

A CONTINUOUS WORD RECOGNITION SYSTEM

based on an

ACOUSTIC - PHONETIC APPROACH:

an application to a small vocabulary

in ITALIAN

M. G. DI BENEDETTO

P. MANDARINI

# description of the vocabulary

● Ten digits

+ 7 control words

digits

uno

due

tre

quattro

cinque

sei

sette

otto

nove

zero

control words

comando

invio

aiuto

annullo

precedente

successivo

conferma

'uno

d'ue

tr'e

kw'at:ro

ts'inkwe

s'eʎ

s'ɛt:e

'ɔt:ɔ

n'ɔve

dz'ɛro

Kom'ando

inv'io

a'juto

an:'ul:ɔ

pretsed'ɛ\*nte

sutʃ:es'ivo

komf'erɔ

# description of the vocabulary

phonemes present in the vocabulary

- vowels : [i, e, ε, a, o, ɔ, u]  
all vowels of the Italian vowel system
- stops : present [d] [p, t, k]  
miss [b, g]
- affricates : present [tʃ, dʒ]  
miss [tʃ, dʒ]
- fricatives : present [f, s] [v]  
miss [z]
- liquids : present [l, r]  
miss [ʎ]
- nasals : present [n, m]  
miss [ɲ]  
+ three allophones of [n] : [n̠, n̠̟, n̠̟̟]  
before dental [t, d] → [n̠̟]  
before velars [k, g] → [n̠̟̟]  
before bilabials [f, v] → [n̠̟̟̟]
- semiconsonants : present [j, w]  
miss NONE

## description of the vocabulary

- stops [d] word initial position  
[p] in cluster pr  
[t] geminate intervocalic position  
" in cluster tr  
[t] in cluster tr
- affricates [tʃ] word initial position  
intervocalic position  
geminate in intervocalic position
- fricatives [s] word initial position  
geminate in intervocalic position  
[f] after m  
[v] intervocalic position
- liquids [l] geminate intervocalic  
[r] in clusters tr pr  
intervocalic  
in group erm
- nasals [n] intervocalic, initial, geminate intervocalic  
[m] intervocalic and after [r]

## description of the vocabulary

- the words can be pronounced in sequences of one to three words -

### only one word

annullo

precedente

successivo

confermo

### one to three words

digits

comento

invio

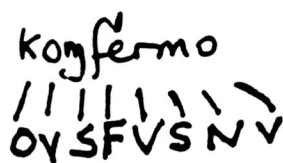
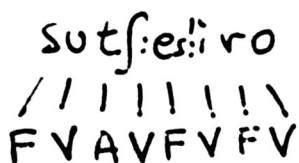
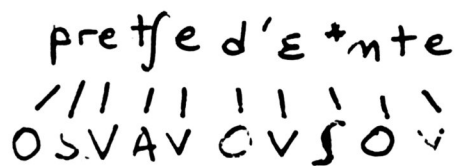
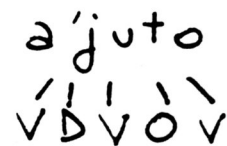
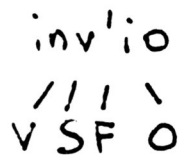
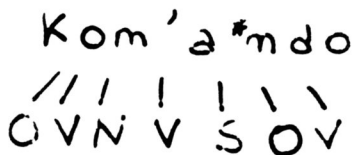
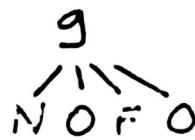
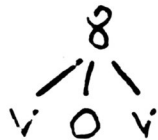
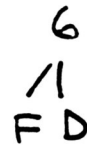
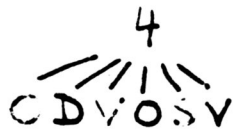
7 broad phonetic classes

1. Vowel-like	D	
2. Vowel	V	
3. Sonorant	S	
4. intervocalic nasal	N	
5. stop	O	
6. affricate	A	
7. fricative	F	

vowel-like	[e <sup>y</sup> ] [w, j]	
vowel	[i, e, ε, a, ɔ, o, u]	
sonorant	[r, l] [ŋ, m, *n]	
intervocalic nasal	[n, m]	
stop	[d, t, k, p]	
affricate	[dz, tʃ]	
fricative	[s, f] [v]	

description of the morphology

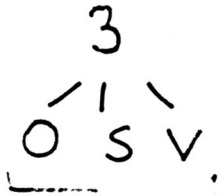
BROAD PHONEMIC CATEGORIES





description of the vocabulary

PREFIX RULE NOT SATISFIED



OSVAVOVS<sup>o</sup>V  
| | | | | | | |  
p r e t s e d ' e \* n t e

VNV<sup>o</sup>SV<sup>o</sup>  
| | | | | | | |  
a h : ' u l : o



pronounced

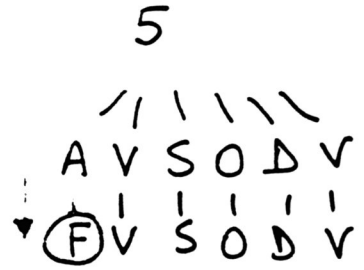
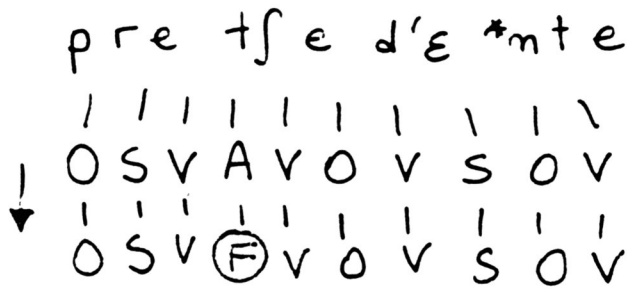
ONLY

in isolation -

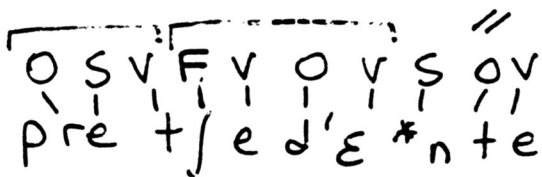
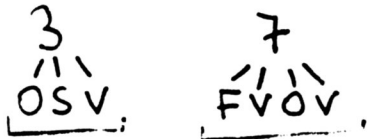
description of the vocabulary

differences in pronunciation

tʃ → ʃ



➔ prefix rule not satisfied



↖ pronounced only in isolation -

	p	b	f	v	t	d	ts	dz	s	z	k	g	tʃ	dʒ	ʃ	m	n	ɲ	ɺ	ɺ	r	i	e	ɛ	a	ɔ	o	u	j	w
vocalic	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
consonantal	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+
nasal	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
compact	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
diffuse																														
grave	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+
acute																														
tense																														
voiced	-	+	~	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+
continuant	-	-	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
strident			-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-

Description of the vocabulary 57

Distinctive Phonetic Features

- 75 sequences formed by one to three words belonging to the vocabulary
- Sequences selected in order to include several types of boundaries
  - Vowel - Vowel
  - Vowel - Consonant
  - consonant - vowel
- ten speakers : five female and five male repeated the 75 sequences twice
- recording in a sound-proof room using high-quality equipment
- speech materials sent on the telephone channel using artificial mouth and telephone-quality speech recorded - different telephone lines of the same telephone trunk.

- speech material digitalized : (toll quality)
  - low-pass filter 4.5 kHz
  - sampling frequency 10 kHz
  - 12 bits A/D

➔ each sentence in a file (stored on diskettes)

ILS compatible

- speech material of five speakers  
MANUALLY segmented using the following tools :
  - acoustic waveform
  - spectrum every 10 msec -

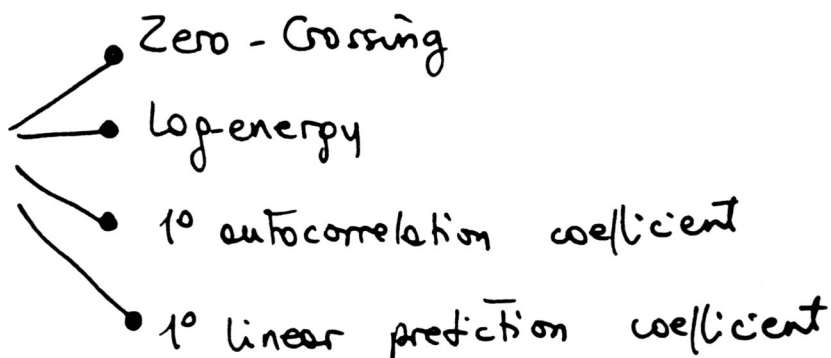
based on the criteria described in :

"A Bayesian - Adaptive decision method for the V/UV/S classification of segments of a speech signal"

IEEE - ASSP, vol 35, n° 4, 1987 -

by Bruno, Di Benedetto, Giulio, Mandorini -

- 4 parameters :



- definition of a-priori probabilities of the classes V-UV-S  $P(V)$   $P(UV)$   $P(S)$

- definition of transition probabilities between classes  $P(V|UV)$   $P(UV|S)$   $P(S|UV)$   
 $P(UV|V)$   $P(V|S)$   $P(S|V)$   
 $P(V|V)$   $P(UV|UV)$   $P(S|S)$

PRELIMINARY RESULTS

- training on two speakers - 10ms (100samples)

	V	UV	S
zero crossing	low ~9	high ~40	medium ~20
log-energy	high ~45	medium ~29	low ~19
1 <sup>o</sup> autoc. coef	high 0.95	low 0.25	medium 0.65
1 <sup>o</sup> LP coeff.	high 1	low -0.6	medium .1

- average values (mean vectors)  
variances (covariance matrices)  
computed for each class -

on the same speakers of the Training

in \ out	V	UV	S
V	0.89	0.07	0.04
UV	=	0.79	0.21
S	=	=	1

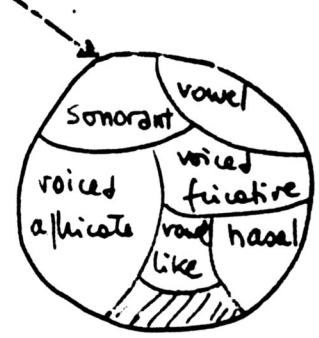
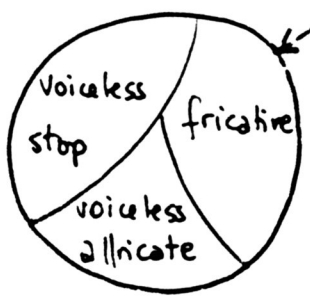
error types

in \ out	V	UV	S
V		[d][v] [dz][r]	0.04 [v]
UV			0.14 [k] in initial position 0.04 [f] first frame 0.03 [t] or [p] in tr/pr
S			

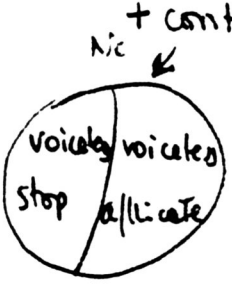


recognition strategy

is the segment  
+ voiced or? voiced



is the segment  
+ continuant



is the segment [+vocalic]  
or [+nasal] or [-vocalic] and  
[-consonantal]



is the segment  
+ grave



is the segment  
+ continuant



is the segment  
- consonantal



is the segment  
+ strident



is the segment  
- strident



is the segment  
+ nasal

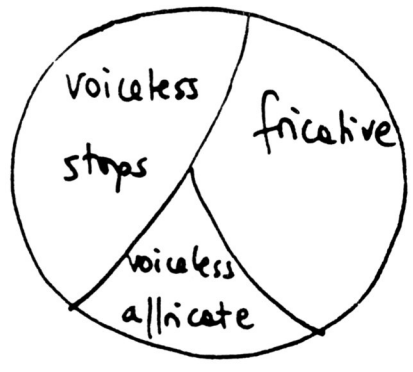


experimentation

- voiced segments

V-UV-S Classifier produces S and UV segments which

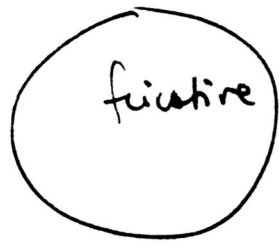
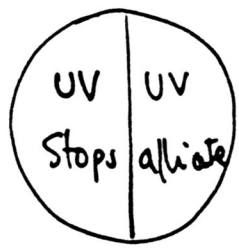
are - voiced.



+ continuant

no

yes

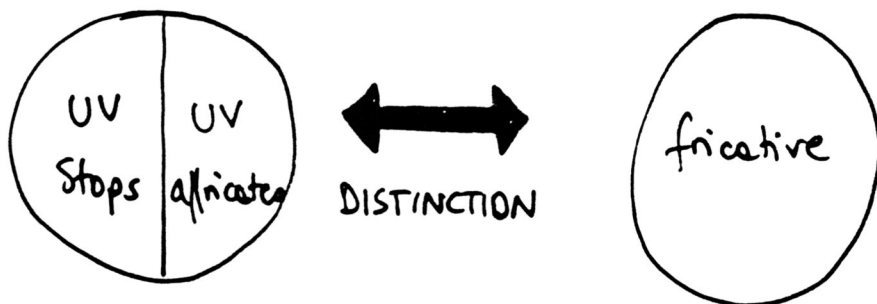


+ grave

no

yes





- sequence

V-S-UV-V  $\Rightarrow$  stop or affricate

V-UV-V  $\Rightarrow$  fricative

⚠ rule not valid for UV in initial sequence position -



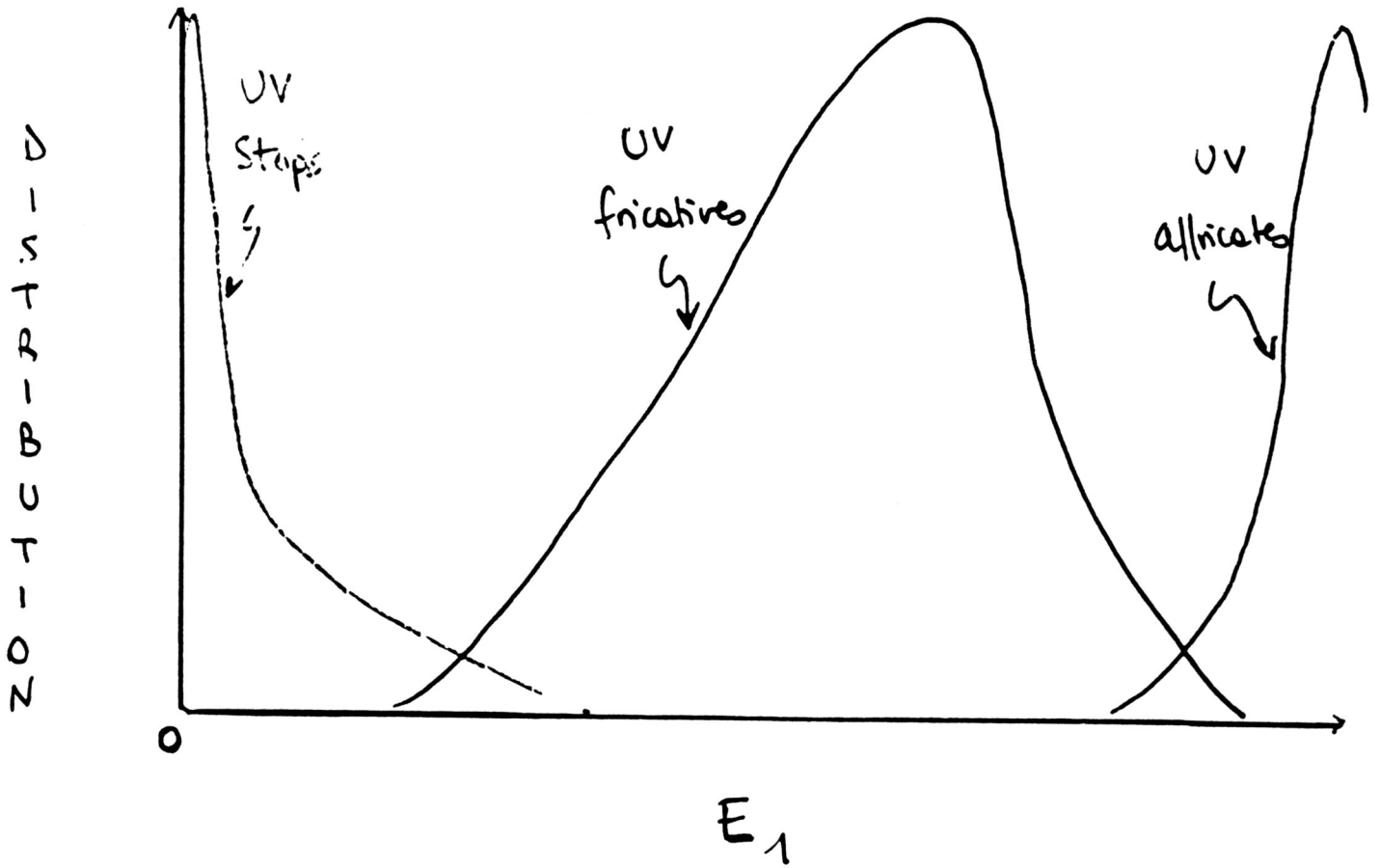
- energy parameter

$$E_1 = \frac{\text{Energy between } 2500 \div 3300 \text{ Hz}}{\text{Energy between } 300 \div 3300 \text{ Hz}}$$

⚡  
telephone channel bandwidth.

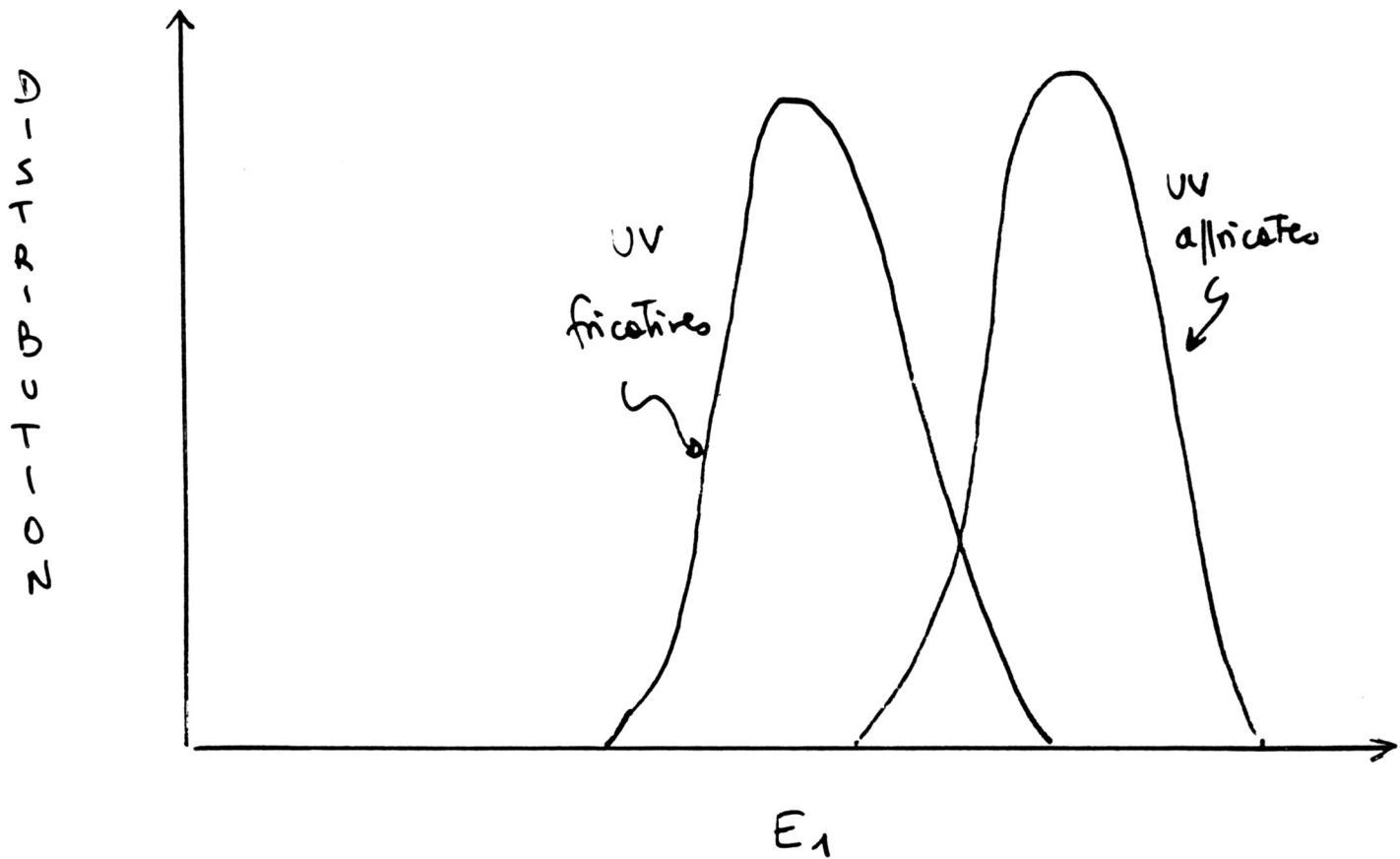
distinction + continuant (fricative)  
 - " (affricate + stop -

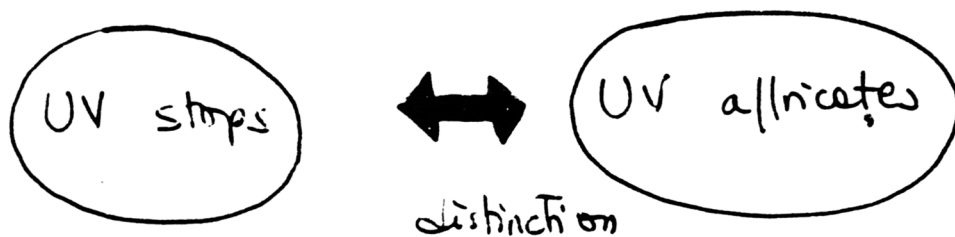
parameter :  $E_1 = \frac{\text{energy between } 2500 \div 3300}{\text{energy between } 300 \div 3300}$



distinction affricates - fricatives

$E_2$  = energy between 1600 and 3300 Hz





- $E_1$  as seen previously
- burst duration

UV stops average burst duration

25 msec

UV allicates average burst duration

60 msec

## Conclusion

- application of the acoustic-phonetic approach to a small vocabulary in Italian connected words
- data-base existing and digitalized
- data-base segmented
- preliminary results on:
  - distinction V-UV-S
  - distinction stops-allocates-fricatives
  - distinction stops-allocates

## Future Work

- extension to all classes
- " to other speakers
- " to telephone quality speech -