

TEMA 2: LA SEÑAL DE VOZ.

MODELADO Y ANÁLISIS

1.- INTRODUCCIÓN

La señal de voz es una señal natural, con bastante complejidad respecto a señales artificiales

Aspectos históricos importantes: PCTE, Vocoder

Problemas básicos relativos a la señal de voz:

- codificación: representar la señal de voz como un conjunto de bits → compromiso entre calidad y velocidad
- sintesis: producir voz a partir de texto
- reconocimiento: obtener el texto correspondiente a una señal de voz.
Necesidad de análisis de voz + contexto (gramática)
- ayudas a discapacitados en producción (obtener voz)

"Señal de voz" se emplea siempre T_x , R_x , y la propia señal

para conocer qué aspectos de la voz son importantes para el éxito del receptor

Comportamiento del oído humano = analizadores de espectros

→ responde a las amplitudes de los espectros,
y tanto a la fase.

Además, la oscila es aproximadamente exponencial

→ encuadreacion: capacidad de los padres distinguir
tomas cercanas a un toco de gran amplitud.

• aparato fonador: sistema realimentado (voces b que hablan)

(*2.1*)

fuerza/potencia { - diafragma
- pliegues
- musculos

partes correspondientes a

producción de
sonido

cuerdas vocales: vibración que se
abre y cierra periódicamente
(señal sin continuidad)

cuerdas resonantes: confinar la
señal

(*2.2*) → transparencia "ingenieril"

(*2.3*) → las cuerdas se cierran sobre el paso del
aire debido a la presión

• frecuencia fundamental o "pitch": frecuencia de vibración de las

cuerdas vocales. Depende de:

- sexo: hombre con fundamental más baja ($\approx 50-250 \text{ Hz}$) que
mujer ($\approx 190-590 \text{ Hz}$)

- tensión variable sobre las cuerdas vocales

• cuidad resonante: vibrante, para difundir la señal deseada

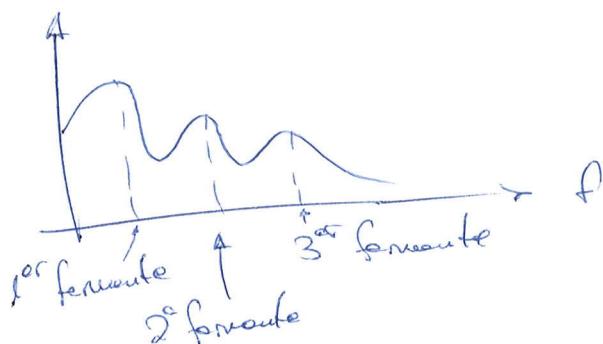
⇒ señal co-estimulativa (para los procesos de interés)

\Rightarrow uso de modelos no estacionarios, para superar la falta de
para obtener los espectros

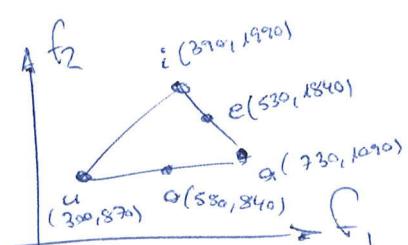
Solución: considerar segmentos de tiempo fijos \Rightarrow análisis localizado
de fuentes

- vez constituida por fuentes, cuya caracterización da lugar a la informa-
ción (vez = sistema de comunicación digital). Las fuentes se caracterizan
por su espectro y dentro de este, por sus formantes.

ejemplo: /a/



Dos los vocales, basta con los 2 primeros formantes \Rightarrow triángulo vocalico
(frecuencias para el hueco)



Se intenta obtener un alfabeto universal de fuentes.

ej: el castellano tiene pocas fuentes

el sonido se genera mediante:

- ondas vocales (vibración)

- mecanismos al deformar el tracto vocal. ej: /f/

Aparato fonador en continuo movimiento de voz no estacionario,
pero fuentes con estacionariedad

• Fuentes:

- sordos vocales (de voz) / sonores (los cuerdas vocales libres)
/ sordos

- vocales: sonidos sonoros propios (cuerdas vocales + confluencia), tienen más energía (o mucha estacionariedad), bastante estacionariedades, se pueden mantener tiempo. Se caracteriza bastante bien por sus 2 primeras fuentes

- diphongs: fuentes, no se pueden considerar 2 vocales juntas ya que la sordedad es distinta.

• Fuentes:

• Ligados o semi-vocales: sonores, con menos energía, se caracterizan por ej: /l/, /r/

• fricativas: totalmente sordos (/s/, /f/) o con algo de sonoro (/v/, /z/, ...). Poco energía, difíciles de medir.

• occlusivas: cierre + apertura del tracto vocal

sordos: /p/, /t/, ...

sonores: /b/, /d/, ...

No estacionario: parte de cierre + parte de señal

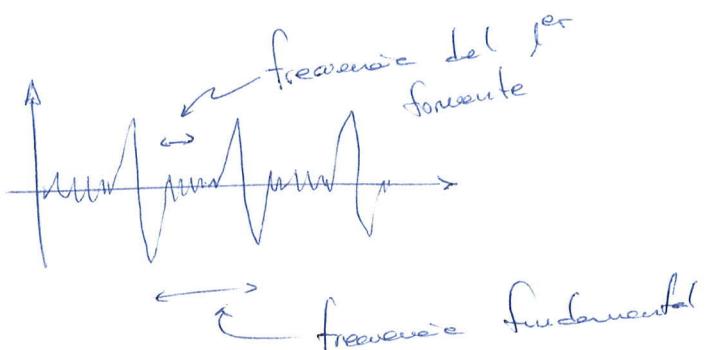
• africadas o semiocclusivas: no sonores (/ch/)

• nasales: se identifican bastante bien por el uso de la cavidad nasal (velo del paladar levantado cuerdas vocales bajas y nasal). Con las 2 cavidades en paralelo aparecen zonas en el espectro.

Ej: /u/, /ui/, ...

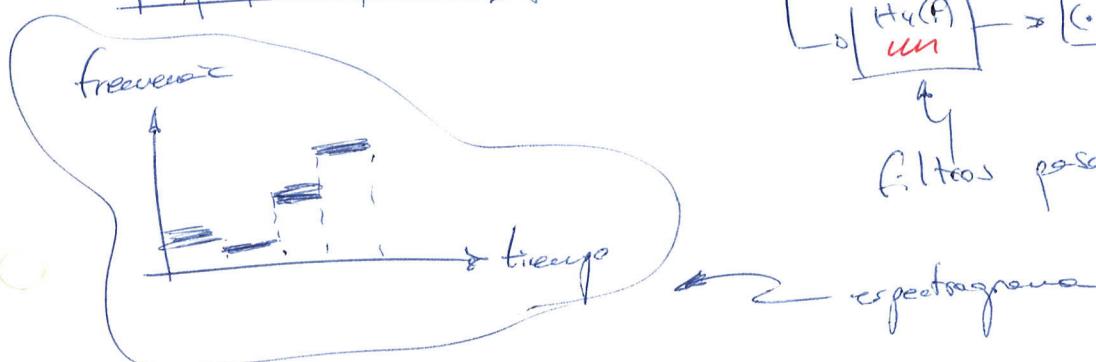
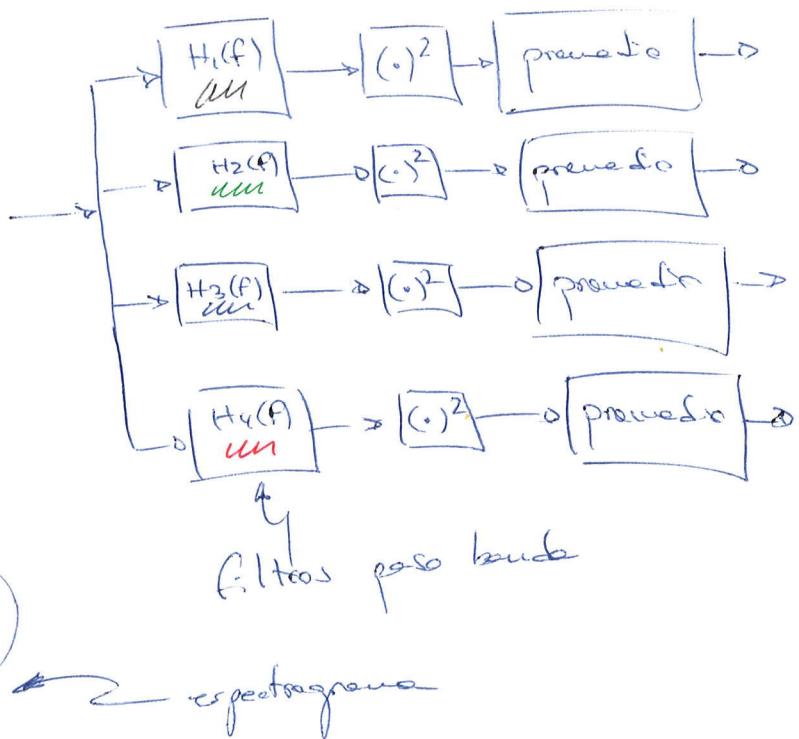
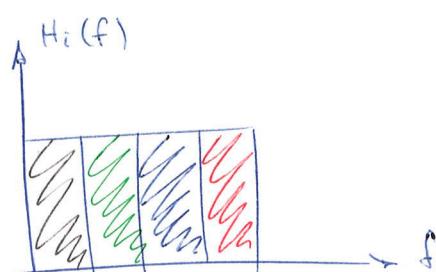
- Frecuencias sonoras no se identifican con el pitch
- Podemos variar dicho frecuencia: entonación, resalte de la voz, canto, ...

(x 2.6*) \rightarrow voces:

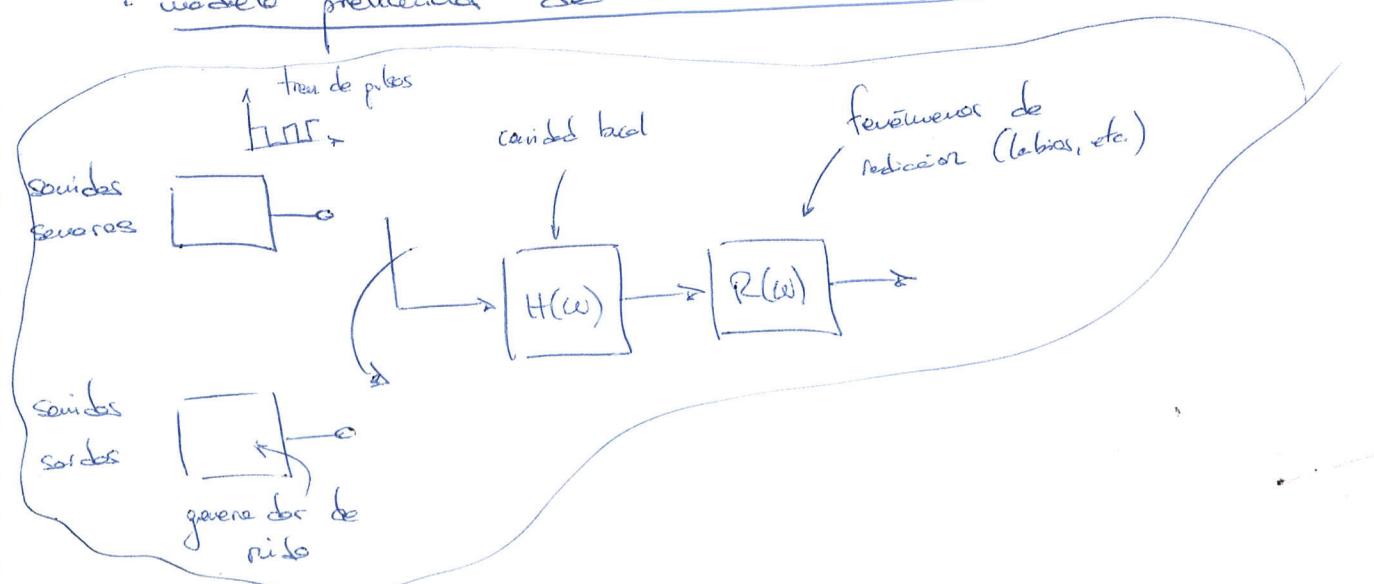


- espectrográmu: intento de análisis espectral de señales no estacionarias.

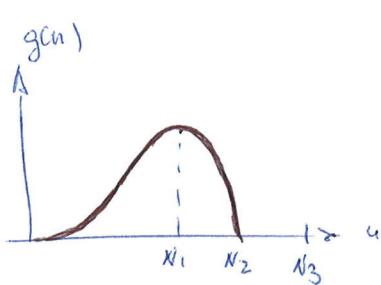
(x 2.7*)



- modelos preliminares de la señal de voz:



Pulso usado para modelar: "pulso de Rosenberg"



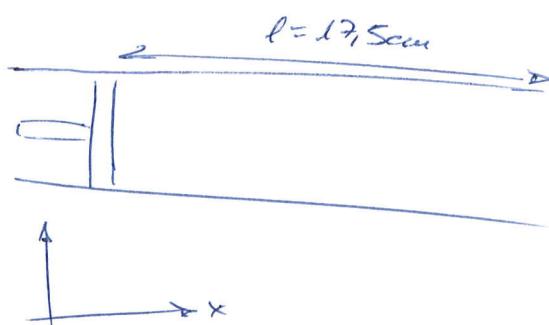
(*2.4*)

$$\left. \begin{array}{l} g(u) = \frac{1}{2} \left(1 - \cos \frac{\pi u}{N_1} \right) \quad 0 \leq u \leq N_1 \\ g(u) = \cos \pi \frac{u - N_1}{2N_2} \quad N_1 \leq u \leq N_1 + N_2 \\ g(u) = 0 \quad N_1 + N_2 < u \leq N_3 \end{array} \right\}$$

• modelo del tracto vocal: complejo, leyes de la física, etc.

Para cada frecuencia: hipótesis:

- Sección constante con el tiempo
- paredes sin pérdidas, no absorben energía
- consideramos cada pluma en la propagación de la onda por el tracto
- lo consideramos recto:



tubo acústico de sección circular constante

variables:

- velocidad: $v(x, t)$ m/s

- presión: $p(x, t)$

- densidad del gas: ρ

- área de la sección del tubo: A

- velocidad de propagación: $c = 20,1 \sqrt{T}$ (T en °K)

- velocidad volumétrica: $u(x, t) = A \cdot v(x, t)$ m³/s

ecuaciones:

$$-\frac{\partial u}{\partial x} = \frac{A}{\rho c^2} \frac{\partial p}{\partial t}$$

conservación de la masa

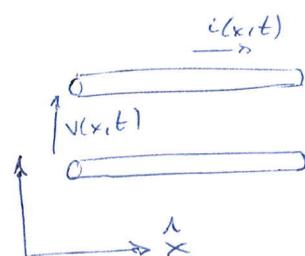
$$-\frac{\partial p}{\partial x} = \frac{\rho}{A} \frac{\partial u}{\partial t}$$

conservación del momento

(tubo uniforme \Rightarrow pérdidas)

análogas a las líneas de transmisión sin pérdidas:

$$\left. \begin{aligned} -\frac{\partial i}{\partial x} &= C \frac{\partial v}{\partial t} \\ -\frac{\partial v}{\partial x} &= L \frac{\partial i}{\partial t} \end{aligned} \right\}$$



L: inductancia / v.l. longitud

de esta forma:

gasos

línea fx

$$v \longleftrightarrow i$$

$$p \longleftrightarrow v$$

$$\begin{aligned} \frac{A}{\rho c^2} &\longleftrightarrow C && \text{el acústico se usa:} \\ \frac{p}{A} &\longleftrightarrow L && \boxed{C = \frac{A}{\rho c^2}} \quad \boxed{L = \frac{p}{A}} \end{aligned}$$

C: compresibilidad

L: inductancia acústica

$$u(x,t) = u^+(t - \frac{x}{c}) - u^-(t + \frac{x}{c})$$

$$p(x,t) = Z_0 \left(u^+(t - \frac{x}{c}) + u^-(t + \frac{x}{c}) \right)$$

donde

$$\boxed{Z_0 = \frac{\rho c}{A}}$$

impedancia acústica

ondas armónicas: entrada/excitación \Rightarrow señal exponencial compleja

$$u(0,t) = u(0)e^{j2\pi ft} = e^{j2\pi ft} \quad \begin{cases} \text{(condición inicial } u(0)=1) \\ \text{(considerando también } p(l,t)=0) \end{cases}$$

$$u(x,t) = u(x) e^{j2\pi ft} \quad \begin{cases} \text{condiciones de contorno:} \\ \quad \quad \quad u(0)=1 \\ \quad \quad \quad p(l)=0 \end{cases}$$

$$p(x,t) = P(x) e^{j2\pi ft}$$

tubo abierto

$$\begin{aligned} -\frac{\partial u(x)}{\partial x} &= \frac{A}{pc^2} p(x) e^{j2\pi f t} \quad \left\{ \begin{array}{l} \frac{\partial^2 u(x)}{\partial x^2} + \frac{4\pi^2 f^2}{c^2} u(x) = 0 \\ u(x) = 0 \end{array} \right. \quad \begin{array}{l} \text{ecuaciones} \\ \text{de onda} \end{array} \\ -\frac{\partial p(x)}{\partial x} &= \frac{P}{A} u(x) e^{j2\pi f t} \quad \left\{ \begin{array}{l} \frac{\partial^2 p(x)}{\partial x^2} + \frac{4\pi^2 f^2}{c^2} p(x) = 0 \\ p(l) = 0 \end{array} \right. \end{aligned}$$

Soluciones: $u(x) = a_1 e^{j\frac{2\pi f}{c}x} + a_2 e^{-j\frac{2\pi f}{c}x}$

$p(x) = a_3 e^{j\frac{2\pi f}{c}x} + a_4 e^{-j\frac{2\pi f}{c}x}$

dnde $\boxed{V = \sqrt{-\frac{4\pi^2 f^2}{c^2}} = j\frac{2\pi f}{c}}$

$$(\dots) = \pi \left\{ \begin{array}{l} p(x) = j\frac{2\pi f}{c} \frac{\sin \frac{2\pi f}{c}(l-x)}{\cos \frac{2\pi f}{c} l} \\ u(x) = \frac{\cos \frac{2\pi f}{c}(l-x)}{\cos \frac{2\pi f}{c} l} \end{array} \right.$$

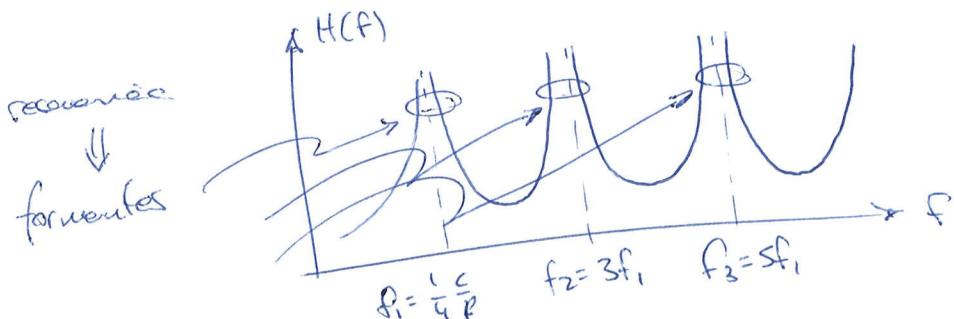
soluciones al exponer con exponente complejo

$$\boxed{\begin{array}{c} \text{---} \\ x=0 \quad x=l \end{array}}$$

$$u(0,t) = e^{j2\pi ft} \quad u(l,t) = U(l) e^{j2\pi ft}$$

responde:

$$\boxed{H(f) = \frac{U(l)}{U(0)} = \frac{1}{\cos \frac{2\pi f}{c} l}}$$



s: $\frac{2\pi f}{c} l = \frac{\pi}{2} + k\pi \Rightarrow$ polos

$$\Rightarrow f_i = \frac{1}{4} \frac{c}{l}$$

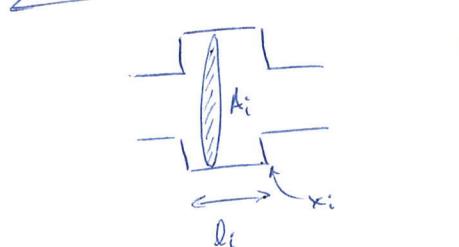
→ esto es para tubo recto, abierto, sin pérdidas

- paralelismo con líneas de transmisión terminadas a abierto.
- en la práctica siempre hay pérdidas \Rightarrow no será invertible
- señal excitadora: $S_e(f) \geq S_o(f) = S_e(f) |H(f)|^2$
- Esta $H(f)$ es propia de sistemas distribuidos

Caso v.i = potencia \Rightarrow potencia constante = U.P

- caracterización de tubos: tenemos ahora que caracterizar de tubos de diferentes longitudes y secciones

$\xrightarrow{\text{sistema LTI} \Rightarrow \text{resonancia en frecuencia}}$



Para cada tubo "i", tenemos:

$$\begin{aligned} \bullet u_i(x_i, t) &= u_i^+(t - x/c) - u_i^-(t + x/c) \\ \bullet P_i(x_i, t) &= 2i \left(u_i^+(t - x/c) + u_i^-(t + x/c) \right) \end{aligned}$$

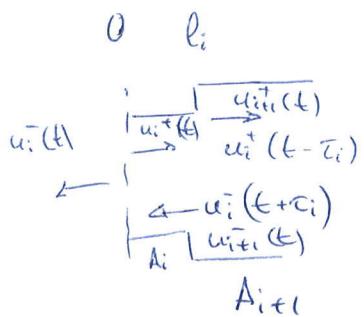
Nos interesa especialmente lo que pasa en las

transiciones \Rightarrow continuidad:

$$\begin{aligned} \bullet u_i(x_i, t) &= u_{i+1}(x_i, t) \\ \bullet P_i(x_i, t) &= P_{i+1}(x_i, t) \end{aligned}$$

Tiempo de propagación de la onda en un tubo:

$$\boxed{\tau_i = \frac{l_i}{c}}$$

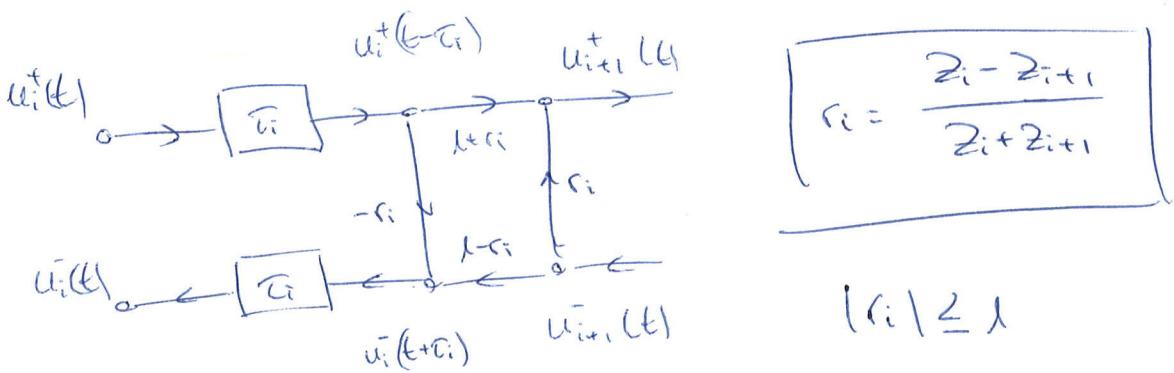


continuidad $u_i^+(t, l_i) = u_i^+(t - \tau_i) - u_i^-(t + \tau_i) = u_{i+1}^-(t, 0) = u_{i+1}^+(t) - u_{i+1}^-(t)$

$$P_i(t, l_i) = 2(u_i^+(t - \tau_i) + u_i^-(t + \tau_i)) = P_{i+1}(t, 0) = 2u_{i+1}^+(t) + u_{i+1}^-(t)$$

$$Z_i = \frac{PC}{A_i}$$

Representación de la onda en el espacio mediante grafos



El diagrama de grafos se puede expresar matricialmente:

$$\begin{bmatrix} u_{i+1}^+(t) \\ u_i^-(t+\tau_i) \end{bmatrix} = \begin{bmatrix} 1+r_i & r_i \\ -r_i & 1-r_i \end{bmatrix} \begin{bmatrix} u_i^+(t-\tau_i) \\ u_{i+1}^-(t) \end{bmatrix}$$

(ondas solientes)

(ondas entrantes)

Sistema completo \Rightarrow constancia de varios tubos

Necesitamos saber ademas que ocurre al principio y al final de la constancia completa

- Comienzo: glorios
- Final: lobios

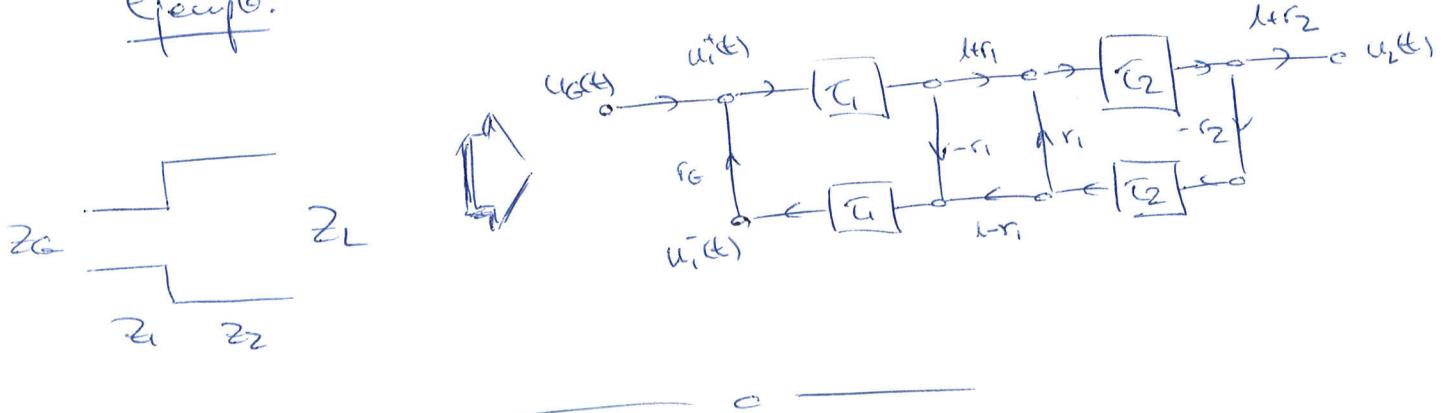
Lobios: Si Z_L es real, no hay onda reflejada

$$Z_L = \frac{\rho c}{A_L}$$
$$\Rightarrow L \rightarrow \infty$$

Glorios:

$$u_i(t) = u_i^+(t) - u_i^-(t)$$
$$P_i(t) = Z_1(u_i^+(t) + u_i^-(t))$$
$$u_i(t) = u_G(t) - \frac{P_i(t)}{Z_G}$$
$$\Rightarrow u_i^+(t) = \frac{1+r_G}{2} u_G(t) + r_G u_i^-(t)$$

Ejemplo:



Ya tenemos ecuaciones suficientes para calcular la respuesta en frecuencia γ_D ($\approx 2.19 - 2.20\omega$)

Lo hemos hecho suponiendo que son todos sistemas LTI, aunque la ganancia puede variar (órgano de impedancia variable), pero considerarlo complicaría innecesariamente el análisis.

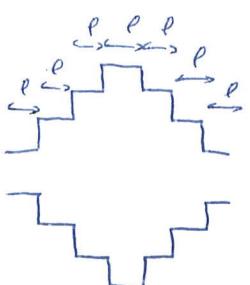
Ejercicio:

Responste en frecuencia de 2 tubos conestador.

Solucion: $H(f) = \frac{\frac{1}{2}(\lambda + r_G)(\lambda + r_1)(\lambda + r_2)}{\lambda + r_1 r_2 \gamma_2^2 + r_1 r_G \gamma_1^2 + r_G r_2 \gamma_1 \gamma_2}$

$$\gamma_i = e^{-j2\pi f t_i}$$

Si se particulariza para el caso de un modelo de tubos, todos con la misma longitud l , la respuesta en frecuencia que se obtiene es periódica de periodo $1/2c$



$$H(f) = H\left(f + \frac{k}{2c}\right) \quad c = \frac{l}{\lambda}$$

$H(f)$ periódica \Rightarrow sistema discreto

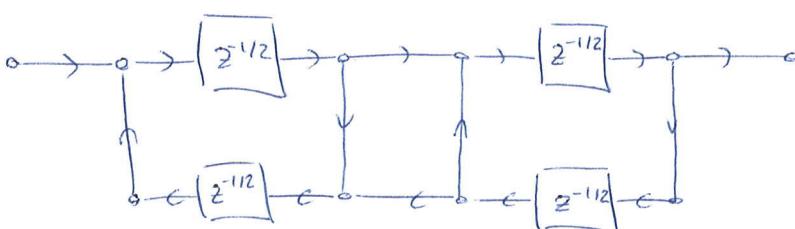
Ritmo de muestreo o periodo para este sistema:

$$\boxed{T_m = 2c}$$

Retardo de una muestra discreta: $T_m \rightarrow 2^{-l}$

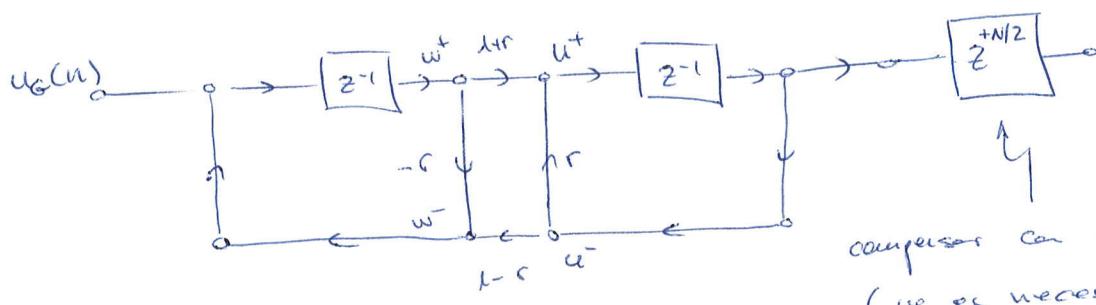
$$\text{Retardo } \tau = \frac{T_m}{2} \rightsquigarrow 2^{-1/2}$$

gráfico para la estructura digital



medie muestra
⇒ interpolador (muestreo)

En cada bucle cerrado tiene un z^{-1} , podemos agruparlos:



compensar con un adelanto
(no es necesario)

Sistema discreto equivalente a 2 filos muestra los
de longitud constante.

Resposta en frecuencia: ($\times 2.22\%$)

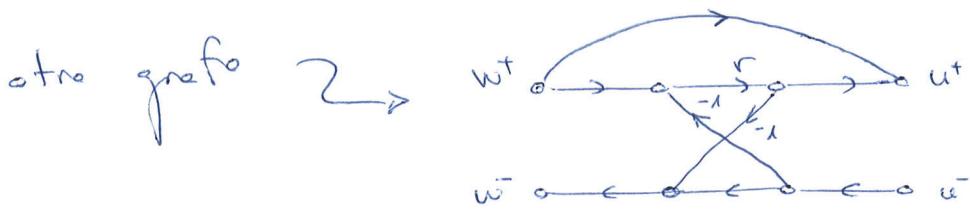
- Consideraciones sobre el modelo discreto:

- las operaciones que se realizan en un trinomial se pueden simplificar:

$$u^+ = (1+r)w^+ + ru^- = w^+ + r(w^+ + u^-)$$

$$w^- = -ru^+ + (1-r)u^- = u^- - r(w^+ + u^-)$$

reducir el número de multiplicaciones



• relación entre σ, L, N, c si $L = 17,5 \text{ cm}$, $c = 35 \cdot 10^3 \text{ cm/s}$

$$\ell = \frac{L}{N} \quad T = \frac{\ell}{c} = \frac{L}{N \cdot c} = \frac{T_m}{2}$$

$$\frac{N}{2} = \frac{L}{c} \frac{\ell}{T_m} = \frac{L}{2} 10^{-3} \text{ fm} = 10^{-3} B$$

$B = \frac{1}{2} f_m$ = banda que se está midiendo

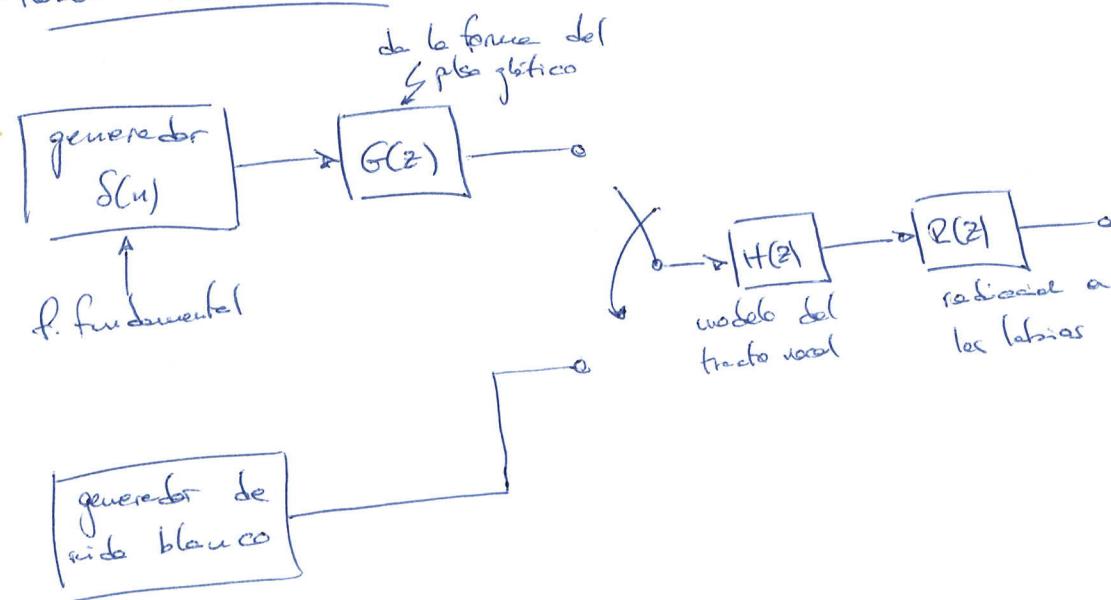
$$\boxed{\frac{N}{2} = B \{ \text{kHz} \}}$$

N / número de tubos

✓ número de picos \Rightarrow n° de picos de resonancia

Cada pico está a 1kHz del otro aproximadamente

-Modelo más completo:

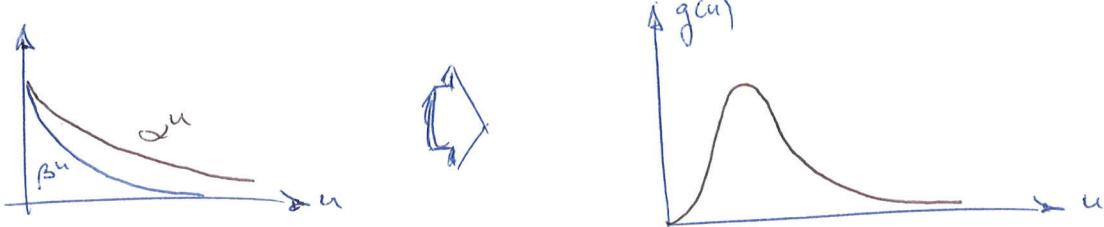


$g(u) = \underline{\text{pulso f\'etico}}$:

- Rosenberg

- de duraci\'on infinita y tipo exponencial:

$$\boxed{g(u) = (\alpha^u - \beta^u) u(u)} \quad \lambda > \alpha > \beta > 0$$



Ventaje de este tipo de pulsos: la respuesta impulsiva del sistema que los genere s\'olo tiene polos

- El modelo de tr\'ocho vocal se suele hacer con sistemas s\'olo polos. Los sonidos vocales tienen cu\~nas, pero se pueden modelar, con un orden suficiente, con sistemas s\'olo polos.

$$H(z) = \frac{1}{1 - \sum_{k=1}^N a_k z^{-k}}$$

- $R(z) = 1 - z_0 z^{-1} \quad z_0 \approx 1 \quad z_0 < 1 \quad (\underline{\text{paso alto}})$

Normalmente se modela s\'olo con $H(z)$, obteniendo $G(z)$ y $R(z)$.

2.- MÉTODOS DE ANÁLISIS LOCALIZADO (STP)

Se imponean díbles para analizar la señal de vez, y generales a el análisis de señales no estacionarias.

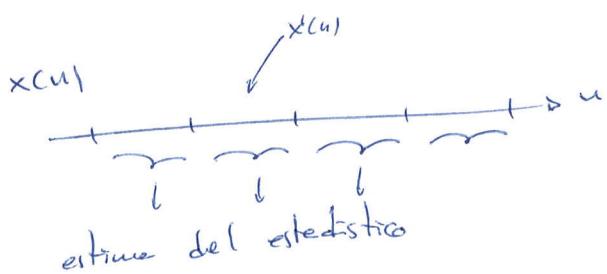
STP = Short Time Processing

Se basan en la dificultad de hacer análisis estadístico de señales no estacionarias. Se supone que la señal, en un segmento corto de duración, se comporta como una señal estacionaria.

El tamaño del segmento es determinante:

- desviación grande: señal muy variable sin estacionariedad
- desviación corta: señal con pocas muestras (mejor resolver temporal), estimación del estadístico corta.

La idea es basar segmentos en los que el estadístico no varía demasiado, pero en el que entra varios períodos de frecuencia fundamental.



• los segmentos no tienen porque ser contiguos, e incluso pueden solaparse

A veces puede interesar aplicar averiación a los segmentos:

$$S(u) = x'(u) w(u) \xleftarrow{TF} S(f) = \sum f_i * w(f_i)$$

y se extraen los estadísticos de $S(u)$

Usualmente: rectangular y Hamming

Para el análisis de la ventana se hay cinque reglas, y básicamente dependen del estadístico que queremos obtener. Requisitos:

- 1.- Que el estadístico no varíe mucho en el tiempo
- 2.- Que la ventana nos permita hacer una estimación suficientemente precisa
- 3.- Tener un número de ventanas suficiente para poder ver una propiedad como la frecuencia fundamental: habrá que ver más de 1 período

En voz, tiempos típicos de ventana van entre 3 y 30 ms

Duración de un foueuco típico: 80 ms

La duración del segmento puede estar determinada también por el tipo de sonido:

- silencio \rightarrow muy rápido
- vocal \rightarrow bastante estacionario

Aplicaciones típicas:

- Detección del comienzo y el final de los palabros y las frases
- Clasificar el tipo de sonido: silencio, silbido, sonido, ...
- Estimar la frecuencia fundamental o los fonemaes
- Estimar los parámetros del modelo

El análisis basado se puede distinguir en los dominios temporal y frecuencial.

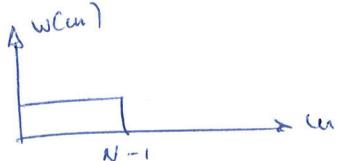
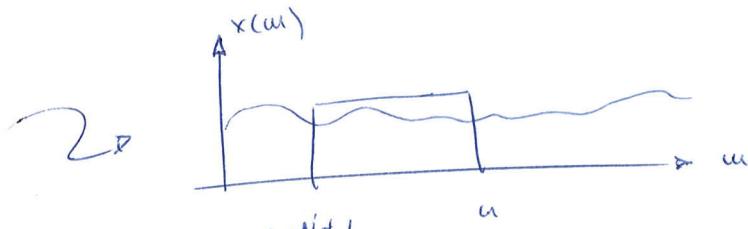
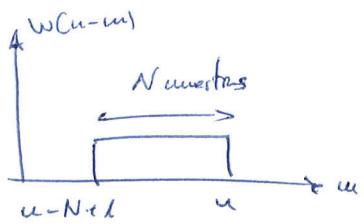
2.1.- ANÁLISIS LOCALIZADO EN EL TIEMPO

Características que veremos:

- Energía local o localizada
 - Magnitud local
 - Ceros por cero
 - Autocorrelación localizada
- } todas tienen una dependencia temporal
→ se subdivide x
- ENERGÍA LOCALIZADA:** ("energía localizada en el instante u ")

$$E_u = \sum_{m=-\infty}^{\infty} |x(mu) w(u-m)|^2$$

Ventana de tamaño finito, N :



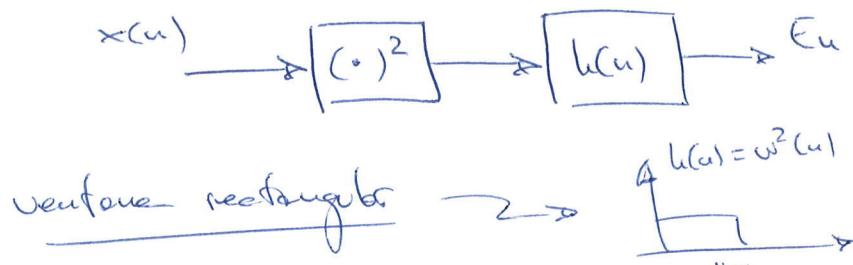
Se calcula la energía de las muestras que van de $u-N+1$ a u .

En una señal no estacionaria, E_u irá variando según el tiempo.

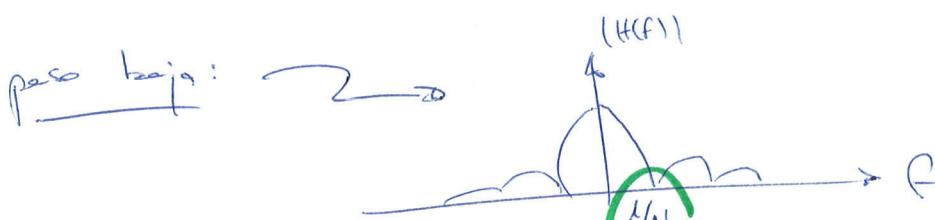
Conseguir ventanas de duración finita:

$$\left[E_u = \sum_{m=u-N+1}^u |x(mu) w(u-m)|^2 = \sum_{m=u-N+1}^u x^2(mu) w^2(u-m) \right]$$

$$h(u) = w^2(u) \rightarrow [E_u = x^2(u) * h(u)] \Rightarrow \text{filtro FIR}$$

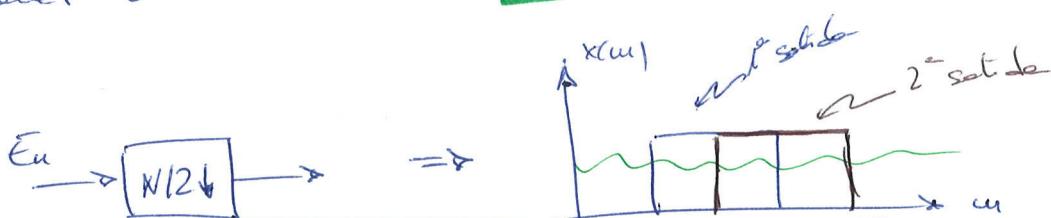


Rectangular y Hamming:
referido $\frac{N-1}{2}$ muestras
ca. fase lineal



En el paso bajo, tanto más ancho mayor sea N.

Las señales paso bajo adentro se cruzan, y en este caso es usual hacer un decaimiento en $N/2$:



se puede reconstruir la señal original tras el decaimiento

\Rightarrow no hay que computar todos los muestras en el filtro FIR

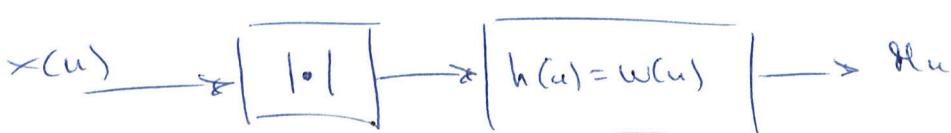
MAGNITUD LOCALIZADA:

$$H_u = \sum_{m=-\infty}^{\infty} |x(m)| w(u-m)$$

es interesante a veces cuando se usa aritmética en punto fijo, ya que reduce el rango dinámico:

$$\frac{\max |H_u|}{\min |H_u|} = \sqrt{\frac{\max |E_u|^2}{\min |E_u|^2}}$$

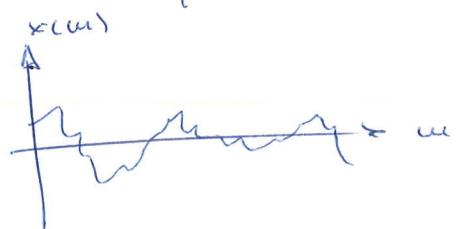
En la señal de voz, E_u puede variar varios decenas de dB



- CRUCES POR CERO:

Sigue para llevar una cuenta de las veces que una señal pasa por cero. Consideraré muy bien ciertas señales (ej: \approx sinusoidal)

Se utiliza porque es fácil de calcular y además es significativa del tipo de trama de voz que se está analizando.



$$Z_u = \sum_u \frac{1}{2} |\text{sign}(x(u)) - \text{sign}(x(u-1))| w(u-u)$$

(de el número de cruces por cero en la ventana.)

Para que la cuenta sea exacta, la ventana debe ser rectangular de valor $1/N_c$ y nos dará el número de cruces por cero por cada muestra.

Se puede obtener el número de cruces por cero de una sucesión:

$f_{su} \rightarrow$ frecuencia de muestras

$f_0 \rightarrow$ frecuencia de la sucesión

$$\frac{f_{su}}{f_0} = \# \text{ de muestras por periodo de la sucesión}$$

En un periodo de la sucesión hay 2 cruces por cero:

$$\frac{2}{f_{su}/f_0} = \frac{2f_0}{f_{su}} = \text{tasa de cruces por cero por cada muestra}$$

(para una sucesión)

Se puede intentar estimar su frecuencia viendo sus cruces por cero.

- Sonidos sonoros: Ent, Znt

- Sonidos sordos: Ent, Znt

- Ruido de fondo: Ent, Znt

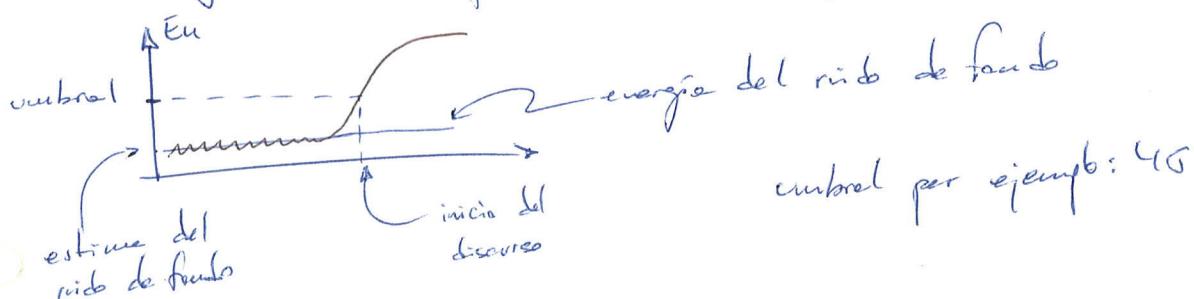
Importante para distinguir sonidos
sordos del ruido ambiente

APLICACIÓN: DISCRIMINACIÓN VOZ/SILENCIO EN UN DISCURSO HABLADO:

Determinar donde el factor rebota y da de la voz silencio

Caso 1: grabación en estudio

Se puede hacer de manera que sea fácil determinar dónde hay silencios a partir de la energía localizada.



ambient por ejemplo: 4G

No va a ser muy largo, o perderemos señalización temporal.
Esto presenta problemas en ambientes ruidosos.

Ejemplos / energia local \rightarrow ($\times 2.22 \times$)
 magnitud local \rightarrow ($\times 2.23 \times$)
 cruces por cero \rightarrow ($\times 2.23 \times$)

| ~ 10 ms de duración
| de la ventilación

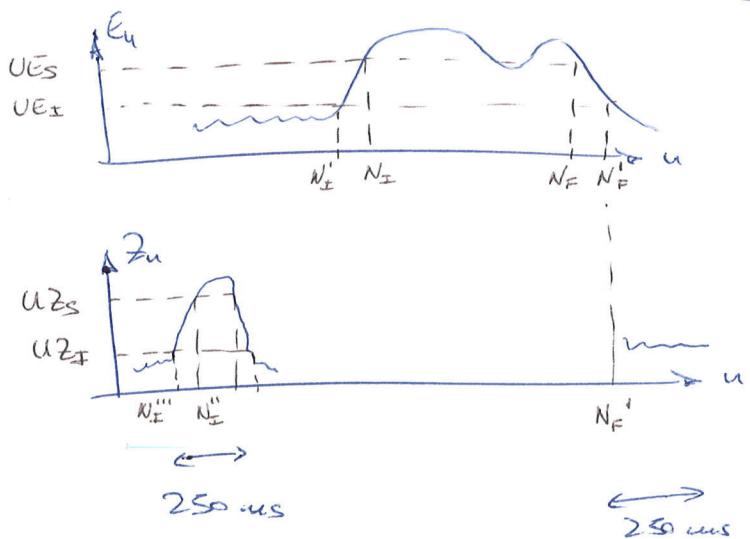
Caso 2: ambiente ruidoso

Fuentes con poca energía pueden confundirse con el ruido ambiente, por lo que la energía localizada no es suficiente.

Se usan la magnitud o la energía total y los cruces por cero para distinguir entre los problemas.

	(Energía)	(Componentes espectrales)
	E_u	Z_u
Sonidos suaves	Alta	Pocos
Sonidos sordos	Baja	Muchos
Ruido de fondo	Baja	~

→ mayor tasa de cruces por cero que el ruido de fondo



La estímulación de las magnitudes E_u y Z_u del ruido de fondo, nos permite establecer vibraciones de energía U_E (Superior e inferior) y de cruces por cero $U2$ (Sup. e inf.)

Superior → muy conservador, el ruido de fondo nunca lo alcanza

$N_I^, N_F^$ instantes de inicio y final, muy conservadores

Ahora se puede trazar el intervalo inferior que el ruido puede superar, pero ya teniendo conocimiento de dónde hay polabro con calzado → $N_I^{'}, N_F^{'}$

Aún queda decir que los sonidos en baja energía

⇒ cruces por cero

$U2_S$ → muy conservador también

Se observan los cruces por cero desde el comienzo de la palabra hasta que transcurren 10 segmentos de voz de 25ms (250ms) $\Rightarrow N_f''$

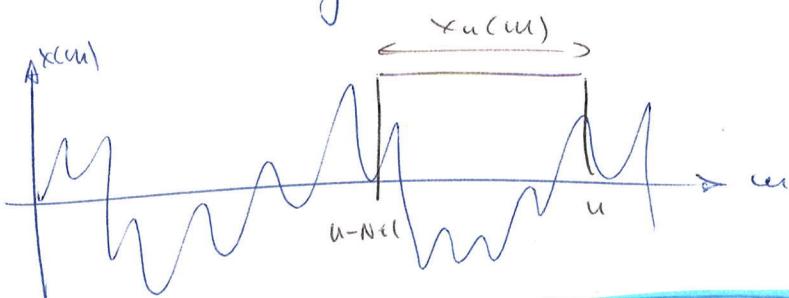
$$U_{Z_f} \geq N_f''$$

Al final, se observan 250ms desde N_f' y si no hay cruces por cero por encima del umbral superior, se \Rightarrow por lo tanto ese instante como final de palabra (que lo supera U_{Z_f})

Primero se mira la energía, ya que las palabras tienen vocales, que son de alta energía \Rightarrow energía más significativa que cruces por cero

- Autocorrelación Localizada:

Es el estimador seguido de la autocorrelación del segmento sin normalizar



$$R_u(k) = N r_{x_u}(k) = \sum_m x_u(m) x_u(m+k) = X_u(k) * X_u(-k)$$

Suponemos $k \geq 0$ siempre

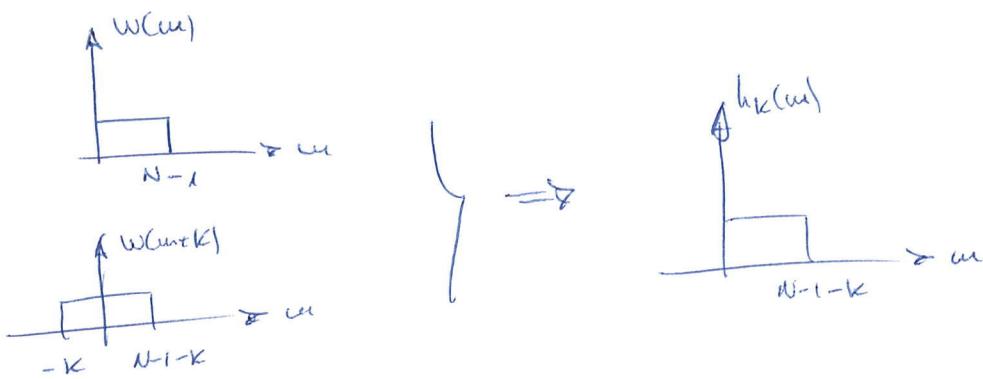
$$x_u(u) = x(u) w(u-u) \quad \left. \begin{array}{l} \\ \end{array} \right\} \Rightarrow \text{TF } \left. \begin{array}{l} R_u(k) \end{array} \right\} = X_u(f) * X_u(-f) = |X_u(f)|^2$$

$X_u(f) = \text{TF } \left. \begin{array}{l} x(u) \end{array} \right\}$
transformada de Fourier localizada

señales reales: $X_u(f) = X_u^*(f)$

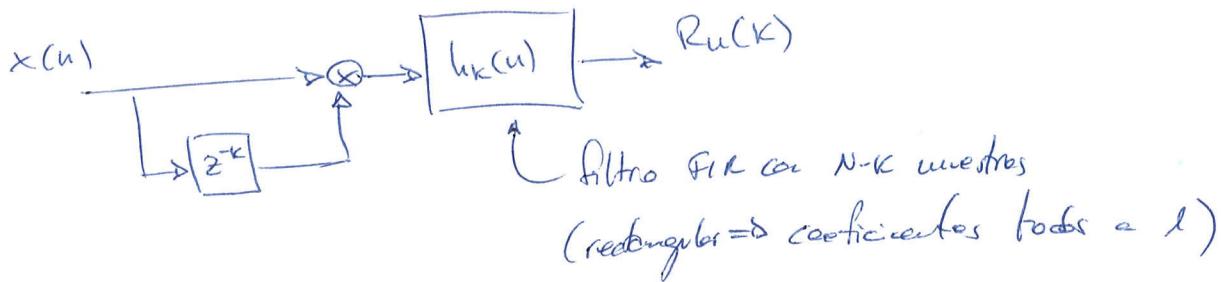
$$\begin{aligned} R_u(k) &= \sum_m x(m) w(m-u) x(m-k) w(m-m+k) = \left\{ h_k(m) = w(m) w(m+k) \right\} = \\ &= \sum_m x(m) x(m-k) h_k(m-u) \end{aligned}$$

Si tenemos la ventana rectangular:

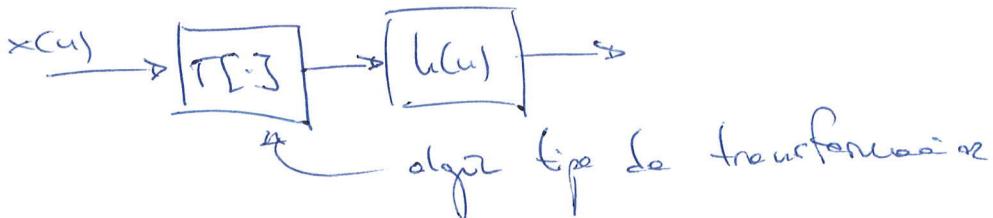


$$\Rightarrow R_{u,k} = (x(u) \cdot x(u-k)) * h_k(u)$$

Podemos considerar que la autocorrelación (realizada es:



Todos estos weightides responden a este estructura:



	$h(u)$	$T\Sigma J$
E_u	$w^2(u)$	$(\cdot)^2$
p_{lu}	$w(u)$	$ \cdot $
Z_u	$w(u) = 1/N$	$\frac{1}{2}(\text{sign}(x(u)) - \text{sign}(x(u-1)))$
$R_{u,k}$	$w(u) w(u+k)$	$x(u) \cdot x(u-k)$

- Problemas de la autocorrelación localizada:

La duración del filtro $R_u(k)$ depende de k ($N-k$ muestras)

$$E[R_u(k)] = (N-|k|)\gamma_x(u)$$

• La extensión de la esperanza de la autocorrelación localizada decrece con k

Sinal periódico \Rightarrow autocorrelación periódica en k

Una de las aplicaciones más típicas de la autocorrelación localizada es la estimación de la frecuencia fundamental de señales sonoras (señales periódicas), en las que la autocorrelación presenta picos periódicos.
(* 2.25 *)

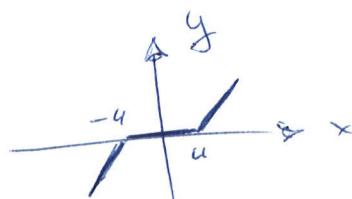
Los picos decrescen porque:

• la señal no es físicamente periódica

• al estimar la autocorrelación llevas visto que su esperanza decrece al crecer el índice k (enfocando tríangular de la autocorrelación)

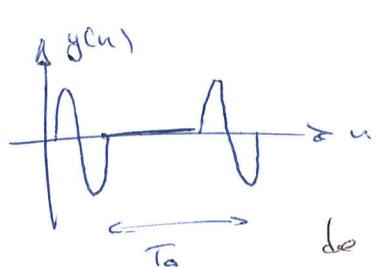
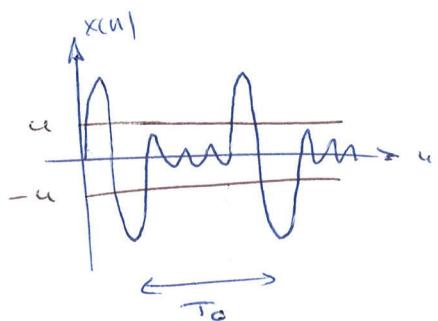
La ventana tiene que ser suficientemente grande para observar varios períodos de la señal, al menos 2.

- Central differencing (recolección central): es un dispositivo que hace una transformación no lineal sin memoria



Para valores de x pequeños, hasta un umbral u , la salida es cero, y a partir de él, es proporcional a x

(* 2.24*) \Rightarrow se reduce la frecuencia del fundamental
(la periodicidad de la señal)



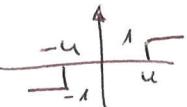
ventanas pequeñas \Rightarrow picos de $R_{xy}(k)$

intervalos grandes, se pierden confundir con el fundamental \Rightarrow (* 2.25*, 4.26c)

También pasa con fuentes que varían rápidamente

Reactor central de 3 niveles \Rightarrow simplifica el cálculo
de la autocorrelación y no degrada la detección del
pínteh

(Rob-Schlosser p. 154)



- Auto-correlación localizada modificada:

$$\hat{R}_{xy}(k) = \sum_m x_u(m) y_{u-k}(m-k)$$

Se miden 2 segmentos distintos
 $x_u(m)$ y $y_{u-k}(m)$

✓ se pierde la simetría de la autocorrelación

✗ no es una estimación de la autocorrelación, ya que son segmentos distintos:

- } x pierde la simetría
- } x pierde las propiedades de su TF

(* 2.26*)

⊗ \Rightarrow los picos no doblan

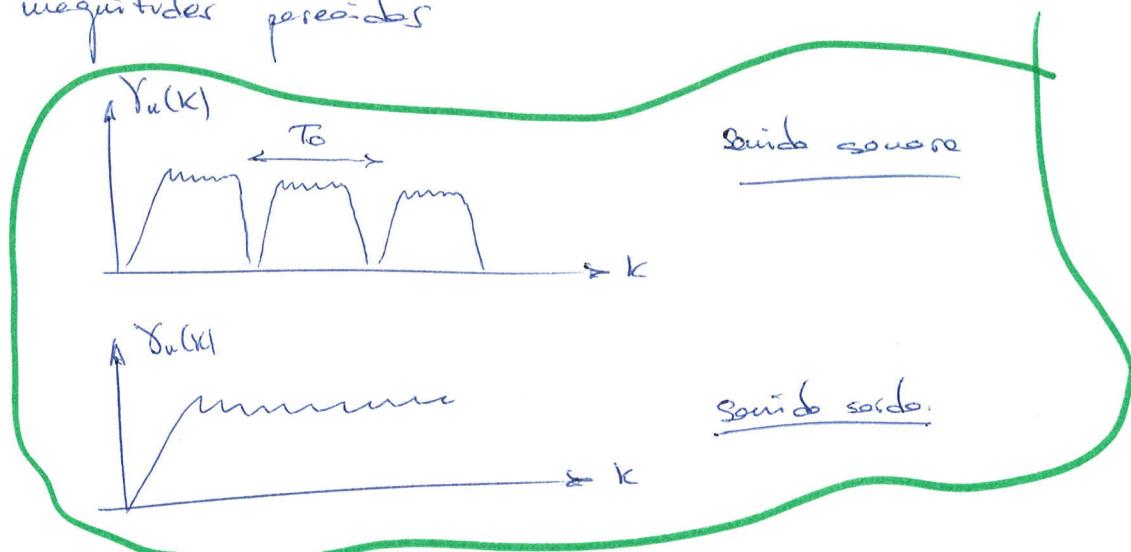
\Rightarrow es una correlación cruzada de los dos segmentos

Autocorrelación \Rightarrow multiplicar \Rightarrow constante
 (en punto fijo)

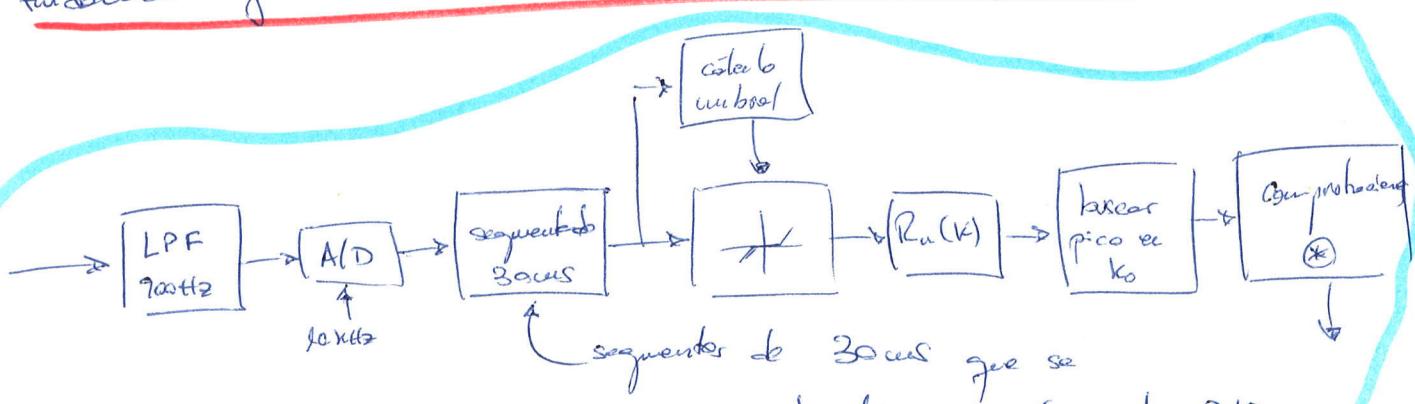
- DIFERENCIA DE MAGNITUD LOCAL:

$$Y_u(k) = \sum_m |x_{ul}(m) - x_{u-k}(m)|$$

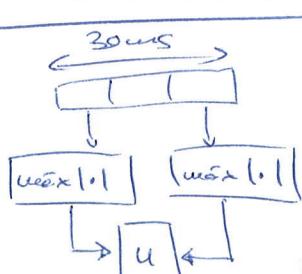
Restamos 2 segmentos separados k , por lo que si k coincide con el periodo, tenemos un mínimo y si no tenemos magnitudas periódicas



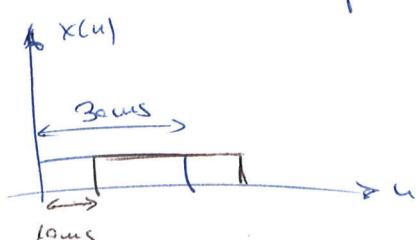
Con lo que hemos visto, hay un esquema que permite calcular la frecuencia fundamental y determinar si un sonido es sordo o suave:



- Cálculo del umbral:



$$U = 0,68 \sin(m) \quad \{ \max_1, \max_2 \}$$



Comprobaciones:

$$\left. \begin{array}{l} R(k_0) > 0,3 R(0) \\ \text{ko valor adecuado} \Rightarrow \text{sínd. sonoro} \\ (\text{periodo fundamental} \\ \text{pequeño}) \\ \text{sólido} = k_0 \end{array} \right\}$$

Si se cumple una de ellas \Rightarrow sínd. sonoro

2.2.- ANÁLISIS LOCALIZADO EN LA FRECUENCIA

Magnitud estrella:

- TRANSFORMADA LOCALIZADA DE FOURIER (STFT = Short Time Fourier Transform):

Es la TF del segmento $x_u(u) = x(u) w(u-u)$

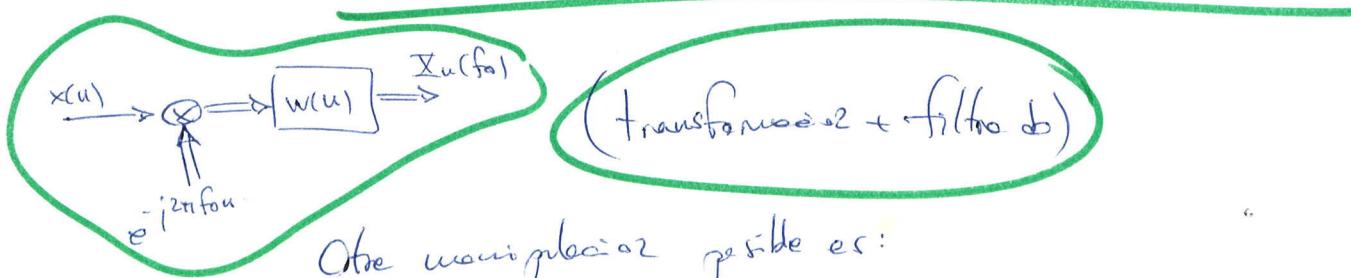
$$\boxed{\begin{aligned} X_u(f) &= \text{TF } \{x_u(u)\} = \sum_u x(u) e^{-j2\pi f u} \\ &= \sum_u x(u) w(u-u) e^{-j2\pi f u} \end{aligned}}$$

- Propiedades:

a) Vemos antes que $\text{TF } \{R_u(k)\} = |\bar{X}_u(f)|^2$

b) Si fijamos una cierta frecuencia: $f = f_0$,

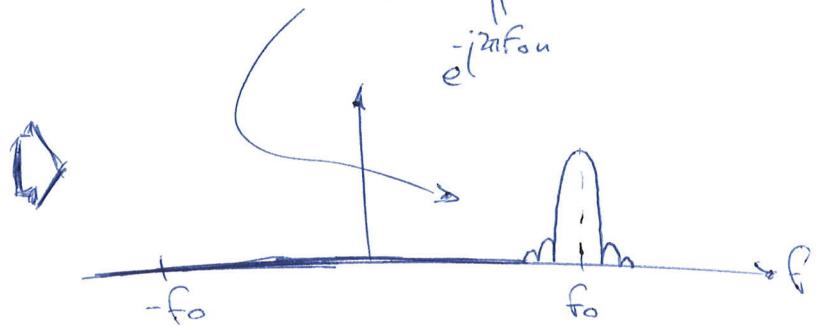
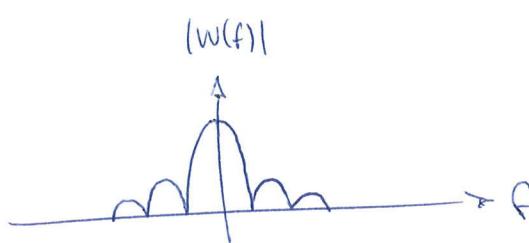
$$X_u(f_0) = \sum_u (x(u) e^{-j2\pi f_0 u}) w(u-u) = (x(u) e^{-j2\pi f_0 u}) * w(u)$$



Otra manipulación posible es:

$$X_u(f_0) = e^{-j2\pi f_0 u} \sum_u x(u) w(u-u) e^{j2\pi f_0 (u-u)} =$$

$$= x(u) * \left(w(u) e^{j2\pi f_0 u} \right) \Rightarrow x(u) \xrightarrow{w(u) e^{j2\pi f_0 u}} \otimes \Rightarrow X_u(f_0)$$



1 base pasa bajo las frecuencias en torno a f_0 y filtra

2 al revés

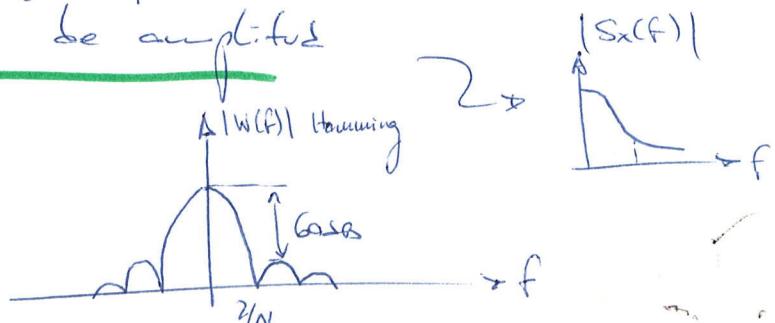
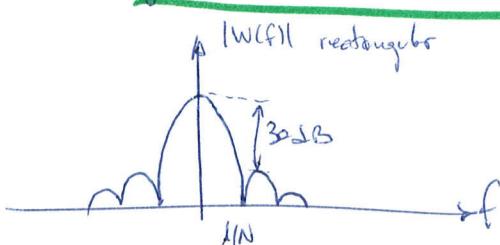
$x_u(u)$ podemos obtener mediante transformada inversa:

$$x_u(u) = \mathcal{F}^{-1} \{ X_u(f) \} = \int_{-1/2}^{1/2} X_u(f) e^{j2\pi fu} df$$

$$x_u(u) = x(u)w(u-u) \Rightarrow x(u)w(0) = x_u(u)$$

$$x(u) = \frac{1}{w(0)} x_u(u) = \frac{1}{w(0)} \int_{-1/2}^{1/2} X_u(f) e^{j2\pi fu} df$$

En el análisis espectral se suele preferir la ventana de Hamming para hacer el análisis localizado de Fourier, ya que el espectro de la señal de voz tiene componentes en diferentes bandas con diferentes amplitudes



Hannning tiene menor resolución, pero mayor atenuación de los lóbulos secundarios.

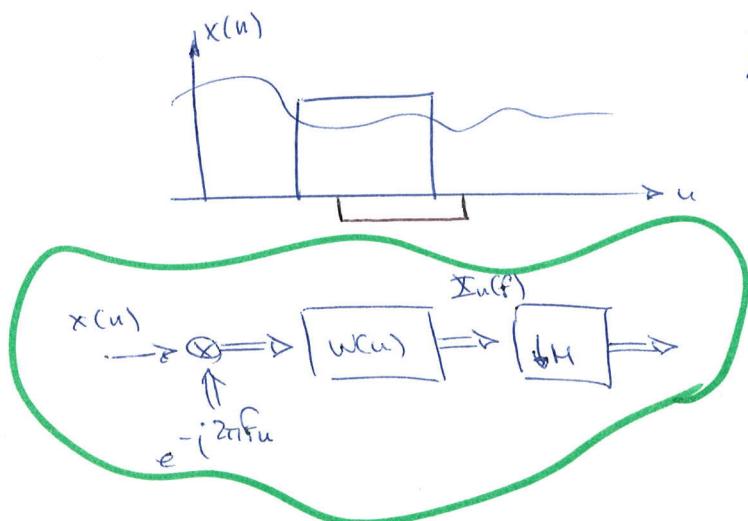
El tamaño de la ventana será suficientemente grande para ver las características periódicas, como la detección del pitch (freg. fundamental)

- Detección del pitch: N grande

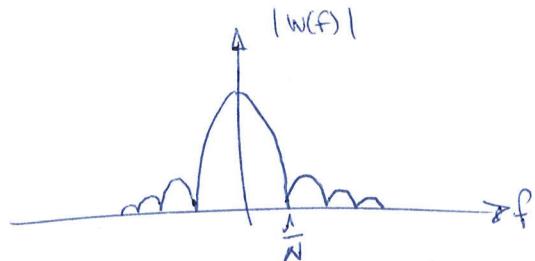
frecuencia de muestreo: $f_m = 10 \text{ kHz} \Rightarrow N = 220$
 $\Rightarrow 72 \mu\text{s}$

$$B_w = 45 \text{ Hz}$$

- buena resolución en frecuencia
- poca resolución temporal
 $(\star 2.28 - 2.31 \star)$



¿rápidos de desplazamiento de la ventana?



paso bajo \Rightarrow admite desplazos
 \Rightarrow desplazar más de 1 muestra

Como el FPR no es ideal, habrá aliasing, que consideremos despreciable si el ancho de banda es suficiente.

• Rectangular: $B = 1/N \rightarrow f_m' = \frac{2}{N} f_m \Rightarrow \frac{N}{2} \downarrow \quad (M = \frac{N}{2})$

• Hannning: $B = \frac{2}{N} \downarrow \rightarrow f_m' = \frac{4}{N} f_m \Rightarrow \frac{N}{4} \downarrow \quad (M = \frac{N}{4})$

Comprobación: \mathcal{F} (en el tiempo / frecuencia continua) \Rightarrow una señal constante

$$x_u(u) \xrightarrow{\text{TF}} X_u(f)$$

(duración N muestras \Rightarrow admite DFT)

$$x_u(u) \xleftarrow{\text{DFT}_N} X_u(k) = X_u(f) \Big|_{f=\frac{k}{N}} \quad k \in \mathbb{Z}$$

$$(\text{DFT}_N \equiv \text{st-DFT})$$

$$\boxed{X_u(k) = \sum_u x_u(u) e^{-j\frac{2\pi}{N}ku}} \quad | \quad k = 0 : N-1$$

| st-DFT

- BANCO DE FILTROS:

Es una interpretación de la st-DFT

• IDFT: $x_u(u) \xleftarrow{\text{DFT}} X_u(k)$

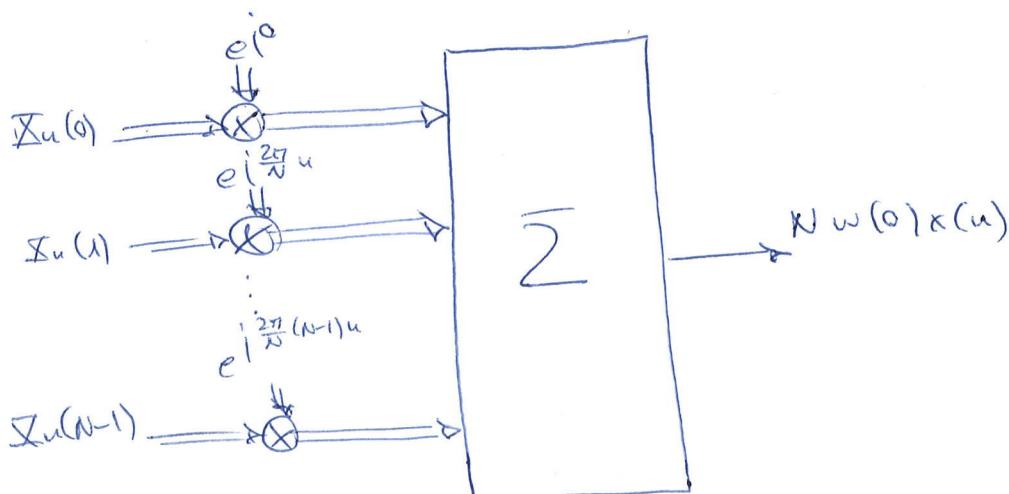
$$x_u(u) = \text{DFT}^{-1} \{ X_u(k) \} = \frac{1}{N} \sum_{k=0}^{N-1} X_u(k) e^{j\frac{2\pi}{N}ku}$$

$$x_u(u) = x(u) w(u-u)$$

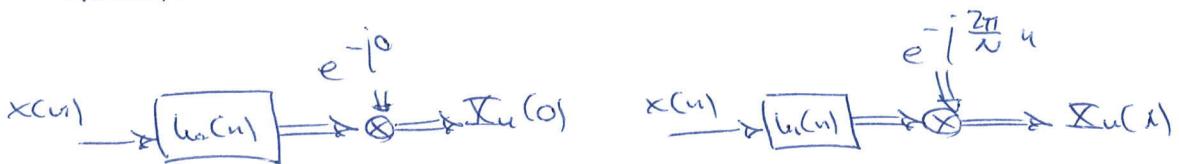
$$\hookrightarrow x_u(u) = x(u) w(u) = \frac{1}{N} \sum_{k=0}^{N-1} X_u(k) e^{j\frac{2\pi}{N}ku}$$

$\underbrace{x(u) \cdot N \cdot w(0)}$

Se hace una modulación de $X_u(k)$:



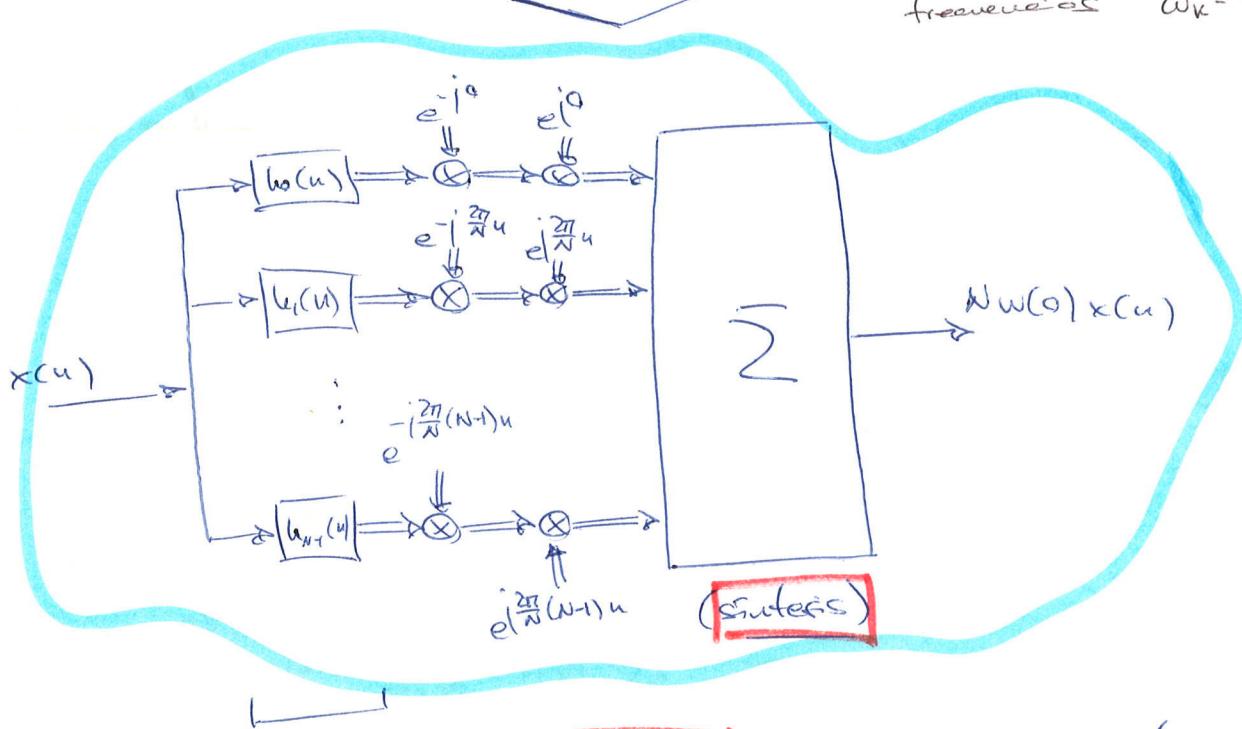
Transformada localizada a una frecuencia:



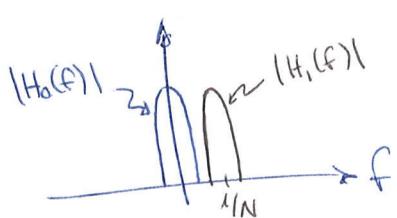
$$h_k(u) = w(u) e^{j \frac{2\pi}{N} k u}$$

Ventana (longitud N)

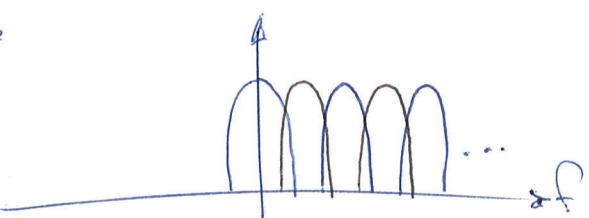
$\Rightarrow X_u(f)$ se muestra en las frecuencias $\omega_k = \frac{2\pi k}{N}$ $k=0:N-1$



banco de filtros (análisis) \Rightarrow hace un análisis de la señal



la señal se divide en bandas



aplicación: codificación subbandas, para quitar componentes innecesarios en voz y audio

Ejemplo: Estimador de energía local

$$e(u) = \sum_{m=u-N+1}^u x^2(m)$$

a) Buscar una expresión recursiva para este estimador

$$e(u) = e(u-1) + \text{algo} \quad \swarrow a(u)$$

$$\begin{aligned} a(u) &= e(u) - e(u-1) = \sum_{m=u-N+1}^u x^2(m) - \sum_{m=u-N}^{u-1} x^2(m) \\ &= x^2(u) - x^2(u-N) \end{aligned}$$

$$\Rightarrow \boxed{e(u) = e(u-1) + x^2(u) - x^2(u-N)}$$

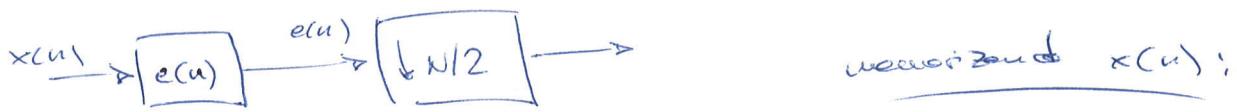
b) Lo mismo pero usando una ventana de Hamming

$$e(u) = \sum_{m=u-N+1}^u x^2(m) w^2(u-m)$$

$$\begin{aligned} a(u) &= e(u) - e(u-1) = \sum_{m=u-N+1}^u x^2(m) w^2(u-m) - \sum_{m=u-N}^{u-1} x^2(m) w^2(u'-m) = \\ &= x^2(u) w^2(0) - x^2(u-1-N) w^2(N-1) + \sum_{m=u-N+1}^{u-1} x^2(m) (w^2(u-m) - w^2(u-1-m)) \end{aligned}$$

→ sumar $N+1$ términos → es más complejo usar una recursión con una ventana no rectangular que usar la expresión directa

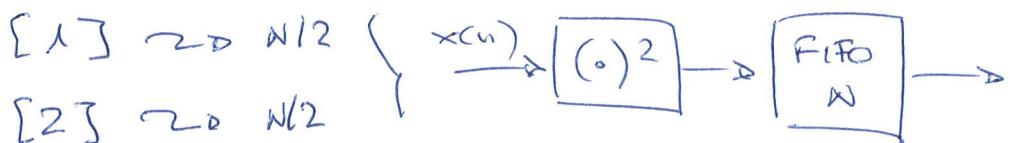
c) Estimación del coste computacional en el caso inicial para la expresión directa [1] y la recursiva [2] si la estimación de energía se calcula cada $N/2$ muestras de la entrada (dividiendo por $N/2$). Claro, teniendo en cuenta sólo los productos



[1] \Rightarrow en cada cálculo requiere N productos

[2] \Rightarrow 2 productos por iteración, $N/2$ veces para la señal
cada $N/2$ muestras $\Rightarrow N$ productos

Otra opción es memorizar $x^2(u)$



d) Tenemos un estimador de la autocorrelación local definido como:

$$r(k; u) = \frac{1}{N-|k|} \sum_{m=u-N+1+|k|}^u x(m)x(u-|k|) \quad (\text{inseguido})$$

$$r(0; u) = \frac{1}{N} \sum_{m=u-N+1}^u x^2(m) = \frac{e(u)}{N}$$

e) Calcular el sesgo y la varianza de $e(u)$:

$$\left. \begin{aligned} x(u) &\leftarrow \text{AWGN} \\ E[x^2(m)x^2(u)] &= \bar{x}^4 (\delta(m-u) + 1) \\ \sigma_x^2 &= \text{Var}(x(u)) = \bar{x}^2 \end{aligned} \right\}$$

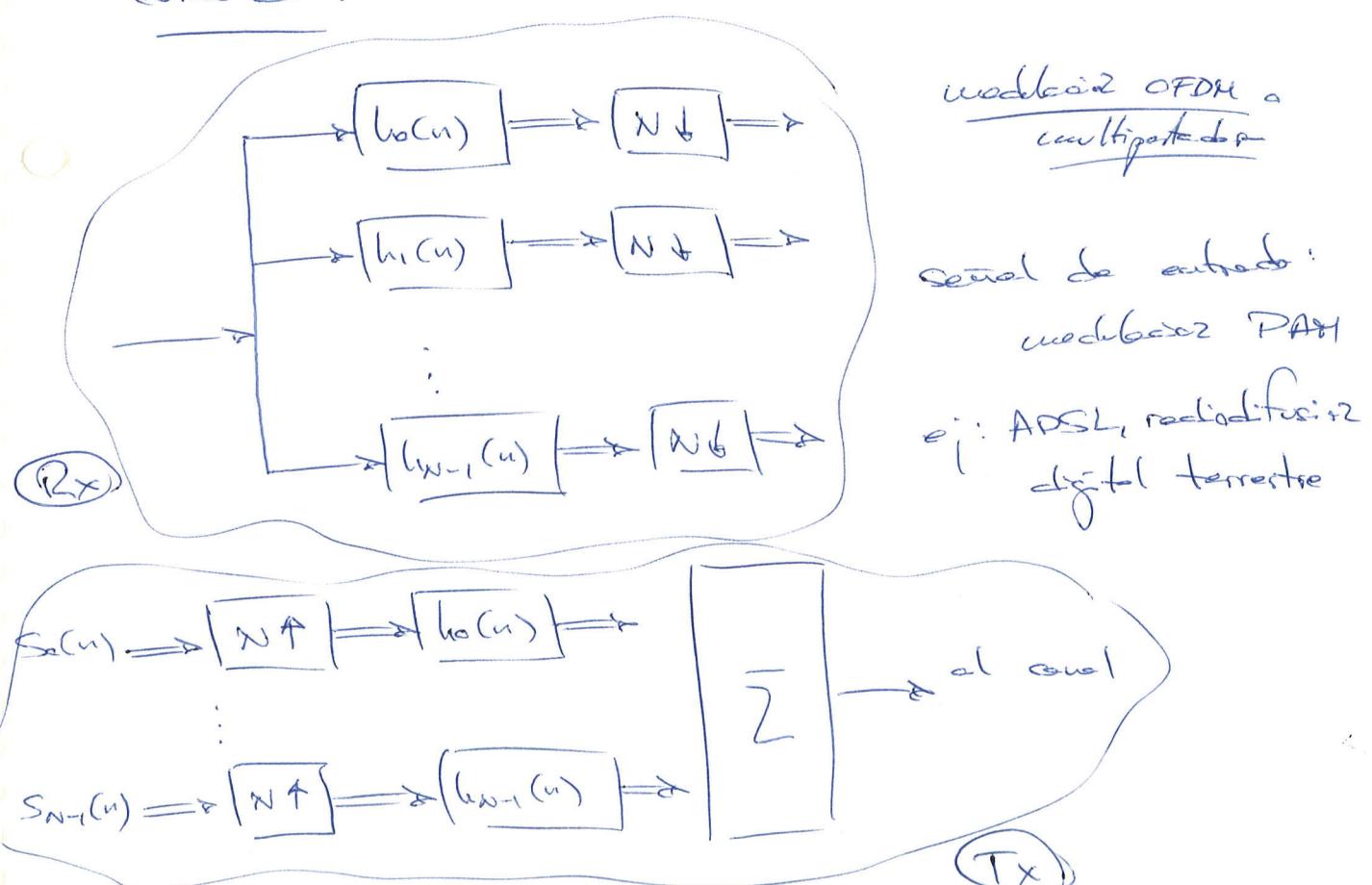
$$E[e(u)] = \sum_{m=u-N+1}^u E[x^2(m)] = N\bar{x}^2 \quad \text{es inseguido}$$

$$\text{Var}(e(u)) = E[e^2(u)] - E^2[e(u)] = E[e^2(u)] - N^2\bar{x}^4$$

$$\begin{aligned}
 E[e^2(u)] &= E\left[\sum_{m=u-N+1}^u \sum_{l=u-N+1}^u x^2(m)x^2(l)\right] = \\
 &= \sum_m \sum_l E[x^2(m)x^2(l)] = \\
 &= \sum_m \sum_l \sigma_x^4 (S(m-l)+1) = \\
 &= \sum_m \sigma_x^4 + \sum_m \sum_l \sigma_x^4 = (N+N^2) \sigma_x^4
 \end{aligned}$$

$$\Rightarrow \boxed{\text{Var}(e(u)) = N \sigma_x^4}$$

Curiosidad:





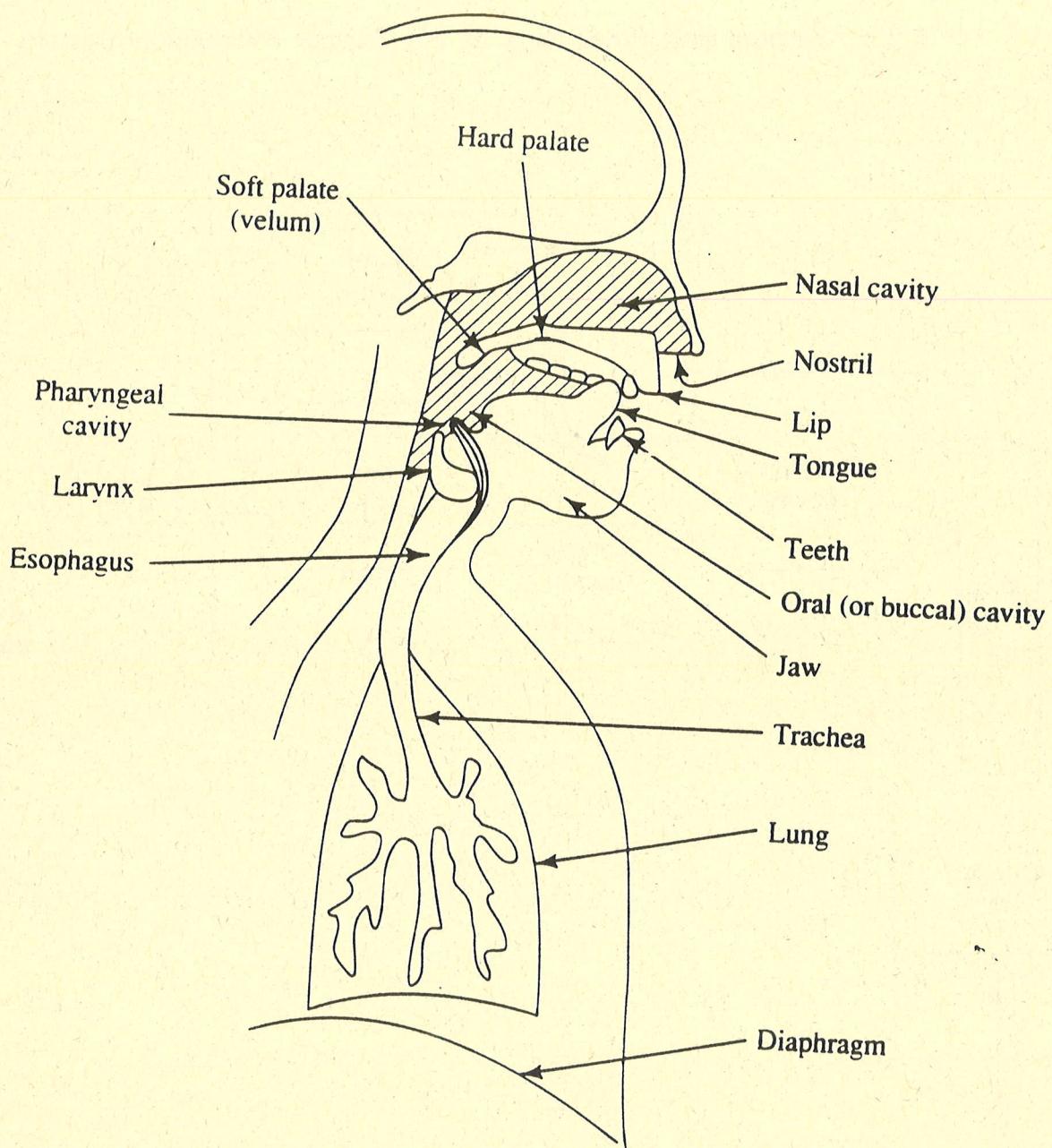


FIGURE 2.1. A schematic diagram of the human speech production mechanism.

2.2 / Anatomy and Physiology of the Speech Production System 103

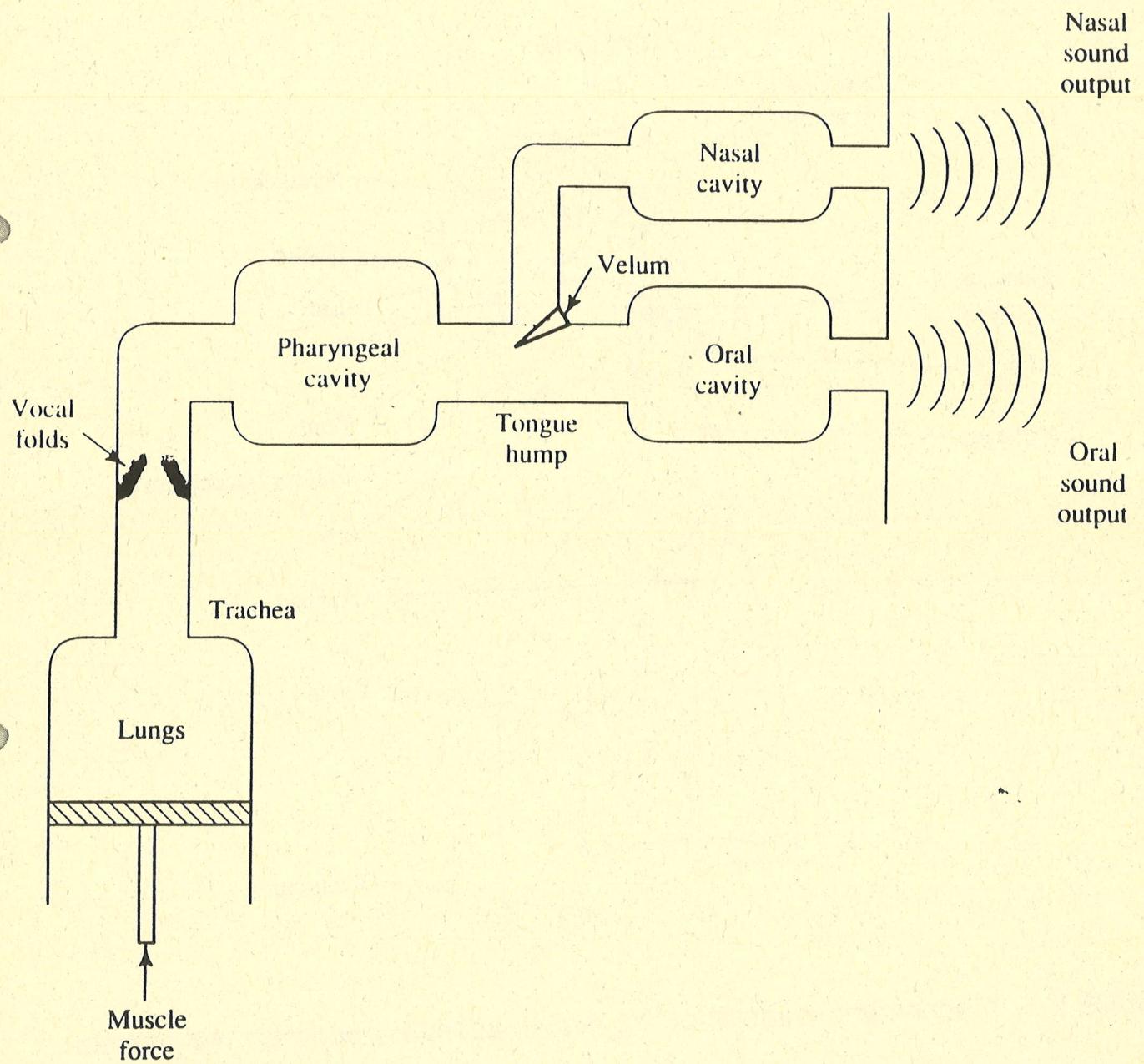


FIGURE 2.2. A block diagram of human speech production.

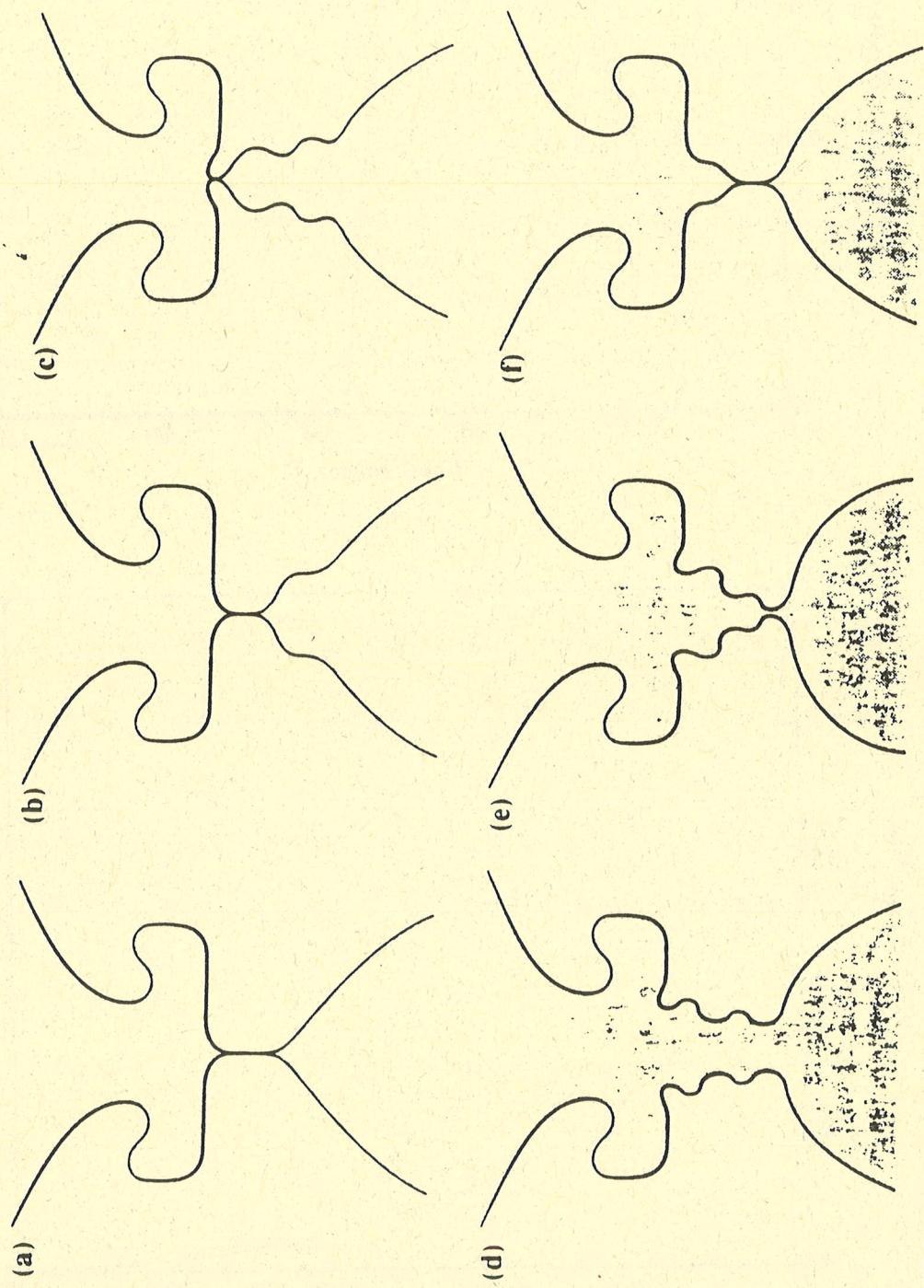
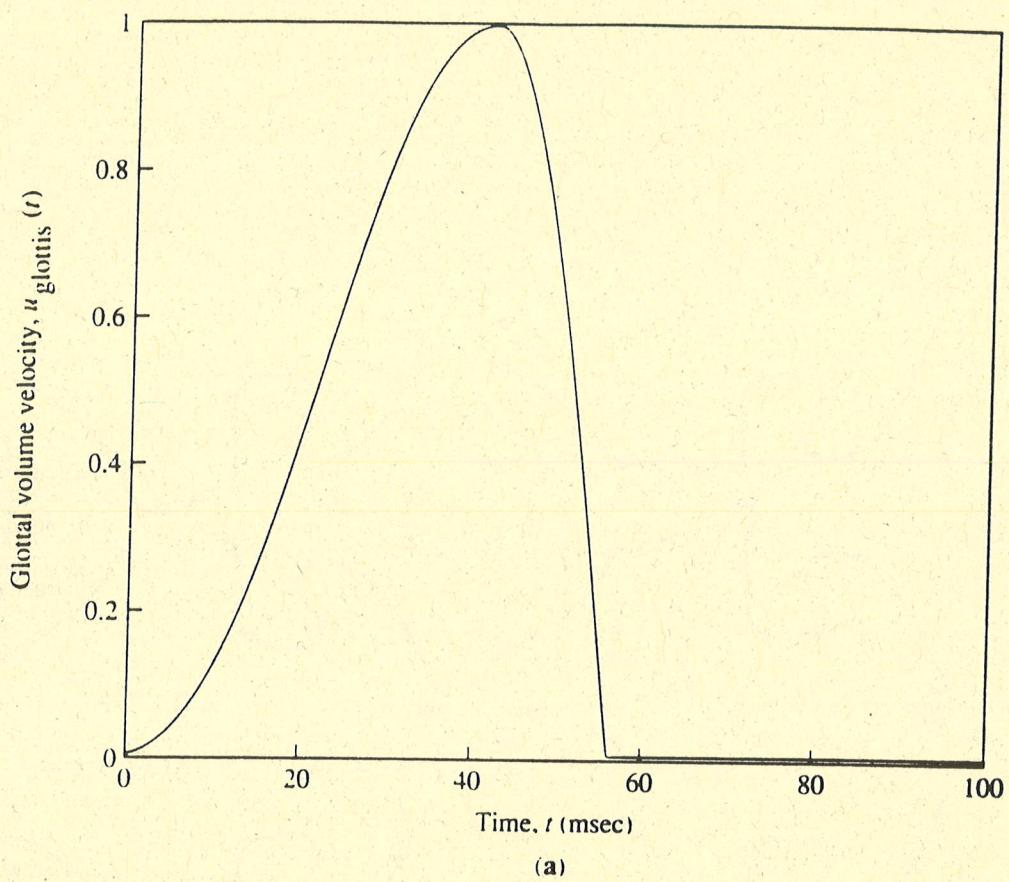
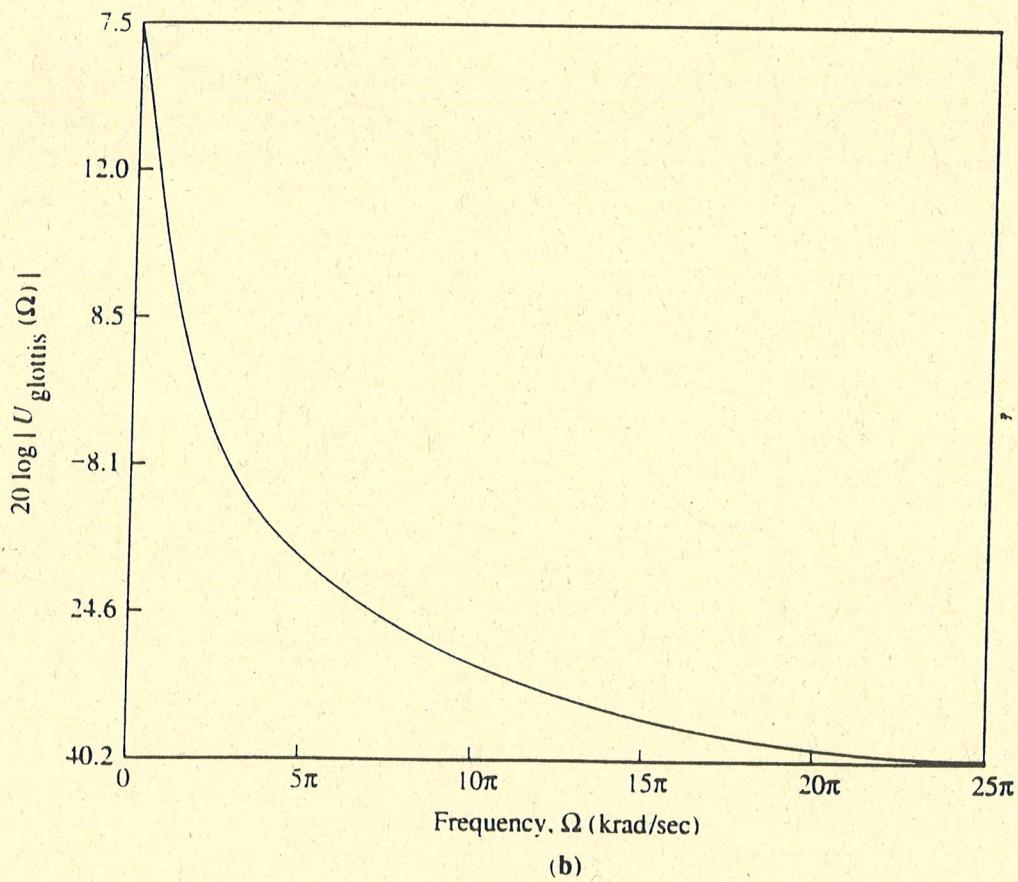


FIGURE 2.7. A sequence of cross sections of the larynx illustrating a complete phonation cycle. After Vennard (1967).

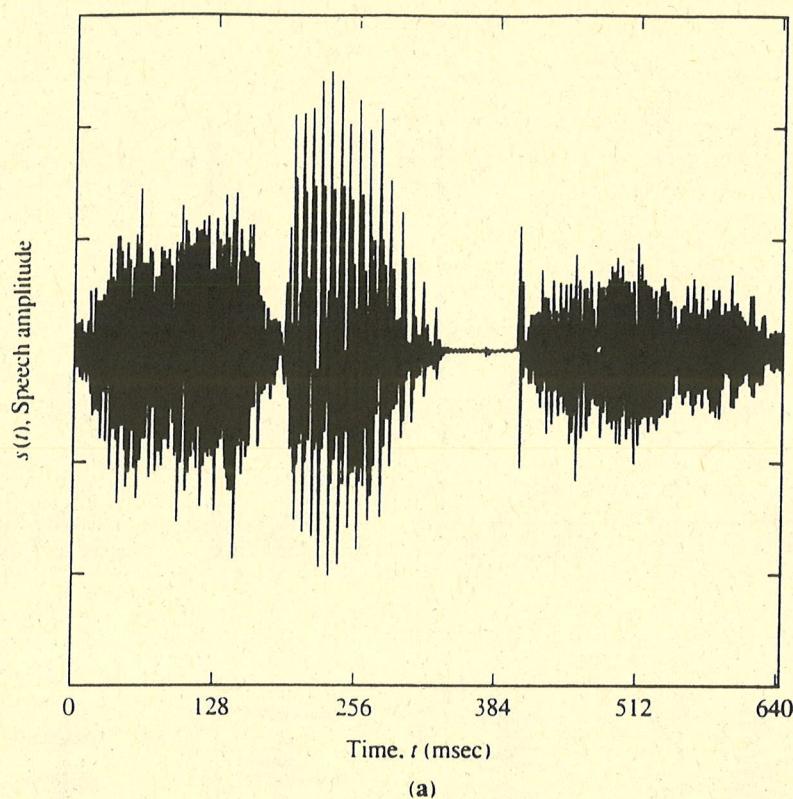


(a)

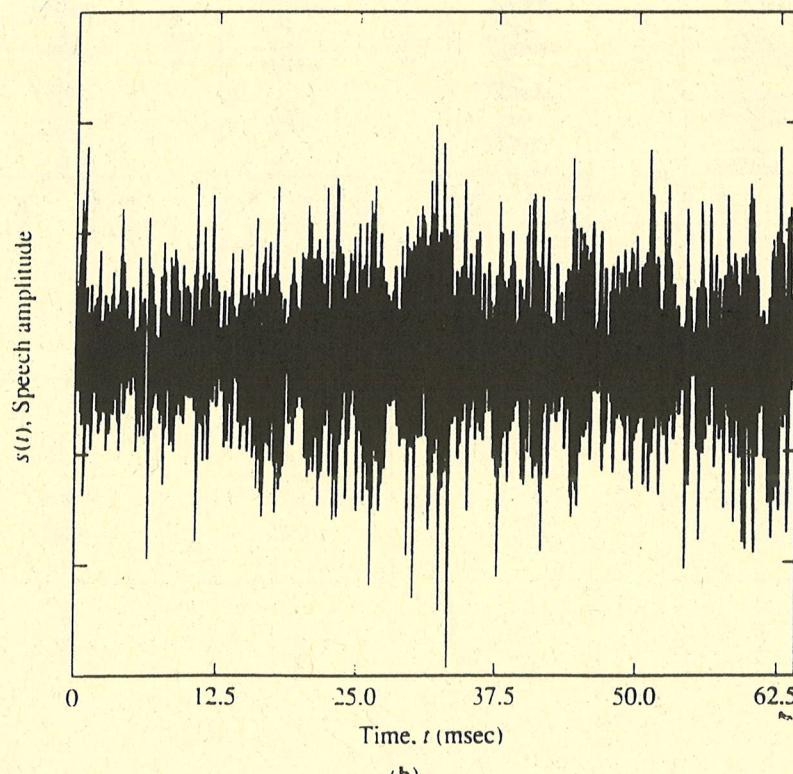


(b)

FIGURE 2.8. (a) Time waveform of volume velocity of the glottal source excitation. (b) Magnitude spectrum of one pulse of the volume velocity at the glottis.

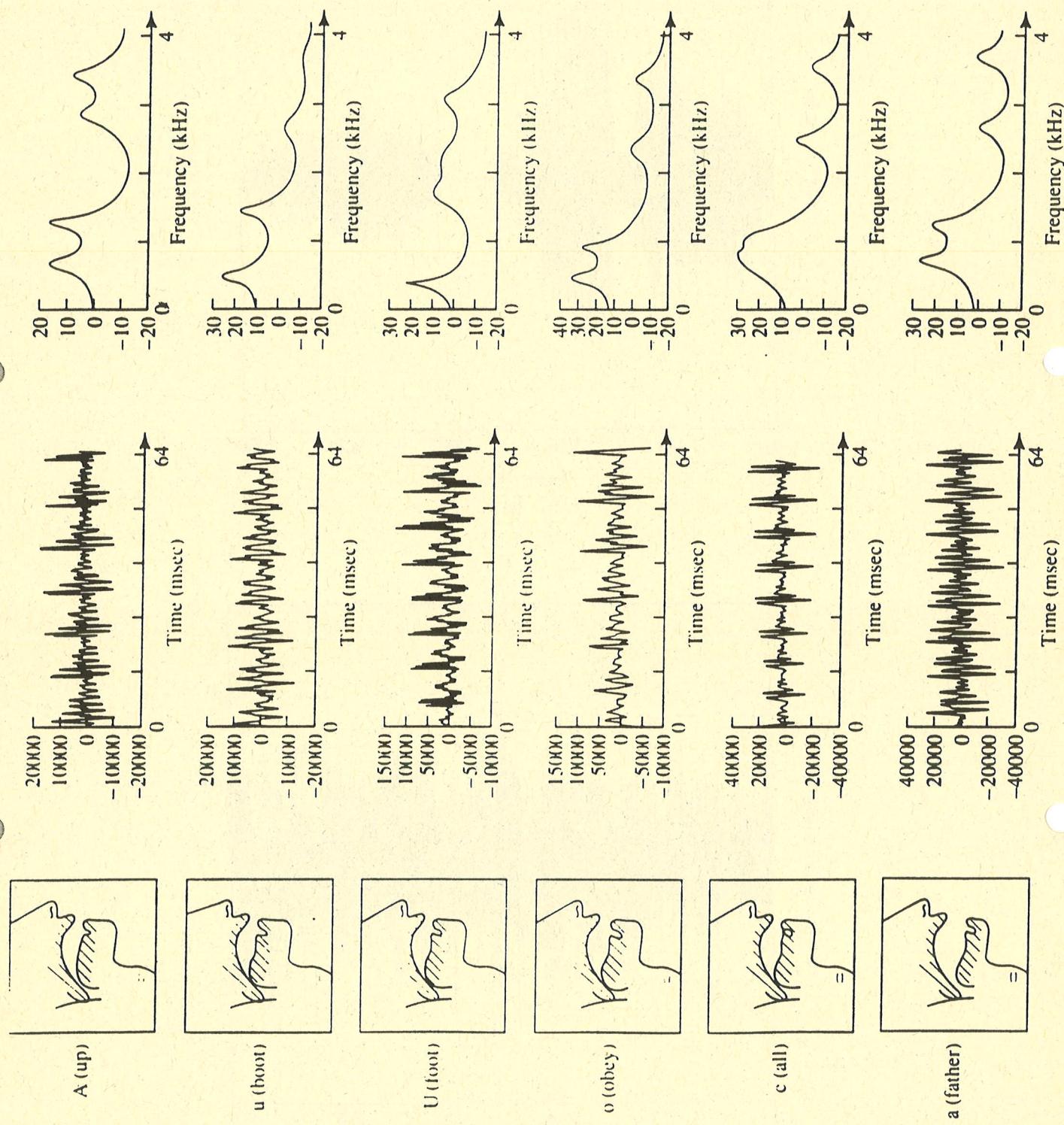


(a)



(b)

FIGURE 2.4. (a) A speech signal of the word "six" spoken by a male speaker; (b) blowup of a frame of the steady-state region of the initial /s/ sound; (c) blowup of the vowel /ɪ/.



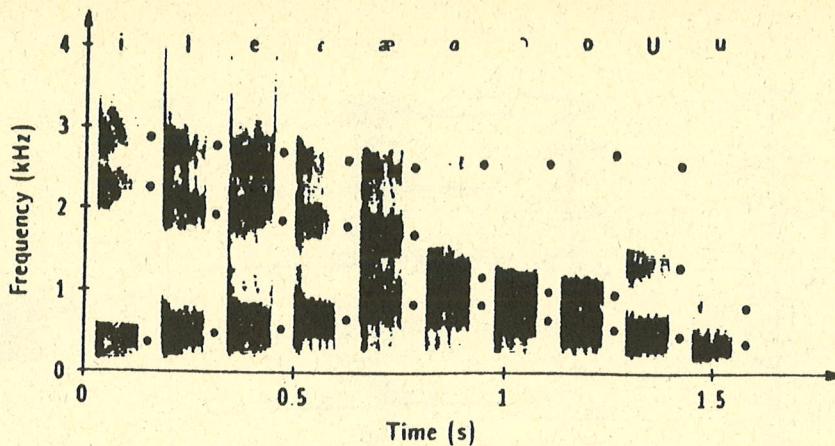


Fig. 3.10 Spectrogram of short sections of English vowels from a male speaker. Formants for each vowel are noted by dots.

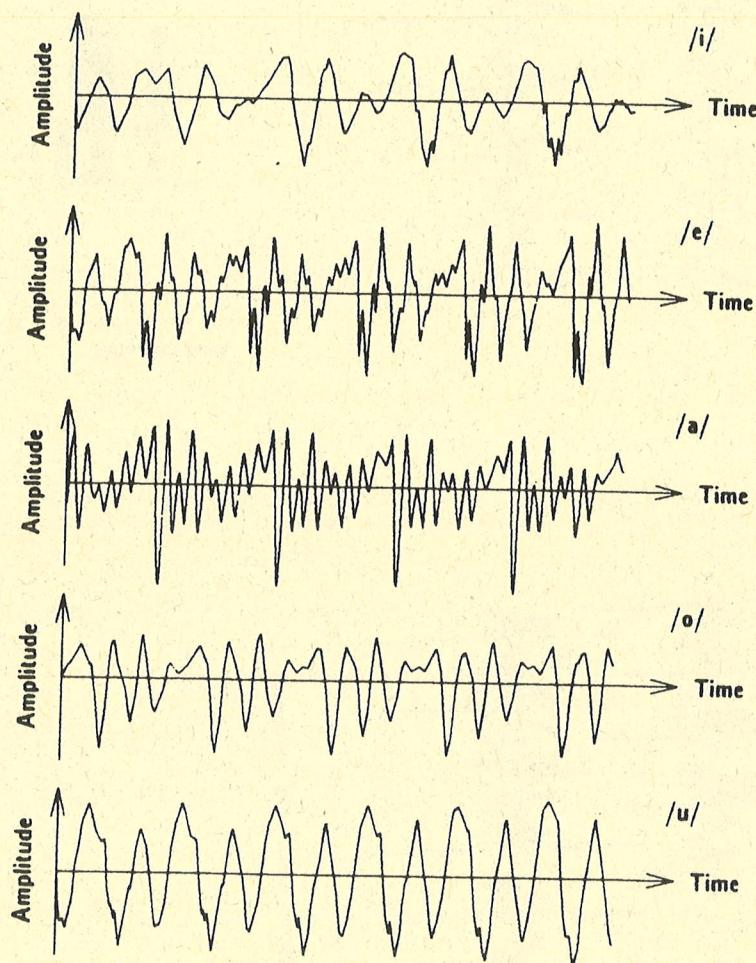


Fig. 3.12 Typical acoustic waveforms for five English vowels. Each plot shows 40 ms of a different vowel, which comprises about 5–6 pitch periods for this speaker. Note the quasi-periodic nature of such voiced speech as well as the varying spectral content for different vowels.

		/i/	/ɪ/	/ɛ/	/æ/	/ɑ/	/ɔ/	/ʊ/	/ə/	/ʌ/	/ɒ/
F1	male	270	390	530	660	730	570	440	300	640	490
	female	310	430	610	860	850	590	470	370	760	500
F2	male	2290	1990	1840	1720	1090	840	1020	870	1190	1350
	female	2790	2480	2330	2050	1220	920	1160	950	1400	1640
F3	male	3010	2550	2480	2410	2440	2410	2240	2240	2390	1690
	female	3310	3070	2990	2850	2810	2710	2680	2670	2780	1960

Table 3.2 Average formant frequencies (in Hz) for English vowels by adult male and female speakers. (After Peterson and Barney [17].)

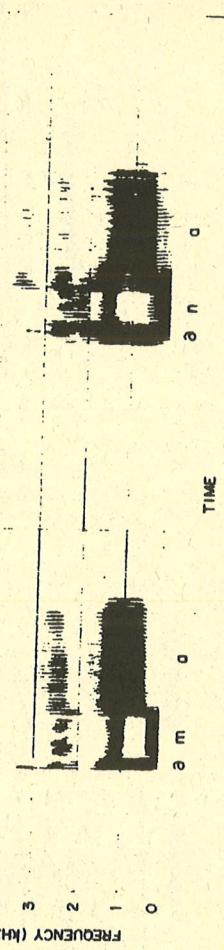


Fig. 3.8 Acoustic waveforms and spectrograms for utterances /UH-M-A/ and /UH-N-A/.

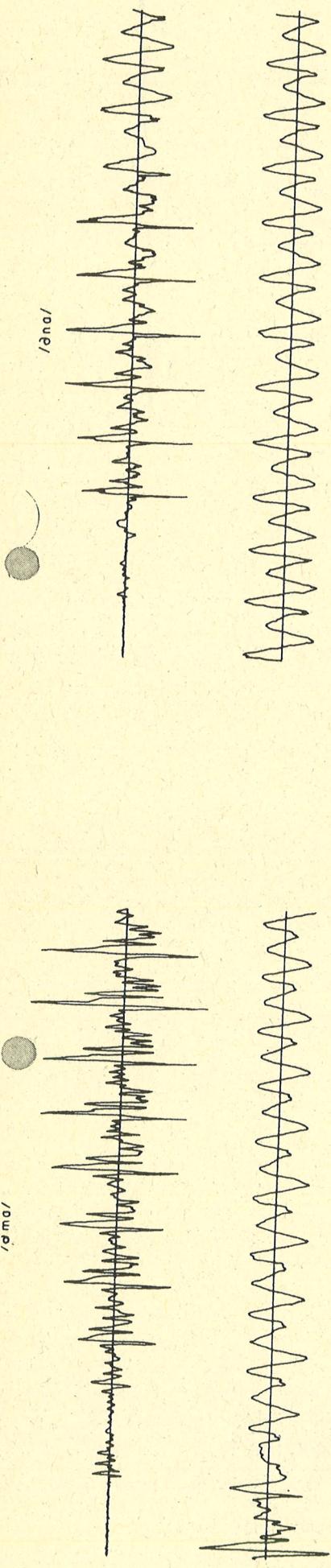
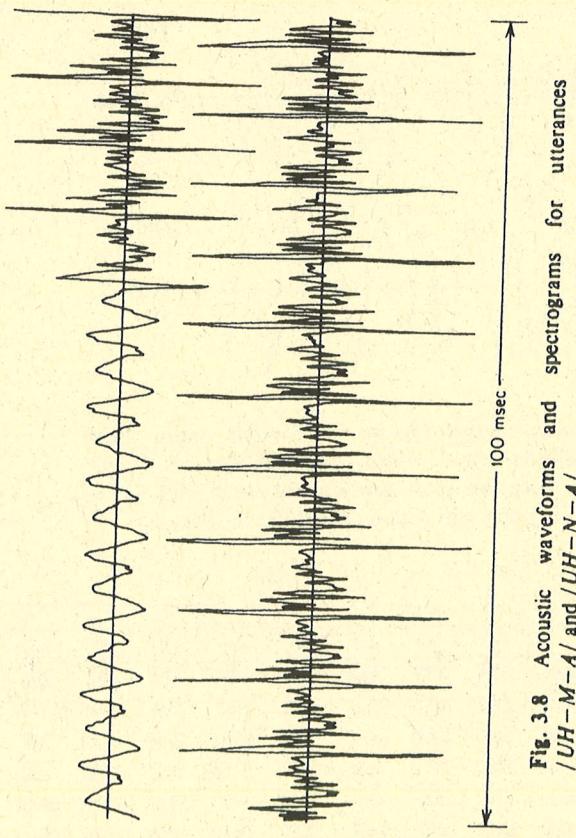


Fig. 3.8 (Continued)

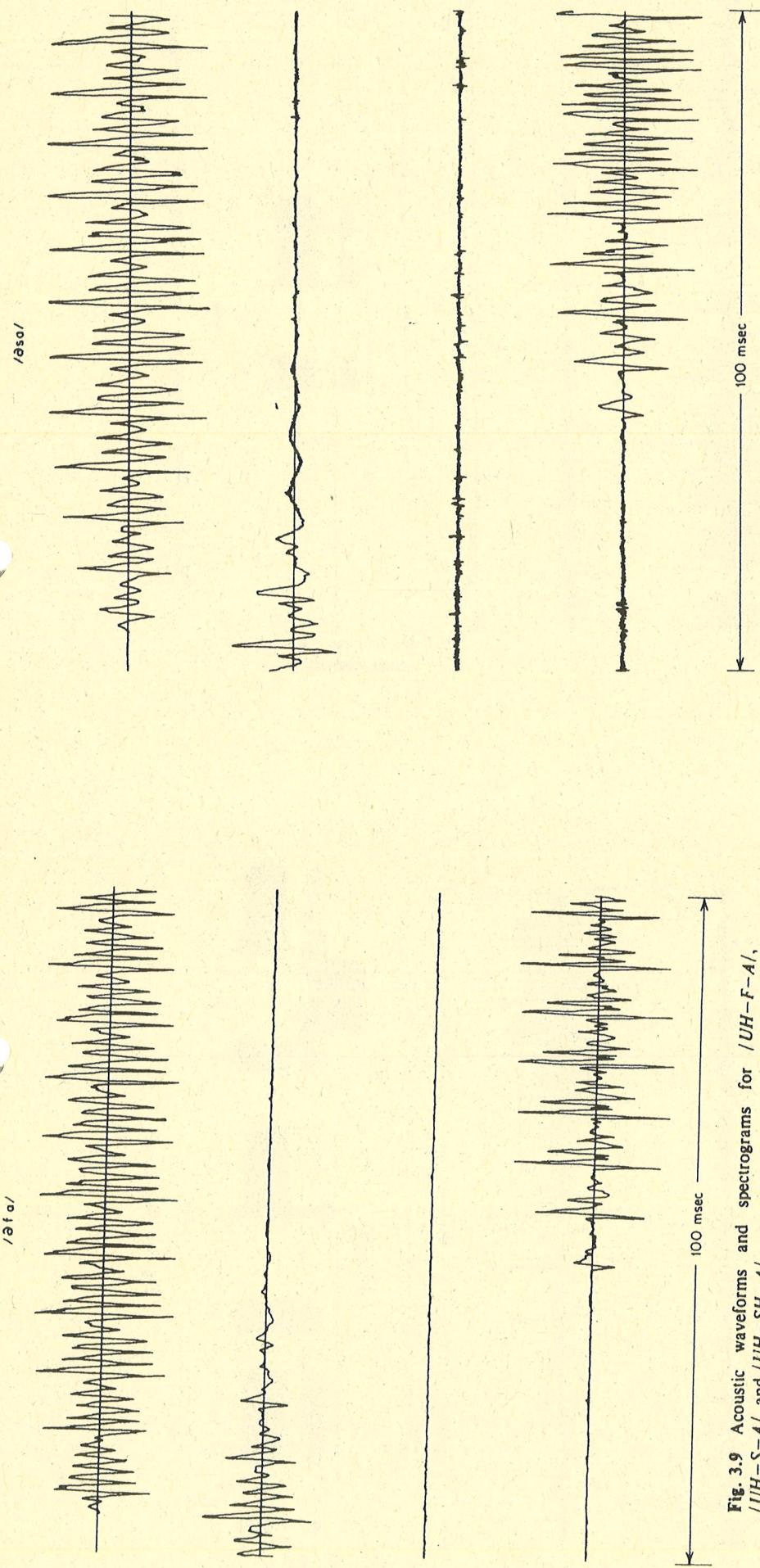


Fig. 3.9 Acoustic waveforms and spectrograms for /əsəv/,
/ətəv/, /UH-S-A/, and /UH-SH-A/.

Fig. 3.9 (Continued)

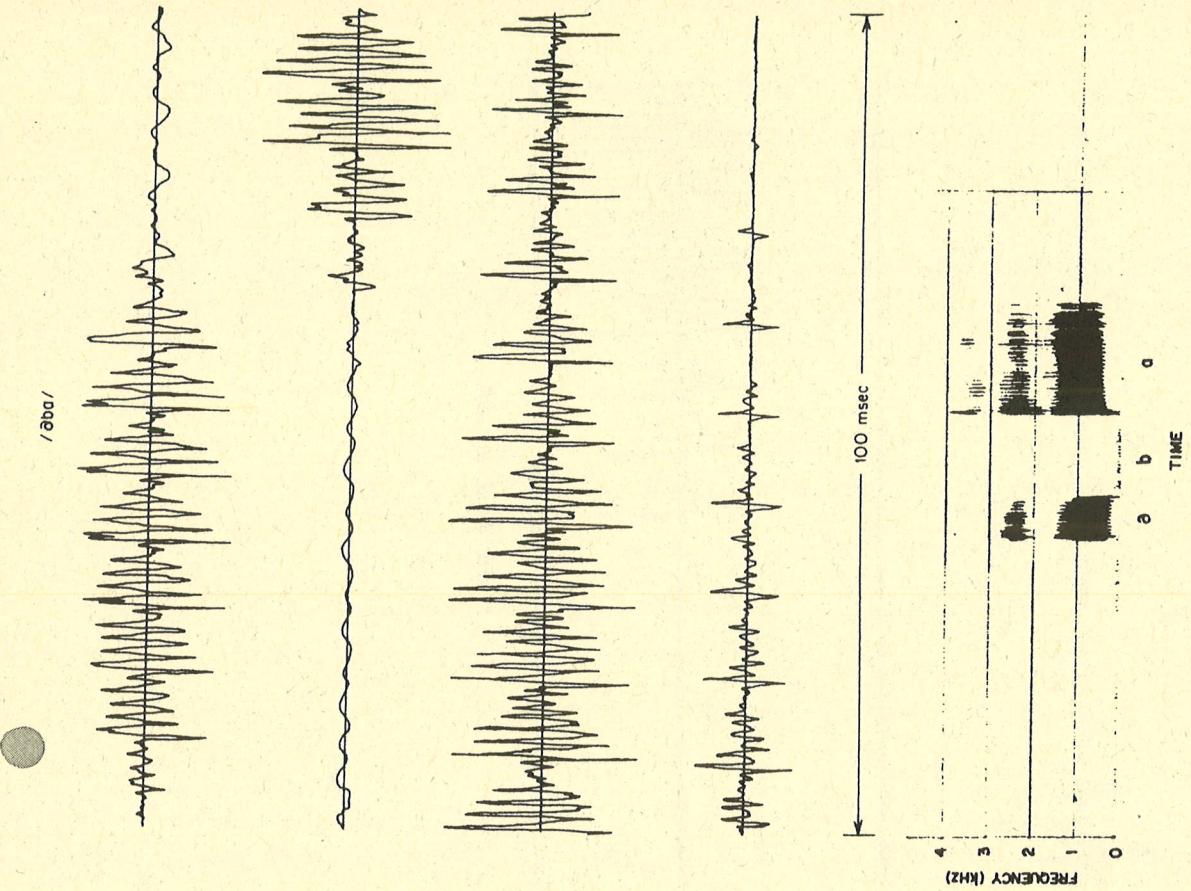


Fig. 3.11 Acoustic waveform and spectrogram for utterance /əʃəʊ-B-A/.

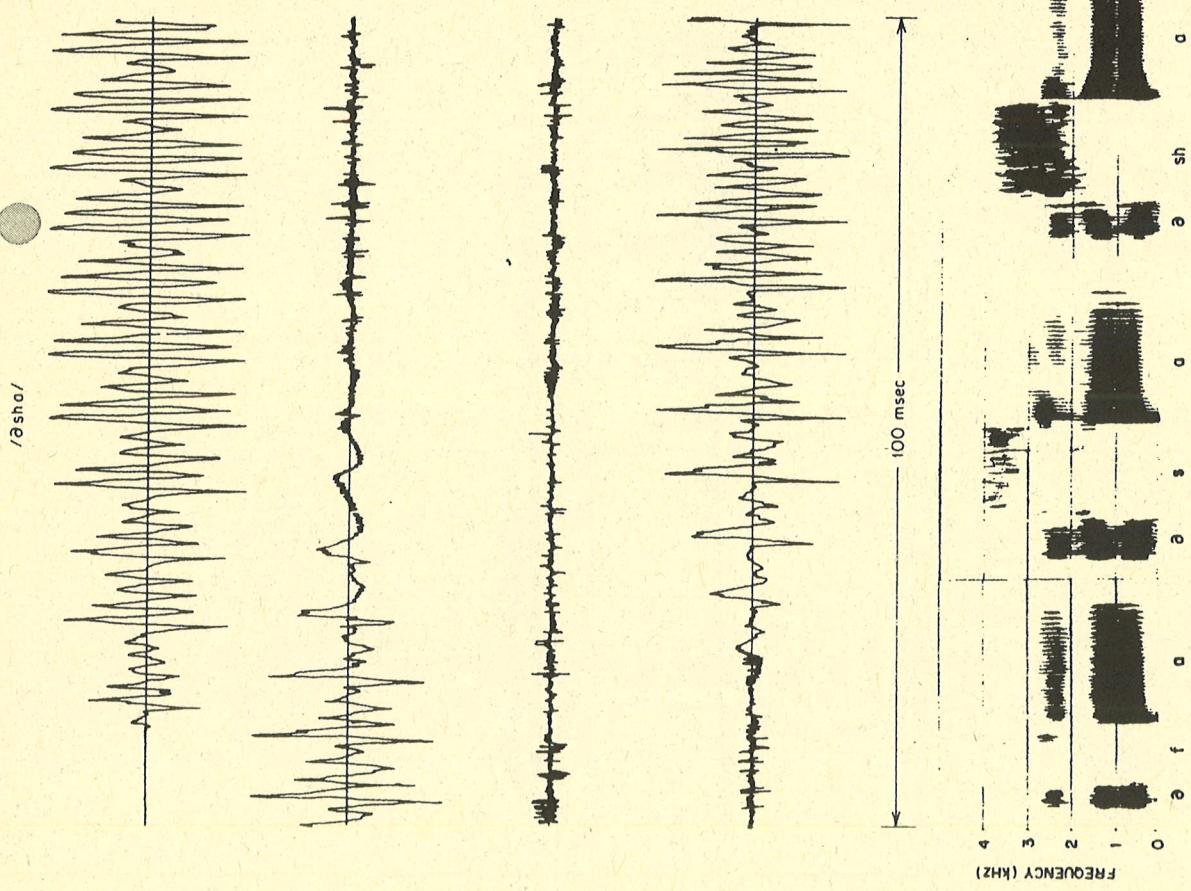
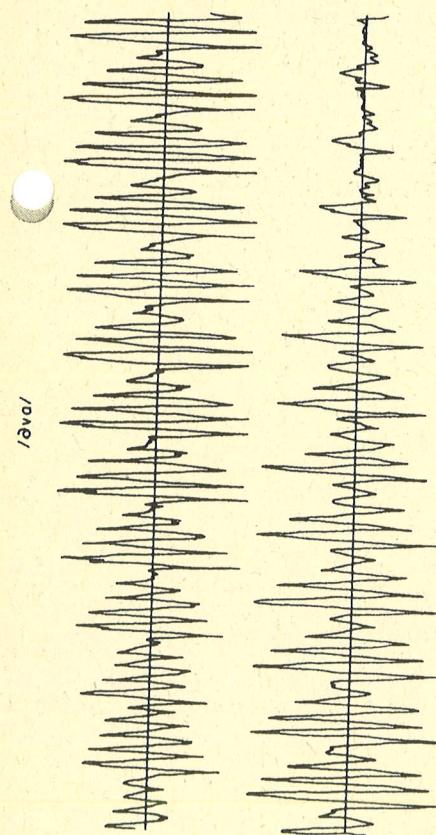


Fig. 3.9 (Continued)



/əzəv/

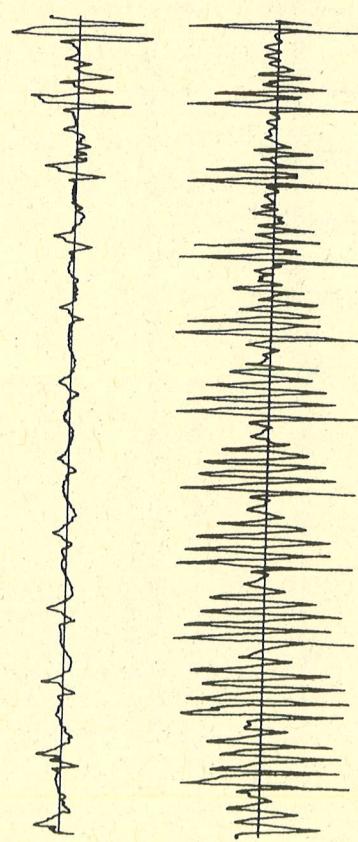
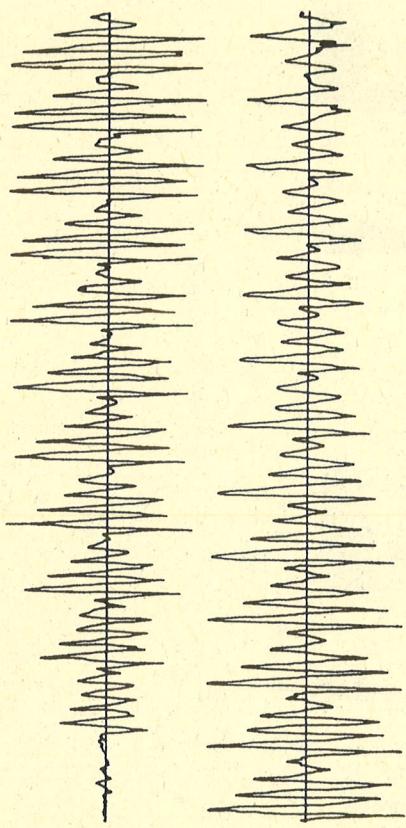


Fig. 3.10 Acoustic waveforms and spectrograms for utterances /UH-V-A/ and /UH-ZH-A/.

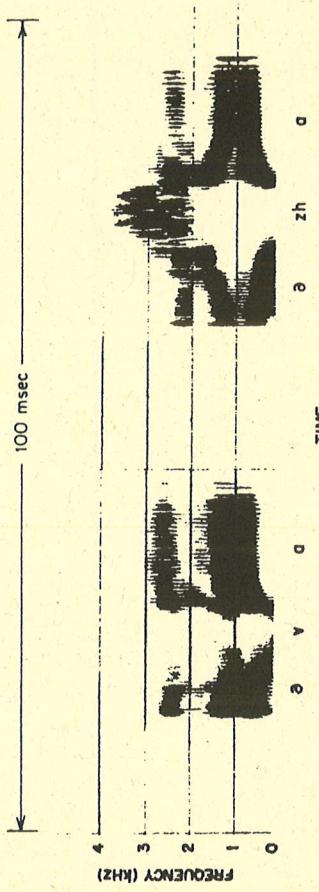


Fig. 3.10 (Continued)

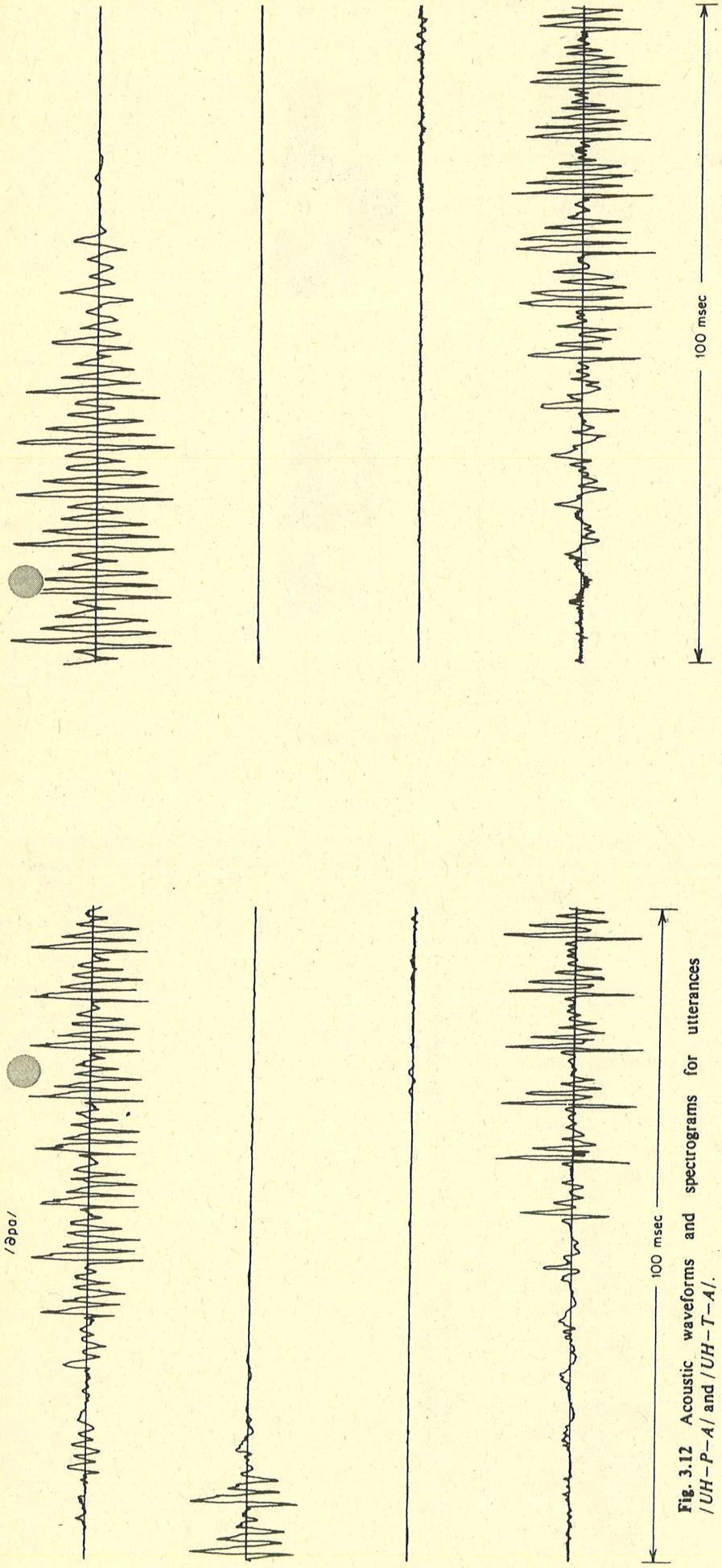


Fig. 3.12 Acoustic waveforms and spectrograms for utterances /UH-P-A/ and /UH-T-A/.

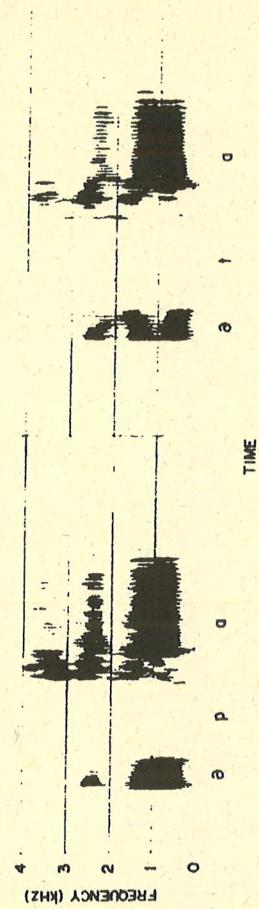


Fig. 3.12 (Continued)

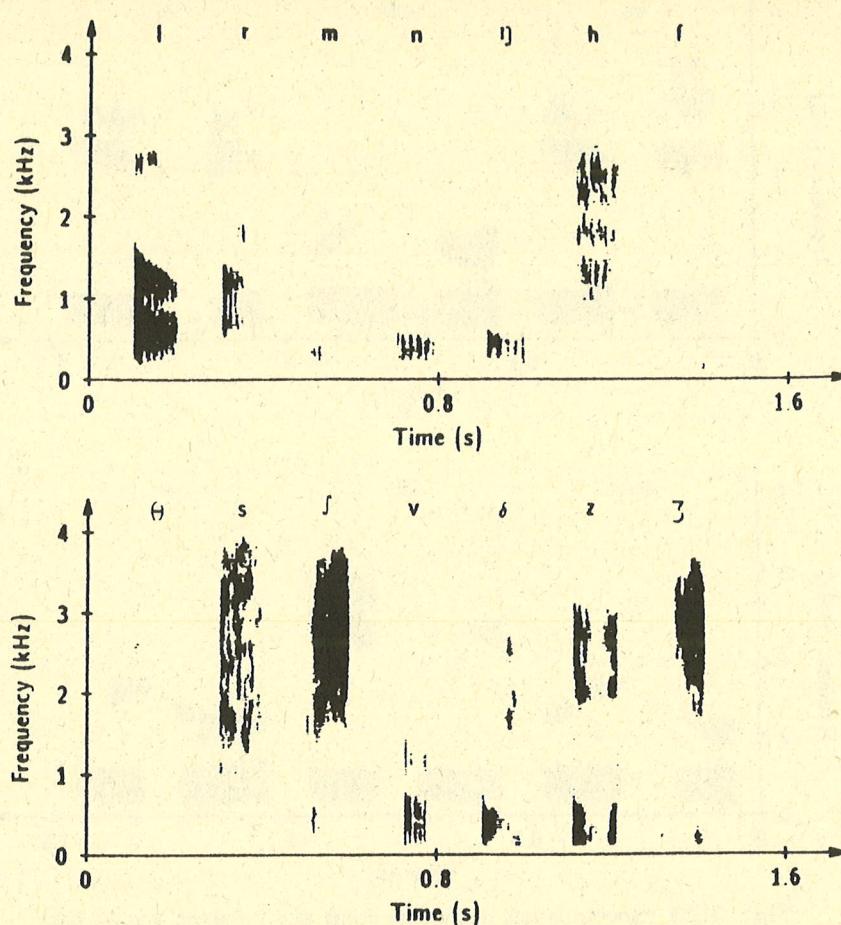


Fig. 3.18 Spectrograms of 14 English steady-state consonants /l,r,m,n,ŋ,h,f,
 θ ,s,ʃ,v,ð,z,ʒ/.

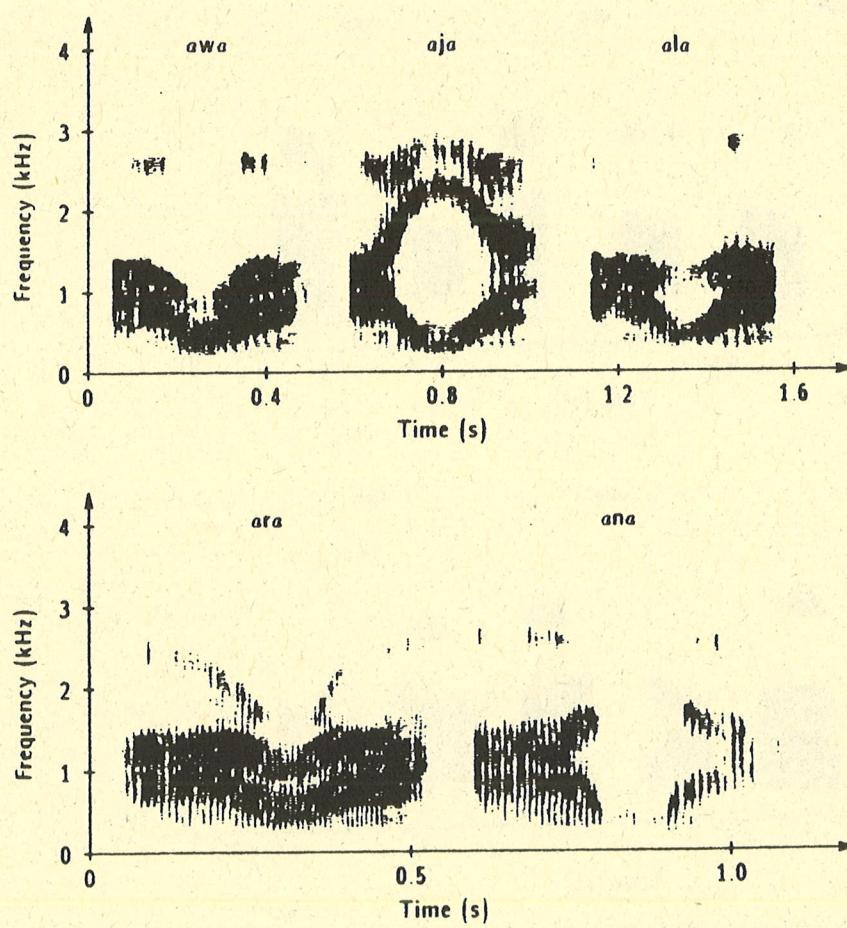


Fig. 3.19 Spectrograms of English sonorants in vocalic context:
/əwə, əjə, ələ, ərə, ənə/.

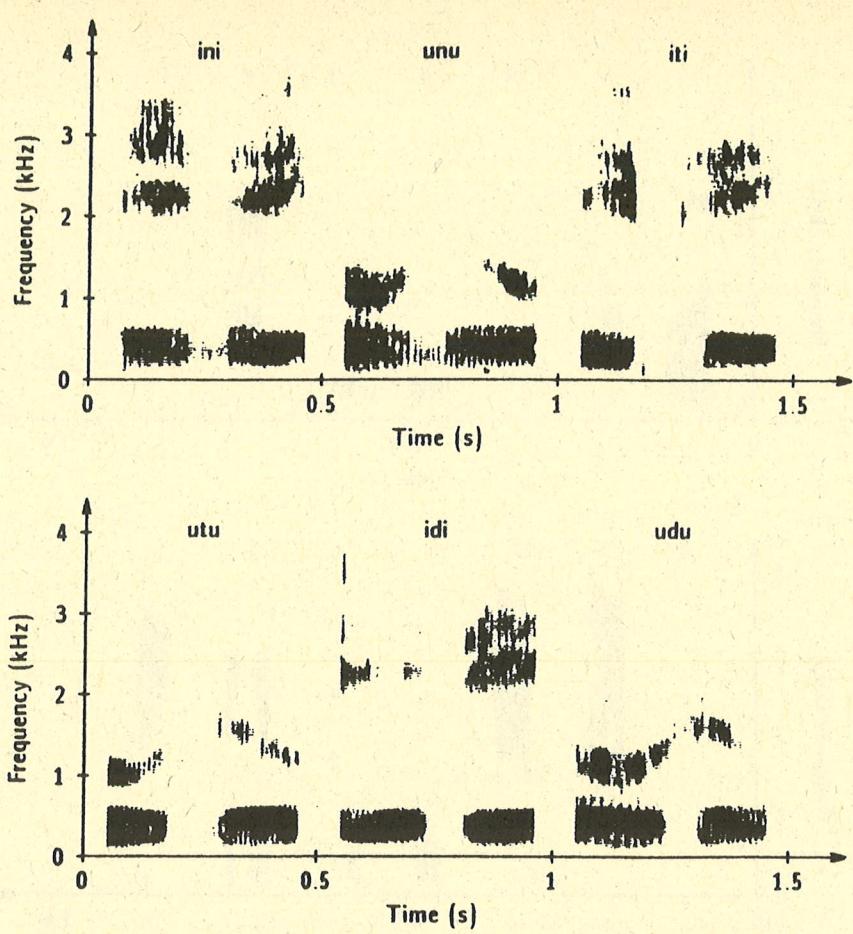


Fig. 3.20 Spectrograms of English stops and nasals in vocalic context: /ini.unu,iti,utu,idi,udu/.

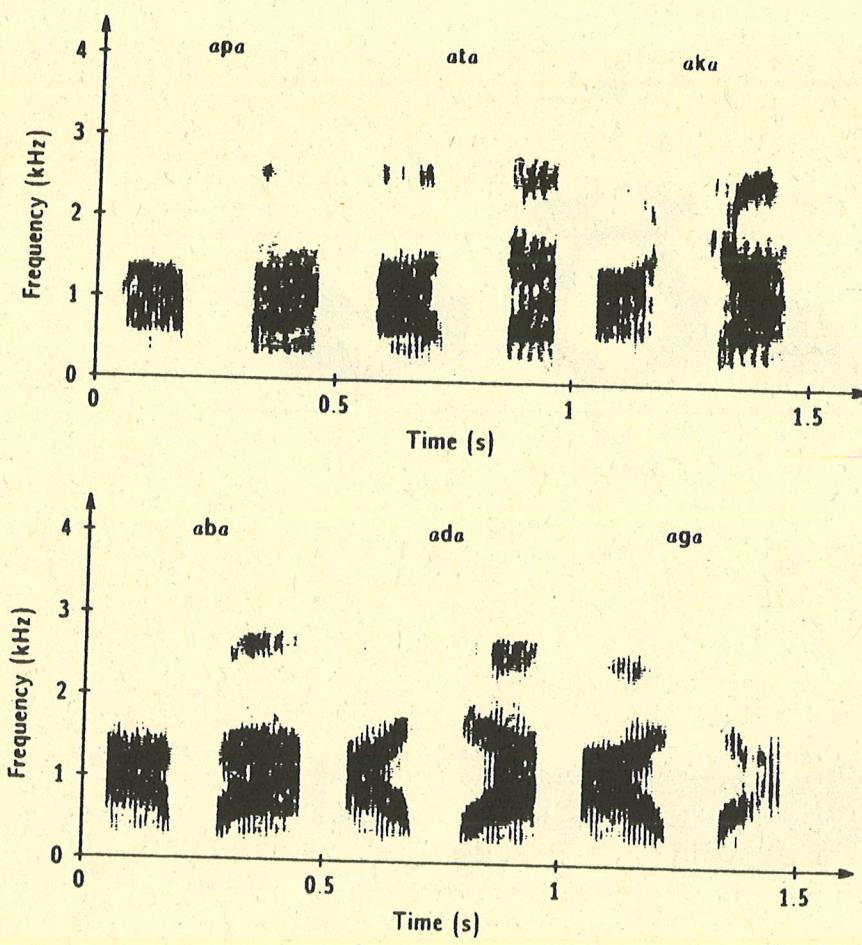


Fig. 3.21 Spectrograms of the six English stops in vocalic context: /apə,atə,akə,abə,adə,agə/.

EVERY SALT BREEZE COME FROM THE SEA - AER

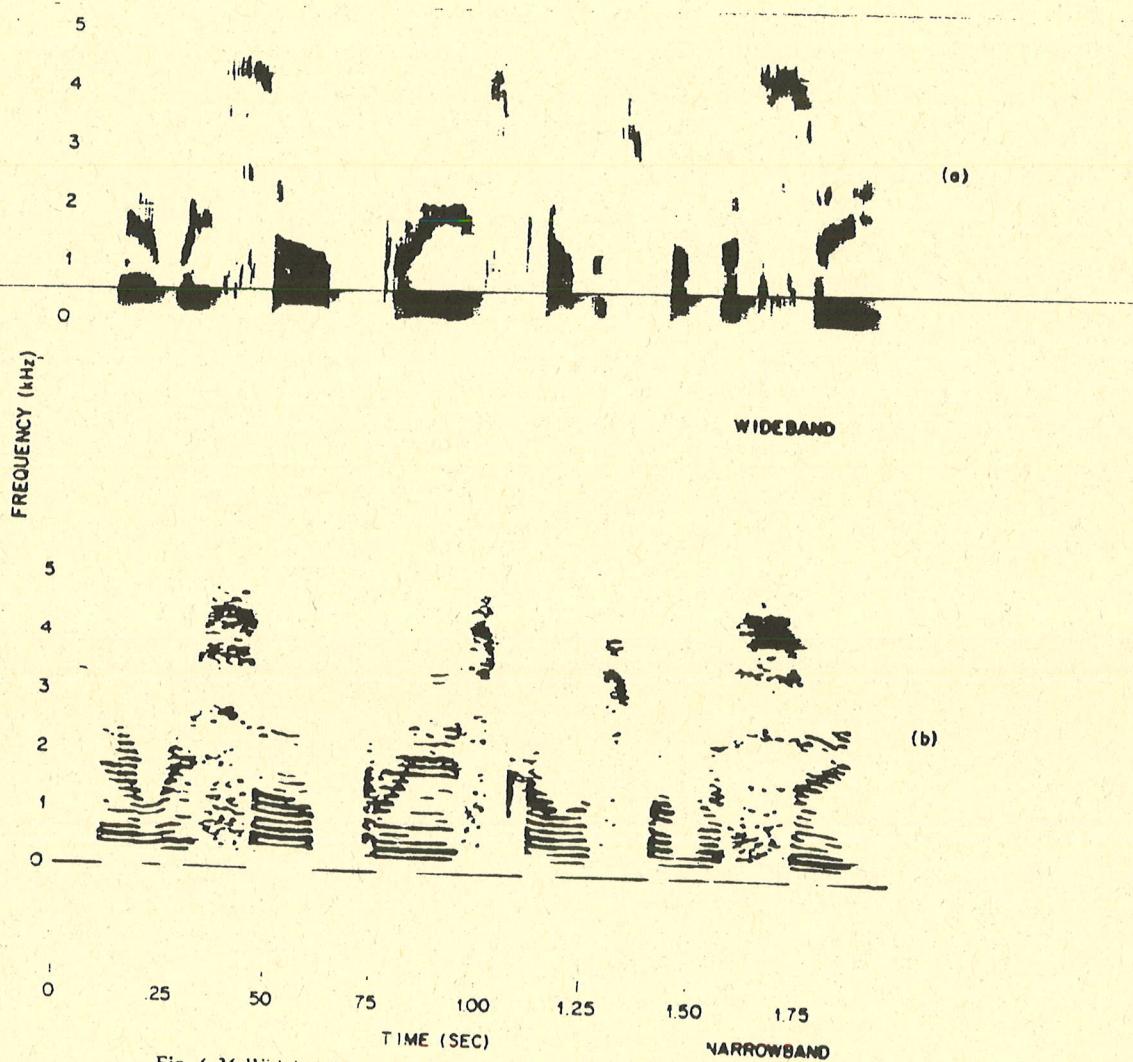


Fig. 6.36 Wideband and narrowband spectrograms of a sentence

Tubo acústico de sección constante.

Ecuaciones

$$-\frac{\partial u}{\partial x} = \frac{A}{\rho \cdot c^2} \frac{\partial p}{\partial t} \quad (\text{Conservación de la masa})$$

$$-\frac{\partial p}{\partial x} = \frac{\rho}{A} \frac{\partial u}{\partial t} \quad (\text{Conservación del momento})$$

$$-\frac{\partial i}{\partial x} = C \frac{\partial v}{\partial t}$$

$$-\frac{\partial v}{\partial x} = L \frac{\partial i}{\partial t}$$

Analogía

$u \leftrightarrow i$	$\frac{A}{\rho c^2} \leftrightarrow C$	Compresibilidad
$p \leftrightarrow v$	$\frac{\rho}{A} \leftrightarrow L$	Inductancia acústica

Solución (Superposición onda directa y reflejada)

$$u(x,t) = u^+ \left(t - \frac{x}{c} \right) - u^- \left(t + \frac{x}{c} \right)$$

$$p(x,t) = Z_0 \left[u^+ \left(t - \frac{x}{c} \right) + u^- \left(t + \frac{x}{c} \right) \right]$$

donde Z_0 es la impedancia acústica característica (del medio):

$$Z_0 = \frac{\rho \cdot c}{A}$$

Análisis armónico

$$u(0,t) = \exp(j2\pi ft) \quad (\text{Excitación})$$

$$u(x,t) = U(x) \cdot \exp(j2\pi ft)$$

$$p(x,t) = P(x) \cdot \exp(j2\pi ft)$$

$$p(l,t) = 0 \Rightarrow P(l) = 0 \quad (\text{Final del tubo abierto})$$

Ecuaciones

$$\left. \begin{aligned} -\frac{\partial U(x)}{\partial x} &= \frac{A}{\rho \cdot c^2} j \cdot 2\pi f \cdot P(x) \\ -\frac{\partial P(x)}{\partial x} &= \frac{\rho}{A} j \cdot 2\pi f \cdot U(x) \end{aligned} \right\} \Rightarrow \begin{aligned} \frac{\partial^2 U}{\partial x^2} + \frac{4\pi^2 f^2}{c^2} U &= 0 \\ \frac{\partial^2 P}{\partial x^2} + \frac{4\pi^2 f^2}{c^2} P &= 0 \end{aligned}$$

Soluciones

$$\begin{aligned} U(x) &= a_1 \exp[\gamma \cdot x] + a_2 \exp[-\gamma \cdot x] \\ P(x) &= a_3 \exp[\gamma \cdot x] + a_4 \exp[-\gamma \cdot x] \end{aligned} \quad \gamma = j \cdot \frac{2\pi f}{c}$$

$$\text{Condiciones de contorno: } U(0) = 1; \quad P(l) = 0$$

$$U(x) = \frac{\cos\left(\frac{2\pi f}{c}(l-x)\right)}{\cos\left(\frac{2\pi f}{c}l\right)}$$

$$P(x) = j \cdot Z_0 \frac{\sin\left(\frac{2\pi f}{c}(l-x)\right)}{\cos\left(\frac{2\pi f}{c}l\right)}$$

Modelo de tubos concatenados.**Ecuaciones en la sección iésima**

$$u_i(x,t) = u_i^+ \left(t - \frac{x}{c} \right) - u_i^- \left(t + \frac{x}{c} \right)$$

$$p_i(x,t) = Z_i \left[u_i^+ \left(t - \frac{x}{c} \right) + u_i^- \left(t + \frac{x}{c} \right) \right] ; \quad Z_i = \frac{\rho c}{A_i}$$

Condiciones de continuidad en la transición i / (i+1)

$$u_i(l_i, t) = u_{i+1}(0, t) \Rightarrow u_i^+(t - \tau_i) - u_i^-(t + \tau_i) = u_{i+1}^+(t) - u_{i+1}^-(t)$$

$$p_i(l_i, t) = p_{i+1}(0, t) \Rightarrow Z_i \left[u_i^+(t - \tau_i) + u_i^-(t + \tau_i) \right] = Z_{i+1} \left[u_{i+1}^+(t) + u_{i+1}^-(t) \right]$$

donde: $\tau_i = \frac{l_i}{c}$

Y expresando las ondas salientes en función de las entrantes:

$$u_{i+1}^+(t) = (1 + r_i)u_i^+(t - \tau_i) + r_i u_{i+1}^-(t)$$

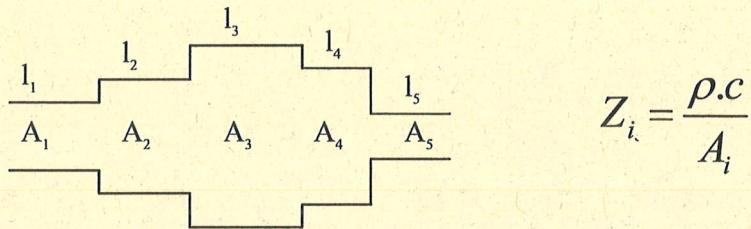
$$u_i^-(t + \tau_i) = -r_i u_i^+(t - \tau_i) + (1 - r_i) u_{i+1}^-(t)$$

O en forma vectorial:

$$\begin{bmatrix} u_{i+1}^+(t) \\ u_i^-(t + \tau_i) \end{bmatrix} = \begin{bmatrix} 1 + r_i & r_i \\ -r_i & 1 - r_i \end{bmatrix} \begin{bmatrix} u_i^+(t - \tau_i) \\ u_{i+1}^-(t) \end{bmatrix}$$

donde: $r_i = \frac{Z_i - Z_{i+1}}{Z_i + Z_{i+1}} = \frac{A_{i+1} - A_i}{A_{i+1} + A_i}$

Modelo de tubos concatenados.



Ondas directas y reflejadas en la transición de los tubos i / i+1

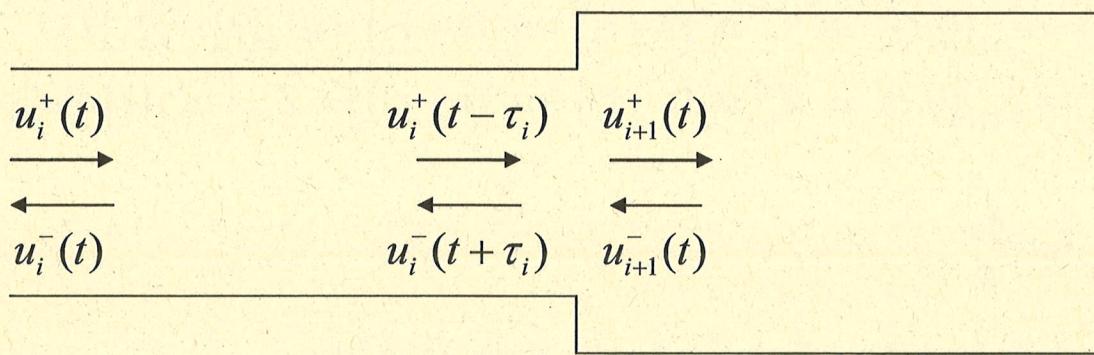
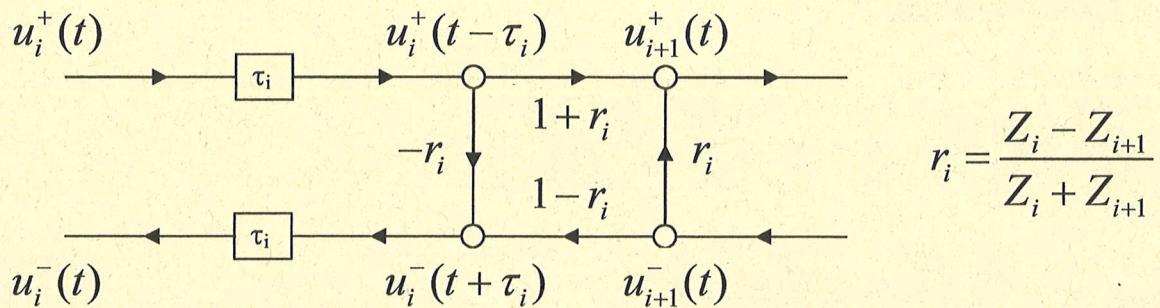


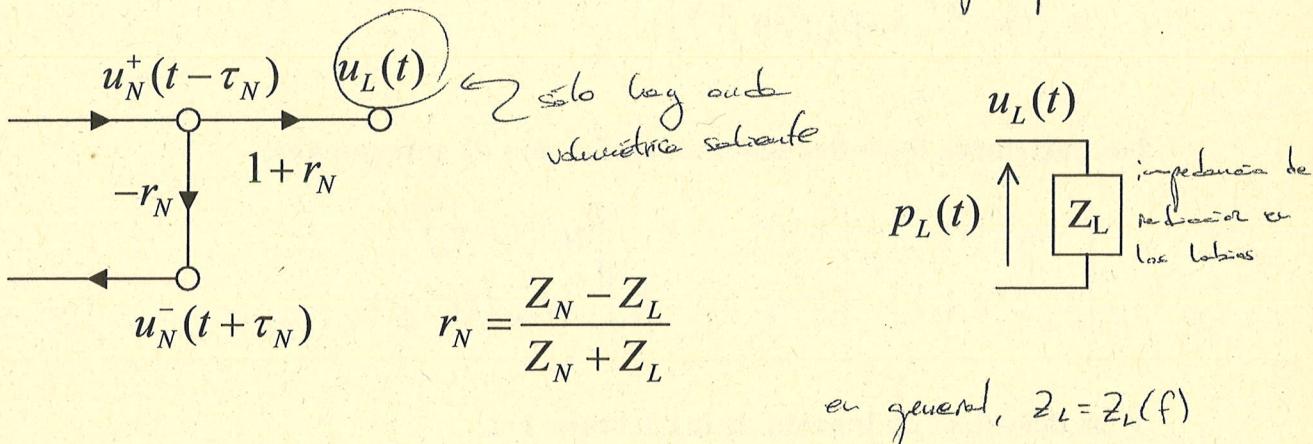
Diagrama de flujo en la transición de los tubos i / i+1



Condiciones de contorno: Radiación en los Labios.

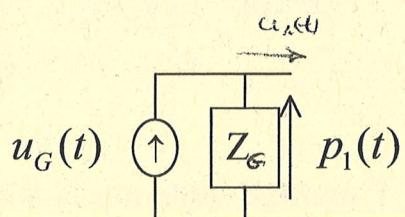
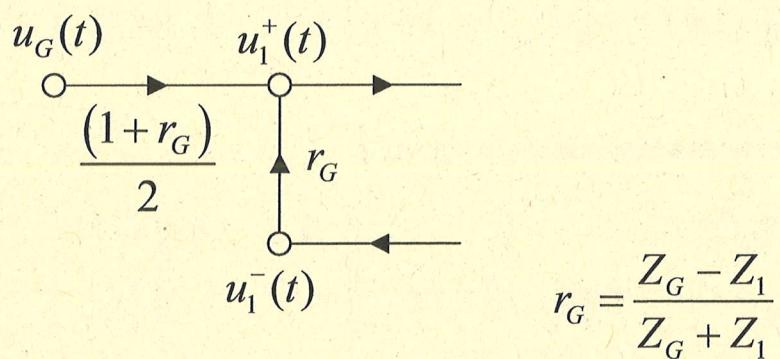
- Se asume que el medio exterior no refleja nada.
- La radiación se modela como una impedancia Z_L

*✓ todo la energía que se emite va
fuera y no vuelve → sistema
exterior muy amplio*



Condiciones de contorno: Glotis.

- La generación del sonido se modela como un generador de corriente con impedancia Z_G en paralelo.



Respuesta en frecuencia**Análisis armónico:**

$$\left. \begin{array}{l} u_G(t) = U_G \exp(j2\pi ft) \\ u_L(t) = U_L \exp(j2\pi ft) \end{array} \right\} \Rightarrow H(f) = \frac{U_G}{U_L}$$

Las relaciones de ondas salientes en función de entrantes es:

$$\begin{bmatrix} u_{i+1}^+(t) \\ u_i^-(t + \tau_i) \end{bmatrix} = \begin{bmatrix} 1+r_i & r_i \\ -r_i & 1-r_i \end{bmatrix} \begin{bmatrix} u_i^+(t - \tau_i) \\ u_{i+1}^-(t) \end{bmatrix}$$

Y las del tubo i en función de las del tubo i+1:

$$\begin{bmatrix} u_i^+(t - \tau_i) \\ u_i^-(t + \tau_i) \end{bmatrix} = \frac{1}{1+r_i} \begin{bmatrix} 1 & -r_i \\ -r_i & 1 \end{bmatrix} \begin{bmatrix} u_{i+1}^+(t) \\ u_{i+1}^-(t) \end{bmatrix}$$

Que para señales armónicas resulta:

reflexos para el

$$\begin{bmatrix} \exp(-j2\pi f\tau_i) & 0 \\ 0 & \exp(j2\pi f\tau_i) \end{bmatrix} \begin{bmatrix} U_i^+ \\ U_i^- \end{bmatrix} = \frac{1}{1+r_i} \begin{bmatrix} 1 & -r_i \\ -r_i & 1 \end{bmatrix} \begin{bmatrix} U_{i+1}^+ \\ U_{i+1}^- \end{bmatrix}$$

Y multiplicando por la matriz de exponentiales inversa:

$$\begin{bmatrix} U_i^+ \\ U_i^- \end{bmatrix} = \frac{1}{1+r_i} \begin{bmatrix} \exp(j2\pi f\tau_i) & 0 \\ 0 & \exp(-j2\pi f\tau_i) \end{bmatrix} \begin{bmatrix} 1 & -r_i \\ -r_i & 1 \end{bmatrix} \begin{bmatrix} U_{i+1}^+ \\ U_{i+1}^- \end{bmatrix}$$

Para la Glotis y los Labios tenemos las relaciones:

$$U_G = \frac{2}{1+r_G} \begin{bmatrix} 1 & -r_G \end{bmatrix} \begin{bmatrix} U_1^+ \\ U_1^- \end{bmatrix}$$

desarrolla con el resto

$$\begin{bmatrix} U_N^+ \\ U_N^- \end{bmatrix} = \frac{1}{1+r_N} \begin{bmatrix} \exp(j2\pi f \tau_N) & 0 \\ 0 & \exp(-j2\pi f \tau_N) \end{bmatrix} \begin{bmatrix} 1 \\ -r_N \end{bmatrix} U_L$$

Y concatenando todas las relaciones obtenemos:

$$U_G = \overbrace{\frac{2}{1+r_G} \begin{bmatrix} 1 & -r_G \end{bmatrix}}^{\Gamma_i} \prod_{i=1}^{N-1} \left\{ \overbrace{\frac{1}{1+r_i} \cdot \underline{\underline{\Gamma_i \cdot R_i}}}^{\underline{\underline{\Gamma_i \cdot R_i}}} \right\} \overbrace{\frac{1}{1+r_N} \begin{bmatrix} 1 \\ -r_N \end{bmatrix}}^{\Gamma_N} U_L$$

Siendo:

$$\underline{\underline{\Gamma_i}} = \begin{bmatrix} \exp(j2\pi f \tau_i) & 0 \\ 0 & \exp(-j2\pi f \tau_i) \end{bmatrix}; \quad \underline{\underline{R_i}} = \begin{bmatrix} 1 & -r_i \\ -r_i & 1 \end{bmatrix}$$

La expresión de la respuesta en frecuencia resulta ser:

$$H(f) = \frac{U_L}{U_G} = \frac{(1+r_G) \prod_{i=1}^N \{1+r_i\}}{2 \begin{bmatrix} 1 & -r_G \end{bmatrix} \prod_{i=1}^{N-1} \{\underline{\underline{\Gamma_i \cdot R_i}}\} \cdot \underline{\underline{\Gamma_N}} \begin{bmatrix} 1 \\ -r_N \end{bmatrix}}$$

Respuesta en frecuencia (Modelo de tiempo discreto)**Haciendo el cambio:**

$$\exp(j2\pi f\tau) = z^{\frac{1}{2}}$$

en la expresión de la respuesta en frecuencia de los tubos concatenados obtenemos la función del sistema del filtro discreto:

$$\left\{ H(z) = \frac{(1 + r_G) \prod_{i=1}^N \{1 + r_i\} \cdot z^{-\frac{N}{2}}}{2[1 - r_G] \prod_{i=1}^{N-1} \{\underline{\mathbf{Q}_i}\} \begin{bmatrix} 1 \\ -r_N \cdot z^{-1} \end{bmatrix}} \right\}$$

donde la matriz $\underline{\mathbf{Q}_i}$ es:

$$\underline{\mathbf{Q}_i} = z^{-\frac{1}{2}} \cdot \underline{\Gamma} \cdot \underline{\mathbf{R}_i} = z^{-\frac{1}{2}} \cdot \begin{bmatrix} z^{\frac{1}{2}} & 0 \\ 0 & z^{-\frac{1}{2}} \end{bmatrix} \cdot \begin{bmatrix} 1 & -r_i \\ -r_i & 1 \end{bmatrix} = \begin{bmatrix} 1 & -r_i \\ -r_i z^{-1} & z^{-1} \end{bmatrix}$$

Se trata de un sistema de solo polos (N):

(verde 6 A.R.)

$$\left\{ H(z) = \frac{b \cdot z^{-\frac{N}{2}}}{1 - \sum_{k=1}^N a_k \cdot z^{-k}} \right\}$$

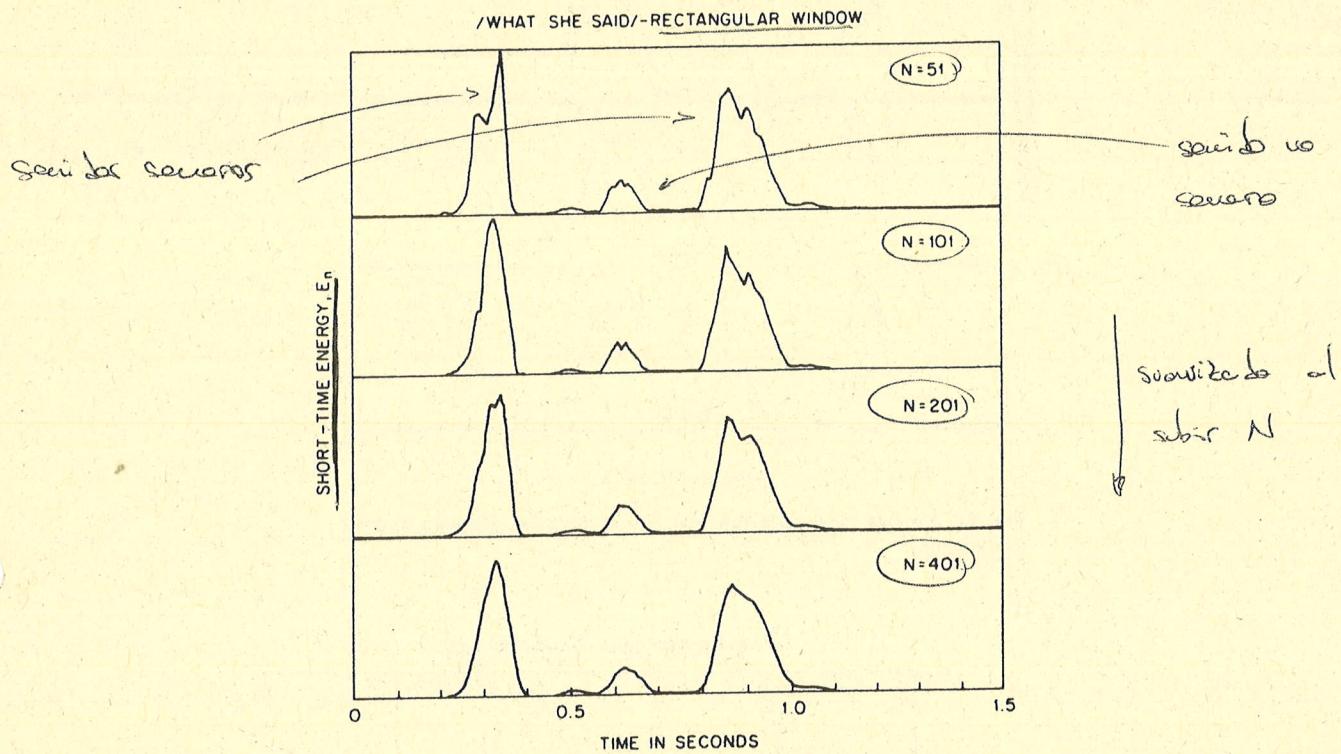
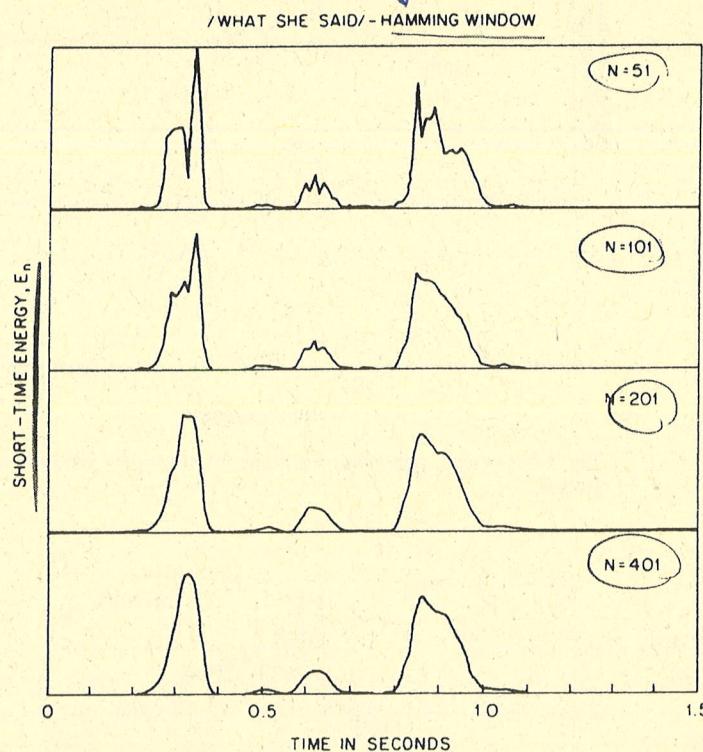
ENERGÍA LOCALIZADA

Fig. 4.6 Short-time energy functions for rectangular windows of various lengths.

ventana de Hamming: $w(u) =$

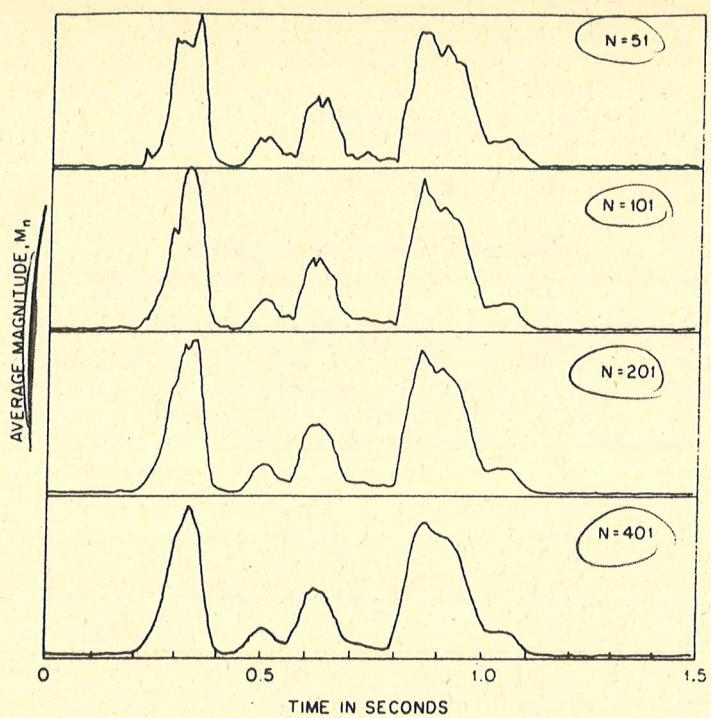
$$\begin{cases} 0,54 - 0,46 \cos \frac{2\pi u}{N-1} & (0 \leq u \leq N-1) \\ 0 & \text{resto} \end{cases}$$



Mejor ancho de banda
para igual duración de
ventana que la
rectangular

Fig. 4.7 Short-time energy functions for Hamming windows of various lengths.

MAGNITUD
LOCAL



II-1
Picos y valles / se han
aceptado como en la
energía local
↓
menos dificultad para
representar en punto fijo

Fig. 4.8 Average magnitude functions for rectangular windows of various lengths.

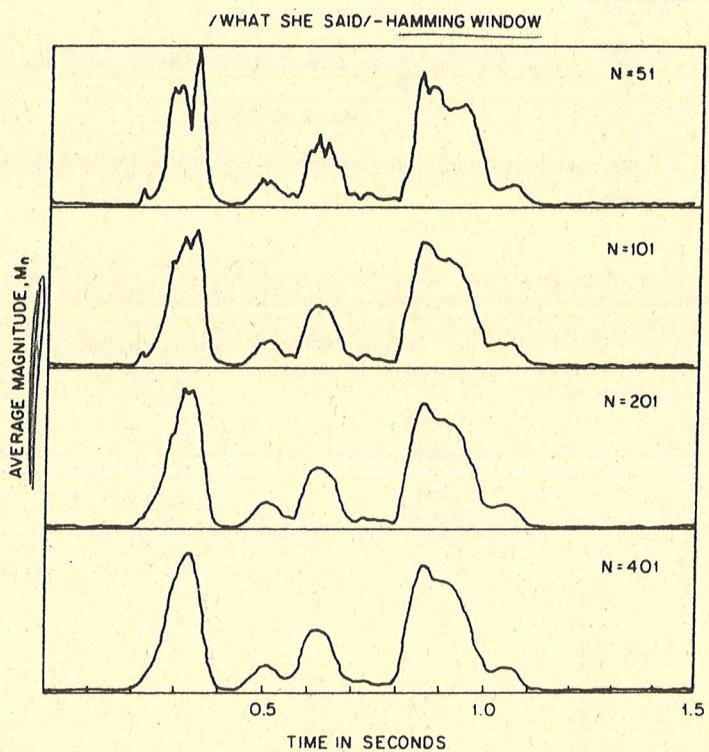
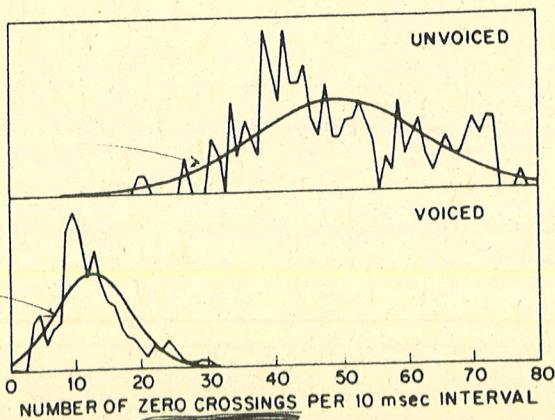


Fig. 4.9 Average magnitude functions for Hamming windows of various lengths.

CROZES PER CERO

se intenta aproximar
por una curva



frecuencias de conversación

frecuencias de conversación

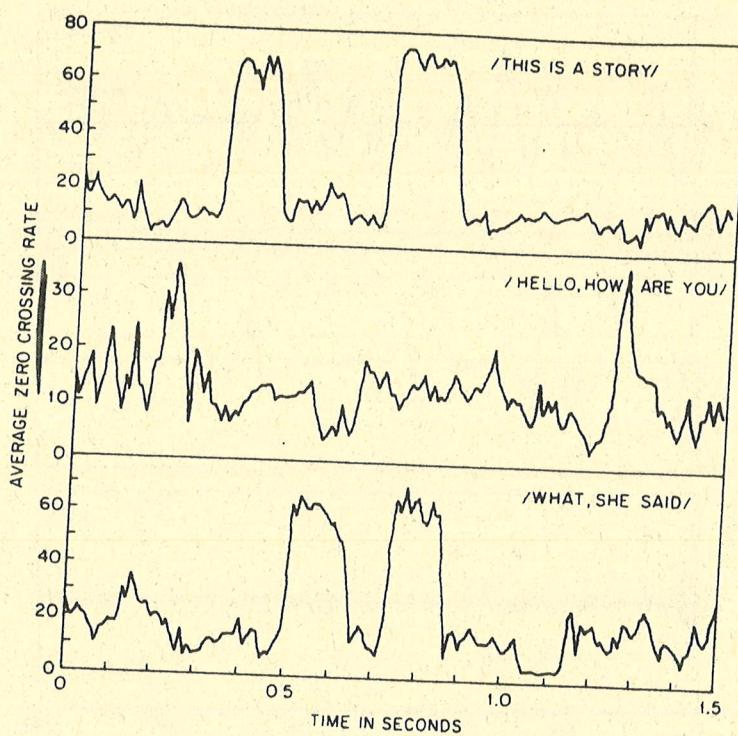


Fig. 4.12 Average zero-crossing rate for three different utterances.

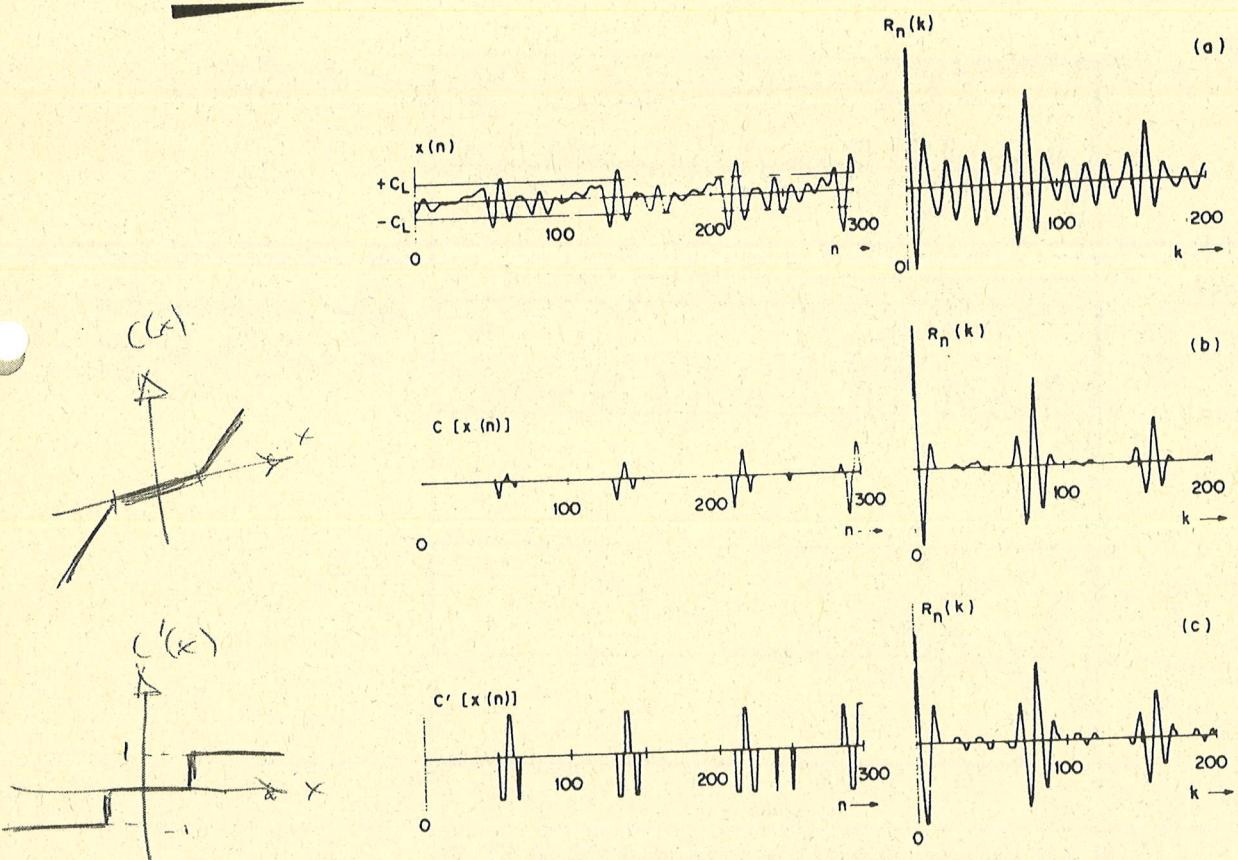


Fig. 4.33 Example of waveforms and correlation function: (a) no clipping; (b) center clipped; (c) 3-level center clipped. (All correlation functions normalized to 1.0.) (After Rabiner [18].)

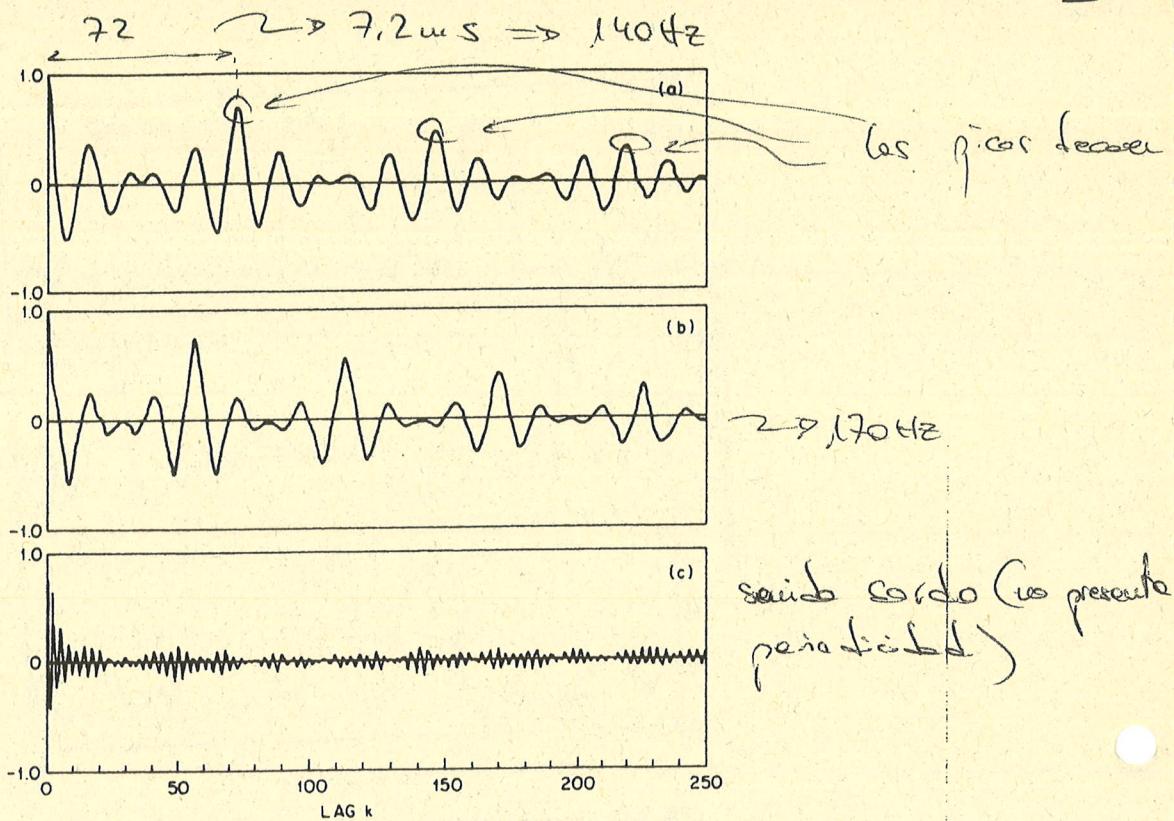


Fig. 4.24 Autocorrelation function for (a) and (b) voiced speech; and (c) unvoiced speech, using a rectangular window with $N = 401$.

Señal de voz muestra $\Rightarrow 140 \text{ Hz}$

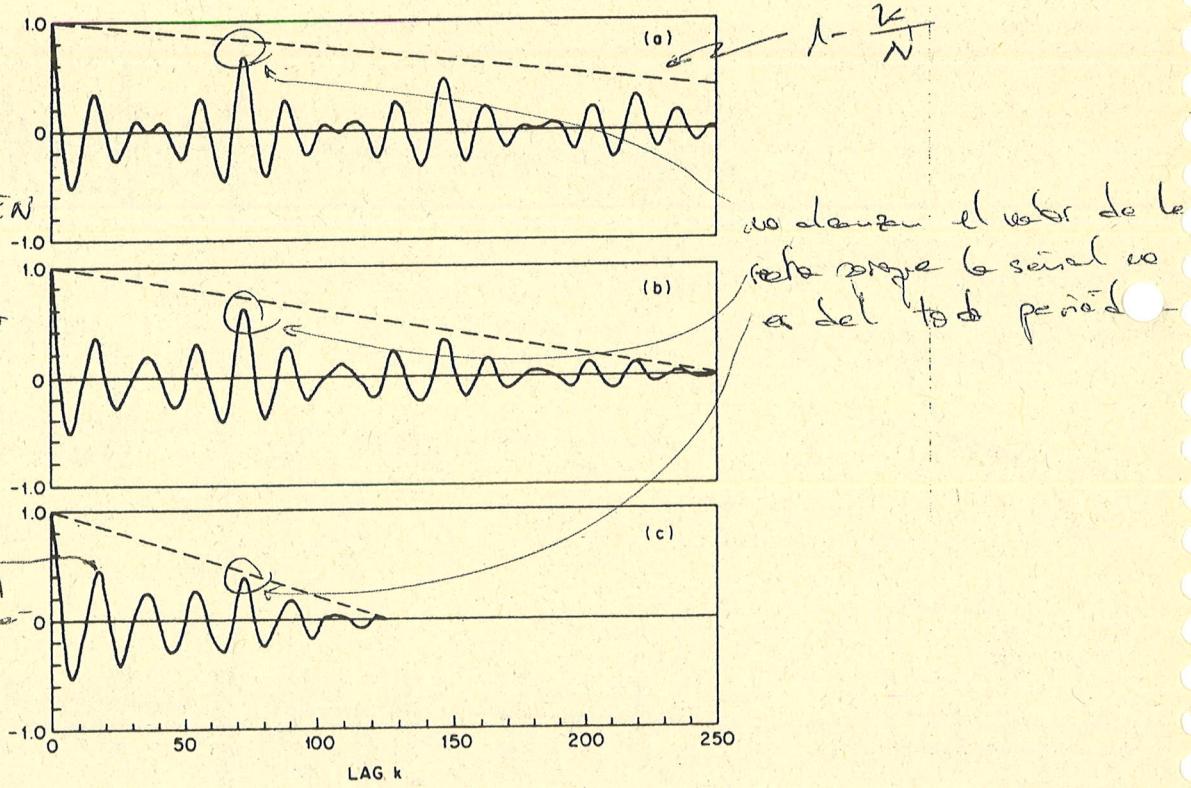


Fig. 4.26 Autocorrelation function for voiced speech with (a) $N = 401$; (b) $N = 251$; and (c) $N = 125$. Rectangular window used in all cases.

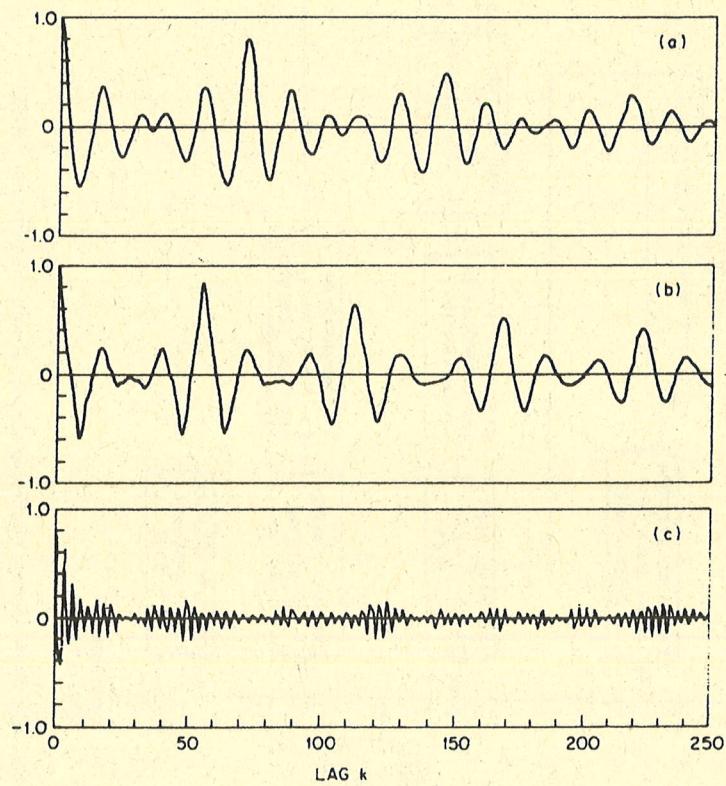
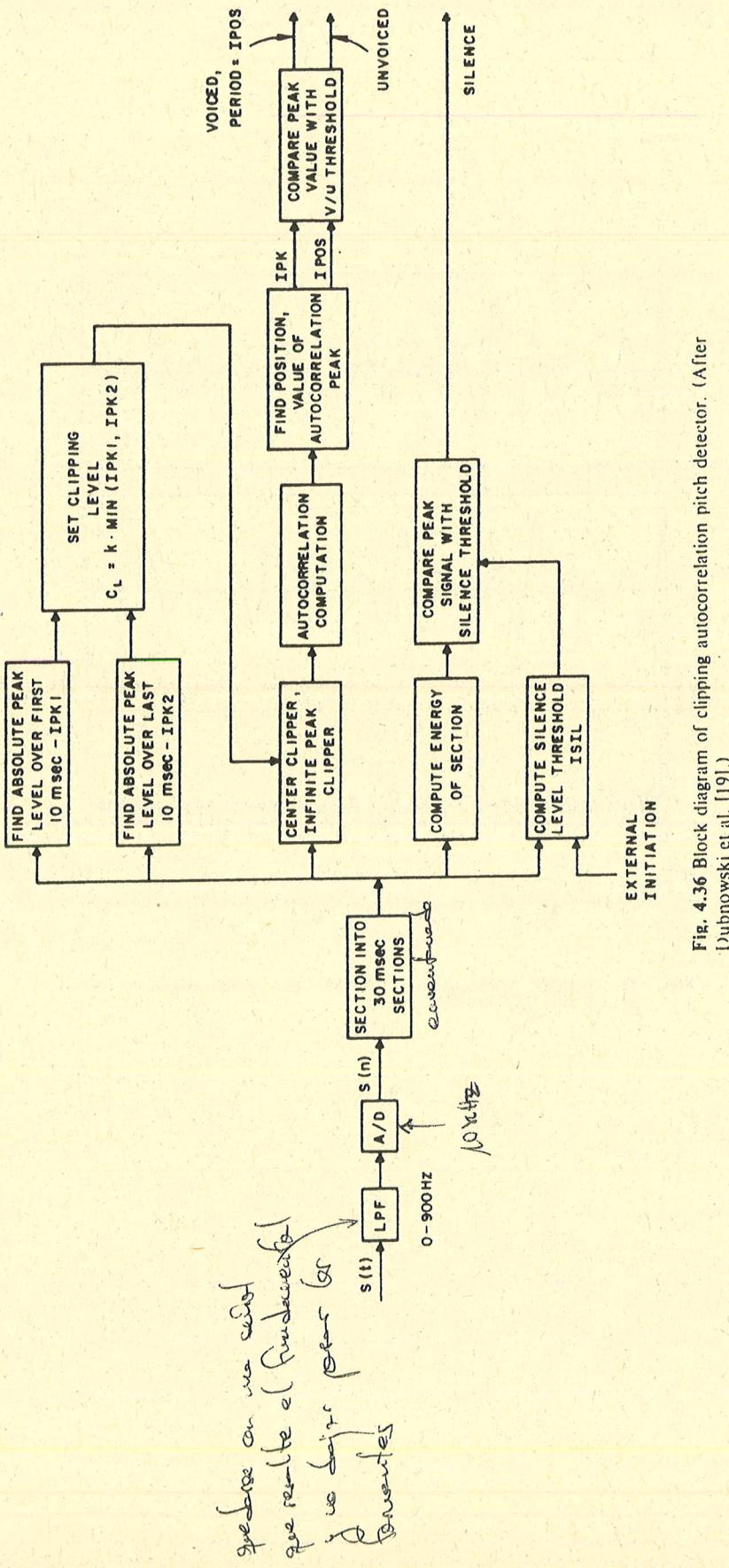


Fig. 4.28 Modified autocorrelation function for speech segments of Fig. 4.24 with $N = 401$.



SOUNDS SOUNDS

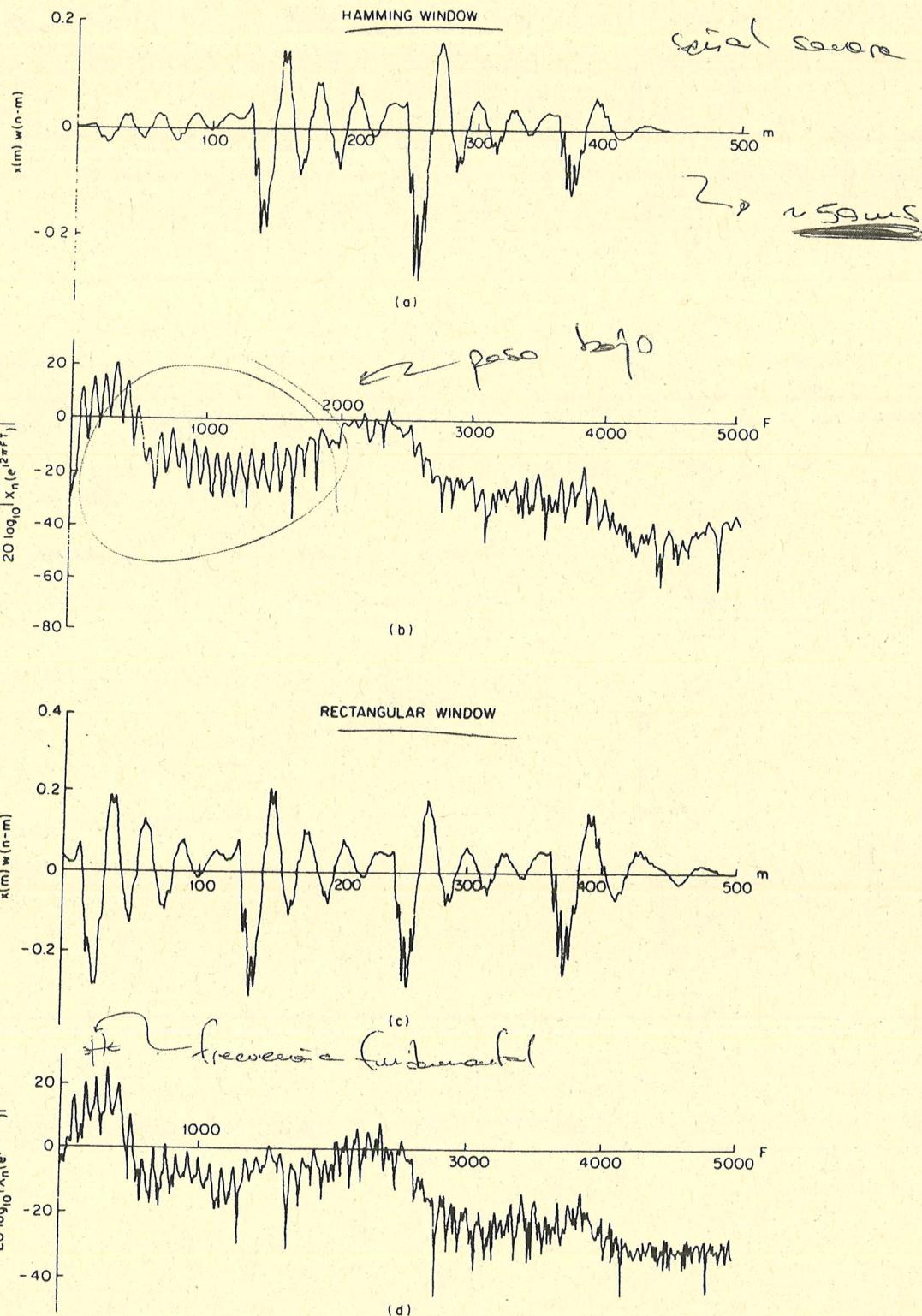
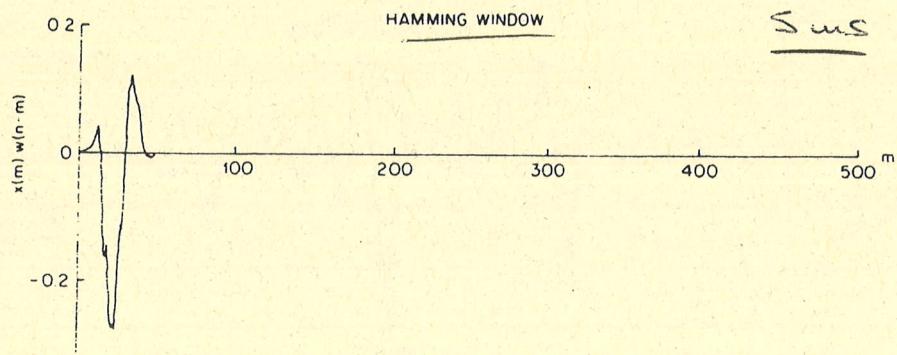


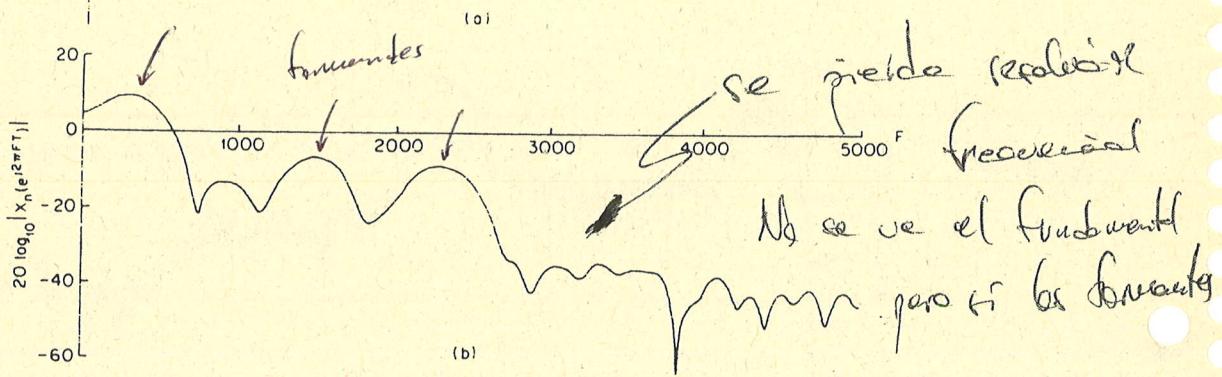
Fig. 6.2 Spectrum analysis for voiced speech using a 50 msec (a,b) Hamming window; (c,d) rectangular window. Parts (a) and (c) show time waveforms; parts (b) and (d) show corresponding spectra.

SEÑAL (SABRO)

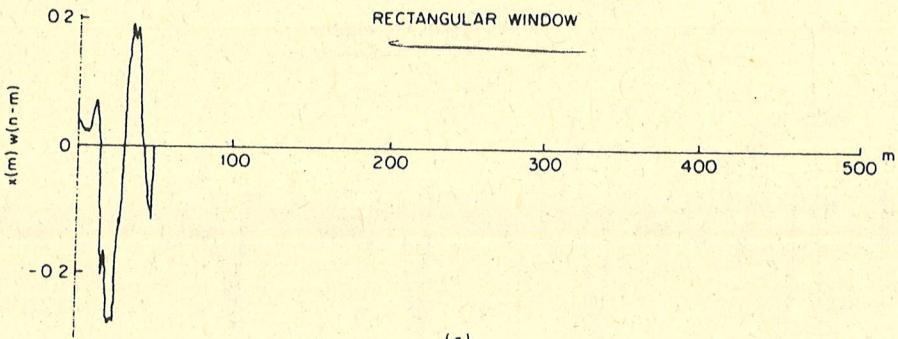
tiempo



frecuencia



tiempo



frecuencia

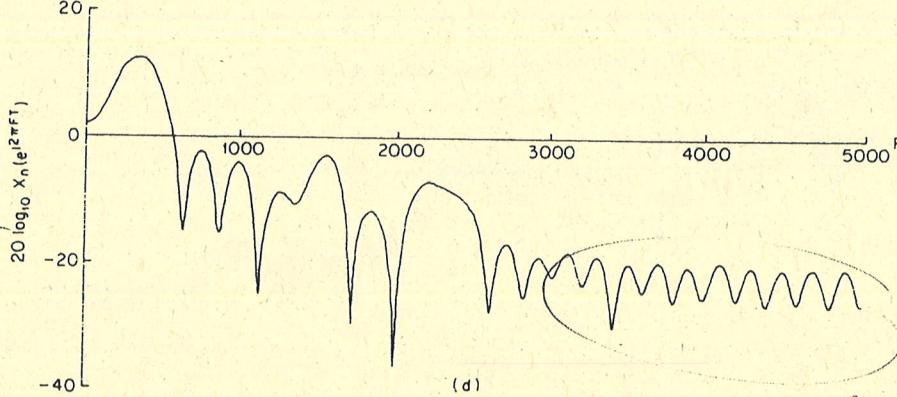


Fig. 6.3 Spectrum analysis of voiced speech using a 5 msec (a,b) Hamming window; (c,d) rectangular window. Parts (a) and (c) show time waveforms; parts (b) and (d) show corresponding spectra.

bolas
biteriales

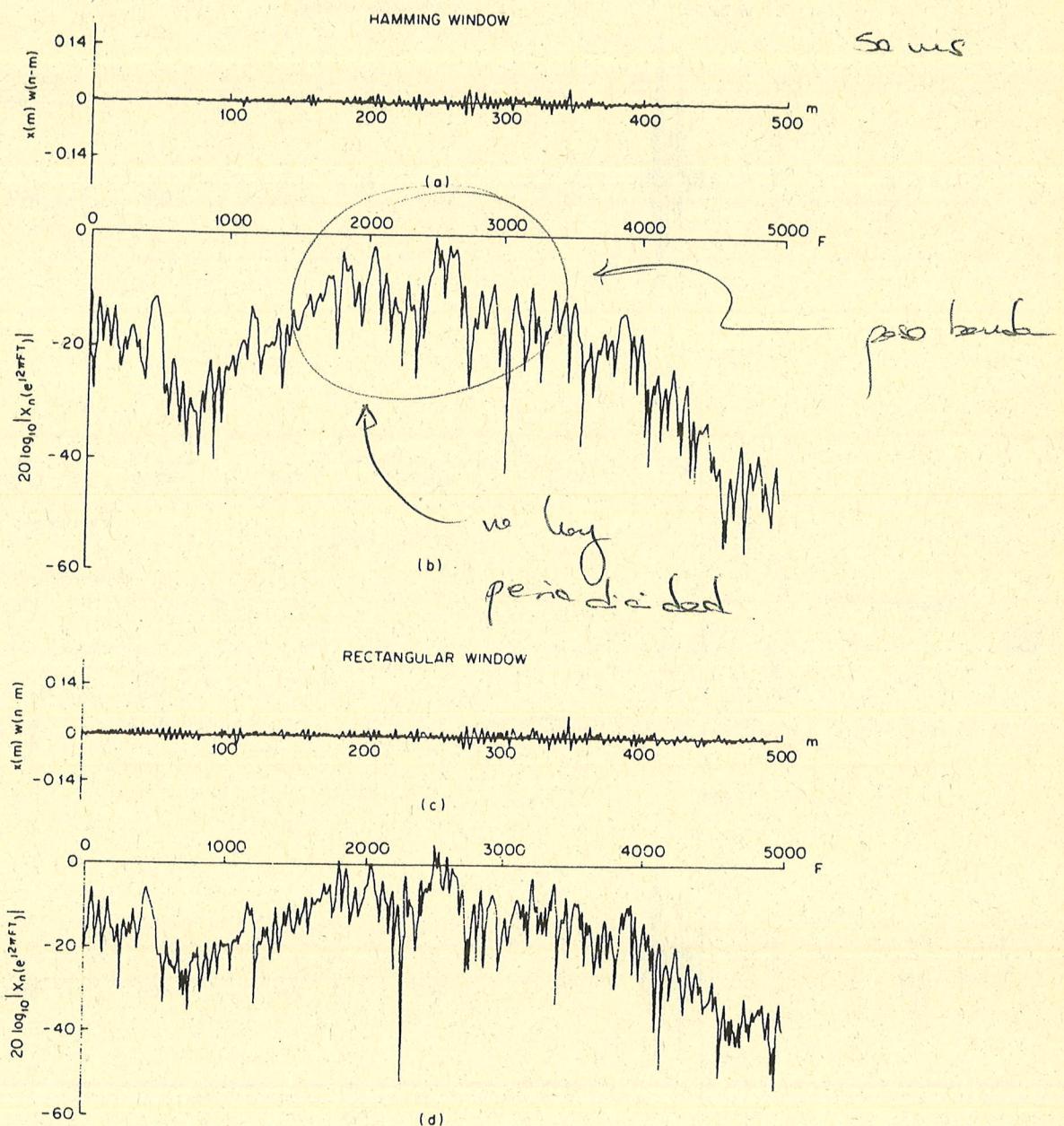


Fig. 6.4 Spectrum analysis of unvoiced speech using a 50 msec (a,b) Hamming window; (c,d) rectangular window. Parts (a) and (c) show time waveforms; parts (b) and (d) show corresponding spectra.

SONIDOS FRICATIVOS Y SONOROS

II-25

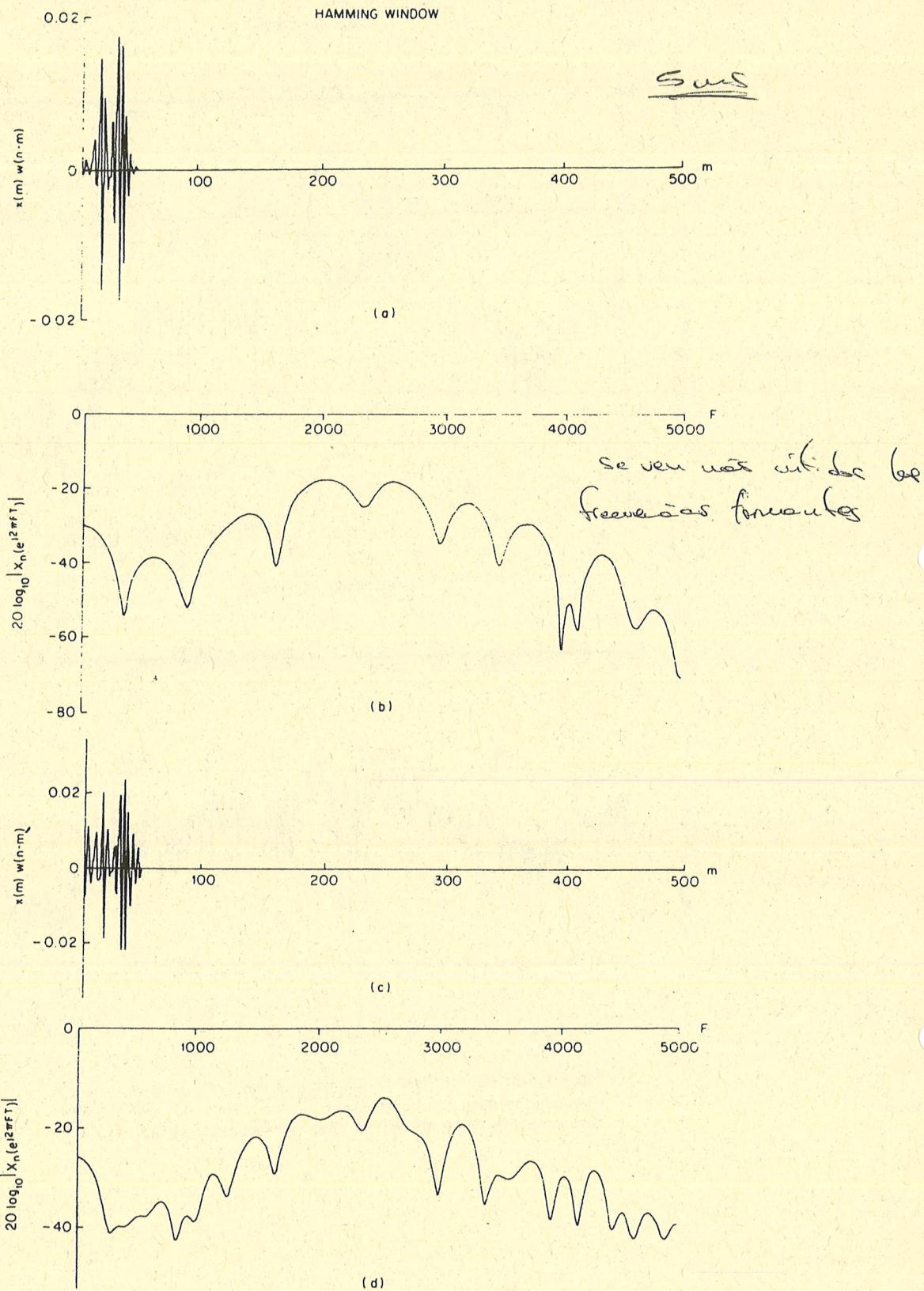


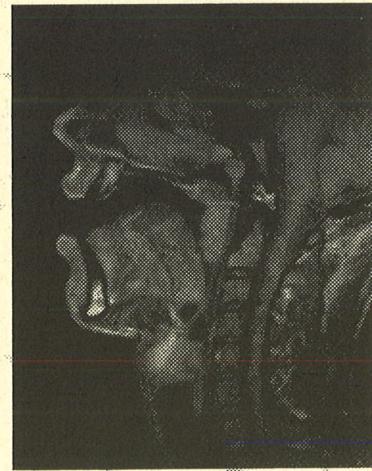
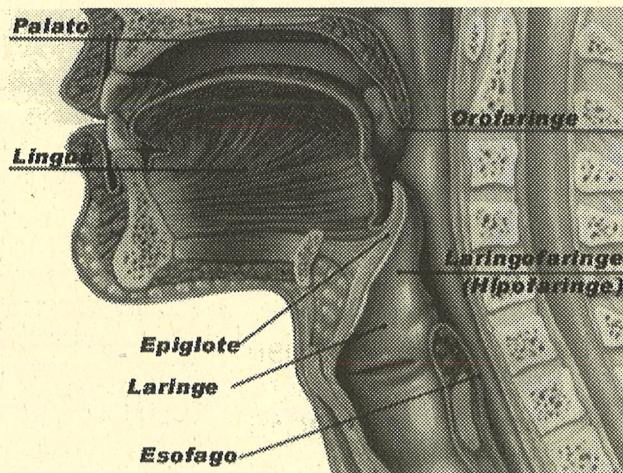
Fig. 6.5 Spectrum analysis of unvoiced speech using a 5 msec (a,b) Hamming window; (c,d) rectangular window. Parts (a) and (c) show time waveforms; parts (b) and (d) show corresponding spectra.

[La señal de voz]

Tratamiento Digital de la Señal II

Enrique Nava

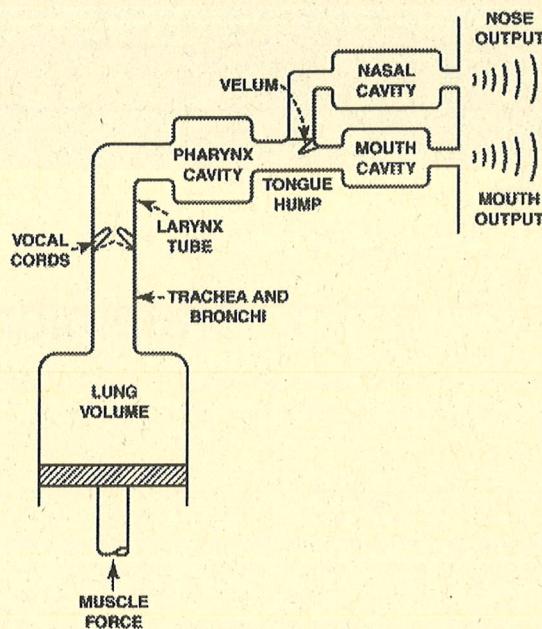
[La voz: fisiología]



Fonemas

■ $s(t) = s(p, v; t)$

sonido = ondas de presión acústica acopladas al movimiento de un fluido



Flanagan: "Speech Analysis and Perception", Springer-Verlag 1965

Fonemas

■ Alfabeto fonético internacional (IPA)

THE INTERNATIONAL PHONETIC ALPHABET (revised to 1993)

CONSONANTS (PULMONIC)

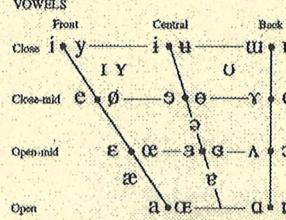
	Bilabial	Lateral/bilabial	Bilabial	Alveolar	Postalveolar	Retroflex	Palatal	Vocal	Uvular	Pharyngeal	Glottal
Phrasal	p b		t d		t̪ d̪ c̪ f̪	k g	q q̪ G̪				?
Nasal	m	n̪	n		ɳ̪	p̪	ɳ̪				N
Trill	R		r								R
Tap or Flap			t̪		t̪						
Plosive	ɸ β	f v	θ ð	s z	ʃ ʒ	s z	ç ɿ	x ɻ	χ ɬ	h ɬ̪	h ɬ̪
Lateral plosive				t̪ l̪							
Approximant		v	i	ɿ	j	w̪					
Lateral approximant			l̪		ɿ	ɬ̪					

Where symbols appear in pairs, the one to the right represents a voiced consonant. Shaded areas denote articulations judged impossible.

CONSONANTS (NON-PULMONIC)

Clicks	Voiced implosives	Ejectives
O Bilabial	b̪	Bilabial
D Dental	d̪	Dental/alveolar
F Postalveolar	f̪	Palatal
G Palatoalveolar	g̪	Vocal
G Alveolar lateral	g̪	Uvular

VOWELS



Where symbols appear in pairs, the one to the right represents a rounded vowel.

OTHER SYMBOLS

ʍ	Voiceless labio-velar fricative
w̪	Voiceless labio-velar approximant
ɥ	Voiceless lateral-palatal approximant
χ	Voiceless epiglottal fricative
ɸ	Voiceless glottal fricative
ʔ	Epiglottal plosive

χ ɬ̪ Alveolo-palatal fricative

ɿ Alveolar lateral flap

ʃ ɻ Sibilants

ts Alveolo-palatal sibilants

Affricates and double articulations can be represented by two symbols joined by a bar if necessary.

kp ts

SUPERSEGMENTALS

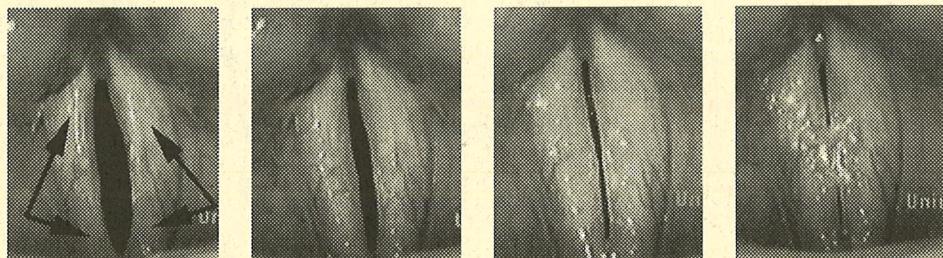
Primary stress	Secondary stress	Lengthening	TONE & WORD ACCENTS
↑	↓	ε or ε̄	Level
Long	ei	ε̄	High
Half-long	ē	ε̄	Mid
Extra-short	ɛ̄	ε̄	High+mid
Syllable break	ə, ekt̪	ε̄	Low
Mixis (beat) group	ε̄	ε̄	Extra-low
Major (intonation) group	↓	ε̄	Rising-falling
Linking (between a word)	↑	Upstep	Global rise
			Global fall

DIACRITICS

Diacritics may be placed above a symbol with a descender, e.g. ī		
Vocaliss n̪ d̪	Breathy voiced b̪ a̪	Dental t̪ d̪
Voiced § t̪	Creaky voiced b̪ a̪	Apical t̪ d̪
Aspirated t̪̪ d̪̪	Languidized t̪̪ d̪̪	Laminal t̪̪ d̪̪
More rounded ɔ̄	Labialized t̪ʷ d̪ʷ	Nasalized ɛ̄
Less rounded ɔ̄	Palatalized t̪̄ d̪̄	Nasal release d̄
Advanced ɥ̄	Volarized t̪̄ d̪̄	Lateral release d̄
Reduced ī	Pharyngealized t̪̄ d̪̄	No middle release d̄
Centralized ē̄	Velarized t̪̄ d̪̄	Velarized or pharyngealized t̄
Mid-centralized ē̄	Raised ɛ̄ (↓ = voiced alveolar fricative)	
Syllabic ɬ̄	Lowered ɛ̄ (↓ = voiced bilabial approximant)	
Non-syllabic ɬ̄	Advanced Tongue Root ɬ̄	
Rhoticity ð̄	Retracted Tongue Root ɬ̄	

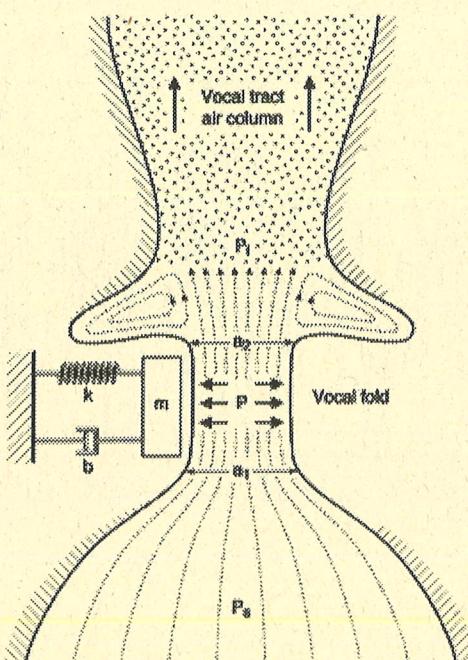
Fonemas: tipos de sonidos

- Sonoros: las cuerdas vocales vibran:
 - Sonidos casi-periódicos (vocales)



- No sonoros (sordos)
 - Cuerdas vocales abiertas, sonidos “tipo ruido” /s/

Sonido sonoro



Frecuencia fundamental

$$F_0 = \frac{1}{2L} \sqrt{\frac{\sigma}{\rho}}$$

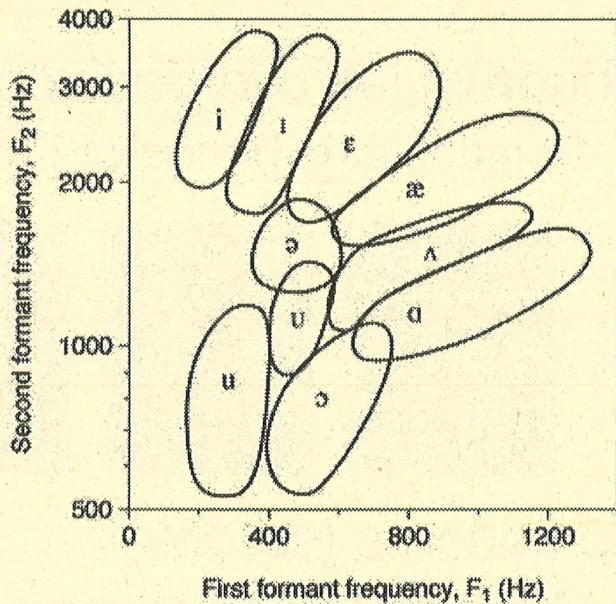
Tensión longitudinal

Longitud de las cuerdas

Densidad del tejido

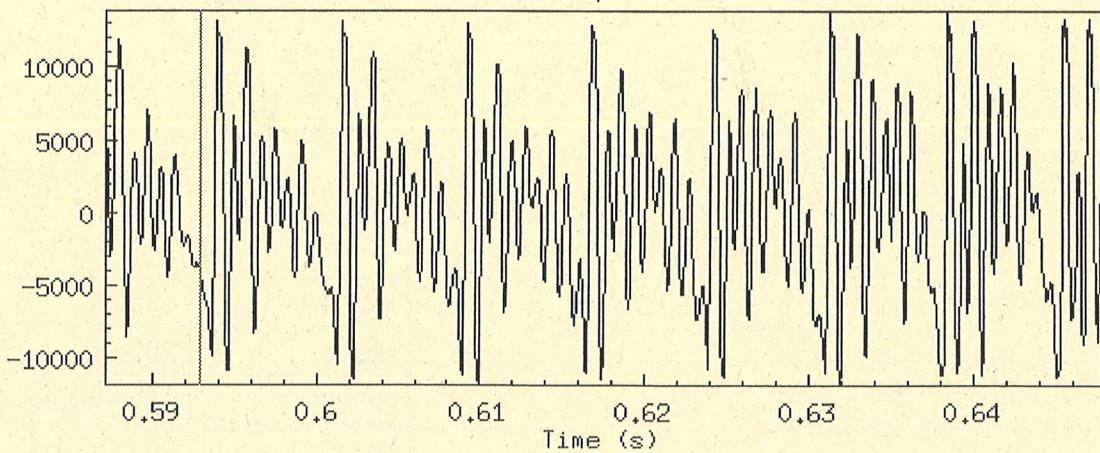
Hombre	125 Hz
Mujer	200 Hz
Niño	250-400 Hz
Bebé	500 Hz

Vocales



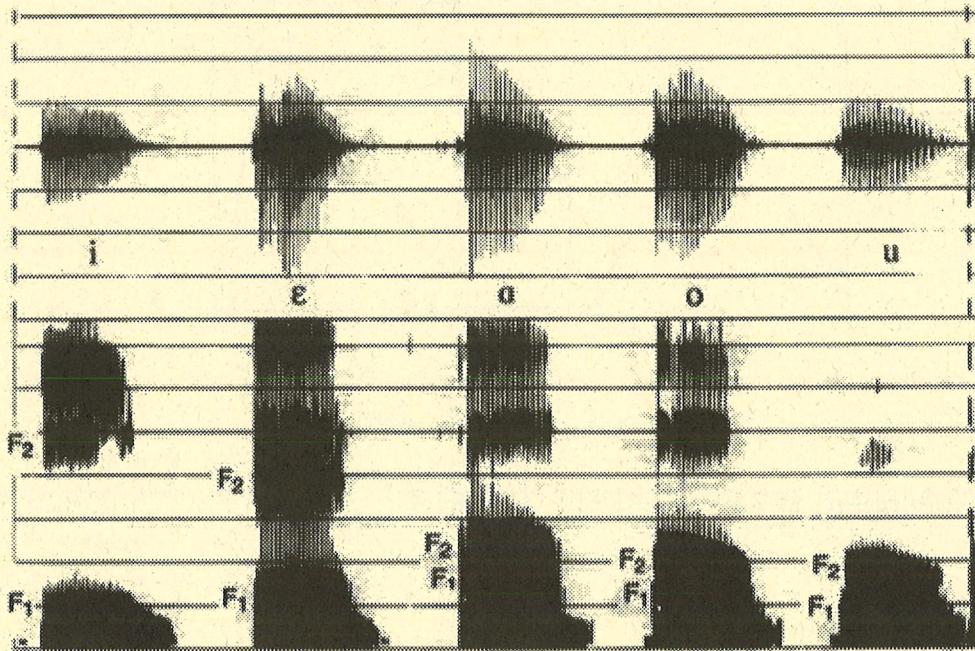
Peterson, Barney 1952

Vocales



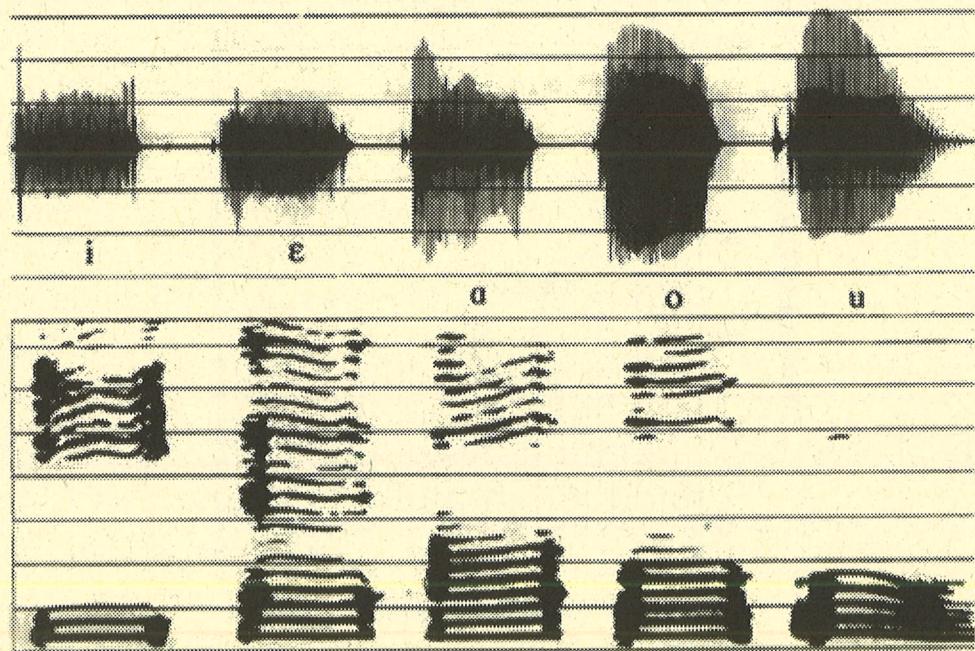
[Espectrograma]

Time waveform and broadband spectrogram of [i e a o u].

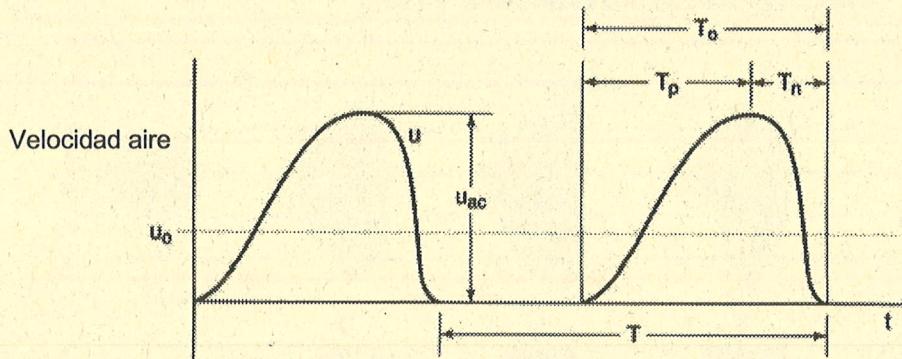


[Espectrograma]

Time waveform and narrowband spectrogram of [i e a o u].



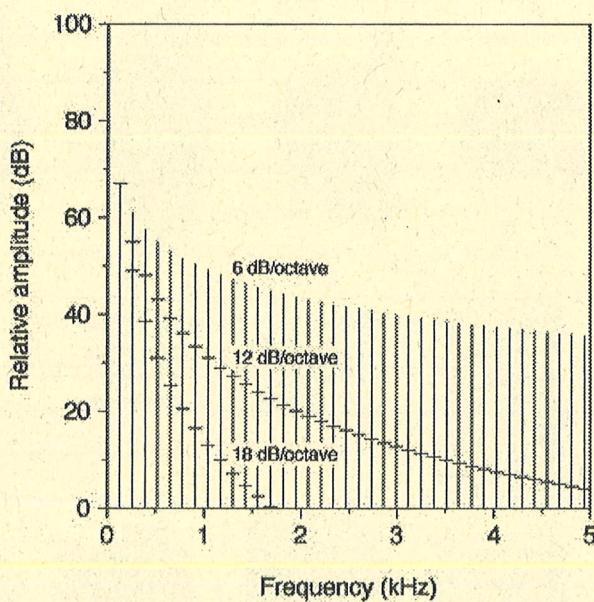
Pulsos glotales



Rosenberg 1971:

$$u[n] = \begin{cases} \frac{1}{2} \left(1 - \cos\left(\frac{\pi n}{T_p}\right) \right) & 0 \leq n \leq T_p \\ \cos\left(\pi \frac{n - T_p}{2T_n}\right) & T_p \leq n \leq T_p + T_n \\ 0 & \text{resto} \end{cases}$$

Efecto de radiación



Modelo de voz

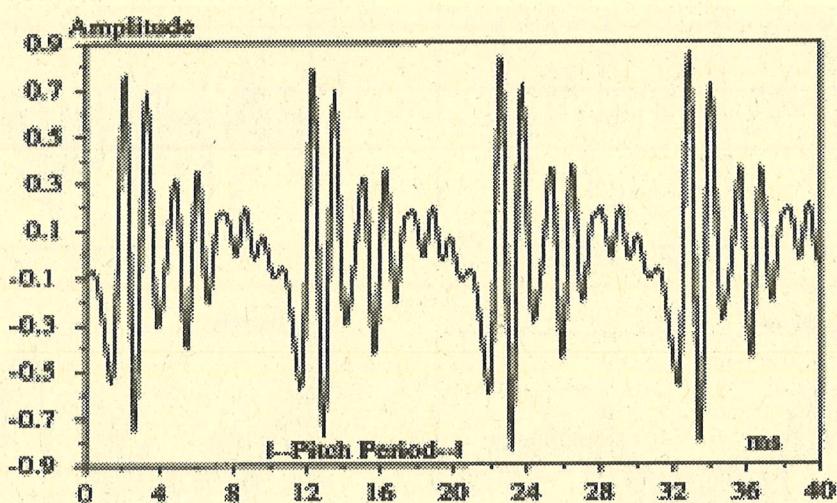


FIGURE 2.3

Time-domain waveform of a short segment of voiced speech, x-axis units in ms, y axis is relative amplitude of sound pressure.

Modelo de voz

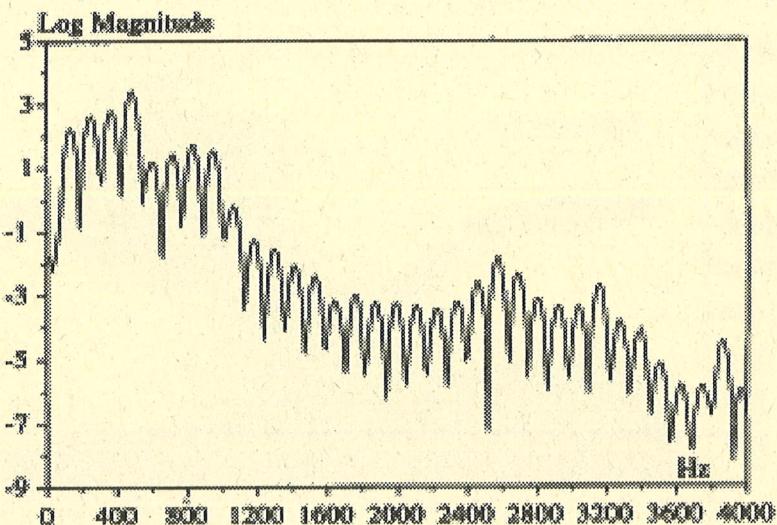


FIGURE 2.4

Log magnitude spectrum of a short segment of voiced speech, X axis unit in Hz.

Fonemas

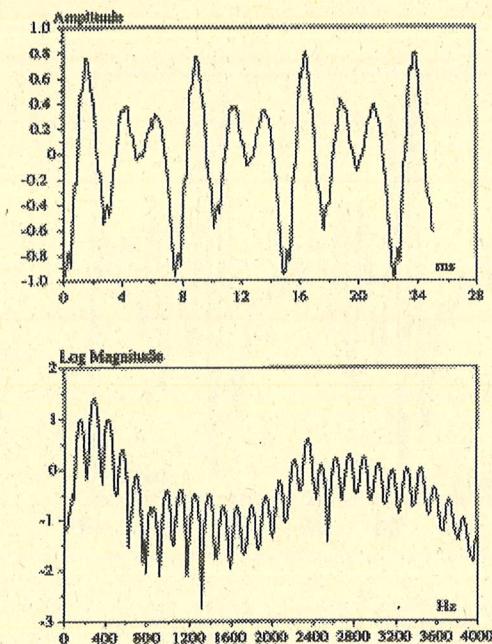


FIGURE 2.5
Time waveform and log magnitude spectrum of /ɪ/, as in the word "bit."

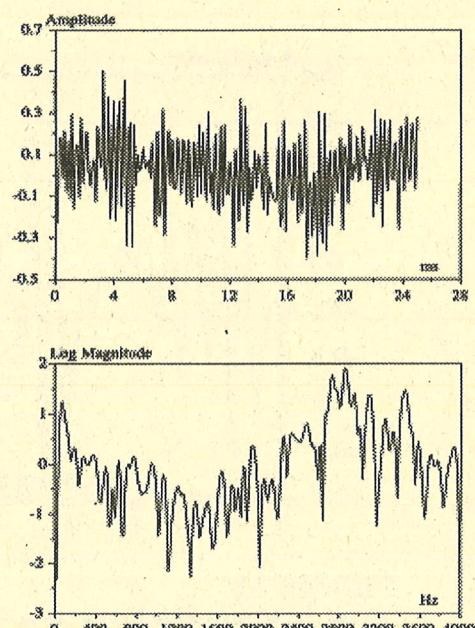


FIGURE 2.7
Time waveform and log magnitude spectrum of /sh/, as in the beginning of the word "shop."

Sonogramma

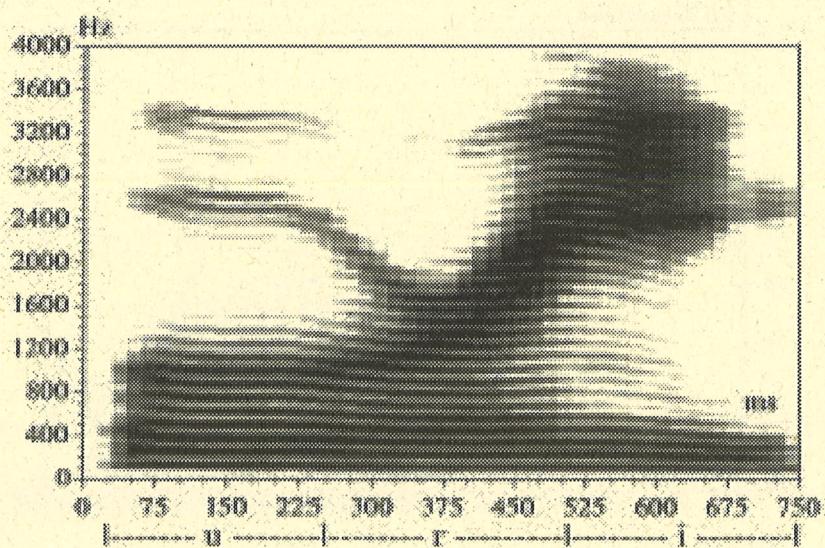
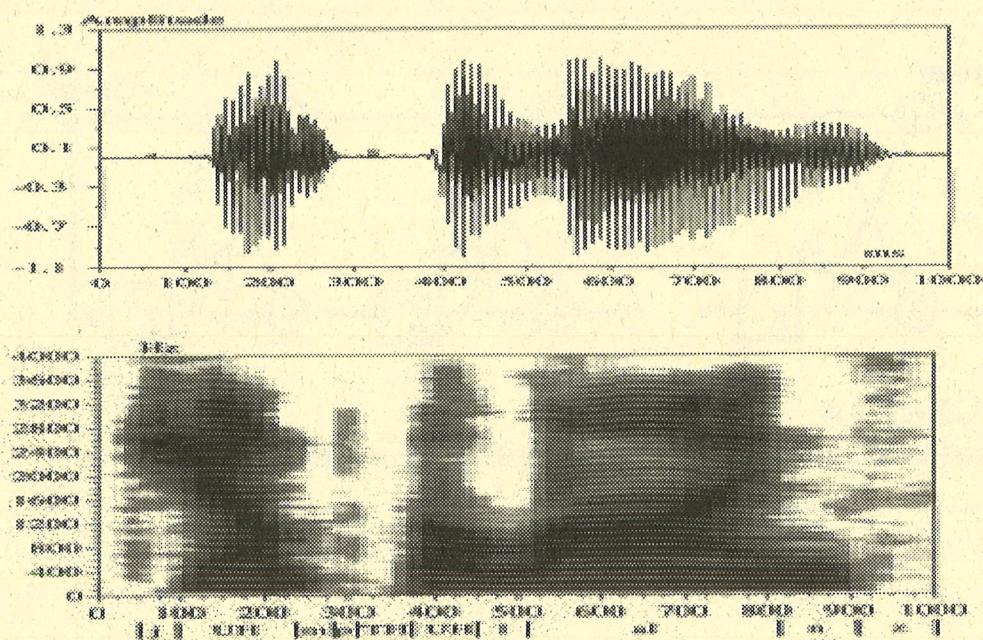


FIGURE 11
Spectrogram of nonword utterance /u-r-i/.

Sonogramas



Modelo de tubos

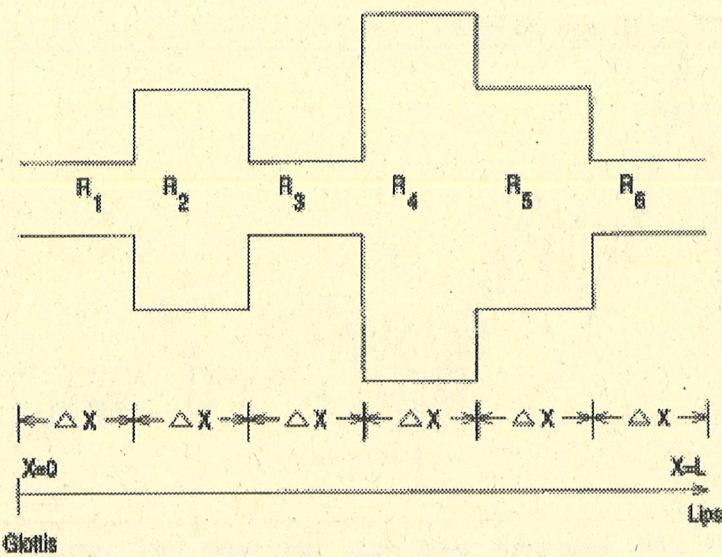


FIGURE 4.3
Multiple concatenated tube model.

Modelo de tubos

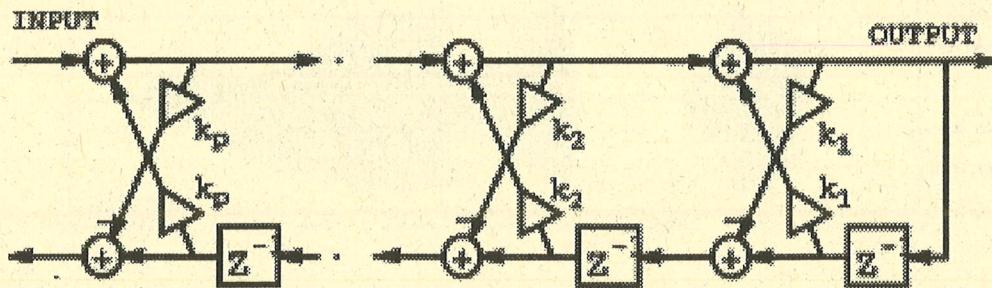


FIGURE 4.4
Lattice filter realization of multiple-tube model.

Modelo de tubos

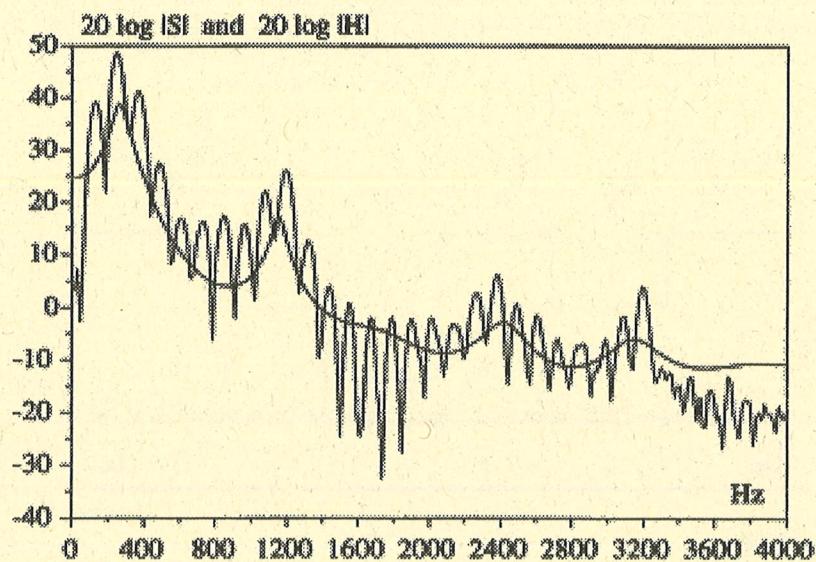


FIGURE 4.6
Log magnitude of DFT and LP spectra for a segment of voiced speech.

Modelo de tubos

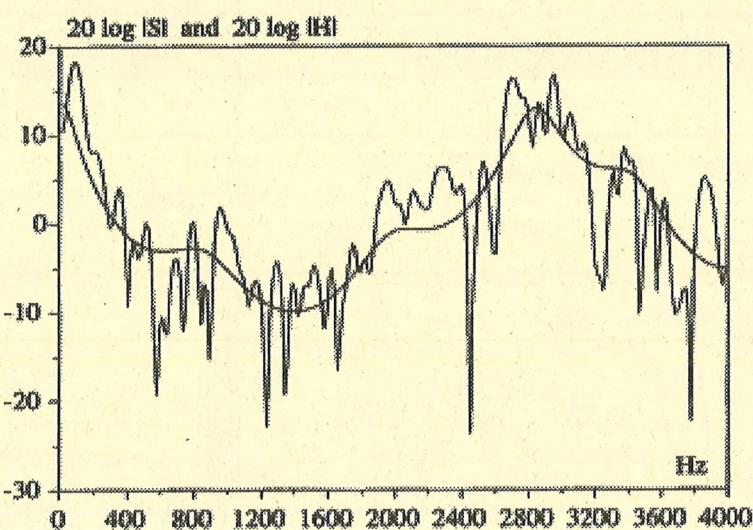


FIGURE 4.7
Log magnitude of DFT and LP spectra for a segment of unvoiced speech.

Detección del tono fundamental

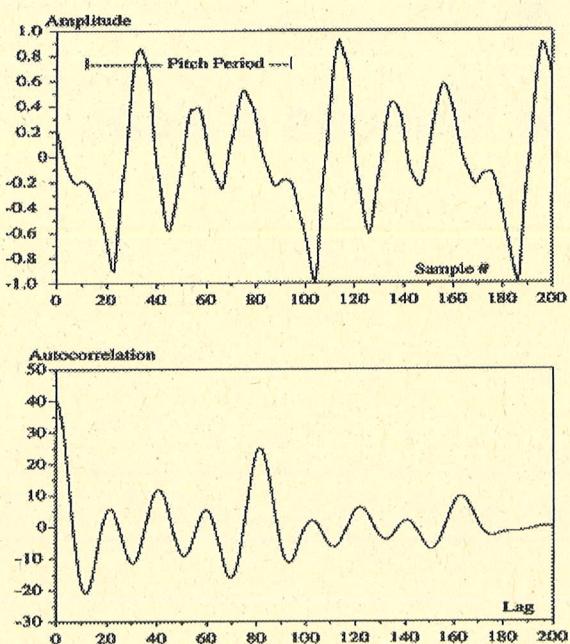


FIGURE 5.1
Time-domain waveform and autocorrelation of a short segment of voiced speech.

Detección del tono fundamental

Análisis cepstral

$$\text{Cepstrum}(d) = \text{IFFT}(\log_{10}|\text{FFT}(z(n))|)$$

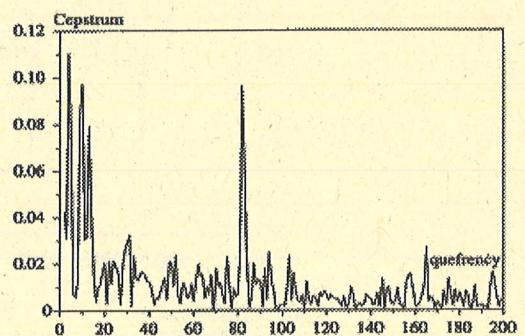
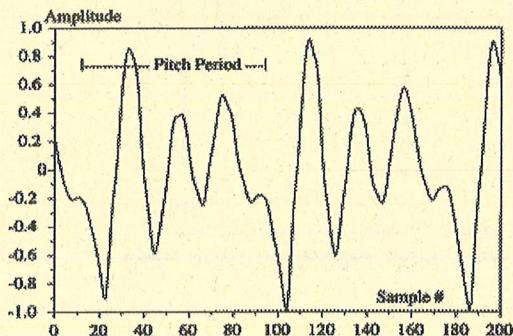
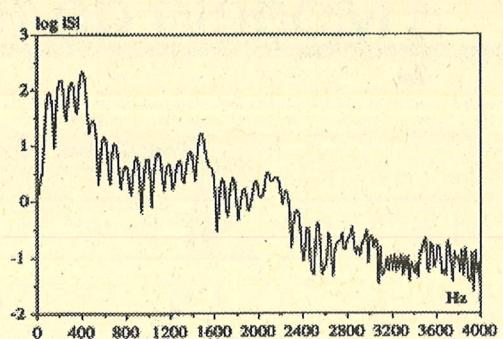


FIGURE 5.6
Log magnitude of DFT and cepstrum of speech segment of Figure 5.1.

Percepción humana

Bandas críticas

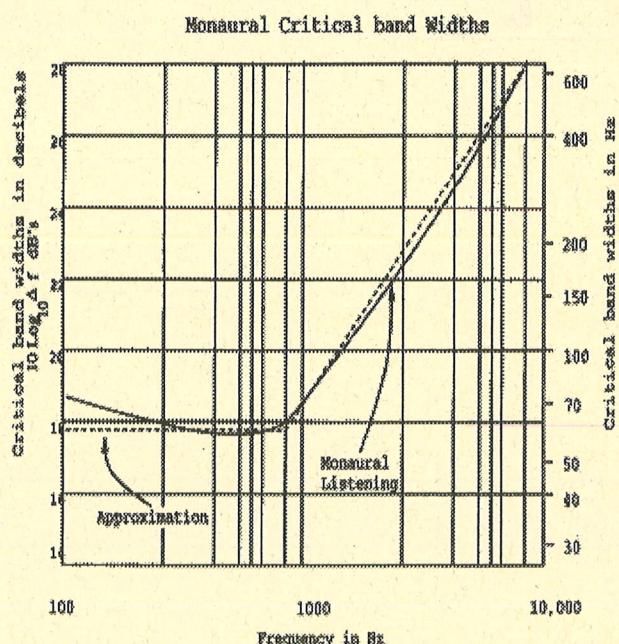


FIGURE 6.1
Frequency width of critical bands as a function of the band center frequency.

Escalas de frecuencia

Critical Band No. (Barks)	Frequency (Hz)	Mels
1	20-100	0-150
2	100-200	150-300
3	200-300	300-400
4	300-400	400-500
5	400-510	500-600
6	510-630	600-700
7	630-770	700-800
8	770-920	800-950
9	920-1080	950-1050
10	1080-1270	1050-1150
11	1270-1480	1150-1300
12	1480-1720	1300-1400
13	1720-2000	1400-1550
14	2000-2320	1550-1700
15	2320-2700	1700-1850
16	2700-3150	1850-2000
17	3150-3700	2000-2150
18	3700-4400	2150-2300
19	4400-5300	2300-2500
20	5300-6400	2500-2700
21	6400-7200	2700-2850
22	7200-9500	2850-3050

Table 6.1 The relationship between the frequency units: Barks, Hertz, and Mels.

Audibilidad

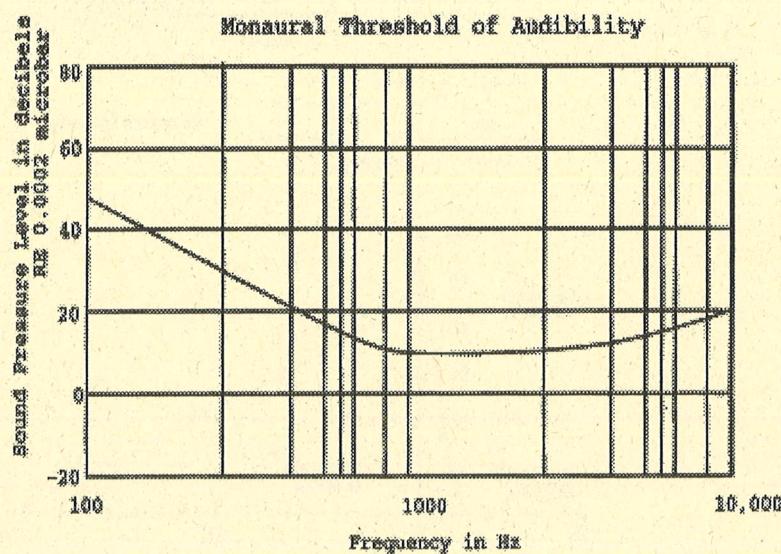


FIGURE 6.3
Threshold of audibility for a pure tone in silence.

Enmascaramiento espectral

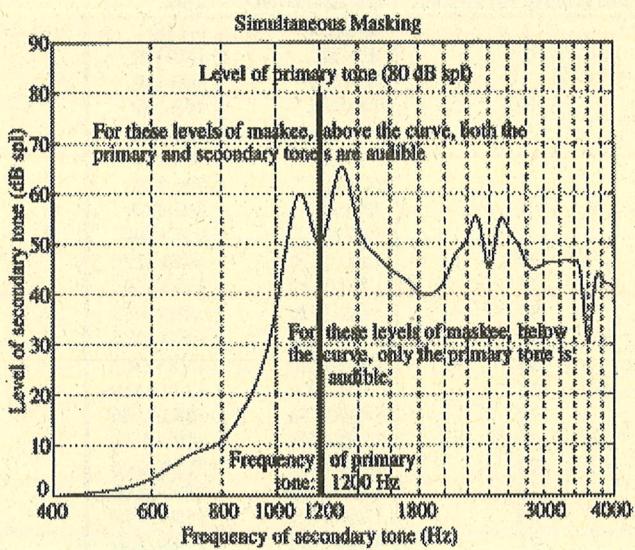
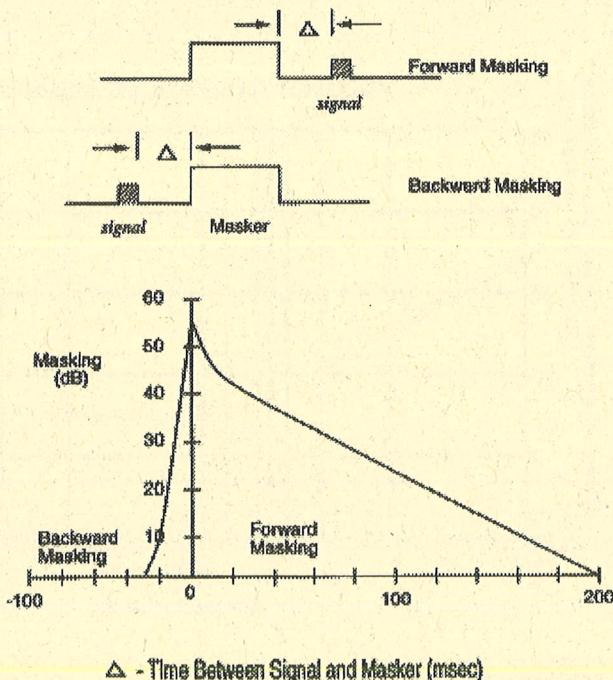


FIGURE 6.4
Simultaneous masking in frequency of one tone on another tone
(data adapted from [81]).

Enmascaramiento temporal



Δ - Time Between Signal and Masker (msec)

FIGURE 6.5
Illustration of the effect of temporal masking.

Enmascaramiento psicoacústico

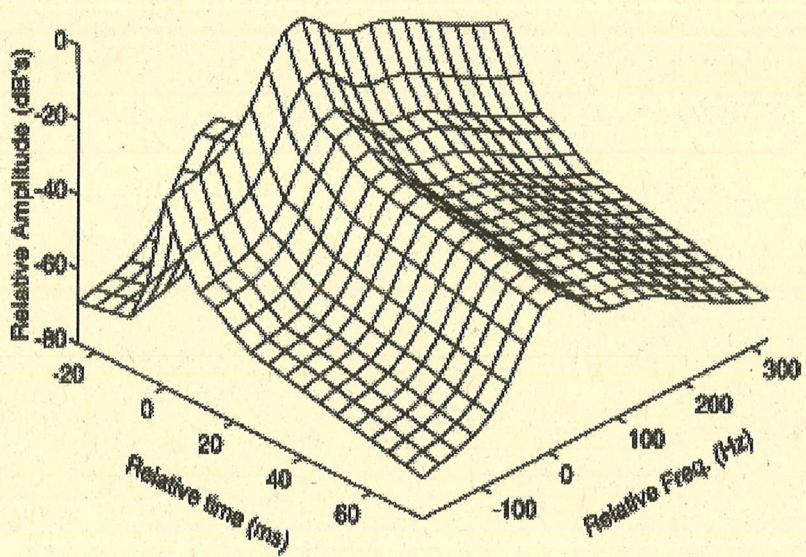


FIGURE 12.3
Psycho-acoustic masking data, both temporal and frequency.

