

Extracción de Características en el Procesamiento Digital de una Señal para el Mejoramiento del Reconocimiento Automático de Habla usando Wavelets

Jorge Luis Guevara Díaz

Universidad Nacional de Trujillo, Escuela de Informática,
Trujillo, Perú
jorge.jorjasso@gmail.com

and

Juan Orlando Salazar Campos

Universidad Nacional de Trujillo, Escuela de Informática,
Trujillo, Perú
josco_orlando@hotmail.com

Que veremos?

1. Introducción
2. Trabajos Previos
3. Procesamiento de la Señal
4. Coeficientes MFCC
5. Transformada Wavelet
6. Extracción de características usando wavelets
7. Experimentos y Resultados
8. Conclusiones

1. Introducción

Podríamos conversar con las maquinas como lo hacemos con los humanos?

hay pamele ya te
formateé hace dos
dias



y como te iba
contando necesito
una formateada



1. Introducción

- Speech Recognition

¿Cómo hacer que las computadoras puedan convertir a texto la palabra hablada?

Problemas:

algoritmos de bajo costo computacional
extracción de “buenas” características
correcta clasificación

1. Introducción

- Extracción de características

Complejidad computacional

Cual es la mejor representación de características?

Reducción de la dimensionalidad

conjunto mas pequeño que contenga la información mas esencial presente en los atributos originales

$$\mathbf{x} = (x_1, \dots, x_n) \mapsto \phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_d(\mathbf{x})), \quad d < n,$$

2. Trabajos Previos

2. Trabajos Previos

- Diversas Técnicas

Bandos de Energias de la Trasformada de Fourier

LPC Coeficientes de Prediccion Lineal [Atal, and Schroeder]

LPC-Cepstrum [Atal, and Schroeder] [Bogert and Tukey]

PLP Coeficientes de Predicción Lineal Perceptuales

MFCC Coeficientes Cepstrales en Frecuencia Mel [Davis and Mermelstein]

Propuestas Basada en Wavelets

3. Procesamiento Digital de la Señal de **Habla**

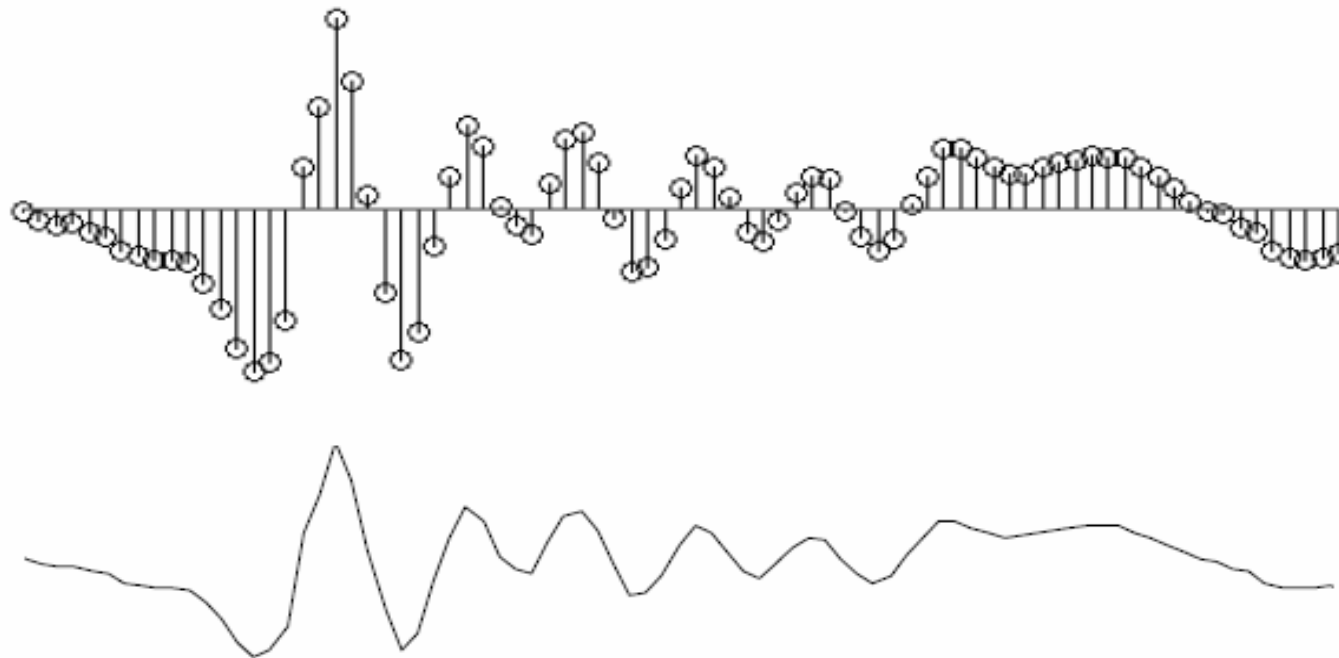
3. Procesamiento Digital de la Señal de Habla

- Diversos algoritmos para procesar la señal digital de habla en una computadora
- Ejemplo:
Eliminación de ruido, análisis de frecuencias, etc

3. Procesamiento Digital de la Señal de Habla

Capturar la señal analógica y digitalizarla para poder usarla en la computadora

$$x[n] = x_0(nT).$$



3. Procesamiento Digital de la Señal de Habla

3.1 Transformada de Fourier

$$X(\omega) = \int_{-\infty}^{\infty} x(t) e^{-j\omega t} dt$$

Donde :

$$e^{j\phi} = \cos \phi + j \sin \phi$$

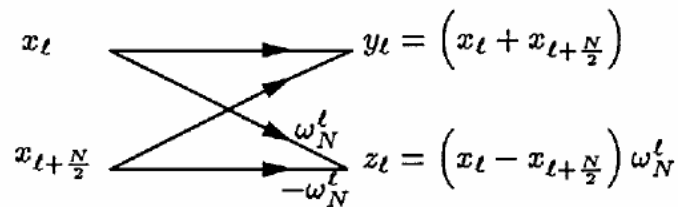
Transformada Discreta de Fourier.

$$X[k] = \sum_{n=0}^{N-1} x[n] e^{-j2\pi nk/N}$$

Complejidad computacional : $O(n^2)$

Transformada Rápida de Fourier.

Algoritmo radix-2 con diezmado en frecuencia y reordenamiento de la salida de bits mezclados, cuya complejidad es $O(n \log n)$.

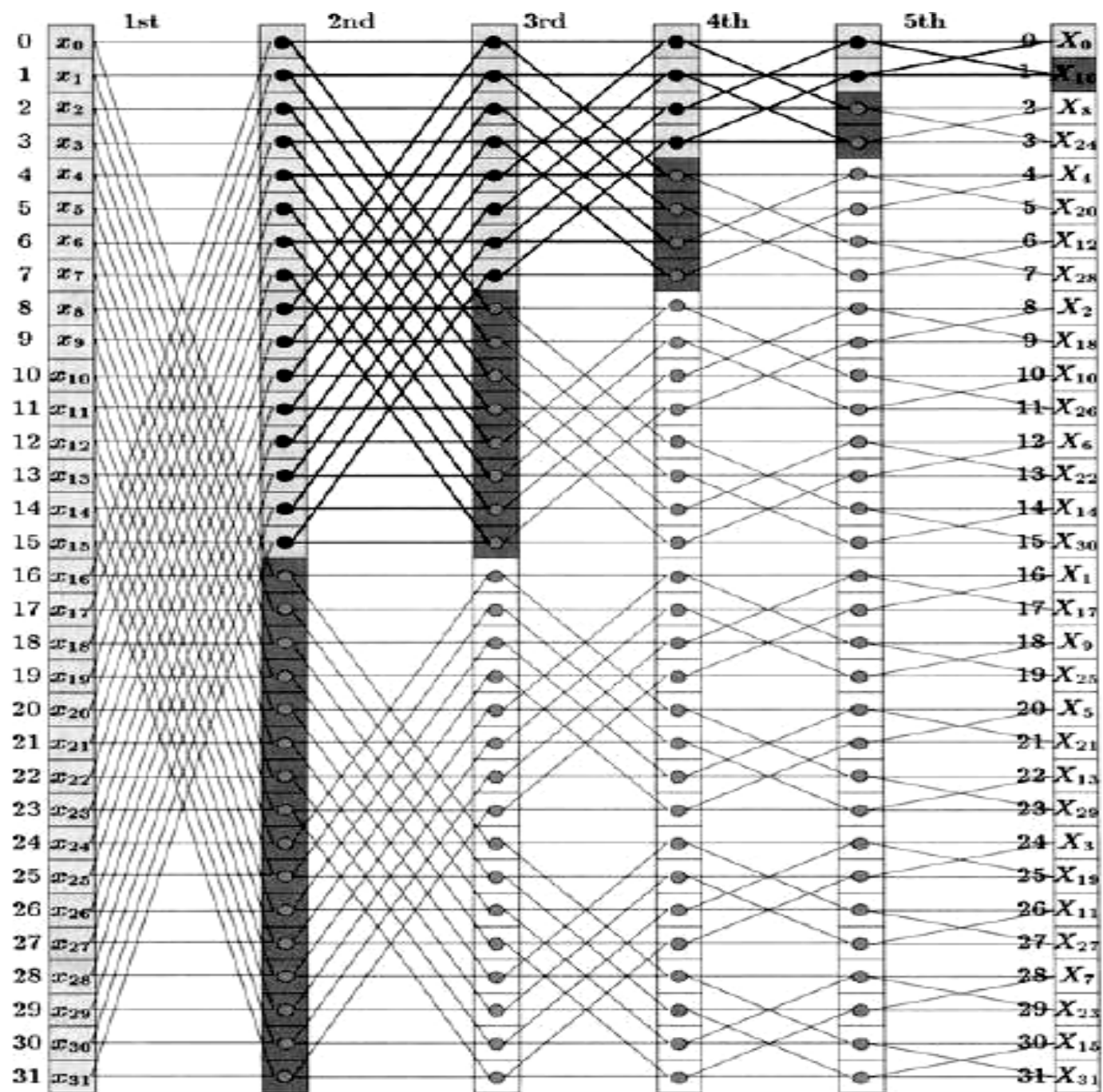


Entrada :

x_0	x_1	x_2	x_3	x_4	x_5	x_6	x_7
$a[0]$	$a[1]$	$a[2]$	$a[3]$	$a[4]$	$a[5]$	$a[6]$	$a[7]$

Salida :

X_0	X_4	X_2	X_6	X_1	X_5	X_3	X_7
$a[0]$	$a[1]$	$a[2]$	$a[3]$	$a[4]$	$a[5]$	$a[6]$	$a[7]$



Complejidad computacional.

$$T[n] = \begin{cases} 2T(n/2) + Cn & \text{si } n \geq 2 \\ 0 & n = 1 \end{cases}$$

Resolviendo la ecuación de recurrencia se tiene:

$$O(n \log n)$$

3. Procesamiento Digital de la Señal de Habla

3.2 Ventaneamiento

Se puede cortar la señal por partes para un análisis más cómodo

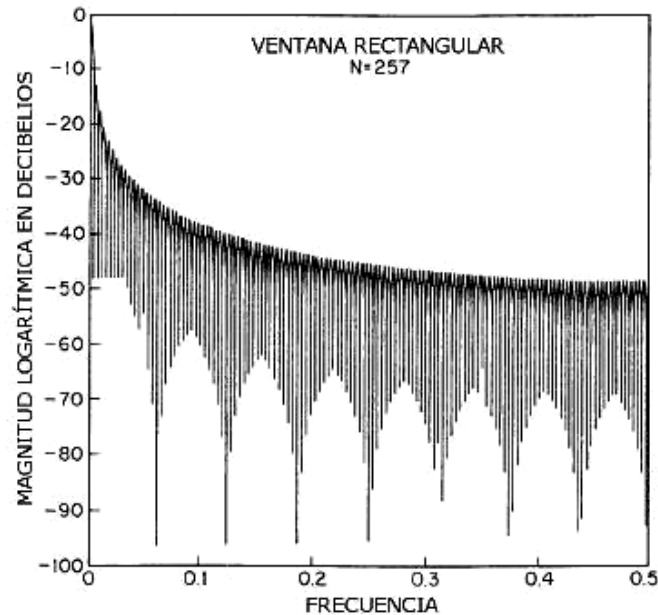
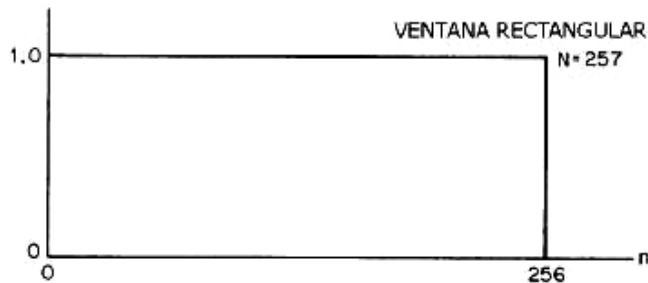
Idea: utilizar ventanitas

Problema : ¿Qué tipo de ventana usar?

3. Procesamiento Digital de la Señal de Habla

- Caso ventana rectangular

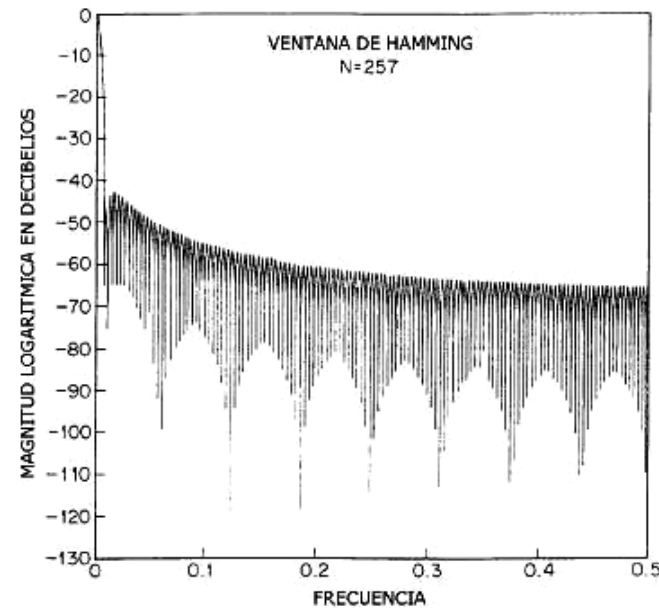
$$w[n] = 1, \quad 0 \leq n \leq N - 1$$



3. Procesamiento Digital de la Señal de Habla

- Caso ventana Hamming

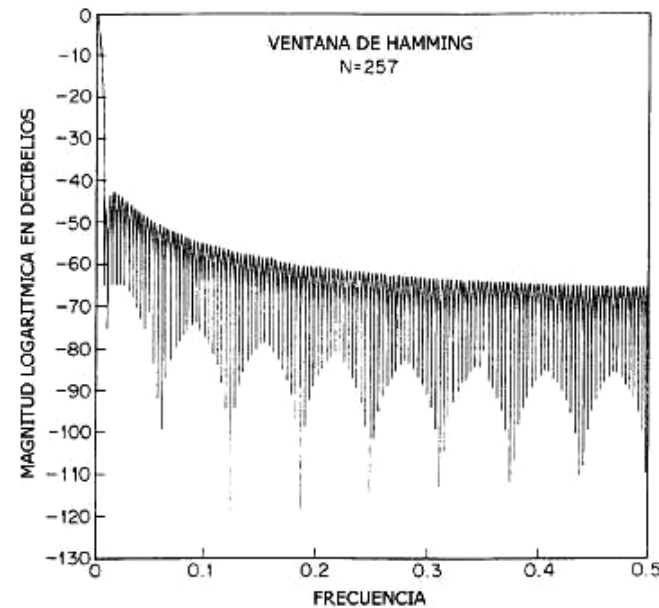
$$w[n] = 0.54 - 0.46\cos\left(\frac{2\pi n}{N-1}\right), \quad 0 \leq n \leq N-1$$



3. Procesamiento Digital de la Señal de Habla

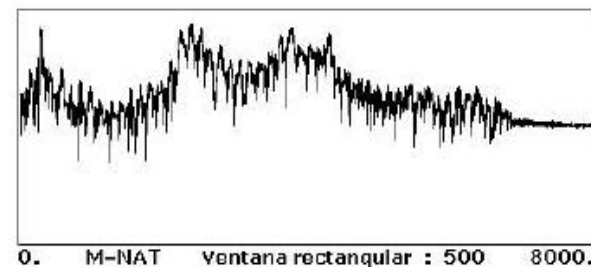
- Por que hamming? Caso ventana rectangular

$$w[n] = 0.54 - 0.46 \cos\left(\frac{2\pi n}{N-1}\right), \quad 0 \leq n \leq N-1$$



3. Procesamiento Digital de la Señal de Habla

Comparación de Ventanas.



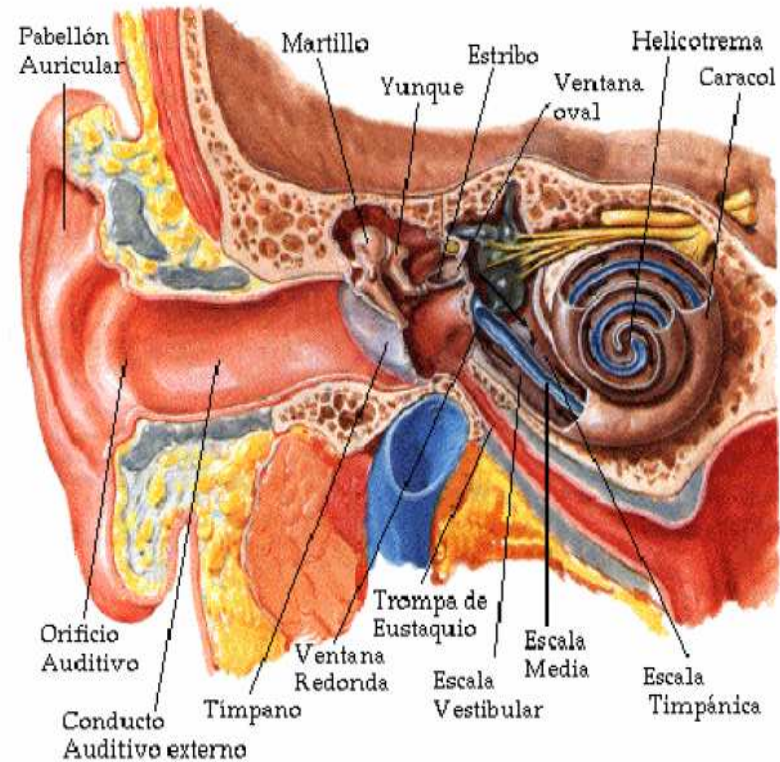
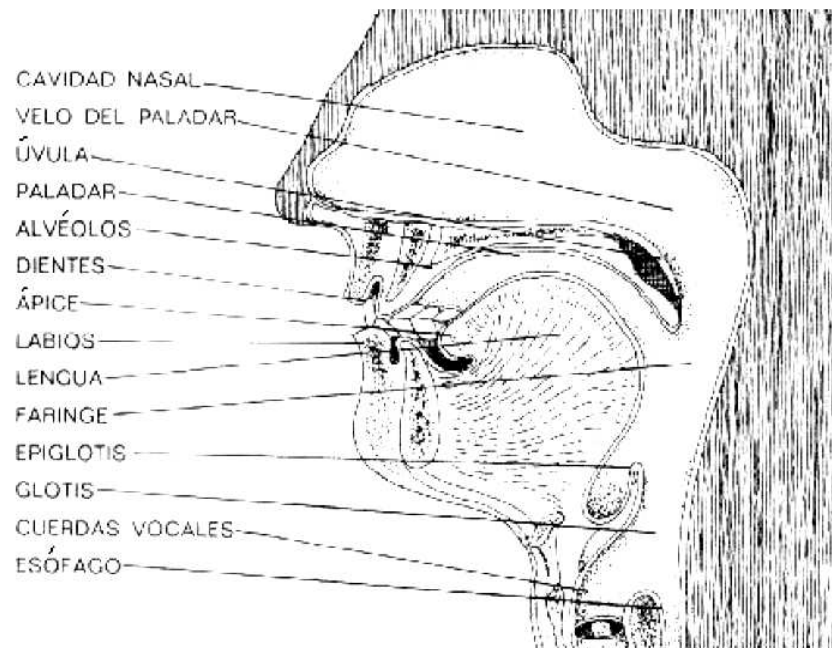
4. Coeficientes MFCC

4. Coeficientes MFCC

- Es un método (el más famoso) para extracción de características
- La idea esta inspirada en un modelo biológico
- Usa Transformada de Fourier
- Complejidad Computacional $O(n \log n)$

4. Coeficientes MFCC

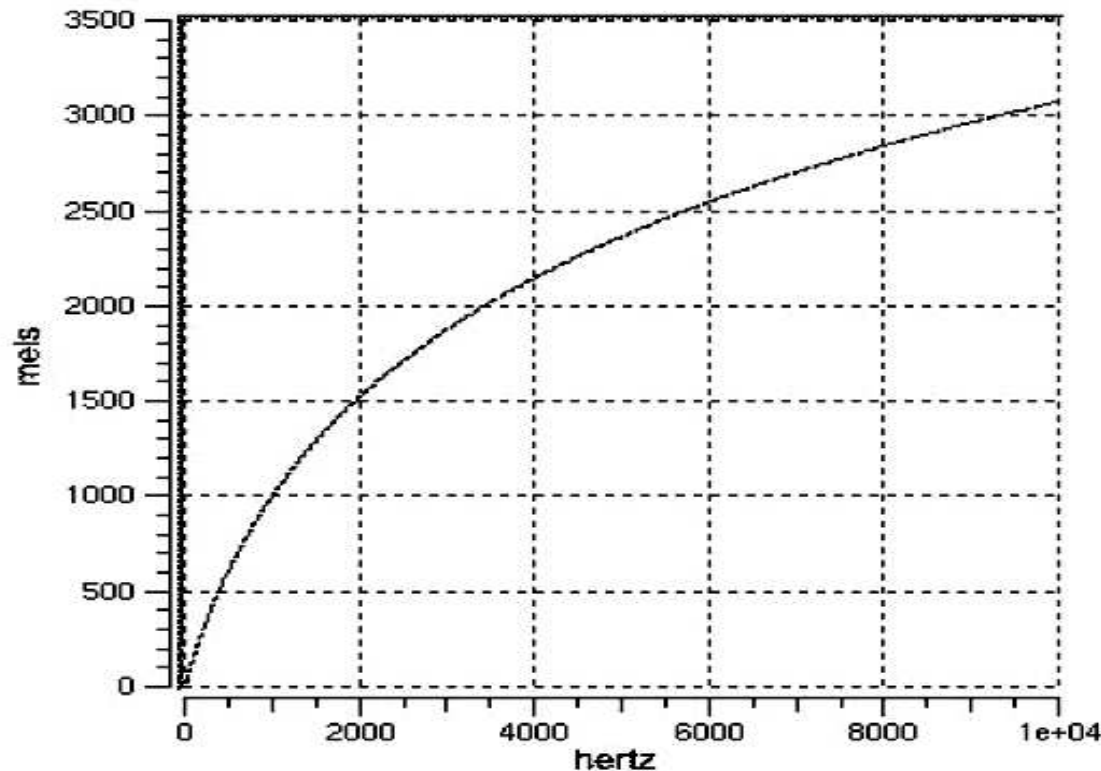
Producción y Percepción del Habla



4. Coeficientes MFCC

Frecuencia Mel.

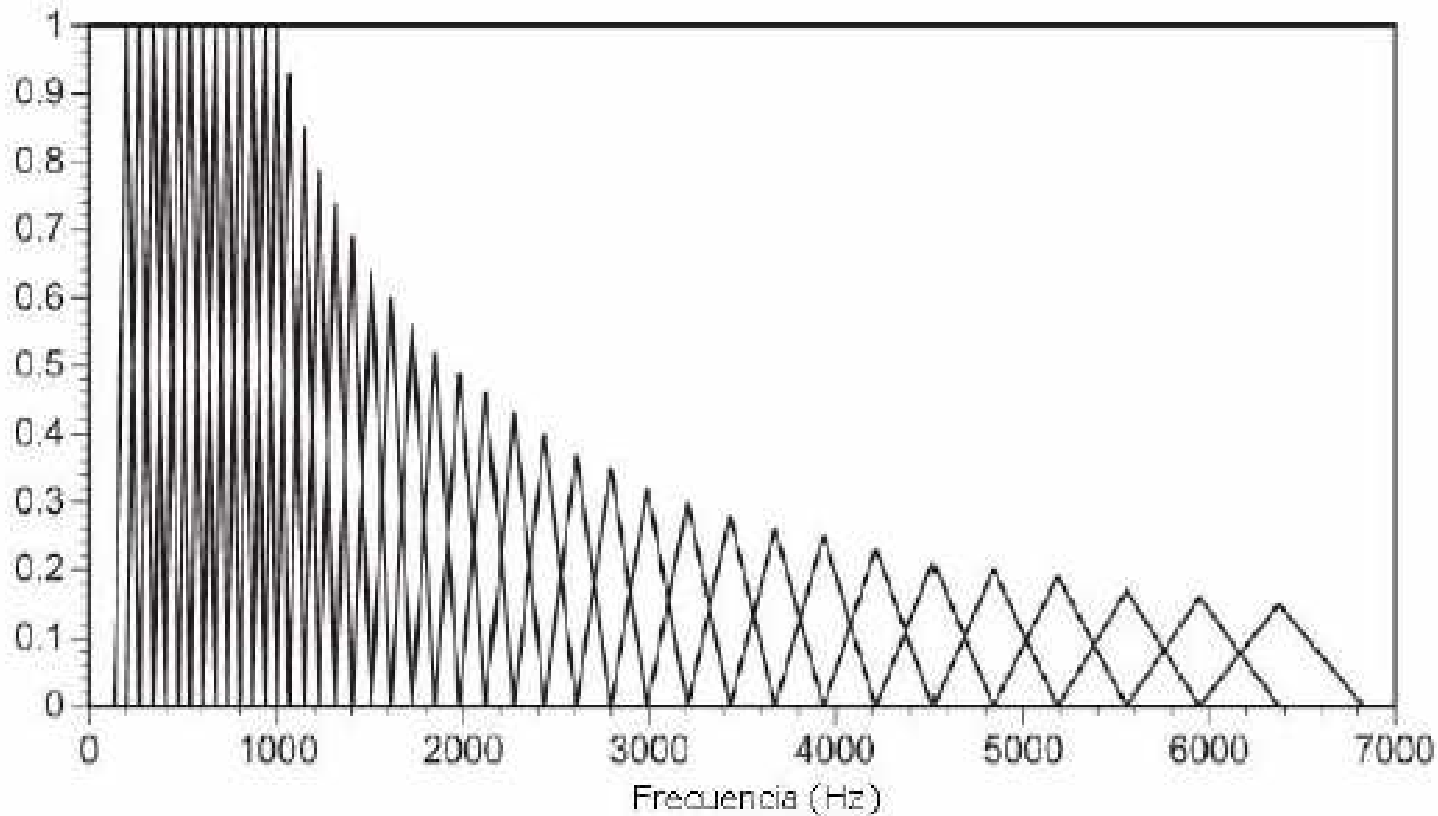
Es una escala basada en como oímos, y se ha construido , a través de experimentos fisiológicos.



$$\beta(f) = 1125 \ln\left(1 + \frac{f}{700}\right)$$

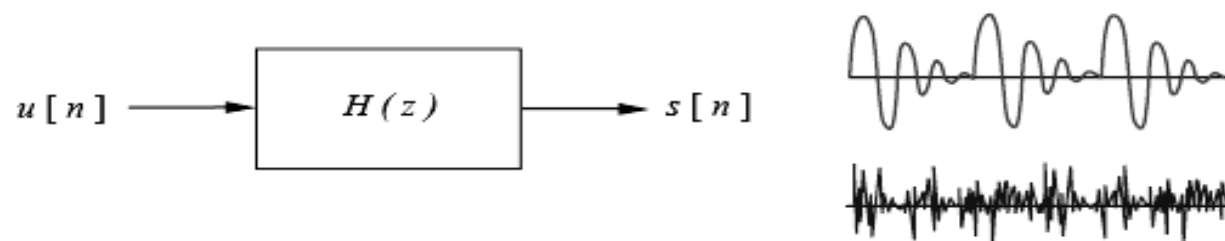
4. Coeficientes MFCC

Frecuencia Mel.



Cepstrum.

Si imaginamos la señal de voz como producto de la convolución del aire que fluye de nuestros pulmones y varios filtros correspondientes al tracto vocal.

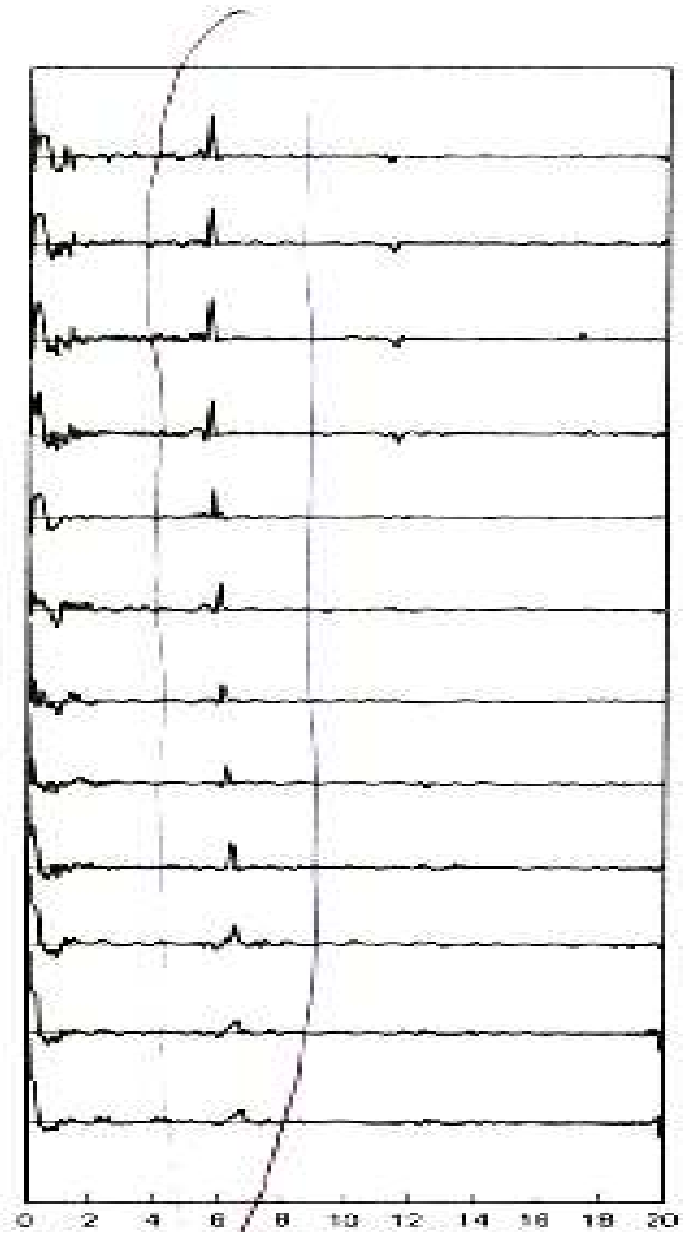


$$x[n] = e[n] * h[n]$$

$$\hat{x}[n] = \hat{e}[n] + \hat{h}[n]$$

Objetivo: Desconvolucionar la señal de voz

$$\hat{x}[n] = \frac{1}{2\pi} \int_{-\pi}^{\pi} \ln |X(e^{j\omega})| e^{j\omega n} d\omega$$



4. Coeficientes MFCC

- Algoritmo

- Se hace un análisis por frames de la señal

$$\chi^m[n] = \chi[n - mF]\omega[n].$$

Con una ventana Hamming

$$\omega[n] = 0,54 - 0,46 \cos \frac{2\pi n}{N}.$$

4. Coeficientes MFCC

- Algoritmo

- Se aplica una Transformada de Fourier a cada Frame (Transformada Corta de Fourier) con un algoritmo rápido $O(n \log n)$

$$X_m(e^{j\omega}) = \sum \chi_m[n] e^{-j\omega n} = \sum \omega[m - n] \chi[n] e^{-j\omega n}$$

En nuestro caso un algoritmo Radix-2 con decimación en frecuencia y reordenamiento de bits mezclados

4. Coeficientes MFCC

- Algoritmo

- Se traspa de la escala de frecuencias a la escala Mel, mediante un ventaneamiento con ventanas triangulares (bins)

$$H = \begin{cases} 0 & \text{si } k < f(m-1) \\ \frac{k-f(m-1)}{f(m)-f(m-1)} & f(m-1) \leq k \leq f(m) \\ \frac{k-f(m-1)}{f(m+1)-f(m)} & f(m) \leq k \leq f(m+1) \\ 0 & k > f(m+1) \end{cases}$$

$$f(m) = \frac{N}{F_s} \beta^{-1} \left(\beta(f_1) + m \frac{\beta(f_h) - \beta(f_1)}{M+1} \right)$$

4. Coeficientes MFCC

- Algoritmo

- Se obtiene el Cepstrum de las frecuencias en escala Mel

$$S(m) = \ln\left(\sum |X(k)| H_m(k)\right), 0 < m < M$$

- Finalmente se una transformada de Coseno II es calculada

$$c(m) = \sum S(m) \cos\left(\pi n \left(\frac{m + \frac{1}{2}}{M}\right)\right)$$

5. Trasformada Wavelet

5. Transformada Wavelet

“La Transformada Wavelet es una herramienta matemática que corta los datos, funciones o operadores en diferentes componentes de frecuencia y estudia cada componente a una resolución ubicada a esa escala.”

Ingrid Daubechies

Ten Lectures of Wavelets

5. Transformada Wavelet

transformada wavelet continua de una función f está dada por:

$$(T^{wav} f)(a, b) = |a|^{-\frac{1}{2}} \int \delta t f(t) \psi\left(\frac{t-b}{a}\right)$$

la familia de wavelets se puede construir dilatando y trasladando

$$\psi^{a,b}(x) = |a|^{-\frac{1}{2}} \psi\left(\frac{t-b}{a}\right)$$

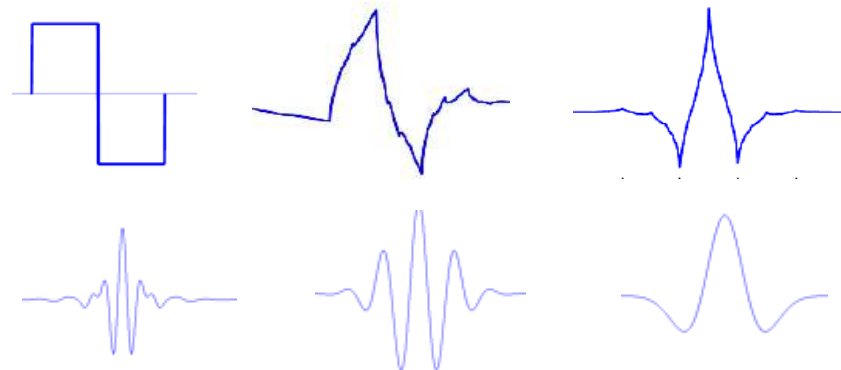
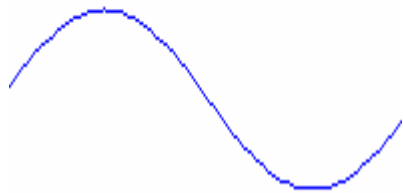
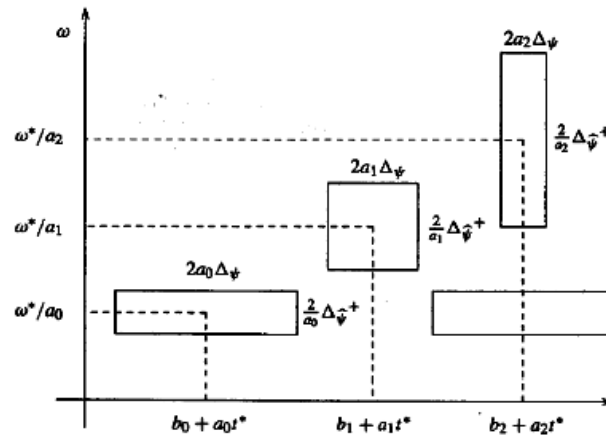
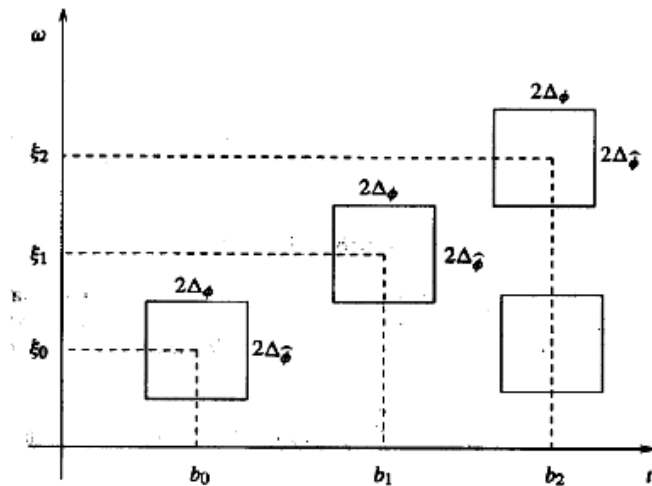
$$f = C_{\psi}^{-1} \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \frac{\delta a \delta b}{a^2} \psi(T^{wav} f)(a, b) \psi^{a,b}$$

5. Trasformada Wavelet

T. Fourier vs T. Wavelet

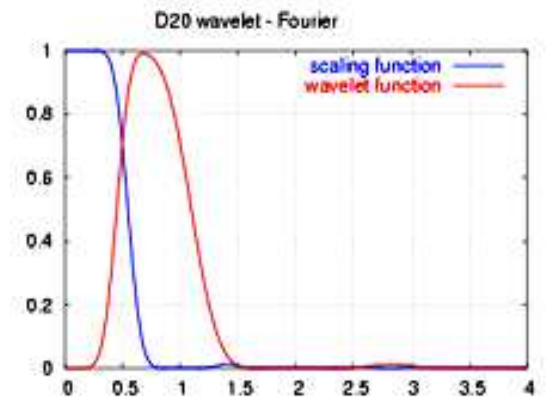
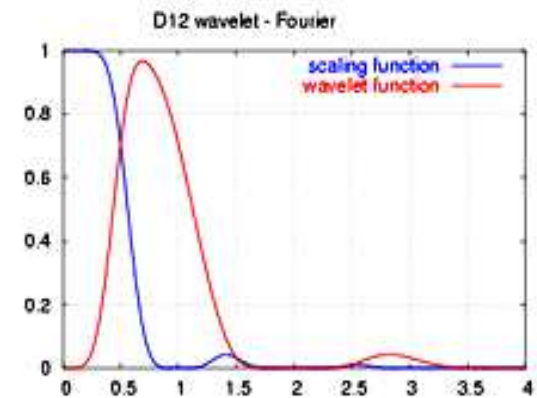
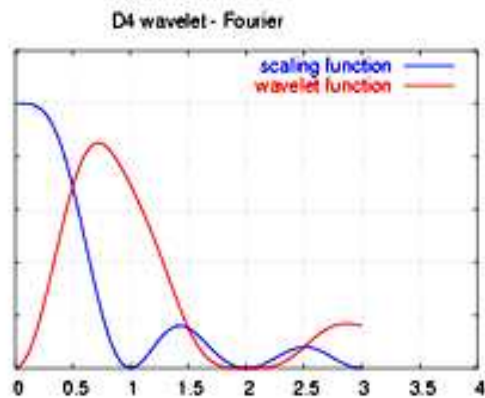
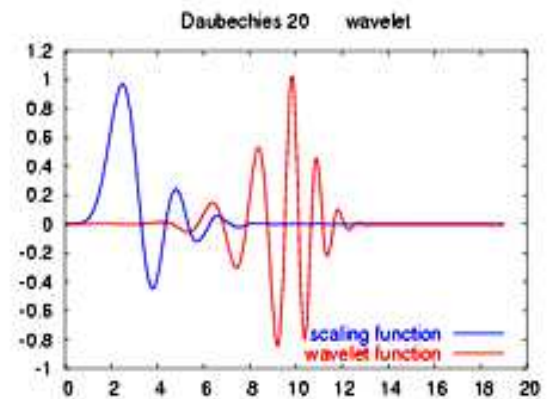
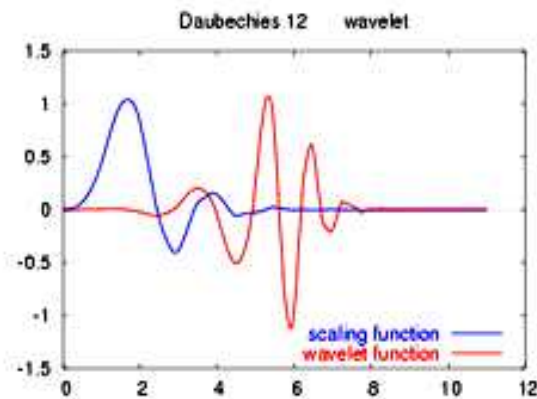
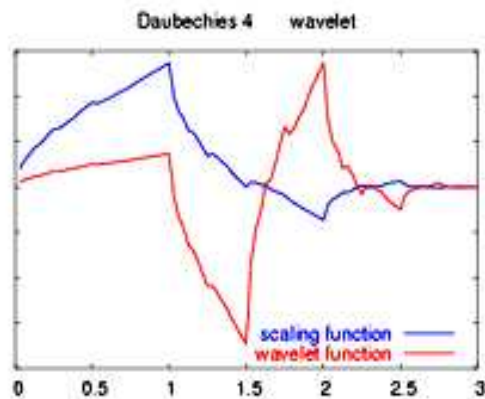
$$T^{win}(w, t) = \int \delta s f(s) g(s - t) e^{-i\omega s}$$

$$CWT(a, b) = \frac{1}{\sqrt{a}} \int x(t) \psi\left(\frac{t - b}{a}\right) \delta t$$



5. Transformada Wavelet

Wavelets en el Dominio de la Frecuencia



5. Trasformada Wavelet

Wavelets Discretas.

$$\psi^{m,n}(x) = a_0^{-\frac{m}{2}} \psi(a_0^{-m}(x - nb_0a_0^m))$$

$$\psi^{m,n}(x) = a_0^{-\frac{m}{2}} \psi(a_0^{-m}x - nb_0)$$

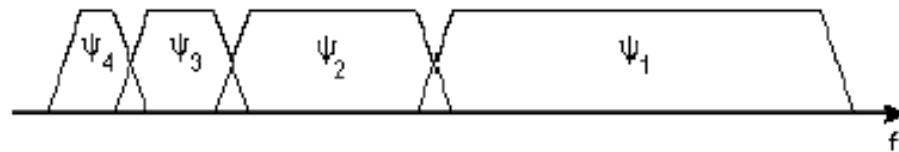
en particular si escogemos $a_0 = 2$ y $b_0 = 1$ entonces:

$$\psi^{m,n}(x) = 2^{-\frac{m}{2}} \psi(2^{-m}x - n)$$

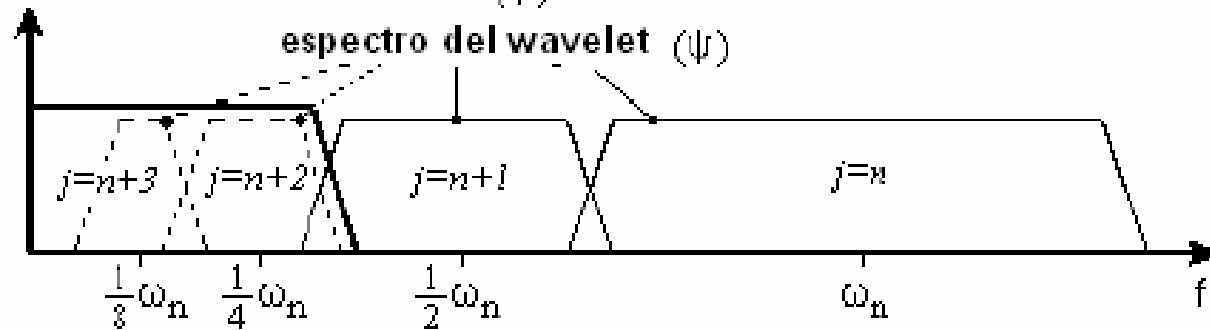
5. Transformada Wavelet

Filtro Pasa Banda

$$H(f(at)) = \frac{1}{|a|} H\left(\frac{\omega}{a}\right)$$

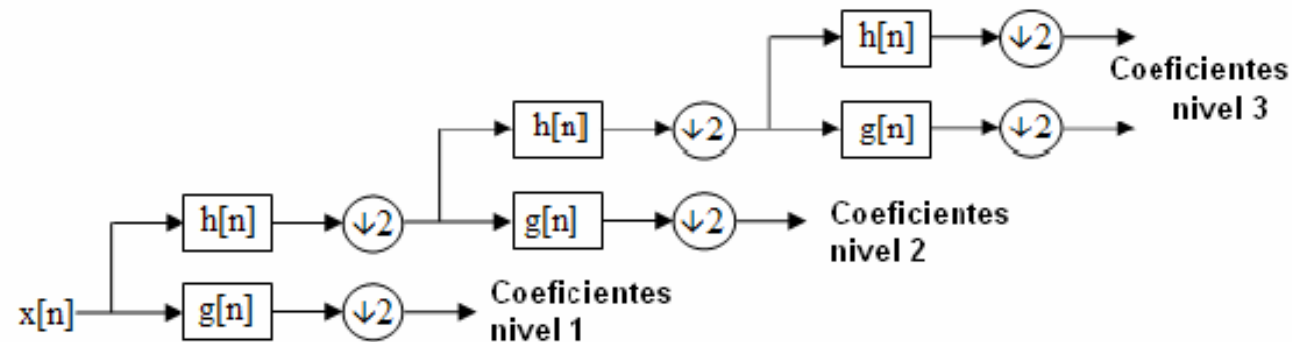


Espectro de la función escala (φ)



5. Transformada Wavelet

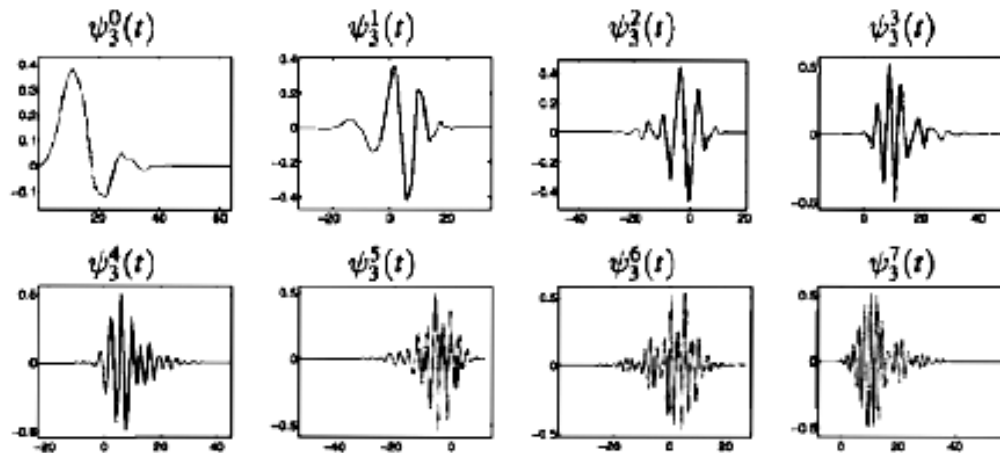
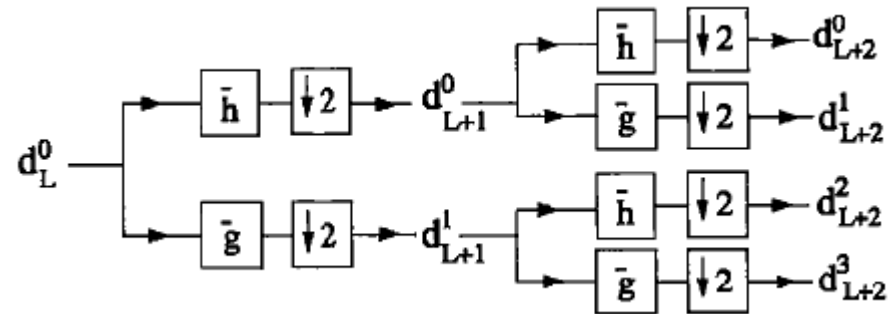
Algoritmo de Banco de Filtros Iterativo



$$T(n) = 2cn - 2c$$

La Transformada wavelet con banco de filtros tiene una complejidad de $O(n)$.

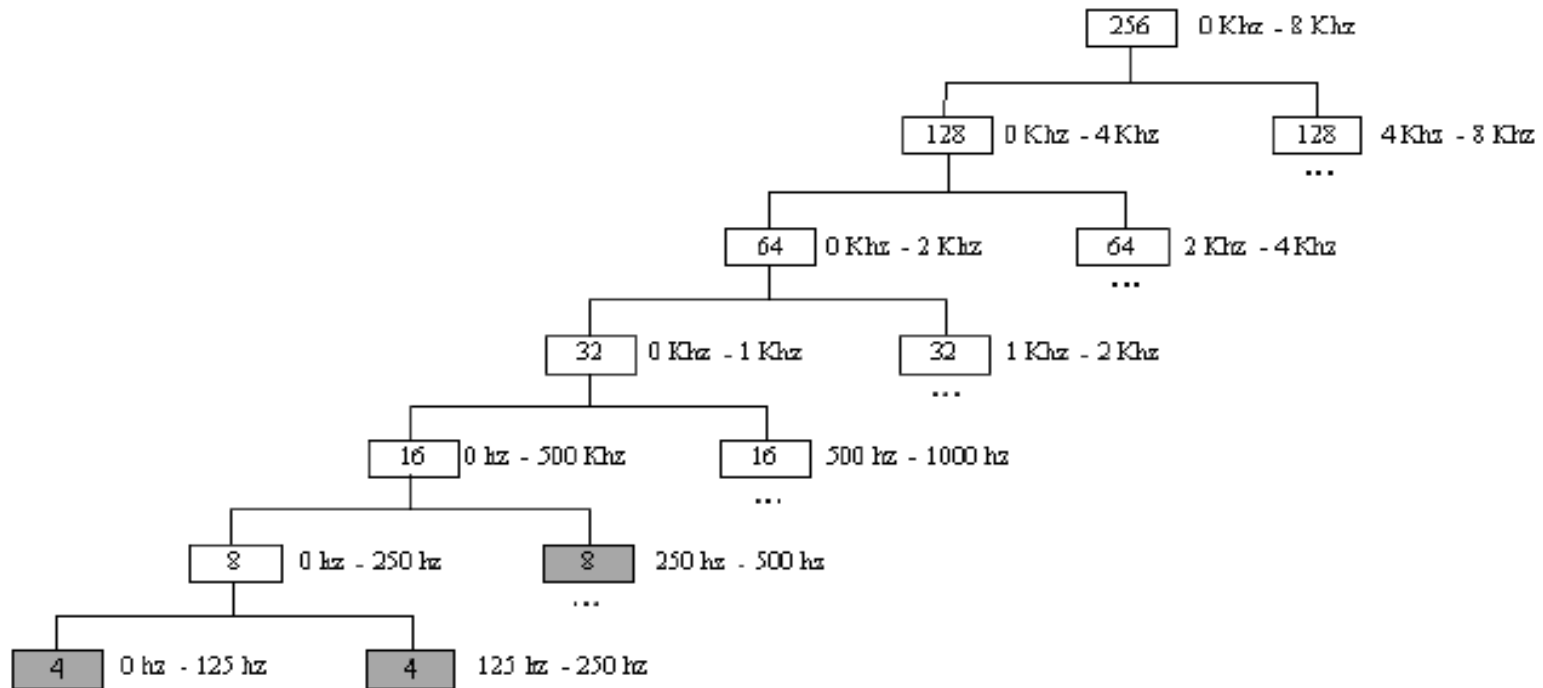
Wavelets Packets



6. Extracción de Características **usando Wavelets**

Extracción de características con Wavelets

Arbol de descomposición.



Extracción de Características con Wavelets

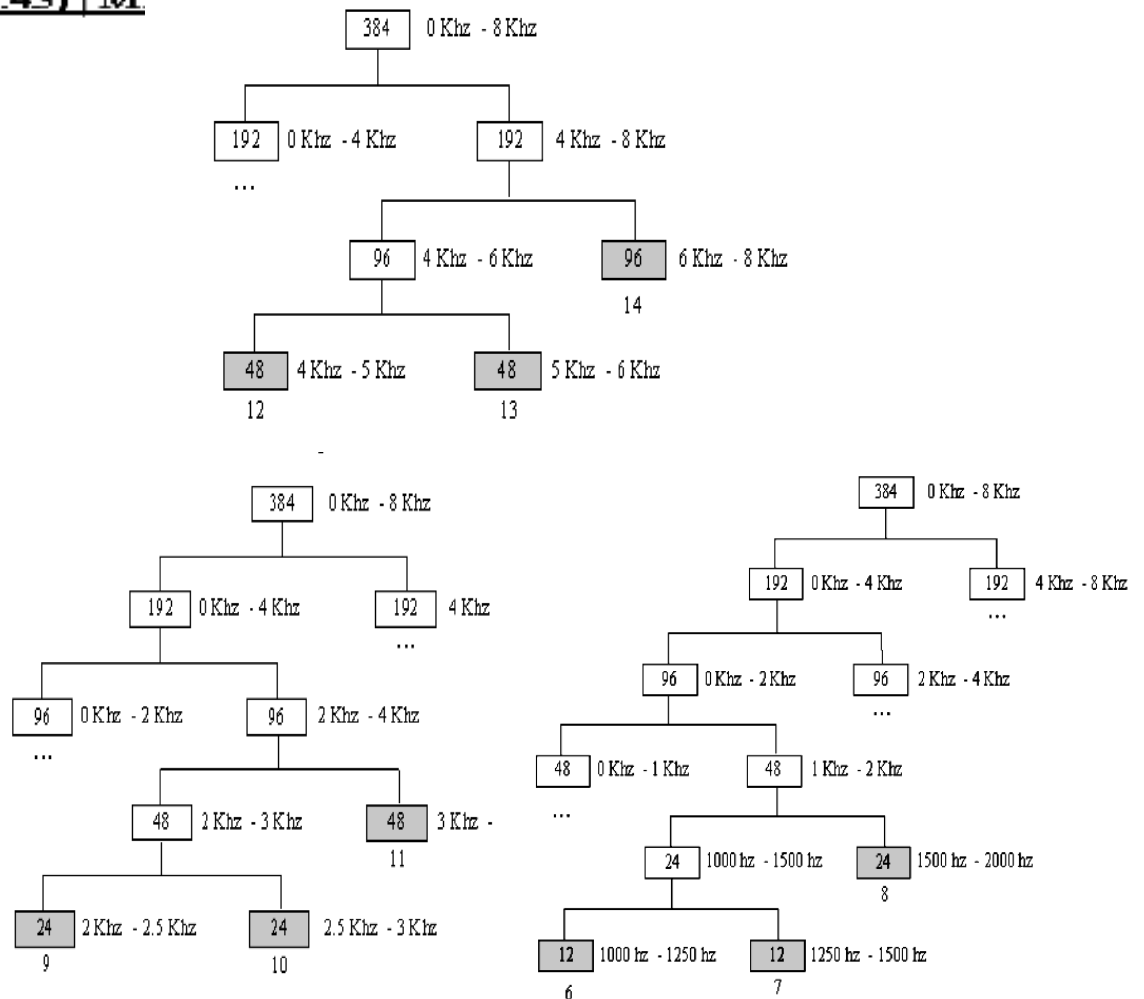
Packet Perceptuales

Análisis con Wavelet Packet paso (202.45) | M

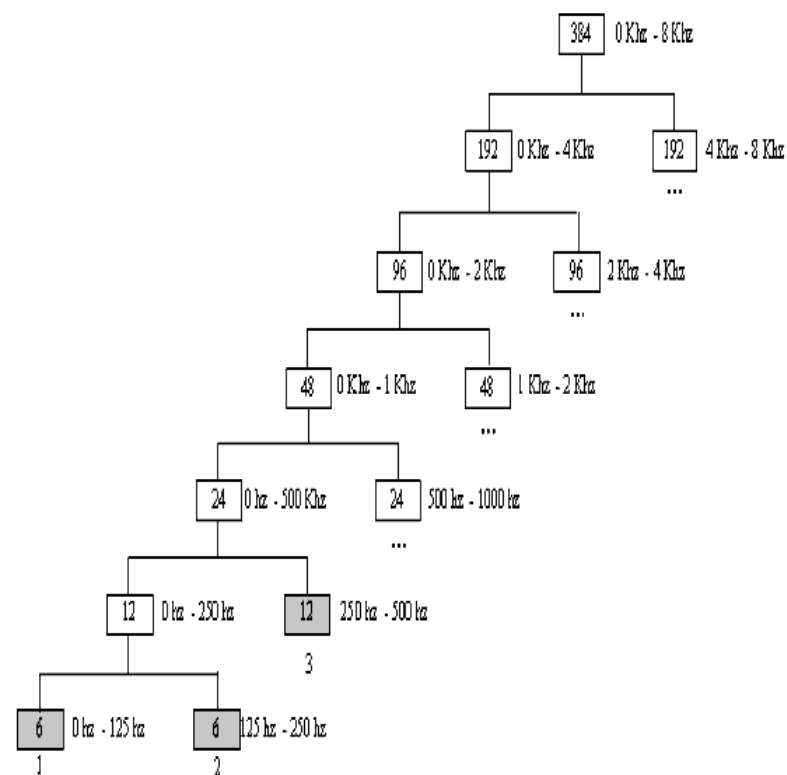
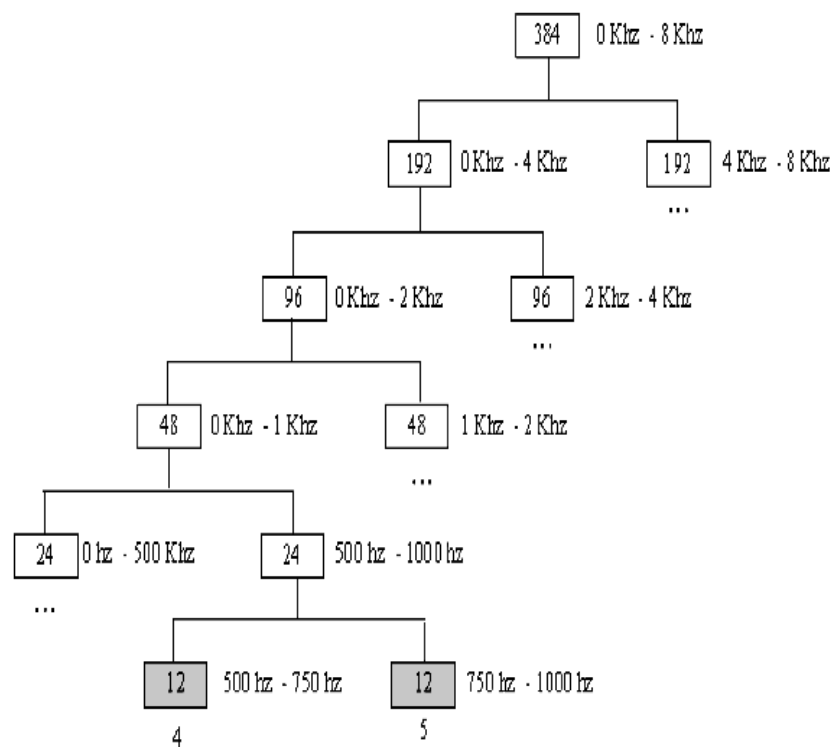
Equivalencias escala MEL a Frecuencias

	MEL	Frecuencias
1	0	0
2	202.10	135
3	404.21	300
4	606.31	500
5	808.42	735
6	1010.53	1030
7	1212.63	1375
8	1414.73	1750
9	1616.85	2235
10	1818.95	2830
11	2021.06	3500
12	2223.16	4333
13	2425.27	5310
14	2627.38	6500
	2829.48	8000

Arbol de descomposición.



Extracción de Características con Wavelets Packet Perceptuales



Filtros Usados

Haar

$$h(n) = \left[\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right]$$

Wavelets de Daubechies 4

$$h(n) = \left[\frac{1 + \sqrt{3}}{4\sqrt{2}}, \frac{3 + 3\sqrt{3}}{4\sqrt{2}}, \frac{3 - 3\sqrt{3}}{4\sqrt{2}}, \frac{1 - \sqrt{3}}{4\sqrt{2}} \right]$$

Wavelets de Daubechies 6

$$h(n) = [0,3326, 0,8068, 0,4598, -0,1350, -0,0854, 0,0352]$$

Wavelet Coiflet 6

$$h(n) = \left[\frac{1 - \sqrt{7}}{16\sqrt{2}}, \frac{5 + \sqrt{7}}{16\sqrt{2}}, \frac{14 + 2\sqrt{7}}{16\sqrt{2}}, \frac{14 - 2\sqrt{7}}{16\sqrt{2}}, \frac{1 - \sqrt{7}}{16\sqrt{2}}, \frac{-3 + \sqrt{7}}{16\sqrt{2}} \right]$$

Siendo los filtros g correspondientes: $g_n = (-1)^n h_{-n+1}$

Extracción de características

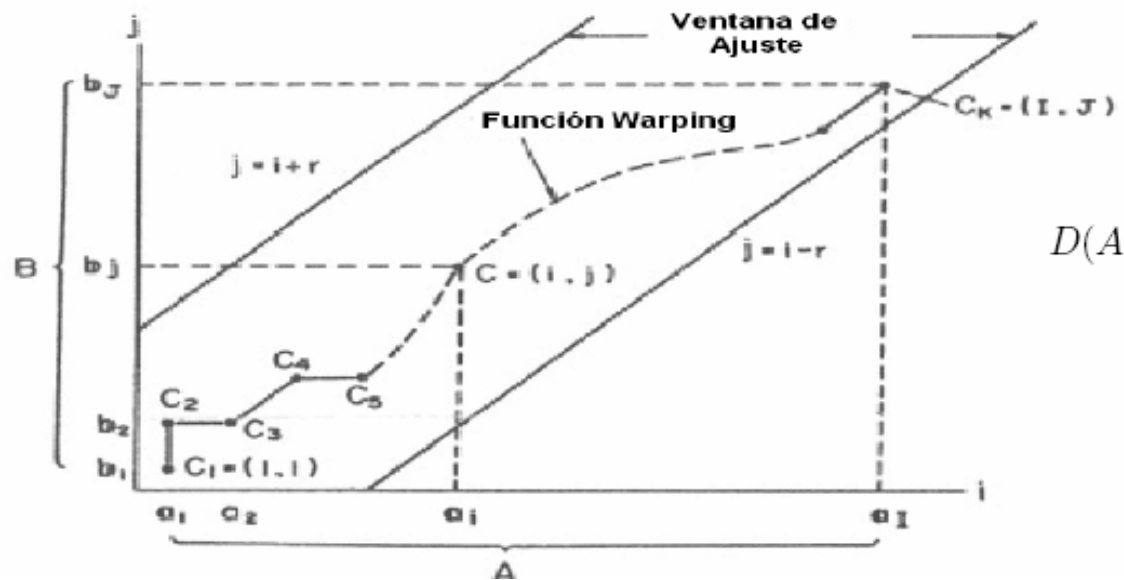
- Cálculo de energías por nivel de resolución aproximadamente igual a la escala Mel

$$E_i = \frac{\sum_{j=1}^N (W_i^p f(j))^2}{N_i}$$

- Aplicación de un “Cepstrum”

$$F(i) = \sum_{n=1}^N \log E_n \cos\left(\frac{i\left(\frac{n-1}{2}\right)}{N}\right)$$

Dynamic Time Warping DTW



$$D(A, B) = \min \left[\frac{\sum_{k=1}^K d(c(k)) \cdot w(k)}{\sum_{k=1}^K w(k)} \right]$$



Experimentos y Resultados

Resultados Wavelets O(n)

DISTRIBUCION DE 900 PALABRAS PRONUNCIADAS POR 60 PERSONAS SEGUN SU IDENTIFICACION Y NO IDENTIFICACION POR LA COMPUTADORA CON EL METODO DE FOURIER Y EL MÉTODO PROPUESTO BASADO EN LAS WAVELET HAAR.

TRUJILLO 2006.

MÉTODO WAVELET HAAR	MÉTODO DE FOURIER		TOTAL
	Palabras identificadas	Palabras no identificadas	
- Palabras identificadas	262	13	275
- Palabras no identificadas	172	453	625
TOTAL	434	466	900

χ^2_{MIN} : Prueba ji-cuadrado de Mc Nemar para datos correlacionados

$\chi^2_{Mc} = 136.65$	$p < 0.01$
------------------------	------------

Eficiencia relativa de aciertos = $275 / 434 \times 100$

EF :63.36%

DISTRIBUCION DE 900 PALABRAS PRONUNCIADAS POR 60 PERSONAS SEGUN SU IDENTIFICACION Y NO IDENTIFICACION POR LA COMPUTADORA CON EL METODO DE FOURIER Y EL MÉTODO PROPUESTO BASADO EN LAS WAVELETS DE DAUBECHIES 6 .

TRUJILLO 2006.

MÉTODO WAVELET DAUB 6	MÉTODO DE FOURIER		TOTAL
	Palabras identificadas	Palabras no identificadas	
- Palabras identificadas	430	93	523
- Palabras no identificadas	4	373	377
TOTAL	434	466	900

χ^2_{MIN} : Prueba ji-cuadrado de Mc Nemar para datos correlacionados

$\chi^2_{Mc} = 81.66$	$p < 0.01$
-----------------------	------------

Eficiencia relativa de aciertos = $523 / 434 \times 100$

ER-aciertos = 120.51

Resultados Wavelets Packets $O(n \log n)$

DISTRIBUCION DE 900 PALABRAS PRONUNCIADAS POR 60 PERSONAS SEGUN SU IDENTIFICACION Y NO IDENTIFICACION POR LA COMPUTADORA CON EL METODO DE FOURIER Y EL MÉTODO PROPUESTO BASADO EN LAS WAVELETS PACKET CON EL METODO DE FOURIER Y EL METODO PROPUESTO BASADO EN LAS DUABECHIES 6.

TRUJILLO 2006.

		MÉTODO DE FOURIER		TOTAL
MÉTODO WAVELET PACKET DAUB6		Palabras identificadas	Palabras no identificadas	
-	Palabras identificadas	425	220	645
-	Palabras no identificadas	9	246	255
TOTAL		434	466	900

χ^2_{Mn} : Prueba ji-cuadrado de Mc Nemar para datos correlacionados

$\chi^2_{Mn} = 194.41$	$p < 0.01$
------------------------	------------

Eficiencia relativa de aciertos = $645 / 434 \times 100$

ER-aciertos = 148.62

DISTRIBUCION DE 900 PALABRAS PRONUNCIADAS POR 60 PERSONAS SEGUN SU IDENTIFICACION Y NO IDENTIFICACION POR LA COMPUTADORA CON EL METODO DE FOURIER Y EL METODO PROPUESTO BASADO EN LAS WAVELETS PACKET PERCEPTUAL CON DAUBECHIES 4.

TRUJILLO 2006.

		MÉTODO DE FOURIER		TOTAL
MÉTODO WAVELET PACKET PERCEPTUAL		Palabras identificadas	Palabras no identificadas	
-	Palabras identificadas	420	183	603
-	Palabras no identificadas	14	283	297
TOTAL		434	466	900

χ^2_{Mn} : Prueba ji-cuadrado de Mc Nemar para datos correlacionados

$\chi^2_{Mn} = 144.98$	$p < 0.01$
------------------------	------------

Eficiencia relativa de aciertos = $603 / 434 \times 100$

ER-aciertos = 138.94

Comparaciones con MFCC

<i>Método</i>	<i>Tasa aceptación</i>	<i>Error</i>
Coeficientes Cepstrales en Escala Mel	85.32%	14.68%
Wavelet Haar	34.47%	65.53%
Wavelet Daubechies 4	51.79%	48.21%
Wavelet Daubechies 6	61.32 %	38.68%
Wavelet Coiflets 6	55.46 %	44.54%
Wavelet Packet Perceptuales Walsh	55.79 %	44.21 %
Wavelet Packet Perceptuales Daubechies 4	69.14 %	30.86%
Wavelet Packet Perceptuales Daubechies 6	74.43 %	25.57%
Wavelet Packet Perceptuales Daubechies 4 (22)	71.21 %	28.79%

Tabla Datos obtenidos utilizando la técnica de DTW como reconocedor, con distancia Chebyshev y con Slope Constrain $P=1$. Se observa la mejor performance en las Wavelet Packet Perceptuales Daubechies 6, y la mas pobre en las Wavelet Haar

Tasa de Reconocimiento por Palabra

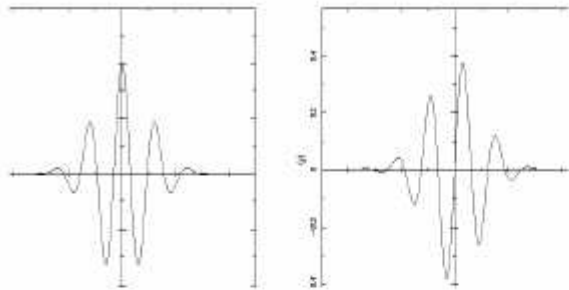
	<i>MFCC</i>	<i>W. Haar</i>	<i>W. Db₄</i>	<i>W. Db₆</i>	<i>W. Coif₆</i>	<i>WP Walsh</i>	<i>WP Db₄</i>	<i>WP Db 6</i>	<i>WP Perc.</i>
<i>Arriba</i>	86%	17%	41%	48%	34%	31%	66%	72%	62
<i>Cerrar</i>	100%	45%	79%	86%	76%	97%	100%	97%	97%
<i>Coger</i>	90%	17%	17%	31%	17%	28%	52%	69%	62%
<i>Cuatro</i>	97%	48%	28%	31%	24%	62%	76%	83%	83%
<i>Dos</i>	86%	28%	38%	48%	48%	59%	69%	69%	72%
<i>Eliminar</i>	72%	24%	34%	52%	48%	24%	48%	76%	45%
<i>Error</i>	97%	24%	66%	93%	83%	93%	97%	97%	93%
<i>Hola</i>	83%	38%	48%	72%	55%	38%	52%	59%	52%
<i>Izquierda</i>	86%	31%	48%	62%	59%	24%	59%	72%	69%
<i>Pez</i>	62%	24%	52%	41%	66%	62%	59%	59%	55%
<i>Salir</i>	86%	66%	79%	79%	72%	76%	76%	79%	76%
<i>Terminar</i>	62%	24%	24%	31%	28%	41%	55%	55%	55%
<i>Tres</i>	83%	24%	52%	59%	45%	41%	66%	62%	66%
<i>Tres</i>	59%	24%	45%	59%	41%	48%	55%	59%	62%
<i>Uno</i>	76%	31%	79%	83%	86%	76%	72%	66%	59%

8. Conclusiones

Conclusiones.

- El mejoramiento del espectro se da gracias al análisis tiempo frecuencia de las wavelets.
- Una extracción de características usando solamente la Transformada de Fourier no brinda buenos resultados.
- Los wavelets pueden ser utilizados alternativamente, para el procesamiento digital de la señal de habla.
- La complejidad computacional de los algoritmos de extracción de características usando las wavelets y las wavelets packets es de $O(n)$ y de $O(n \log n)$ respectivamente.
- La ventaja de utilizar wavelets radica, en la variedad de funciones wavelet que se puede escoger.
- Las wavelets que mejor funcionan, son aquellos que tienen su espectro parecido a un filtro paso de banda ideal.

Proyecciones

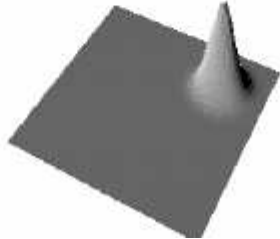
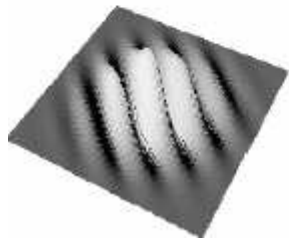
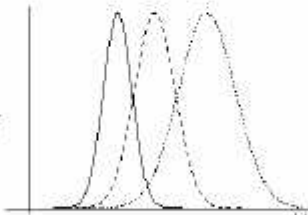
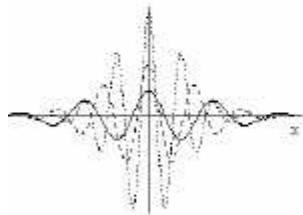


$$\psi(x) = \pi^{-\frac{1}{4}} \left(e^{-i\omega x} - e^{-\frac{\pi^2 \alpha^2}{4}} \right) e^{-\frac{x^2}{\alpha^2}}$$

Parte real e imaginaria de las Wavelets de Morlet

Dominio del Tiempo

Dominio de la Frecuencia



<i>Método</i>	<i>Tasa aceptación</i>	<i>Error</i>
Coefficientes Cepstrales en Escala Mel	85.32%	14.68%
Wavelet Continuo de Morlet	87.32%	12.68%

$O(n \log n)$.

Referencias

- [1] ABOUFADEL, E. A wavelets approach to voice recognition. *Grand Valley State University* (2001).
- [2] ATAL, AND SCHROEDER. Predictive coding of speech signals. *Report of the 6th Int. Congress on Acoustics, Tokio, Japan* (1968).
- [3] BAUN, L., AND EAGON, J. Perceptual linear predictive analysis of speech. *RBulletin of American Mathematical Society, 1967, 73, pp. 360-363* (1968).
- [4] DAUBECHIES, I. *Ten lectures on wavelets*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1992.
- [5] GROSSMANN, A., AND MORLET, J. Decomposition of hardy functions into square integrable wavelets of constant shape. *SIAM Journal on Mathematical Analysis* 15, 4 (1984), 723–736.
- [6] GUEVARA, J. L. *Lorito, speech recognition software*, v 1.0 ed. Universidad Nacional de Trujillo, Trujillo, Enero 2007.
- [7] GUEVARA, J. L., AND SALAZAR, J. O. *Extracción de Características en el Procesamiento Digital de una Señal para el Mejoramiento del Reconocimiento Automático de Habla usando Wavelets*. Tesis, Trujillo, Enero 2007.
- [8] HERMANSKY, H. A an inequality with applications to statistical estimation for probabilistic functions of markov processes and to a model for ecology. *J. Acoust. Soc. Am.* (2005).
- [9] HUANG, X., ACERO, A., AND HON, H.-W. *Spoken Language Processing: A Guide to Theory, Algorithm and System Development*. Prentice Hall PTR, April 2001.
- [10] M. SIAFARIKAS, TODOR GANCHEV, N. F. Objective wavelet packet features for speaker verification, 2000.

Software de pruebas

LORITO version 3.14

