

卒業論文

MP3 エンコーダのハフマン符号器を用いたエ
ンコード処理の高速化

東海大学情報理工学部ソフトウェア開発工学科

6ADF1118 岡部裕之

2010年1月31日

東海大学 情報理工学部
ソフトウェア開発工学科
学生番号：6ADF1118
氏名：岡部 裕之

題：MP3 エンコーダのハフマン符号器を用いたエンコード処理の高速化

1、本研究の意味

本研究は、1989年に開発された圧縮音声情報規格「MPEG Audio-Layer 3 (以下MP3.)」のエンコード処理の高速化をはかることを目的とした研究である。この研究の最終目標は、リアルタイムで音声情報をエンコードすることである。この研究が完成すれば、高音質である情報をより遠くに、より早く届けることが出来る他、インターネットの音楽配信コンテンツでの活躍など、様々な分野において活躍することが予想される。

2、MP3とは

1970年代初頭、ドイツのエルランゲン-ニュルンベルク大学のディター・ザイター (Dieter Seiter) 教授が、電話回線を通して音楽を圧縮転送するという課題を研究したことがMP3の起源となる。後、MP1 (MPEG Audio-Layer 1) 規格とMP2 (MPEG Audio-Layer 2) 規格を通して、それらにハフマン符号化アルゴリズムを取り入れた規格が1989年に誕生した。まだ、明確な名前や規格として確立していなかった。そして、パーソナルコンピュータが一般に普及し始めた1995年に規格にMP3という名前がついた。衛星放送の音声送信フォーマットにも選ばれ

た。その後、インターネットが普及することにより電子情報をより遠くに届けることが容易になったことで、CD音源をMP3に圧縮し他人に配信する行為が流行る。また、MP3はCD音源の10分の1のファイルサイズで、CD音源と同じ程度の音質を維持できるため、小型のポータブルオーディオ再生機が出来る。現代では、アップルコンピュータの「i-Pod」やソニーの「ウォークマン」を代表する所謂ポータブルオーディオプレイヤーが人気となっている。それらが再生することが出来るオーディオ規格の中心はもちろんMP3である。また、携帯電話にオーディオ再生機能を付けたものまである。このように、MP3は一般大衆に深く浸透している。

なぜMP3は普及したのだろうか。MP3が開発された時代は、インターネットの回線は電話回線であり、非常に転送速度が遅く、インターネットの利用者数も少なかった。Windows95が出てきた頃、Microsoft社の「Internet Explorer」がOSに標準で付属されたのをきっかけに、インターネット利用者数が増え、個人でホームページを持つ人が増えた。そして、そのホームページで自分の作曲した音楽を配布するという行為が可能となった。この頃は、MP3よりもMIDIという規格で配布する形が主流だった。MIDIとは、楽譜とそれを演奏する演奏者を電子化したようなもので、これを再生する特別な再生装置で音を鳴らすと言う方式だった。ファイルサイズは1メガバイトもなく、通信速度が遅い電話回線でもそんなに時間をかけずにダウンロードできた。しかし、インターネットの通信速度が上がっていくにつれて、より大きなファイルサイズのデータをダウンロードすることが可能となり、ISDN回線の普及により、

国際電話料金をかけずに海外のインターネットサイトに接続できるようになったことから、利用方法の幅が増えた。そして、自分の気に入った曲の CD を MP3 に変換し、自分のインターネットサイトで公開し、その曲を広めようという行為が出来るようになった。そして、小型で持ち運べる MP3 再生装置が開発され、従来のテープ、CD、MD を再生する所謂「ウォークマン」を置き換えていった。代表的なものにアップルコンピュータの「i-Pod」がある。この「i-Pod」が、さらに MP3 を広めることになる。

i-Pod の登場により、ポータブルオーディオが変わった。今までは、買って来た（もしくは CD レンタルショップ等で借りて来た）CD から MD やテープに録音し、それを専用のポータブル再生機器にて、音楽をどこでも楽しむと言うスタイルだった。この場合、コンポネントステレオや何かしらの録音機を通じて、録音媒体に録音することになる。MD の場合、CD の PCM から AAC という圧縮規格に変換されて MD に記録される。この時、MP3 と同じで可視聴音域外の音をカットする方式を取っているため、音質はすこし劣化する。MD から MD へ録音することを繰り返すと、ノイズが激しくなり、データは劣化する。テープはデジタルではなくアナログである。こちらも、何回もダビング（テープからテープへの録音、もしくは録画）を繰り返すと、データが劣化してノイズが乗るようになる。この二つの方式では、分計算で録音できる容量が決まっている。MD の場合、1 倍速（標準）モードでは、多くて 80 分、4 倍速でその 4 倍の時間、8 倍速では 1 倍速の 8 倍の時間。ただし、音質は倍速にすればするほど、劣化する。テープも同じである。MD の場合は、圧

縮レートを上げることにより、一曲のデータ容量を少なくしている。テープでは、倍速録音をする時、テープの回転をその分遅くすることにより、録音時間を長くしている。

一方、i-Pod はパソコンに USB 端子で繋ぎ、専用のソフトウェアでパソコンの HDD 上にある音楽ファイルを i-Pod 本体に転送するだけで、どこでも音楽を楽しむことが出来る。扱えるファイルは MP3 だけではなく、WAVE ファイルや AAC も扱うことが出来る。MP3 は CD と同等程度の音質で 10 分の 1 程度のファイルサイズまで圧縮することが出来る。MD とは違い、専用の録音機が必要でなく、パソコンを使って簡単に録音ができ、録音している曲が最後まで終わるのを局の時間と同じ時間待つ必要もない。CD アルバムを MP3 に圧縮するのに長くても 10 分で済む。さらに、i-Pod 本体のメモリ容量も大きかった。今のようには 160 ギガバイトなんて大きな容量ではないものの、MP3 を 100 曲入れるには十分な容量だった。CD を沢山持っている人は、その分の MD やテープを用意する必要がなくなったのである。これ一つで全部収まってしまうのである。

しかし、MP3 ファイルを取り扱うことが出来ると言うことは様々な問題を起こした。特に、一番大きな問題は、著作権の侵害である。前述したとおり、当時はインターネット上に MP3 ファイルを公開するという行為が行われていた。しかし、この行為は音楽のアーティストへの著作権を侵害する行為として問題になっていた。i-Pod が出る前までは、パソコンを使い音楽を聴くスタイルをとる人たちが MP3 を多く利用していた。音楽を配信、配布する人たちも利用していた。更に、「WinMX」「Winny」などの P2P (Peer To Peer) 通信

ソフト（所謂、交換ソフト、共有ソフト）を利用して、不正に人の著作物を交換、共有する行為が蔓延していた。音楽も例外ではなかった。i-Pod が世界中に広まり始めた頃、このソフトを利用して音楽ファイルをゲットしようとする人々が増えた。無料で好きなアーティストの曲が入手できるソフトがある、という噂が広まった結果、i-Pod も売れていった。MP3 はその波に乗り、世の中に浸透していった。デジタルデータであり、パソコンで簡単に扱えると言うことは、コピーしても劣化しないデータである。何度コピーしても、劣化しない。MD やテープから比べればかなり便利になった。

MP3 は曲の著作権だけでなく、エンコーダそのものにも特許権がある。MP3 は圧縮方式にハフマン符号を利用している。このハフマン符号化する部分は、非常に複雑であり、過去の多くの MP3 エンコーダは ISO 仕様書規格通りに作られた「dist10」という MP3 エンコーダのハフマン符号化アルゴリズムを元に作っていた。しかし、このハフマン符号化アルゴリズムの特許を取った会社が現れたためインターネットでフリーソフトとして配布する形を取っていた多くの MP3 エンコーダは特許の使用料金を払うことが出来ない、または裁判所からの警告や会社からの警告によって、撤退することを余儀なくされてしまった。このハフマン符号化アルゴリズムの特許を回避して MP3 にエンコードするためには、特許方式とは違うやりかたでエンコードするほかない。この違うやり方で MP3 にエンコードすることを成功させたのが「LAME (LAME Ain't an MP3 Encoder = “LAME は MP3 エンコーダではない”）」である。実は、MP3 エンコーダには特に決められたエンコード方法

は仕様で定められていないのである。つまり最初からエンコードする方法は、自由だったのだ。MP3 デコーダで読むことができるファイルにしまえば、それはもう MP3 なのである。特に LAME は“力技”ともいえるやり方で、特許を回避しているだけでなく、他のエンコーダよりもいづれか音質が良いという評判がある。さらに、エンコード時に音質を調整できるオプションや、ショートブロックのみ、ロングブロックのみでエンコードできるなど、様々なオプションが用意されている。現在出回っている殆どの MP3 はこの LAME でエンコードされている。また、LAME エンコーダのエンコード速度を早くした「午後のコーダ」というエンコーダもある。「午後のコーダ」は一時期、日本のインターネット上で爆発的な人気を誇っていた。GUI 環境で LAME を使ったエンコードが出来たことが人気になった理由の一つであると考えられる。DLL (Dynamic Link Library) 版もあり、他のソフトウェアで MP3 をエンコードする、又はデコードする場合、この DLL 版を組み込んでいたソフトウェアが多かった。午後のコーダでエンコードした MP3 ファイルをバイナリエディタで開くと、いたるところに「GOGO」という文字が埋め込まれているのが特徴である。この二つのエンコーダは、コンパイル済みのものを配布せず、コンパイルする前のソースコードとコンパイラを配布している。こうすることで、研究目的のオープンソースソフトウェアという建前で、特許を回避している。

MP3 には特徴がある。前述したハフマン符号化による圧縮もその一つと言える。MP3 は人間には聞き取ることの出来ない高い音域を削ることにより、更にファイルサイズを小さ

くしている。しかし、高音域を削ることは、音質に対して悪影響を及ぼす。たとえば、音が妙に狭く聞こえるとか、迫力がなくなるとか、爽快感がなくなるとか。それをある程度防止する方法として、心理聴覚音響モデルという人間の錯覚を利用する方法がある。

3、MP3 の仕様

前述したとおり、MP3 はエンコードの仕方を仕様で決めていない。デコードできれば良いとされている。ISO 仕様書にはエンコード時に、デコードする時に必要な情報の項目が書かれている。ISO 仕様書通りに作られた Dist10 でエンコードした MP3 ファイルは、ISO 仕様書に書かれている情報の項目の順番にビットを出力しているため、仕様書をみれば人でも解読できる。なので、私は仕様書と照らし合わせながら Dist10 でエンコードした MP3 ファイルを手で解読していった。時間はかかったものの、何とか基本の構造は理解することが出来た。

さて、あなたが Dist10 でエンコードした MP3 ファイルを期待に胸を躍らせながら、バイナリエディタを使い開いた時、そこには意味不明な数字とアルファベットがずらりと並んでいるだろう。そしてあなたは驚愕し、恐れ、逃げたくなるだろう。しかし、その文字列にはそれぞれ意味があることを忘れてはいけない。森を見て木を見ず。その木にはいろいろな情報が書かれていて、あなたを MP3 エンコードという目的に導いてくれる。そして、あなたは右手にペンと消しゴムを、左手には「ISO/IEC 11172-3」と書かれた一冊の地図とノートを持ち、果敢にその文字列の森に挑むのだろう(笑)。正直、この作業には相

当な気力と根気が必要である。しかし、もしあなたが仕様書を見ただけで理解が出来、C言語なりのプログラムが書けるのならば、その仕様書にそってプログラムを作り、そのプログラムを使い解読するほうが労力は大分減るだろうとコメントしておく。ちなみに、文字列の森と述べたものをビットストリームと呼ぶ。本当は文字列ではなく、バイナリコードであったりする。

さて、まずは MP3 ファイルの基本構造から説明する。MP3 は「フレーム」という大きなくりの中に「グラニューール」という中くらのくり、そして「チャンネル」という小さなくりによって構成されている(図1)。

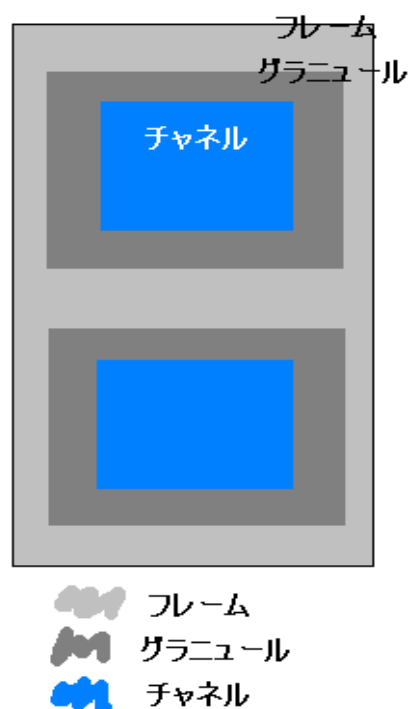


図1、MP3 の1フレームの構成図

チャンネルはモノラルの場合は1つ、それ以外は2つある。グラニューールは必ず2つある。このグラニューールとチャンネルの情報を持った部分がフレームとなる。このフレームをたくさん並べたものが MP3 である。1フレーム含

まれるバイト数をNとすると次の計算でNを
求めることが出来る。

$$N = 144 \times \frac{\text{bitrate}}{\text{sampling_frequency}}$$

今、ここに bitrate が 128kb/s で
sampling_frequency が 44,1kHz の固定レ
ートのMP3データがあるとする。上の式を使
い計算するとN=418 (四捨五入)。つまり、
1 フレーム中に 418 バイトの情報があると言
うことである。この 418 バイトの情報の中に
元の波形にもどすための情報が含まれている。

フレームの一番最初にはフレームの初めを
標している同期語「Syncword」がある。この
同期語を使い、フレームの最初を探す。同期
語は 12bit の”1”で構成されている。つまり、
1 が 12 個並んでいるだけである。バイナリエ
ディタであれば、16 進数なので “FFF” とな
る。この同期語のあとに、デコードに必要な
情報が記述されている。

同期語を含む最初の 4 バイトのデータをヘ
ッドと言う。このヘッドには連続して「同期
語」「ID」「レイヤ情報」「プロテクション(CRC
検査)の有無」「ビットレート情報」「サンプ
リング周波数」「パディングビット」「プラ
イベートビット」「モード情報」「モードエク
ステンション」「コピーの許可」「オリジナル
データかコピーデータか」「エンファシス」と
いう順に記述されている(表1)。詳細は仕様
書を参照してほしい。

表1、ヘッダ情報

Header()		
情報	ビット数	意味
Syncword	12bit	同期語
ID	1bit	
Layer	2bit	レイヤ情報
Protection_bit	1bit	CRCの有無
Bitrate_index	4bit	ビットレート
Sampling_frequency	2bit	サンプリング周波数
Padding_bit	1bit	パディングビットの有無
Private_bit	1bit	
Mode	2bit	モード
Mode_extention	2bit	
copyright	1bit	コピーの許可
original/copy	1bit	オリジナルかコピーか
emphasiss	2bit	

このヘッダの節ではそのフレームの MP3 デ
ータがどのようなサンプリング周波数で、1
秒にどれくらいのビットを使っているか、レ
イヤはいくつか、モード(ステレオとかモノ
ラルとか)はどれでエンコードしているのか
等、音のモデルを記している。このヘッダ情
報を元に、MP3 データを PCM などの音の波
形データにデコードしていく。逆に、エンコ
ードする時はこれらの値を最初に決めておき、
それに沿ってエンコードする。

ヘッダ情報の次には CRC 巡回冗長検査の
ビットが並ぶ。しかし、このビット列はヘッ
ダの「Protection_bit」が0の場合のみ出力さ
れる。出力された場合、そのフレームにおけ
る巡回冗長検査の値を書き込む。検査方法は
CRC-16 を使用する。ビット長は 16 ビット。
検査の対象のビットはヘッダ情報の 4 バイト
分と、付加情報(ヘッダの後ろの 32 バイト分)
である。CRC の検査方法は図2の通りである。
この CRC 検査がない場合、予期せぬエラーで
ビット情報に欠損が出た場合、デコードする
際ひずみが出てしまう可能性が高い。CRC が
あれば、ある程度のエラーは回復できるので、

CRC は付加しておくほうが良い。

CRC 検査ビット列の次には、オーディオデータ (Audio_data0) が続く。この部分で初めてグラニューールとチャンネルに触れる。このオーディオデータ節では MP3 データのデコードに必要な情報が記述されている。オーディオデータの最初には「main_data_begin」という 9bit のデータがある。main_data_begin とは、メインデータの開始位置を示す値である。同期語からこの値のバイト分だけ、前のフレームのメインデータ節にこのフレームのメインデータが含まれることになる。実は、MP3 は 1 フレームに入るメインデータの量が、1 フレームに入れられる分より少なかった場合、その次のメインデータを、今のメインデータにの後ろに入れることが可能なのである。こうすることにより、無駄を省き、ファイルサイズをより小さくしている。それを実現するための情報がこの main_data_begin である (図 3)。

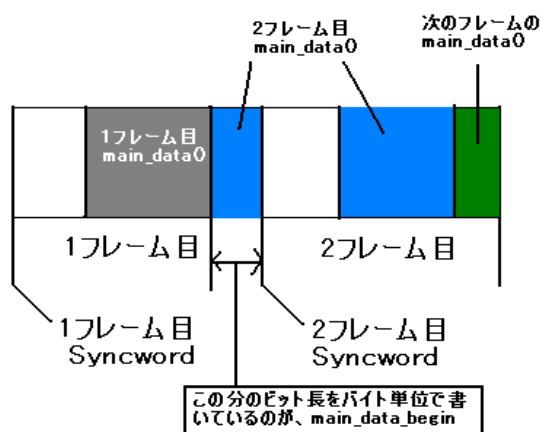


図 3、main_data_begin の説明

Main_data_begin の次は「Private_bit」が続く。この Private_bit は意味を持たない情報である。仕様書によると「将来にわたり、この

規格では規定しない」らしい。しかし、この Private_bit はモノラルモードと、それ以外のモードではビット長が変わってしまうので注意が必要だ。モノラルモードでは 5 ビット、それ以外のモードでは 3 ビットと決められている。

このオーディオデータ節の中には、MP3 の特徴とも言えるスケールファクタ (倍率) を定める節も記される。Private_bit の次に来る「Scfsi[ch][scfsi_band] (ch はチャンネル、scfsi_band はスケールファクタバンド。同倍率選択情報群。)」がそれである。MP3 で言うスケールファクタとは、音の音量を決める係数である。この値が大きければ大きいほど、音量が大きくなる。Scfsi ではスケールファクタの適用方法を、チャンネルごとにスケールファクタバンドで設定している。スケールファクタバンドは、1 チャンネルにつき 4 つある。なので、ステレオの場合は 8 つあることになる。短いブロックを使用する場合、Scfsi は 0 とする。

倍率情報が定まったところで、次からはグラニューールとチャンネルごとに情報が分かれる。モノラルのデータの場合、2 つのグラニューールに、それぞれ 1 チャンネル分のデータが書かれている。それ以外のモードの場合は、2 つのグラニューールにそれぞれ 2 チャンネル分のデータが入っている。一番最初の情報は「Part2_3_length[gr][ch]」である。Part2_3_length ではハフマン符号のビット長を記してある 12 ビットの情報である。1 グラニューール中に含まれるハフマン符号の長さはそれぞれ長さが違う。この値を読み、どこまでがこのグラニューールのどこのチャンネルのメインデータかを確認することが出来る。次に続くのが「Big_value[gr][ch]」である。

ビット長は9ビット。**Big_value** 領域は、ハフマン符号表を用いて0からナイキスト周波数までの全周波数帯域をいくつかの領域に分割し、ハフマン符号化した部分である。この節の**Big_value** は、その**Big_value** 領域のデータの個数を表している。**Big_value** のハフマンコードを符号表と照らし合わせると、**x** と **y** の2つの値を得ることが出来、その後ろに続くビット列からその2値がプラスかマイナスかを知ることが出来る。1 グラニュールにおける**Big_value** の個数は576個と決められている。3つの短いブロックの場合、1つの窓に最大で192個の**Big_value** を入れることが出来る。

その次は「**Global_gain[gr][ch]**」である。量子化ステップ情報を対数量子化して付加する付加情報変数である。この値は、逆量子化する時に使用する。次の「**Scalefac_compress[gr][ch]**」はスケールファクタの同倍率群のビット長を決める情報である。仕様書には、**Scalefac_compress** 表が附属されている（表2）。同倍率群は**Block_type**（後述）によって、範囲が定められる。表2の**Slen1**、**Slen2** は同倍率群の前半と後半のビット長の値である。**Block_type** が0、1、又は3のとき**Slen1** は同倍率群0～10に対する倍率のビット数、**Slen2** は同倍率群11～20に対する倍率のビット数。**Block_type** が2でかつ**Mixed_block_flag**（後述）が0のとき、**Slen1** は同倍率群0～5に対する倍率のビット数。**Slen2** は同倍率群6～11に対する倍率のビット数。**Block_type** が2でかつ**Mixed_block_flag** が1のとき、**Slen1** は同倍率群0～7（長い窓）および3～5（短い窓）に対する倍率のビット数。**Slen2** は同倍率群6～11に対する倍率のビット数

である。

表2、**Scalefac_compress** 表

Scalefac_compress[gr][ch]	slen1	slen2
0	0	0
1	0	1
2	0	2
3	0	3
4	3	0
5	1	1
6	1	2
7	1	3
8	2	1
9	2	2
10	2	3
11	3	1
12	3	2
13	3	3
14	4	2
15	4	3

その次は「**Window_switching_flag[gr][ch]**」が続く。**Window_switching_flag** はそのブロックが普通の窓（**Block_type**==0）以外を使用することを示す1ビットの値。**Window_switching_flag** が設定された時、他のいくつかの変数を特に指定しない限り、次に示す値に設定する。

Region0_count = 7

(**Block_type** == 1, **Block_type** == 3 又は **Block_type** == 2かつ **mixed_block_flag** == 1のとき。)

Region0_count = 8

(**Block_type** == 2かつ **mixed_block_flag** == 1のとき。)

Region1_count = 36

Big_value 領域の残りの全ての値は、領域1（**Region1**）に含む。**Window_switching_flag** が設定されていない（0である）時、

Block_type は 0 とする。Region0_count、Region1_count については、後述する。

Window_switching_flag の次の情報は「Block_type[gr][ch]」である。この情報は、対応するグラニューールのブロックタイプ（窓タイプ）を示す。ブロックタイプには開始ブロック、3 つの短いブロック、停止ブロックの 3 つと、予約されたブロックタイプがあり、主に 3 つの短いブロックを利用してエンコードする（表 3）。開始ブロックは、主に MP3 データの 1 フレーム目（開始フレーム）で使用される。短いブロックは、高い周波数の音をエンコードする際、長いブロックではプリエコーがかかってしまうため、3 つに短く切り、時間解像度をよくして、プリエコーの発生を抑えるためにできた。停止ブロックは、最後のフレームで使われる。

表 3、Block_type

Block type[gr][ch]	ブロックタイプ
0	予約
1	開始ブロック
2	3つの短いブロック
3	停止ブロック

Block_type および mixed_block_flag は、そのブロックの数値（情報）の並び方、変換の回数についての情報を与える。Window_switching_flag == 1 のとき、mixed_block_flag は低域側の多重位相フィルタ分割帯域を普通のブロックタイプを使って符号化するかどうかを示す。長いブロックの場合（Block_type が 2 でない場合、または Block_type == 2 であり、Mixed_block_flag == 1 であり、かつ低域側分割帯域の場合）IMDCT は 18 個の入力ごとに 36 個の出力を生成する。その出力に Block_type に応じた窓をかけ、その前半を前ブロックの後半とかさ

ねあわせ処理をする。得られたベクトルを多重位相フィルタバンク合成部の一つの帯域の入力とする。短いブロックの場合（Mixed_block_flag == 1 であり、Block_type == 2 であり、かつ高域側分割帯域の場合、又は mixed_block_flag == 0 でありかつ Block_type == 2 の場合）、それぞれ 12 個の出力値を生成する変換を 3 回実行する。この 3 ベクトルに窓を掛け、互いに重ねあわせ処理をする。得られたベクトルの両端にそれぞれ 6 個の“0”を連結し、長さ 36 のベクトルを作り出し、それを長いブロックに用いた変換の出力と同様に処理をする（図 4）。Block_type の次は、「Mixed_block_flag[gr][ch]」の情報が入っている。Mixed_block_flag は高域側で使用する窓と異なる窓で低域側を変換することを示す。

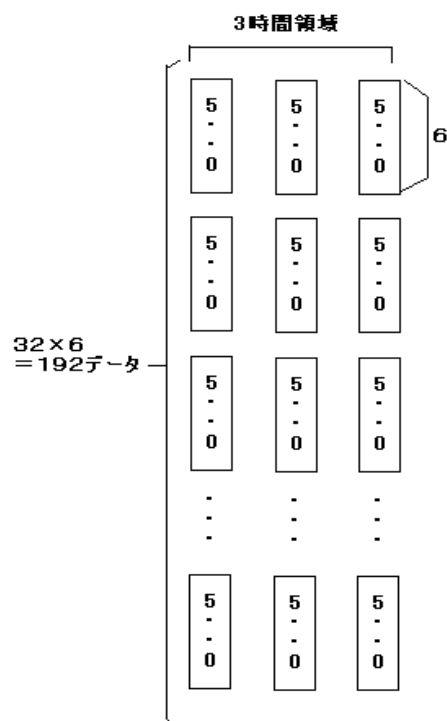


図 4、3 つの短いブロックの 1 ブロック

Mixed_block_flag == 0 のとき全てのブロックに Block_type[gr][ch]が示す変換を行う。Mixed_block_flag == 1 のとき、最低域側の2個の多重位相分割帯域に対応する周波数成分を普通の窓(Block_type == 0) で変換し、残りの30個の分割帯域を Block_type[gr][ch]に応じて変換する。Mixed_block_type が設定された場合、そのグラニューールには短いブロックと長いブロックが混在していることになる。Mixed_block_type の次は「Table_select[gr][ch][region]」が入っている。この情報は、最大量子化値および信号の局所的な統計的性質に応じて異なるハフマン符号表を使用する。合計32種の符号表があり、この中から適切な符号表を使い符号化する。この Table_select で選択された符号表は、Big_value 領域のみで使用される。MP3におけるハフマン符号化は、通常ハフマン符号化とは違い、もともとコード表が決まっている。この表の中から音質とファイルサイズを見て適切な表を選択しなければならない。この部分で、複雑な計算をする。この部分にかかる時間を短縮することにより、エンコード速度の向上につながると考えている。このコード表には Big_value 領域で得られる x と y の2つの値が記されている。

次に入る情報は、「subblock_gain[gr][ch][window]」である。あるサブブロックについて global_gain[gr][ch]からの偏位量(4のべき乗による量子化)を示す。Block_type == 2 (短い窓)の場合だけ使用する。復号器においてサブブロックの数値を以下の式で除する

$$4^{(\text{subblock_gain}[\text{gr}][\text{ch}][\text{window}])}$$

Subblock_gain の次には、

「region0_count[gr][ch]」が入る。ハフマン符号器の能力を高めるためにスペクトルを更に分割する。これは、Big_value で表される領域を更に再分割することに相当する。この再分割は誤り耐性改善および符号化率向上を目的として行う。領域0、1又は2の3領域を用いる。各領域を最大量子化値および信号の局所的性質に応じて、異なるハフマン符号表で符号化する。region0_count には、count1 領域(後述)の前半の領域の終点の同倍率群番号が書かれている。Count1 領域は、仕様書付属の同倍率群表と、ハフマン符号表 A、およびハフマン符号表 B(後述)を使い変換するため、この値が必要となる。領域の境界はスペクトルの同倍率群の分割点と一致させる必要がある。値、region0_count は領域0 (region0) の同倍率群数より1少ない値とする。短いブロックの場合、各同倍率群表は各短い窓に対し、1回ずつ、合計3回数える。したがって、region0_count の値が8の時、領域1を同倍率群表3からはじめることを示す。Block_type == 2かつMixed_block_flag == 1 のとき、そのグラニューールの同倍率群総数は $12 \times 3 = 36$ とする。Block_type == 2かつMixed_block_flag == 1 のとき、同倍率群の総数は、 $8 + 9 \times 3 = 35$ とする。Block_type が2でないとき、同倍率群の総数は21とする。Region0_count の次には、「Region1_count[gr][ch]」が入っている。Region0と同じく、Count1領域のスペクトルの分割点を示している。領域内の同倍率群数より少ない数とする。Block_type == 2 のとき、異なる時間スロットを表す同倍率群は別々に数える。

その次は「preflag[gr][ch]」である。Preflag は量子化値の高域を更に増幅する簡便な方法

である。Preflag を設定すると、仕様書の付属の表に示された値を倍率に加算する。これは、逆量子化した倍率に表の値を乗算することと投下である。Block_type == 2 (短いブロック) の時、Preflag は使用しない。

次は、「Scalefac_scale[gr][ch]」である。Scalefac_scale に応じて、倍率は、2 又は√2 のステップ幅で対数量子化する。次の表に各ステップ幅に対する逆量子化の式で使用する、倍率への乗数を示す (表 4)。

表 4、scalefac_scale

scalefac_scale[gr][ch]	scalefac_multiplier
0	0.5
1	1

次は、Count1table_select[gr][ch]である。大きさ 1 以下の量子化値からなる四つ組みの領域に対して、使用可能な 2 種のハフマン符号表のうちの一つを選択する。この 2 種ハフマン符号表が付属書の符号表 A、符号表 B である。Count1table_select で、Count1 領域が符号表 A か符号表 B を使っているかが分かる (表 5)。

表 5、Count1table_select

count1table_select[gr][ch]	選択するハフマン符号表
0	付属書B(規定)表 B.7(A)
1	付属書B(規定)表 B.7(B)

これらの情報が、ステレオの場合、1 グラニューールに 2 チャンネル分。グラニューールは 1 フレームに 2 つあるので、合計 4 つ分出力されている。以上、ここまでの情報が Audio_data である。

Audio_data ではデコードに必要な情報が記載されていた。次の情報は、倍率に関する情報の「main_data0」である。仕様書によれば「量子化雑音を周波数領域で重み付けするために、倍率を用いる。量子化雑音の重み付けが正しければ、人間はその量子化雑音を知覚できない。」と書いてある。つまり、MP3 の一つの特徴である、人間が聞くことの出来ない高音域を削るということをこの倍率を用いて実現しているということである。レイヤ 1 とレイヤ 2 とは異なり、レイヤ 3 の倍率は量子化した信号の局所的最大値について、何も規定しない。レイヤ 3 では倍率は複合器で数値の集まりに対する割り算の除数を得るためだけに使用する。レイヤ 3 の場合、その集まりはいくつかの周波数成分にまたがっている。これらの集まりを同倍率群と呼び、可能な限り臨界帯域に一致するように選択する。

Scalefac_compress 表は、倍率 0 ~ 10 が 0 ~ 15 (最大長 4 ビット) の範囲にあること、および、倍率 11 ~ 21 は 0 ~ 7 (最大 3 ビット) の範囲にあることを示す。スケールファクタのビット長は 0 ~ 4 ビットで、Scalefac_compress 表から得られる Slen1 と Slen2 がそのままビット長になる。スケールファクタの符号化に使用するビット数は Part2_length と呼び、次のように計算する。

• block_type == 0, 1 又は 3 (長いブロックに対して)

$$\text{Part2_length} = 11 \times \text{slen1} + 10 \times \text{slen2}$$

• block_type == 2 (短いブロック) かつ mixed_block_flag == 0 に対して

$$\text{Part2_length} = 18 \times \text{slen1} + 18 \times \text{slen2}$$

• block_type == 2 (短いブロック) かつ mixed_block_flag == 1 に対して

$$\text{Part2_length} = 17 \times \text{slen1} + 18 \times \text{slen2}$$

Audio_data の次は、「Huffmancodebit()」つまり、ハフマン符号である。このハフマン符号の中に、音の元となるデータが詰まっている。Huffmancodebit の中には、Big_value 領域と、Count1 領域と、ZERO 領域という3つの領域がある(図5)。Big_value 領域では32個のハフマン符号表のうち1つの表を選んでエンコードされた値が詰まっている。これらの個数は2の倍数である。この領域ではxとyの2値が得られる。また、2値のほかにlinbitsと言う値も得られる。linbitsは各は不満符号表に書かれているので、その値の分のビットを読む。Count1 領域では符号表Aか符号表Bでエンコードされた値が入っている。この領域では、v、w、x、yの4値が得られる。いずれも、-1、0、1のどれか。個数は4の倍数でなければならない。ZERO 領域は0の値が入っている領域である。多くの場合、前のフレームのmain_data0が入っているため、ビットストリームには出力されていない。

Big_value 領域は次のパラメタで構成されている。

- Hcod [|x|] [|y|]
ハフマン符号表への引数。この引数を使って、xとyの2値を得ることが出来る。
- hlen [|x|] [|y|]
値x、yのハフマン符号長符号表への引数。要は、ハフマン符号のビット長である。
- signx
もし、xの値がある場合、xの値が正の値か、負の値かを示す。0の場合、正。1の場合負。
- signy

もし、yの値がある場合、yの値が正の値か、負の値かを示す。他はsignxと同じ。

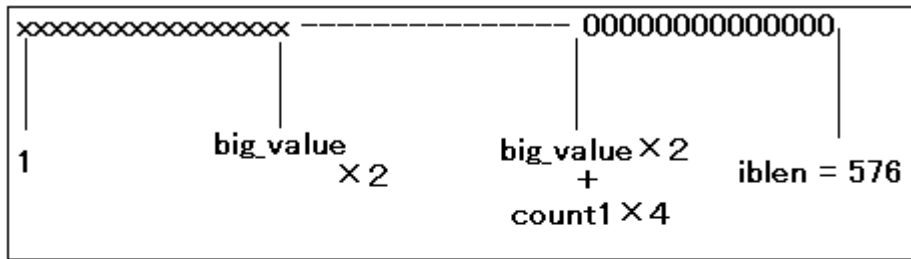
- linbits
linbitsx 又は linbitsy がある時、そのビットの長さ。

Count1 領域は次のパラメタで構成されている。

- Hcod [|v|] [|w|] [|x|] [|y|]
ハフマン符号表への引数。Big_value 領域とは違い、v、w、x、yの4値を得ることが出来る。
- signv
もし、vの値がある場合、vの値が正の値か、負の値かを示す。0の場合、正。1の場合、負。
- signw
もし、wの値がある場合、wの値が正の値か、負の値かを示す。他はsignvと同じ。

あとはBig_valueと同じく、signx、signyがある。意味も同じである。

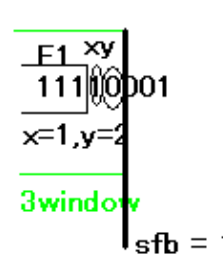
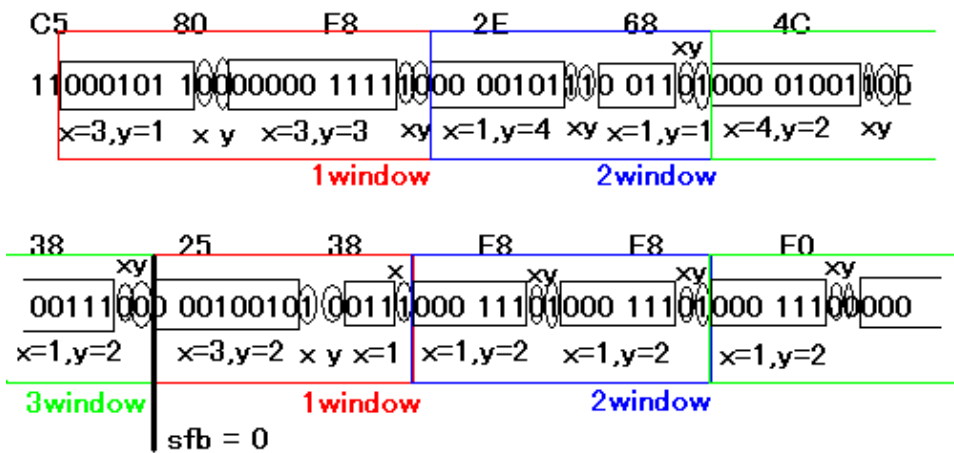
これらから得られた数値を逆量子化(IMDCT変換)することにより、音の波形を得ることが出来る。デコードするときBlock_type == 2(短いブロック)の場合、Big_value 領域は同倍率群表を使い、デコードする(図6)。以上の内容が、MP3の1フレームのビットストリーム構造になっている。



値 000 は全て0。0の個数は2の倍数
 値 --- は-1, 0, 1のいずれか。これらの個数は4の倍数。
 値 xxx は制限なし。これらの個数は2の倍数。
 iblen は 576。

1フレーム中に576個のデータがある。

図5、Huffmancodebit の中身



sfb = 0 → start = 0, end = 4 → 4 - 0 = 4 なので2回
 sfb = 1 → start = 4, end = 8 → 8 - 4 = 4 なので2回

 Layer = 3, サンプリング周波数 = 44.1kHz, ビットレート128kb/s
 block_type == 2 (短いブロック)
 scalefac_compress = 1010 → slen1 = 2, slen2 = 3
 window_switching_flag = 1
 table_select[0][0] = 7
 table_select[0][1] = 2

図6、short block におけるデコード

4、MP3のデコード

ビットストリーム構造が分かったので早速本物のMP3をデコードしてみよう。まずエンコードするWAVEファイルを用意する。エンコーダはDist10を使用する。ここで忘れてはいけないことがある。Dist10はWindowsで標準とされているPCMやRIFF形式のWAVEファイルだと、エンコードできないので、注意が必要だ。Dist10がエンコードできるファイル形式はAIFFとraw形式のファイルでなければならない。私は「TiMidity++」というソフトウェアシンセサイザがAIFFファイルを出力できるので、こちらにて再生時間が短めなMIDIファイルをAIFFファイルに出力した。ソースとなるファイルが用意できたので、Dist10を使い、MP3にエンコードする。サンプリング周波数は44.1kHz、ビットレートは128kb/s、モードはステレオ。プロテクションは使用しないでエンコードする。ここで注意。Dist10には聴覚心理モデルが2つある(正しくはMP3エンコーダ。ISOの仕様書にも、同様に2つの聴覚真理モデルの説明が記載されている)。聴覚心理モデル1は、Layer3には使えない。Dist10は、デフォルトで使用する聴覚心理モデルは「聴覚心理モデル1」なので、エンコード時にオプションで「聴覚心理モデル2」に指定しなければならない。このように細かくオプションを決めてあげないと、きちんとエンコードできない。出力されるMP3ファイルの拡張子は最初「~.mpg」になっている。エンコードし終わったら、一度再生して音を聞いてみよう。成功していれば、ソースのファイルとほぼ同じ音が聞こえるはずだ。しかし、失敗していると、変なノイズがのっていたり、「ザー」というノイズしか乗っていない場合がある。再

生するソフトウェアはWindows環境では「SCMPX」がお勧めである。SCMPXは、古くからフリーソフトとして配布されているMP3を再生できるソフトウェアである。エンコードも可能である。LAMEエンコーダとは違い、ある程度ISO仕様書にそったMP3ファイルを生成するようである。自前のエンコーダとデコーダを持っており、CDから直接MP3にエンコードすることが、CD2WAVEというソフトウェアの外部ツールとして使用することで、可能だった。何故、私がSCMPXを奨めるかというと、実はDist10でエンコードしたMP3ファイルをWindows Media Playerで再生しようとしたところ、どれも再生できずにフリーズしてしまった。(私のミスかもしれない。)しかし、RIFF形式のWAVEファイルをエンコードしたMP3ファイルは再生できた。けれども雑音しか入っていなかった。他の知っているMP3を再生できるプレイヤーソフトも使ってみたが、やはり雑音だけだったり、再生出来ずに終了した。MP3ファイルが壊れているのかもしれないと思い、他の音源もエンコードしたものの、結果は同じだった。しかし、SCMPXならばとりあえずは再生可能だったので、これはきっとMP3ファイルなのだろう、と仮定して作業を進めた(教授曰く、「Windows Media PlayerがISOの仕様を満たしていない可能性も否定できない。」)そういえば昔、Windows98に付属していたMedia Playerか、レコーダだったかに、MP3エンコード機能がついていたのを見つけて、期待に胸を膨らませながら、僕が好きだったある曲をエンコードしたら、プツプツというノイズと、ブロックノイズだらけで、がっくりした記憶がある……。ちなみに、サンプリング周波数は44.1kHz、ビットレ

ト 128kbps、ステレオ。記憶があいまいだけれども、とにかくノイズだらけだったのは記憶に残っている)。エンコードが終わり、音も確認し終わったら、次はバイナリエディタで MP3 ファイルを開く。Windows 環境ならば「Stirling」、Linux 環境ならば「bvi」がお勧めである。私は現在 Windows 環境であるので Stirling を使用する。さて、バイナリエディタで先ほど作成した MP3 ファイルを開くと図 7 のようになる。これだけを見ても何がなんだか分からないはず。これを見ただけで分かたれば、あなたには才能があるのかもしれない (Dist10 以外のエンコーダでエンコードすると、ID3 タグという情報タグが一番上、もしくは一番下についている可能性が高い。すると、この部分には生の文字が入っているため、バイナリエディタの右の文字列表示に、目で読める部分があったりする。Dist10 には ID3 タグがない。なので、仕様書と照らし合わせて解読するには、非常に良い)。まず、バイナリエディタの見かたを説明していく。一番上の黄色い部分は、列数を意味している。バイナリエディタは 16 進数で表示されるので、最後の「0F」は 16 列目ということになる。左の灰色の部分は、白い部分に書かれた数値のアドレス番号を表している。なので、2 行目の 1 番最初のアドレスは 17 なので、F (16) から 1 繰り上がって「10」となる。さて、このバイナリエディタを使って、MP3 をデコードしていきたいと思う。MP3 の一番最初は「同期語 (syncword)」の 1 が 12 ビット並んでいるのであるから、16 進数では「FFF」となるはずである。バイナリエディタで開いた MP3 の一番最初は「FFF」である。よってここがヘッダの始点となる。ヘッダは始点から 4 バイト分であるので、0 から 4 番

までのアドレスを読む。バイナリコードでは「FF FB 92 00」。これを 2 進数 1 ビット単位に直すと「1111 1111 1111 1011 1001 0010 0000 0000」となる。これをデコードする。最初に述べた 12 ビットは同期語である。

ADDRESS	00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F	0123456789ABCDEF
00000000	FF	FB	92	00	00	00	02	68	24	53	05	3D	00	02	68	06	羽....h\$\$.=...h.
00000010	2A	30	AC	30	00	0C	2A	41	62	18	D3	80	09	7F	48	2C	*0y0...*Ab.モ...H,
00000020	43	1A	70	01	1E	81	08	07	20	19	00	C8	12	01	70	3A	C.p.....ネ..p:
00000030	26	60	39	0D	03	A1	40	C9	A8	6F	C1	A0	28	01	40	14	&`9...@ノioチ.(.@.
00000040	06	82	88	88	2E	7B	B8	B8	B8	B8	B9	E6	20	A0	A2	22	.h...{クククク.「”
00000050	22	3B	BF	BB	BB	A2	7C	22	22	3B	B9	62	E2	E2	E2	EF	”;ツサ ””;外戸靴
00000060	C1	10	3C	30	00	03	EC	1E	1E	1E	1E	19	03	FF	FC	44	チ.<0.....鮭
00000070	0F	45	42	E0	00	00	62	11	88	05	98	40	02	E8	A2	7D	.EB...b...蓮.睡]
00000080	18	62	EC	62	0D	72	29	8C	AE	36	24	03	40	20	0D	04	.b.r)鍵6\$.@ ..
00000090	43	05	94	58	BD	7A	F5	EB	CC	CC	CC	CF	D6	2C	58	B1	C.囊又z・フマヨ,X7
000000A0	63	94	A5	17	DE	F7	BD	EF	7B	D2	94	A5	29	4A	52	F7	
000000B0	5E	BC	CC	CC	CC	FD	FA	52	94	A5	27	6E	B1	CA	52	94	
000000C0	A5	29	7B	DE	F7	6D	7B	E0	E0	20	AC	10	04	01	01	07)['・[珥 ヲ.....
000000D0	FC	F8	3E	0F	89	FF	FF	38	20	16	43	1C	2F	42	6C	25	・>....8 .C./BI%
000000E0	20	9B	25	C8	C4	C1	BC	15	F1	80	22	92	4E	6C	48	93	
000000F0	C5	84	48	8D	46	A3	C2	82	42	C1	D2	84	82	88	54	B9	3孝]ツ・ヲp・ケ
00000100	87	A1	11	B0	F2	17	35	CE	2E	40	C3	E4	99	4E	52	EC	
00000110	6B	2B	A8	E9	C8	9B	64	5E	75	0C	67	DC	F4	9D	D0	FB	
00000120	6E	95	F7	FB	FF	FF	FF	FF	FF	FF	FF	FF	FF	FF	FF	FF	庖.....
00000130	FF	FF	FF	FF	FF	FF	FF	FF	FF	FF	C5	90	EF	11	A1	36并撰..6
00000140	12	90	4D	92	E4	62	60	DE	0B	58	C0	11	49	27	36	24	.信弟b`°.Xタ.I'6\$
00000150	49	E2	C2	26	8D	47	47	85	04	85	83	A5	09	05	10	A9	[籠&宏G...ウ
00000160	73	0F	42	23	61	E4	2E	6B	9C	5C	81	87	C9	32	9C	A5	s.B#a..k彌∞ノ2倍
00000170	D9	99	5D	47	4E	44	DB	22	F3	A8	63	3E	EA	97	E8	7D	焼克GNDa”・c>語聞
00000180	B7	D7	DF	EF	FF	FF	FF	FF	FF	FF	FF	FF	FF	FF	FF	FF	キ°.....
00000190	FF	FF	FF	FF	FF	FF	FF	FF	FF	FF	00	00	14	56	B6	9B	
000001A0	75	4B	FF	FB	92	00	04	00	03	11	90	5C	56	1C	40	02	K.羽.....申V.@.
000001B0	60	32	0B	8A	C3	88	00	4B	DA	41	6D	58	A1	00	09	80	`2.甘..KvAmX....
000001C0	48	2D	6B	14	20	01	AE	FF	FF	C5	80	F8	2E	68	4C	38	H-k. .ョ..ナ...hL8
000001D0	25	38	E8	E0	F1	A3	55	26	A6	D4	EE	18	CC	C4	49	18	%8鞆・U&ヲ..ガI.
000001E0	89	22	15	E5	47	32	1E	2E	6C	EC	AF	B1	8C	D9	68	CF	.”.薑2..I・ア雇h7
000001F0	7A	33	1D	9D	16	F5	F6	DE	CE	49	DF	47	73	31	ED	BB	z3...・°ホI°Gs1峻
00000200	B3	ED	3B	BB	B2	5A	7A	EA	5D	A4	AF	FE	BD	2B	F6	F7	ウ.;サイz鶴.ツ.ス+・
00000210	FD	3D	7F	FD	FF	FF	FF	FF	FF	FF	FF	FF	FF	FF	FF	FF	.=.....

図 7、バイナリエディタで開いた MP3 データ

ここからは、ISO/IEC11172-3つまりMP3のISO仕様書を読みながら説明する。もし手元にISO仕様書がある場合は、見て確認しながら読むとより理解できるはずだ。さて、解読作業に入る。同期語の次のビットは「ID」を示している。このIDが1であれば、MP3である。2進数コードは「1」なので、MP3である。次は「レイヤ」だ。レイヤは2ビットなので2ビット分読むと「01」だ。レイヤに01が入っている場合、レイヤは3になるので、このデータはレイヤ3でエンコードされていることが分かる。次はProtection_bitの1ビットだ。値は「1」が入っているので、このフレームのデータにはCRC巡回冗長検査が出力されていない。次はbitrate_indexの4ビットだ。「1001」が入っている。ISO仕様書の表を見ると、レイヤ3の1001は、128kbit/sなので、このデータは一秒に128kビットのビットレートになっている。次はsampling_frequencyの2ビット。値は「00」が入っている。仕様書の表ではsampling_frequencyが「00」のときは標準化周波数（サンプリング周波数）は44.1kHzになっている。次は、モードだ。2ビット。値は「00」仕様書によれば、00はステレオである。次はmode_extentionの2ビット。こちらは、「00」である。残りの4ビットはcopyrightの1ビット、original/copyの1ビット、emphasisの2ビットである。全て0なので著作権なし、複製されたデータ、エンファシスはなし。となっている。最初に設定したとおりのMP3データが出力されているのが分かる。Protection_bitは1なので、CRCデータはない。次は、Audio_data()だ。Audio_dataはモノラルモード以外では、ヘッダの4倍と以降の32バイト分のデータがこれ

に相当するのでそれを読む。全て読むと分かりづらくなるので、0グラニューール0チャンネルのみ読む。そこからあとは、図8、図9、図10、図11に記す。Audio_data()の最初はmain_data_beginの9ビットである。このフレームは、最初のフレームなので、main_data_beginは0になっている。Private_bitはステレオなので3ビットあればいい。次のScfsiはステレオモードなので2チャンネル分取らなければならない。Scfsi一つ分で1ビット、1チャンネルに4つのscfsi_bandがあるので合計8ビット。Scfsiの次は、part2_3_lengthにはいる。こちらは12ビットである。値は「001001101000」だ。これを10進数に直すと616になる。よってハフマン符号のビット数は616ビットである。次はBig_valueの9ビット。値は「001001000」。10進数に直すと72。よって、72個のBig_valueがhuffmancodebitのbig_value領域に入っていることになる。次は、Global_gainの8ビットだ。そして、scalefac_compressの4ビットが続く。Scalefac_compressは「0000」の値が入っているのでslen1とslen2は両方とも0。このフレームは音が出ないと予測される。次はwindow_switching_flagの1ビット。値は1。Window_switching_flagが設定されているので、次はblock_typeの2ビットになる。こちらは、「01」が入っているので開始ブロックであることが分かる。次は、mixed_block_flagの1ビット。「0」が入っている。その次のtable_selectは、window_switching_flagが1なので領域は2つ。領域1のtable_selectは01111、領域2のtable_selectは01000。なので、領域1ではハフマン符号表15、領域2ではハフマン符号表8を使う。

Subblock_gain は 3 つの窓分用意している。1 つの窓で 4 ビット使うので、12 ビット分用意されている。このフレームではどれも 0 が入っている。次は Preflag、scalefac_scale、count1table_select の 3 ビット。Preflag は 0、scalefac_scale も 0。Count1table_select は 1 が入っている。Count1table_select が 1 なので、count1 領域では、符号表 B を使う。他の情報は図 8 から図 11 を参照。これらの情報をもとに、ハフマン符号をデコードする。0 グラニュール、0 チャネル目のハフマン符号は、開始ブロックなので、長いブロックである。よって Huffmancodebit() は big_value と count1 領域、0 領域で構成されている。今、Big_value 領域はさらに二つの領域に分けることが出来る。前半の領域は Region0 で、後半の領域は Region1 だ。Table_select で選択された符号表は、符号表 15 と符号表 8 なので、Region0 は符号表 15、Region1 は符号表 8 でデコードする。ここまでのデータが、シングルチャネル (モノラル) の場合、1 グラニュール、1 チャネルの 2 つ、それ以外のモード (ステレオ等) の場合、2 グラニュール 2 チャネルの 4 つあれば、そのあとに続くデータは、MP3 のメインデータ (量子化された値) なので、とりあえず ISO で定められた仕様を満たしたデータになる。

さて、Audio 情報を読みどいう音の情報が、どれくらいの音質で入っているかがわかったので、今度は音そのものを符号化圧縮した Main_data () と Huffmancodebit () を解読していく。私、実はこの項について正しい答えが出せていない。理由は、エンコーダもデコーダも作れず、これであっているのかどうか、確かめられなかったからである。唯一、私が確かめられるエンコーダとデコーダは、

Dist10 だけだった。なので、この項は Dist10 の仕様書と Dist10 のソースコードを頼りに、私なりに解釈したデコード方法である。よって確証性が無い。

さて、本題に入る。まず Audio_data 節と Huffmancodebit 節は図 12 のようなになっていることを確認してほしい。こちらは、仕様書にも同じような図があるので、正しいはずである。解釈が間違っていなければのはなしではあるが。モノラルの場合、big_value 領域には 1 グラニュールにつき 1 チャネル分のデータがあり、それ以外 (ステレオなどの 2 チャネル必要とする) の場合、1 グラニュールにつき 2 チャネル分ある。また、倍率に対応する 2 つの領域に区切られており、前半の領域を region0 領域、後半の領域を region1 領域と呼ぶことにする。Count1 領域は、「-1,0,1」のいずれかの値が入っている領域で、「v,w,x,y」の 4 値得ることが出来る。最後の赤く囲まれた部分は、次のフレームのデータで、次のフレームの main_data_begin に入っている数のバイト数分のビットが入っている (つまり main_data_begin × 8 ビット入っている)。1 つのグラニュール分に big_value 領域と count1 領域、そして出力されない 0 領域を含めて 576 データある。Big_value 領域と count1 領域を合わせて 576 データに満たない場合、0 領域がある。0 領域を出力しないことで、ファイルサイズを小さくしている。

まず、main_data 節の解読に入りたい。この節では、倍率の細かい指定をしている。倍率は scalefac_l (長いブロック用) と scalefac_s (短いブロック用) の二種類があり各ブロックタイプに応じて使い分ける。倍率のビット長は scalefac_compress 表を使い、図 14 の式で求めることが出来る。倍率のビ

ット長を `part2_length` と呼びこの値は、ブロックタイプと `scalefac_compress` に応じて変化する。`part2_3_length` とは違うものなので注意。この `part2_length` の分だけ倍率が1グラニューールの1チャンネル、あるということになる。さて、図14の式を使いこのフレームの倍率を求める。0グラニューール、0チャンネルの `scalefac_compress` には0が入っている。`slen1` は0、`slen2` は0。よって、`part2_length` は0。0グラニューール0チャンネルの倍率は出力されていないことになる。同様に、0グラニューール1チャンネル、1グラニューール0チャンネル、1グラニューール1チャンネルも `scalefac_compress` が0なので、このフレームのビットストリーム上には倍率に関するビットは出力されていない。

さて、`big_value` 領域の解読に入る。`Big_value` 領域を解読するには、仕様書付属のハフマンコード表を利用して解読していくのは前述したとおりである。`Big_value` 領域に適用できる符号表は0-31の符号表である。では、0グラニューールの `table_select` はどうなっているかと言うと、「`table_select[0][0][0] 01111`」「`table_select[0][0][1] 01000`」なので `region0` 領域には符号表15を、`region1` 領域には符号表8を適用してデコードする。0グラニューール0チャンネル目の `part2_3_length` は616、`big_value` は72が入っているので、0グラニューール0チャンネル目の `huffmancodebit` 節は616ビットあり、そのうち `big_value` 領域は72組ある。その後ろに `count1` 領域、0領域と続く。`big_value` 領域は図13のような構造になっている。符号表に書かれたハフマン符号と、MP3のビットストリーム上に書かれているコードとを照らし合わせて、デコードしていく。ハフマン符号表とビットストリ

ームから得られたハフマン符号から、`x,y` の2値を取り出すことができ、ハフマン符号の後に続く2つの値が `x,y` が正の値か負の値かを示している。0グラニューール0チャンネル目はブロックタイプが開始ブロックなので長いブロックである。そして `region0_count` が7、`region1_count` は36なので、仕様書付属の同倍率群表を使い、`big_value` 領域を `region0` 領域と `region1` 領域で分割する。72組データが取れた時点で、`big_value` 領域は終了する。そして、`count1` 領域に入る。`Count1` 領域では符号表Aと符号表Bのどちらかを選択し符合する。どちらの符号表を使えば良いかは、`count1table_select` に書いてあるので、それをみて判断する。0グラニューール0チャンネル目の `count1table_select` は1が入っているので、符号表Bを使いデコードする。616ビット目で、`count1` 領域が終わり、0グラニューール0チャンネルのハフマン符号ビットのデコードが完了する。このやり方を残りの3つにも適用すると、1フレームのハフマン符号のデコードが完了する。

上記から得た値を逆量子化器に入れ、得られた値を合成フィルタバンクに入力し、並べ替える。逆量子化に使用する公式は図15に示す。このとき `isi` はバッファ指標 `i` におけるハフマン複合値とし、合成フィルタバンクへの入力値を `xri` とする。図15の `pretab` はプリアンファシスのことである。210と言う値はシステム定数で、出力値を適切に倍率変換するために用いる。短いブロックを使用した場合、逆倍率変換したデータ `xr[scf_band][window][freq_line]` はIMDCT演算の前に分割帯域順 `xr[subband][window][freq_line]` に並べ替えないといけない。

逆量子化が終わると、次はステレオ処理が入る。ステレオ処理では、和差ステレオモード（ミッドサイドステレオモード）、共包絡ステレオモード（ジョイントステレオモード）の2種類を決定する。和差ステレオモードにおいて、正規化された和チャンネルの M_i と S_i の値を左右チャンネル値 L_i と R_i の代わりに伝送する。 L_i と R_i は次の式で復元する。

$$L_i = \frac{M_i + S_i}{\sqrt{2}}$$

$$S_i = \frac{M_i - S_i}{\sqrt{2}}$$

共包絡ステレオは音声圧縮規格が、普通のステレオと、和差ステレオを、ブロックごとに分けて音声圧縮をかける時の方法のことで、この方法でエンコードすると、和差ステレオだけでエンコードするよりも、効率がよく圧縮できる。MP3では共包絡ステレオで符号化された各同倍率群は sb （スケールファクタ・バンド）にたいして次の（1）～（5）を実行する。

- （1） 共包絡ステレオ位置 $is_pos_{sb}[sfb]$ を右チャンネルの倍率から読む。
- （2） もしくは、 $is_pos_{sb} == 7$ なら（3）～（5）は実行しない。（ is_pos_{sb} は無効）
- （3） $is_ratio = \tan\left(is_pos_{sb} \times \frac{\pi}{12}\right)$
- （4） 現在の同倍率群 sb 内の全ての i に対して、次の式に示す計算を行う。

$$L_i = L_i \times \frac{is_ratio}{1 + is_ratio}$$

- （5） 現在の同倍率群 sb 内のすべての i に対して、次の式に示す計算を行う。

$$R_i = L_i \times \frac{1}{1 + is_ratio}$$

ステレオ相関の次は、折り返しひずみ削減である。ロングブロックに適応する処理である。ロングブロックの場合、ショートブロックよりも長い時間の音が入っているため、折り返しひずみという雑音が入りやすいので、ここでノイズを取り除き、PCM 標本に悪影響を与えないようにしている。ショートブロックには適応しない。

折り返しひずみ削減の次は、IMDCT を使い、PCM 標本を作る。IMDCT は逆修正離散コサイン変換と言い、エンコーダで使われている MDCT 処理の逆をする処理である。IMDCT の式を図 1 6 に示す。 n は窓掛けされた標本値の個数とする（短いブロックに対する n は 12 で、長いブロックに対する n は 36 とする。）短いブロックの場合は、3 個のブロックをそれぞれ別々に変換する。長いブロックの場合はそのブロックごとに変換する。IMDCT で掛ける窓はブロックタイプに応じて異なる窓を掛ける。各ブロックタイプに対応する窓の形は図 1 7、図 1 8 の通りである。そして、36 個の値を持つブロックの前半分を直前ブロックの後半分に重ね合わせて加算する。現ブロックの後半分は次のブロックで使用するために保存する。重ねあわせ加算された出力値は 32 の多重位相分割帯域ごとに 18 個の時間標本値とする。最初の時間標本値から順に 0～17 の番号をつけ、分割帯域も 0～31 の番号を付ける。多重位相フィルタバンクに入力する前に全ての奇数分割帯域の奇数時間標本値に -1 を乗じる。

Audio_data()

main_data_begin	9bit	0000 0000 0
private_bit	3bit	000

scfsi[ch][scfsi_band]

scfsi[0][0]	0
scfsi[0][1]	0
scfsi[0][2]	0
scfsi[0][3]	0
scfsi[1][0]	0
scfsi[1][1]	0
scfsi[1][2]	0
scfsi[1][3]	0

Ogranule,0channel

part2_3_length[0][0]	12bit	0010 0110 1000
big_values[0][0]	9bit	0010 0100 0
global_gain[0][0]	8bit	1010 0110
scalefac_compress[0][0]	4bit	0000 -> slen1=0, slen2=0
window_switching_flag[0][0]	1bit	1

block_type[0][0]	2bit	01
mixed_block_flag[0][0]	1bit	0

table_select[gr][ch][region] region < 2

table_select[0][0][0]	5bit	01111 -> 符号表15
table_select[0][0][1]	5bit	01000 -> 符号表8

subblock_gain[gr][ch][window]

subblock_gain[0][0][0]	3bit	000
subblock_gain[0][0][1]	3bit	000
subblock_gain[0][0][2]	3bit	000

preflag[0][0]	1bit	0
scalefac_scale[0][0]	1bit	0
count1 table_select[0][0]	1bit	1 -> count1 領域の符号表はB

図8、Audio_data その1

0granule,1channel

part2_3_length[0][1]	12bit	0011 0100 0000 -> 832bit
big_value[0][1]	9bit	0011 0001 0 -> 98個
glbal_gain[0][1]	8bit	1010 0011
scalefac_compress[0][1]	4bit	0000 -> slen1=0,slen2=0
window_switching_flag[0][1]	1bit	1

block_type[0][1]	2bit	01
mixed_block_flag[0][1]	1bit	0

table_select[gr][ch][region] region < 2

table_select[0][1][0]	5bit	11000 -> 符号表24
table_select[0][1][1]	5bit	01100 -> 符号表12

subblock_gain[gr][ch][window] window < 3

subblock_gain[0][1][0]	3bit	000
subblock_gain[0][1][1]	3bit	000
subblock_gain[0][1][2]	3bit	000

preflag[0][1]	1bit	0
scalefac_scale[0][1]	1bit	0
count1table_select[0][1]	1bit	0 -> count1領域符号表はA

図9、Audio_data その2

1granule,0channel

part2_3_length[1][0]	12bit	0011 0000 1010 -> 778bit
big_value[1][0]	9bit	1001 0000 0 -> 288個
glbal_gain[1][0]	8bit	1011 0001
scalefac_compress[1][0]	4bit	0000 -> slen1=0,slen2=0
window_switching_flag[1][0]	1bit	1

block_type[1][0]	2bit	10 -> 3つの短いブロック
mixed_block_flag[1][0]	1bit	0

table_select[gr][ch][region] region < 2

table_select[1][0][0]	5bit	01101 -> 符号表13
table_select[1][0][1]	5bit	00111 -> 符号表7

subblock_gain[gr][ch][window] window < 3

subblock_gain[1][0][0]	3bit	000
subblock_gain[1][0][1]	3bit	000
subblock_gain[1][0][2]	3bit	000

preflag[1][0]	1bit	0
scalefac_scale[1][0]	1bit	0
count1table_select[1][0]	1bit	1 -> count1領域符号表はB

図10、Audio_data その3

1granule,0channel

part2_3_length[1][1]	12bit	0010 1111 1110 -> 766bit
big_value[1][1]	9bit	1001 0000 0 -> 288個
glbal_gain[1][1]	8bit	1011 0001
scalefac_compress[1][1]	4bit	0000 -> slen1=0, slen2=0
window_switching_flag[1][1]	1bit	1

block_type[1][1]	2bit	10 -> 3つの短いブロック
mixed_block_flag[1][1]	1bit	0

table_select[gr][ch][region] region < 2		
table_select[1][1][0]	5bit	01101 -> 符号表13
table_select[1][1][1]	5bit	00111 -> 符号表7

subblock_gain[gr][ch][window] window < 3		
subblock_gain[1][1][0]	3bit	000
subblock_gain[1][1][1]	3bit	000
subblock_gain[1][1][2]	3bit	000

preflag[1][1]	1bit	0
scalefac_scale[1][1]	1bit	0
count1table_select[1][1]	1bit	1 -> count1領域符号表はB

図 1 1、Audio_data その4

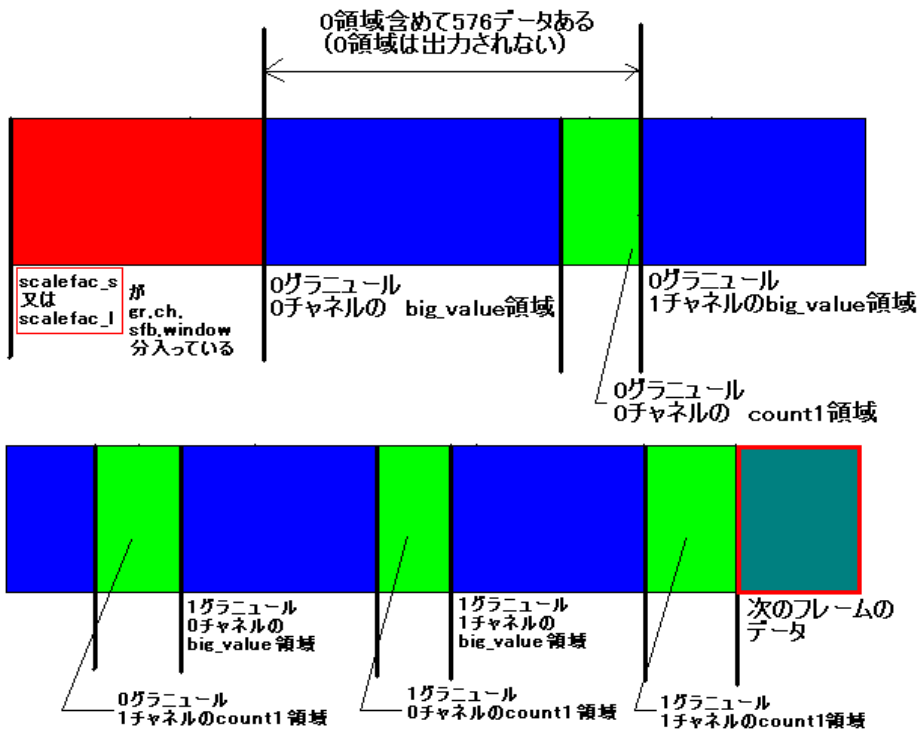
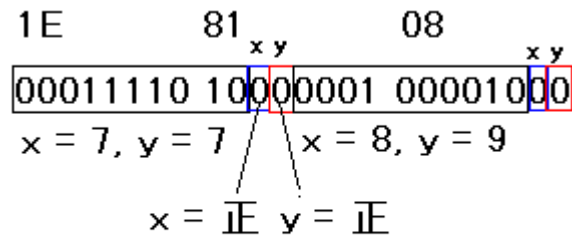


図 1 2、main_data と Huffmancodebit のビットストリーム



- 符号表15
- ブロックタイプは開始ブロック。

図 1 3、big_value 領域の構造 (long_block)

- block_type == 0, 1 又は3(長いブロック)に対して
part2_length = 11 × slen1 + 10 × slen2 → 21個のデータ
- block_type == 2(短いブロック)かつmixed_block_flag == 0に対して
part2_length = 18 × slen1 + 18 × slen2 → 36個のデータ
- block_type == 2(短いブロック)かつmixed_block_flag == 1に対して
part2_length = 17 × slen1 + 18 × slen2 → 35個のデータ

図 1 4、part2_length

短いブロックの逆量子化公式

$$x_{r_i} = \text{Sign}(is_i) \times |is_i|^{\frac{4}{3}} \times 2^{\frac{1}{4}(\text{global_gain}[\text{gr}][\text{ch}][\text{window}] - 210 - 8 + \text{subblock_gain}[\text{gr}][\text{ch}][\text{window}])} \\ \times 2^{-(\text{scalefac_multiplier} \times \text{scalefac_s}[\text{gr}][\text{ch}][\text{sfb}][\text{window}])}$$

長いブロックの逆量子化公式

$$x_{r_i} = \text{Sign}(is_i) \times |is_i|^{\frac{4}{3}} \times 2^{\frac{1}{4}(\text{global_gain}[\text{gr}][\text{ch}][\text{window}] - 210)} \\ \times 2^{-(\text{scalefac_multiplier} \times (\text{scalefac_l}[\text{sfb}][\text{ch}][\text{gr}] + \text{preflag}[\text{gr}] \times \text{pretab}[\text{sfb}]))}$$

図 1 5、逆量子化公式

$$X_i = \sum_{k=0}^{\frac{n}{2}-1} X_k \cos\left(\frac{\pi}{2n}\left(2i+1+\frac{n}{2}\right)(2k+1)\right) \quad i = 0, \dots, n-1$$

図 1 6、IMDCT 式

・ロングブロック

(a) block_type = 0 (普通のブロック)

$$z_i = x_i \sin\left(\frac{\pi}{36}\left(i + \frac{1}{2}\right)\right) \quad i = 0, \dots, 35$$

$$= x_i \sin\left(\frac{\pi}{36}\left(i + \frac{1}{2}\right)\right) \quad i = -, \dots, 35$$

(b) block_type = 1 (開始ブロック)

$$z_i \begin{cases} x_i \sin\left(\frac{\pi}{36}\left(i + \frac{1}{2}\right)\right) & i = 0, \dots, 17 \\ x_i & i = 18, \dots, 23 \\ x_i \sin\left(\frac{\pi}{12}\left(i - 18 + \frac{1}{2}\right)\right) & i = 24, \dots, 29 \\ 0 & i = 30, \dots, 35 \end{cases}$$

(c) block_type = 3 (停止ブロック)

$$z_i \begin{cases} 0 & i = 0, \dots, 5 \\ x_i \sin\left(\frac{\pi}{12}\left(i - 6 + \frac{1}{2}\right)\right) & i = 6, \dots, 11 \\ x_i & i = 12, \dots, 17 \\ x_i \sin\left(\frac{\pi}{36}\left(i + \frac{1}{2}\right)\right) & i = 18, \dots, 35 \end{cases}$$

図17、ロングブロックの窓

(d) block_type = 2(短いブロック)

3個の短いブロックは別々に窓掛けする。

$$y_i^{(j)} = x_i^{(j)} \sin\left(\frac{\pi}{12}\left(i + \frac{1}{2}\right)\right) \quad \begin{array}{l} i = 0, \dots, 11, \\ j = 0, \dots, 2 \end{array}$$

窓掛けされた短いブロックは重ねあわせ加算し連結しなければならない。

$$z_i \left\{ \begin{array}{ll} 0 & i = 0, \dots, 5 \\ (1) & i = 6, \dots, 11 \\ y_{i-6}^{(1)} & \\ (2) + (2) & i = 12, \dots, 17 \\ y_{i-6}^{(2)} \ y_{i-12}^{(2)} & \\ (2) + (3) & i = 18, \dots, 23 \\ y_{i-12}^{(2)} \ y_{i-18}^{(3)} & \\ (3) & i = 24, \dots, 29 \\ y_{i-18}^{(3)} & \\ 0 & i = 30, \dots, 35 \end{array} \right.$$

図18、ショートブロックの窓

5、MP3 のエンコード方法

次に、エンコードの方法を説明する。MP3 は仕様でエンコード方法が定められていないので、自由にエンコードすることが出来る。エンコードの流れは、大まかに言うと 1、量子化、2、ハフマン符号化、3、付加情報出力、4 ビットストリームに出力。この手順でエンコードできるはずである。この中でも、一番計算が大変なのが 2 のハフマン符号化で、32 個ある符号表から、音質とファイルサイズの両方を見て、最適な符号表を選び出さなければならない。エンコードの方法について、仕様書では「1 例」としてエンコードの方法の説明が載っているその説明は以下の通りである。

MP3 は聴覚心理アルゴリズムを用いている。聴覚心理アルゴリズムの主要部分は図 14 に示す。

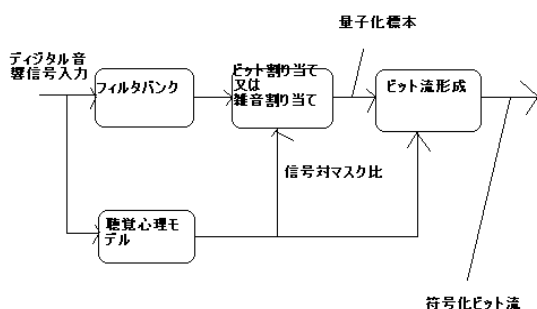


図 14、MP3 の符号化方法

フィルタバンクでは多重位相フィルタバンクおよび混成多重位相フィルタバンク（MDCT フィルタバンク）で各々の時間および周波数上での特殊な写像を行う。フィルタバンクから出力した標本を量子化する。また、このフィルタバンクはレイヤ 3 の場合、入力信号に依存した 6×32 個または 18×32 個の周

波数帯域の解像度を持つ。

聴覚心理モデルは MP3 の場合、モデル 2 を当てる。最終出力は、分割帯域のグループに対する信号対マスク比（SMR）とする。

ビット割り当て又は雑音割り当てでは、フィルタバンクの出力標本および SMR を考慮しながら、雑音割り当てしていく。

ビット流形成では、量子化されたフィルタバンクの出力、雑音割り当て、およびほかの必要な付加情報の符号化およびフォーマットを行う。レイヤ 3 の場合はハフマン符号も行う。ビット割り当てまたは雑音割り当ての方法は、レイヤ 1 およびレイヤ 2 ではビット割り当て処理として、レイヤ 3 では雑音割り当てループとして実現する。いずれの場合も、結果を量子化パラメタおよび量子化された標本として、ビット流形成に伝達する。

ビット流形成では、分割帯域の各々の信号に対してのハフマン符号を出力する。ここまでの結果をビットストリームに出力することで、MP3 の 1 フレームが出来上がる。これを繰り返すことにより、一曲分の MP3 が出来上がる。

また、符号器の入力側に遮断周波数 2 ~ 10 kHz の高域通過フィルタを含んでいることが望ましい。これによって、もっとも低域側の分割帯域に対して、必要以上に多くビットを割り当てることを防ぎ、結果的に全体の音質を向上させることになる。と、仕様書には書かれている。

6、エンコードの速度の向上

MP3 がエンコードする中で、一番時間がかかる場所が、ハフマン符号器である。なぜならば、32 個もあるハフマン符号表から一つだ

け適切な符号表を選ばなければならないからである。この部分にはかなり複雑な計算が必要であり、計算ステップ数が多くなることからコンピュータで計算するにも、その分だけ時間がかかってしまう。このハフマン符号表を選び、ハフマン符号するハフマン符号器のアルゴリズムを改良することにより、エンコード速度を向上させることが出来ると考えられる。しかし、残念ながらまだ実際に動くエンコーダを作っていないため、実証することができない。

7、MP3の今後

インターネットの普及と技術向上によって、コンピュータの技術も向上した昨今、音楽と言うものは昔みたいにテレビやラジオで聞くよりも、パソコンやポータブル機器で聞くほうが主流になりつつある。いままで CD ショップへ行かなければ、手に入れることの出来なかった、人気アーティストの新盤が、今ではインターネットを使い自宅のパソコンでボタン一つ購入できる。しかも、その場で、電線を通して手元にやってくる。そんな時代になった今、MP3の技術は大きく社会に貢献している。MP3がなければ、このようなミュージックライフはやってこなかっただろう。今では、MP3のほかに、VorbisのOGG、i-Podなどでよく使われるATRAC、5.1チャンネルのオーディオデータが作成できるAAC、Microsoftが作ったWMAなど、MP3以外の圧縮オーディオ規格は沢山出てきている。最近では、MP3がやっている人間には聞こえない高音域を削ると言うことをやめて、可逆性圧縮（情報の劣化が無い圧縮）が出来る規格まで存在する。一時期、景気がよかった頃は、

大手のオーディオメーカーから、小さな町工場まで、ポータブルオーディオの新商品をこぞって出していたのに、今ではそのオーディオブームも下火になり、不況の風が追い討ちを掛けるように、この業界に襲い掛かっているようだ。しかし、圧縮音源はポータブルオーディオだけが活躍する舞台ではない。パソコンで視聴できる動画の音源部分に使われているのだ。「youtube」や「ニコニコ動画」、「veoh」などの動画サイトでも、この圧縮音源は使われている。

MP3に限らず、圧縮音源はこれからもきっと、進化し続けるだろう。現にMP3PROという新しい規格が誕生しようとしている。このMP3PROはMP3を可逆性圧縮に改良したものになりそうである。可逆性圧縮にしたほうが、従来の方式より音質面において有利ではある。しかし、ファイルサイズは少し大きくなる。

また、音楽の再生の面だけでなく、録音の面でもMP3は非常に役に立っている。今では、小型のICレコーダが、MP3に直接録音できるくらいまで、技術が向上している。

8、感想

研究を本格的に始めたのが今年の10月でした。非常に遅いスタートでした。ここまで来るのに、こんなにも時間がかかるとは思わず、非常に後悔しております。MP3に関する資料が少なく、量子化などで使われる数式が難しく、それを解読するのも時間もかかり……。もっと早く取り組むべきでした。すみませんでした。

9、参考文献

[1]、MP3 の復元

<http://www.ctrl.nara-k.ac.jp/~idlab/2004member/koba/mp3.doc>

[2]、日本工業標準調査会：データベース JIS
詳細表示- JISX4323

<http://www.jisc.go.jp/app/pager?id=31259&%23jps.JPSH0090D:JPSO0023:/JPS/JPSO0090.jsp=&AKKNB vJISJISNO=X4323>