



UNIVERSITÀ
DEGLI STUDI DI TRIESTE



Psicoacustica

A.Carini – Elettronica per l'audio e l'acustica

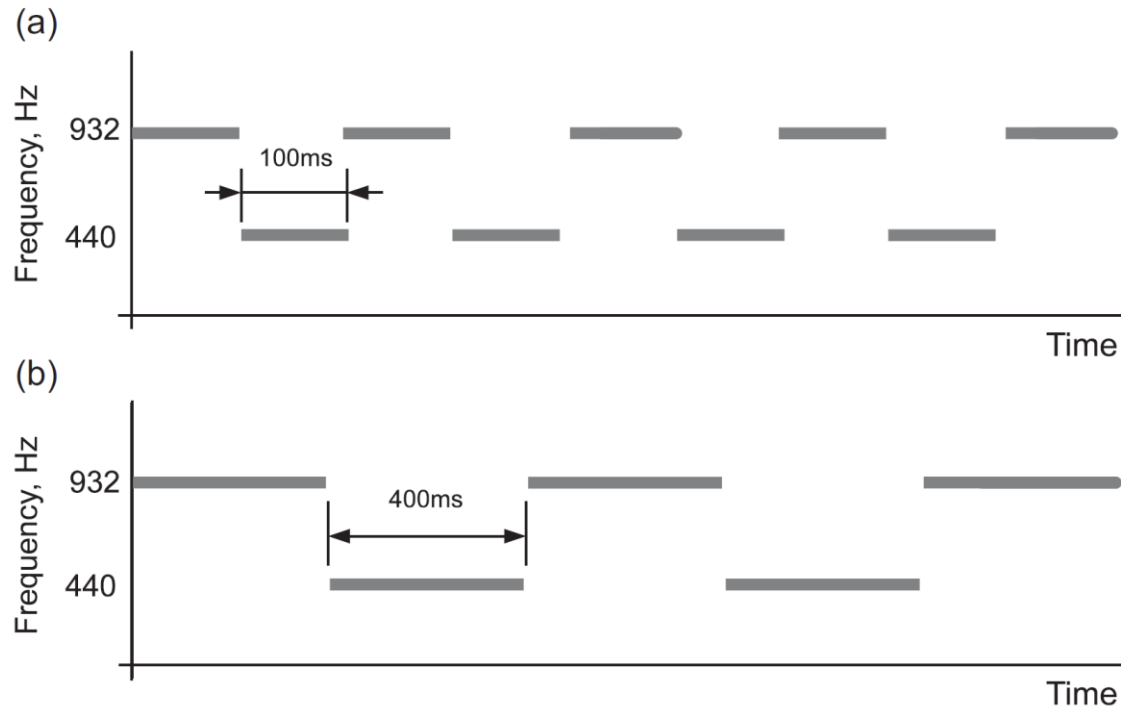
Analisi della scena sonora

- L'analisi della scena sonora descrive il processo uditivo che viene usato per trattare e interpretare combinazioni di suoni complessi.
- Quando sulla stessa scena ci sono più suoni, l'orecchio riesce a sintonizzarsi sul singolo suono che vuole ascoltare estraendolo dalla scena.
- Ci sono diversi meccanismi che intervengono nell'analisi di una scena sonora:
 - La prossimità,
 - La chiusura,
 - Il destino comune,
 - La buona continuazione.

La prossimità

- Suoni in prossimità sono quelli vicini in termini di ampiezza, pitch, durata, timbro.
- Il cervello ha la tendenza a classificare suoni in prossimità come appartenenti alla stessa sorgente.

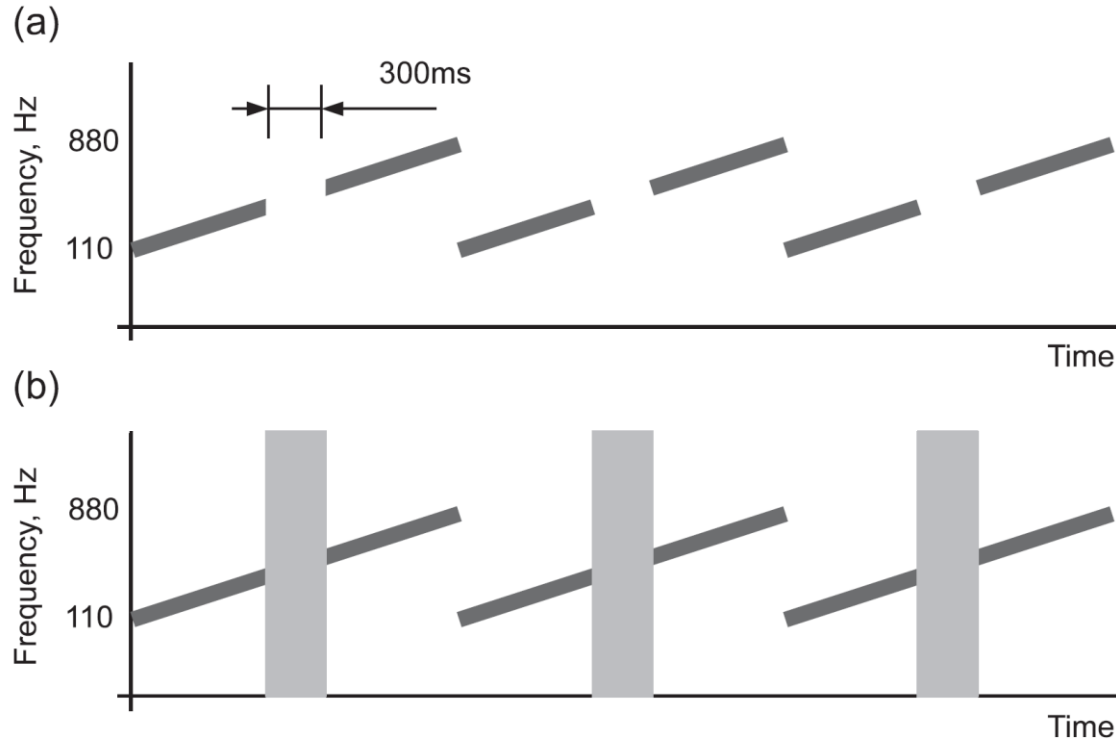
La prossimità - esempio



Chiusura

- Il principio di chiusura consente al cervello di riempire eventuali «buchi» per soddisfare le sue aspettative rispetto alla continuità o all'origine del suono.
- Il fenomeno è anche noto con il nome di induzione uditoria: il cervello non deduce la continuazione del suono mancante ma induce/crea un rimpiazzo per coprire i buchi.

Chiusura - esempio



Destino comune

- Quando gruppi di toni o rumori partono, si arrestano, fluttuano insieme, vengono normalmente interpretati come parte di un suono combinato avente una sorgente comune.

Destino comune - esempio

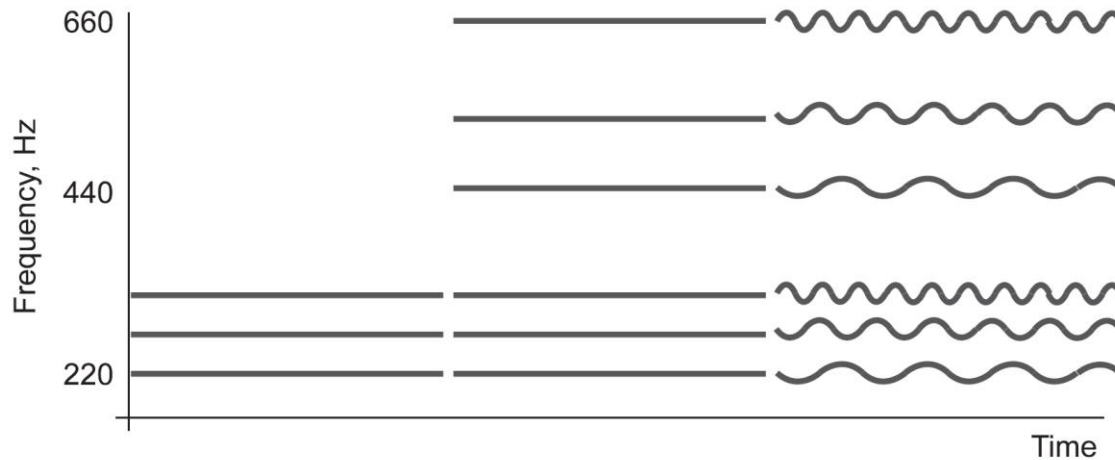
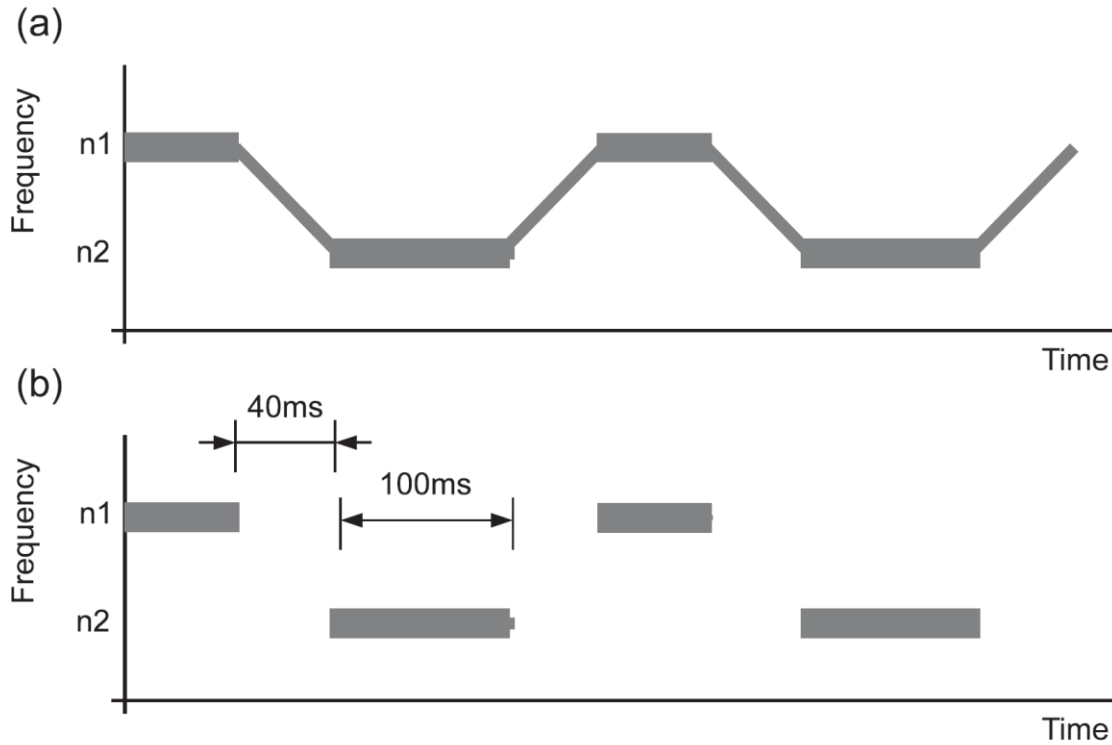


Figure 5.3 Illustration of common fate by comparing three chords. The first generates a pleasing note from three related sinewaves of different frequency. The second adds a first harmonic to each of the three fundamentals. The third modulates each of the three tones plus their first harmonic with unrelated modulating frequencies.

La buona continuazione

- In genere in natura i suoni non iniziano, si fermano o cambiano di frequenza istantaneamente. Ci sarà normalmente un «attacco» dell'ampiezza all'inizio, un decadimento alla fine. La frequenza istantanea slitterà da una frequenza all'altra più o meno lentamente.
- In presenza di suoni complessi, il cervello tende a classificare suoni connessi da queste variazioni graduali come provenienti dalla stessa sorgente. Al contrario, suoni che non sono legati vengono ritenuti provenienti da sorgenti diverse.
- Il fenomeno è detto della «buona continuazione» ma più propriamente dovrebbe essere chiamato della connettività dei suoni.

La buona continuazione - esempio



Modellazione psicoacustica

- La registrazione di un'onda sonora «fisica» contiene elementi che sono molto rilevanti per un ascoltatore e elementi che non lo sono affatto.
- Esempio: il mascheramento post-stimolo

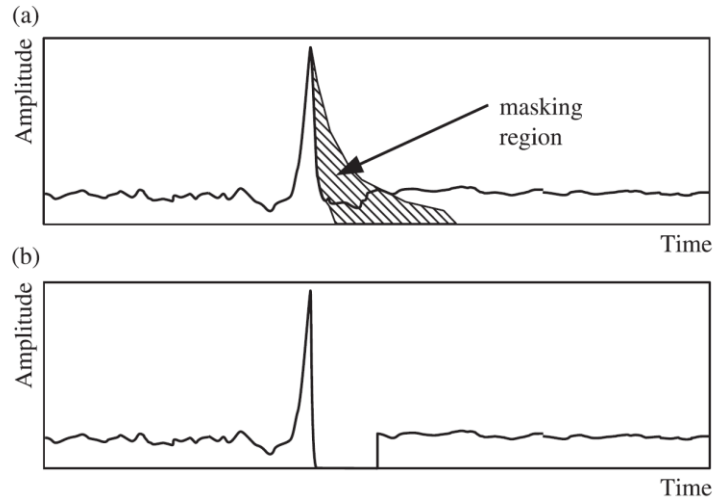


Figure 5.5 Illustration of post-stimulatory masking in the time domain detected (a) and inaudible information stripped out (b).

Applicazioni della psicoacustica

- Compressione di segnali audio ad alta fedeltà.
- Compressione della voce.
- Steganografia audio (data hiding o watermarking).
- Cancellazione attiva del rumore.
- Miglioramento dell'intellegibilità della voce.
- Sistemi di riconoscimento della voce.
- Machine hearing.
- Sintesi del suono e della voce.

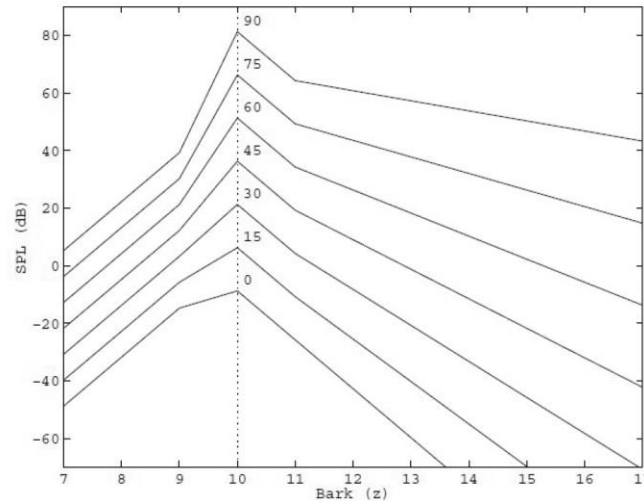
- La maggior parte di questi sistemi impiega:
 - Un modello che predice il mascheramento,
 - Una pre-enfasi basata su equal-loudness (curve isofoniche),
 - Metodi basati sulla discriminazione delle frequenze (scala MEL o Bark).

Applicazioni della psicoacustica

- Il mascheramento di frequenze simultanee è molto utilizzato in quanto facilmente modellabile.
- I modelli di mascheramento considerano l'effetto che hanno toni o rumori passabanda e, per quanto opinabile, applicano i risultati a suoni complessi.
- Si assume che i suoni complessi possano essere scomposti in un insieme di toni o rumori passabanda, ognuno dei quali produce un effetto di mascheramento, con il mascheramento complessivo dato dalla somma dei contributi separati.
- Il mascheramento reale dipende dall'esatta ampiezza del suono che giunge all'orecchio. Questa informazione in genere non è disponibile. Si fa riferimento alle ampiezze relative dei suoni e si assume un livello approssimato d'ascolto.

Modello a bande critiche

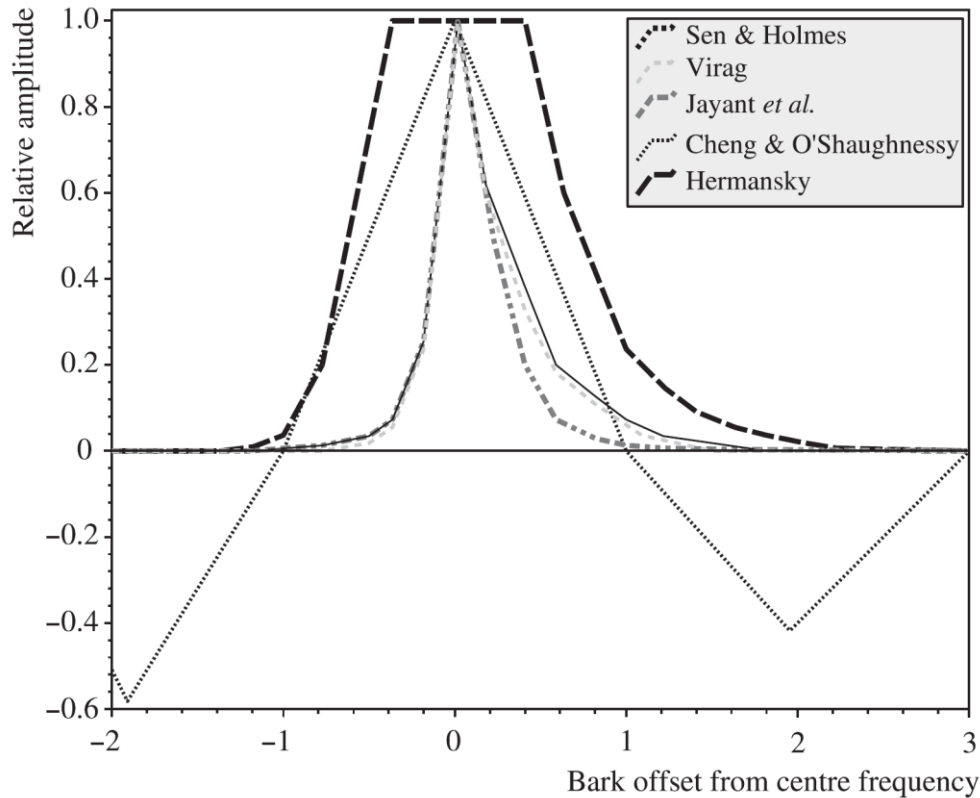
- Il suono in una banda critica influenza pure le bande critiche vicine. Dato un forte suono a una certa frequenza la sua influenza tende a diminuire con la distanza in frequenza. Ciascun suono influenza gli altri suoni in relazione alla loro ampiezza e frequenza.



(mascheramento con toni nell'MPEG)

Prototype spreading functions at $z=10$ as a function of masker level

Critical band spreading functions



Frequenza centrale di
1 kHz e 70 dB SPL

“It is a rather crude approximation of what is known about the shape of auditory filter»
Hermansky

Critical band spreading functions

Listing 5.1 spread_hz.m

```
1 function band=spread_hz(hz_array,hz_c)
2     %hz_array is an array of Hz frequencies
3     %hz_c is the current centre frequency in Hz
4     barkc=f2bark(hz_c);
5     band=zeros(size(hz_array));
6
7     for hi=1:length(hz_array)
8         barki=f2bark(hz_array(hi));
9         barkd=barki-barkc; ← barkd=barkc-barki;
10        if barkd >= -2.5 & barkd <=-0.5
11            band(hi)=10^(1*(barkd+0.5));
12        elseif barkd > -0.5 & barkd <0.5
13            band(hi)=1;
14        elseif barkd >= 0.5 & barkd <=1.3
15            band(hi)=10^-(2.5*(barkd-0.5));
16        end
17    end
```


Critical band spreading functions

```
F = [0:512] * 4000 / 512;  
spread = spread_hz(F, 1200);  
plot(F, spread, 'o-')
```

Mel-based spreading functions

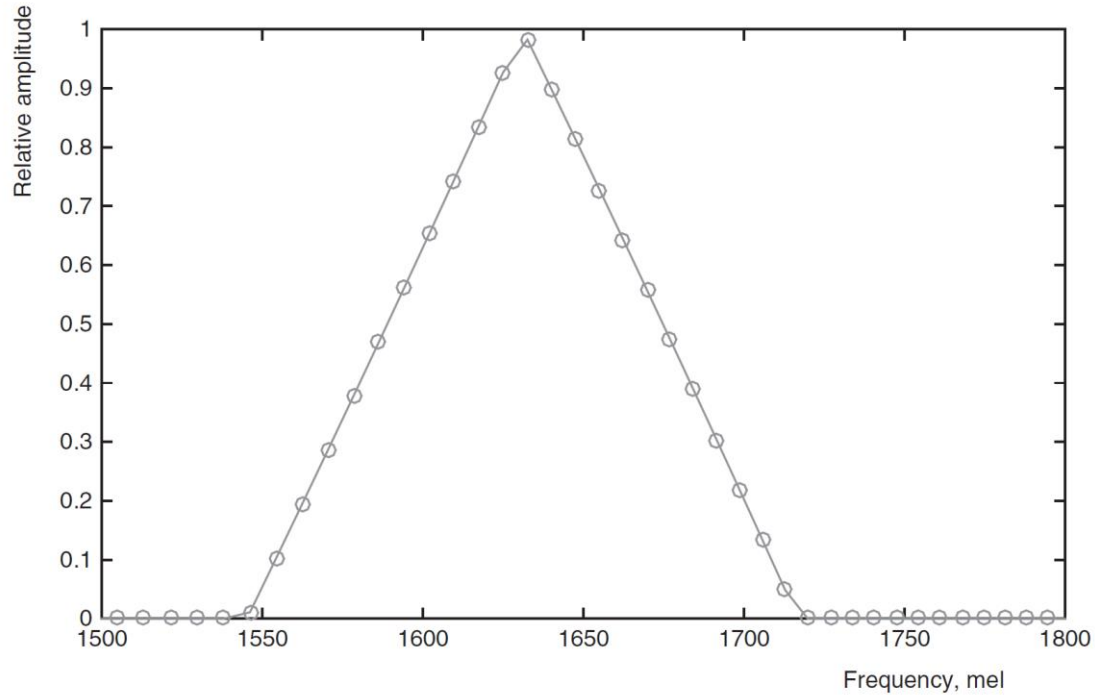


Figure 5.7 The triangular mel-based spreading function.

Mel-based spreading functions

Listing 5.2 spread_mel.m

```
1 function band=spread_mel(hz_points,hz_c,hz_size,hz_max)
2 %hz_array is an array spaced in Hz
3 %hz_c is the current index
4 band=zeros(1, hz_size);
5 hz1=hz_points(max(1,hz_c-1));           %start
6 hz2=hz_points(hz_c);                   %middle
7 hz3=hz_points(min(length(hz_points),hz_c+1)); %end
8 %-----
9 for hi=1:hz_size
10     hz=hi*hz_max/hz_size;
11     if hz > hz3
12         band(hi)=0;
13     elseif hz>=hz2
14         band(hi)=(hz3-hz)/(hz3-hz2);
15     elseif hz>=hz1
16         band(hi)=(hz-hz1)/(hz2-hz1);
17     else
18         band(hi)=0;
19     end
20 end
```

Mel-based spreading functions

```
mmax=f2mel(4000);  
melarray=[0:mmax/255:mmax]; %256 elements  
hzarray=mel2f(melarray);  
[idx idx]=min(abs(hzarray-1200));  
spread=spread_mel(hzarray,idx,100,4000);  
plot(f2mel([1:100]*4000/100),spread,'o-')
```

```
mmax=f2mel(4000);  
melarray=[0:mmax/24:mmax];  
hzarray=mel2f(melarray);  
[idx idx]=min(abs(hzarray-1200));  
spread=spread_mel(hzarray,idx,256,4000);  
plot(f2mel([1:256]*4000/256),spread,'o-')
```

Critical band filter-banks

- Le spreading functions viste sono usate per modellare l'effetto di una singola banda critica. Normalmente sono usate per consentire la rappresentazione di suoni in modo percettivamente rilevante, mediante un vettore corrispondente all'eccitazione delle singole bande critiche.
- L'eccitazione viene calcolata a partire dallo spettro del segnale (dalla DFT) calcolando il contributo a ciascuna banda critica.
- Si lavora o nel dominio Bark (ad es. con funzione di Hermansky) o nel dominio Mel con funzioni triangolari.

Critical band filter-banks

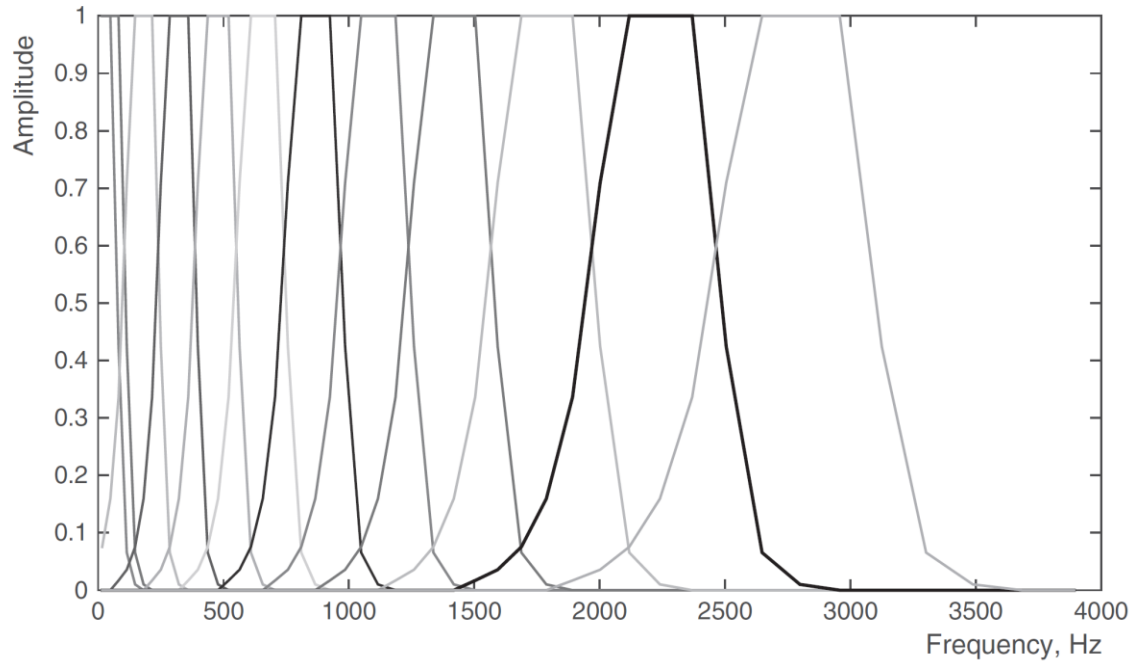


Figure 5.8 A plot of 12 critical-band spreading functions, based on Equation (5.1), over a 48-element array.

Critical band filter-banks

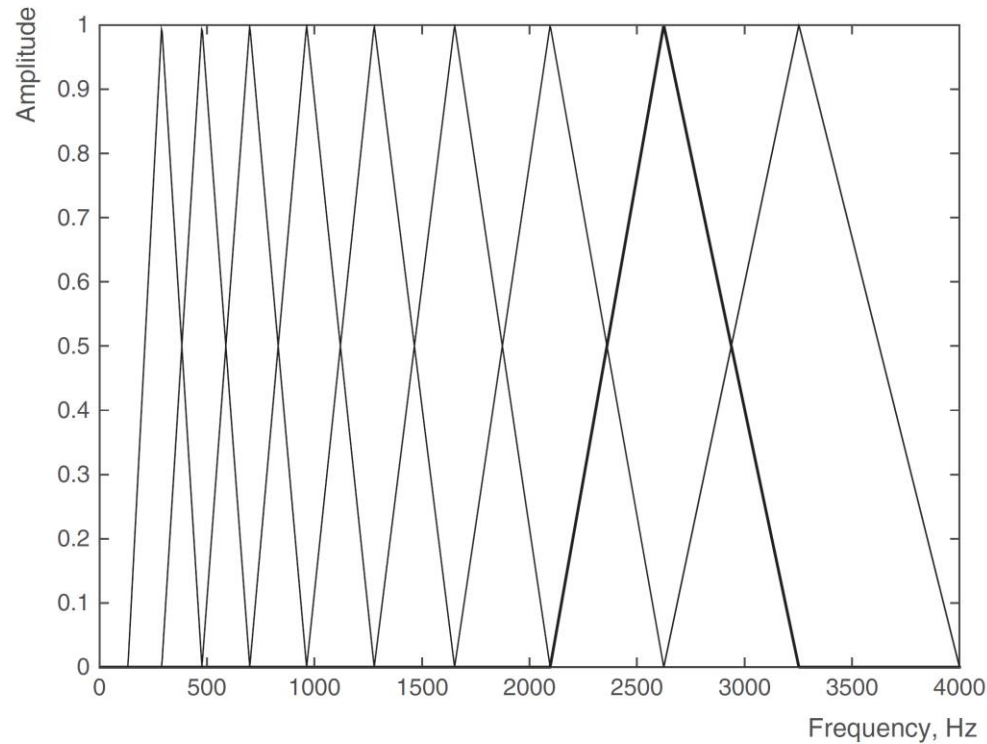


Figure 5.9 A plot of nine mel frequency critical-band spreading functions.

Critical band filter-banks

- Consideriamo l'implementazione di un modello psicoacustico nell'ipotesi che si voglia esaminare un segnale audio per determinare la soglia di mascheramento dovuta ai suoni presenti nel segnale stesso.
- Tre passi:
 - Analisi spettrale
 - Warping a bande critiche e convoluzione con funzione di spreading.
 - Conversione / pre-enfasi equal-loudness (es. pesatura con legge A).

Spectral analysis

```
S=fft(seg.*hamming(1,256));  
S=S(1:128);
```

Equal-loudness pre-emphasis

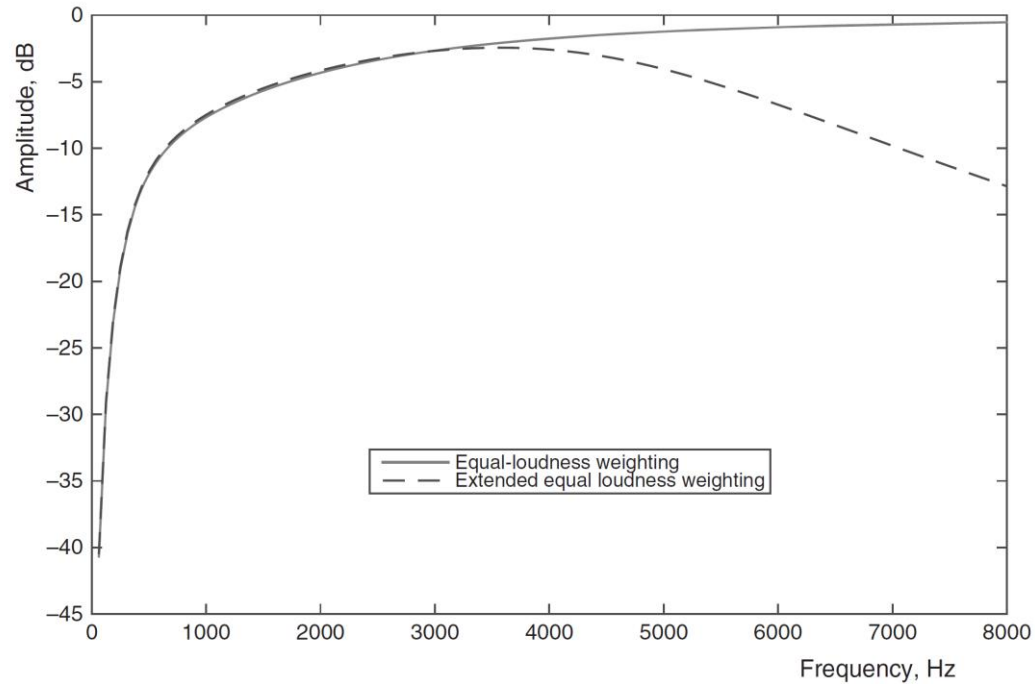


Figure 5.10 The shape of the equal-loudness emphasis curves of Equations (5.5) and (5.4).

Equal-loudness pre-emphasis

$$E(\omega) = \frac{\omega^4(\omega^2 + 56.8 \times 10^6)}{(\omega^2 + 6.3 \times 10^6)^2(\omega^2 + 0.38 \times 10^9)}. \quad <4 \text{ kHz}$$

$$E(\omega) = \frac{\omega^4(\omega^2 + 56.8 \times 10^6)10^{27}}{(\omega^2 + 6.3 \times 10^6)^2(\omega^2 + 0.38 \times 10^9)(\omega^6 + 9.58 \times 10^{26})}. \quad >4 \text{ kHz}$$

Listing 5.3 equal_loudness_preemph.m

```
1 function eql=equal_loudness_preemph(hz,varargin)
2 %hz is frequency in Hz
3 %if second argument given, use extended form
4 w=hz*2*pi; %convert to abs angular freq.
5 w2=w.^2; %squared version
6 if ~isempty(varargin)
7     eql=1e27*((568e5+w2).*(w2.^2))./(((63e5+w2).^2)
8         .* (38e7+w2).*(w.^6)+958e24));
9 else
10    eql=((568e5+w2).*(w2.^2))./(((63e5+w2).^2).*(38e7+w
11        .^2));
12 end
```

Hermansky-style model

Alla normalizzazione con legge A si fa in genere seguire un' enfasi secondo legge di potenza

$$\mathcal{P}(\omega) = \{E(\omega)\Theta(\omega)\}^{0.33}.$$

Listing 5.4 percept_model.m

```
1 function [x, xf]=percept_model(seg, N, Fs)
2 % Map audio frame seg to N-point 0:(Fs/2)Hz percept
   model. N typically 40 at Fs = 8 kHz
3 b_low=0; %Bark span lower limit
4 b_top=f2bark(Fs/2); %Bark span upper limit
5 bdiv=(b_top-b_low)/N; %Bark resolution
6 %Define an array of centre frequencies
7 xb=b_low+bdiv/2:bdiv:b_top-bdiv/2;
8 xf=bark2f(xb); %Convert to Hz
9 S=abs(fft(seg));
10 S=S(1:length(S)/2);
11 F=[1:length(S)]*(Fs/2)/length(S);
12 x=zeros(1,N);
13 for xi=1:N
14     bark=xb(xi);
15     hz=bark2f(bark);
16     %compute spreading function
17     spr=spread_hz(F,hz);
18     %compute summed influence
19     x_sum=sum(spr.*S')*equal_loudness_preemph(hz);
20     %intensity loudness power law
21     x(xi)=(x_sum)^0.33;
22 end
```

MFCC model

Listing 5.6 mfcc_model.m

```
1 function cc=mfcc_model(seg, N, M, Fs)
2 % Do FFT of audio frame seg, map to M MFCCs
3 % from 0 Hz to Fs/2 Hz, using N filterbanks
4 % typical values N=26,M=12,Fs=8000,seg~20ms
5 m_low=0; %mel span lower limit
6 m_top=f2mel(Fs/2); %mel span upper limit
7 mdiv=(m_top-m_low)/(N-1); %mel resolution
8 %Define an array of centre frequencies
9 xm=m_low:mdiv:m_top;
10 %Convert this to Hz frequencies
11 xf=mel2f(xm);
12 %Quantise to the FFT resolution
13 xq = floor((length(seg)/2 + 1)*xf/(Fs/2));
14 %Take the FFT of the speech...
15 S=fft(seg);
16 S=abs(2*(S.*S)/length(S));
17 S=S(1:length(S)/2);
18 F=[1:length(S)]*(Fs/2)/length(S);
19 %Compute the mel filterbanks.m
20 x1=zeros(1,N);
21 for xi=1:N
22     band=spread_mel(xf,xi,length(S),Fs/2);
23     x1(xi)=sum(band.*S');
24 end
25 x=log(x1);
26 %Convert to MFCC using loop (could use matrix)
27 cc=zeros(1,M);
28 for xc=1:M
29     cc(xc)=sqrt(2/N)*sum(x.*cos(pi*xc*( [1:N]-0.5)/N));
30 end
```

Mascheramento di rumori

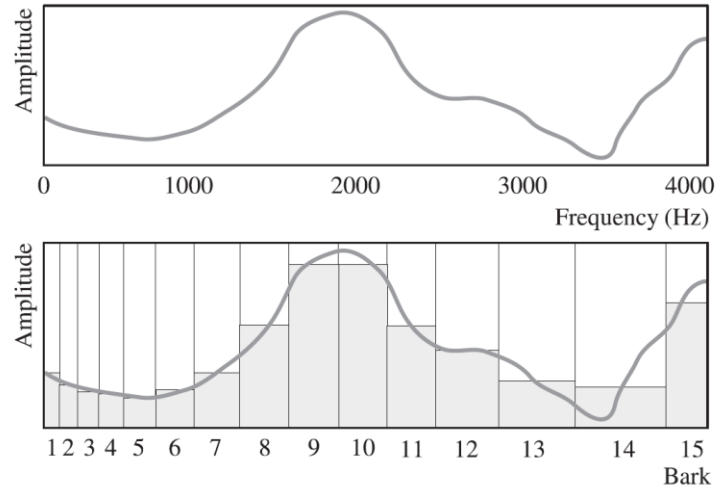


Figure 5.13 In-band masking level derived and plotted (bottom) from an example spectrum (top) for 15 critical bands of constant Bark width.

Vedere:

- Ian Vince McLoughlin, “Speech and Audio Processing”- Cambridge Univesity Press (2016)
 - Cap. 5