



31

—

-

, -

-

,

"

2013

4



31

.

,
Automatic)

.(Speech Recognition



3

(diphone)

(occurrences)

.

(triphone)

"It is crucial that all phonemes are represented in the speech corpus in **sufficiently high numbers.**" (Gibbon, Moore, and Winski (eds.), 1997)



?

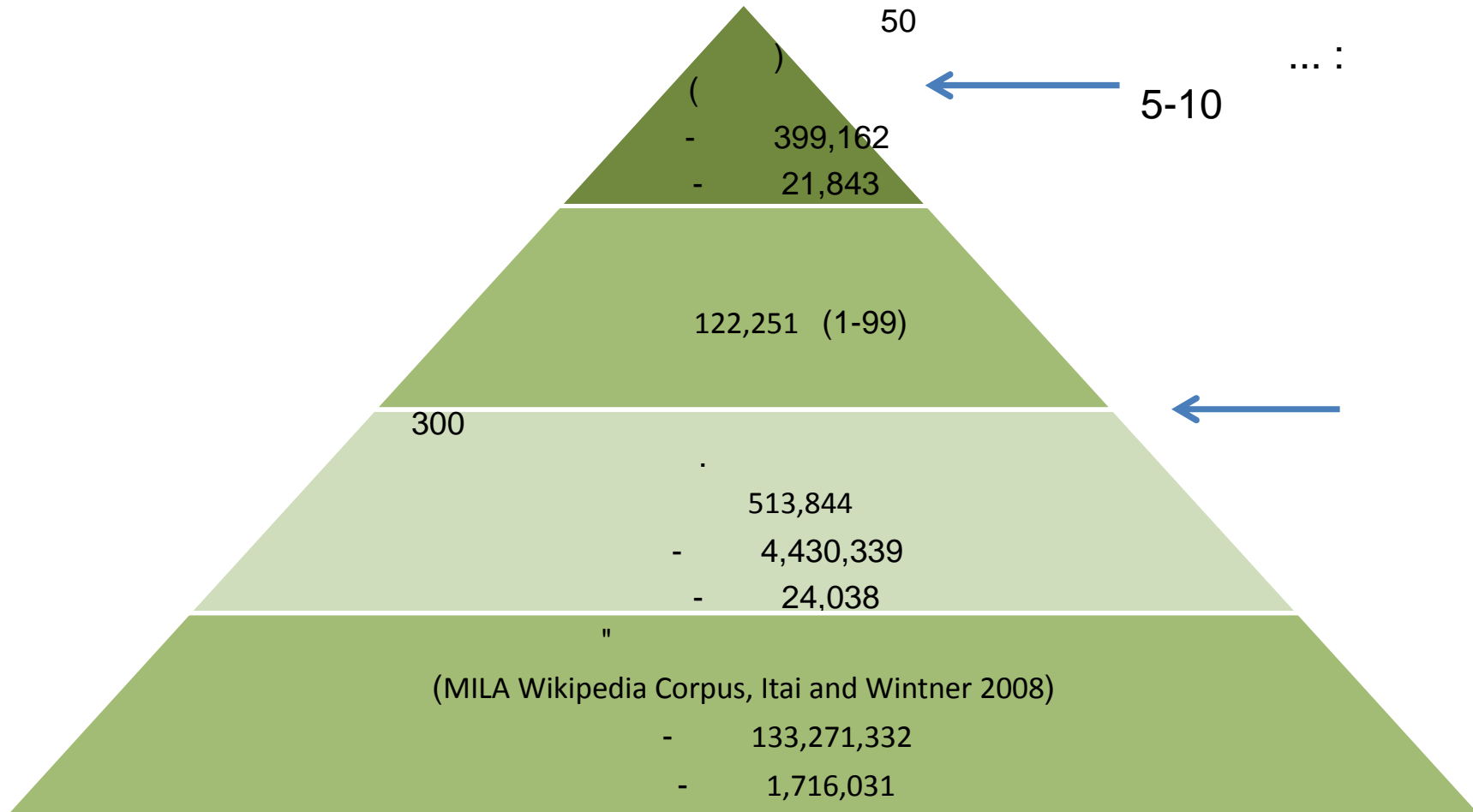
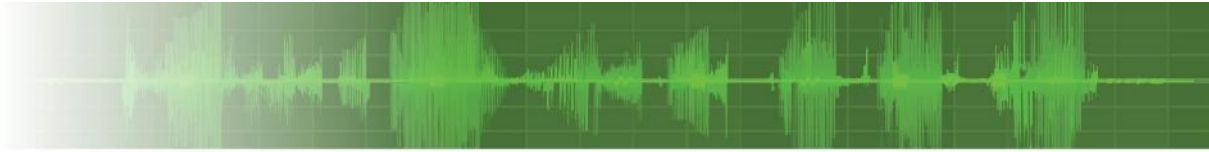
(phonetically balanced)

(phonetically rich)

· —
—
, ●
—
●



:"	"		.1
.		= (G2P)	.2
-	,TYPES-		.3
.	() TOKENS	
	TYPES -		.4
.		,	
.()		.5





(,)

z v k 1) ("	t h a (10,641)	p 1 " ' "	h a (55,250)	(279)	a (344,100)	50



:(,)

z v k (1)	t h a (11,176)	p (1)	h a (55,250)	(279)	a (344,100)	50 ()
z z o (1)	a n i (934)	z s (1)	h a / ma (2,000+)	(7)	A (20,553)	-) COSIH (2011) ()



(60) TIMIT

Sentence Type	Sentences	Speakers	Total	Sentences/ Speaker
Dialect (SA)	2	630	1260	2
Compact (SX)	450	7	3150	5
Diverse (SI)	1890	1	1890	3
Total	2342	630	6300	10



.[pda pde pdo pdu]

[p d i]

⋮ ●

—

100

⋮ ●

,

)

"

—

(

·

" ...

,

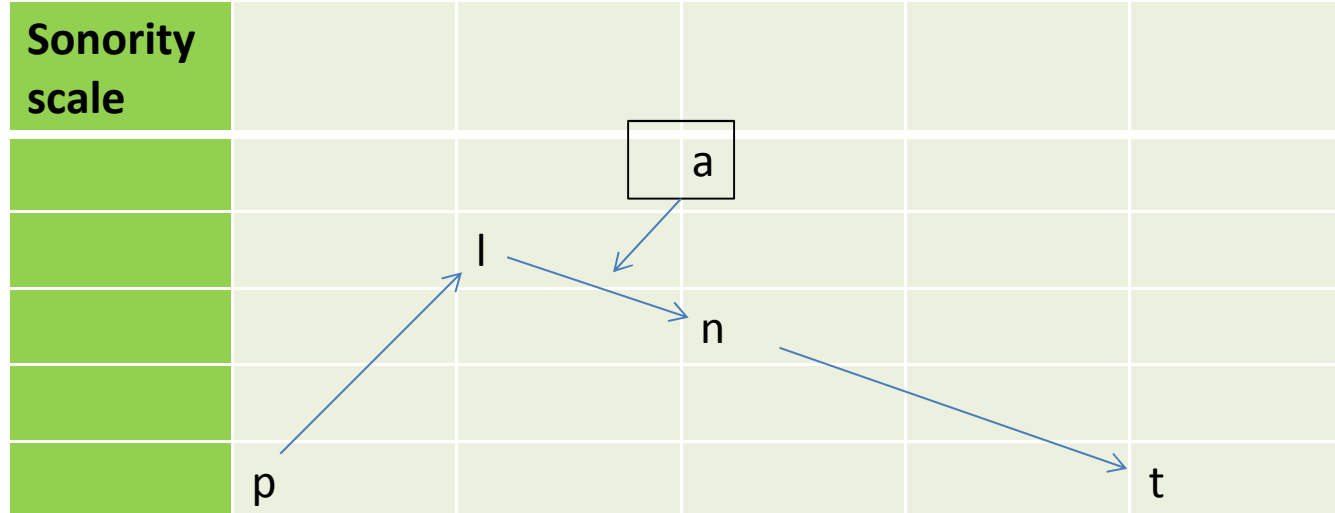
,

(12, "

)



(") (CCC)V(CCC) :
 (SSP – Sonority Sequencing Principle) –
 !' ' [stira] , ' [pli/ a] –





(CCVCC)

:

		stress	unstress	stress	unstress
		p		b	
onset (XY)	p	"petel	po"lin	bpVCC	bpVCC
onset (YX)		Rits"pa	ipa"ron	pbVCC	pbVCC
coda (XY)				CCVbp	CCVbp
coda (YX)				CCVpb	CCVpb
onset (XY)	b	pbVCC	pbVCC	"beReX	be"ReX
onset (YX)		bpVCC	bpVCC	niS"ba	"?aba
coda (XY)		CCVpb	CCVpb
coda (YX)		CCVbp	CCVbp
		s	
onset (XY)	f	fsVCC	fsVCC
onset (YX)		"sfinks	sfaRa"di sfi"na sfat ha"jam
coda (XY)		CCVfs	CCVfs
coda (YX)		CCVsf	CCVsf
onset (XY)	v	vsVCC	vsVCC
onset (YX)		svVCC	svi"va
coda (XY)		CCVvs	CCVvs
coda (YX)		CCVsv	CCVsv



:

	a	e	i	o	u
p d V	+	+	-	+	+
p g V	+	-	+	-	-
p n V	+	+	+	+	-
p z V	+	+	-	-	-
p S V	+	+	+	-	+
p X V	+	-	-	-	-
p ts V	+	-	+	-	+
b X V	+	+	+	-	-
t X V	+	+	+	-	+
V jt	+	+	-	+	+
d X V	+	-	+	-	-
D XT V	+	+	+	+	+
k z V	+	+	+	+	-
k X V	+	-	-	-	-
g X V	-	-	-	-	-



. :
 , [] , - •
 ' " [vue]
 " .
 ' " [wwo] - •
 : , [rfg] -
 " " " -
 () ,
 " " " , •



▪

50-

•

•

▪



todaRaba



Total No. of Word Types: 28
Total No. of Word Tokens: 117,536

1	20553	a
2	14669	e
3	8988	i
4	7270	m
5	7100	l
6	6897	o
7	6221	t
8	4643	n
9	4331	R
10	4165	h
11	3997	S
12	3666	k
13	3548	X
14	3252	j
15	2748	v
16	2607	u
17	2505	z
18	2474	b
19	2075	d
20	1724	s
21	1210	ts
22	989	g
23	960	f
24	772	p
25	78	w
26	49	tS
27	38	dZ
28	7	Z



Total No. of Word Types: 28

Total No. of Word Tokens: 21,843,070

1	3898936	a	
2	2458405	e	
3	1625533	i	:
4	1374107	m	:
5	1239149	l	
6	1192523	t	
7	1122641	o	
8	976604	R	
9	877492	h	
10	817748	n) 2010
11	744567	X	
12	735492	S	
13	664370	k	(300
14	626049	j	
15	584213	u	
16	523874	b	
17	498836	v	
18	462244	d	
19	337849	s	
20	236129	z	
21	234633	ts	
22	230061	f	
23	175947	g	
24	150020	p	
25	31042	dZ	
26	21571	w	
27	2762	tS	
28	273	Z	