



Department of Electrical Engineering
Technion – Israel Institute of Technology



Artificial Bandwidth Extension of Band Limited Speech Based on Vocal Tract Shape Estimation

Itai Katsir

MSc. Research

Supervised by Prof. David Malah
and Prof. Israel Cohen

30-Nov-11



Outline

- Introduction
- Methods of BWE
- Proposed BWE Algorithm
- Performance Evaluation
- Conclusion



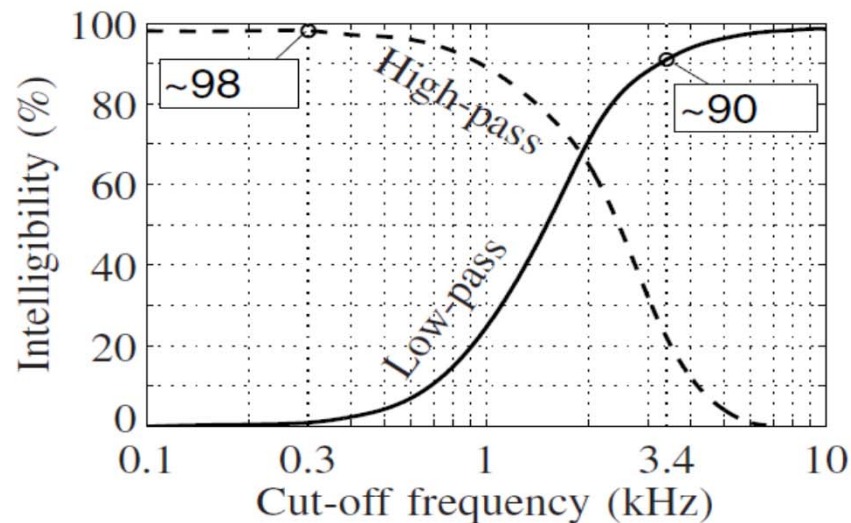
Outline

- **Introduction**
- Methods of BWE
- Proposed BWE Algorithm
- Performance Evaluation
- Conclusion

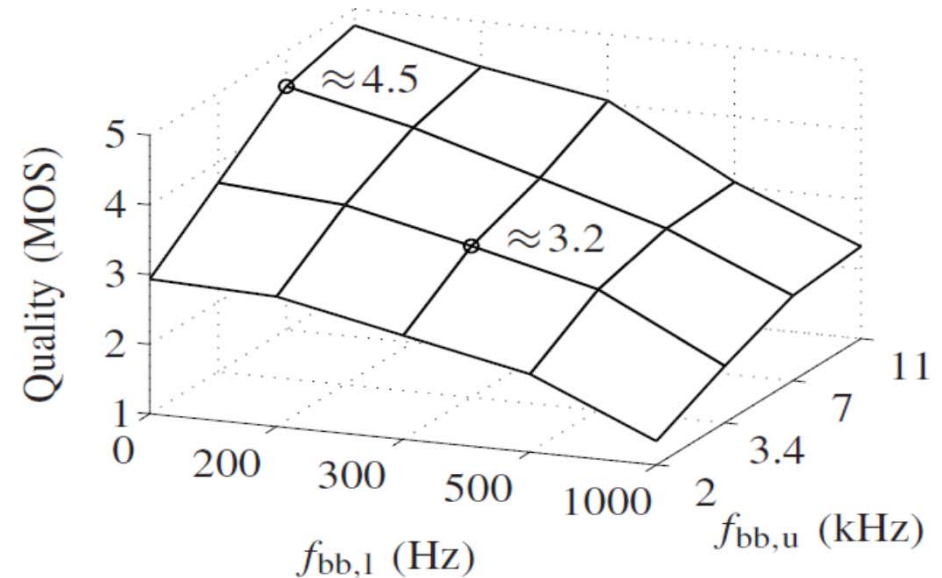


Introduction (1/8)

- **Fact** – Growing consumer demand for HD media → high quality speech communication.
- **Problem** – Today's analog telephone and PSTN limit the speech to narrowband (NB) frequency range of about 300-3400Hz → lower speech quality compared to wideband (WB) speech of range 50-7000Hz.



(a) Intellig. of meaningless syllables



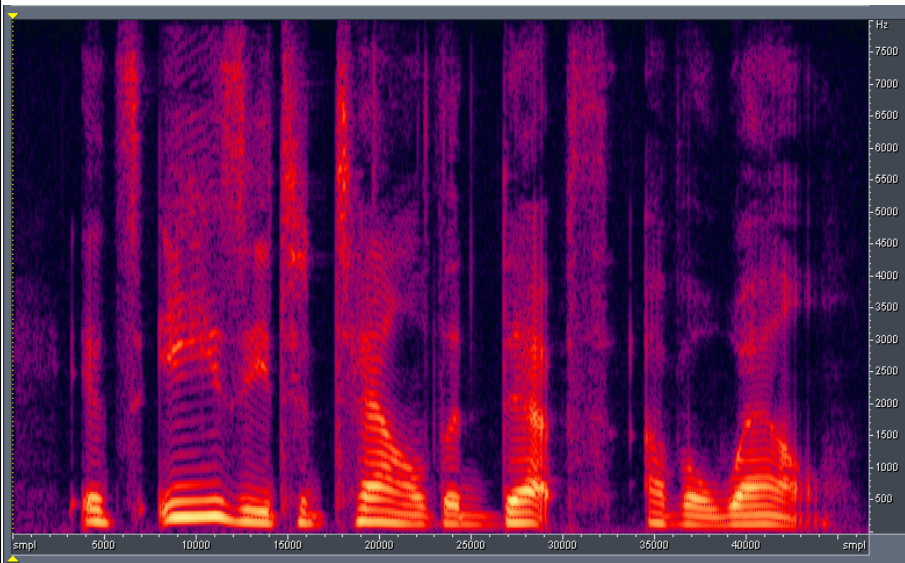
(b) Subjective speech quality

Introduction (2/8)

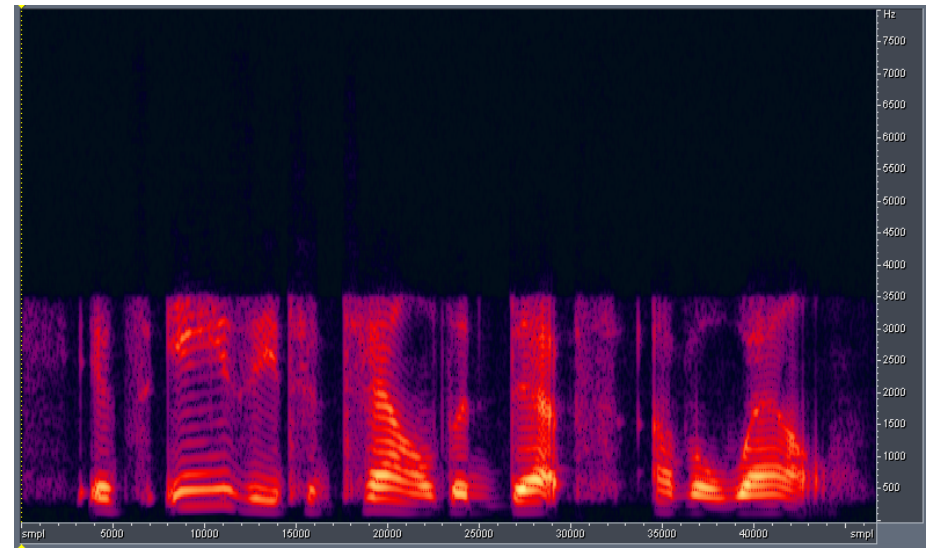


Wideband vs. Narrowband Speech

Wideband (WB)



Narrowband (NB)



Introduction (3/8)



Wideband vs. Narrowband Speech

- “Seed, Feed” – Spoken at different bandwidths.
- What do they sound like?
 - Reference: G.A. Miller and P.E. Nicely, “An analysis of perceptual confusions among some English consonants” Lincoln Laboratory, MIT, 1955 (*J. Acoust. Soc. Amer.* Vol. 27, pp. 338-352)

300-3400 [kHz]



50-5000 [kHz]



50-7000 [kHz]



- The same sentence – “Seed, Feed, Seed” in all cases !
- Spectrograms...

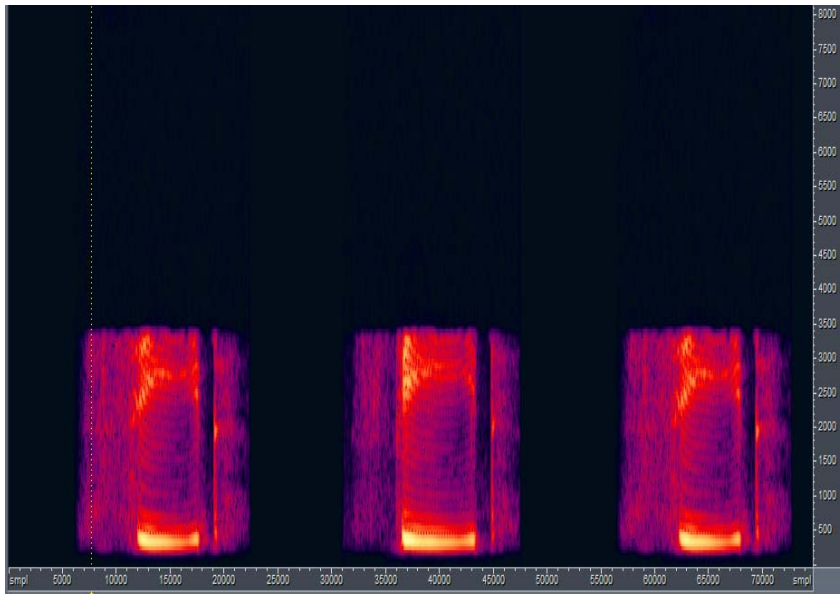
Introduction (4/8)



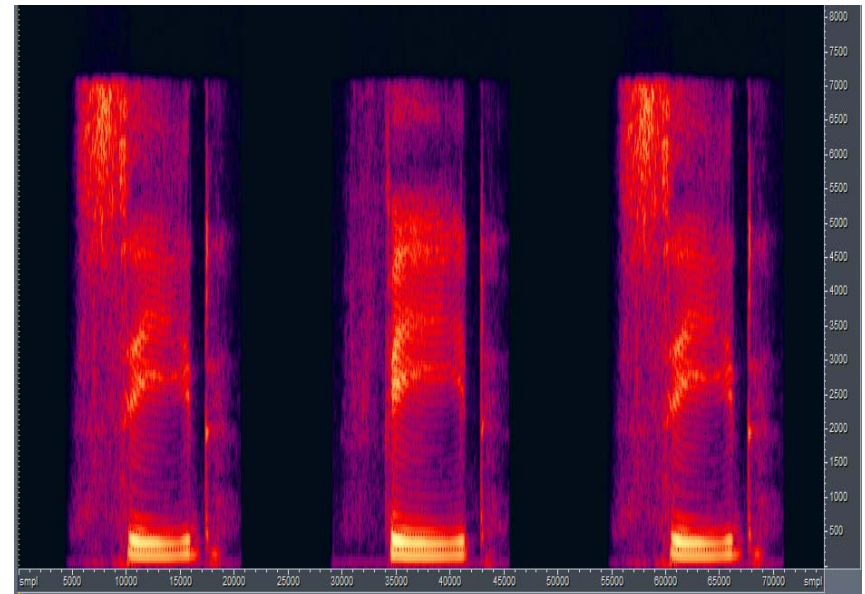
Wideband vs. Narrowband

- Spectrograms

300-3400 [kHz]



50-7000 [kHz]

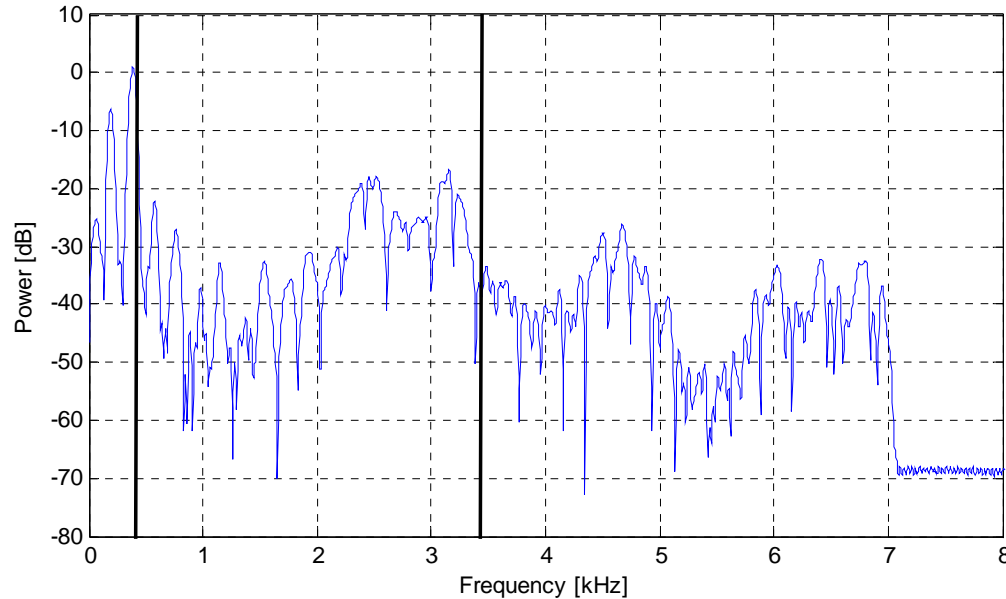


Introduction (5/8)



Voiced sounds

- Most of the energy is present in the low frequencies -> filtering out below 300 Hz affects **speech naturalness**.

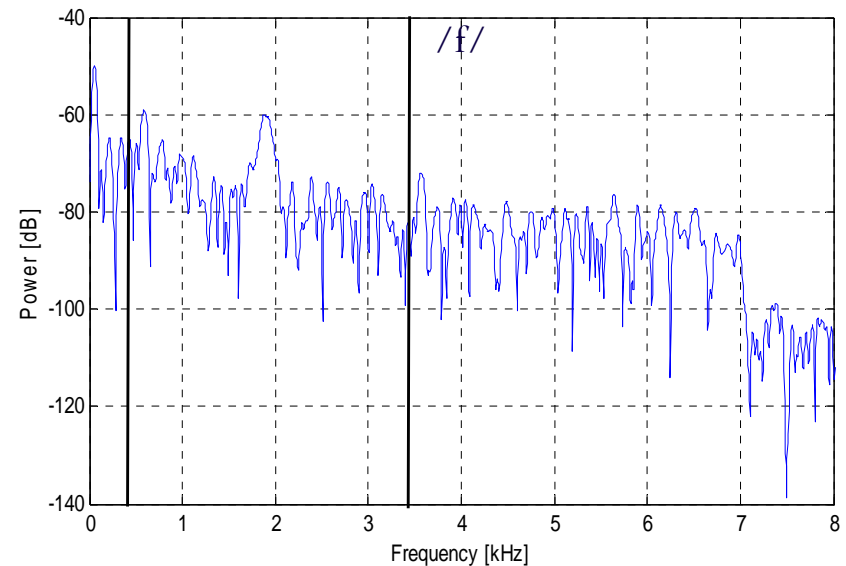
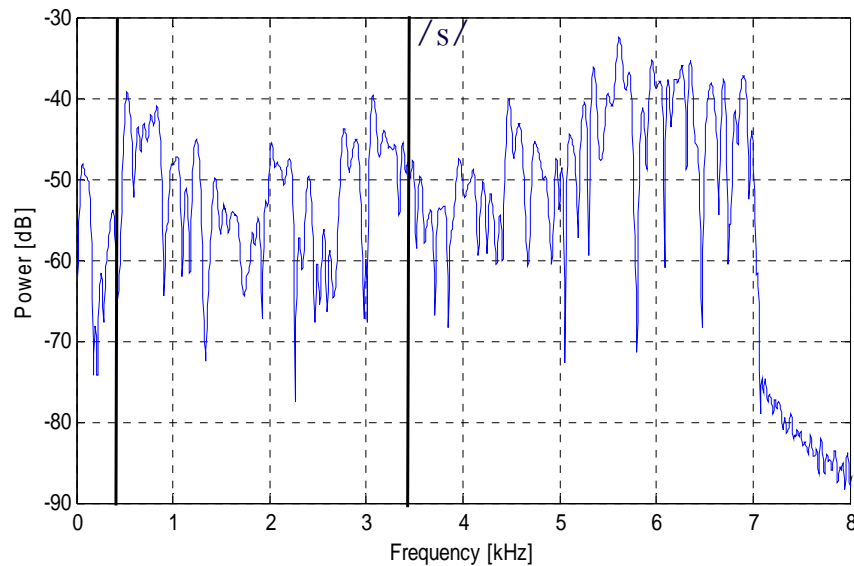


Introduction (6/8)



Unvoiced sounds

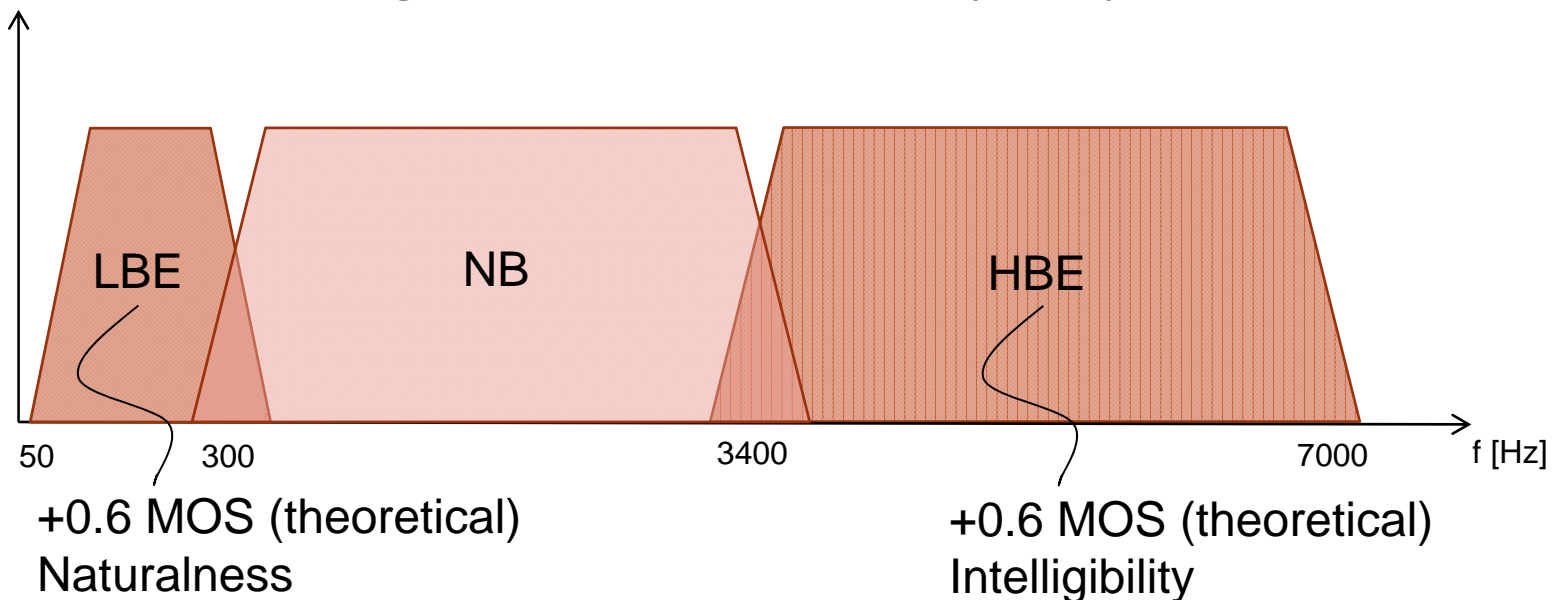
- Important energy above 3.4 kHz -> filtering out affects **speech intelligibility** such as differentiation between “s” and “f”.



Introduction (7/8)



- **Solution** – Artificially extend speech bandwidth to achieve speech quality enhancement and “listening effort” reduction.
 - 3.4-7kHz – Higher intelligibility and quality
 - 0-0.3kHz – Higher naturalness and quality

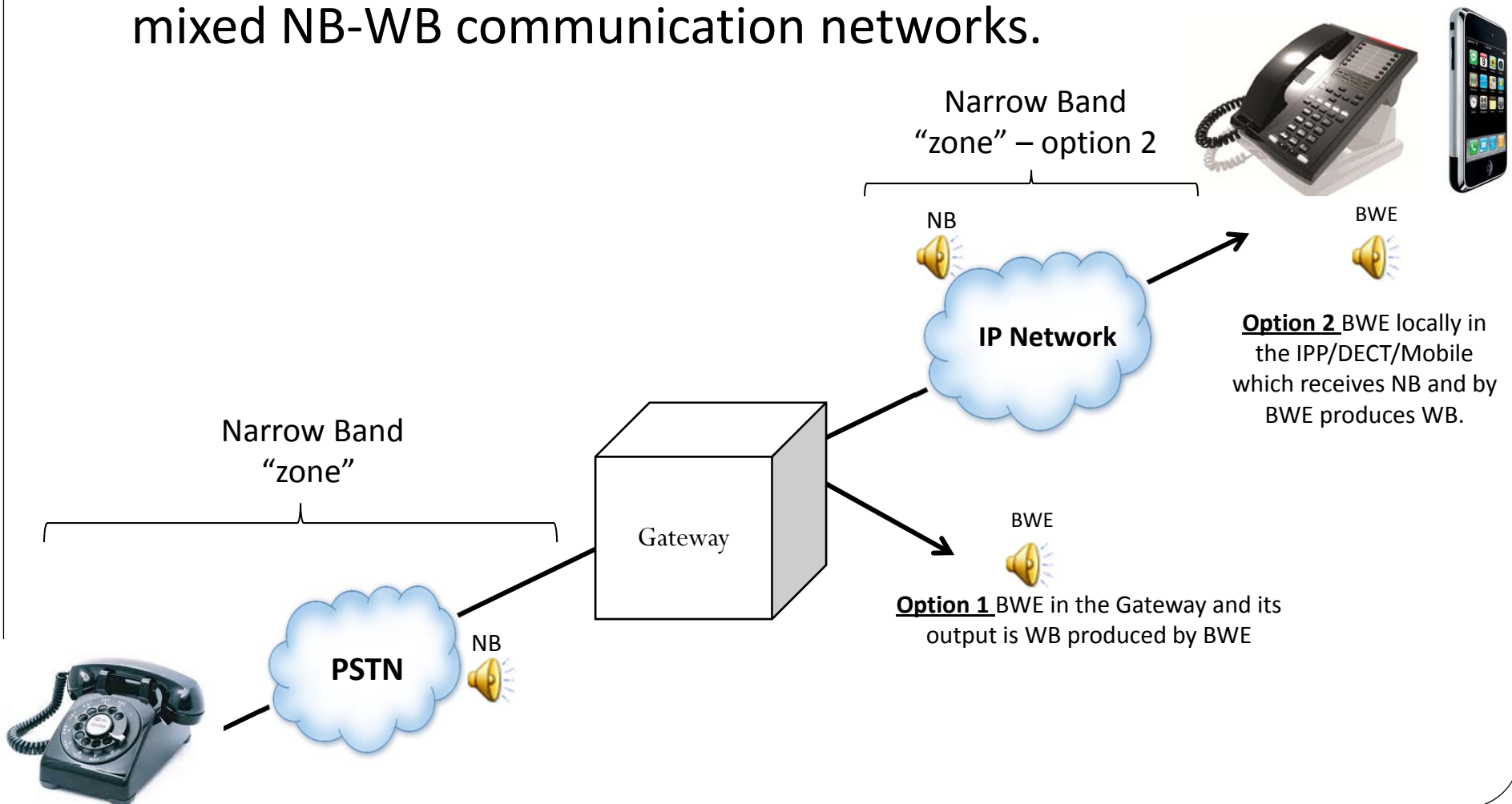


- +0.3 MOS - state of the art published BWE algorithms.

Introduction (8/8)



- **Application** – In the transition time to full WB communication networks, BWE can be used in mixed NB-WB communication networks.





Outline

- Introduction
- **Methods of BWE**
- Proposed BWE Algorithm
- Performance Evaluation
- Conclusion

Methods of BWE (1/8)



Model-less methods (Non parametric)

- Extend the bandwidth of the input narrowband speech signal directly, i.e., without any signal analysis.
- Utilize rather simple signal processing techniques (filtering and resampling), in time or frequency domain.

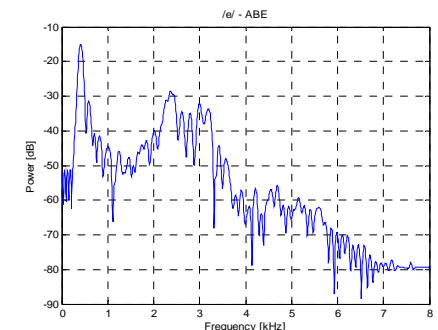
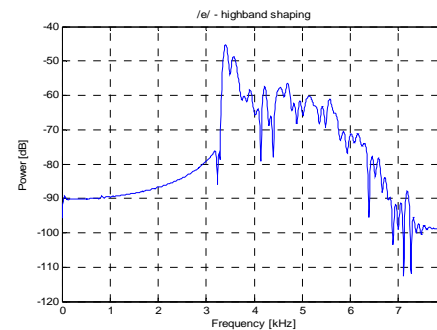
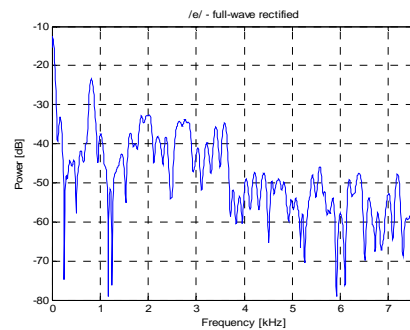
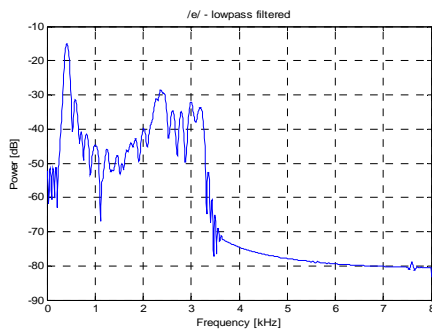
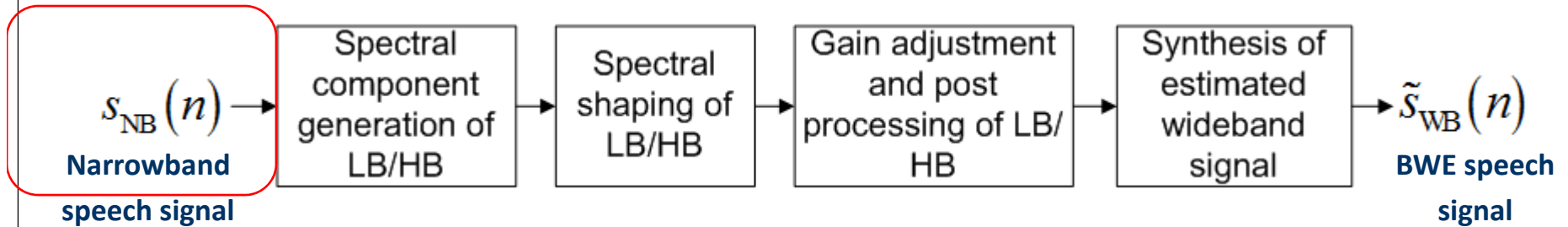
Model-based methods (Parametric)

- Estimate WB speech parameters from NB parameters.
- Rely on state-of-the-art signal processing techniques taken from pattern recognition, estimation theory, signal classification, etc.

Methods – Model-less BWE (2/8)



General scheme



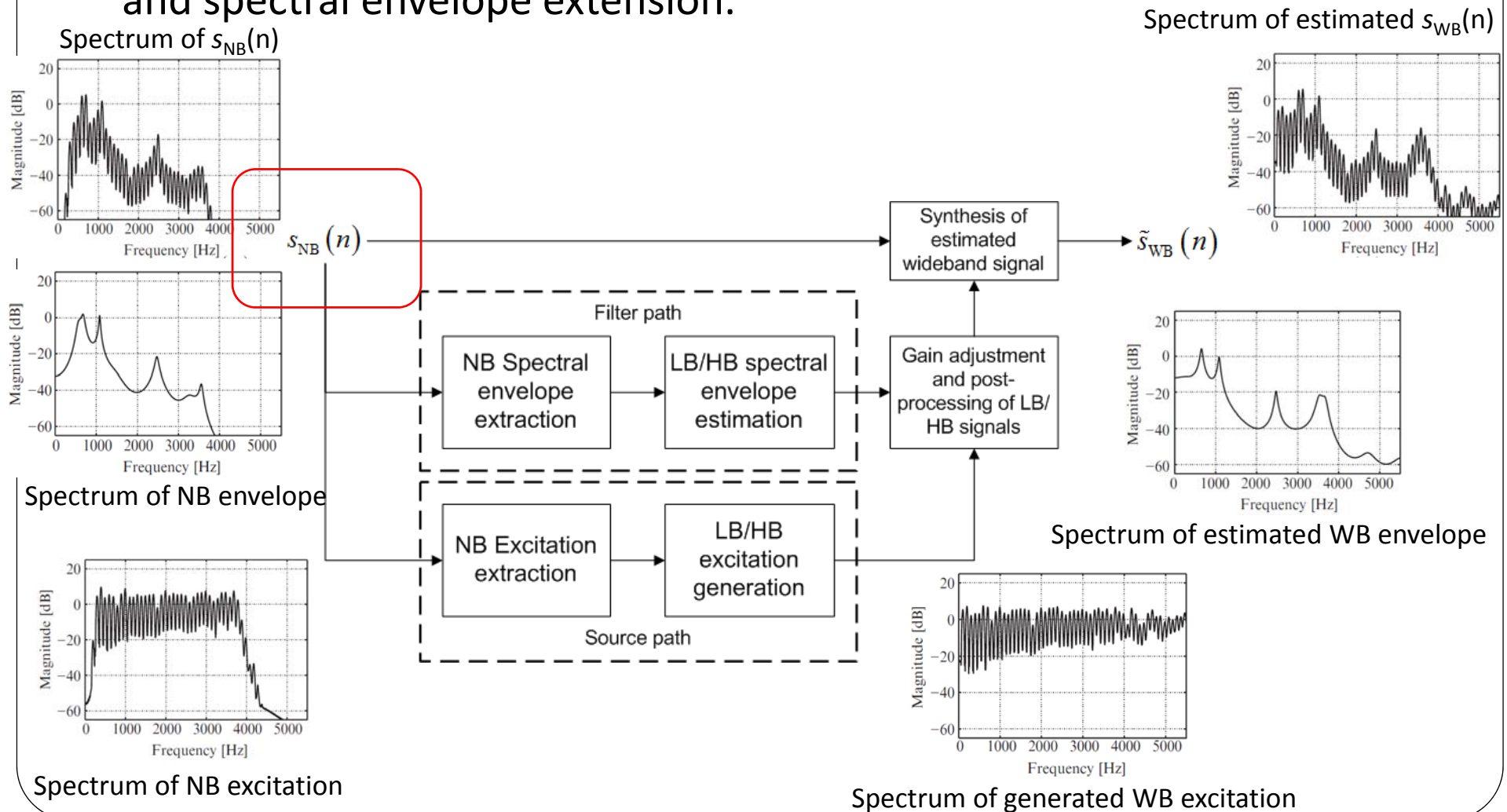
- Advantages
 - Low complexity as compared to parametric methods.
- Drawback
 - Lower quality as compared to parametric methods.

Methods – Model Based BWE (3/8)



General scheme

- Allows separate and independent algorithms for excitation extension and spectral envelope extension.

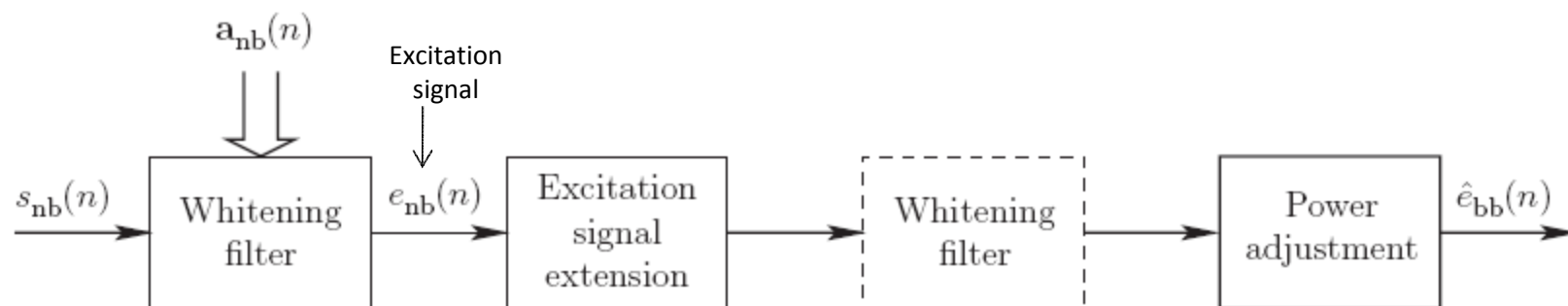


Methods – Model Based BWE (4/8)



Excitation BWE - General scheme

- Generation of WB excitation from NB excitation and NB scalar parameters.
- Various methods for wideband excitation generation.
 - Spectral shifting
 - Non-linear operators
 - Function generation

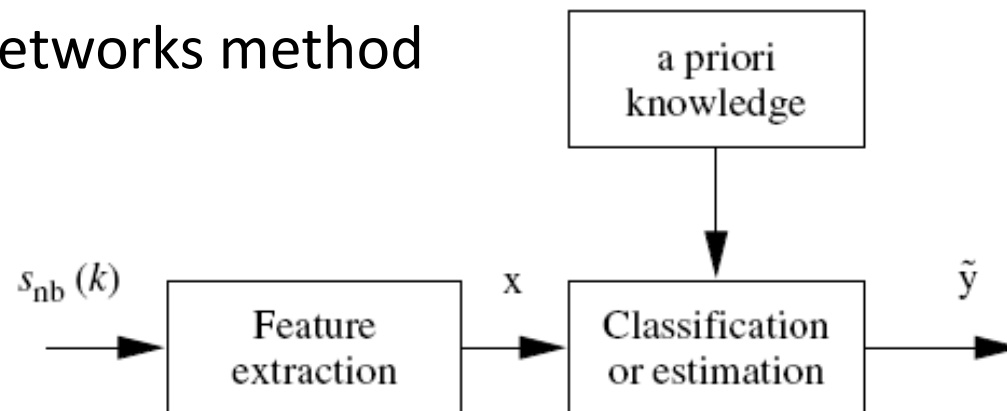


Methods – Model Based BWE (5/8)



Spectral envelope BWE - General scheme

- Estimate wideband spectral envelope \tilde{y} from narrowband signal features x (narrowband spectral envelope and voiced/unvoiced scalar features).
- Various methods for wideband spectral envelope estimation.
 - Linear or piece-wise linear mapping
 - Codebook mapping
 - Bayesian estimation based on Gaussian Mixture Models (GMMs) or Hidden Markov Models (HMMs)
 - Neural networks method



Methods – Existing Algorithms (6/8)



	Energy shift	Spectral shaping	Gain adjustment	Results
Nokia, Laaksonen [2008] (model-less)	Spectral folding	Frequency domain filtering using off-line trained control points	Normalized frequency domain filter	0.2 points improvement in MOS score for AMR coded signals
Ericsson, Gustafsson [2006] (model-based)	Spectral copy	Frequency domain filtering using estimated formants peaks	linear mapping based on off-line training	Evaluation of absence of distortion annoyance (using MOS scale)
Motorola, Ramabadran [2008] (model-based)	Full wave rectifier and noise generation	frequency domain filtering using 6, energy based, off-line trained spectral envelopes	mapping using off-line trained control constants	0.25 points improvement in MOS score
NTT – Quasi WB		Extending signal up to 6kHz		

Methods – Existing Algorithms (7/8)



	Energy shift	Spectral shaping	Gain adjustment	Results
Kornagel [2006] (model-based)	Spectral copy	LPC based codebook	narrowband energy equalization of estimated wideband signal to original signal	Used objective cepstral distance
Nilsson and Kleijn [2002] (model-based)	Repeated spectral folding of 2-3kHz excitation band and smooth transition to white noise	Estimate spectral envelope using GMM based mapping	Estimate gain using GMM based mapping	Used subjective degree of artifacts survey
Jax and Vary [2001] (model-based)	Modulation of narrowband excitation	Estimate spectral envelope and gain using HMM based mapping	Estimate gain using HMM based mapping	Used objective RMS LSD

Methods – Existing Algorithms (8/8)



Summary

- The major challenges are to estimate the high-band gain for unvoiced sounds and the spectral envelope for voiced sounds.
- BWE algorithms focus on tuning the extended bandwidth to minimize possible artifacts, by sophisticated spectral envelope and gain estimation techniques.
- Jumpy behavior of spectral envelope estimation, from one frame to another, in time, and from NB to HB shapes, in frequency, causes noticeable artifacts, like high frequency whistling.
- Power adjustment of synthesized signal to match narrow-band signal is crucial for artifact removal.
- Low-band frequencies estimation allows major improvement of BWE speech quality, but it is much more sensitive to erroneous estimation compared to high-band estimation.



Outline

- Introduction
- Methods of BWE
- **Proposed BWE Algorithm**
- Performance Evaluation
- Conclusion

Proposed BWE Algorithm (1/26)



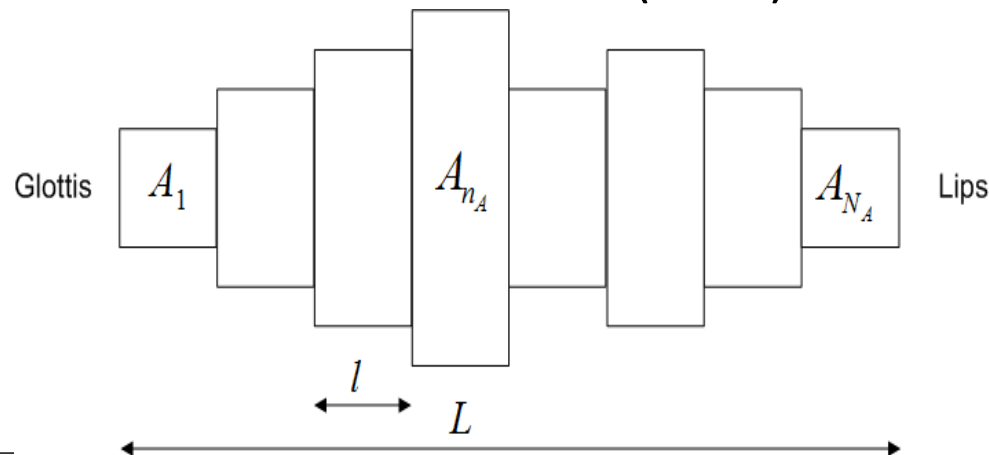
Vocal Tract Modeling

- The proposed algorithm tries to estimate the speaker physical vocal tract shape.
- Atal and Wakita showed the equivalence of acoustic tube model and the linear prediction (LP) model under certain conditions [Atal, 1970; Wakita, 1973].
- The M^{th} order filter transfer function derived through LP is equivalent to the transfer function of an acoustical tube made up of M equal length sections of variable areas.
- This is referred to as Vocal Tract Area Function (VTAF).

$$M = f_s \frac{2L}{c}$$

c - sound velocity

$$A_{n_A} = \frac{1 + r_{n_A}}{1 - r_{n_A}} A_{n_A+1}$$

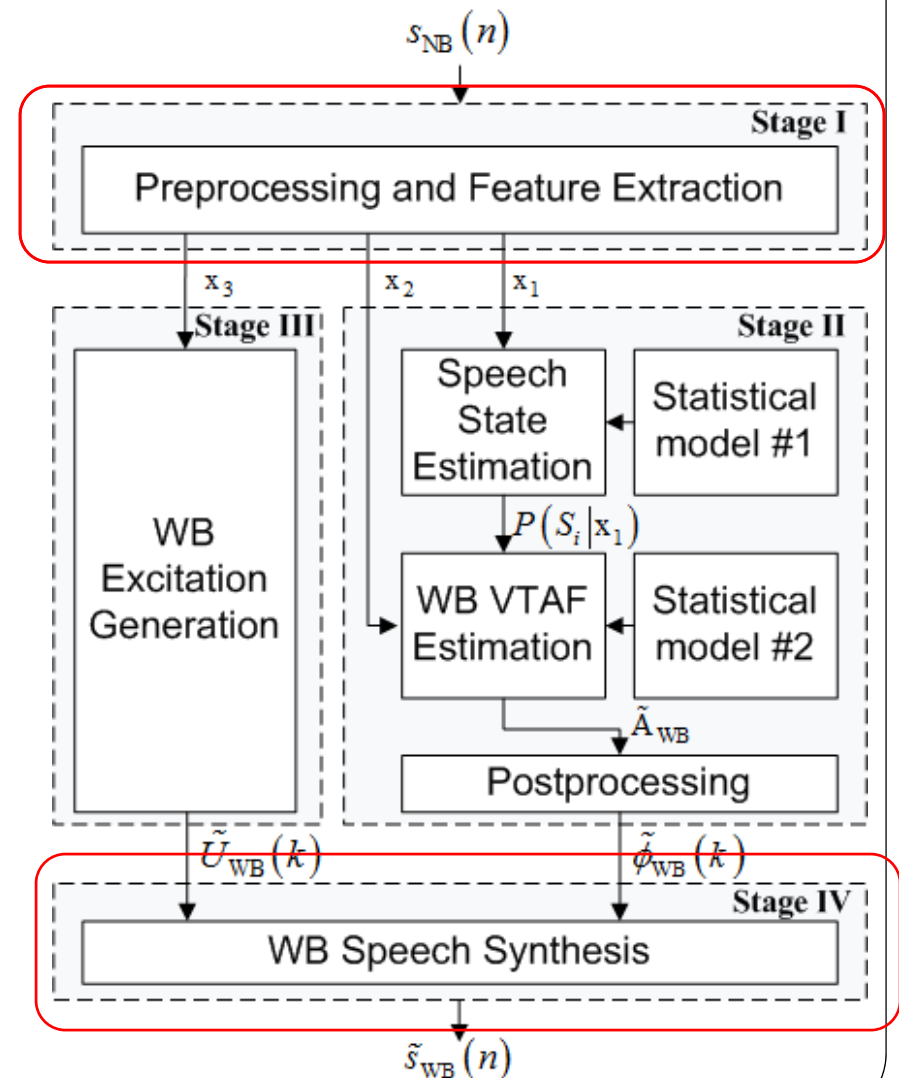


Proposed BWE Algorithm (2/26)



Algorithm stages:

- I. NB signal preprocessing and features extraction
- II. HB spectral envelope estimation and postprocessing
- III. WB excitation generation
- IV. Wideband signal synthesis

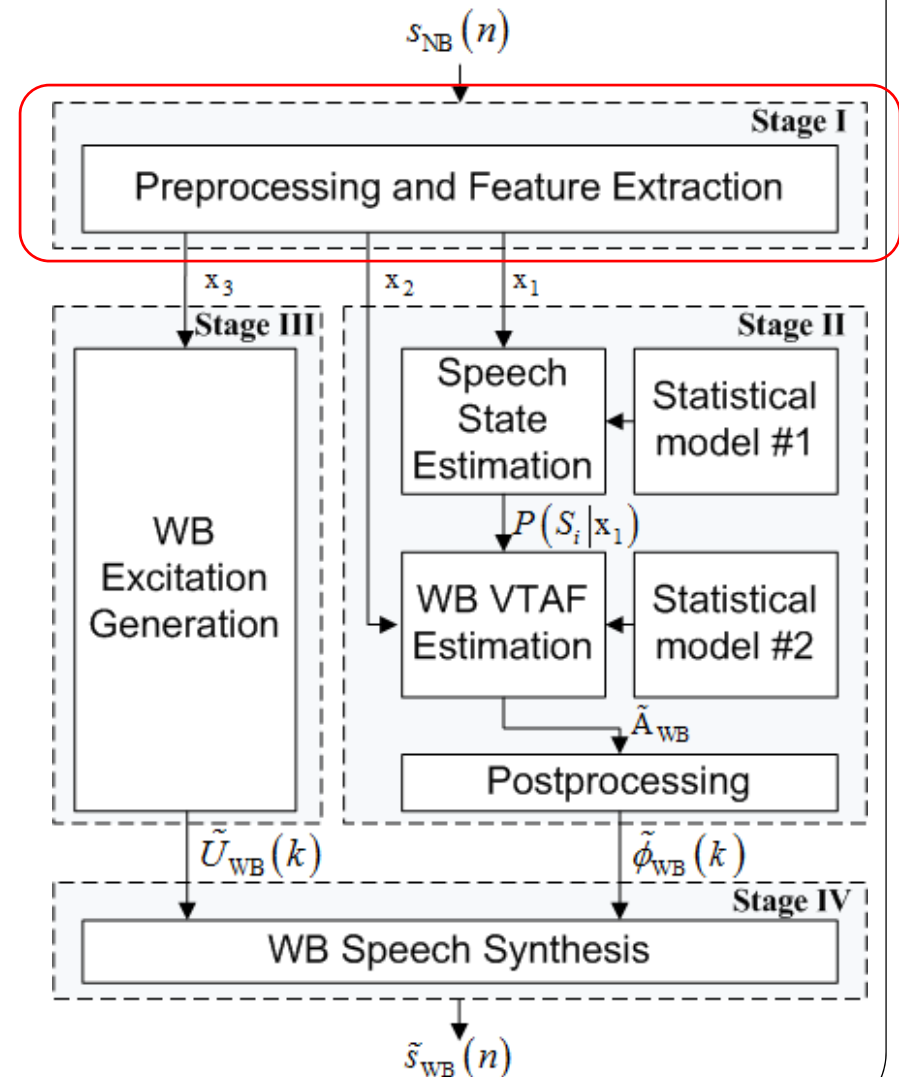


Proposed BWE Algorithm (2/26)



Algorithm stages:

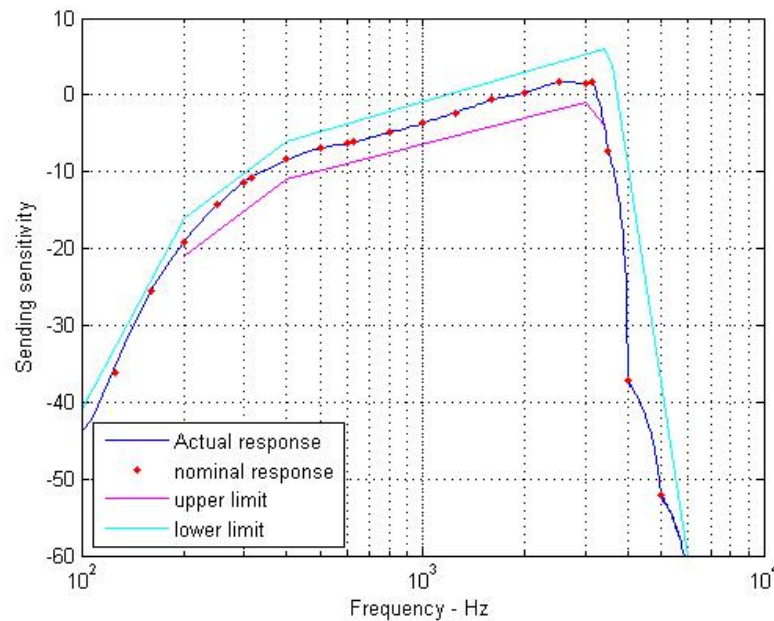
- I. NB signal preprocessing and features extraction
 - x_1 - NB spectral envelope estimation (MPEG spectral centroid, spectral flatness, spectral slope and normalized energy)
- II. WB excitation generation
 - x_2 - NB VTAF for WB VTAF estimation
- III. Wideband signal synthesis
 - x_3 - NB excitation for WB excitation generation



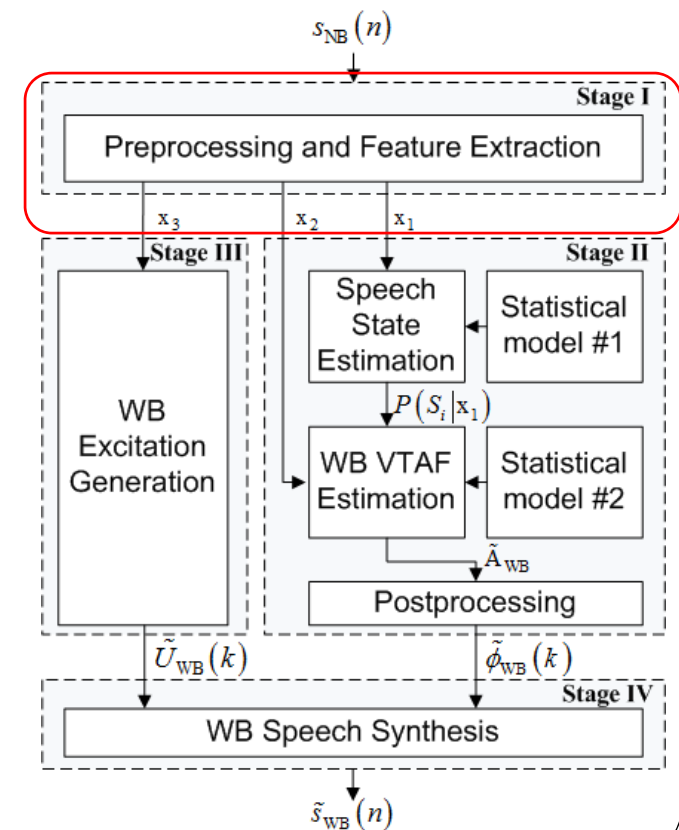
Preprocessing and Features Extraction (3/26)

Incentives

- Compensate for the IRS filter response.
- Extract features that allow good classification and estimation.



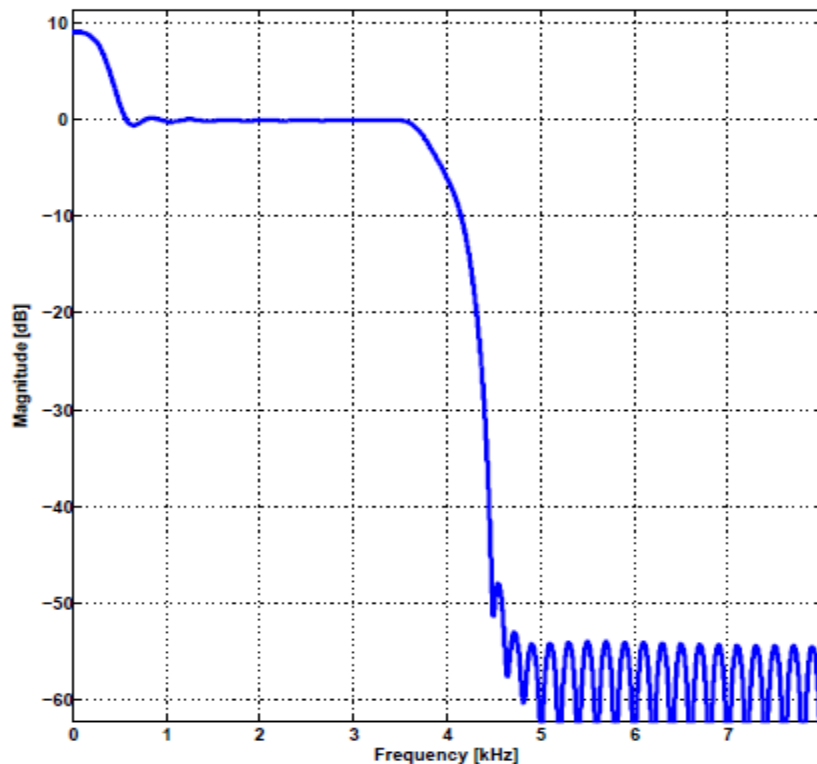
IRS send filter



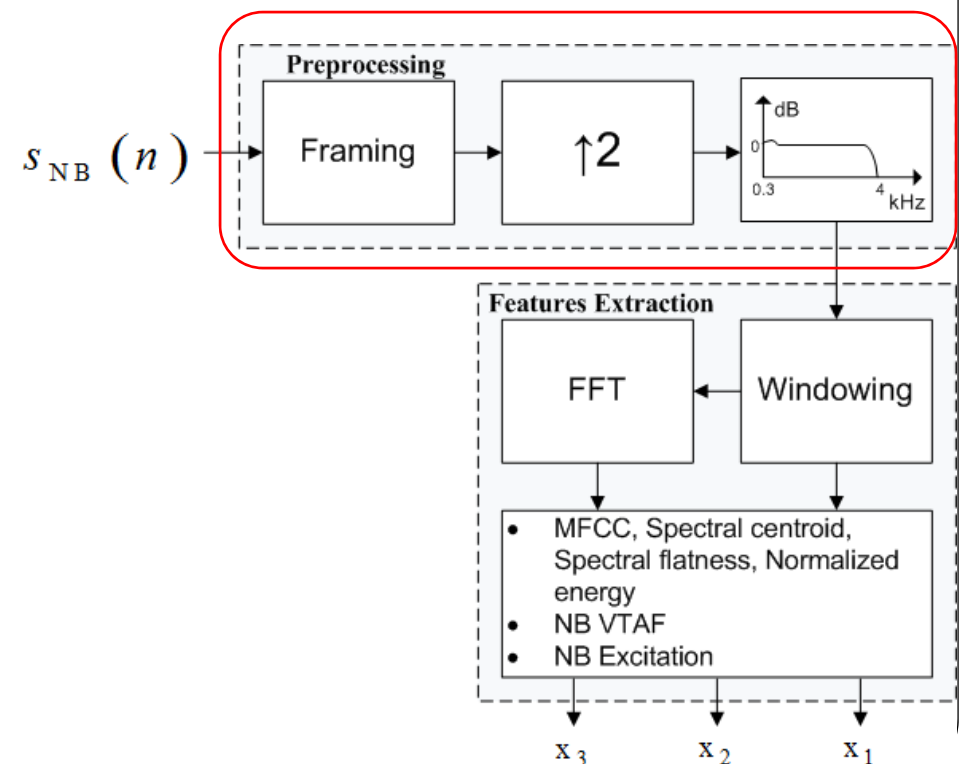
Preprocessing and Features Extraction (4/26)

Preprocessing

- NB signal interpolation.
- Equalization – 10dB boost to compensate for IRS filter at 300 Hz:



Equalizer filter



Preprocessing and Features Extraction (5/26)

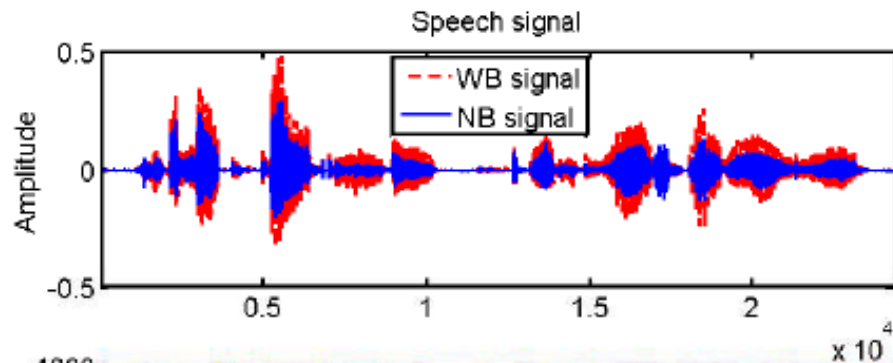
- Feature selection for BWE algorithm depends on the following considerations:
 - Computational complexity
 - Mutual information, $I(\mathbf{x}; \mathbf{y})$, between narrow- and high-band parameters.
 - Separability, $\zeta(\mathbf{x})$, the quality of particular feature set \mathbf{x} , for a classification problem. A higher separability value indicates a better suitability for classification.



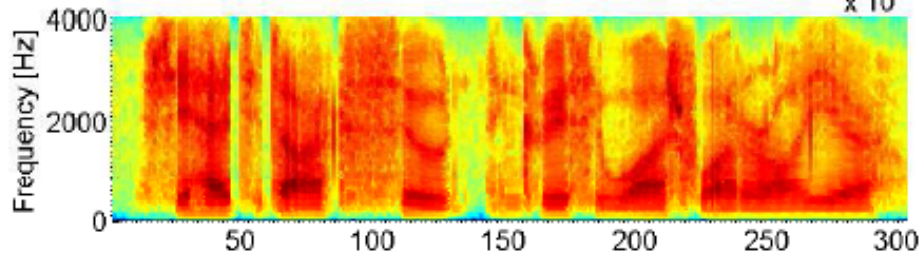
Preprocessing and Features Extraction (6/26)

Feature vector \mathbf{x}	dim \mathbf{x}	Towards high frequencies		Towards low frequencies	
		$I(\mathbf{x}; \mathbf{y})$ [bit/frame]	$\zeta(\mathbf{x})$ (16 classes)	$I(\mathbf{x}; \mathbf{y})$ [bit/frame]	$\zeta(\mathbf{x})$ (16 classes)
ACF	10	2.6089	1.6349	2.7530	2.3977
LPC	10	2.3054	1.5295	2.1100	1.7901
LSF	10	2.3597	1.5596	2.2125	2.5817
LPC-cepstrum	10	2.2401	1.4282	2.1778	2.3879
Cepstrum	10	2.3075	1.5483	1.9398	2.5473
MFCC	10	2.3325	2.2659	3.0771	6.6142
ACF (1)	1	0.7514	1.1237	0.7324	1.1065
ACF (pitch period)	1	0.4450	0.4058	0.5441	0.6745
Frame energy	1	0.9285	1.0756	1.3968	4.2328
Gradient index	1	0.8011	1.2520	0.5403	0.6983
Zero-crossing rate	1	0.7453	1.0795	0.7456	1.1685
Pitch period	1	0.2451	0.0530	0.4823	0.1122
Local kurtosis	1	0.2037	0.0225	0.2979	0.0809
Spectral centroid	1	0.7913	1.0179	0.6630	0.9276
Spectral flatness	1	0.4387	0.3538	0.4201	0.4648

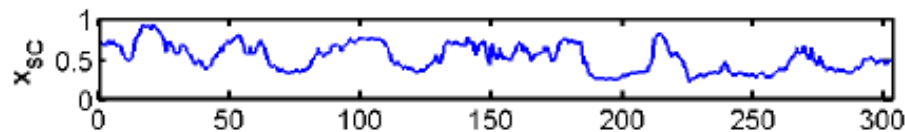
Preprocessing and Features Extraction (7/26)



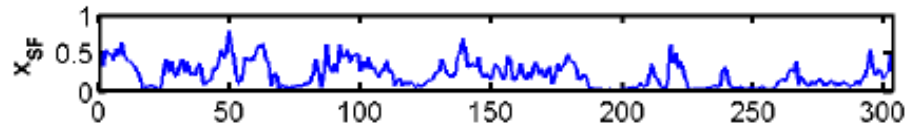
Speech signal



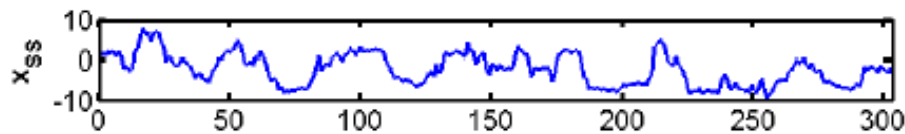
Spectrogram



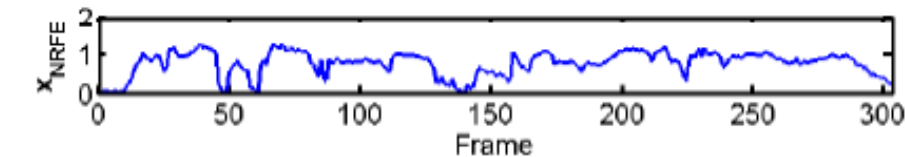
Spectral centroid



Spectral flatness



Spectral slope



Normalized relative
frame energy



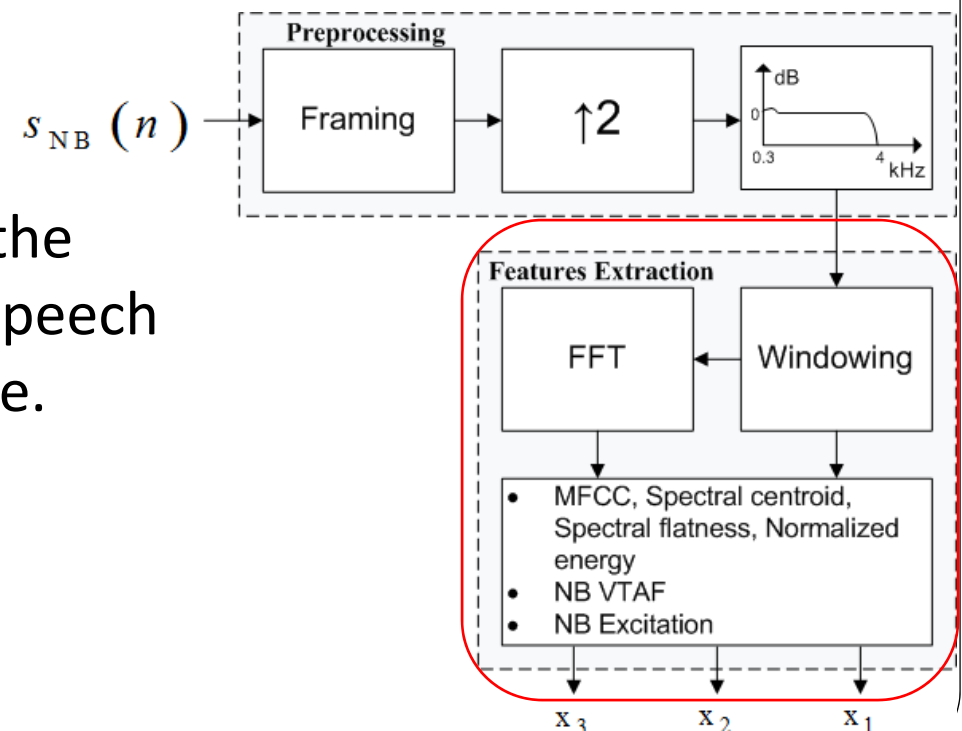
Preprocessing and Features Extraction (8/26)

Feature extraction

- x_1 – Frequency based features for phoneme estimation (MFCC, spectral centroid, spectral flatness, spectral slope and normalized energy).
- x_2 – NB VTAF calculated from the reflection coefficients

$$A_{n_A} = \frac{1 + r_{n_A}}{1 - r_{n_A}} A_{n_A+1}$$

- x_3 – NB excitation calculated in the frequency domain by dividing speech signal with the spectral envelope.

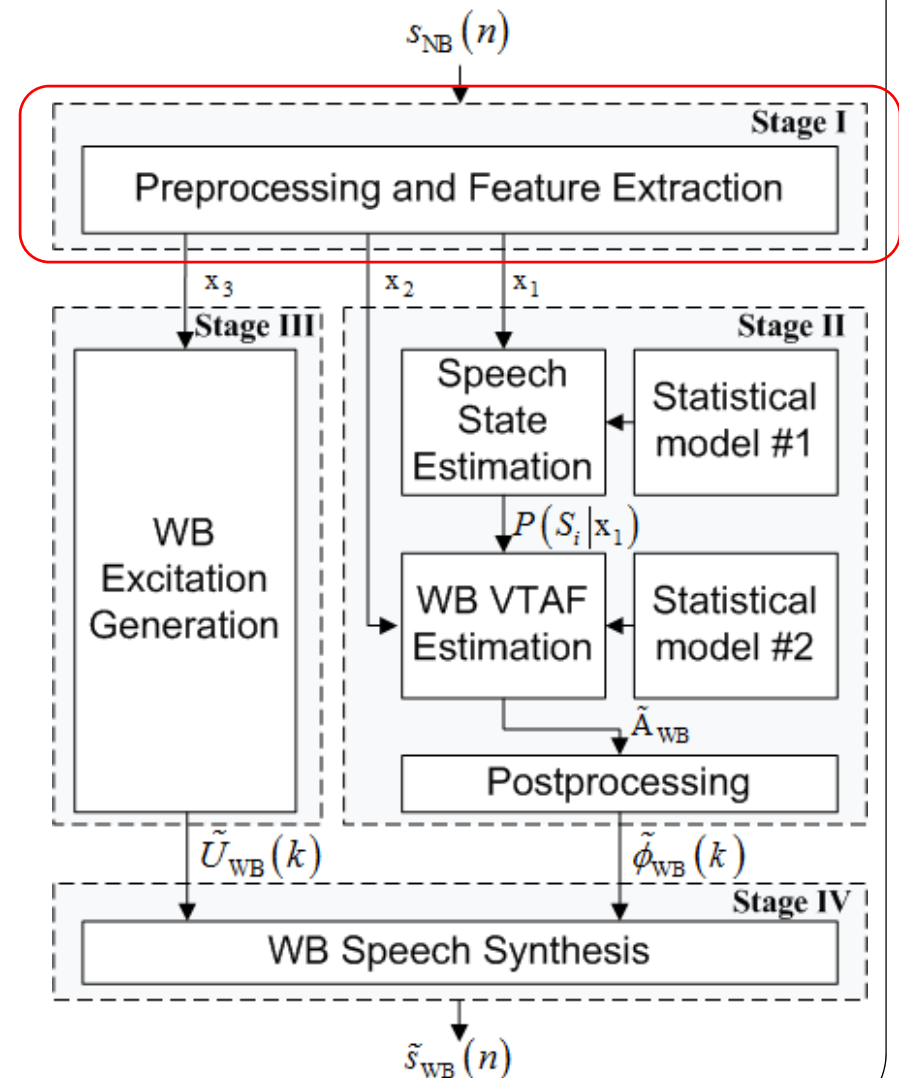


Proposed BWE Algorithm (9/26)



Algorithm stages:

- I. NB signal preprocessing and features extraction
- II. HB spectral envelope estimation and postprocessing
- III. WB excitation generation
- IV. Wideband signal synthesis

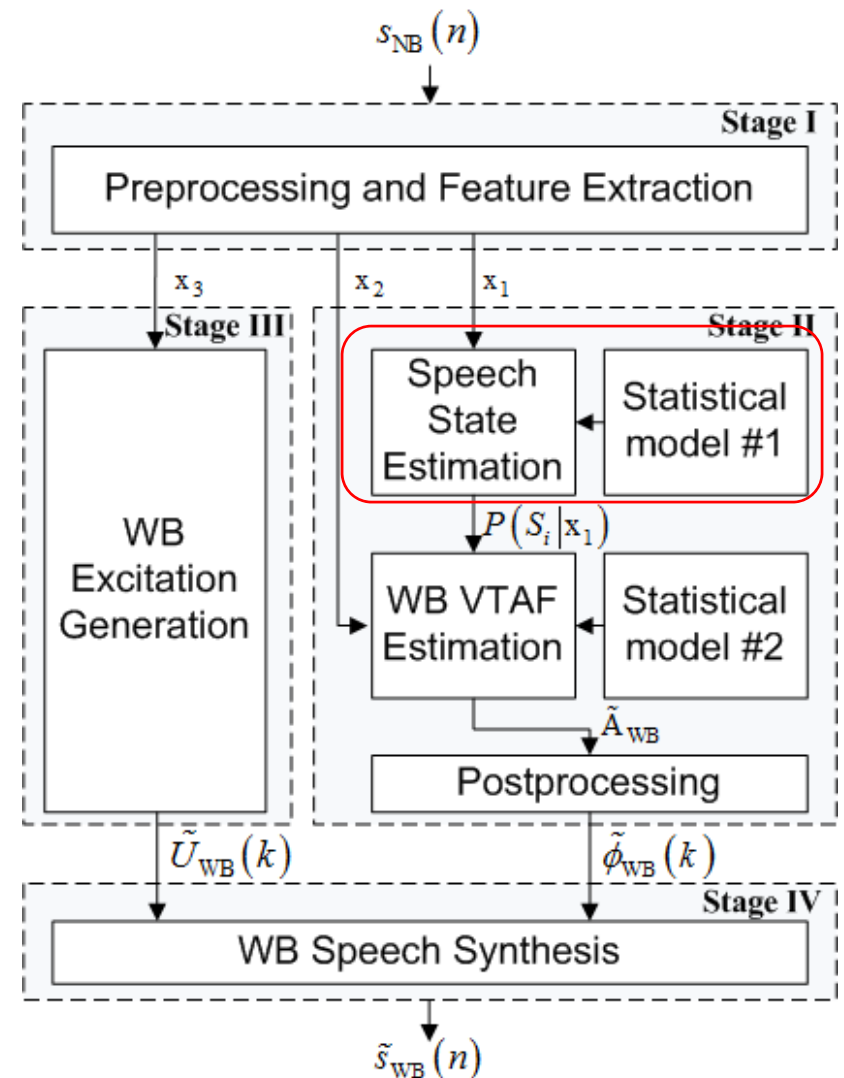


Proposed BWE Algorithm (9/26)



Algorithm stages:

- I. NB signal preprocessing and features extraction
- II. HB spectral envelope estimation and postprocessing
- III. WB excitation generation
- IV. Wideband signal synthesis



WB Spectral Envelope Estimation (10/26)

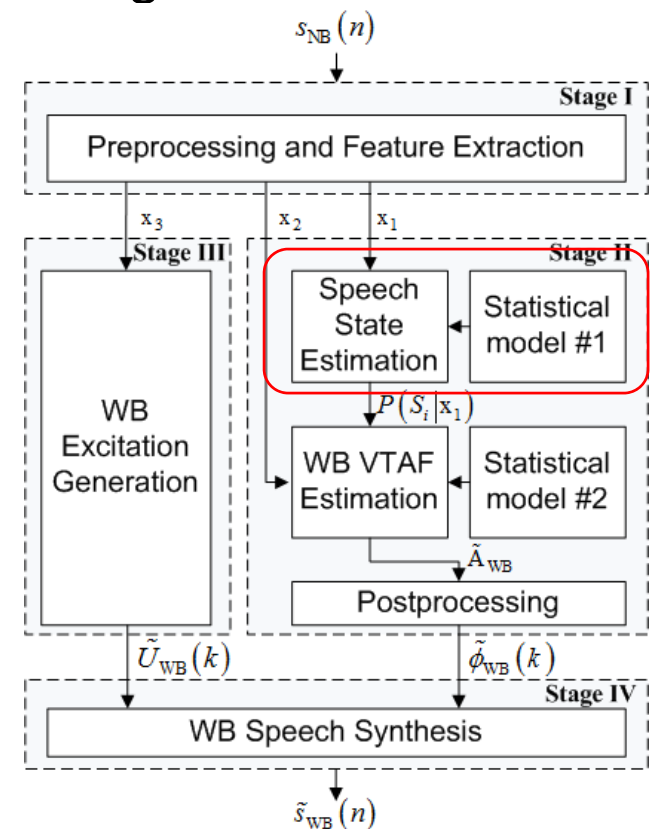
Phoneme estimation: using HMM to estimate each speech frame linguistic state.

- Off-line process: using TIMIT transcription to build HMM statistical model using phoneme based states. Calculating the following PDFs:

➤ $p(S_i)$ - Initial probability of each state.

➤ $p(S_i(m)|S_j(m-1))$ - Transition probability of the Markov chain from state j to state i .

➤ $p(x_1|S_i)$ - Observation probability for each state. Approximated by GMM parameters using the EM algorithm.



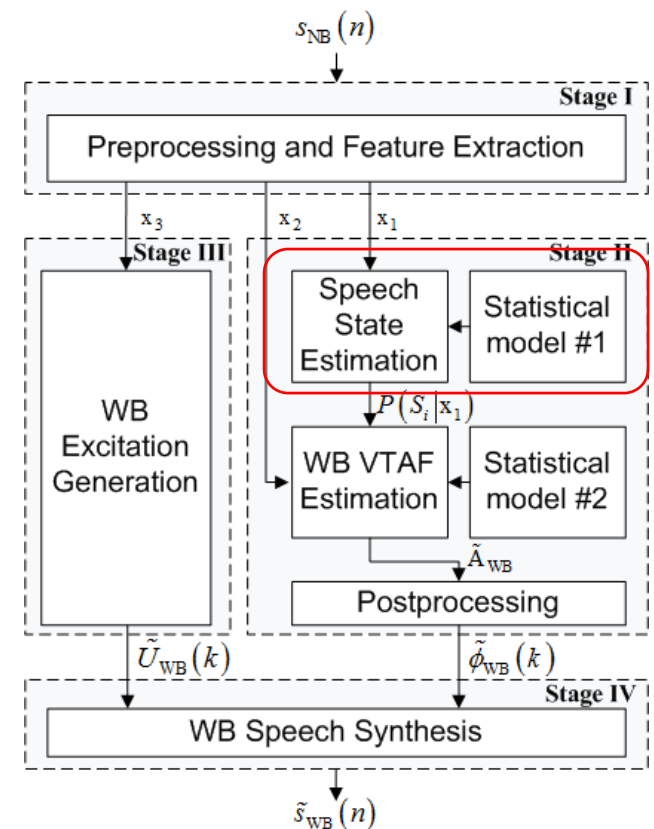
WB Spectral Envelope Estimation (11/26)

Phoneme estimation: using HMM to estimate each speech frame linguistic state.

- On-line process: making a decision on current frame state (phoneme) by maximizing the a-posteriori PDF:

$$p(S_i(m)|X_1(m)) = p(x_1(m)|S_i(m)).$$

$$\sum_{j=1}^{N_s} p(S_i(m)|S_j(m-1)) p(S_j(m-1)|X_1(m-1))$$



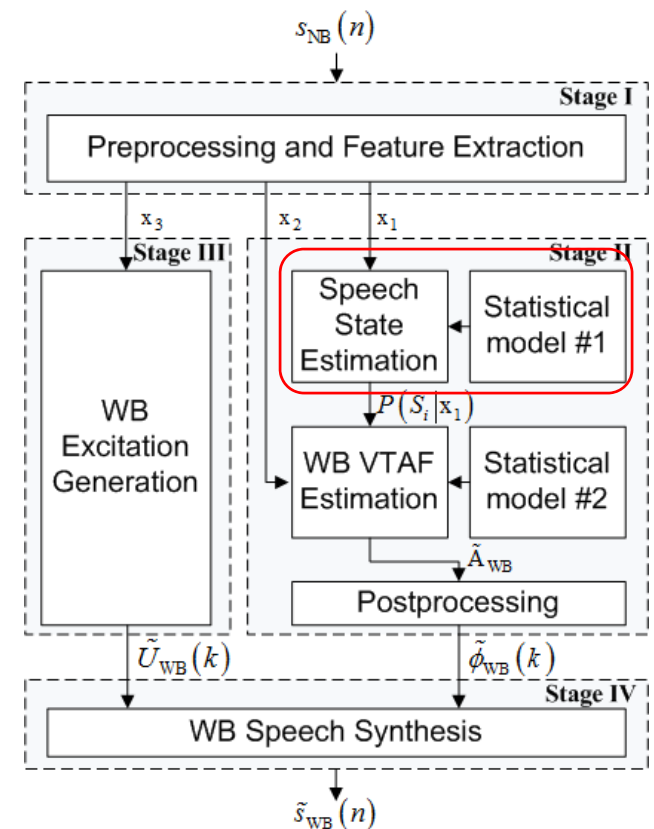
WB Spectral Envelope Estimation (12/26)

Estimate WB VTAF: using codebook matching of calculated NB VTAF to WB VTAF.

- Off-line process: for each speech state, clustering of N_{CB} WB VTAF using vector quantization of speech frames training set.
- On-line process: finding closest WB VTAF codeword to extracted NB VTAF using Euclidean distance.

$$\tilde{A}_{WB}^{S_i} = A_{WB}^{S_i} (j^{opt})$$

$$j^{opt} = \arg \min_{j=1}^{N_{CB}} \left\| \log(A_{NB}) - \log(A_{WB}^{S_i}(j)) \right\|_2^2$$

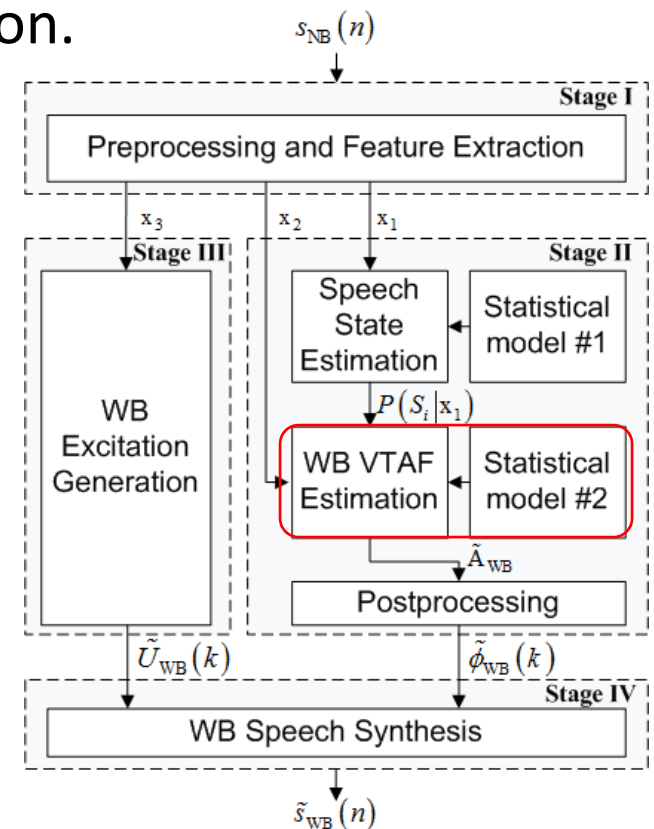


WB Spectral Envelope Estimation (13/26)

Postprocessing

- Reduce artifacts due to erroneous estimation in first two steps.
- Reduce artifacts due to erroneous state estimation by using N_{best} highest probability states for VTAF estimation.

$$\tilde{\mathbf{A}}_{\text{WB}} = \mathbf{C} \cdot \left(p_1 \cdot \tilde{\mathbf{A}}_{\text{WB}}^{S_{i_1}} + \dots + p_{N_{\text{best}}} \cdot \tilde{\mathbf{A}}_{\text{WB}}^{S_{i_{N_{\text{best}}}}} \right)$$

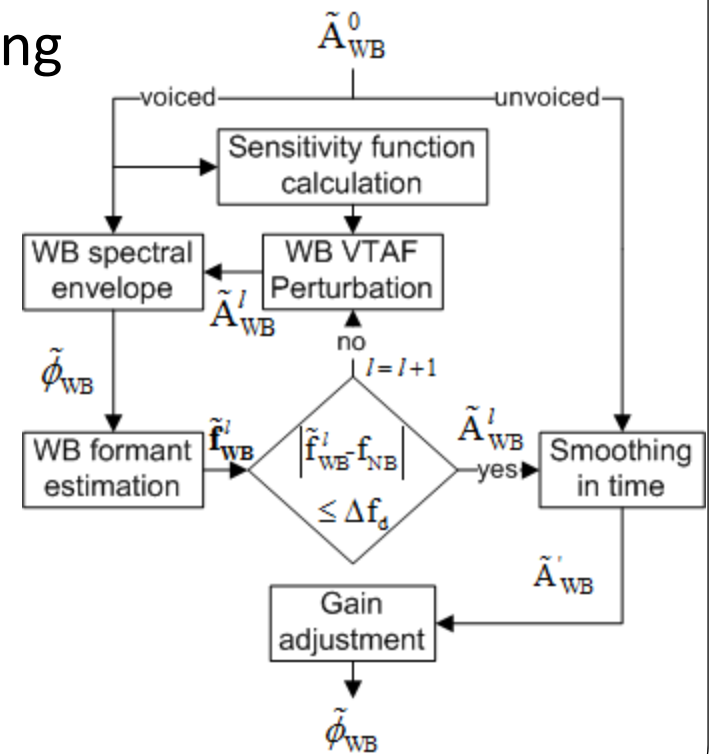


WB Spectral Envelope Estimation (14/26)

Postprocessing

- Reduce artifacts due to erroneous WB VTAF estimation.
- Estimated WB envelope fit to NB envelope by tuning formant frequencies of estimated WB VTAF to allow better gain adjustment to NB envelope. Iterative tuning by VTAF perturbation.
- Iterative VTAF perturbation based on the sensitivity function:

$$\frac{\Delta f_{n_f}}{f_{n_f}} = \sum_{n_A}^{N_A} S_{n_f, n_A} \frac{\Delta A_{n_A}}{A_{n_A}}$$



WB Spectral Envelope Estimation (15/26)

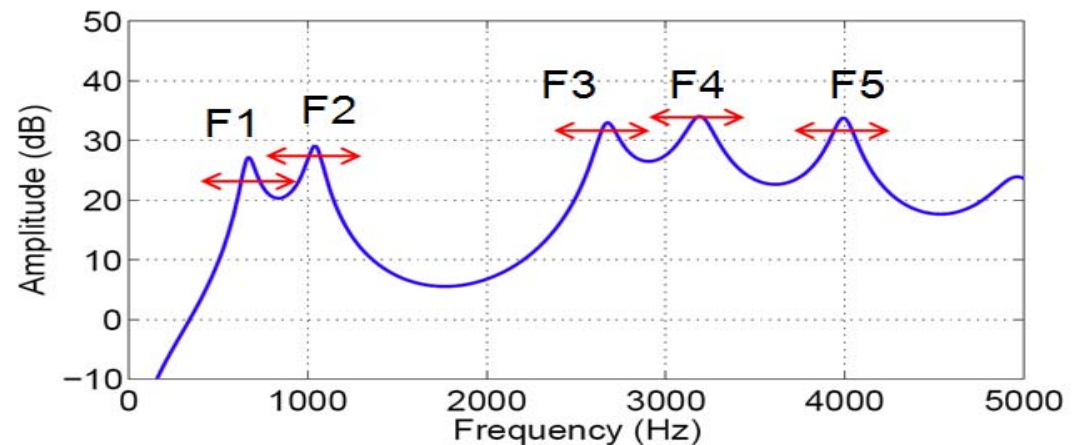
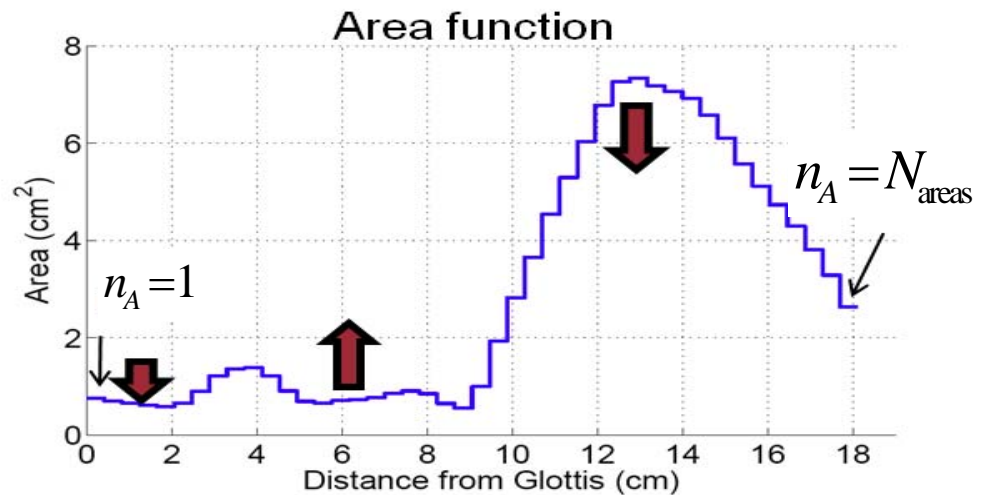
Vocal Tract Sensitivity Function

Relate changes in area to changes in formant frequencies

$$\frac{\Delta f_{n_f}}{f_{n_f}} \Leftrightarrow \frac{\Delta A_{n_A}}{A_{n_A}}$$

Sensitivity function

$$\frac{\Delta f_{n_f}}{f_{n_f}} = \sum_{n_A}^{N_A} S_{n_f, n_A} \frac{\Delta A_{n_A}}{A_{n_A}}$$

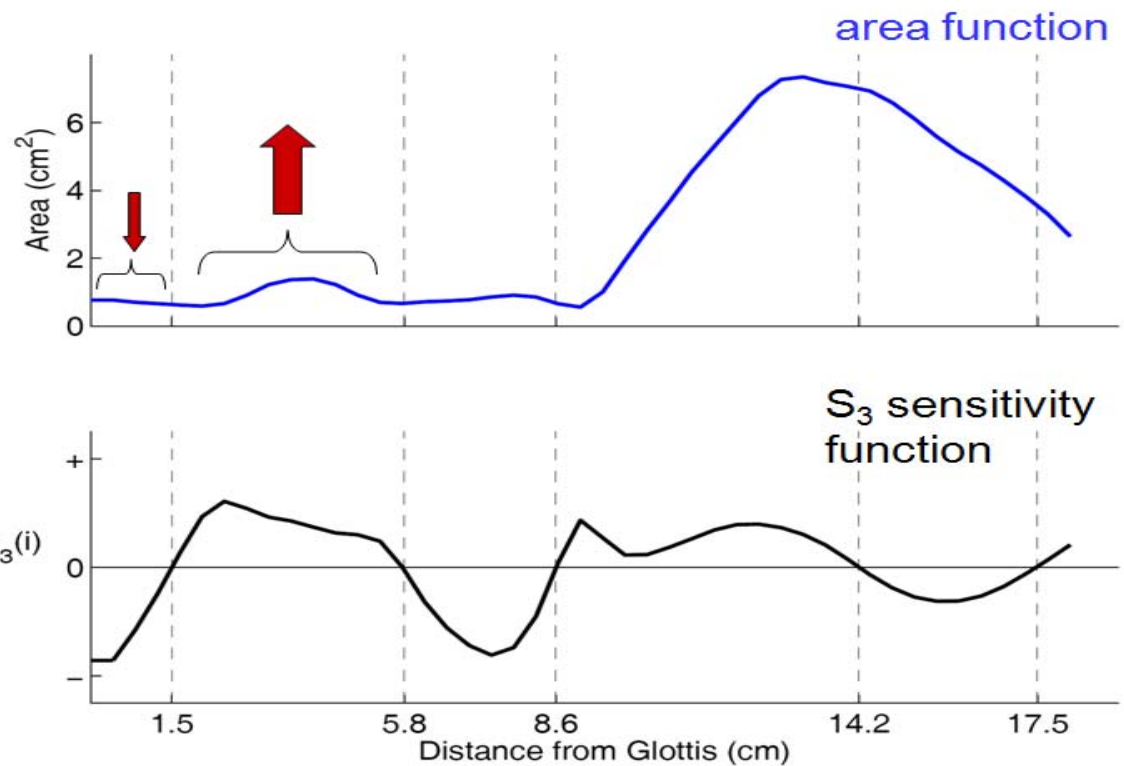
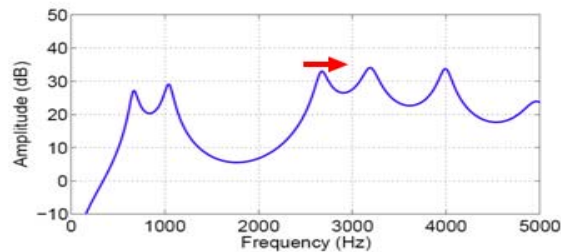


WB Spectral Envelope Estimation (16/26)

Vocal Tract Sensitivity Function

$$\frac{\Delta f_3}{f_3} = \sum_{n_A}^{N_A} S_{3,n_A} \frac{\Delta A_{n_A}}{A_{n_A}}$$

$\Delta a(i)$	$S_3(i)$	ΔF_3
+	+	+
+	-	-
-	+	-
-	-	+



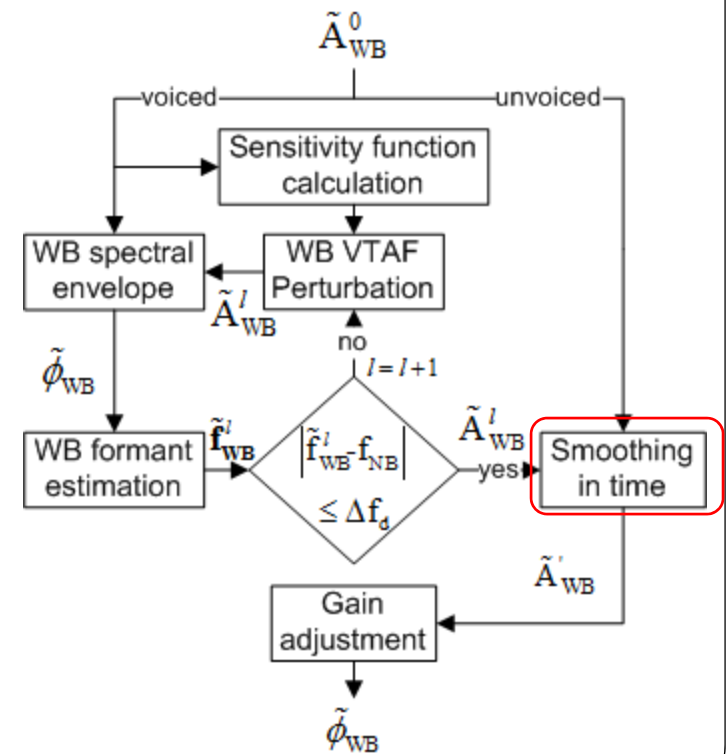
WB Spectral Envelope Estimation (17/26)

Postprocessing

- Stopping condition for iterative process is formant frequencies difference.
- Smoothing in time of tuned estimated VTAF.

$$\tilde{A}'_{\text{WB}}(m) = \beta \cdot \tilde{A}'_{\text{WB}}(m-1) + (1-\beta) \cdot \tilde{A}_{\text{WB}}(m)$$

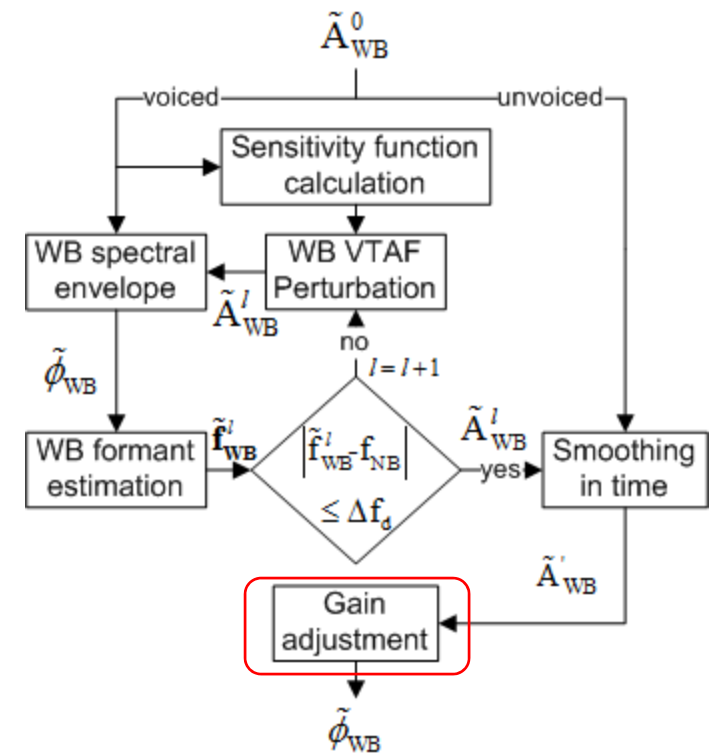
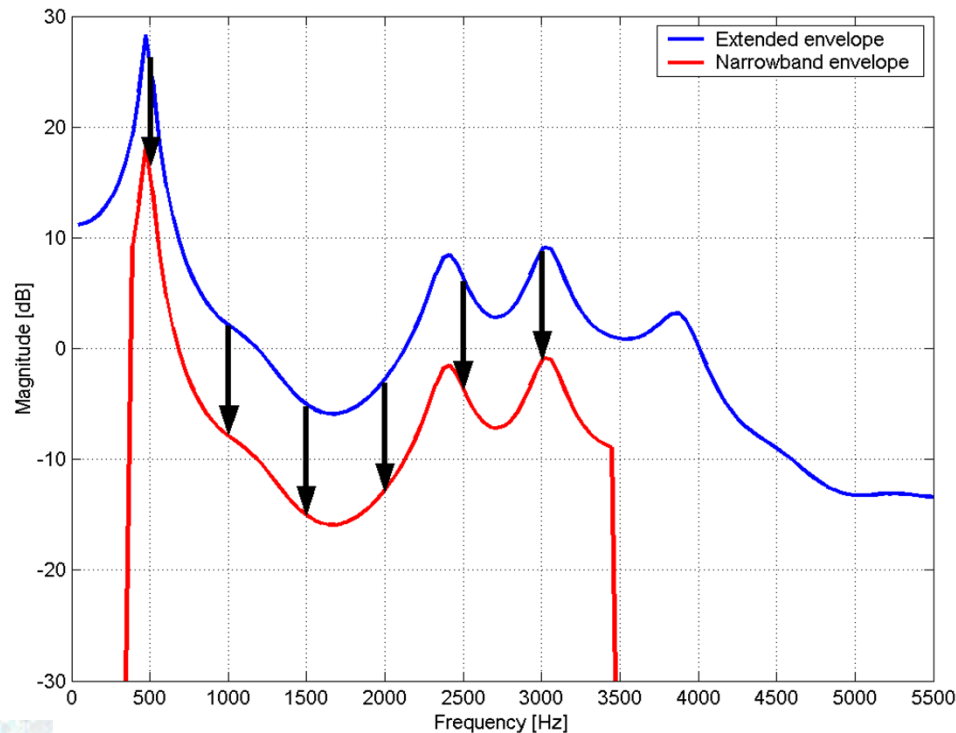
- Converting WB VTAF to WB spectral envelope.



WB Spectral Envelope Estimation (18/26)

Postprocessing

- Gain adjustment of estimated WB spectral envelope to match the input NB spectral envelope.

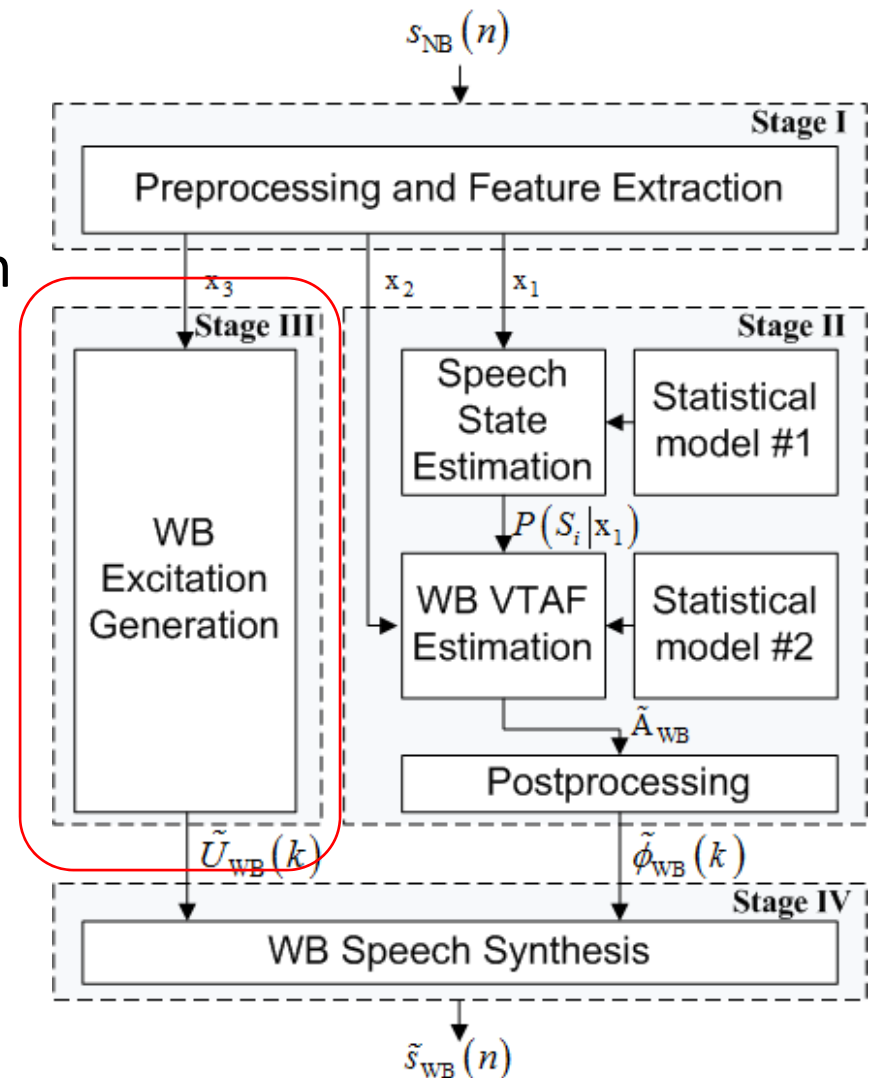


Proposed BWE Algorithm (19/26)



Algorithm stages:

- I. NB signal preprocessing and features extraction
- II. HB spectral envelope estimation and postprocessing
- III. **WB excitation generation**
- IV. Wideband signal synthesis



WB Excitation Generation (20/26)



Excitation generation - Shifting / modulation approaches:

- Fixed spectral translation:

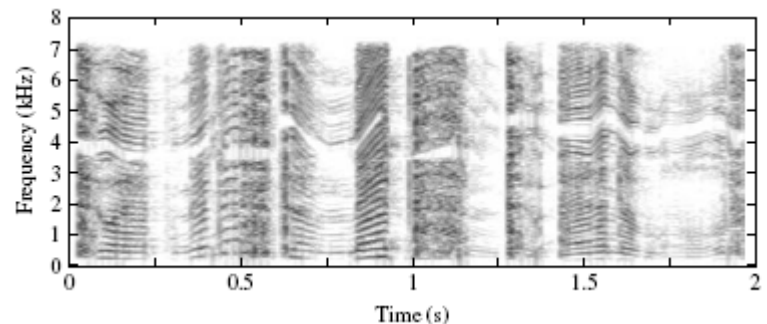
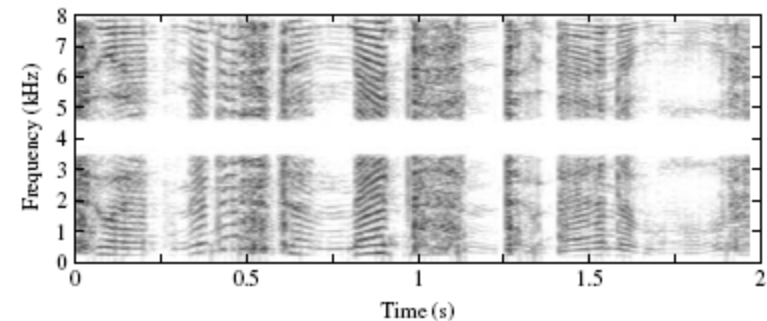
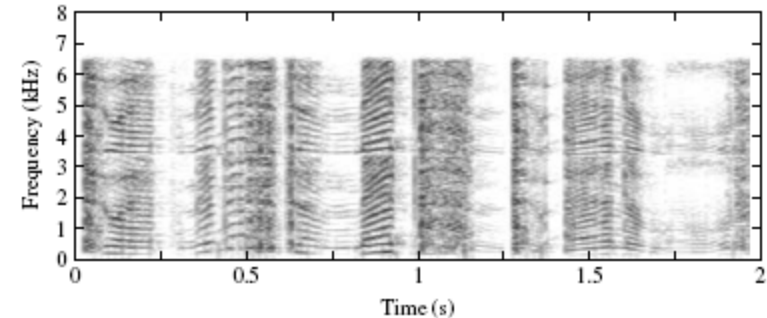
$$\Omega_m = 2\pi \frac{3.4\text{kHz}}{f_s}$$

- Spectral folding:

$$\Omega_m = 2\pi \frac{8\text{kHz}}{f_s}$$

- Pitch adaptive modulation:

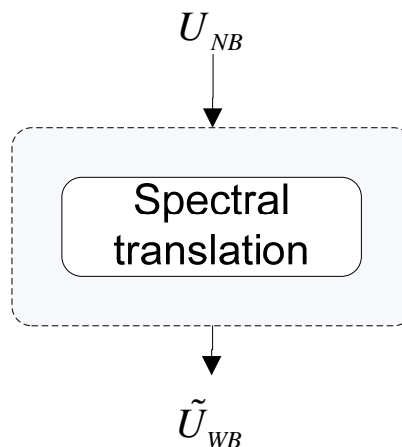
$$\Omega_m = 2\pi \frac{(nF_0)\text{kHz}}{f_s}, n = \max_n nF_0 \leq 3.4\text{kHz}$$



WB Excitation Generation (21/26)



- Extract narrowband (NB) excitation using analyzed LPC.
- Using fixed spectral translation method.
- Advantages
 - Fill all high-band frequencies.
 - Does not change narrowband excitation.
 - Low computational complexity.
- Disadvantages
 - Does not keep harmonic structure of voiced excitation.

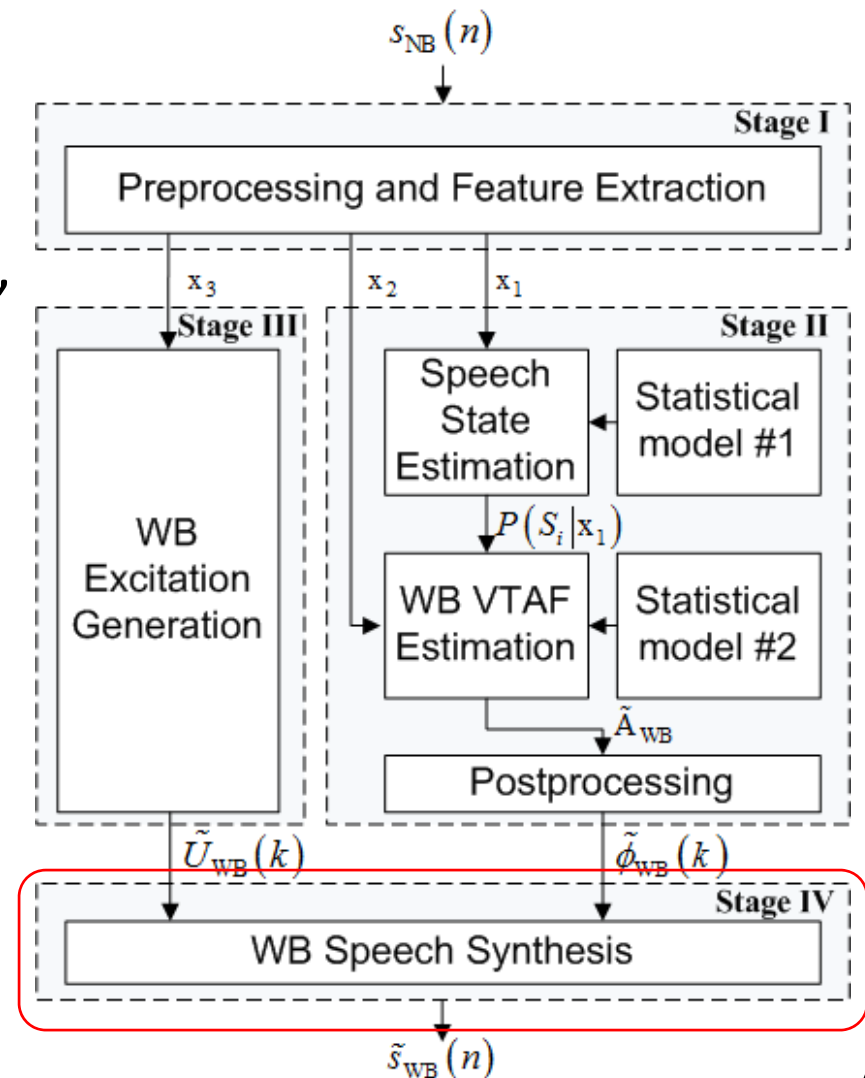


Proposed BWE Algorithm (22/26)



Algorithm stages:

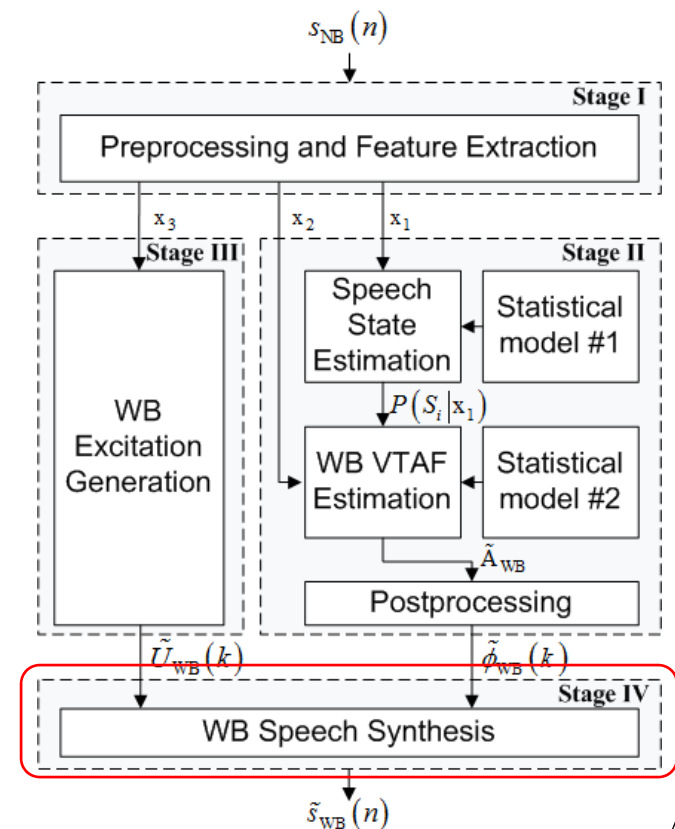
- I. NB signal preprocessing and features extraction
- II. HB spectral envelope estimation, and post processing
- III. WB excitation generation
- IV. **Wideband signal synthesis**



Wideband Speech Synthesis (23/26)



- Make sure not to change received narrow-band signal.
- Combine the estimated high-band signal and the narrow band signal.
 - Summation in time domain.
 - Concatenation in frequency domain.



Wideband Speech Synthesis (24/26)

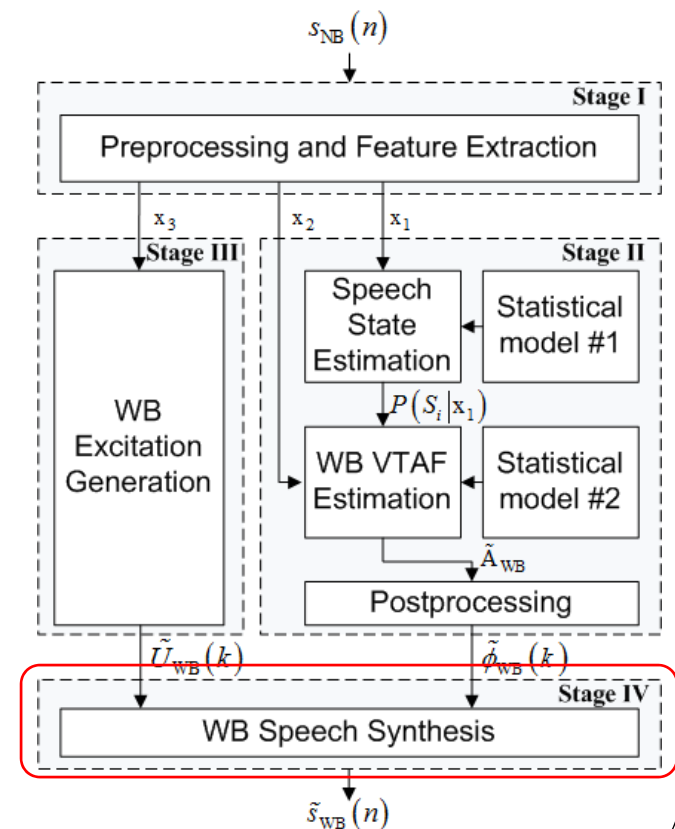


- Frequency domain synthesis.
- Calculate high-band estimated signal using:

$$\tilde{S}_{hb}(k) = \tilde{U}_{hb}(k) \cdot \tilde{\phi}_{hb}(k)$$

- Concatenating narrowband and high band signal in frequency.

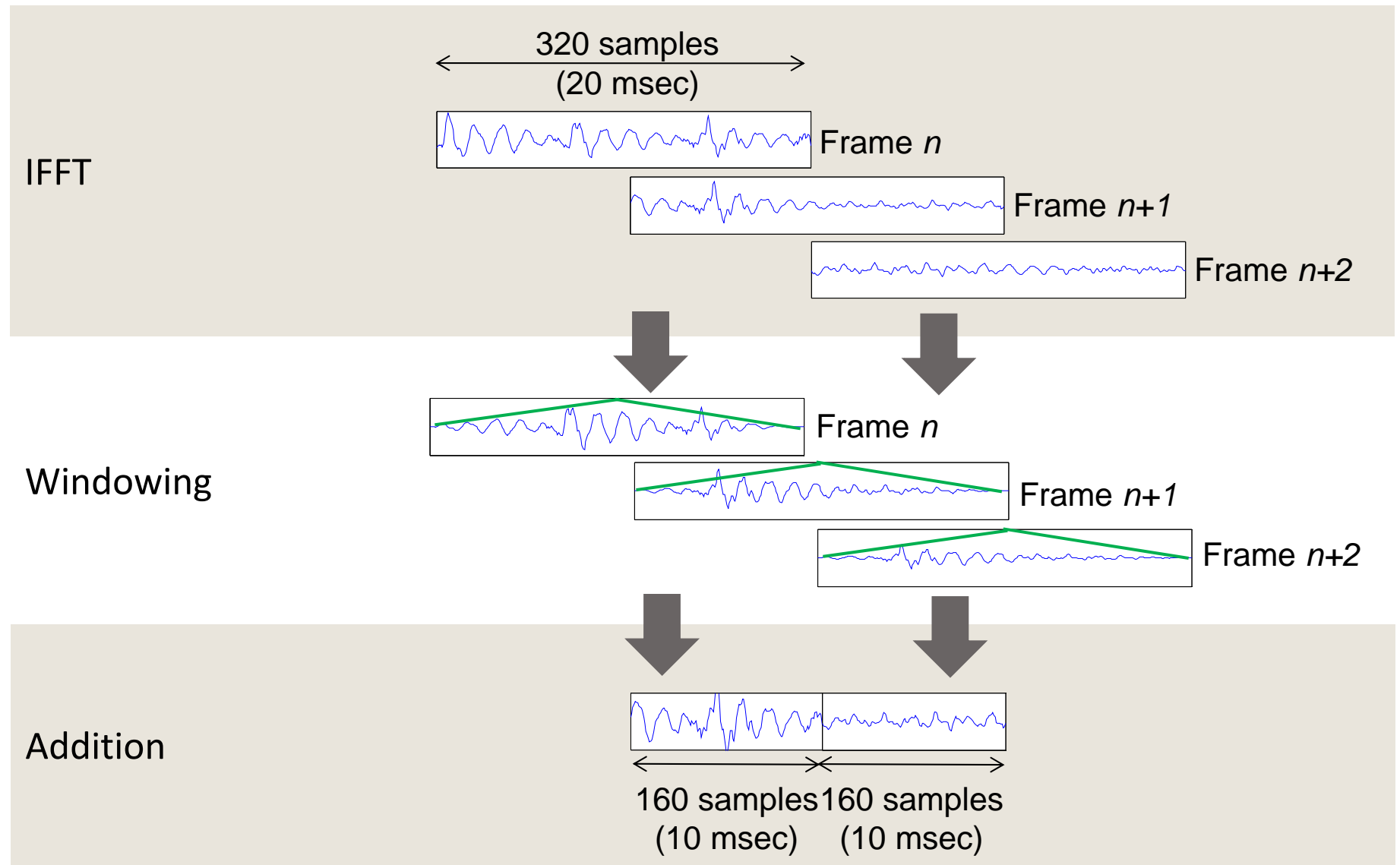
$$\tilde{S}_{wb}(k) = \begin{cases} S_{nb}(k), & 0 < k < 3.5/4kHz \\ \tilde{S}_{hb}(k), & 3.5/4 < k < 8kHz \end{cases}$$



Wideband Speech Synthesis (25/26)



Overlap add - Synthesize 20msec frames, 10msec overlap

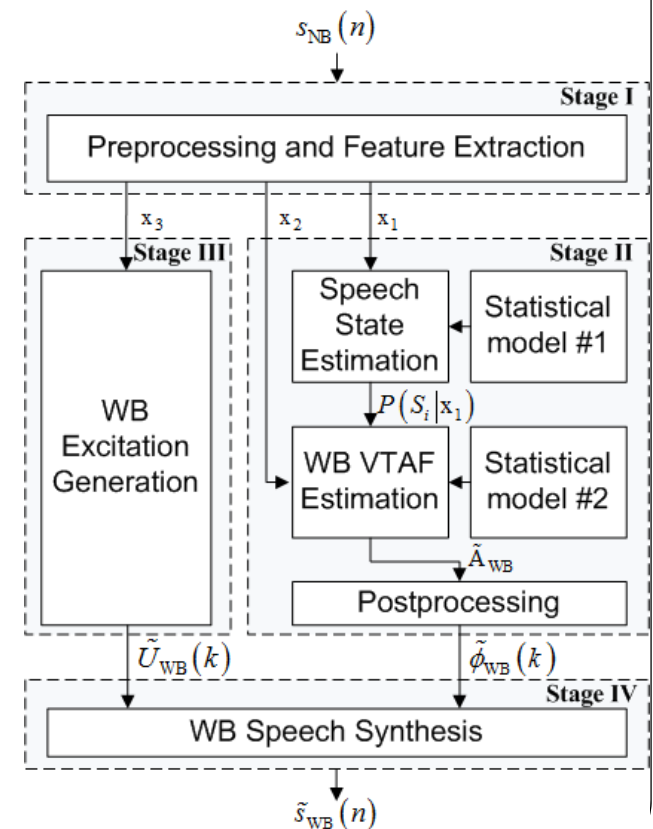


Proposed BWE Algorithm (26/26)



Summary

- Separate extension of excitation and spectral envelope.
- signal equalization of NB attenuated band – 50-300 Hz.
- Spectral envelope estimation
 - **phoneme estimation** and **speaker vocal tract area function (VTAF) estimation.**
 - **Iterative tuning of estimated VTAF** to reduce possible artifacts.
- HB Excitation estimation by spectral shifting.
- Overlap-add synthesis for better time transition.





Outline

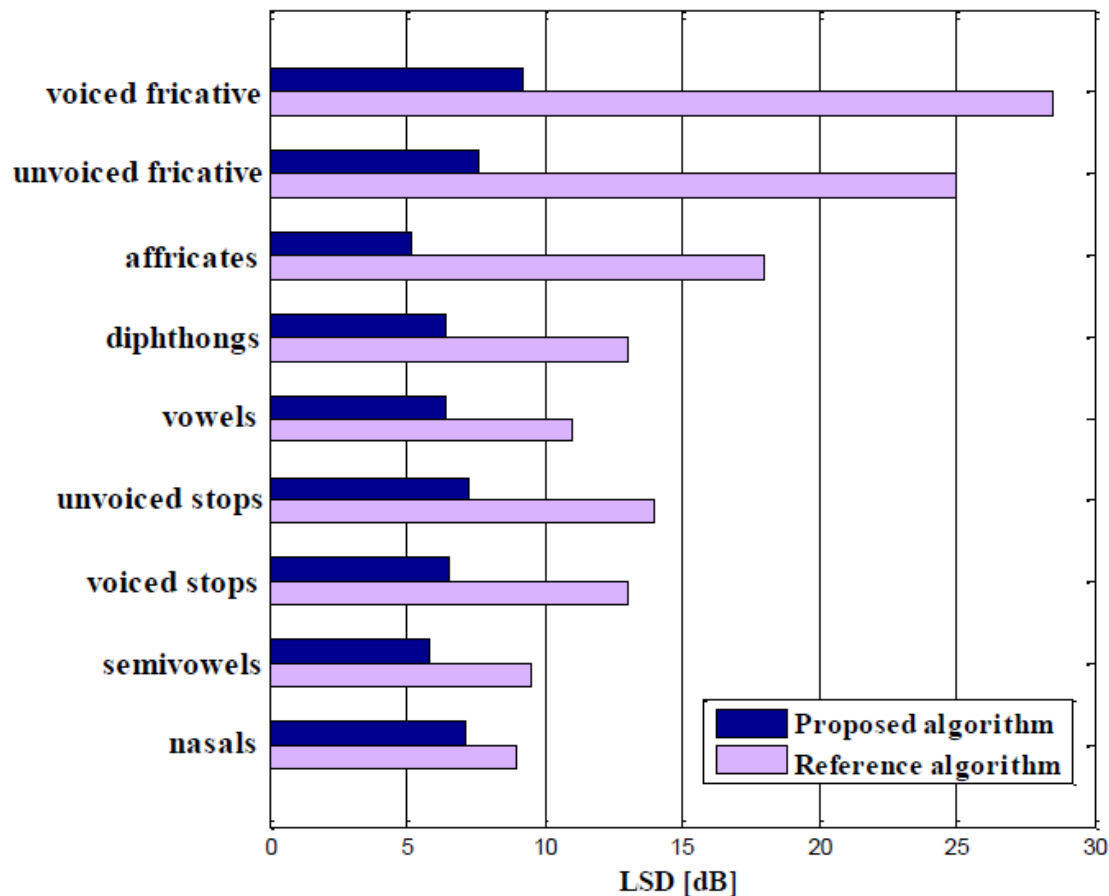
- Introduction
- Methods of BWE
- Proposed BWE Algorithm
- **Performance Evaluation**
- Conclusion

BWE Performance Evaluation (1/8)



Objective quality evaluation - Log Spectral Distance for phone category

- Reference algorithm: *Evaluation of an Artificial Speech Bandwidth Extension Method in Three Languages* [Pulakka et al., 2008].



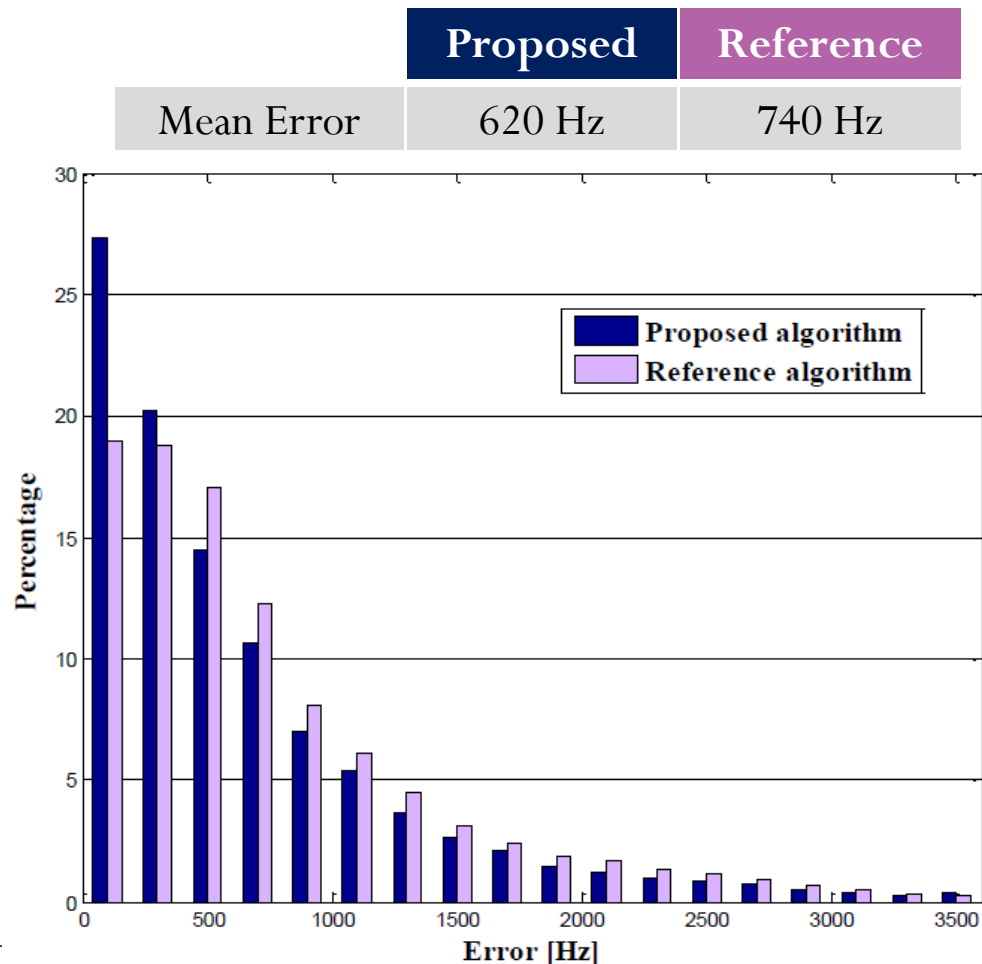
$$\text{LSD} = \sqrt{\frac{1}{k_{\text{high}} - k_{\text{low}} + 1} \sum_{k=k_{\text{low}}}^{k_{\text{high}}} \left[10 \log_{10} \frac{P(k)}{\tilde{P}(k)} \right]^2}$$

BWE Performance Evaluation (2/8)



Objective quality evaluation - Histogram of estimated formant frequencies error

- Reference algorithm: *Low-Complexity Feature-Mapped Speech Bandwidth Extension* [Gustafsson et al., 2006].



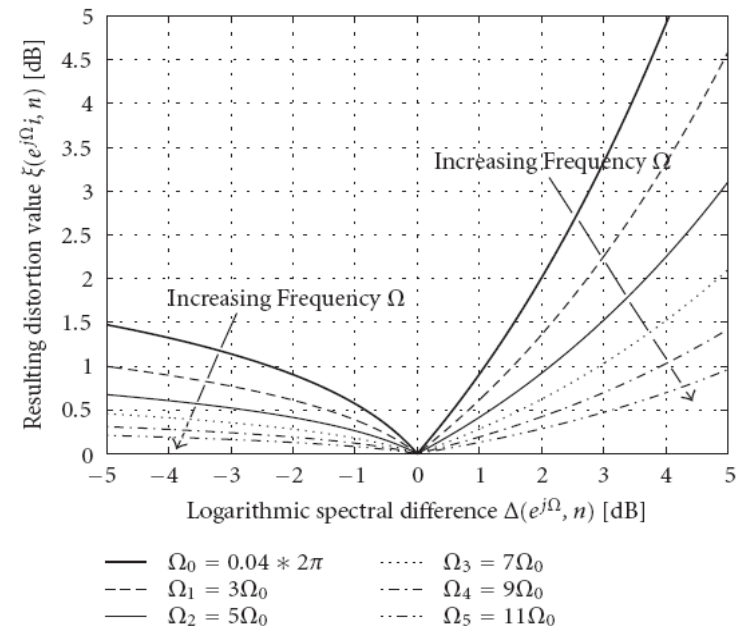
$$Error = \left| \tilde{f}_{HB} - f_{HB} \right|$$

BWE Performance Evaluation (3/8)



Objective quality evaluation - Spectral distortion measure of estimated spectral envelope with and without the iterative process

- Spectral Distortion Measure (SDM)
 - Non-symmetric distortion measure.
 - Takes the human auditory system into account.
 - penalizes spectral overestimation higher than underestimation.



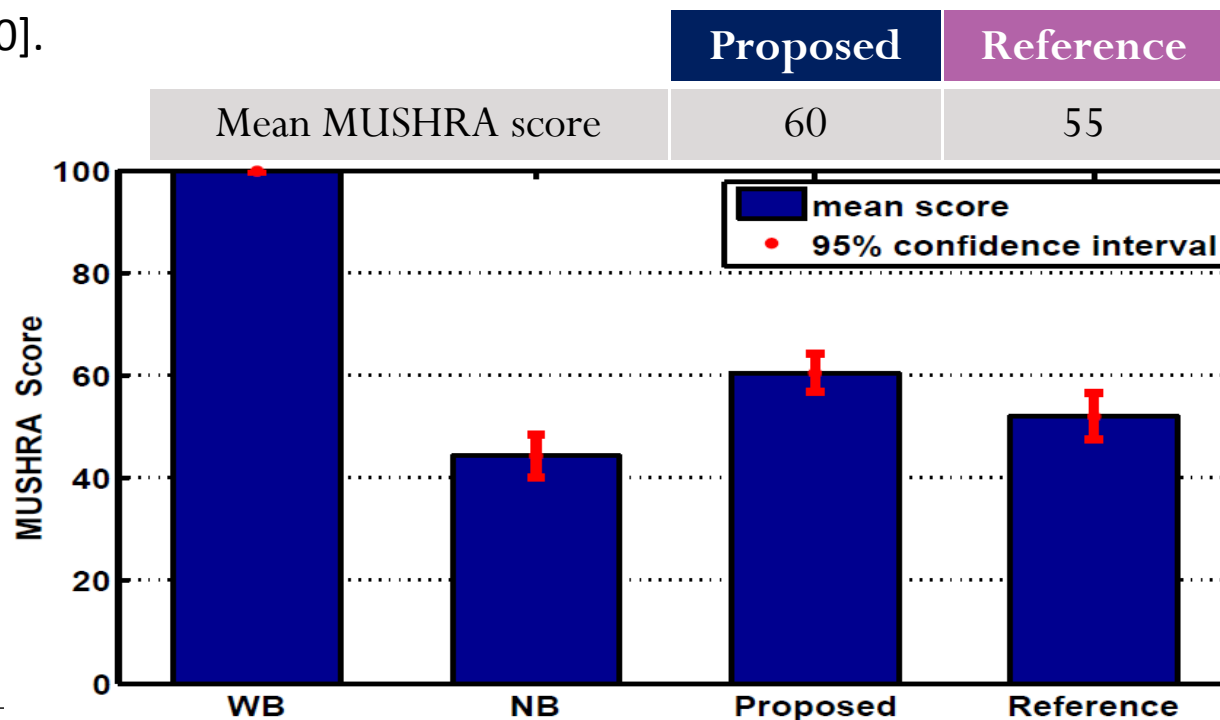
Measured	SDM [dB]	LSD [dB]
Without iterative process	13.6380	9.9759
With iterative process	9.8889	9.9057

BWE Performance Evaluation (4/8)



Subjective quality evaluation

- **MUSHRA (Multiple Stimuli with Hidden Reference and Anchor)** – ranks several speech samples for score between 0-100.
- Recommendation ITU-R BS.1116-1.
- Test setup: 11 listeners, 6 different sentences (English, 3 male, 3 female) – each with WB reference signal, NB anchor signal, proposed BWE signal and a reference BWE signal from Geiser [Based on Jax et al., 2003. modified in 2010].



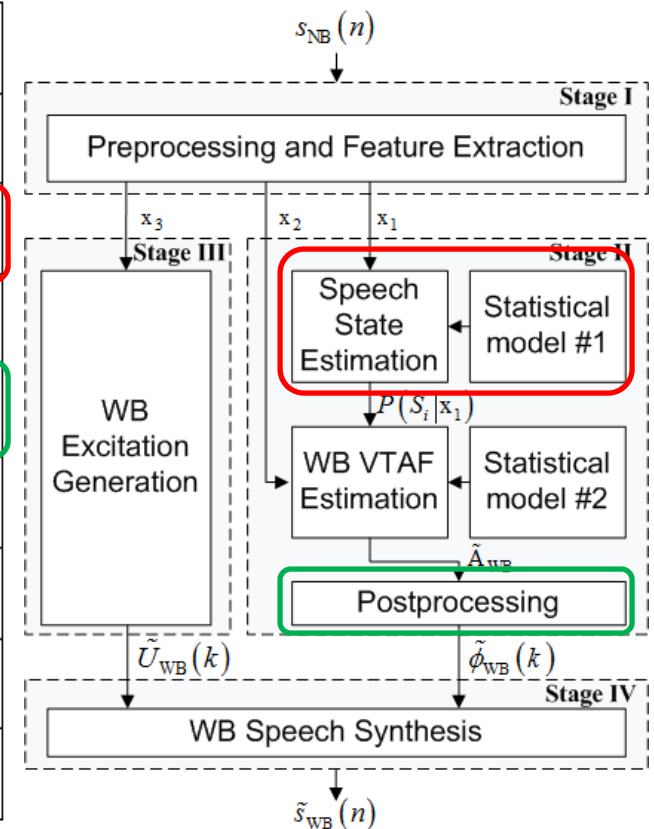
BWE Performance Evaluation (5/8)



Complexity evaluation

- Calculating Matlab processing time of each major processing block.
- The table indicate the average processing time of a 20 msec speech frame.

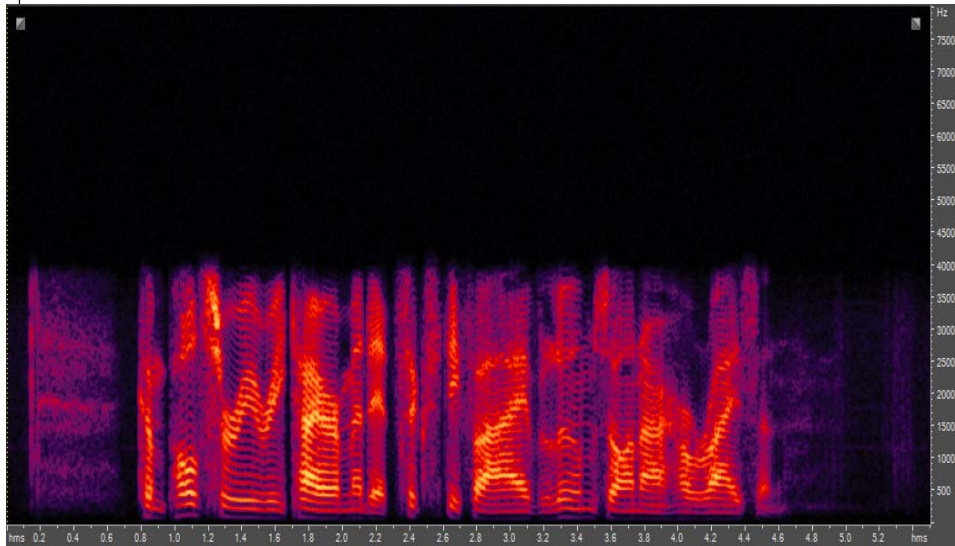
Algorithm Processing Block	Computation Time [msec]
Preprocessing and feature extraction	1.27
State estimation	19.39
WB VTAF estimation	0.59
Postprocessing (iterative process)	7.69
Postprocessing (gain adjustment)	0.36
WB excitation generation	0.04
WB speech synthesis	0.57
Total	29.91



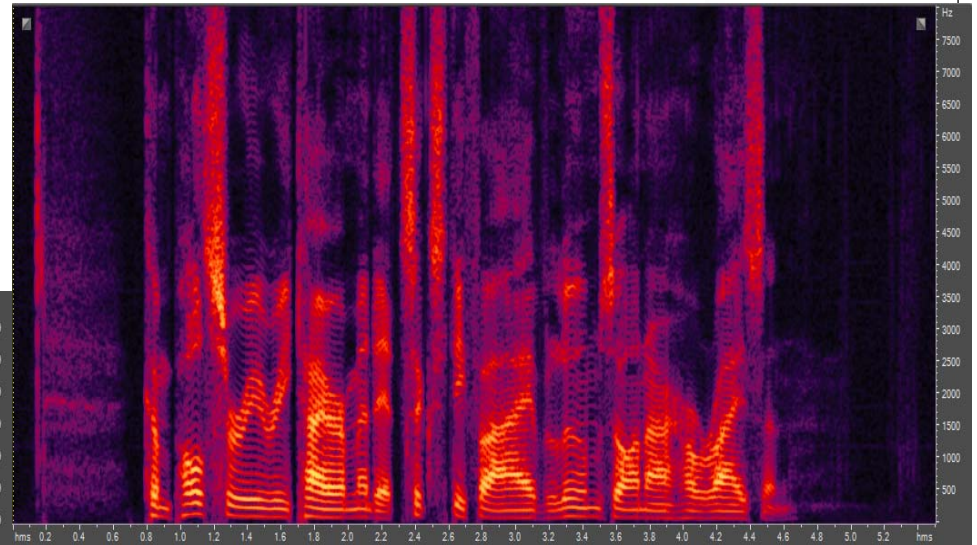
BWE Performance Evaluation (6/8)



Male spectrograms



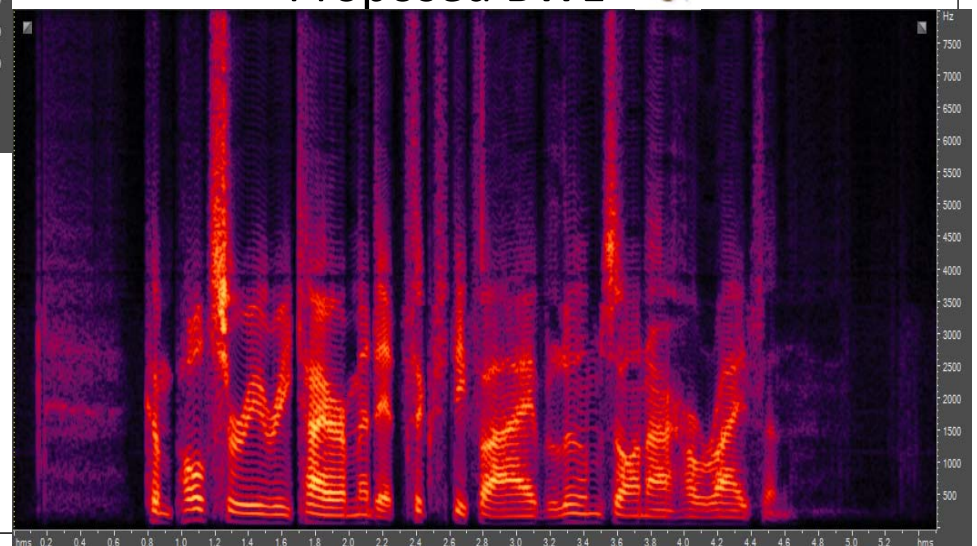
Narrowband



Wideband



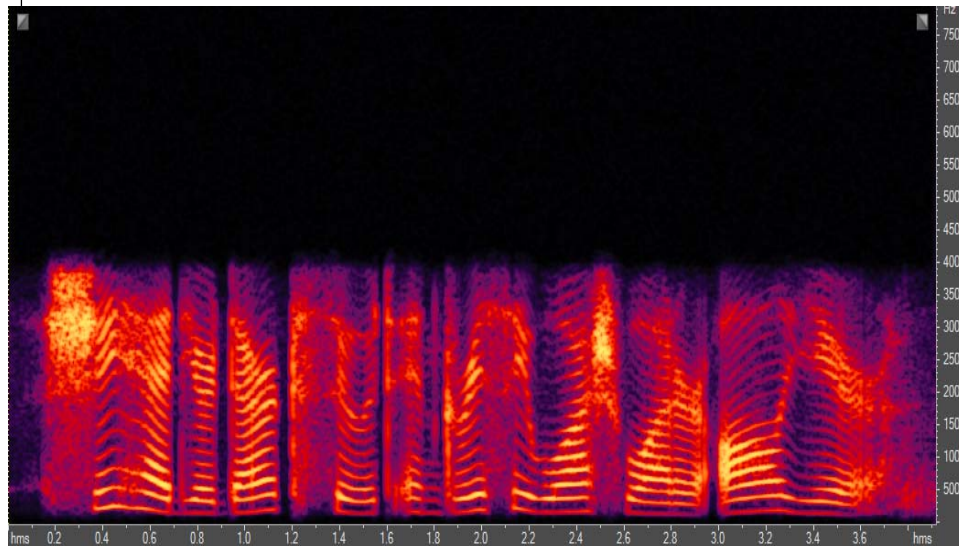
Proposed BWE



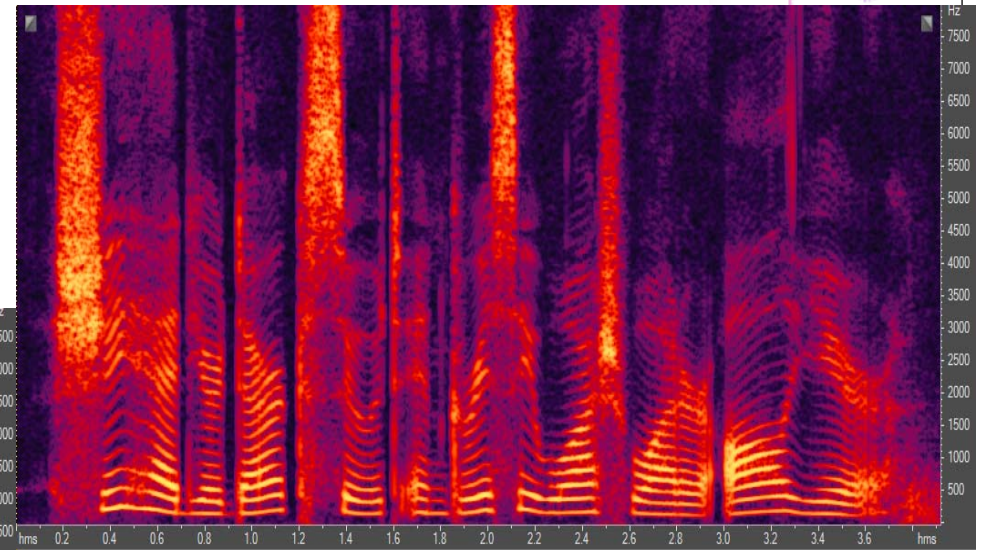
BWE Performance Evaluation (7/8)



Female spectograms

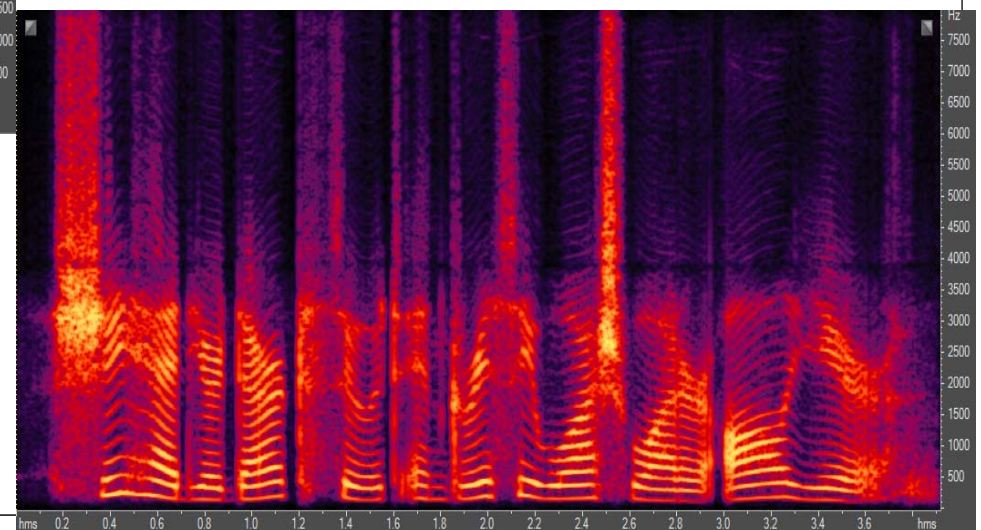


Narrowband 🗣️



Wideband 🗣️

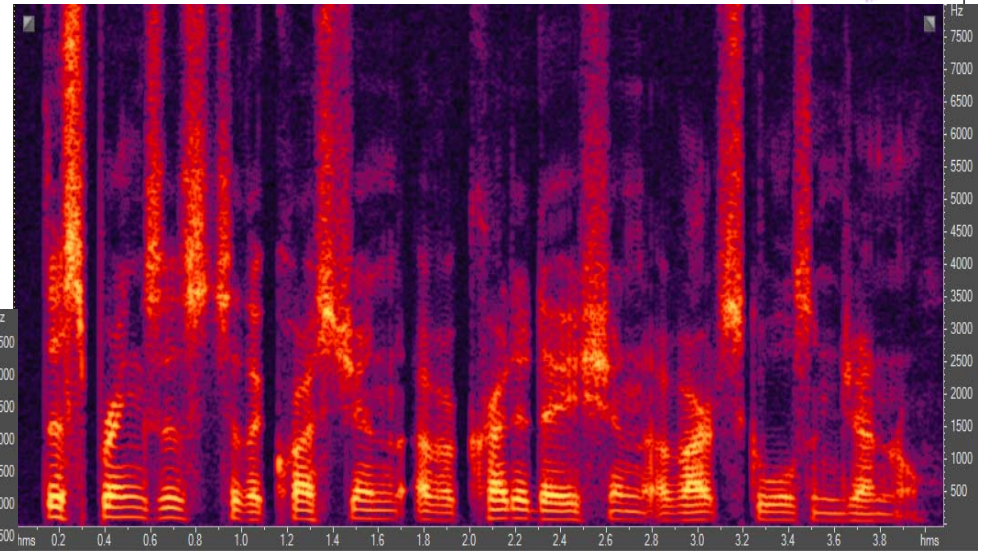
Proposed BWE 🗣️



BWE Performance Evaluation (8/8)

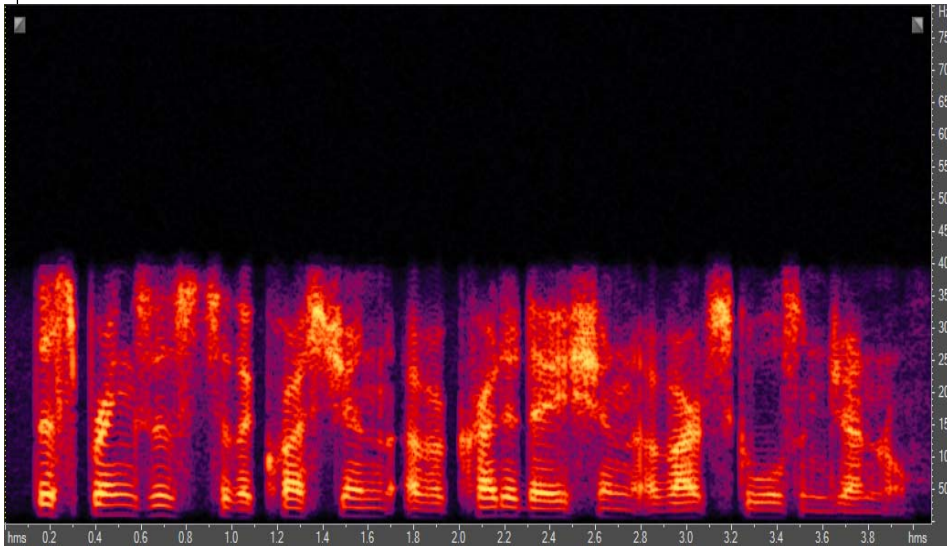
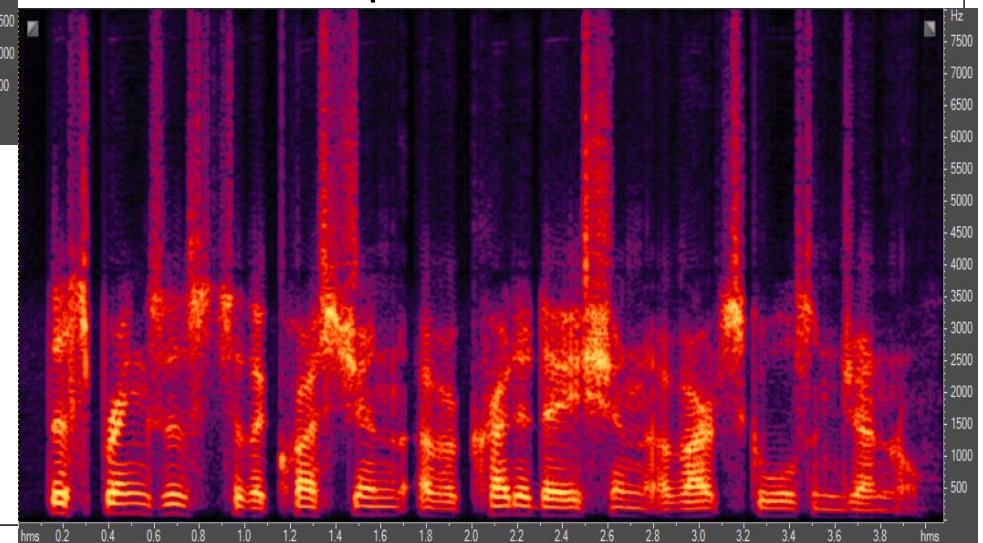


Male spectrograms



Wideband 🗣️

Proposed BWE 🗣️



Narrowband 🗣️



Outline

- Introduction
- Methods of BWE
- Proposed BWE Algorithm
- Performance Evaluation
- **Conclusion**



Conclusion (1/2)

- Proposed BWE algorithm innovations
 - Phonetic content estimation for each speech frame.
 - WB VTAF estimation for specific speaker using codebook search.
 - Iterative tuning of estimated WB VTAF for better gain adjustment.
- Algorithm advantages
 - Phonetic-based estimation reduces estimation error for unvoiced frames.
 - Iterative tuning at postprocessing reduces estimation error artifacts.
- Algorithm shortcomings
 - Using VTAF is lacking in modeling nasal and unvoiced sounds.
 - HMM-based phoneme estimation and online sensitivity function calculation at the postprocessing step, result in high algorithm complexity.

Conclusion (2/2)



- Future Work
 - Reduce algorithm complexity by using sensitivity functions tables for each WB VTAF codeword.
 - Use the postprocessing iterative procedure for better refinement and control of estimated spectral envelope by high-band formants tuning to past estimated high-band formants.
 - Change spectral envelope estimation model for unvoiced and nasal sounds.
 - Check algorithm robustness to background noise and to different languages.



Artificial Bandwidth Extension of Band Limited Speech Based on Vocal Tract Shape Estimation

Thank You

References - General



- “Bandwidth Extension of Speech Signals”, Bernd Iser, Wolfgang Minker, Gerhard Schmidt, Springer 2008.
- “Audio Bandwidth Extension”, Erik Larsen, Ronald M. Aarts, wiley 2004.
- “Artificial Bandwidth Extension of Narrowband Speech”, Master’s Thesis, Nels Rohde, Svend Aage Vedstesen, Aalborg University 2007.
- “Toward Wideband Speech by Narrowband Speech Bandwidth Extension: Magic Effect or Wideband Recovery?”, Doctor’s Thesis, Gilles MIET, Maine University 2001.
- “Bandwidth Extension of Speech Signals: A Catalyst for the Introduction of Wideband Speech Coding?”, Peter Jax and Peter Vary, RWTH Aachen University, IEEE Communications Magazine, May 2006.

References - Papers



- "Quality Enhancement of Band Limited Speech by Filtering and Multirate Techniques", H. Yasukawa, International Conference on Spoken Language Processing, ICSLP '94,27.7, pp. 1607-1610, Sept. 1994.
- "Adaptive filtering for broad band signal reconstruction using spectrum extrapolation", H. Yasukawa, The 7th IEEE 1996 Digital Signal Processing Workshop, 1996.
- "Bandwidth expansion of speech based on vector quantization of the mel-frequency cepstral coefficients", N. Enbom, W. B. Kleijn, in Proc. IEEE Workshop on Speech Coding, Porvoo, Finland, pp. 171–173, 1999.
- "Techniques for artificial bandwidth extension of telephone speech", Ulrich Kornagel, Signal Processing , Vol. 86 , Nr. 6, p. 1296—1306, 2006.
- "Low-complexity feature-mapped speech bandwidth extension", H. Gustafsson, U. A. Lindgren, I. Claesson, IEEE Trans. On Audio, Speech & Language Processing, vol. 14, pp. 577-588, 2006.
- "On artificial bandwidth extension of telephone speech", P. Jax and P. Vary, Signal Processing, vol. 83, no. 8, pp. 1707–1719, 2003.
- "Evaluation of an artificial speech bandwidth extension method in three languages", H. Pulakka, L. Laaksonen, M. Vainio, J. Pohjalainen, P. Alku, IEEE Trans. Audio Speech Lang. Process. 16 (6) (August 2008) 1124–1137.
- "Estimation of vocal-tract shapes from acoustical analysis of the speech wave. The state of the art", H. Wakita. IEEE Transactions on Acoustics, Speech and Signal Processing. ASSP-27(3), June 1979.
- "Direct estimation of the vocal tract shape by inverse filtering of acoustic speech waveforms", H. Wakita. IEEE Transactions on Audio and Electroacoustics, AU-21(5)L417. October 1973.

Introduction – Methods of BWE



Non-model BWE:

- Nonlinear based approach [Yasukawa 96]
 - The highband content is extracted by nonlinear waveform rectification in time domain:

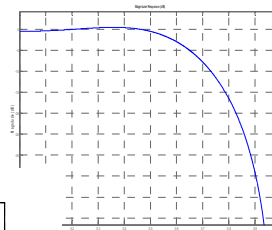
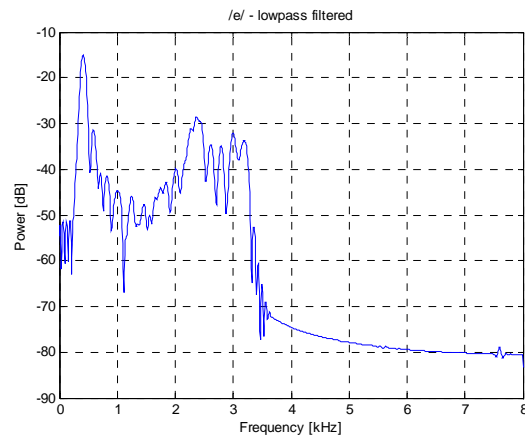
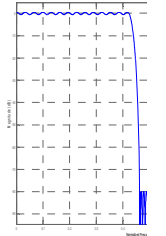
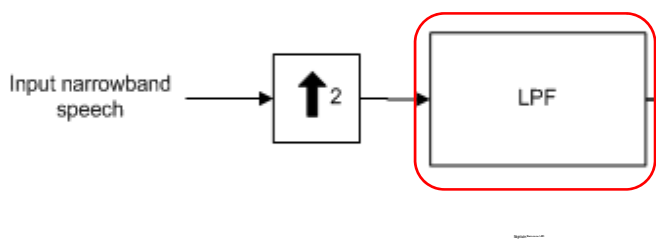
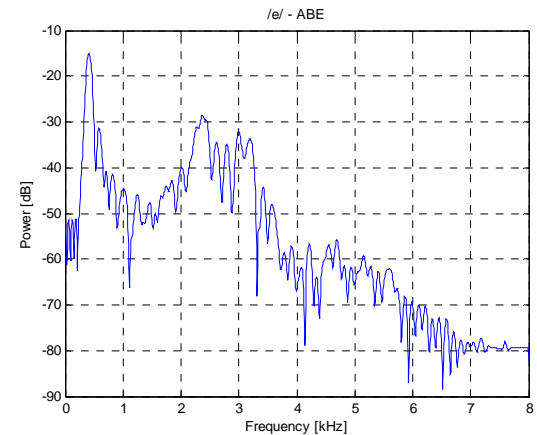
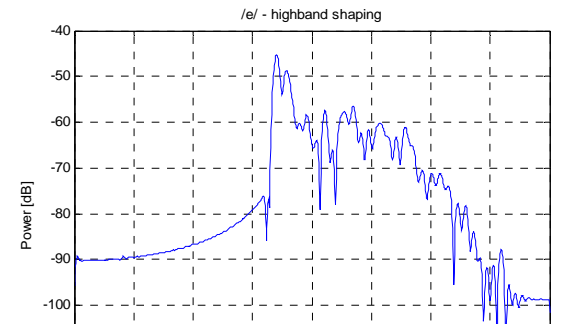
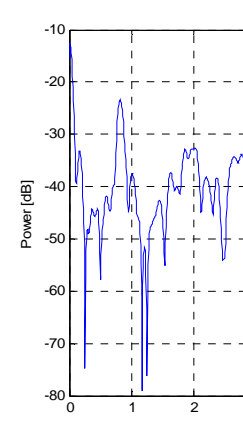
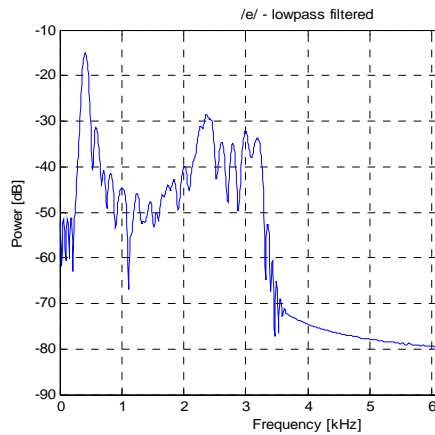
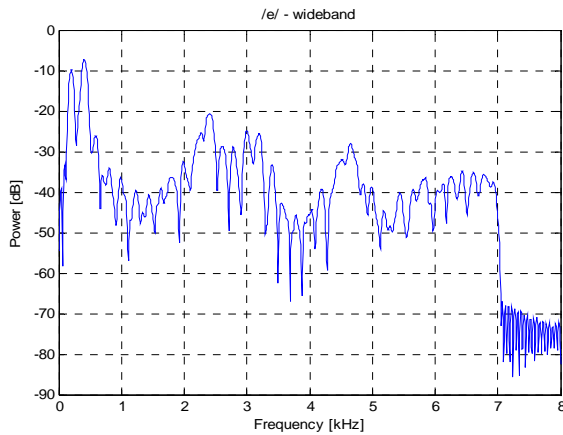
$$y(t) = \frac{[(1+\alpha) \cdot |x(t)| + (1-\alpha) \cdot x(t)]}{2}, 0 \leq \alpha \leq 1$$

$\alpha=0$ corresponds to half-wave rectification

$\alpha=1$ corresponds to full-wave rectification

- Highband is attenuated and shaped
- Wideband signal is generated as the summation of the lowpassed upsampled narrowband signal and the shaped highband signal

Introduction – Methods of BWE



Introduction – Algorithms Review

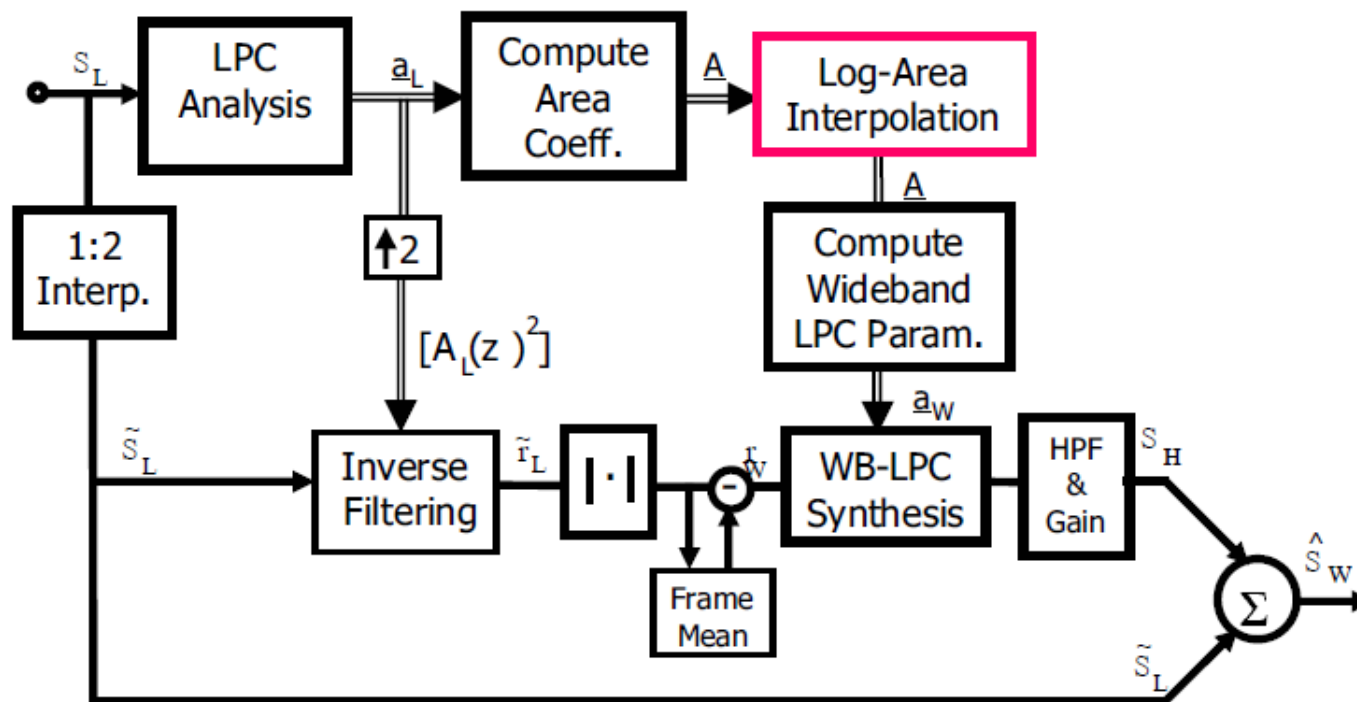


	Energy shift	Spectral shaping	Gain adjustment	Results
Kabal (model based)	Narrowband equalizer and modulation of band-limited signal	Estimate spectral envelope using GMM based mapping of LSF/MFCC	Estimate gain using GMM based mapping	Used objective LSD
Epps (model based)	Sinusoids in the pitch harmonic frequencies	Estimate spectral envelope using codebook mapping	narrowband energy equalization of estimated wideband signal to original signal in 3-3.5kHz	Used objective SD
Yasukawa (model based)	Non-linear / spectral folding using multirate analysis	Constant shaping filter	Level adjustment based on off line training	Used subjective evaluation



Introduction - Malah's BWE Algorithm

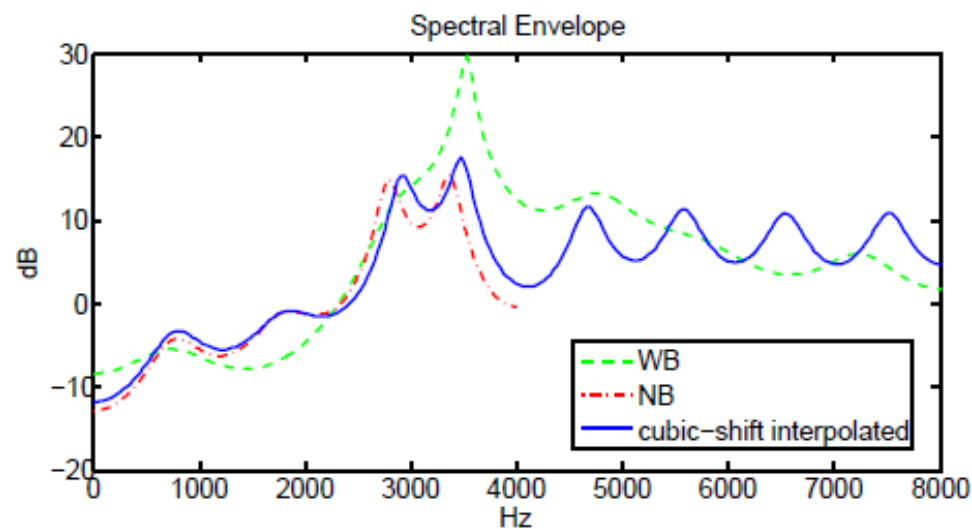
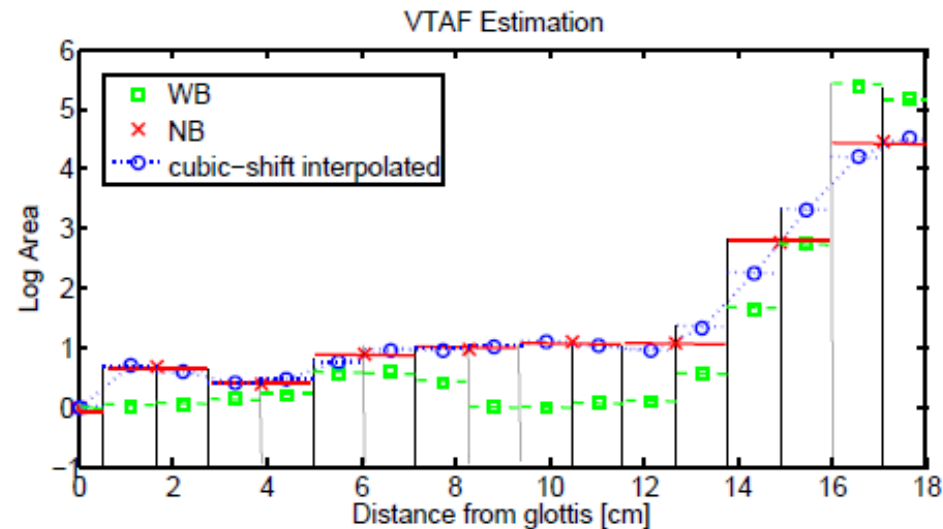
- Excitation extension
 - Using waveform rectification on the NB excitation
- Spectral envelope extension
 - Derive area/log area coefficients from LPC
 - Shift interpolate coefficients with cubic-spline function
 - Derive LPC from interpolated area/log area coefficients



Introduction - Malah's BWE Algorithm



- Example of Malah's BWE algorithm result

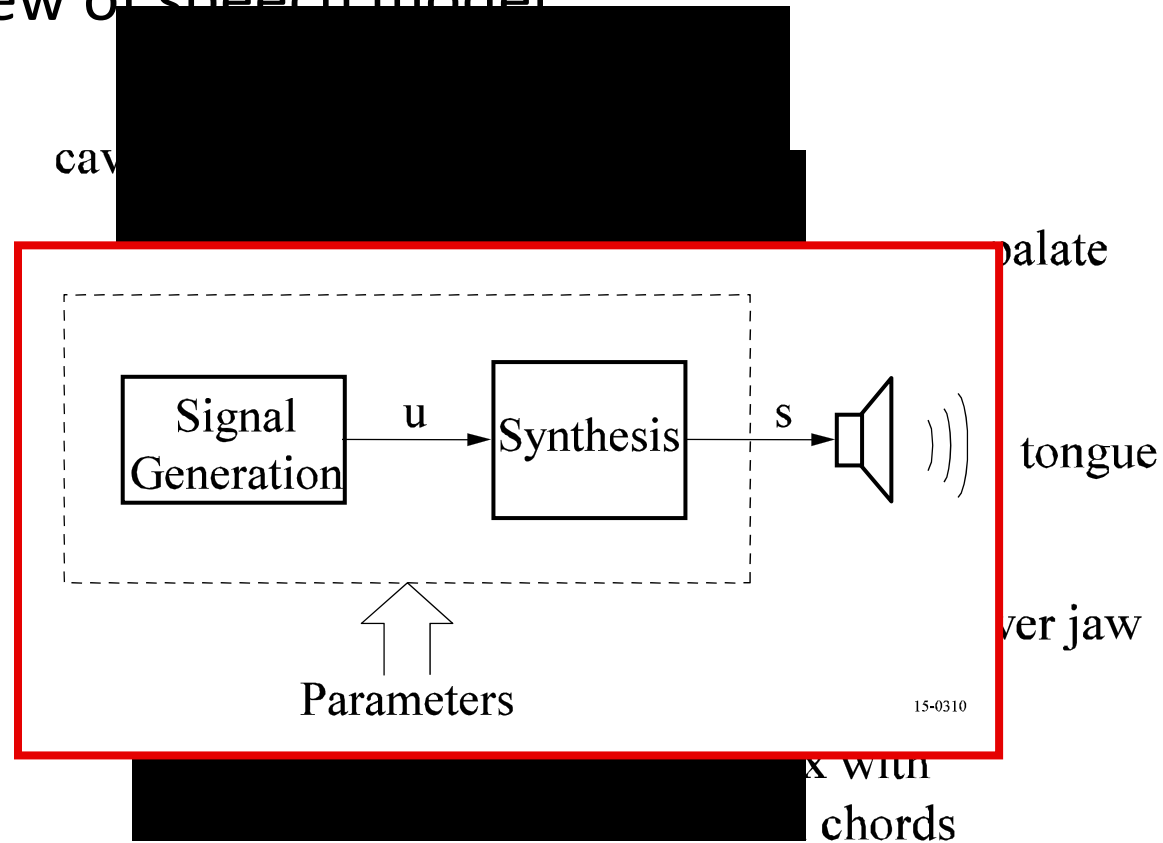


Introduction – Methods of BWE



Model-based BWE:

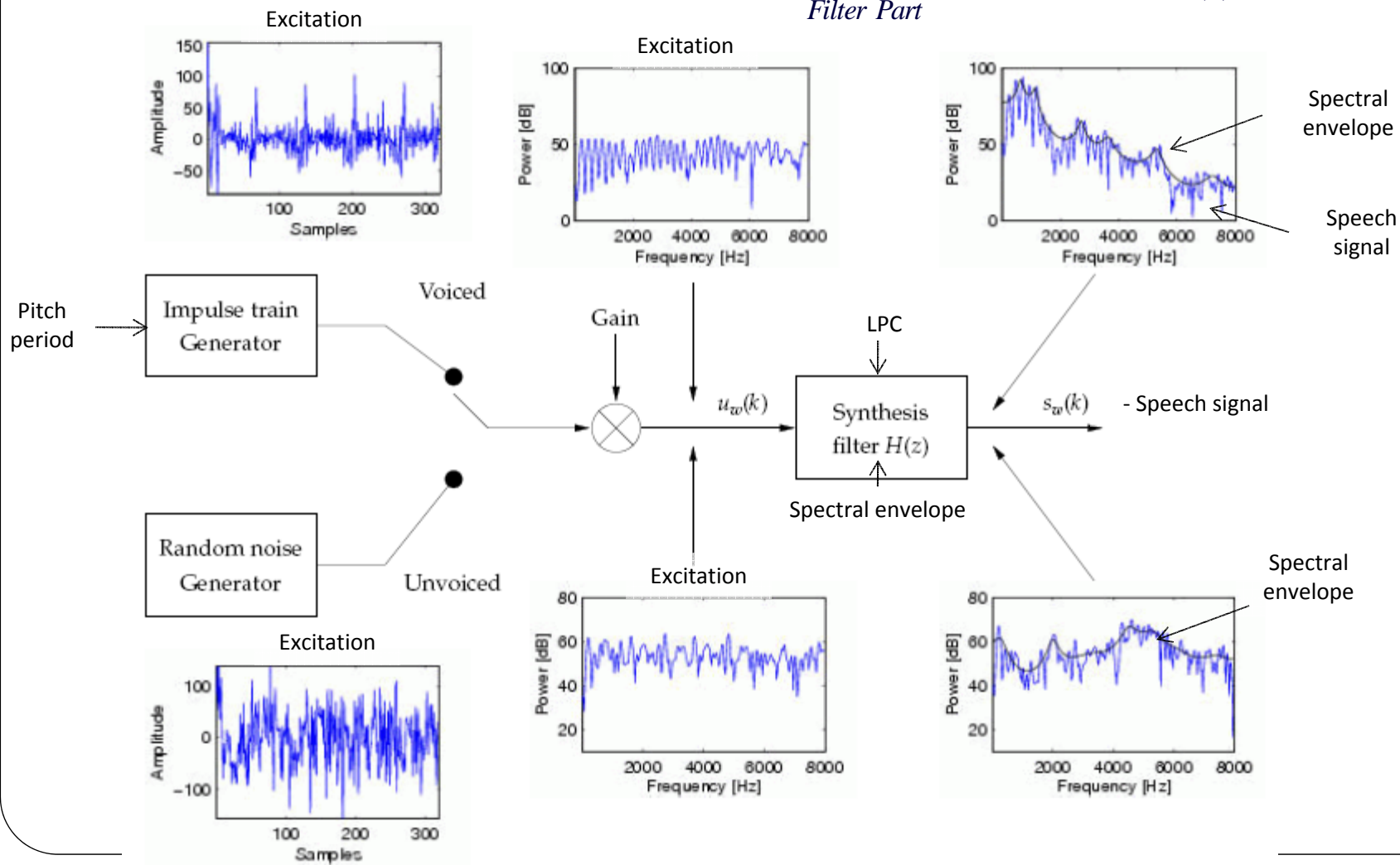
- Short review of speech model



Introduction – Methods of BWE



- Source-filter model:
$$x(n) = \underbrace{\sum_{k=1}^P a_k \cdot x(n-k)}_{\text{Filter Part}} + \underbrace{G \cdot u(n)}_{\text{Source Part - } e(n)}$$



Introduction – Methods of BWE



Source-filter model parameters:

- Spectral envelope:
 - LPC/AR - Linear Prediction Coefficients / Autoregressive coefficients
 - LSF/LSP – Line Spectral (Frequency/Pair) Coefficients
 - Cepstral coefficients
 - MFCC – Mel Frequency Cepstral Coefficients
 - Reflection coefficients
 - Area (Log Area) coefficients
- Excitation signal:
 - Pitch frequency
 - Voicing degree
 - Gain

Introduction - Model Based BWE



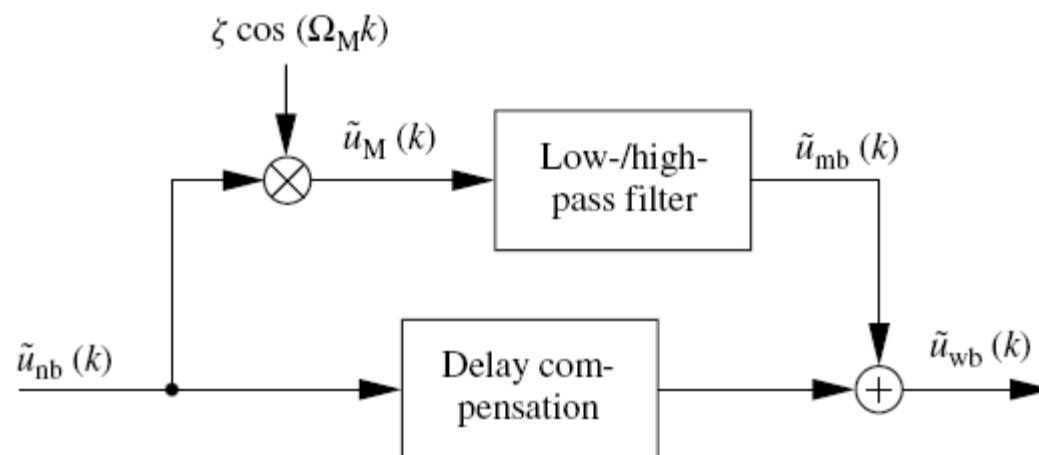
- Model based methods that need training vs. no training
 - Training – estimation of the wideband parameters using a-priori information based on large databases training.
 - Various training techniques
 - VQ (LBG).
 - Statistical approaches (GMM/HMM based).
 - Pattern recognition techniques (Neural networks).
 - Mostly used for spectral envelope extension.
 - Mostly gives higher quality as compared to no training methods.
 - No training – estimation of the wideband parameters using signal processing methods on the narrowband parameters.
 - Interpolation, Modulation, Nonlinear methods (etc).
 - Mostly used for excitation extension.
 - simpler and more robust to speech variability.

Introduction - Model Based BWE



Excitation generation - Shifting / modulation approaches:

- Shifting the spectrum of the narrowband excitation into the upper part of the spectrum
- Different possibilities for modulation frequency Ω_m
- Band pass filter of the extended band excitation



Introduction - Model Based BWE



Excitation generation - Shifting / modulation approaches:

- Spectral fixed translation:

$$\Omega_m = 2\pi \frac{3.4\text{kHz}}{f_s}$$

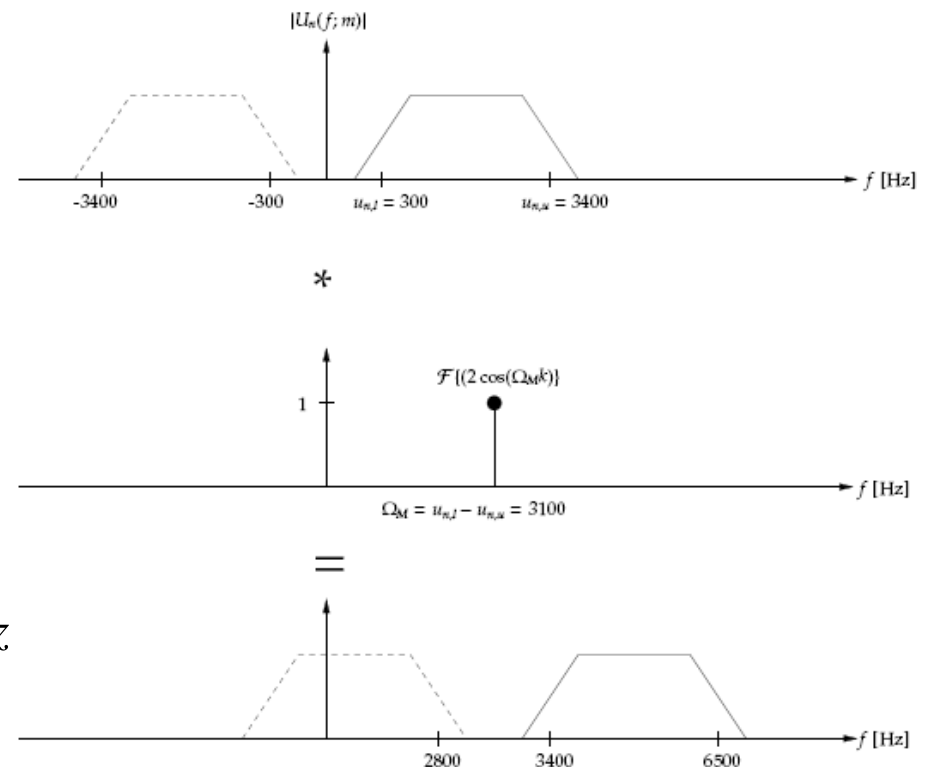
- Spectral folding:

$$\Omega_m = 2\pi \frac{8\text{kHz}}{f_s}$$

- Pitch adaptive modulation:

$$\Omega_m = 2\pi \frac{(nF_0)\text{kHz}}{f_s}, n = \max_n nF_0 \leq 3.4\text{kHz}$$

Example:



Introduction - Model Based BWE



Linear and piecewise-linear mapping

- The wideband spectral envelope is derived from the observed feature vector (NB feature vector) - $\tilde{y} = W^T x$.
- The transformation matrix W is derived using least-squares by offline training with true wideband speech database $W = (X^T \cdot X)^{-1} X^T \cdot Y$
 - Y - wideband speech spectral envelope database
 - X - narrowband speech features database

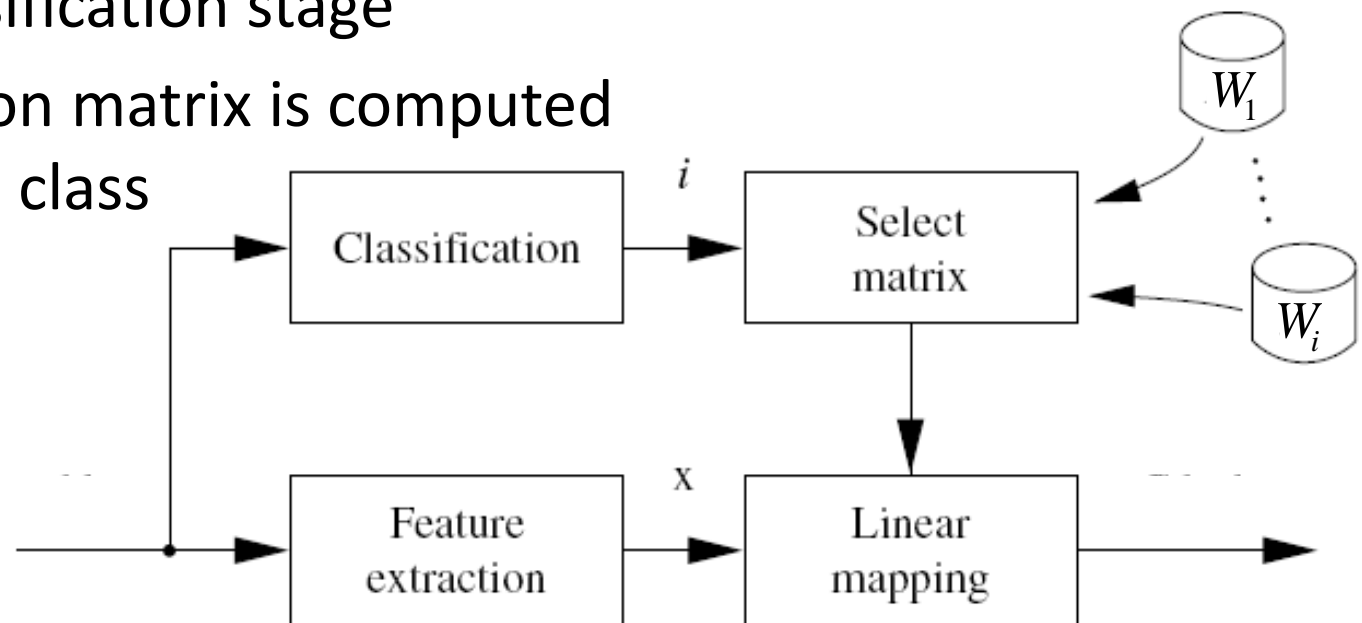
Introduction - Model Based BWE



Piecewise-linear mapping:

- The main problem of linear mapping is that linear transformation model is too simple to describe the true relationship between narrowband feature vector and the wideband spectral envelope
- The direct linear mapping approach is extended by a preceding classification stage
- A transformation matrix is computed offline for each class

$$\tilde{y} = W^T_i x$$

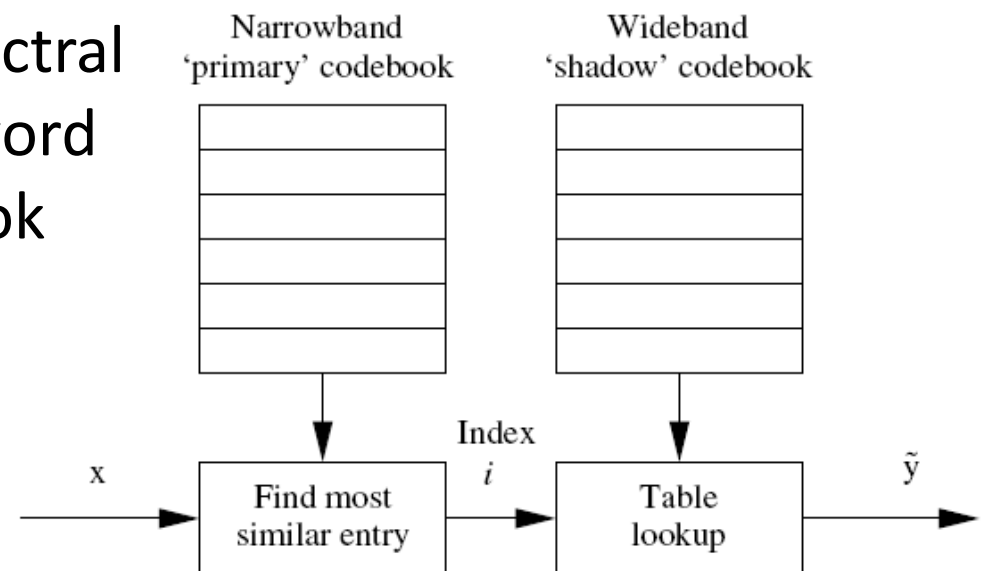


Introduction - Model Based BWE



Codebook mapping approach:

- Based on a pair of coupled codebooks that contain spectral envelopes representations of the NB and WB speech respectively
- The feature vector of the NB speech is extracted and compared with a pre-trained codebook
- The most similar NB codeword is selected
- The estimate \tilde{y} of the WB spectral envelope is simply the codeword of the “shadow” WB codebook which is assigned to the “primary” NB codeword



Introduction - Model Based BWE



Codebook mapping approach – (cont'd)

- Training phase
- Primary codebook:
 - VQ is usually utilized (LBG)
 - The quantization mapping Q is defined such as to minimize the error criterion $d(x, \hat{x}_i)$ between the input vectors x and the codewords \hat{x}_i
 - $Q(x) = \arg \min_{\hat{x}_i \in C_x} d(x, \hat{x}_i)$
- Shadow codebook :
 - All the entries in this codebook are averaged for each NB cluster
 - $\hat{y}_i = \arg \min_{\hat{y}_i \in \mathbb{R}^d} E\{d(y, \hat{y}_i) | x \in Y_i\}$

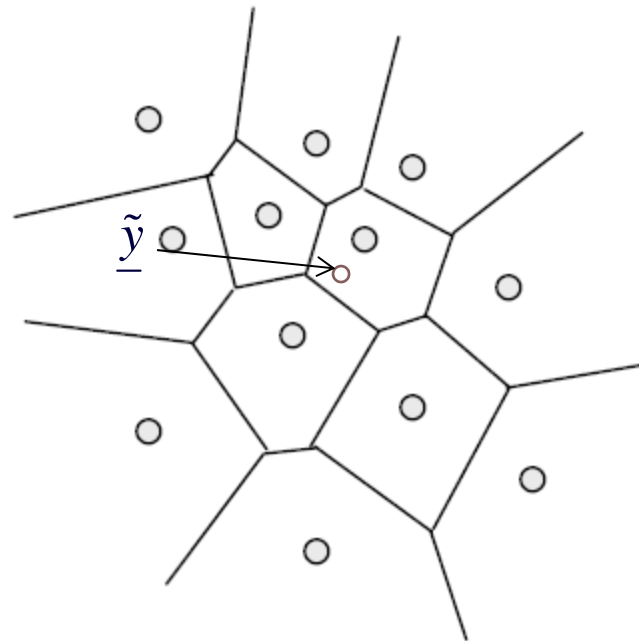
Introduction - Model Based BWE



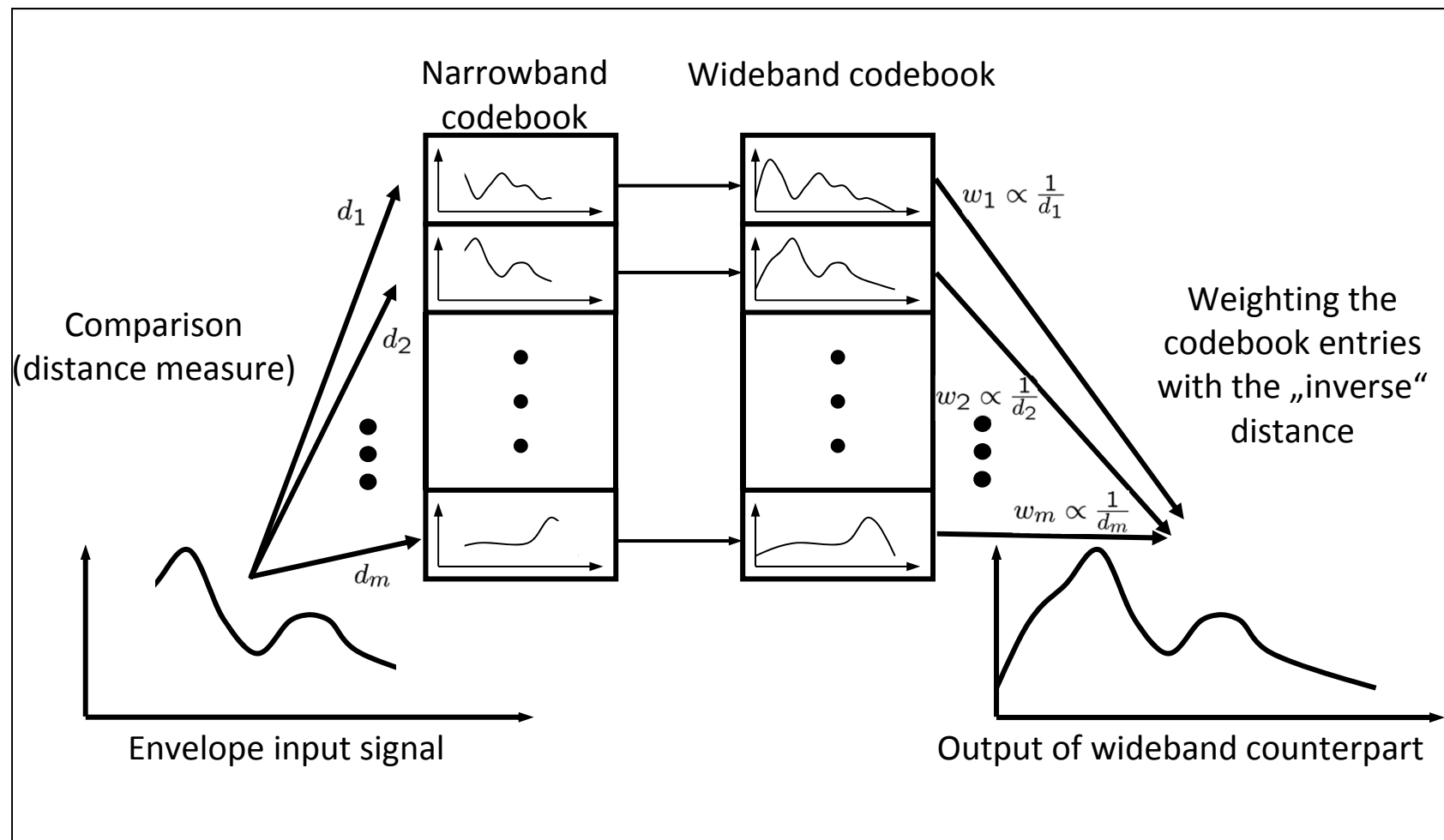
Codebook mapping approach – (cont'd)

- To improve performance, the estimate \tilde{y} can be determined by a weighted sum of most probable codewords

- $$\tilde{y} = \sum_{i=1}^{N_i} \omega_i \hat{y}_i$$



Introduction - Model Based BWE



Introduction - Model Based BWE



GMM based approach

- Taking into consideration more sophisticated statistical behavior between x and y , we must find a more exact model of their joint PDF

- Thus, by formulating a column vector $z = \begin{bmatrix} x^T & y^T \end{bmatrix}^T$

- The joint PDF can be approximated by GMM.

$$p(x, y) \approx \tilde{p}(x, y) = \tilde{p}(z) = \sum_{l=1}^L \rho_l N(z; \mu_{z,l}, V_{z,l})$$

- Where,

$$N(z; \mu_{z,l}, V_{z,l}) = \frac{\sqrt{\det A_{z,l}}}{2\pi^{\frac{\dim x + \dim y}{2}}} \exp\left(-\frac{1}{2}(z - \mu_{z,l})^T A_{z,l} (z - \mu_{z,l})\right), \mu_{z,l} = \begin{bmatrix} \mu_{x,l}^T & \mu_{y,l}^T \end{bmatrix}^T$$

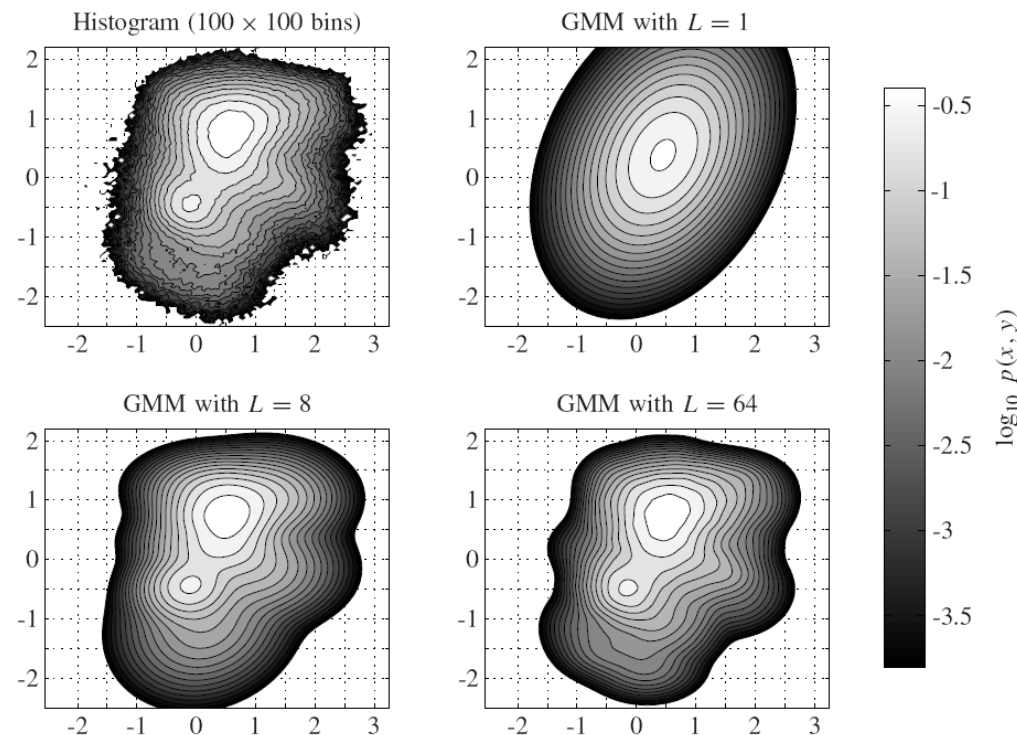
- $A_{z,l} = V_{z,l}^{-1} = \begin{bmatrix} V_{xx,l} & V_{xy,l} \\ V_{yx,l} & V_{yy,l} \end{bmatrix}^{-1}$ and $0 \leq \rho_l \leq 1, \sum_{l=1}^L \rho_l = 1$

Introduction - Model Based BWE



GMM based approach

- Models parameters are computed solving MMSE of the form $D_{MSE}(y, \tilde{y}|x) = E\{\|y - \tilde{y}\|^2 | x\}$ utilizing EM algorithm.

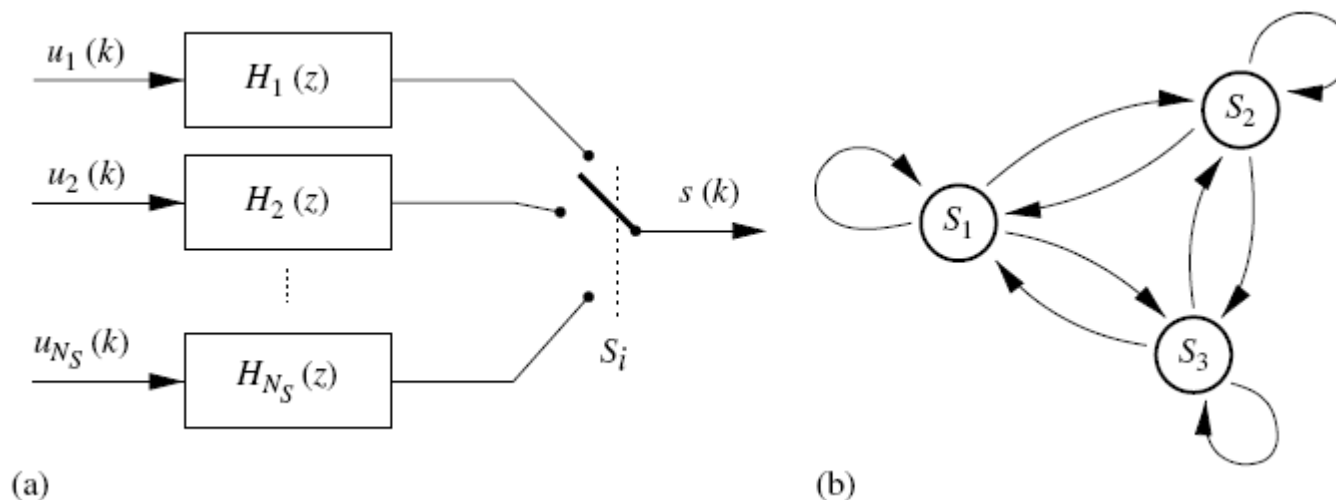


Introduction - Model Based BWE



HMM based approach

- The HMM is able to model hidden information, e.g. how a speech sequence evolves over time. Therefore it utilizes information about previous frames to estimate the extension-band.



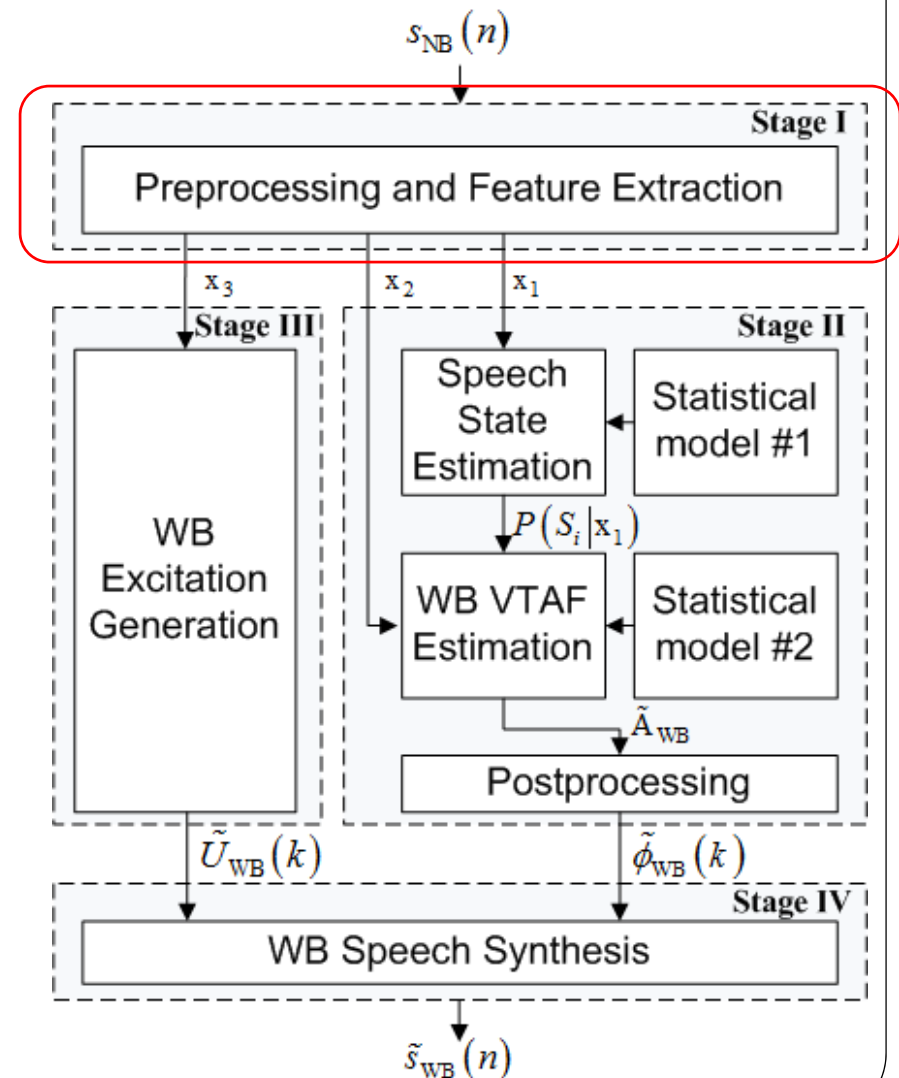
- (a) Hidden Markov model of the process of speech generation. The AR filters $H_i(z)$ and excitation signals $u_i(k)$ represent typical (wideband) speech sounds for each state.
- (b) State transition diagram of an ergodic first-order Markov chain with $N_S = 3$ states

Proposed BWE Algorithm (2/26)



Algorithm stages:

- I. NB signal preprocessing and features extraction
- II. HB spectral envelope estimation and postprocessing
- III. WB excitation generation
- IV. Wideband signal synthesis



Preprocessing and Features Extraction



Separability

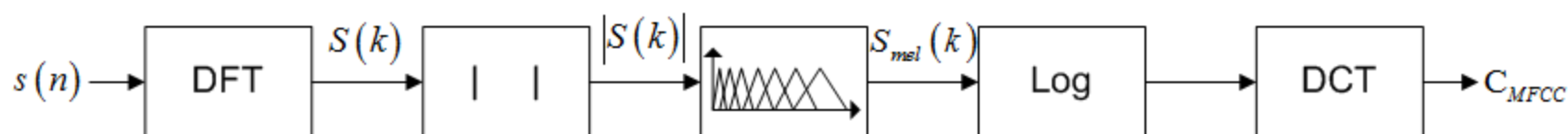
- Calculated from labeled set of training data (known class per each \mathbf{x})
- \mathbf{X}_i the set of feature vectors \mathbf{x} assigned to i^{th} class with $N_{\mathbf{X}_i}$ vectors
- N_s number of classes; $N_m = \sum_{i=1}^{N_s} N_{\mathbf{X}_i}$ number of feature vectors in DB
- Within-class covariance matrix
$$\mathbf{V}_x = \frac{1}{N_m} \sum_{i=1}^{N_s} \sum_{\mathbf{x} \in \mathbf{X}_i} (\mathbf{x} - \mu_i)(\mathbf{x} - \mu_i)^T$$
- Between-class covariance matrix
$$\mathbf{B}_x = \sum_{i=1}^{N_s} \frac{N_{\mathbf{X}_i}}{N_m} (\mu_i - \mu)(\mu_i - \mu)^T$$
- Where
$$\mu_i = \frac{1}{N_{\mathbf{X}_i}} \sum_{\mathbf{x} \in \mathbf{X}_i} \mathbf{x}, \quad \mu = \sum_{i=1}^{N_s} \frac{N_{\mathbf{X}_i}}{N_m} \mu_i$$
- Now, the separability is defined as:
$$\zeta(\mathbf{x}) = \text{tr}(\mathbf{V}_x^{-1} \mathbf{B}_x)$$
- A large separability value indicates a better suitability of the corresponding feature vector for classification and estimation

Preprocessing and Features Extraction



Extracted speech signal features

- MFCC – represent the short term power spectrum of a sound, based on a linear cosine transform of a log power spectrum on a non-linear mel-scale of frequency



- Spectral centroid – indicates where most of the power of a speech frame is spectrally located. It is generally high for unvoiced sounds

$$\mathbf{x}_{SC} = \frac{\sum_{k=0}^{N_{FFT}/2} k |S(k)|}{\left(\frac{N_{FFT}}{2} + 1\right) \sum_{k=0}^{N_{FFT}/2} |S(k)|}$$

Preprocessing and Features Extraction



Extracted speech signal features

- Spectral flatness – indicates the tonality of the speech signal. A high value indicates unvoiced sounds. A low value indicates voiced sounds.
- Spectral slope – increases during unvoiced sounds and plosives
- Normalized frame energy – calculated by normalizing the short term energy with the long term energy to receive independent measure to different speakers

WB Spectral Envelope Estimation



- Post-processing and gain adjustment
 - Reduce artifact due to erroneous estimation
 - WB spectral envelope shape fit by formant frequencies tuning to allow better gain adjustment to NB spectral envelope. Iterative tuning by VTAF perturbation
 - Iterative VTAF perturbation using the sensitivity function

$$\frac{\Delta f_{n_f}}{f_{n_f}} = \sum_{n_A}^{N_A} S_{n_f, n_A} \frac{\Delta A_{n_A}}{A_{n_A}}$$

- Stopping condition for iterative process is formant frequencies difference
- Time smoothing of tuned estimated VTAF

$$\tilde{A}'_{WB}(m) = \beta \cdot \tilde{A}'_{WB}(m-1) + (1-\beta) \cdot \tilde{A}_{WB}(m)$$

- Converting WB VTAF to WB spectral envelope
- Gain adjustment of Estimated WB spectral envelope to calculated NB spectral envelope

