

COMPRESIA AUDIO.

- **Semnalul necomprimat:**
 - Frecvența maximă percepută de urechea umană este de aprox. 20kHz;
 - frecvența de eșantionare este de minim 40 kHz;
 - cuantizare cu 16 biti/eșantion;
 - pentru un semnal stereo calitate CD (eșantionat cu 44,1 kHz) rezultă o rată de transmisiune pentru semnalul necomprimat de 1.4 Mbps.
- **Metodele de compresie fără pierderi** (Huffman, LZW, etc.) în general nu funcționează bine pentru compresia audio.
- **Metode de compresie cu pierderi:**
 - **Silence Compression**
 - detectează zonele de “liniște”, asemănătoare cu codarea run-length;
 - **Adaptive Differential Pulse Code Modulation (ADPCM)**
 - în CCITT G.721 -- 16 sau 32 kbiți/sec.
 - codează diferența între două eșantioane consecutive;
 - adaptează pasul de cuantizare așa încât să se micșoreze varianța (puterea) zgomotului de cuantizare.
 - se obține o compresie de aproximativ 4:1.

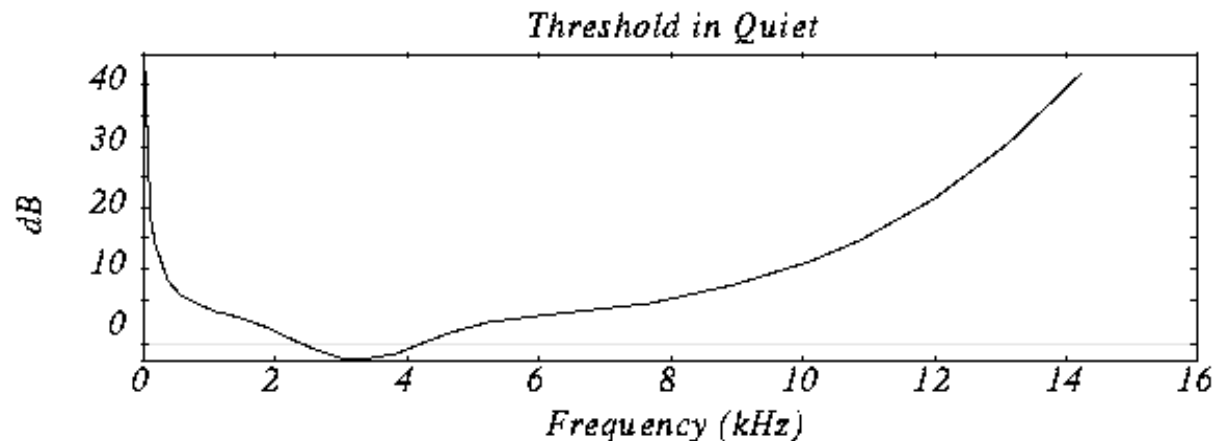
- **Linear Predictive Coding (LPC)**
 - se transmit, conform modelului vorbirii, parametrii de model ai corzilor vocale, laringelui, cavității bucale.
 - sună ca și vorbirea sintetizată pe calculator.
 - rată de 2.4kbiti/sec.
- **Code Excited Linear Predictor (CELP)**
 - efectuează LPC, dar transmite și termenul de eroare
 - calitate de audio-conferință la o rată de 4,8 kbiți/sec.
- Codarea audio poate fi făcută în:
 - **TIMP**
 - complexitate redusă;
 - necesită mai mult de 10 biți/eșantion pentru păstrarea calității;
 - **FRECVENȚĂ**
 - se poate obține o calitate înaltă cu numai 3 biți/eșantion;
 - se utilizează codarea în subbenzi și prin transformări;
- Pentru obținerea unor rate de compresie mari toate metodele de codare se bazează pe percepția audio umană (**PSIHOACUSTICĂ**).

Auzul si vocea umană

- Domeniul audibil este între 20 Hz și 20 kHz, cel mai sensibil la frecvențe de la 2 la 4 kHz.
- Dinamica auzului (încet la tare) e de aproximativ 96 dB.
- Vocea are domeniul normal de frecvență între 500 Hz și 2 kHz
- Fonemele sonore (m, v, l) au frecvențe joase.
- Fonemele insonore (f, s) au frecvențe înalte.

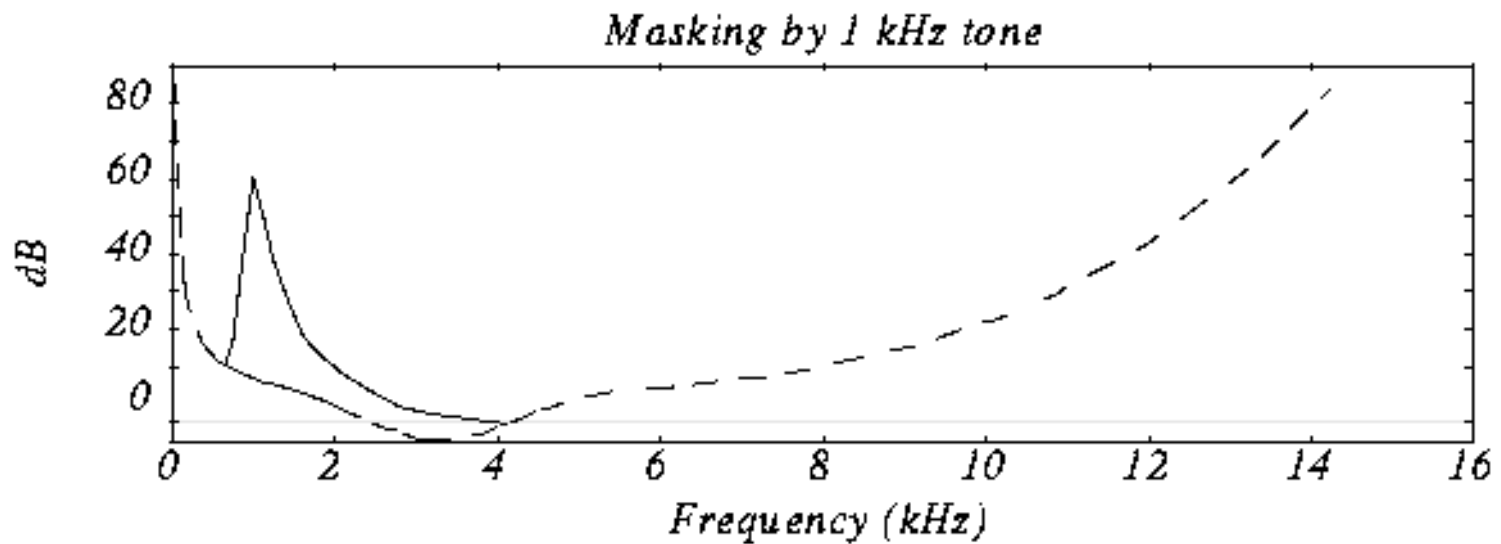
Cât de sensibilă este urechea umana?

- Experiment: O persoană ascultă un semnal de 1 kHz într-o cameră liniștită. Se reduce nivelul semnalului până când acesta nu se mai aude. La fel se reprezintă pentru toată gama de frecvențe audio și rezultă **curba de mascare în liniște**:

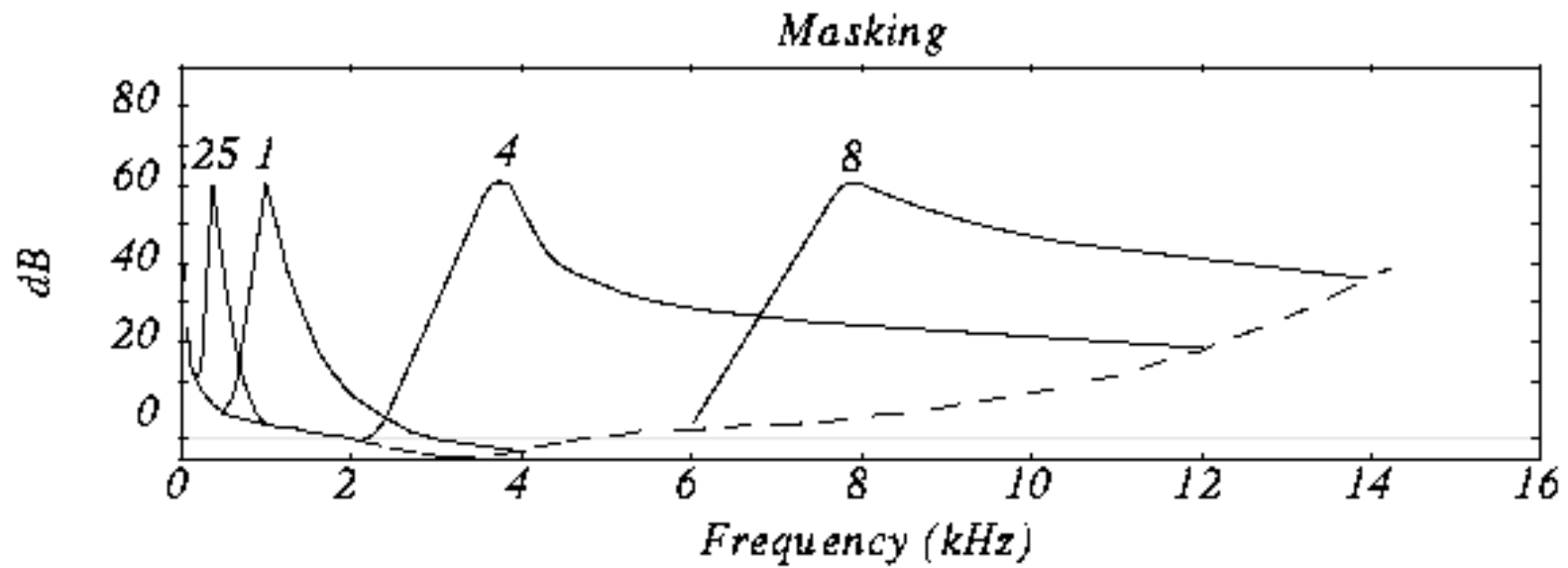


Mascarea în frecvență

- Experiment: Se asculta un ton de 1 kHz (**ton de mascare**) la un nivel fixat (60 dB).
- Se asculta un **ton de test** cu nivel variabil până când acesta începe să se audă.
- Se variază frecvența semnalului de test în jurul lui 1 kHz.



- Se repetă experimentul pentru mai multe frecvențe ale **tonului de mascare** obținându-se **curbele pragului de mascare în frecvență**.

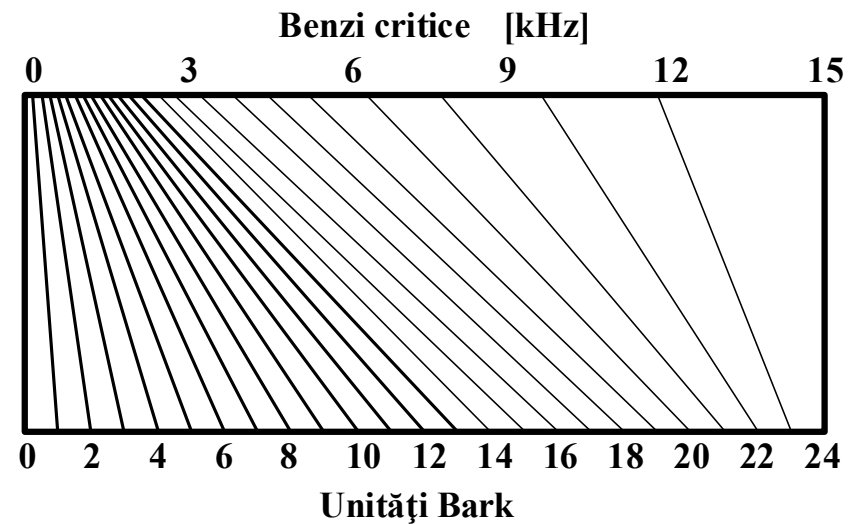


Benzi critice

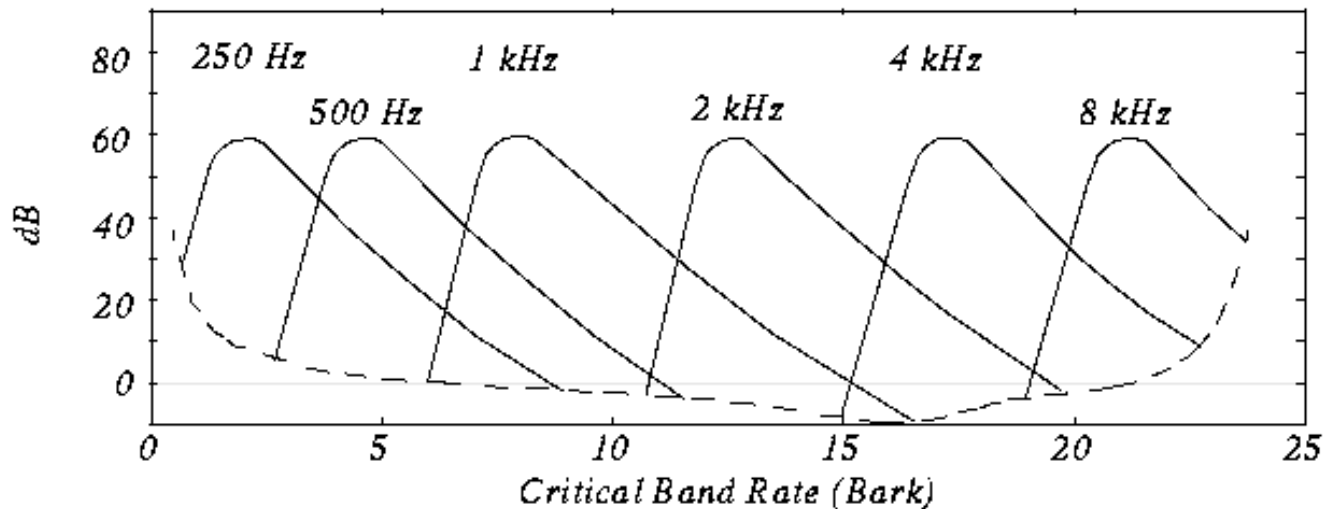
- Măsură uniformă de percepție a frecvenței neproportională cu lățimea curbei de mascare.
- Aproximativ 100 Hz pentru frecvențe de mascare <500 Hz, crește din ce în ce mai mult peste 500 Hz.
- Lățimea benzii se numește mărimea benzii critice.

Bark

- O altă unitate de măsură pentru frecvență (după Barkhausen).
- 1 Bark = lățimea unei benzi critice.
- Pentru frecvențe < 500 Hz, $f/100$
- Pentru frecvențe > 500 Hz, $9+4 \cdot \log_2(f/1000)$

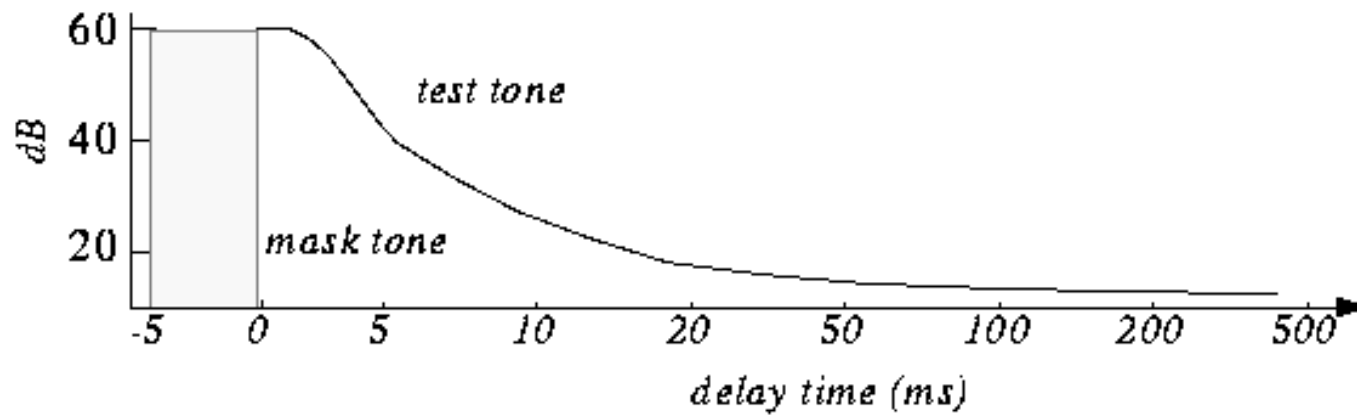


- Pragurile de mascare reprezentate în funcție de banda critică:

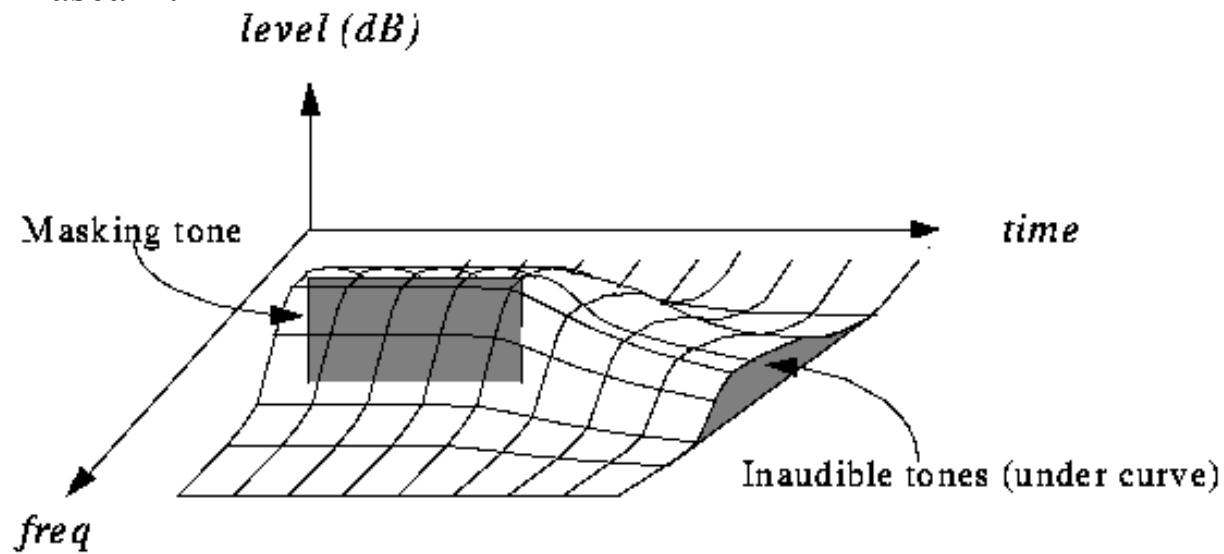


Mascare temporală

- Dacă ascultăm un sunet puternic, apoi acesta se oprește, trebuie să treacă un timp scurt pentru ca să putem auzi un sunet slab în apropiere.
 - Experiment: Se ascultă un **ton de mascare** de 1 kHz, 60 dB și un **ton de test** de 1.1kHz, 40 dB. Tonul de test nu se poate auzi (e mascat).
 - Se oprește **tonul de mascare**, apoi, după o scurtă întârziere, se oprește tonul de test.
 - Se ajustează întârzierea la durată minimă la care tonul de test mai poate fi auzit (aprox. 5 ms).
 - Se repetă cu niveluri diferite ale tonului de test.



- Se încearcă alte frecvențe pentru tonul de test (durata tonului de mascare rămâne constantă).
- Efectul total al mascării:

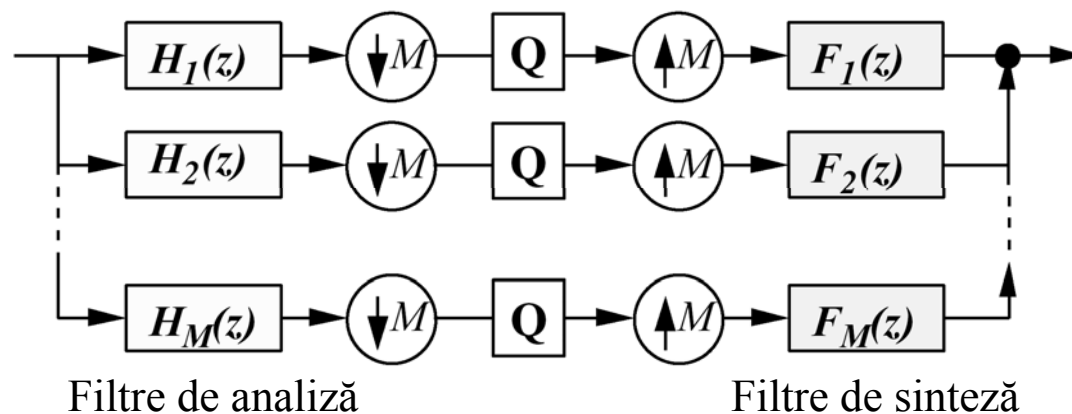


Concluzii:

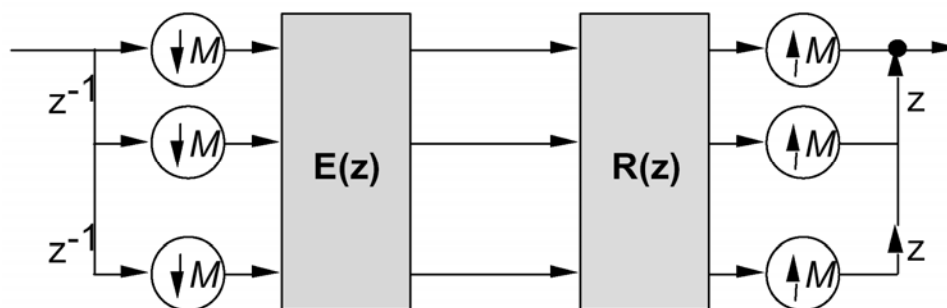
- Dacă avem un ton puternic (de exemplu 1 kHz), atunci tonurile mai slabe, apropiate în frecvență și timp, sunt mascate.
- Comparațiile se fac pe scala benzilor critice (dimensiunea mascării e aproximativ o bandă critică).
- Există doi factori de mascare: mascare în frecvență și mascare temporală.

Cum se poate folosi mascarea în compresia audio?

- Un semnal mascat de altul mai puternic este comparabil cu zgomotul de cuantizare.
- Funcția de mascare oferă distorsiunea maximă acceptabilă pentru fiecare bandă critică.
- Codorul determină mascarea din fiecare bandă cauzată de semnalele din benzile apropiate.
- Dacă puterea în bandă este sub pragul de mascare aceasta nu se codează.
- Altfel, se determină numărul de biți necesari pentru cuantizarea fiecărui coeficient astfel încât zgomotul introdus de cuantizare este sub pragul de mascare. (1 bit de cuantizare introduce 6 dB zgomot).

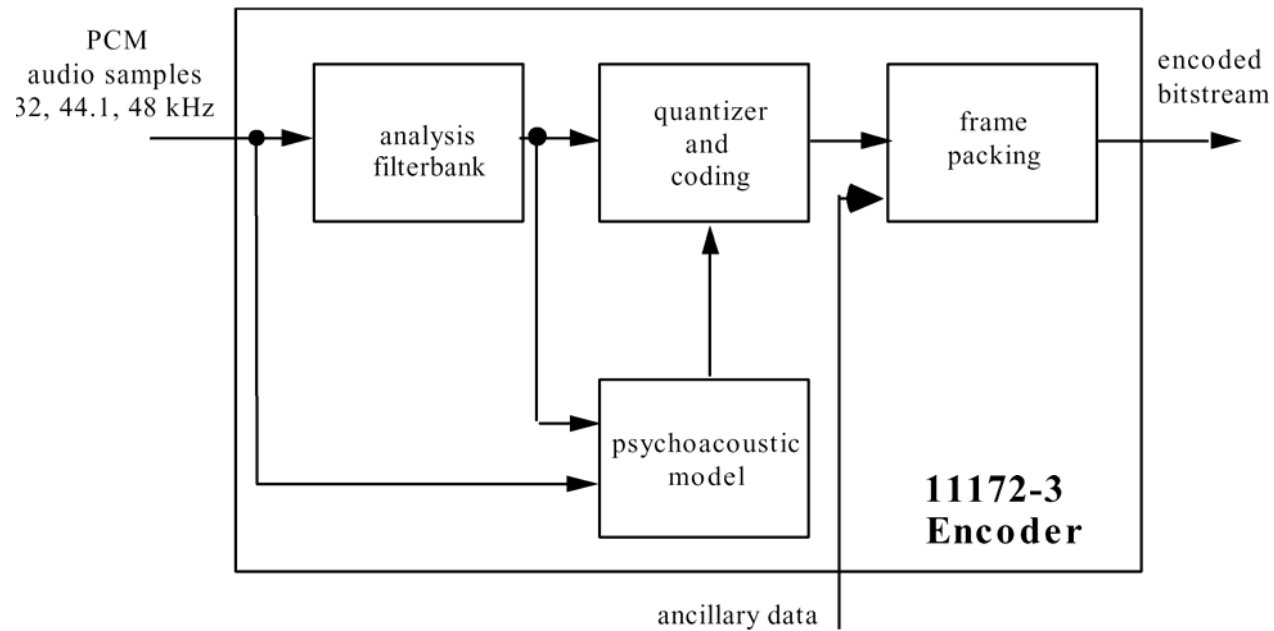


- Benzile de frecvență percepute de ureche nu sunt uniforme ci logaritmice.
- Bancul de filtre de analiză ar trebui să aproximeze benzile critice.
- Minimizarea ratei de biți în limitele date de mascare conduce la o compresie audio optimă.
- Se poate folosi pentru analiza în subbenzi transformata cosinus dacă $E(z)$ este matricea DCT și $R(z)$ este matricea IDCT.



Codarea MPEG-1 audio

- Standardul ISO/IEC 11172-3 elaborat între 1988 și 1991.
- Este primul standard de compresie audio de înaltă calitate.
- Codează semnale audio cu frecvențele de eșantionare de 32, 44.1 și 48kHz.
- Rata de bit comprimată pentru un semnal de calitate CD-audio stereo este între 64kbiți/s și 256kbiți/s față de 1.4Mbiți/s necomprimat.
- Schema bloc a codorului:

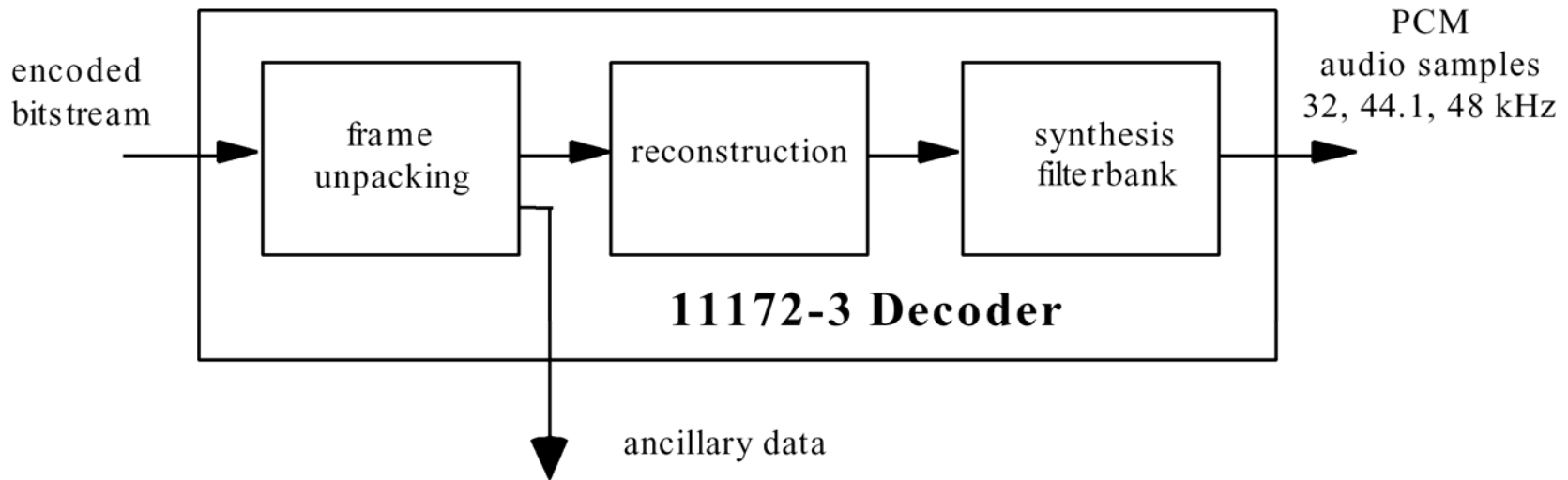


- Codorul analizează componentele spectrale ale semnalului audio cu ajutorul bancului de filtre sau a transformării MDCT (layer 3).
- Aplica un model psihoacustic pentru a estima nivelul minim de zgomot.
- Se furnizează SMR (Signal-to-Mask Ratio) pentru alocarea biților sau a zgomotului.
- Se formează fluxul de biți după cum urmează:

Header 32 biți	CRC 16 biți	Audio Data	Anciliary Data
--------------------------	-----------------------	-----------------------	---------------------------

- Header-ul
 - Syncword (12 biți)
 - Layer code (2 biți) reprezentând layererele I, II și III
 - Bit-rate index (4 biți) indexul debitului utilizat (diferă pentru fiecare layer în parte)
 - Frecvența de eșantionare (2 biți) poate fi 48, 44.1 și 32kHz
 - Padding bit
 - Mod (2 biți) stereo, joint stereo, unu sau două canale

- Schema bloc a decodorului



- Standardul MPEG audio include 3 layere diferite corespunzător diverselor aplicații, cu creșterea complexității codorului dar și a performanțelor (calitatea sunetului raportată la rata de bit).
- Layerele sunt compatibile în sensul ierarhic (layerul N poate decoda fluxul de date codate în layerul N și în toate layerele inferioare).
- Toate layerele au aceeași structură de bază.

- Layer 1
 - de la 32 kbps pînă la 448 kbps
 - rata de compresie 1:4
- Layer 2
 - de la 32 kbps pînă la 384 kbps
 - rata de compresie 1:6..8
- Layer 3
 - de la 32 kbps pînă la 320 kbps
 - rata de compresie 1:10..12

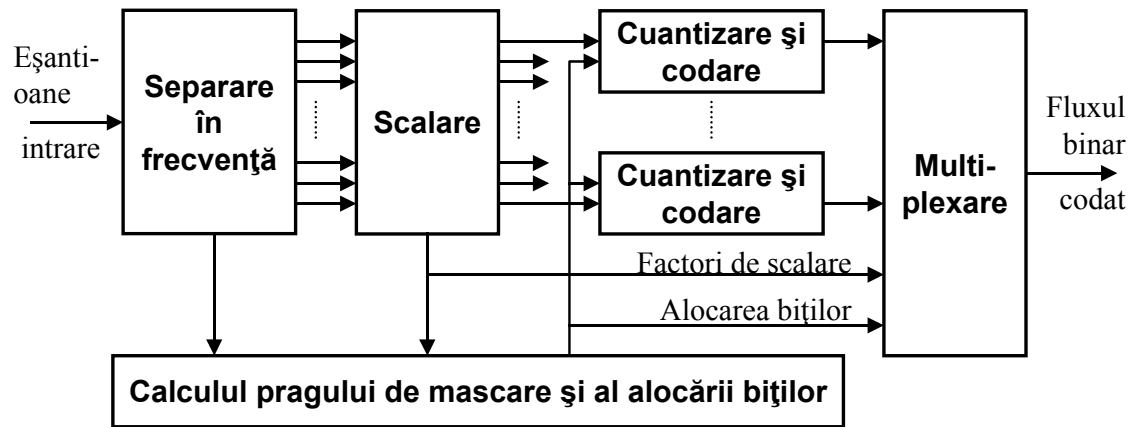
- Layer I - algoritmul de bază pentru codarea audio
 - Bancul de filtre împarte semnalul audio în 32 de subbenzi de frecvență cu lățime egală. Acestea nu corespund cu benzile critice.
 - Codorul calculează pragul de mascare folosind benzile critice.
 - Eroarea care apare la reconstrucție nu este mare.
 - Cadrul este format din 12 eşantioane/subbandă.
 - Conține un model psihoacustic pentru determinarea adaptivă a alocării biților și pentru cuantizare.
 - Domeniile de aplicație includ înregistrarea digitală pe bandă sau disc.

- Layer II - algoritm îmbunătățit față de layer I
 - Îmbunătățirea constă într-o codare suplimentară a alocării biților, a factorilor de scalare și o structură diferită a cadrului.
 - Codorul formează 3 blocuri cu 12 eșantioane/bloc și 32 de subbenzi (1152 eșantioane).
 - Se transmite un tip de alocare a biților și maxim 3 factori de scalare pentru 3 blocuri (câte un factor de scalare pentru fiecare bloc).
 - Aplicații în studiourile profesionale (radiodifuziune, înregistrări), telecomunicații, multimedia etc.
- Layer III - cea mai bună compresie
 - Crește complexitatea codorului/decodorului.
 - Conține un banc de filtre hibrid (filtre plus MDCT- modified discrete cosine transform).
 - Se obține o rezoluție mai bună în frecvență prin utilizarea MDCT.
 - Două lungimi ale blocului MDCT: 36 eșantioane și 12 eșantioane.
 - Aplicații în telecomunicații pe canale de banda îngustă ISDN, mp3 și alte aplicații cu debit foarte redus.

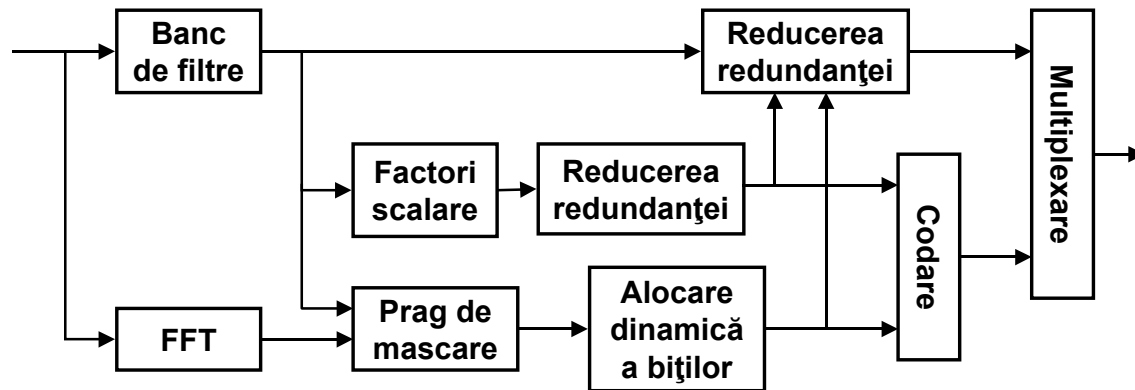
Algoritmi propuși:

- ASPEC (Audio Spectral Perceptual Entropy Coding): codare cu transformate cu suprapunerea blocurilor;
- ATAC (Adaptive Transform Aliasing Cancellation): codare cu transformate fara suprapunerea blocurilor;

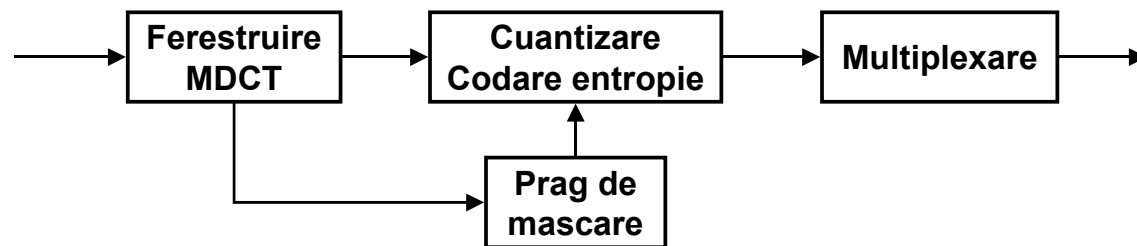
- SB/DPCM (Subband Coding and DPCM): codare pe subbenzi cu mai puțin de 8 subbenzi;
- MUSICAM (Masking-pattern Universal Subband Integrated Coding and Multiplexing): codare pe subbenzi cu mai mult de 8 subbenzi;



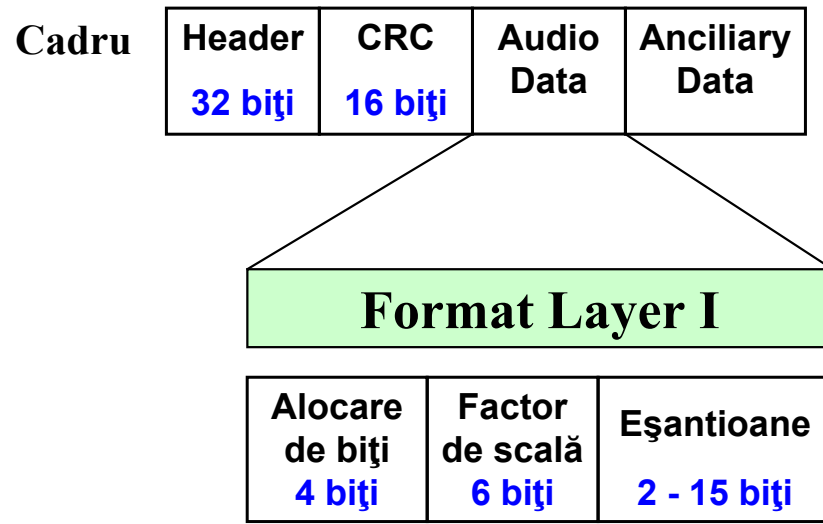
- Eșantioanele audio sunt mapate în frecvență printr-o transformare sau cu un banc de filtre.
- Coeficienții audio din domeniul de frecvență sunt normați cu un factor de scalare detreminat din pragul de mascare al răspunsului psihoacustic.
- Codorul MUSICAM



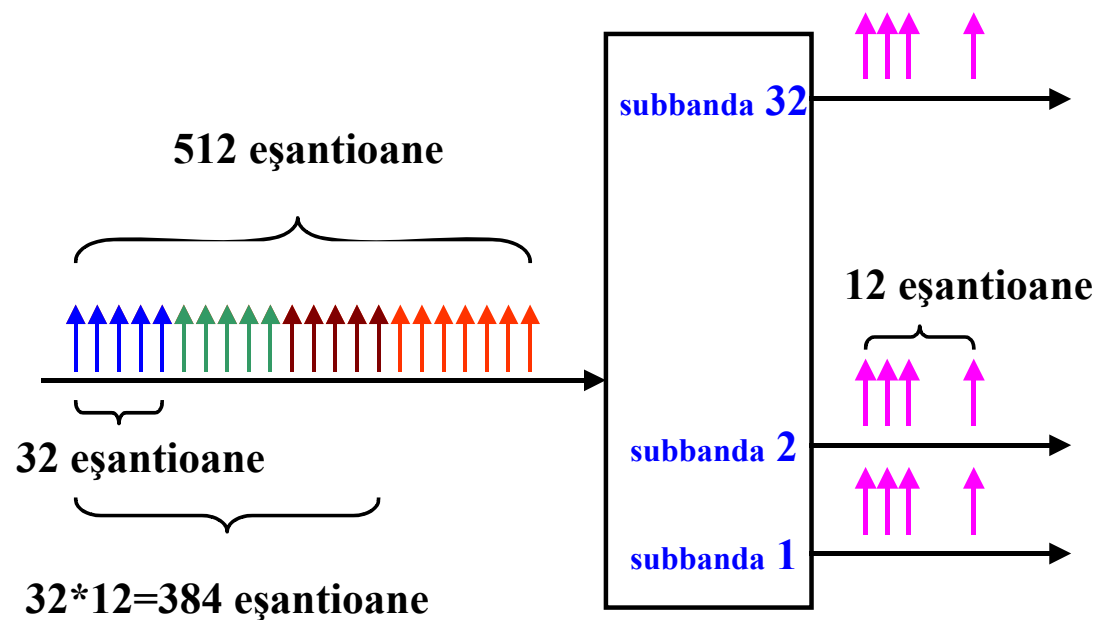
- Filtrele polifazice au complexitate de calcul redusă iar faza liniară permite reconstrucția perfectă.
- Un semnal eșantionat cu 48kHz este împărțit în 32 de subbenzi, fiecare bandă având lățimea de 0.75kHz.
- Semnalele de subbandă sunt împărțite în cadre digitale de 12 eșantioane succesive (8 ms).
- Intervalul de eșantionare în fiecare subbandă este de 2/3 ms.
- Pragul de mascare se calculează dintr-un estimat pe termen scurt al densității spectrale de putere prin medierea transformatei FFT.
- Calculul se repetă la fiecare 24 ms.
- Lățimea constantă a subbenzilor nu coincide cu benzile critice.
- După calculul puterii zgomotului de mascare, biții se alocă cuantizoarelor minimizându-se NMR.
- Factorii de scalare pot fi calculați folosind cuantizarea adaptivă așa încât eșantioanele să fie între [-1,1].
- Factorii de scalare au redundanță mare și pot fi codați, urmând a fi transmiși împreună cu informația de alocare a biților în fluxul de date.
- Codorul ASPEC



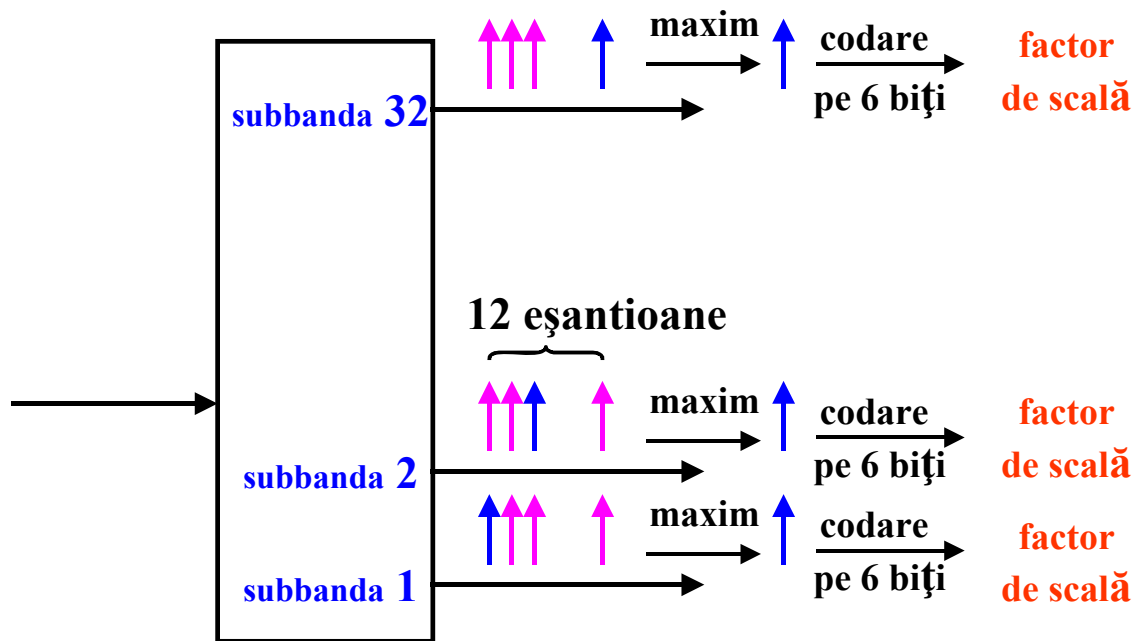
- Pentru separarea în frecvență se utilizează MDCT.
- Eșantioanele sunt ferestruite pentru limitarea alierii în domeniul timp.
- MDCT împreună cu subeșantionarea creează dintr-un bloc de $2N$ eșantioane, N coeficienți în frecvență.
- Calculul pragului de cuantizare:
 - Este calculată energia semnalului în domeniul frecvență (amplitudine și fază);
 - Se calculează energia în fiecare bandă critică. Această energie dă pragul neîmprăștiat.
 - Împrăștierea e calculată cu o funcție de împrăștiere.
 - În final se calculează entropia perceptuală utilizată pentru estimarea numărului de biți necesari pentru blocul curent.
- Datele cuantizate sunt codate cu cod Huffman.
- Factorii de scalare și alocarea biților sunt multiplexati în fluxul de date.
- **MPEG Layer I**
 - Filtrarea în subbenzi;
 - Modelare psihoacustică;
 - Scalare și alocarea biților;
 - Cuantizare și codare
 - Formarea fluxului de date



- **Filtrarea în subbenzi**
- Se folosește un banc de filtre pentru a transforma semnalul audio din domeniul timp în frecvență.
- Filtrele împart semnalul inițial în 32 de benzi de frecvență echidistante cu frecvență de eșantionare $F_s/32$.



- Pentru fiecare subbandă se calculează maximul (în modul) pentru fiecare set de 12 eșantioane.
- Factorul de scalare se alege dintr-un tabel și este valoarea imediat superioară maximului găsit.
- Se codează indexul factorului de scalare din tabel, pe 6 biți pentru fiecare subbandă.
- Acesta se transmite doar dacă a fost alocat benzii un număr nenul de biți.



- **Modelare psihoacustică**
- Layer I suporta atât modelul psihoacustic I cât și modelul psihoacustic II.
- Totuși, modelul psihoacustic I este suficient pentru Layer I, care implică un FFT de 512 elemente.
- SMR (signal-to-mask ratio) se determină din modelul psihoacustic folosit.
- **Modelul psihoacustic I**

- Calculul FFT în paralel cu filtrarea în subbenzi compensează lipsa de selectivitate a filtrelor în zona de joasă frecvență. FFT este de 512 eșantioane pentru layer I și de 1024 eșantioane pentru layer II.
- Se cunoaște pragul de mascare în liniște.
- Se extrag din spectrul de putere FFT componentele tonale și netonale deoarece ele influențează pragul de mascare în benzile critice.
- Componentele tonale sunt cele care respecta relațiile:

$$\begin{aligned} power_x(i-j) < power_x(i) - 7 \leq power_x(i+j) \quad j \in \{2,3,6\} \\ power_x(i-1) < power_x(i) \leq power_x(i+1) \end{aligned}$$

- Se elimina componentele vecine componentelor tonale.
- Se elimina componentele tonale și netonale care sunt sub pragul de mascare în liniște.
- Dacă mai multe componente tonale sunt la distanță mai mică de 0.5 Bark se păstrează maximul lor.
- Calculul pragului global de mascare (în dB):

$$LT_G(i) = 10 \log_{10} \left[10^{LT_q(i)/10} + \sum_{j=1}^m 10^{LT_{tm}(j,i)/10} + \sum_{j=1}^n 10^{LT_{nm}(j,i)/10} \right]$$

unde LT_q este pragul în liniște, iar LT_{tm} și LT_{nm} sunt pragurile de mascare datorate componentelor tonale și netonale.

- Pragul global de mascare minim din subbanda n se utilizează pentru determinarea raportului semnal-mascare (SMR):

$$SMR_{sb}(n) = L_{sb}(n) - LT_{\min}(n) \text{ dB}$$

unde $L_{sb}(n)$ este nivelul semnalului în subbanda n .

- Se calculează SMR pentru fiecare subbandă.
- **Modelul psihoacustic II**
- Dimensiunea FFT și a ferestrei Hann poate fi variată. Layer III calculează modelul de două ori în paralel cu FFT de 192 și de 576 esantioane (bloc scurt / lung).
- Se consideră o funcție de împrăștiere între benzile critice vecine bazată pe mascarea temporală (sunetele se “sting” în timp iar curba de mascare este influențată de sunetele precedente).
- Pragul audibil final se calculează prin convoluția energiei împrăștiate și a energiei parțiale inițiale.

- SMR e calculat ca raport între energia parțială e_{part} și nivelul zgomotului n_{part} :

$$SMR_n = 10 \log_{10} (e_{part_n} / n_{part_n})$$

- **Alocarea biților**

- Conceptul de bază în alocarea biților este minimizarea **MNR** din cadru cu constrângerea ca numărul total de biți utilizați să nu depășească numărul de biți disponibili în cadru B_f . B_f se calculează cu formula:

$$B_f = \frac{Bit \ rate}{f_s} \cdot 384 \text{ biti / cadru}$$

- Procedura de alocare de biți e iterativă și pornește din starea “zero bit allocation”.
- Întâi se calculează “mask-to-noise ratio” **MNR** care se obține cu formula:

$$\text{MNR} = \text{SNR} - \text{SMR} \text{ (dB)}$$

unde :

SNR se găsește în tabelul următor

SMR este furnizat de modelul psihoacustic.

Biți	Codul	Număr de niveluri	SNR (dB)
0	0000	0	0.00
2	0001	3	7.00
3	0010	7	16.00
4	0011	15	25.28
5	0100	31	31.59
6	0101	63	37.75
7	0110	127	43.84
8	0111	255	49.89
9	1000	511	55.93
10	1001	1023	61.96
11	1010	2047	67.98

12	1011	4095	74.01
13	1100	8191	80.03
14	1101	16383	86.05
15	1110	32767	92.01
invalid	1111	-	-

- **MNR** arată diferența dintre eroarea de cuantizare și măsurarea perceptuală.
- Eșantioanele audio pot fi comprimate de MNR ori.
- De aceea minimul MNR din fiecare subbandă e determinat la fiecare iterație.
- Procedura iterativă se repetă până când MNR e minimizat și numărul de biți folosiți pentru cele 4 componente se apropie de numărul de biți disponibili.
- Biții marginali calculați la fiecare iterație B_{mg} , pot fi calculați ca:

$$B_{mg} = B_{tav} - (b_{bal} + b_{scf} + b_{spl} + b_{anc})$$

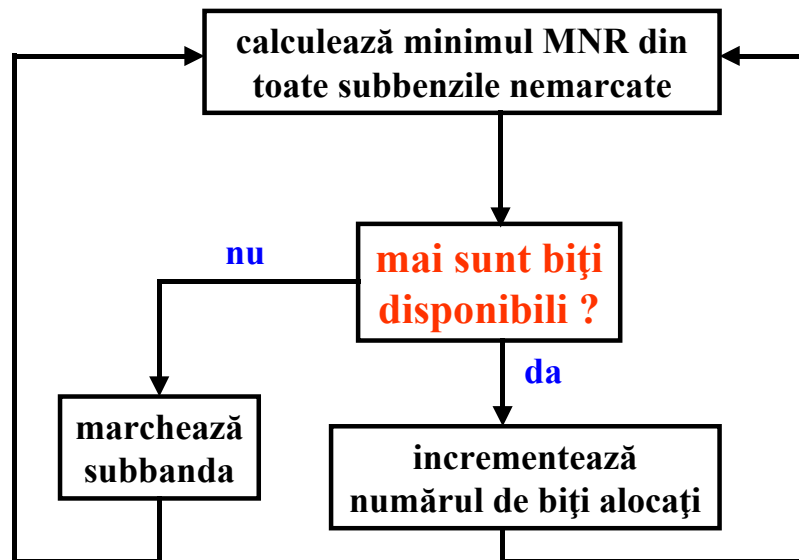
unde:

b_{bal} - numărul de biți de alocare

b_{scf} - numărul de biți pentru factorul de scală

b_{spl} - numărul de biți pentru eșantioane

b_{anc} - numărul de biți pentru “ancillary data”



- **Cuantizarea și codarea**

- Eșantioanele de subbandă sunt codate și cuantizate de un cuantizor uniform cu o reprezentare simetrică față de 0.
- Fiecare eșantion de subbandă S_i este normat la factorul de scală și cuantizat utilizând formula :

$$S_{qi} = \left(A \left(\frac{S_i}{scf} \right) + B \right) \Big|_N$$

- Coeficienții A și B sunt tabelați.

Număr de	A	B
----------	---	---

niveluri		
3	0.750000000	-0.250000000
7	0.875000000	-0.125000000
15	0.937500000	-0.062500000
31	0.968750000	-0.031250000
63	0.984375000	-0.015625000
127	0.992187500	-0.007812500
255	0.996093750	-0.003906250

- **Fluxul de biți**

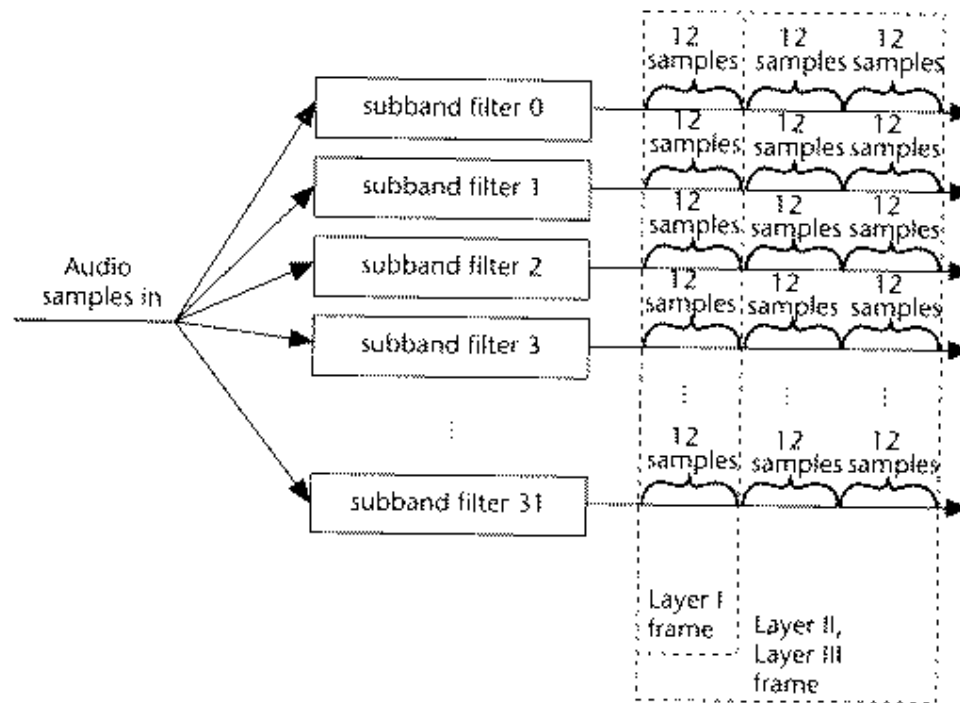
- Informația codată din subbenzi e multiplexată în cadre. Această operație nu presupune o codare suplimentară.
- Un cadru este compus dintr-un număr întreg de sloturi pentru a ajusta fluxul mediu de biți.
- În Layer I un slot are 32 de biți în timp ce în Layer II și III un slot are 8 biți.
- Numărul de sloturi dintr-un cadru se obține împărțind B_f la numărul de biți dintr-un slot.
- Dacă frecvența de eșantionare este 44.1 kHz numărul de sloturi nu este întreg. În asemenea cazuri cadrul trebuie ajustat prin adăugarea de biți (padding). Astfel numărul de sloturi dintr-un cadru poate fi N sau N+1.

Ex:

$F_s=44.1$ kHz, 114.84 cadre/sec, 1 cadru=8.70ms

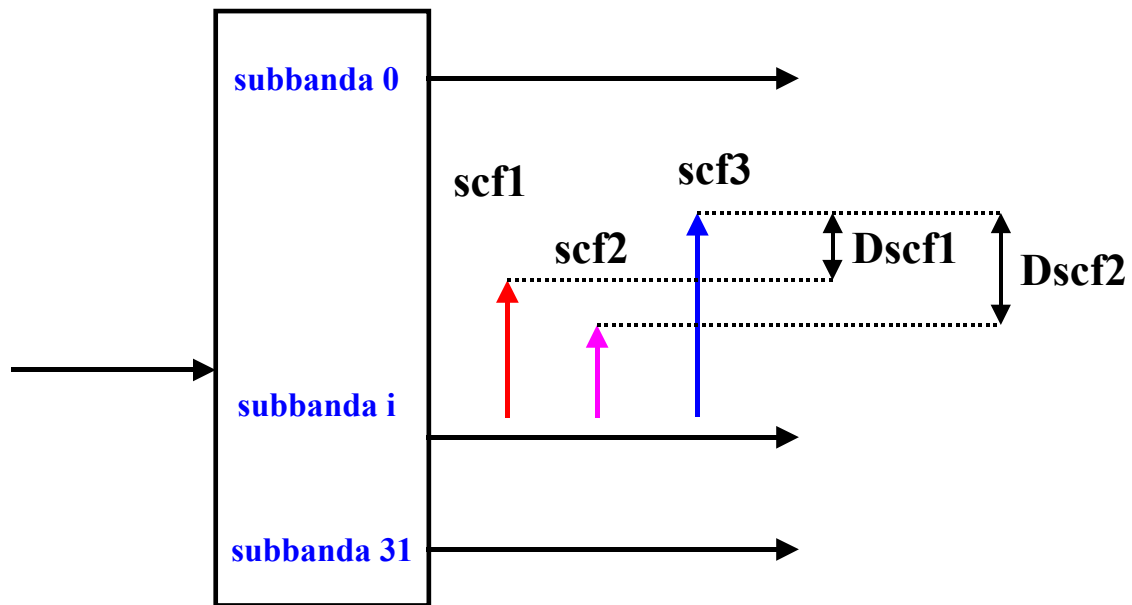
rezultă 17.41 sloturi => 18 sloturi

- **MPEG Layer II**
- Layer II urmărește în principiu aceleași reguli de codare și decodare ca și Layer I.
- Principala diferență este ca Layer II introduce corelație între subbenzi. Layer II conține informații pentru 1152 de esantioane (3 x 12 x 32 esantioane = 1152 de esantioane).
- În fluxul de date apare și un selector al factorului de scală.



- Layer II suportă atât modelul psihoacustic I cât și modelul psihoacustic II.

- Modelul psihoacustic I implică un FFT de 1024 esantioane iar modelul II 512 esantioane.
- SMR din fiecare subbandă se determină din modelul psihoacustic folosit.
- **Codarea factorilor de scalare**
- Se poate folosi aceeași analiză și sinteză a filtrelor ca în cazul Layer I.
- În Layer II un cadru conține 36 (3 x 12) esantioane de subbandă (12 granule) și 3 factori de scală pe subbandă.
- Cele două diferențe se obțin din cei trei factori de scală după cum urmează:
$$D_{scf1} = scf3 - scf1$$
$$D_{scf2} = scf3 - scf2$$



- Fiecare diferență este clasificată în una din cele 5 clase după cum urmează:

Clasa	Condiția
1	$Dscf_i \leq -3$
2	$-3 < Dscf_i < 0$
3	$Dscf_i = 0$
4	$0 < Dscf_i < 3$
5	$Dscf_i \geq 3$

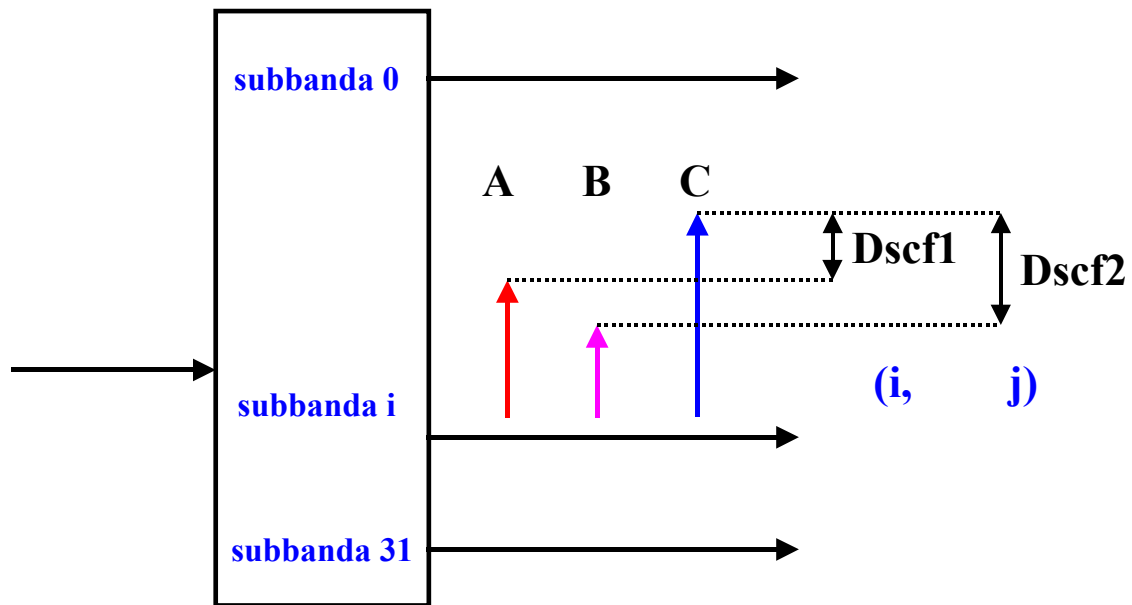
- Cele doua clase ce rezultă corespund transmisiei unui pattern (cei trei factori de scală care trebuie transmiși).
- Redundanța e redusă cu prețul codării informației de selectare a factorului de scală (2 biți).

(Clasa1,Clasa2)	Pattern transmis	Factor de scală selectat
(1,1), (1,5), (4,5), (5,1),(5,5)	123	0
(1,2), (1,3), (5,2),(5,3)	122	3
(1,4), (5,4)	133	3
(2,1), (2,5), (3,5)	113	1
(2,2), (2,3),(3,1),(3,2),(3,3)	111	2
(2,4)	444	2
(3,4), (4,4)	333	2
(4,1), (4,2), (4,3)	222	2

- Biții de selecție a factorului de scală reprezintă numărul și poziția factorilor de scală din fiecare subbandă.

scfsi	Factori de scală	Factor de scală
-------	------------------	-----------------

	codafi	decodat
0 (00)	3	scf1, scf2, scf3
1 (01)	2	primul □□scf1 și scf2 al doilea □ scf3
2 (10)	1	scf=scf1=scf2=scf 3
3 (11)	2	primul □□scf1 al doilea □ scf2 siscf3



Ex: Presupunem că 3 factori de scală A, B, C sunt obținuți într-o subbandă.

Clasa	Factori de scală transmiși	scfs i	Factori de scală decodați
(1,1)	ABC	00	ABC
(1,3)	AB	11	ABB
(3,2)	A	10	AAA

- **Alocarea biților**

- **SMR** din modelul psihoacustic este folosit pentru a obține **MNR** și operația iterativă este similară cu cea din Layer I, incluzând și câmpul de selecție a factorului de scală.

- **Cuantizarea și codarea:**

- Același algoritm folosit la cuantizare în Layer I se aplică și aici.
- Trei eșantioane succesive (1 granulă) sunt codate ca un singur cuvânt de cod.

- La decodare se va folosi următorul algoritm ($s(0)$, $s(1)$ și $s(2)$ sunt cele 3 eșantioane codate):

```

for i=0 to 2
s(i)=(code) MOD (număr de nivele)
code=(code) DIV (număr de nivele)

```

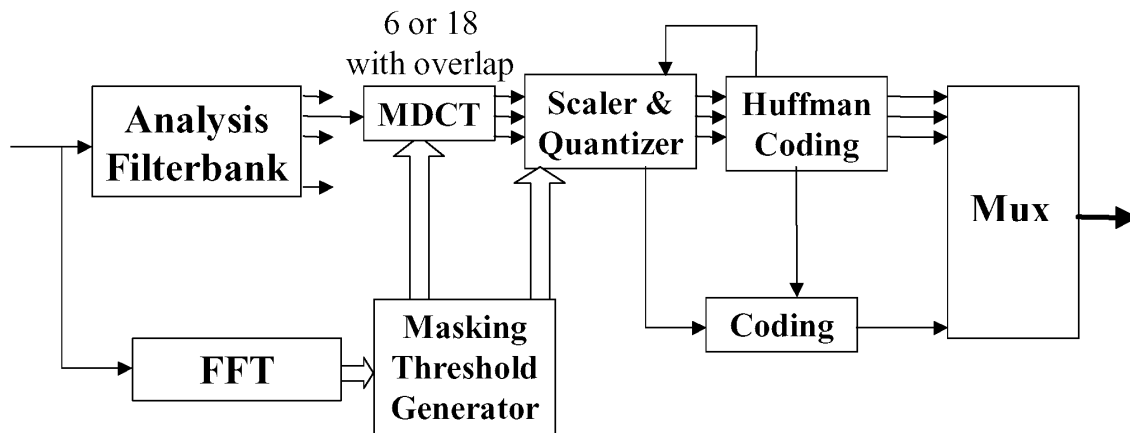
- Cele trei eșantioane sunt decuantizate după formula:

$$S(i)=C(S''(i)+D)$$

Unde C și D sunt constante tabelate.

- **Formarea fluxului de biți:**

- Aceleași operații ca și în Layer I sunt efectuate. Această operație nu presupune o codare suplimentară.
- În Layer II un slot are 8 biți.
- Aceiași algoritmi de padding se aplică și aici.
- **MPEG Layer III**
- Codarea în MPEG Layer III e mult mai sofisticată decât cea din Layer I/II.
- Cu ajutorul unui banc de filtre hibride se obține o mai bună rezoluție în frecvență.
- Filtrele hibride sunt obținute prin cascada filtrelor polifazice de analiză (folosite și în Layer I și II) cu operația MDCT (Modified DCT).
- Modelul perceptual combină calculul energiei cu FFT și cu bancul de filtre.
- Ieșirile modelului perceptual sunt valorile pragului de mascare echivalent cu valoarea acceptată a zgomotului în fiecare bandă.
- Benzile de frecvență sunt egale cu benzile critice.
- Cuantizarea nu mai este uniformă, se introduce codarea entropică, se introduc mai multe bucle pentru modelul psihoacustic și pentru alocarea de biți.
- Codarea Huffman se face în funcție de statistica semnalului muzical alegându-se tabelul de codare optim.



- Ferestrele definite pentru MDCT sunt pentru blocuri lungi și scurte.
- Pentru blocuri lungi (N=36) formula este:

$$h(k) = x(k) \sin\left(\frac{\pi}{N}\left(k + \frac{1}{2}\right)\right) \quad k=0, 1, \dots, 35, N=36$$

- Pentru blocuri scurte se aplică aceeași formulă doar că N=12.
- Comutarea între blocuri nu e instantanee. Pentru aceasta se definesc ferestre de tranziție (lung => scurt și scurt => lung).
- Decizia de comutare se ia din curba de mascare obținută din estimatul entropiei psihoacustice. Dacă valoarea entropiei psihoacustice (PE) depășește un anumit nivel (PE>1800) atunci se va trece la blocul scurt.

- **Transformarea Cosinus Modificată (MDCT)**

- Următoarea ecuație se folosește pentru a obține $N/2$ coeficienți S_i din N eșantioane de intrare x_k :

$$S_i = \sum_{k=0}^{N-1} x_k \cos\left(\frac{\pi}{2N}\left(2k+1+\frac{N}{2}\right)(2i+1)\right)$$

unde : $i=0, 1, \dots, \frac{N}{2}-1$

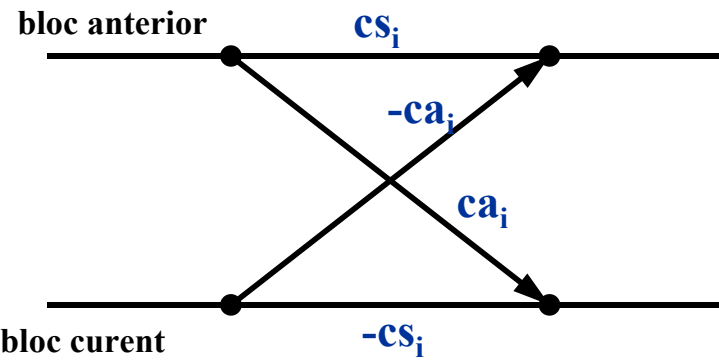
- N poate fi 12 pentru blocuri scurte și 36 pentru blocuri lungi.
- Transformarea MDCT inversă are expresia:

$$x_k = \sum_{i=0}^{\frac{N}{2}-1} S_i \cos\left(\frac{\pi}{2N}\left(2k+1+\frac{N}{2}\right)(2i+1)\right)$$

unde $k=0, 1, \dots, N-1$

- **Reducerea efectului de aliere**

- Calculul de reducere a alierii se face atât în codor cât și în decodor.
- Numai blocurilor lungi li se aplică această procedură.
- Transformarea MDCT dă 18 coeficienți din 36 de eșantioane de intrare. Între 2 seturi de 18 coeficienți se aplică un operator fluture ca în figura următoare.



unde $i=0, 1, \dots, 7$ iar c_{si} și c_{ai} se calculează cu formulele:

$$cs_i = \frac{1}{\sqrt{1+c_i^2}} \quad ca_i = \frac{c_i}{\sqrt{1+c_i^2}}$$

- Cei 8 coeficienți c_i sunt tabelați:

i	c_i
0	-0.6
1	-0.535
2	-0.33
3	-0.185
4	-0.095
5	-0.0041
6	-0.0142
7	-0.0037

- **Cuantizarea și codarea**

- Cuantizorul MPEG Layer III este neliniar. Legea de cuantizare este de forma:

$$Q\left(\alpha \cdot x^{\frac{3}{4}}\right)$$

- La decodare va trebui efectuată operația inversă adică ridicarea la puterea 4/3.
- Codorul Huffman este utilizat pentru codare entropică.
- Procesul de găsimă a câștigului și factorilor de scalare optimi pentru un bloc, rata de bit și ieșirea modelului perceptual este realizat în două cicluri iterative prin analiză-sinteză.
- Ciclul interior (ciclul de rată):
 - Codul Huffman alocă valorilor cuantizate mici (cele mai frecvente) cuvinte de cod de lungime minimă.
 - Dacă numărul de biți rezultat depășește numărul de biți disponibili pentru codarea unui bloc de date, aceasta se poate ajusta prin modificarea câștigului global care rezultă într-un pas de cuantizare mai mare, ceea ce conduce la valori cuantizate mai mici.
 - Operația este repetată cu diferiți pași de cuantizare până când cererea de biți pentru codarea Huffman este suficient de mică.
- Ciclul exterior (ciclul de control al zgomotului):
 - Pentru a dimensiona zgomotul de cuantizare în funcție de pragul de mascare, se aplică un factor de scalare fiecărei benzi.
 - Sistemul pornește cu un factor de scalare 1.

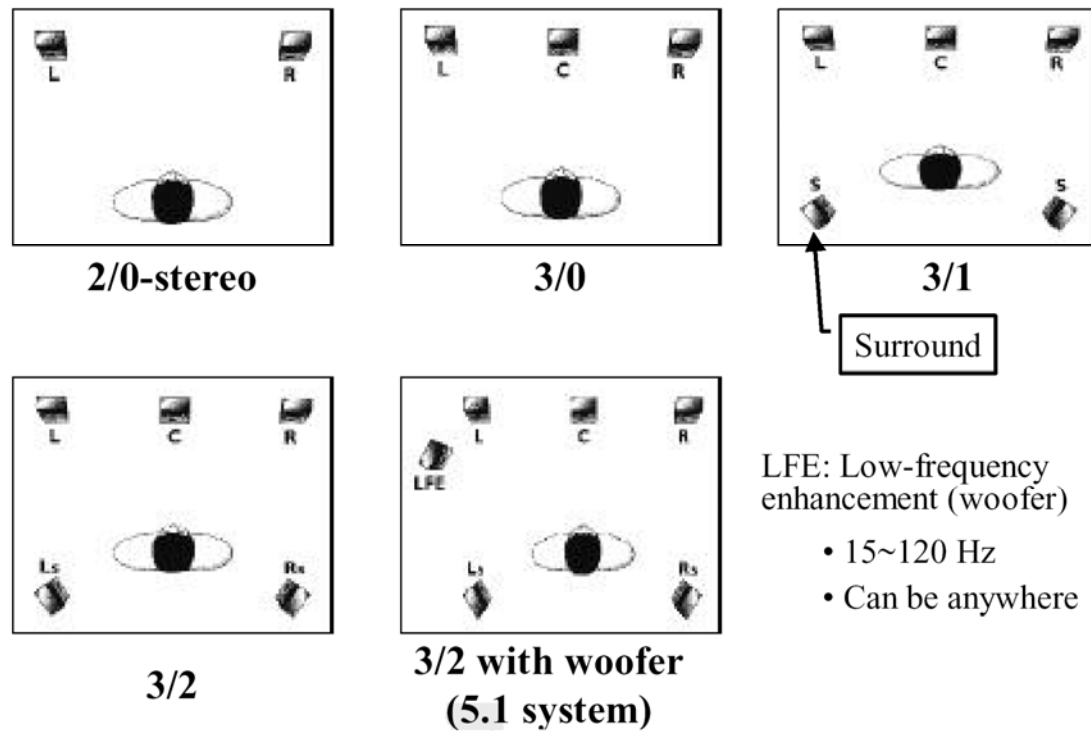
- Dacă zgomotul de cuantizare într-o bandă depășește pragul de mascare (zgomotul permis), factorul de scalare pentru această bandă este ajustat pentru a reduce zgomotul de cuantizare.
- Deoarece pentru a reduce zgomotul de cuantizare sunt necesari mai mulți pași de cuantizare deci o rată de bit mai mare, ciclul interior de rată este repetat de fiecare dată când se modifică factorii de scalare.
- Ciclul exterior este repetat până când zgomotul (calculat ca diferența între valorile spectrale originale și cuantizate) este sub pragul de mascare.

- **Codarea semnalului stereo.**

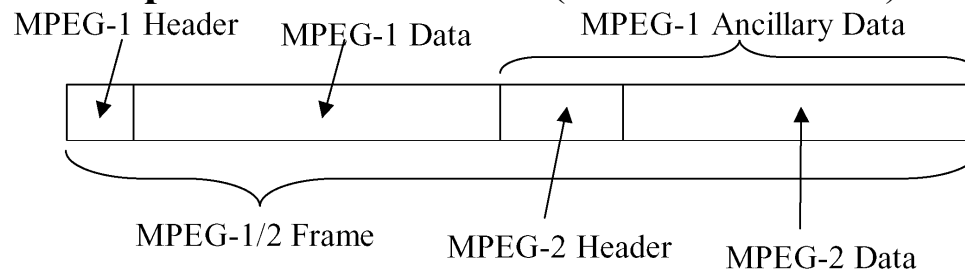
- MPEG-1 audio codează atât cu semnal mono cât și stereo.
- Sunt patru moduri de codare: mono, stereo, două canale separate și joint stereo.
- O tehnică de codare eficientă a semnalului stereo se numește *joint stereo coding*:
 - **Codarea stereo a intensității** exploatează redundanța din semnalele stereofonice bazată pe perceperea la frecvențe mai mari de 2kHz numai a anvelopei energiei canalelor drept și stâng.
 - **Codarea MS(middle/side) stereo** exploatează redundanța din semnalele stereofonice bazată pe codarea sumei și diferenței dintre canalele drept și stâng.

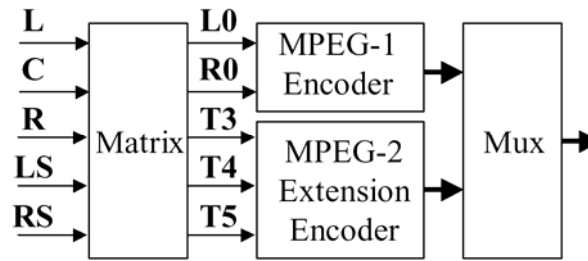
- **MPEG-2 AUDIO**

- Permite și codarea semnalelor cu frecvențe mai mici de eșantionare: 16, 22 și 24kHz.
- Realizează o analiză în frecvență cu rezoluție mărită.
- Include codorul MPEG-1 (Layer I, II și III)
- Codare multicanal:
 - Permite codarea a 2 până la 5 canale: sunet surround sau coloana sonoră pentru mai multe limbi

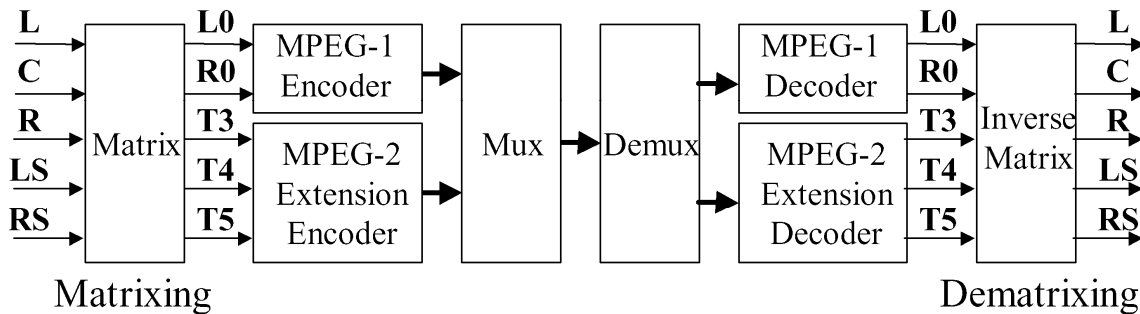


- **Compatibilitatea MPEG audio.**
- Compatibilitate directă (forward):
 - Un decodor nou poate decoda un flux de biți creat de un codor mai vechi.
 - Se poate obține relativ ușor.
- Compatibilitate inversă (backward):
 - Un decodor mai vechi poate decoda un flux de biți creat de un codor nou, cel puțin parțial.
 - Limitează eficiența codării.
- Codorul audio MPEG-2 **compatibil în sens invers** (ISO/IEC 13818-3):



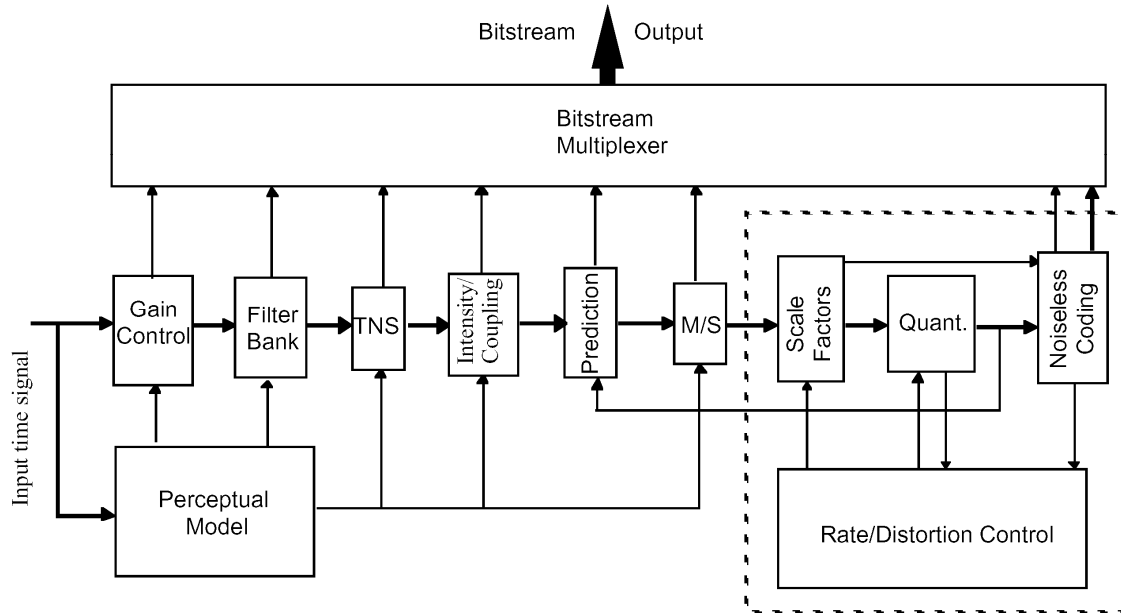


$$\begin{cases} L0 = \alpha(L + \beta \cdot C + \delta \cdot LS) \\ R0 = \alpha(R + \beta \cdot C + \delta \cdot RS) \end{cases} \quad \alpha = \frac{1}{1+\sqrt{2}}; \beta = \delta = \frac{1}{\sqrt{2}} \text{ or } \alpha = 1; \beta = \delta = 0$$



- Codarea **Non Backward Compatible (NBC)**
- **MPEG-2 Advanced Audio Coding (AAC)** ISO/IEC 13818-7 (Aprilie 1997).
- Rata de codare: 320-384 kbiți/s pentru 5 canale, 64 kbiți/canal.
- Semnal codat NBC la 320kbiți/s are aceeași calitate ca semnalul codat BC la 640kbiți/s.

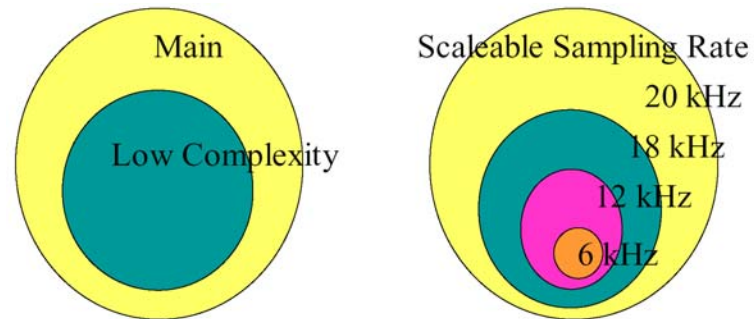
- Permite codarea multicanal: 1-48 canale audio, 0-16 canale LFE (low frequency enhancement), 0-16 canale de date.
- Aceeași structură (codare perceptuală pe subbenzi) ca la MPEG-1 cu unele îmbunătățiri.



- Îmbunătățiri
 - Banc de filtre cu rezoluție mărită (MDCT în 1024 sau 128 puncte) cu răspuns la impuls micșorat la 5.3 ms (față de 18.6 ms la Layer III) reduce distorsiunile de tip pre-echo (zgomotul de cuantizare se aude înaintea muzicii care îl produce).
 - Cuantizarea dependentă de evoluția în timp a semnalului (Temporal noise shaping TNS).
 - Predicție inversă în subbenzi oferă o codare eficientă a semnalelor tonale.

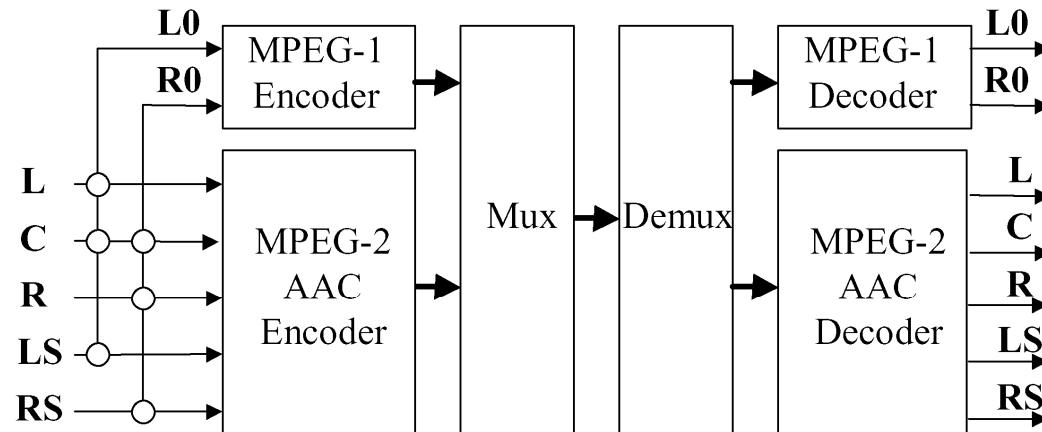
- Codare stereo Middle/Side și de intensitate mai flexibilă reduce rata de bit.
- Codare Huffman cu tabele de codare pe fiecare bloc al codorului.

- **Profiluri MPEG-2 AAC**



- Profilul principal
 - Cea mai bună calitate, complexitate maximă
 - MDCT în 1024 sau 128 puncte
- Profilul de complexitate redusă
 - Fără predicție și TNS
- Profil cu frecvența de eșantionare scalabilă
 - Complexitatea și frecvența de eșantionare sunt scalabile
 - Folosește filtre hibride ca la MPEG-1 Layer III
 - Fără predicție și intercorelare canal

- Pentru a obține compatibilitate în sens invers dar cu o rată de bit mai mare se poate folosi schema (Simulcast):



- **MPEG-4 AUDIO**

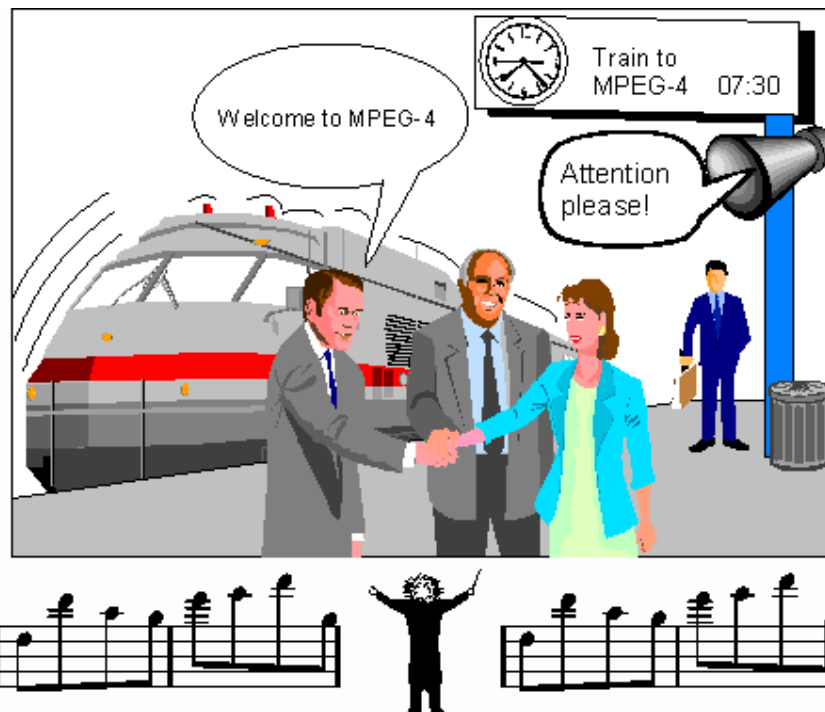
- MPEG-4 Audio integrează codarea audio sintetizată și naturală.
- Partea de codare sintetizată cuprinde realizarea muzicii și vorbirii definite simbolic. Include sisteme MIDI și Text-to-Speech. În plus, sunt incluse tehnici de localizare 3-D a sunetului, permițând crearea unor medii de sunet artificiale folosindu-se surse artificiale și naturale.

- **Codarea audio naturală**

- pentru debite între 2 kbiti/s și 64 kbiti/s.
- trei tipuri de codecuri:
 - un **codec parametric** pentru cele mai mici debite

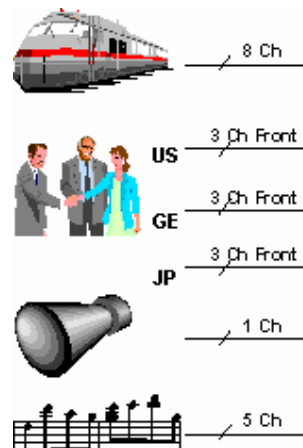
- un **codec CELP** (Code Excited Linear Predictive) pentru debite medii:
- **codecuri timp-frecvență** (TF) incluzând MPEG-2 AAC și Cuantizare Vectorială.
- Sunt oferite facilități pentru o gamă largă de aplicații de la vorbirea inteligibilă la audio-multicanal de înaltă calitate.
- În MPEG-4 sunt incluse funcții adiționale
 - controlul vitezei la redare.
 - modificarea înălțimii sunetului.
 - înlăturarea erorilor.
 - scalabilitatea.
- **Obiecte audio MPEG-4**
- MPEG-4 definește obiectele audio ca obiecte “realistice”.
- Un obiect audio “real-world” poate fi definit ca o entitate semantică audibilă (vocea unor vorbitori, instrumente muzicale etc.).
- Acesta poate fi înregistrat cu un microfon (înregistrare mono) sau cu mai multe microfoane în direcții diferite (înregistrare multicanal).
- Obiectele audio pot fi grupate sau mixate împreună dar nu pot fi (ușor) descompuse în sub-obiecte.
- Un singur obiect audio poate fi reprezentat pe unu sau mai multe canale audio, dacă definim canalele audio ca informația pentru poziția unei boxe. De exemplu un flux audio MPEG-1 poate fi un obiect audio în MPEG-4. Acest obiect poate conține un canal (mono) sau 2 canale (stereo etc.)

- **Exemple de aplicații tipice pentru MPEG-4 Audio**
- **Cântă N-1 Obiecte Audio**
 - Transmiterea a cinci semnale multicanal care reprezintă cinci instrumente ale unui cvintet. Ascultătorul poate asculta numai patru instrumente deoarece vrea să cânte el la al cincilea instrument.
- **Servicii de difuzare în mai multe limbi**
 - Cei ce urmăresc programele sportive sunt frecvent distrași de vocea comentatorului. MPEG-4 permite un “mix-minus” stil de prezentare unde să fie incluse toate sunetele, mai puțin vocea comentatorului.
 - Alternativ, într-un serviciu multi-limbi, poate fi inclus unul din comentariile în limbi străine.
- **Filme**
 - O scenă la gară dintr-un film poate conține de exemplu patru tipuri de obiecte audio:



- **Obiectul conversație:**
 - Vocea 'welcome' este cu siguranță cea mai importanta informație.
 - Vorbirea este întotdeauna localizată în fața ascultătorului.
 - Această conversație poate fi de asemenea disponibilă în mai multe limbi.
- **Obiectul fundal:**
 - Trenul va veni din depărtare spre centrul scenei, va trece de ascultător și va dispăre în spatele lui.

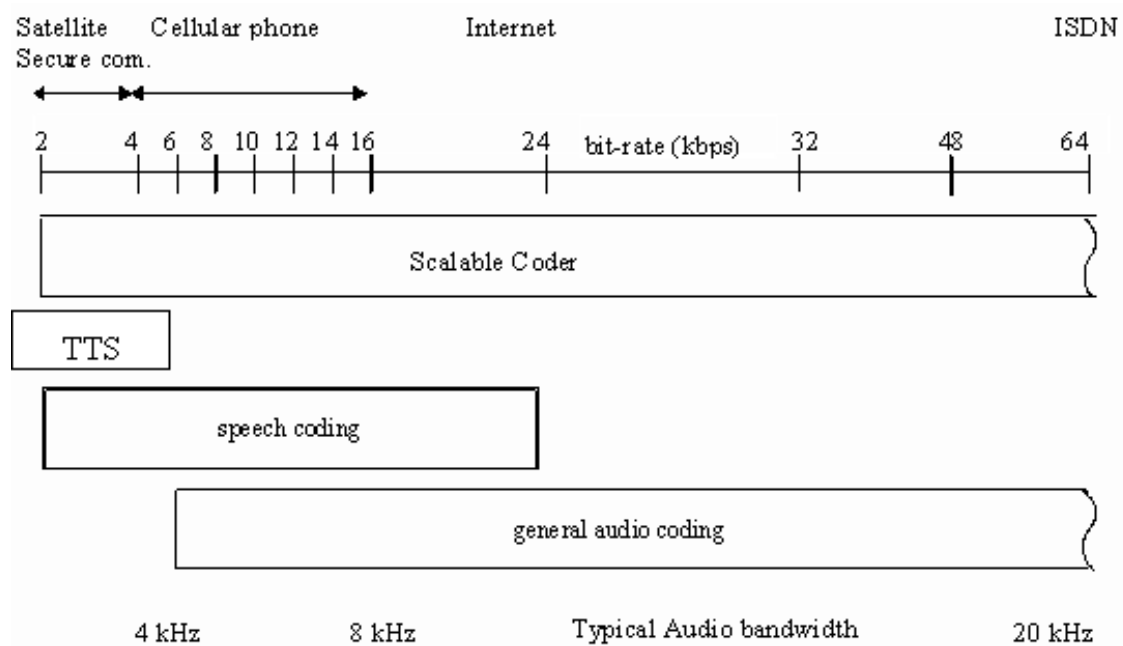
- În plus canalul pentru efecte de joasă frecvență va produce un zgomot de huruit.
- Deși includerea acestui obiect este dorită, el poate fi exclus în cazul unei conexiuni cu debit foarte redus.
- **Obiectul anunț:**
 - Pentru anunț este suficient de transmis vorbire cu calitate redusă.
 - Pot fi generate ușor unele efecte pseudo 3D și de ecou la prezentarea scenei.
- **Muzica de fundal:**
 - Orchestra poate fi codată cu MPEG-2 multicanal și fluxul de biți poate fi folosit fără necesitatea recodării.



- **Obiecte audio multi-limbă**
- Pentru o producție internațională mai mult de un obiect conversație este necesar.

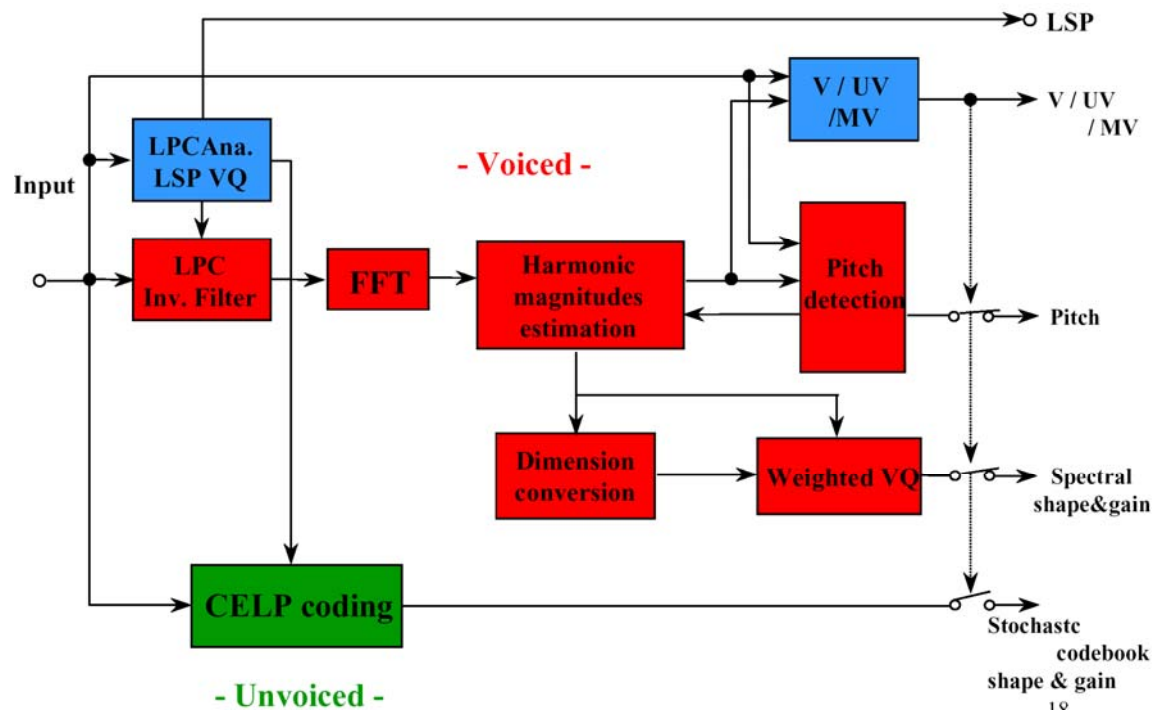
- Același obiect audio din scenă poate exista în mai multe limbi.
- Fiecare limbă este un obiect audio separat, va fi codată cu un codor independent, și va fi selectată la cerere în decodor.

- **Codarea obiectelor audio**
- Codarea MPEG-4 a obiectelor audio oferă tehnici pentru reprezentarea sunetelor naturale și pentru sunetele sintetizate pe baza descrierii structurii.
- Reprezentarea pentru sunetele sintetizate poate deriva dintr-un șir de date sau așa numita descriere de instrument și prin codarea parametrică pentru a furniza efecte ca reverberația și spațializarea.
- Această reprezentare avantajează compresia și alte funcții cum ar fi scalabilitatea și redarea la diferite viteze.
- MPEG-4 standardizează codarea audio naturală pentru debite între 2 kbiți/s și 64 kbiți/s.



- Pentru obținerea celei mai bune calități posibile pentru toate debitele și să ofere și funcții suplimentare, în standard au fost incluse trei tipuri de structuri de codare:
- **Tehnici de codare parametrică (HVXC),**
 - Codare voce cu 8 kHz frecvență de eșantionare la rate de bit foarte mici (între 2 – 4 kbiți/s).
 - **Scalabilitatea ratei de bit:** Este posibilă decodarea la 2kbiți/s dintr-un flux de bit codat cu 4kbiți/s.
 - **Variația vitezei de redare și a pitch-ului:** Utilă pentru căutarea în baze de date de vorbitori.
 - Sunt combinate două tipuri de scheme de codare: una pentru segmente vocale și alta pentru segmente nevocale.

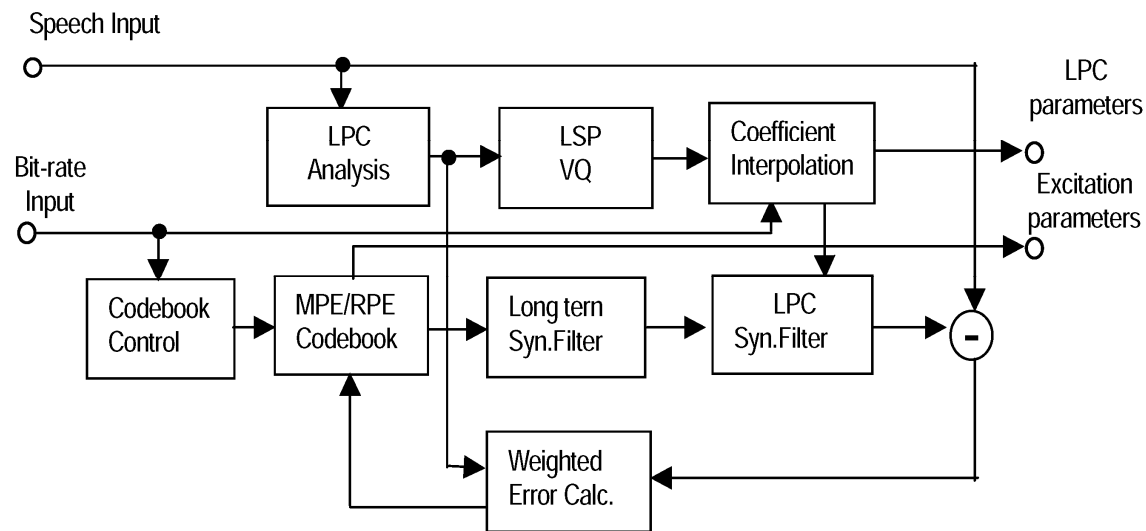
- Voce: Informația de fază este eliminată la reprezentarea spectrului de putere a erorii de predicție a filtrului LPC.
- Nevocal: Parametrii consoanelor sunt obținuți cu codorul CELP.



- **Tehnici de codare Code Excited Linear Predictive (CELP).**

- Codarea vorbirii la debite medii între 6 –24 kbiți/s.

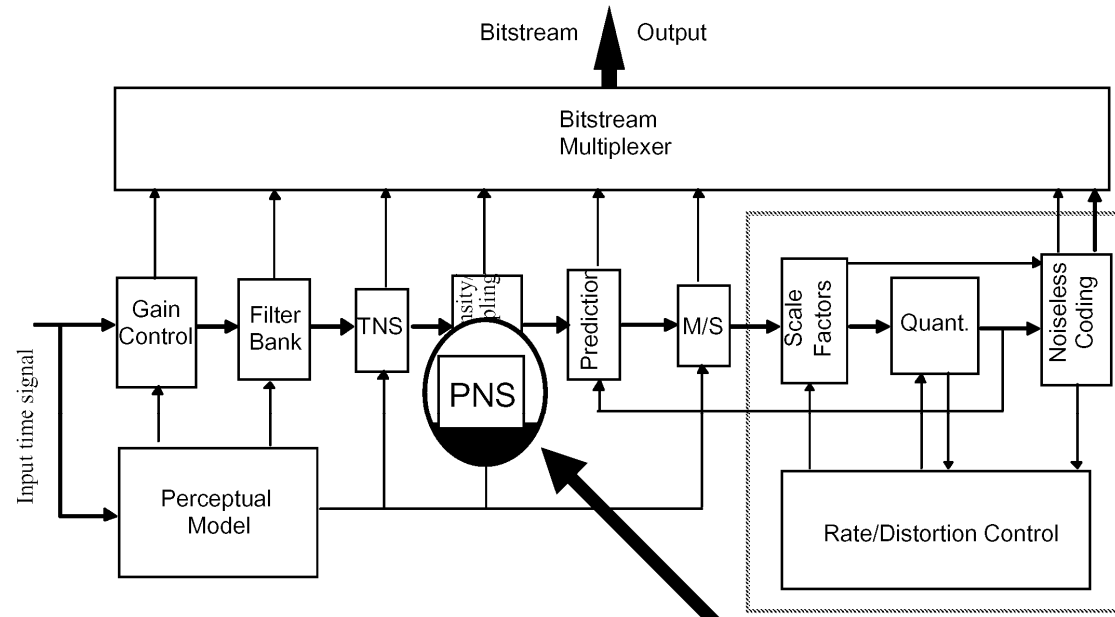
- În această zonă, două frecvențe de eșantionare, 8 și 16 kHz, sunt folosite pentru vorbirea de bandă îngustă și bandă largă.
- Banda îngustă: 3,85-12,2 kbps, pentru cadre de 10-40 ms.
- Bandă largă: 10,9-23,8 kbps, pentru cadre de 10-20 ms.



- **Tehnici de codare timp-frecvență (T/F),**

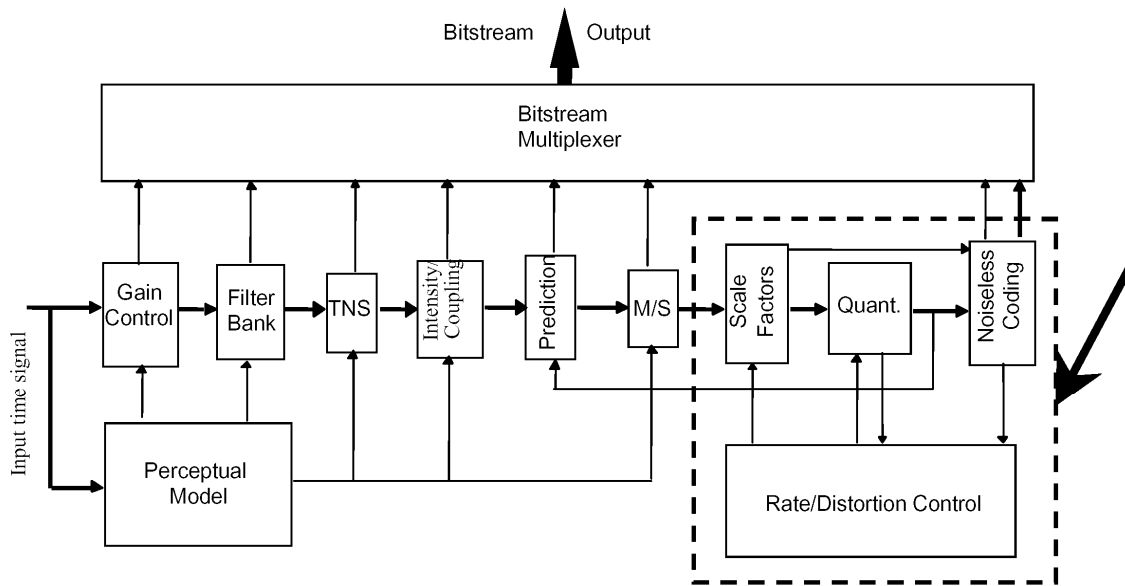
- Pentru debite peste 16 kbiți/s semnale audio.
- Se folosesc în principal codoarele TwinVQ și AAC.
- Frecvențele de eșantionare sunt peste 8 kHz.
- Extensii la AAC:

- Substituția zgomotului perceptual (PNS)



- Codarea parametrică a semnalelor asemănătoare zgomotului se folosește în codarea vorbirii (consoane).
- **Perceptual Noise Substitution** (PNS) permite o codare selectivă a frecvențelor pentru semnale similare zgomotului.
- Componentele ca de zgomot se detectează în funcție de factorul de scalare al benzii.
- Coeficienții spectrali corespunzători nu sunt cuantizați și codați. În loc de aceștia se transmite un flag de înlocuire cu zgomot și puterea totală a benzii substitute.
- Decodorul generează semnal pseudo aleator cu puterea echivalentă a coeficienților spectrali.

- **Predicție pe termen lung**
- Semnalele tonale necesită precizie la codare mai mare decât semnalele similare zgomotului (netonale).
- Componentele tonale sunt predictibile
- Predicția fiecărui coeficient spectral se face în MPEG-2 AAC cu un predictor invers adaptiv. Acesta are complexitate mare (50% din complexitatea decodării).
- În MPEG-4 se folosește Long Time Predictor (LTP) cunoscut în codarea vorbirii.
- Acesta are complexitate redusă (cu 50% mai mică față de MPEG-2 la aceleași performanțe)
- **Codecul TwinVQ** (Transform-Domain Weighted Interleave Vector Quantization)
- Codare audio la rate de bit extrem de mici (6-8 kbiți/s)
- Codoarele CELP nu se comportă bine la codarea muzicii.
- La rata dorită se obțin 0,5 biți pe componenta de frecvență!
- Selectează vectorul codat controlat de modelul perceptual.
- Este complet integrat în MPEG-4 AAC.
- Folosește aceeași reprezentare spectrală ca și codorul AAC.
- Folosește facilitățile MPEG-4 (LTP, TNS, joint stereo)



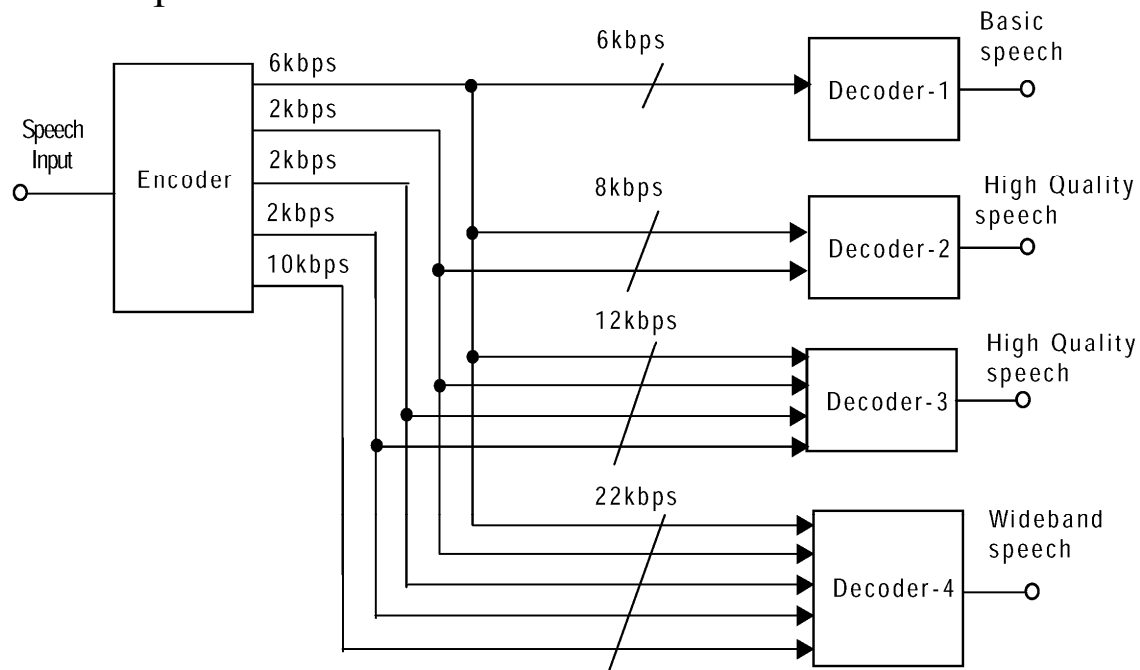
- **Structura TwinVQ:**

- Normalizarea coeficienților spectrali:
 - Anvelopa LPC (curba globală a spectrului)
 - Codarea componentelor periodice (componente armonice)
 - Codarea curbei după scara bark.
- Cuantizarea Vectorială (VQ)
 - Întreșeserea coeficienților spectrali în sub-vectori
 - Cuantizarea vectorială se face cu două seturi de cuvinte de cod.

• **Scalabilitatea codorului audio MPEG-4**

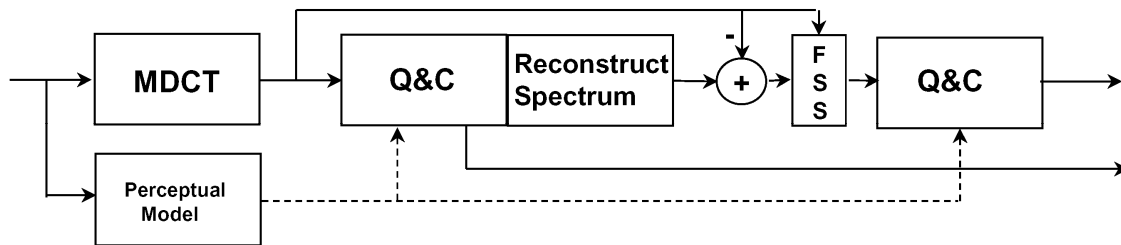
- Există mai multe tipuri de scalabilitate:

- **Scalabilitatea debitului** permite unui flux de biți să fie partiționat într-un flux cu debit mai mic care să poată fi încă decodat într-un semnal inteligibil. Partiționarea poate fi efectuată fie în timpul transmisiei sau la decodor.

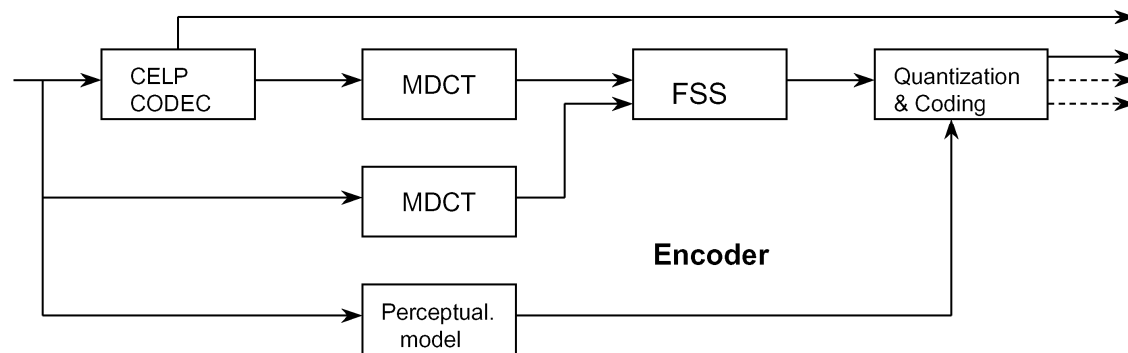


- **Scalabilitatea benzii de frecvență** este un caz particular al scalabilității debitului, unde o parte din fluxul de biți reprezintă o parte din spectrul de frecvență care poate fi ignorat în timpul transmisiunii sau la decodare.

- **Scalabilitatea complexității codorului** permite ca codoare de complexitate diferită să genereze fluxuri de biți valide și inteligibile.
- **Scalabilitatea complexității decodorului** permite ca un flux de biți să fie decodat de decodoare cu diferite niveluri de complexitate.
- Scalabilitatea funcționează cu unele din tehnicile MPEG-4, dar poate fi aplicată și unei combinații de tehnici (de exemplu cu Twin VQ ca layer de bază și AAC pentru layere extinse).
- Exemplu: Codarea semnalului eroare de cuantizare al unui modul AAC sau TwinVQ ca intrare într-un al doilea modul cuantizare/codare în frecvență.



- Exemplu: Combinarea cu codor CELP:



- **Codarea audio sintetizată**

- **Codarea Text To Speech (TTS)**

- Codoarele TTS asigură un debit între 200 biți/s și 1.2 kbiți/s și permit ca să se genereze o vorbire sintetizată inteligibilă, primind la intrare text sau text și parametrii prozodici (conturul înălțimii, durata fonemelor etc.)
- MPEG-4 oferă o interfață standard pentru operarea unui codor TTS și nu standardizează un anume sintetizor TTS.
- Sunt incluse următoarele funcționalități:
 - Sinteza vorbirii folosind prozodia vorbirii originale.
 - Controlul sincronizării buzelor cu informația despre foneme.
 - Pauză, reluare, derulare înainte/înapoi.
 - Suport pentru limbi străine și dialecte pentru text.

- Suport pentru simboluri de foneme internaționale, și suport pentru specificarea vârstei, sexului, debitului verbal al vorbitorului.

- **Sinteza după partitură**

- Tehnicile de Structurare Audio decodează datele de intrare și produc sunete.
- Această decodare este condusă de un limbaj special de sinteza numit SAOL (Structured Audio Orchestra Language), standardizat ca parte a MPEG-4.
- Acest limbaj e utilizat pentru a defini o “orchestră” alcătuită din “instrumente” (provenite din fluxul de biți și nu fixate în terminal) care creează și procesează data de control.
- Un instrument este o mică rețea de primitive de procesare de semnal care poate emula sunete specifice ca ale instrumentelor acustice naturale.
- Rețeaua de procesare a semnalului poate fi implementată hardware sau software și include generarea și procesarea sunetelor și manipularea sunetelor pre-stocate.
- MPEG-4 nu standardizează o metoda de sinteză ci mai degrabă o metodă de descriere a sintezei.
- Orice metodă curentă sau viitoare poate fi descrisă în SAOL, inclusiv sinteza wavetable, FM, aditivă, modelare psihică și granulară, precum și metode hibride non-parametrice.
- Controlul sintezei este desăvârșit prin extragerea “partiturii” sau “scenariului” din fluxul de biți.
- O partitură este un set de comenzi în timp care invocă diferite instrumente la momente de timp specifice, fiecare contribuind la interpretarea globală a muzicii sau la generarea efectelor sonore.

- Descrierea partiturii, integrată într-un limbaj numit SASL (Structured Audio Score Language), poate fi folosită pentru a crea sunete noi și de a include informații adiționale de control pentru modificarea sunetului existent.
- Aceasta permite compozitorului un control mai fin asupra sunetului final sintetizat.
- Pentru sinteza care nu necesită un control așa de fin, se poate utiliza protocolul MIDI pentru controlul orchestrei.
- Controlul fin împreună cu definirea de instrumente proprii, permite generarea unor sunete pornind de la simple efecte audio cum ar fi zgomot de pași sau de uși închise, până la simularea sunetelor naturale cum ar fi ploaia sau de la muzica cântată pe instrumente convenționale până la sunete integral sintetizate pentru efecte audio complexe sau muzica futuristă.
- Pentru terminale cu mai puține facilități și pentru aplicații care nu necesită o sinteză atât de sofisticată, un “wavetable bank format” (SASBF) este standardizat.
- Cu acest format pot fi extrase eșantioane de sunet care vor fi folosite în sinteza wavetable, de asemenea și procesări simple cum ar fi: filtre, reverberații și efecte de cor.
- În acest caz, complexitatea de calcul pentru procesul de decodare poate fi determinată exact, examinându-se fluxul de biți.
- **Efecte audio speciale**
- Decodorul bazat pe Structurarea audio/Efecte permite la decodare un flux de date care să includă atât canalele audio decodate cât și parametrii necesari pentru controlul efectelor (desfășurarea lor în timp etc.)

- Efectele sunt în esență descrieri de instrumente “speciale” servind procesoarelor de efecte aplicate asupra fluxului de intrare.
- Procesarea de efecte include reverberatoare, spațializatoare, mixere, limitatoare, controlul dinamicii, filtre, flangere, coruri și efecte hibride.
- Avându-se în vedere aceste facilități, se poate realiza pe lângă compoziția muzicală, organizarea altor tipuri de audiții cum ar fi voce, efecte sonore și ambianță generală.