

# CS578 - SPEECH SIGNAL PROCESSING

## LECTURE : HARMONIC AND QUASI-HARMONIC MODELS OF SPEECH

George P. Kafentzis



University of Crete, Computer Science Dept., Speech Signal Processing Lab  
kafentz@csd.uoc.gr  
(based on work from Prof. Stylianou and Dr. Pantazis)

Univ. of Crete

- ① FIRST WORKS ON SPEECH DECOMPOSITION...
- ② INTRODUCTION TO HNMs
- ③ ANALYSIS
  - Frequency
  - Maximum Voiced Frequency
  - Amplitudes and Phases
    - Error Function - for  $\text{HNM}_1$
    - Least Squares - for  $\text{HNM}_1$
  - Residual
- ④ SYNTHESIS
- ⑤ ENERGY MODULATION FUNCTION
- ⑥ TOWARDS QUASI-HARMONICITY
- ⑦ THANKS
- ⑧ REFERENCES

# OUTLINE

- 1 FIRST WORKS ON SPEECH DECOMPOSITION...
- 2 INTRODUCTION TO HNMs
- 3 ANALYSIS
  - Frequency
  - Maximum Voiced Frequency
  - Amplitudes and Phases
    - Error Function - for  $\text{HNM}_1$
    - Least Squares - for  $\text{HNM}_1$
  - Residual
- 4 SYNTHESIS
- 5 ENERGY MODULATION FUNCTION
- 6 TOWARDS QUASI-HARMONICITY
- 7 THANKS
- 8 REFERENCES

Mentioning just a few works for speech analysis...

- Multi-Band Excitation Vocoder (Griffin et al.1988 [1])
  - $S(\omega) = H(\omega)E(\omega)$
  - $E(\omega)$  is represented by an  $f_0$ , a V/UV decision for each harmonic, and the phase of each voiced harmonic
  - Parameters are estimated by comparing the original vs the synthetic speech spectrum
  - Voiced portion is synthesized in time domain while unvoiced part is synthesized in frequency domain

- Multi-Band Excitation Vocoder (Griffin et al.1988 [1])
  - $S(\omega) = H(\omega)E(\omega)$
  - $E(\omega)$  is represented by an  $f_0$ , a V/UV decision for each harmonic, and the phase of each voiced harmonic
  - Parameters are estimated by comparing the original vs the synthetic speech spectrum
  - Voiced portion is synthesized in time domain while unvoiced part is synthesized in frequency domain

- Multi-Band Excitation Vocoder (Griffin et al.1988 [1])
  - $S(\omega) = H(\omega)E(\omega)$
  - $E(\omega)$  is represented by an  $f_0$ , a V/UV decision for each harmonic, and the phase of each voiced harmonic
  - Parameters are estimated by comparing the original vs the synthetic speech spectrum
  - Voiced portion is synthesized in time domain while unvoiced part is synthesized in frequency domain

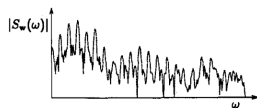
- Multi-Band Excitation Vocoder (Griffin et al.1988 [1])
  - $S(\omega) = H(\omega)E(\omega)$
  - $E(\omega)$  is represented by an  $f_0$ , a V/UV decision for each harmonic, and the phase of each voiced harmonic
  - Parameters are estimated by comparing the original vs the synthetic speech spectrum
  - Voiced portion is synthesized in time domain while unvoiced part is synthesized in frequency domain



- Multi-Band Excitation Vocoder (Griffin et al.1988 [1])
  - $S(\omega) = H(\omega)E(\omega)$
  - $E(\omega)$  is represented by an  $f_0$ , a V/UV decision for each harmonic, and the phase of each voiced harmonic
  - Parameters are estimated by comparing the original vs the synthetic speech spectrum
  - Voiced portion is synthesized in time domain while unvoiced part is synthesized in frequency domain

# BACKGROUND

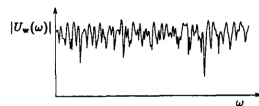
## Multi-band Excitation Vocoder (Griffin et al.1988 [1])



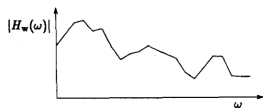
(a)



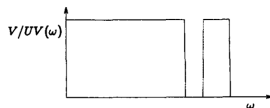
(c)



(e)



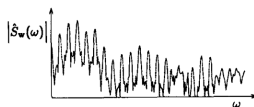
(b)



(d)



(f)



(g)

- Sinusoids + band-pass random signals (Abrantes et al.1991 [2])
  - Completely avoids V/UV decision
  - Harmonically related sinusoids model the voiced parts
  - Random band-pass signals model the unvoiced parts
    - White noise filtered by a group of band-pass filters (filterbank) with center frequencies  $k\omega_s$

- Sinusoids + band-pass random signals (Abrantes et al.1991 [2])
  - Completely avoids V/UV decision
  - Harmonically related sinusoids model the voiced parts
  - Random band-pass signals model the unvoiced parts
    - White noise filtered by a group of band-pass filters (filterbank) with center frequencies  $k\omega_s$

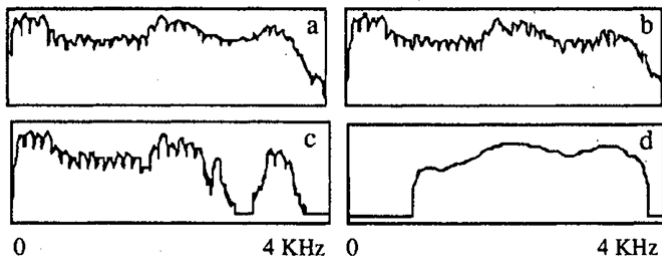
- Sinusoids + band-pass random signals (Abrantes et al.1991 [2])
  - Completely avoids V/UV decision
  - Harmonically related sinusoids model the voiced parts
  - Random band-pass signals model the unvoiced parts
    - White noise filtered by a group of band-pass filters (filterbank) with center frequencies  $k\omega_0$

- Sinusoids + band-pass random signals (Abrantes et al.1991 [2])
  - Completely avoids V/UV decision
  - Harmonically related sinusoids model the voiced parts
  - Random band-pass signals model the unvoiced parts
    - White noise filtered by a group of band-pass filters (filterbank) with center frequencies  $k\omega_s$

- Sinusoids + band-pass random signals (Abrantes et al.1991 [2])
  - Completely avoids V/UV decision
  - Harmonically related sinusoids model the voiced parts
  - Random band-pass signals model the unvoiced parts
    - White noise filtered by a group of band-pass filters (filterbank) with center frequencies  $k\omega_s$

# BACKGROUND

Sinusoids + band-pass random signals (Abrantes et al.1991 [2])





- Periodic + Aperiodic Decomposition (Yegnayarayana et al.1995 [3])
  - The LP residual signal is used as an approximation to the excitation
  - V/UV analysis is used
  - Frequency regions of harmonic and noise components in the spectral domain
  - An iterative algorithm is proposed which reconstructs the aperiodic component in the harmonic regions
  - The periodic component is obtained by subtracting the reconstructed aperiodic component signal from the residual signal.

- Periodic + Aperiodic Decomposition (Yegnayarayana et al.1995 [3])
  - The LP residual signal is used as an approximation to the excitation
  - V/UV analysis is used
  - Frequency regions of harmonic and noise components in the spectral domain
  - An iterative algorithm is proposed which reconstructs the aperiodic component in the harmonic regions
  - The periodic component is obtained by subtracting the reconstructed aperiodic component signal from the residual signal.

- Periodic + Aperiodic Decomposition (Yegnayarayana et al.1995 [3])
  - The LP residual signal is used as an approximation to the excitation
  - V/UV analysis is used
  - Frequency regions of harmonic and noise components in the spectral domain
  - An iterative algorithm is proposed which reconstructs the aperiodic component in the harmonic regions
  - The periodic component is obtained by subtracting the reconstructed aperiodic component signal from the residual signal.

- Periodic + Aperiodic Decomposition (Yegnayarayana et al.1995 [3])
  - The LP residual signal is used as an approximation to the excitation
  - V/UV analysis is used
  - Frequency regions of harmonic and noise components in the spectral domain
  - An iterative algorithm is proposed which reconstructs the aperiodic component in the harmonic regions
  - The periodic component is obtained by subtracting the reconstructed aperiodic component signal from the residual signal.

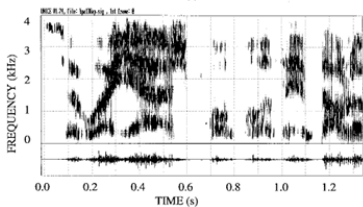
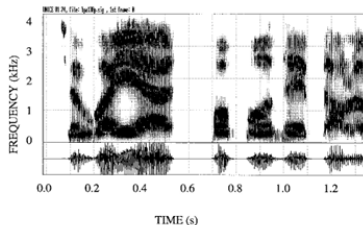
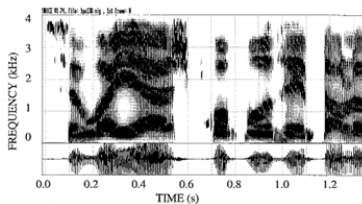
# BACKGROUND

- Periodic + Aperiodic Decomposition (Yegnayarayana et al.1995 [3])
  - The LP residual signal is used as an approximation to the excitation
  - V/UV analysis is used
  - Frequency regions of harmonic and noise components in the spectral domain
  - An iterative algorithm is proposed which reconstructs the aperiodic component in the harmonic regions
  - The periodic component is obtained by subtracting the reconstructed aperiodic component signal from the residual signal.

- Periodic + Aperiodic Decomposition (Yegnayarayana et al.1995 [3])
  - The LP residual signal is used as an approximation to the excitation
  - V/UV analysis is used
  - Frequency regions of harmonic and noise components in the spectral domain
  - An iterative algorithm is proposed which reconstructs the aperiodic component in the harmonic regions
  - The periodic component is obtained by subtracting the reconstructed aperiodic component signal from the residual signal.

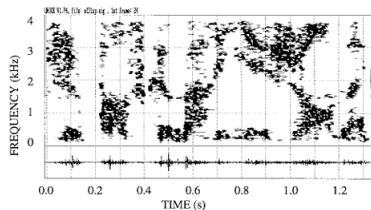
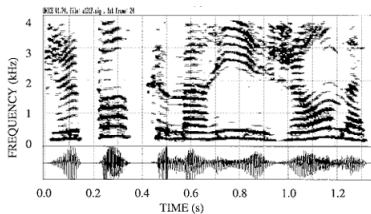
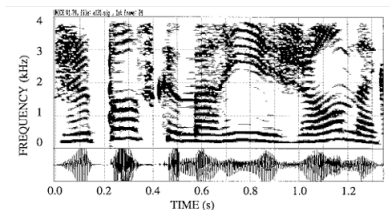
# BACKGROUND

Periodic + Aperiodic Decomposition (Yegnayarayana et al.1995 [3])



# BACKGROUND

Periodic + Aperiodic Decomposition (Yegnayarayana et al.1995 [3])





# WHY DECOMPOSE?

Decomposing speech into (quasi)periodic and non-periodic part has many applications in:

- Speech modification
- Speech coding
- Pathologic voice detection (i.e., HNR ...)
- Psychoacoustic research

# WHY DECOMPOSE?

Decomposing speech into (quasi)periodic and non-periodic part has many applications in:

- Speech modification
- Speech coding
- Pathologic voice detection (i.e., HNR ...)
- Psychoacoustic research

# WHY DECOMPOSE?

Decomposing speech into (quasi)periodic and non-periodic part has many applications in:

- Speech modification
- Speech coding
- Pathologic voice detection (i.e., HNR ...)
- Psychoacoustic research

# WHY DECOMPOSE?

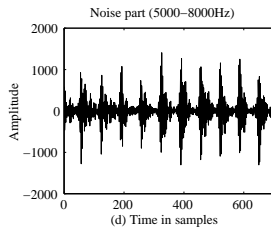
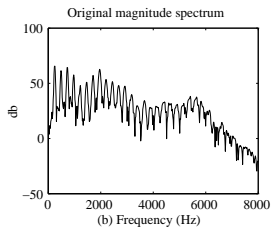
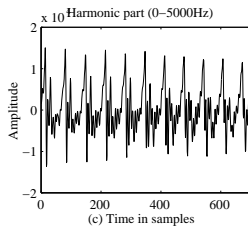
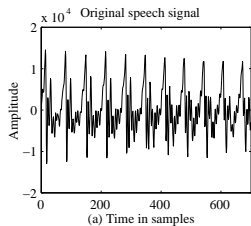
Decomposing speech into (quasi)periodic and non-periodic part has many applications in:

- Speech modification
- Speech coding
- Pathologic voice detection (i.e., HNR ...)
- Psychoacoustic research

# OUTLINE

- 1 FIRST WORKS ON SPEECH DECOMPOSITION...
- 2 INTRODUCTION TO HNMs
- 3 ANALYSIS
  - Frequency
  - Maximum Voiced Frequency
  - Amplitudes and Phases
    - Error Function - for  $\text{HNM}_1$
    - Least Squares - for  $\text{HNM}_1$
  - Residual
- 4 SYNTHESIS
- 5 ENERGY MODULATION FUNCTION
- 6 TOWARDS QUASI-HARMONICITY
- 7 THANKS
- 8 REFERENCES

# MOTIVATION FOR HNM



# BRIEF OVERVIEW OF HNM

- HNM (Stylianou 1995 [4]) is a pitch-synchronous harmonic plus noise representation of the speech signal.
- Speech spectrum is divided into a low and a high band delimited by the so-called *maximum voiced frequency*
- The *lower* band of the spectrum (below the maximum voiced frequency) is represented solely by harmonically related sine waves.
- The *upper* band is modeled as a noise component modulated by a time-domain amplitude envelope.
- HNM allows high-quality copy synthesis and prosodic modifications.

# BRIEF OVERVIEW OF HNM

- HNM (Stylianou 1995 [4]) is a pitch-synchronous harmonic plus noise representation of the speech signal.
- Speech spectrum is divided into a low and a high band delimited by the so-called *maximum voiced frequency*
- The *lower* band of the spectrum (below the maximum voiced frequency) is represented solely by harmonically related sine waves.
- The *upper* band is modeled as a noise component modulated by a time-domain amplitude envelope.
- HNM allows high-quality copy synthesis and prosodic modifications.



# BRIEF OVERVIEW OF HNM

- HNM (Stylianou 1995 [4]) is a pitch-synchronous harmonic plus noise representation of the speech signal.
- Speech spectrum is divided into a low and a high band delimited by the so-called *maximum voiced frequency*
- The *lower* band of the spectrum (below the maximum voiced frequency) is represented solely by harmonically related sine waves.
- The *upper* band is modeled as a noise component modulated by a time-domain amplitude envelope.
- HNM allows high-quality copy synthesis and prosodic modifications.

# BRIEF OVERVIEW OF HNM

- HNM (Stylianou 1995 [4]) is a pitch-synchronous harmonic plus noise representation of the speech signal.
- Speech spectrum is divided into a low and a high band delimited by the so-called *maximum voiced frequency*
- The *lower* band of the spectrum (below the maximum voiced frequency) is represented solely by harmonically related sine waves.
- The *upper* band is modeled as a noise component modulated by a time-domain amplitude envelope.
- HNM allows high-quality copy synthesis and prosodic modifications.

# BRIEF OVERVIEW OF HNM

- HNM (Stylianou 1995 [4]) is a pitch-synchronous harmonic plus noise representation of the speech signal.
- Speech spectrum is divided into a low and a high band delimited by the so-called *maximum voiced frequency*
- The *lower* band of the spectrum (below the maximum voiced frequency) is represented solely by harmonically related sine waves.
- The *upper* band is modeled as a noise component modulated by a time-domain amplitude envelope.
- HNM allows high-quality copy synthesis and prosodic modifications.

- Harmonic part:

$$h(t) = \sum_{k=-L(t)}^{L(t)} A_k(t) e^{j2\pi k f_0(t) t}$$

where  $A_k(t)$  and  $f_0(t)$  are the instantaneous complex amplitude and real frequency, respectively

- Noise part:

$$n(t) = e(t) [v(\tau, t) \star g(t)]$$

where  $e(t)$ ,  $v(\tau, t)$ ,  $g(t)$  are a time envelope, an estimation of the PSD (filter), and white gaussian noise, respectively

- Speech:

$$s(t) = h(t) + n(t)$$

- Harmonic part:

$$h(t) = \sum_{k=-L(t)}^{L(t)} A_k(t) e^{j2\pi k f_0(t) t}$$

where  $A_k(t)$  and  $f_0(t)$  are the instantaneous complex amplitude and real frequency, respectively

- Noise part:

$$n(t) = e(t) [v(\tau, t) \star g(t)]$$

where  $e(t)$ ,  $v(\tau, t)$ ,  $g(t)$  are a time envelope, an estimation of the PSD (filter), and white gaussian noise, respectively

- Speech:

$$s(t) = h(t) + n(t)$$

- Harmonic part:

$$h(t) = \sum_{k=-L(t)}^{L(t)} A_k(t) e^{j2\pi k f_0(t) t}$$

where  $A_k(t)$  and  $f_0(t)$  are the instantaneous complex amplitude and real frequency, respectively

- Noise part:

$$n(t) = e(t) [v(\tau, t) \star g(t)]$$

where  $e(t)$ ,  $v(\tau, t)$ ,  $g(t)$  are a time envelope, an estimation of the PSD (filter), and white gaussian noise, respectively

- Speech:

$$s(t) = h(t) + n(t)$$

# MODELS FOR PERIODIC PART

- HNM<sub>1</sub>: Sum of exponential functions without slope

$$h_1[n] = \sum_{k=-L(n_a^i)}^{L(n_a^i)} a_k(n_a^i) e^{j2\pi k f_0(n_a^i)(n-n_a^i)}$$

- HNM<sub>2</sub>: Sum of exponential function with complex slope

$$h_2[n] = \Re \left\{ \sum_{k=1}^{L(n_a^i)} A_k(n) e^{j2\pi k f_0(n_a^i)(n-n_a^i)} \right\}$$

where

$$A_k(n) = a_k(n_a^i) + (n - n_a^i) b_k(n_a^i)$$

with  $a_k(n_a^i)$ ,  $b_k(n_a^i)$  to be complex numbers (amplitude and slope respectively).  $\Re$  denotes taking the real part.

# MODELS FOR PERIODIC PART

- HNM<sub>1</sub>: Sum of exponential functions without slope

$$h_1[n] = \sum_{k=-L(n_a^i)}^{L(n_a^i)} a_k(n_a^i) e^{j2\pi k f_0(n_a^i)(n-n_a^i)}$$

- HNM<sub>2</sub>: Sum of exponential function with complex slope

$$h_2[n] = \Re \left\{ \sum_{k=1}^{L(n_a^i)} A_k(n) e^{j2\pi k f_0(n_a^i)(n-n_a^i)} \right\}$$

where

$$A_k(n) = a_k(n_a^i) + (n - n_a^i) b_k(n_a^i)$$

with  $a_k(n_a^i)$ ,  $b_k(n_a^i)$  to be complex numbers (amplitude and slope respectively).  $\Re$  denotes taking the real part.



# MODELS FOR PERIODIC PART

- HNM<sub>3</sub>: Sum of sinusoids with time-varying real amplitudes

$$h_3[n] = \sum_{k=0}^{L(n_a^i)} a_k(n) \cos(\varphi_k(n))$$

where

$$\begin{aligned} a_k(n) &= c_{k0} + c_{k1} (n - n_a^i)^1 + \dots + c_{kp} (n - n_a^i)^{p(n)} \\ \varphi_k(n) &= \epsilon_k + 2\pi k\zeta (n - n_a^i) \end{aligned}$$

where  $a_k(n)$ ,  $\phi_k(n)$  are real functions of discrete time and  $p(n)$  is the order of the amplitude polynomial, which is, in general, a time-varying parameter.

## RESIDUAL (NOISE) PART

The non-periodic part is just the *residual* signal obtained by subtracting the periodic-part (harmonic part) from the original speech signal in the time-domain

$$r[n] = s[n] - h[n]$$

where  $h[n]$  is either  $h_1[n]$ ,  $h_2[n]$ , or  $h_3[n]$  (harmonic part of HNM<sub>1</sub>, HNM<sub>2</sub>, and HNM<sub>3</sub>, respectively).

# OUTLINE

- 1 FIRST WORKS ON SPEECH DECOMPOSITION...
- 2 INTRODUCTION TO HNMs
- 3 ANALYSIS
  - Frequency
  - Maximum Voiced Frequency
  - Amplitudes and Phases
    - Error Function - for  $\text{HNM}_1$
    - Least Squares - for  $\text{HNM}_1$
  - Residual
- 4 SYNTHESIS
- 5 ENERGY MODULATION FUNCTION
- 6 TOWARDS QUASI-HARMONICITY
- 7 THANKS
- 8 REFERENCES

# INITIAL FUNDAMENTAL FREQUENCY

- Get an initial estimation of fundamental frequency  $f_0$  [5]
- Determine the voicing of the frame using normalized error over first four harmonics:

$$E = \frac{\int_{0.7f_0}^{4.3f_0} (|S(f)| - |\tilde{S}(f)|)^2}{\int_{0.7f_0}^{4.3f_0} |S(f)|^2}$$

where  $\tilde{S}(f)$  is a synthetic DFT-based spectrum using the initial  $f_0$  estimation

- If  $E < T$ , where  $T$  an appropriate threshold (e.g.  $-15$  dB), then frame is voiced, else it is labeled as unvoiced

# INITIAL FUNDAMENTAL FREQUENCY

- Get an initial estimation of fundamental frequency  $f_0$  [5]
- Determine the voicing of the frame using normalized error over first four harmonics:

$$E = \frac{\int_{0.7f_0}^{4.3f_0} (|S(f)| - |\tilde{S}(f)|)^2}{\int_{0.7f_0}^{4.3f_0} |S(f)|^2}$$

where  $\tilde{S}(f)$  is a synthetic DFT-based spectrum using the initial  $f_0$  estimation

- If  $E < T$ , where  $T$  an appropriate threshold (e.g.  $-15$  dB), then frame is voiced, else it is labeled as unvoiced

# INITIAL FUNDAMENTAL FREQUENCY

- Get an initial estimation of fundamental frequency  $f_0$  [5]
- Determine the voicing of the frame using normalized error over first four harmonics:

$$E = \frac{\int_{0.7f_0}^{4.3f_0} (|S(f)| - |\tilde{S}(f)|)^2}{\int_{0.7f_0}^{4.3f_0} |S(f)|^2}$$

where  $\tilde{S}(f)$  is a synthetic DFT-based spectrum using the initial  $f_0$  estimation

- If  $E < T$ , where  $T$  an appropriate threshold (e.g.  $-15$  dB), then frame is voiced, else it is labeled as unvoiced

# MAXIMUM VOICED FREQUENCY - MVF

- The MVF  $F_M$  is determined frame-wise from the speech spectrum
- Starting from the frequency  $f_c$  of the maximum spectral peak,  $A_m$ , in  $[f_0/2, 3f_0/2]$ , spectral peak values are collected around that maximum peak, along with their frequencies
- The range of collection is  $R_{search} = [f_c - f_0/2, f_c + f_0/2]$
- Determine peak frequencies  $f_i$  in  $R_{search}$ , and the corresponding amplitudes,  $A(f_i)$  and cumulative amplitudes  $A_c(f_i)$
- Cumulative amplitude  $A_c(f)$  is the sum of all spectral peak values from previous valley to following valley

# MAXIMUM VOICED FREQUENCY - MVF

- The MVF  $F_M$  is determined frame-wise from the speech spectrum
- Starting from the frequency  $f_c$  of the maximum spectral peak,  $A_m$ , in  $[f_0/2, 3f_0/2]$ , spectral peak values are collected around that maximum peak, along with their frequencies
- The range of collection is  $R_{search} = [f_c - f_0/2, f_c + f_0/2]$
- Determine peak frequencies  $f_i$  in  $R_{search}$ , and the corresponding amplitudes,  $A(f_i)$  and cumulative amplitudes  $A_c(f_i)$
- Cumulative amplitude  $A_c(f)$  is the sum of all spectral peak values from previous valley to following valley



# MAXIMUM VOICED FREQUENCY - MVF

- The MVF  $F_M$  is determined frame-wise from the speech spectrum
- Starting from the frequency  $f_c$  of the maximum spectral peak,  $A_m$ , in  $[f_0/2, 3f_0/2]$ , spectral peak values are collected around that maximum peak, along with their frequencies
- The range of collection is  $R_{search} = [f_c - f_0/2, f_c + f_0/2]$
- Determine peak frequencies  $f_i$  in  $R_{search}$ , and the corresponding amplitudes,  $A(f_i)$  and cumulative amplitudes  $A_c(f_i)$
- Cumulative amplitude  $A_c(f)$  is the sum of all spectral peak values from previous valley to following valley

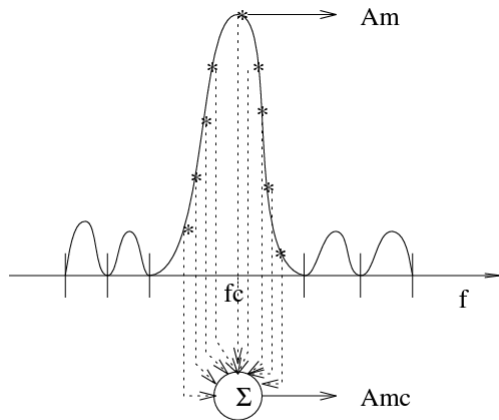
# MAXIMUM VOICED FREQUENCY - MVF

- The MVF  $F_M$  is determined frame-wise from the speech spectrum
- Starting from the frequency  $f_c$  of the maximum spectral peak,  $A_m$ , in  $[f_0/2, 3f_0/2]$ , spectral peak values are collected around that maximum peak, along with their frequencies
- The range of collection is  $R_{search} = [f_c - f_0/2, f_c + f_0/2]$
- Determine peak frequencies  $f_i$  in  $R_{search}$ , and the corresponding amplitudes,  $A(f_i)$  and cumulative amplitudes  $A_c(f_i)$
- Cumulative amplitude  $A_c(f)$  is the sum of all spectral peak values from previous valley to following valley

# MAXIMUM VOICED FREQUENCY - MVF

- The MVF  $F_M$  is determined frame-wise from the speech spectrum
- Starting from the frequency  $f_c$  of the maximum spectral peak,  $A_m$ , in  $[f_0/2, 3f_0/2]$ , spectral peak values are collected around that maximum peak, along with their frequencies
- The range of collection is  $R_{search} = [f_c - f_0/2, f_c + f_0/2]$
- Determine peak frequencies  $f_i$  in  $R_{search}$ , and the corresponding amplitudes,  $A(f_i)$  and cumulative amplitudes  $A_c(f_i)$
- Cumulative amplitude  $A_c(f)$  is the sum of all spectral peak values from previous valley to following valley

# MAXIMUM VOICED FREQUENCY - MVF



**Fig. 1.** Cumulative amplitude definition

# MAXIMUM VOICED FREQUENCY - MVF

- Compute the average cumulative amplitude for all  $f_i$ :  $\bar{A}_c(f_i)$
- Pass  $f_c$  through the *voicing test* (see next slide)
- Search for the maximum spectral peak in  $[f_c + f_0/2, f_c + 3f_0/2]$ , and find new  $f_c$
- Repeat the steps until  $f_c \leq f_s/2$ .
- Determine voiced and unvoiced spectral areas
- Maximum voiced frequency  $M_F$  is the maximum frequency of the last voiced spectral area.

# MAXIMUM VOICED FREQUENCY - MVF

- Compute the average cumulative amplitude for all  $f_i$ :  $\bar{A}_c(f_i)$
- Pass  $f_c$  through the *voicing test* (see next slide)
- Search for the maximum spectral peak in  $[f_c + f_0/2, f_c + 3f_0/2]$ , and find new  $f_c$
- Repeat the steps until  $f_c \leq f_s/2$ .
- Determine voiced and unvoiced spectral areas
- Maximum voiced frequency  $M_F$  is the maximum frequency of the last voiced spectral area.

# MAXIMUM VOICED FREQUENCY - MVF

- Compute the average cumulative amplitude for all  $f_i$ :  $\bar{A}_c(f_i)$
- Pass  $f_c$  through the *voicing test* (see next slide)
- Search for the maximum spectral peak in  $[f_c + f_0/2, f_c + 3f_0/2]$ , and find new  $f_c$
- Repeat the steps until  $f_c \leq f_s/2$ .
- Determine voiced and unvoiced spectral areas
- Maximum voiced frequency  $M_F$  is the maximum frequency of the last voiced spectral area.

# MAXIMUM VOICED FREQUENCY - MVF

- Compute the average cumulative amplitude for all  $f_i$ :  $\bar{A}_c(f_i)$
- Pass  $f_c$  through the *voicing test* (see next slide)
- Search for the maximum spectral peak in  $[f_c + f_0/2, f_c + 3f_0/2]$ , and find new  $f_c$
- Repeat the steps until  $f_c \leq f_s/2$ .
- Determine voiced and unvoiced spectral areas
- Maximum voiced frequency  $M_F$  is the maximum frequency of the last voiced spectral area.



# MAXIMUM VOICED FREQUENCY - MVF

- Compute the average cumulative amplitude for all  $f_i$ :  $\bar{A}_c(f_i)$
- Pass  $f_c$  through the *voicing test* (see next slide)
- Search for the maximum spectral peak in  $[f_c + f_0/2, f_c + 3f_0/2]$ , and find new  $f_c$
- Repeat the steps until  $f_c \leq f_s/2$ .
- Determine voiced and unvoiced spectral areas
- Maximum voiced frequency  $M_F$  is the maximum frequency of the last voiced spectral area.

# MAXIMUM VOICED FREQUENCY - MVF

- Compute the average cumulative amplitude for all  $f_i$ :  $\bar{A}_c(f_i)$
- Pass  $f_c$  through the *voicing test* (see next slide)
- Search for the maximum spectral peak in  $[f_c + f_0/2, f_c + 3f_0/2]$ , and find new  $f_c$
- Repeat the steps until  $f_c \leq f_s/2$ .
- Determine voiced and unvoiced spectral areas
- Maximum voiced frequency  $M_F$  is the maximum frequency of the last voiced spectral area.

# VOICING TEST

Voicing Test:

- If

$$\frac{A_c}{\bar{A}_c(f_i)} > 2$$

or

$$|A - \max \{A(f_i)\}| > 13 \text{ dB}$$

then

- if  $f_c$  is *really* close to the closest harmonic  $lf_0$ , then
- declare  $f_c$  as voiced frequency. Otherwise, declare  $f_c$  as unvoiced frequency.

# VOICING TEST

Voicing Test:

- If

$$\frac{A_c}{\bar{A}_c(f_i)} > 2$$

or

$$|A - \max \{A(f_i)\}| > 13 \text{ dB}$$

then

- if  $f_c$  is *really* close to the closest harmonic  $lf_0$ , then
- declare  $f_c$  as voiced frequency. Otherwise, declare  $f_c$  as unvoiced frequency.

# VOICING TEST

Voicing Test:

- If

$$\frac{A_c}{\bar{A}_c(f_i)} > 2$$

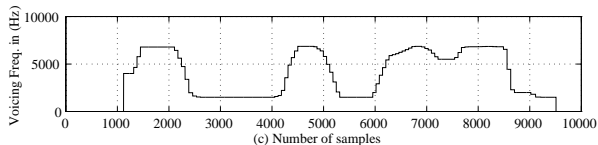
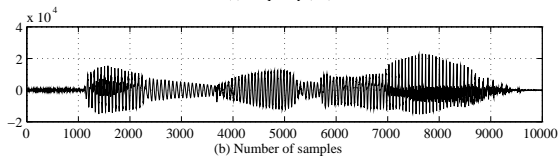
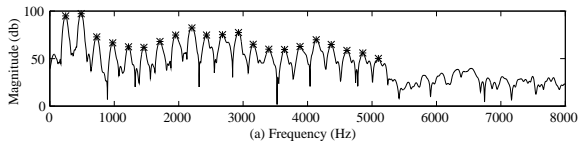
or

$$|A - \max \{A(f_i)\}| > 13 \text{ dB}$$

then

- if  $f_c$  is *really* close to the closest harmonic  $lf_0$ , then
- declare  $f_c$  as voiced frequency. Otherwise, declare  $f_c$  as unvoiced frequency.

# MAXIMUM VOICED FREQUENCY EXAMPLE

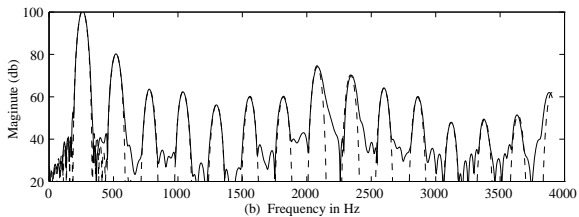
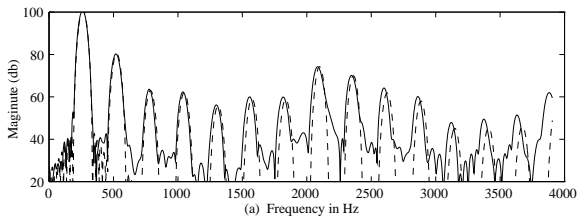


# FUNDAMENTAL FREQUENCY REFINEMENT

Using the initial  $f_0$  value and the  $L$  detected voiced frequencies  $f_i$ , then the refined fundamental frequency,  $\hat{f}_0$  is defined as the value that minimizes the error:

$$E(\hat{f}_0) = \sum_{i=1}^L |f_i - i \cdot \hat{f}_0|^2$$

# REFINEMENT FREQUENCY EXAMPLE





# AMPLITUDES AND PHASES ESTIMATION

Having  $f_0$  estimated for voiced frames, amplitudes and phases are estimated by minimizing the criterion:

$$\epsilon = \sum_{n=n_a^i-N}^{n_a^i+N} w^2[n](s[n] - \hat{h}[n])^2$$

where  $n_a^i = n_a^{i-1} + P(n_a^{i-1})$ , and  $P(n_a^{i-1})$  denotes the pitch period at  $n_a^{i-1}$ .

- for HNM<sub>1</sub> and HNM<sub>2</sub>, this criterion has a quadratic form and is solved by inverting an over-determined system of linear equations.
- For HNM<sub>3</sub>, however, a non-linear system of equations has to be solved.

# AMPLITUDES AND PHASES ESTIMATION

Having  $f_0$  estimated for voiced frames, amplitudes and phases are estimated by minimizing the criterion:

$$\epsilon = \sum_{n=n_a^i-N}^{n_a^i+N} w^2[n](s[n] - \hat{h}[n])^2$$

where  $n_a^i = n_a^{i-1} + P(n_a^{i-1})$ , and  $P(n_a^{i-1})$  denotes the pitch period at  $n_a^{i-1}$ .

- for HNM<sub>1</sub> and HNM<sub>2</sub>, this criterion has a quadratic form and is solved by inverting an over-determined system of linear equations.
- For HNM<sub>3</sub>, however, a non-linear system of equations has to be solved.

# REFORMULATE THE ERROR FUNCTION - FOR HNM<sub>1</sub>

Cost function:

$$\epsilon(a_{-L}, \dots, a_L, f_0) = \frac{1}{2} \sum_{n=-N}^N (e[n])^2 = \frac{1}{2} \mathbf{e}^h \mathbf{e}$$

where

$$e[n] = w[n](s[n] - h[n])$$

or

$$\mathbf{e} = [e[-N], e[-N+1], \dots, e[N]]^T$$

# REFORMULATE THE ERROR FUNCTION - FOR HNM<sub>1</sub>

$$\epsilon(\mathbf{a}) = \frac{1}{2}(\mathbf{s} - \mathbf{E}\mathbf{a})^h \mathbf{W}^2 (\mathbf{s} - \mathbf{E}\mathbf{a})$$

where

$$\mathbf{a} = [a_{-L}, \dots, a_0, \dots, a_L]^T$$

and

$$\mathbf{E} = \begin{bmatrix} e^{j2\pi(-L)\hat{f}_0(-N)/f_s}, & \dots & e^{j2\pi L\hat{f}_0(-N)/f_s} \\ e^{j2\pi(-L)\hat{f}_0(-N+1)/f_s}, & \dots & e^{j2\pi L\hat{f}_0(-N+1)/f_s} \\ \vdots & \vdots & \vdots \\ e^{j2\pi(-L)\hat{f}_0 N/f_s}, & \dots & e^{j2\pi L\hat{f}_0 N/f_s} \end{bmatrix}^T \quad (2L+1 \times 2N+1)$$

# LEAST SQUARES - FOR HNM<sub>1</sub>

- Setting:

$$\frac{\partial \epsilon(\mathbf{a})}{\partial \mathbf{a}} = 0 \implies \mathbf{E}^h \mathbf{W}^2 \mathbf{E} \mathbf{a} - \mathbf{E}^h \mathbf{W}^2 \mathbf{s} = 0$$

- Solution:

$$\mathbf{a}_{LS} = (\mathbf{E}^h \mathbf{W}^2 \mathbf{E})^{-1} \mathbf{E}^h \mathbf{W}^2 \mathbf{s}$$

- Properties:

- Asymptotically efficient even when the noise is colored.
- Rather fast,  $O(L(N+L))$ .
- Assumes no errors in  $E$  matrix.

# LEAST SQUARES - FOR HNM<sub>1</sub>

- Setting:

$$\frac{\partial \epsilon(\mathbf{a})}{\partial \mathbf{a}} = 0 \implies \mathbf{E}^h \mathbf{W}^2 \mathbf{E} \mathbf{a} - \mathbf{E}^h \mathbf{W}^2 \mathbf{s} = 0$$

- Solution:

$$\mathbf{a}_{LS} = (\mathbf{E}^h \mathbf{W}^2 \mathbf{E})^{-1} \mathbf{E}^h \mathbf{W}^2 \mathbf{s}$$

- Properties:

- Asymptotically efficient even when the noise is colored.
- Rather fast,  $O(L(N+L))$ .
- Assumes no errors in  $E$  matrix.

# LEAST SQUARES - FOR HNM<sub>1</sub>

- Setting:

$$\frac{\partial \epsilon(\mathbf{a})}{\partial \mathbf{a}} = 0 \implies \mathbf{E}^h \mathbf{W}^2 \mathbf{E} \mathbf{a} - \mathbf{E}^h \mathbf{W}^2 \mathbf{s} = 0$$

- Solution:

$$\mathbf{a}_{LS} = (\mathbf{E}^h \mathbf{W}^2 \mathbf{E})^{-1} \mathbf{E}^h \mathbf{W}^2 \mathbf{s}$$

- Properties:

- Asymptotically efficient even when the noise is colored.
- Rather fast,  $O(L(N + L))$ .
- Assumes no errors in  $E$  matrix.

# LEAST SQUARES - FOR HNM<sub>1</sub>

- Setting:

$$\frac{\partial \epsilon(\mathbf{a})}{\partial \mathbf{a}} = 0 \implies \mathbf{E}^h \mathbf{W}^2 \mathbf{E} \mathbf{a} - \mathbf{E}^h \mathbf{W}^2 \mathbf{s} = 0$$

- Solution:

$$\mathbf{a}_{LS} = (\mathbf{E}^h \mathbf{W}^2 \mathbf{E})^{-1} \mathbf{E}^h \mathbf{W}^2 \mathbf{s}$$

- Properties:

- Asymptotically efficient even when the noise is colored.
- Rather fast,  $O(L(N + L))$ .
- Assumes no errors in  $E$  matrix.



# LEAST SQUARES - FOR HNM<sub>1</sub>

- Setting:

$$\frac{\partial \epsilon(\mathbf{a})}{\partial \mathbf{a}} = 0 \implies \mathbf{E}^h \mathbf{W}^2 \mathbf{E} \mathbf{a} - \mathbf{E}^h \mathbf{W}^2 \mathbf{s} = 0$$

- Solution:

$$\mathbf{a}_{LS} = (\mathbf{E}^h \mathbf{W}^2 \mathbf{E})^{-1} \mathbf{E}^h \mathbf{W}^2 \mathbf{s}$$

- Properties:

- Asymptotically efficient even when the noise is colored.
- Rather fast,  $O(L(N + L))$ .
- Assumes no errors in  $E$  matrix.

# LEAST SQUARES - FOR HNM<sub>1</sub>

- Setting:

$$\frac{\partial \epsilon(\mathbf{a})}{\partial \mathbf{a}} = 0 \implies \mathbf{E}^h \mathbf{W}^2 \mathbf{E} \mathbf{a} - \mathbf{E}^h \mathbf{W}^2 \mathbf{s} = 0$$

- Solution:

$$\mathbf{a}_{LS} = (\mathbf{E}^h \mathbf{W}^2 \mathbf{E})^{-1} \mathbf{E}^h \mathbf{W}^2 \mathbf{s}$$

- Properties:

- Asymptotically efficient even when the noise is colored.
- Rather fast,  $O(L(N + L))$ .
- Assumes no errors in  $E$  matrix.

# AVOIDING ILL-CONDITIONING

- For  $\text{HNM}_1$  there is no problem if window length is twice the local pitch period
- Same thing for  $\text{HNM}_2$
- For  $\text{HNM}_3$  stands the same in case the maximum voiced frequency is less than  $3/4$  of the sampling frequency and order of amplitude polynomial is 2

# AVOIDING ILL-CONDITIONING

- For  $\text{HNM}_1$  there is no problem if window length is twice the local pitch period
- Same thing for  $\text{HNM}_2$
- For  $\text{HNM}_3$  stands the same in case the maximum voiced frequency is less than  $3/4$  of the sampling frequency and order of amplitude polynomial is 2

# AVOIDING ILL-CONDITIONING

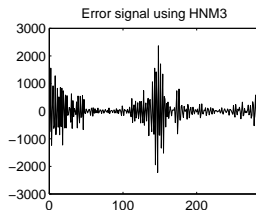
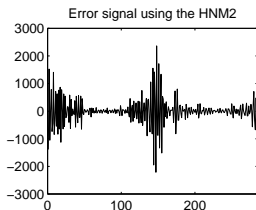
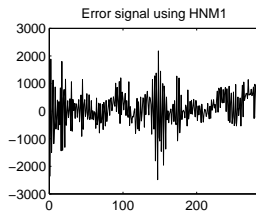
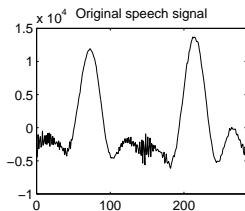
- For  $\text{HNM}_1$  there is no problem if window length is twice the local pitch period
- Same thing for  $\text{HNM}_2$
- For  $\text{HNM}_3$  stands the same in case the maximum voiced frequency is less than  $3/4$  of the sampling frequency and order of amplitude polynomial is 2

# RESIDUAL SIGNAL

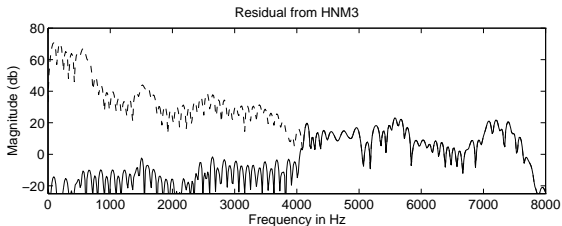
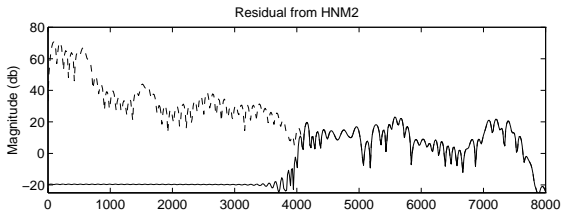
The residual signal  $r[n]$  is estimated by

$$\hat{r}[n] = s[n] - \hat{h}[n]$$

# TIME DOMAIN CHARACTERISTICS OF $\hat{r}[n]$

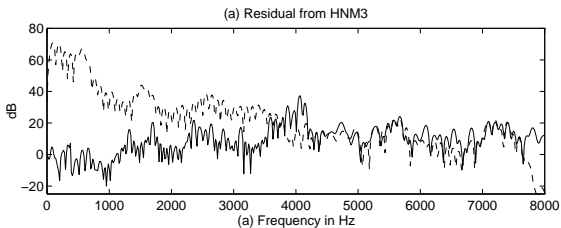
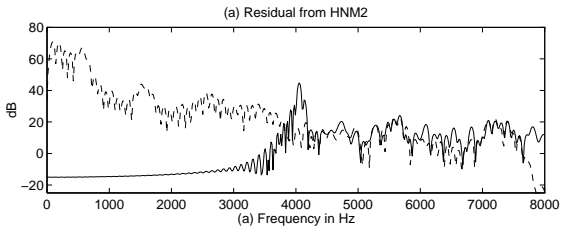


# SPECTRAL DOMAIN CHARACTERISTICS OF $\hat{r}[n]$

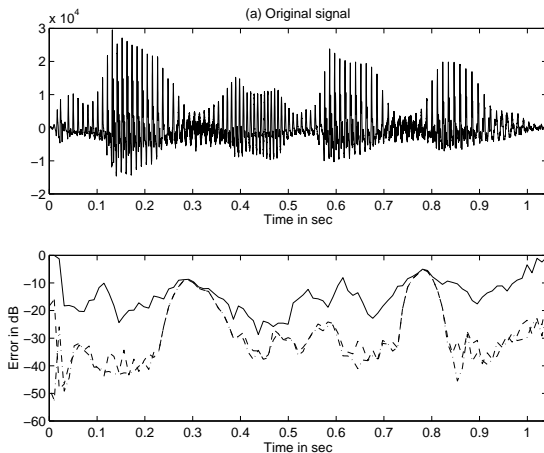




# ... AND AFTER ADDING NOISE



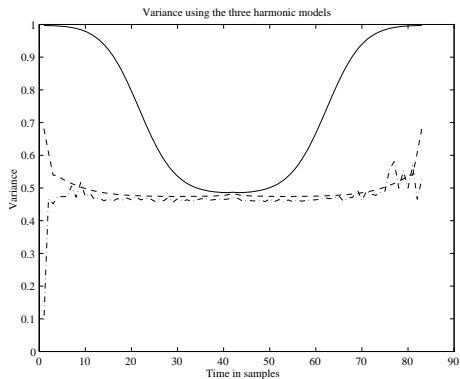
# MODELING ERROR



# VARIANCE OF THE RESIDUAL SIGNAL

The variance of the residual signal is given as:

$$E(\mathbf{r}\mathbf{r}^h) = \mathbf{I} - \mathbf{W}\mathbf{P}(\mathbf{P}^h\mathbf{W}^h\mathbf{W}\mathbf{P})^{-1}\mathbf{P}^h\mathbf{W}^h$$



# MODELING THE RESIDUAL SIGNAL

- Full bandwidth representation using a low-order (10th) AR filter
- Time-domain characteristics of the residual signal are modeled using deterministic functions

# OUTLINE

- 1 FIRST WORKS ON SPEECH DECOMPOSITION...
- 2 INTRODUCTION TO HNMs
- 3 ANALYSIS
  - Frequency
  - Maximum Voiced Frequency
  - Amplitudes and Phases
    - Error Function - for  $\text{HNM}_1$
    - Least Squares - for  $\text{HNM}_1$
  - Residual
- 4 SYNTHESIS
- 5 ENERGY MODULATION FUNCTION
- 6 TOWARDS QUASI-HARMONICITY
- 7 THANKS
- 8 REFERENCES

# FOR ALL HNMs

- $n_s^i \longleftrightarrow n_a^i$
- For the periodic part: Overlap-and-Add
- For the stochastic (noise) part):
  - Instead of AR coefficients we use reflection coefficients
  - Sample-by-sample filtering of Gaussian noise using normalized lattice filtering
  - Modulation in time with a deterministic function (i.e., triangular)

# FOR ALL HNMs

- $n_s^i \longleftrightarrow n_a^i$
- For the periodic part: Overlap-and-Add
- For the stochastic (noise) part):
  - Instead of AR coefficients we use reflection coefficients
  - Sample-by-sample filtering of Gaussian noise using normalized lattice filtering
  - Modulation in time with a deterministic function (i.e., triangular)

# FOR ALL HNMS

- $n_s^i \longleftrightarrow n_a^i$
- For the periodic part: Overlap-and-Add
- For the stochastic (noise) part):
  - Instead of AR coefficients we use reflection coefficients
  - Sample-by-sample filtering of Gaussian noise using normalized lattice filtering
  - Modulation in time with a deterministic function (i.e., triangular)



# FOR ALL HNMs

- $n_s^i \longleftrightarrow n_a^i$
- For the periodic part: Overlap-and-Add
- For the stochastic (noise) part):
  - Instead of AR coefficients we use reflection coefficients
  - Sample-by-sample filtering of Gaussian noise using normalized lattice filtering
  - Modulation in time with a deterministic function (i.e., triangular)

# FOR ALL HNMs

- $n_s^i \longleftrightarrow n_a^i$
- For the periodic part: Overlap-and-Add
- For the stochastic (noise) part):
  - Instead of AR coefficients we use reflection coefficients
  - Sample-by-sample filtering of Gaussian noise using normalized lattice filtering
  - Modulation in time with a deterministic function (i.e., triangular)

# FOR ALL HNMS

- $n_s^i \longleftrightarrow n_a^i$
- For the periodic part: Overlap-and-Add
- For the stochastic (noise) part):
  - Instead of AR coefficients we use reflection coefficients
  - Sample-by-sample filtering of Gaussian noise using normalized lattice filtering
  - Modulation in time with a deterministic function (i.e., triangular)

# FOR $\text{HNM}_1$ SPECIFICALLY

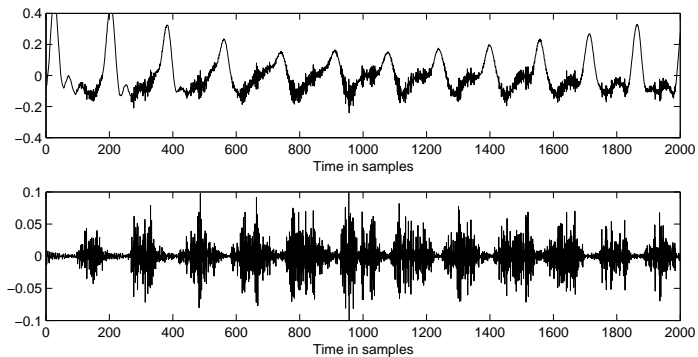
for Periodic part (as an alternative to OLA)

- Direct frequency matching
- Linear amplitude interpolation
- Linear phase interpolation using average pitch value

# OUTLINE

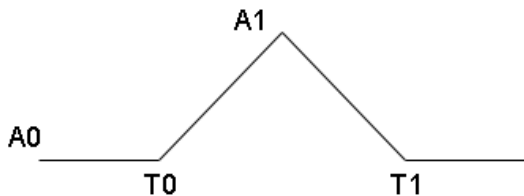
- 1 FIRST WORKS ON SPEECH DECOMPOSITION...
- 2 INTRODUCTION TO HNMs
- 3 ANALYSIS
  - Frequency
  - Maximum Voiced Frequency
  - Amplitudes and Phases
    - Error Function - for  $\text{HNM}_1$
    - Least Squares - for  $\text{HNM}_1$
  - Residual
- 4 SYNTHESIS
- 5 ENERGY MODULATION FUNCTION
- 6 TOWARDS QUASI-HARMONICITY
- 7 THANKS
- 8 REFERENCES

# AGAIN ON THE ENERGY MODULATION



# SO FAR, MAINLY

So far we mainly use the Triangular Envelope:



There are many ways to obtain the “envelope” of a signal, as:

- Hilbert Transform (analytic signal)
- Low-pass local energy (energy envelope):

$$e[n] = \frac{1}{2N+1} \sum_{k=-N}^N |r[n-k]|$$

where  $r[n]$  denotes the residual signal.



There are many ways to obtain the “envelope” of a signal, as:

- Hilbert Transform (analytic signal)
- Low-pass local energy (energy envelope):

$$e[n] = \frac{1}{2N+1} \sum_{k=-N}^N |r[n-k]|$$

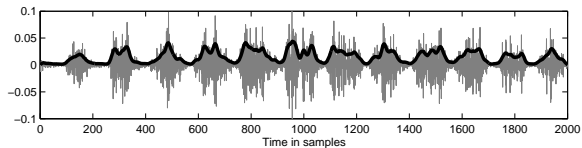
where  $r[n]$  denotes the residual signal.

We may also use the Hilbert envelope, computed as:

$$\tilde{e}_H[n] = \sum_{k=L-M+1}^L a_k e^{2\pi k(f_0/f_s)n}$$

# EXAMPLE OF ENERGY ENVELOPE

Example of Energy Envelope, with  $N = 7$



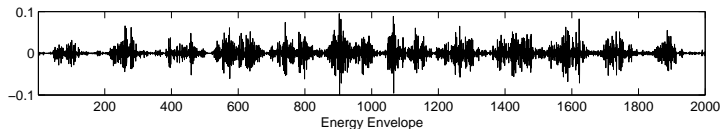
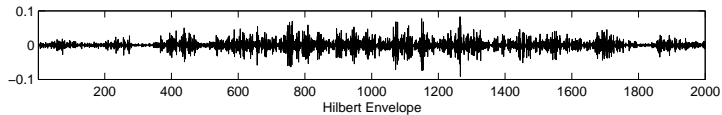
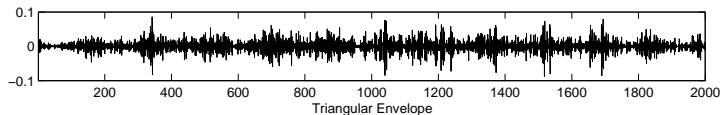
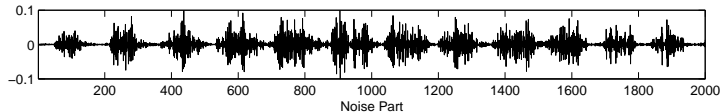
# ENERGY ENVELOPE

The energy envelope can be efficiently parameterized with a few Fourier coefficients:

$$\hat{e}[n] = \sum_{k=-L_e}^{L_e} A_k e^{j2\pi k(f_0/f_s)n}$$

where  $L_e$  is set to be 3 to 4

# LOOKING AT TIME DOMAIN PROPERTIES



# RESULTS FROM LISTENING TEST I

	Triangular	No pref.	Hilbert
Male	8 (8.3%)	43 (44.8%)	45 (46.9%)
Female	40 (41.7%)	47 (48.9%)	9 (9.4%)

	Hilbert	No pref.	Energy
Male	22 (22.9%)	47 (49.0%)	27 (28.1%)
Female	22 (22.9%)	54 (56.3%)	20 (20.8%)

	Energy	No pref.	Triangular
Male	43 (44.8%)	50 (52.0%)	3 (3.2%)
Female	16 (16.7%)	67 (69.8%)	13 (13.5%)

TABLE: Results from the listening test for the English sentences.

## RESULTS FROM LISTENING TEST II

	Triangular	No pref.	Hilbert
Male	10 (10.4%)	47 (49.0%)	39 (40.6%)
Female	8 (8.3%)	71 (74.0%)	17 (17.7%)

	Hilbert	No pref.	Energy
Male	11 (11.5%)	58 (60.4%)	27 (28.1%)
Female	13 (13.5%)	58 (60.4%)	25 (26.1%)

	Energy	No pref.	Triangular
Male	42 (43.7%)	48 (50.0%)	6 (6.3%)
Female	16 (16.7%)	68 (70.8%)	12 (12.5%)

TABLE: Results from the listening test for the French sentences.

# OUTLINE

- 1 FIRST WORKS ON SPEECH DECOMPOSITION...
- 2 INTRODUCTION TO HNMs
- 3 ANALYSIS
  - Frequency
  - Maximum Voiced Frequency
  - Amplitudes and Phases
    - Error Function - for  $\text{HNM}_1$
    - Least Squares - for  $\text{HNM}_1$
  - Residual
- 4 SYNTHESIS
- 5 ENERGY MODULATION FUNCTION
- 6 TOWARDS QUASI-HARMONICITY
- 7 THANKS
- 8 REFERENCES



# ESTIMATING SINUSOIDAL PARAMETERS

- Sinusoidal representation for a speech/signal frame:

$$x(t) = \left( \sum_{k=-K}^K a_k e^{j2\pi f_k t} \right) w(t)$$

- Methods:
  - FFT-based methods (i.e., QIFFT [Abe et al., 2004-05, [6] [7]])
  - Subspace methods
  - Least Squares (LS) method

- Frequency mismatch:

$$\hat{f}_k = f_k + \eta_k$$

- How to deal with that?
- You will discuss more advanced sinusoidal models in the following lecture! :)

# ESTIMATING SINUSOIDAL PARAMETERS

- Sinusoidal representation for a speech/signal frame:

$$x(t) = \left( \sum_{k=-K}^K a_k e^{j2\pi f_k t} \right) w(t)$$

- Methods:

- FFT-based methods (i.e., QIFFT [Abe et al., 2004-05, [6] [7]])
- Subspace methods
- Least Squares (LS) method

- Frequency mismatch:

$$\hat{f}_k = f_k + \eta_k$$

- How to deal with that?
- You will discuss more advanced sinusoidal models in the following lecture! :)

# ESTIMATING SINUSOIDAL PARAMETERS

- Sinusoidal representation for a speech/signal frame:

$$x(t) = \left( \sum_{k=-K}^K a_k e^{j2\pi f_k t} \right) w(t)$$

- Methods:
  - FFT-based methods (i.e., QIFFT [Abe et al., 2004-05, [6] [7]])
  - Subspace methods
  - Least Squares (LS) method
- Frequency mismatch:

$$\hat{f}_k = f_k + \eta_k$$

- How to deal with that?
- You will discuss more advanced sinusoidal models in the following lecture! :)

# ESTIMATING SINUSOIDAL PARAMETERS

- Sinusoidal representation for a speech/signal frame:

$$x(t) = \left( \sum_{k=-K}^K a_k e^{j2\pi f_k t} \right) w(t)$$

- Methods:
  - FFT-based methods (i.e., QIFFT [Abe et al., 2004-05, [6] [7]])
  - Subspace methods
  - Least Squares (LS) method
- Frequency mismatch:

$$\hat{f}_k = f_k + \eta_k$$

- How to deal with that?
- You will discuss more advanced sinusoidal models in the following lecture! :)

# ESTIMATING SINUSOIDAL PARAMETERS

- Sinusoidal representation for a speech/signal frame:

$$x(t) = \left( \sum_{k=-K}^K a_k e^{j2\pi f_k t} \right) w(t)$$

- Methods:
  - FFT-based methods (i.e., QIFFT [Abe et al., 2004-05, [6] [7]])
  - Subspace methods
  - Least Squares (LS) method
- Frequency mismatch:

$$\hat{f}_k = f_k + \eta_k$$

- How to deal with that?
- You will discuss more advanced sinusoidal models in the following lecture! :)

# ESTIMATING SINUSOIDAL PARAMETERS

- Sinusoidal representation for a speech/signal frame:

$$x(t) = \left( \sum_{k=-K}^K a_k e^{j2\pi f_k t} \right) w(t)$$

- Methods:
  - FFT-based methods (i.e., QIFFT [Abe et al., 2004-05, [6] [7]])
  - Subspace methods
  - Least Squares (LS) method
- Frequency mismatch:

$$\hat{f}_k = f_k + \eta_k$$

- How to deal with that?
- You will discuss more advanced sinusoidal models in the following lecture! :)

# ESTIMATING SINUSOIDAL PARAMETERS

- Sinusoidal representation for a speech/signal frame:

$$x(t) = \left( \sum_{k=-K}^K a_k e^{j2\pi f_k t} \right) w(t)$$

- Methods:
  - FFT-based methods (i.e., QIFFT [Abe et al., 2004-05, [6] [7]])
  - Subspace methods
  - Least Squares (LS) method
- Frequency mismatch:

$$\hat{f}_k = f_k + \eta_k$$

- How to deal with that?
- You will discuss more advanced sinusoidal models in the following lecture! :)

# ESTIMATING SINUSOIDAL PARAMETERS

- Sinusoidal representation for a speech/signal frame:

$$x(t) = \left( \sum_{k=-K}^K a_k e^{j2\pi f_k t} \right) w(t)$$

- Methods:
  - FFT-based methods (i.e., QIFFT [Abe et al., 2004-05, [6] [7]])
  - Subspace methods
  - Least Squares (LS) method
- Frequency mismatch:

$$\hat{f}_k = f_k + \eta_k$$

- How to deal with that?
- You will discuss more advanced sinusoidal models in the following lecture! :)



# OUTLINE





- 1 FIRST WORKS ON SPEECH DECOMPOSITION...
- 2 INTRODUCTION TO HNMs
- 3 ANALYSIS
  - Frequency
  - Maximum Voiced Frequency
  - Amplitudes and Phases
    - Error Function - for  $\text{HNM}_1$
    - Least Squares - for  $\text{HNM}_1$
  - Residual
- 4 SYNTHESIS
- 5 ENERGY MODULATION FUNCTION
- 6 TOWARDS QUASI-HARMONICITY
- 7 THANKS
- 8 REFERENCES

THANK YOU  
for your attention




# OUTLINE

- 1 FIRST WORKS ON SPEECH DECOMPOSITION...
- 2 INTRODUCTION TO HNMs
- 3 ANALYSIS
  - Frequency
  - Maximum Voiced Frequency
  - Amplitudes and Phases
    - Error Function - for  $\text{HNM}_1$
    - Least Squares - for  $\text{HNM}_1$
  - Residual
- 4 SYNTHESIS
- 5 ENERGY MODULATION FUNCTION
- 6 TOWARDS QUASI-HARMONICITY
- 7 THANKS
- 8 REFERENCES

# REFERENCES I

-  D. Griffin and J. Lim, "Multiband-excitation vocoder," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-36, pp. 236–243, Feb 1988.
-  A. Abrantes, J. Marques, and I. Transcoso, "Hybrid sinusoidal modeling of speech without voicing decision," *Eurospeech-91*, pp. 231–234, 1991.
-  B. Yegnanarayana, C. d'Alessandro, and V. Darsinos, "An iterative algorithm for decomposition of speech signals into periodic and aperiodic components," *IEEE Trans. Speech and Audio Processing*, vol. 6, no. 1, 1998.
-  Y. Stylianou, *Harmonic plus Noise Models for Speech, combined with Statistical Methods, for Speech and Speaker Modification*. PhD thesis, Ecole Nationale Supérieure des Télécommunications, Jan 1996.

# REFERENCES II

-  W. Hess, *Pitch determination of Speech Signals: Algorithmes and Devices*.  
Berlin: Springer, 1983.
-  M. Abe and J. S. III, "CQIFFT: Correcting Bias in a Sinusoidal Parameter Estimator based on Quadratic Interpolation of FFT Magnitude Peaks," Tech. Rep. STAN-M-117, Stanford University, California, Oct 2004.
-  M. Abe and J. S. III, "AM/FM Estimation for Time-varying Sinusoidal Modeling," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Processing*, (Philadelphia), pp. III 201–204, 2005.

