



Ecole Polytechnique
Electronique, composants & systèmes

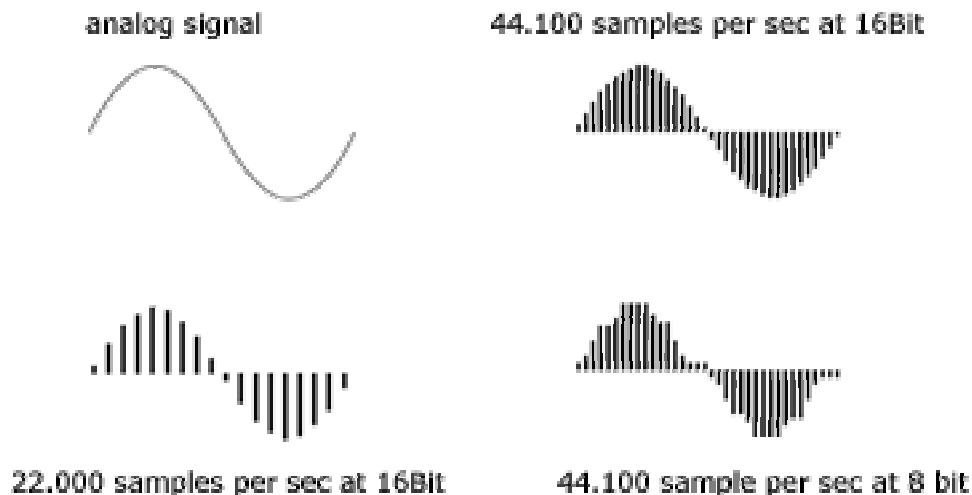
Le son numérique

La compression audio
MPEG Layer 3

Guillaume Rincé
X96

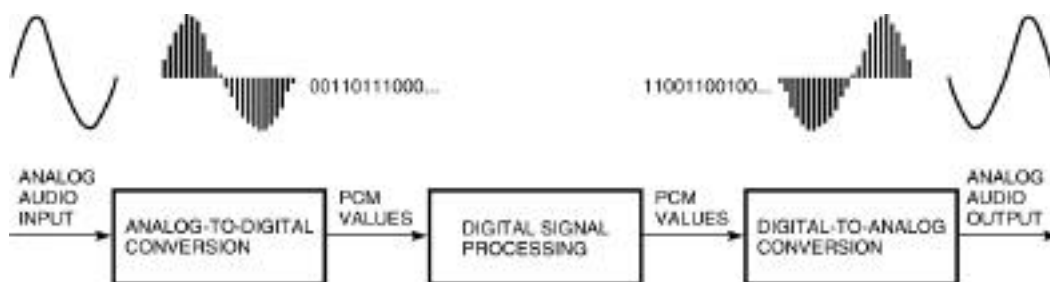
I Le principe de la conversion numérique

La conversion numérique consiste à convertir un signal analogique en une série de bits. Cette conversion consiste à prélever un certain nombre d'échantillon à une fréquence dite « fréquence d'échantillonnage » puis à les coder sur un certain nombre de bits. Selon le résultat souhaité, on définit la fréquence d'échantillonnage et le nombre de bits. Le schéma ci-dessous illustre l'influence de ces deux facteurs.



Une chaîne complète de traitement numérique comporte :

- Conversion analogique / numérique
- Un éventuel traitement des données numériques (stockage, filtrage, ...)
- Conversion numérique / analogique

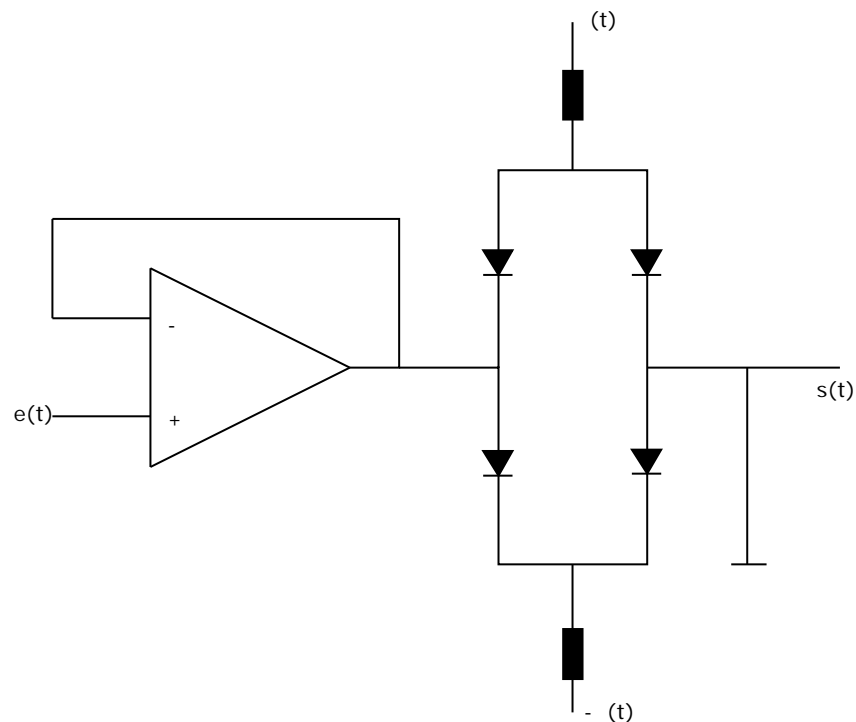


I.1 L'échantillonnage du signal

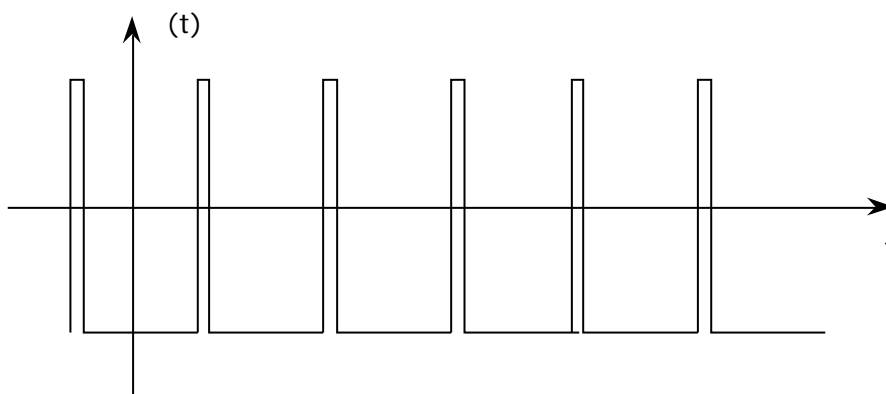
L'échantillonnage du signal analogique est une étape très importante de la conversion numérique car le choix de la fréquence d'échantillonnage détermine la qualité du son converti. L'oreille humaine perçoit les sons d'une fréquence inférieure à 20 kHz. Par conséquent, à partir d'une fréquence d'échantillonnage supérieure à 40 kHz – le facteur 2 provient des problèmes de repliement traités

plus loin – toutes les fréquences audibles seront restituées. Si l'on veut coder la voix humaine dont les fréquences sont inférieures à 8 kHz, il faudrait en théorie échantillonner à 16 kHz. En réalité, la voix est compréhensible même si on détruit les hautes fréquences, c'est pourquoi l'échantillonnage de qualité téléphonique s'effectue en général à 8 kHz.

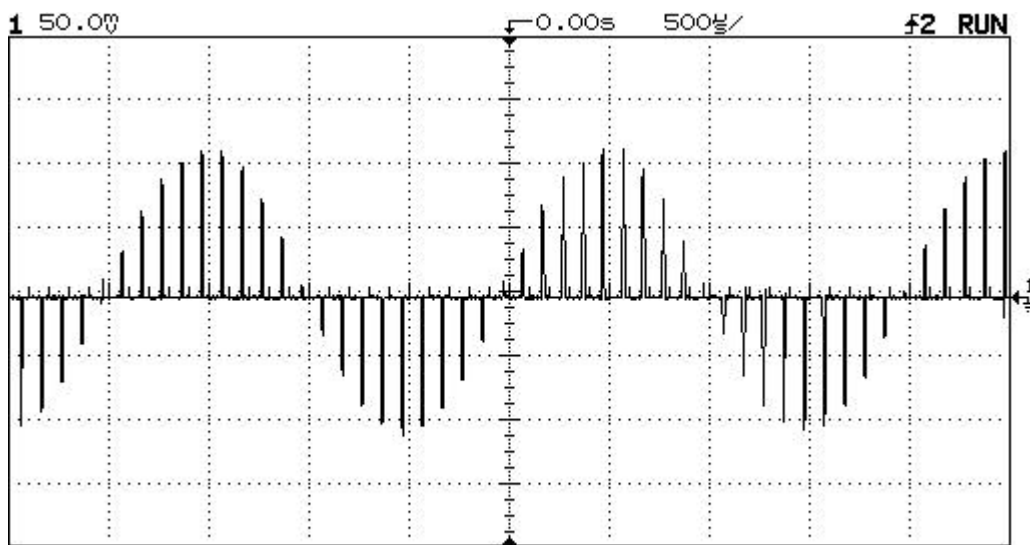
Le principe de l'échantillonnage est de prendre des échantillons du signal analogique à la fréquence d'échantillonnage. Le dispositif ci-dessous permet d'échantillonner le signal analogique $e(t)$.



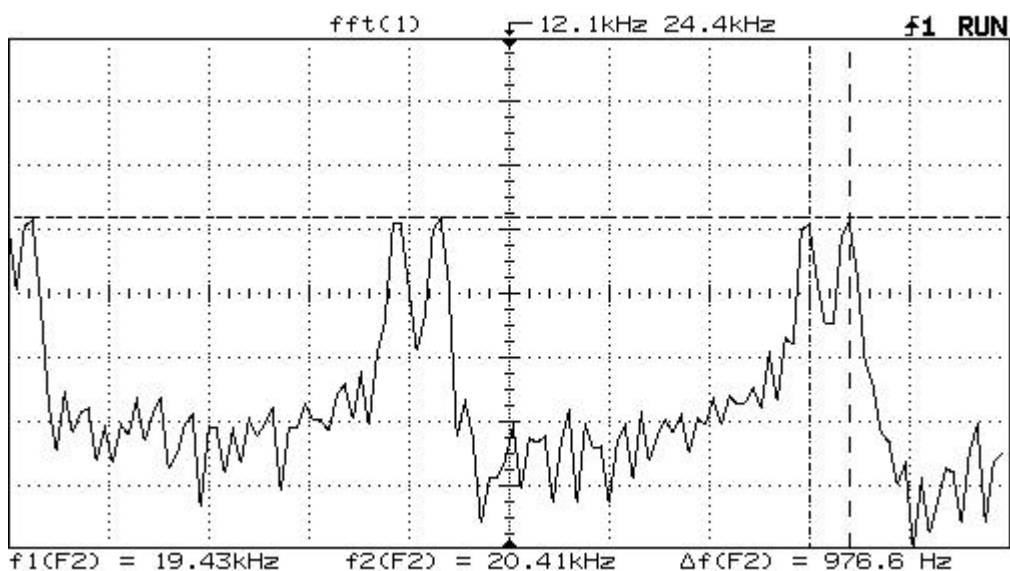
Le signal (t) est un signal rectangulaire de fréquence la fréquence d'échantillonnage. Plus la largeur des raies est faible, plus l'échantillonnage sera de bonne qualité.



On prend pour signal d'entrée $e(t)$ un signal sinusoïdal à 500 Hz. Le signal $s(t)$ d'échantillonnage a une fréquence de 10 kHz. En sortie $s(t)$ on obtient le signal ci-dessous.



On échantillonne à une fréquence 20 fois plus rapide que la fréquence du signal d'entrée, on a donc 20 échantillons par période du signal d'entrée. Il est aussi intéressant de regarder la transformée de Fourier du signal de sortie.



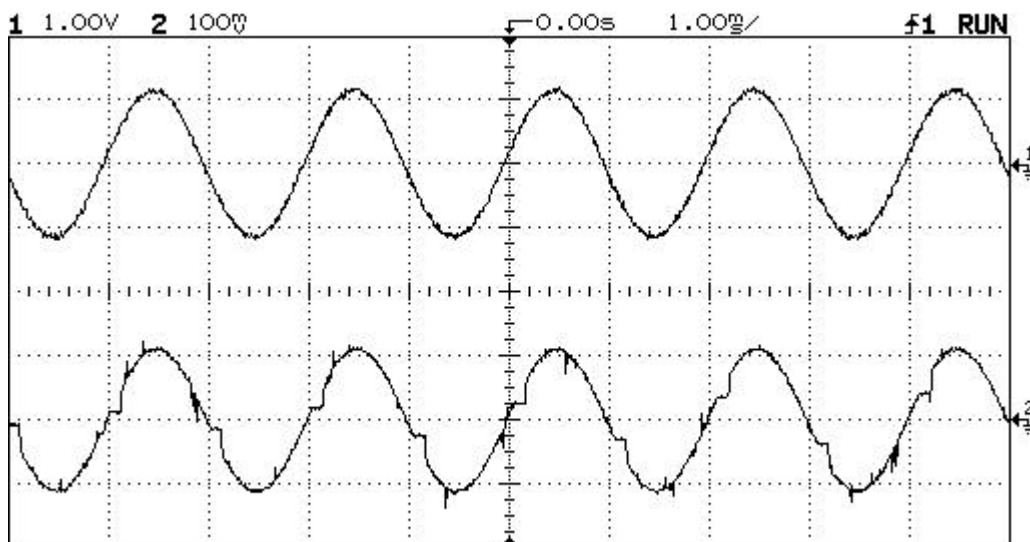
Si l'on visualisait la transformée de Fourier du signal d'entrée, on trouverait une raie pure à la fréquence de ce signal, 500 Hz. Dans la transformée du signal de sortie, on retrouve tous les 10 kHz, la fréquence d'échantillonnage, deux pics situés relativement par rapport aux multiples de la fréquence d'échantillonnage à -500 Hz et à +500 Hz, la fréquence du signal d'entrée.

Cette transformée de Fourier permet de mettre en évidence les deux problèmes posés par l'échantillonnage. Tout d'abord, on voit ici l'intérêt du facteur 2 qu'il est nécessaire d'avoir entre la fréquence maximale du signal d'entrée et la fréquence d'échantillonnage. En effet, si on ne respecte pas cette règle, les pics autour de deux multiples consécutifs de la fréquence d'échantillonnage se recouvrent. Pour reconstruire le signal d'origine, il faut pouvoir filtrer le signal échantillonné pour isoler une série de pics. Or, il n'existe pas de filtre parfait, il faut donc qu'il y ait un certain espace entre la fin d'une série de pic et la suivante pour que le filtre puisse couper. Par exemple, dans le cas de l'échantillonnage qualité CD, on échantillonne à 44,1 kHz, ce qui laisse une plage de 2 kHz entre deux séries de pics consécutifs.

1.2 Le blocage

Le blocage consiste à transformer le signal non pas en signal échantillonné, mais en signal crénelé. En fait, entre deux prises d'un échantillon, le signal bloqué reste à la valeur de l'échantillon qui vient d'être pris.

Cette fonction d'échantillonnage blocage peut être réalisée à l'aide du convertisseur analogique numérique AD 1380.



Le signal sinusoïdal de l'entrée se trouve transformé en signal crénelé, la largeur de chaque créneau dépendant de la fréquence d'échantillonnage choisi. Le signal ainsi obtenu est maintenant prêt à être quantifié et codé.

1.3 La quantification et le codage

La quantification du signal est l'étape qui consiste à quantifier l'amplitude de chaque échantillon. Il y a plusieurs manières de procéder, on peut en particulier choisir des pas de quantification fixes ou variables. Nous nous limiterons ici à des pas fixes, c'est à dire à une quantification dite « linéaire ».

Le codage binaire consiste à coder l'amplitude quantifiée en un nombre binaire. La relation suivante permet de déterminer le nombre de bits nécessaires au codage. (On note S le rapport signal sur bruit en décibels).

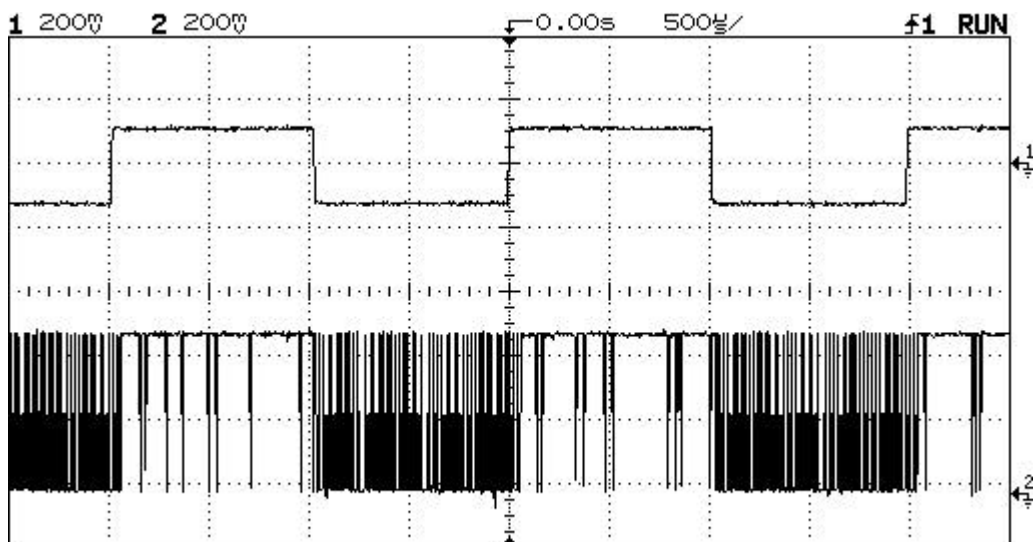
$$S = 20 \log_0 2 N = 6,02N$$

Il y a plusieurs manières de coder en binaire, le tableau suivant décrit les trois manières les plus utilisés.

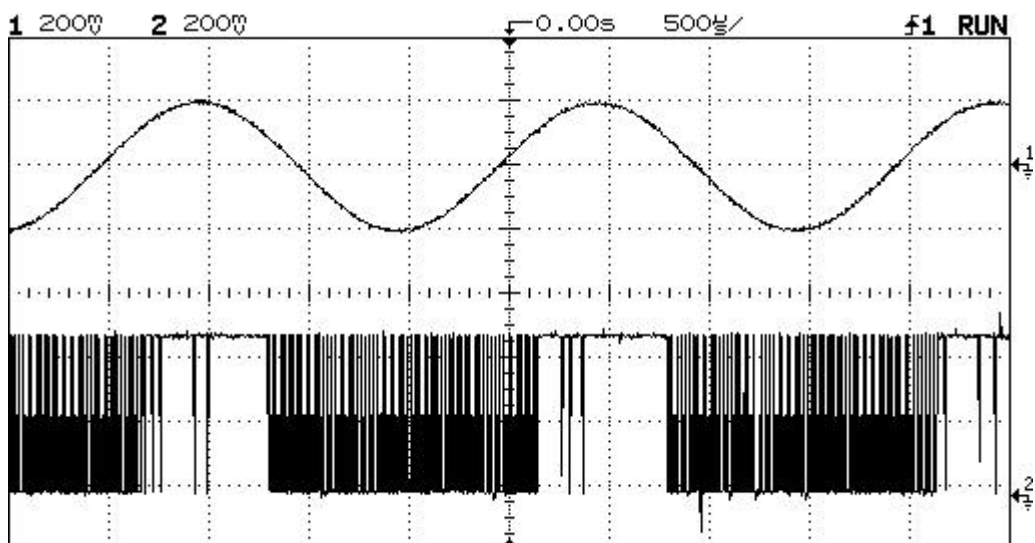
Code binaire naturel		Code binaire avec offset		Code avec complément à 2	
7	111	3	111	3	011
6	110	2	110	2	010
5	101	1	101	1	001
4	100	0	100	0	000
3	011	-1	011	-1	111
2	010	-2	010	-2	110
1	001	-3	001	-3	101
0	000	-4	000	-4	100

Le troisième code, le code complément à 2, est le plus utilisé, c'est en particulier celui du CD. Il consiste à utiliser le MSB – Most Significant Byte – pour coder le signe. La valeur absolue est alors codée en code binaire naturel si la valeur est positive (MSB = 0) ou en code binaire complément à 2^n-1 si la valeur est négative (MSB = 1).

La courbe suivante montre le MSB d'un signal créneau à 500 Hz échantillonné à 40 kHz codé par le convertisseur analogique numérique AD 1380 (L'AD 1380 code en complément à 2, le MSB est donc le bit de signe).



La courbe suivante montre le MSB d'un signal sinusoïdal à 500 Hz échantillonné à 40 kHz codé par le convertisseur analogique numérique AD 1380.



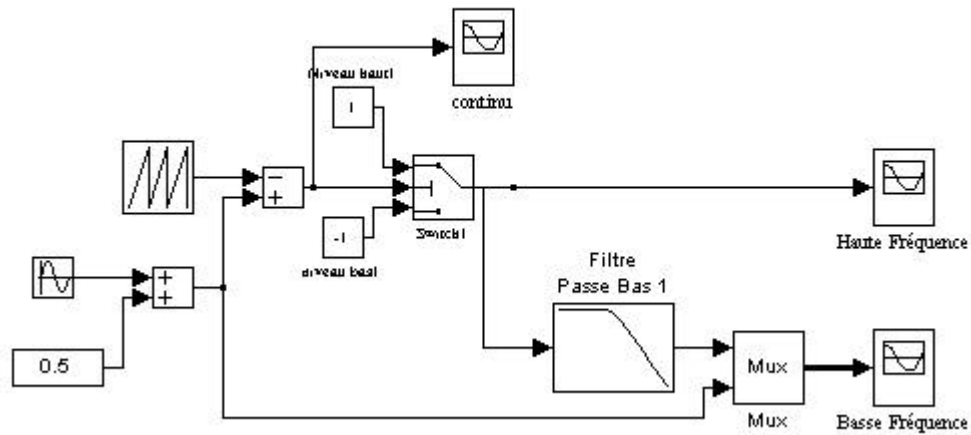
On peut remarquer que le MSB est plus souvent à 0 qu'à 1. Ceci est dû au fait que le signal sinusoïdal d'entrée n'est pas bien centré en amplitude autour de zéro.

1.4 Le décodage

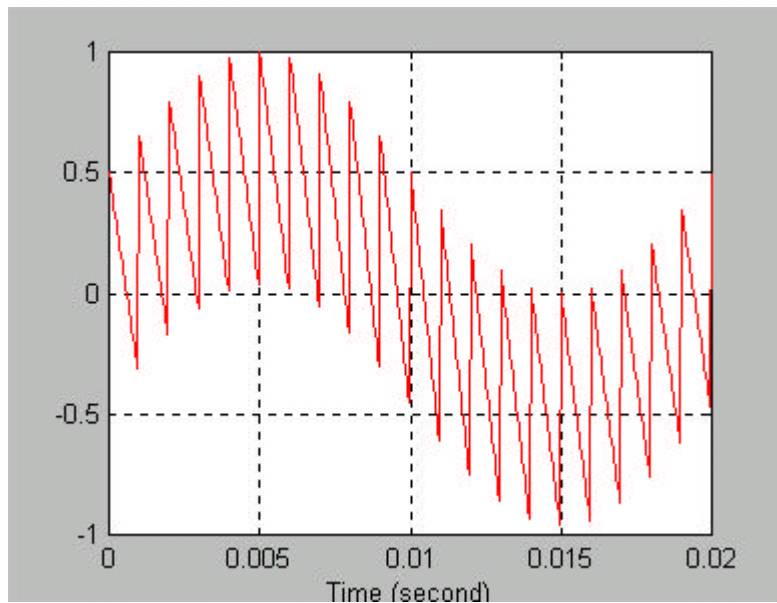
Le décodage est l'étape de restitution du signal analogique à partir du signal numérique codé. Il existe plusieurs méthodes, l'une des plus utilisées est celle de la conversion « one-bit ». Une simulation sous MatLab permet d'étudier

le fonctionnement d'un convertisseur numérique analogique de ce type. Le schéma sous Simulink est donné ci-dessous.

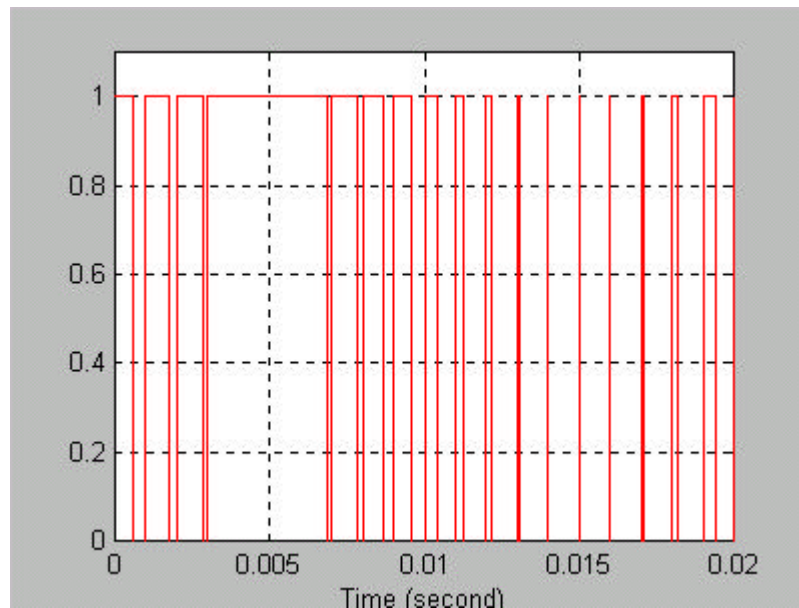
Principe des convertisseurs One-Bit



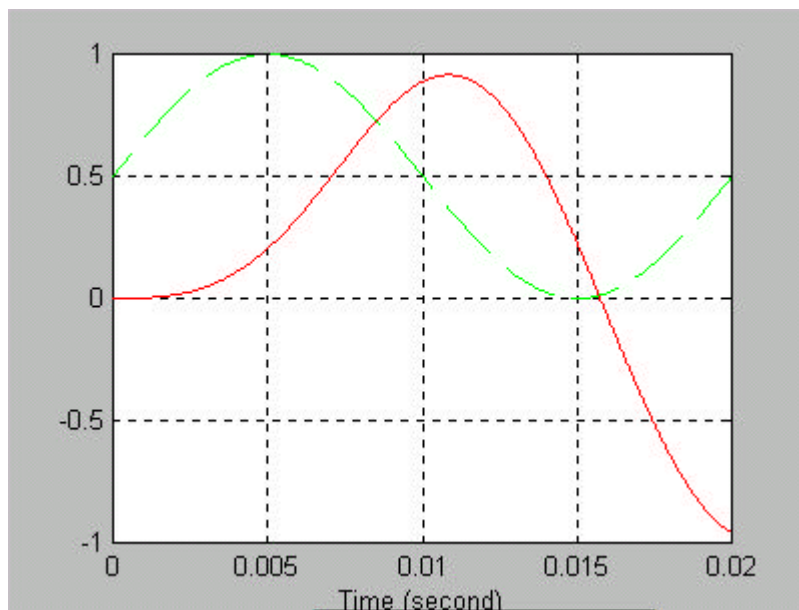
Le signal en entrée du switch est représenté par la courbe suivante :



Le signal à la sortie haute fréquence est représenté par la courbe suivante :

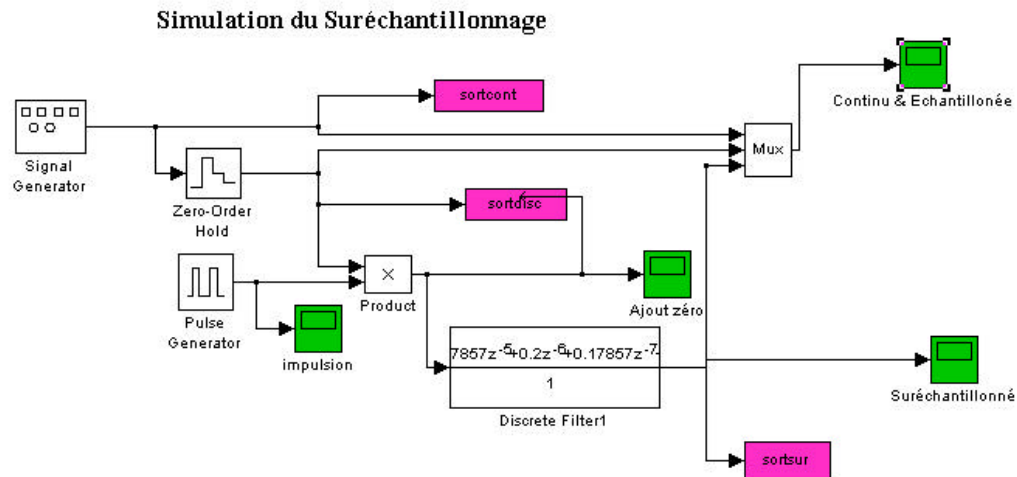


Le signal d'entrée (vert) et le signal de sortie basse fréquence (rouge) sont représentés sur le graphique ci-dessous :

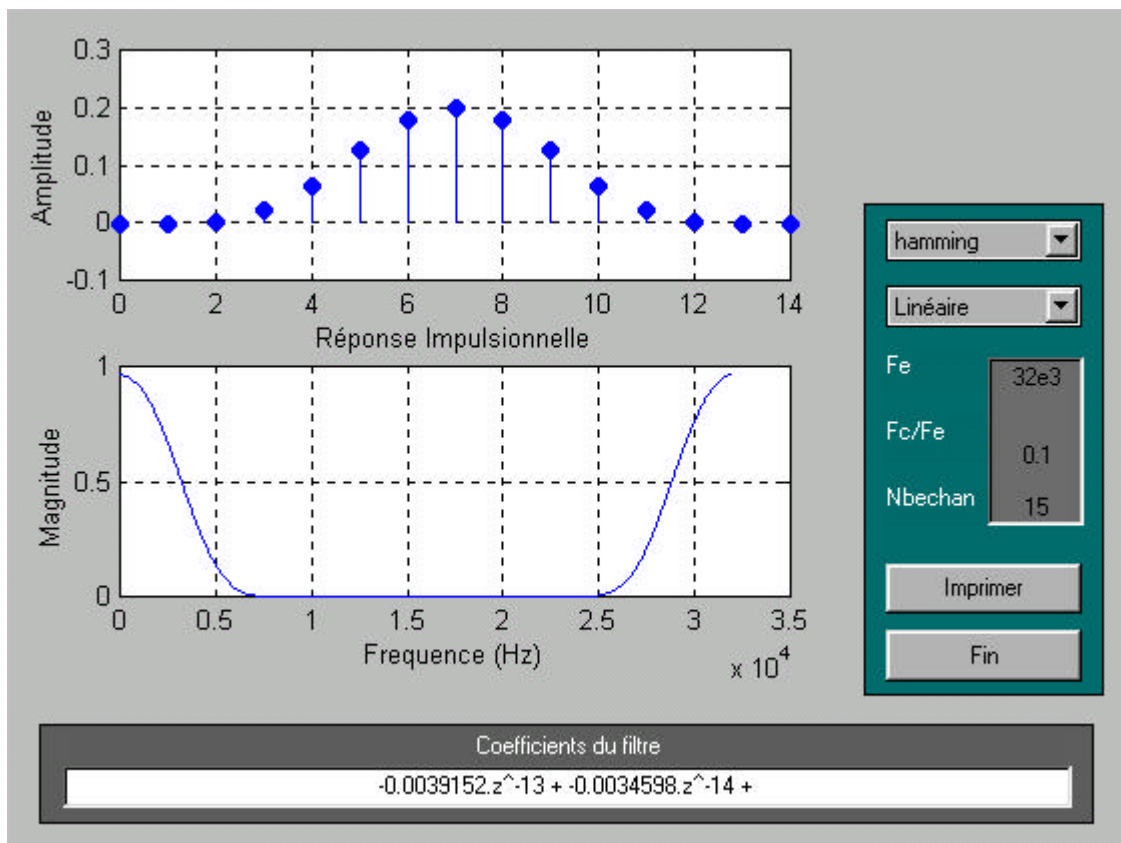


1.5 Le suréchantillonnage

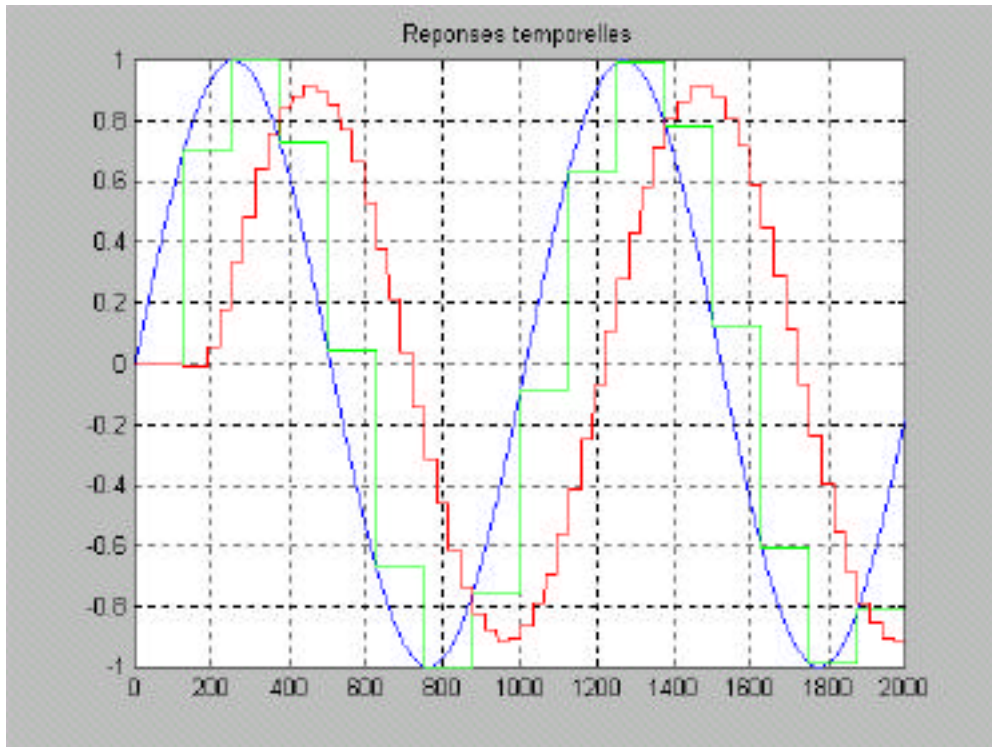
Le suréchantillonnage utilise le principe de Shannon et permet de retrouver les valeurs du signal de départ entre deux échantillons du signal échantillonné. Une simulation sous MatLab permet d'illustrer le suréchantillonnage. Le schéma sous Simulink de la simulation est le suivant :



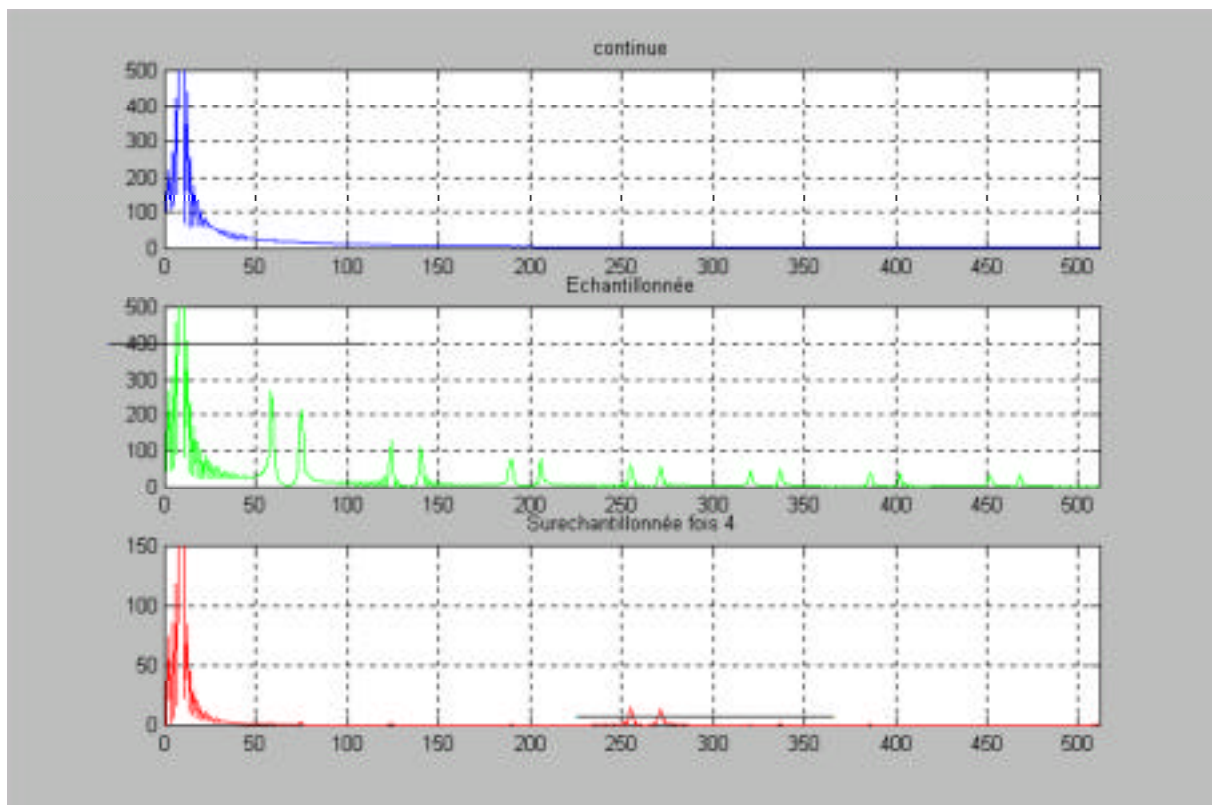
Les caractéristiques du filtre utilisé sont les suivantes :



Le graphique ci-dessous représente le signal d'entrée (bleu), le signal échantillonné (vert) et le signal surechantillonné (rouge).

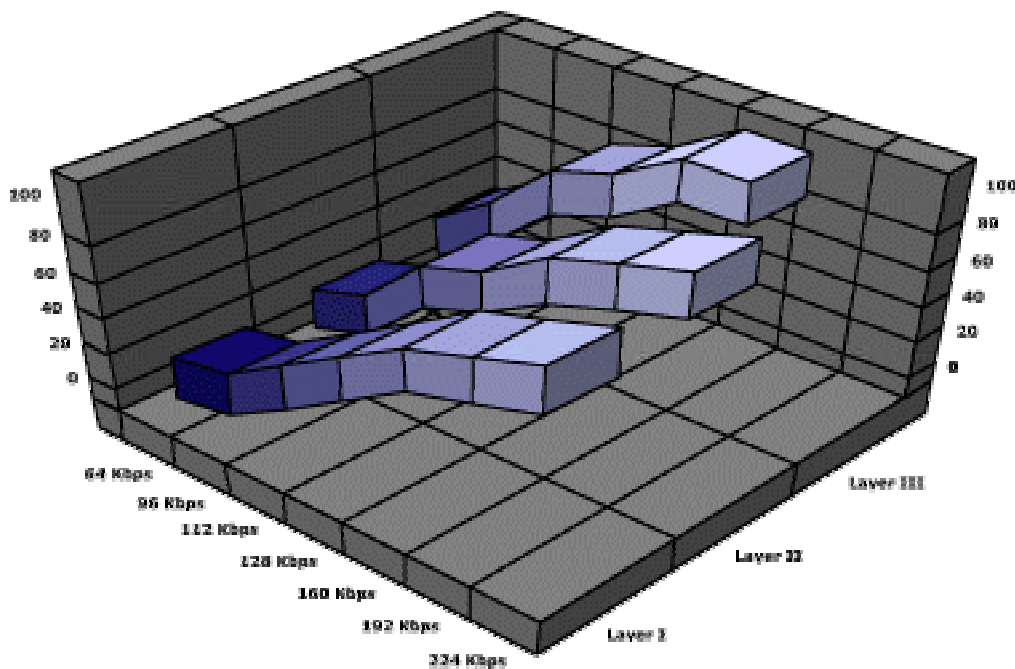


Le graphique ci-dessous montre les transformées de Fourier des trois signaux précédents.

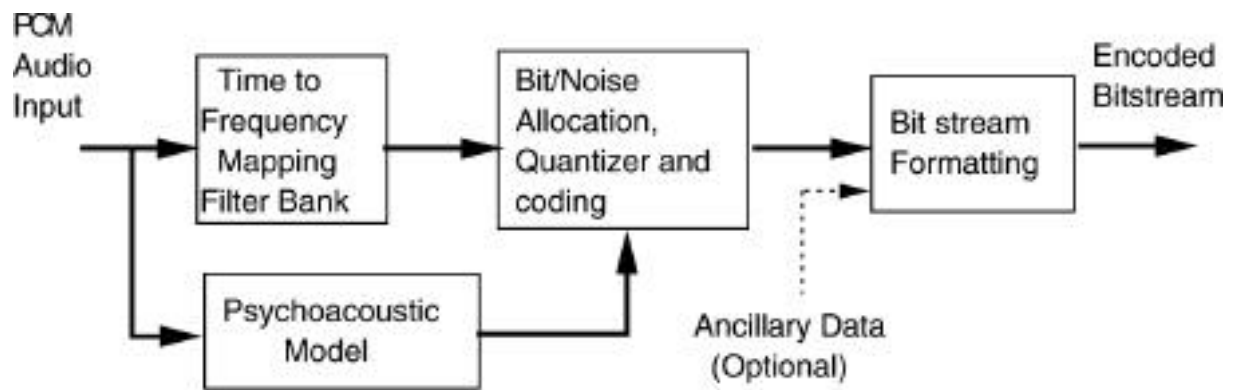


II La compression audio MPEG Layer 3

Le standard de compression audio MPEG a été développé par le « Motion Picture Experts Group » (MPEG) pour la compression du son de qualité « haute fidélité ». L'adoption de ce standard date de 1992. Il s'agit d'un standard à la fois rigide et ouvert. Il est rigide dans sa syntaxe de codage pour permettre la compatibilité au décodage, mais il est ouvert en ce qui concerne la manière d'encoder le signal. Le standard propose trois modes de compression, les layer 1, 2 et 3. Le layer 1 est le plus simple, il permet des débits de l'ordre de 128 kbits/s. Il est par exemple utilisé par Philips dans son système de cassette DCC. Le layer 2 offre une complexité intermédiaire, il est par exemple utilisé dans les CD multimédias ainsi que dans les Vidéo Disc. Le layer 3 est le plus complexe, mais il offre d'excellentes performances tant en taux de compression qu'en qualité du son restitué. Ces trois layer ont une qualité très importante, ils permettent un décodage temps réel. Le layer 3 est ainsi utilisé dans un baladeur. Les performances de ces trois layer différents sont regroupés dans le graphique ci-dessous.

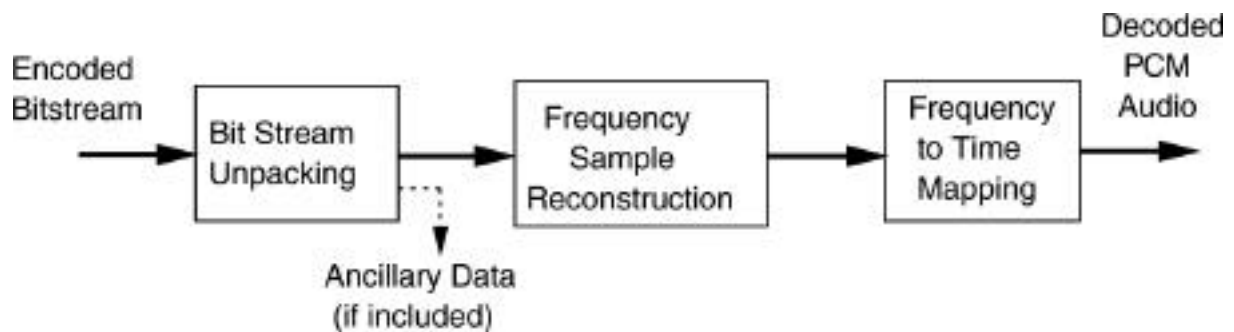


La compression audio MPEG Layer 3 se décompose en deux phases. La première est la phase de calcul du modèle psycho-acoustique, la seconde est celle, plus classique, du codage numérique du signal audio.



MPEG/Audio Encoder

La décompression audio MPEG Layer 3 consiste à effectuer les étapes inverses afin de reconstruire un signal analogique.



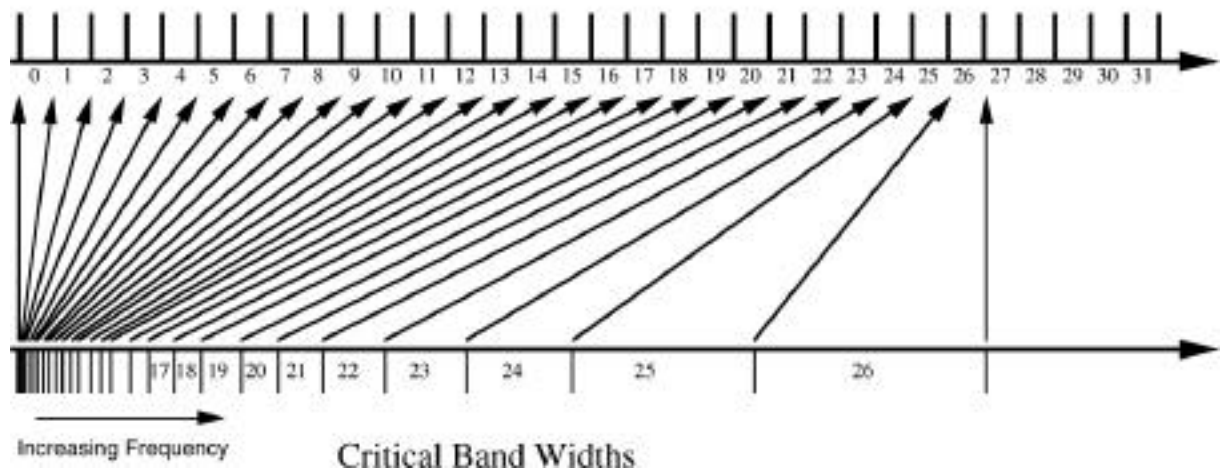
MPEG/Audio Decoder

II.1 Filtrage en sous-bandes

Cette étape consiste à découper la plage de fréquence 0 - 24 kHz en 32 sous-bandes de largeur 750 Hz. On la réalise en deux temps. Tout d'abord, on filtre le signal avec des filtres passe-bande adaptés. Ensuite, les signaux résultants de ces filtrages sont décimés.

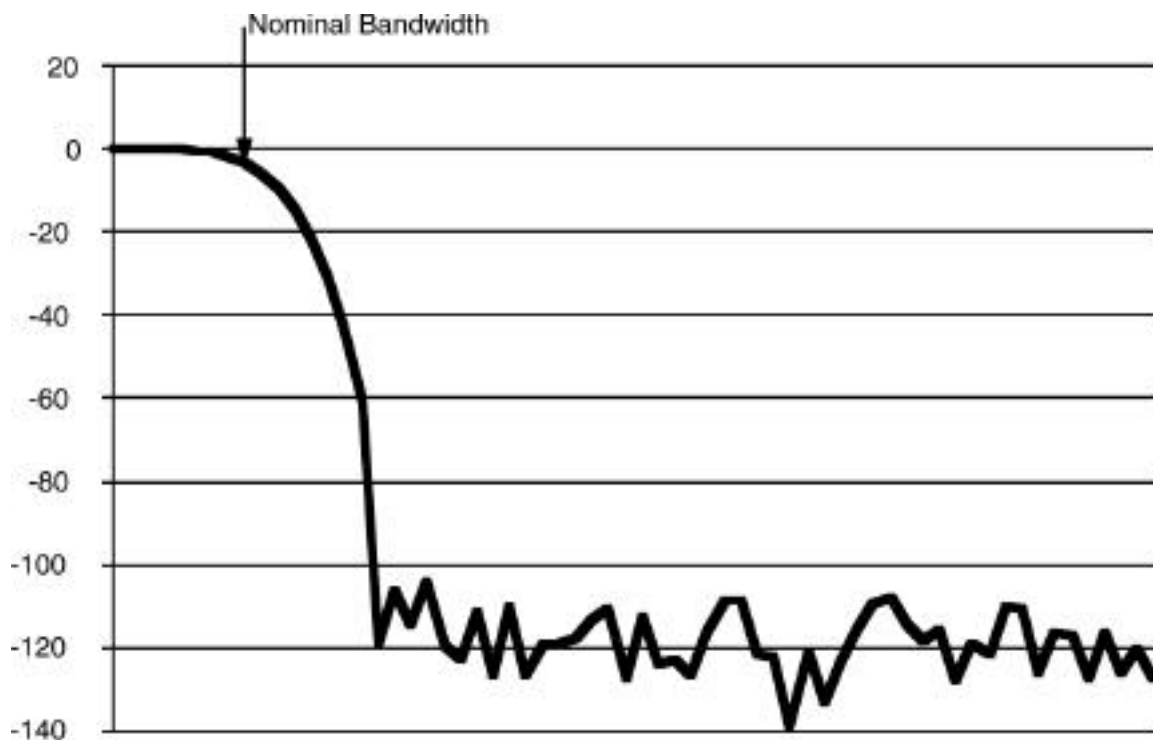
Le fait de découper en 32 bandes de même largeur la plage de fréquence est un compromis. En effet, ce découpage ne reflète pas les performances de l'oreille humaine. Les performances de l'oreille humaine sont définies sur des sous-bandes critiques.

MPEG/Audio Filter Bank Bands



D'autre part, le découpage en bande par filtrage est destructif, toutefois, l'erreur introduite est faible et inaudible.

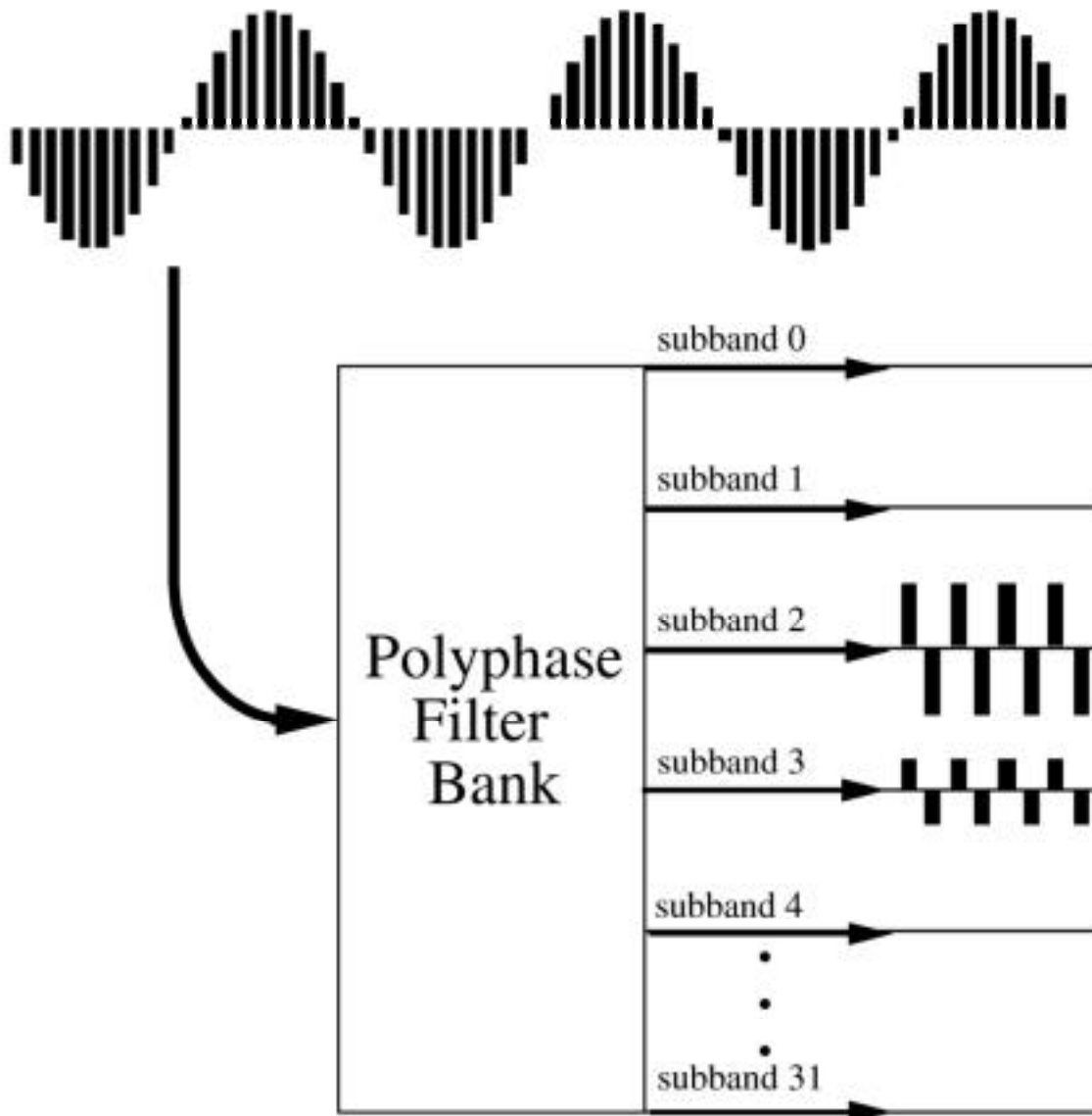
Le graphique ci-dessous montre la réponse en fréquence d'un filtre basé sur les travaux de Rothweiler, Polyphase Quadrature Filters – a New Subband Coding Technique (1983), utilisé pour le découpage en sous-bande.



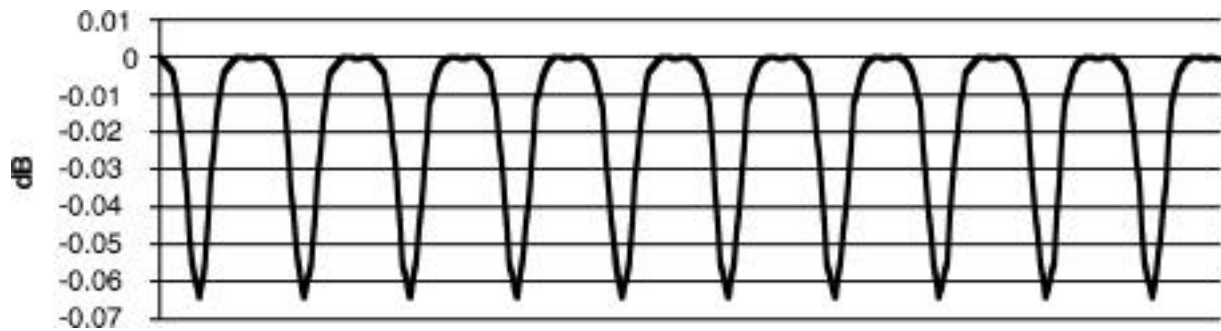
On remarque que ces filtres ne coupent de manière très rapide à la fréquence de coupure. Ceci implique un aliasing important puisque le signal est découpé en 32 sous-bandes. Ainsi, si l'on considère une fréquence proche d'une

borne d'une plage de fréquence, on peut retrouver de l'énergie dans deux sous-bandes consécutives comme le montre le schéma suivant.

Input audio: 1,500 Hz sinewave sampled at 32 kHz, 64 of 256 samples shown

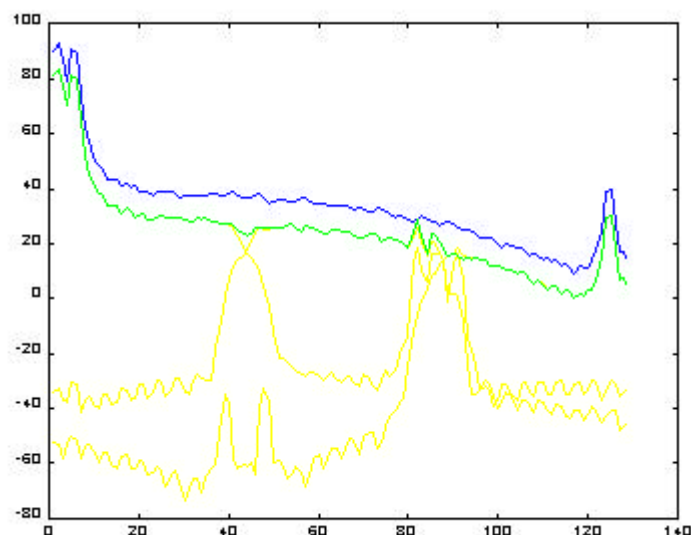


Toutefois, même si ce filtrage est destructif, les erreurs qui en résultent sont négligeables. Ainsi, on ne perd pas plus de 0,07 dB comme le montre le graphique ci-dessous.



Une simulation sous MatLab permet d'illustrer l'étape de découpage en sous-bande. Le graphique ci-après montre les transformées de Fourier de :

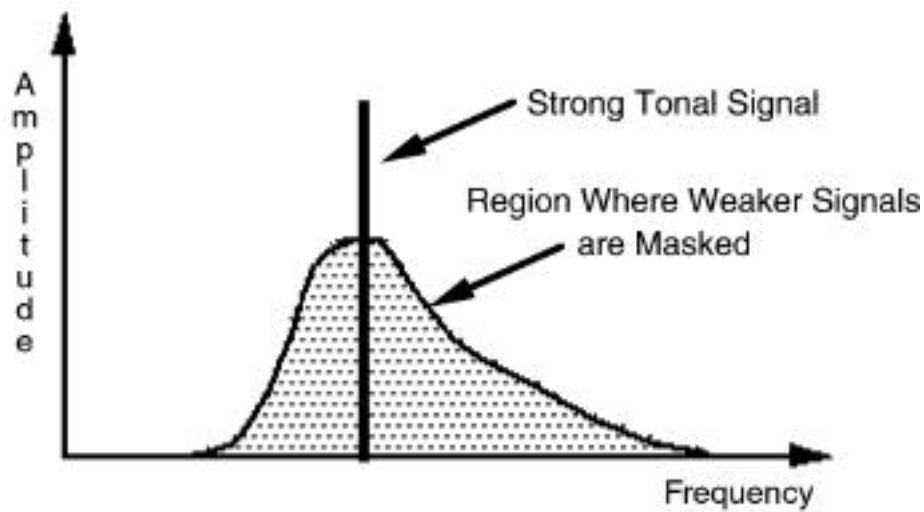
- Le signal initial en bleu
- Les signaux filtrés et décimés dans chaque sous-bande de fréquence (ici, pour simplifier, on ne considère que 3 sous-bandes au lieu de 32) en jaune
- La somme des signaux filtrés et décimés en vert



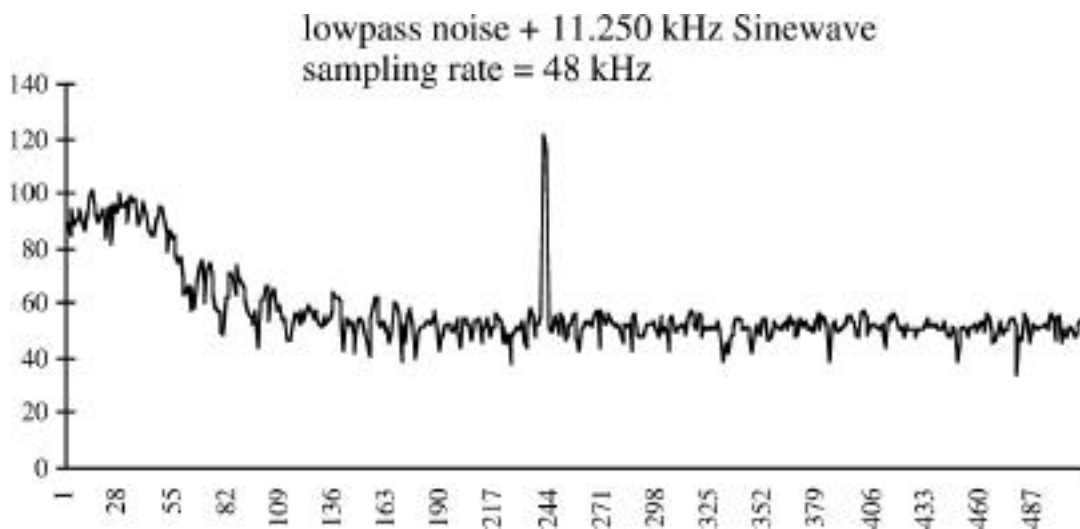
II.2 Le calcul du modèle psycho-acoustique

L'objectif est de déterminer pour chaque bloc de fréquence la courbe de masquage dynamique afin de pouvoir ensuite calculer le bruit de quantification maximal que l'on peut injecter dans chacune des sous-bandes et donc le nombre de bits nécessaires au codage du signal pour chacune des sous-bandes.

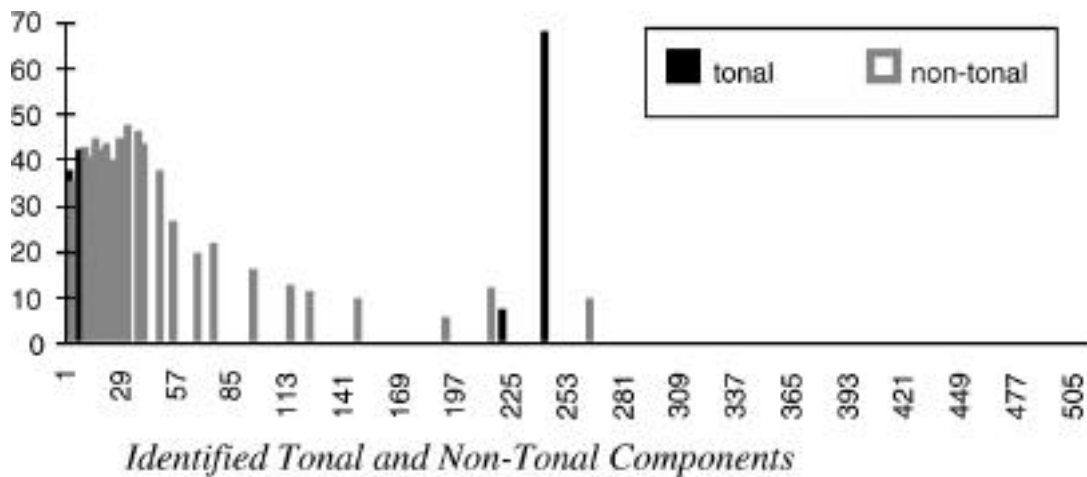
Le principe consiste à dire qu'une raie fréquentielle de forte intensité masque ces voisines de plus faible intensité. Les raies voisines ne sont alors plus perçues par l'oreille humaine.



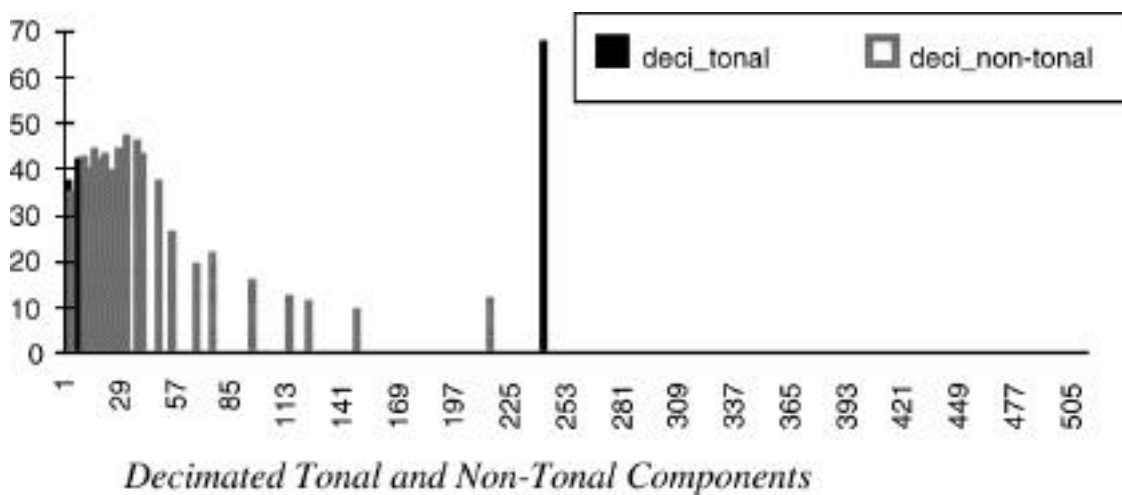
Regardons comment mettre en application ce principe. On s'intéresse à un signal d'entrée composé d'une sinusoïde à 11,25 kHz et de bruit que l'on échantillonne à 48 kHz.



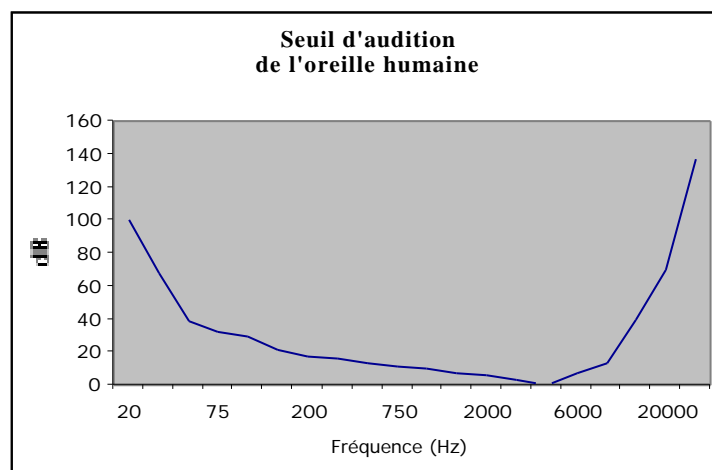
Dans un premier temps, on extrait les composantes tonales et non-tonales.



Ensuite, on décime les tonales et les non-tonales, c'est-à-dire qu'on élimine celles qui sont masquées par leurs voisines proches.



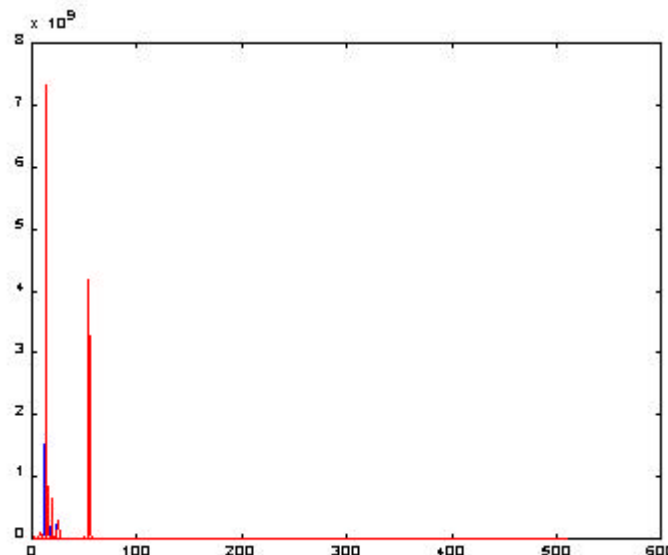
On en déduit alors la courbe de masquage dynamique en tenant compte de la réponse statique de l'oreille et du calcul des tonales et des non-tonales.



La simulation réalisée sous MatLab permet d'illustrer le calcul du modèle psycho-acoustique en six étapes successives :

- Analyse du signal par transformée de Fourier
- Calcul du niveau acoustique dans chaque sous-bande
- Détermination des composantes tonales et sous-tonales
- Calcul de la courbe de masquage globale
- Détermination du nombre de bits nécessaires au codage dans chaque sous-bande

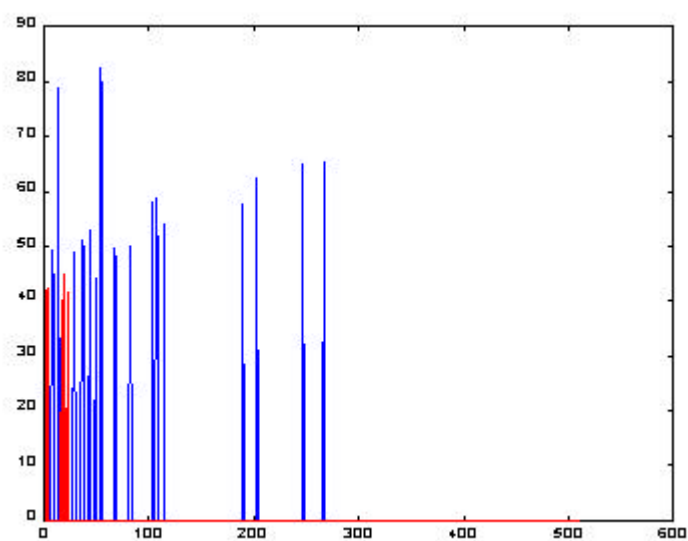
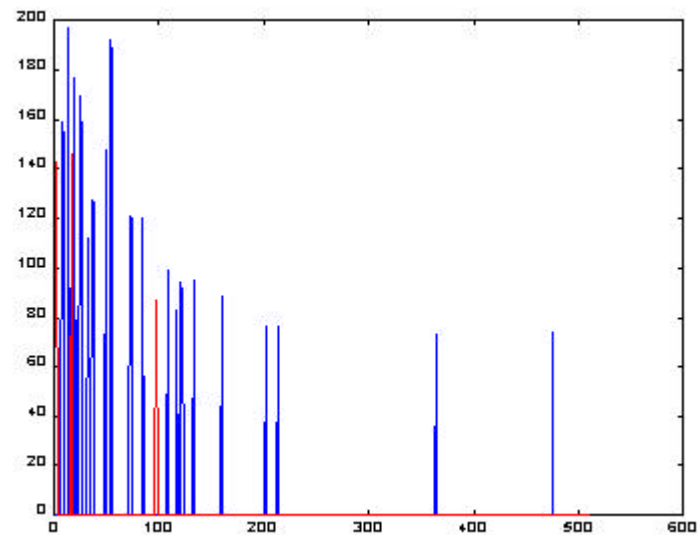
Le graphique ci-après représente en bleu la transformée de Fourier d'un échantillon (j'appelle échantillon le son dans l'une des 32 sous-bandes de fréquence) et en rouge les maximums locaux. Il faut noter que la courbe des maximums locaux vient masquer celle de la transformée de Fourier de l'échantillon.

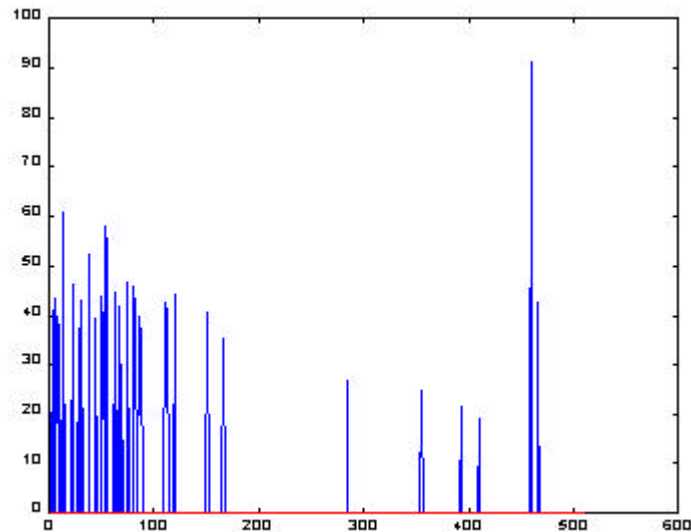


L'étape suivante consiste à extraire les tonales des maximums locaux. Un maximum est une tonale si son amplitude est suffisamment grande (écart supérieur à 7 dB) par rapport aux amplitudes de ses voisines, sachant que le nombre de voisines pris en compte dépend de la fréquence du maximum considéré.

Il faut ensuite extraire les non-tonales. Pour cela, il faut diviser le spectre de fréquence en 26 bandes en utilisant une échelle en Barks. Cette échelle est une échelle non linéaire directement liée au fonctionnement interne de l'oreille. Pour chaque bande critique dans laquelle il n'existe pas de composante tonale, on place au centre une raie fictive dite non-tonale contenant la somme de toutes les raies de la bande.

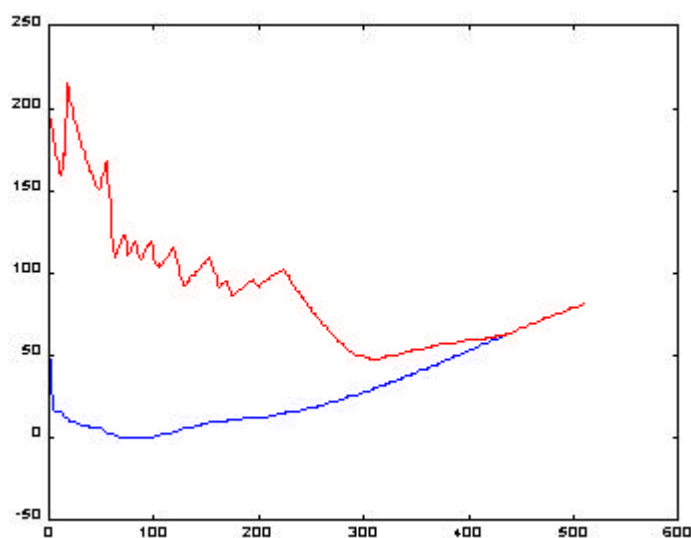
Les graphiques ci-après représentent en bleu les tonales et en rouge les non-tonales dans trois bandes de fréquence d'un même son.





Ensuite, on décime les composantes de masquage. Dans un premier temps, on ne conserve que les masqueurs tonales ou non-tonales dont l'intensité est supérieure au seuil d'audition statique de l'oreille. Ensuite on supprime sélectivement les composantes de masquage dont la distance est inférieure à 0,5 bark.

On peut alors calculer la courbe de masquage dynamique de l'oreille, sachant que chaque tonale et chaque non-tonale masque une partie du son autour d'elle. On obtient alors une courbe comme celle ci-dessous où est représenté en bleu le seuil d'audition statique de l'oreille humaine et en rouge la courbe de masquage dynamique pour un échantillon donné.



On peut maintenant calculer le rapport signal sur bruit et en déduire le nombre de bits nécessaires au codage pour chaque sous-bande.

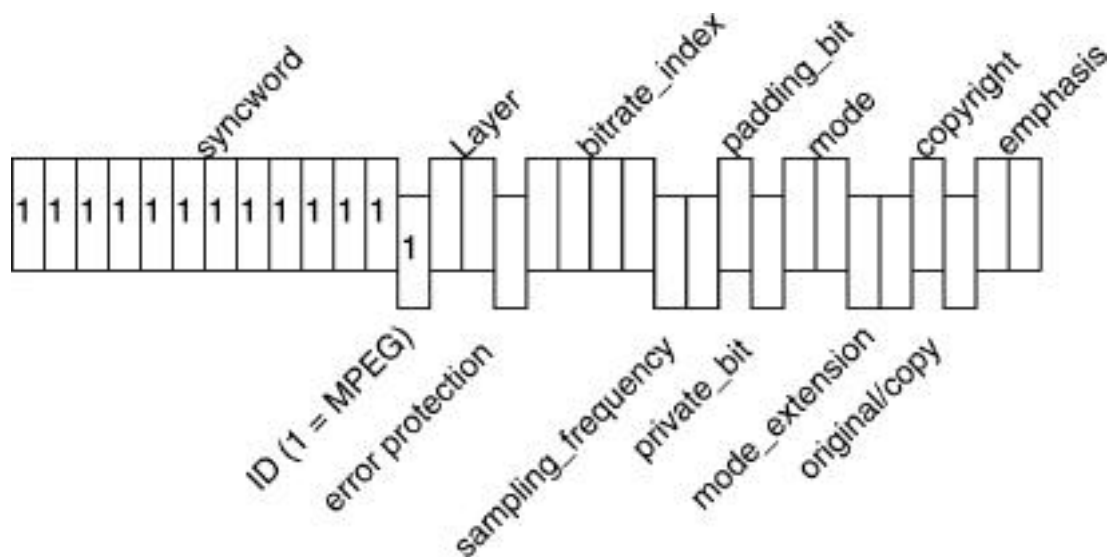
II.3 La quantification et le codage MPEG Layer 3

Les fichiers MPEG ont une forme bien spécifique.

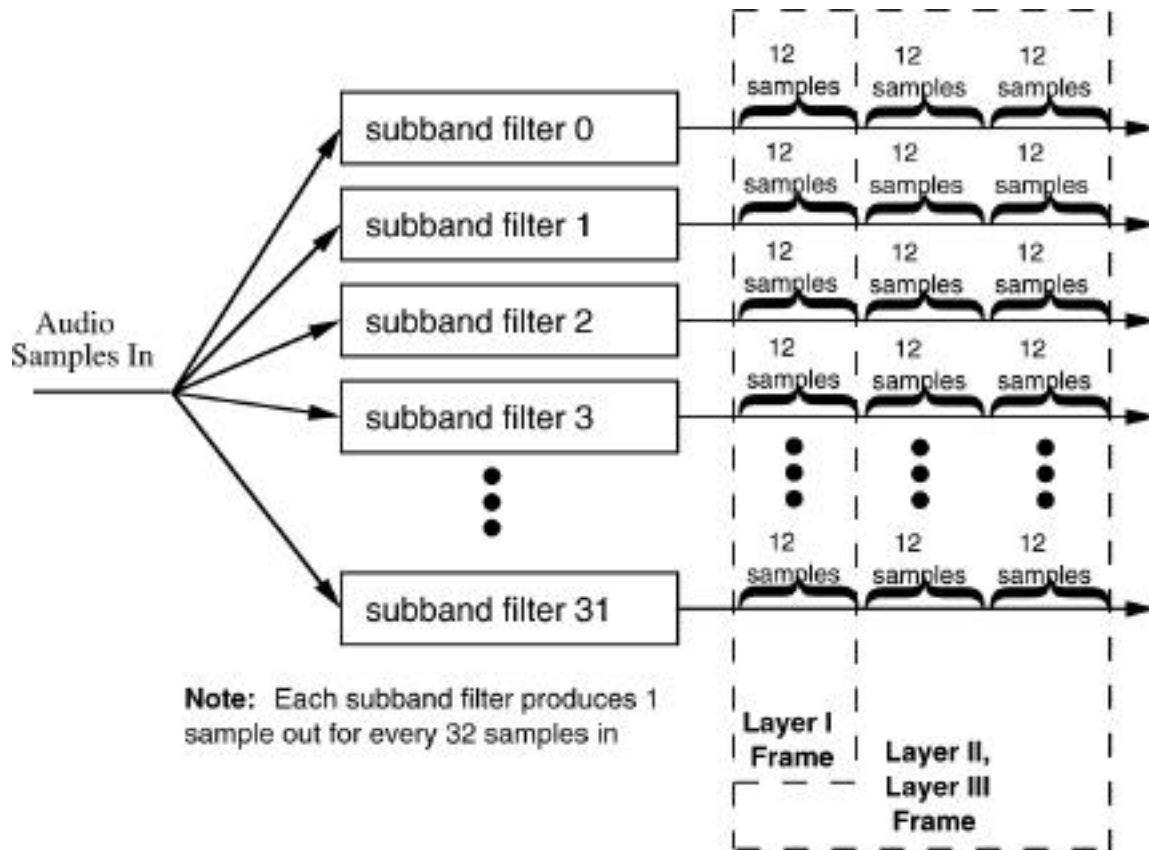
Header (32)	CRC (0,16)	Side Information (136,256)	Main Data; not necessarily linked to this frame. See figure 22
----------------	---------------	-------------------------------	-------------------------------------------------------------------

The Frame Format of a Layer III Bitstream

Leur en-tête est le suivant :



Dans la norme MPEG Layer 3, on code 36 échantillons pour chaque sous-bande. Ces échantillons sont regroupés en 3 blocs de 12 échantillons.

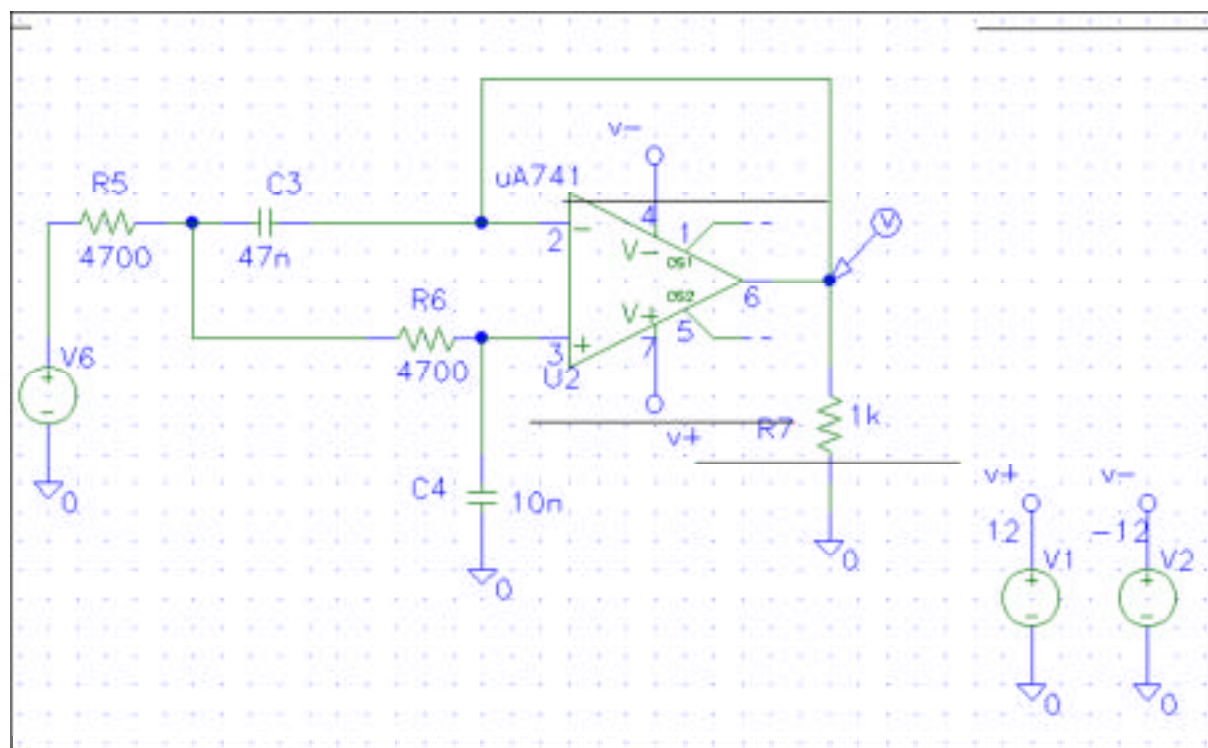


Grouping of Subband Samples for Layer I and Layer II

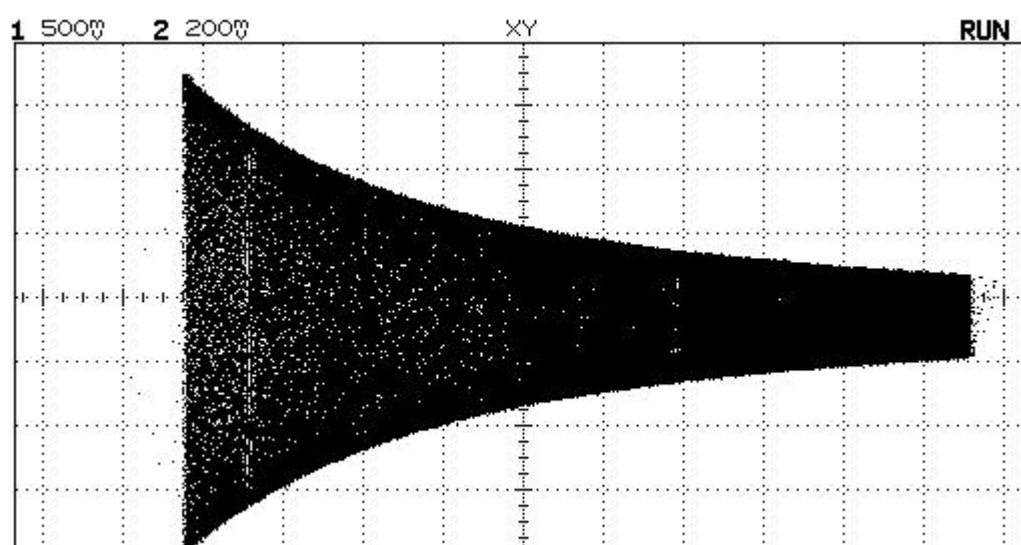
Pour chaque bloc de 12 échantillons, on repère le maximum de la valeur absolue de l'amplitude du signal que l'on appelle facteur d'échelle (SCF – « Scale Factor »). Ce facteur d'échelle permet de ramener le signal dans la plage $[-1,1]$ et de le quantifier en utilisant le nombre de bits calculé à l'étape précédente. Le signal est donc prêt à être codé. Pour que le codage soit réversible, il faut à la fois coder les amplitudes ramenées sur $[-1,1]$ et le facteur d'échelle.

Annexe

Un exemple de filtre passe-bas susceptible d'être utilisé dans une chaîne de conversion analogique numérique :



Le graphique ci-dessous donne la réponse en fréquence d'un tel filtre :



Bibliographie

1. Le son numérique – Yvan Bonnassieux – Ecole Polytechnique
2. Codage Musicam – Yvan Bonnassieux – Ecole Polytechnique
3. Digital Audio Compression – Davis Yen Pan – Digital Technical Journal
4. A Tutorial on MPEG / Audio Compression – Davis Yen Pan – Motorola