    Џ џ*  
 ђе    џе  
 đ    dž  
 [dz]    [dʒ]

Ж ж    З з

*Ж ж    З з*  
 же    зе  
 ž    z  
 [ʒ]    [z]

Л л    Љ љ

*Л л    Љ љ*  
 ле    ље  
 l    lj  
 [l]    [ɭ]

Н н    Њ њ

*Н н    Њ њ*  
 не    ње  
 n    nj  
 [n]    [ɲ]

Ћ ћ    Џ џ    Ч ч

*Ћ ћ    Џ џ    Ч ч*  
 ће    џе    че  
 ć    c    č  
 [tɕ]    [tʂ]    [tʃ]

С с    Ш ш

*С с    Ш ш*  
 се    ша  
 s    š  
 [s]    [ʃ]

# Greek alphabet with IPA phonemes

Letter	Name	Sound	
		Ancient <sup>[7]</sup>	Modern <sup>[8]</sup>
Α α	alpha, άλφα	[a] [a:]	[a]
Β β	beta, βήτα	[b]	[v]
Γ γ	gamma, γάμμα	[g], [ŋ] <sup>[ex 1]</sup>	[ɣ] ~ [j], [ŋ] <sup>[ex 2]</sup> ~ [ɲ] <sup>[ex 3]</sup>
Δ δ	delta, δέλτα	[d]	[ð]
Ε ε	epsilon, έψιλον		[e]
Ζ ζ	zeta, ζήτα	[zd] <sup>A</sup>	[z]
Η η	eta, ήτα	[ɛ:]	[i]
Θ θ	theta, θήτα	[tʰ]	[θ]
Ι ι	iota, ιώτα	[i] [i:]	[ç], <sup>[ex 4]</sup> [j], <sup>[ex 5]</sup> [ɲ] <sup>[ex 6]</sup>
Κ κ	kappa, κάππα	[k]	[k] ~ [c]
Λ λ	lambda, λάμδα		[l]
Μ μ	mu, μυ		[m]

# Greek alphabet with IPA phonemes

Letter	Name	Sound	
		Ancient <sup>[7]</sup>	Modern <sup>[8]</sup>
Ν ν	nu, νυ		[n]
Ξ ξ	xi, ξι		[ks]
Ο ο	omicron, όμικρον		[o]
Π π	pi, πι		[p]
Ρ ρ	rho, ρώ		[r]
Σ σ/ς <sup>[note 1]</sup>	sigma, σίγμα	[s]	[s] ~ [z]
Τ τ	tau, ταυ		[t]
Υ υ	upsilon, ύψιλον	[y] [y:]	[i]
Φ φ	phi, φι	[p <sup>h</sup> ]	[f]
Χ χ	chi, χι	[k <sup>h</sup> ]	[x] ~ [ç]
Ψ ψ	psi, ψι		[ps]
Ω ω	omega, ωμέγα	[ɔ:]	[o]

# Greek Digraphs and letter combinations with their phonemes

Combination	Pronunciation	Devoiced pronunciation
⟨αυ⟩	[av]	[af]
⟨ευ⟩	[ev]	[ef]
⟨ηυ⟩	[iv]	[if]
⟨μπ⟩	[b]	
⟨ντ⟩	[d]	
⟨τζ⟩	[dz]	
⟨τσ⟩	[ts̥]	
γγ	[ŋg], [ŋʝ], [ŋɣ]	
γκ	[ŋg], [ŋʝ], [g], [ɣ]	
ει	[i]	
οι	[i]	
αι	[e]	



# Albanian Alphabet with IPA phonemes

A	B	C	Ç	D	Dh	E	Ë	F	G	Gj	H	I	J	K	L	LI	M
a	b	c	ç	d	dh	e	ë	f	g	gj	h	i	j	k	l	ll	m
ä	b	t͡s	t͡ʃ	d	ð	e, ε	ə, ʌ, ɜ	f	g	ɟ	h	i	j	k	l	ɫ	m

# Albanian Alphabet with IPA phonemes

N	Nj	O	P	Q	R	Rr	S	Sh	T	Th	U	V	X	Xh	Y	Z	Zh
n	nj	o	p	q	r	rr	s	sh	t	th	u	v	x	xh	y	z	zh
n	ɲ	o, ɔ	p	c	r	r	s	ʃ	t	θ	u	v	d͡ž	d͡ž̄	y	z	ʒ

# Some phonemes does not exist in every language

- Serbian has 23 phoneme classes
- Greek has 27 phoneme classes
- Albanian has ... phoneme classes

# Some phonemes does not exist in every language (GR vs ALB)

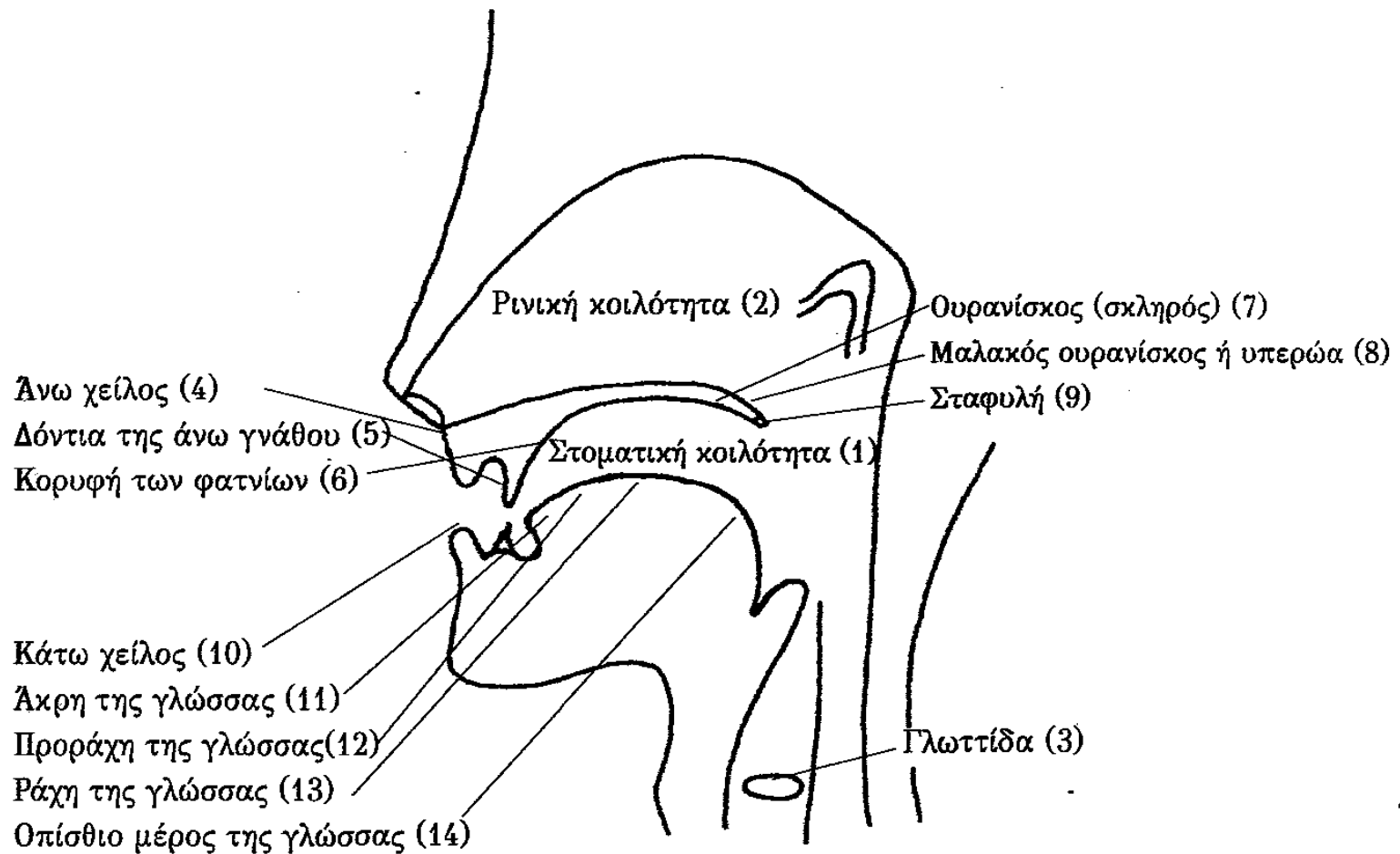
- GR cap – GR low – ALB – IPA
- ALB, GR, IPA
- Δ – δ – Dh – ð
- Dhjaku
- Διάκου
- /ðjaku/
- Γ – γ – Ø – γ
- Approximate with grafo
- γράφω
- /ɣrafo/
- ΟΥ – ου – U – u
- Pule
- Πούλε
- /puɫe/
- Ø – Ø – Υ – y
- ylber
- Approximate with Ιλμπερ
- /ylbeɪ/

Other GR Phonemes not existing in Albanian: Γ γ, Ξ ξ, Χ χ, Ψ ψ

# Some phonemes does not exist in every language (GR vs SRB)

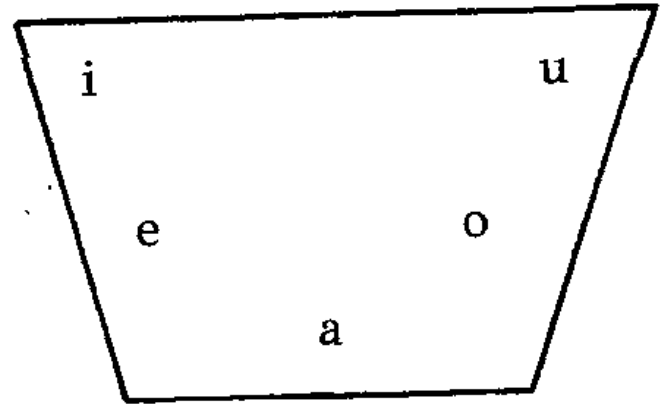
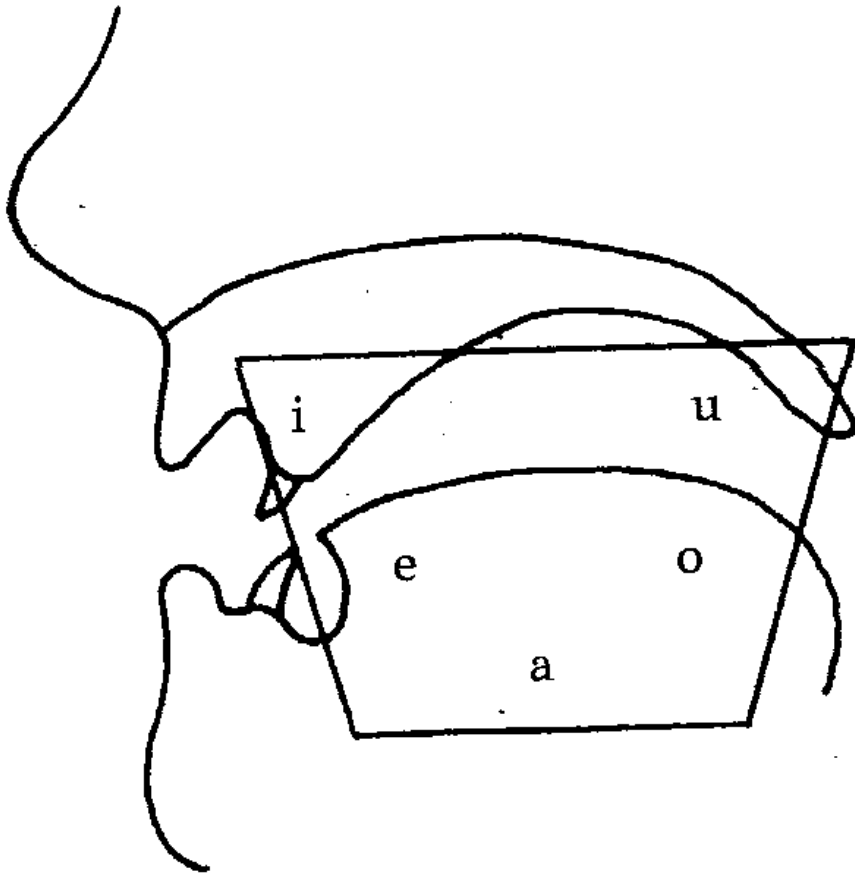
- Phonemes not existing in Serbian:
  - Γ γ
  - Δ δ
  - Θ θ
  - Ξ ξ
  - Ψ ψ
- Phonemes not existing in Greek:
  - J j

# Organs of humans relevant with phonology



picture from [3]

# Vowels of the Greek Language



picture from [4]

# Phonetic Alphabets

- Why we need phonetic alphabets?
  - To be able to represent graphically all the phonemes exists in every human language
  - To be able to represent with the same symbol a single phoneme that is represented with different letters in different languages
  - To solve the restrictions of the written alphabets
    - γέρος (/jeros/ )
    - γαρίδα (/ɣariða/ )
- How many Exists ?
  - 2, IPA and SAMPA



# IPA

- The **International Phonetic Alphabet (IPA)** is an alphabetic system of phonetic notation based primarily on the Latin alphabet.
- It was devised by the International Phonetic Association in the late 19th century as a standardized representation of the sounds (phonemes) of spoken language.
- The IPA is used by lexicographers, foreign language students and teachers, linguists, speech-language pathologists, singers, actors, constructed language creators and translators

# More about IPA

- <http://www.internationalphoneticalphabet.org/ipa-sounds/ipa-chart-with-sounds/>
- [https://www.internationalphoneticassociation.org/sites/default/files/IPA\\_Kiel\\_2015.pdf](https://www.internationalphoneticassociation.org/sites/default/files/IPA_Kiel_2015.pdf)
- [https://www.internationalphoneticassociation.org/sites/default/files/IPA2005\\_3000px.png](https://www.internationalphoneticassociation.org/sites/default/files/IPA2005_3000px.png)

# IPA excerpt

## CONSONANTS (PULMONIC)

	Bilabial	Labiodental	Dental	Alveolar	Postalveolar	Retroflex
Plosive	p b			t d		ʈ ɖ
Nasal	m	ɱ		n		ɳ
Trill	ʙ			r		
Tap or Flap		ɸ		ɾ		ɽ
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ
Lateral fricative				ɬ ɮ		
Approximant		ʋ		ɹ		ɻ
Lateral approximant				l		ɭ

Where symbols appear in pairs, the one to the right represents a voiced consonant.



# SAMPA

- The **Speech Assessment Methods Phonetic Alphabet (SAMPA)** is a computer-readable phonetic script using 7-bit printable ASCII characters, based on the International Phonetic Alphabet (IPA).
- It was originally developed in the late 1980s for six European languages by the EEC ESPRIT information technology research and development program.
- As many symbols as possible have been taken over from the IPA; where this is not possible, other signs that are available are used, e.g. [a] for schwa (IPA [ə]), for the vowel sound found in French *deux* (IPA [ø]), and for the vowel sound found in French *neuf* (IPA [œ]).

# More about SAMPA

- SAMPA at UCL
- <http://www.phon.ucl.ac.uk/home/sampa/index.html>
- SAMPA for Greek
- <http://www.phon.ucl.ac.uk/home/sampa/greek.htm>

# Greek Sampa 1/3

## Vowels

i, e, a, o, u

(see examples below)

## Consonants

### plosives

Symbol	Gloss	Transcription	Orthography
p	I drink	"pino	πίνω
b	I can	bo"ro	μπορώ
t	then	"tote	τότε
d	I dress	"dino	ντύνω
k	I do	"kano	κάνω
g	I throw down	gre"mizo	γκρεμίζω

### affricates

ts	pocket	"tsepi	τσέπη
dz	freeloader	dzaba"dzis	τζαμπαιζής

# Greek Samba 2/3

## fricatives

f	I depart	"fevGo	φεύγω
v	evening	"vraDi	βράδυ
T	I want	"Telo	θέλω
D	route	"Dromos	δρόμος
s	house	"spiti	σπίτι
z	heat	"zesti	ζέστη
x	time	"xronos	χρόνος
G	milk	"Gala	γάλα

## nasals

m	I speak	mi"lo	μιλώ
n	water	ne"ro	νερό
N	cucumber	an"guri	αγγούρι

([N] = allophone of /n/)



# Greek Sampa 3/3

## liquids

l	flower	lu"luDi	λουλούδι
r	clothes	"ruxa	ρούχα

## semivowel

j	I read	Dja"vazo	διαβάζω
---	--------	----------	---------

## (palatals)

c	and	ce	καί
	I sleep	ci"mame	κοιμάμαι
	I frighten	"scazo	σκιάζω
gj	angel	"angjelos	άγγελος
	bad luck	"gjinja	γκίνια
	screech-owl	"gjonis	γκιώνης
C	hand	"Ceri	χέρι
	winter	Ci"monas	χειμώνας
	snow	"Coni	χιόνι
jj	old man	"jjeros	γέρος
	circle	"jjiros	γύρος
	son	jjos	γιός

# Greek terms with IPA

- [https://en.wiktionary.org/wiki/Category:Greek\\_terms\\_with\\_IPA\\_pronunciation](https://en.wiktionary.org/wiki/Category:Greek_terms_with_IPA_pronunciation)
- At 30-4-2018 there were 3,481 items
- Items are organized per Greek letter: Αα Ββ Γγ Δδ Εε Ζζ Ηη Θθ Ιι Κκ Λλ Μμ Νν Ξξ Οο Ππ Ρρ Σσ Ττ Υυ Φφ Χχ Ψψ Ωω
- For each item there are:
  - Modern Greek: θάλασσα
  - Pronunciation with IPA: [ˈθalasa]
  - Usage examples
  - Declension
  - Related terms (some kind of thesaurus)

# Albanian terms with IPA

- [https://en.wiktionary.org/w/index.php?title=Category:Albanian\\_terms\\_with\\_IPA\\_pronunciation](https://en.wiktionary.org/w/index.php?title=Category:Albanian_terms_with_IPA_pronunciation)
- At 30-4-2018 there were 1,023 items
- Items are organized per Albanian letter: A B C Ç D Dh E Ē F G Gj H I J K L LI M N Nj O P Q R Rr S Sh T Th U V X Xh Y Z Zh
- For each item there are:
  - Albanian word: ylber
  - Pronunciation with IPA: /ylbɛɹ/
  - English Translation: rainbow
  - Declension

# Albanian Language Phonology - Vowels

- Albanian Language Phonology is simply
- Two tables of symbols are enough to transcribe from Albanian to IPA and vice versa

IPA	Description	Written as	English approximation
i	Close front unrounded vowel	i	seed
ɛ	Open-mid front unrounded vowel	e	bed
a	Open central unrounded vowel	a	father, Spanish casa
ə	Schwa	ë	about, the
ɔ	Open-mid back rounded vowel	o	law
y	Close front rounded vowel	y	French tu, German über
u	Close back rounded vowel	u	boot

# Albanian Language Phonology – Consonants 1/3

IPA	Description	Written as	English approximation
m	Bilabial nasal	m	man
n	Alveolar nasal	n	not
ɲ	Palatal nasal	nj	~onion
ŋ	Velar nasal	ng	bang
p	Voiceless bilabial plosive	p	spin
b	Voiced bilabial plosive	b	bat
t	Voiceless alveolar plosive	t	stand
d	Voiced alveolar plosive	d	debt
k	Voiceless velar plosive	k	scar
g	Voiced velar plosive	g	go

# Albanian Language Phonology – Consonants 2/3

$\widehat{ts}$	Voiceless alveolar affricate	c	hats
$\widehat{dz}$	Voiced alveolar affricate	x	goods
$\widehat{tʃ}$	Voiceless postalveolar affricate	ç	chin
$\widehat{dʒ}$	Voiced postalveolar affricate	xh	jet
$\widehat{cç}$	Voiceless palatal affricate	q	~china (RP)
$\widehat{ɟʝ}$	Voiced palatal affricate	gj	~gem (RP)
f	Voiceless labiodental fricative	f	far
v	Voiced labiodental fricative	v	van
θ	Voiceless dental fricative	th	thin
ð	Voiced dental fricative	dh	then

# Albanian Language Phonology – Consonants 3/3

s	Voiceless alveolar fricative	s	son
z	Voiced alveolar fricative	z	zip
ʃ	Voiceless postalveolar fricative	sh	show
ʒ	Voiced postalveolar fricative	zh	vision
h	Voiceless glottal fricative	h	hat
r	Alveolar trill	rr	Spanish perro
r	Alveolar tap	r	Spanish pero
l	Alveolar lateral approximant	l	lean
ɫ	Velarized alveolar lateral approximant	ll	ball
j	Palatal approximant	j	yes

# Albanian Language Phonology

## Not so simple - Not context free

- Before *q* and *gj*, the *n* is always pronounced /ɲ/ but it's not reflected in the orthography. That means:  
...nq... → /... ɲcç.../ and not /...ncç.../
- Next example (from [9])  
ngjashëm → [ɲʝaʃəm]  
does not follow the tables.  
According to tables it should be transcribed [ɲʝaʃəm]



# Do you pronounce the same way the letter on top in every following word

- **Q – q**  
Suflaqe  
Xaxiq  
Qirici  
Qepë  
Qumeshtur  
Qentër

- **Ç – ç**  
çift  
çakmak

- **C – c**  
cigare

- **GJ – gj**  
Gjizë  
Gjumë  
Gjasht  
Gjermane  
Gjrokaster

- **XH – xh**  
Xhina  
Xhiola  
Maxhelaku  
Xheni

- **X – x**  
Xulja  
Xeni

# Can you explain the differences

- **Gjrokaster** → something between Τζίροκαστερ and Γκίροκαστερ
- **Xhrokaster** → Ντζίροκαστερ or Τζίροκαστερ
- This is a practical understanding of the problem we want to resolve: Transliterate in another alphabet (not in IPA) understood by the user

# Greek Language Phonology

- It is more difficult.
- There are a lot of context sensitive rules
- γέρος → /jeros/  
while  
γαρίδα → /ɣariða/
- αγγαρεία → /aŋgaria/  
αγγελία → /aŋjelja/  
εγγόνι → /eŋgoni/  
έγγραφο → /eŋɣrafo/

# Orthographic Transcription

- **Orthographic transcription** is a transcription method that employs the standard spelling system of each target language.
- Examples of orthographic transcription are "Pushkin" and "Pouchkine", respectively the English and French orthographic transcriptions of the surname "Пу́шкин" in the name Алекса́ндр Пу́шкин (Alexander Pushkin).
- Thus, each target language (English and French) transcribes the surname according to its own orthography.

# Elaborating Orthographic Transcription

We can form training sets or corpus like:

Albanian word (1)	Orthographic Transcription Greek (2)	Phonetic Transcription (IPA) (3)	Translation Greek (4)
Bagazh	μπαγκάζ	[bagaʒ]	βαλίτσα
Dashuroj	ντασσουρόι	[daʃuroj]	αγαπάω
dymbëdhjetë	ντουμπαδιέτ(α)	[dymbə'ðjet(ə)]	Δώδεκα
Derr	ντερ	der	γουρούνι
Buzëqeshje	μπουζατσέσχ ιε	/buzə'ceʃje/, [bus'ceʃjɛ]	γελώ
Kolloface	κολοφάτσε	[kɔʎɔ'fatse]	λουκάνικο
Kuptoj	κουπτόι	[kup'tɔj]	καταλαβαί νω
Gjumë	τζιούμ(α)	/ɣumə/ /dʒumə/	κοιμάμαι

# Problems we can solve

For a Greek User

Usage	Input	Output
A	(4) Γουρούνι	(1), (2)* derr, ντερ
B	(1) Derr	(4), (2) γουρούνι, ντερ

We assume that a simple translation dictionary exist that offers:  
 $(4) \rightarrow 1$  and  $(1) \rightarrow (4)$ .

How A can be done:

1st step:  $(4) \rightarrow (1)$ ; from translation dictionary

2nd step:  $(1) \rightarrow (2)$ ; can already exist (in training set),

Otherwise  $(1) \rightarrow (3) \rightarrow (2)$ ; two steps are needed

So we need Transcribers for **Albanian to IPA (3)** and **IPA to Greek (2)**

\* numbers in parentheses are columns from previous table

# Problems we can solve

Similar usages can supported for the Albanese User

Usage	Input	Output
C	kolloface	λουκάνικο, lukaniko
D	λουκάνικο	kolloface, lukaniko

We need Transcribers for **Greek to IPA** and **IPA to Albanian**

# Machine Learning of Phonological Rules for Greek Transcription

- Is it possible to create some program that learns how to transcribe from Greek to IPA and from IPA to Greek?
- Yes, We think so,
- We need some resources, like “Wiktionary, Greek terms with IPA pronunciation” and some Algorithm.
- We start with Greek, because it has many rules (dependency of contexts) for phonological transcription.
- The algorithm should uncover (mechanically learn, mine) these rules.



# First step

- Consider words having the same number of Greek letters and IPA symbols in the transcription
- Example: πορτοκάλι → /portokali/
- With the one by one correspondence we can conclude:
  - π transcribe to p
  - ο transcribe to o
  - ρ transcribe to r
  - ...

# Fist Step after many words

$\alpha = \{$ a=323, k=4, p=2, n=1, s=1, e=1 $\},$	$\beta = \{$ v=67 $\},$	$\gamma = \{$ $\chi = 51,$ j=20, $\eta = 5,$ $\Upsilon = 1,$ g=1, j=1 $\},$
---	-------------------------------	--

# First step erroneous results

- Can have some erroneous results
- Example:  
εγκέφαλος → /eŋgefalos/  
conclude: γ transcribe to ŋ  
it is one of the 5 cases (out of 79) we have found  
where γ transcribe to ŋ
- Another example:  
ευθανασία → /efθanasia/  
conclude υ transcribe to f  
it is one of the 9 cases (out of 58) we have found  
where υ transcribe to f

# Step two

- We keep only transcriptions having a percentage above a predefined threshold
- For 20% transcription table is reduced to:

$\alpha = \{$ a=323, $\},$	$\beta = \{$ v=67 $\},$	$\gamma = \{$ $\gamma = 51,$ j=20, $\},$
----------------------------------	-------------------------------	---

# Third Step

- Consider words having one more Greek letter than the symbols in the IPA transcription
- Example: ουρανός → /uranos/
- With the one by one correspondence and **respecting** the results of the **second step**, the algorithm can conclude:
  - ou transcribes to u
  - ει transcribes to i
  - οι transcribes to i

# Third Step after many words

$OU = \{$ $u = 1,$ $\},$	$EL = \{$ $i = 1,$ $\},$	$OL = \{$ $i = 1,$ $\},$	$αL = \{$ $e = 1,$ $ε = 1,$ $\},$
--------------------------------	--------------------------------	--------------------------------	--

# Fourth Step

- Consider words having one less Greek letter than the symbols in the IPA transcription
- Example: ψάρια → /psaria/
- With the one by one correspondence and **respecting** the results of the **second step**, the algorithm can conclude:
  - ψ transcribes to ps
  - ξ transcribes to ks

# Fourt Step after many words

$\xi = \{$ ks = 10, $\},$	$\psi = \{$ ps = 15, $\},$
---------------------------------	----------------------------------



# Fifth Step

- Consider words having not resolved in previous steps
- GR letters can be +1 | +2 | -1 | -2 relatively to IPA symbols in transcription
- Example are:
  - ηλιόλουστος → /iɫɔlustɔs/
  - θηλιά → θiɫa
  - γιάννα → /jɔna/
  - γιαούρτι → /jaurti/
  - γράμμα → /ɣrama/
  - γέννηση → /jɛnisi/
  - μελισσοκομία → /melisokomia/
- With the one by one correspondence and **respecting** the results of **all previous steps**, the algorithm can conclude interesting valid transcriptions:

# Fifth Step after many words

$\lambda\iota=\{$ $\lambda=2,$ $\},$	$\mu\mu=\{$ $m=1,$ $\},$	$\nu\nu=\{$ $n=2,$ $\},$	$\sigma\sigma=\{$ $s=1,$ $\},$	$\gamma\iota=\{$ $j=2,$ $\},$
--	--------------------------------	--------------------------------	--------------------------------------	-------------------------------------

# Protected couples

- There are couples of letters that correspond phonetically to IPA couples of symbols.
- It is wrong to split the couple and consider each letter separately.
- These couples sometimes are also depended to their context (usually previous and next letter).
- These letters should be examined by the next step. For this reason, the **operator** of the Algorithm should have **declare these couples** in order the words having them not to be considered by previous steps.
- Such couples we call protected
- For the Greek language we suggest:  
γγ, γκ, τσ, τζ, μπ, ντ, αυ, ευ
- Also the same stand for some triangles:  
ντσ, ντζ

# Sixth Step

- The Algorithm considers only words not resolved by previous steps.
- It tries to find correspondences respecting all the previous findings and resolving the protected couples (and triangles).
- Given:
  - καλιαρντά → /kaldarda/
  - Μέτσοβο → /metsovo/
  - τσιμπούκι → /tsimbuci/
  - μπαμπάς → /babas/
  - ...
  - Εύβοια → /evia/
  - Ευγενία → /evjenia/
  - ευθανασία → /efthanasia/
  - αυγό → /avgo/
  - αυτοκίνητο → /aftocinito/

# Sixth Step after many words

$\nu\tau=\{$ d=1, $\},$	$\tau\sigma=\{$ ts=1, $\},$	$\mu\pi=\{$ mb=2, b=2, $\},$	$\tau\zeta=\{$ ts=1, dz=2, $\},$	$\gamma\gamma=\{$ $\eta g=2,$ $\eta j=1,$ $\eta\chi=1,$ $\},$	$\gamma\kappa=\{$ $\eta g=2,$ $\eta j=1,$ g=2, j=1, $\},$	$\epsilon\upsilon=\{$ e=1, ev=1, ef=1, $\},$
-------------------------------	-----------------------------------	---------------------------------------	---	---	--	--

# Contextual data

- The unification of all previous tables (step 2 to 6) are the set of transcription rules
- However, there are ambiguity cases. For example when the Greek word contains  $\kappa$  when it is transcribed to  $k$  and when it is transcribed to  $c$ ?
- The algorithm should also learn disambiguation of rule usage.
- To do this, the algorithm should keep the contexts.
- Next we see Greek couple  $\alpha\upsilon$  with contextual data:

$$\alpha\upsilon = \left\{ \begin{array}{l} av=1, -\alpha\upsilon\gamma \\ af=1, -\alpha\upsilon\tau \end{array} \right\},$$

- with more data can become:

$\alpha\upsilon=\{$

$av=8, -\alpha\upsilon\gamma, -\alpha\upsilon\lambda, -\alpha\upsilon\nu, \mu\alpha\upsilon\rho, \tau\alpha\upsilon\rho, \rho\alpha\upsilon\lambda, \beta\alpha\upsilon\alpha, \kappa\alpha\upsilon\lambda,$   
 $af=4, -\alpha\upsilon\tau, \epsilon\alpha\upsilon\tau, \nu\alpha\upsilon\pi,$

$\},$

- Can be generalized to:

$\alpha\upsilon=\{$

$av=8, \text{ when next letter is one of } \gamma, \lambda, \nu, \rho, \alpha,$   
 $af=4, \text{ when next letter is one of } \tau, \pi,$

$\},$

- And with more data can be generalized to:

$\alpha\upsilon=\{$

$av=x_1, \text{ when next letter is vowel or voiced consonant,}$   
 $af=x_2, \text{ when next letter is voiceless consonant,}$

$\},$

# Other rules we expect to mine

- $\mu\pi = \{$   
     $mb = r_1$ , after vowel  
     $b = r_2$ , in the beginning of word  
    or after consonant  
     $\}$ ,
- $\nu\iota = \{$   
     $\eta = t_1$ , when next letter is vowel  
    otherwise  $n_i$   
     $\}$ ,



# More problems

- A problem remains: How to orthographically transcribe when the target language does not have equivalent phoneme.
- Possible solution: use patterns from some other target language the user speaks (or english)
- Example from Greek to Albanian+English  
γέρος → **y**-eros (like y in **y**-ellow)

# References

- [1] Wikipedia, Phonology  
<https://en.wikipedia.org/wiki/Phonology>
- [2] Wikipedia, Phoneme  
<https://en.wikipedia.org/wiki/Phoneme>
- [3] Marina Nespou, ΦΩΝΟΛΟΓΙΑ, Chapter 2, Schema 2, ISBN: 960-378-083-9
- [4] Marina Nespou, ΦΩΝΟΛΟΓΙΑ, Chapter 2, Schema 3, ISBN: 960-378-083-9
- [5] Wikipedia, International Phonetic Alphabet  
[https://en.wikipedia.org/wiki/International\\_Phonetic\\_Alphabet](https://en.wikipedia.org/wiki/International_Phonetic_Alphabet)

# References

- [6] IPA chart with sounds  
<http://www.internationalphoneticalphabet.org/ipa-sounds/ipa-chart-with-sounds/>
- [7] IPA revised 2015  
[https://www.internationalphoneticassociation.org/sites/default/files/IPA\\_Kiel\\_2015.pdf](https://www.internationalphoneticassociation.org/sites/default/files/IPA_Kiel_2015.pdf)
- [8] IPA 2005  
[https://www.internationalphoneticassociation.org/sites/default/files/IPA2005\\_3000px.png](https://www.internationalphoneticassociation.org/sites/default/files/IPA2005_3000px.png)
- [9] Albanian Language Phonology  
[https://en.wikipedia.org/wiki/Albanian\\_language#Phonology](https://en.wikipedia.org/wiki/Albanian_language#Phonology)
- [10] Wikipedia SAMPA  
<https://en.wikipedia.org/wiki/SAMPA>

# References

- [11] SAMPA  
<http://www.phon.ucl.ac.uk/home/sampa/index.html>
- [12] SAMPA for Greek  
<http://www.phon.ucl.ac.uk/home/sampa/greek.htm>
- [13] Wiktionary, Greek terms with IPA pronunciation  
[https://en.wiktionary.org/wiki/Category:Greek\\_terms\\_with\\_IPA\\_pronunciation](https://en.wiktionary.org/wiki/Category:Greek_terms_with_IPA_pronunciation)
- [14] Wiktionary, Albanian terms with IPA pronunciation  
[https://en.wiktionary.org/wiki/Category:Albanian\\_terms\\_with\\_IPA\\_pronunciation](https://en.wiktionary.org/wiki/Category:Albanian_terms_with_IPA_pronunciation)
- [15] Wikipedia, Orthographic Transcription  
[https://en.wikipedia.org/wiki/Orthographic\\_transcription](https://en.wikipedia.org/wiki/Orthographic_transcription)
- [16 ] Wikipedia, Serbian Cyrillic alphabet  
[https://en.wikipedia.org/wiki/Serbian\\_Cyrillic\\_alphabet](https://en.wikipedia.org/wiki/Serbian_Cyrillic_alphabet)

# References

- [17] Omniglot, Serbian  
<https://www.omniglot.com/writing/serbian.htm>
- [18] Wikipedia, Greek alphabet  
[https://en.wikipedia.org/wiki/Greek\\_alphabet](https://en.wikipedia.org/wiki/Greek_alphabet)
- [19] Wikipedia, Albanian alphabet  
[https://en.wikipedia.org/wiki/Albanian\\_alphabet](https://en.wikipedia.org/wiki/Albanian_alphabet)

# Machine Learning of Phonetic Transcription Rules for Greek

Nikitas N. Karanikolas

Dept. of Informatics and Computer Engineering

University of West Attica, Athens, Greece

[nnk@teiath.gr](mailto:nnk@teiath.gr)

<http://users.uniwa.gr/nnk/>