



ELSEVIER

Available online at www.sciencedirect.com

SCIENCE @ DIRECT®

Speech Communication xxx (2006) xxx–xxx

SPEECH
COMMUNICATIONwww.elsevier.com/locate/specom

Effect of voice quality on frequency-warped modeling of vowel spectra [☆]

Pushkar Patwardhan *, Preeti Rao

Department of Electrical Engineering, ACRE, Indian Institute of Technology Bombay, Powai, Mumbai 400 076, India

Received 19 July 2005; received in revised form 7 December 2005; accepted 6 January 2006

Abstract

The perceptual accuracy of an all-pole representation of the spectral envelope of voiced sounds may be enhanced by the use of frequency-scale warping prior to LP modeling. For the representation of harmonic amplitudes in the sinusoidal coding of voiced sounds, the effectiveness of frequency warping was shown to depend on the underlying signal spectral shape as determined by phoneme quality. In this paper, the previous work is extended to the other important dimension of spectral shape variation, namely voice quality. The influence of voice quality attributes on the perceived modeling error in frequency-warped LP modeling of the spectral envelope is investigated through subjective and objective measures applied to synthetic and natural steady sounds. Experimental results are presented that demonstrate the feasibility and advantage of adapting the warping function to the signal spectral envelope in the context of a sinusoidal speech coding scheme.

Keywords: Voice quality; Spectral envelope modeling; Frequency warping; All-pole modeling; Partial loudness

1. Introduction

A successful low bit rate speech coding algorithm requires a good model for the speech signal together with an effective parameter quantization algorithm. A popular model for low bit rate speech coding has been the sinusoidal model, an important example of which is the multiband excitation (MBE) model (Griffin and Lim, 1988). In the case of voiced speech, the parameters of the model are the funda-

mental frequency (or pitch), and the amplitudes and phases of the harmonics. The harmonic amplitudes represent the product of the source excitation and vocal tract spectra. At low bit rates, estimated phases are usually dispensed with, and the accurate representation of the pitch and harmonic amplitudes becomes critical to the perceptual quality of decoded speech. Steady vowel sounds are particularly sensitive to harmonic amplitude representation errors.

The quantization of harmonic amplitudes is most demanding on the bit allotment in sinusoidal coding, and methods for efficient quantization have been an important topic of research. A widely used method for the quantization of the amplitudes of the harmonics is based on the modeling of a spectral

[☆] A portion of this work was presented at the International Conference on Spoken Language Technology-04, New Delhi.

* Corresponding author. Tel.: +91 22 25346665; Mob: +91 9819073984; fax: +91 22 25723806.

E-mail address: pushkar@ee.iitb.ac.in (P. Patwardhan).

45 envelope fitted to the harmonic peaks (MacAulay
46 and Quatieri, 1995). The spectral amplitudes are
47 then reconstructed from samples of the modeled
48 spectral envelope at harmonic frequencies. Represent-
49 ing the spectral envelope by the coefficients of
50 an all-pole filter enables the use of one of the many
51 efficient quantization methods available in the
52 speech coding literature. The order of the all-pole
53 model has a significant effect on the accuracy of
54 the modeled spectral amplitudes. While the all-pole
55 representation of the spectral envelope is expected
56 to capture local resonances accurately, capturing
57 additional features such as overall spectral tilt and
58 spectral zeros due to nasality typically lead to an
59 increase in the number of poles required for an ade-
60 quate approximation. Further, for similar spectral
61 envelope, low pitched sounds require higher LP
62 model order for similar perceived quality levels
63 (Champion et al., 1994; Rao and Patwardhan,
64 2005).

65 In the interest of achieving low bit rates, how-
66 ever, it is necessary to keep model order as low as
67 possible. Frequency-scale warping before all-pole
68 modeling of the spectral envelope is a widely used
69 method to improve the perceptual accuracy of mod-
70 eling for a given model order. Frequency-scale
71 warping leads to a more accurate representation of
72 the low frequency spectrum at the cost of increased
73 errors in the high frequency region. Although per-
74 ceptual scales such as the Bark-scale and its variants
75 have been widely used in LP modeling of speech
76 spectra, our recent work (Rao and Patwardhan,
77 2005) on synthetically generated steady vowel
78 sounds (using fixed excitation source parameters)
79 indicated that the performance of frequency warp-
80 ing depended to a great extent on the nature of
81 the underlying sound spectrum. It was observed that
82 front vowels such as [e] and [i] in fact were better
83 modeled without Bark-scale warping. This was
84 explained by the low frequency first formant struc-
85 ture of these vowels which fails to mask high fre-
86 quency distortion adequately. In a subjective
87 experiment using natural speech, it was found that
88 the comparative behaviour of modeling error under
89 different warping conditions in the case of the non-
90 front vowels was inconsistent across speakers indi-
91 cating a further dependence on speaker voice qual-
92 ity as reflected in the overall spectral envelope.

93 In the present study we attempt to extend our
94 previous work by exploring aspects of voice quality
95 that could influence the spectrum modeling error.
96 The influence of voice quality on perceived model-

ing error in frequency-warped all-pole modeling of
the spectral envelope is studied via the framework
of MBE model analysis–synthesis.

While the overall spectral envelope for voiced
speech is determined by the glottal source wave-
form, vocal tract transfer function and lip radiation,
it is the glottal source waveform that most directly
affects the relative strengths of the low and high fre-
quency harmonics. The scope of the present study is
restricted to variations in laryngeal voice qualities
for they are closely linked to gross differences in
spectral envelope. An investigation based on subjec-
tive and objective experimental evaluation is pre-
sented. Finally, the applicability of the results to
improving speech quality in a low bit rate sinusoidal
coder is discussed.

2. Frequency-warped LP modeling of discrete spectra

A discrete spectrum, characterized by a funda-
mental frequency and the amplitudes of the compo-
nents at harmonic frequencies, can be represented
by the coefficients of an all-pole model. A smooth
spectral envelope is first derived to fit the harmonic
amplitudes using a suitable interpolation method
such as linear interpolation of log amplitudes (Her-
mankys et al., 1985; Rao and Patwardhan, 2005).
The power spectrum obtained from the interpolated
envelope is used to compute the autocorrelation
function via the inverse DFT. The Levinson–
Durbin algorithm is applied to obtain the LP coeffi-
cients. The spectral amplitudes can be recovered by
sampling the reconstructed spectral envelope repre-
sented by the all-pole coefficients at the harmonic
frequency locations. Frequency-scale warping may
be incorporated in the spectral envelope modeling
by mapping the input harmonic frequencies to cor-
responding warped frequency locations by means of
a warping function based on a chosen perceptual
scale. The log-linear interpolation of the discrete
spectral amplitudes (now non-uniformly spaced in
frequency) is carried out to obtain a densely sam-
pled spectral envelope as detailed in Section 3 of
(Rao and Patwardhan, 2005). Uniformly spaced
samples at 20 Hz interval are found to be adequate
for the LP modeling of narrowband speech spectra
in the sinusoidal coding context.

In this work, we use the framework of MBE nar-
row band coding of speech to evaluate frequency-
warped LP modeling of voiced sounds. The MBE
analysis–synthesis model offers a convenient frame-

work for the evaluation of spectral envelope modeling (Molyneux et al., 1998; Rao and Patwardhan, 2005) although the results are applicable more generally to other vocoders. Voiced regions are modeled by harmonics of a fundamental frequency in the MBE vocoder, and unvoiced bands by spectrally shaped random noise. Thus parameters of the MBE speech model consist of the fundamental frequency, band voicing decisions, and the harmonic amplitudes (Griffin and Lim, 1988).

MBE analysis involves the use of high resolution DFT in an analysis-by-synthesis loop for the accurate determination of pitch, voicing and spectral amplitudes for each input frame of speech (typically 20 ms duration). In the case of fully voiced speech, synthesis is achieved by summing of sinusoids each corresponding to one harmonic. Adjacent frames are combined using either overlap-add or the interpolation of phase depending on the extent of pitch variation (Griffin and Lim, 1988). With unquantised parameters, synthesized speech of very high quality is obtained, particularly in voiced regions. Spectral envelope-interpolated LP modeling is utilized to achieve the low bit rate coding of the harmonic spectral amplitudes in an MBE model based speech coder. In the present work, we investigate some aspects of the performance of frequency-scale warping in obtaining improved quality at low LP model order. Based on the considerations discussed in (Rao and Patwardhan, 2005) the LP model order selected for the speech quality evaluation experiments is 10. The test and evaluation framework is

depicted Fig. 1 which shows the generation of reference and test signals used in the quality comparisons.

3. Voice quality and its spectral correlates

Voice quality refers to the auditory impression a listener gets upon hearing the speech of another talker. Voice quality is determined by the articulators of the vocal tract as well as the characteristics of the vocal folds. For speech, the vocal folds have a predominant importance. We refer to this aspect of voice quality that results from differences in vocal cord vibratory patterns, or laryngeal voice quality (Childers and Lee, 1991). Periodic glottal excitation is a characteristic of “modal” voices. The perceptual and spectral variations within modal voices can be attributed to the differences in the timing parameters of the glottal pulse. Since we are concerned with the spectral representation of steady periodic sounds, we do not consider those voice types that are characterized by voicing aperiodicities such as aspiration noise and pitch/amplitude jitter.

3.1. Modal voice production parameters

The speech signal is generated by excitation of vocal tract caused by the modulated flow of air generated by the lungs. The volume of air passing through the vocal cords, also known as glottal volume flow is the excitation signal which is periodic in the case of voiced speech. The glottal pulse shape

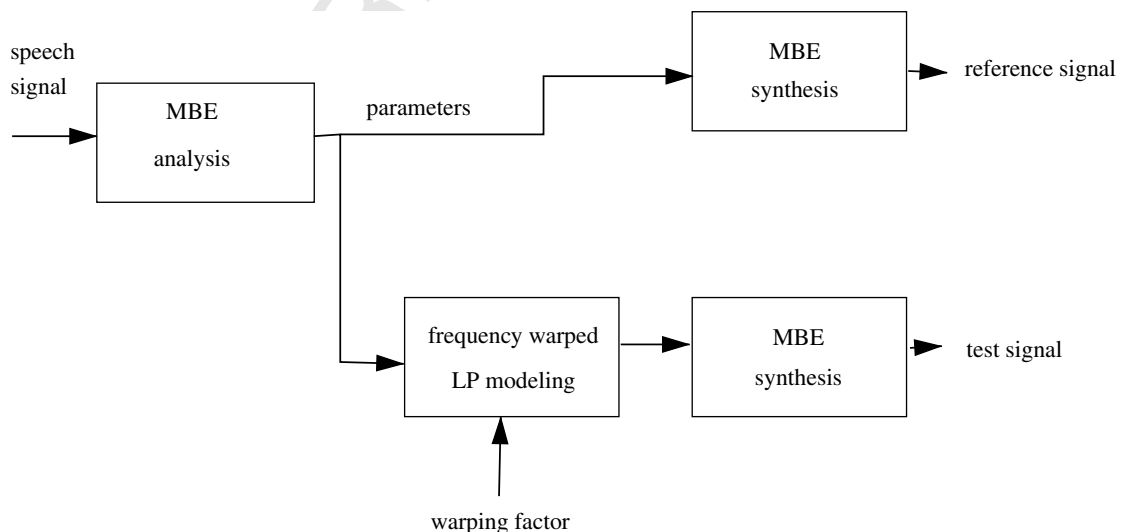


Fig. 1. Generation of MBE modeled test and reference signals.

207 has been extensively studied and modeled for differ- 228
 208 ent voice qualities. A typical glottal pulse period 229
 209 (glottal volume flow) and its derivative (Childers, 230
 210 2000; Childers and Lee, 1991) are shown in 231
 211 Fig. 2(a) and (b).

212 The parameter t_p denotes the instant of the max- 232
 213 imum in the glottal flow waveform. The parameter 233
 214 t_e is the instant of the maximum negative differenti- 234
 215 ated glottal flow. The t_a is the time constant of the 235
 216 exponential curve of the second segment of the LF 236
 217 model, and parameter t_c is the instant at which com- 237
 218 plete glottal closure is reached. The spectral shape 238
 219 of the glottal pulse derivative is dependent on the 239
 220 values of the timing parameters. Modal voices, com- 240
 221 mon among male speakers, are characterized by the 241
 222 nearly complete closure of the vocal folds. Fig. 2(c) 242
 223 shows glottal pulse derivative waveforms corre- 243
 224 sponding to three different sets of timing parameters 244
 225 and Fig. 2(d) shows the corresponding spectra. An 245
 226 abrupt closure of the vocal folds (low t_a) creates 246
 227 strong high frequency harmonics while a relatively

gradual closure of vocal folds results in spectra with 228
 strong low frequency harmonics and weak high fre- 229
 quency harmonics. The perceptual correlate of the 230
 relative strengths of the low and high frequency har- 231
 monics in the glottal source spectrum is the bright- 232
 ness of the voice. “Dark” voices have relatively 233
 weak high harmonics while “bright” voices are 234
 characterized by flatter spectra. Apart from the 235
 abruptness in the closure of the vocal folds, the sym- 236
 metry of the open phase of the glottal pulse as cap- 237
 tured by the “pulse skew” also affects the spectrum. 238
 It influences the relative amplitudes of the low fre- 239
 quency harmonics (especially the first and second 240
 harmonics) (Doval and d’Alessandro, 1997). 241

3.2. Spectral correlates 242

Spectral shape can provide useful cues to relevant 243
 aspects of voice quality. Some of the important cues 244
 that reflect the characteristics of glottal signal are as 245
 follows: 246

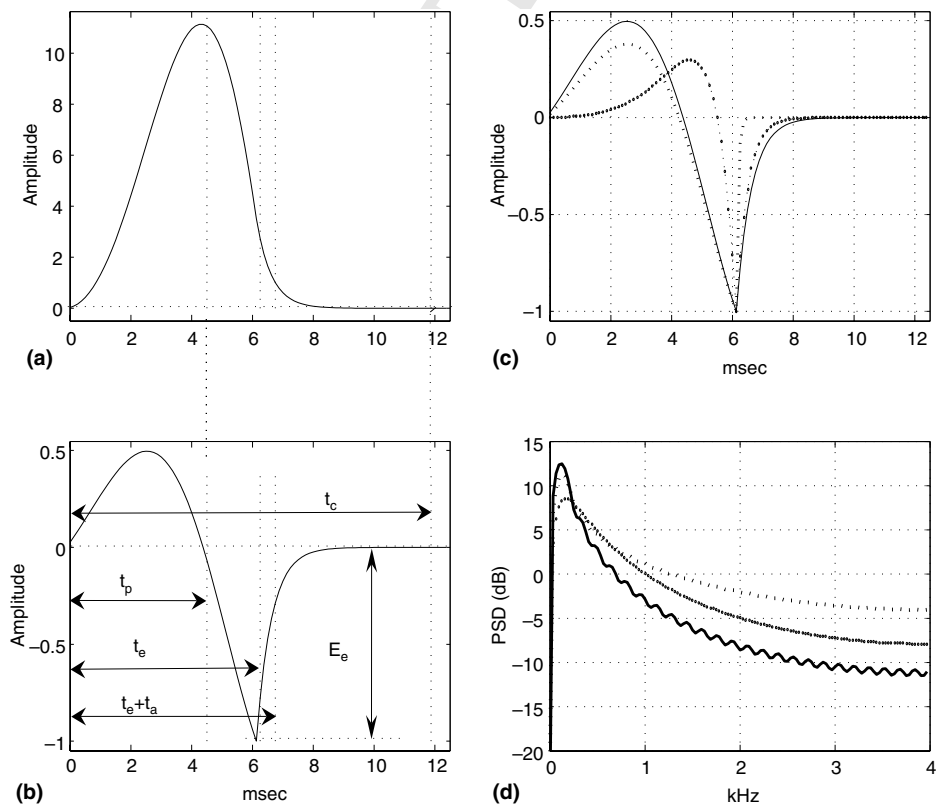


Fig. 2. Temporal and spectral structure of glottal pulse. (a) Glottal flow, (b) glottal flow derivative, (c) glottal flow derivative for three different sets of timing parameters, (d) spectra corresponding to (c).

247 H1–A3: This is the ratio of amplitude of first harmonic (H1) relative to that of the third-
 248 formant spectral peak (A3), and has
 249 been used by Hanson (1997) to charac-
 250 terize spectral tilt. The middle and high
 251 frequency components are directly
 252 affected by the duration of glottal pulse.
 253 An abrupt closure in the glottal cycle
 254 results in relatively strong middle and
 255 high frequency components, i.e. A3 is
 256 typically higher than what it would be
 257 with a gradual glottal closure. Fig. 3
 258 shows the comparison of harmonic spec-
 259 tra for vowel [ɑ] for dark and bright
 260 sounds. It clearly shows the differences
 261 in slopes of overall spectral envelopes.
 262 H1–H2: This is the ratio of amplitudes of the first
 263 two harmonics. It is affected by the glot-
 264 tal pulse skew as well as the open quo-
 265 tient (OQ) of the glottal pulse. The
 266 glottal pulse skew is measured by the
 267 speed quotient (SQ) which is computed
 268 as ratio of opening interval to the closing
 269 interval. While the open quotient is ratio
 270 of open glottal interval to the pitch per-

272 iod. A higher SQ indicates lower H1–H2
 273 (Doval and d’Alessandro, 1997), while
 274 higher OQ implies higher H1–H2.

275 Roll-off: Roll-off is an indicator of shape of the
 276 spectrum. This is defined as frequency
 277 within which 85% of the total accumu-
 278 lated magnitude is concentrated (Burred
 279 and Lerch, 2004). The roll-off is deter-
 280 mined by the largest DFT bin “ R ”,
 281 which satisfies

$$\sum_{k=N1}^{R-1} |X(k)| \leq 0.85 \sum_{k=N1}^{N2} |X(k)| \quad (1) \quad 283$$

284 where $X(\cdot)$ is the DFT spectrum of the in-
 285 put frame. $N1$ and $N2$ define the DFT
 286 bins corresponding to the range of fre-
 287 quencies over which the roll-off is com-
 288 puted. Based on the values of $N1$ and
 289 $N2$, R can represent the roll-off in any
 290 frequency region of interest. The roll-off
 291 computed over the full frequency spec-
 292 trum (0, 4000 Hz) captures the overall
 293 slope of the spectrum. For right-skewed
 294 spectra the value of roll-off turns to be
 295 high while for left-skewed spectra the

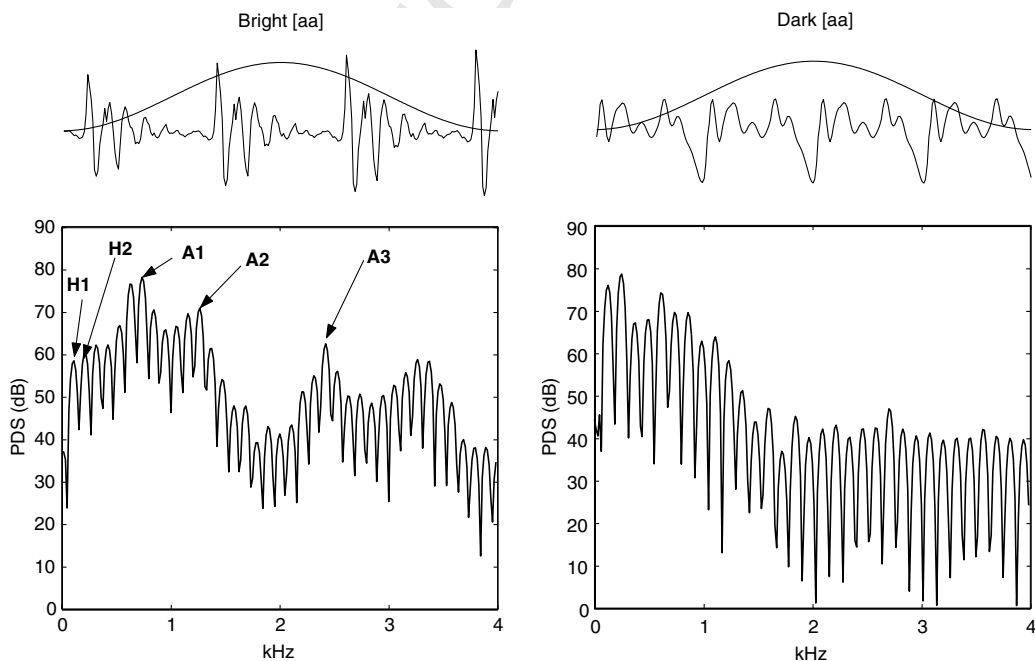


Fig. 3. Temporal waveforms and spectra of “dark” and “bright” vowel [ɑ]. Acoustic parameters labeled in “bright” [ɑ] are: H1 – amplitude of first harmonic, H2 – amplitude of second harmonic, A1 – amplitude of first formant, A2 – amplitude of second formant, A3 – amplitude of third formant.

296 roll-off is lower. For bright vowels we expect
297 a high roll-off.
298

299 It must be remarked that perceived voice quality
300 may also be influenced by recording conditions. The
301 spectral measures described in this section would
302 reflect both glottal waveform characteristics as well
303 as the transfer function of the recording setup.

304 4. Experimental evaluation

305 Experiments were designed to investigate the
306 influence of voice quality on the modeling error
307 from frequency-warped LP modeling of the spectral
308 envelope as presented in Section 2. Vowels uttered
309 in low pitched modal voices of different brightness
310 were obtained by synthesis as well as extraction
311 from natural speech.

312 4.1. Test set generation

313 Natural and synthetic samples corresponding to
314 eight distinct vowels, as shown in Table 1, were gener-
315 ated for the experiment. Synthesis of vowels
316 allowed the manipulation of the glottal timing
317 parameters which in turn enabled control over the
318 voice quality. The synthetic vowels were generated
319 using articulatory synthesis with the articulation
320 parameters estimated from provided target speech
321 based upon an analysis-by-synthesis approach
322 (Childers, 2000). For the synthesis of each of the
323 vowels the estimated articulatory parameters
324 together with an excitation signal, separately gener-
325 ated using an LF model with fixed parameters, are
326 used in a period-by-period synthesis. The synthe-
327 sized voice quality is dependent on the LF model
328 parameters (t_o , t_p , t_a , t_e , t_c and E_e) (Childers,
329 2000). By varying the timing parameters it is possi-
330 ble to generate variations in the spectral slope (and
331 consequently in the perceived brightness) of the

voice. The four sets of glottal pulse timing param- 332
eters used in the experiment are shown in Table 2. 333
Four instances of each vowel corresponding to each 334
set of glottal parameters (i.e. a total of 32 synthetic 335
sounds) were generated. Based on the resulting 336
overall spectral slope, we term these as “dark”, 337
“weak-dark”, “weak-bright” and “bright”. Each 338
vowel was generated at 120 Hz pitch (i.e. 339
 $t_o = 8.3$ ms) for a duration of 350 ms. The start 340
and end of the sound were tapered to avoid abrupt 341
transitions at the boundaries. 342

The set of natural vowels consisted of four 343
instances of each vowel (pitch ranging between 344
80 Hz and 120 Hz) taken from slowly uttered words 345
with neutral intonation from a set of six male 346
(Indian) speakers. The sounds were selected such 347
that two instances of each vowel were of bright 348
quality and two of dark quality as judged by listen- 349
ing as well as visual examination of the spectral 350
slope. The ends were tapered to eliminate abrupt 351
transitions. This resulted in a total set of 32 vowel 352
sounds (four instances of eight distinct vowels 353
selected in all across six speakers). In a few instances 354
of the vowel sounds (i, ɔ, u) it was found that the 355
actual vowel durations were in between 250 ms 356
and 280 ms and insufficient for judgment. The dura- 357
tion of such sounds were extended to reach the min- 358
imum duration of 350 ms by repeating the first two 359
periods at the start of the sounds and final two peri- 360
ods at the end of file. This artificial extension was 361
not detectable upon listening. The sounds were ana- 362
lyzed by the MBE model algorithm to estimate the 363
pitch and spectral amplitudes. The spectral ampli- 364
tudes thus obtained for each 20 ms speech frame 365
were modeled with 10th order frequency-warped 366
LP coefficients with a chosen frequency warping fac- 367
tor. The synthesis was carried out by standard sinu- 368
soidal synthesis method (MacAulay and Quatieri, 369
1995) using the spectral amplitudes obtained by fre- 370
quency-warped all-pole model approximation. It 371

Table 1
Set of vowels used in the experiment

Vowel ID	IPA symbol	Typical word	ARPABET (symbol)
1	ɑ	“guard”	[aa]
2	æ	“cat”	[ae]
3	ɔ	“orchid”	[ao]
4	ʌ	“cut”	[uh]
5	ou	“lotus”	[ow]
6	u	“boot”	[uw]
7	i	“beet”	[iy]
8	eɪ	“mate”	[ey]

Table 2

Glottal waveform timing parameters as a % of t_o , (Childers, 2000) used to generate the synthetic vowels with four voice qualities

Voice quality	Timing parameters		
	t_p (%)	t_e (%)	t_a (%)
Dark	35	50	15
Weak-dark	40	50	15
Weak-bright	40	50	2
Bright	45	50	2

$t_o = 8.3$ ms; $t_c = 0.9t_o$; $E_e = 40$.

372 was then compared with a reference sound synthe-
 373 sized using originally estimated spectral amplitudes.
 374 There were three test sounds for each of the 32 ref-
 375 erence sounds: LP modeled sound without fre-
 376 quency warping ('U'), LP modeled with a mild
 377 version of Bark-scale warping ('M') and LP mod-
 378 eled with Bark-scale warping ('B') based on the fre-
 379 quency warping scales described in (Rao and
 380 Patwardhan, 2005).

381 4.2. Subjective test

382 A subjective listening experiment was set up to
 383 compare the perceived qualities of the LP modeling
 384 with different warping factors. Six normal hearing
 385 listeners participated in the test. The test material
 386 was presented to the subject at normal listening lev-
 387 els through high quality head-phones connected to a
 388 PC sound card in a quiet room. For each of the test
 389 sounds the following "trio" presentation format
 390 was used: "reference-test-reference" where 250 ms
 391 silence separated the sounds. This format made
 392 the degradation of the test with respect to the refer-
 393 ence sound relatively easy to detect. The reference
 394 sound was also separately available for listening.

395 The subjects were asked to rank (using ranks 1, 2
 396 and 3) the relative perceived degradations of the test
 397 sounds U, M and B with respect to the correspond-
 398 ing reference sound for each of the vowel sounds in
 399 test set. Rank 1 would correspond to the least deg-
 400 radation. An undetectable difference would result in
 401 a suitable tied ranking. Subjects were allowed to lis-
 402 ten to the test trios associated with a given refer-
 403 ence sound any number of times before making a decision. Each
 404 listener did the test using the same set of items in dif-
 405 ferent orders on three separate occasions. Since it
 406 was found that the subjective ranks were nearly con-
 407 stant across listeners and trials, an overall ranking
 408 order was derived for each test vowel by combining
 409 the numerical ranks across listeners and trials as
 410 follows.

411 A numerical value is assigned to each rank in a
 412 single trial as follows: rank 1 = 2 points, rank
 413 2 = 1 point, rank 3 = 0 point. In case of a tie, equal
 414 number of points are given to both (i.e. two points
 415 each for first position or one point each for second
 416 position). The points are summed for each item
 417 across listeners to get the "total score". The total
 418 scores are then used to derive an overall subjective
 419 ranking. While the data within each item (i.e. refer-
 420 ence sound) can be compared, this is not so across
 421 items because the relative degradations were rated

from best to worst (ranked) only within each item 422
 set. E.g. a high score in item 1 may indicate percep- 423
 tually transparent modeling, while the same score in 424
 item 2 may indicate a small but clearly audible 425
 degradation. 426

4.3. Objective measurement of degradation 427

428 To quantify the perceived modeling error in each 428
 of the vowel sounds, we used partial loudness (PL), 429
 a psychoacoustical distance measure with demon- 430
 strated correlation with subjective judgement of 431
 spectral degradation in steady vowel sounds (Moore 432
 et al., 1997; Rao and Patwardhan, 2005). The spec- 433
 trum modeling error is treated as the signal whose 434
 loudness is to be estimated in the presence of a back- 435
 ground masker (the reference sound). The reference 436
 sound is intended to take the role of the background 437
 noise and the modeled sound that of the signal plus 438
 background noise (reference sound plus the model- 439
 ing error). The linear spectral distortion is treated 440
 as additive noise with power spectrum given by 441
 the difference between reference and modeled power 442
 spectra (Rao et al., 2001). The PL of the spectrum 443
 modeling error is computed for each frame and 444
 averaged across frames to obtain the PL value for 445
 the entire sample. 446

447 The objective rankings of relative degradation 447
 across U/M/B warping conditions were derived 448
 based upon the computed distortion measure for 449
 each of the reference-test sound pairs. The perfor- 450
 mance of an objective distance measure in predict- 451
 ing subjective judgments may be evaluated by 452
 computing a measure of correlation between corre- 453
 sponding objective and subjective rankings. The 454
 Spearman's correlation coefficient (Miller, 1989) is 455
 a suitable measure since it makes minimal assump- 456
 tions about the data. 457

4.4. Sentence level test 458

459 The experiments described so far involved quality 459
 judgments on isolated vowels. In order to obtain a 460
 perspective on the usefulness of the observations 461
 in the context of a speech coding application, a sen- 462
 tence-level listening test was designed. Phonetically 463
 balanced sentences as well as sentences which pre- 464
 dominantly contained one or another vowel were 465
 collected for a subjective listening experiment. 466
 Low pitched voices from a set of male speakers were 467
 selected. Table 8 lists the 15 sentences along with an 468
 indication of the dominant phonetic content as well 469

470 as perceived brightness of the voice. The sentences
 471 were uttered by speakers whose voices would be
 472 clearly classified as bright sounding or dark sound-
 473 ing in neutral speech.

474 Each of the sentences was subjected to MBE
 475 analysis and synthesis. The reference sentence was
 476 taken to be the one synthesized with the estimated
 477 spectral amplitudes while the test sentences were
 478 synthesized from 10th order LP modeled spectral
 479 amplitudes with and without Bark-scale warping
 480 to obtain “B” and “U” versions, respectively. Eight
 481 subjects were presented with the reference and the
 482 two corresponding test sentences, and asked to
 483 choose the one perceptually most similar to the refer-
 484 ence sentence. Based on a spectral shape measure
 485 described in Section 5.3, test sentences constructed
 486 by switching the warping parameter between “U”
 487 and “B” (depending on the nature of the frame)
 488 were also included in the listening test.

489 **5. Results and discussion**

490 Tables 3 and 4 summarise the results of subjec-
 491 tive and objective ranking of the experiment involv-
 492 ing the natural vowels. Tables 5 and 6 are the
 493 corresponding results for the synthetic vowels.

494 The best rank (rank 1) implies the least perceived
 495 degradation among the three test sounds, while the
 496 worst (rank 3) implies perceived degradation was

497 maximum among the three test sounds. Higher
 498 score, and lower PL value, implies lower perceived
 499 degradation. The last column contains Spearman’s
 500 rank correlation (Miller, 1989) which is an indicator
 501 of how close to each other the subjective and objec-
 502 tive ranks are. In this section, we comment on sev-
 503 eral aspects of the obtained experimental results.

504 *5.1. Subjective ranking experiment*

505 We first discuss the observations on frequency-
 506 warped all-pole modeling of natural vowels. From
 507 Table 3, which provides the subjective and objective
 508 evaluation of the reference-U, reference-M and refer-
 509 ence-B pairs of the dark sounding vowels, we
 510 observe that the reference-B pair is ranked better
 511 than the reference-U for all the vowels except for
 512 the vowels [eɪ] and [i]. This is consistent with the
 513 experimental results of (Rao and Patwardhan,
 514 2005). From Table 4 we observe that in the case
 515 of the “bright” vowels the reference-B pair is always
 516 degraded as compared to the reference-U pair. This
 517 is evident from the fact that it has always been
 518 ranked subjectively lower than the reference-U pair.
 519 This is contrary to what has been generally accepted
 520 in the literature. The reference-M pair is seen to be
 521 typically ranked between the other two warping
 522 scales. In the case of the synthetic vowels, we see
 523 from Table 5 that the dark vowels satisfy the same

Table 3
 Objective and subjective results on degradation due to spectral envelope modeling of natural vowels with dark voice quality

ID	IPA sym.	Sample	Subjective results						Objective results						Correlation coefficient
			Scores			Ranks			PL			Ranks			
			U	M	B	U	M	B	U	M	B	U	M	B	
1	ɑ	1	2	14	19	3	2	1	0.77	0.28	0.25	3	2	1	1
		2	5	20	28	3	2	1	0.86	0.14	0.09	3	2	1	1
2	æ	1	3	22	24	3	2	1	1.34	0.45	0.57	3	1	2	0.5
		2	4	22	23	3	2	1	0.77	0.23	0.23	3	1	1	0.75
3	ɔ	1	8	25	28	3	2	1	0.51	0.25	0.24	3	2	1	1
		2	7	20	27	3	2	1	1.28	0.82	0.61	3	2	1	1
4	ʌ	1	7	15	25	3	2	1	0.44	0.33	0.33	3	1	1	0.75
		2	5	23	22	3	1	2	1.29	1.01	0.98	3	2	1	0.5
5	ou	1	9	26	18	3	1	2	0.59	0.13	0.14	3	1	2	1
		2	4	20	27	3	2	1	1.28	0.17	0.15	3	2	1	1
6	u	1	1	26	18	3	1	2	0.85	0.11	0.13	3	1	2	1
		2	2	15	27	3	2	1	0.99	0.68	0.31	3	2	1	1
7	i	1	30	15	0	1	2	3	0.64	0.83	0.94	1	2	3	1
		2	24	21	4	1	2	3	0.38	0.32	0.42	2	1	3	0.5
8	eɪ	1	28	16	3	1	2	3	0.15	0.19	0.36	1	2	3	1
		2	24	20	3	1	2	3	0.09	0.16	0.2	1	2	3	1

All vowels were modeled with 10th LP model order with each of the three warping conditions: unwarped (U), mild-Bark-scale warped (M) and Bark-scale warped (B). “1”, “2” indicate samples from different voice qualities. The last column shows the Spearman’s correlation coefficient between subjectively and objectively measured ranks.

Table 4
Same as Table 3 but for natural vowels with bright voice quality

ID	IPA sym.	Sample	Subjective results						Objective results						Correlation coefficient
			Scores			Ranks			PL			Ranks			
			U	M	B	U	M	B	U	M	B	U	M	B	
1	ɑ	1	30	11	5	1	2	3	0.34	0.93	0.6	1	3	2	0.5
		2	30	15	0	1	2	3	0.11	0.35	0.49	1	2	3	1
2	æ	1	30	17	2	1	2	3	0.31	0.57	0.57	1	2	2	0.75
		2	25	23	8	1	2	3	0.43	0.24	0.48	2	1	3	0.5
3	ɔ	1	26	8	12	1	3	2	0.59	0.64	0.94	1	2	3	0.5
		2	14	7	5	1	2	3	0.04	0.06	0.14	1	2	3	1
4	ʌ	1	28	13	4	1	2	3	0.33	0.31	0.38	2	1	3	0.5
		2	29	16	0	1	2	3	0.43	0.49	0.78	1	2	3	1
5	ou	1	30	9	6	1	2	3	0.15	0.39	0.34	1	3	2	0.5
		2	28	16	1	1	2	3	0.29	0.16	0.39	2	1	3	0.5
6	u	1	30	14	1	1	2	3	0.43	1.04	1.02	1	3	2	0.5
		2	30	15	0	1	2	3	0.21	0.92	0.77	1	3	2	0.5
7	i	1	29	19	3	1	2	3	0.11	0.07	0.49	2	1	3	0.5
		2	30	15	0	1	2	3	0.1	0.37	0.83	1	2	3	1
8	er	1	28	17	2	1	2	3	0.3	0.3	0.37	1	1	3	0.75
		2	27	10	8	1	2	3	0.4	0.39	0.52	2	1	3	0.5

Table 5
Objective and subjective results on degradation due to spectral envelope modeling of synthetic vowels with dark voice quality

ID	IPA sym.	Sample	Subjective results						Objective results						Correlation coefficient
			Scores			Ranks			PL			Ranks			
			U	M	B	U	M	B	U	M	B	U	M	B	
1	ɑ	1	5	19	11	3	1	2	0.46	0.09	0.11	3	1	2	1
		2	6	14	13	3	1	2	0.11	0.10	0.09	3	2	1	0.5
2	æ	1	0	19	11	3	1	2	0.70	0.08	0.06	3	2	1	0.5
		2	3	19	11	3	1	2	0.29	0.06	0.04	3	2	1	0.5
3	ɔ	1	3	19	14	3	1	2	0.82	0.11	0.12	3	1	2	1
		2	5	18	15	3	1	2	0.32	0.10	0.13	3	1	2	1
4	ʌ	1	9	13	11	3	1	2	0.67	0.29	0.44	3	1	2	1
		2	5	13	14	3	2	1	0.39	0.32	0.28	3	2	1	1
5	ou	1	4	11	16	3	2	1	0.42	0.13	0.12	3	2	1	1
		2	2	13	15	3	2	1	0.25	0.17	0.13	3	2	1	1
6	u	1	7	12	17	3	2	1	0.38	0.10	0.12	3	1	2	0.5
		2	0	16	14	3	1	2	0.36	0.10	0.12	3	1	2	1
7	i	1	16	16	4	1	1	3	0.42	0.21	0.27	3	1	2	-0.25
		2	16	15	2	1	2	3	0.22	0.18	0.35	2	1	3	0.5
8	er	1	18	14	2	1	2	3	0.19	0.15	0.24	2	1	3	0.5
		2	20	10	0	1	2	3	0.06	0.12	0.15	1	2	3	1

All vowels were modeled with 10th LP model order with each of the three warping conditions: unwarped (U), mild-Bark-scale warped (M) and Bark-scale warped (B). “1”, “2” indicate samples from different voice qualities. The last column shows the Spearman’s correlation coefficient between subjectively and objectively measured ranks.

524 trend as the dark natural vowels (that is reference-B
525 pair has been ranked better except for [i] and [er]).
526 Table 6 for the bright vowels, shows some inconsis-
527 tencies in [ɔ] and [u] while all other vowels show the
528 expected superiority of reference-U. An explanation
529 for the inconsistencies is provided in Section 5.3.

5.2. Correlation with partial loudness

530
531 Tables 3–6 indicate an overall high positive cor-
532 relation between subjective ranks and partial loud-
533 ness based ranks. In fact, the relative positions of
534 non-warped and Bark-warped in the subjective pref-
535 erence are correctly captured by the objective dis-

Table 6

Same as Table 5 but for synthetic vowels with bright voice quality

ID	IPA sym.	Sample	Subjective results						Objective results						Correlation coefficient
			Scores			Ranks			PL			Ranks			
			U	M	B	U	M	B	U	M	B	U	M	B	
1	ɑ	1	15	14	5	1	2	3	0.06	0.07	0.15	1	2	3	1
		2	15	12	6	1	2	3	0.05	0.06	0.15	1	2	3	1
2	æ	1	2	18	10	3	1	2	0.13	0.04	0.09	3	1	2	1
		2	15	13	3	1	2	3	0.08	0.03	0.09	2	1	3	0.5
3	ɔ	1	1	19	12	3	1	2	0.43	0.05	0.20	3	1	2	1
		2	8	18	12	3	1	2	0.28	0.06	0.23	3	1	2	1
4	ʌ	1	13	14	5	2	1	3	0.25	0.34	0.5	1	2	3	0.5
		2	14	13	6	2	1	3	0.31	0.52	0.43	1	3	2	0.5
5	ou	1	12	12	8	2	1	3	0.15	0.17	0.22	1	2	3	0.5
		2	14	13	6	1	2	3	0.19	0.16	0.27	2	1	3	0.5
6	u	1	3	14	19	3	2	1	0.70	0.16	0.17	3	1	2	0.5
		2	2	15	13	3	1	2	0.55	0.20	0.27	3	1	2	1
7	i	1	18	14	2	1	2	3	0.19	0.09	0.56	2	1	3	0.5
		2	12	17	2	2	1	3	0.37	0.32	0.75	2	1	3	1
8	eɪ	1	20	10	2	1	2	3	0.05	0.13	0.22	1	2	3	1
		2	20	10	7	1	2	3	0.05	0.14	0.26	1	2	3	1

536 tance. By and large, the only inconsistencies are
537 with respect to the relative position of mild-Bark-
538 warped test sounds.

539 The partial loudness is an indication of how audi-
540 ble is the spectral distortion. It is the sum of the par-
541 tial specific loudness contributions of the distortion
542 distributed across the signal spectrum. In frequency-
543 warped LP modeling, there is an improved spectral
544 match in the low frequency region at the cost of
545 increased errors in the high frequency region.
546 Whether frequency warping improves the overall
547 perceptual accuracy of the fit will depend on
548 whether the increase in the partial loudness of the
549 high frequency distortion is exceeded by the
550 decrease in the partial loudness of the low frequency
551 distortion. Fig. 4 shows, for a particular frame of a
552 dark vowel sound [ɑ] the comparison between Bark-
553 warped and non-warped LP modeling in terms of
554 spectral amplitude distortion and partial loudness
555 distribution. On comparing Fig. 4(a) and (b), we
556 note the reduced low frequency distortion and
557 increased high frequency distortion introduced by
558 the Bark-scale frequency warping. Further, a com-
559 parison of the partial loudness distributions of
560 Fig. 4(c) and (d) indicates that the reduction in the
561 low frequency distortion contributes more to
562 decreasing audible error than the effect of the
563 increase in the high frequency distortion. The cor-
564 rect prediction of this perceptual effect by the partial
565 loudness model indicates that the high frequency

566 spectral distortion is masked to a great extent by
567 the relatively strong low frequency signal spectral
568 components. This explanation is consistent with
569 the opposite behaviour of the front vowels [i] and
570 [eɪ] since their low first formant causes reduced
571 spread of masking to the high frequency region. In
572 the corresponding plots of Fig. 5 for a frame of a
573 bright vowel [ɑ], we note the similar effect. That
574 is, the reduction in the low frequency distortion
575 due to Bark-scale warping does not compensate
576 adequately for the increased high frequency distor-
577 tion as observed in the partial loudness distributions
578 of Fig. 5(c) and (d). The higher “loudness” of the
579 high frequency error can be attributed to the insuf-
580 ficient masking from the low frequency components
581 of the reference signal which are only comparable in
582 strength to the high frequency components.

583 5.3. Relation to spectral cues

584 Although the partial loudness is able to predict
585 the subjectively preferred warping condition accu-
586 rately in most cases, it is desirable to have a less
587 computationally complex measure for use in prac-
588 tice. The simple spectral cues of Section 3.2 capture
589 the overall spectral balance in some way or other
590 and therefore merit investigation.

591 Fig. 6 plots the H1–A3 of each of the set of 32
592 natural vowels plus 32 synthetic vowels used in the
593 experiments. The H1–A3 is measured for a repre-

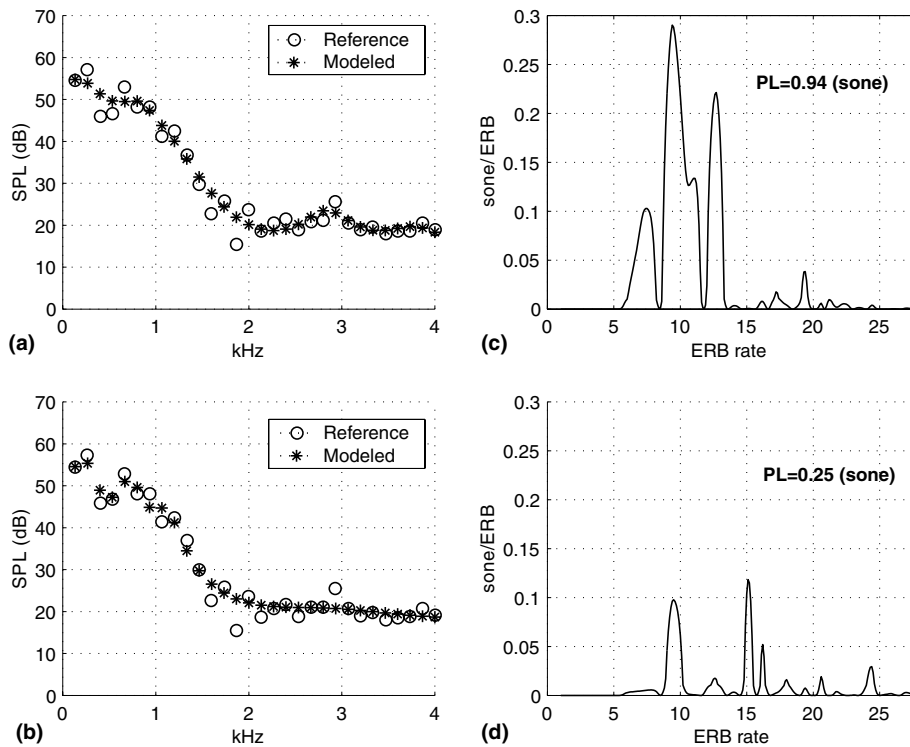


Fig. 4. Modeling of “dark” natural [a] with pitch = 105 Hz, without and with Bark-scale frequency warping at LP model order = 10. (a) and (b) Spectral harmonic amplitudes of reference (O) and modeled (*) sounds are connected with a dotted line to show the spectral envelopes without warping and with Bark-scale warping respectively. (c) and (d) Partial loudness distribution of spectral distortion in unwarped and Bark-warped modeling, respectively.

594 tentative frame picked from center of steady vowel
 595 sound. The bright sounding vowels fall at the lower
 596 end of the y -axis relative to the corresponding vow-
 597 els with dark voice quality.

598 Based on the results of Tables 3 and 4, the vowel
 599 instances that have been ranked better with Bark-
 600 scale frequency-warped all-pole modeling are shown
 601 encapsulated by a triangle while the ones without
 602 warping before modeling are shown encapsulated
 603 by a square. Fig. 6 clearly shows that the sounds
 604 with higher H1–A3 benefit more from Bark-scale
 605 warping whereas the ones with lower H1–A3 benefit
 606 less (but for the expected cases of [i] and [e]). In
 607 fact, from the data on natural vowels, it appears
 608 that H1–A3 = 10 dB may be considered a rough
 609 threshold for prediction on whether Bark-warping
 610 before LP modeling would help to improve per-
 611 ceived quality. From Fig. 6(b) we note that the syn-
 612 thetic vowels [u] and [ɔ] in our set had H1–A3 values
 613 clustered at the higher end of the scale which may
 614 explain why they do not show any variation in the
 615 subjectively preferred warping condition. It is seen
 616 that while H1–A3 captures the overall spectral tilt,

617 it is insensitive to the relative distribution of spectral
 618 energy in the low frequency region, and hence can-
 619 not distinguish the low first formant characteristic
 620 of the vowels [i] and [e].

621 In the view of these above observations, it is
 622 desirable to have a measure that describes not only
 623 the overall frequency concentration but also that in
 624 a specific frequency region. The spectral roll-off
 625 defined in Section 3.2 was adapted for the purpose.
 626 Experiments revealed that the spectral roll-off value
 627 of (1) evaluated in the region (0, 1500 Hz) was a
 628 good indicator of the lowness of the first formant
 629 as well as of the largeness of the gap between the
 630 first and second formant. That is, this measure is
 631 low only in cases of front high and front-mid vowels
 632 where there is large gap between F1 and F2. Further
 633 the spectral roll-off computed in the reverse direc-
 634 tion in the same interval was a good indicator of
 635 the low frequency skew as represented by H1–H2.
 636 Based on these considerations and experimentally
 637 determined thresholds, the following roll-off based
 638 rule was applied to determine whether warping
 639 should be enabled:

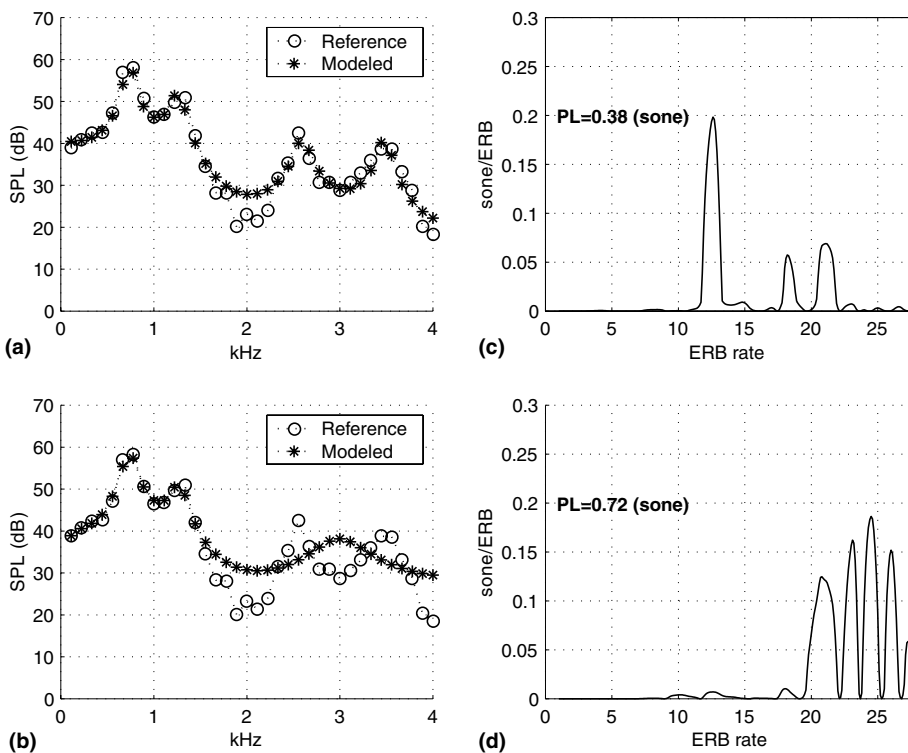


Fig. 5. Modeling of “bright” natural [ɑ] with pitch = 105 Hz, without and with Bark-scale frequency warping at LP model order = 10. (a) and (b) Spectral harmonic amplitudes of reference (O) and modeled (*) sounds are connected with a dotted line to show the spectral envelopes without warping and with Bark-scale warping respectively. (c) and (d) Partial loudness distribution of spectral distortion in unwarped and Bark-warped modeling, respectively.

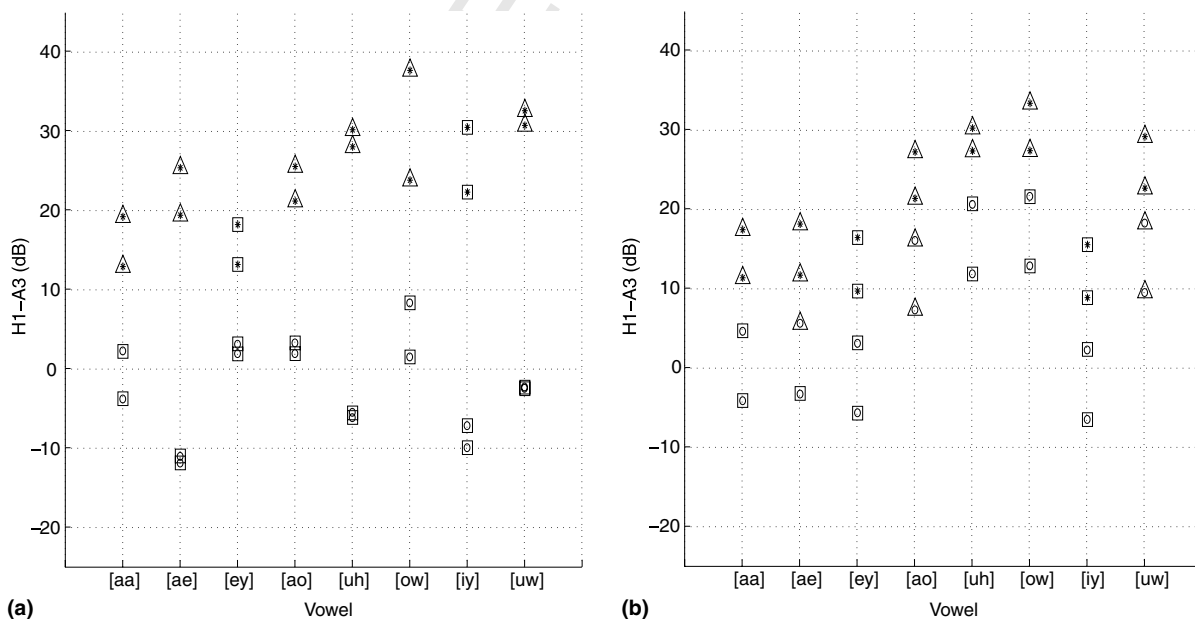


Fig. 6. Spectral cue H1–A3 along with subjective preference of quality between Bark-warped and unwarped modeling (a) natural vowels and (b) synthetic vowels. (Δ) Bark-scale warping preferred. (□) Unwarped preferred, “o”: Vowels with “bright” voice quality, “*”: Vowels with “dark” voice quality.

if (SR < 1700 Hz) and (PSR > 400 Hz) and
 (RPSR < 350 Hz), use Bark warping.

SR is the roll-off computed (0, 4000 Hz), PSR is the
 roll-off computed in (0, 1500 Hz), RPSR is the roll-
 off computed in (1500, 0 Hz). Table 7 shows the nor-
 malized correlation of the objectively predicted pre-
 ferred warping condition with the subjectively
 preferred condition for each vowel for the three
 objective measures: PL, H1–A3 and roll-off based.
 We see that the roll-off based rule does better than
 H1–A3 and is comparable to PL while being signif-
 icantly less complex computationally.

5.4. Some comments on nasalised vowels

Nasalisation of vowels is known to be associated
 with a change in the first formant region of the spec-
 trum that may be modeled with an added pole and
 zero, and therefore expected to be difficult to model
 with a low order LP model. Since a prominent fea-
 ture of nasal vowels lies in the low frequency spec-
 trum, it is expected that Bark-scale warping before
 the LP modeling would help to retain the perception
 of nasality. Informal listening experiments with arti-
 ficially nasalized natural vowels support the above
 observation. To nasalize the vowels, a filter was
 applied that had a zero and a pole located above
 the first formant (Feng and Castelli, 1996). A sepa-
 ration of 100 Hz between F1 (frequency of first
 formant), FNZ (frequency of zero) and FNP (fre-
 quency of pole) is sufficient to introduce perceptible
 nasality in the vowel. In a listening experiment with
 nasalised vowels, when asked to rank overall qual-
 ity, however, with and without Bark-scale warping
 before LP modeling, subjects often judged that
 modeling without warping was less degraded com-

pared with modeling with Bark-scale warping, espe-
 cially in the case of bright vowels. In the case of
 dark vowels too, high frequency distortion from
 warping was sometimes more significant than in
 the case of the corresponding non-nasalised vowel.
 This may be attributed to lowered masking from
 reduced low frequency components due to the pres-
 ence of a spectral zero. Overall there was less consis-
 tency between subjects and also less correlation
 between subjective and objective rankings (based
 on the partial loudness measure) compared with
 the results on the non-nasalised vowels. This may
 be explained by the existence of conflicting percep-
 tual effects from the two different spectral cues
 namely of loss of nasality, and high frequency spec-
 trum distortion. In such cases, it is expected that
 higher level cognition would come into play in a
 subject's judgement of rank, a situation that cannot
 be predicted from a low level distance measure
 based on a model of the peripheral auditory system.

5.5. Sentence level test

Table 8 shows the preferred warping condition
 (as selected between non-warped, “U” and Bark-
 warped, “B”) in terms of overall subjective quality
 on sentences. We note that in the case of the dark-
 voiced sentences, subjects preferred the Bark-
 warped condition except when the phonemes [eɪ]
 and [i] dominated. In the case of the bright voices,
 the subjects preferred modeling without warping in
 all cases. These observations are consistent with
 the isolated phoneme results discussed earlier. For
 the speech coding application, this suggests the pos-
 sibility of improving overall perceived quality by
 making available a limited set of warping factors
 for dynamic selection based on frame-level spectral
 characteristics. Information on the selected warping
 factor can be conveyed to the decoder via one or
 two bits depending on the number of distinct warp-
 ing factors available.

In the second part of the experiment instead of
 switching the frames manually, the prediction rule
 of Section 5.3 was implemented and warping condi-
 tion was switched automatically frame by frame. A
 subjective listening test comparing the dynamically
 switched warped sentences with the preferred fixed
 warping condition of Column 5 of Table 8 revealed
 that in all cases the subjects rated the former either
 better or indistinguishable from the latter. This is
 shown in Column 6 of Table 8. The percentage of
 frames that were selectively switched to Bark-scale

Table 7

Normalised correlation values between subjectively observed and
 objectively predicted preferred warping condition (between non-
 warped and Bark-warped only) for each vowel set

Vowel	Natural			Synthetic		
	PL	H1–A3	Roll-off	PL	H1–A3	Roll-off
ɑ	1	1	1	1	1	1
æ	1	1	1	1	0.75	0.5
ɔ	1	1	1	1	0.75	1
ʌ	1	1	1	1	0.5	0.5
ou	1	1	1	1	0.5	0.5
u	1	1	1	1	0.75	1
i	1	0.5	1	0.75	0.75	1
eɪ	1	0.5	1	1	0.75	1

Table 8

Sentences used in the subjective evaluation of preferred frequency warping at the sentence level

Sr. no.	Sentences	Voice quality	Predominant phoneme	Subjective preference		% Switched to B
1	All the balls were brought from shopping mall	Dark	[ɔ]	B	S	49
2	By and large he was un-harmed	Dark	[ɑ]	B	S	42
3	Lord Paul was tall	Dark	[ɔ]	B	N	71
4	Early bird earns a worm	Dark	[ʌ]	B	N	57
5	Can that man carry those pans	Dark	[æ]	B	S	39
6	Draw each graph on new axis	Dark	Balanced	B	S	35
7	They may say the same in Spain	Dark	[eɪ]	U	N	15
8	His meal is eel meat	bright	[i]	U	N	0
9	We were away to walla walla	Dark	[eɪ], [i], [ɑ]	U	S	55
10	They all agree that the essay is barely intelligible	Bright	Balanced	U	N	17
11	Thick glue oozed out of the tube	Bright	[i], [ɔ]	U	N	4
12	Dont ask me to carry an oily rag like that	Bright	Balanced	U	N	4
13	A muscular abdomen is good for your back	Bright, nasal	Balanced	U	N	7
14	Withdraw only as much money as you need	Dark, nasal	Balanced	B	N	48
15	Withdraw only as much money as you need	Bright, nasal	Balanced	U	N	13

The first sub-column of the “subjective preference” indicates the warping condition preferred by listeners between two fixed warping conditions (“U” and “B”). The second sub-column indicates the subjective preference when the automatic warping condition prediction rule was used. “S”: automatically switched preferred, “N”: unable to differentiate and “O”: fixed warping preferred.

724 warped modeling is shown in the last column of
725 Table 8. We find that, for sentences no. 7, 8, 10,
726 11, 12, 13 and 15, the automatic algorithm keeps
727 most frames unwarped. These sentences are charac-
728 terized by either bright voice or presence of [eɪ] and
729 [i], which are better modeled without warping.

730 At this point, it is relevant to comment on the use
731 of high frequency pre-emphasis commonly used to
732 improve LP modeling of speech spectra (Markel
733 and Gray, 1976). Pre-emphasis serves to suppress
734 the spectral tilt to an extent and thus improve the
735 LP modeling of the formants. However it is sug-
736 gested that pre-emphasis be avoided for unvoiced
737 speech, while in the case of voiced speech, the pre-
738 emphasis factor should ideally be signal dependent
739 (Markel and Gray, 1976). This is also borne out by
740 our own observations where we found that while
741 pre-emphasis improves the modeling of voices with
742 high spectral tilt (such as breathy, female voices), it
743 degrades sounds with low energy at low frequencies
744 [also noted by (Wong et al., 1980)] including sound
745 spectra characterised by low H1–H2. Pre-emphasis
746 was not used in the present study. It is anticipated
747 that the selective use of pre-emphasis combined with
748 frequency-scale warping based on the considerations
749 presented in this paper can help to achieve further

improvements in the perceptual accuracy of low 750
order LP modeling of the speech spectral envelope. 751

6. Conclusion 752

Frequency-warping according to a perceptual 753
scale is often applied to improve the perceptual 754
accuracy of low order LP modeling of the speech 755
spectral envelope in sinusoidal speech coding. 756
Understanding the factors that influence the subjective 757
perception of spectral envelope modeling errors 758
can be useful in improving the attained speech qual- 759
ity. Steady vowel sounds are particularly sensitive to 760
spectrum envelope modeling errors. Experimental 761
investigations of the relative improvement in per- 762
ceived quality from the use of different warping 763
functions on natural and synthetic steady vowels 764
have been presented. Contrary to what is generally 765
accepted, it is found that whether the widely used 766
Bark-scale frequency warping is effective depends 767
on the vowel and voice quality attributes of the sig- 768
nal spectrum. Dark voices with their steep spectral 769
slopes were found to benefit from Bark-scale warp- 770
ing but not so the relatively bright voices. The 771
exceptions to this are the front-mid and front-low 772
vowels [eɪ] and [i]. These vowels are observed to typ- 773

ically degrade when modeled with Bark-scale frequency warping relative to the reproduced quality without warping. In summary, our studies based on observations presented in this paper and a previous paper (Rao and Patwardhan, 2005) suggest that the vowels [ei] and [i] dominate over the dark/bright voice quality differences.

The subjective ratings of relative perceived degradation were closely predicted by an objective distance measure based on the partial loudness of the spectrum distortion. This suggested that frequency masking plays a role in determining whether the improved low frequency spectrum match from Bark-scale warping is sufficient to compensate for the accompanying high frequency spectrum distortion. Only when the low frequency components of the spectrum are relatively strong and sufficiently spread is there a clear benefit from Bark-scale warping. This condition was shown to be closely linked to the spectral slope as quantified by measures such as H1–A3 and roll-off which can therefore act as useful cues in predicting the suitability of a warping function for a particular sound. A rule has been proposed to predict the subjectively preferred warping condition. The rule is based on a combination of spectral roll-off values computed in different regions of the spectrum.

A sentence-level listening test confirmed the results of the isolated vowel experiments and also demonstrated the effectiveness of the warping prediction rule. For the speech coding application, this suggests the possibility of improving overall perceived quality by making available a limited set of warping factors for dynamic selection based on frame-level spectral characteristics. Information on the selected warping factor can be conveyed to the decoder via one or two bits depending on the number of distinct warping factors available. Alternatively, a simple but slightly less effective solution is to use fixed warping according to a mild version of the Bark-scale.

Audio demonstrations may be accessed at http://www.ee.iitb.ac.in/~prao/speech_comm_1/index.html.

7. Uncited references

Klatt and Klatt (1990) and Hanson and Chuang (1999).

Appendix A. Results on synthetic vowels

A subjective ranking experiment similar to that carried out for natural set of vowels was carried

out on the synthetic vowels. The results appear in Tables 5 and 6. The trend nearly matches with the natural vowels, except for the dark vowel [IY], which may be attributed to the difficulty of carefully controlling the output of the articulatory synthesizer for this particular phoneme. (We found that the articulatory synthesizer produced a harmonic envelope which became less smooth on the higher frequency end for this phoneme. The partial loudness based ranking was not consistent with the subjective test in such cases probably due to the deviation from the model assumptions.)

References

- Burred, J., Lerch, A., 2004. Hierarchical automatic audio signal classification. *J. Audio Eng. Soc.* 52 (8), 724–739.
- Champion, T., MacAulay, R., Quatieri, J., 1994. High-order all-pole modeling of the spectral envelope. In: *Proceedings of the IEEE International Conference on Acoustics, Speech, Signal Processing*, pp. 529–532.
- Childers, D., 2000. *Speech Processing and Synthesis Toolboxes*. John Wiley and Sons.
- Childers, D., Lee, C., 1991. Vocal quality factors: analysis, synthesis and perception. *J. Acoust. Soc. Am.* 90 (5), 2394–2411.
- Doval, B., d’Alessandro, C., 1997. Spectral correlates of glottal waveform models: an analytic study. In: *Proceedings of the IEEE International Conference on Acoustics, Speech, Signal Processing*, pp. 1295–1298.
- Feng, G., Castelli, E., 1996. Some acoustic features of nasal and nasalized vowels: a target for vowel nasalization. *J. Acoust. Soc. Am.* 99 (6), 3694–3706.
- Griffin, D., Lim, J., 1988. Multiband excitation vocoder. *IEEE Trans. Acoust. Speech Signal Process.* 36 (8), 1223–1235.
- Hanson, H., 1997. Glottal characteristics of female speakers: acoustic correlates. *J. Acoust. Soc. Am.* 101 (1), 466–481.
- Hanson, H., Chuang, E., 1999. Glottal characteristics of male speakers: acoustic correlates and comparison with female data. *J. Acoust. Soc. Am.* 106 (2), 1064–1077.
- Hermansky, H., Hanson, B., Wakita, B., Fujisaki, H., 1985. Linear predictive modeling of speech in modified spectral domain. In: *Digital Processing of Signals in Communications*, pp. 55–63.
- Klatt, D., Klatt, L., 1990. Analysis, synthesis and perception of voice quality variations among female and male talker. *J. Acoust. Soc. Am.* 87 (2), 820–855.
- MacAulay, R., Quatieri, T., 1995. *Sinusoidal Coding*, in *Speech Coding and Synthesis*. Elsevier, Amsterdam.
- Markel, J., Gray, A., 1976. *Linear Prediction of Speech*. Springer Verlag, Berlin.
- Miller, J., 1989. *Correlation*, in *Statistics for Advanced Level*. Cambridge University Press.
- Molyneux, D., Parris, C., Sun, X., Cheetham, B., 1998. Comparison of spectral estimation techniques for low bit-rate speech coding. In: *Proceedings of the International Conference on Spoken Language Processing*, pp. 946–949.

- 877 Moore, B., Glasberg, B., Baer, T., 1997. Model for prediction of 883
878 thresholds, loudness and partial loudness. *J. Audio Eng. Soc.* 884
879 45 (4), 224–240. 885
880 Rao, P., Patwardhan, P., 2005. Frequency warped modeling of 886
881 vowel spectra: dependence on vowel quality. *Speech Com-* 887
882 *mun.* 47, 322–335. 888
- Rao, P., van Dinther, R., Veldhuis, R., Kohlrausch, A., 2001. A 883
measure for predicting audibility discrimination thresholds 884
for spectral envelope distortions in vowel sounds. *J. Acoust.* 885
Soc. Am 109 (4), 2085–2097. 886
- Wong, R., Hsiao, C., Markel, J., 1980. Spectral mismatch due to 887
preemphasis in LPC analysis/synthesis. *IEEE Trans. Acoust.* 888
Speech Signal Process. ASSP 28 (2), 263–264. 889
890

UNCORRECTED PROOF