

# COMPRESSION OF ACOUSTIC INVENTORIES USING ASYNCHRONOUS INTERPOLATION

Alexander B. Kain and Jan P. H. van Santen

Center for Spoken Language Understanding  
OGI School of Science & Engineering at OHSU  
20000 NW Walker Road, Beaverton, OR 97006, USA

## ABSTRACT

A compression method is proposed that takes advantage of a powerful property of acoustic unit inventories: In the appropriate acoustic space, units that share a (context-dependent or -independent) phoneme label must be close to a vector *phoneme template* associated with the phoneme. The method approximates units by interpolation between templates. The interpolation operation involves two asynchronous weight functions operating on the template. One is associated with spectral peak locations, the second with spectral balance. This enables approximating transitions such as [i:] → [v], in which formant movement precedes frication onset. The algorithm guarantees smooth concatenation points.

## 1. INTRODUCTION

Recent advances in hardware technologies have introduced fast processors on hand-held devices, such as personal digital assistants and mobile phones. However, what has prevented the most common form of text-to-speech (TTS) synthesis, *concatenative* TTS, from becoming ubiquitous on these devices is the large amount of memory that is required for storing a system’s *acoustic inventory*, which contains the necessary acoustic units. Moreover, recent research has shown that substantial voice quality improvements can be obtained when inventories are used that contain not only phoneme pair units, or *diphones*, but also longer units, entire words, and even phrases. Thus, there is a quality/size trade-off in concatenative TTS. It follows that compression technologies are needed for any storage-limited device, even if the cost of storage continues to decline.

Over the last several decades, the area of speech coding has created many technologies for speech compression, and some of these have been used for compressing acoustic units [1]. However, these technologies have not produced the required compression ratios. This is no surprise, because speech coding applications differ in several key respects from synthesis:

- In an acoustic inventory (AI), units are spoken by the same speaker under acoustically constant and noise free conditions. By contrast, speech coding applications are designed to work under a variety of acoustic conditions and cannot be fine-tuned for a specific speaker.

---

This research was conducted with support from NSF Grants 0117911 “Making Dysarthric Speech Intelligible” and 0082718 “Modeling Degree of Articulation for Speech Synthesis”. We are grateful to Gilead Cohen and Johan Wouters for helpful discussions, and to Mike Macon for inspiring us to work on this topic.

- The AI encoder can be arbitrary complex, since encoding takes place only once, off-line, on a large, static, random access data set. Speech compression algorithms on the other hand must usually encode and decode in real-time, working on very short signals.
- A high quality AI has the *close acoustic match* property, by which is meant that, in an appropriate acoustic feature space, units that share a (context-dependent or -independent) phoneme label (joinable units) are close to a common feature vector, or *template*, associated with the common phoneme. For example, the /p/→/i:/ and /i:/→/tS/ diphone units have been selected from a larger corpus of recorded speech to be maximally compatible in the sense that their concatenation does not lead to audible spectral discontinuities.
- Additional information is available in the inventory, such as pitch marks and acoustic segment boundaries, which can be used advantageously by the compression algorithm.

In this paper, we propose a special purpose compression algorithm that achieves a high rate of compression while eliminating spectral discontinuities. We implement and evaluate the algorithm using a diphone inventory, but the approach can be extended to operate on larger and variable size unit inventories as well.

## 2. METHOD

### 2.1. Core idea

The core idea of our method is to approximate a diphone by interpolating between a left and a right phoneme template. The general idea of using interpolation for compression is not new, as exemplified by Atal’s Temporal Decomposition method [2], a general purpose compression method in which speech frames are approximated by linear combinations of *basis vectors*. However, the proposed method is fundamentally different in that (i) the templates correspond to known phonetic labels instead of arbitrary (and numerous) frames, and (ii) it involves *two asynchronous, non-linear* interpolation operations. The latter is based on the realization that during transitions between certain phonetic segments, articulators and their acoustic manifestations often change asynchronously. For example, a transition from a high vowel to a /v/, as in the word “leave”, shows a lowering of the second formant that starts well before the onset of frication. Simple interpolation between an /i:/ and /v/ would create an unnatural transition in which either the frication onset is too early and smooth or in which the formant movement is too late and sudden. We need a method that can track asynchronous formant and spectral balance changes.

We represent a diphone in the form of  $\mathbf{D}_{n,m}^{l \rightarrow r}$ , which is the  $n^{\text{th}}$  frequency component of the complex spectrum of a short segment

of speech at a discrete time point, or frame,  $m$  of the original di-  
phone with phoneme labels  $l$  to  $r$ . In our specific method, we let

$$\mathbf{D}_{n,m}^{l \rightarrow r} \approx \hat{\mathbf{D}}_{n,m}^{l \rightarrow r} = \mathbf{T} \left( \mathbf{P}^l, \mathbf{P}^r, \alpha_m^{l \rightarrow r}, \beta_m^{l \rightarrow r} \right)_n \quad (1)$$

$$n = 1, 2, \dots, N \quad m = 1, 2, \dots, M$$

where  $\hat{\mathbf{D}}_{n,m}^{l \rightarrow r}$  is the compressed di-  
phone, and  $\mathbf{T}$  is the *transition function*.  $\mathbf{T}$  takes as arguments *phoneme templates*  $\mathbf{P}^l$  and  $\mathbf{P}^r$ , as  
well as *transition weights*  $\alpha_m^{l \rightarrow r}$  and  $\beta_m^{l \rightarrow r}$ . A phoneme template is  
defined as the representative complex spectrum of a typical rendi-  
tion of a phoneme, whereas transition weights govern the temporal  
evolution of the transition.

Expanding the right hand side of Eq. 1, we define

$$\begin{aligned} & \mathbf{T} \left( \mathbf{P}^l, \mathbf{P}^r, \alpha_m^{l \rightarrow r}, \beta_m^{l \rightarrow r} \right)_n \\ &= \beta_m^{l \rightarrow r} \cdot \mathbf{M} \left( \mathbf{P}^l, \mathbf{v}_m^{l \rightarrow r} \right)_n \oplus \left( 1 - \beta_m^{l \rightarrow r} \right) \cdot \mathbf{M} \left( \mathbf{P}^r, \mathbf{v}_m^{l \rightarrow r} \right)_n \end{aligned} \quad (2)$$

$$0 \leq \beta_m^{l \rightarrow r} \leq 1$$

where  $\mathbf{v}_m^{l \rightarrow r}$  is the  $K$ -dimensional *feature parameter state* at frame  
 $m$ , defined as

$$\mathbf{v}_{k,m}^{l \rightarrow r} = \alpha_m^{l \rightarrow r} \cdot \mathbf{V} \left( \mathbf{P}^l \right)_k + \left( 1 - \alpha_m^{l \rightarrow r} \right) \cdot \mathbf{V} \left( \mathbf{P}^r \right)_k \quad (3)$$

$$k = 1, 2, \dots, K \quad 0 \leq \alpha_m^{l \rightarrow r} \leq 1.$$

The function  $\mathbf{V}(\mathbf{P})$  returns the feature parameters of a vocal tract  
model for  $\mathbf{P}$ , whereas  $\mathbf{M}(\mathbf{P}, \mathbf{v}_m^{l \rightarrow r})$  is a speech modification func-  
tion such that

$$\mathbf{V} \left( \mathbf{M} \left( \mathbf{P}, \mathbf{v}_m^{l \rightarrow r} \right) \right) = \mathbf{v}_m^{l \rightarrow r}.$$

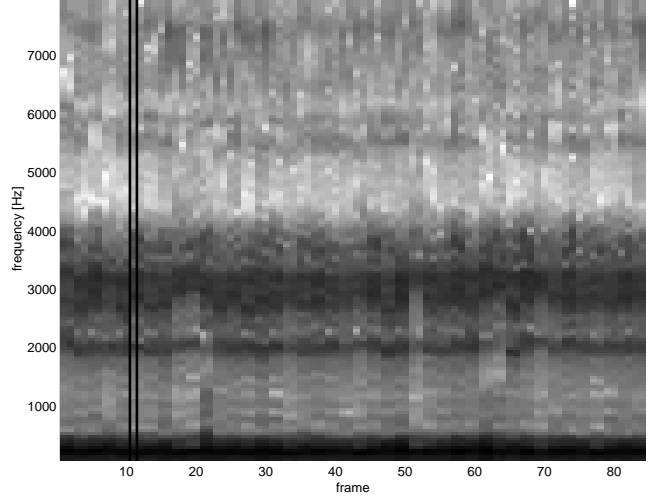
In other words,  $\mathbf{M}$  modifies  $\mathbf{P}$  to have features  $\mathbf{v}_m^{l \rightarrow r}$ . We use  
the operator  $\oplus$  to indicate addition taking place separately in the  
log magnitude and unwrapped phase domains, with subsequent re-  
combination.

Examining Eq. 3, we observe that  $\mathbf{v}_m^{l \rightarrow r}$  is calculated by inter-  
polating between feature parameters of the two phoneme template  
using weight  $\alpha_m^{l \rightarrow r}$ . The result is used in Eq. 2 to modify both  
phoneme templates accordingly. Finally, the two modified spectra  
are interpolated using weight  $\beta_m^{l \rightarrow r}$ . Thus, the value of  $\alpha$  will cor-  
relate with changes in tongue position or place of articulation and  
height, while  $\beta$  will correlate with spectral balance or the man-  
ner of articulation. Because the values of the two weights evolve  
asynchronously, we call our method *asynchronous interpolation compression (AIC)*. We will now take a closer look at the various  
components of the transition function.

## 2.2. Estimating phoneme templates

Given an acoustic inventory, the initial task is to obtain phoneme  
templates  $\mathbf{P}$  for each phoneme of the database. In order to inter-  
polate between the templates, it is necessary to fix the frequencies of  
each component of  $\mathbf{P}$ .

The simplest approach to estimation is to select the appropriate  
speech locations manually. An automatic approach, taken in this  
paper, consists of collecting one or more *boundary frames* from all  
diphones that involve a specific phoneme. Boundary frames are  
frames that start or terminate a di-  
phone and are typically located  
in the center of phonemes. For example, to estimate the phoneme



**Fig. 1.** Log-magnitude spectrogram of concatenated boundary  
frames for the phoneme /i:/. The phoneme template  $\mathbf{P}^{i:}$ , defined  
as the centroid, is marked by the rectangle. The close acoustic  
match property of the acoustic inventory is evident by the similar-  
ity across frames.

template for /i:/, we collect left frames of diphones that begin with  
/i:/ as well as right frames of diphones that end in /i:/. Then, we  
set the template equal to the complex spectrum of the frame that  
has been identified as the centroid in the log-magnitude spectrum  
domain (see Fig. 1). This may not be the optimal solution; how-  
ever, we do not expect significant improvements if the acoustic  
inventory has the close acoustic match property.

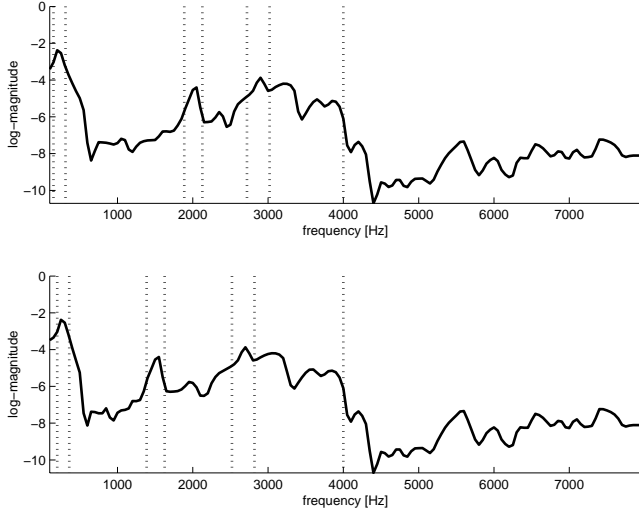
In certain diphones, acoustic events occur that are spectrally  
dissimilar to either phoneme template. An obvious example are  
diphthongs, such as in the di-  
phone /aI/→/@/ . In this case, unit-  
internal templates will have to be used and interpolated in se-  
quence; in our example the templates would be  $\mathbf{P}^{a(al)}$ ,  $\mathbf{P}^{I(al)}$ ,  
and  $\mathbf{P}^{@/}$ , where the first and last templates are (joinable) phoneme  
templates, but the middle template is used only in conjunction with  
its original di-  
phone (see also Sec. 4).

## 2.3. Speech model and modification function

The particular implementations of  $\mathbf{V}$  and  $\mathbf{M}$  depend on the choice  
of speech model. In this work, we use the concept of “spectral  
landmark” features. Since it is possible to pre-compute  $\mathbf{V}(\mathbf{P})$ ,  
we manually label the first three formant frequencies and band-  
widths of each phoneme, resulting in a sorted list of frequencies  
of the form  $\{F_1-B_1, F_1+B_1, F_2-B_2, F_2+B_2, F_3-B_3, F_3+B_3, C\}$ ,  
where  $C$  is a constant upper limit frequency beyond which  
no spectral modification is desired. When formants were not visi-  
ble, we used formant frequencies from locus theory [3].

During synthesis,  $\mathbf{M}(\mathbf{P}, \mathbf{v}_m^{l \rightarrow r})$  has access to both the original  
features  $\mathbf{V}(\mathbf{P})$  and the desired features  $\mathbf{v}_m^{l \rightarrow r}$ , using them to cal-  
culate a non-uniform sampling of the original frequency locations.  
To obtain the final spectrum, the magnitude and unwrapped phases  
are interpolated at the new frequencies. The original spectral bal-  
ance is mostly preserved by this frequency warping (see Fig. 2).

The advantage of this method is that of computational effi-  
ciency — important for synthesis on a computationally limited de-  
vice. However, the quality of modified speech spectra may be im-



**Fig. 2.** Log-magnitude spectrum of the phoneme template  $\mathbf{P}^{i:/}$  before (top) and after modification (bottom). The dashed lines represent vocal tract feature parameters  $V(\mathbf{P})$  and  $\mathbf{v}_m^{l \rightarrow r}$ , respectively.

proved upon by using parts of the spectral modification method suggested by Wouters et al [4].

#### 2.4. Estimating transition weights

To guarantee smooth concatenation at the diphone boundaries, we impose conditions  $\alpha_1 = \beta_1 = 0$  and  $\alpha_M = \beta_M = 1$ . It then follows that

$$\hat{\mathbf{D}}_{n,1}^{l \rightarrow r} = \mathbf{P}_n^l \quad \hat{\mathbf{D}}_{n,M}^{l \rightarrow r} = \mathbf{P}_n^r.$$

Thus the spectrum of two joinable diphones will be identical at their left and right boundaries, respectively, even if the acoustic units in the original inventory do not exhibit the close acoustic match property perfectly. Additionally, it is reasonable to require  $\alpha_{m-1} \leq \alpha_m$  and  $\beta_{m-1} \leq \beta_m$ .

Given phoneme templates  $\mathbf{P}$  and their feature parameters  $V(\mathbf{P})$ , we can now optimally estimate the transition weights  $\alpha$  and  $\beta$  to approximate a diphone  $\mathbf{D}$  by using a dynamic time warping (DTW) algorithm. Because  $\alpha$  and  $\beta$  can evolve independently, we construct the three-dimensional *transformation cube*

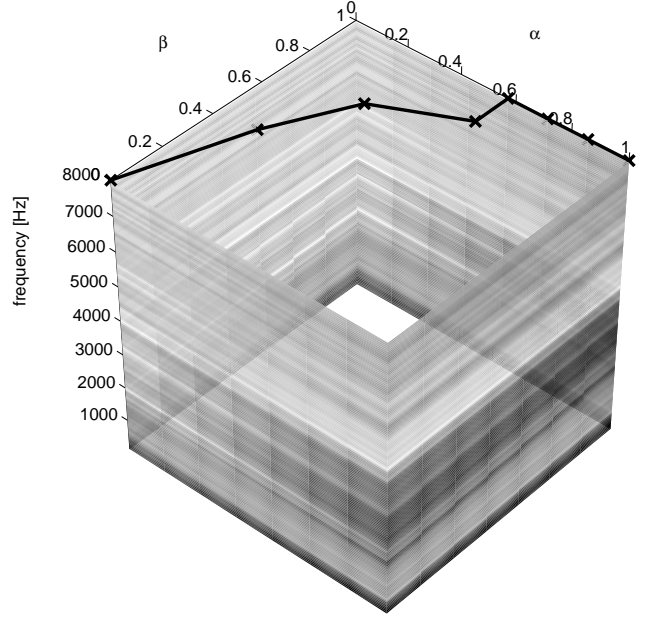
$$\mathbf{Q}_{n,a,b}^{l \rightarrow r} = \mathbf{T}(\mathbf{P}^l, \mathbf{P}^r, a/A, b/B)_n \quad (4)$$

$$a = 0, 1, \dots, A \quad b = 0, 1, \dots, B.$$

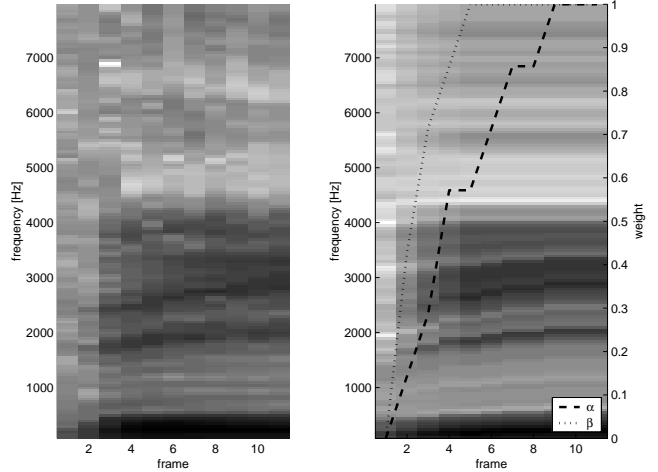
Now we align  $\mathbf{Q}$  to the template  $\mathbf{D}$  and the resulting DTW path is used to set values for  $\alpha$  and  $\beta$  (see Figures 3 and 4). The constraints above are realized by restricting the starting and ending points of the path, as well as using appropriate local path constraints. While we allowed for large jumps in  $\beta$  (allowing for naturally occurring sudden changes in spectral balance), we reduced the maximum rate of change for  $\alpha$  (since vocal tract feature parameters evolve relatively slower).

### 3. EVALUATION

Currently, our diphone inventory has not been re-labeled yet to take into account unit-internal templates. Therefore, we do not include diphthongs, plosives, or affricates in our evaluation. As reference, we use the ‘‘MWM’’ database of the OGIresLPC package



**Fig. 3.** Transformation cube  $\mathbf{Q}$  for the diphone  $\mathbf{D}^{D/- \rightarrow i:/}$  (only four sides are shown). Each column is a log-magnitude spectrum, the result of a particular choice of  $\alpha$  and  $\beta$  in Eq. 4. The left front side corresponds to a change in  $\alpha$  only, equivalent to a spectral cross-fade between  $\mathbf{P}^{D/}$  and  $\mathbf{P}^{i:/}$ . The right front side corresponds to a change in  $\beta$  only, which represents a continuous modification to  $\mathbf{P}^{i:/}$ . The DTW path, superimposed onto the top side, is constrained to start at the left corner and end at the right corner. The transition weight values for each frame of the diphone are associated with the crosses in the path.



**Fig. 4.** Log-magnitude spectrogram of the original diphone  $\mathbf{D}^{D/- \rightarrow i:/}$  (left) and compressed diphone  $\hat{\mathbf{D}}^{D/- \rightarrow i:/}$  with superimposed transition weights (right). The spectral frames of the compressed diphone correspond to columns in the transformation cube that were traced out by the DTW path. The distinct role of the two transitions weights is evident: While  $\alpha$  effects relatively slow formant movements,  $\beta$  causes a relatively fast change of spectral balance, as is commonly associated with fricative-vowel transitions.

Coder	SD	SEGSNR	SPD	Size	Ratio
Wave	-	-	5.9 (5.2)	6,519	1:1
G.722.1	4.1 (3.8)	16 (20)	5.9 (5.0)	638	1:10
G.722.2	7.3 (7.0)	4 (8)	4.9 (4.5)	197	1:33
AIC	8.7 (9.8)	1 (14)	0.0 (0.0)	57	1:114

**Table 1.** Mean spectral distortion (SD) and segmental signal-to-noise ratio (SEGSNR) across frames against the reference, and discontinuity distortion (SPD) between joinable frames of various coders. Variances are shown in parentheses. The size column refers to the total size of the compressed database in Kbytes, which includes other, uncompressed information such as pitch marks.

(<http://cslu.cse.ogi.edu/tts>), sampled at 16 kHz. We compare the reference against the ITU-T recommendation G.722.1 for speech coding at 24 Kbit/s [5], G.722.2 at 6.6 Kbit/s [6], and against our AIC approach.

A first evaluation involves frame-by-frame comparison of the reference with the coded speech, using standard spectral distortion (SD) and segmental signal-to-noise ratio (SEGSNR) measures. For each compression scheme, we encoded and decoded (i.e. synthesis without any prosody modification) the recorded diphones. For fricatives, the AIC coder randomizes the phases to avoid tonal artifacts. A second evaluation focuses on the spectral discontinuity (SPD) between the last and first frame of joinable diphones, using standard SD. By construction, spectral discontinuity is zero for the AIC method.

The results are shown in Table 1. While the size of AIC is an order of magnitude smaller than the G.722.2 coder, it also performs worse. However, it completely eliminated any spectral discontinuities that are present in the reference and other coders. Of course these objective measures, while informative, need to be validated with perceptual experiments. For example, the SEGSNR is sensitive to phase differences, yet it is unknown what the perceptual significance of these differences is.

Note that the AIC coding scheme is not fully size-optimized (e.g. there is no compression of phoneme templates, nor did we parameterize the weight trajectories or use parameter sharing between units of the same class; see Section 4).

#### 4. EXTENSIONS

The proposed method is the initial result of a longer-term project, which has the following goals:

**Application to longer units.** To apply the method to longer units, more unit-internal templates are needed. A key issue is to what degree these unit-internal templates can be shared by different units. If they cannot be shared, then the proposed method still provides compression because a template requires the same amount of storage as a single frame, and the total number of frames in a unit substantially exceeds the total number of phonemes in a unit; presumably, the latter is roughly equal to the number of templates required to generate a unit. If templates can be shared, larger savings will result. This extension will enable hand-held devices to benefit from progress in long-unit based TTS.

**Parameterized weight functions.** The weights shown in Figure 4 are typical in that they can be approximated by a sigmoid function characterized by a slope and a location parameter. Such a scheme can achieve a further substantial increase in compression.

**Parameter sharing.** The weight trajectories may be shared by units belonging to certain classes, such as Vowel→Nasal transi-

tions. Further savings would be obtained if the temporal patterns of all units in a given class have been normalized.

**Voice adaptation.** Changing the TTS voice can be accomplished by changing the phoneme templates while keeping the original transition weights, which appear to be highly speaker independent. This would permit the generation of personalized TTS systems using a minimum of recordings. This can be crucial for applications such as voice transformation for individuals with Dysarthria, where a coupled ASR-TTS system would be used to transform the Dysarthric speech into more intelligible speech. Individuals with Dysarthria are generally not able to create the amount of recordings needed for even a small diphone inventory.

**Enhanced spectral control.** Our hybrid Formant/Concatenative synthesizer has spectral control while maintaining naturalness. This enables development of synthesis modules that affect speaking rate, articulation effort, etc.

#### 5. CONCLUSION

A compression method has been proposed that takes advantage of properties of acoustic inventories, in particular the property that the end and start points of joinable units must be close in acoustic space. This immediately suggests using some type of interpolation scheme, but the asynchronous character of transitions between phonetic segments makes it impossible to use any straightforward interpolation method. The proposed method captures transition asynchrony by using a two-step process in which spectrum peak locations and spectral balance are asynchronously manipulated.

The method can be considered as a hybrid of traditional formant or rule based systems such as MITTalk and concatenative systems: It uses recorded speech, but generates transitions via (trained) interpolation functions.

#### 6. REFERENCES

- [1] A. Acero, "A mixed-excitation frequency domain model for time-scale pitch-scale modification of speech," in *Proc. of the Int. Conf. on Spoken Language Processing*, Sydney, Australia, December 1998.
- [2] B. Atal, "Efficient coding for lpc parameters by temporal decomposition," in *Proc. Int. Conf. Acoust., Speech and Signal Processing*, 1983, pp. 81–84.
- [3] D. Klatt, "Review of text-to-speech conversion for english," *J. Acoust. Soc. Am.*, vol. 82, no. 3, pp. 737–793, September 1987.
- [4] J. Wouters and M. Macon, "Control of spectral dynamics in concatenative speech synthesis," *IEEE Trans. Speech and Audio Proc.*, vol. 9, no. 1, pp. 30–38, January 2001.
- [5] ITU-T, *G.722.1 7 KHz Audio coding at 24 Kbit/s and 32 Kbit/s for hands free operation in systems with low frame loss*, September 1999.
- [6] ITU-T, *G.722.2 Wideband coding of speech at around 16 Kbit/s using adaptive multi-rate wideband (AMR-WB)*, January 2002.