

Compression des signaux audio

Gaël Mahé

présentation de Sonia Larbi et Gaël Mahé

Université Paris Descartes

Mastère TICV - janvier 2017

Plan

- 1 Introduction
- 2 Numérisation
- 3 Quantifications
- 4 Perception auditive
- 5 Codage perceptif
- 6 Codage entropique
- 7 Quelques codeurs
 - Historique
 - Codeurs MPEG-1

Plan

- 1 Introduction
- 2 Numérisation
- 3 Quantifications
- 4 Perception auditive
- 5 Codage perceptif
- 6 Codage entropique
- 7 Quelques codeurs
 - Historique
 - Codeurs MPEG-1

Problématique

Objectif = transformer un signal audio en flux binaire

- 1 Echantillonner
- 2 Quantifier
- 3 Comprimer

Codage PCM

Représentation la plus simple d'une paire de signaux stéréo de qualité CD :

- Echantillonnage à 44.1kHz
- Chaque échantillon représenté avec 16 bits

⇒ Débit résultant : 1.4 Mbit/s !

Nécessité d'un codage efficace

Décrire l'information avec le moins de bits possibles :

- 1 Pour stocker plus d'information sur un même support.
- 2 Pour transmettre l'information sur des canaux à débit limité.

Différentes gammes de qualité

	F_e (kHz)	R (bits)	Débit kbit/s	Débit usuel kbit/s
Bande téléphonique	8	13	104	4 à 64
Bande "élargie"	16	14	224	16 à 64
Bande FM (stéréo)	32	16	1024	64 à 192
Bande Hi-Fi (CD stéréo)	44.1	16	1411	64 à 192
Qualité studio (5.1)	48	16	3840	384

Un bon codeur : un compromis...

Sur le débit

- Le codeur doit proposer un débit adapté au canal de transmission.
- Dans certains cas, le codeur doit pouvoir adapter son débit (codeur hiérarchique).

Sur la complexité

Codeurs embarqués dans des systèmes aux ressources de calcul limitées (Lecteur mp3, PDA, téléphone portable) :

- Diffusion : le codeur peut être plus complexe que le décodeur.
- Communications : faible complexité du codeur et du décodeur.
- De nombreuses architectures ne disposent pas d'une unité de calcul en virgule flottante !

Un bon codeur : un compromis... (2)

Sur le retard de reconstruction

Un retard algorithmique peut être induit par les traitements par blocs ou fenêtres.

- Indifférent pour les applications de diffusion (pre-buffering).
- Gênant pour les applications conversationnelles :
 - Délai $> 150\text{ms}$ désagréable
 - Au delà de 400 ms, perte de l'interactivité.

Sur la résistance aux erreurs de transmission

- Codes correcteurs d'erreur.
- Resynchronisation : si une section du flux binaire est altérée, le codeur doit pouvoir "sauter" jusqu'à la prochaine section saine
⇒ Découpage du flux en "trames" ne dépendant pas les unes des autres.

Un critère essentiel : la qualité

Définitions variées

- Parole : Intelligibilité + reconnaissance du locuteur.
- Musique : Transparence,
i.e. pas de différence perceptive entre le signal original et le signal codé-décodé.

Tests subjectifs de qualité : ITU-R BS.1116

- "Test avec triple stimulus et référence dissimulée".
- Enregistrements courts.
- Signal original : A, signal codé : B. Trois signaux sont présentés. Deux configurations : AAB ou ABA.
- On demande au sujet de repérer le signal codé et de le noter sur une échelle de 1 à 5.
- Traitement statistique des résultats.

Plan

- 1 Introduction
- 2 Numérisation**
- 3 Quantifications
- 4 Perception auditive
- 5 Codage perceptif
- 6 Codage entropique
- 7 Quelques codeurs
 - Historique
 - Codeurs MPEG-1

Echantillonnage (1)

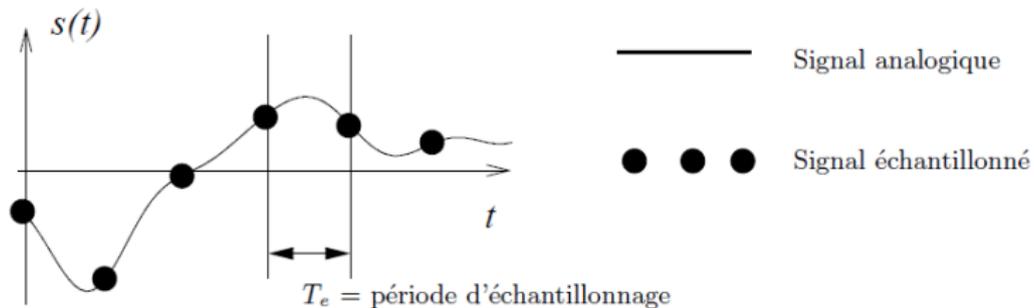


FIGURE – Echantillonnage d'un signal analogique.

Echantillonnage (2)

Comment choisir T_e ?

Théorème de Shannon

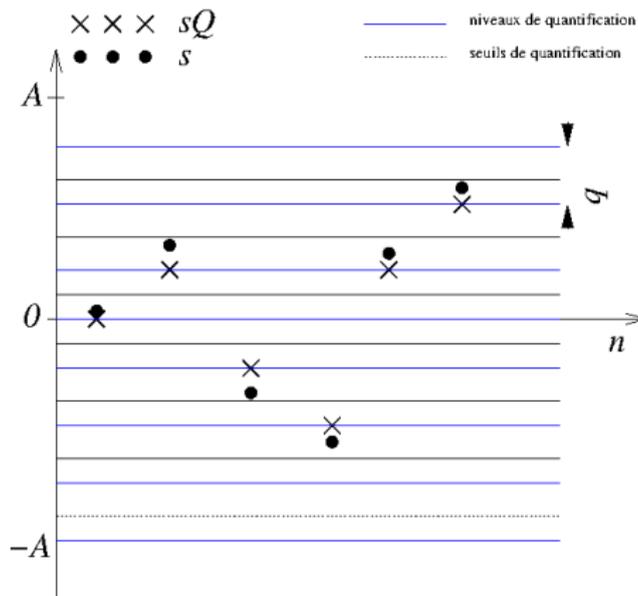
- Soit un signal de spectre à support borné $[-B; B]$.
- Si $\nu_e = 1/T_e > 2B$,
alors l'échantillonnage est sans perte d'information.

Echantillonnage (3)

Application	Fréquence max B	ν_e
téléphone	3400 Hz	8000 Hz
FM	15 kHz	32 kHz
Hifi	22 kHz	44,1 kHz

TABLE – Fréquences d'échantillonnage selon les applications audio.

Quantification



► Erreur de quantification :

$$q[n] = s[n] - Q(s[n])$$

Comment minimiser cette erreur ?

Plan

- 1 Introduction
- 2 Numérisation
- 3 Quantifications**
- 4 Perception auditive
- 5 Codage perceptif
- 6 Codage entropique
- 7 Quelques codeurs
 - Historique
 - Codeurs MPEG-1

Définitions

Formulation

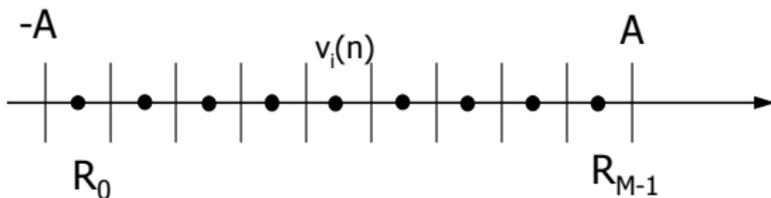
Quantificateur scalaire, cas général : $Q : \mathbb{R} \rightarrow \{0, 1, \dots, M - 1\}$

En pratique, signaux bornés : $Q : [-A, A] \rightarrow \{0, 1, \dots, M - 1\}$

Quantificateur régulier (uniforme)

Caractérisé par :

- Une partition de $[-A, A]$ en M intervalles égaux : $\bigcup R_i = [-A, A]$ et de largeur $\frac{2A}{M}$
- M représentants v_i , avec $v_i \in R_i$.
 $\{v_0, v_1, \dots, v_{M-1}\}$ forme le "codebook" ou dictionnaire.



Définitions (2)

Processus de quantification

Quantification : $Q : x \in [-A, A] \rightarrow i$ tel que $x \in R_i$

Déquantification : $Q^{-1} : i \in \{0, \dots, M-1\} \rightarrow v_i$

Quantif. puis déq. : $q : x \in [-A, A] \rightarrow y \in \{v_0, \dots, v_{M-1}\}$

Mesure de l'erreur $\epsilon(x) = q(x) - x$

Mesure de distorsion possible : l'EQM (erreur quadratique moyenne)

$$\begin{aligned} \sigma_Q^2 &= E[|q(X) - X|^2] = \int_{-A}^A p(x) |q(x) - x|^2 dx \\ &= \sum_{i=0}^{M-1} \int_{R_i} p(x) |v_i - x|^2 dx \end{aligned}$$

où $p(x)$ est la densité de probabilité de X .

Minimisation de la puissance de l'erreur

Partition optimale pour $\{v_i\}$ donné

La minimisation de la puissance de l'erreur σ_Q^2 donne : $R_i = \{x/\forall j, ||v_i - x|| \leq ||v_j - x||\}$ (régions de Voronoï / règle du plus proche voisin).

Dans notre cas unidimensionnel, $R_i = [\frac{v_{i-1}+v_i}{2}, \frac{v_i+v_{i+1}}{2}]$

Représentants optimaux pour $\{R_i\}$ donné : règle du centroïde (centre de gravité)

$$v_i = \frac{\int_{R_i} xp(x)dx}{\int_{R_i} p(x)dx}$$

Dans le cas d'une distribution uniforme, si $R_i = [a, b]$, $v_i = \frac{a+b}{2}$.

Minimisation de la puissance de l'erreur (détails)

- 1) soit $\{v_0, \dots, v_{M-1}\}$ un dictionnaire donné. La meilleure partition est alors $R_i = \{x/\forall j, ||v_i - x|| \leq ||v_j - x||\}$
 Soit t_i la valeur définissant la limite entre les partitions R_i et R_{i+1} , la minimisation de σ_Q^2 par rapport à t_i donne :

$$\frac{\partial}{\partial t_i} \left[\int_{t_{i-1}}^{t_i} (x - v_i)^2 p(x) dx + \int_{t_i}^{t_{i+1}} (x - v_{i+1})^2 p(x) dx \right] = 0 \Rightarrow (t_i - v_i)^2 p(t_i) - (t_i - v_{i+1})^2 p(t_i) = 0$$

$$\text{soit } t_i = \frac{v_i + v_{i+1}}{2}$$

- 2) Etant donnée une partition $R_0, \dots, R_i, \dots, R_{M-1}$, les représentants optimaux sont donnés par la règle du centroïde (ou centre de gravité) de la partie de la DDP placée dans la région R_i : $v_i = \frac{\int_{x \in R_i} x p(x) dx}{\int_{x \in R_i} p(x) dx}$ (1).

La minimisation de σ_Q^2 par rapport à v_i donne :

$$\frac{\partial}{\partial v_i} \int_{x \in R_i} p(x) (v_i - x)^2 dx = 0 \Rightarrow 2 \int_{x \in R_i} p(x) (v_i - x) dx = 0 \Rightarrow (1)$$

Rappel Intégrale de Riemann

$f(t)$ continue sur $[a, b]$, alors 1) l'intégrale de Riemann de f existe, i.e. $\int_a^b f(t) dt$ existe, 2) soit $F(t) = \int_a^t f(x) dx, \forall t \in [a, b] \Rightarrow \frac{\partial}{\partial t} F(t) = f(t), \forall t \in [a, b]$

Minimisation de la puissance de l'erreur (2)

Recherche d'un quantificateur optimal

Algorithme itératif de Lloyd-Max : optimisation alternée des représentants et de la partition.

Très sensible à l'initialisation. Ne converge que vers un minimum local.

⇒ Quantificateur uniforme à M pas sur $[-A, A]$

Optimal sous condition de distribution uniforme.

- $\Delta = \frac{2A}{M}$
- $v_i = -A + \Delta(\frac{1}{2} + i)$
- $R_i = [-A + \Delta i, -A + \Delta(i + 1)]$

Quantificateur uniforme

Puissance de l'erreur de Q (X de distribution uniforme)

$$\begin{aligned}
 \sigma_Q^2 = E[|\epsilon(x)|^2] &= \sum_{i=0}^{M-1} \int_{R_i} p(x) |v_i - x|^2 dx \\
 &= \sum_{i=0}^{M-1} \int_{-A+\Delta i}^{-A+\Delta(i+1)} \frac{1}{2A} |x + A - \Delta(\frac{1}{2} + i)|^2 dx \\
 &= \frac{1}{2A} \sum_{i=0}^{M-1} \int_{-\Delta/2}^{\Delta/2} x^2 dx = \frac{M}{2A} \frac{\Delta^3}{12} = \frac{\Delta^2}{12} = \frac{1}{3} \frac{A^2}{M^2}
 \end{aligned}$$

Note

$$\sigma_X^2 = \int_{-\infty}^{\infty} x^2 p_X(x) dx = \int_{-A}^A x^2 \frac{1}{2A} dx = \frac{A^2}{3}$$

Quantificateur uniforme (détails)

Puissance de l'erreur de Q (X de distribution uniforme)

$$\begin{aligned}
 \sigma_Q^2 &= \sum_{i=0}^{M-1} \int_{R_i} p(x)(v_i - x)^2 dx \\
 &= \int_{x \in R_0} \frac{1}{2A} (v_0 - x)^2 dx + \dots + \int_{x \in R_{M-1}} \frac{1}{2A} (v_{M-1} - x)^2 dx \\
 &= \frac{1}{2A} \sum_{i=0}^{M-1} \int_{-A+\Delta i}^{-A+\Delta(i+1)} (x + A - \Delta(\frac{1}{2} + i))^2 dx
 \end{aligned}$$

En posant $t = x + A - \Delta(\frac{1}{2} + i)$, on obtient

$$\sigma_Q^2 = \frac{1}{2A} \sum_{i=0}^{M-1} \int_{-\Delta/2}^{\Delta/2} t^2 dt$$

Quantificateur uniforme (2)

Rapport signal à bruit de quantification

Quantification sur M pas avec $M = 2^k$ ($k =$ nombre de bits par échantillon).

$$\begin{aligned} RSB &= 10 \log \frac{E[X^2]}{E[(q(X) - X)^2]} = 10 \log \frac{\sigma_X^2}{\sigma_Q^2} \\ &= 10 \log M^2 = 20 \log 2^k \approx 6k \text{ dB} \end{aligned}$$

Règle des "6dB par bit"

Le rapport signal à bruit de quantification (uniforme) augmente de 6dB lorsque le niveau de quantification augmente d'un bit.

Quantificateur uniforme (3)

- La formule $RSB = 6k$ n'est valable que pour un signal de distribution uniforme
- Ce n'est pas le cas des signaux audio...
- On montre que :

$$RSB_{\text{dB}} = 6k + 4,76 - f_c$$

avec f_c le *facteur de crête* :

$$f_c \simeq \begin{cases} 13 \text{ dB} & \text{pour la musique multi-instruments} \\ 18 \text{ dB} & \text{pour la parole} \end{cases}$$

- La règle des "6dB par bit" reste valable

Quantification non uniforme

Quantificateurs non uniformes

On a supposé que x était uniformément distribué sur $[-A, A]$.
C'est rarement le cas !

Quantification optimale

$$\sigma_Q^2 = E[|\epsilon(x)|^2] = \sum_{i=0}^{M-1} \int_{R_i} p(x) |v_i - x|^2 dx$$

→ Adapter la répartition des v_i et des R_i à p pour minimiser σ_Q^2 :
Plus de v_i là où les valeurs du signal sont plus fréquentes

Pourquoi est-ce délicat d'utiliser des quantificateurs optimaux ?

- La ddp du signal n'est jamais connue précisément
- σ_Q^2 peut augmenter fortement si la ddp réelle s'éloigne de la ddp théorique
- ▶ Quantificateur logarithmique ("loi A") pour la parole

Quantification vectorielle

La même chose dans \mathbb{R}^N !

- Dictionnaire optimal pour une partition donnée : centroïdes.
- Partition optimale pour un dictionnaire donné : régions de Voronoï.

Par contre, pour l'optimisation alternée, l'algorithme de Lloyd-Max est trop instable.

Application aux signaux unidimensionnels

- 1 On groupe les échantillons à coder en paquets de N échantillons
 $X(m) = [x_{1+(m-1)N} \dots x_{mN}]$.
- 2 On quantifie chaque paquet à l'aide d'un quantificateur vectoriel.

Quantification vectorielle (2)

Formalisation

On appelle quantificateur vectoriel de dimension N et de taille M l'application :

$$Q : \mathbb{R}^N \rightarrow D \text{ avec } D = \{\vec{v}_1 \dots \vec{v}_M\} \text{ où } \vec{v}_i \in \mathbb{R}^N.$$

- L'espace \mathbb{R}^N est partitionné en M régions définies par :
 $\mathcal{R}^i = \{X : Q(X) = \vec{v}_i\}$.
- D est le dictionnaire, \vec{v}_i un représentant.
- Mesure de distorsion possible : **distance euclidienne** :

$$d(X, \vec{v}_i) = \frac{1}{N} \|X - \vec{v}_i\|^2.$$

- Pour une résolution k (nombre de bit par échantillon) et une dimension N des vecteurs, taille du dictionnaire : $M = 2^{kN}$.
- k n'est plus nécessairement entier. Il suffit que kN le soit.

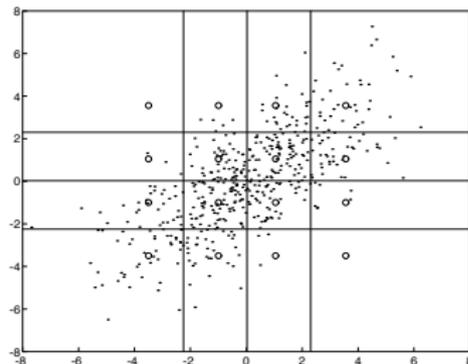
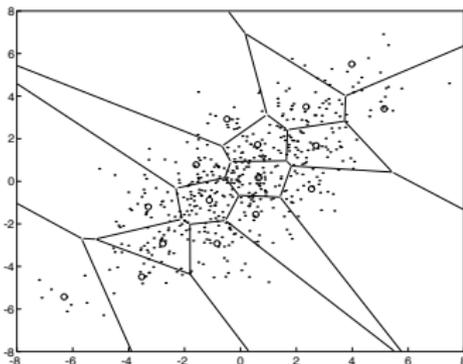
Application de la QV aux signaux audio

Intérêt

Les échantillons du signal audio sont corrélés, les vecteurs $X(m)$ n'occupent qu'une région limitée de l'espace \mathbb{R}^N . Le QV s'adapte mieux à la forme du nuage de points.

Exemple : QV avec $N = 2$ et $M = 16$, donc $k = 2$, comparé à un QS avec une partition $M = 4$ représentée dans le plan \mathbb{R}^2 .

- le QS impose une occupation identique relativement aux 2 axes,
- le QV permet une grande liberté dans le choix de la partition,
- la partition n'est pas nécessairement composée de régions rectangulaires,
- plus le signal est corrélé et plus le gain du QV est important.



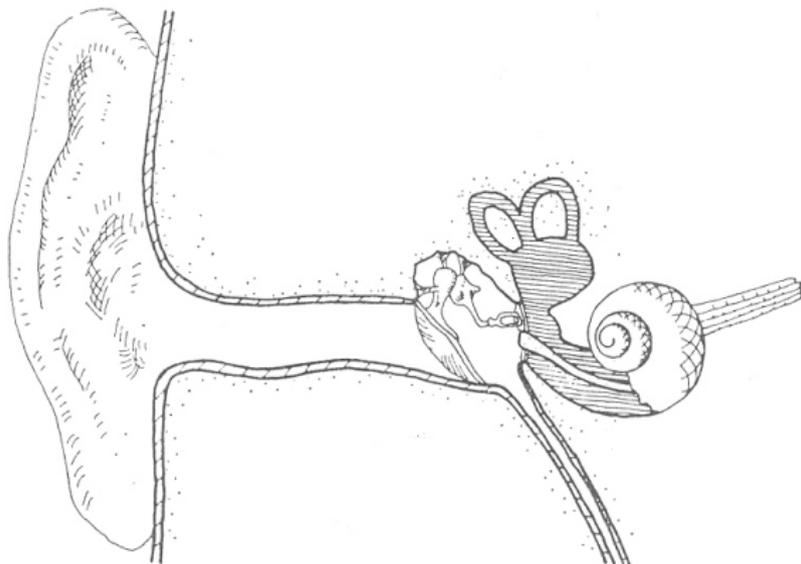
Quelle erreur de quantification est-elle acceptable ?

- Dans tous les cas (quantification uniforme ou optimale, scalaire ou vectorielle), erreur d'approximation \rightarrow bruit
- Critère d'acceptabilité = audibilité ou désagément de ce bruit
- ▶ Tenir compte de la perception auditive

Plan

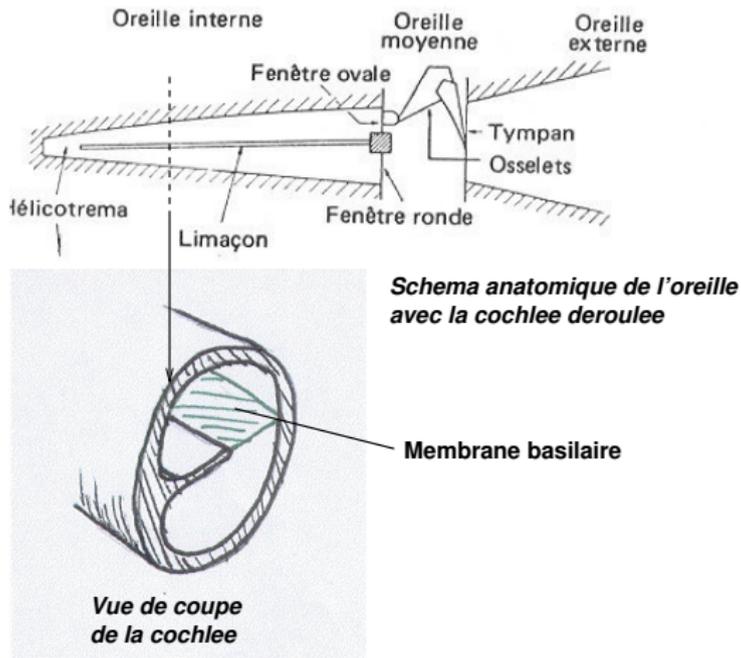
- 1 Introduction
- 2 Numérisation
- 3 Quantifications
- 4 Perception auditive**
- 5 Codage perceptif
- 6 Codage entropique
- 7 Quelques codeurs
 - Historique
 - Codeurs MPEG-1

Vue générale



La cochlée

Cochlée = cavité de 32 mm de long, enroulée sur elle-même (2,5 tours), baignée par le liquide lymphatique.



Une transformation fréquence-espace

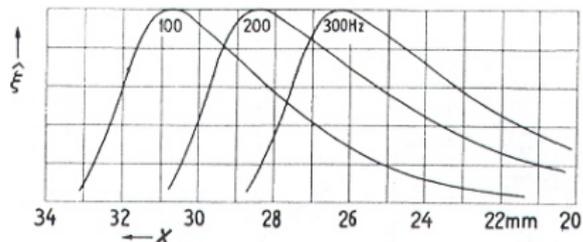


FIGURE – Amplitude des oscillations de la membrane basilaire selon la distance à la fenêtre ovale et la fréquence.

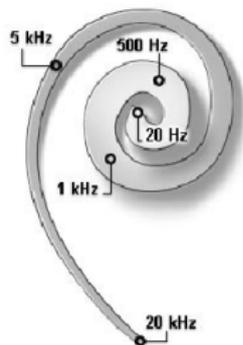


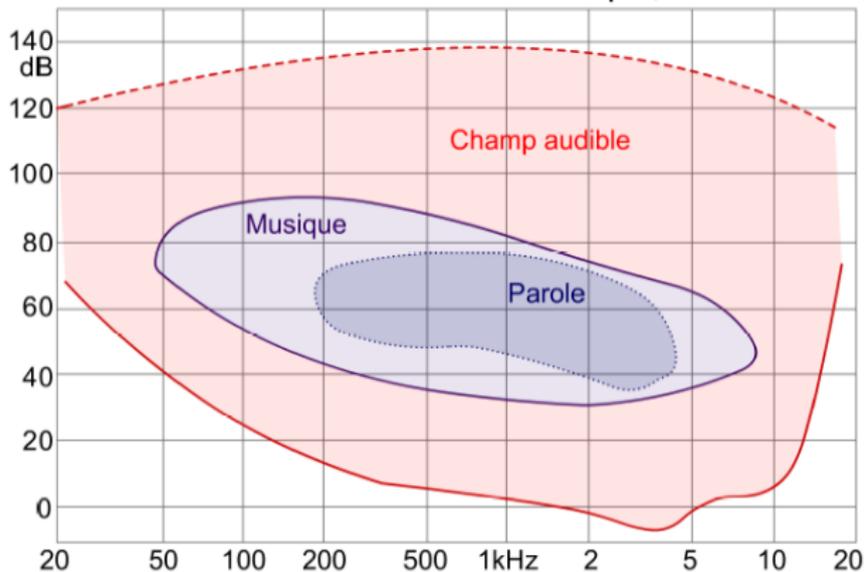
FIGURE – Tonotopie de la cochlée.

Le champ audible

Niveau sonore en échelle ajustée des dB :

$$L_{\text{dB}} = 10 \log \left(\frac{I}{10^{-12}} \right)$$

avec I = intensité de l'onde acoustique, en W/m^2



Perception de la puissance sonore

Pourquoi ?

- Seuil d'audition d'un son pur autour de 1000 Hz = +3 dB
- 2 sons purs, à 800 et 1000 Hz, sont audibles si leurs puissances $> 3\text{dB}$
- 2 sons purs, à 1000 et 1010 Hz, sont audibles si leurs puissances $> 0\text{dB}$

Interprétation

La force sonore perçue par l'oreille autour d'une fréquence f résulte d'une intégration de la puissance du son dans une bande fréquentielle Δf appelée **bande critique**

- Pour $f < 500$ Hz, $\Delta f = 100$ Hz
- Pour $f > 500$ Hz, Δf proportionnel à f

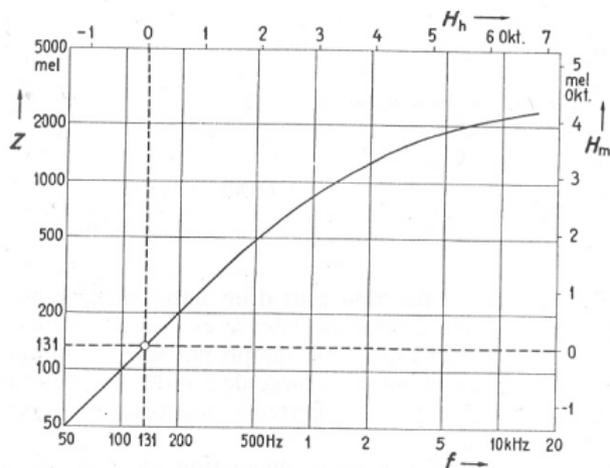
Perception de la hauteur

● Définition physique :

- Pour un son pur, hauteur = fréquence
- Pour un son harmonique, hauteur = fréquence fondamentale

● Echelle perceptive : tonie (en mels)

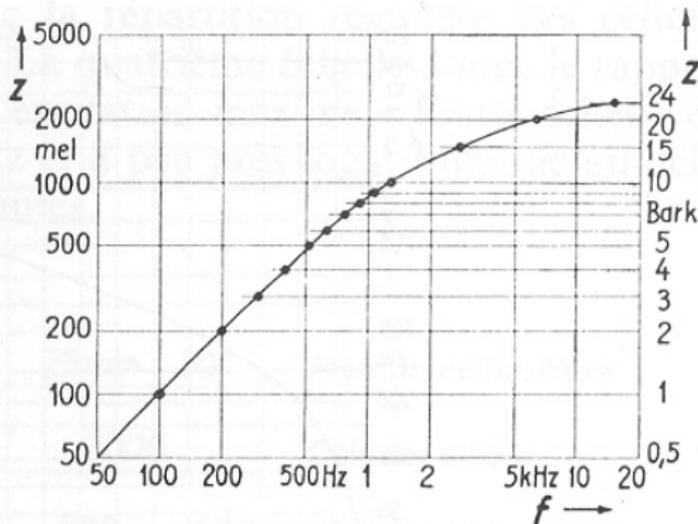
Un son A perçu comme 2 fois plus aigu qu'un son B
a une tonie 2 fois supérieure



Relation tonie / fréquence
pour un son pur

Bandes critiques et hauteur

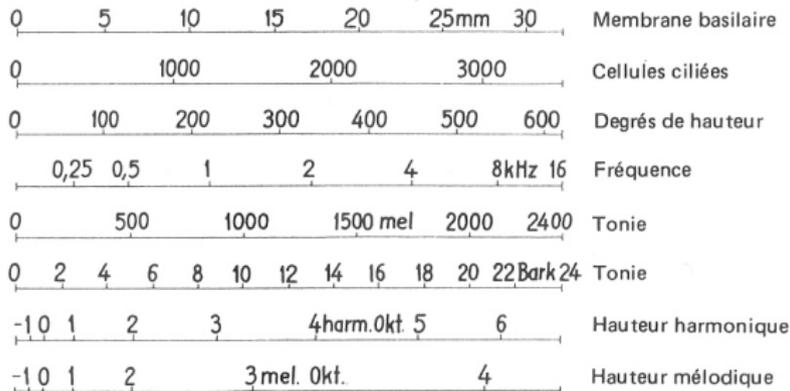
- Décomposition de la bande audible en 24 bandes critiques adjacentes.
- Tracé des points de coordonnées (f_{max} de la i^{e} BC, i) :



→ échelle mel = échelle des bandes critiques.
Nouvelle unité de tonie : 1 Bark = 1 BC = 100 mels

Grandeurs physiques, perceptives et physiologiques

Echelles de hauteur et membrane basilaire :

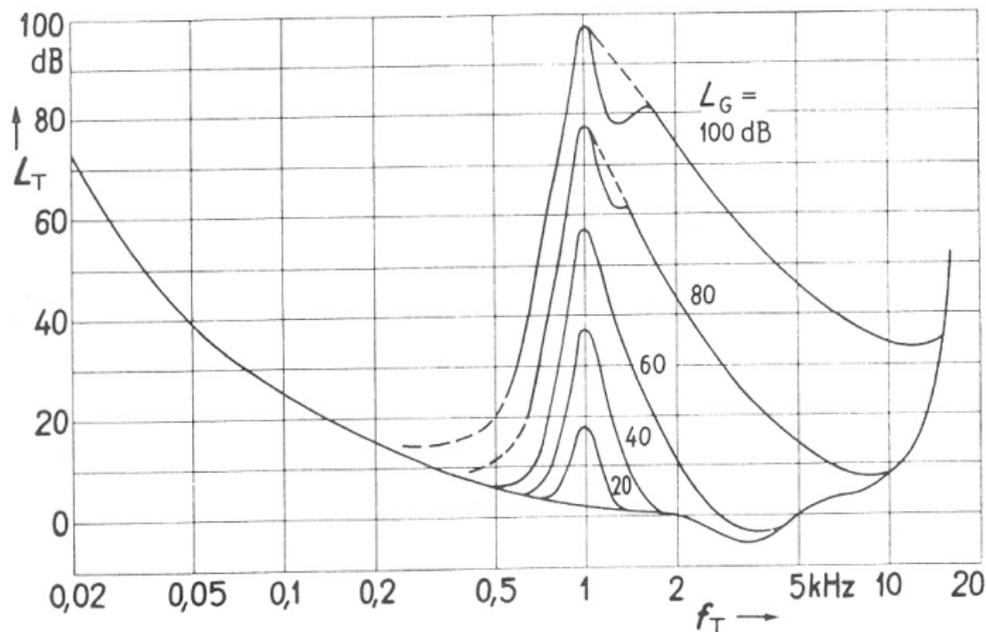


Relation linéaire

- tonie / lieu d'excitation principal :
1 Bark \simeq 1,3 mm de membrane
- degrés de hauteur perceptible / tonie / nombre de cellules ciliées :
1 degré \simeq 3,9 mel \simeq 6 cellules

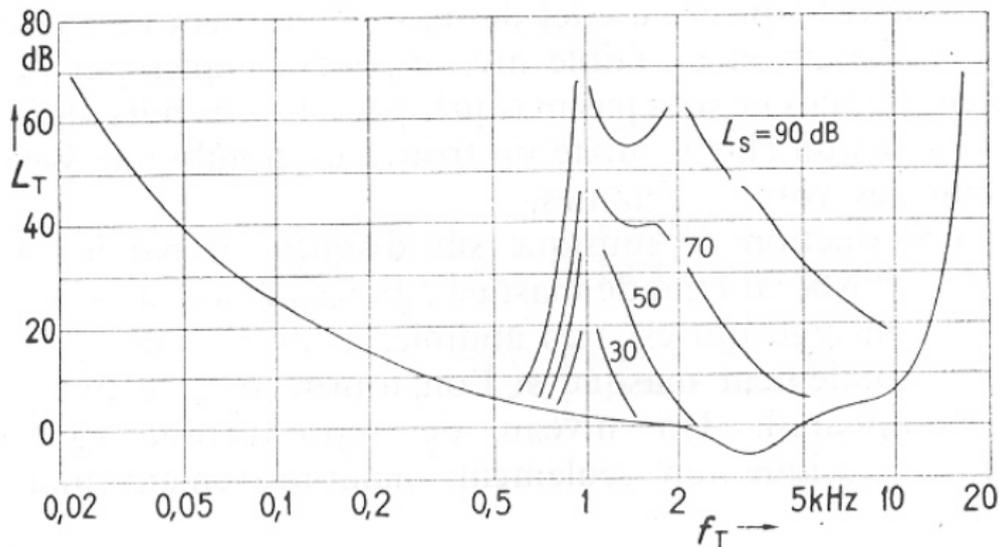
Masquage simultané (1)

Seuil d'audition en présence d'un bruit de bande étroite (160 kHz) de fréquence centrale 1 kHz :

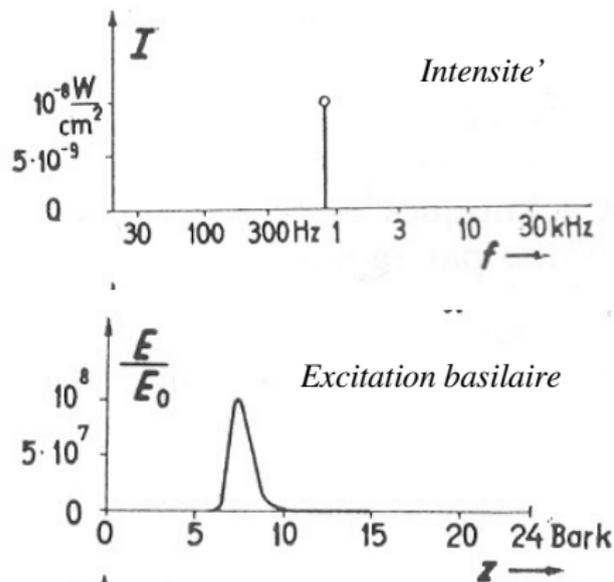


Masquage simultané (2)

Seuil d'audition en présence d'un son pur de fréquence 1 kHz :



Masquage simultané : pourquoi ?

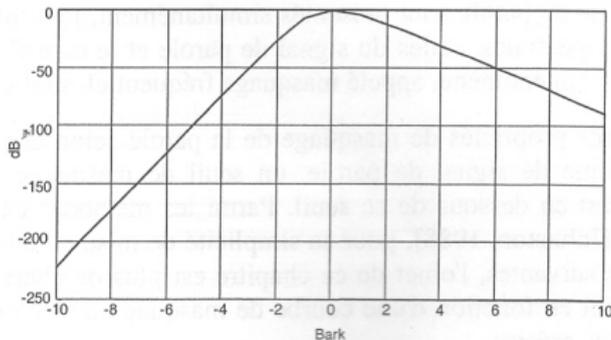


Ex :

Pour un son pur de 1 kHz,
l'excitation s'étend sur plu-
sieurs bandes critiques
et inhibe les cellules ciliées
correspondantes

Seuil de masquage d'un son complexe : méthode de Johnston

- 1 Energie par bande critique : $B(i) = \int_{f_{\min}(i)}^{f_{\max}(i)} |X(f)|^2 df$
- 2 Convolution par la fonction d'étalement : excitation $E(i) = B(i) * f_E(i)$



- 3 Seuil de masquage $S(i) = E(i) - O(i)$, avec

$$O(i) = \begin{cases} 5,5 \text{ dB} & \text{pour du bruit} \\ 14,5 + i \text{ dB} & \text{pour une tonale} \end{cases}$$

→ $O(i) = \alpha(15.5 + i) + (1 - \alpha)5.5$
avec α coef de tonalité

MPEG : Modèle psychoacoustique I

Rôle

- Calculer le seuil de masquage à partir de la d.s.p.
- Complexité moyenne.

Cas simple : sinusoïde ou bruit bande étroite

- $M_{dB}(f_2; f_1, \sigma^2) = \sigma_{dB}^2 + a(f_1) + E_{dB}(f_2 - f_1, \sigma^2)$
- Pour une raie tonale : $a(f_B) = -6.025 - 0.275 \times f_B$
- Pour un bruit bande étroite : $a(f_B) = -2.025 - 0.175 \times f_B$
- Fonction d'étalement :

$$E_{dB}(f_B, \sigma_{dB}^2) = \begin{cases} 17(f_B + 1) - 0.4 \sigma_{dB}^2 - 6 & \text{si } -3 \leq f_B < -1 \\ (0.4 \sigma_{dB}^2 + 6)f_B & \text{si } -1 \leq f_B < 0 \\ -17f_B & \text{si } 0 \leq f_B < 1 \\ -(17 - 0.15 \sigma_{dB}^2)(f_B - 1) - 17 & \text{si } 1 \leq f_B < 8 \end{cases}$$

MPEG : Modèle psychoacoustique I (2)

Procédure

- 1 Identifier les maxima locaux dans le spectre.
- 2 Classifier chaque maximum en tonal / non-tonal selon la différence d'amplitude avec ses voisins.
- 3 Ne garder qu'un maximum par bande critique (composantes).
- 4 Calculer la courbe de masquage pour chaque composante.
- 5 Additionner l'ensemble des courbes de masquage.

MPEG : Modèle psychoacoustique II

Procédure

- 1 Découper le spectre en partitions (1/3 bande critique).
- 2 Calculer une proportion tonal/non-tonal dans chaque partition.
- 3 Calculer une mesure d'imprédictibilité en fonction des spectres précédents (les composantes stables sont plus difficiles à masquer).
- 4 Calculer la courbe de masquage dans chaque partition, fonction de la mesure d'imprédictibilité et du rapport tonal/non tonal.
- 5 Additionner les courbes de masquage pour chaque partition.

Complexité élevée

Beaucoup plus de composantes à calculer et sommer !

Plan

- 1 Introduction
- 2 Numérisation
- 3 Quantifications
- 4 Perception auditive
- 5 Codage perceptif**
- 6 Codage entropique
- 7 Quelques codeurs
 - Historique
 - Codeurs MPEG-1

Comment réduire le débit binaire des signaux audio ?

(1)

Le problème

Plus le nombre de bits par échantillon est faible, plus le signal est bruité...

Les solutions précédentes

Quantification non-uniforme ?

- Q. loi A (voix) → on passe de 12 à 8 bits par échantillon à qualité constante

Quantifier des blocs d'échantillons ? (Q vectorielle)

- Lourd à mettre en oeuvre, le quantificateur optimal dépend du signal.
- Erreur de quantification audible dès que le taux de compression $> 35\%$.

Quantifier la composante imprédictible du signal ? (codage prédictif)

- Modèle auto-régressif : $s(n) = e(n) - \sum_i a_i s(n - i)$
- Au lieu de coder les $s(n)$, on code les $e(n)$ et les a_i (qui varient lentement)

Comment réduire le débit binaire des signaux audio ? (2)

Le "défaut" commun de ces méthodes

Exploitent les propriétés du signal audio, mais pas celles de son récepteur (l'oreille)
→ réduction du débit de 50% au mieux

Approche perceptive

- Au lieu de minimiser la puissance de l'erreur de quantification σ_Q^2 ,
Viser une erreur de quantification non audible
- **Objectif = minimiser le débit binaire sous contrainte d'inaudibilité du bruit de quantification.**

Une première idée

Rapport signal à bruit :

$$RSB = \frac{P_s}{P_q} = \frac{E[s^2]}{E[q^2]}$$

En décibels

- puissance en échelle des dB ajustée :

$$P^{\text{dB}} = 10 \log(P/10^{-12})$$

- $RSB_{\text{dB}} = 10 \log(RSB) = P_s^{\text{dB}} - P_q^{\text{dB}}$

Or :

$$RSB_{\text{dB}} = 6k + 4,76 - f_c$$

avec

$$f_c \simeq \begin{cases} 13 \text{ dB} & \text{pour la musique multi-instruments} \\ 18 \text{ dB} & \text{pour la parole} \end{cases}$$

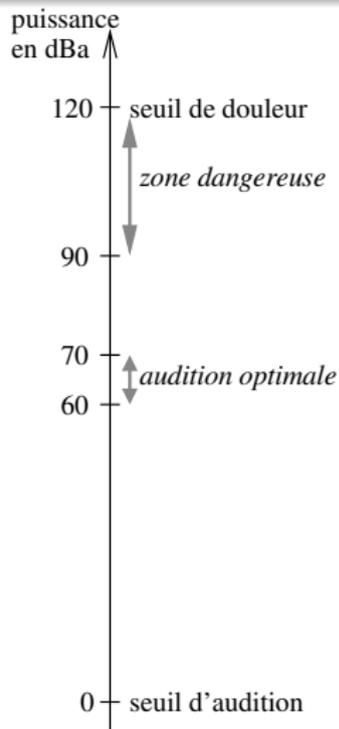


FIGURE – Echelle des puissances

Nécessité d'une compression plus forte

Qualité	ν_e (kHz)	k (bit/ech)	Débit (kbit/s)	Débit usuel (kbit/s)
Téléphone	8	13	104	4 à 64
Tel. bande élargie	16	14	224	16 à 64
FM stéréo	32	16	1024	64 à 192
HiFi (CD stéréo)	44,1	16	1411	64 à 192
Qualité studio (5.1)	48	16	3840	384
Qualité "parfaite"	96	24	13824	1000

TABLE – Débits nominal et compressé usuel selon la qualité audio souhaitée.

Idée : exploiter le masquage fréquentiel

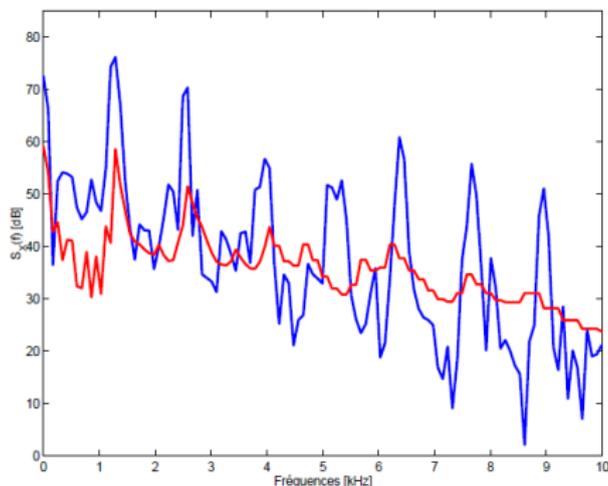


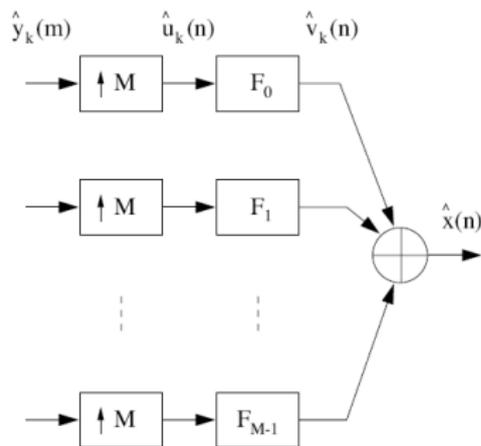
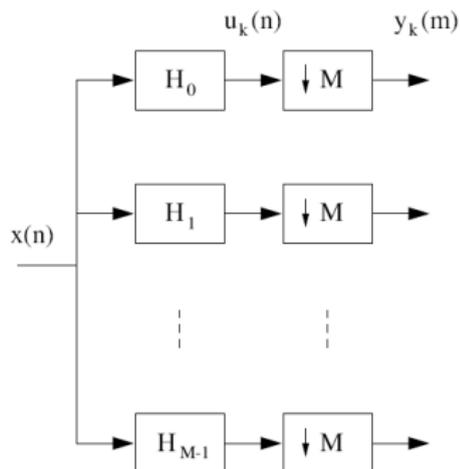
FIGURE – Seuil de masquage d'un son complexe.

Si spectre du bruit de codage $<$ seuil de masquage,
alors la dégradation est inaudible.

Comment reformer le spectre du bruit de quantification

- Le bruit de quantification a un spectre d'amplitude plat
- Idée : décomposer le signal à coder en plusieurs signaux à bande étroite et utiliser un quantificateur différent pour chaque signal de bande
→ spectre en escalier qui « suit » à peu près le seuil de masquage
- Solution : bancs de filtres

Traitement multicanal



Banc de filtres d'analyse : on applique en parallèle M filtres au signal x , et on code chacun des signaux de sous-bandes y_k .

Idée : au lieu de coder un signal échantillonné à f_e , on code M signaux échantillonnés à $\frac{f_e}{M}$. Le nombre d'échantillons à coder est préservé.

Banc de filtres modulés

Définition

On suppose que H est un filtre passe-bas idéal de fréquence de coupure réduite $\frac{1}{4M}$.
 On forme des filtre $H_k(f)$ par translation dans le domaine fréquentiel du filtre prototype H :

$$H_k(f) = H\left(f - \frac{2k+1}{4M}\right) + H\left(f + \frac{2k+1}{4M}\right)$$

$$h_k(n) = 2h(n) \cos\left(2\pi \frac{2k+1}{4M}n\right)$$

On montre que la condition de reconstruction parfaite est vérifiée.

En pratique...

On ne peut pas construire un filtre passe-bas idéal de longueur N finie.
 $\Rightarrow h(n)$ est une approximation d'un passe-bas idéal,
 \Rightarrow la reconstruction n'est pas exactement parfaite
 \Rightarrow MAIS les distorsions introduites par les bancs d'analyse et de synthèse sont négligeables devant celles introduites par la quantification.

Banc de filtres modulés (2)

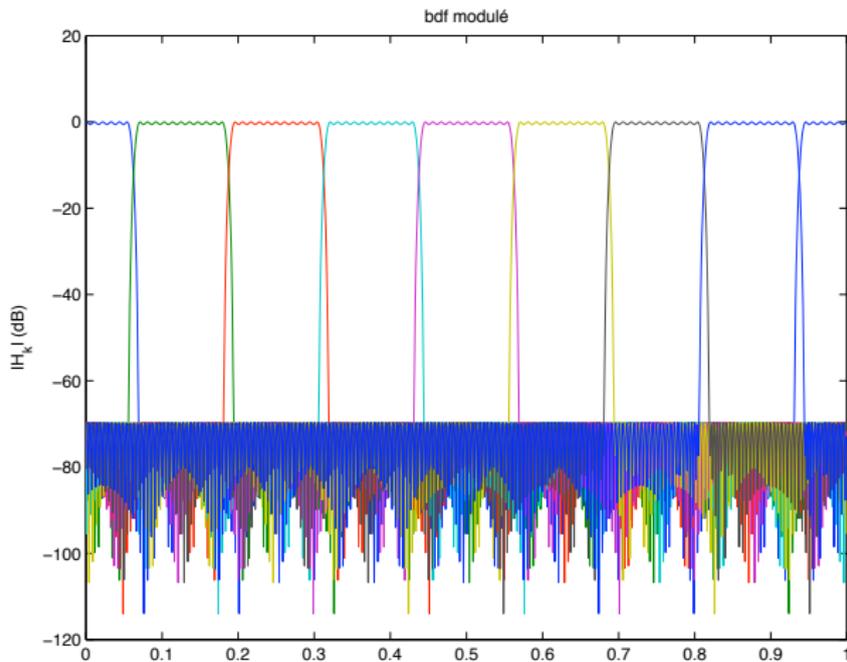


Schéma général d'un codeur perceptif

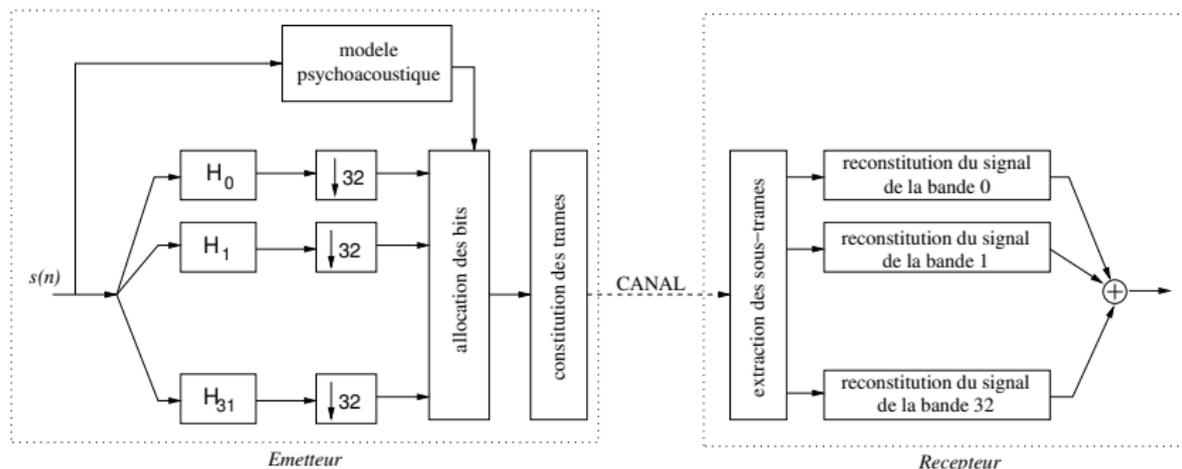
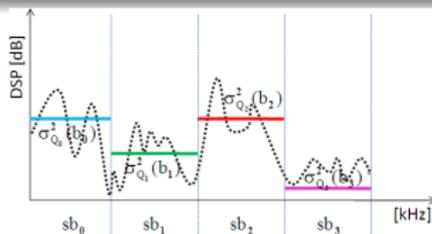


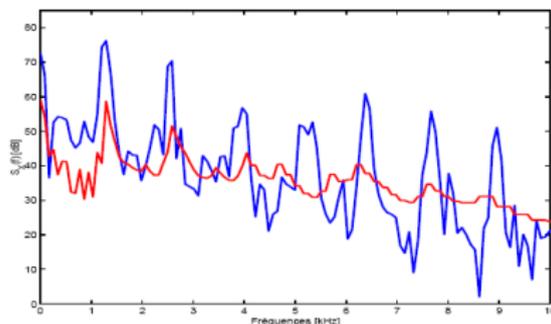
FIGURE – Schéma de principe des codecs (codeur-décodeur) perceptuels en sous-bandes

Comment allouer les bits ? (1)

- 0 bit sur les sous-bandes où $S_x(f) < M(f)$
- Sur les autres bandes, nombre de bits de quantification selon le rapport entre $S_x(f)$ et $M(f)$

Sur cet exemple :

suppression de $\approx 50\%$ des composantes
grâce au masquage fréquentiel
→ taux de compression élevé



Comment allouer les bits ? (2)

2 contraintes pour fixer le nombre de bit k_i par échantillon sur la bande i :

- 1 $\forall i$, la dsp du bruit est constante et doit être inférieure au seuil de masquage :
 $Q_i^{\text{dB}} < M_i^{\text{dB}}$, avec

- $Q_i^{\text{dB}} = S_i^{\text{dB}} - RSB_i^{\text{dB}}$
- et $RSB_i^{\text{dB}} = 6k_i - 4,76 + f_c$

$$\blacktriangleright \forall i, k_i > \frac{S_i^{\text{dB}} - M_i^{\text{dB}} - 4,76 + f_c}{6}$$

- 2 $\sum_i k_i <$ une valeur dépendant du débit binaire fixé

Comment allouer les bits ? (2)

2 contraintes pour fixer le nombre de bit k_i par échantillon sur la bande i :

- 1 $\forall i$, la dsp du bruit est constante et doit être inférieure au seuil de masquage :

$Q_i^{\text{dB}} < M_i^{\text{dB}}$, avec

- $Q_i^{\text{dB}} = S_i^{\text{dB}} - RSB_i^{\text{dB}}$
- et $RSB_i^{\text{dB}} = 6k_i - 4,76 + f_c$

$$\blacktriangleright \forall i, \quad k_i > \frac{S_i^{\text{dB}} - M_i^{\text{dB}} - 4,76 + f_c}{6}$$

- 2 $\sum_i k_i <$ une valeur dépendant du débit binaire fixé

Si impossible, on minimise $\sum_i Q_i$ (puissance totale du bruit) sous la contrainte sur $\sum_i k_i$ liée au débit.

Plan

- 1 Introduction
- 2 Numérisation
- 3 Quantifications
- 4 Perception auditive
- 5 Codage perceptif
- 6 Codage entropique**
- 7 Quelques codeurs
 - Historique
 - Codeurs MPEG-1

Principes (1)

Principe

Dans tous les cas (quantification scalaire « de base, vectorielle, codage perceptif...), l'étape de quantification associe à chaque réel (ou bloc de réels) un entier $i \in \{0, 1, \dots, M - 1\}$. Comment représenter cet entier dans un flux binaire ?

Solution naïve

Code à longueur fixe : si $M = 2^b$, associer à i sa représentation binaire sur b bits.

Oui mais...

Le même nombre de bits est utilisé pour chaque valeur de i . On aimerait que les valeurs de i les plus courantes (les plus probables) soient codées avec moins de bits.

Comment construire un code à longueur variable décodable et minimisant le débit ?

Principes (2)

Hypothèses

- Source $X \rightarrow$ symboles $x_i \in \{x_1, \dots, x_m\}$ **alphabet fini**
- Source **sans mémoire**,
i.e. chaque symbole est indépendant des précédents.

Notion d'information

- **Information** portée par un symbole x_i :

$$I(x_i) = \log_2 \left(\frac{1}{P(x_i)} \right) = -\log_2 P(x_i)$$

Unité de mesure : le **bit**

- Information moyenne ou **entropie** :

$$H(X) = E[I(x_i)] = \sum_{i=1}^m P(x_i) I(x_i) = - \sum_{i=1}^m P(x_i) \log_2 P(x_i)$$

Notion d'information

- **Information** portée par un symbole x_i :

$$I(x_i) = \log_2 \left(\frac{1}{P(x_i)} \right) = -\log_2 P(x_i)$$

Unité de mesure : le **bit**

- Information moyenne ou **entropie** :

$$H(X) = E[I(x_i)] = \sum_{i=1}^m P(x_i) I(x_i) = - \sum_{i=1}^m P(x_i) \log_2 P(x_i)$$

Soit une source X générant m symboles.

- $0 \leq H(X) \leq \log_2 m$
- Si $\exists x_i$ tel que $P(x_i) = 1$, alors $H(X) = 0$
- Si les m symboles sont équiprobables, ...

Décodabilité

Chaque symbole $x_i \rightarrow$ mot binaire de n_i éléments binaires

► Comment segmenter le flux binaire en une suite de symboles ?

Si le code est un **code préfixe**,
alors le décodage est

- **unique**
- **instantané**

Ex : $x_1 = 0, x_2 = 10, x_3 = 101, x_4 = 110, x_5 = 1110, x_6 = 1111$

Ce n'est pas une condition nécessaire !

Ex : $x_1 = 0, x_2 = 01, x_3 = 011, x_4 = 011$

Comment vérifier la condition du préfixe ?

Efficacité d'un code

- **Objectif** : minimiser la longueur moyenne L des mots de code.

$$L = \sum_{i=1}^m P(x_i) n_i$$

- **Théorème du codage** : $L \geq H(X)$ et il existe un code à décodage unique et instantané tel que :

$$H(X) \leq L < H(X) + 1$$

- **Efficacité** du code : $\eta = H(X)/L$
= nombre de bits d'information par élément binaire

Ex : codage de Huffman

Algorithme :

- Initialisation : chaque symbole = 1 groupe
- Répéter
 - ① ordonner les groupes de symboles selon leurs probabilités
 - ② fusionner les 2 groupes les moins probables en 1 seul.
 - ▶ probabilité du groupe = somme des 2 probabilités.

Tant que nombre de groupes > 1

Ex : codage de Huffman

Algorithme :

- Initialisation : chaque symbole = 1 groupe
- Répéter
 - ① ordonner les groupes de symboles selon leurs probabilités
 - ② fusionner les 2 groupes les moins probables en 1 seul.
 - ▶ probabilité du groupe = somme des 2 probabilités.

Tant que nombre de groupes > 1

- ▶ Arbre binaire avec 1 symbole sur chaque feuille.

Ex : codage de Huffman

Algorithme :

- Initialisation : chaque symbole = 1 groupe
- Répéter
 - ① ordonner les groupes de symboles selon leurs probabilités
 - ② fusionner les 2 groupes les moins probables en 1 seul.
 - ▶ probabilité du groupe = somme des 2 probabilités.

Tant que nombre de groupes > 1

- ▶ Arbre binaire avec 1 symbole sur chaque feuille.

On affecte 1 aux branches hautes, 0 aux branches basses

- ▶ code préfixe par construction

Plan

- 1 Introduction
- 2 Numérisation
- 3 Quantifications
- 4 Perception auditive
- 5 Codage perceptif
- 6 Codage entropique
- 7 Quelques codeurs**
 - Historique
 - Codeurs MPEG-1

MPEG-1 Audio (1991)

Généralités sur la norme MPEG-1

- Norme internationale ISO 11172.
- Codage de l'image animée et du son associé pour les supports de stockage jusqu'à 1.5 Mbit/s.
- Norme pour la vidéo, l'audio, et les tests de conformité.
- Normalise le décodeur mais pas le codeur.
- fréquences d'échantillonnage audio : 48kHz, 44.1kHz, 32kHz.

MPEG-1 Audio : Couches de complexité croissante

Couche 1 : Studio, diffusion par satellite.

- Débit recommandé : 192 kbit/s.
- Retard théorique de codage/décodage : 19ms.

Couche 2 : DAB (radio numérique). Diffusion.

- Débit recommandé : 128 kbit/s.
- Retard théorique de codage/décodage : 35ms.

Couche 3 : "MP3"

- Débit recommandé : 96 kbit/s.
- Retard théorique de codage/décodage : 59ms.

MPEG-2 Audio (1993)

Extension de MPEG-1 (TV numérique, DVD, vidéoconférence)

- Support des basses fréquences d'échantillonnage (16, 22.05 et 24 kHz).
- Support du multicanal (5.1).

MPEG-2 AAC (Advanced Audio Coder)

Evolution de la couche 3. Popularisé par Apple (iTunes, iPod).

- Transparent à 64kbit/s/ch.
- Longues fenêtres pour les parties stationnaires, fenêtres courtes pour les parties transitoires.
- Plus large gamme de fréquences d'échantillonnage (de 8 à 96 kHz).
- Multicanaux (de 1 à 48 canaux).

MPEG-4 (1999)

Notion de scène sonore

Approche "Orientée objet" de la description de contenus sonores.

- Représentation de sources diverses : parole et musique naturelles et codées, ou synthétiques.
- Sources combinées pour former des "scènes sonores".

Famille de codeurs parole/musique

- Parole de 6 à 24 kbits/s : codeur UIT G.729.
- Parole ou musique de 2 à 24 kbit/s : codeur paramétrique (HVXC, HILN).
- Musique de 16 à 64 kbit/s : MPEG-2 AAC.

MPEG-4 (2)

Génération de sons synthétiques

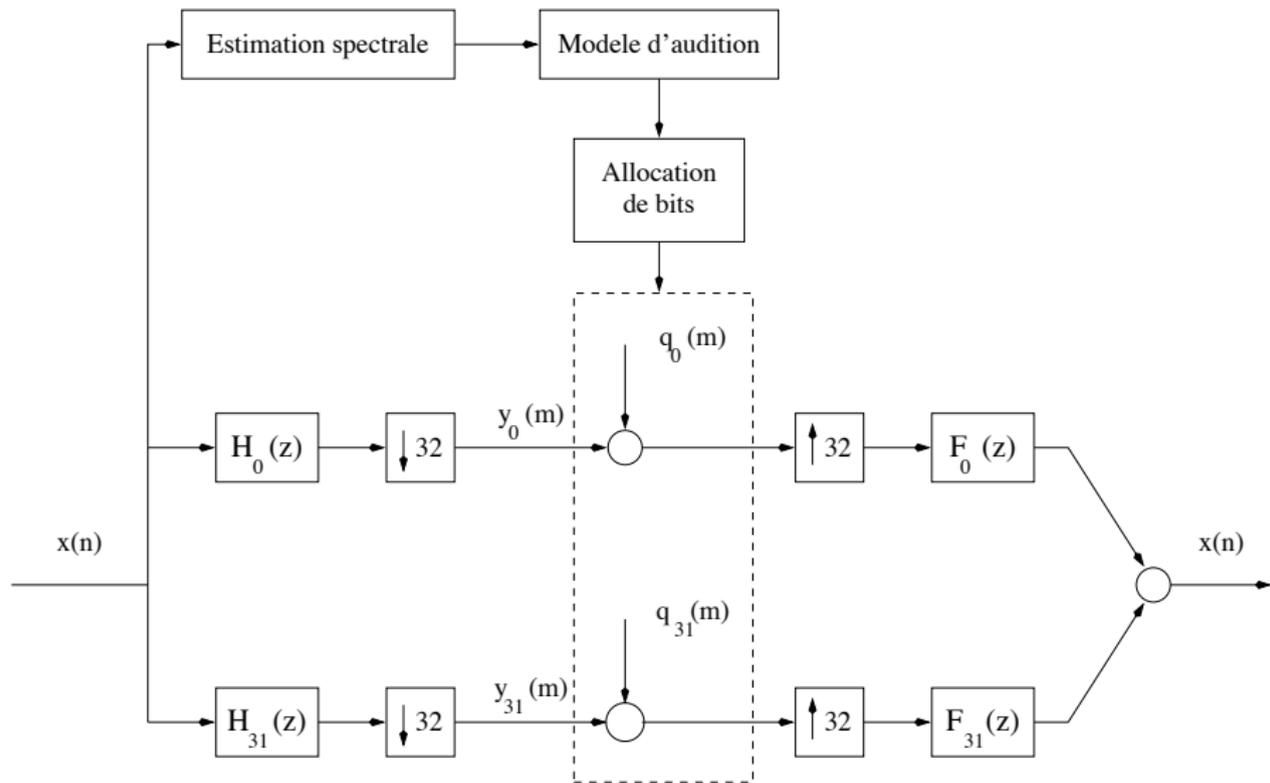
- Synthèse vocale.
- Langage de représentation de musique et de synthèse Structured Audio Orchestra Language, Structured Audio Sample Bank (type MIDI + SoundFont).
<http://web.media.mit.edu/~eds/mpeg4/>

⇒ Ce ne sont pas des codeurs. Il faut que le contenu ait été préparé (à partir du texte prononcé ou de la partition).

Scènes sonores

- Localisation de chaque source dans l'espace.
- Paramètres de modification des sources.
- Spatialisation et réverbération.

MPEG-1 Couche 1 : principe



MPEG-1 Couche 1 : transformation

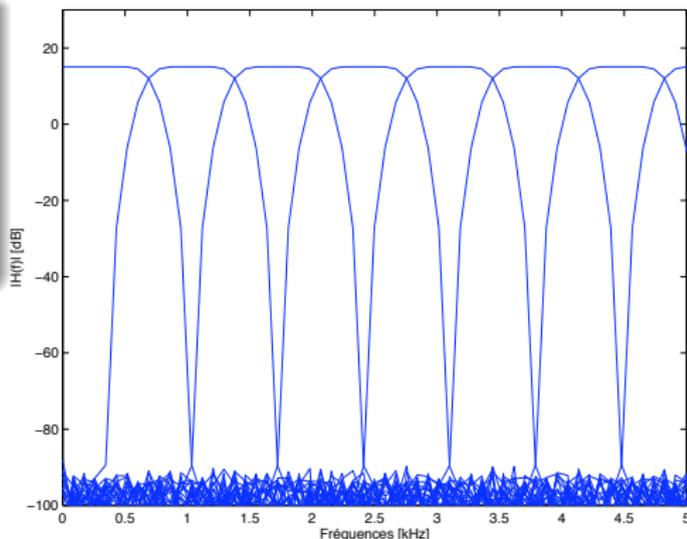
Banc de $M = 32$ filtres modulés

$$H_k(f) = H(f - \frac{2k+1}{4M}) + H(f + \frac{2k+1}{4M}), h_k(n) = 2h(n)\cos(2\pi\frac{2k+1}{4M}n + \phi_k).$$

Filtres de longueur $N = 512$.

Décimation

- Sous-échantillonnage de chaque bande par un facteur M .
- $M \ll N$, donc bonne résolution temporelle mais mauvaise résolution fréquentielle.



MPEG-1 Couche 1 : quantification (1)

Constitution des trames :

- Sur chaque voie i , échantillons groupés par 12 :

$$x_i = [x_i(0), x_i(1) \dots x_i(11)]$$

- Normalisation :

$$x'_i = x_i / \arg \min(|x_i|) \in [-1, 1]$$

- Chaque $x'_i(k)$ codé sur n_i e.b. selon le modèle psycho-acoustique

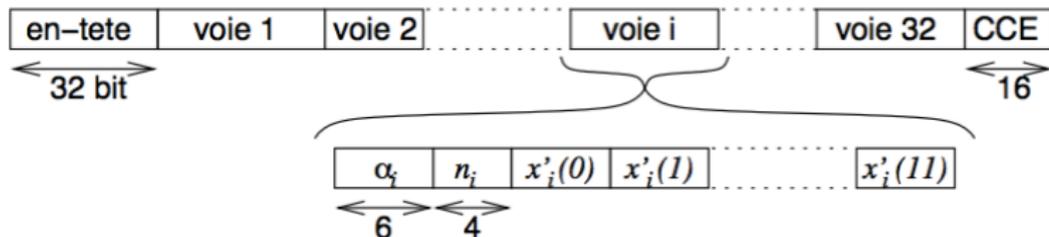
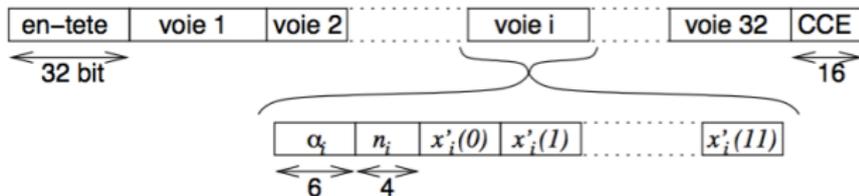


FIGURE – Trame MPEG-1 (correspondant à 384 échantillons de signal).

MPEG-1 Couche 1 : quantification (2)



Débits

Si l'on souhaite obtenir un débit de 96 kbit/s pour un signal à 44.1kHz,

- on dispose pour chaque bloc de 384 échantillons de :

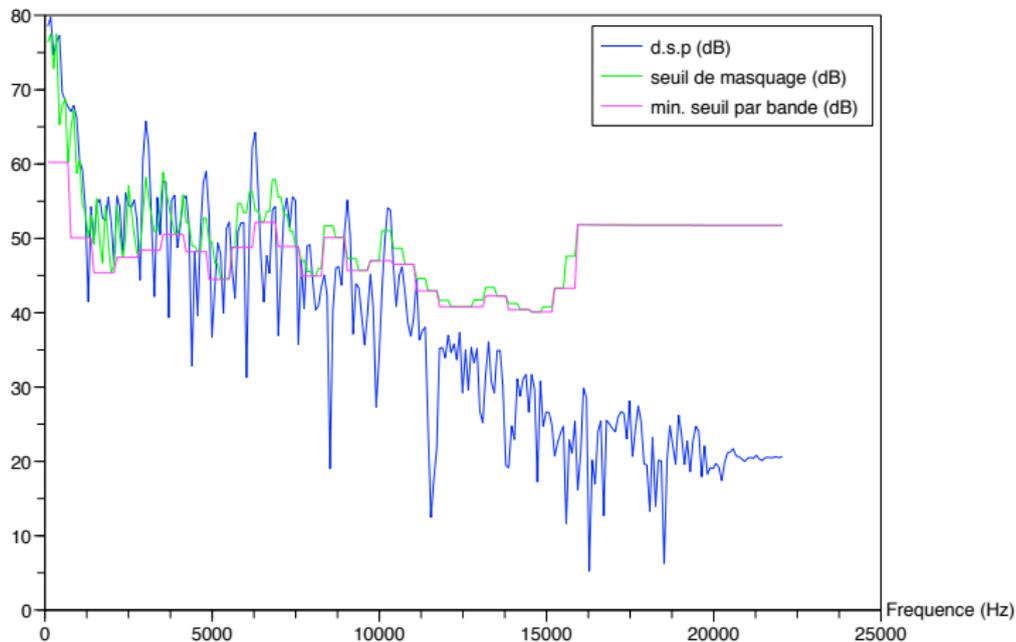
$$\frac{96000 \text{ bit/s}}{44100 \text{ ech/s} / 384 \text{ ech/trame}} = 828 \text{ bits/trame}$$

- 32 bits de tête + 16 bits de queue + (6+4) bits pour α_i et n_i dans chaque sous-trame = 368 bits.
- Il reste 460 bits pour coder les 384 coefficients

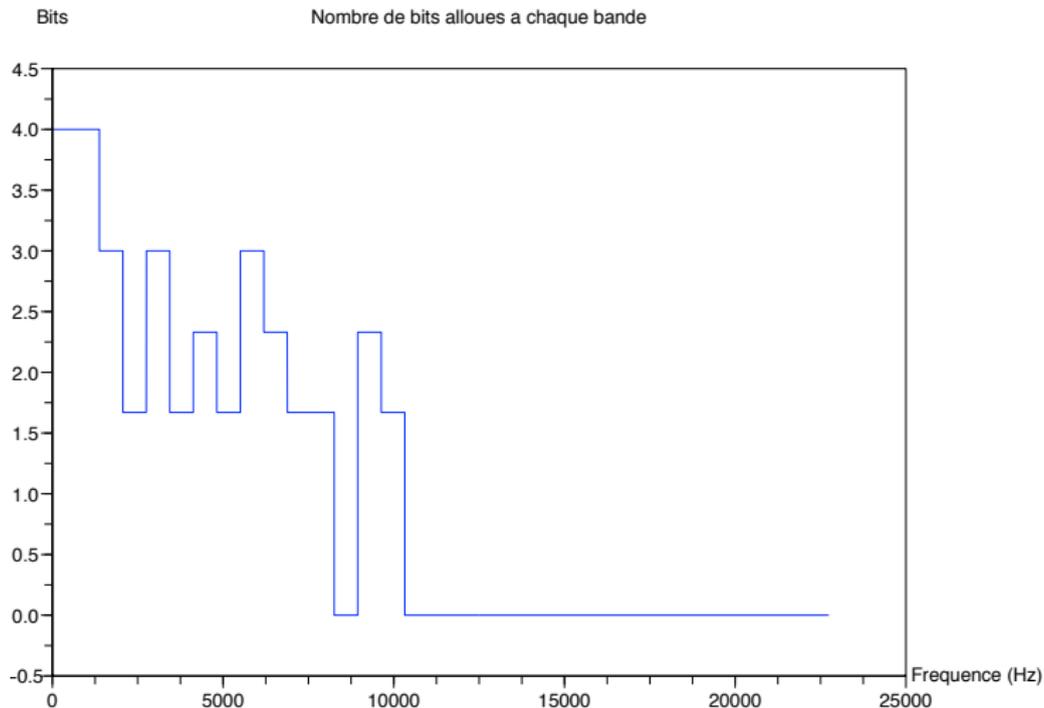
MPEG-1 Couche 1 : quantification (3)

Amplitude (dB)

dsp et seuils de masquage



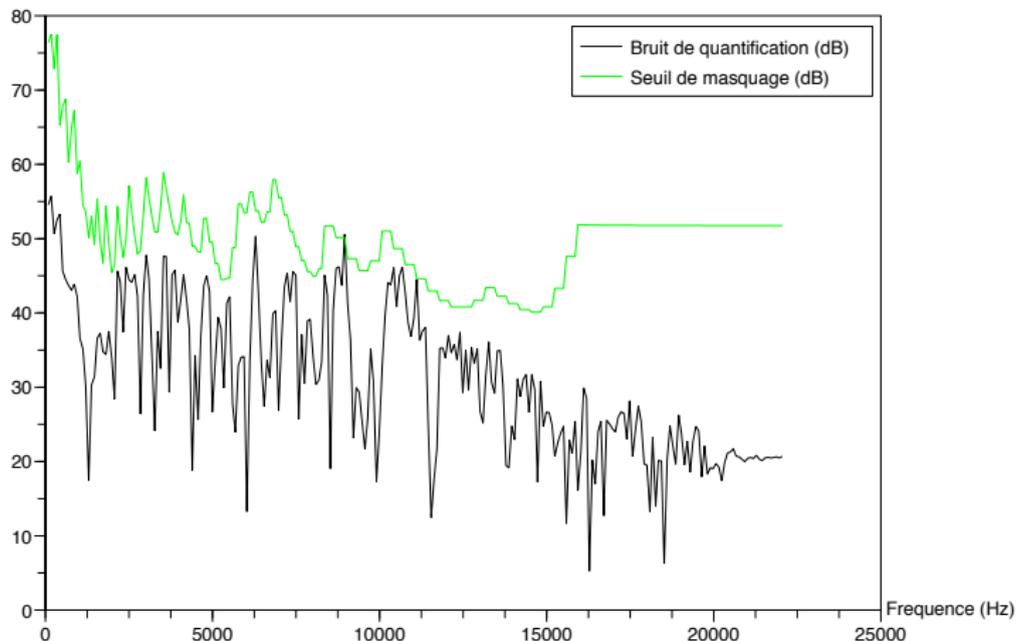
MPEG-1 Couche 1 : quantification (3)



MPEG-1 Couche 1 : quantification (3)

Amplitude (dB)

Effet de la quantification



MPEG-1 Couche 2

Amélioration

- Modèle psychoacoustique plus fin.
- Autorise plus de pas de quantification par sous-bandes.
- Ne transmet pas les facteurs d'échelle quand ceux-ci ne changent pas.

MPEG-1 Couche 3 : Principe général

Spécificités du codeur

- Décomposition hybride banc de filtres + MDCT (La MDCT est appliquée aux signaux de sous-bande).
- Taille des fenêtres variable (courtes ou longues).
- Quantification non uniforme.
- Elimination de la redondance (Codage de Huffman) lors du codage des coefficients.
- "Réservoir de bits" : permet d'allouer plus de bits sur les transitoires, et de compenser en allouant moins de bits sur les parties stationnaires.

MPEG-1 Couche 3 : MDCT par bandes

Décomposition en 32 bandes

Même banc de filtres que pour la couche 1.

Dans chaque bande

- 1 Groupement des échantillons en fenêtres de 12 ou 36.
- 2 Transformée en cosinus discrète modifiée de la fenêtre

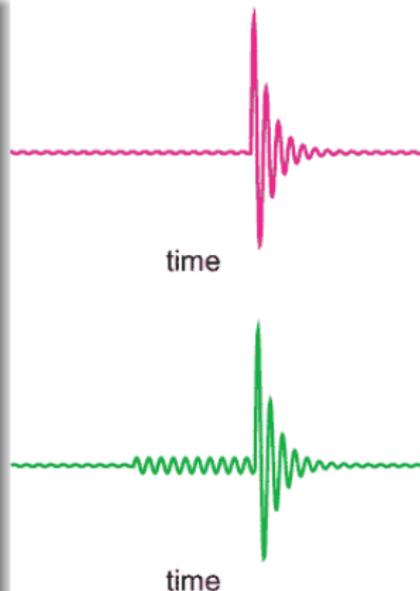
$$X_i = \sum_{k=0}^{2n-1} x_k \cos \left[\frac{\pi}{n} \left(i + \frac{1}{2} \right) \left(k + \frac{n+1}{2} \right) \right]$$

- 3 \Rightarrow Blocs de 576 ou 192 coefficients X_i .

Pourquoi une taille de fenêtre variable ?

Phénomène de pré-echo

- 1 Instrument en fond à niveau faible.
- 2 Une nouvelle note à l'attaque très forte se présente en fin de fenêtre.
 - ⇒ Beaucoup de bits sont alloués pour la gamme de fréquence occupée par cette note.
 - ⇒ Peu de bits sont alloués pour la gamme de fréquence occupée par l'instrument en fond.
 - ⇒ Augmentation du bruit de fond sur l'ensemble de la fenêtre.
 - ⇒ Effet "d'aspiration" avant l'attaque.



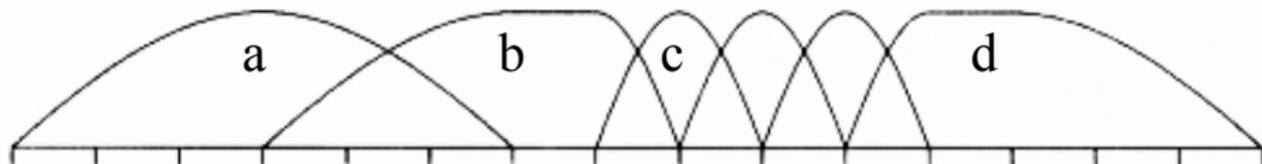
Solution

Compromis temps/fréquence

Si l'on se contente de prendre des petites fenêtres (comme c'est le cas avec la couche 1), la décomposition fréquentielle n'est pas fine et on ne peut pas finement allouer des bits.

"Window switching"

Utiliser une fenêtre longue (36 coefficients), et passer à une fenêtre courte (12 coefficients) lorsqu'un transitoire est détecté.



MPEG-1 Couche 3 : Quantification et codage

Quantification non uniforme

- $q(X_i) = \text{sign}(X_i) \text{round} \left[\left(\frac{|X_i|}{b_i} \right)^{\frac{3}{4}} - 0.0946 \right]$: adaptée à la distribution des X_i .

Classement des coefficients par fréquence croissante

- Basses fréquences : quelques coefficients, valeurs fortes.
- Mediums : beaucoup de coefficients, valeurs proches de zéro.
- Hautes fréquences : valeurs presque toujours nulles.

⇒ Partition des coefficients en 3 régions

- Région des zéros : non codée.
- Région des faibles valeurs : groupées en quadruplets.
- Région des grandes valeurs : groupées par paires, utilisation de plusieurs tables de Huffman adaptées.

MPEG-1 Couche 3 : Différents codeurs

Pas de norme pour le codeur !

Ne sont spécifiés que le format du flux binaire et le décodeur.

Conséquence

Les codeurs peuvent offrir différents compromis complexité/qualité :

- Pour une meilleure qualité : optimisation des facteurs d'échelle plus avancée, avec plus d'itérations.
- Au contraire, un codeur peut ne pas chercher à optimiser la qualité et n'effectuer que l'optimisation de débit en utilisant des facteurs d'échelle par défaut. Codeurs "sales" à faible complexité (Baladeur MP3 avec mémo vocal, PDA).