

未来の没入型オーディオ・ システム、コンピュータ・ ビジョンが音声の再生を支える

著者: Santosh Singh、コンスーマ・システム・アプリケーション担当シニア・エンジニア **Aravind Navada、** コンスーマ・システム・アプリケーション・アジア担当ディレクタ

概要

民生向けのエンターテインメントの分野では、没入型の体験 (エクスペリエンス) がより強く求められるようになっていま す。つまり、ユーザが物理的な現実と区別できない形でコンテ ンツを楽しめるようにすることが重視されています。より優れ た没入感を実現する上では、音声が重要な役割を果たすことに なります。そのような背景を踏まえ、アナログ・デバイセズは 近い将来、音声を再生するための新たな手法が活用されるよう になると考えています。その手法は、ビジョン・インテリジェ ンスを活用したオーディオ・システムとして具現化されるはず です。そうしたシステムを開発するには、人間の脳が音声を処 理して局所化する方法について深く理解する必要があります。 技術的に言えば、最先端のToF (Time of Flight) イメージャ とクラス最高のDSPが重要な役割を果たします。それらの適 切な組み合わせにより、次世代の没入型オーディオ・システム を開発するための理想的なプラットフォームが実現されます。

あらゆる新世代の民生向けエンターテインメント機器について語 る際には、例外なく没入感という言葉が使われます。それは具体 的には、どういう意味なのでしょうか。1999年に「マトリック ス」という映画がヒットしました。その登場人物であるモーフィ アスは、ネオに対し、「匂いを嗅ぎ、味を感じ、物体に触れると いう行為は現実のものなのか」と尋ねます。その上で、ネオがそ れまで現実だと思っていたことは、人間の感覚を欺くためにコン ピュータによって作られた虚構にすぎなかったという事実を示し ます。これが、まさに没入感という言葉が意味するところです。 すべての人工的な没入体験は、そのような世界が目指すべきゴー ルだと捉えて開発されています。

体験の中に本当に入りこんでいるような感覚が得られるようにす るためには、音声そのものと、それをどのように体験させるのか ということが重要になります。それらは、高揚感を演出するため

の非常に重要な要素です。音声は、脳の中の原始的な反応として 始まり、任意の状況に対する私たちの最初の反応へと発展してい きます。脳は音声を活用することにより、周囲の環境や状況を表 す明確なイメージを形成します。人工的に作られた没入体験を脳 に信じ込ませることによって意図した没入感を与える上で、音声 は非常に重要な役割を果たします。

音声を再生する技術は、長い年月の間に大きな進化を遂げまし た。当初使われていたのは基本的なモノラル・オーディオ・シス テムで、オーディオ用のチャンネルは1つしかありませんでした。 それが現在では、非常に高度なサラウンド・サウンド・システム が使われるようになっています。ホーム・シアターでよく用いら れるのは、5.1ch (6チャンネル)、7.1ch (8チャンネル) といっ た最小構成のシステムです。それに対し、映画館のスクリーン向 けには、64以上のチャンネルを備える大規模なシステムが提供 されています。ただ、そうしたシステムにおいて、音声に対する 空間的な感覚と精度は、スピーカの数とそれらの位置によって制 限されます。

次世代の没入型オーディオ・システムは、音声の再生方法に新た な手法を適用することによって実現されるはずです。その手法を 実現するためには、人間の脳が音声を処理して局所化する方法に 関する深い理解が必要になります。新たなシステムでは、多数の スピーカをリスナーの周りに配置する必要はありません。そのよ うなことを行わなくても、リスナーの周囲360度に対応した没入 感あふれる音声体験がホーム・シアターにもたらされます。その ような新たなシステムを開発するためには、リスナーとそのリス ニング環境に関する高度な知識を有していなければなりません。 そうした知識が不足していた場合、没入型のオーディオ・システ ムに求められる要件を満たすことは難しいでしょう。では、次世 代の没入型オーディオ・システムを実現するためには、具体的に は何が必要なのでしょうか。その答えは、音声の再生技術にビ ジョン・インテリジェンスを融合させるというものです。











現実の世界で音声を自然に耳にするとき、人間の脳は、わずか2 つのオーディオ信号に基づいて音源に関する空間的な手がかりを 導き出します。2つのオーディオ信号の一方は左耳、もう一方は 右耳に届きます。これは、人間の両眼による視覚系の仕組みによ く似ています。その仕組みとは、左目と右目の視界を比較するこ とによって、脳内で奥行き感を生成するというものです。それと 同様に、人間の脳は左耳と右耳に届く音声を処理し、その振幅と 遅延時間を比較することによって、音源の位置を大まかに割り出 します。この能力は、人類の進化の過程で徐々に発達していきま した。人類が自然界で生き延びるためには、そのような能力が不 可欠だったからです。

自然な聴覚体験を再現することを目的とした手法に、バイノーラ ル・オーディオ再生 (Binaural Audio Reproduction) というも のがあります。これは、信号処理を利用することにより、現実の 世界で音を聴いた場合と同じ状態が得られるように、左耳と右耳 のそれぞれに向けて2つのオーディオ信号を生成するというもの です(図1)。しかし、この手法を具現化するのは非常に難しく、 様々な問題が発生します。

まず、バイノーラル・オーディオ再生に使用可能な音声を録音す るにはどうすればよいのでしょうか。最も簡単な方法は、2つの マイクを実際に人の左右の外耳道に1つずつ配置し、その位置で 音声信号を録音するというものです。この方法をバイノーラル録 音と呼びます。音声の再生は、ヘッドフォンを介してリスナーの 耳に向けて行われます。では、この手法は狙いどおりに機能する のでしょうか。その答えは、ある条件が満たされるならば「イエ ス」です。その条件とは、「同じ人を対象として音声の録音と再 生が行われる」というものです。言い換えれば、録音時と再生時 の対象人物が異なると、狙いどおりには機能しないということで す。その原因は、私たちの脳が音声を局所化する仕組みにあり ます。私たちの頭/耳介/体は、周波数領域において特定のシ グネチャを放置することにより、脳による音声の局所化プロセス を補助するという形で影響を及ぼします。どのシグネチャが対象 になるのかは人によって異なります。その違いは、頭部伝達関数 (HRTF: Head-related Transfer Function) によって表されます。 バイノーラルの手法を適切に機能させるには、音声を再生する際、 リスナーの耳において、HRTFが音声に与える影響を正確に再現 しなければなりません。

HRTFについては、リスナーごとに測定を行ってパーソナライズ する必要があります。つまり、汎用的なソリューションによって 対応できるものではありません。ある人が、別の人のHRTFを使 用して生成された音声を聞かされたとします。その場合、人が音 声を局所化する能力は著しく低下します。このことは、複数の研 究によって確認されています 1, 2, 3。

ラウド・スピーカによってバイノーラル・オーディオを再生する 際には、更に難易度の高い課題が浮上します。まず、複数のス ピーカからの音声信号は互いに干渉を引き起こします。この現象 は、クロストーク効果と呼ばれています(図2)。また、リスニン グ環境によっては、リスナーの耳に届く前に、音声に対して望ま しくない影響が及ぶ可能性もあります。

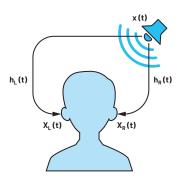


図 1. 音源 x(t) からの音声を聞く際の仕組み。 $X_L(t)$ は左耳に届く オーディオ信号、 $X_R(t)$ は右耳に届くオーディオ信号です。

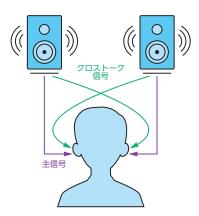


図2. ステレオ・スピーカにおける クロストーク効果

目標とするのは、自然な聴覚体験を正確に再現することです。ス ピーカのクロストーク、HRTFのパーソナライズの必要性、リス ニング環境の影響は、その実現を阻む主要な要因です。このよう なバイノーラル・オーディオ再生が抱える課題を解決するために 役立つものがあります。それは、リスナーとリスニング環境に関 連して必要になるすべての詳細情報を捉えることが可能なビジョ ン・システムです。

例えば、コンピュータ・ビジョンのアルゴリズムにデータを供給 するカメラを構築したとします。それを利用すれば、リスニング 環境の3次元構造の詳細(音声を聞いている部屋の形状、各所 表面の幾何学的な測定データ、物体の存在)を捉えることがで きます。そして、その情報(洞察)を利用すれば、リスニング環 境が音声に及ぼす影響を計算することが可能になります。また、 音声再生システムにおいて適切な係数を備えるフィルタを使用す れば、その望ましくない影響を排除することができます。実は、 この種のシステムは、ホーム・シアターで利用するオーディオ向 けのものとしては既に存在していました。従来は、キャリブレー ションを実行する際、マイクを使用して音声に対する部屋の影響 を把握するということが行われていました。ただ、この手法では、 リスナーが音声を聞く位置が変わったり、部屋に構造的な変化が 生じたりした場合、処理をやり直す必要がありました。

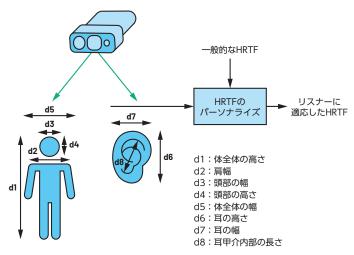
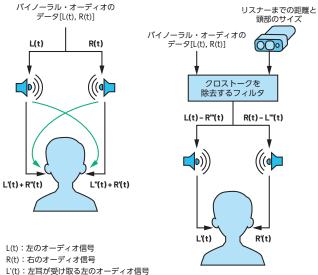


図3. 身体の測定値に基づく HRTFのパーソナライズ

ビジョン・システムを利用すれば、身体の位置を測定したり、構 造に関する詳細情報を把握したりすることができます⁴。それに より、空間的な手がかりを正確に把握し、それをレンダリングす ることで、HRTFをパーソナライズするための計算を行うことが できます(図3)。例えば、スピーカに対するリスナーの頭の位置 と頭部のサイズに関する情報を使用し、クロストークを除去する アルゴリズムを適用するといった処理を行います。それにより、 ラウド・スピーカからのバイノーラル・オーディオをリアルタイ ムにレンダリングすることができます。結果として、リスナーが 部屋の中を動き回ったとしても、理想的な音声体験が得られるよ うになります (図4)。

ビジョン・システムを利用する場合には、ある一般的な問題が浮 上します。それは、ユーザのプライバシーの侵害をどのようにし て防ぐのかというものです。例えば、エッジの専用プロセッサを 使用して、ビジョン・データの分析を行うのであれば、ユーザの プライバシーが侵害されることはありません。その場合、ビジョ ン・システムのカメラでキャプチャされた情報はリアルタイムに 処理されます。それらのデータを保存したり、別のリモート・マ シンに転送したりする必要はありません。そのため、プライバ シーの問題は生じないということです。

バイノーラル・オーディオに 対応してない 従来型のシステム クロストークを除去する バイノーラル・オーディオ対応のシステム



L"(t): 右耳が受け取る左のオーディオ信号(クロストーク・オーディオ)

L'''(t):右耳のL''(t)を除去するために生成された左のクロストーク除去信号

R'(t): 右耳が受け取る右のオーディオ信号

R"(t): 左耳が受け取る右のオーディオ信号 (クロストーク・オーディオ) R'"(t): 左耳のR"(t)を除去するために生成された右のクロストーク除去信号

> 図4. クロストークを除去するための実装方法。 フリー・フィールド・スピーカ・システムによるバイノーラル・オーディオの再生を実現します。

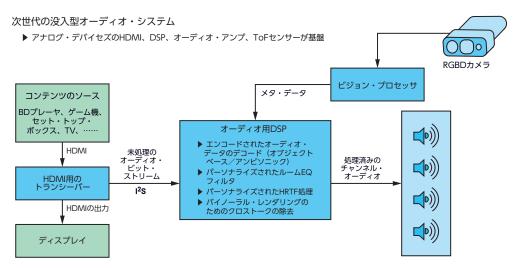


図5. 次世代の没入型オーディオ・システム

ここまでに説明したように、次世代の没入型オーディオ・システ ムを構築するには、ビジョンとオーディオの融合が必要です。つ まり、ハードウェア・プラットフォームをベースとして、図5の ようなシステムを構築しなければなりません。アナログ・デバイ セズの場合、最新のマルチコアSHARC® DSPと最先端のToFイ メージャを主要な構成要素としてプラットフォームを構成してい ます。

アナログ・デバイセズの [ADSP-SC598] は、2個のSHARCコ アと1個のArm® Cortex®-A55コアを搭載したSoC (System on Chip)です。大容量のオンチップ・メモリと、外付けメモリ 用のDDRインターフェースを備えています。また、遅延を小さ く抑えつつ、メモリに関連する負荷の大きい演算に対応すること が可能です。同製品は、真の没入型オーディオ・システム向け のものとして理想的なプラットフォームだと言えます(図6)。

ADSP-SC598であれば、SHARC DSPが備える演算リソースを使 用して負荷のパーティショニングを実現できます。例えば、オー ディオ・データのデコードに関連する処理を1つ目のSHARCコ アで実行し、オーディオを再生するための後処理とパーソナライ ズを2つ目のSHARCコアで実行するといった具合です。また、 Arm Cortex-A55は、様々な制御の処理に使用できます⁶。図5 に示したビジョン・システムは、RGBカメラと深度カメラを組み 合わせるか、またはスタンドアロンの深度カメラによって実現で きます。アナログ・デバイセズは解像度が1MPのToFイメージャ 「ADSD3100」を提供しています。この製品は、mmのレベルの 分解能で深度マップをキャプチャし、様々な照明の条件下で機能 するように設計されています。これを使用すれば、パーソナライ ズ用のアルゴリズム(クロストークの除去、ルーム・イコライゼー ション [EQ]、HRTFのパーソナライズなど) の適用対象となる 非常に精度の高い幾何学的な測定データが得られます。

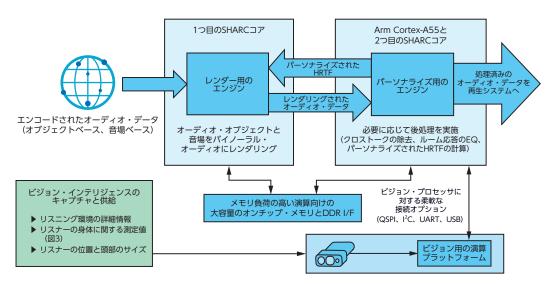


図6. 次世代の没入型オーディオ・システムにおける パーティショニング

アナログ・デバイセズは、深度の検出を担うToFイメージャとし てADSD3100を採用したToFモジュール「ADTF3175」も提 供しています。その解像度は1MPで、視野角(FOV: Field of View) は75°×75°です。ADTF3175には、ADSD3100用のレ ンズと光学的バンドパス・フィルタに加え、光学系/レーザー・ ダイオードとドライバ/フォトディテクタを備える赤外照射光 源、フラッシュ・メモリ、ローカル用の電源電圧を生成するレギュ レータが統合されています。このモジュールは、複数の範囲と分 解能のモードに対して完全にキャリブレーションされています。

完全な深度検出システムでは、ADTF3175によって得られる未 処理のイメージ・データを、ホスト側のシステム・プロセッサ または深度ISP (Image Signal Processor) によって外部で処 理することになります。ADTF3175から出力されるイメージ・ データは、MIPI CSI-2 (Mobile Industry Processor Interface Camera Serial Interface 2) に対応する4レーンのトランスミッ タ・インターフェースを介してホスト・システムに引き渡されま す。モジュールのプログラミングと制御には、4線式のSPI (Serial Peripheral Interface) とI²Cのシリアル・インターフェースを使 用します。

アナログ・デバイセズは、開発プラットフォーム [EVAL-MELODY-8/9」、評価用ボード「EV-2159X/SC59x-EZKIT」、 Eclipseベースのエディタ・ツール「CrossCore® Embedded Studio」も提供しています。これらを利用すれば、ADSP-SC598 を使用するアプリケーションのリアルタイムの配備とデバッグを 直ちに開始することができます⁷。

Melodyプラットフォームは、AVR (AV Receiver) やサウンド バーのアプリケーションに適したシグナル・チェーン・ソリュー ションです。この統合型のソリューションには、ビデオ、DSP、 オーディオ、電力、ソフトウェアの面でアナログ・デバイセズの 中でもクラス最高レベルのコンポーネントが盛り込まれていま す。これを利用すれば、例えば年に1度のアップグレードのスケ ジュールに合わせて、最新技術を採用した製品を素早く市場に投 入することができます⁸。

図7に示したのは、次世代の没入型オーディオ・システムを実現 するためのハードウェア・プラットフォームです。ご覧のように、 ToFモジュールであるADTF3175をビジョン用の演算プラット フォームに接続し、そのプラットフォームをMelodyのボードに 接続しています。なお、RGBカメラとADTF3175を組み合わせ れば、ビジョンに関する高度な分析を行うためのRGBD(色+深 度) カメラを構成することが可能です。

まとめ

アナログ・デバイセズは、DSP、HDMI用のトランシーバー、D 級アンプ、ToFイメージャなどを含むポートフォリオを有してい ます。現実世界の音声と区別がつかない音声を再生可能な真の 没入型オーディオ・システムの実現に向けて、全力で取り組みを 行っています。

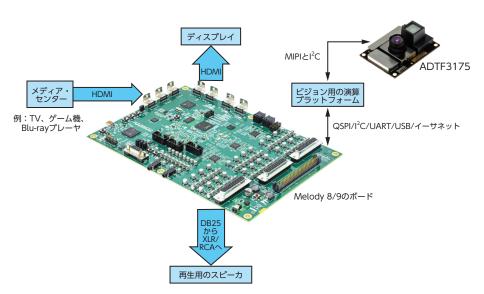


図7. アナログ・デバイセズのプラットフォームによって構成した 没入型のオーディオ・システム

参考資料

¹ Philipp Paukner, Martin Rothbucher, Klaus Diepold Sound Localization Performance Comparison of Different HRTF-Individualization Methods (異なるHRTFのパーソナライズ手 法による音声の局所化性能の比較) L Technische Universität München、2014年4月

² Parham Mokhtari, Ryouichi Nishimura, Hironori Takemoto. [Toward HRTF Personalization: An Auditory-Perceptual Evaluation of Simulated and Measured HRTFs (HRTFのパーソナライズに向けて:HRTFのシミュレーション 結果と実測値の聴覚知覚評価)」International Conference on Auditory Display (ICAD)、2008年7月

³ Jenny Claudia、Christoph Reuter [Usability of Individualized Head-Related Transfer Functions in Virtual Reality: Empirical Study With Perceptual Attributes in Sagittal Plane Sound Localization(仮想現実におけるパーソナ ライズされたHRTFの有用性:矢状面の音声局所化における知覚 属性による実証研究)」JMIR Serious Games、2020年9月

⁴ Geon Woo Lee, Hong Kook Kim [Personalized HRTF Modeling Based on Deep Neural Network Using Anthropometric Measurements and Images of the Ear (身 体測定と耳の画像を用いたディープ・ニューラル・ネットワー クに基づくパーソナライズされたHRTFのモデル)」Applied Sciences、2018年11月

⁵ Sanket Nayak, Mitesh Moonat [EE-436: Using ADSP-SC59x/2159x High Performance FIR/IIR Accelerators (ADSP-SC59x/2159xによる高性能のFIR/IIRアクセラレータ)」 Analog Devices、2022年6月

⁶ 「ADSP-SC59x/2159x SHARCシリーズがポートフォリオの性 能を2倍に向上し、複雑なオーディオ・アプリケーション向けプ ラットフォームを実現」Analog Devices

⁷ [CrossCore® Embedded Studio] Analog Devices

⁸ [ホーム・シアターとゲーム] Analog Devices

著者について

Santosh Singhは、2016年にインドのビルラ工科大学で電 子工学と通信工学の学士号を取得しました。インターンを 経て、アナログ・デバイセズに入社。デジタル設計者とし てキャリアをスタートし、最先端の技術を活用した民生市 場向けの様々な製品を担当しました。現在は、シニア・シ ステム・アプリケーション・エンジニアとして、コンスー マ・ビジネス・ユニットに所属。アプリケーションと技術の ギャップを埋めるべく、システムの主要な課題を解決する ことに主眼を置いて業務に取り組んでいます。

Aravind K. Navadaは、アナログ・デバイセズで、アジア 全域のお客様を対象とするコンスーマ・システム・アプリ ケーション・チームを統括しています。20年以上にわた り、ICの設計に従事。様々なミックスド・シグナルIoT製 品や民生向けのSoC製品を市場に投入してきました。現在 は、製品のターゲットを最大化することを目指し、民生市 場向けのプラットフォームとソリューションの開発を担当 しています。インドのバンガロールにあるヴィスエソラヤ 国立工科大学 (UVCE) で電子工学と通信工学の学士号を 取得。米国テキサス大学ダラス校でマイクロエレクトロニ クスに関する修士号も取得しました。

EngineerZone® オンライン・サポート・コミュニティ

アナログ・デバイセズのオンライン・サポート・コミュ ニティに参加すれば、各種の分野を専門とする技術者と の連携を図ることができます。難易度の高い設計上の問 題について問い合わせを行ったり、FAQを参照したり、 ディスカッションに参加したりすることが可能です。

△ ADI EngineerZone[®]

Visit ez.analog.com

*英語版ソート・リーダーシップ記事はこちらよりご覧いただけ ます。



VISIT ANALOG.COM/JP