

**VILLE DE LIÈGE**

**Institut de Technologie  
Enseignement de Promotion sociale**

Année académique 2021 – 2022

**Développement d'un codec audio AAC :  
optimisation de l'algorithme MDCT  
pour l'architecture ARM**

Étudiante :

**Laura Binacchi**

Lieu de stage :

**EVS Broadcast Equipment**

Rue du Bois Saint-Jean 13, 4102 Ougrée

Maître de stage :

**Bernard Thilmant**

Software Engineer

Épreuve intégrée présentée pour l'obtention du diplôme de  
**BACHELIER.E EN INFORMATIQUE ET SYSTÈMES**  
**FINALITÉ : INFORMATIQUE INDUSTRIELLE**

## Table des matières

<b>Introduction</b>	<b>1</b>
<b>1 EVS Broadcast Equipment</b>	<b>2</b>
1.1 Présentation d'EVS et du département R&D	2
1.2 Le serveur XT	2
<b>2 L'encodage audio numérique : généralités</b>	<b>4</b>
2.1 Le son	4
2.2 La numérisation d'un signal	4
<b>3 Les codecs audio</b>	<b>5</b>
3.1 Définition d'un codec	5
3.2 Les codecs MPEG	5
3.3 Les modèles psychoacoustiques	6
<b>4 Le codec AAC</b>	<b>8</b>
4.1 Fonctionnement de l'encodeur AAC	8
4.2 Le bloc MDCT	8
<b>5 Environnement de développement</b>	<b>9</b>
<b>6 Algorithmes MDCT de référence</b>	<b>10</b>
6.1 Description mathématique	10
6.2 Implémentations des MDCT de référence en <i>floating point</i> et <i>fixed point</i>	10
6.3 Validation des algorithmes de référence	11
<b>7 Algorithme MDCT basé sur la FFT</b>	<b>12</b>
7.1 Optimisations attendues	12
7.2 Implémentation de la MDCT basée sur la FFT de la librairie FFTW3	12
7.3 Validation	12
<b>8 Intégration de la librairie Ne10</b>	<b>13</b>
8.1 Choix de la librairie	13
8.2 Implémentation de la MDCT basée sur la FFT Ne10 en <i>float 32</i>	13
8.3 Validation	14
8.4 Performances	14
<b>9 Algorithme MDCT en arithmétique fixed point</b>	<b>15</b>
9.1 Arithmétique fixed point	15
9.2 Améliorations attendues	15
9.3 Implémentation de la MDCT basée sur la FFT Ne10 en arithmétique <i>fixed point</i>	15
9.4 Performances	15
9.5 Arithmétique fixed point	16

---

<b>10 Optimisations à l'architecture ARM</b>	<b>16</b>
10.1 Spécificités de l'architecture ARMv8 . . . . .	16
10.2 Utilisation des fonctions Neon SIMD (intrinsic) . . . . .	16
<b>11 Analyse des résultats</b>	<b>17</b>
11.1 Validation des données . . . . .	17
11.2 Gain en performance . . . . .	17
11.3 Perte de précision . . . . .	18
<b>12 Améliorations possibles</b>	<b>19</b>
<b>Conclusion</b>	<b>20</b>
<b>Références</b>	<b>20</b>

## Table des figures

1	Logo de la société EVS Broadcast Equipment[1]	2
2	Vues avant et arrière (en configuration IP) de l'XT-VIA[2]	3
3	Vue simplifiée d'un codec MPEG basé sur un modèle psychoacoustique[3]	5
4	Effets de masque dans le domaine fréquentiel[4]	7
5		14

## **Remerciements**

## Introduction

Développement d'une solution de software embarqué sur processeur ARM pour encodage audio AAC optimisé aux applications d'EVS :

- Prise de connaissance de l'encodage AAC et de l'environnement EVS qui utilise ce type de format ;
- Prise de connaissance des résultats des optimisations possibles du modèle psycho-acoustique développé par EVS ;
- Développement du code en C ou Assembleur pour l'encodage AAC sur plateforme ARM ;
- Test du système et documentation de son implémentation.

Ce travail commencera par une présentation d'EVS et du département dans lequel s'est déroulé mon stage. Parmi les nombreux produits d'EVS, seul le serveur XT sera brièvement présenté puisque c'est spécifiquement pour ce dernier que le codec AAC est développé et optimisé.

Quelques notions théoriques indispensables à la compréhension du travail pratique seront ensuite développées avec une section consacrée au son et sa numérisation et une autre consacrée aux codecs MPEG, à leur fonctionnement et en particulier au fonctionnement du bloc MDCT de l'encodeur AAC.

# 1 EVS Broadcast Equipment

## 1.1 Présentation d'EVS et du département R&D

Mon stage s'est déroulé au sein de la société EVS Broadcast Equipment dont la figure 1 représente le logo. EVS est une entreprise d'origine liégeoise devenue internationale. Fondée en 1994 par Pierre L'Hoest, Laurent Minguet et Michel Counson, EVS compte aujourd'hui plus de 600 employés dans plus de 20 bureaux à travers le monde mais son siège principal se situe toujours à Liège.



FIGURE 1 – Logo de la société EVS Broadcast Equipment[1]

EVS est devenu leader dans le monde du broadcast avec ses serveurs permettant l'accès et la diffusion instantanée des données audiovisuelles enregistrées sur ses serveurs. L'entreprise est également célèbre pour ses ralentis instantanés. Ces technologies sont utilisées pour la production live des plus importants événements sportifs dans le monde : le matériel EVS est notamment utilisé pour la retransmission des Jeux Olympiques depuis 1998.

Plus de 50% des employés d'EVS travaillent en recherche et développement afin de répondre au marché du broadcast en constante évolution. Outre ses solutions techniques innovantes, EVS se différencie de ses concurrents par la proximité entretenue avec les clients en leur proposant des solutions à l'écoute de leurs besoins et en leur offrant un service de support de qualité.

C'est en R&D, dans l'équipe Hardware-Firmware, que s'est déroulé mon stage. Sous la direction de Justin Mannesberg, cette équipe se compose d'une vingtaine d'employés spécialisés en développement embarqué et en développement FPGA. La situation particulière dans laquelle s'est déroulé mon stage, en pleine pandémie de Covid et alors que tous les employés étaient confinés, ne m'a pas permis d'interagir avec beaucoup de membres de l'équipe et ni de pouvoir observer leur travail. Bernard Thilmant (Software Engineer dans l'équipe Hardware-Firmware) a cependant réussi à m'apporter le soutien nécessaire à la bonne réalisation de mon stage : il m'a permis de m'initier au C++, m'a aidée à ne pas me perdre dans les concepts parfois complexes de l'encodage audio et m'a aidée à apporter la rigueur scientifique nécessaire à la réalisation de mon travail. J'ai également pu bénéficier de l'expertise technique de Frédéric Lefranc (Principal Embedded System Architect dans l'équipe Hardware-Firmware) ainsi que du suivi de Justin Mannesberg (Manager de l'équipe Hardware-Firmware).

## 1.2 Le serveur XT

EVS développe et commercialise de nombreux produits allant des serveurs de production aux interfaces permettant d'exploiter des données audiovisuelles ou de monitorer des systèmes de production[2]. Le serveur de production live XT est un des produits emblématiques d'EVS. Il permet de stocker de grandes quantités de données audiovisuelles et d'y accéder en temps réel afin de répondre aux besoins de la production en live. Par exemple, la remote LSM (*Live Slow Motion*) permet d'accéder aux contenus des serveurs XT afin de créer les ralentis pour lesquels EVS est célèbre dans le monde.



FIGURE 2 – Vues avant et arrière (en configuration IP) de l'XT-VIA[2]

Le serveur XT a connu plusieurs versions : XT, XT2, XT2+, XT3 et enfin l'XT-VIA. L'XT-VIA (cf figure 2), la plus récente version du serveur XT, en quelques informations clés[2] :

- offre un espace de stockage de 18 à 54 TB, soit plus de 130h d'enregistrement en UHD-4K ;
- dispose de 2 à plus de 16 canaux selon le format choisi : 2 canaux en UHD-8K (4320p), 6 canaux en UHD-4K (2160p) et plus de 16 canaux en FHD and HD (720p, 1080i, 1080p) ;
- permet une configuration hybride de ses entrées et sorties en IP (10G Ethernet SFP+, 100G en option, ST2022-6, ST2022-7, ST2022-8, ST2110, NMOS IS-04, IS-05, EMBER+, PTP) ou SDI (1.5G-SDI, 3G-SDI et 12G-SDI) ;
- supporte de nombreux formats d'encodage vidéo : UHD-4K (XAVC-Intra et DNxHR), HD/FHD (XAVC-I, AVC-I, DNxHD et ProRes), PROXY (MJPEG et H264) ;
- peut enregistrer 192 canaux audio non compressés et supporte les standards AES et MADI ;
- offre de nombreuses possibilités de connexion avec du matériel EVS ou non.

C'est pour la dernière génération du serveur XT, l'XT-VIA, que le codec AAC est développé. La compression avec perte de données de ce codec permet d'optimiser l'espace occupé par les données audio sans en altérer la qualité perçue. Outre la qualité audio, les performances de l'encodage sont importantes à prendre en compte pour permettre l'enregistrement de plusieurs canaux en parallèle tout en conservant un traitement de l'information qui tienne le temps réel. L'optimisation des performances doit tenir compte de l'architecture de l'XT-VIA : l'architecture ARM Neon remplace l'architecture Intel x86 de ses prédécesseurs avec des différences importantes dans les fonctions intrinsèques.



## **2 L'encodage audionumérique : généralités**

### **2.1 Le son**

### **2.2 La numérisation d'un signal**

### 3 Les codecs audio

#### 3.1 Définition d'un codec

Un codec est un procédé logiciel composé d'un encodeur (*coder*) et d'un décodeur (*decoder*)[5]. Un codec audio permet donc, d'une part, de coder un signal audio dans un flux de données numériques et, d'autre part, de décoder ces données afin de restituer le signal audio.

Les codecs sont dits avec perte (*lossy*) ou sans perte (*lossless*). Le PCM est par exemple un codec sans perte puisqu'il encode la totalité des informations sonores dans la bande de fréquences humainement audible. Ce type de codec permet de conserver la qualité de l'audio mais nécessite en contrepartie un espace de stockage conséquent, même avec une compression des données.

Afin de réduire l'espace de stockage nécessaire, les codecs avec perte permettent de supprimer une partie des données audio. C'est le cas des codecs définis par les normes MPEG dont fait partie le codec AAC.

#### 3.2 Les codecs MPEG

MPEG (*Moving Picture Experts Group*) désigne une alliance de différents groupes de travail définissant des normes d'encodage, de compression et décompression et de transmission de média audio, vidéo et graphiques[6]. Le groupe est actif depuis 1988 et a produit depuis de nombreuses normes.

Les codecs audio qui implémentent les normes MPEG ont pour point commun d'être des codecs avec perte de données basés sur un modèle psychoacoustique. Le premier est le MP3, défini par la norme MPEG-1 Layer-3 ISO/IEC 11172-3 :1993. Le codec AAC est conçu en 1997 pour remplacer le MP3. Il est défini par les normes MPEG-2 partie 7 ISO/IEC 13818-7 :2006[7] et MPEG-4 partie 3 ISO/IEC 14496-3 :2019[8].

Les normes MPEG définissent les grandes lignes de l'encodage et du décodage et le format du conteneur mais pas l'implémentation du codec qui peut de ce fait être plus ou moins performant. Les codecs MPEG sont typiquement composés des blocs suivants :

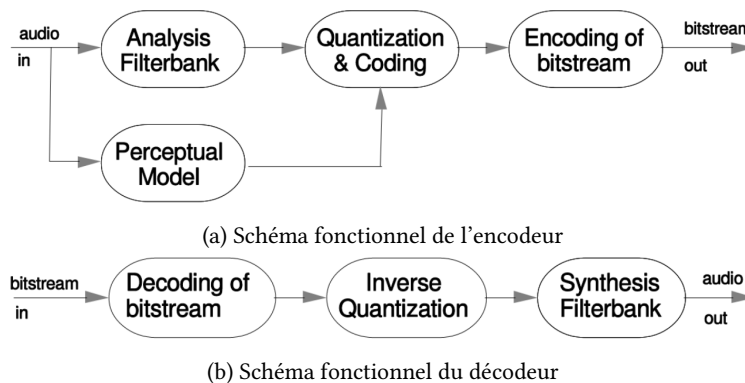


FIGURE 3 – Vue simplifiée d'un codec MPEG basé sur un modèle psychoacoustique[3]

L'encodeur est composé des blocs suivants :

**filter bank** la banque de filtres décompose le signal temporel d'entrée en différentes composantes fréquentielles

**perceptual model** le modèle psychoacoustique utilise le signal temporel et/ou sa décomposition fréquentielle pour éliminer les données audio dont l'absence ne nuira pas à la qualité perçue à l'écoute ()

**quantization and coding** la quantification attribue une valeur numérique aux données du spectre de fréquences : elle sont typiquement codées avec une méthode entropique qui peut être optimisée avec le modèle psychoacoustique

**encoding of bitstream** les données sont formatées en un flux contenant typiquement le spectre de fréquences codé et des informations supplémentaires permettant l'encodage

Le décodeur a un fonctionnement inverse : le flux de données est décodé (**decoding of bitstream**), les composantes fréquentielles du signal sont retrouvées par l'opération inverse à la quantification (**inverse quantization**) et ces sous-bandes fréquentielles sont finalement rassemblées pour reconstituer le signal temporel (**synthesis filter bank**).

Le fonctionnement du décodeur ne sera pas plus développé dans ce travail car le bloc MDCT fait partie de la banque de filtres de l'encodeur. Le fonctionnement spécifique de l'encodeur AAC sera par contre détaillé dans la section 4.

### 3.3 Les modèles psychoacoustiques

Le section précédente a défini les codecs MPEG comme étant basés sur un modèle psychoacoustique. La psychoacoustique est une branche de la psychophysique qui étudie la manière dont l'oreille humaine perçoit le son[9]. Cette discipline permet d'améliorer la compression d'un signal audio en éliminant les sons qui sont captés par un microphone mais qui ne peuvent pas être perçus par l'oreille humaine et les avancées dans cette discipline permettent de développer des encodeurs audio de plus en plus performants. Les codecs basés sur un modèle psychoacoustique sont toujours des codecs avec perte puisqu'une partie des informations auditives sera définitivement perdue, ce qui ne nuit pour autant pas à la qualité perçue du son.

L'encodage audionumérique tient déjà compte des seuils de fréquences humainement audibles pour limiter les données audio enregistrées : nous l'avons vu dans la section 2.2, aucun son n'est perçu en-deça de 20Hz ou au-delà de 20kHz. La psychoacoustique permet de mieux dessiner la limite entre ce qui est humainement audible ou non afin d'éliminer un maximum des informations non pertinentes et ainsi augmenter le facteur de compression des données : le facteur de compression des codecs MPEG est environ 15 fois supérieur à celui du CD[4].

Les effets de masque sont au centre des différents modèles psychoacoustiques utilisés pour la compression audio. L'enjeu afin d'obtenir le meilleur taux de compression est de calculer le plus finement possible les seuils de masquage, i.e. la limite entre les informations pertinentes et celles qui peuvent être éliminées. Les effets de masques dans le domaine fréquentiel (*spectral masking effects*) sont parmi les plus utilisés mais il en existe d'autres, e.g. dans le domaine temporel. La figure suivante représente différents effets de masque du domaine fréquentiel :

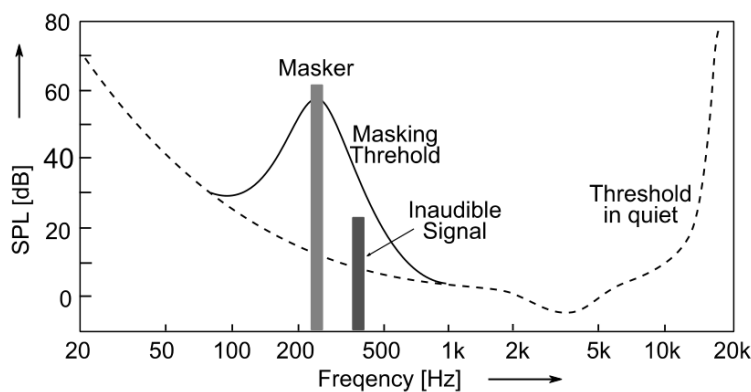


FIGURE 4 – Effets de masque dans le domaine fréquentiel[4]

Les lignes représentent le seuil

**Threshold in quiet** la ligne en pointillé représente le seuil d'audibilité dans le calme, indépendamment de tout autre élément qui pourrait interférer

#### Masking threshold

Le calcul du seuil de masquage tient compte[4] :

- des effets de masquage monophonique de la perception non-linéaire des fréquences, plus fine pour les basses fréquences ;
- des effets de masque dans le temps ;
- de l'effet de masque dans une bande de fréquence ou entre bandes de fréquences ;
- l'impact de la tonalité sur le masquage
- masking over time

## **4 Le codec AAC**

### **4.1 Fonctionnement de l'encodeur AAC**

Le codec AAC (Advanced Audio coding) est défini, défini par la norme . Les recherches en psychoacoustique ont permis de développer un algorithme d'encodage plus performant pour l'AAC que pour le MP3 : il permet d'encoder moins données audio tout en gardant la même qualité perçue au décodage[3].

### **4.2 Le bloc MDCT**

## 5 Environnement de développement

Plateformes : passage de Intel à ARM nécessaire à cause de la lib Ne10 et impossibilité de maintenir une version de référence Intel

Développement remote sur RPI + photo raspberry

CentOS7 parce que utilisé sur tout le matériel linux EVS

Projet en C++ mais qui ressemble fort à du C car implémentation d'un codec -> améliorations possibles : juste du C?

Le build du projet se fait grâce à un fichier CMake à la racine du projet. Ce fichier présenté dans l'annexe ?? permet d'appeler les CMake du projet *audio\_encoding* (projet contenant les MDCT et ses tests) et de la librairie *Ne10*. Contrairement à l'usage, ces deux projets sont construits par défaut en mode *release* et non en mode *debug* afin de ne pas risquer de faire tourner les tests de performance en mode *debug*. Les commandes du CMake générant les exécutables de tests seront fournies en annexe à la suite du code source de ces tests.

## 6 Algorithmes MDCT de référence

La première étape de ce travail consiste à développer des algorithmes de référence afin de valider les différentes MDCT implémentées par la suite. Ces algorithmes de référence sont développés en algorithmique flottante sur base de la formule mathématique de la MDCT. Ils permettent de générer des spectres de fréquence en *float* ou en *integer* afin de valider les données de sortie des MDCT optimisées.

### 6.1 Description mathématique

La transformation effectuée par la MDCT est donnée par l'équation suivante[10] :

$$X_k = \frac{2}{\sqrt{2N}} \sum_{n=0}^{2N-1} x_n \cos \left[ \frac{\pi}{N} \left( N + \frac{1}{2} + \frac{N}{2} \right) \left( k + \frac{1}{2} \right) \right]$$

$X_k$  avec  $k \in [0, N[$  pour une fenêtre d'entrée de  $2N$  échantillons

$x_n$  avec  $n \in [0, 2N[$  : la fenêtre d'entrée

$F : \mathbb{R}^{2N} \rightarrow \mathbb{R}^N$  la MDCT est une fonction linéaire qui pour  $2N$  nombres réels en entrée produit  $N$  nombres réels en sortie

La MDCT a été implémentée avec une fenêtre d'entrée de  $2N = 1024$  échantillons. Le bloc de sortie, i.e. le spectre de fréquences de la fenêtre d'entrée, aura donc une taille de 512. Ces valeurs, utilisées à de très nombreux endroits du code, sont rassemblées dans le header `mdct_constants.h` présenté dans l'**annexe B**. Ce fichier contient également d'autres valeurs précalculées sur base de la taille de la fenêtre d'entrée.

La section suivante présente deux implémentations simples de cette formule. Ces implémentations ne pourraient pas être utilisées sans avoir été optimisées car elles seraient beaucoup trop lentes pour un codec qui doit tenir le temps réel sur plusieurs canaux. La complexité de implémentation de cette formule est de  $O(N^2)$  opérations (où  $N$  est la taille de la fenêtre d'entrée). Cette complexité peut être ramenée à  $O(N \log N)$  opérations par une factorisation récursive. La complexité peut également être diminuée en se basant sur une autre transformation, e.g. une DFT (*Discrete Fourier Transform*) ou une autre DCT (*Discrete Cosine Transform*) : la complexité sera alors de  $O(N)$  opérations de pre- et post-processing en plus de la complexité de la DFT ou de la DCT choisie[10]. C'est cette dernière option qui a été retenue pour ce travail.

### 6.2 Implémentations des MDCT de référence en *floating point* et *fixed point*

La formule mathématique de la MDCT a été implémentée très simplement en algorithmique flottante avec la possibilité d'obtenir le spectre de fréquences codés en *float*, *double* ou *int32*. L'objectif de ces implémentations est de pouvoir valider les spectres de fréquence calculés par les implémentations optimisées de la MDCT. Les MDCT de référence serviront également à mesurer la précision des MDCT optimisées.

La première implémentation de l'équation de la MDCT est présentée dans l'**annexe C.1**. La fonction développée permet de faire ses calculs et d'obtenir un résultat aussi bien en *float* (32 bits) qu'en *double* (64 bits) grâce à l'utilisation d'un *template*. Le signal temporel a la même précision (*float* ou *double*) que le spectre généré.

La seconde fonction de référence est présentée dans l'**annexe C.2**. Elle permettra de vérifier les résultats des implémentations optimisées en *fixed point*. Tous les calculs ne sont pas fait en algorithmique entière : la fonction fait les mêmes calculs que la fonction de référence en algorithmique flottante (uniquement en *double* cette fois pour garder le plus de précision possible) et transtype le résultat final dans un *integer* de 32 bits qui correspond à une notation Qx.15 signée.

### 6.3 Validation des algorithmes de référence

Validation avec un code d'exemple qui correspond à un signal sinusoïdal (single tone) -> annexes : génération d'un signal single tone ??+ code qui sort les données. Code pas ici mais renvoi à la section sur la validation des données -> pour le float, on regarde ce qui sort, pour le integer, on vérifie avec calcul q15

Présentation des résultats sous forme de données brutes ou de graphique.

Explication de la lecture des résultats, calcul des bandes de fréquences représentées. Mise en évidence qu'on a bien une seule composante fréquentielle comme attendu pour un signal single tone

+ les résultats auraient été plus précis avec la fonction de fenêtre -> voir améliorations possibles



## 7 Algorithme MDCT basé sur la FFT

### 7.1 Optimisations attentues

Choix de l'optimisation de la MDCT basé sur sur une DCT : la FFT -> complexité de  $O(N \log N)$  ( $N$  taille de la fenetre) +  $O(N)$  opérations de pré et de post processing.

But : utiliser une FFT déjà optimisée pour n'avoir à optimiser "que" les opérations de pré et de post processing  
Exemple trouvé sur le site DSP related -> Annexe? Citation?

### 7.2 Implémentation de la MDCT basée sur la FFT de la librairie FFTW3

Code développé à partir de cet exemple en annexe. Il est basé sur la librairie FFTW3 qu'on ne peut pas garder car ne permet pas de travailler en integer -> sera amené à être retravaillé.

Explications des paramètres du code :

- fenêtre de 1024
- pre-twiddle -> 256
- FFT => FFT avec une fenêtre d'entrée réduite et donc une complexité réduite
- post-twiddle -> 512

### 7.3 Validation

Validé avec un single tone signal + l'algo de référence.

Code d'exemple qui teste l'algo de référence et l'algo FFTW3 avec les mêmes données d'entrée pour comparaison.  
(on fait la différence pour la précision? ce n'est peut être pas pertinent de la faire déjà)

Présentation des résultats sous forme de données brutes ou de graphiques.

A ce niveau, pour valider, j'ai aussi essayé de faire la IMDCT mais dû à l'overlap, sur une seule fenetre, on ne sait pas reconstruire le signal -> pas possible de valider comme ça

## 8 Intégration de la librairie *Ne10*

### 8.1 Choix de la librairie

La librairie *FFTW3* utilisée pour l'itération précédente de la MDCT ne propose pas de FFT en algorithmique entière. Le passage à une autre librairie était donc nécessaire et le choix s'est porté sur la librairie *Ne10* qui propose différentes FFT en *fixed point*.

Le projet *Ne10* propose toute une série de fonctions mathématiques et physiques de base ainsi que des fonctions de traitement de signal et de traitement d'image. La librairie est spécifiquement développée pour les architectures ARM possédant les opérations SIMD Neon (ARMv7 et ARMv8-A)[11].

*Ne10* propose à la fois des fonctions développées en *plain C* et des fonctions optimisées avec les instructions SIMD Neon : les deux types de fonctions seront utilisées pour développer une MDCT optimisée et pour conserver une MDCT de référence en *plain C*. Maintenir une MDCT *plain C* permettra d'avoir une référence pour la mesure des performances de l'algorithme *fixed point* mais pourrait aussi s'avérer utile pour une utilisation de l'encodeur AAC sur une architecture ARM ne possédant pas les instructions Neon.

L'utilisation de la librairie *Ne10* est soumise à la licence *3-Clause BSD*, licence permissive qui permet un usage commercial des produits intégrant la librairie et qui ne contraint pas à en distribuer le code source[12].

### 8.2 Implémentation de la MDCT basée sur la FFT *Ne10* en *float 32*

La librairie *Ne10* s'installe simplement en suivant les instructions données par la documentation : clone du projet GitHub, run du CMake et build du projet[11]. La librairie ne peut cependant être installée que sur une plateforme Linux, Android ou iOS reposant sur une architecture ARM. À partir de ce moment, il n'est donc plus possible de maintenir une implémentation de référence de la MDCT pour une architecture Intel.

*Ne10* propose des algorithmes de FFT *real to complex* et *complex to complex* en *floating point* (32 bits) ou en *integer* (32 bits et 16 bits). L'objectif est évidemment de passer toute la MDCT en *integer* mais pour un premier test de la librairie, l'algorithme de la section précédente a tout d'abord été repris en remplaçant la FFT de *FFTW3* par la fonction `ne10_fft_c2c_1d_float32_neon` de *Ne10* : FFT en *complex to complex* en *float 32*, i.e. l'entrée et la sortie de la FFT sont représentées sous forme de tableaux de nombres complexes codés en *float* sur 32 bits.

L'annexe ?? présente le code de cette implémentation. L'annexe ?? montre que ce code est construit sur le même modèle que le code utilisant la librairie *FFTW3*. La classe contient la configuration de la FFT et les tableaux contenant les données d'entrée, les données de sortie et les facteurs de twiddling. La classe définit trois fonctions publiques : le constructeur, le destructeur et la fonction MDCT.

L'annexe ?? montre l'initialisation de la MDCT dans le constructeur de la classe. Le constructeur initialise les tableaux de facteurs de twiddling de la même manière que l'algorithme basé sur la FFT de *FFTW3*. La configuration de la FFT de *Ne10* se fait conformément au code d'exemple donné par la documentation de la librairie[11] avec en paramètre la taille de la FFT qui correspond au quart de la taille de la fenêtre d'entrée.

Le destructeur présenté à l'annexe ?? permet de libérer la mémoire allouée pour la FFT en appelant la fonction adéquate de la librairie *Ne10*.

La fonction MDCT présentée dans l'annexe ??, comme pour l'algorithme précédent :

- effectue les opérations de pre-processing ou pre-twiddling : ce sont les mêmes que celle de la MDCT *FFTW3* en arithmétique *floating point* sur 32 bits;
- appelle l'algorithme de FFT : la FFT de *Ne10* prend en paramètres les tableaux contenant les données d'entrée et de sortie, la configuration de la FFT et un *integer* à 0 pour réaliser la FFT ou à 1 pour réaliser l'opération inverse;
- effectue les opérations de post-processing ou post-twiddling : ici aussi les mêmes que celles de la MDCT *FFTW3* en arithmétique *floating point* sur 32 bits.

L'algorithme développé ici ne diffère donc pas de l'algorithme présenté à la section précédente. Son développement est trivial mais il permet de tester et de valider le fonctionnement de la librairie *Ne10*.

### 8.3 Validation

L'utilisation de la librairie *Ne10* est validée par comparaison du spectre de fréquence qu'elle génère avec les spectres générés par la MDCT de référence en *double* et par la MDCT basée sur *FFTW3* également en *double*. Les MDCT sont appelées en *double* pour plus de précision. La comparaison aurait pu se faire sur base de tous les algorithmes de MDCT en *float* sur 32 bits mais le choix qui a été fait ici permet en plus de vérifier que la perte de précision entre le *float* et le *double* soit acceptable.

Le code source permettant de comparer les trois MDCT en *floating point* est présenté à l'annexe ?. Il sera expliqué en détail dans la section 11.1 consacrée à la validation des données des MDCT. Compilé avec les commandes CMake de l'annexe ?, le code produit un exécutable permettant de comparer les spectres de fréquence produits par les trois MDCT *floating point* en les affichant en console.

TODO

FIGURE 5

En redirigeant les données sorties en console vers un fichier texte, il est possible d'exploiter ces données sous forme de graphique. (voir figure 5) TODO : explication du graphique

TODO : présentation de la précision

### 8.4 Performances

Une fois l'utilisation de la librairie *Ne10* validée, Mesure de la différence de performance entre différentes FFT : code en annexe

FFT *Ne10* f32 plain C average run time : 5632.06 ns standard deviation : 4.38439e-09 ns

FFT *Ne10* f32 Neon average run time : 3115.18 ns standard deviation : 2.81769e-06 ns

FFT *Ne10* i32 plain C average run time : 10723.1 ns standard deviation : 1.77145e-06 ns

FFT *Ne10* i32 Neon average run time : 3455.78 ns standard deviation : 5.35823e-07 ns

FFT *Ne10* i16 plain C average run time : 9557.52 ns standard deviation : 1.69145e-07 ns

FFT *Ne10* i16 Neon average run time : 2007.07 ns standard deviation : 2.58836e-07 ns

FFT *FFTW3* f32 plain C average run time : 4641.64 ns standard deviation : 2.69836e-07 ns

## 9 Algorithmes MDCT en arithmétique fixed point

Le passage d'une algorithmique flottante à une algorithmique entière est une des optimisations envisagées par l'analyse préalable à ce travail. Le bénéfice attendu est double : l'algorithmique entière est généralement plus rapide et un bloc MDCT en algorithmique entière permettrait d'économiser des opérations de transtypage des données entières en *float* à l'entrée de la MDCT et inversement en sortie.

### 9.1 Arithmétique fixed point

donner quelques exemple de remplacement d'une opération floating point en fixed point  
trivial en float (reprendre le code de dsp) mais pour l'implémentation en arithmétique entière, attention à mettre dans les bonnes ranges  
La notation *fixed point* peut se faire

### 9.2 Améliorations attendues

Le code de DSP related est en floating point -> passage en fixed point : opérations sur des entiers plus performantes + meilleure intégration (éviter le passage artificiel de l'entier au flottant et inversement)

### 9.3 Implémentation de la MDCT basée sur la FFT Ne10 en arithmétique *fixed point*

L'annexe ?? présente l'implémentation de la MDCT *Ne10 fixed point* en *plain C*. La classe `mdct_ne10_i32_c` a la même structure que l'itération précédente de la MDCT. Le header présenté dans l'annexe ?? montre que seul le type de certaines données a changé :

- Le tableau de facteurs de *twiddling* passe du `float` au `int16_t` ;
- La fonction MDCT prend en paramètres un signal temporel en `int16_t` et renvoie un spectre en `int32_t` au lieu des tableaux de `float`.

Le passage de la FFT en

Explication du code en annexe en accord avec les explications de la première subsection sur le fixed point

### 9.4 Performances

Explication du code en annexe qui permet de mesurer les performances (avec code Timer) performances de la FFT FFTW3 en f32 pour comparaison : average run time : 4259.95 ns standard deviation : 2.07873e-06 ns

Démonstration des résultats

Explication des résultats pas attendus : à ce niveau, résultats pas intéressants car + d'opérations qu'en flottant et ARM a des fonctions d'algo flottante

Pour soutenir cette hypothèse : Test de la librairie et des performances des différentes FFT : performances équivalentes en 32 bits (floating ou fixed) et deux fois plus rapides en 16 bits

Pour comparaison, le temps d'exécution de la FFT de la librairie *FFTW3* a également été mesuré, en *float 32* puisque la FFT de *Ne10* est en *float 32*. Le code de la

Pendant mon stage, j'ai aussi essayé l'algo en Q15 mais échec à cause du manque de documentation des fonctions (quelle headroom prévoir?) -> on n'avait pas seulement une perte de précision mais des dépassements -> pas admissible d'avoir des données fausses

## 9.5 Arithmétique fixed point

# 10 Optimisations à l'architecture ARM

## 10.1 Spécificités de l'architecture ARMv8

Globalement reprend les tutos ARM

## 10.2 Utilisation des fonctions Neon SIMD (intrinsic)

Les opérations SIMD permettent de faire plusieurs opérations en une fois où l'algo normal n'en fait qu'une à la fois. L'algo SIMD permet de faire plusieurs modifications à la fois -> il faut ranger les données de manière à pouvoir l'appliquer facilement (fonction fenêtre -> tri des données? )

Plus les données sur lesquelles on travaille sont petites, plus on peut faire d'opérations en parallèle -> il faut voir si la perte de précision est ok. -> mettre une représentation graphique, c'est beaucoup plus simple à comprendre. D'où l'intérêt de passer à du 16 bits, plutôt que de rester en 32 bits et il faut absolument éviter le 64 bits (aucun intérêt).

Attention au flag pour la compilation + au header (accès aux intrinsics pas activé par défaut)

## 11 Analyse des résultats

### 11.1 Validation des données

Les sections précédentes ont montré que les différentes itérations de la MDCT développées ont été validées à chaque étape. Cette validation consiste essentiellement en une vérification manuelle des données de sortie de la MDCT : avec un signal sinusoïdal connu en entrée, il est facile de vérifier que l'analyse fréquentielle ne contient bien qu'une seule composante fréquentielle, que la vérification se fasse en lisant les données brutes à la sortie de l'algorithme ou par une analyse graphique de celle-ci.

Ces tests auraient pu être améliorés en automatisant la vérification, e.g. en générant une fois les données de référence attendues pour permettre le développement d'un code de test qui compare automatiquement les données de références avec les données à vérifier. En effet, devoir relancer les tests et vérifier les données à chaque fois qu'une modification est faite dans le code peut s'avérer laborieux et mettre en place des tests automatiques aurait permis de gagner un temps précieux.

### 11.2 Gain en performance

La mesure des performances a pour but de valider le bloc MDCT avant de l'intégrer au codec AAC. Le cahier des charges du stage ne contenait pas d'objectif à atteindre en terme de performance, ni absolu (e.g. un temps d'exécution maximal à respecter dans des conditions données), ni relatif (e.g. gagner un certain pourcentage de performances par rapport à une MDCT de référence).

L'objectif en terme de temps d'exécution n'étant pas fini, il a été décidé de tenter de gagner le maximum de performances sur le temps de mon stage. Le critère de réussite est dès lors d'obtenir des performances au moins équivalentes pour la version finale de la MDCT que pour ses versions moins optimisées. Le temps d'exécution de la MDCT *fixed point* doit évidemment être inférieur au temps d'exécution de l'algorithme de référence puisque celui-ci ne contient aucune optimisation. Ce temps peut toutefois être équivalent au temps d'exécution de la MDCT *floating point* : à performances équivalentes, l'algorithme *fixed point* rendra tout de même l'encodeur AAC plus performant en économisant les conversions *integer-floating point* à l'entrée et à la sortie du bloc MDCT. En effet, les données sont reçues par la MDCT en *integer* et devront être traitées en *integer* par le bloc de quantification à la sortie de la MDCT.

Le temps d'exécution de la MDCT a été mesuré sur base du code de l'annexe ?? compilé par les commandes CMake présentée à l'annexe ???. Le code permet de générer plusieurs exécutables en fonction de la variable d'environnement définie. Par défaut, l'exécutable permet de mesurer le temps d'exécution de l'algorithme de référence en , plusieurs exécutables sont générés, chacun d'eux permettant d'exécuter l'algorithme MDCT un certain nombre de fois (donné en paramètre à l'exécution). La fréquence du signal sinusoïdal utilisé en entrée de la MDCT peut également être donnée en paramètre mais elle ne varie pas entre les différentes exécutions afin de ne pas introduire d'aléatoire dans les mesures.

L'annexe En fonction du paramètre défini à la création des exécutables, le code exécuté sera celui de la MDCT de référence en *double floating point* (paramètre -DREF) ou de la MDCT basée sur la FFT de la librairie Ne10 en *floating point* (paramètre -DFLOATING\_POINT) ou en *fixed point* en *plain C* (paramètre -DFLOATING\_POINT)

perfs de la MDCT basée sur la FFT de FFTW3 en f32 pour comparaison : average run time : 7686.56 ns standard deviation : 6.61219e-07 ns

### **11.3 Perte de précision**

## 12 Améliorations possibles

- Fonction fenêtre intégrée aux opérations de pre twiddling
- Quantification intégrée au post twiddling
- tests automatisés
- code en C
- tests de performances plus poussés avec comparaison avec un algo existant



## Conclusion

Sur base du cahier des charges de début de stage, il a été décidé que mon travail s

## Références

- [1] EVS Website, “Page d’accueil d’EVS Broadcast Equipment.” [<https://evs.com>], consulté le 21 avril 2022.
- [2] EVS Website, “Page de présentation des produits commercialisés par EVS Broadcast Equipment.” [<https://evs.com/products>], consulté le 21 avril 2022.
- [3] K. Brandenburg, “Mp3 and aac explained,” *AES 17<sup>th</sup> International Conference on High Quality Audio Coding*, 1999.
- [4] J. Herre and S. Dick, “Psychoacoustic models for perceptual audio coding—a tutorial review,” *Applied Sciences*, vol. 9, p. 2854, 07 2019.
- [5] Wikipedia, “Codec.” [<https://en.wikipedia.org/wiki/Codec>], consulté le 2 mai 2022.
- [6] Wikipedia, “Moving picture experts group.” [[https://en.wikipedia.org/wiki/Moving\\_Picture\\_Experts\\_Group](https://en.wikipedia.org/wiki/Moving_Picture_Experts_Group)], consulté le 30 mars 2022.
- [7] “Information technology – Generic coding of moving pictures and associated audio information – Part 7 : Advanced Audio Coding (AAC),” standard, International Organization for Standardization, 2006. [<https://www.iso.org/standard/43345.html>].
- [8] “Information technology – Coding of audio-visual objects – Part 3 : Audio,” standard, International Organization for Standardization, 2019. [<https://www.iso.org/standard/76383.html>].
- [9] Wikipedia, “Psychoacoustics.” [<https://en.wikipedia.org/wiki/Psychoacoustics>], consulté le 2 mars 2022.
- [10] Wikipedia, “Modified discrete cosine transform.” [[https://en.wikipedia.org/wiki/Modified\\_discrete\\_cosine\\_transform](https://en.wikipedia.org/wiki/Modified_discrete_cosine_transform)], consulté le 17 septembre 2021.
- [11] Project Ne10 Website, “Documentation du projet Ne10.” [<http://projectne10.github.io/Ne10/doc/>], consulté le 9 mai 2022.
- [12] “The 3-Clause BSD License.” [<https://opensource.org/licenses/BSD-3-Clause>], consulté le 9 mai 2022.

## Liste des annexes

<b>A</b>	<b>CMake principal</b>	<b>I</b>
<b>B</b>	<b>Valeurs constantes des MDCT</b>	<b>II</b>
<b>C</b>	<b>Algorithmes de référence</b>	<b>III</b>
C.1	MDCT de référence en <i>float</i>	III
C.2	MDCT de référence en <i>integer</i>	III
<b>D</b>	<b>Génération d'un signal sinusoïdal</b>	<b>IV</b>
D.1	Génération d'un signal sinusoïdal en <i>float</i>	IV
D.2	Génération d'un signal sinusoïdal en <i>integer</i>	IV
<b>E</b>	<b>Implémentation de la MDCT basée sur la FFT de <i>FFTW3</i></b>	<b>V</b>
E.1	Header	V
E.2	Constructeur	V
E.3	Destructeur	VI
E.4	Fonction MDCT	VI
<b>F</b>	<b>Validation de la MDCT <i>FFTW3</i> en <i>float 32</i></b>	<b>VIII</b>
F.1	Code source	VIII
F.2	Compilation	IX
<b>G</b>	<b>Implémentation de la MDCT basée sur la FFT de <i>Ne10</i> en <i>floating point</i></b>	<b>X</b>
G.1	Header	X
G.2	Constructeur	X
G.3	Destructeur	XI
G.4	Fonction MDCT	XI
<b>H</b>	<b>Mesure des performances des FFT de <i>Ne10</i></b>	<b>XIII</b>
H.1	Code source	XIII
H.2	Compilation	XV
<b>I</b>	<b>Mesure des performances des FFT de <i>FFTW3</i></b>	<b>XVII</b>
I.1	Code source	XVII
I.2	Compilation	XVIII
<b>J</b>	<b>Implémentation de la MDCT basée sur la FFT de <i>Ne10</i> en <i>fixed point</i></b>	<b>XIX</b>
J.1	Header	XIX
J.2	Constructeur	XIX
J.3	Destructeur	XX
J.4	Fonction MDCT	XX

<b>K</b>	<b>Implémentation de la MDCT basée sur la FFT de <i>Ne10</i> en <i>fixed point</i> avec optimisations <i>Neon</i></b>	<b>XXII</b>
K.1	Header . . . . .	XXII
K.2	Constructeur . . . . .	XXII
K.3	Destructeur . . . . .	XXIII
K.4	Fonction MDCT . . . . .	XXIII

## A CMake principal

Fichier CMake principal placé à la racine du projet. Il permet de compiler :

- le projet *audio\_encoding* contenant les différentes MDCT et leurs tests : les commandes CMake de ce sous projet sont présentées dans les annexes suivantes sous le code qu'elles permettent de compiler;
- la librairie *Ne10* : les variables `NE10_LINUX_TARGET_ARCH` et `GNULINUX_PLATFORM` sont initialisées conformément aux recommandations de la documentation pour la compilation de la librairie.

```
cmake_minimum_required(VERSION 3.13)

set(NE10_LINUX_TARGET_ARCH armv7)
set(GNULINUX_PLATFORM ON)
if (CMAKE_BUILD_TYPE STREQUAL "DEBUG")
    set(BUILD_DEBUG ON)
endif (CMAKE_BUILD_TYPE STREQUAL "DEBUG")

add_subdirectory(audio_encoding)
add_subdirectory(Ne10)
```

## B Valeurs constantes des MDCT

Le fichier `mdct_constants.h` rassemble les valeurs constantes des MDCT pour une fenêtre d'entrée de 1024 échantillons.

```
// Sampling frequency: 48kHz
#define FS 48000

// Window length and derived constants
#define MDCT_WINDON_LEN 1024
#define MDCT_M (MDCT_WINDON_LEN > > 1) // spectrum size
#define MDCT_M2 (MDCT_WINDON_LEN > > 2) // fft size
#define MDCT_M4 (MDCT_WINDON_LEN > > 3)
#define MDCT_M32 (3 * (MDCT_WINDON_LEN > > 2))
#define MDCT_M52 (5 * (MDCT_WINDON_LEN > > 2))
```

## C Algorithmes de référence

### C.1 MDCT de référence en *float*

Algorithme de référence basé sur la formule mathématique de la MDCT. Le template permet de réaliser les calculs en *float* ou en *double*.

```
#include <cmath>

#include "mdct_constants.h"

template<typename FLOAT>
void ref_float_mdct(FLOAT *time_signal, FLOAT *spectrum)
{
    FLOAT scale = 2.0 / sqrt(MDCT_WINDON_LEN);
    FLOAT factor1 = 2.0 * M_PI / static_cast<FLOAT>(MDCT_WINDON_LEN);
    FLOAT factor2 = 0.5 + static_cast<FLOAT>(MDCT_M2);
    for (int k = 0; k < MDCT_M; ++k)
    {
        FLOAT result = 0.0;
        FLOAT factor3 = (k + 0.5) * factor1;
        for (int n = 0; n < MDCT_WINDON_LEN; ++n)
        {
            result += time_signal[n] * cos((static_cast<FLOAT>(n) + factor2) * factor3);
        }
        spectrum[k] = scale * result;
    }
}
```

### C.2 MDCT de référence en *integer*

Algorithme de référence basé sur la formule mathématique de la MDCT. Le spectre est calculé en *double* puis converti en *integer* sur 32 bits en représentation Q15.

```
#include <cassert>

#include "ref_mdct.h"

void ref_int_mdct(int16_t *time_signal, int32_t *spectrum)
{
    double scale = sqrt(MDCT_WINDON_LEN) / 2.0; // MDCT scale (2/sqrt(WIN_LEN)) + Q15 scale
    double factor1 = 2.0 * M_PI / MDCT_WINDON_LEN;
    double factor2 = 0.5 + MDCT_M2;
    for (int k = 0; k < MDCT_M; ++k)
    {
        double result = 0.0;
        double factor3 = (k + 0.5) * factor1;
        for (int n = 0; n < MDCT_WINDON_LEN; ++n)
        {
            result += time_signal[n] * cos((n + factor2) * factor3);
        }
        assert(round(result * scale) == static_cast<int32_t>(round(result * scale)));
        spectrum[k] = static_cast<int32_t>(round(result / scale));
    }
}
```

## D Génération d'un signal sinusoïdal

### D.1 Génération d'un signal sinusoïdal en *float*

Code de génération d'un signal sinusoïdal en *float* ou en *double*.

```
#include <cmath>

template<typename FLOAT>
void sin_float(FLOAT *out, int n_samples, double amplitude,
              double frequency, double phase_shift, int sampling_frequency)
{
    FLOAT omega = 2.0 * M_PI * frequency / static_cast<FLOAT>(sampling_frequency);
    for (int i = 0; i < n_samples; ++i)
    {
        out[i] = amplitude * sin(static_cast<FLOAT>(i) * omega + phase_shift);
    }
}
```

### D.2 Génération d'un signal sinusoïdal en *integer*

La génération du signal sinusoïdal en *integer* fait appel à la génération du signal sinusoïdal en *double* avant de convertir le résultat en *integer* (représentation Q15).

```
#include <cstring>

#include "sin_wave.h"

void sin_int(int16_t *out, int n_samples, double amplitude,
            double frequency, double phase_shift, int sampling_frequency)
{
    double scale = 1.0;
    if (abs(amplitude) < 1.0) scale *= amplitude;

    double *temp_sin = static_cast<double *>(malloc(n_samples * sizeof(double)));
    memset(temp_sin, 0, n_samples * sizeof(double));

    sin_float<double>(temp_sin, n_samples, scale, frequency, phase_shift, sampling_frequency);

    for (int i = 0; i < n_samples; ++i)
    {
        out[i] = static_cast<int16_t>(temp_sin[i] * pow(2.0, 15.0));
    }
}
```

## E Implémentation de la MDCT basée sur la FFT de *FFTW3*

### E.1 Header

Header de la classe `mdct_fftw3_f32` : MDCT basée sur la FFT de la librairie *FFTW3* en *float* (32 bits). La classe contient les structures de données `fft_in` et `fft_out`, le tableau de facteurs de `twiddle` utilisé pour le pre- et le post-processing et la configuration de la FFT (`fft_plan`). L'implémentation des fonctions de ce header est présentée dans les annexes suivantes.

```
#include <fftw3.h>

#include "mdct_constants.h"

class fftw3_mdct_f32
{
    private:
        fftwf_plan fft_plan;           // FFT configuration
        fftwf_complex *fft_in;         // FFT input buffer
        fftwf_complex *fft_out;        // FFT output buffer
        float twiddle[MDCT_M];

    public:
        fftw3_mdct_f32();
        ~fftw3_mdct_f32();
        void mdct(float *time_signal, float *spectrum);
        void imdct(float *spectrum, float *time_signal);
};
```

### E.2 Constructeur

Initialisation de la MDCT dans le constructeur de la classe `mdct_fftw3_f32` :

- Le tableau de `twiddle` est initialisé en *float* sur 32 bits;
- La FFT de *FFTW3* est initialisée en une dimension (pour l'audio) avec la taille de la FFT réduite à un quart de la taille de la fenêtre d'entrée par le pre-processing et avec l'option `FFTW_MEASURE` plus lente à l'initialisation mais qui permet d'optimiser le temps d'exécution de la FFT;
- Les tableaux contenant les données d'entrée (`fft_in`) et de sortie (`fft_out`) de la FFT sont alloués dynamiquement avec la fonction de *FFTW3* et ils sont passés en paramètre à la configuration de la FFT.

```
#include <cmath>

fftw3_mdct_f32::fftw3_mdct_f32()
{
    float alpha = M_PI / (8.f * MDCT_M);
    float omega = M_PI / MDCT_M;
    float scale = sqrt(sqrt(2.f / MDCT_M));

    for (int i = 0; i < MDCT_M2; ++i)
    {
        float x = omega*i + alpha;
        twiddle[2*i] = scale * cos(x);
        twiddle[2*i+1] = scale * sin(x);
    }
}
```



```

fft_in = (fftwf_complex *)fftwf_malloc(sizeof(fftwf_complex) * MDCT_M2);
fft_out = (fftwf_complex *)fftwf_malloc(sizeof(fftwf_complex) * MDCT_M2);
fft_plan = fftwf_plan_dft_1d(MDCT_M2, fft_in, fft_out, FFTW_FORWARD, FFTW_MEASURE);
}

```

### E.3 Destructeur

Destructeur de la classe `mdct_fftw3_f32` qui permet de libérer la mémoire allouée aux tableaux d'entrée et de sortie de la FFT et à sa configuration avec les fonctions appropriées fournies par la librairie *FFTW3*.

```

fftw3_mdct_f32::~fftw3_mdct_f32()
{
    fftwf_destroy_plan(fft_plan);
    fftwf_free(fft_in);
    fftwf_free(fft_out);
}

```

### E.4 Fonction MDCT

Implémentation de l'algorithme de MDCT basé sur la FFT de la librairie *FFTW3* :

- Initialisation du tableau d'entrée de la FFT : les opérations de *pre-twiddling* permettent de réduire la fenêtre d'entrée de la FFT;
- Appel de la fonction FFT de *FFTW3*;
- Calcul du spectre de fréquence : les opérations de *post-twiddling* permettent de calculer le spectre à partir des données de sortie de la FFT et des facteurs de *twiddle*.

```

void fftw3_mdct_f32::mdct(float *time_signal, float *spectrum)
{
    float *cos_tw = twiddle;
    float *sin_tw = cos_tw + 1;

    /* odd/even folding and pre-twiddle */
    float *xr = (float *)fft_in;
    float *xi = xr + 1;

    for (int i = 0; i < MDCT_M2; i += 2)
    {
        float r0 = time_signal[MDCT_M32-1-i] + time_signal[MDCT_M32+i];
        float i0 = time_signal[MDCT_M2+i] - time_signal[MDCT_M2-1-i];

        float c = cos_tw[i];
        float s = sin_tw[i];

        xr[i] = r0*c + i0*s;
        xi[i] = i0*c - r0*s;
    }

    for(int i = MDCT_M2; i < MDCT_M; i += 2)
    {
        float r0 = time_signal[MDCT_M32-1-i] - time_signal[-MDCT_M2+i];
        float i0 = time_signal[MDCT_M2+i] + time_signal[MDCT_M52-1-i];

        float c = cos_tw[i];
    }
}

```

```

        float s = sin_tw[i];

        xr[i] = r0*c + i0*s;
        xi[i] = i0*c - r0*s;
    }

    /* complex FFT of size MDCT_M2 */
    fftwf_execute(fft_plan);

    /* post-twiddle */
    xr = (float *)fft_out;
    xi = xr + 1;

    for (int i = 0; i < MDCT_M; i += 2)
    {
        float r0 = xr[i];
        float i0 = xi[i];

        float c = cos_tw[i];
        float s = sin_tw[i];

        spectrum[i] = -r0*c - i0*s;
        spectrum[MDCT_M-1-i] = -r0*s + i0*c;
    }
}

```

## F Validation de la MDCT *FFTW3* en *float 32*

### F.1 Code source

Test de la MDCT basée sur la FFT de *FFTW3* avec un signal d'entrée sinusoïdal à 200Hz :

- Génération et affichage d'un signal sinusoïdal à 200Hz;
- Calcul et affichage du spectre de fréquences de ce signal;
- Opération inverse de la MDCT et affichage du signal temporel calculé à partir du spectre.

```
#include <iomanip>
#include <iostream>

#include <cstring>

#include "mdct_constants.h"
#include "fftw3_mdct_f32.h"
#include "sin_wave.h"

/**
 * @brief MDCT algorithm calling the FFT of the fftw3 library
 * Code based on https://www.dsprelated.com/showcode/196.php
 */
int main(void)
{
    float time_in[MDCT_WINDON_LEN];           // input time signal
    sin_float(time_in, MDCT_WINDON_LEN, 0.9, 200.0, 0.0, FS);

    float time_out[MDCT_WINDON_LEN];           // output time signal (generated by the IMDCT)
    memset(time_out, 0, MDCT_WINDON_LEN*sizeof(float));

    float spectrum[MDCT_M];                    // frequency spectrum
    memset(spectrum, 0, MDCT_M*sizeof(float));

    fftw3_mdct_f32 fftw3_mdct;
    fftw3_mdct.mdct(time_in, spectrum);
    fftw3_mdct.imdct(spectrum, time_out);

    for (int i = 0; i < MDCT_WINDON_LEN; ++i)
    {
        std::cout << "time_in[" << std::setw(4) << i << "]" << std::setw(12) << time_in[i]
                    << " | _time_out[" << std::setw(4) << i << "]" << std::setw(12) << time_out[i]
                    << std::endl;
    }
    std::cout << std::endl;

    for (int i = 0; i < MDCT_M; ++i)
    {
        std::cout << "spectrum[" << std::setw(4) << i << "]"
                    << std::setw(12) << spectrum[i] << std::endl;
    }
    std::cout << std::endl;

    return 0;
}
```

## F.2 Compilation

Commandes CMake permettant de compiler le code d'exemple.

```
# MDCT using the fftw3 library f32
add_executable(fftw3_mdct_f32 test/validation/fftw3_example.cpp
               src/fftw3_mdct_f32.cpp src/sin_wave.cpp)
target_link_libraries(fftw3_mdct_f32 fftw3f)
```

## G Implémentation de la MDCT basée sur la FFT de *Ne10* en *floating point*

### G.1 Header

Header de la classe `mdct_ne10_f32_c` : MDCT basée sur la FFT de la librairie *Ne10* en *float* (32 bits). La classe contient les structures de données `fft_in` et `fft_out`, le tableau de facteurs de `twiddle` utilisé pour le pre- et le post-processing et la configuration de la FFT (`cfg`). L'implémentation des fonctions de ce header est présentée dans les annexes suivantes.

```
#pragma once

#include "mdct_constants.h"
#include "NE10.h"

class ne10_mdct_f32_c
{
private:
    ne10_fft_cfg_float32_t cfg; // Ne10 configuration
    ne10_fft_cpx_float32_t fft_in[MDCT_M2]__attribute__((aligned(16))); // Ne10 FFT input buffer
    ne10_fft_cpx_float32_t fft_out[MDCT_M2]__attribute__((aligned(16))); // Ne10 FFT output buffer
    float twiddle[MDCT_M]__attribute__((aligned(16))); // twiddle factors

public:
    ne10_mdct_f32_c();
    ~ne10_mdct_f32_c();
    void mdct(float *time_signal, float *spectrum);
};
```

### G.2 Constructeur

Initialisation de la MDCT dans le constructeur de la classe `mdct_ne10_f32_c` :

- Le tableau de `twiddle` est initialisé en *float* sur 32 bits;
- La configuration de la FFT de *Ne10* est initialisée en *complex to complex* en *float 32* avec en paramètre la taille de la fenêtre de la FFT réduite à un quart de la taille de la fenêtre d'entrée.

```
ne10_mdct_f32_c::ne10_mdct_f32_c()
{
    float alpha = M_PI / (8.0 * static_cast<float>(MDCT_M));
    float omega = M_PI / static_cast<float>(MDCT_M);
    float scale = sqrt(sqrt(2.0 / static_cast<float>(MDCT_M)));
    for (int i = 0; i < MDCT_M2; ++i)
    {
        float x = omega * i + alpha;
        twiddle[2*i] = static_cast<float>(scale * cos(x));
        twiddle[2*i+1] = static_cast<float>(scale * sin(x));
    }

    cfg = ne10_fft_alloc_c2c_float32_c(MDCT_M2);
}
```

### G.3 Destructeur

Destructeur de la classe `mdct_ne10_f32_c` qui permet de libérer la mémoire allouée à la configuration de la FFT avec la fonction appropriée de la librairie *Ne10*.

```
ne10_mdct_f32_c::~ne10_mdct_f32_c()
{
    ne10_fft_destroy_c2c_float32(cfg);
}
```

### G.4 Fonction MDCT

Implémentation de l'algorithme de MDCT basé sur la FFT de la librairie *Ne10* en *float 32* et en *plain C* :

- Initialisation du tableau d'entrée de la FFT : les opérations de *pre-twiddling* permettent de réduire la fenêtre d'entrée de la FFT;
- Appel de la fonction FFT de *Ne10*;
- Calcul du spectre de fréquence : les opérations de *post-twiddling* permettent de calculer le spectre à partir des données de sortie de la FFT et des facteurs de *twiddle*.

```
void ne10_mdct_f32_c::mdct(float *time_signal, float *spectrum)
{
    // pre-twiddling
    float *cos_tw = twiddle;
    float *sin_tw = cos_tw + 1;
    for (int i = 0; i < MDCT_M2; i += 2)
    {
        float r0 = time_signal[MDCT_M32-1-i] + time_signal[MDCT_M32+i];
        float i0 = time_signal[MDCT_M2+i] - time_signal[MDCT_M2-1-i];

        float c = cos_tw[i];
        float s = sin_tw[i];

        fft_in[i/2].r = r0*c + i0*s;
        fft_in[i/2].i = i0*c - r0*s;
    }

    for (int i = MDCT_M2; i < (MDCT_M); i += 2)
    {
        float r0 = time_signal[MDCT_M32-1-i] - time_signal[-MDCT_M2+i];
        float i0 = time_signal[MDCT_M2+i] + time_signal[MDCT_M52-1-i];

        float c = cos_tw[i];
        float s = sin_tw[i];

        fft_in[i/2].r = r0*c + i0*s;
        fft_in[i/2].i = i0*c - r0*s;
    }

    // FFT
    ne10_fft_c2c_1d_float32_c(fft_out, fft_in, cfg, 0);

    // post-twiddling
    for (int i = 0; i < (MDCT_M); i += 2)
    {
        float r0 = fft_out[i/2].r;
```

```

    float i0 = fft_out[i/2].i;

    float c = cos_tw[i];
    float s = sin_tw[i];

    spectrum[i] = -r0*c - i0*s;
    spectrum[(MDCT_M)-1-i] = -r0*s + i0*c;
}

```

## H Mesure des performances des FFT de *Ne10*

### H.1 Code source

Code permettant de tester la vitesse d'exécution moyenne de différentes FFT proposées par la librairie *Ne10*. La moyenne est calculée sur 10.000.000 exécutions avec des données aléatoires en entrée de la FFT différentes pour chaque exécution. Les variables de préprocesseur définies à la compilation permettent sur base du même code de mesurer le temps d'exécution moyen avec écart type :

- de la FFT *complex to complex* en *float 32* en *plain C* ou avec les optimisations Neon;
- de la FFT *complex to complex* en *integer 32* en *plain C* ou avec les optimisations Neon;
- de la FFT *complex to complex* en *integer 16* en *plain C* ou avec les optimisations Neon.

```
#include <iomanip>
#include <iostream>
#include <limits>

#include <cmath>
#include <cstring>

#include "mdct_constants.h"
#include "Timers.h"
#include "NE10.h"

#ifdef F32          // 32 bits floating point arithmetic

#define INPUT_RANGE      1.8
#define INPUT_DATA       ne10_fft_cpx_float32_t
#define OUTPUT_DATA      ne10_fft_cpx_float32_t
#define FFT_CONFIG       ne10_fft_cfg_float32_t
#define DESTROY_CONFIG   ne10_fft_destroy_c2c_float32

#ifdef NEON
#define ALLOC_CONFIG      ne10_fft_alloc_c2c_float32_neon
#define PERFORM_FFT       ne10_fft_c2c_1d_float32_neon
#else
#define ALLOC_CONFIG      ne10_fft_alloc_c2c_float32_c
#define PERFORM_FFT       ne10_fft_c2c_1d_float32_c
#endif

#elif I32          // 32 bits fixed point arithmetic

#define INPUT_RANGE      std::numeric_limits<int16_t>::max()*2
#define INPUT_DATA       ne10_fft_cpx_int32_t
#define OUTPUT_DATA      ne10_fft_cpx_int32_t
#define FFT_CONFIG       ne10_fft_cfg_int32_t
#define DESTROY_CONFIG   ne10_fft_destroy_c2c_int32

#ifdef NEON
#define ALLOC_CONFIG      ne10_fft_alloc_c2c_int32_neon
#define PERFORM_FFT       ne10_fft_c2c_1d_int32_neon
#else
#define ALLOC_CONFIG      ne10_fft_alloc_c2c_int32_c
#define PERFORM_FFT       ne10_fft_c2c_1d_int32_c
#endif

#endif
```



```

#else                                // 16 bits fixed point arithmetic

#define INPUT_RANGE                    std::numeric_limits<int16_t>::max()*2
#define INPUT_DATA                     ne10_fft_cpx_int16_t
#define OUTPUT_DATA                    ne10_fft_cpx_int16_t
#define FFT_CONFIG                      ne10_fft_cfg_int16_t
#define ALLOC_CONFIG                   ne10_fft_alloc_c2c_int16
#define DESTROY_CONFIG                 ne10_fft_destroy_c2c_int16

#ifdef NEON
#define PERFORM_FFT                    ne10_fft_c2c_1d_int16_neon
#else
#define PERFORM_FFT                    ne10_fft_c2c_1d_int16_c
#endif

#endif

#define RUNS                          10000000
#define FFT_SCALE_FLAG                0

int main()
{
    // print which FFT will be tested
#ifdef F32
#ifdef NEON
        std::cout << "FFT_Ne10_f32_Neon" << std::endl;
#else
        std::cout << "FFT_Ne10_f32_plain_C" << std::endl;
#endif

#elif I32
#ifdef NEON
        std::cout << "FFT_Ne10_i32_Neon" << std::endl;
#else
        std::cout << "FFT_Ne10_i32_plain_C" << std::endl;
#endif
#else
#ifdef NEON
        std::cout << "FFT_Ne10_i16_Neon" << std::endl;
#else
        std::cout << "FFT_Ne10_i16_plain_C" << std::endl;
#endif
#endif

    // feed the random
    srand(static_cast<unsigned>(time(0)));

    // inititalize the configuration
    FFT_CONFIG cfg = ALLOC_CONFIG(MDCT_M2);

    // start the loop executing the FFTs
    int64_t *runtimes = static_cast<int64_t *>(malloc(RUNS * sizeof(int64_t)));
    for (int i = 0; i < RUNS; ++i)
    {
        // initialize an empty spectrum
        OUTPUT_DATA spectrum[MDCT_M2]__attribute__((aligned(16)));
        memset(&spectrum, 0, (MDCT_M2)*sizeof(OUTPUT_DATA));
    }

```

```

    // generate random input data
    INPUT_DATA time_signal[MDCT_M2]__attribute__((aligned(16)));
    for (int i = 0; i < MDCT_M2; ++i)
    {
        time_signal[i].r = INPUT_RANGE * rand() / RAND_MAX - INPUT_RANGE / 2;
        time_signal[i].i = INPUT_RANGE * rand() / RAND_MAX - INPUT_RANGE / 2;
    }

    // perform the FFT and measure the run time
    EvsHwLGPL::CTimers timer;
    timer.Start();
#ifdef F32
    PERFORM_FFT(time_signal, spectrum, cfg, 0);
#else
    PERFORM_FFT(time_signal, spectrum, cfg, 0, FFT_SCALE_FLAG);
#endif
    timer.Stop();
    runtimes[i] = timer.GetTimeElapsed();
}

// clean
DESTROY_CONFIG(cfg);

// compute the average
double avg = 0.0;
for (int i = 0; i < RUNS; ++i) avg += static_cast<double>(runtimes[i]);
avg = avg / static_cast<double>(RUNS);
std::cout << "average_run_time:_" << avg << "_ns" << std::endl;

// compute the standard deviation
double dev = 0.0;
for (int i = 0; i < RUNS; ++i) dev += static_cast<double>(runtimes[i]) - avg;
dev = dev * dev / static_cast<double>(RUNS);
dev = sqrt(dev);
std::cout << "standard_deviation:_" << dev << "_ns" << std::endl;

return 0;
}

```

## H.2 Compilation

Commandes CMake utilisées pour générer les exécutables permettant de mesurer le temps d'exécution de différentes FFT proposées par la librairie *Ne10*. En fonction des variables de préprocesseur définies, les exécutables suivants sont générés :

- `run_fft_f32_c` est généré si la variable `F32` est définie pour mesurer le temps d'exécution de la FFT *float 32 plain C*;
- `run_fft_i32_c` est généré si la variable `I32` est définie pour mesurer le temps d'exécution de la FFT *integer 32 plain C*;
- `run_fft_i16_c` est généré par défaut pour mesurer le temps d'exécution de la FFT *integer 16 plain C*;
- `run_fft_f32_neon` est généré si les variables `F32` et `NEON` sont définies pour mesurer le temps d'exécution de la FFT *float 32* avec optimisations *Neon*;
- `run_fft_i32_neon` est généré si les variables `I32` et `NEON` sont définies pour mesurer le temps d'exécution de la FFT *integer 32* avec optimisations *Neon*;

- `run_fft_i16_neon` est généré si la variable `NEON` est définie pour mesurer le temps d'exécution de la FFT *integer 16* avec optimisations *Neon*.

```
# Ne10 FFT performance (float32 plain C)
add_executable(run_ne10_fft_f32_c test/performance/run_ne10_fft.cpp src/Timers.cpp)
target_compile_definitions(run_ne10_fft_f32_c PUBLIC -DF32)
target_link_libraries(run_ne10_fft_f32_c NE10)

# Ne10 FFT performance (int32 plain C)
add_executable(run_ne10_fft_i32_c test/performance/run_ne10_fft.cpp src/Timers.cpp)
target_compile_definitions(run_ne10_fft_i32_c PUBLIC -DI32)
target_link_libraries(run_ne10_fft_i32_c NE10)

# Ne10 FFT performance (int16 plain C)
add_executable(run_ne10_fft_i16_c test/performance/run_ne10_fft.cpp src/Timers.cpp)
target_compile_definitions(run_ne10_fft_i16_c PUBLIC -DI16)
target_link_libraries(run_ne10_fft_i16_c NE10)

# Ne10 FFT performance (float32 with neon optimizations)
add_executable(run_ne10_fft_f32_neon test/performance/run_ne10_fft.cpp src/Timers.cpp)
target_compile_definitions(run_ne10_fft_f32_neon PUBLIC -DF32 -DNEON)
target_link_libraries(run_ne10_fft_f32_neon NE10)

# Ne10 FFT performance (int32 with neon optimizations)
add_executable(run_ne10_fft_i32_neon test/performance/run_ne10_fft.cpp src/Timers.cpp)
target_compile_definitions(run_ne10_fft_i32_neon PUBLIC -DI32 -DNEON)
target_link_libraries(run_ne10_fft_i32_neon NE10)

# Ne10 FFT performance (int16 with neon optimizations)
add_executable(run_ne10_fft_i16_neon test/performance/run_ne10_fft.cpp src/Timers.cpp)
target_compile_definitions(run_ne10_fft_i16_neon PUBLIC -DI16 -DNEON)
target_link_libraries(run_ne10_fft_i16_neon NE10)
```

# I Mesure des performances des FFT de *FFTW3*

## I.1 Code source

Code permettant de tester la vitesse d'exécution moyenne de la FFT en *float 32* de la librairie *FFTW3*. La moyenne est calculée sur 10.000.000 exécutions avec des données aléatoires en entrée de la FFT différentes pour chaque exécution.

```
#include <iomanip>
#include <iostream>

#include <cmath>

#include <fftw3.h>

#include "mdct_constants.h"
#include "Timers.h"

#define RUNS 10000000

int main()
{
    // print which FFT will be tested
    std::cout << "FFT_FFTW3_f32_plain_C" << std::endl;

    // feed the random
    srand(static_cast<unsigned>(time(0)));

    // start the loop executing the FFTs
    int64_t *runtimes = static_cast<int64_t *>(malloc(RUNS * sizeof(int64_t)));
    for (int i = 0; i < RUNS; ++i)
    {
        // initialize an empty spectrum
        fftwf_complex *fft_out = (fftwf_complex *)fftwf_malloc(sizeof(fftwf_complex) * MDCT_M2);

        // generate random input data
        fftwf_complex *fft_in = (fftwf_complex *)fftwf_malloc(sizeof(fftwf_complex) * MDCT_M2);
        float *x = (float *)fft_in;
        for (int i = 0; i < MDCT_M2; ++i)
        {
            x[i] = 1.8f * rand() / RAND_MAX - 1.8f / 2.0f;
        }

        // initialize the configuration
        fftwf_plan fft_plan = fftwf_plan_dft_1d(MDCT_M2, fft_in, fft_out,
            FFTW_FORWARD, FFTW_MEASURE);

        // perform the FFT and measure the run time
        EvsHwLGPL::CTimers timer;
        timer.Start();
        fftwf_execute(fft_plan);
        timer.Stop();
        runtimes[i] = timer.GetTimeElapsed();

        // clean
        fftwf_destroy_plan(fft_plan);
        fftwf_free(fft_in);
        fftwf_free(fft_out);
    }
}
```

```

    // compute the average
    double avg = 0.0;
    for (int i = 0; i < RUNS; ++i) avg += static_cast<double>(runtimes[i]);
    avg = avg / static_cast<double>(RUNS);
    std::cout << "average_run_time:_" << avg << "_ns" << std::endl;

    // compute the standard deviation
    double dev = 0.0;
    for (int i = 0; i < RUNS; ++i) dev += static_cast<double>(runtimes[i]) - avg;
    dev = dev * dev / static_cast<double>(RUNS);
    dev = sqrt(dev);
    std::cout << "standard_deviation:_" << dev << "_ns" << std::endl;

    return 0;
}

```

## I.2 Compilation

Commandes CMake utilisées pour générer l'exécutable permettant de mesurer le temps d'exécution de la FFT *float32* de la librairie *FFTW3*.

```

# FFTW3 FFT performance (float32)
add_executable(run_fftw3_fft_f32 test/performance/run_fftw3_fft_f32.cpp src/Timers.cpp)
target_link_libraries(run_fftw3_fft_f32 fftw3f)

```

## J Implémentation de la MDCT basée sur la FFT de *Ne10* en *fixed point*

### J.1 Header

Header de la classe `mdct_ne10_i32_c` : MDCT basée sur la FFT de la librairie *Ne10* en *integer* (32 bits). La classe contient les structures de données `fft_in` en représentation Q1.15 et `fft_out` en Q9.15, le tableau de facteurs de `twiddle` utilisé pour le pre- et le post-processing et la configuration de la FFT (`cfg`). L'implémentation des fonctions de ce header est présentée dans les annexes suivantes.

```
#pragma once

#include "mdct_constants.h"
#include "NE10.h"

class ne10_mdct_i32_c
{
private:
    ne10_fft_cfg_int32_t cfg; // Ne10 configuration
    ne10_fft_cpx_int32_t fft_in[MDCT_M2] __attribute__((aligned(16))); // Ne10 FFT input buffer
    // Q1.15
    ne10_fft_cpx_int32_t fft_out[MDCT_M2] __attribute__((aligned(16))); // Ne10 FFT output buffer
    // Q9.15
    int16_t twiddle[MDCT_M] __attribute__((aligned(16))); // MDCT twiddle factors

public:
    ne10_mdct_i32_c();
    ~ne10_mdct_i32_c();
    void mdct(int16_t *time_signal, int32_t *spectrum);
};
```

### J.2 Constructeur

Initialisation de la MDCT dans le constructeur de la classe `mdct_ne10_i32_c` :

- Le tableau de `twiddle` est initialisé en *double* puis converti en *integer* (représentation Q15);
- La configuration de la FFT de *Ne10* est initialisée en *complex to complex* en *integer 32* avec en paramètre la taille de la fenêtre de la FFT réduite à un quart de la taille de la fenêtre d'entrée.

```
ne10_mdct_i32_c::ne10_mdct_i32_c()
{
    // initialize the twiddling factors
    double alpha = M_PI / (8.0*MDCT_M);
    double omega = M_PI / MDCT_M;
    double scale = sqrt(sqrt(2.0 / static_cast<double>MDCT_M));
    for (int i = 0; i < MDCT_M2; ++i)
    {
        double x = omega * i + alpha;
        twiddle[2*i] = static_cast<int16_t>(cos(x)*scale*pow(2.0, 15.0));
        twiddle[2*i+1] = static_cast<int16_t>(sin(x)*scale*pow(2.0, 15.0));
    }

    // initialize the Ne10 FFT configuration
    cfg = ne10_fft_alloc_c2c_int32_c(MDCT_M2);
}
```

### J.3 Destructeur

Destructeur de la classe `mdct_ne10_i32_c` qui permet de libérer la mémoire allouée à la configuration de la FFT avec la fonction appropriée de la librairie *Ne10*.

```
ne10_mdct_i32_c::~ne10_mdct_i32_c()
{
    ne10_fft_destroy_c2c_int32(cfg);
}
```

### J.4 Fonction MDCT

Implémentation de l'algorithme de MDCT basé sur la FFT de la librairie *Ne10* en *integer 32* et en *plain C* :

- Initialisation du tableau d'entrée de la FFT : les opérations de *pre-twiddling* permettant de réduire la fenêtre d'entrée de la FFT sont faites en algorithmique *fixed point*;
- Appel de la fonction FFT de *Ne10*;
- Calcul du spectre de fréquence : les opérations de *post-twiddling* permettant de calculer le spectre à partir des données de sortie de la FFT et des facteurs de *twiddle* sont faites en algorithmique *fixed point*.

```
void ne10_mdct_i32_c::mdct(int16_t *time_signal, int32_t *spectrum)
{
    // pre-twiddling
    // fft_in = (Q1.15 + Q1.15) * Q1.15/4 + (Q1.15 + Q1.15) * Q1.15/4
    //          1/4 Q1.30 + 1/4 Q1.30 + 1/4 Q1.30 + 1/4 Q1.30 -> Q1.30
    //          >>7 -> Q1.23 + 8 bits reserved for the FFT
    int16_t *cos_tw = twiddle;
    int16_t *sin_tw = cos_tw + 1;
    for (int i = 0; i < MDCT_M2; i += 2)
    {
        int32_t r0 = static_cast<int32_t>(time_signal[MDCT_M32-1-i]) + time_signal[MDCT_M32+i];
        int32_t i0 = static_cast<int32_t>(time_signal[MDCT_M2+i]) - time_signal[MDCT_M2-1-i];

        int16_t c = cos_tw[i];
        int16_t s = sin_tw[i];

        fft_in[i/2].r = (((r0*c)+64)>>7) + (((i0*s)+64)>>7);
        fft_in[i/2].i = (((i0*c)+64)>>7) - (((r0*s)+64)>>7);
    }

    for (int i = MDCT_M2; i < MDCT_M; i += 2)
    {
        int32_t r0 = static_cast<int32_t>(time_signal[MDCT_M32-1-i]) - time_signal[-MDCT_M2+i];
        int32_t i0 = static_cast<int32_t>(time_signal[MDCT_M2+i]) + time_signal[MDCT_M52-1-i];

        int16_t c = cos_tw[i];
        int16_t s = sin_tw[i];

        fft_in[i/2].r = (((r0*c)+64)>>7) + (((i0*s)+64)>>7);
        fft_in[i/2].i = (((i0*c)+64)>>7) - (((r0*s)+64)>>7);
    }

    // perform the FFT
    ne10_fft_c2c_1d_int32_c(fft_out, fft_in, cfg, 0, 0);
}
```

```

// post-twiddling
// spectrum = Q9.23>>8 * Q1.15/4 + Q9.23>>8 * Q1.15/4
//           = Q9.15 * Q1.15/4 + Q9.15 * Q1.15/4
//           = Q10.30/4 + Q10.30/4
//           = Q11.30/4
//           = Q9.30 >> 15 = Q9.15
for (int i = 0; i < MDCT_M; i += 2)
{
    int32_t r0 = fft_out[i/2].r;
    int32_t i0 = fft_out[i/2].i;

    int16_t c = cos_tw[i];
    int16_t s = sin_tw[i];

    spectrum[i] = ((((-r0+128)>>8)*c+16384)>>15) - (((i0+128)>>8)*s+16384)>>15);
    spectrum[MDCT_M-1-i] = ((((-r0+128)>>8)*s+16384)>>15) + (((i0+128)>>8)*c+16384)>>15);
}
}

```



## K Implémentation de la MDCT basée sur la FFT de *Ne10* en *fixed point* avec optimisations *Neon*

### K.1 Header

Header de la classe `mdct_ne10_i32_neon` : MDCT basée sur la FFT de la librairie *Ne10* en *integer* (32 bits) optimisée par l'utilisation des opérations SIMD Neon. La classe contient les structures de données `fft_in` en représentation Q1.15 et `fft_out` en Q9.15, les tableaux de facteurs de twiddle utilisé pour le pre- et le post-processing et la configuration de la FFT (`cfg`). Contrairement aux autres implémentations, les facteurs de twiddle non sont pas rassemblés dans un seul tableau. Les tableaux de facteurs de *pre-twiddling* et de *post-twiddling* sont séparés car ils sont utilisés en 16 bits pour le *pre-twiddling* et en 32 bits pour le *post-twiddling*. Chacun de ses tableaux est séparé en deux car pour que chaque moitié puisse être initialisée dans un ordre qui facilite l'utilisation des opérations SIMD. L'implémentation des fonctions de ce header est présentée dans les annexes suivantes.

```
#pragma once

#include <arm_neon.h>
#include "mdct_constants.h"
#include "NE10.h"

class ne10_mdct_i32_neon
{
private:
    ne10_fft_cfg_int32_t cfg; // Ne10 configuration
    ne10_fft_cpx_int32_t fft_in[MDCT_M2]__attribute__((aligned(16))); // Ne10 FFT input buffer
    ne10_fft_cpx_int32_t fft_out[MDCT_M2]__attribute__((aligned(16))); // Ne10 FFT output buffer
    int16_t pretwiddle_start[MDCT_M2]__attribute__((aligned(16))); // pre-twiddle factors
    int16_t pretwiddle_end[MDCT_M2]__attribute__((aligned(16))); // second half is stored
    // in reversed order
    int32_t posttwiddle_start[MDCT_M2]__attribute__((aligned(16))); // post-twiddle factors
    int32_t posttwiddle_end[MDCT_M2]__attribute__((aligned(16))); // second half is stored
    // in reversed order

public:
    ne10_mdct_i32_neon();
    ~ne10_mdct_i32_neon();
    void mdct(int16_t *time_signal, int32_t *spectrum);
};
```

### K.2 Constructeur

Initialisation de la MDCT dans le constructeur de la classe `mdct_ne10_i32_neon` :

- Le tableau de twiddle est initialisé en *double* puis converti en *integer* (représentation Q15 en 16 bits pour le *pre-twiddling* et en 32 bits pour le *post-twiddling*) : la première moitié des tableaux est rangée à l'endroit dans les tableaux `pretwiddle_start` et `posttwiddle_start` tandis que la seconde est rangée à l'envers dans les tableaux `pretwiddle_end` et `posttwiddle_end`;
- La configuration de la FFT de *Ne10* est initialisée an *complex to complex* en *integer 32* avec en paramètre la taille de la fenêtre de la FFT réduite à un quart de la taille de la fenêtre d'entrée avec la fonction adaptée pour l'exécution d'une FFT optimisée avec les instructions SIMD Neon.

```

ne10_mdct_i32_neon::ne10_mdct_i32_neon()
{
    double alpha = M_PI / (8.0*MDCT_M);
    double omega = M_PI / MDCT_M;
    double scale = sqrt(sqrt(2.0 / static_cast<double>MDCT_M));
    for (int i = 0; i < MDCT_M4; ++i)
    {
        double start = omega * (i) + alpha;
        double end = omega * (i+MDCT_M4) + alpha;
        double cos_start = cos(start);
        double sin_start = sin(start);
        double cos_end = cos(end);
        double sin_end = sin(end);
        pretwiddle_start[2*i] = static_cast<int16_t>(cos_start*scale*pow(2.0, 15.0));
        pretwiddle_start[2*i+1] = static_cast<int16_t>(sin_start*scale*pow(2.0, 15.0));
        pretwiddle_end[MDCT_M2-2*i-2] = static_cast<int16_t>(cos_end*scale*pow(2.0, 15.0));
        pretwiddle_end[MDCT_M2-2*i-1] = static_cast<int16_t>(sin_end*scale*pow(2.0, 15.0));
        posttwiddle_start[2*i] = static_cast<int32_t>(cos_start*scale*pow(2.0, 31.0));
        posttwiddle_start[2*i+1] = static_cast<int32_t>(sin_start*scale*pow(2.0, 31.0));
        posttwiddle_end[MDCT_M2-2*i-2] = static_cast<int32_t>(cos_end*scale*pow(2.0, 31.0));
        posttwiddle_end[MDCT_M2-2*i-1] = static_cast<int32_t>(sin_end*scale*pow(2.0, 31.0));
    }

    cfg = ne10_fft_alloc_c2c_int32_neon(MDCT_M2);
}

```

### K.3 Destructeur

Destructeur de la classe `mdct_ne10_i32_c` qui permet de libérer la mémoire allouée à la configuration de la FFT avec la fonction appropriée de la librairie *Ne10*.

```

ne10_mdct_i32_neon::~ne10_mdct_i32_neon()
{
    ne10_fft_destroy_c2c_int32(cfg);
}

```

### K.4 Fonction MDCT

Implémentation de l'algorithme de MDCT basé sur la FFT de la librairie *Ne10* en *integer 32* avec utilisation des instructions SIMD Neon :

- Initialisation du tableau d'entrée de la FFT : les opérations de *pre-twiddling* permettant de réduire la fenêtre d'entrée de la FFT sont faites en algorithmique *fixed point*. L'utilisation des fonctions SIMD permet d'effectuer quatre opérations en parallèle afin de réduire le temps d'exécution. Les facteurs de *pre-twiddling* en 16 bits transforment le signal d'entrée en 16 bits en un tableau d'entrée de la FFT en 32 bits;
- Appel de la fonction FFT de *Ne10* optimisée par l'utilisation des instructions SIMD Neon;
- Calcul du spectre de fréquence : les opérations de *post-twiddling* permettant de calculer le spectre à partir des données de sortie de la FFT et des facteurs de *twiddle* sont faites en algorithmique *fixed point*. Les fonctions SIMD permettent d'effectuer deux ou quatre opérations en parallèle afin de réduire le temps d'exécution.

```

void ne10_mdct_i32_neon::mdct(int16_t *time_signal, int32_t *spectrum)
{
    // see the twiddling_loops.nlsx file for more details

    // for i from 0 to 254, step 2

    // r[ 0 -> 127, pas 1] = time_signal[ 767 -> 513, pas 2] * c[ 0 -> 127, pas 1]
    //                        + time_signal[ 768 -> 1022, pas 2] * c[ 0 -> 127, pas 1]
    //                        + time_signal[ 256 -> 510, pas 2] * s[ 0 -> 127, pas 1]
    //                        + -time_signal[ 255 -> 1, pas 2] * s[ 0 -> 127, pas 1]

    // i[ 0 -> 127, pas 1] = time_signal[ 256 -> 510, pas 2] * c[ 0 -> 127, pas 1]
    //                        + -time_signal[ 255 -> 1, pas 2] * c[ 0 -> 127, pas 1]
    //                        + time_signal[ 767 -> 513, pas 2] * s[ 0 -> 127, pas 1]
    //                        + -time_signal[ 768 -> 1022, pas 2] * s[ 0 -> 127, pas 1]

    // fft_in[i/2].r = time_signal[M32-1-i] * cos_tw[i] + time_signal[M32+i] * cos_tw[i]
    //                  + time_signal[M2+i] * sin_tw[i] + (-time_signal[M2-1-i]) * sin_tw[i]

    // fft_in[i/2].i = time_signal[M2+i] * cos_tw[i] + (-time_signal[M2-1-i]) * cos_tw[i]
    //                  + (-time_signal[M32-1-i]) * sin_tw[i] + (-time_signal[M32+i]) * sin_tw[i]

    // for i from 510 to 256, step 2

    // r[255 -> 128, pas 1] = time_signal[ 257 -> 511, pas 2] * c[255 -> 128, pas 1]
    //                        + -time_signal[ 254 -> 0, pas 2] * c[255 -> 128, pas 1]
    //                        + time_signal[ 766 -> 512, pas 2] * s[255 -> 128, pas 1]
    //                        + time_signal[ 769 -> 1023, pas 2] * s[255 -> 128, pas 1]

    // i[255 -> 128, pas 1] = time_signal[ 766 -> 512, pas 2] * c[255 -> 128, pas 1]
    //                        + time_signal[ 769 -> 1023, pas 2] * c[255 -> 128, pas 1]
    //                        + -time_signal[ 257 -> 511, pas 2] * s[255 -> 128, pas 1]
    //                        + time_signal[ 254 -> 0, pas 2] * s[255 -> 128, pas 1]

    // fft_in[i/2].r = time_signal[M32-1-i] * cos_tw[i] + (-time_signal[-M2+i]) * cos_tw[i]
    //                  + time_signal[M2+i] * sin_tw[i] + time_signal[M52-1-i] * sin_tw[i]

    // fft_in[i/2].i = time_signal[M2+i] * cos_tw[i] + time_signal[M52-1-i] * cos_tw[i]
    //                  + (-time_signal[M32-1-i]) * sin_tw[i] + time_signal[-M2+i] * sin_tw[i]

    for (int i = 0; i < MDCT_M2; i += 8)
    {
        // tx.val[0] -> odd indexes
        // tx.val[1] -> even indexes
        int16x4x2_t t1 = vld2_s16(time_signal+MDCT_M2+i);
        int16x4x2_t t2 = vld2_s16(time_signal+MDCT_M2-8-i);
        int16x4x2_t t3 = vld2_s16(time_signal+MDCT_M32+i);
        int16x4x2_t t4 = vld2_s16(time_signal+MDCT_M32-8-i);

        t2.val[0] = (int16x4_t)vrev64_s32((int32x2_t)vrev32_s16(t2.val[0]));
        // reverse the t2 even values: 0 2 4 6 -> 6 4 2 0
        t2.val[1] = (int16x4_t)vrev64_s32((int32x2_t)vrev32_s16(t2.val[1]));
        // reverse the t2 odd values: 1 3 5 7 -> 7 5 3 1
        t4.val[0] = (int16x4_t)vrev64_s32((int32x2_t)vrev32_s16(t4.val[0]));
        // reverse the t4 even values: 0 2 4 6 -> 6 4 2 0
        t4.val[1] = (int16x4_t)vrev64_s32((int32x2_t)vrev32_s16(t4.val[1]));
        // reverse the t4 odd values: 1 3 5 7 -> 7 5 3 1
    }
}

```

```

// x_tw.val[0] -> cos twiddle
// x_tw.val[1] -> sin twiddle
int16x4x2_t start_tw = vld2_s16(pretwiddle_start+i);
int16x4x2_t end_tw = vld2_s16(pretwiddle_end+i);

// start.val[0] -> real part
// start.val[1] -> imaginary part
int32x4x2_t start;
start.val[0] = vshrq_n_s32(
    vaddq_s32(
        vaddq_s32(
            vmull_s16(t4.val[1], start_tw.val[0]),
            vmull_s16(t3.val[0], start_tw.val[0])),
        vaddq_s32(
            vmull_s16(t1.val[0], start_tw.val[1]),
            vmull_s16(vneg_s16(t2.val[1]), start_tw.val[1]))),
    7);
start.val[1] = vshrq_n_s32(
    vaddq_s32(
        vaddq_s32(
            vmull_s16(t1.val[0], start_tw.val[0]),
            vmull_s16(vneg_s16(t2.val[1]), start_tw.val[0])),
        vaddq_s32(
            vmull_s16(vneg_s16(t4.val[1]), start_tw.val[1]),
            vmull_s16(vneg_s16(t3.val[0]), start_tw.val[1]))),
    7);

// end.val[0] -> real part
// end.val[1] -> imaginary part
int32x4x2_t end;
end.val[0] = vshrq_n_s32(
    vaddq_s32(
        vaddq_s32(vmull_s16(t1.val[1], end_tw.val[0]),
            vmull_s16(vneg_s16(t2.val[0]), end_tw.val[0])),
        vaddq_s32(
            vmull_s16(t4.val[0], end_tw.val[1]),
            vmull_s16(t3.val[1], end_tw.val[1]))),
    7);
end.val[1] = vshrq_n_s32(
    vaddq_s32(
        vaddq_s32(
            vmull_s16(t4.val[0], end_tw.val[0]),
            vmull_s16(t3.val[1], end_tw.val[0])),
        vaddq_s32(
            vmull_s16(vneg_s16(t1.val[1]), end_tw.val[1]),
            vmull_s16(t2.val[0], end_tw.val[1]))),
    7);

// reverse the end part
end.val[0] = (int32x4_t)vrev64q_s32(end.val[0]);
end.val[0] = vcombine_s32(vget_high_s32(end.val[0]), vget_low_s32(end.val[0]));
end.val[1] = (int32x4_t)vrev64q_s32(end.val[1]);
end.val[1] = vcombine_s32(vget_high_s32(end.val[1]), vget_low_s32(end.val[1]));

// store the result
vst2q_s32((int32_t *)fft_in+i, start);
vst2q_s32((int32_t *)(fft_in+MDCT_M2-4-i/2), end);
}

```

```

// perform the FFT
ne10_fft_c2c_1d_int32_neon(fft_out, fft_in, cfg, 0, 0);

// post-twiddling
for (int i = 0; i < MDCT_M2; i += 8)
{
    // load the fft output and reverse the end part
    // fft_out_x.val[0] -> real part
    // fft_out_x.val[1] -> imaginary part
    int32x4x2_t fft_out_start = vld2q_s32((int32_t *)fft_out+i);
    int32x4x2_t fft_out_end = vld2q_s32((int32_t *)fft_out+MDCT_M-8-i);
    fft_out_end.val[0] = (int32x4_t)vrev64q_s32(fft_out_end.val[0]);
    fft_out_end.val[0] = vcombine_s32(vget_high_s32(fft_out_end.val[0]),
                                     vget_low_s32(fft_out_end.val[0]));
    fft_out_end.val[1] = (int32x4_t)vrev64q_s32(fft_out_end.val[1]);
    fft_out_end.val[1] = vcombine_s32(vget_high_s32(fft_out_end.val[1]),
                                     vget_low_s32(fft_out_end.val[1]));

    // load the twiddle factors
    // x_tw.val[0] -> cos twiddle
    // x_tw.val[1] -> sin twiddle
    int32x4x2_t start_tw = vld2q_s32(posttwiddle_start+i);
    int32x4x2_t end_tw = vld2q_s32(posttwiddle_end+i);

    int32x4x2_t spectrum_start;
    spectrum_start.val[0] = vshrq_n_s32(vaddq_s32(
        vqrdmulhq_s32(vnegq_s32(fft_out_start.val[0]), start_tw.val[0]),
        vqrdmulhq_s32(vnegq_s32(fft_out_start.val[1]), start_tw.val[1])), 8);
    spectrum_start.val[1] = vshrq_n_s32(vaddq_s32(
        vqrdmulhq_s32(vnegq_s32(fft_out_end.val[0]), end_tw.val[1]),
        vqrdmulhq_s32(fft_out_end.val[1], end_tw.val[0])), 8);

    int32x4x2_t spectrum_end;
    spectrum_end.val[0] = vshrq_n_s32(vaddq_s32(
        vqrdmulhq_s32(vnegq_s32(fft_out_end.val[0]), end_tw.val[0]),
        vqrdmulhq_s32(vnegq_s32(fft_out_end.val[1]), end_tw.val[1])), 8);
    spectrum_end.val[1] = vshrq_n_s32(vaddq_s32(
        vqrdmulhq_s32(vnegq_s32(fft_out_start.val[0]), start_tw.val[1]),
        vqrdmulhq_s32(fft_out_start.val[1], start_tw.val[0])), 8);

    spectrum_end.val[0] = (int32x4_t)vrev64q_s32(spectrum_end.val[0]);
    spectrum_end.val[0] = vcombine_s32(vget_high_s32(spectrum_end.val[0]),
                                     vget_low_s32(spectrum_end.val[0]));
    spectrum_end.val[1] = (int32x4_t)vrev64q_s32(spectrum_end.val[1]);
    spectrum_end.val[1] = vcombine_s32(vget_high_s32(spectrum_end.val[1]),
                                     vget_low_s32(spectrum_end.val[1]));

    // store the result
    vst2q_s32((int32_t *)spectrum+i, spectrum_start);
    vst2q_s32((int32_t *)spectrum+MDCT_M-8-i, spectrum_end);
}
}

```