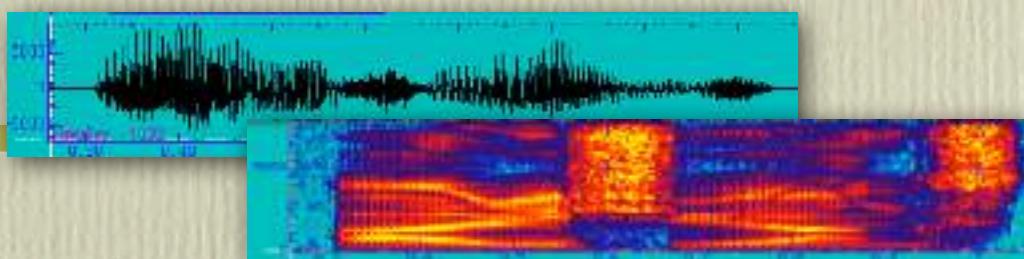


音響音声学

(Topics in Acoustic Phonetics)



峯松 信明

工学系研究科電気系工学専攻

教科書のコラム記事より

コーヒーブレイク

スペクトル包絡に安住していてよいのだろうか？

音声認識にしろ、音声合成にしろ、スペクトル包絡が基本的な音声特徴量として用いられている。人間の聴覚は音声の位相成分に鈍感との知見より、位相スペクトルを削除して振幅スペクトルを抽出し、さらに、音素情報は音声の音高成分とは独立であるため、ピッチハーモニクスを削除してスペクトル包絡特性を抽出している。しかし、包絡特性には音素情報と話者情報とが同居している。音声認識は音声中の音素情報（テキスト情報）を抽出することが目的だが、スペクトル包絡から話者情報を削除した特徴量を定義することはできないのだろうか。

不特性話者音響モデルは話者独立モデルとも呼ばれるが、この独立性は話者性を削除して得られるのではなく、多数の話者からデータを集めることで、話者性を分布の中に隠すことで得られるモデルである。位相やピッチは物理的に削除して独立性を担保するが、話者性は隠すことで独立性を担保している。

幼児の言語獲得は他者の音声活動を模倣することが必要となる（音声模倣）。この場合、話者性までを模倣しようとはしない。声帯模写はしない。話者の違いを超えた模倣をしている。音声模倣をする動物は小鳥、クジラ、イルカなどがあるが、彼らは基本的に音響的な模倣をする（なお、動物は相対音感を持っていないため、移調前後のメロディの同一性の認識が難しい。異なる音は異なるものと

教科書のコラム記事より

して扱っているのだろう)。動物とは異なり、人が他者の発声をコピーして言葉を獲得するとき、話者性には鈍感なのである。しかし技術はそうなっていない。ある話者の音声サンプルを用いて音声合成装置をつくれば、当然その人の声を出力する装置ができる。機能的には声帯模写装置と呼ぶべきである。

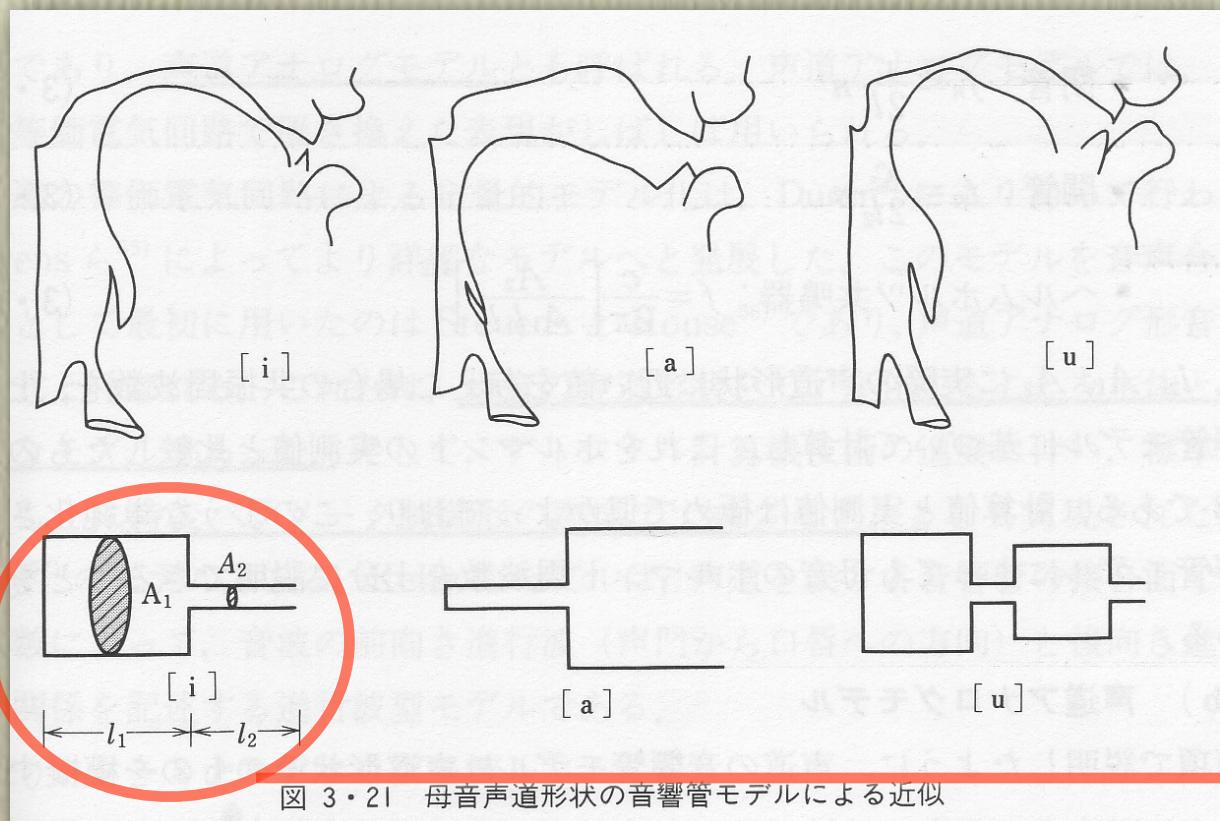
近年、音声工学の分野でもコンテスト形式の研究発表会が多くなってきた。音声合成の場合、blizzard challenge と呼ばれるワークショップが毎年行われている。ここでは合成音の品質を評価する際に、言葉としての自然さ以外に、学習話者の個�性が適切に再現されているか否かも評価対象となる。やはり今の音声合成技術は、声帯模写技術と呼んだほうが適切のように思われる。

言葉を真似る模倣行為が声帯模写的になってしまう場合がある。発達障害の一種、自閉症の中で見られる症状である。このような場合、音声言語の獲得はしばしば難しくなる。音の獲得 ≠ 言語の獲得なのだろう。人間と動物に見られる音情報処理の差異、典型的な言語発達が容易 / 困難な場合に見られる音情報処理の差異を概観すると、話者性を削除して音声を表現する技術の構築が待たれる。確率論は、集めてしまえば消したい要因が消せることを保証するが ($P(a) = \sum_b P(a, b)$)、これはあまりにもナイーブすぎないか。話者性をそぎ落とす一手法として**音声の構造的表象**が提案されている。興味のある読者は文献35) を参照されたい。

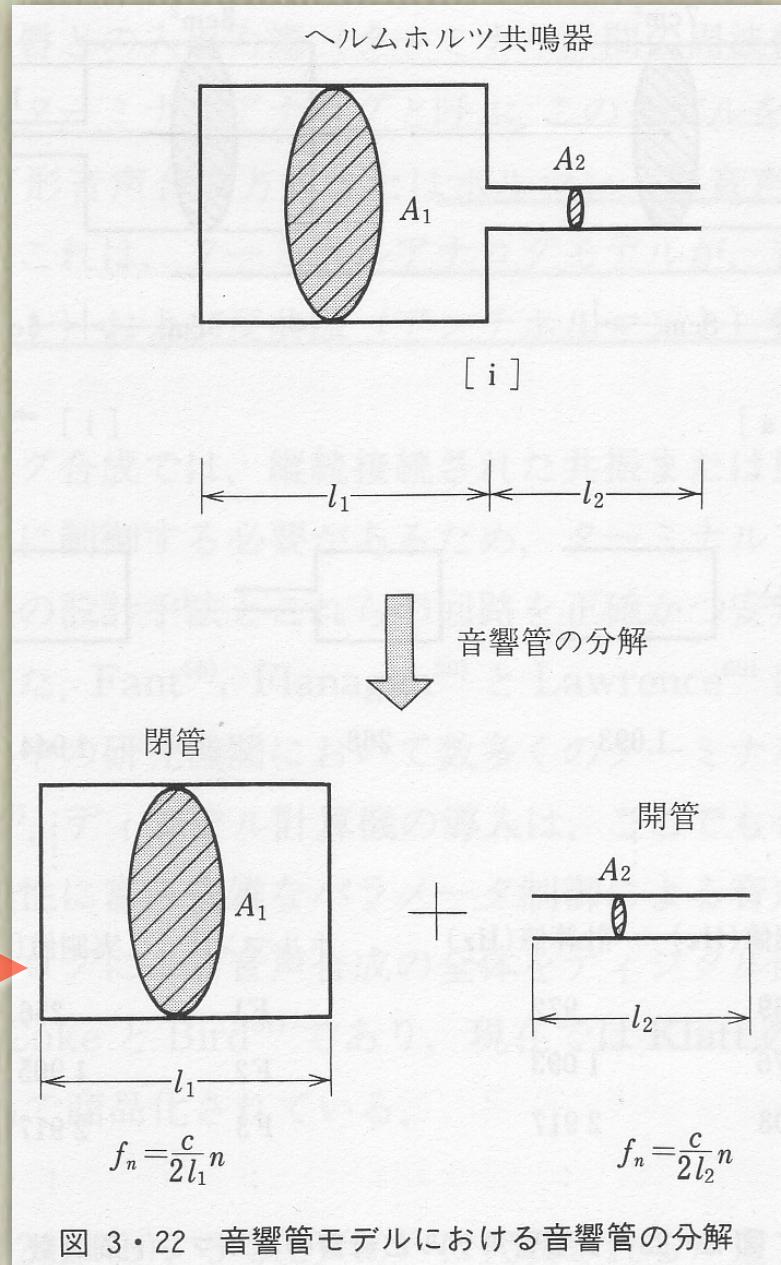
母音=定常波 ~気柱の共鳴現象~

複雑な管になつても原理は同じ

● 定常波の共振周波数を求めて



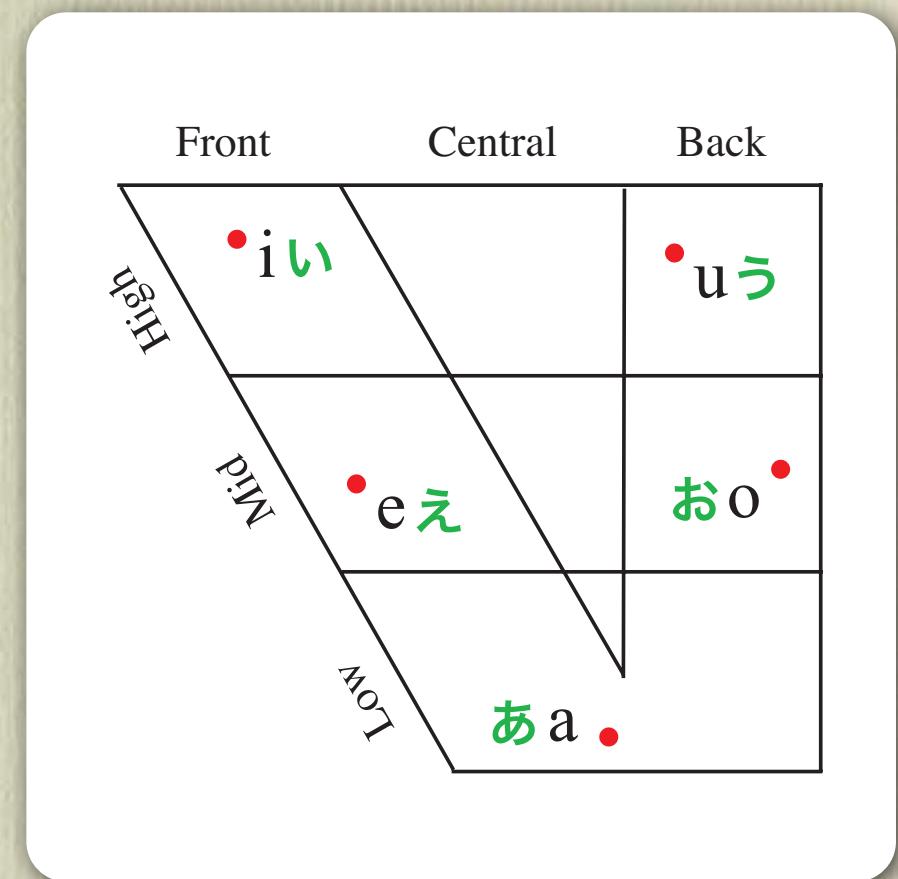
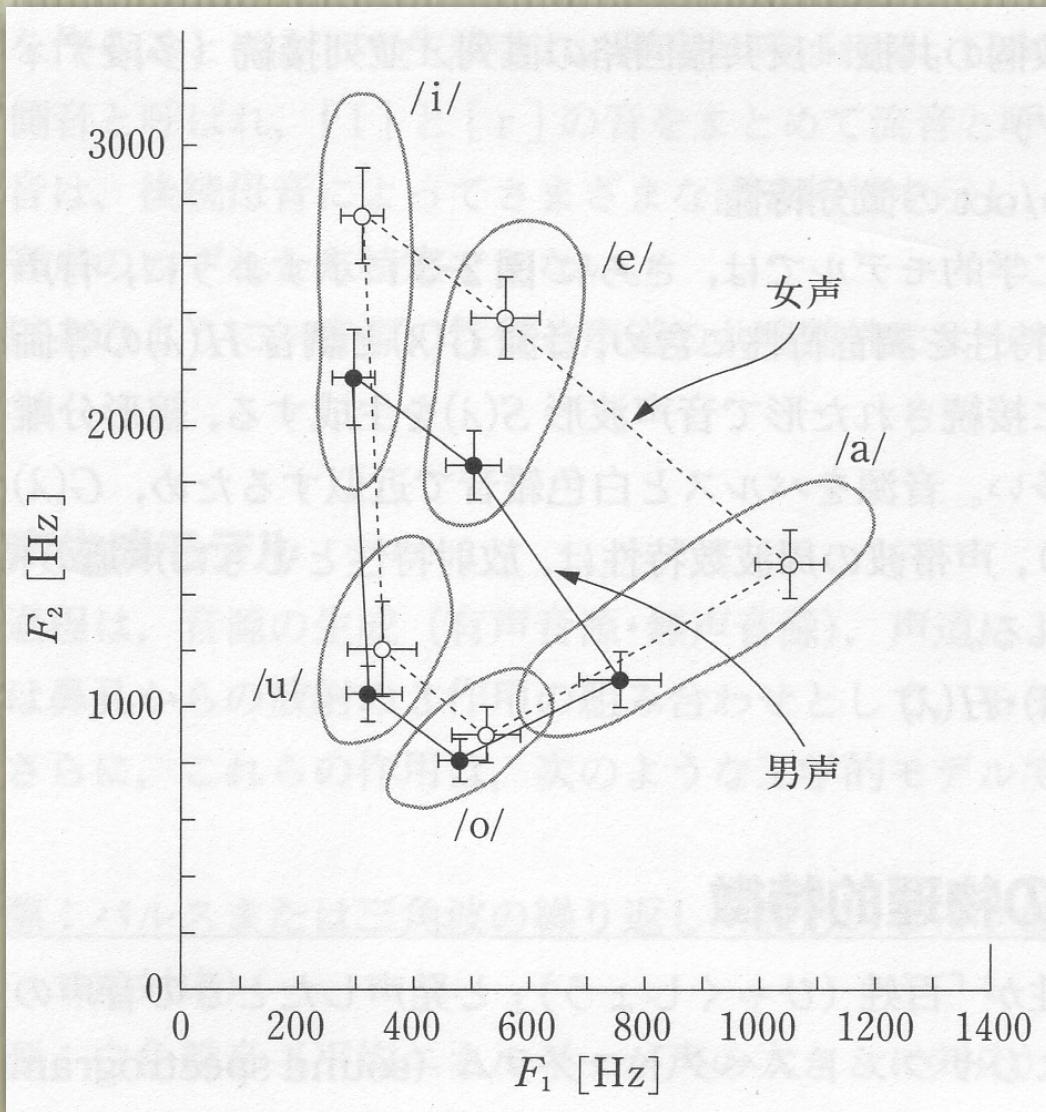
$$f_n = \frac{c}{2l_1} n \quad f_n = \frac{c}{2l_2} n \quad f = \frac{c}{2\pi} \left[\frac{A_2}{A_1 l_1 l_2} \right]^{1/2}$$



話者間における母音の差異

形の違い=長さの違い=共振周波数の違い

「あ」の一部=「お」の一部



音声が運ぶ様々な情報

言語的情報

- 何を話したのか？
- 狭義の言語的情報、語彙、音素
- どのように話したのか？
- パラ言語的情報、意図、発話スタイル、感情



非言語的情報

- 意図的制御は困難であり、不可避的に付与されてしまう情報
- 話者性、年齢、性別、体格、健康状態
- マイク、伝送特性、部屋の音響特性

音声信号：一次元の数値列

- 多様な情報を適切に反映しつつ数値列を生成：音声合成
- 数値列から様々な情報を的確に抽出：音声認識・理解
- 計算機に音声コミュニケーション能力を与えるためには？

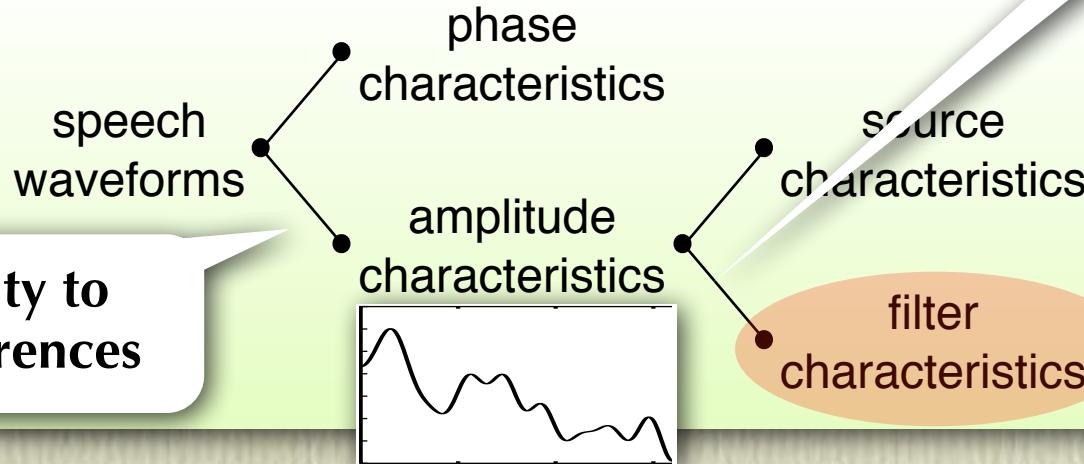
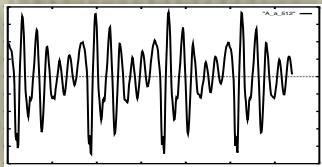
人間のような



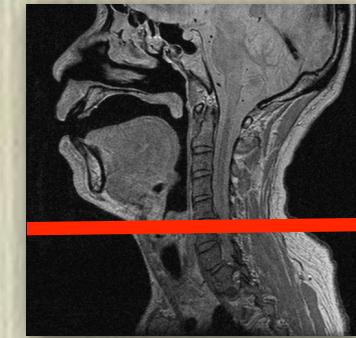
その情報を運ぶ媒体・音響特徴量

二段階の分離に基づく特徴量抽出

Independence bet.
phonemes and pitch



Insensitivity to
phase differences



● スペクトル包絡(σ)は何を運ぶのか？

言・パラ言・非言

● σ の中のある特定の情報のみに着眼したい。

- 当該情報に対応しない特徴量を揃える
- 当該情報に対応しないモデルパラメータを調節
- 確率定義に従って着目しない情報を分布に隠す
- ラベル情報を使って識別的な特徴量へ変換
- 当該情報に直接対応する特徴量を探求する

特徴量正規化

モデル適応

統計的モデル

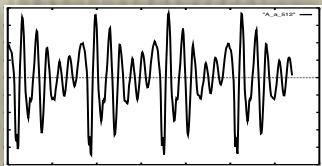
識別的変換

不变量

その情報を運ぶ媒体・音響特徴量

二段階の分離に基づく特徴量抽出

Independence bet.
phonemes and pitch



speech
waveforms

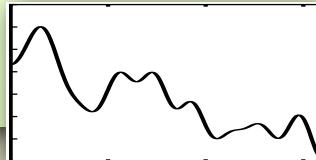
phase
characteristics

amplitude
characteristics

source
characteristics

filter
characteristics

Insensitivity to
phase differences



● スペクトル包絡(o)は何を運ぶのか？

言・パラ言・非言

● 二つの音響モデル $P(o|w)$ と $P(o|s)$

$s = \text{speaker}$
 $w = \text{word}$

● 不特定話者の単語音響モデル

$$P(o|w) = \sum_s P(o, s|w) = \sum_s P(o|w, s)P(s|w) \sim \sum_s \underline{P(o|w, s)}P(s)$$

● テキスト非依存の話者音響モデル

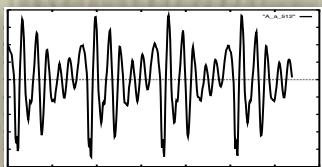
$$P(o|s) = \sum_w P(o, w|s) = \sum_w P(o|w, s)P(w|s) \sim \sum_w \underline{P(o|w, s)}P(w)$$

● 集めてしまえば「確率の定義」が見たくないものを隠してくれる。

その情報を運ぶ媒体・音響特徴量

二段階の分離に基づく特徴量抽出

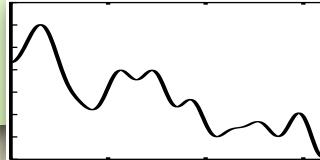
Independence bet.
phonemes and pitch



speech
waveforms

phase
characteristics

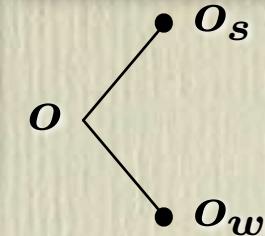
amplitude
characteristics



Insensitivity to
phase differences

source
characteristics

filter
characteristics



スペクトル包絡(o)は何を運ぶのか？

言・パラ言・非言

真の音声の統計的モデル～波形の統計的モデル～

不特定話者・不特定基本周波数・不特定位相の音響モデル

見たくないものは全て「確率の定義」で集めて隠してしまおう。

$$P(o|w) \approx \sum_{s,h,p} P(o|w, s, h, p) P(s) P(h) P(p)$$

s : speaker, h : harmonics, p : phase

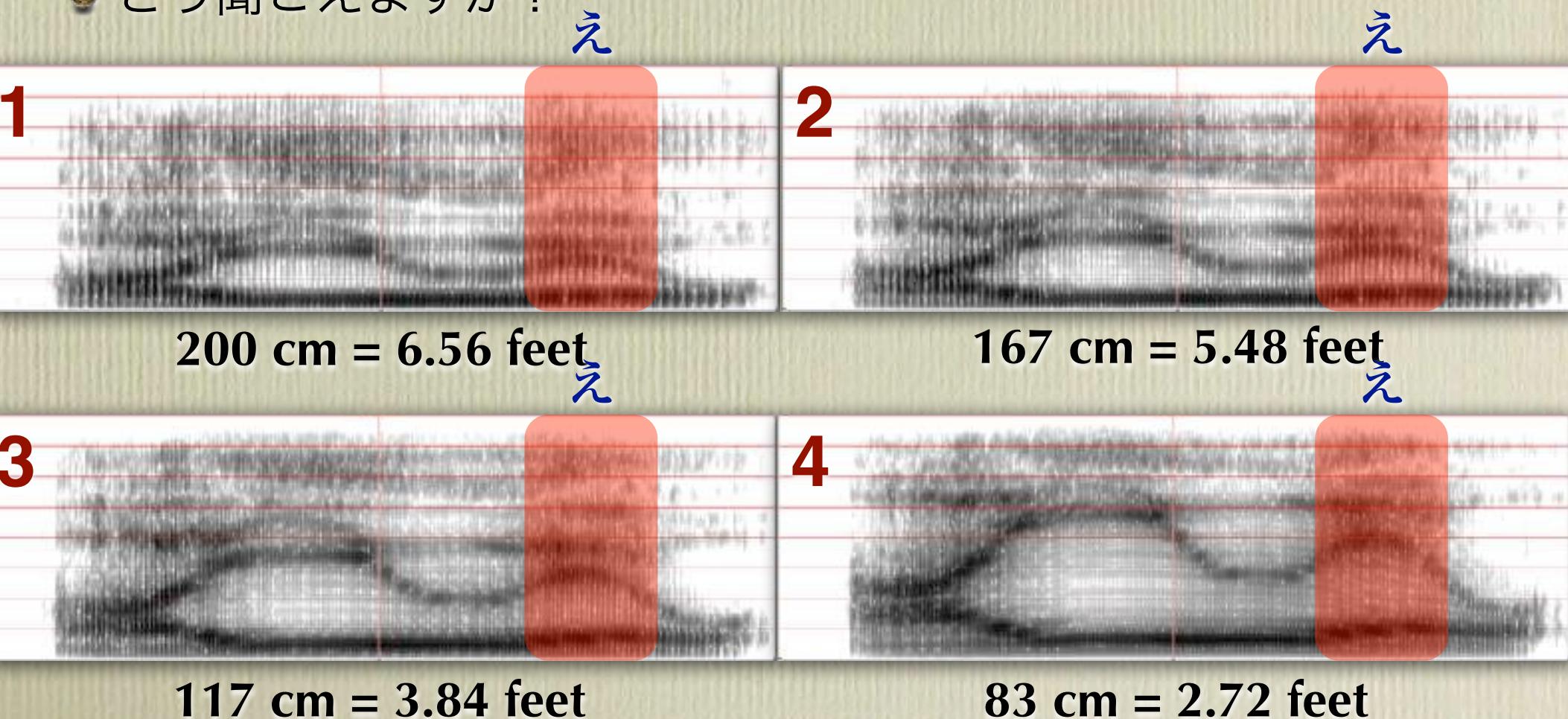
一般的な解決策：各手法の組み合わせ

最終的に性能を最大化する組み合わせを追求する。

音声物理の多様性と音声知覚の不变性

身長（喉の長さ）の違いと声の違い

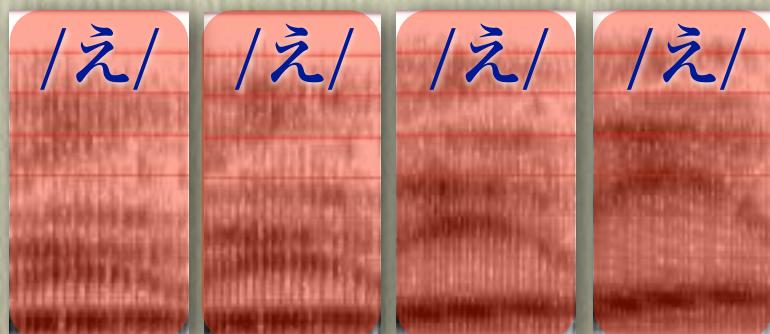
- 分析合成と呼ばれる技術 (STRAIGHT) を用いて音声を変形
- いろんな身長の男声を生成／但しオリジナルは167cm
- どう聞こえますか？



音声工学（科学）のアプローチ

音声認識 = 音声→テキスト変換

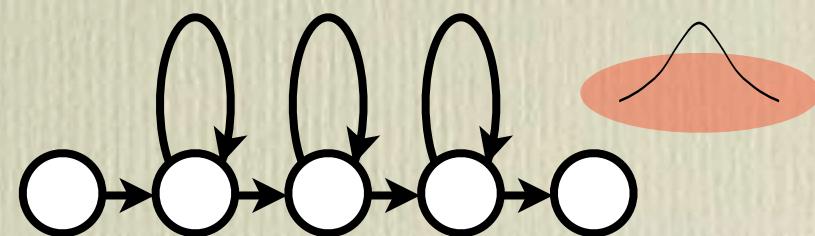
- しかし、個々の音韻の音的実体は様々な音となる
- 性別、年齢、マイク、部屋、伝送系などなど
- IBM の偉業：35万人の音声を収録



.....



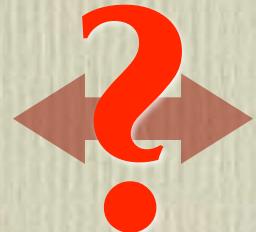
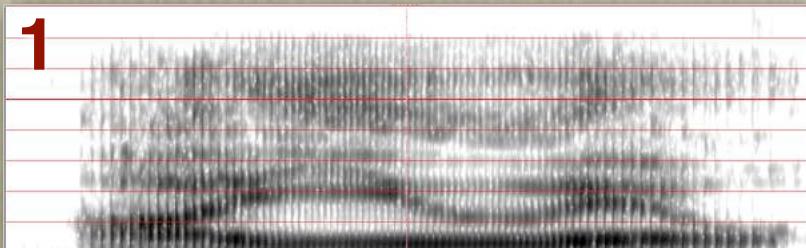
HMM of /え/



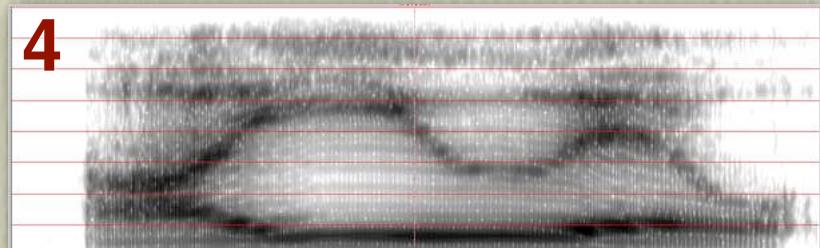
$$P(o|w) = \sum_s P(o, s|w) = \sum_s P(o|w, s)P(s|w) \sim \sum_s P(o|w, s)P(s)$$

この両者の一体何が「同一」なのか？

1



4



音を集めることは本当に必要なのか？

幼児の音響音声的環境 ~Another PoS~

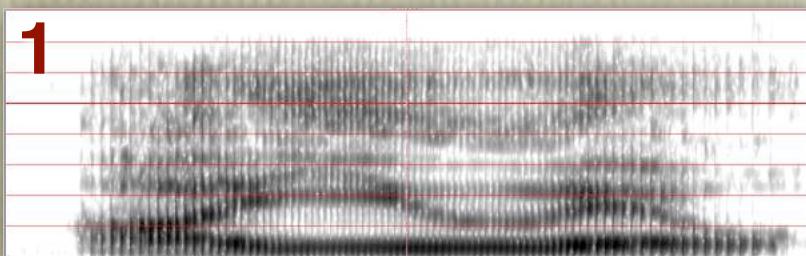
- 大部分は母親と父親の音声（日本の場合は母親ばかり？）
- 話し出せば、人の聞く声は半分は自分の声
- 極端に偏った音声コーパスに基づく超頑健な音声情報処理



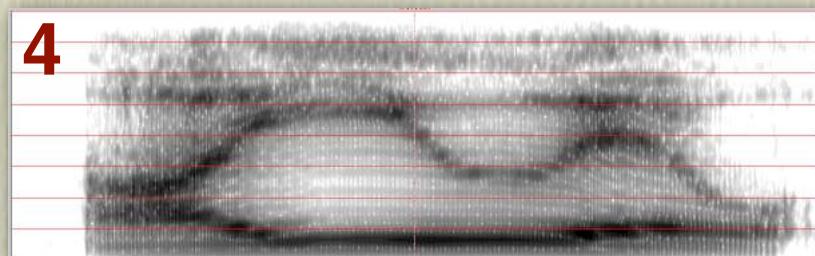
音声物理の多様性と音声知覚の不变性

- 集めずに知覚できる同一性とは一体何なのか？

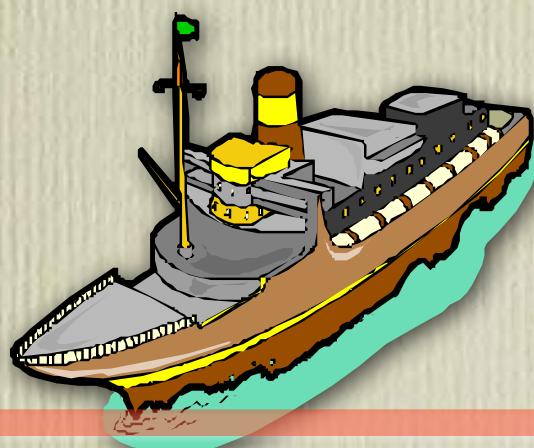
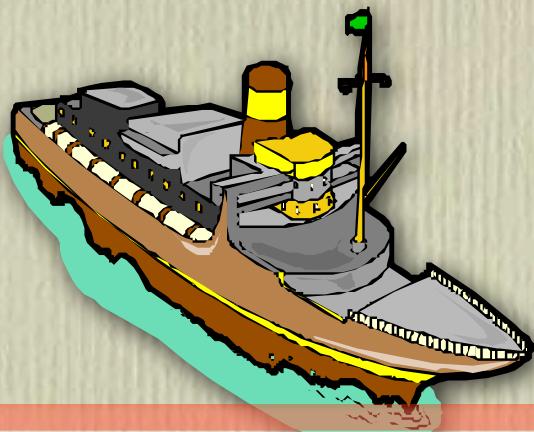
1



4



意地悪な思考実験



子育て奮闘中のお母様方

20/20

音声認識研究者

0/N

幼児は親の声の何を真似ているのか？

音



言



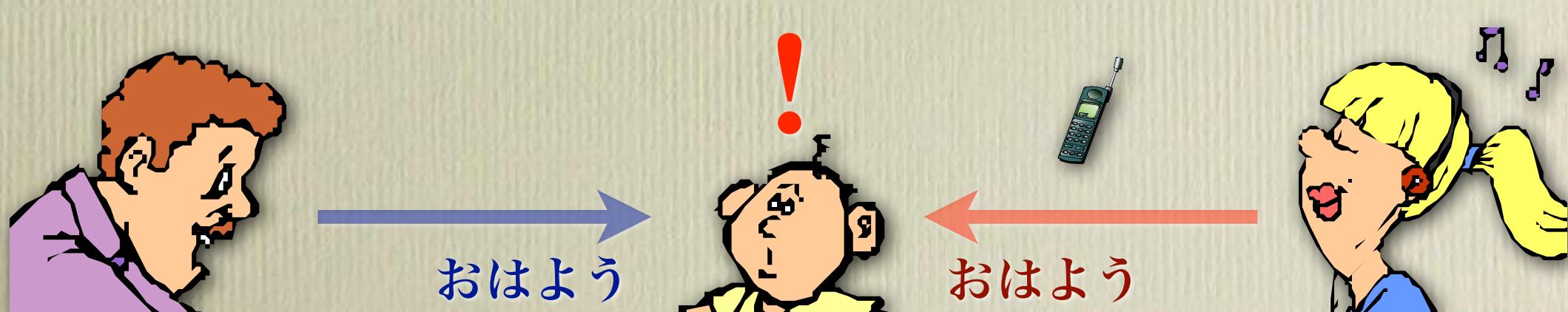
要素から？全体から？

発達心理学の主張 ~全体から入る音声処理~

- 音シンボル（音韻）の意識の定着は小学校入学以降
- 単語音声を要素に分割することが困難でも、音声活動を始める
- 「子供は語全体の音形・ゲシュタルトを獲得してから、語を構成する個々の分節音の獲得へと進む（加藤'03, 早川'06）」
- 語（発話）全体の音形を音にしたもののが音声

幼児の音声音響的環境 ~Another PoS~

- 大部分は母親と父親の音声／人の聞く声は半分は自分の声
- 極端に偏った音声コーパスに基づく超頑健な音声情報処理



?二つの問題?

音声物理の多様性と音声知覚の不变性→another PoS

- [え] は年齢・性別・マイクなどの非言語的要因によって変化
 - 音声認識 (IBM ViaVoice) = 35万人の話者から [え] を集める
 - 幼児 = 大半は母親／父親, そして, 自分の声



九官鳥の音声模倣と幼児の音声模倣

- 九官鳥は声（音）を真似る
 - 良い九官鳥は聞けば飼い主が分かる。「声そのもの」が模倣対象
- 幼児は親の声の何を真似ているのか?
 - 声を真似ようとはしていない。電話声になるようになんかしない。
 - 個々の音をモーラ同定して、それを一つずつ再生する？ それは困難
 - 語全体の音形・語ゲシュタルトをまず獲得。分節音はその後。

200cm

え

80cm

え

発声全体を通して定義できる話者不变な音声の物理表象

ある種の違和感

音声言語工学の構築してきた技術

- 波形素片／スペクトル素片に対するDBの構築
- 波形接続型音声合成／HMM音声合成（**学習用話者の声の合成**）
- 数理統計的手法に基づく音響モデルのパラメータ推定
　　● 数千～数十万人の音声を用いた音響モデルによる不特定話者音声認識
- 圧倒的な計算速度の向上と**大規模クラスター**の構成
- 様々な環境別に構築したシステム群によるパラレルデコーディング



ふと、我が子を見てみる・・・



- 母親の声を一番よく聞いて日本語を獲得したはず・・・
- でも、**母親の声の模倣（声帯模写）** なんて一度もしていない。
- この子の聞く声の多くは、**母親、自分、そして、父親**・・・
- 昨晩、お婆ちゃんに初めて声を聞かせた。電話で。で、会話してた。
- 音声って、**非常にモロい物理現象なんだよな**・・・
- でも、何故か、そのメディアを使うのが**一番楽なんだよな**・・・

教科書のコラム記事より

コーヒーブレイク

スペクトル包絡に安住していてよいのだろうか？

音声認識にしろ、音声合成にしろ、スペクトル包絡が基本的な音声特徴量として用いられている。人間の聴覚は音声の位相成分に鈍感との知見より、位相スペクトルを削除して振幅スペクトルを抽出し、さらに、音素情報は音声の音高成分とは独立であるため、ピッチハーモニクスを削除してスペクトル包絡特性を抽出している。しかし、包絡特性には音素情報と話者情報とが同居している。音声認識は音声中の音素情報（テキスト情報）を抽出することが目的だが、スペクトル包絡から話者情報を削除した特徴量を定義することはできないのだろうか。

不特性話者音響モデルは話者独立モデルとも呼ばれるが、この独立性は話者性を削除して得られるのではなく、多数の話者からデータを集めることで、話者性を分布の中に隠すことで得られるモデルである。位相やピッチは物理的に削除して独立性を担保するが、話者性は隠すことで独立性を担保している。

幼児の言語獲得は他者の音声活動を模倣することが必要となる（音声模倣）。この場合、話者性までを模倣しようとはしない。声帯模写はしない。話者の違いを超えた模倣をしている。音声模倣をする動物は小鳥、クジラ、イルカなどがあるが、彼らは基本的に音響的な模倣をする（なお、動物は相対音感を持っていないため、移調前後のメロディの同一性の認識が難しい。異なる音は異なるものと

教科書のコラム記事より

して扱っているのだろう)。動物とは異なり、人が他者の発声をコピーして言葉を獲得するとき、話者性には鈍感なのである。しかし技術はそうなっていない。ある話者の音声サンプルを用いて音声合成装置をつくれば、当然その人の声を出力する装置ができる。機能的には声帯模写装置と呼ぶべきである。

近年、音声工学の分野でもコンテスト形式の研究発表会が多くなってきた。音声合成の場合、blizzard challenge と呼ばれるワークショップが毎年行われている。ここでは合成音の品質を評価する際に、言葉としての自然さ以外に、学習話者の個�性が適切に再現されているか否かも評価対象となる。やはり今の音声合成技術は、声帯模写技術と呼んだほうが適切のように思われる。

言葉を真似る模倣行為が声帯模写的になってしまう場合がある。発達障害の一種、自閉症の中で見られる症状である。このような場合、音声言語の獲得はしばしば難しくなる。音の獲得 ≠ 言語の獲得なのだろう。人間と動物に見られる音情報処理の差異、典型的な言語発達が容易 / 困難な場合に見られる音情報処理の差異を概観すると、話者性を削除して音声を表現する技術の構築が待たれる。確率論は、集めてしまえば消したい要因が消せることを保証するが ($P(a) = \sum_b P(a, b)$)、これはあまりにもナイーブすぎないか。話者性をそぎ落とす一手法として**音声の構造的表象**が提案されている。興味のある読者は文献35) を参照されたい。

本発表の流れ

● 刺激の物理的多様性とその認知的不变性

- 見え／色み／音高の多様性と自然・進化が編み出した解決方法

● 音声の物理的多様性とその認知的不变性

- 音色の多様性と工学者が編み出した解決方法

● 音声の構造的表象とそれに関する様々な考察

- 常識を覆すことで、違和感の解消を試みてみる。

● 音声の構造的表象と数学的表现と技術的実装

- 体格・性別に不变な音声波形・スペクトルの表現とは？

● 音声の構造的表象を用いた音声アプリケーション

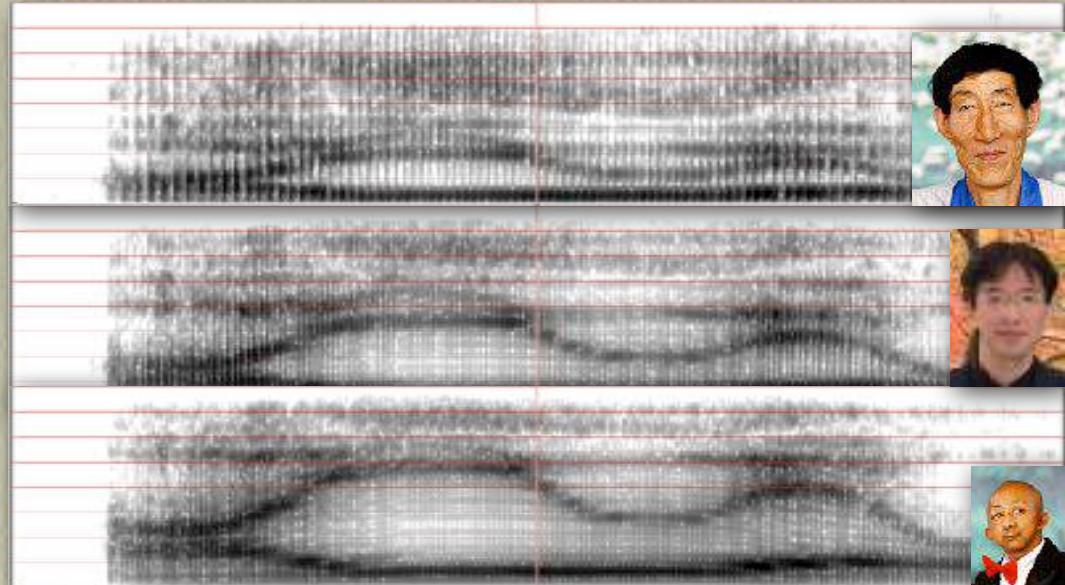
- 音声認識、音声合成、発音分析、etc

● 音声の構造的表象の言語学的妥当性

- 何故、こうしてこなかったのか？　観測技術の功罪？

年齢・性別・体格による音声の変形

巨人と小人の会話は、何故成立するのか？



刺激の物理的多様性とその認知的不变性

感覚受容器が受け取る情報は容易に変貌する

見えの変化

- 視点を変えて見た犬
- 対象との距離を変えて見た像



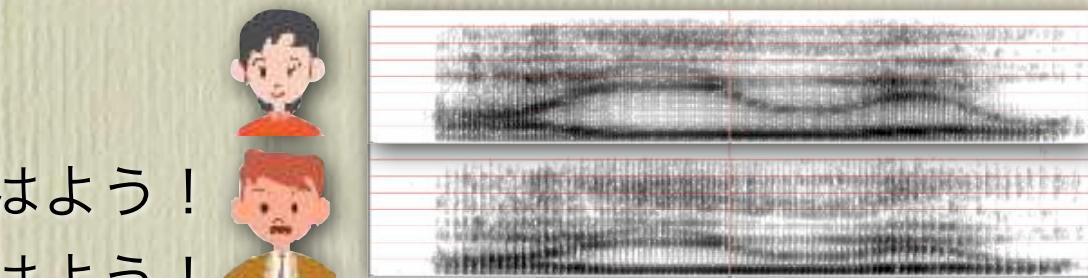
色みの変化

- 朝日の花と夕焼け空の花
- 異なる色眼鏡を通して見た像



音高の変化

- 男性のハミングと女性のハミング
- カラオケでのキーの上げ下げ



音色の変化

- 男性のおはよう！と女性のおはよう！
- 大人のおはよう！と子供のおはよう！

でも、我々は容易に「同一性」を認知できる

刺激の物理的多様性とその認知的不变性

感覚受容器が受け取る情報は容易に変貌する

見えの変化

- 視点を変えて見た犬
- 対象との距離を変えて見た像



色みの変化

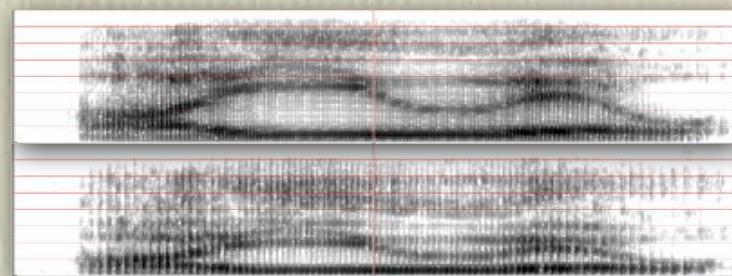
朝日の赤い夕焼け空の赤

静的偏差による刺激変形と 刺激変形に不变な認知様式

カラオケでのキーの上げ下げ

音色の変化

- 男性のおはよう！と女性のおはよう！

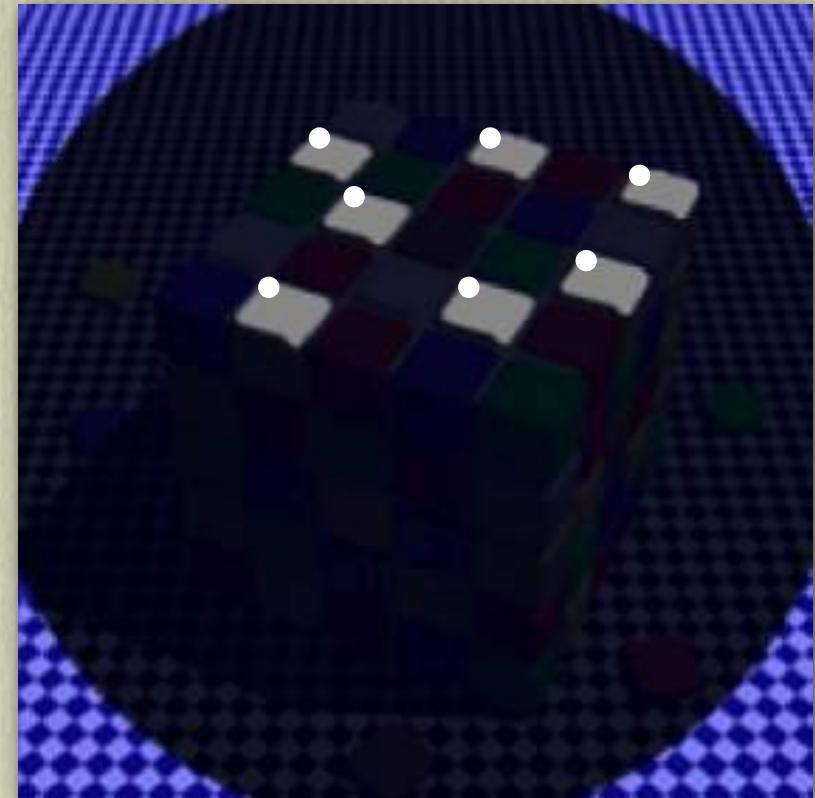
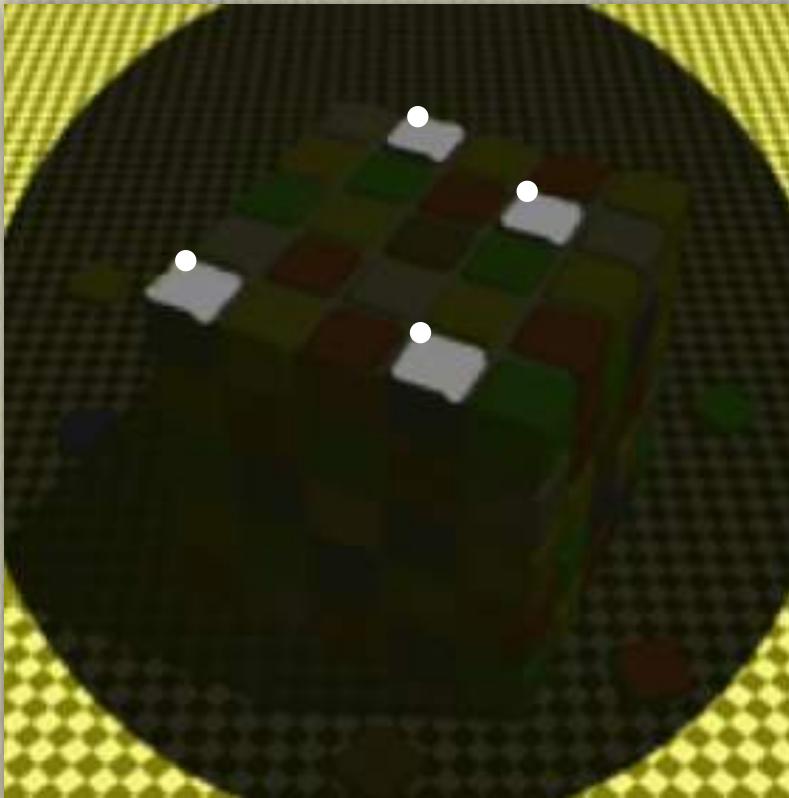


- 大人のおはよう！と子供のおはよう！

でも、我々は容易に「同一性」を認知できる

色みの偏差とその認知的不变性

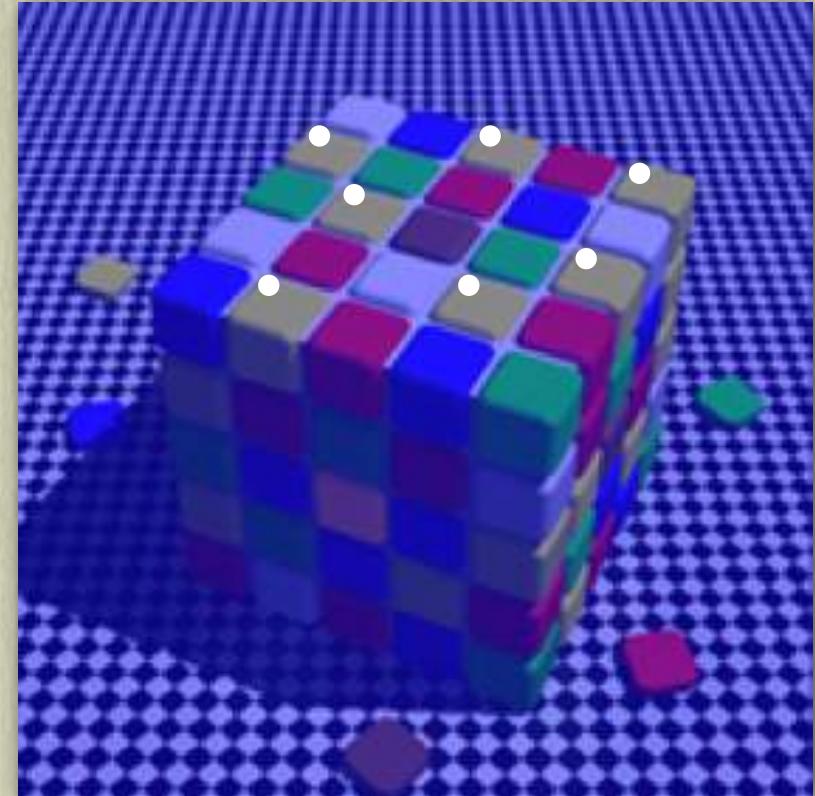
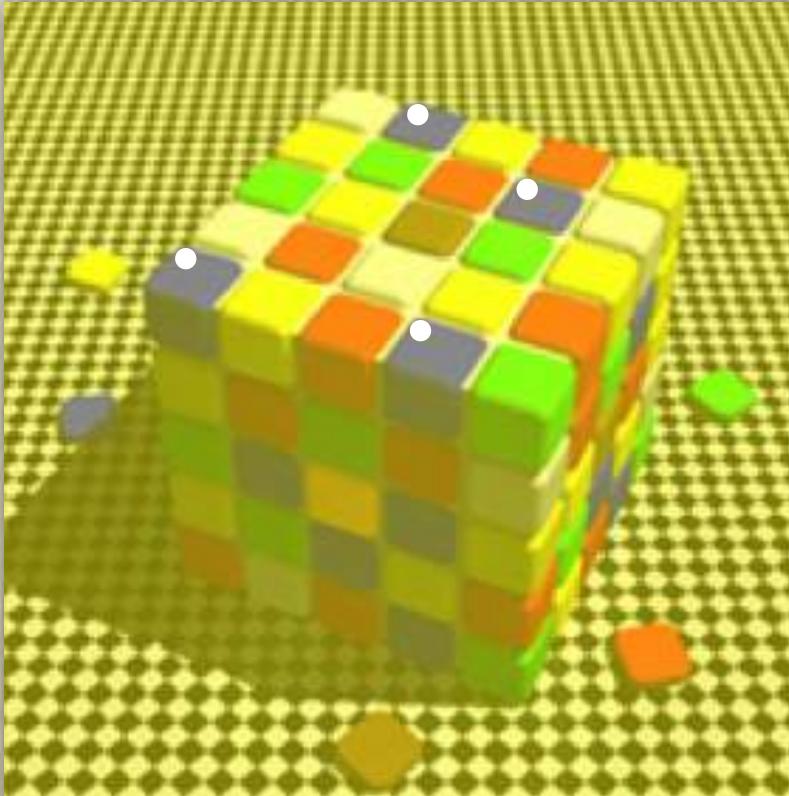
黄・青眼鏡を通して眺めるルービックキューブ[1,2]



- 両者が同一のキューブであることは容易に認知可能
- 異なる色を同一と主張し、同一の色を異なると主張する。
- 各パッチが持つ波長（絶対量）だけではなく、各パッチが他のパッチ群とどのようなコントラストを持つのか、が非常に重要

色みの偏差とその認知的不变性

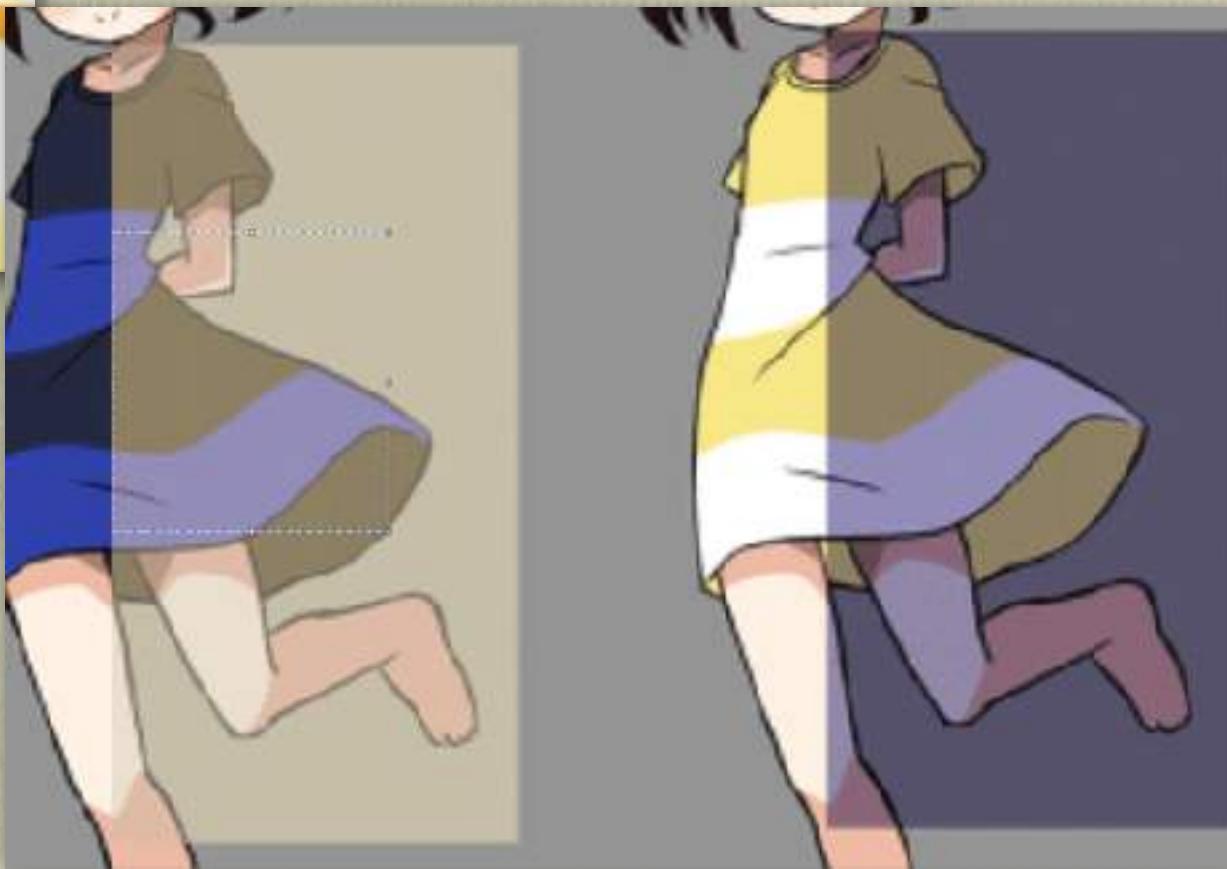
黄・青眼鏡を通して眺めるルービックキューブ[1,2]



- 両者が同一のキューブであることは容易に認知可能
- 異なる色を同一と主張し、同一の色を異なると主張する。
- 各パッチが持つ波長（絶対量）だけではなく、各パッチが他のパッチ群とどのようなコントラストを持つのか、が非常に重要

これ、覚えてますか？

数年前、ネットで大騒ぎになりました。



音高の偏差とその認知的不变性

カラオケでキーを上げ下げして曲を聞く[3,4]

The image shows two musical staves. Staff 1 is in C major (no sharps or flats) and staff 2 is in G major (one sharp). Both staves have a common time signature (indicated by 'c'). The notes are primarily eighth notes. In staff 1, the first note is highlighted with a blue oval and the third note with a red square. In staff 2, the second note is highlighted with a blue oval and the fourth note with a red square.

● 絶対音感者（ドレミは**音名**）

● 1 = ソーミソドーラードドソー, 2 = レーシレソーミーソソレー

● 言語化可能な相対音感者（ドレミは**階名**）

● 1 = ソーミソドーラードドソー, 2 = ソーミソドーラードドソー

● 言語化困難な相対音感者（**ラーラ音感者**）

● 1 = ラーラララーラーラララー, 2 = ラーラララーラーラララー

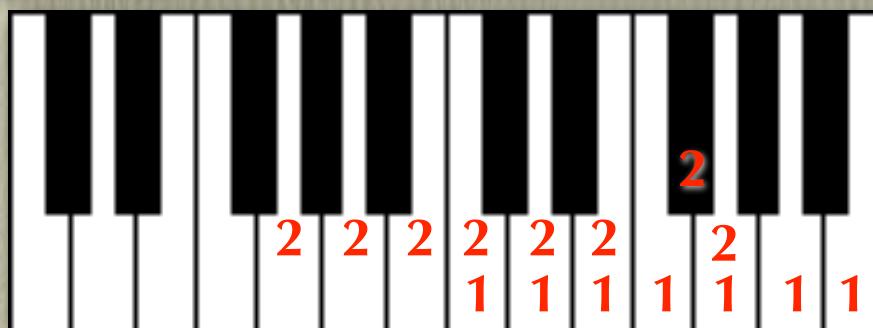
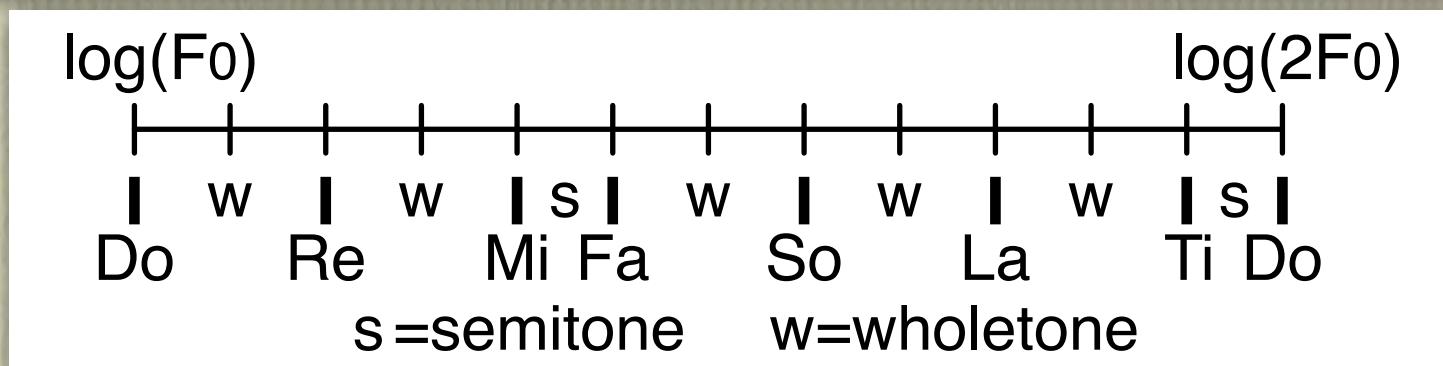
● 異なる音を同一と主張し、同一の音を異なると主張する。

● 各音が持つ基本周波数（絶対量）ではなく、各音が他の音群との
ようなコントラストを持つのか、のみによって決定

音高の偏差とその認知的不变性

カラオケでキーを上げ下げして曲を聞く[3,4]

1
2



- 各音が持つ基本周波数（絶対量）ではなく、各音が他の音群とどのようなコントラストを持つのか、のみによって決定

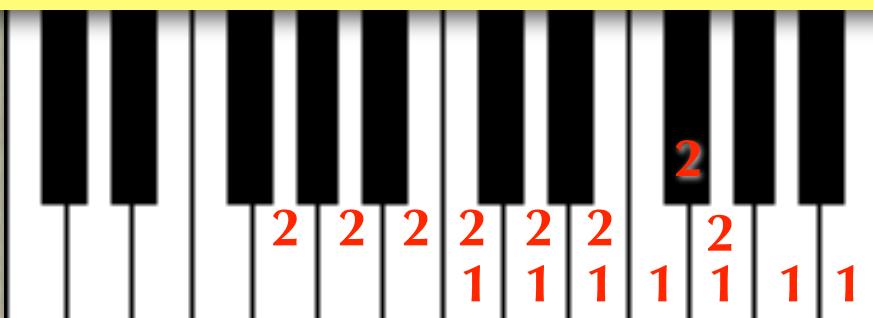
音高の偏差とその認知的不变性

カラオケでキーを上げ下げして曲を聞く[3,4]

1
2

$\log(F_0)$ $\log(2F_0)$

但し、孤立音の同定は不可能
そこにはコントラストが無いから



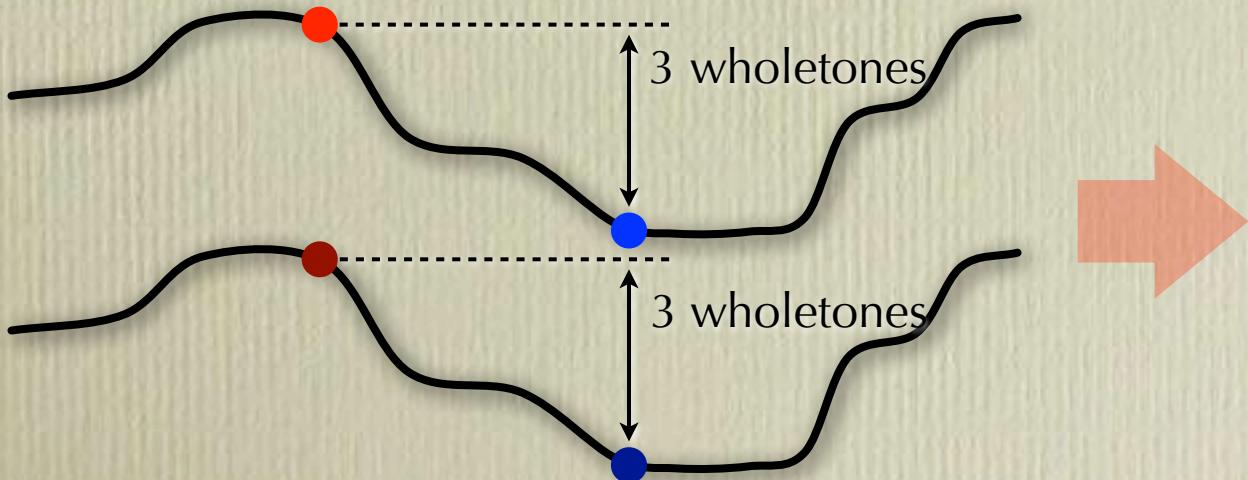
- 各音が持つ基本周波数（絶対量）ではなく、各音が他の音群との
ようなコントラストを持つのか、のみによって決定

Invariant pitch perception against its bias

A melody and its transposed version [Higashikawa'05]

1)  2) 

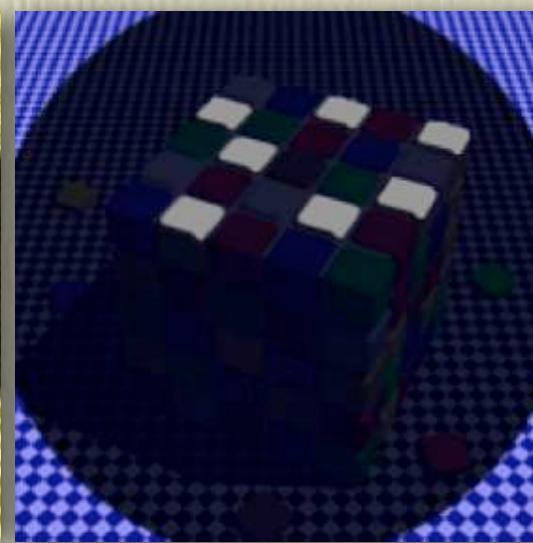
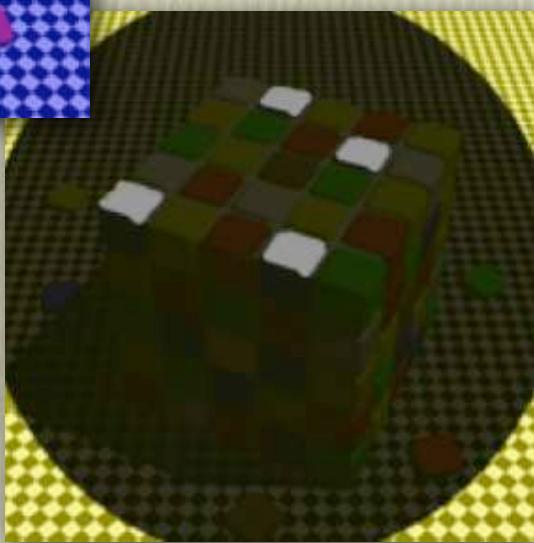
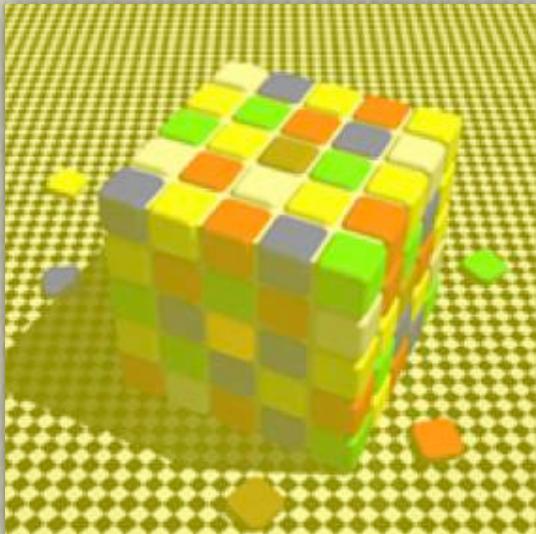
- Listeners with RP can perceive the same sound name sequence.
 - So Mi So Do / Ra Do Do So / So Do Re Mi Re Do / Re
 - The same sound distribution pattern is found in 1) and 2).



● ● and ● ● have to be
fa & ti or ti & fa due to
contrastive constraints.

生物が獲得した静的バイアス除去術

色の恒常的・不变的認知はどこまで遡れるのか？[5]



生物が獲得した静的バイアス除去術

音高の恒常的・不变的認知はどこまで遡れるのか？[6]

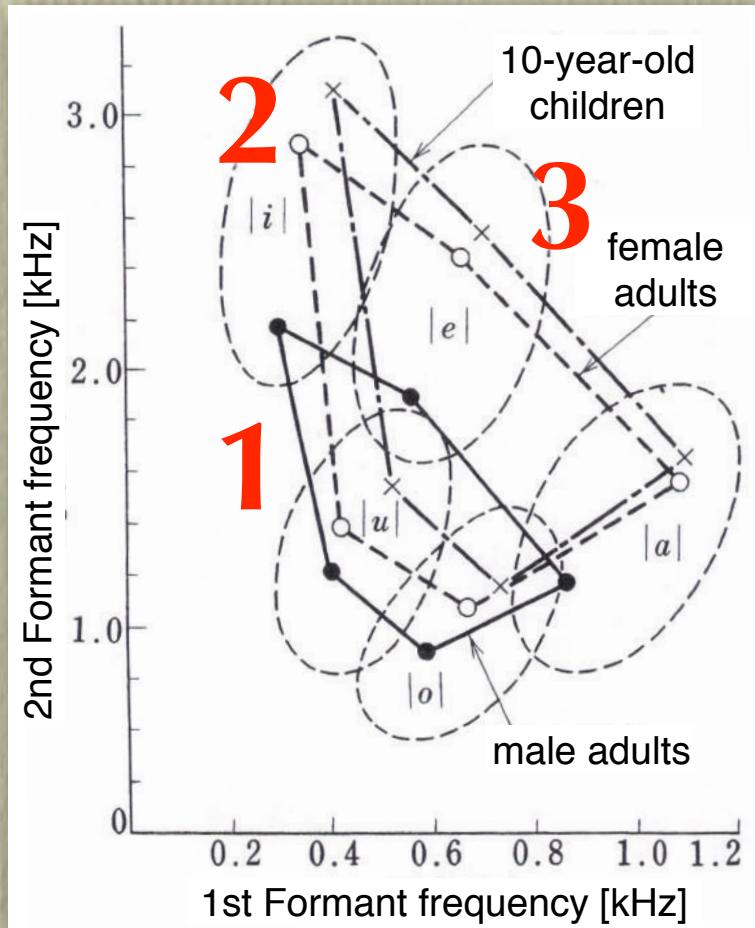


1 = 2



じゃあ、これは？

音色の恒常的・不变的認知はどこまで遡れるのか？

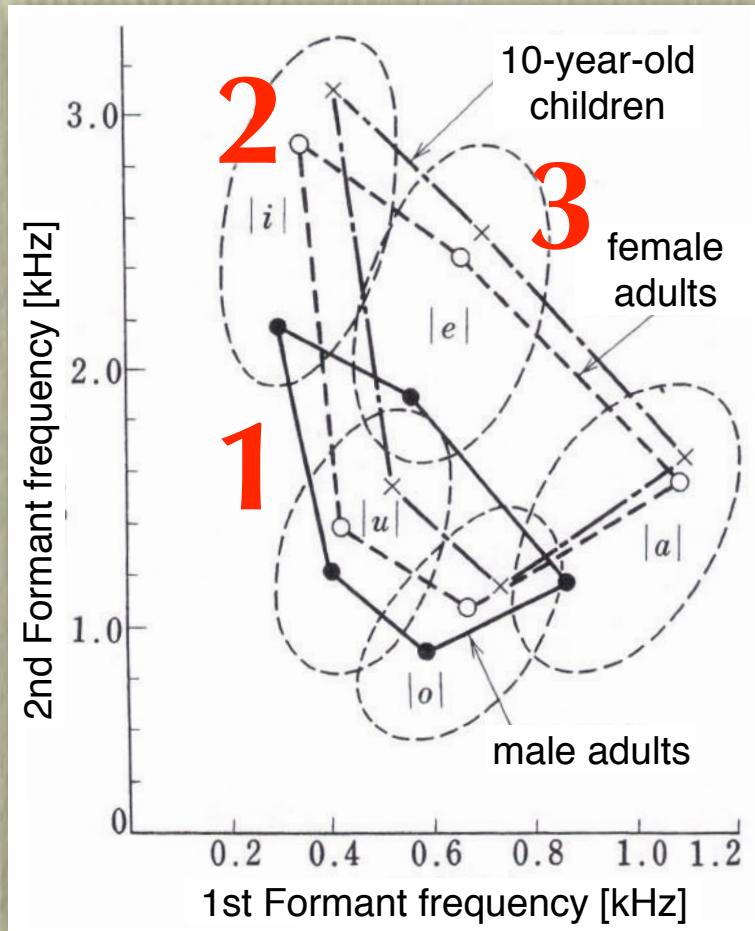


1 = 2 = 3



じゃあ、これは？

音色の恒常的・不变的認知はどこまで遡れるのか？



1 = 2 = 3

?



本発表の流れ

刺激の物理的多様性とその認知的不变性

- 見え／色み／音高の多様性と自然・進化が編み出した解決方法

音声の物理的多様性とその認知的不变性

- 音色の多様性と工学者が編み出した解決方法

音声の構造的表象とそれに関する様々な考察

- 常識を覆すことで、違和感の解消を試みてみる。

音声の構造的表象と数学的表現と技術的実装

- 体格・性別に不变な音声波形・スペクトルの表現とは？

音声の構造的表象を用いた音声アプリケーション

- 音声認識、音声合成、発音分析、etc

音声の構造的表象の言語学的妥当性

- 何故、こうしてこなかったのか？ 観測技術の功罪？

音色の偏差とその認知的不变性

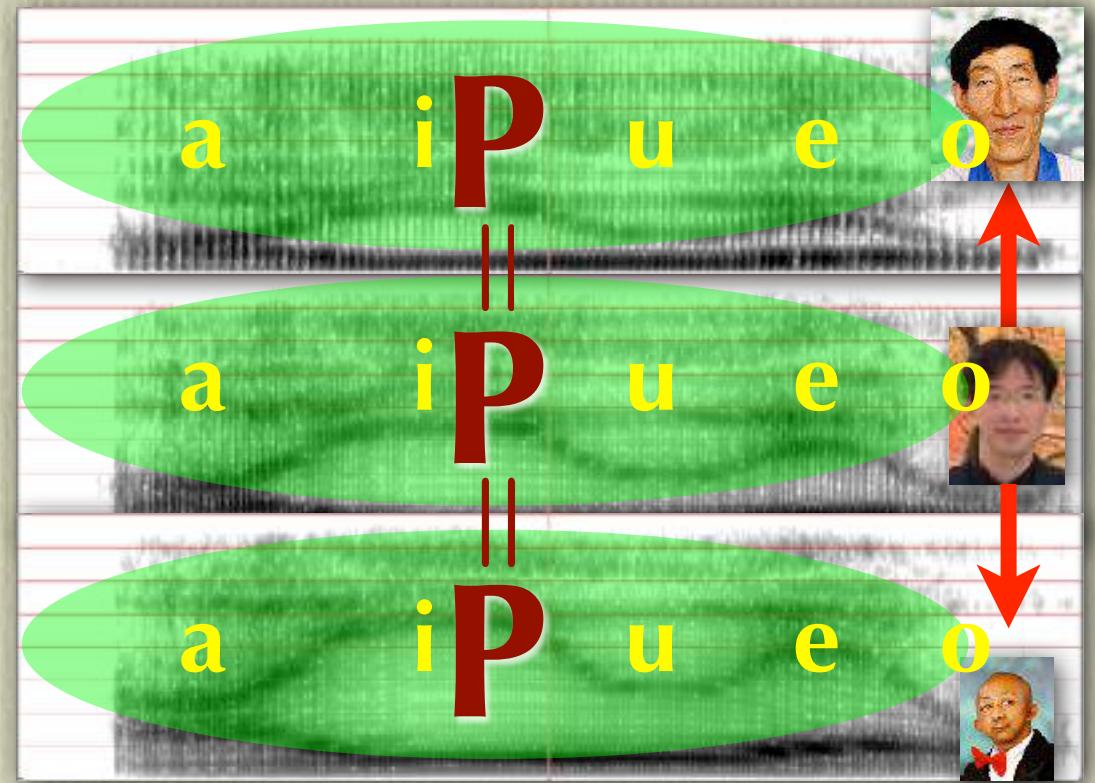
音高の静的偏差を生み出す要因

- 男女の音高偏差 = 声帯の長さ・重さの性差



音色の静的偏差を生み出す要因

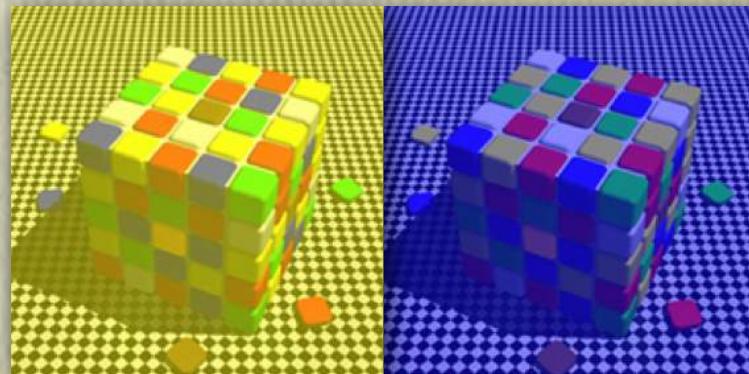
- 男女の音色偏差 = 声道の形状（主に長さ）の性差



音色の偏差とその認知的不变性

色み・音高の恒常・不变的認知

- コントラスト情報に基づく処理が重要
- 要素同定ではなく、コントラスト群から成る全体的パターン処理

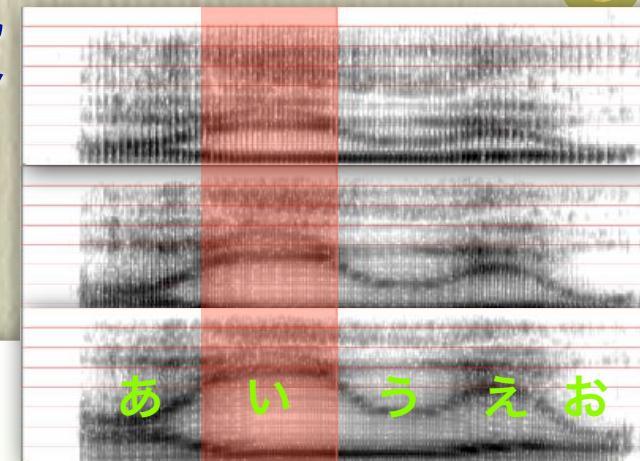


$$P(o|w) \sim \sum_s P(o|w, s)P(s)$$

音色の偏差とその認知的不变性

音色の偏差に対する工学的な常套手段

- 音声ストリームを要素列として表象し、
- 個々の要素の統計モデルを作る。

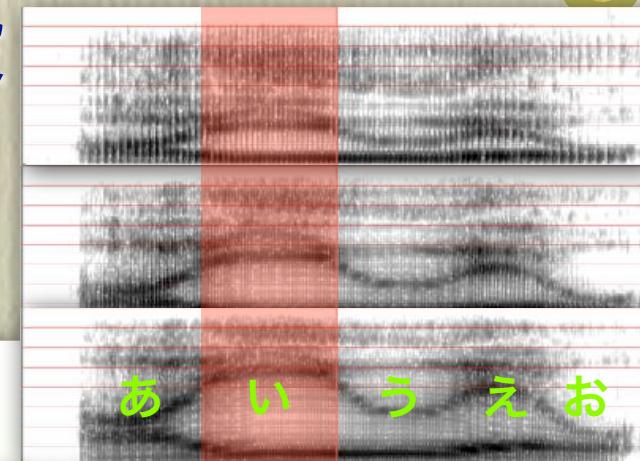


数千～数十万

音色の偏差とその認知的不变性

音色の偏差に対する工学的な常套手段

- 音声ストリームを要素列として表象し、
- 個々の要素の統計モデルを作る。



数十～数十万

幼児の言語獲得と音声模倣

音声模倣＝親の発声行為を子が積極的に模倣する行為

- これを通して幼児は言語を獲得する[7]
- 動物学的には非常に稀な行為。靈長類では人間だけ[8]
- 他の動物では小鳥、クジラ、イルカくらいか[10]

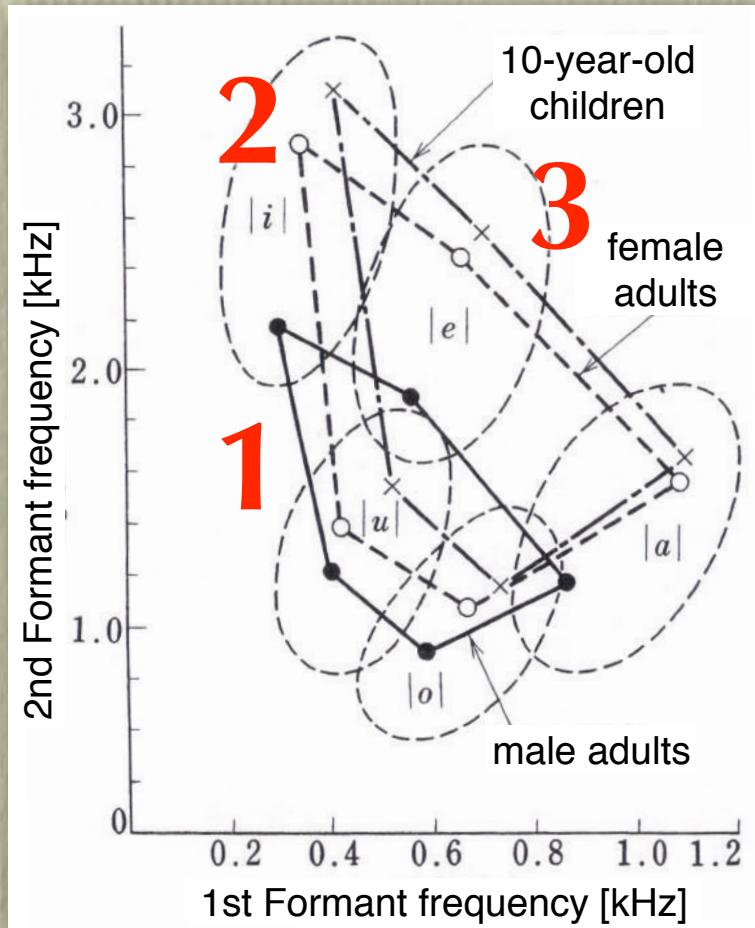
動物の模倣＝声帯模写、ヒトの音声模倣≠声帯模写

- 九官鳥の音声模倣[9]
 - 車、ドア、椅子、犬、猫、音を真似る。人の声も音でしかない。
 - 良い九官鳥を聞くと、飼い主が分かる。
- 幼児の音声模倣
 - 動物学的には奇妙な模倣行為[10]
 - いくら良い子でも、声から父親を割り出せずにお巡りさんは困る。



じゃあ、これは？

音色の恒常的・不变的認知はどこまで遡れるのか？



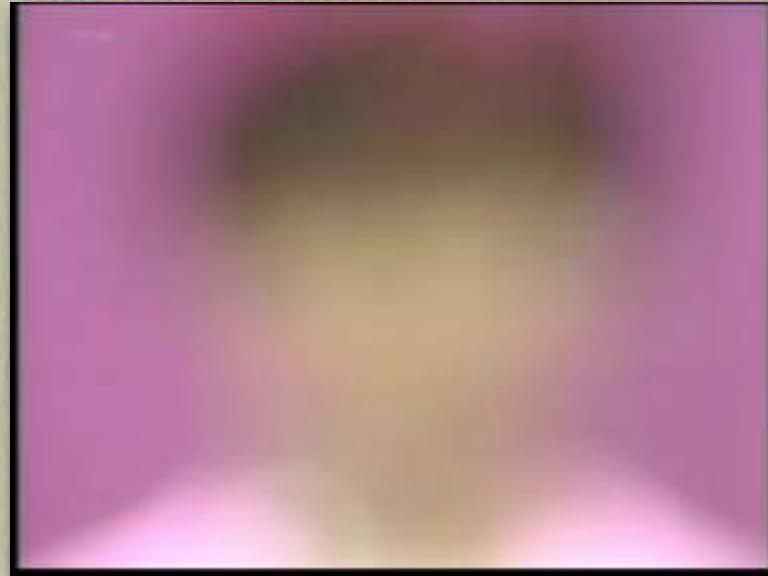
1 = 2 = 3

?



声帯模写と非声帯模写

松田聖子・まねだ聖子・神田沙也加



音声模倣の二面性～音真似と？真似～

親の発声→音韻同定→音韻列→個々の音韻を発声？



×
/おはよう/ →



- 音韻意識（仮名の意識）が希薄／しり取りも出来ない。

発達心理学からの回答

- 幼児は語全体の語形・音形・枠組み・ゲシュタルトを獲得し、その後、個々の分節音（音韻・仮名）を獲得する
- 語ゲシュタルトには話者の情報は含まれない。話者不变量
 - if not, 幼児は動物のように音声模倣することになる。
- 語ゲシュタルトの物理的・音響的定義は何か？
- 親の声と幼児の声の「物理的な共通項」は何か？

?

音色の偏差とその認知的不变性

色み・音高の恒常・不变的認知

- コントラスト情報に基づく処理が重要
- コントラスト群から成る全体的パターン処理が要素同定を可能に



音色の恒常・不变的認知

- コントラスト情報に基づく処理が重要
- コントラスト群から成る全体的パターン処理が要素同定を可能に



本発表の流れ

刺激の物理的多様性とその認知的不变性

- 見え／色み／音高の多様性と自然・進化が編み出した解決方法

音声の物理的多様性とその認知的不变性

- 音色の多様性と工学者が編み出した解決方法

音声の構造的表象とそれに関する様々な考察

- 常識を覆すことで、違和感の解消を試みてみる。

音声の構造的表象と数学的表現と技術的実装

- 体格・性別に不变な音声波形・スペクトルの表現とは？

音声の構造的表象を用いた音声アプリケーション

- 音声認識、音声合成、発音分析、etc

音声の構造的表象の言語学的妥当性

- 何故、こうしてこなかったのか？ 観測技術の功罪？

音高の偏差とその認知的不变性

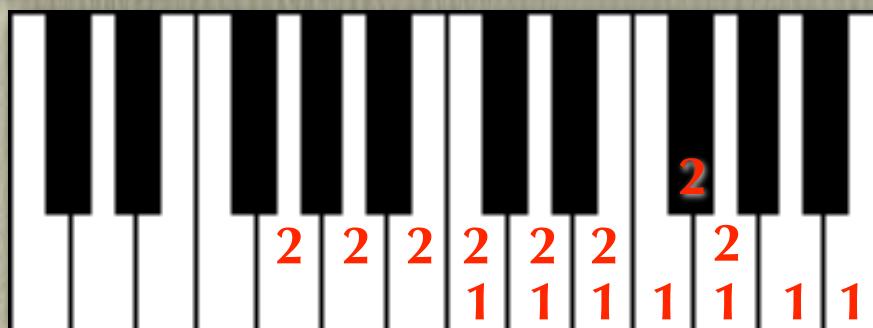
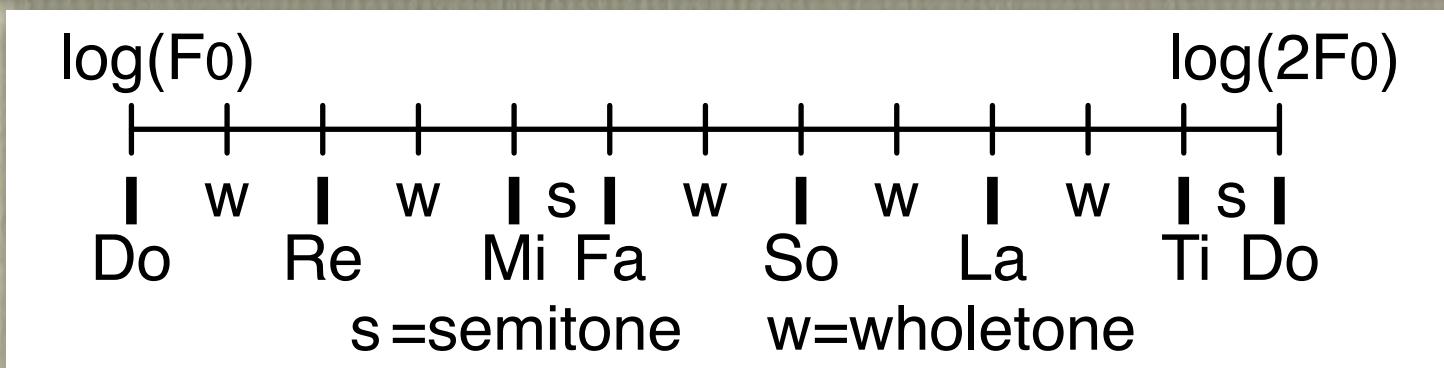
カラオケでキーを上げ下げして曲を聞く[3,4]

The image shows two musical staves. Staff 1 is in C major (no sharps or flats) and staff 2 is in G major (one sharp). Both staves have a common time signature (indicated by 'c'). The notes are primarily eighth notes. In staff 1, the first note is highlighted with a blue oval and the third note with a red square. In staff 2, the second note is highlighted with a blue oval and the fourth note with a red square.

- 絶対音感者（ドレミは**音名**）
 - 1 = ソーミソドーラードドソー, 2 = レーシレソーミーソソレー
- 言語化可能な相対音感者（ドレミは**階名**）
 - 1 = ソーミソドーラードドソー, 2 = ソーミソドーラードドソー
- 言語化困難な相対音感者（**ラーラ音感者**）
 - 1 = ラーラララーラーラララー, 2 = ラーラララーラーラララー
- 異なる音を同一と主張し, 同一の音を異なると主張する。
- 各音が持つ基本周波数（絶対量）ではなく, **各音が他の音群とのようなコントラストを持つのか**, のみによって決定

音高の偏差とその認知的不变性

カラオケでキーを上げ下げして曲を聞く[3,4]



各音が持つ基本周波数（絶対量）ではなく、各音が他の音群との
ようなコントラストを持つのか、のみによって決定

音高の偏差とその認知的不变性

- ## カラオケでキーを上げ下げして曲を聞く[3,4]

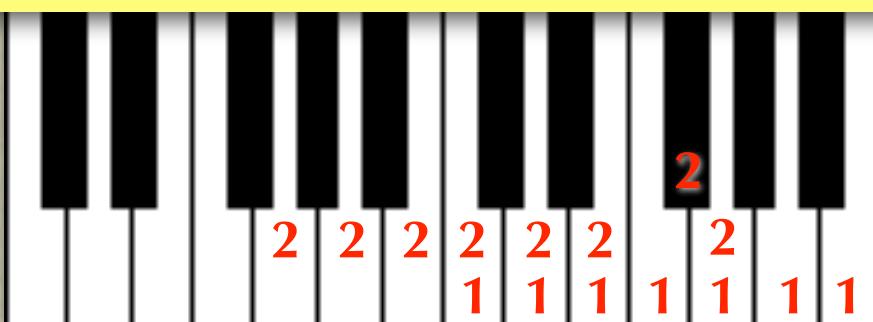
1

2

$\log(F_0)$

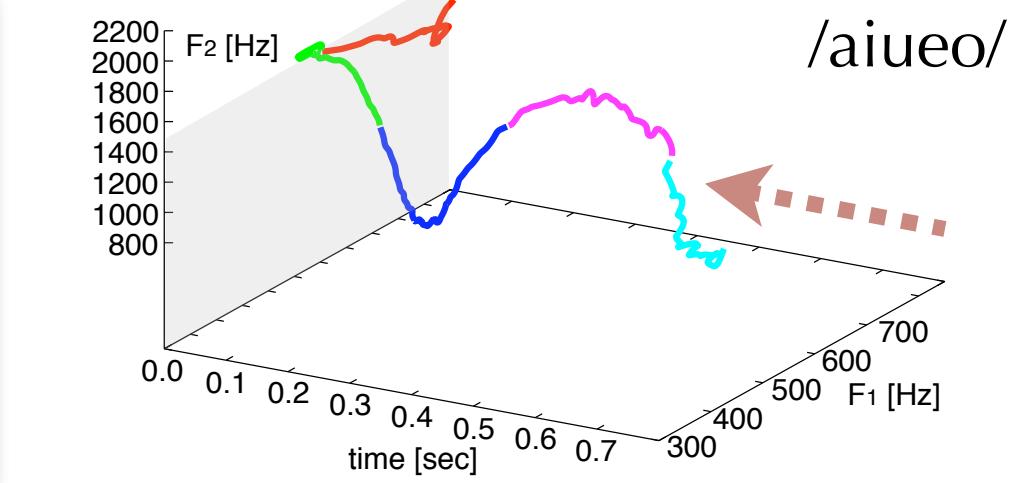
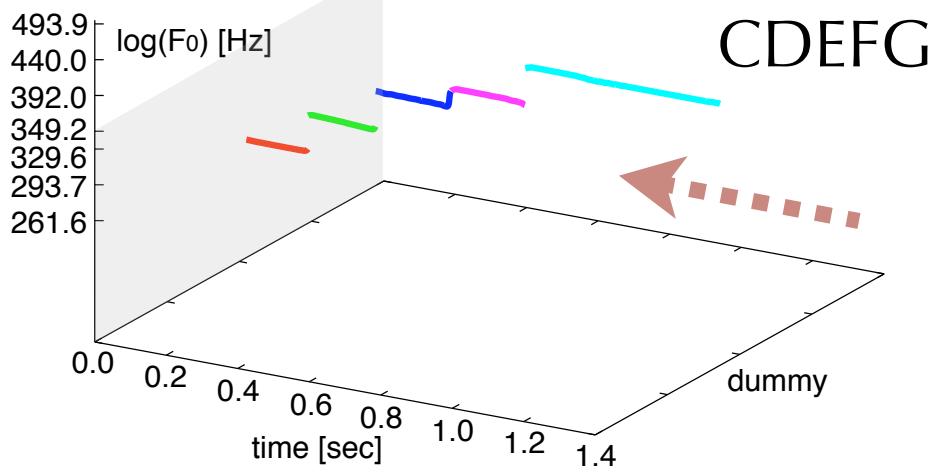
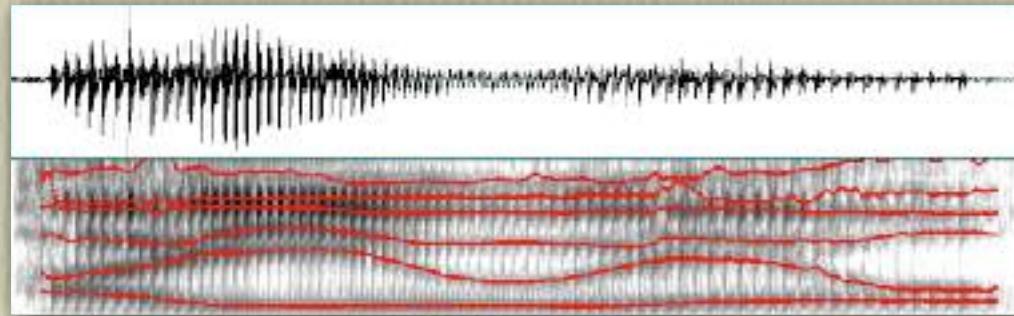
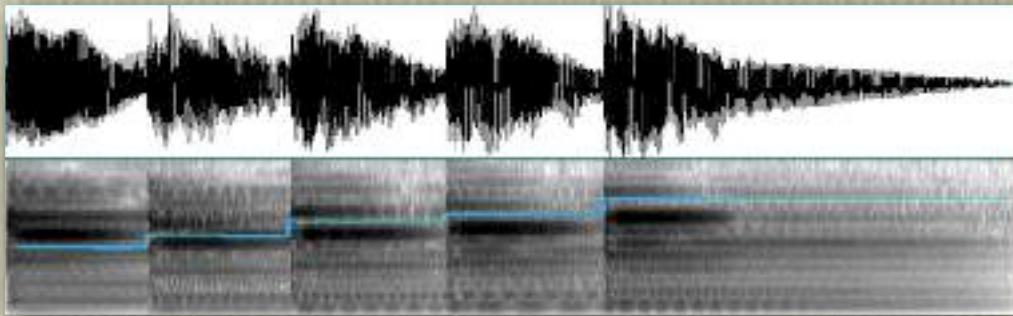
$\log(2F_0)$

但し、孤立音の同定は不可能
そこにはコントラストが無いから



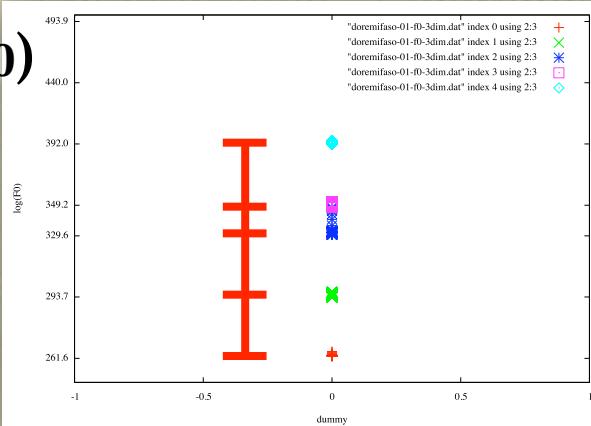
- 各音が持つ基本周波数（絶対量）ではなく、各音が他の音群との
のようなコントラストを持つのか、のみによって決定

音声の構造的表象／音色の相対音感



音高の動的变化パターン

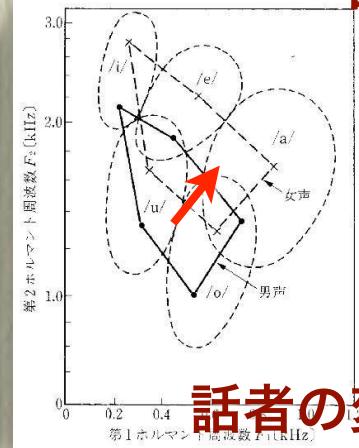
$\log(F_0)$



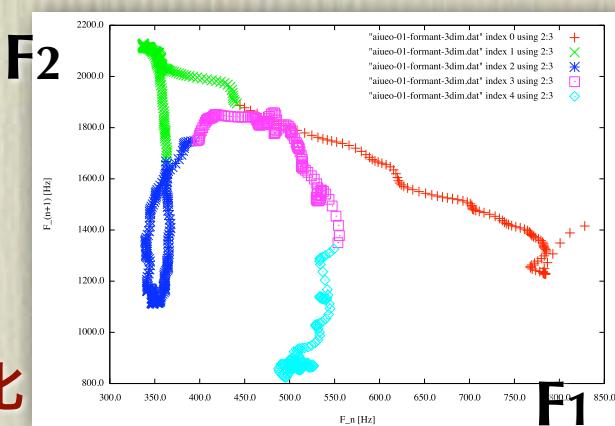
調の変化

音色の動的变化パターン

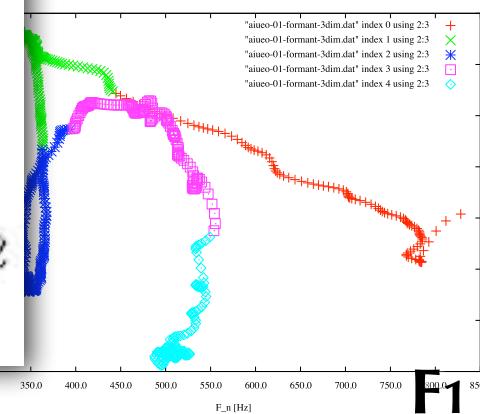
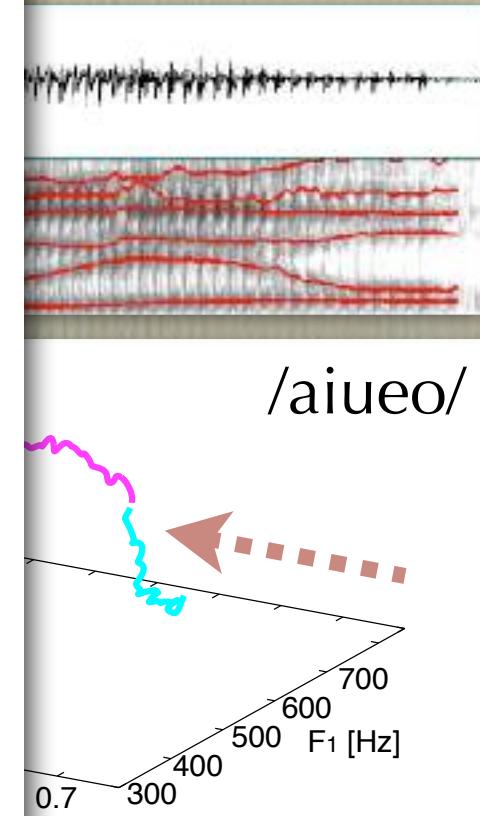
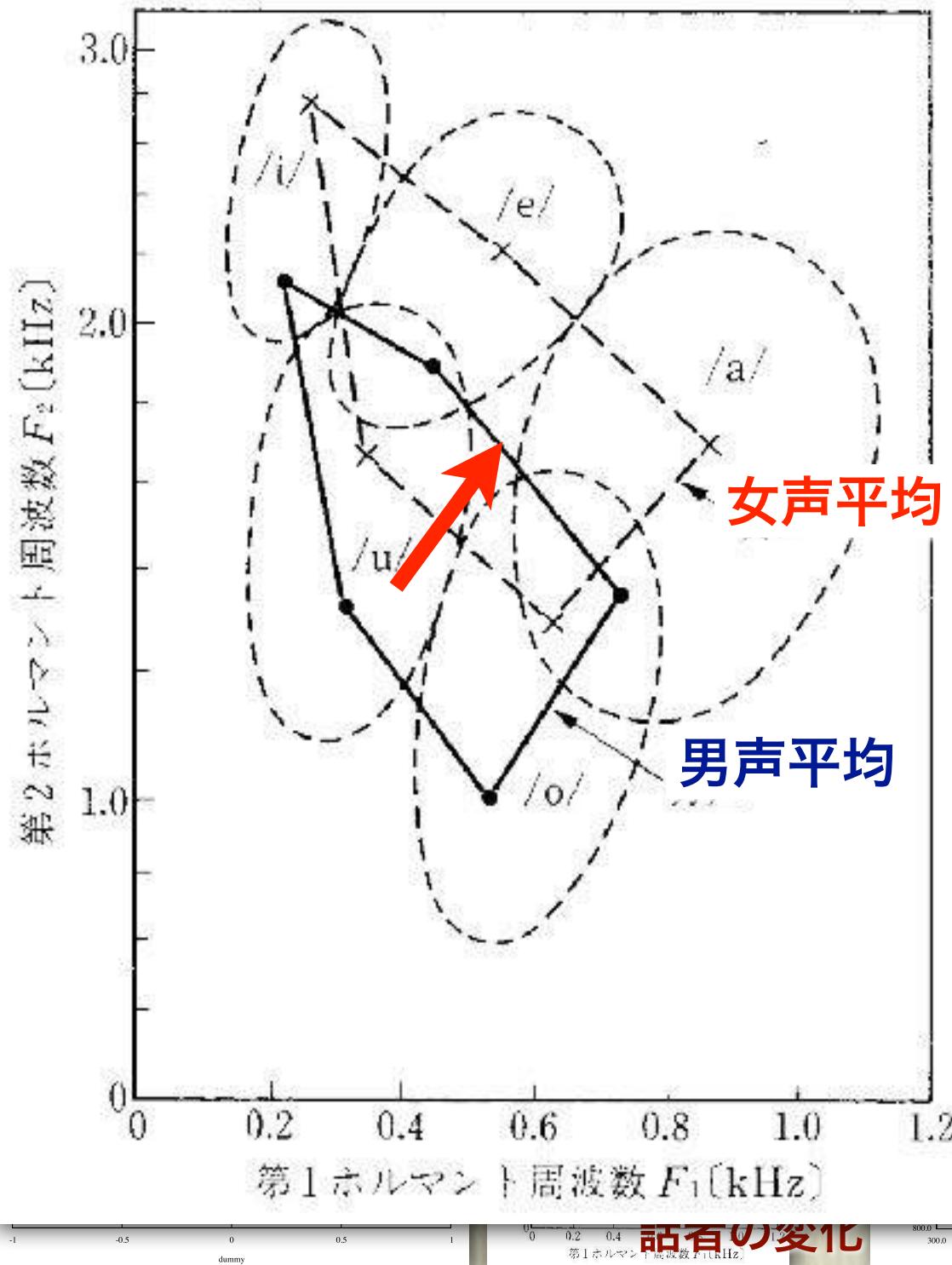
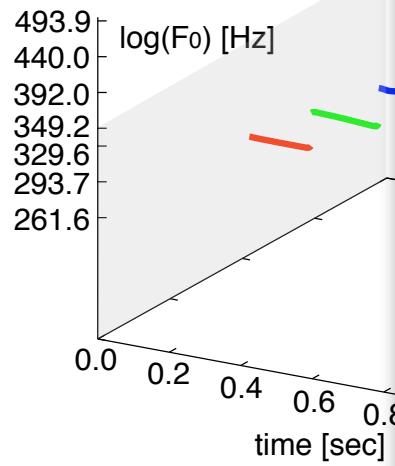
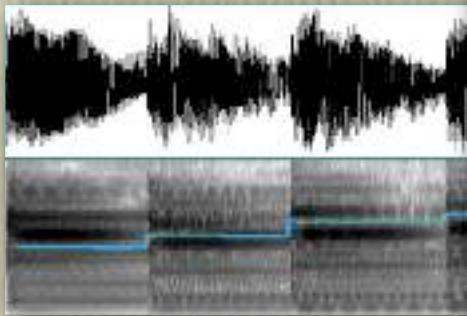
F_2



話者の変化

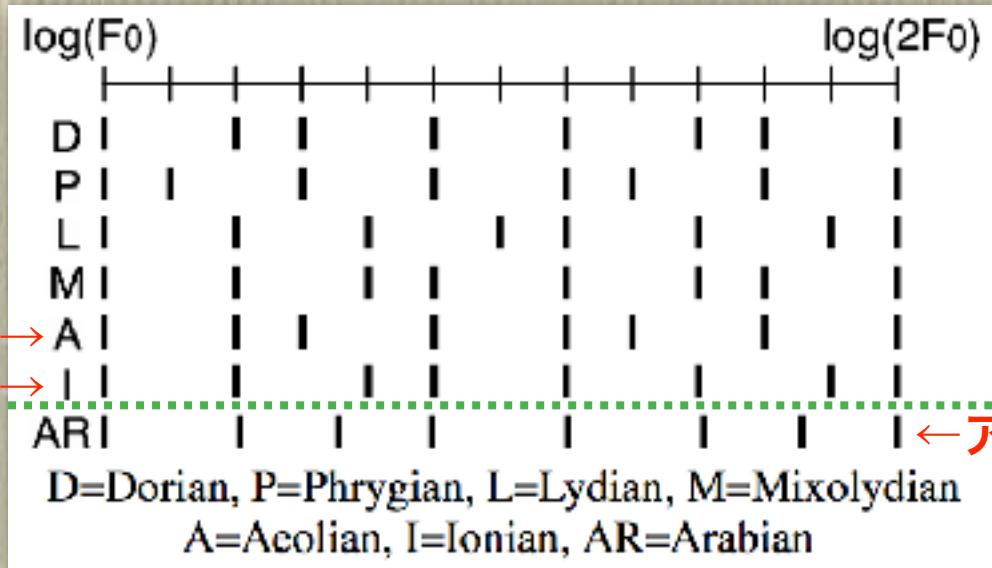


音声の構造的音色 / 音色の相対音感



音声の構造的表象／音色の相対音感

音楽における調不变の音配置とその変種

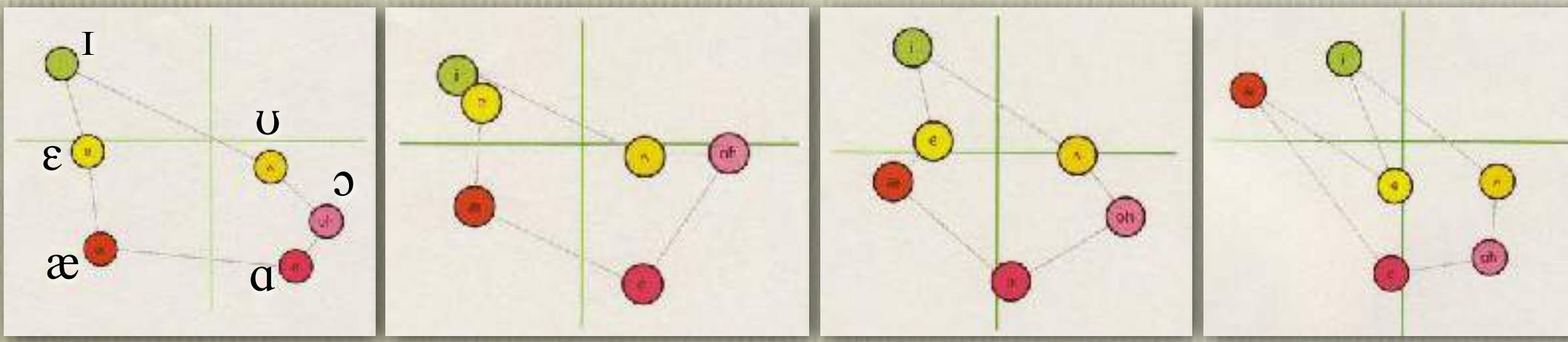


- 西洋音楽 = 5全音 + 2半音
- 種々の配置 = 教会音楽
- 民族音楽には半音以外の配置

أهلاً وسهلاً



音声における話者不变の音配置とその変種 = 欧米の方言



Williamsport, PA

Chicago, IL

Ann Arbor, MI

Rochester, NY

相対音感者による転調部分のドレミ同定

曲の途中で調が変わるとドレミ同定はどうなる？

- 絶対音感者 ~音名としてのドレミを使って書き起こす~
 - 何ら問題無く、ドレミ同定する。
 - 時には転調したことに気付かないことすらある。
- 相対音感者 ~階名としてのドレミを使って書き起こす~
 - 転調すると、とたんに「ドレミ」同定ができなくなる。
 - いわゆるパニクる。その後しばらくして、まだ出来るようになる。
 - アラビア犬のワルツ→ドレミが聞こえる所と聞こえない所がある。

発話の途中で話者が変わるとモーラ同定はどうなる？

- 話者が変わろうが、同定率が落ちない人
 - 音色の絶対音感者（音色の絶対的特性に基づいて判断する）？
- 話者が変わると、同定率が落ちる人
 - 音色の相対音感者（音色の相対的特性に基づいて判断する）？
 - 話者変化による音色の変化を音韻の変化と捉える人がいる？

興味深いサイト

絶対音感ある人に30の質問

http://www.100q.net/100/question.cgi?que_no=51

The screenshot shows a web browser window with the URL www.100q.net/100/question.cgi?que_no=51. The title bar says "絶対音感ある人に30の質問". The main content area displays a list of 30 questions, each with a question number, name, and answer.

[Home]	[ReLoad]
絶対音感ある人に30の質問	
[715] heyjoe	---
[714] レイン	---
[713] Julio	---
[712] vista	---
[711] ざるそば	---
[710] ただは	---
[709] 優々	☒ --
[708] あいり	---
[707] 唯。	---
[706] あじさい	---
[705] ろはん	---
[704] mao	---
[703] 岡崎汐	---
[702] とも	---

Q1 氏名前と、この質問の回答日を教えてください。
vistaです。2014.04.19

Q2 年齢・性別をお願いします。
16。女。

Q3 今のお仕事を教えてください。
学生。

Q4 初めて音楽の手ほどきを受けたのは何歳?そのときの楽器は?
4歳のときにピアノを。

Q5 あなたの絶対音感は先天性?それとも後天性?
たぶん先天性

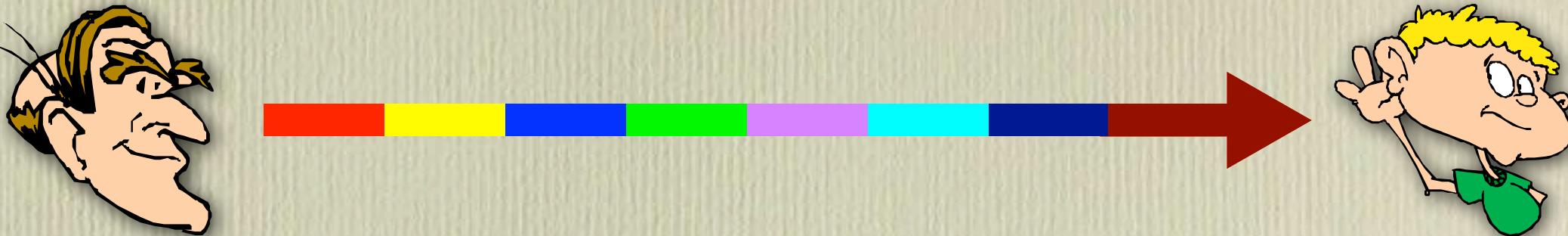
Q6 ピアノのA(ラ)を44Hzでとらえていますか?
441Hz

Q7 暗譜は得意ですか?
はい

Q8 移調は得意ですか?
大好き

話者がコロコロ変わる音声の知覚

話者性が時間軸に沿って変化する音声



- もし全体的表象が使用されていれば、同定率は低下するはず。

音声刺激の作成

- HMM合成（男性アナウンサー7名／ATR503文）
- メルケプストラム（0～24次元），7状態5分布
- 無意味8モーラ列 ($F0=LHHHLLLL$, 4型)
 - 促音, 撥音, 拗音, 濁音, 半濁音などのモーラは使用せず。全43種類
- 話者性変化のタイミング
 - 8/4/2/1モーラ, 1音素, 1分布 (5人／音素)

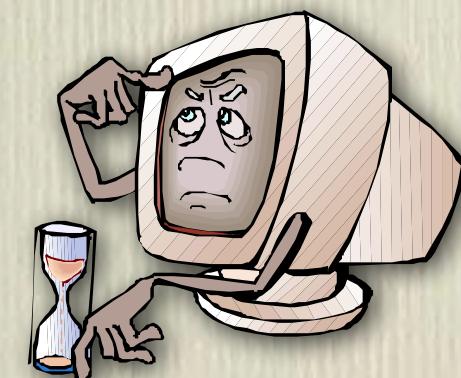
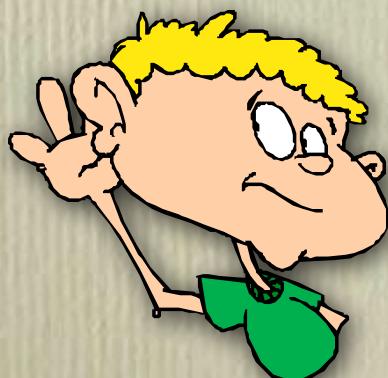
話者がコロコロ変わる音声の知覚

実験手順

- 話者性変化間隔6種類 × 25音声 = 150音声刺激
- Web上の音声ファイルをクリックにてヘッドフォン提示
- 聴取回数 = 2回
- 8モーラであること、一部モーラの欠落は事前教示

三種類の被験者

- 音声研究に従事する大学院生（合成音評価実験経験有）5名
- 法学部大学生（合成音評価実験経験無）3名
- HVite+CSRC不特定話者音響モデル+8モーラ認識文法



話者がコロコロ変わる音声の知覚

実験結果に対する四つの予測

- 前後音との関係に基づく聴取 = 音声をより大きな単位で知覚
- 話者性変化頻度の上昇と共にモーラ同定率は劣化
- 音声をより小さな単位で知覚 = 分析的な聴取
- 話者性変化頻度によらず一定のモーラ同定率
- 音色の相対音感の存在を本当に知らない唯一の実体
- 不特定話者音響モデル = 一定のモーラ同定率
- 話者性頻度が極めて大 = 話者性の差の表出・知覚が困難
- 話者性変化頻度の極端な上昇によってモーラ同定率は向上

不特定話者音声の例

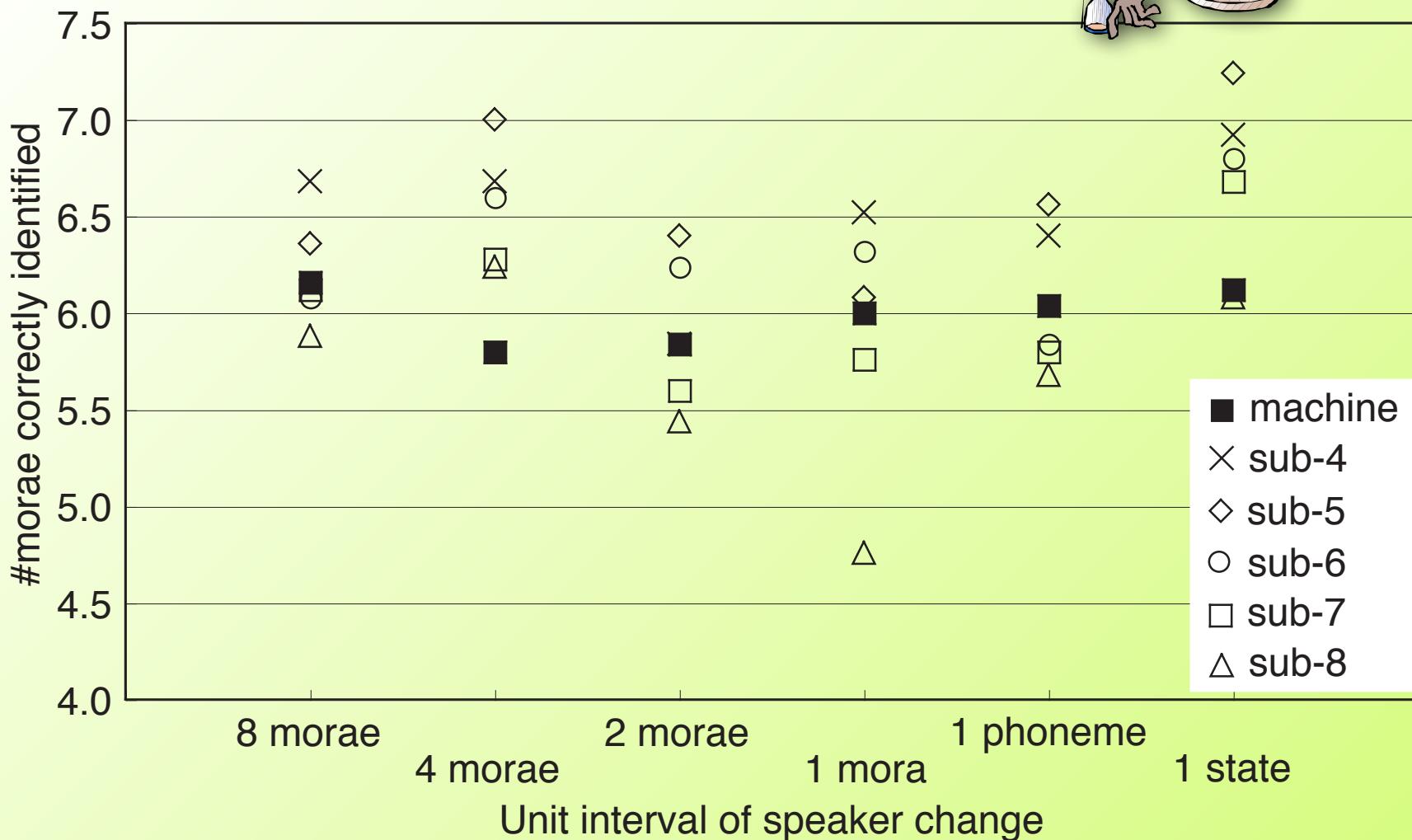
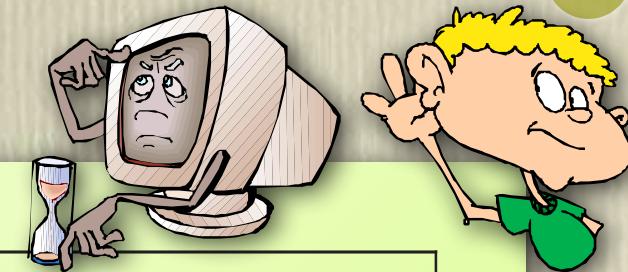
- 1モーラ単位での不特定話者音声
- 話者性の変化を音韻の変化として認知する可能性

レカルツヌチオキ



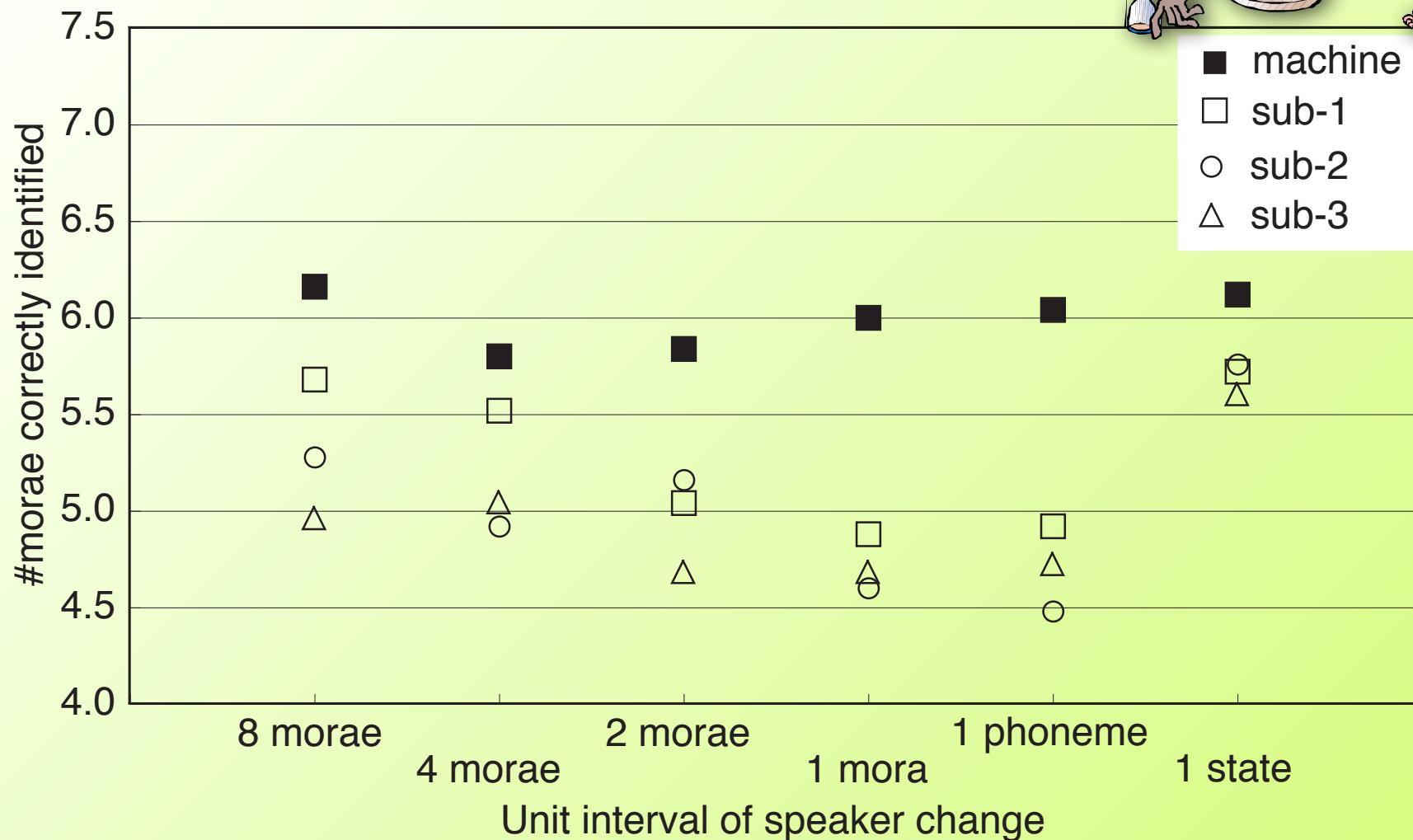
話者がコロコロ変わる音声の知覚

実験結果（計算機と音声研究者）



話者がコロコロ変わる音声の知覚

実験結果（計算機と法学部学生）



音高の偏差とその認知的不变性

カラオケでキーを上げ下げして曲を聞く[3,4]

The image shows two musical staves. Staff 1 is in C major (no sharps or flats) and staff 2 is in G major (one sharp). Both staves have a common time signature (indicated by 'c'). The notes are primarily eighth notes. In staff 1, the first note is highlighted with a blue oval and the third note with a red square. In staff 2, the second note is highlighted with a blue oval and the fourth note with a red square.

- 絶対音感者（ドレミは**音名**）
 - 1 = ソーミソドーラードドソー, 2 = レーシレソーミーソソレー
- 言語化可能な相対音感者（ドレミは**階名**）
 - 1 = ソーミソドーラードドソー, 2 = ソーミソドーラードドソー
- 言語化困難な相対音感者（**ラーラ音感者**）
 - 1 = ラーラララーラーラララー, 2 = ラーラララーラーラララー
- 異なる音を同一と主張し, 同一の音を異なると主張する。
- 各音が持つ基本周波数（絶対量）ではなく, **各音が他の音群とのようなコントラストを持つのか**, のみによって決定

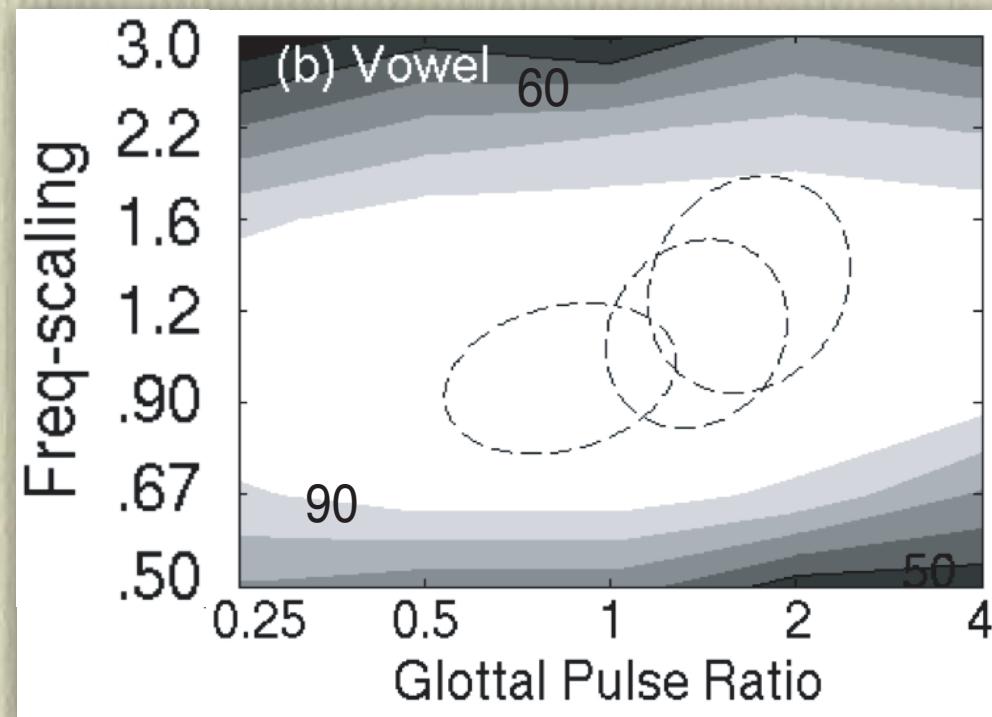
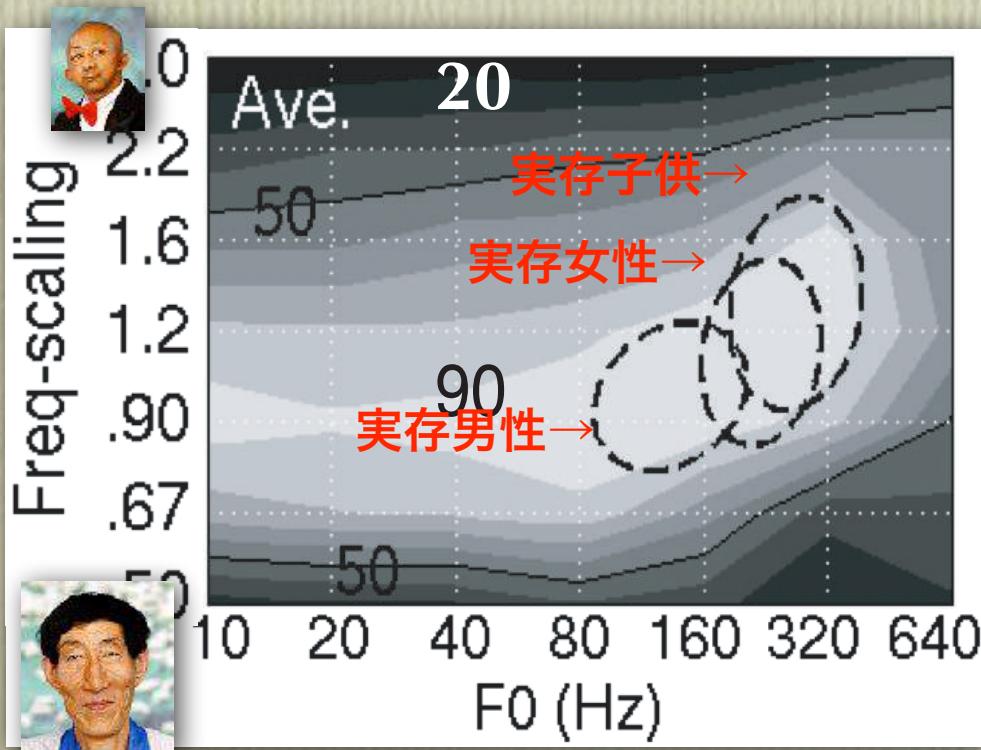
音声の構造的表象／音色の相対音感

言語化できる相対音感者が出来ないこと

- 孤立的に提示された音をドレミ同定することは出来ない。
- 孤立的に提示された音を母音同定できない人などいるのか？

巨人＆小人の音声を使った母音同定・単語同定実験

- 孤立母音の同定は困難になる[18]
- でも、無意味語でよいので単語音声にすると書き起こせる[19]



音声の構造

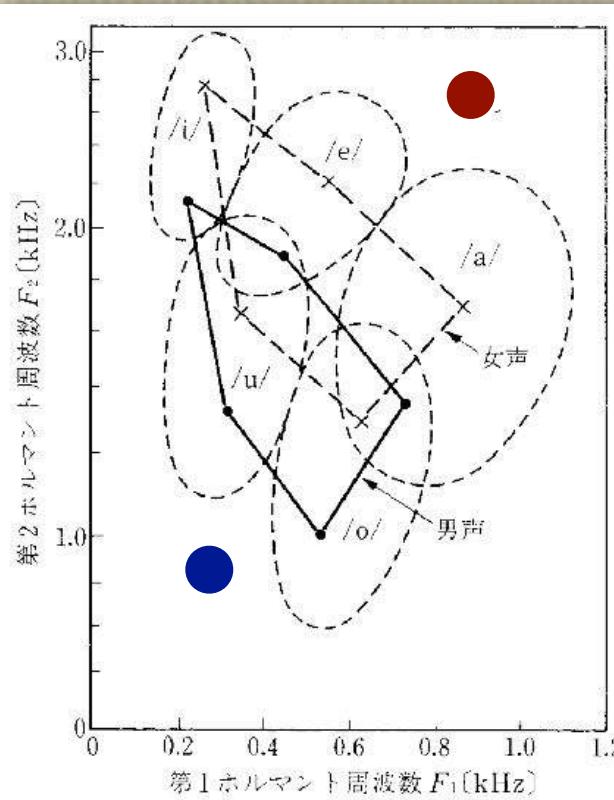
の相対音感

言語化できる相対

- 孤立的に提示された
- 孤立的に提示された

巨人&小人の音声

- 孤立母音の同定は困
- でも、無意味語でよ

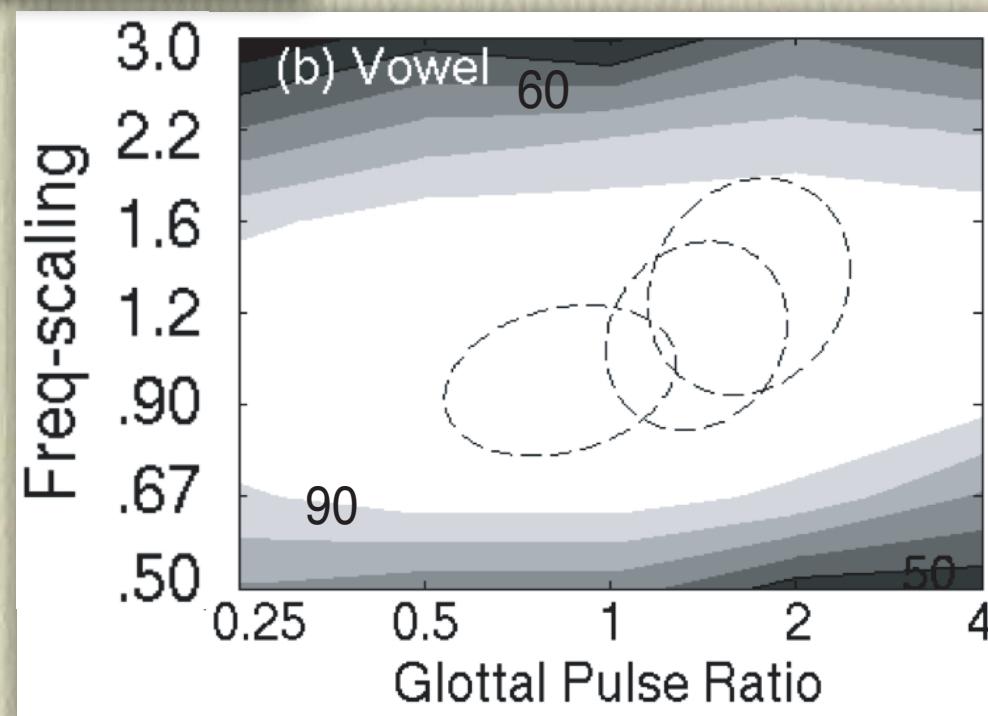
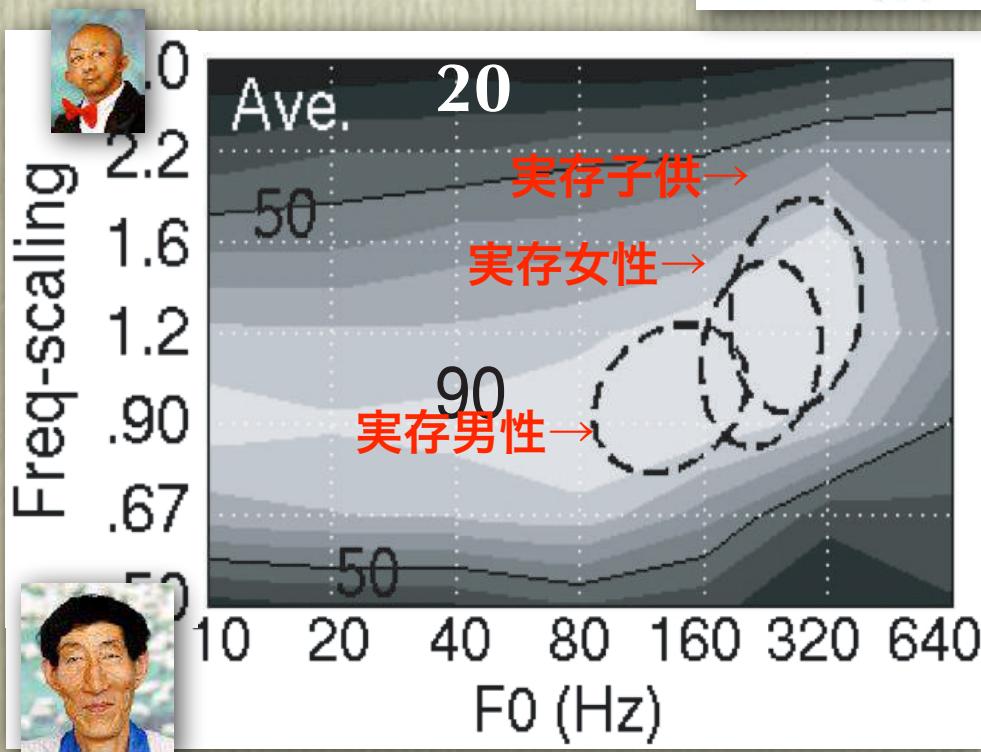


こと

ことは出来ない。
い人などいるのか？

・単語同定実験

ると書き起こせる[19]



音声の構造的主角 / 产角の相対音感

言語化できる相対音感

- 孤立的に提示される相対音感
- 孤立的に提示される相対音感

巨人&小人の音声

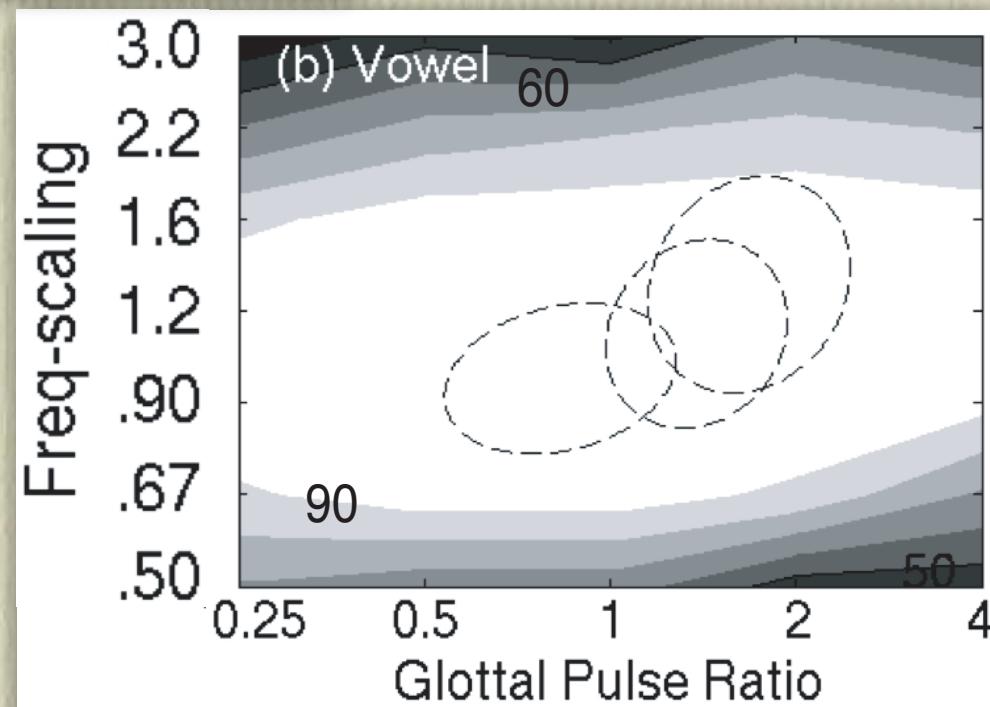
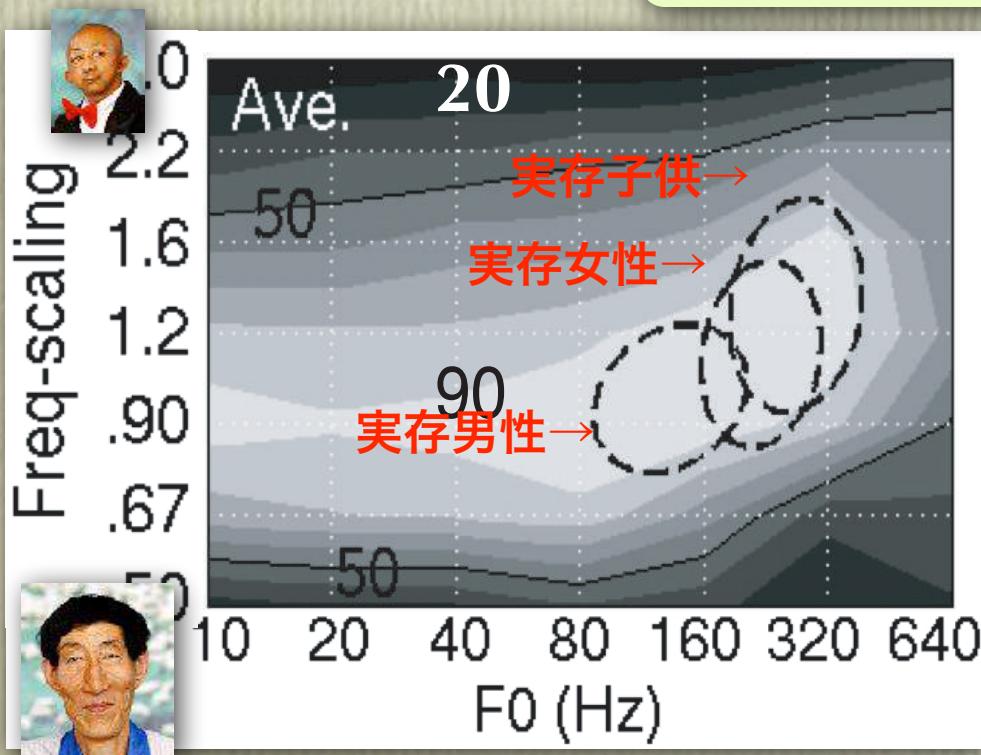
- 孤立母音の同定は可能
- でも、無意味語では



こと
とは出来ない。
人などいるのか？

単語同定実験

と書き起こせる[19]



音声の構造的主角 / 產角の相対音感

言語化できる相対音感

- 孤立的に提示されると認識される
- 孤立的に提示されると認識される

巨人＆小人の音声

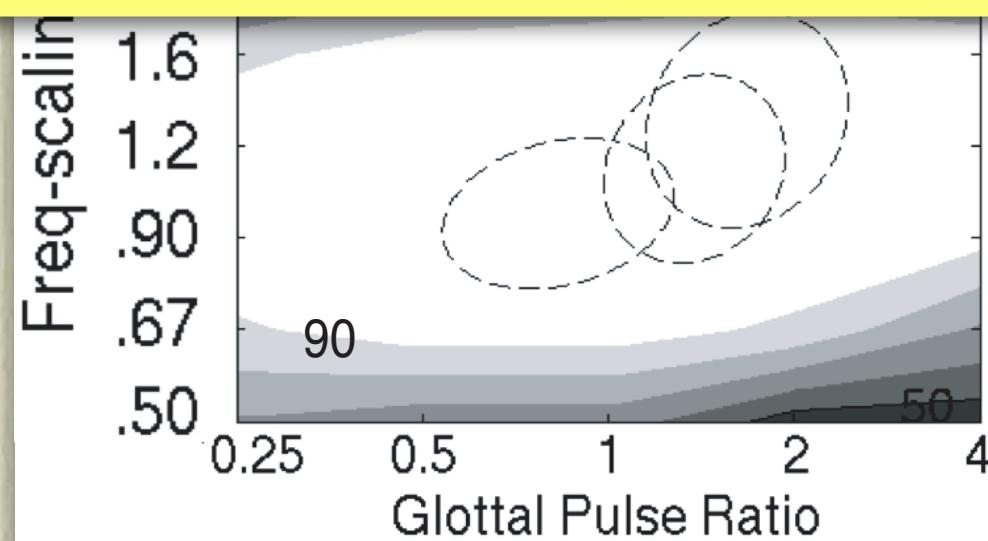
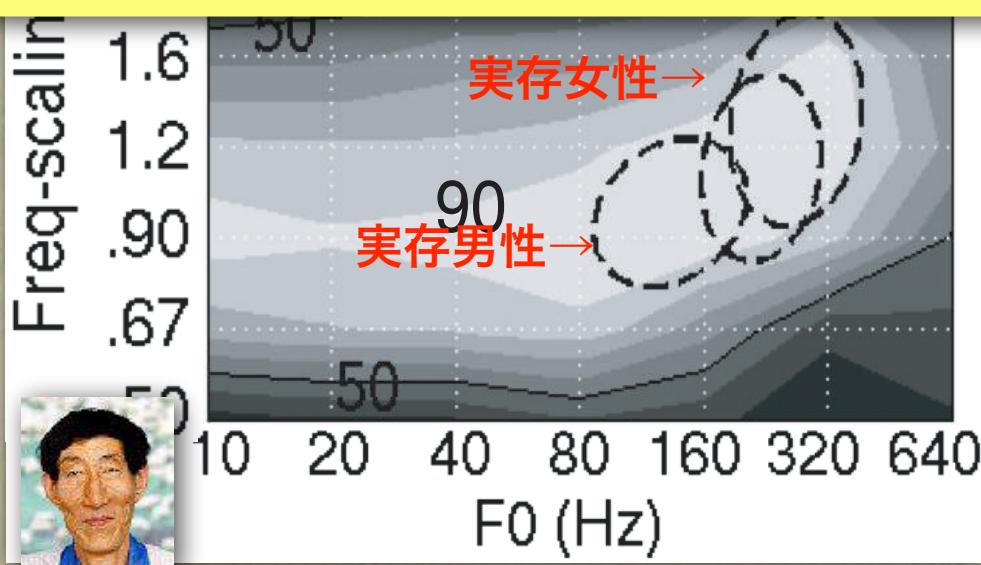
- 孤立母音の同定は可能



こと
とは出来ない。
人などいるのか？

単語同定実験

孤立提示された音を音韻同定する能力は
音声言語運用には不要なのかもしれない



音高の偏差とその認知的不变性

カラオケでキーを上げ下げして曲を聞く[3,4]

The image shows two musical staves. Staff 1 is in C major (no sharps or flats) and staff 2 is in G major (one sharp). Both staves have a common time signature (indicated by 'c'). The notes are primarily eighth notes. In staff 1, the first note is highlighted with a blue oval and the third note with a red square. In staff 2, the second note is highlighted with a blue oval and the fourth note with a red square.

● 絶対音感者（ドレミは**音名**）

● 1 = ソーミソドーラードドソー, 2 = レーシレソーミーソソレー

● 言語化可能な相対音感者（ドレミは**階名**）

● 1 = ソーミソドーラードドソー, 2 = ソーミソドーラードドソー

● 言語化困難な相対音感者（**ラーラ音感者**）

● 1 = ラーラララーラーラララー, 2 = ラーラララーラーラララー

● 異なる音を同一と主張し、同一の音を異なると主張する。

● 各音が持つ基本周波数（絶対量）ではなく、各音が他の音群との
ようなコントラストを持つのか、のみによって決定

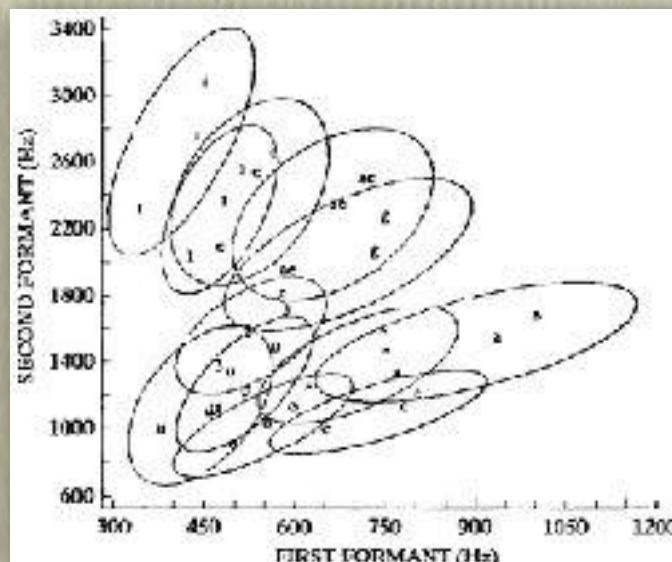
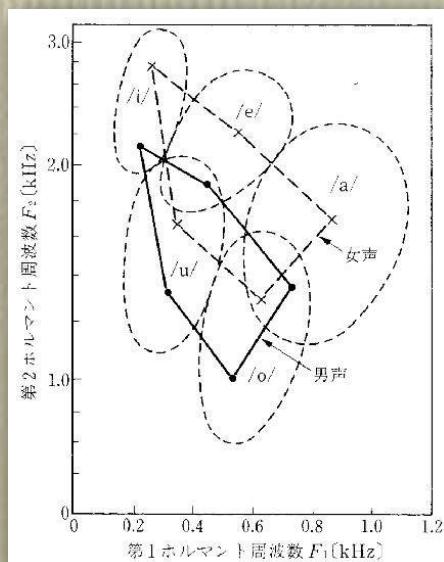
音声の構造的表象／音色の相対音感

言語化困難な相対音感者（ラーラ音感者）

- 次に示すメロディーの3番目の音を覚えて下さい。その後、別のメロディーを提示します。同じ音が出て来たら挙手しなさい。
- メロディーをシンボル列に変換できないので、困難な問い合わせとなる。

言語化困難な音声の相対音感者（幼児的な成人？）

- 次に示す発声の3番目の音を覚えて下さい。その後、別の発声を提示します。同じ音が出て来たら挙手しなさい。
- 発声をシンボル列（音韻列）に変換できなければ、困難な問い合わせとなる



英語圏には十分な教育を受けているが、読み書きに苦労する人が多く存在しなければならない？

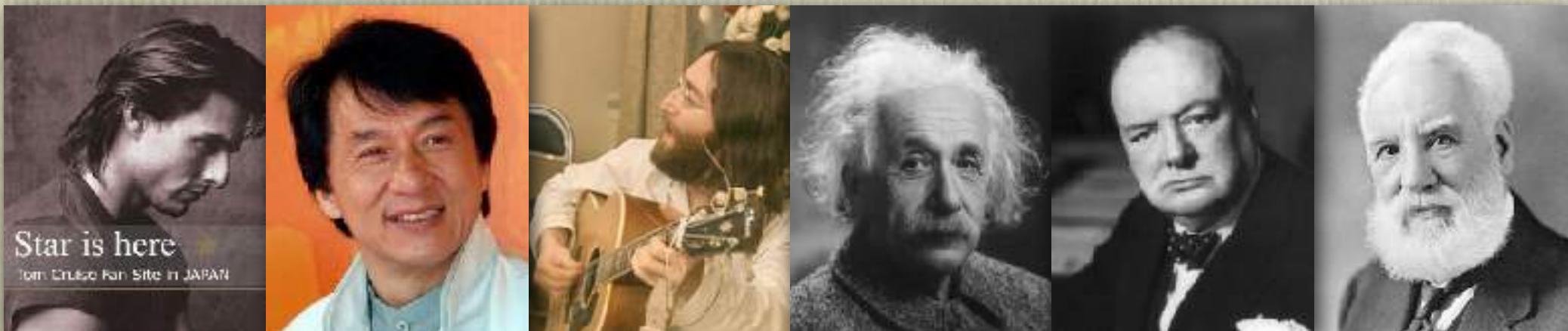
音声の構造的表象／音色の相対音感

言語化困難な相対音感者（ラーラ音感者）

- 次に示すメロディーの3番目の音を覚えて下さい。その後、別のメロディーを提示します。同じ音が出て来たら挙手しなさい。
- メロディーをシンボル列に変換できないので、困難な問い合わせとなる。

言語化困難な音声の相対音感者（幼児的な成人？）

- 次に示す発声の3番目の音を覚えて下さい。その後、別の発話を提示します。同じ音が出て来たら挙手しなさい。
- 発話をシンボル列（音韻列）に変換できなければ、困難な問い合わせとなる



ディスレクシア（読字障害・難読症）

音声の構造

言語化困難な相対

- 次に示すメロディーをメロディーを提示します。
- メロディーをシンボル列

言語化困難な音声

- 次に示す発声の3音を提示します。同じ音
- 発話をシンボル列



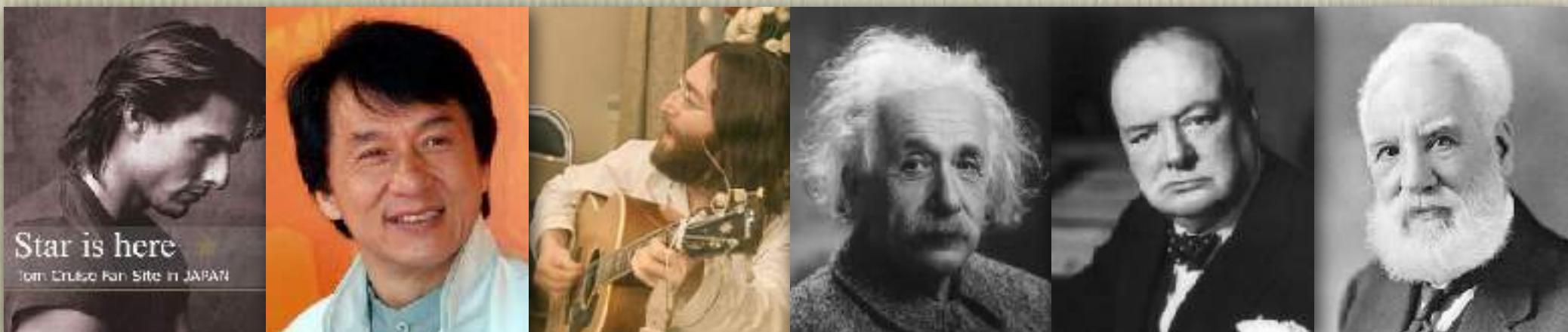
の相対音感

感者)

て下さい。その後、別の文を読んだら挙手しなさい。
これが、困難な問い合わせとなる。

児的な成人?)

い。その後、別の発話を読みなさい。
これが、困難な問い合わせとなる



ディスレクシア（読字障害・難読症）

とある作文

「あ」という声を聞いて母音「あ」と同定する能力は音声言語運用に必要か?

話し言葉の音声

第4章

「あ」という声を聞いて母音「あ」と同定する 能力は音声言語運用に必要か?

——音声認識研究からの一つの提言——

峯松 信明

▼はじめに　～何、この変なタイトル?～

タイトルを見て、多くの読者が首を傾げていることだろう。しかし、十一頁の本記事を読み終えた時に、ほとんどの読者に私の意図は通じるもの、と考えている。そう。「あ」という声を聞いて、それを有限個の音カテゴリーの一つとしての母音「あ」であると同定する能力は、音声言語運用の必要条件ではない。との主張を本稿では展開する(文献1)(文献2)。

そんな馬鹿な、と思われるかもしれない。こんな実験を考えてみよう。身長300cmの巨人と50cmの小人に孤立母音を発声してもらう。通常音声学の教科書には、F₁・F₂

の母音図が出ている(図1参照)。複数の男性/女性のサンプルから、凡そ男性の各母音はこの領域、女性の各母音はこの領域にあるといった図である。フォルマント周波数(共鳴周波数)は声道長に依存するため、身長が50cm、300cmという架空の大人を想定した場合、彼らの母音は、通常知られている領域の外に存在する。そのような母音でも、現在の音声分析・再合成技術を使えば非常に高品質な音声として生成できる。さて、聞いたことのない母音音声を孤立提示されて、読者は同定できるだろうか?

「音声言語は流暢だし雄弁。頭は良いのかもしない。でも何故か本が読めない、手紙が書けない。そういう成人が米国や英国に多かつたりしませんか? えと、教育を受けていないとか、そういう事ではなく、彼らの認知特性として文字言語が何故か難しい……」

「先生、ディスレクシアってご存知なんですか? 特に音韻性のやつ。」

「でいすれ……何ですかそれ?」

「変だな。先生、今、自分でディスレクシアの説明してたじやないですか。」

四一年間の人生の中で、あれほど口をあんぐり開けたことは無い。顎が外れるかと思った。これは実話である。私は彼ら(文献15)の存在を、音声の物理学に基づいて予言していた。

日本語学4月号, p.187-197,
明治書院(2008)

とある反応

音声工学からディスレクシアの仕組みに迫る論文 に、心の底から感動する

1月、当ブログに、あの東京大学からアクセスが集中した日がありました。

なんだろうと思い、リンク元をたどってみると、

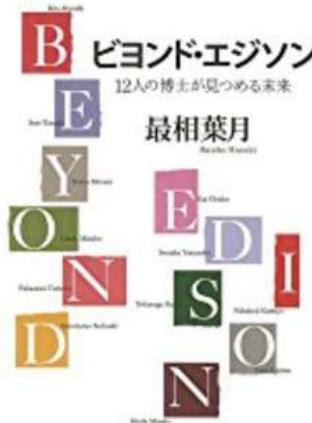
東大工学部・峯松信明教授の「音響音声学」という授業で、
「ディスレクシアであることの利点」が参考資料にあがっていました。

峯松研究室のサイトに、

**～「あ」という声を聞いて母音「あ」と同定する能力は音
声言語に必要か～**

という論文が置いてあったので、読んだところ・・・・

これまで読んだ「ディスレクシアの日本語と英語の出方の差」を説明している
どの論文よりも、激しく腑に落ちるものでした！！



(←『ビヨンド・エジソン』という本に、
同じ内容がより一般向けに書かれています。)

峯松論文では、まず音楽の「絶対音感」と「相
対音感」の違いを説明します。

絶対音感者とは、どんな音を聞いても、音程が
わかる人です・・・①

いわゆる耳コピ（音楽を聴いて譜面に起こす）
ができる人です。

こういう人たちは、音を流れとしてではなく、
孤立的にしかとらえられません（その証拠に、
絶対音感がある人は、カラオケで転調されると

興味深い思考実験を一つ

一卵性双生児が生まれた直後に両親が離婚した・・・

- 一人ずつ引き取られた。
- 彼らは10年後どんな発音をしているのだろうか？

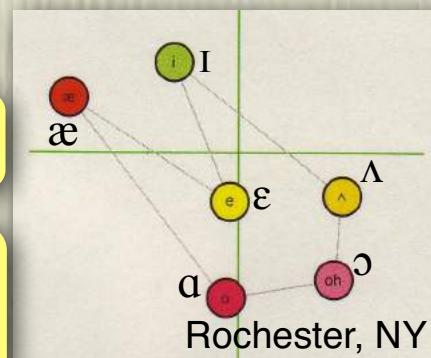
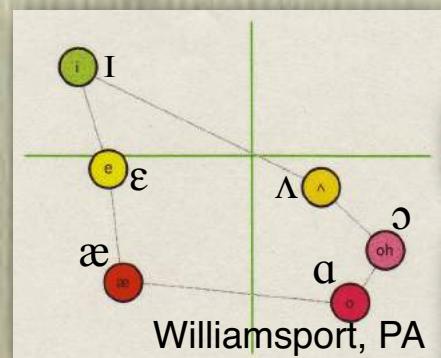


声道形状の性差＝音色の差異



方言差＝音色の差異

九官鳥は音を真似る
幼児は音の体系を真似る



興味深い思考実験を一つ

一卵性双生児が生まれた直後に両親が離婚した・・・

- 一人ずつ引き取られた。
- 彼らは10年後どんな発音をしているのだろうか？

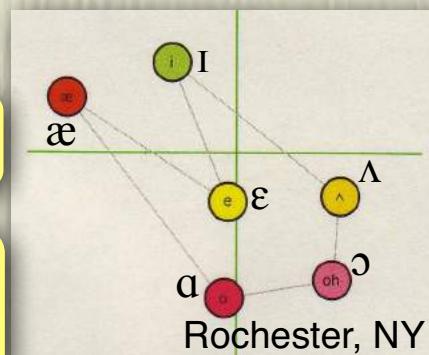
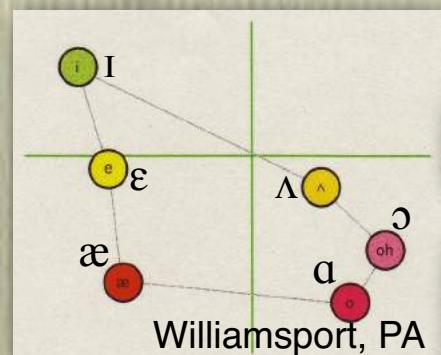


声道形状の性差＝音色の差異



方言差＝音色の差異

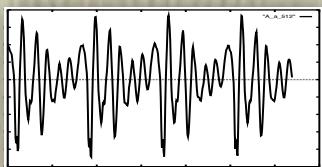
幼児が学ぶものを学ぶ機械
幼児が無視するものは無視する機械



その情報を運ぶ媒体・音響特徴量

二段階の分離に基づく特徴量抽出

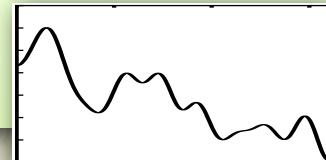
Insensitivity to pitch differences



Insensitivity to phase differences

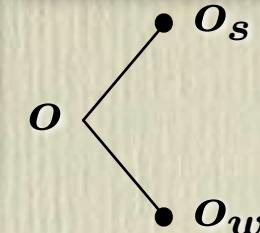
speech waveforms
phase characteristics

amplitude characteristics



source characteristics

filter characteristics



スペクトル包絡(O)は何を運ぶのか？

言・パラ言・非言

真の音声の統計的モデル～波形の統計的モデル～

不特定話者・不特定基本周波数・不特定位相の音響モデル

見たくないものは全て「確率の定義」で集めて隠してしまおう。

$$P(o|w) \approx \sum_{s,h,p} P(o|w, s, h, p) P(s) P(h) P(p)$$

s : speaker, h : harmonics, p : phase

一般的な解決策：各手法の組み合わせ

最終的に性能を最大化する組み合わせを追求する。

幼児の言語獲得と音声模倣

音声模倣＝親の発声行為を子が積極的に模倣する行為

- これを通して幼児は言語を獲得する[7]
- 動物学的には非常に稀な行為。靈長類では人間だけ[8]
- 他の動物では小鳥、クジラ、イルカくらいか[10]

動物の模倣＝声帯模写、ヒトの音声模倣≠声帯模写

- 九官鳥の音声模倣[9]
 - 車、ドア、椅子、犬、猫、音を真似る。人の声も音でしかない。
 - 良い九官鳥を聞くと、飼い主が分かる。
- 幼児の音声模倣
 - 動物学的には奇妙な模倣行為[10]
 - いくら良い子でも、声から父親を割り出せずにお巡りさんは困る。



幼児は親の声の何を真似ているのか？

音



~~お+は+よ+う~~

言



音色の偏差とその認知的不变性

色み・音高の恒常・不变的認知

- コントラスト情報に基づく処理が重要
- コントラスト群から成る全体的パターン処理が要素同定を可能に



音色の恒常・不变的認知

- コントラスト情報に基づく処理が重要
- コントラスト群から成る全体的パターン処理が要素同定を可能に



本発表の流れ

刺激の物理的多様性とその認知的不变性

- 見え／色み／音高の多様性と自然・進化が編み出した解決方法

音声の物理的多様性とその認知的不变性

- 音色の多様性と工学者が編み出した解決方法

音声の構造的表象とそれに関する様々な考察

- 常識を覆すことで、違和感の解消を試みてみる。

音声の構造的表象と数学的表現と技術的実装

- 体格・性別に不变な音声波形・スペクトルの表現とは？

音声の構造的表象を用いた音声アプリケーション

- 音声認識、音声合成、発音分析、etc

音声の構造的表象の言語学的妥当性

- 何故、こうしてこなかったのか？ 観測技術の功罪？

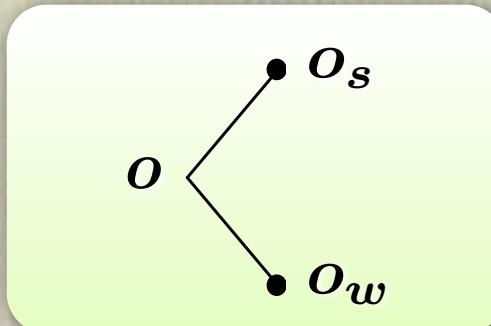
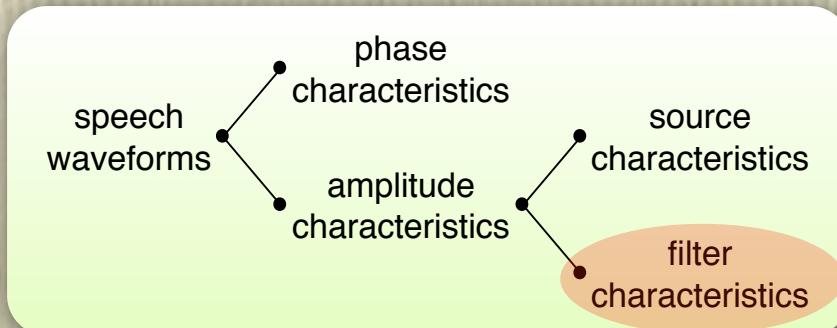
音声模倣とその技術的実装

まねだ聖子・松田聖子・神田沙也加



学習話者そっくりの声色で読み上げる技術＝音声合成

- Blizzard Challengeでは学習話者の個人性の再現も採点対象[7]
- 黒柳徹子を使えば、黒柳徹子の声になる。



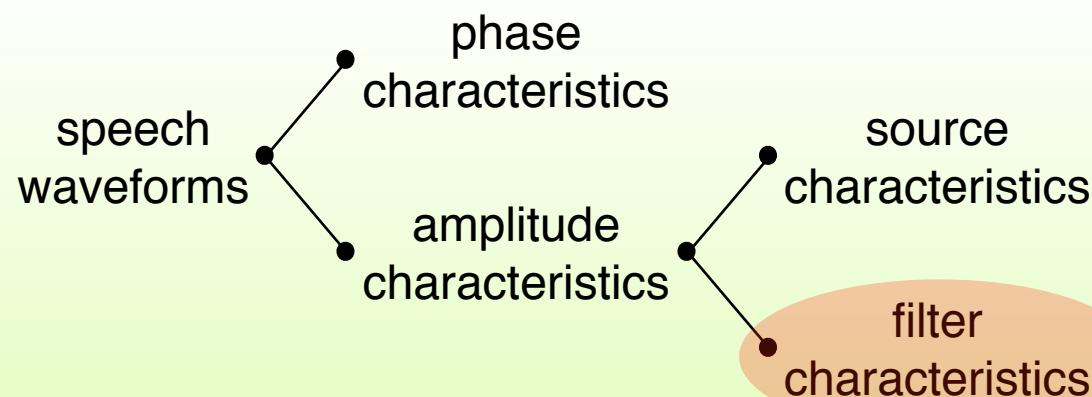
他者の音声の模倣＝声帯模写となる方々

(重度) 自閉症者に見られる音声模倣＝声帯模写

- 七色の声を持つ中村メイ子の声をそっくり真似る[8]
- 相手そっくりの声を模倣する[11]
- 車, 電車, などの音響音の模倣[12]
- 移調してしまうと, その曲だと認識してくれない[8]
- 母親の声は理解できるが, それ以外は難しい[13]

「言語＋非言語」が同居したままの音声の捉え方

- 音声コミュニケーションに困難を抱える場合が多い[18]



とある自閉症者の訴え

とある自閉症者（アスペルガー症候群）の手記

- 「発達障害当事者研究」（綾屋紗月、熊谷晋一郎著）[9]
- 「外国語の発音練習」「カラオケ」が難しい。
- どうしても、先生／職業歌手の声帯模写をしてしまう。
- みんなの真似は真似じゃない。だって、声色違うじゃない。
- 「自分の声でいいんだよ」と言われるけど。
- 「そもそも、私の声って何なの？いつの私の声のこと言ってるの？」



とある自閉症者の訴え

・ものまね歌合戦って、面白いですか？

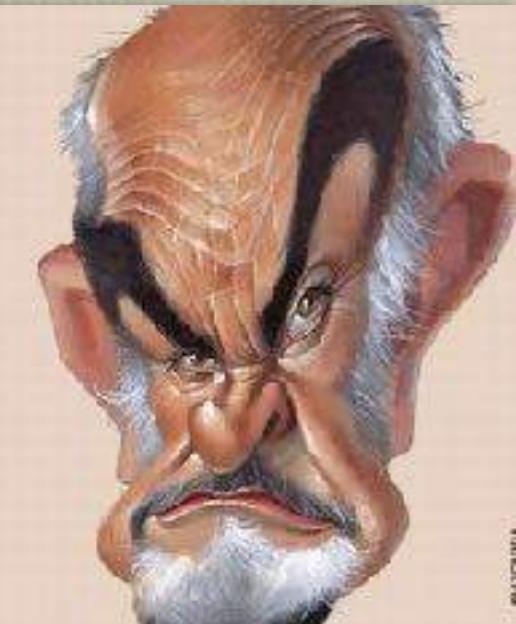
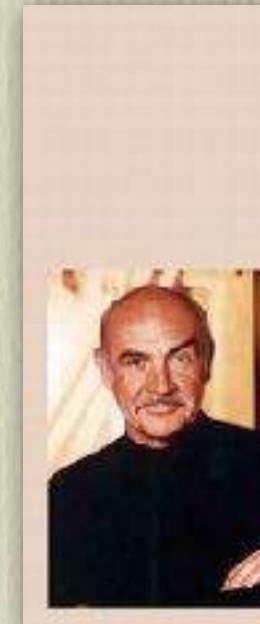
- 何が面白いのか、さっぱり理解できません。

・この似顔絵、似てるって分かりますか？

- 分かりません。こっちの方が似てると思いますが。

・綾屋さんの言語活動の主メディア

- 手話と文字言語



とある自閉症者の訴え

・ものまね歌合戦って、面白いですか？

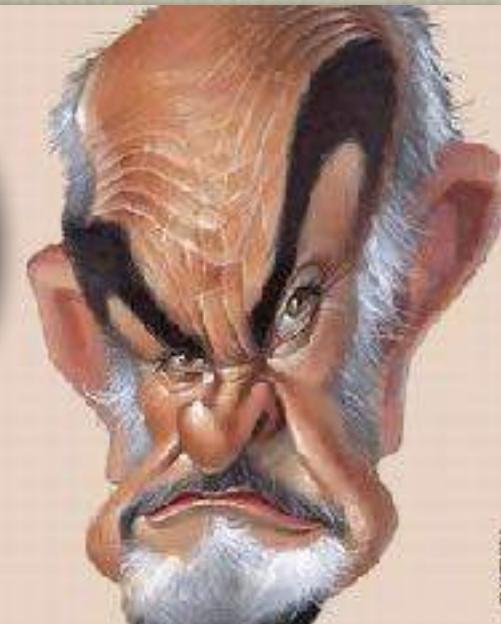
○ 何が面白いのか、さっぱり理解できません。

・この似顔絵、似てるって分かりますか？

○ 分かりません。こっちの方が似てると思いませんが。

・綾屋さんの言語活動の主メディア

○ 手話と文字言語



自閉症の方々に見られる症状

とある web より

自閉症の特徴の強みと弱み

強み→① 具体的なことをよく理解し、記憶する。

② 目で見て認知したり記憶する視覚的な認識・記憶力がいい。

③ 決まったパターンのくり返しに強い。

④ 好きなことへの集中力。

弱み→① 曖昧なこと、抽象的なことに弱い。

(一つひとつ的情報はキャッチしていても、それらの相互関係がつかみにくい。
目に見えないこと、経験していないことを想像することが難しい。)

② 時間の見通しをたてるのが苦手。

(物事の終わりがわかりにくい。いつもの流れが変更されると、わからなくなる。)

③ 状況を認識すること。

(人の表情、しぐさ雰囲気などが理解しにくく、人の感情がわかりにくい。
怒られているのに嬉しがったり、ほめられているのに知らん顔など・・・。)

④ 話し言葉への理解、自分からのコミュニケーションが難しい。

(言葉が出てもオウム返しになるなど。)

⑤ 感覚刺激に対して特異な反応をする。

(感覚刺激に対して過敏だったり鈍感だったりする。感覚刺激が一度にたくさん入りすぎてしまう。特定の感覚刺激に苦痛を感じる。)

幼児の言語獲得と音声模倣

音声模倣＝親の発声行為を子が積極的に模倣する行為

- これを通して幼児は言語を獲得する[7]
- 動物学的には非常に稀な行為。靈長類では人間だけ[8]
- 他の動物では小鳥、クジラ、イルカくらいか[10]

動物の模倣＝声帯模写、ヒトの音声模倣≠声帯模写

- 九官鳥の音声模倣[9]
 - 車、ドア、椅子、犬、猫、音を真似る。人の声も音でしかない。
 - 良い九官鳥を聞くと、飼い主が分かる。
- 幼児の音声模倣
 - 動物学的には奇妙な模倣行為[10]
 - いくら良い子でも、声から父親を割り出せずにお巡りさんは困る。



自閉症・・絶対的記憶・・動物・・??

動物の情報処理と自閉症者情報処理

- Dr. Temple Grandin (アスペルガー&動物学者)

- 「動物感覺」[17]

- 動物と自閉症者情報処理類似性を主張

- 局所的／具体的／実体的 \longleftrightarrow 全体的／抽象的／概念的



(定型発達を遂げた) 人間が有する特異的な能力?

- 音を用いた情報伝達において、情報同一性は如何に確保できる？

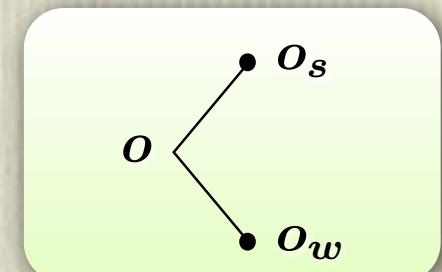
- 動物、重度自閉症者、現在の音声認識（情報分離が困難）

- 情報（メッセージ）の同一性 = 音響的特徴 (o) の同一性

- 定型発達を遂げた人間

- 情報の同一性を確保するために、音の同一性は必要でなくなった種では、音のどの側面は同一なのか？

- 語全体の語形・音形・ゲシュタルト



生物が獲得した静的バイアス除去術



音高の恒常的・不变的認知はどこまで遡れるのか？[6]

1

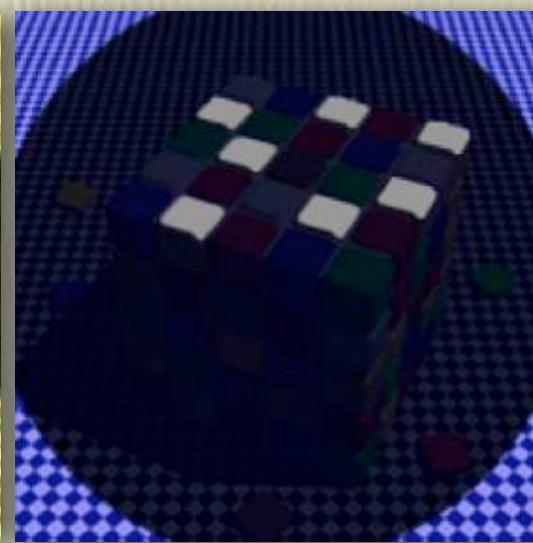
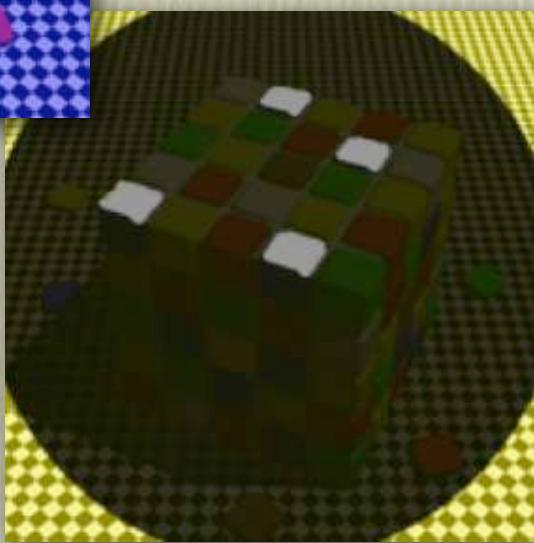
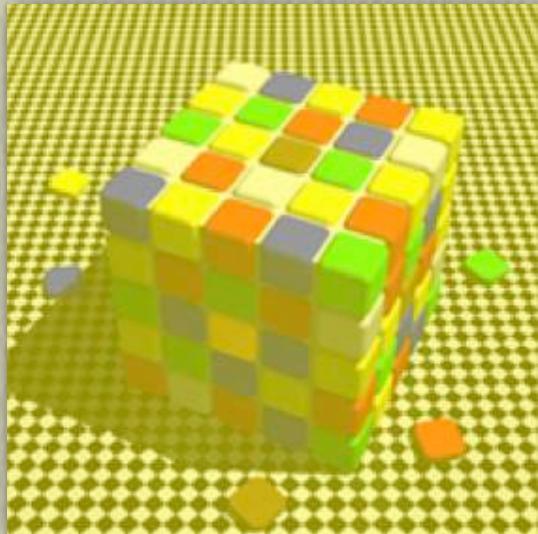
2

$$1 = 2$$



生物が獲得した静的バイアス除去術

色の恒常的・不变的認知はどこまで遡れるのか？[5]



彼女と会ってきました



第六章 言葉の不思議を探究する ······ 119

音声工学者・峯松信明と動物科学者テンプル・グランティンの自閉症報告

彼女と会ってきました



絶対音感
ABSOLUTE PITCH

最相葉月
SASEKO TSUYOSHI

五嶋みどり、五嶋龍、
千住真理子、矢野顕子……

一流音楽家、科学者ら200人以上に証言を求める、
「絶対音感」の謎を探り
音楽の魅力の本質に迫る
ノンフィクション! 文庫決定版

新潮文庫
の
新刊

第六章 言葉の不思議を探究する ······ 119

音声工学者・峯松信明と動物科学者テンプル・グランティンの自閉症報告

こそっと隠してあります

gavo.t.u-tokyo.ac.jp

最相葉月: わたしの偉人伝 第4回 言葉の不思議を探求する | ボブアビーチ

イラスト: 嶋井沙奈美

わたしの偉人伝

最相葉月

人は人に騙される——。子どもの頃に「エジソン伝」「野口英世伝」などを読み、医学や科学の道を志した人も多いのではないでしょうか。本連載では、各分野で活躍している方に、感銘を受けた伝記・評伝を学びいただきながら、現在のお仕事への想いや研究内容について伺っていきます。

企画監修

第4回 言葉の不思議を探求する

音声工学者・峯松信明さんと動物科学者テンブル・グランティン

峯松信明さんは、音声工学を研究する中で自閉症と出会い、彼らの音声認知について育声の物理的な側面から考察しています。峯松さんが衝撃を受けた『動物感覺』は、自閉症者であり動物科学者であるテンブル・グランティンが著したノンフィクション。今回は、グランティンと実際に会ったときのエピソードや音声認識研究について伺いました。

言葉とは何だろう。それはおそらく、人間とは何かを問うのと同じだけの時間、問われ続けてきた根源的な問いだろう。なぜ人間だけが言葉をもつのか、なぜ子どもは文法など何も知らないから言葉を話し始めるのか、言語の起源は何か……等々、言葉にまつわる疑問は尽きない。哲学や心理学、文化人類学、言語学など多様な角度から研究されてきたテーマである言葉に、工学と物理学の観点からアプローチしているのが、東京大学大学院工学系研究科准教授の峯松信明さんである。柏キャンパスから本郷キャンパスの工学部2号館に移転中の実験室には、運びこまれたばかりのコンピュータがカバーをかけられたまま並ぶ。峯松さんの研究室もまだ真っ新しい気もない状態だが、研究は大きな展開を見せていく。

自閉症への関心

本当の音声の絶対音感者？

とある自閉症児が書いた本



僕はお母さんの言うことならすべてわかります。それは、第1に安心感、第2に言葉のリズムや高低が良くわかっていること、第3に話の予測がつきやすいためでしょう。

どこにいてもどんなときでも、僕がわかる言葉は、お母さんだけです。
僕は、どうして今まで言葉が理解できないのか、わかりませんでした。他のみんなが指示されたことにすぐに反応できて、その通りに動けることが不思議でした。
僕には聞こえないのです。
音は聞こえているけれど、意味になつて頭の中に入つてこないのです。話しているのが本人だとわかれば、慣れれば言つていることはわかります。でも、同じ人でも場所や状況が違うと、その人だということがわからないのです。

興味深いサイト

絶対音感ある人に30の質問

http://www.100q.net/100/question.cgi?que_no=51

The screenshot shows a web browser window with the URL www.100q.net/100/question.cgi?que_no=51. The title bar says "絶対音感ある人に30の質問". The main content area displays a list of 30 questions, each with a question number, name, and answer.

[Home]	[ReLoad]
絶対音感ある人に30の質問	
[715] heyjoe	---
[714] レイン	---
[713] Julio	---
[712] vista	---
[711] ざるそば	---
[710] ただは	---
[709] 優々	☒ --
[708] あいり	---
[707] 唯。	---
[706] あじさい	---
[705] ろはん	---
[704] mao	---
[703] 岡崎汐	---
[702] とも	---

Q1 氏名前と、この質問の回答日を教えてください。
vistaです。2014.04.19

Q2 年齢・性別をお願いします。
16。女。

Q3 今のお仕事を教えてください。
学生。

Q4 初めて音楽の手ほどきを受けたのは何歳?そのときの楽器は?
4歳のときにピアノを。

Q5 あなたの絶対音感は先天性?それとも後天性?
たぶん先天性

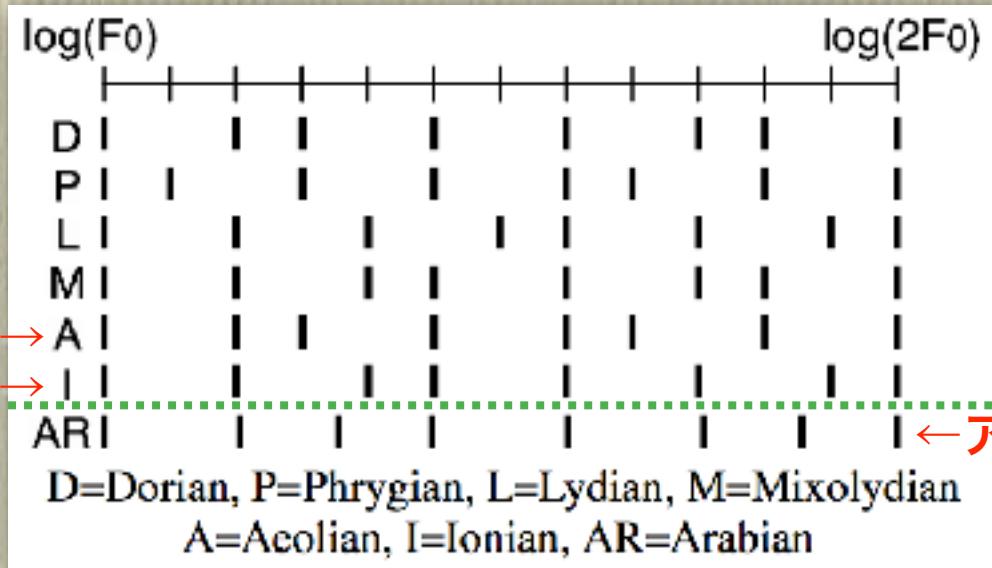
Q6 ピアノのA(ラ)を44Hzでとらえていますか?
441Hz

Q7 暗譜は得意ですか?
はい

Q8 移調は得意ですか?
大好き

音声の構造的表象／音色の相対音感

音楽における調不变の音配置とその変種

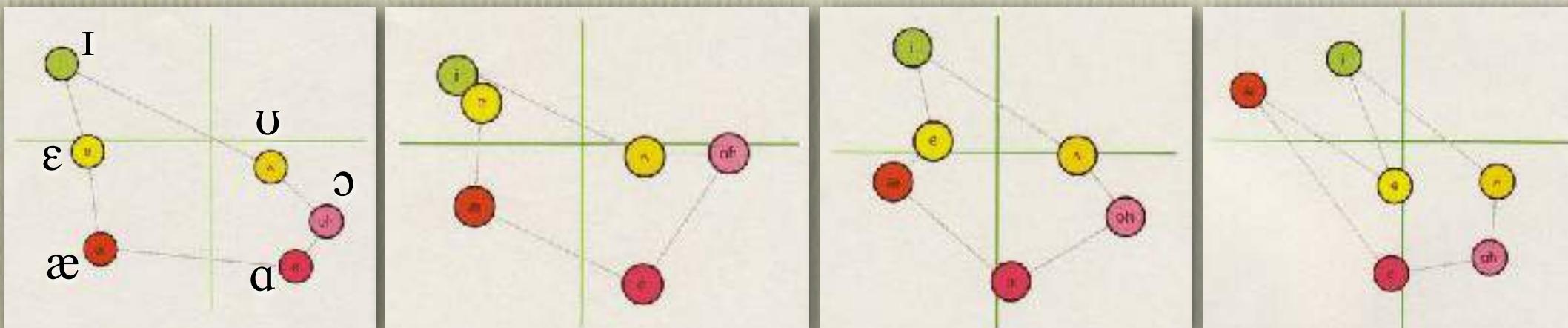


- 西洋音楽 = 5全音 + 2半音
- 種々の配置 = 教会音楽
- 民族音楽には半音以外の配置

أهلاً وسهلاً



音声における話者不变の音配置とその変種 = 欧米の方言



Williamsport, PA

Chicago, IL

Ann Arbor, MI

Rochester, NY

なるほどね・・・



自閉症と音声

第5回「自閉症と音声」研究会 開催案内

12月4th, 2017 | Published in [開催情報](#)

主催：電子情報通信学会 発達障害支援研究会（ADD）

共催：東京大学 奉松・瀬藤研音声技術セミナー

開催趣旨：

自閉症当事者や関係者と音声研究者など理系研究者の間の相互交流の場として、下記の研究会（第5回）を計画しています。専門的議論をするよりも、多領域の間で疑問や注文など、誤解を恐れずに出し合い、相互理解を深め、将来の協力の芽となることを目指しています。

今回は話題提供として自閉症（自閉スペクトラム症）と音声発話などの現象に関する話題に加えて、発達障害の支援技術について計4件お話しいただき、会場からの活発な発言を期待したいと思います。

市川 紘（ADD顧問、千葉大・名誉教授/早大・招聘研究員/工学院大・客員研究員）

記

日時：2018年1月28日（日）

13時～18時

会場：東京人学本郷キャンパス 工学部2号館10階 電気系会議室5

アクセスマップ

http://www.u-tokyo.ac.jp/campusmap/map01_02_j.html

工学部2号館

http://www.u-tokyo.ac.jp/campusmap/cam01_04_03_j.html

参加費：無料です。事前登録など必要ありません。直接会場にお越しください。

話題提供：

1. 自閉症は津軽弁を話さない

松本敏治（教育心理支援教室・研究所 ガジュマルつがる）

<http://add.shimane-u.ac.jp>

2. 自閉症スペクトラム児のフィクショナルナラティブにおける発話特徴の検討

ここで、シラバスを再掲

2017年度 人文社会系研究科 21170104 音響音声学（1） 峰松 信明

複数型

本授業では高校で物理を履修しなかった学生を対象に、音声の物理的・音響的側面について分かり易く解説する。音声は音、即ち、空気（酸素・窒素・二酸化炭素など）の振動現象でしかない。しかし、その振動現象を鼓膜が捉えると、言語メッセージ、意図、感情、更には話者の健康状態など、様々な情報を我々は知覚できる。一体、空気振動のどこにこれらの豊富な情報が隠れているのだろうか？

音響音声学（1）では、音の基礎物理から始め、音声を音響的に眺めるために必要な基礎知識を提供すると共に、音刺激に対するインターフェースである聴覚の処理についても学ぶ。

音響音声学（2）では、スマホで有名になった音声認識（音声テキスト変換）や音声合成（テキスト音声変換）についても、その基礎知識を提供する。その後、言語獲得、外国語学習、言語障害、更には言語の起源に関する様々な話題も提供する。音声の音響的側面についての知識が身に付くと、これら様々な言語現象に対して、従来とは違った視点で議論を展開できる可能性があることを示す。

なお、音響音声学（1）、（2）で通年の授業となるが、年明けてからの5コマが一番面白い講義となるはずである。

（1）は文系学生でも十分理解できる内容だと自負している。（2）の技術的な内容をなんとか（概要だけでも）理解できれば、一番面白い最後の5コマに辿り着ける、そういう通年授業の構成となっている。是非頑張って欲しい。

本発表の流れ

刺激の物理的多様性とその認知的不变性

- 見え／色み／音高の多様性と自然・進化が編み出した解決方法

音声の物理的多様性とその認知的不变性

- 音色の多様性と工学者が編み出した解決方法

音声の構造的表象とそれに関する様々な考察

- 常識を覆すことで、違和感の解消を試みてみる。

音声の構造的表象と数学的表現と技術的実装

- 体格・性別に不变な音声波形・スペクトルの表現とは？

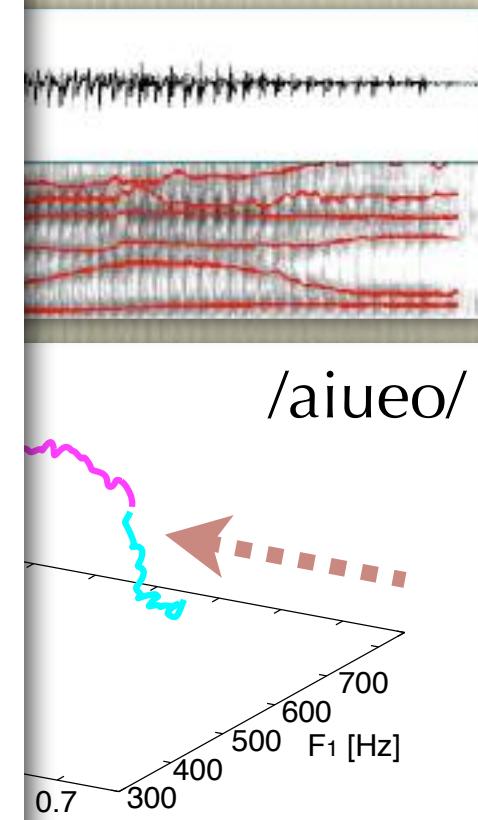
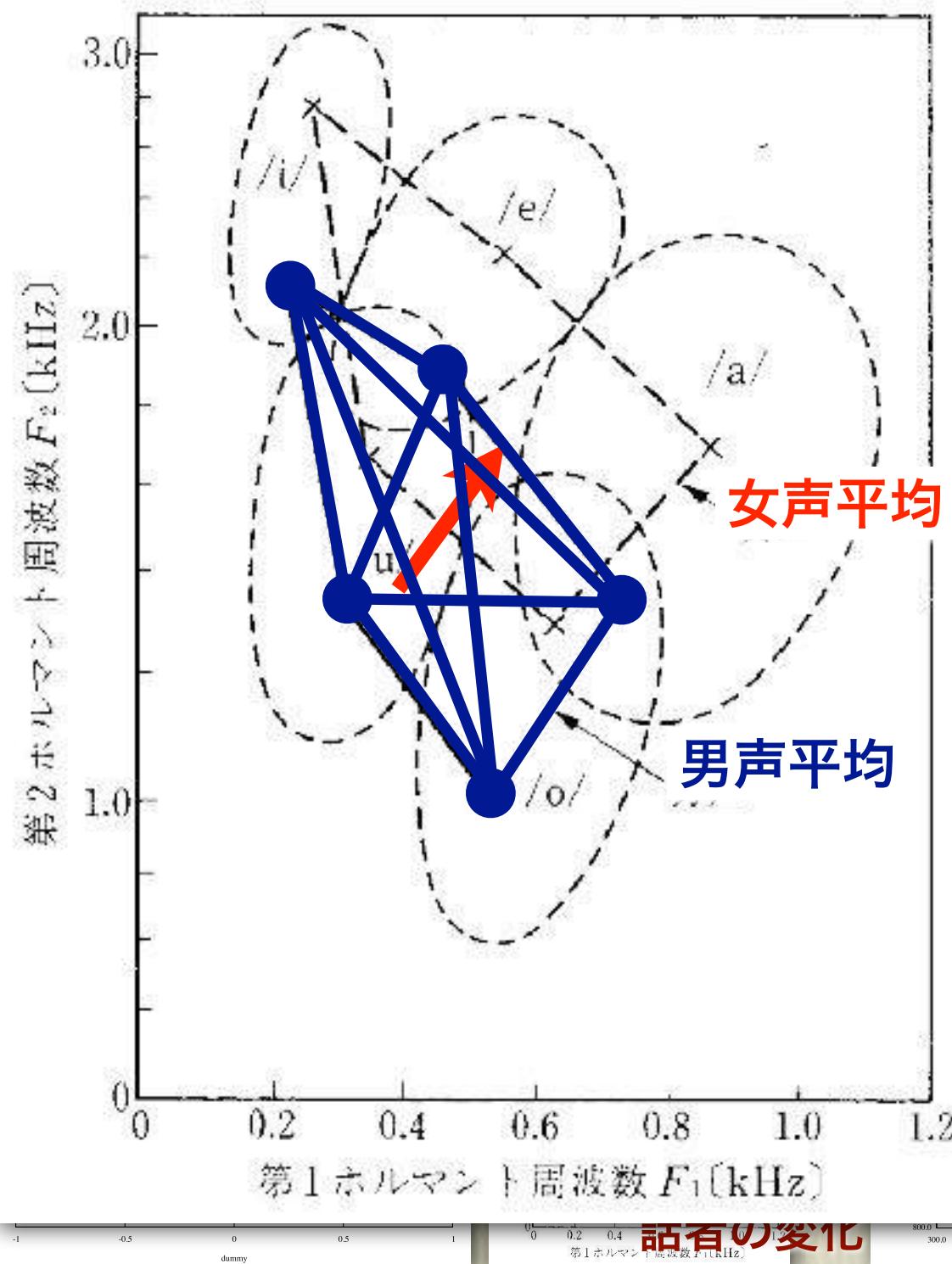
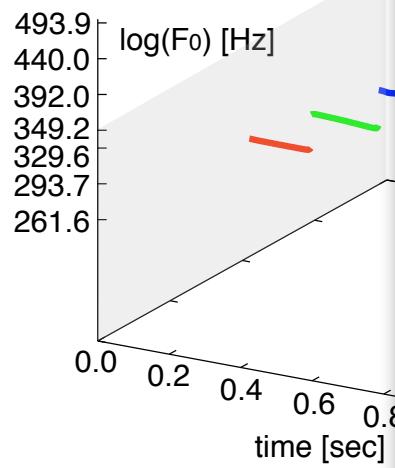
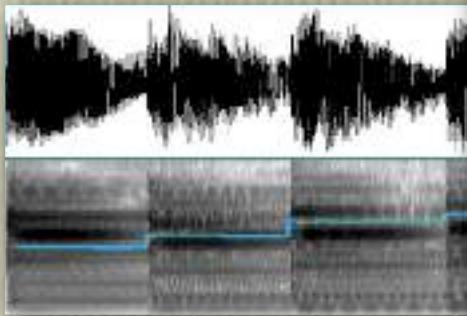
音声の構造的表象を用いた音声アプリケーション

- 音声認識、音声合成、発音分析、etc

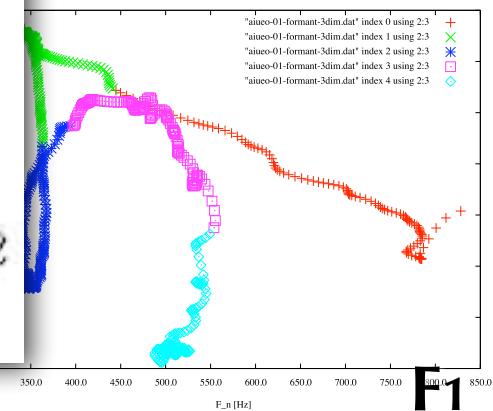
音声の構造的表象の言語学的妥当性

- 何故、こうしてこなかったのか？ 観測技術の功罪？

音声の構造的音色 / 音色の相対音感

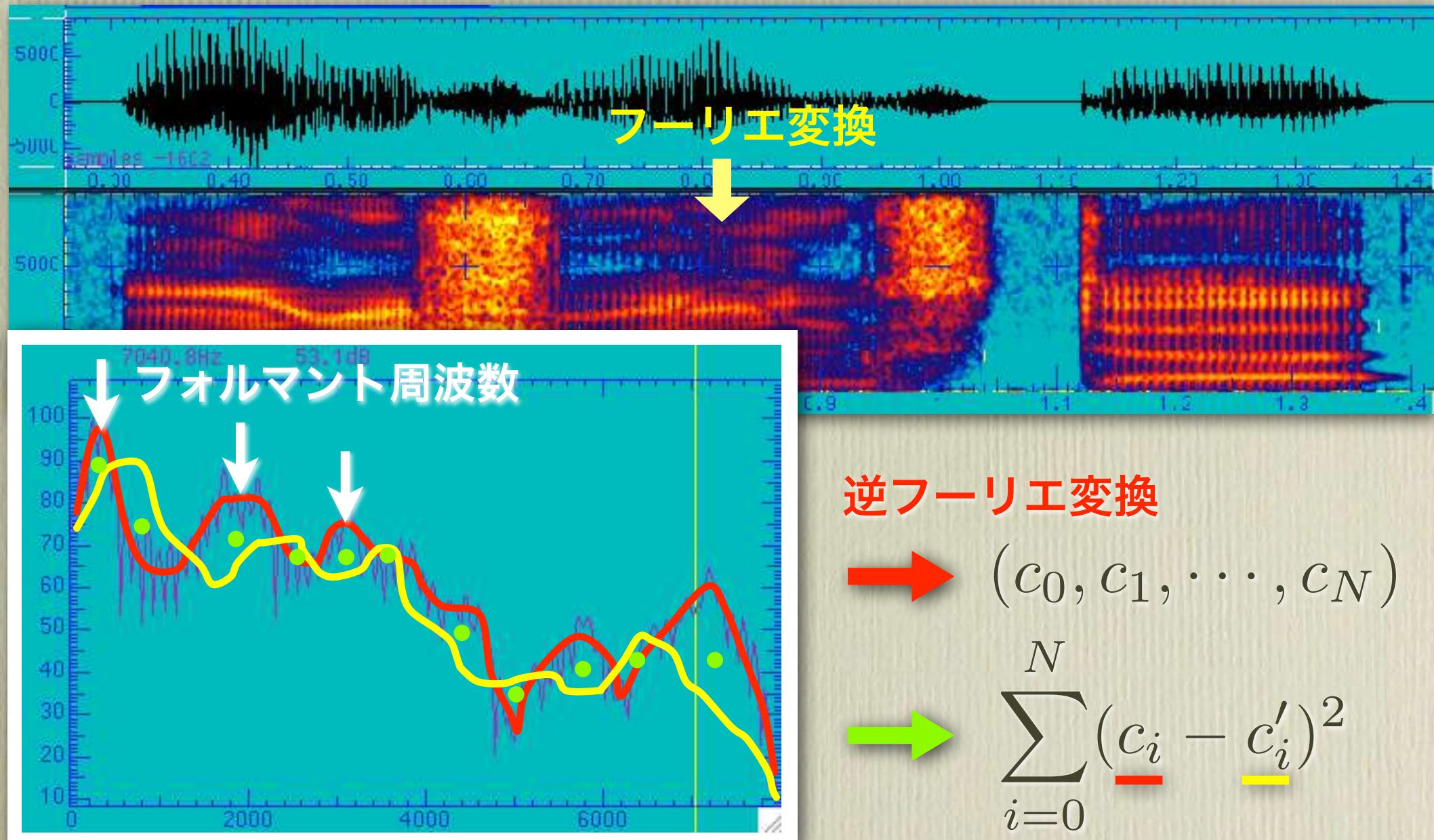


的変化パターン



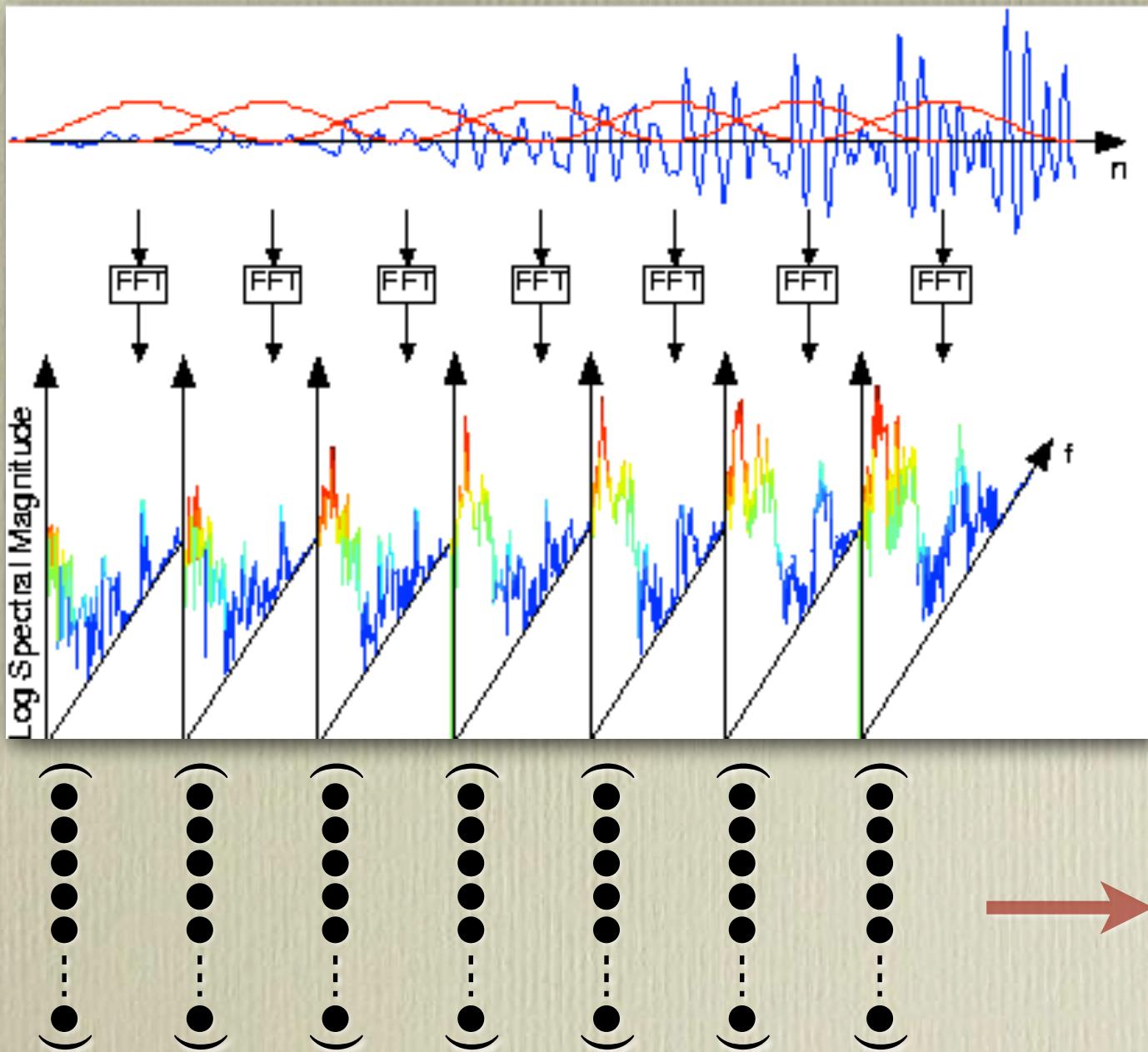
スペクトル包絡の効率的なベクトル表現

音声波形 → スペクトラム → ケプストラム



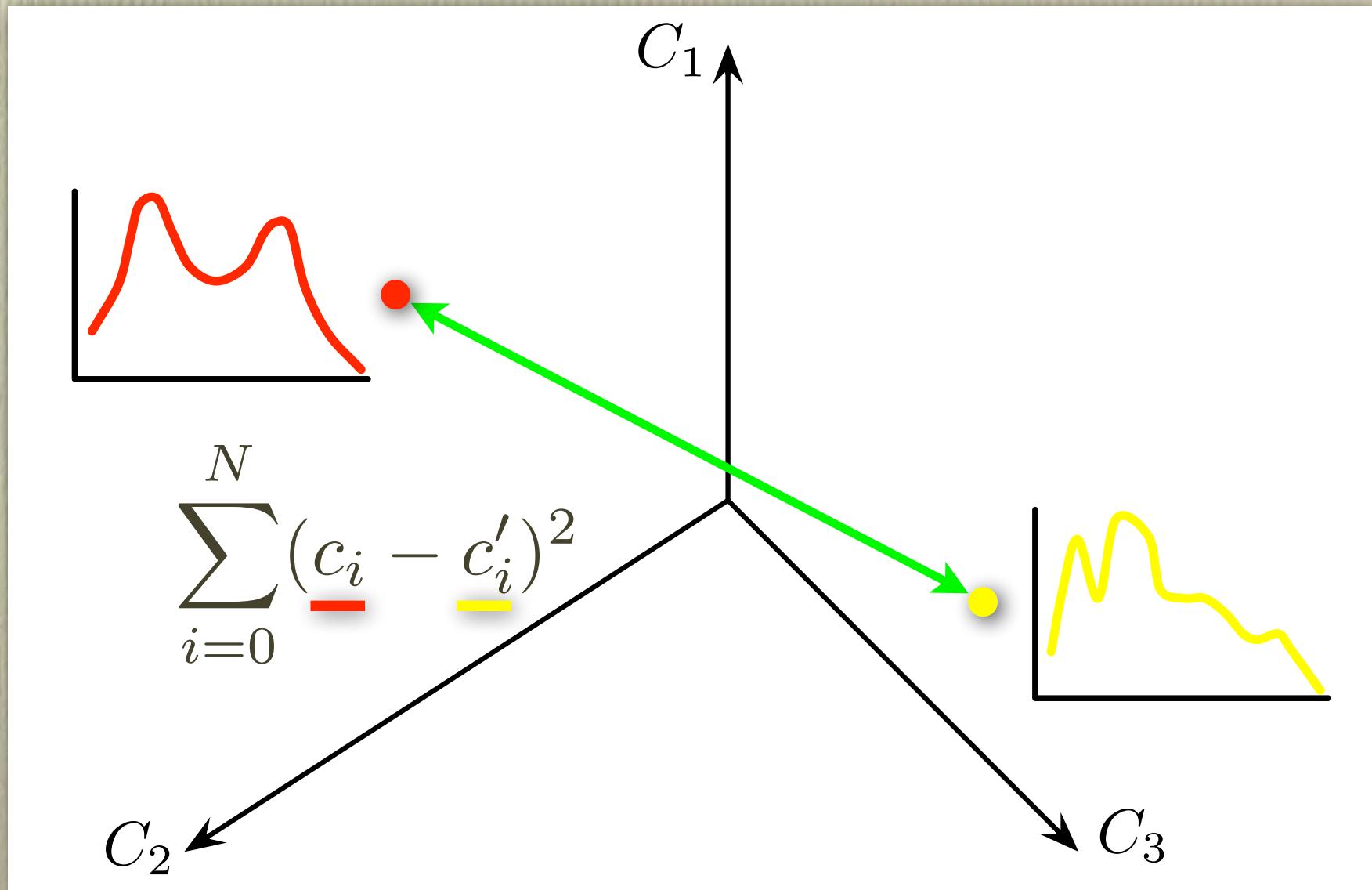
スペクトル包絡の効率的なベクトル表現

音声波形 → スペクトラム → ケプストラム



スペクトル包絡の効率的なベクトル表現

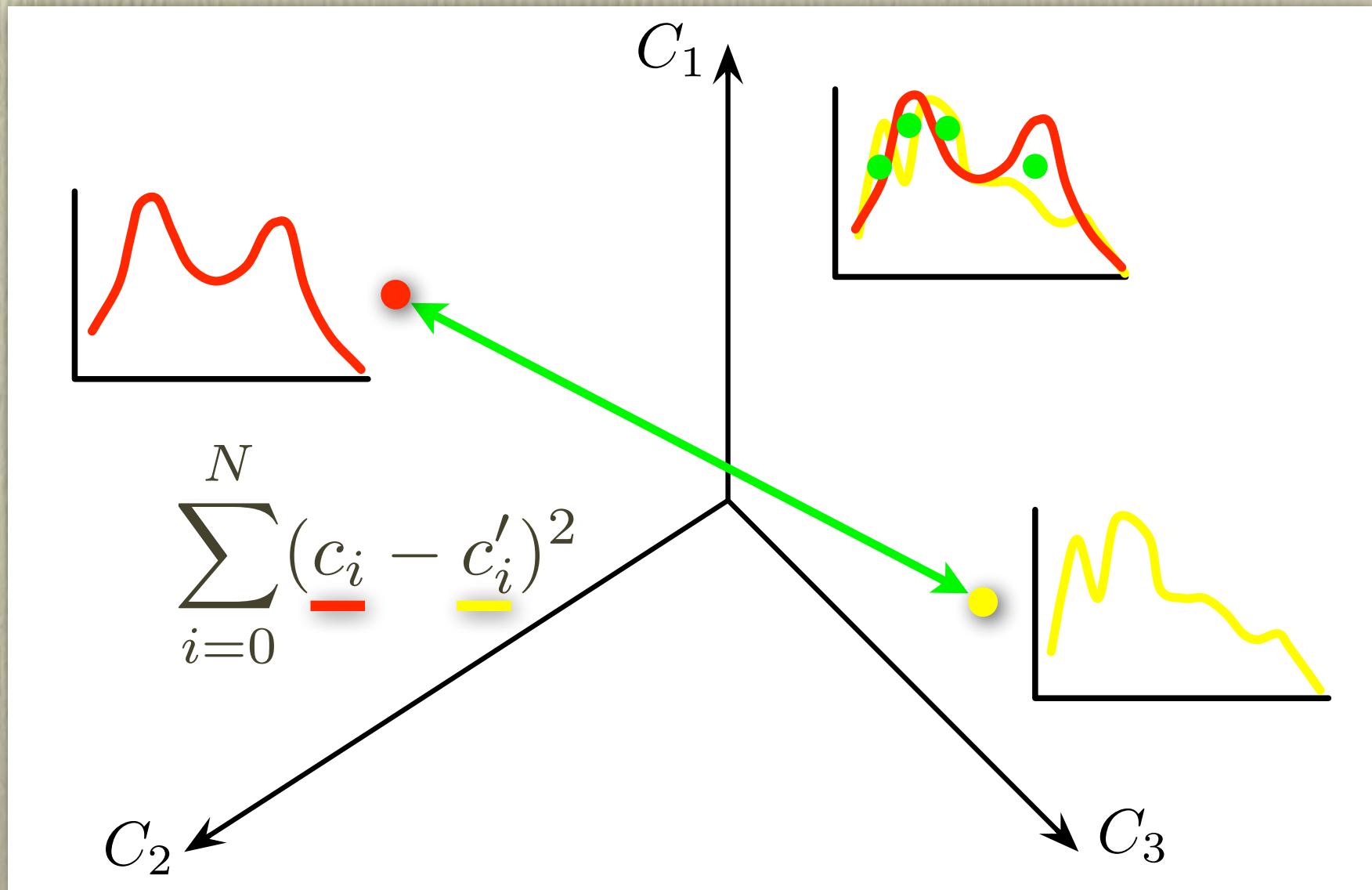
ケプストラム空間における「点」と「点間距離」



点=スペクトル包絡, 点間距離=スペクトル間差異

スペクトル包絡の効率的なベクトル表現

ケプストラム空間における「点」と「点間距離」

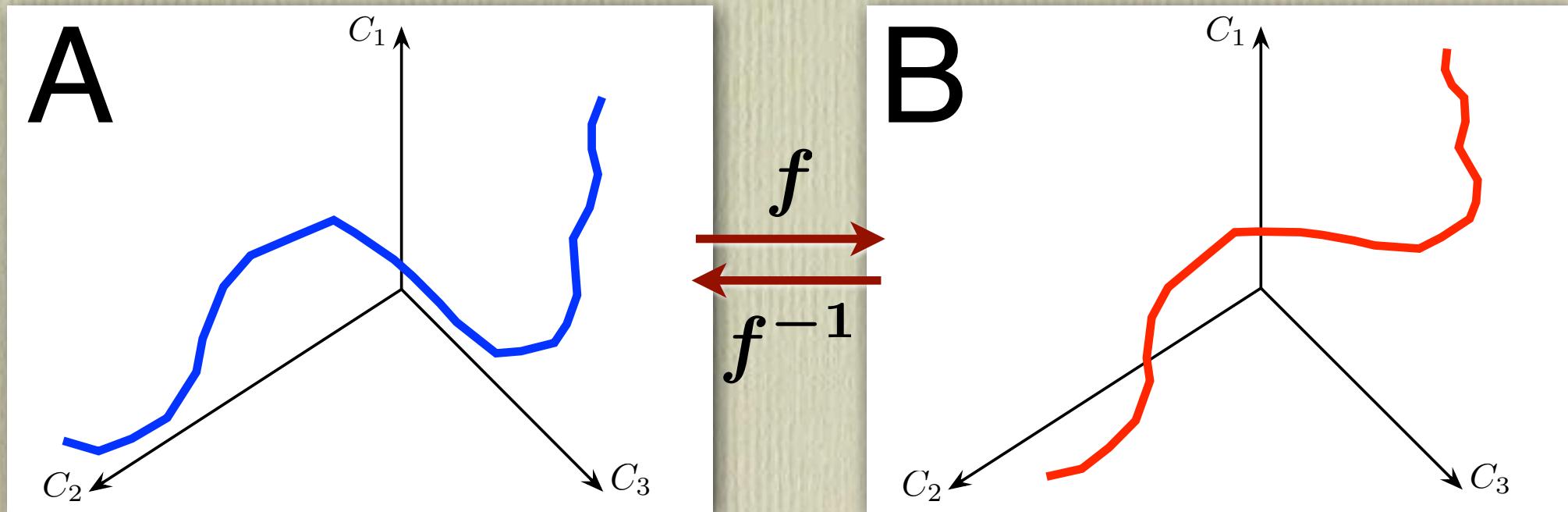


点=スペクトル包絡、点間距離=スペクトル間差異

変換不变な音響量の数学的探求

話者の違い = 空間写像 (話者・声質変換)

- 話者 A の音響空間 \leftrightarrow 話者 B の音響空間



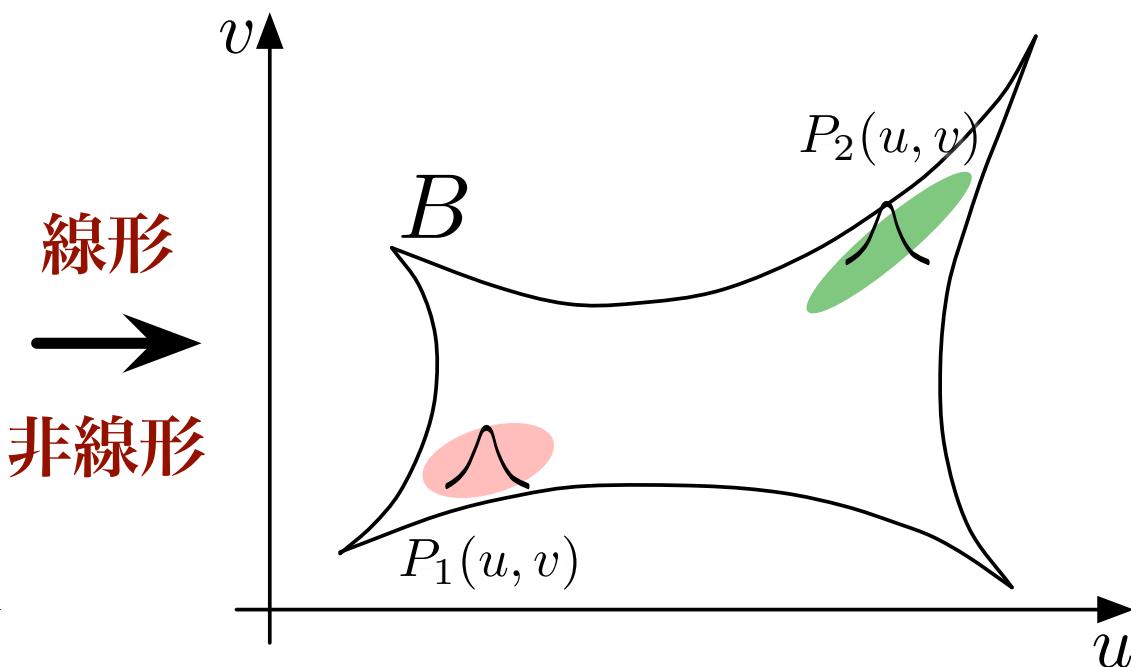
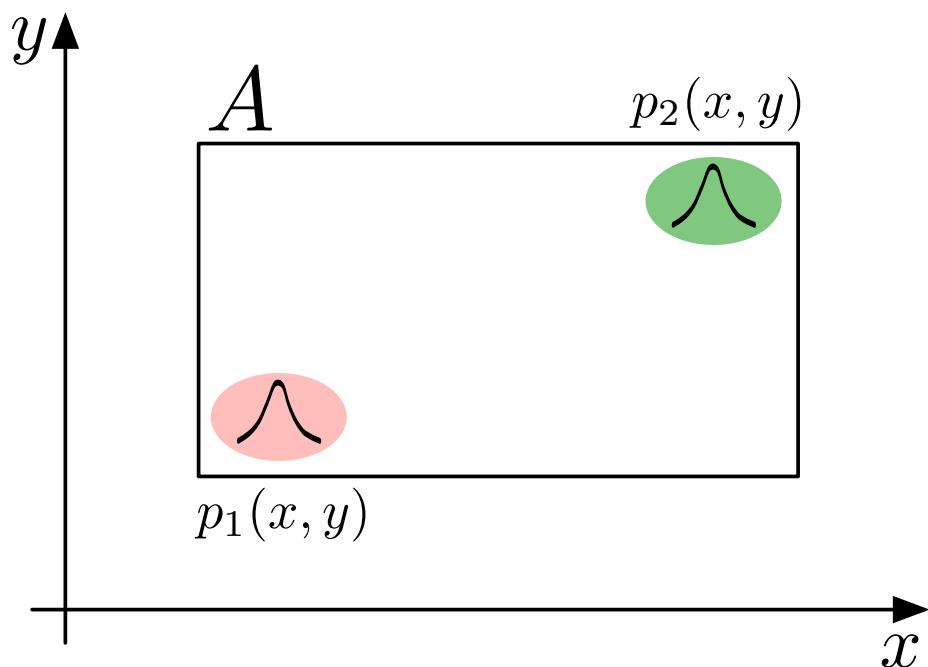
話者 A の声を 65 億全ての話者の声へ変形する

- 65 億 \times 65 億 の写像関数が定義可能
- 話者不变のコントラスト = 写像不变のコントラスト
- 任意の写像に対して不变なるコントラスト量は存在するのか？

変換不变な音響量の数学的探求

二人の話者空間（一对一対応）における不变音響量

- 音響事象を点ではなく、分布として表現する。
- 空間Aの分布 p は空間Bの分布 P へと写像される。
- p と P は異なる物理特性を持つ（[あ] と [あ]）
- 両空間において不变な物理量はどこに？不变コントラストは存在する？



変換不变な音響量の数学的探求

変数変換と積分

- 一変数: $x = x(t)$ ($x_1 = x(t_1), x_2 = x(t_2)$)

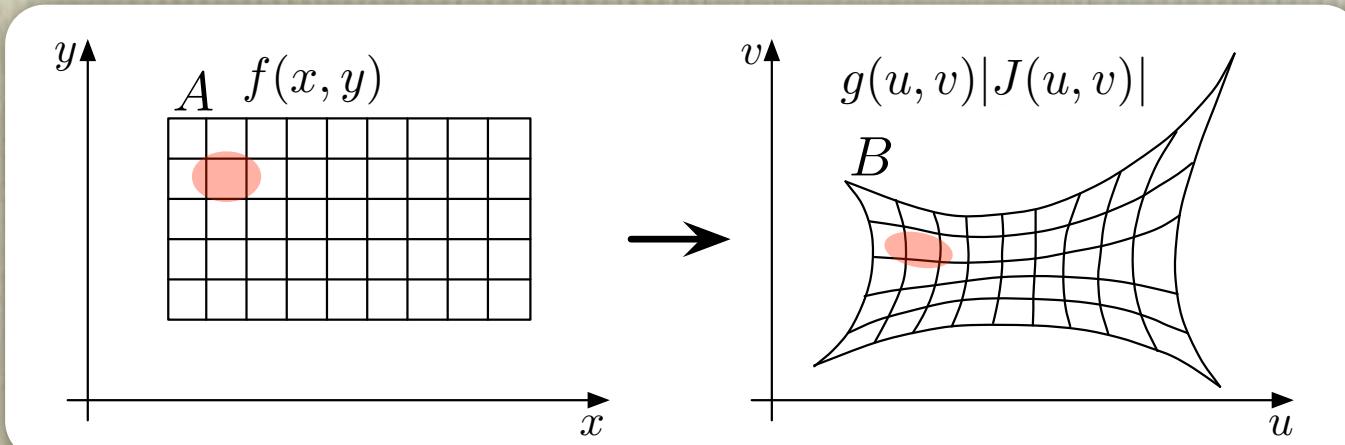
$$\int_{x_1}^{x_2} f(x) dx = \int_{t_1}^{t_2} f(x(t)) \frac{dx(t)}{dt} dt = \int_{t_1}^{t_2} g(t) x'(t) dt$$

- 二変数: $x = x(u, v)$, $y = y(u, v)$

$$\begin{aligned} x &= 3u + 2v - 5 \\ y &= 4u + 5v + 3 \end{aligned}$$

$$\iint_A f(x, y) dxdy = \iint_B f(x(u, v), y(u, v)) |J(u, v)| du dv$$

$$= \iint_B g(u, v) |J(u, v)| du dv \quad J(u, v) \equiv \frac{\partial(x, y)}{\partial(u, v)} \equiv \det \begin{bmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{bmatrix}$$



変換不变な音響量の数学的探求

変数変換と確率密度分布関数

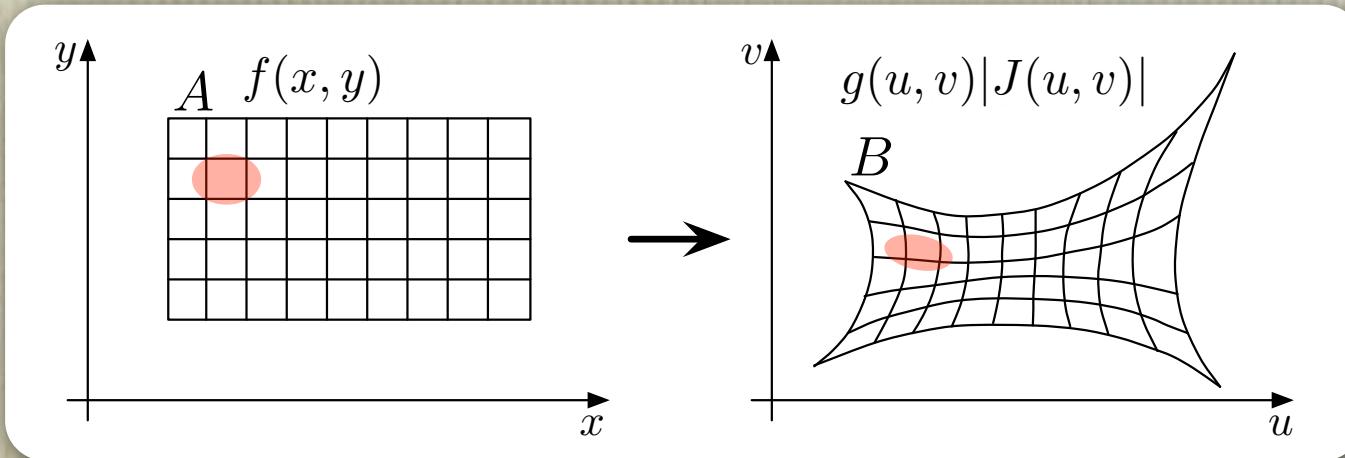
- 一変数: $x = x(t)$ ($x_1 = x(t_1), x_2 = x(t_2)$)

$$1.0 = \int_{x_1}^{x_2} p(x) dx = \int_{t_1}^{t_2} p(x(t)) \frac{dx(t)}{dt} dt = \int_{t_1}^{t_2} q(t) x'(t) dt$$

- 二変数: $x = x(u, v)$, $y = y(u, v)$

$$1.0 = \iint_A f(x, y) dxdy = \iint_B f(x(u, v), y(u, v)) |J(u, v)| du dv$$

$$= \iint_B g(u, v) |J(u, v)| du dv \quad J(u, v) \equiv \frac{\partial(x, y)}{\partial(u, v)} \equiv \det \begin{bmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{bmatrix}$$



変換不变な音響量の数学的探求

バタチャリヤ距離 (=分布間距離尺度の一つ)

● 座標変換による式の変形 $x = x(u, v), y = y(u, v)$

$$BD(p_1(x, y), p_2(x, y))$$



$$= -\log \iint \sqrt{p_1(x, y)p_2(x, y)} dx dy$$

$$= -\log \iint \sqrt{q_1(u, v)q_2(u, v)} |J(u, v)| dx dy$$

$$= -\log \iint \sqrt{q_1(u, v)|J(u, v)| \cdot q_2(u, v)|J(u, v)|} du dv$$

$$= -\log \iint \sqrt{P_1(u, v)P_2(u, v)} du dv$$

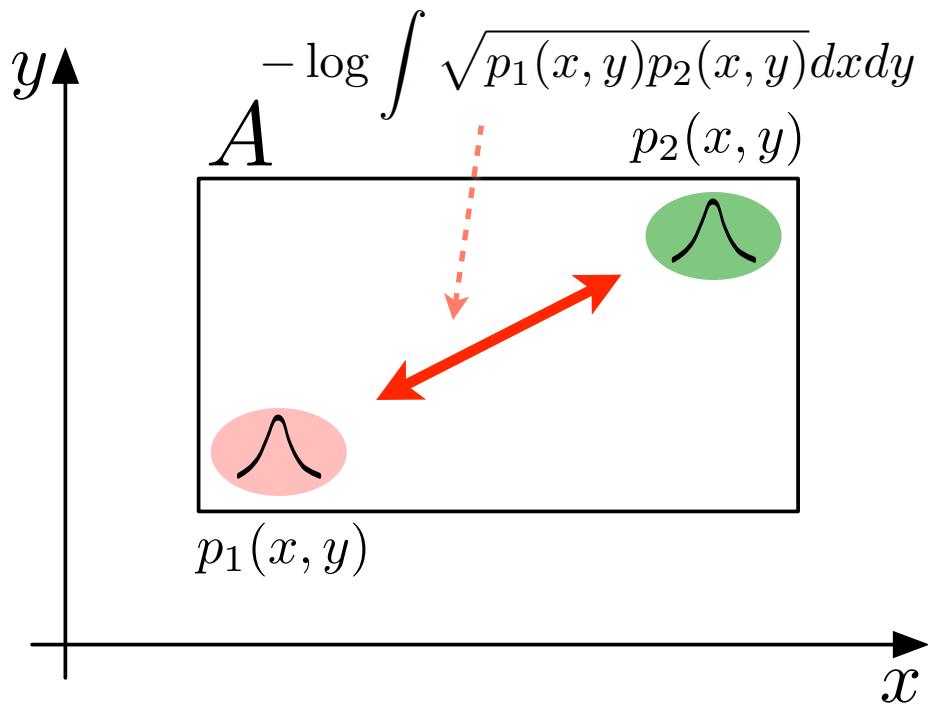
$$= BD(P_1(u, v), P_2(u, v))$$

$$q_1(u, v) = p_1(x(u, v), y(u, v)), \quad J = \text{Jacobian}$$

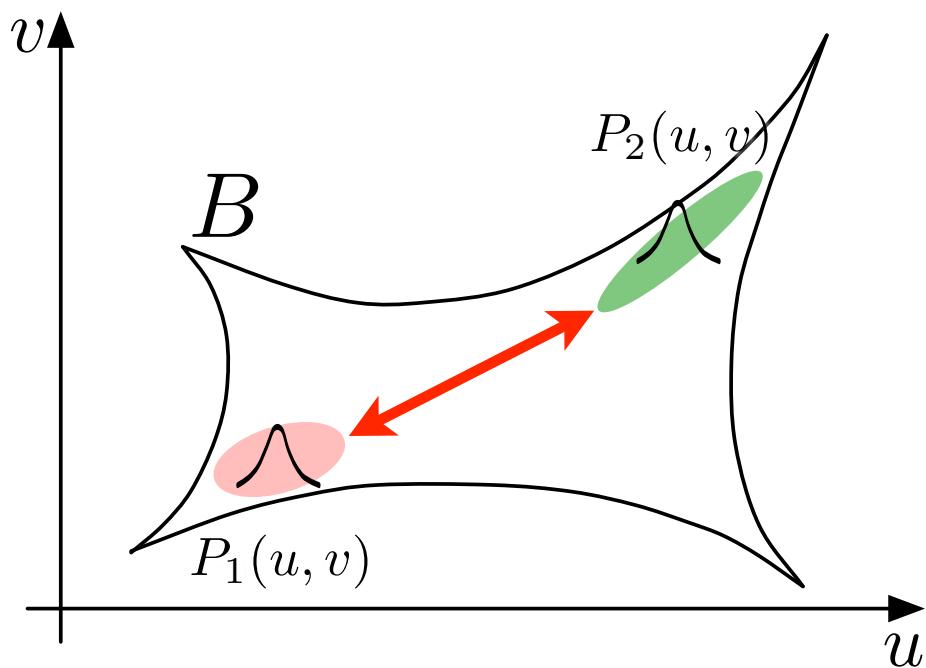
変換不变な音響量の数学的探求

二人の話者空間（一对一対応）における不变音響量

- 音響事象を点ではなく、分布として表現する。
- 空間Aの分布 p は空間Bの分布 P へと写像される。
- p と P は異なる物理特性を持つ（[あ] と [あ]）
- 両空間において不变な物理量はどこに？不变コントラストは存在する？
- 各事象は可変、しかし、少なくともバタチャリヤ距離は不变。



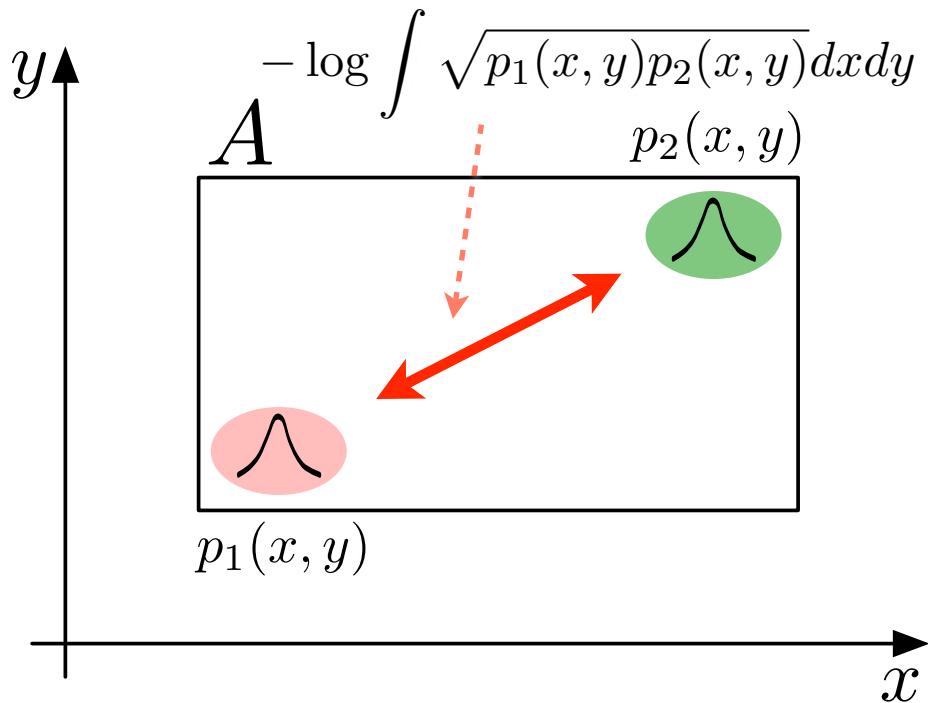
線形
→
非線形



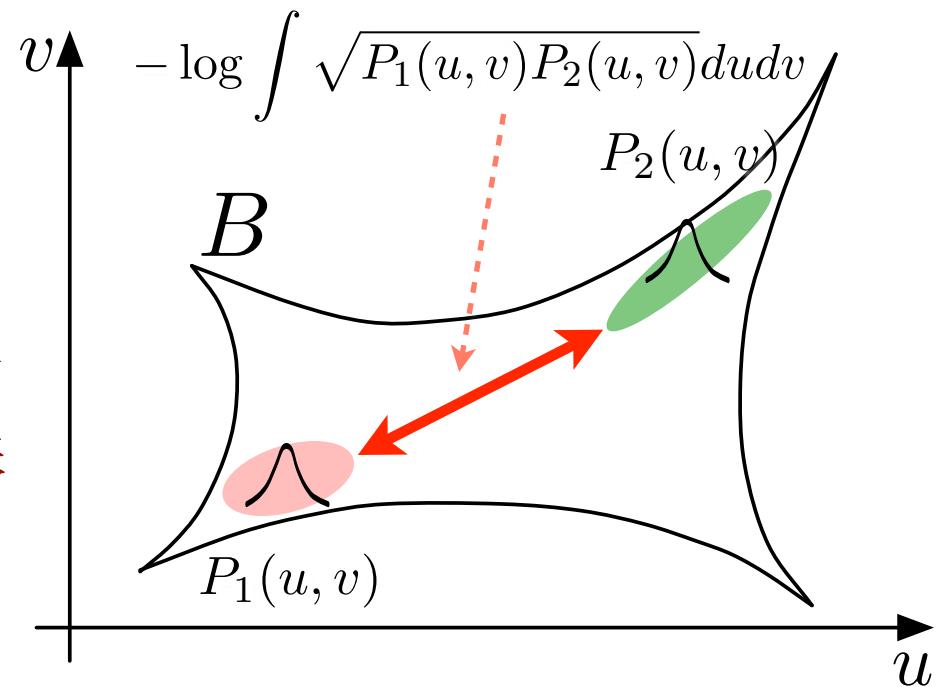
変換不变な音響量の数学的探求

二人の話者空間（一对一対応）における不变音響量

- 音響事象を点ではなく、分布として表現する。
- 空間Aの分布 p は空間Bの分布 P へと写像される。
- p と P は異なる物理特性を持つ（[あ] と [あ]）
- 両空間において不变な物理量はどこに？不变コントラストは存在する？
- 各事象は可変、しかし、少なくともバタチャリヤ距離は不变。



線形
→
非線形



変換不变な音響量の数学的探求

変換不变量の一般式はあるのか？

◆ f-divergence 不変性の十分性

$$f_{div}(p_1, p_2) = \int p_2(x)g\left(\frac{p_1(x)}{p_2(x)}\right) dx$$

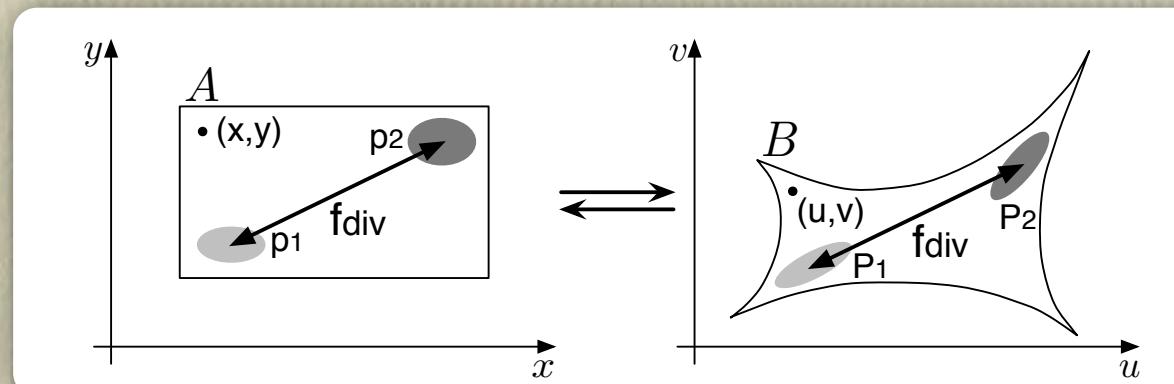
$$g(t) = t \log(t) \rightarrow f_{div} = \text{KL} - \text{div.} \quad g(t) = \sqrt{t} \rightarrow -\log(f_{div}) = \text{BD}$$

$$f_{div}(p_1, p_2) = f_{div}(P_1, P_2)$$

◆ f-divergence 不変性の必要性

$\int M(p_1(x), p_2(x))dx$ が如何なる可逆&連続の変換に対しても不变

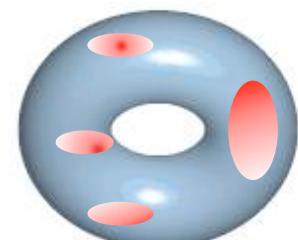
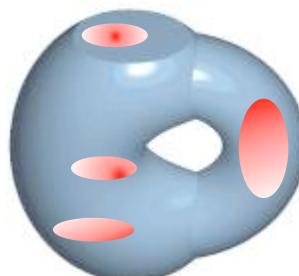
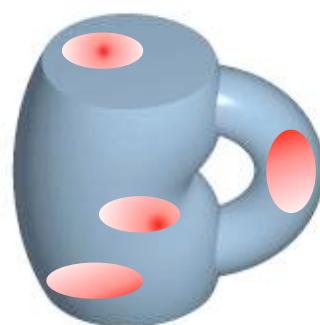
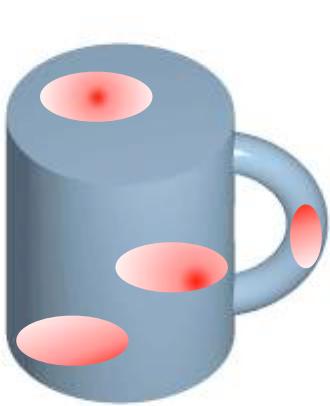
この場合, $M = p_2(x)g\left(\frac{p_1(x)}{p_2(x)}\right)$ であることが必要。



変換不变な音響量の数学的探求

位相幾何学（トポロジー）における不变量

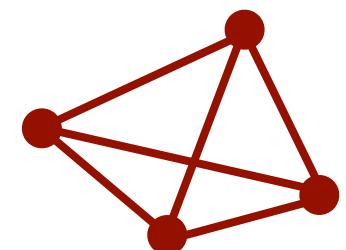
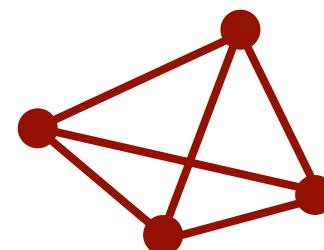
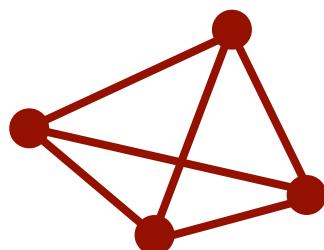
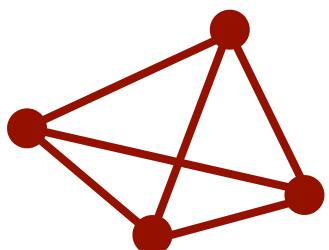
- 連続かつ可逆な任意の変形を施しても不变なる幾何学的性質



変換不变な音響量の数学的探求

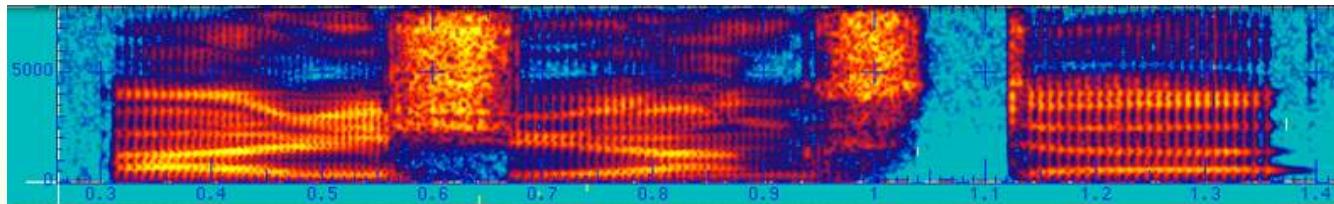
位相幾何学（トポロジー）における不变量

- 連続かつ可逆な任意の変形を施しても不变なる幾何学的性質

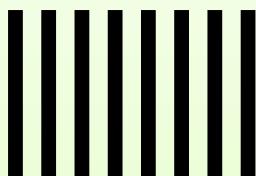
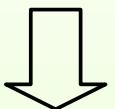


分布間距離群としての音声表象

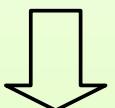
頑健に話者不变な音声表象 = 構造的・全体的表象



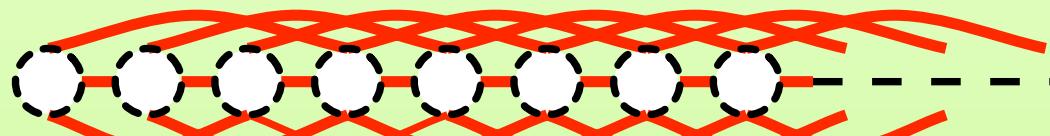
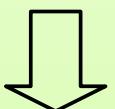
Sequence of spectrum slices



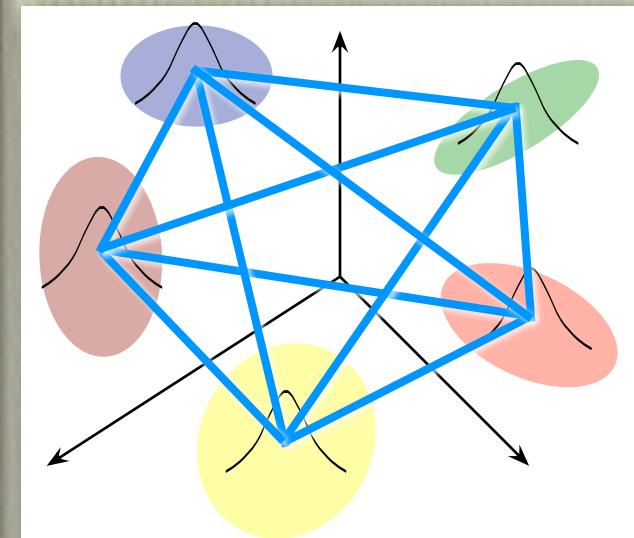
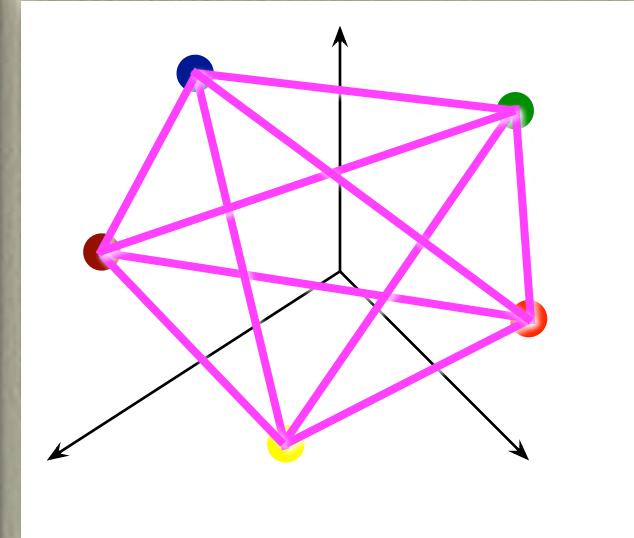
Sequence of cepstrum vectors



Sequence of distributions

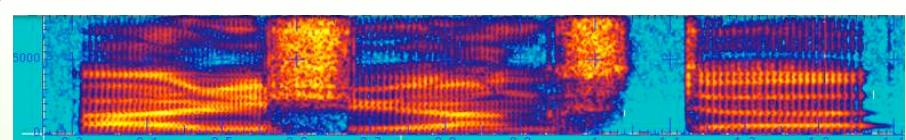
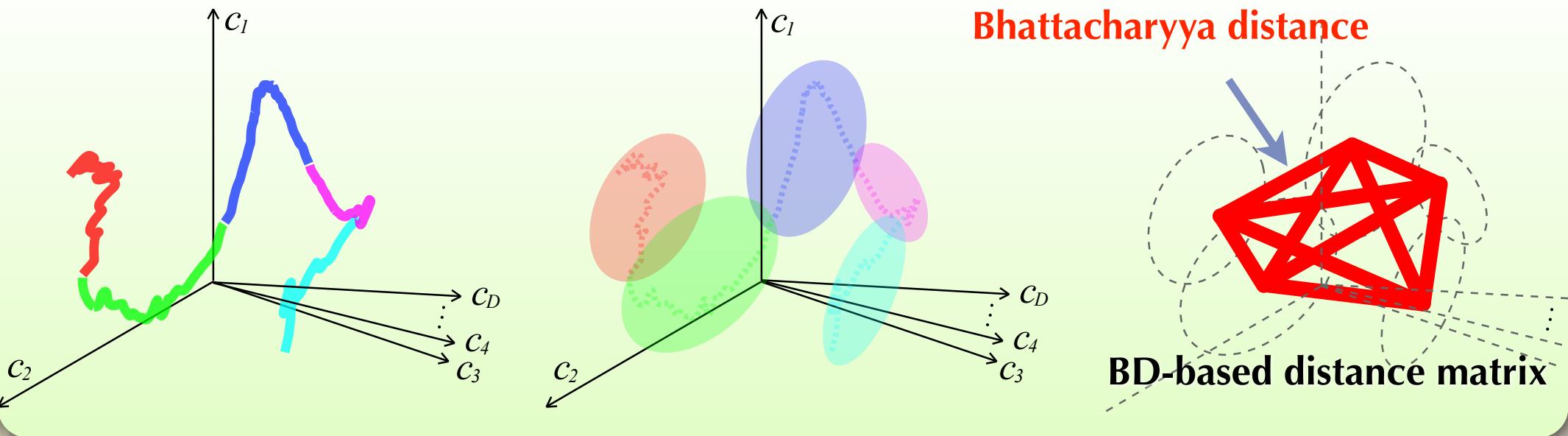


Structuralization by interrelating temporally-distant events



分布間距離群としての音声表象

ケプトラム系列 → 分布系列 → 距離行列



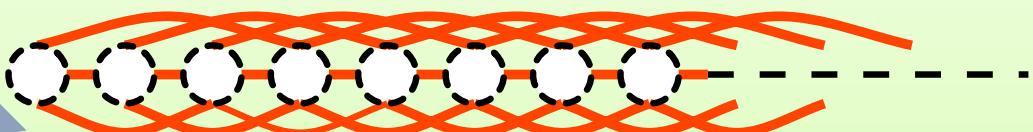
spectrogram (spectrum slice sequence)



cepstrum vector sequence

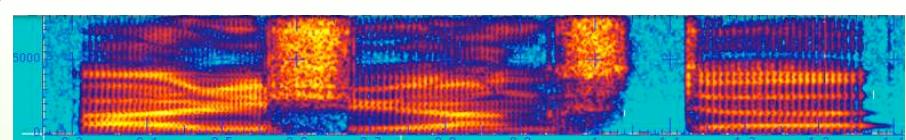
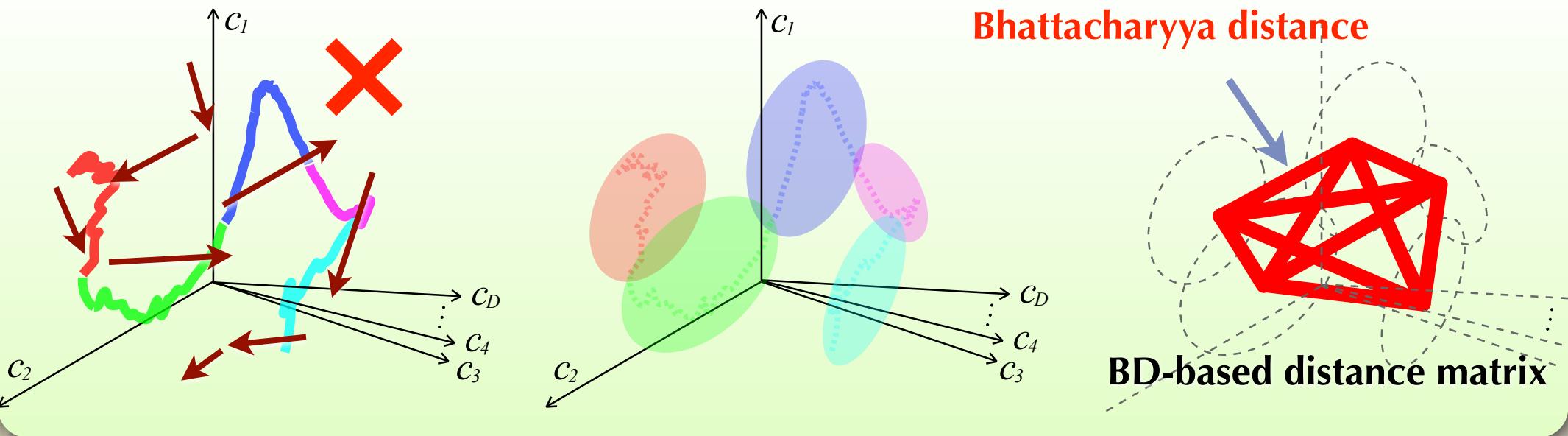


distribution sequence

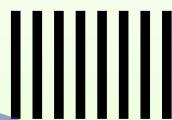


分布間距離群としての音声表象

ケプトラム系列 → 分布系列 → 距離行列



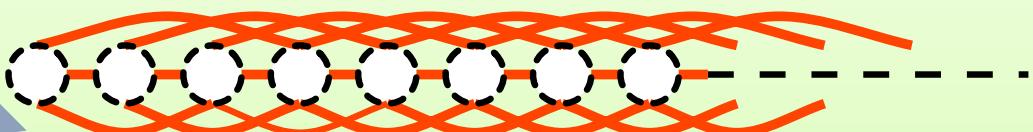
spectrogram (spectrum slice sequence)



cepstrum vector sequence



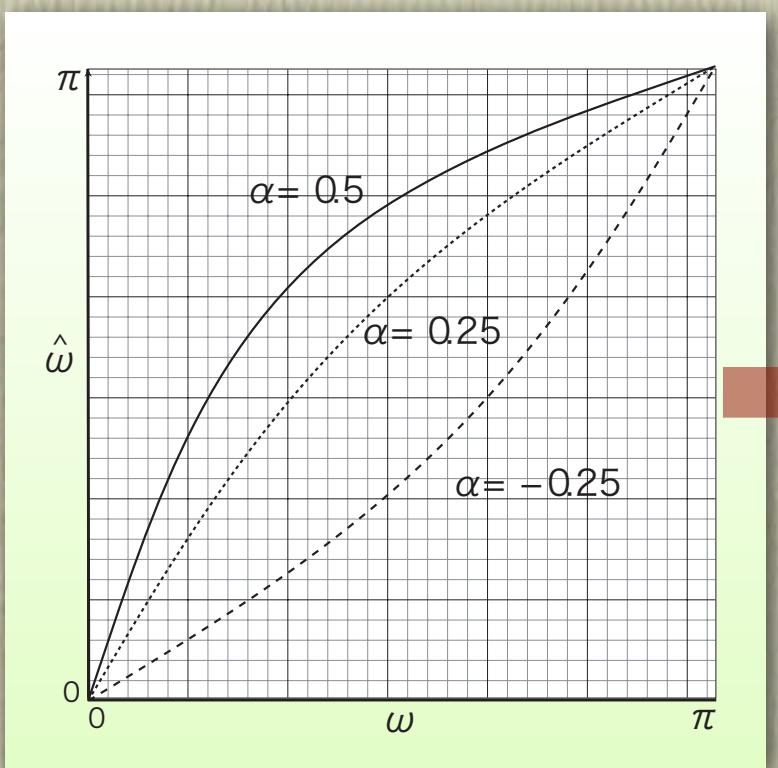
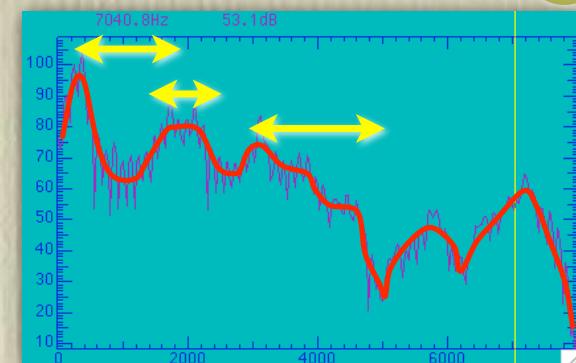
distribution sequence



声道長の変化=行列Aの掛け算

具体的な行列Aの実装は?

- 声道長が伸びる=フォルマントがより低く
- 声道長が縮む = フォルマントがより高く
- スペクトルに対する周波数ウォーピング



$$\hat{c} = (\hat{c}_1 \ \hat{c}_2 \ \hat{c}_3 \ \hat{c}_4 \ \dots)^t$$

$$A = \begin{pmatrix} 1-\alpha^2 & 2\alpha-2\alpha^3 & \dots & \dots \\ -\alpha+\alpha^3 & 1-4\alpha^2+3\alpha^4 & \dots & \dots \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \end{pmatrix}$$

$$c = (c_1 \ c_2 \ c_3 \ c_4 \ \dots)^t.$$

$$a_{ij} = \frac{1}{(j-1)!} \sum_{m=\max(0, j-i)}^j \binom{j}{m} \times \frac{(m+i-1)!}{(m+i-j)!} (-1)^m \alpha^{(2m+i-j)}$$

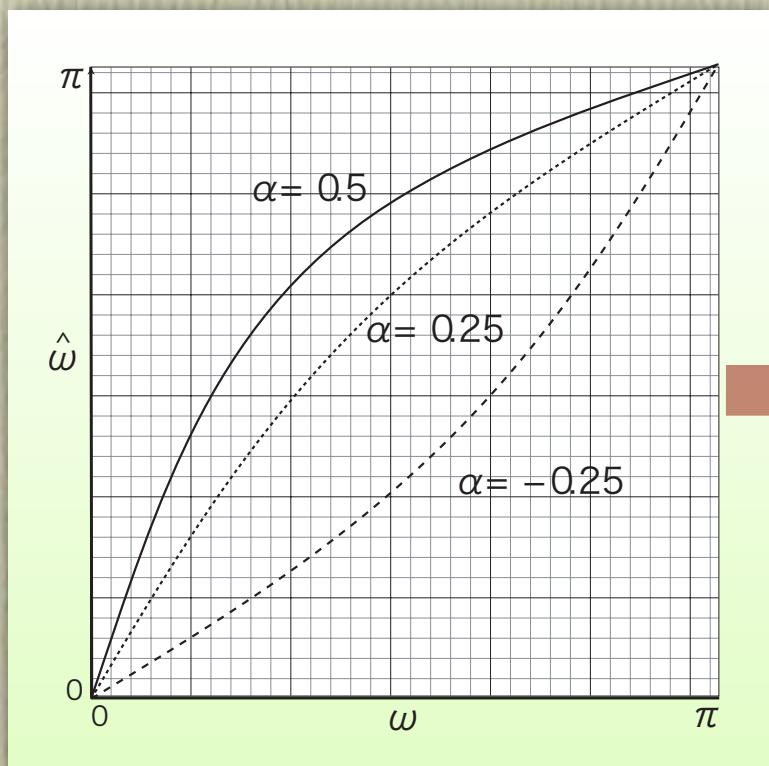
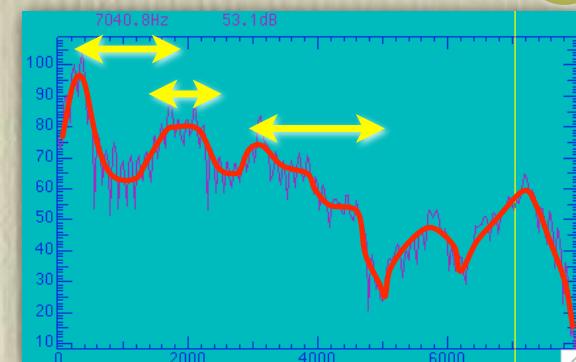
$$\hat{z}^{-1} = \frac{z^{-1} - \alpha}{1 - \alpha z^{-1}}, \quad z = e^{j\omega}, \quad \hat{z} = e^{j\hat{\omega}}$$

$$c' = Ac$$

声道長の変化=行列Aの掛け算

具体的な行列Aの実装は?

- 声道長が伸びる=フォルマントがより低く
- 声道長が縮む = フォルマントがより高く
- スペクトルに対する周波数ウォーピング



➡

$$\hat{\mathbf{c}} = (\hat{c}_1 \ \hat{c}_2 \ \hat{c}_3 \ \hat{c}_4 \ \cdots)^t$$

$$A = \begin{pmatrix} 1-\alpha^2 & 2\alpha-2\alpha^3 & \cdots & \cdots \\ -\alpha+\alpha^3 & 1-4\alpha^2+3\alpha^4 & \cdots & \cdots \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \end{pmatrix}$$

$$\mathbf{c} = (c_1 \ c_2 \ c_3 \ c_4 \ \cdots)^t.$$

$$a_{ij} = \frac{1}{(j-1)!} \sum_{m=\max(0,j-i)}^j \binom{j}{m} \times \frac{(m+i-1)!}{(m+i-j)!} (-1)^m \alpha^{(2m+i-j)}$$

$$\hat{z}^{-1} = \frac{z^{-1} - \alpha}{1 - \alpha z^{-1}}, \quad z = e^{j\omega}, \quad \hat{z} = e^{j\hat{\omega}}$$

$$\mathbf{c}' = A\mathbf{c}$$

行列 A の幾何学的性質

$$\begin{pmatrix} \hat{c}_1 \\ \hat{c}_2 \end{pmatrix} = \begin{pmatrix} 1-\alpha^2 & 2\alpha-2\alpha^3 \\ -\alpha+\alpha^3 & 1-4\alpha^2+3\alpha^4 \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix}$$

$$T = R + O$$

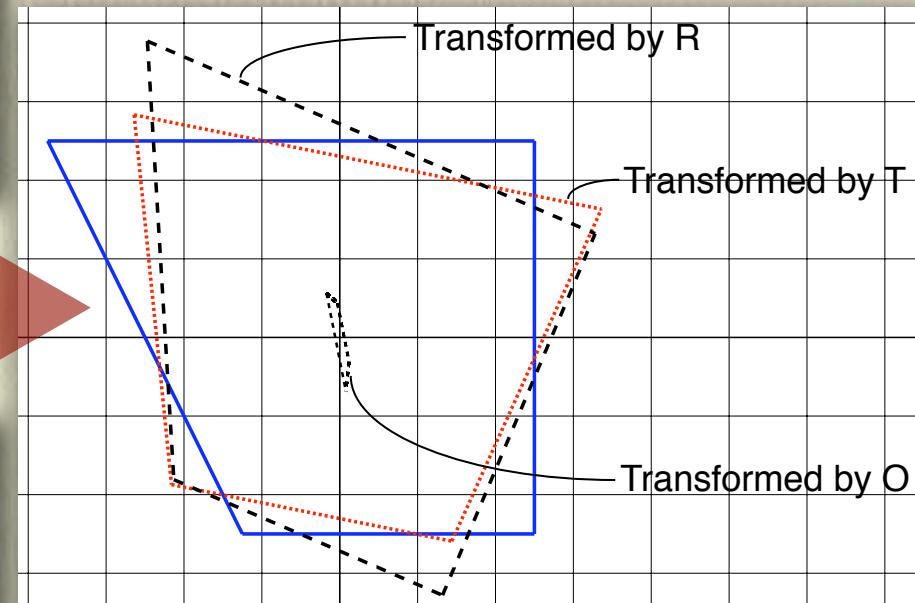
$$R = \begin{pmatrix} 1-2\alpha^2 & 2\alpha(1-\frac{1}{2}\alpha^2) \\ -2\alpha(1-\frac{1}{2}\alpha^2) & 1-2\alpha^2 \end{pmatrix}$$

$$O = \begin{pmatrix} \alpha^2 & -\alpha^3 \\ -\alpha & -2\alpha^2+3\alpha^4 \end{pmatrix}.$$



$$\begin{aligned} R &\simeq \begin{pmatrix} 1-2\alpha^2 & 2\alpha\sqrt{1-\alpha^2} \\ -2\alpha\sqrt{1-\alpha^2} & 1-2\alpha^2 \end{pmatrix} \\ &= \begin{pmatrix} \cos 2\theta & \sin 2\theta \\ -\sin 2\theta & \cos 2\theta \end{pmatrix} (\alpha = \sin \theta) \end{aligned}$$

$$A = \begin{pmatrix} 1-\alpha^2 & 2\alpha-2\alpha^3 & \cdots & \cdots \\ -\alpha+\alpha^3 & 1-4\alpha^2+3\alpha^4 & \cdots & \cdots \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \end{pmatrix}$$



N次元ではどうなる？

行列 A の幾何学的性質

N次元空間における回転行列とは？

$$R^t R = R R^t = I$$

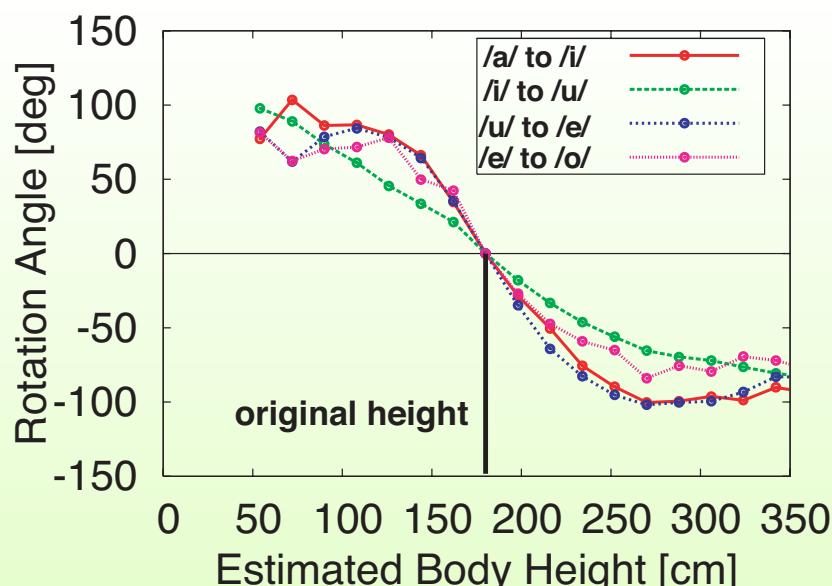
$$\det R = +1.$$



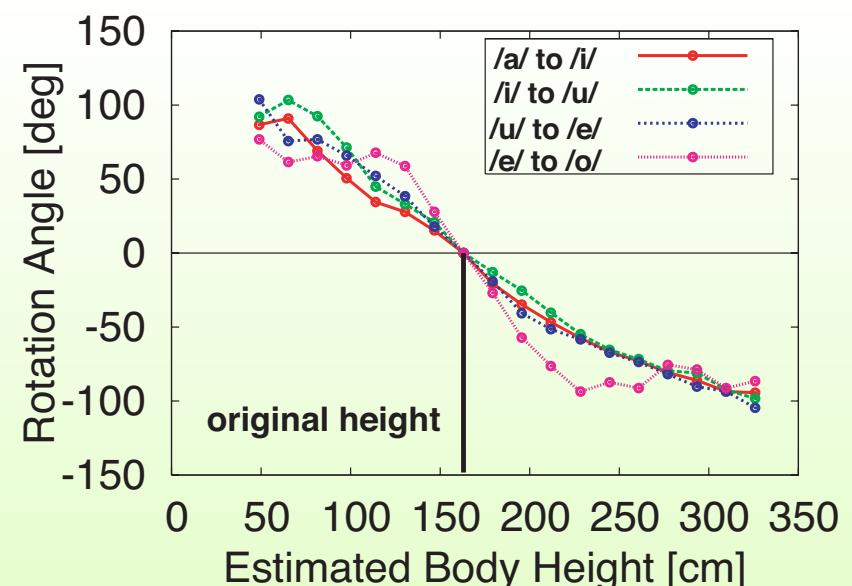
近似的に成立！

$$a_{ij} = \frac{1}{(j-1)!} \sum_{m=\max(0, j-i)}^j \binom{j}{m} \times \frac{(m+i-1)!}{(m+i-j)!} (-1)^m \alpha^{(2m+i-j)}$$

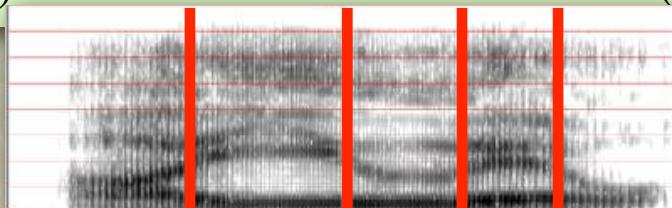
周波数ウォーピングはケプストラムを回転させる！



(a):MFCC (male)

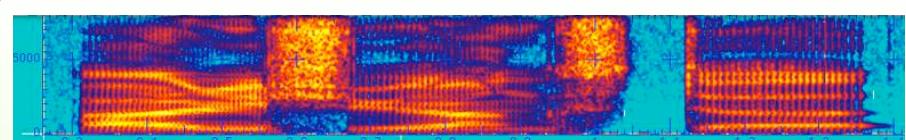
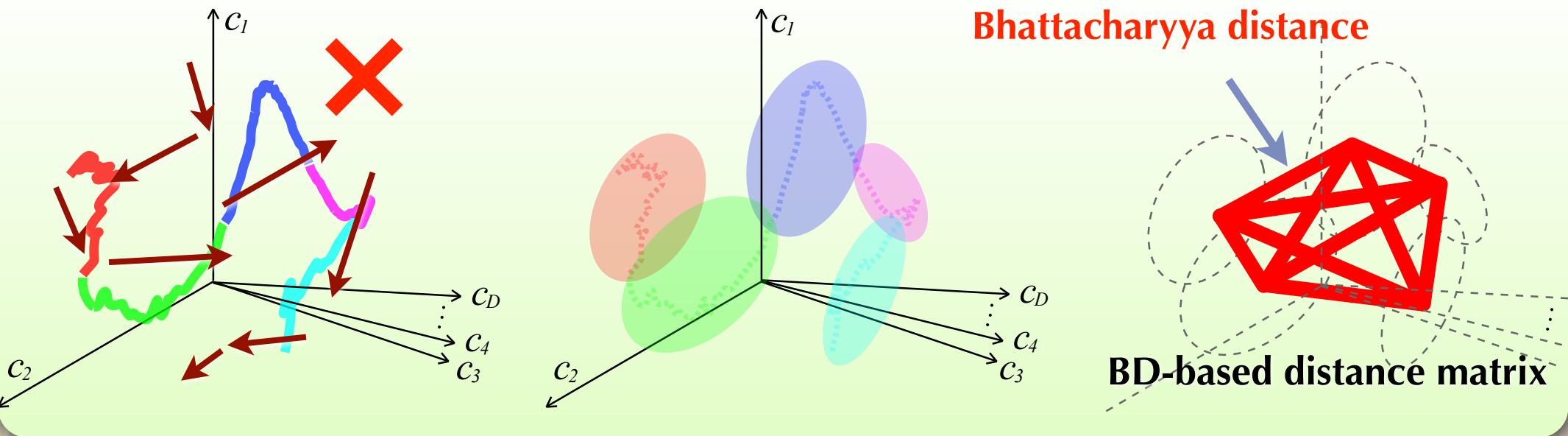


(d):MFCC (female)

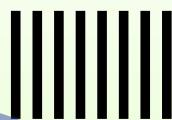


分布間距離群としての音声表象

ケプトラム系列 → 分布系列 → 距離行列



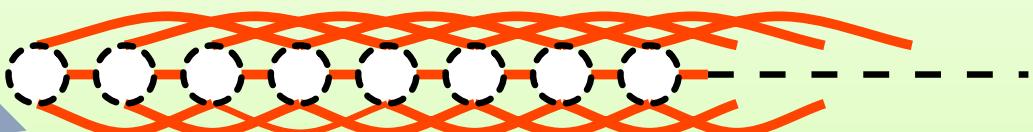
spectrogram (spectrum slice sequence)



cepstrum vector sequence



distribution sequence



音高の相対音感／音色の相対音感

音高＝基本周波数＝1次元の量

- 音の高さ＝直感的に理解しやすい。
 - 単旋律のメロディーを聴いて、音高の動きの様子は把握しやすい。
 - 鼻歌聞かせて「メロディーを描いて」と言えば、指で描ける。
 - 言葉としても「高い↔低い」の対義語で事足りる。
 - 一次元の量だから、その動きの様子を「視覚的に」捉えやすいから？

音色＝周波数軸のエネルギー分布＝多次元の量

- 音の音色＝何か「もやもや」していて掴みどころがない。
 - 「あいうえお」と聞いて、音色の動きの様子を把握できる？
 - その動きを「描いて」と言われても、どう描くべきかすら分らない。
 - 言葉としても「太い↔細い」「しぶい↔若い」など色々。
 - 多次元の量だから、その動きの様子は「視覚的に」捉えられない？
 - 四次元をありありと感覚できる数学家なら捉えられる？
- 隣接音だけでなく、離れた音とのコントラストも必要

面白い事実

Dyslexia であることの利点

空間把握能力

- 空間における物体の形、大きさ、動き、位置、位置関係、及びそれらの相互関係を把握する能力

つながりを把握する能力

- 異なる事物や概念、出来事の相互関係を見抜く力。様々な領域のアプローチやテクニックを使い、物事を様々な視点から見る力

物語を作る能力

- 過去の個人的経験の心的場面を繋ぎ合わせて、過去・現在をリアルに思い出したり、未来をリアルに描く力

未来を予測する能力

- エピソードのシミュレーションを使い、過去や未来の状態を正確に予測する力



音色の偏差とその認知的不变性

色み・音高の恒常・不变的認知

- コントラスト情報に基づく処理が重要
- コントラスト群から成る全体的パターン処理が要素同定を可能に



音色の恒常・不变的認知

- コントラスト情報に基づく処理が重要
- コントラスト群から成る全体的パターン処理が要素同定を可能に

