

## 要素論から全体論へ ～全体から入る音声情報処理への招待～

峯松 信明<sup>†</sup> 西村多寿子<sup>††</sup> 朝川 智<sup>†</sup> 櫻庭 京子<sup>†††</sup> 齋藤 大輔<sup>†</sup>

<sup>†</sup> 東京大学大学院新領域創成科学研究科, <sup>††</sup> 東京大学大学院医学系研究科, <sup>†††</sup> 清瀬市障害者福祉センター  
E-mail:{mine,asakawa,dsk\_saito}@gavo.t.u-tokyo.ac.jp, nt-tazuko@ams.odn.ne.jp, sakuraba@mt.d.biglobe.ne.jp

**あらまし** 一つの言語には通常数十種類の音素 (phoneme) がある。しかし音素の音的実体は前後文脈 (音素環境) などによって多様に變形し, 異音 (allophone) と呼ばれる。音素と比較して種類数も多く, より具体的な音的現象に対応している。しかし奇妙なことに, これら音的事象を記号を用いて記す場合, 性別, 年齢, 収録・伝送機器特性などによる音の變形 (非言語的要因による音響的變形) は一切無視される。その音響的變形が幾ら大きくても, である。音声認識の音響モデリングは, 凡そ, 異音に相当する音事象を triphone としてモデル化しているが, 「非言語的變形の無視」を実装するために, 数万人の話者から, 様々な環境で収録した音サンプル群を統計的にモデル化している。本稿では, 「非言語的變形の無視」の実装は, 集めることではなく, 音事象間の差異を捉えることで可能となることを数学的に示し, 極めて少数の話者の音声で, 不特定話者音声認識が可能であることを示す。提案する枠組みでは, 音的要素をモデル化するのではなく, 音的差異に着眼し, 差異を集めることで構成される全体的な音的構造をモデル化する。  
**キーワード** 変換不変, 音色差異, 相対音感, 音声模倣, 言語障害

## From Reductionism to Holism

— Holistic Approach for Speech Information Processing —

N. MINEMATSU<sup>†</sup>, T. NISHIMURA<sup>††</sup>, S. ASAKAWA<sup>†</sup>, K. SAKURABA<sup>†††</sup>, and D. SAITO<sup>†</sup>

<sup>†,††</sup>The University of Tokyo, <sup>†††</sup>Kiyose-shi Welfare Center for the Handicapped  
E-mail:{mine,asakawa,dsk\_saito}@gavo.t.u-tokyo.ac.jp, nt-tazuko@ams.odn.ne.jp, sakuraba@mt.d.biglobe.ne.jp

**Abstract** A language generally has several tens of phonemes. Acoustic substances of the phoneme depend upon its phonemic environment and the context-sensitive phonemes are called allophones. The number of the allophones in a language is naturally much larger than that of the phonemes. Although the allophones represent finer acoustic differences between linguistic sounds, it is very strange that they completely ignore the acoustic variations in an allophone caused by differences in age, gender, microphone, room, etc. The triphones, which are acoustic models widely used in speech recognition, correspond to the allophones and the ignorance of the acoustic variations caused by the non-linguistic factors are implemented by collecting speech samples from an enormous number of speakers and training statistical acoustic models of the individual allophones. In this paper, it is mathematically shown that the ignorance can be realized not by collecting samples but by capturing timbral differences between two sounds. Then, the possibility of speaker-independent speech recognition only with a very small number of training speakers is experimentally examined. In the proposed framework, what is modeled is not the elementary sound substances, i.e., reductionism, but the holistic sound system exclusively composed of the timbral differences, i.e., holism.

**Key words** transform-invariance, timbral difference, relative sense of sounds, vocal imitation, language disorder

### 1. 「集めること」「合わせること」は必要なのか?

数十万人の話者の声を用いて構築された各異音の統計モデルに基づく音声認識エンジンが市販されている [1]。一方, 言語の唯一のユーザである人間を眺めた場合, 数十万人の声を聞いて初めて頑健な音声情報処理が可能となった個体は恐らく存在しない。子供の言語発達を考えた場合, 乳児の聞く声の大半は母親, 父親の声である。更に自らが話せるようになると, その子の聞く声の半分は自らの声である (speech chain)。話者バランスのとれた大規模音声コーパスを要求する現在の音響モデリング技術は, 些か不可思議な技術体系と言わざるを得ない。

Speaker adaptive training (SAT) という枠組みに従えば [2], 基本的には架空の特定話者音響モデルを用意し (例えば自己聴取音に基づく音響モデル), 音声認識時には, 常時, 話者適応を施す形で音響的照合が可能である。通常の音声コミュニケーションを考えれば, 書き起こしが与えられない教師無し適応である必要があるが, この場合, 認識性能の向上には限界がある。

現在の音声認識技術は, 音響モデリングに関して言えば, 与えられた「音」を学習し, 「音」のモデルを構築し, 新たに与えられた「音」が, 以前与えられた「音 (カテゴリ)」の何れかを推定する技術である。学習話者が一人であれば, 特定話者音響モデルとなり, それを用いた音声合成も可能である [3]。

幼児の言語獲得は、音声模倣という言葉で表現される[4]。他の霊長類には無い行動であると、しばしば指摘される[5]。しかしこの時、幼児は「音」を学習しない。両親の「声」を模倣しようとし、一方九官鳥は「声」を真似る。「音」を真似る。優秀な九官鳥は聞けば飼い主が分かる[6]が、どんなに優秀な幼児を聞いても、父親を当てることはできない。九官鳥は「音」を学習し、(恐らく「音」のモデルを構築し)以前聞いた「音」に反応し、その「音」を鳴管を使って生成する。幼児の模倣を「音声」模倣と呼ぶならば、九官鳥の模倣は「声」模倣である。「声」は音そのものである。では、「音声」とは音の何を指すのだろうか?本節では以下、音声と声を特に区別して記述する。

音声認識系を構築する時に、鳥と飛行機を例にとり「人間を模倣する必要は必ずしも無い」という意見は頻りに耳にする。では、現在の音響モデリングは、幼児と九官鳥、何れに似ているだろうか?間違いなく、九官鳥、である。「音声」認識とは名ばかりで、実は「声」認識というのが、その実体である。

「声」は如何にしたら「音声」となるのか、音声学(音声科学も含まれるであろう)が導いた方程式は下記である。

$$\text{音声} = \sum_{\text{話者, マイク, 場所}} \text{声, 或いは, 音声} = \text{変換関数 } f(\text{声})$$

「集めること」「合わせること」で「音声」になるとの期待である。しかし、幼児は数十万の話者の「声」、及びその書き起こしテキストがなくても、電話越しのお婆ちゃんと会話し始める。そもそも、「変換関数  $f(\text{声})$ 」は「別の人の声」にしかならない。

### 何か、根本的に、間違っているのだろうか?

「声」を音事象の連続体(ストリーム)と考えた場合、個々の音事象は、話者・収録機器の違いで当然変形を被る。音韻意識が未熟な幼児は、両親の「声」を仮名列に落とし、個々の仮名を音にして出すという技は使えない<sup>(注1)</sup>。変形を被った音ストリームに対し、彼らが模倣しているのは「音」ではない。彼らが模倣している、音ストリーム内に符号化されているコンテンツ(つまり音声)を直接的にモデル化したいのであれば、音そのもの(つまり声)に対する音響モデリング技術、即ち、音の生成モデル(generative model)は甚だ不適切である。母親の「おはよう」を模倣しても、父親の「おはよう」を模倣しても、幼児の同じ「オハヨウ」となる。結局、この三者に共通して存在する話者不変の音響現象をモデル化する必要性が生じる。

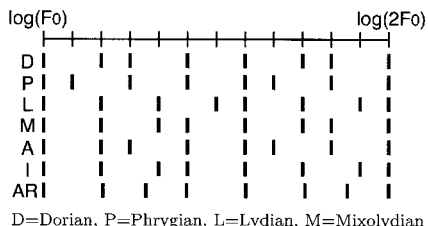
### 2. 「音」の何をモデル化すべきなのか?

声を音韻(音シンボル)列として認知していない幼児が、如何にして母親の「おはよう」と父親の「おはよう」の(類似性ではなく)同一性を感覚するのだろうか?発達心理学は「幼児は単語全体の語形・音形(語ゲシュタルト[4])を獲得し、その後、個々の分節音を獲得する」と主張する[8]。上記はこの話者不変と思いき「語ゲシュタルト」の音響的定義を問うている。筆者らの一部は発達心理学、言語獲得研究者に広く問いかけたが[9]、適切な回答は無い。「惑星」の定義がされないまま議論

(注1): そもそも彼らは「しりとり」を行なうことが困難である[7]。



図1 とあるメロディー(ハ長調)とその移調版(ト長調)



D=Dorian, P=Phrygian, L=Lydian, M=Mixolydian  
A=Aeolian, I=Ionian, AR=Arabian

図2 6種類の古典的教会音階とアラビア音階

を繰り返した天文学と同じである、との意見も得た。その物理的存在は議論せず、存在を仮定した議論を繰り返すのみである。

そもそも、二つの音の同一性を感覚するのに、その二音の物理的同一性が必要なのだろうか?人間は他の霊長類と異なり、全く異なる物理特性を有する二音(ある環境下では)「同一である」と感覚する能力を持っている[10]。相対音感である。

#### 2.1 調不変のドレミ同定 ~言語化可能な相対音感~

図1に示す二つの曲(上曲を移調したものが下曲)をドレミに落とすよう依頼した場合、どのような反応が考えられるだろうか。返答は三通りある。「初めはソーミソドー、次がレーシレソー」と答えた場合、その人は絶対音感者であり、この場合ドレミは音名である。「両方ともソーミソドー」と答えたとすれば、その人は言語化可能な相対音感者であり、この場合ドレミは階名である。「ラーララーとしか歌えません」となった場合、その人は、言語化できない相対音感者である。

言語化可能な相対音感に着眼する。この場合、調を幾ら変えても(カラオケに行ってキーを上げ下げしても)「ソーミソドー、と聞こえてきます」と彼らは主張する<sup>(注2)</sup>。彼らは、何故、音高の異なる音を「ド、と内なる声が聞こえる」と主張する程に、その同一性を感覚するのだろうか?この認知プロセスの必要条件の一つとして、調不変の音階構造(音配置構造)がある[11]。

西洋音楽(平均律)では、1オクターブ( $\log(F_0)$ から $\log(2F_0)$ )に渡る音高帯域)を12個の音程に区分する(12半音)。 $\log(F_0)$ が第1音であれば、 $\log(2F_0)$ は第13音となる。長調と呼ばれる音階は、1オクターブを「全全全全全全」という音程に区分して8音を配置する。これが「ドレミファソラシド」である。上記音程が満たされさえすれば、各音の絶対的な音高には意味はない。個々の音には機能名があり、第1音=主音、第3音=中音、第5音=属音、などと呼ばれ、ドミソ、はそのニックネームである。これが階名の定義である。彼らはこの音の機能・価値を感覚して、ドレミが聞こえてくるのである。移調したところで音配置構造は不変であるため、ドレミ列は変わらない。長調の曲は、オクターブ等価性を前提にすれば、原則的に上記8音で構成されている。極端な場合を考えると、メロディーの中

(注2): 声を出さずに「ソーミソドー」と心の中でつぶやいた時と全く同一と思われる感覚・記憶が、無意識的に再生される、と言う主張である。

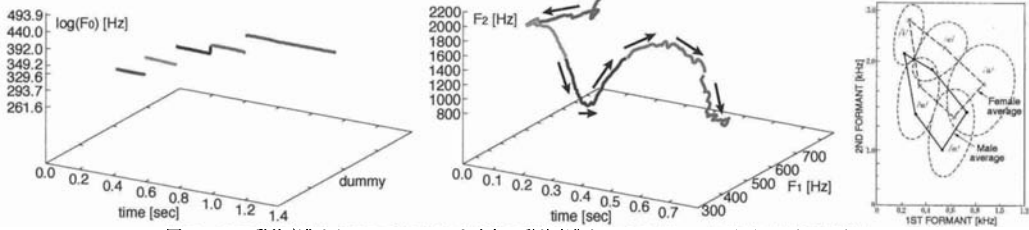


図3  $F_0$ の動的変化としてのCDEFGと音色の動的変化としての/aieuo/,及び、日本語母音図

の任意の2音が、三全音を音程（音高差）として持つ場合、その2音に対して「ファとシ」が聞こえてくる[12]。調不変の音高差異に基づいて要素音の同定を行うのが、言語化できる相対音感者である。彼らが超頑健な要素音同定を行なえるのは、個々の音の絶対的な物理特性など、記憶しないからである。

さて、この音配置構造が崩れるとどうなるのだろうか？古典的教会音楽には、種々の音階がある。図2のイオニア音階、エオリア音階が現代音楽の長調、短調として生き延びている。これらの音階では12半音の原則は守られており、5全音と2半音の配置の違いとなっている。さらに12半音の原則までも壊すとどうなるだろうか？図2にはアラビア音階も示している。12半音では表現できない音が要求されるため、通常のピアノでは再生できない。西洋音楽をアラビア音階で再生した場合、言語化できる相対音感者は「ドレミが聞こえてくるところと、聞こえて来ないところがある」という反応を示す。彼らの言語化は、音配置の様子に依存し、個々の音の音高には全く無関係に行なわれる。しかし逆に、孤立音の言語化は不可能である。メロディーという全体像があって初めて要素音のシンボル化が可能となる。シンボルを並べてメロディーが構成されるのではない。

## 2.2 音高の動的変化と音色の動的変化

主旋律（メロディー）のみを対象とすれば、音楽は  $F_0$ （ピッチ）の動的変化パターンである。音声として母音列のみを対象とすれば、下記に示す様に、これは音色の動的変化パターンである。母音の生成は声道（音響管）の共鳴現象であり、これは、管楽器における音生成と物理的には等価である。即ち「あいうえお」の違いは、声道形状の変化による共鳴現象の変化である。音楽学では音色はしばしば「基本音及び各倍音に対するエネルギー分布（配分）」として定義されるが、これはスペクトル包絡と同値である。結局、音色を表現するための最も簡素な物理パラメータはフォルマント周波数となり、ここでは  $F_1$  と  $F_2$  を考える（十次元のケプストラムを考えても下記の議論は成立する）。なお母音同様、複数の管楽器を  $F_1$ - $F_2$  平面上にプロットし、音色配置を示す場合もある[11]。図3に  $F_0$ の動的変化としてのCDEFG、及び音色の動的変化としての/aieuo/を示す。前者を移調しても、この動的パターンは上下に移動するだけであり、階名同定が要求する音群配置は不変である。一方/aieuo/の動的パターンであるが、日本語母音図（図3）に示すように、音響音声学では、 $F_1$ - $F_2$ 平面上で男声の母音構造を移動すると女声の母音構造に重なると言われる[13]。このような単純な写像で変換できれば、母音構造の話者不変性は容易に実現できるが（即ち二次元の移調＝平行移動）、厳密にはこのような単純な写

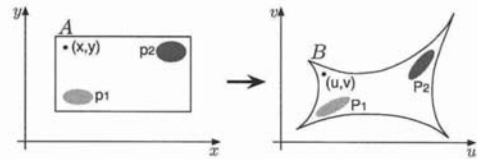


図4 一対一対応関係を有する二つの空間 A と B

像で変換できる訳では無い。音声合成の話者変換技術は、話者Aの音響空間と話者Bの音響空間との対応付け（写像）を精密に定義することで実装されるが[14]、音群構造の不変性は、この両空間における不変構造を要求する。逆に言えば、線形・非線形を問わずあらゆる写像関数に対して、不変なる構造が定義できれば、「音色の相対音感」は議論可能となる。なお、三角形は三辺の長さを規定すればその形状が一意に定まるように、 $N$ 角形の場合、全ての二点間距離（距離行列）を規定すれば、その形状は一意に定まる。即ち、不変なる構造は、不変なる差異（群）の存在を証明することで、立証されることになる。

## 3. 非言語的音響変動不変の音声の構造的表象

### 3.1 2つの空間における頑健な不変量

図4に示す様な、二つの空間AとBを考える。両者には一対一の対応関係があり、空間Aのある点は空間Bの対応点へ写像され、逆もまた成立する。但し、その写像関数は明示的には与えられていない。以下、一般性を失わない範囲で2次元空間を用いて説明する。空間AとBの間に一対一の対応があれば、空間Aの分布  $p_1$  は空間Bの分布  $P_1$  へと写像され、それを  $P_2$  とする。この時、次の等式が常に成立する[15]。

$$\iint_A \sqrt{p_1(x,y)p_2(x,y)} dx dy \equiv \iint_B \sqrt{P_1(u,v)P_2(u,v)} du dv$$

上式は、量子化学の世界では「重なり積分」と呼ばれる量であり<sup>(注3)</sup>、この量に対して  $-\log$  をとったものがパタチャリヤ距離（分布間距離の一つ）である。結局、パタチャリヤ距離は任意の二空間（話者）間で常に等しい。この距離（差異）不変性は、空間写像の種類に依らず、また、カルバックライブラ距離、ヘリンジャ距離でも成立する一般的性質である（頑健な不変性）。

### 3.2 不変事象間距離から普遍的に存在する不変構造へ

頑健に変換不変な距離尺度を用いて、ある発話を変換不変的に表象することを考える。図5に示すように、音声ストリーム

(注3)：この場合、分布は電子雲を指す。任意の二電子雲間の「重なり積分」を全て集めたのが「重なり行列」となる[16]。分子軌道法などで使われる。

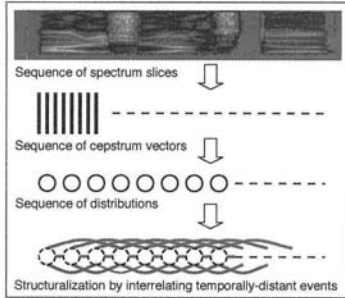


図5 音響事象間の差のみを抽出して構成される不変構造

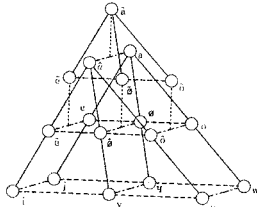


図6 ヤコブソンによるフランス語の母音・準母音構造 [17]

を分布系列へと変換した後に(系列長 =  $N$ )、時間的に離れているものも含め、全ての二分空間距離を求めて  $N \times N$  の距離行列として表象する。この時、個々の音響事象の絶対的な物理特性は全て捨象する。距離行列は一つの幾何学的構造を規定するが、この構造が変換不変となる。この構造は、例えば  $m + 1$  次元の音響パラメータ時空間に存在する音色の動的変化パターンを分布系列化し、各分布を  $m$  次元空間へと射影して得られる分布群が成す構造である。図6にヤコブソンによる仏語の母音・準母音構造を示す [17]。構造音韻論では、このような構造が話者に依らず観測されることを主張するが、筆者らが提唱する音声表象は構造音韻論の物理的・数学的解釈である。

#### 4. 音的実体を全く使用しない構造的音声認識

##### 4.1 連続母音系列発声をタスクとした音声認識

図5に示した、音声の音的実体を一切捨象した物理表象を用いた音声認識を検討した。日本語五母音を入れ替えて構成される連続母音系列発声(語彙数120であるため、 $PP=120$ の孤立単語認識となる)を対象語彙として検討した [18]。

図7にその枠組みを示す。入力音声を構造化し、統計的にモデル化された構造的テンプレートと照合する。この際図8に示す様に、片方の構造を回転及び平行移動して両構造を合わせた上で照合する。提案する構造的表象は変換不変性を有するため、任意の変換関数は、幾何学構造に対して回転或いは平行移動として作用する。例えば、声道長の差異(周波数ウォーピング)は構造の回転として、音響機器特性の差異(伝達関数の掛け算)は構造の平行移動として解釈される<sup>(注4)</sup>。回転&平行移動後の音響スコアは二つの距離行列を用いて計算されるが、これはタンパク質の構造解析などで用いられている手法と同一である。

(注4)：ケプストラム空間では、周波数ウォーピングは行列  $A$  の掛け算 [19]、伝達関数の掛け算はベクトル  $b$  の足し算となるため、最も簡単な話者変換は線形変換  $c' = Ac + b$  となる。この時、 $\times A$  が回転、 $+b$  が平行移動となる。

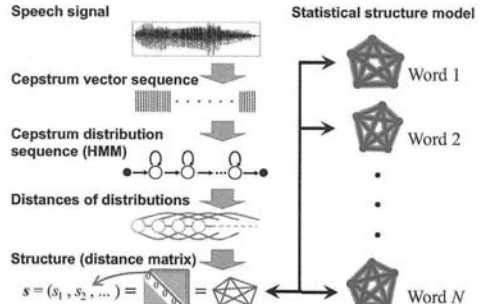


図7 音的実体を用いない構造的音声認識

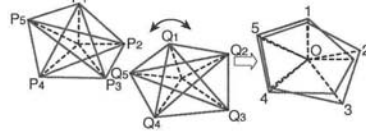


図8 回転及び平行移動を通して行なう音響照合

表1 音的実体を用いない構造的音声認識結果 [%]

|      | HMM(4,130) | HMM(260) | 提案手法 (8) |
|------|------------|----------|----------|
| 単語単位 | 97.4       | 82.1     | 96.6     |
| 母音単位 | 98.8       | 90.4     | 98.6     |

男女計8名に120単語を5回ずつ発声させ、これを用いて120単語の統計的構造モデルを作成した。これとは異なる男女8名に同様の発声を依頼し、評価データとした(合計4,800発声)。結果を表1に示す。学習話者260名、4,130名の不特定話者HMM+CMNの結果も示す。単語単位、母音単位の両性能において、HMM(4,130)とほぼ同等の性能を示している。スペクトル包絡など、音的実体に関する物理量を一切用いず、音色の動きのみを捉えることで、連続発声中の母音の約99%が非常に少ない学習話者数で同定できて「しまう」事実、甚だ驚嘆に値する。声に含まれる言語情報は、音的実体ではなく、音色の動きとして符号化されている、と解釈すべきであろう。

##### 4.2 音高に対する相対処理/音色に対する絶対処理

男女が同一歌詞の歌を歌った時、音高の動的パターンには絶対的な違いがある。男声は低く、女声は高い。これは男性の声帯が長く、重たいために声帯振動周期が長くなるためである。このような純粋に物理的な要因のために男女間の音高差は生じる。よって、両者の動的パターンの同一性を論じる場合、絶対的な音高知覚は役に立たない。極端な絶対音感者は、移調前後で曲の同一性認知が有意に遅れ [20]、困難になる場合もある。

その男女が同一歌詞を読み上げた場合、音色の動的パターンには絶対的な違いがある。男声は太く、女声は細い。これは男性の声道長が長いがために、共鳴周波数が低くなるためである。このような純粋に物理的な要因のために男女間の音色差は生じる。よって、両者の動的パターンの同一性を論じる場合、絶対的な音色知覚は役に立たない、と記したいところであるが、筆者らの知る限り、全ての音声科学・工学の議論は音色に対しては絶対的な処理系を常に指向・構築してきた。筆者らは、この両者の隔たりに強い不自然さ(恣意性)を感覚している。何故、音色に対しては絶対音感ばかりを議論してきたのだろうか？



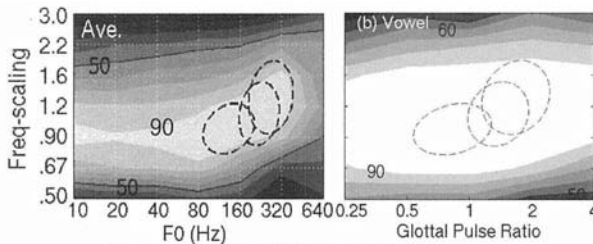


図9 孤立母音及び連続音声声中の母音同定 [21]~[23]

答えは簡単である。孤立音 [あ] を聞いて、それを音韻/あ/であると同定できるからである。これは完全な絶対音感であり、音楽の階名同定とは完全に異なる。この絶対音感を抛り所として、数十万人の音声から統計モデルを構築してきた。ならば聞いてみたい。「孤立音 [あ] を聞いて、音韻/あ/であると同定できる能力は、音声言語の運用に果たして必要なか？」と。

### 5. 母音は音名なのか、階名なのか？

図3に示す母音図から分かるように、日本語の場合、話者による違いを考えても、母音間の重なりはそれほど大きく無い。しかし、この重なりは容易に増加できる。フォルマント周波数は声道長の関数であるため、巨人/小人の声を合成すればよい。通常の領域から外れた孤立母音に対する同定は可能なのだろうか？もしそれが困難であった場合、音の連続ストリームの中にある母音はどのようなのだろうか？孤立母音の場合は困難であるにも拘らず、連続ストリーム中であれば容易である場合、これこそ、音色に関する階名認知として考えることができる。

先行研究にその答えを見ることができる [21]~[23]。図9左が孤立母音に対する同定率、右が無意味4モーラ単語中の母音同定率である<sup>(注5)</sup>。縦軸の値  $y$  に対して、 $170y[\text{cm}]$  が凡そ話者の身長となる。また、右図の横軸の値  $x$  に対して、 $160x[\text{Hz}]$  が基本周波数である。即ち、様々な身長・基本周波数の音声に対する、孤立母音の同定、及び無意味モーラ列中の母音同定の正解率である。図中点線の楕円が3つあるが、これは、実在する男性、女性、子供の領域を示す。全ての提示音声は STRAIGHT による分析合成音声である。孤立母音提示時（絶対的音認知時）は、実際に人間が存在する領域では90%を超えるが、それを越え始めると同定率は下がり、例えば  $65[\text{cm}]$  の小人となると、 $160[\text{Hz}]$  の音声で同定率は約20%となる。これはチャンスレベルであり、母音同定は全く不可能の状態になる。

一方、無意味連続モーラ列中に母音が置かれると、たとんに同定率が上昇する。 $65[\text{cm}]$  の小人ですら、約60%の正解率を呈する。提示単語が有意味語や親密語であれば、正解率は更に上昇するだろう。孤立音の同定はできないが、連続ストリームに対しては、個々の音事象を同定できる。これは、階名認知そのものである。再度聞いてみたい。孤立音を聞いて音韻同定できたとして、それは音声言語運用と関係あるのだろうか？

(注5)：厳密には、観密度データベース [24] における最低観密度単語群である。よって、音楽配列的には正しい日本語である。無意味語と記したのは、上記 DB の開発者が「未知語と考えて差し支えない」と言及しているからである。

言語化できない絶対音感者（ラウラ音感者）は次の要求に難儀する。「次に提示されるメロディーの三番目の音を覚えてください。その後、別のメロディーが提示されます。同一音が出てきたら手を上げてください」音のシンボル（音名/階名）化が出来なければ、この問いは困難である。同様に「次に提示される音声の三番目の音を覚えてください。その後別の音声提示されます。同一音が出てきたら手を上げてください」という問いに難儀するのが発達性ディスレクシアであり、欧米には数多く存在する。音声を音韻（音シンボル）列として認知することが困難であり、その結果、文字の読み書きに苦勞する。語ゲシュタルトに基づく認知プロセスを引きずり、個々の分節音をシンボル認知することが困難である [25]。米国では現在、教科書は音声 CD 添付が義務付けられている [26]。視覚障害を含め、読めない子供が数多く存在するからである。これらの事実を省みた時に、音声ストリームを音シンボル列として認知する能力、孤立音を音シンボルとして同定する能力は、そもそも、音声言語運用の必要条件なのだろうか？幼児にとって必要なのは、母親の「おはよう」と父親の「おはよう」に同一のコンテンツが乗って（符号化されて）いると認知する能力であり、それがどう視覚化されるのか、は楽しい朝食を囲む際に何ら必要ない。音高に対する極端な絶対音感を持つと、移調前後で曲の同一性認知が遅れる。同様に、音色に対する極端な絶対音感を持つと「おはよう」と「おはよう」の同一性認知が困難となるが、自閉症者の一部にその症状は観測される [27]。当然、音声言語（コミュニケーション）は成立しない。彼らの中には、音声模倣ではなく、声模倣を楽しむ者もいる [28]。当然音声言語は無い。

やはり、根本的に、間違っていたのだろうか？

### 6. 音色の動的変化パターンが示す話者依存性

第2.1節において「メロディーという全体像があつて初めて要素音のシンボル（階名）化が可能となる。シンボルを並べてメロディーが構成されるのではない」と書いた。前節の聴取実験は、「音声ストリームという全体像があつて初めて要素音のシンボル化が可能となる。シンボルを並べて音ストリームが構成されるのではない」ことを示唆する。全体が先にあるのか、要素が先にあるのか。言語音群を系（システム）として捉え、各音の（他音群との差異を通して定義される）相対的価値を議論するのが音韻論であり、個々の音を個別に観測し、その絶対的価値を議論するのが音声学である。となれば、(音響)音声学は果たして正しいのだろうか、という問いすら、生まれてくる。

系としての音群構造を考えた場合、声道長差異による周波数ウォーピングはケプストラム  $c$  に対する行列  $A$  の掛け算となる。そしてそれは、構造に対する回転演算子である。従来音声ストリームのダイナミクスを議論する場合、図3に示す、音色の動的変化パターンの各時刻における速度成分（方向成分とその大きさ）を考えていた。しかし上記の議論は、この方向成分は極めて強い声道長（年齢）依存性を持つことを示唆する。

身長約  $165[\text{cm}]$  の男性が発声した/aieuo/の [a] から [i] の遷移区間の中心時刻に着眼する。周波数ウォーピングを施して、

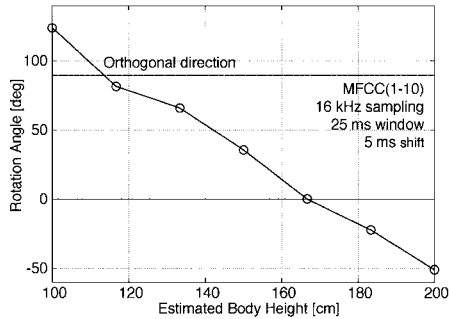


図 10 音色の動的変化パターンの方向成分が示す身長依存性

身長約 100[cm] から 200[cm] の音声を作成し、同一時刻の  $\Delta$  ケプストラムの方向成分の変化を分析した。図 10 は、身長約 165[cm] の方向成分に対する、各身長の方向成分が示す回転角である。身長が増減に伴い回転（やがて直交）するなど、極めて強い身長依存性を呈している。音声認識における音響モデリングは、より詳細な「音」モデリングを模索する傾向にあるが、それは「声」モデルであって、「音声」モデルからは逸脱する。そもそも HMM に代表される「音」の生成モデルは「声」モデルであり、「音声」モデルではない。本稿では、「声」の全体像を捉え、非言語的変形に不変な「音声」モデルを提案している。

## 7. 様々な議論と本稿のまとめ

第 2.1 節において、音階における音配置構造のバリエーションを示した。では、図 3 に示した母音配置構造に対するバリエーションを考えた場合、これは、何に対応するのだろうか？周知のように、これは欧米圏における方言である [29]。

幼児の音声模倣を思考実験と共に再考する。一卵性双生児を出産直後に親が離婚して、父親、母親が一人ずつ養育する場合を考える。10 年後、この双子の発音は（どれほど父親、母親のことを愛していたとしても）片方がより太く、他方がより細くなることは無いだろう<sup>(注6)</sup>。彼らは声（音）を模倣する訳ではない。10 年後彼らの発音は、一つの例外を除いて、非常に類似しているだろう。その例外とは、両親が異なる方言話者であった場合である。この場合、例えば apple の最初の母音 /æ/ は双子の間で異なることは容易に想像できる。同一方言話者の男女の /æ/ の違いは、共鳴周波数の違いである。異なる方言話者の男女の /æ/ の違いも、共鳴周波数の違いである。前者は発音に影響せず、後者は影響する。何故か？結局「幼児が模倣するのは音ではなく、音群の体系である」との説明が最も妥当かつ簡潔である。もし、両方の家庭で九官鳥が飼育されていれば、彼等は「音」を模倣するため、などの議論はもはや不要だろう。

九官鳥は提示された「声」から何を学び、何を模倣するのか？一方幼児は、提示された「声」から何を学び、何を模倣するのか？そして両者の違いは何なのか？これを考えた場合、「音」の生成モデルとしての音響モデルは、一つの例外を除いて、甚だ見当違いの議論であると言わざるを得ない。その例外とは、音声認識研究の目指すゴールが九官鳥シミュレータの構築である

(注6)：但し、発話スタイルに相違が生じることは考えられる。

場合である。この場合は、より詳細な音モデリング技術を目指すべきである。筆者らが目指すのは幼児シミュレータの構築である。彼らが模倣しない側面は積極的にそぎ落とす必要がある。

本稿をここまで読まれた読者に対して一言言いたい。本論を「新しい音声認識論」と感じたのなら、それは全くの誤解である。本論は極めて古典的な音声認識論に過ぎない。近代言語学の祖ソシュールによる一世紀以上も昔の言を示したい。“The important thing in the word is not the sound alone but the phonic differences that make it possible to distinguish this word from all others [30].” 問題は、多くの（恐らく全ての）音声工学者が、波形やスペクトルといった音声の音的実体、即ち「声」に目・耳を奪われて、彼らの言動を物理的、数学的に探求することを怠った点にあると筆者らは考えている。

## 文 献

- [1] <http://tepia.or.jp/archive/12th/pdf/viavoice.pdf>
- [2] T. Anastasakos *et al.*, Proc. ICSLP, vol.2, pp.1137-1140 (1996)
- [3] <http://hts.sp.nitech.ac.jp/>
- [4] 早川, 月刊言語, 35, 9, pp.62-67 (2006)
- [5] W. Gruhn, “The audio-visual system in sound perception and learning of language and music,” Int. Conf. on Language and Music as Cognitive Systems (2007)
- [6] 宮本, 音を作る・音を見る, 森北出版 (1995)
- [7] 原, コミュニケーション障害学, 20, 2, pp.98-102 (2003)
- [8] 加藤, コミュニケーション障害学, 20, 2, pp.84-85 (2003)
- [9] N. Minematsu *et al.*, “Universal and invariant representation of speech,” Proc. Int. Conf. Infant Study (2006) <http://www.gavo.t.u-tokyo.ac.jp/~mine/paper/PDF/2006/ICIS.t2006-6.pdf>
- [10] D. J. Levitin *et al.*, Trends in Cognitive Sciences, vol.9, no.1, pp.26-33 (2005)
- [11] 谷口, 音は心の中で音楽になる, 北大路書房 (2003)
- [12] 東川, 読譜力ー「移動ド」教育システムに学ぶ, 春秋社 (2005)
- [13] 古井, デジタル音声処理, 東海大学出版会 (1985)
- [14] 高橋他, 信学技報, SP-2006-162, pp.13-18 (2007)
- [15] 峯松他, 春音講論, 1-P-12, pp.147-148 (2007)
- [16] 武次他, 早わかり分子起動法, 裳華房 (2003)
- [17] R. Jakobson *et al.*, Notes on the French phonemic pattern, Hunter, N.Y. (1949)
- [18] S. Asakawa *et al.*, “Automatic recognition of connected vowels only using speaker-invariant representation of speech dynamics,” Proc. InterSpeech (2007, accepted)
- [19] M. Pitz, *et al.*, IEEE Trans. Speech and Audio Processing, 13, 5, pp.930-944 (2005)
- [20] 宮崎, 日本音響学会誌, vol.60, no.11, pp.682-688 (2004)
- [21] D. Smith *et al.*, J. Acoust. Soc. Am., 117(1), pp.305-318 (2005)
- [22] 青木他, 秋音講論, 2-P-6, pp.373-374 (2004)
- [23] 林他, 春音講論, 2-Q-27, pp.473-474 (2007)
- [24] 天野他, 日本語の語彙特性, 三省堂 (2000)
- [25] S. Shaywitz, 読み書き障害 (ディスレクシア) のすべて～頭はいいのに本が読めない～, PHP 研究所 (2006)
- [26] 河村, “DAISY を活用したディスレクシアの方々への支援”, 日本障害者リハビリテーション協会セミナー「ディスレクシアの支援・デンマークでの活動から」より (2006)
- [27] 東田他, この地球にすんでいる僕の仲間たちへ, エスコアール出版社 (2005)
- [28] 深見, ひろしくんの本, vol.5, 中川書店 (2006)
- [29] W. Labov *et al.*, Atlas of North American English, Mouton and Gruyter (2005)
- [30] F. D. Saussure, Course in general linguistics, McGraw-Hill Humanities/Social Sciences/Langua (1965)