

基本周波数の周波数変調に対する発声の不随意応答： 純音・複合音を用いた検討

LIAO JIAHUI^{1,a)} 河原 英紀^{2,b)} 松井 淑恵^{1,c)}

概要：聴覚フィードバックは発声基本周波数の調整において重要な役割を果たす。ピッチシフトのある音を聴取した際の発声の応答について、様々な研究が展開してきた。しかし、ランダムな基本周波数の周波数変動に対する発声の不随意応答の分析は難しく、応答の仕組みはまだ明かされていない。河原らはパルスに復元することのできるランダムな系列を用いて基本周波数を変調する方法[1]を考案した。本研究ではこの手法を用いて、基本周波数が周波数変調された純音や複合音を聴取した際の発声実験を実施し、聴取する音の構成が発声の不随意応答にどのように影響するかを調査した。その結果、複合音条件では、基本周波数の周波数変動に対する補償応答が見られたが、純音条件では見られなかった。また、高次の高調波のみ含むミッシングファンダメンタル音はその他の条件より、応答が小さかった。この結果から、基本周波数の周波数変動に対する不随意応答には低次の高調波が影響していることが示唆された。また、発声応答にピッチ知覚が与える影響を確かめるため、発声実験と同じ種類の刺激音に対して、ピッチマッチング実験を行なった。試行ごとのばらつきはあったが、発声応答との関係は明らかにならなかつた。

キーワード：聴覚フィードバック、ピッチ、ピッチシフト、発声、基本周波数変調、応答

Involuntary vocalizations corresponding to modulated fundamental frequency: pilot experiment using pure tone and complex tone

LIAO JIAHUI^{1,a)} KAWAHARA HIDEKI^{2,b)} MATSUI TOSHIE^{1,c)}

Abstract: Auditory feedback plays an important role in the regulation of voice pitch. Various studies have been conducted on vocal responses to pitch-shifted signals. However, it is difficult to analyze voice involuntary response to random modulated fundamental frequency, and the response mechanism has not yet been elucidated. Kawahara et al. devised a method[1] of creating sounds with random fundamental frequency modulation from maximum length sequence and converting them back to pulses. In this study, vocalization experiments were conducted while listening to a tone and a complex tone with momentary modulated fundamental frequency to investigate how the components of the sound being listened to affects the involuntary response of the vocalization. The results showed that the compensatory response to fundamental frequency modulation was observed in the complex tone condition, but not in the pure tone condition. Also, the missing fundamental sound, which contained only higher harmonics, had a smaller response than the other conditions. These results suggest that the involuntary response to random fundamental frequency modulation is influenced by lower-order harmonics. Pitch-matching experiments were also conducted for the same types of stimulus sounds to investigate the effects of pitch perception to the vocalization results. Although there was some variation from trial to trial, the relationship between the vocal response and the perceived pitch was unclear.

Keywords: Auditory feedback, pitch, pitch shift, voice, modulated fundamental frequency, response

1. はじめに

ヒトにはフィードバックされた聴覚情報によって、発声の基本周波数を調整する機構があることが古くから指摘されている。聴覚フィードバック制御による発声ピッチの調整の仕組みを調べるために、ピッチシフトされた音声を聴取した際の発声の変化を観察する実験が多く行われてきた。Elman らは、リアルタイムでピッチがシフトされた自己音声を聞きながら発声する実験 [2] で、ピッチシフトを打ち消すような方向への発声ピッチの変化、すなわち補償応答が見られたことを報告している。このような応答が、随意にコントロールされているものかどうかを調べるために、Hain らは、刺激音のピッチシフトに対する、発声の応答方向（ピッチの上下方向）を教示した実験 [3] を行なった。この実験の結果では、刺激音に 500 ms 以上持続するピッチシフトのある場合、潜時の異なる 2 回の発声応答が見られた。遅い応答は教示した発声のピッチ上下方向と同じ方向へのピッチシフト応答であった。早い応答は応答のピッチ上下方向を教示したにもかかわらず、どれもピッチシフトとは反対方向への発声ピッチのシフトであった。よって、遅い応答は自発的にコントロールできるが、早い応答はコントロールできない、不随意の応答であることが示唆された。また、Zarate らは、ピッチシフト量が異なる二つの条件で、発声応答方向を教示した実験 [4] を行なった。その結果、25 cents のピッチシフト量が小さい方が、教示内容による影響が少なく、不随意であることを示唆した。

これらの実験で使用した刺激音は、ピッチシフトの持続時間が長い。不随意応答を調べるために、持続時間の短いピッチシフトを持つ刺激音を用いる必要がある。Sapir らは、ピッチが正弦波のような（滑らかに）上下変動する刺激音を用いて発声実験 [5] を実施した。その結果、刺激音の基本周波数の周波数変動に合わせた発声ピッチの変動が見られた。しかし、この実験は参加者が少ないとのことと、基本周波数の周波数変動は単一で予測されやすいという課題が残っている。

河原らは、システムのインパルス応答を測定するため、1980 年代から広く用いられていた M 系列（maximum length sequence (MLS)）からできた擬似ランダム信号 [6] を用いて不随意応答の測定実験 [7][8] を行なった。しかし、この方法はハードウェアとソフトウェアの複雑な組み合せが必要であり、また、理論的に不適切な実装が用いられていたという問題がある [9][10]。最近では、河原らは新たに予測不可能な系列を簡単に作成する手法を考案し、CAPRICEP(Cascaded All-Pass filters with RandomIzed

CEnter frequencies and Phase Polarities)[11] で生成する波形を用いて、聴覚-発声基本周波数の制御システムのインパルス応答を分析する手法 [1][12] を開発した。これらの手法を用いて、ランダムな基本周波数の周波数変動のある変調波を生成し、発声音声とともに基本周波数を分析し、基本周波数の周波数変調に対する不随意応答を調査できるようになった。本研究では、この手法を用いて基本周波数の周波数変動のある純音と様々な複合音を作成し、これらを聴取した際の不随意応答を調査した。

2. CAPRICEP による実験刺激と分析方法

CAPRICEP で生成した一個の波形（以下、単位 CAPRICEP）と、その時間反転した波形を畳み込むと、単位インパルスになる。この特性を利用して、ある線形時不变システムに単位 CAPRICEP を畳み込むと、その出力信号からシステムのインパルス応答が求められる。また、この単位 CAPRICEP に ± 1 の重みをつけた直交する信号を、繰り返し配置することで、非線形時不变応答やランダム応答も同時に測定できる。

具体的な手続きは以下の通りである。この直交系列に平滑化処理を加え、周波数変調のための変調信号とする。この変調信号を使って、任意の刺激音を搬送波と見做して周波数変調する。作成した周波数変調音（刺激音）と、この刺激音を聞いて発声した音声を、時間反転した同じ直交系列を用いて、パルス列に復元する。これにより、ランダムの基本周波数の周波数変動をパルス状に変換したものに対する応答、すなわちインパルス応答としての発声の反応を観察できる [13]。発声した音声のはじめと終わりに雑音が入るため、パルス列に復元した信号を全て信頼できる結果ではない。ここで、パルスが最大となった最初の 2 周期分を信頼できる区間として分析に利用した。

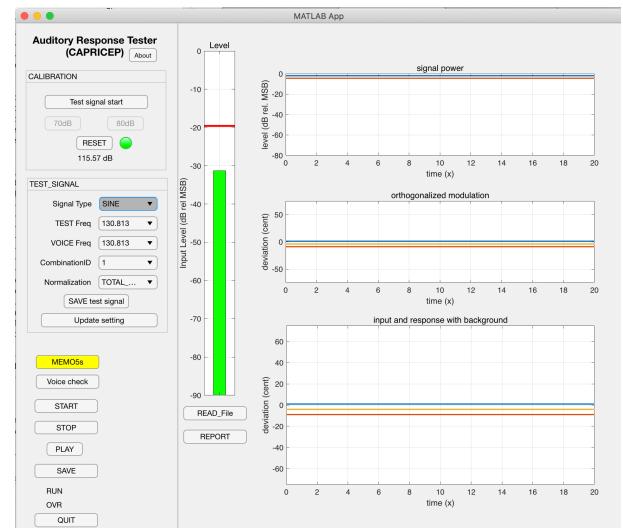


図 1 発声実験の GUI 例。実験者が条件を決め、参加者に合図をしたら「START」ボタン押して実験を開始する。

1 豊橋技術科学大学 Toyohashi University of Technology

2 和歌山大学 Wakayama University

a) liao.jiahui.vy@tut.jp

b) kawahara@wakayama-u.ac.jp

c) tmatsui@cs.tut.ac.jp

3. 発声実験

基本周波数の周波数変動のある刺激音を聴取しながら、発声する実験を行なった。CAPRICEPを利用した分析で、基本周波数の周波数変動に対する発声の応答をインパルス応答として観察した。

3.1 参加者

実験参加者は21歳から24歳の男女9名（男性8名、女性1名）であり、片耳のどちらかで125Hzから8000Hzの聽力レベルが20dBを超えないことを純音聽力検査で確認した。実験参加者には実験について十分な説明を実施し、実験の実施に関する同意を得た。本実験は豊橋技術科学大学の人を対象とする研究倫理審査委員会の承認を受けて行った。

3.2 実験環境

純音聽力検査は聴力検査室（リオン、AT-64）で実施し、音声収録実験は防音室（YAMAHA、AVITECS）で実施した。音声収録実験では、パソコンコンピュータ（Apple、MacBook Air (Retina, 13-inch, 2020)）にオーディオインターフェース（Roland, Rubix24）を接続し、ノイズキャンセリングヘッドホン（SONY, WH-1000XM4）を介して、刺激音を実験参加者に呈示し、実験参加者の発声をマイクスタンドに固定したマイクロホン（Shure, Microflex MX153）から音声を収録した。

3.3 音圧較正

マイクロホンの音圧較正は-3dB/octの特性を有する広帯域雑音を用いて行なった。スピーカー（IK Multimedia, iLoud Micro Monitor）により雑音を流し、マイク付近の音圧レベルが80dBとなるように出力レベルを調節し、この際の較正情報をGUIで記録する。ヘッドホン音圧の較正は、ヘッドホンを人工耳（Brüel & Kjaer, Type 4153）にかぶせてSINES音声条件を流し、サウンドレベルメータ（Brüel & Kjaer, G-4 2250Light）を用いて行い、音圧レベルが85dBであることを確認した。

3.4 実験手続き

本実験は4種類の搬送波を刺激音の条件として用いた：(1) SINE: 基本波となる正弦波のみの純音。(2) SINES: SINEの1-20倍波を足し合わせた複合音。(3) MFND: SINESの基本波を除いた複合音。(4) MFNDH: MFNDの8倍波以下の低い成分を除いた複合音。これらの搬送波に周波数変調で25centsの基本周波数の周波数変動を加えた4種類の刺激音を作成した。刺激音のF0は固定とし、男性はC3の130.813Hz、女性はC4の216.626Hzとした。実験参加者には防音室内でノイズキャンセリングヘッドホン

を装着して刺激音を聴取しながら、マイクロホンに向かって母音の/a/を20s間発声するように指示した。なお、実験者は収録プログラムを操作するために防音室内で収録に同席し、実験状況を観察した。実験で使用したGUIを図1に示す。刺激音の条件はランダム順に呈示し、条件ごと10回ずつ収録した。

3.5 実験結果

実験参加者9名のうち、全体的な発声基本周波数が大幅にずれたため応答の分析が困難となった参加者1名を除いた、男女8名のデータを分析した。

実験参加者全体の発声をパルス復元分析した応答の平均を図2に示す。刺激音のパルスが最大となる2周期を示す。青色は刺激音で、オレンジ色は発声した音声である。横軸はパルス最大となった時刻を0に合わせた2周期分の時間(s)で、縦軸は刺激音のパルスと発声応答の大きさを示す基本周波数(cent)である。観察のしやすさのため、刺激音と発声音の平均基本周波数を0に合わせた。なお、刺激音の大きさは分析によって25cents大きく示しているが、応答の大きさは正しい数値である。複合音条件（SINES, MFND, MFNDH）では、刺激音の基本周波数が変動した後に、発声基本周波数が反対方向にシフトする補償応答が見られた。純音条件（SINE）では補償応答は見られなかった。さらに複合音のうち、MFND条件はSIENS条件より、SINES条件はMFNDH条件より、応答が大きくなる傾向が見られた。

各実験参加者の刺激音の条件ごとの全試行の平均を図3に示す。図の読み方は図2と同様である。青色は刺激音で、その他の色は各参加者が発声した音声を示す。実験参加者の間で特徴が異なることがわかる。MFND条件の方が応答が大きい参加者もいれば、SINES条件の方が応答が大きい参加者もいる。また、応答の大きさにはばらつきもある。

分析結果から応答の大きさを抽出した。パルスが最大となった後の2周期分の信頼区間から、パルスが発生した後の300msの間の最小値を応答とし、パルス発生後の500msから次のパルスの発生前500msを発声の平均とし、応答と発声の差を応答の大きさとした。応答の時刻とその前のパルスの時刻の差を応答の潜時とした。応答の大きさと潜時の抽出例は図4に示す。赤色は応答区間であり、その最小値（青色の点）を応答とみなした。黒色は応答外区間であり、その平均値を発声平均ピッチとした。応答区間の最小値と発声平均ピッチの差を応答の大きさとした。分析対象とした実験参加者の発声応答の大きさに対して、刺激音の条件を要因とする1要因分散分析($\alpha = 0.05$)を行ったところ、刺激音の条件の主効果が認められた [$F(3, 21) = 4.7717, p = 0.0109, \eta^2 = 0.2043$]。しかし、刺激音の条件に対して多重比較を行なった結果、各条件

間の有意差は認められなかった。刺激音の条件における発声応答の平均値はそれぞれ、SINE 条件: 3.819 cents, SINES 条件: 5.313 cents, MFND 条件: 7.950 cents, MFNDH 条件: 4.654 cents であった。

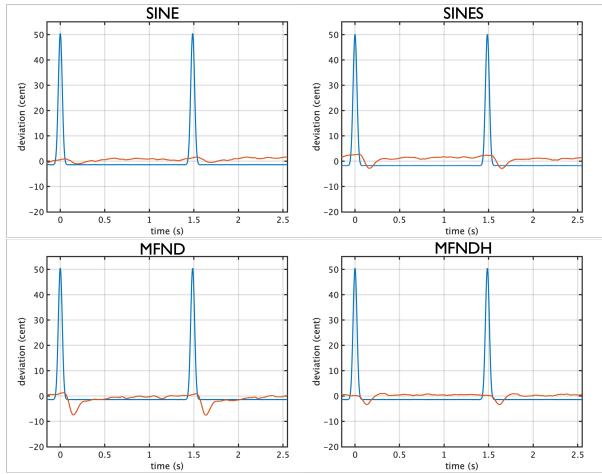


図 2 実験参加者全体の発声を分析した平均結果。刺激音のパルスが最大となる 2 周期を示す。青色は刺激音で、オレンジ色は発声した音声である。横軸はパルス最大となった時刻を 0 に合わせた 2 周期分の時間 (s) で、縦軸は刺激音のパルスと発声応答の大きさを示す基本周波数 (cent) である。複合音条件では、刺激音の基本周波数の周波数変動に対する、発声の補償応答が見られた。

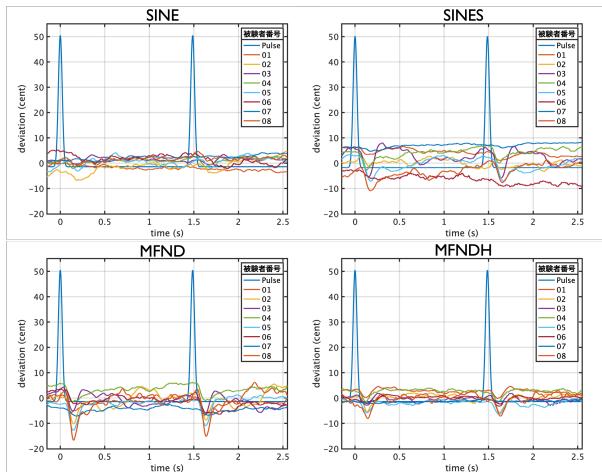


図 3 各実験参加者全体の刺激音の条件ごとの平均結果。青色は刺激音で、その他の色は各参加者が発声した音声を示す。軸の読み方は図 2 と同様である。各実験参加者の応答はばらつきがある。

4. 考察

SINE 条件以外で、刺激音の基本周波数が周波数変調した後に、発声基本周波数が反対方向にシフトする補償応答が見られた。基本周波数の周波数変動に対する応答は、純音より複合音の方が大きい。Sivasankar らの先行研究に

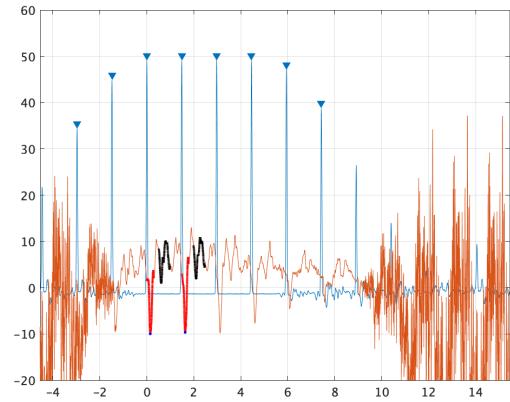


図 4 発声応答の大きさと潜時の抽出例。横軸はパルス最大となった時刻を 0 に合わせた時間 (s) で、縦軸は刺激音のパルスと発声応答の大きさを示す基本周波数 (cent) である。青色が刺激音で、オレンジ色が発声した音声を示す。赤色は応答区間であり、その最小値（青色の点）を応答とみなした。黒色は応答外区間であり、その平均値を発声平均ピッチとした。応答区間の最小値と発声平均ピッチの差を応答の大きさとした。

よると、人の音声と非音声刺激（三角波あるいは純音）のピッチシフトに対する発声応答を調べた [14] 結果、人の音声に対する応答の方が大きく、応答した確率も高いことがわかっている。発声ピッチを調整する聴覚フィードバック制御器は、非音声より人の音声に対して敏感であることが示唆されたこの先行研究と、本研究の結果は一致している。

応答の小さかった純音（SINE 条件）に関しては、純音から知覚されるピッチが曖昧であるため、刺激音の基本周波数変調を検出できなかったことによって本実験と先行研究の結果となった可能性がある。さらに、MFND 条件は SINES 条件より、SINES 条件は MFNDH 条件より、応答が大きい。MFND 条件と SINES 条件の違いは、基本波を含むかどうかであり、MFND 条件と MFNDH 条件の違いは低次の高調波（分解成分）を含むかどうかである。これらの結果から、低次の高調波は基本周波数の周波数変動に対する応答に影響すると考えられる。

参加者は全て歌唱の訓練を受けたことがない学生であり、「どの高さで発声すればいいのかわからない」という感想があった。ある参加者の発声実験結果と発声 F0 の比較を図 5 に示す。この参加者は SINES 条件に比べ、MFND 条件での応答が大きく、試行ごとに発声 F0 のばらつき小さいことが観察できる。

分析結果から応答の大きさを抽出したが、パルスに対する応答なのか、それとも発声の基本周波数が不安定なのかを切り分けることは現段階では難しい。基本周波数の変調に対する応答かどうかを判断する方法を考察し、応答した割合を抽出した上で、応答の大きさと潜時を再度分析する必要があると考えている。

一方で、先行研究での議論や、参加者の内観報告からも、

今回使用した基本周波数の周波数変動のある刺激音から、知覚されるピッチの不明確さや不安定感が、発声応答に関与した可能性はある。そこで、刺激音に対するピッチ知覚が応答とどのような関連性を持つのかをピッチマッチング実験で予備的に調査した。

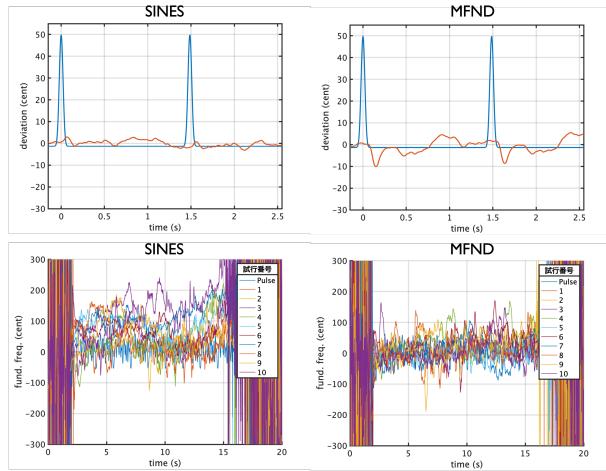


図 5 ある参加者の発声実験分析結果（上パネル）と発声 F0（下パネル）。上の 2 つのパネルは SINES と MFND 条件での発声応答の平均を示し、MFND 条件の方が発声応答が大きい。対して、下の 2 つのパネルは全試行の発声 F0 を示し、MFND 条件の方がばらつきが小さい。

5. ピッチマッチング実験

3 節の実験で用いた刺激音の条件のピッチ知覚を確認するため、ピッチマッチング実験を行なった。3.4 節で述べた基本周波数の周波数変動のある 4 種類の刺激音を用いた。

5.1 実験手続き

刺激音はパーソナルコンピュータにヘッドホンアンプ (FOSTEX, HP-A8) を接続し、実験参加者が装着するヘッドホン (SENNHEISER, HD-650) から刺激音を呈示した。防音室やその他実験機材は 3 と同様である。呈示音圧のレベル較正には、人工耳とサウンドレベルメータを用いた。較正は基本周波数の周波数変動のない SINES 音を使用し、快適に聴取できる 65 dB になるように調整した。実験参加者は 3.1 節で述べた参加者の内男女 5 名である。実験は調整法を用いて、ターゲット音のピッチに合うように、調整音のピッチを調整するタスクとした。刺激音の詳細については、参加者に伏せた。実験参加者は防音室内でヘッドホンを装着して、刺激音を聴取しながら、タスクを行なった。実験者は防音室外で実験の進行をモニターした。

刺激音として、ピッチ調整目標であるターゲット音と、ピッチの調整ができる調整音を呈示した。ターゲット音は、発声実験で聴取した音と同様に、25 cents の基本周波数の周波数変動のある 4 種類の音 (SINE, SINES, MFND,

MFNDH) とした。ターゲット音の F0 は、男性では C3 の 130.813 Hz、女性では C4 の 216.626 Hz を基準に、-75, -50, -25, 0, +25, +50, +75 cents の 7 つとし、これらをターゲットピッチ条件とした。したがって、試行数は 4 種類 × 7 つのターゲットピッチ条件 = 28 個のターゲット音につき 1 試行で、28 試行とした。調整音はヒトの音声と同様の複合音で、基本周波数の周波数変動のない SINES 音とした。調整音の基準となる F0 はターゲット音と同様であるが、ピッチは-100 cents から+100 cents を 5 cents 刻みで調整できる。ターゲット音の呈示順と調整音の初期ピッチはランダムとした。調整音の初期ピッチはターゲット音声のピッチと同じにならないようにした。

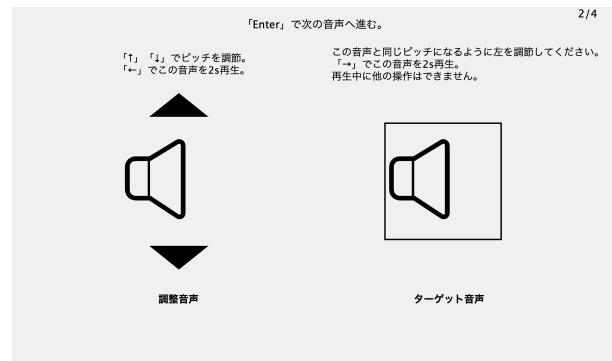


図 6 ピッチマッチング実験の GUI 例。参加者は二つの音声を聞き、右のターゲット音声に合わせるように、キーボードを操作して、左の調整音声のピッチを変更した。

5.2 実験結果

ピッチマッチング実験の結果を図 7 に示す。縦軸は調整音声のピッチとターゲットピッチのずれを cent 単位で示す。横軸はターゲット音声の条件を示す。実線の各色は 4 種類のターゲット音条件ごとの、各参加者のピッチマッチング結果の試行平均を示し、灰色の点線は全参加者のピッチマッチング結果の平均を示す。エラーバーは標準誤差を示す。全参加者の平均はどの種類のターゲット音条件においても、ターゲットピッチからのずれは大きくないが、ばらつきがある。全参加者は SINE 条件において、ターゲットピッチに近い値を選定し、試行ごとのばらつきも小さい。その他条件においては、参加者間に異なる傾向が見られ、試行ごとのばらつきは大きい。

5.3 考察

ピッチマッチング実験での標準偏差が大きいほど、そのターゲット音に対して知覚したピッチは不安定である。この不安定性は、発声応答にも影響すると考えられる。ピッチマッチング実験参加者のみの発声実験の結果を図 8 に示す。ピッチマッチングのばらつきを示す標準偏差に注目すると、参加者 06, 07 の標準偏差が他の参加者より大きい。

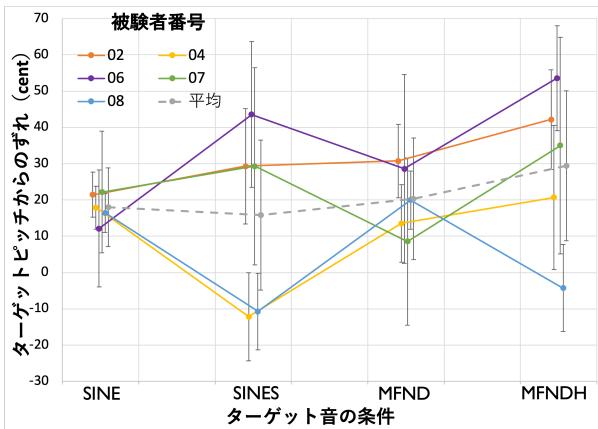


図 7 ピッチマッチング実験の結果。縦軸は調整音声のピッチマッチング結果とターゲットピッチのずれを cent 単位で、横軸はターゲット音声の条件を示す。実線の各色はターゲット音条件ごとの、各参加者のピッチマッチング結果の試行平均を示し、灰色の点線は全参加者のピッチマッチング結果の平均を示す。エラーバーは標準誤差を示す。

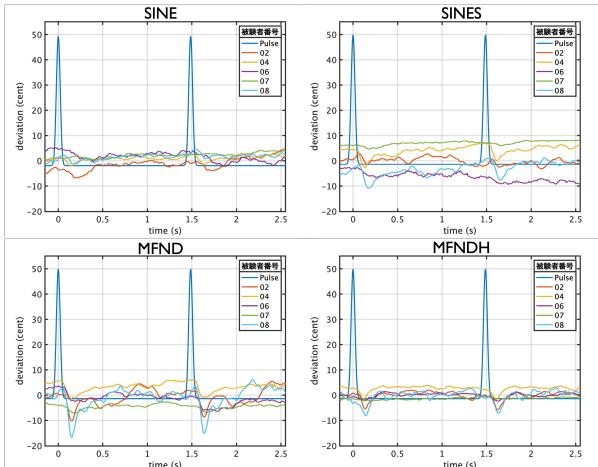


図 8 ピッチマッチング実験参加者のみの発声実験の結果。読み方は図 3 と同様。

そして、発声実験においても、同参加者 2 名ではどの条件においてもはっきりとした応答が見られなかった。しかし、ピッチマッチング実験の試行数が少ないとや、ピッチシフトに対する応答は個人差が存在することなどの要因もあり、ピッチマッチングと発声応答の間に関連性があるかどうかはまだ不明瞭である。

今後の展望として、刺激音の基本周波数の周波数変動を検出できたかどうかは、発声応答に影響すると考えられるため、基本周波数の変動の検知閾を調べることを予定している。また、今回の実験の基本周波数の周波数変動量は、先行研究 [4] で不随意応答が出るであろう 25 cents に設定した。しかし、先行研究では 25 cents と 200 cents しか比較していないため、200 cents より小さい基本周波数の周波数変動に対する応答も不随意である可能性がある。基本周波数の周波数変動の大きさを変えて実験を行い、発声応

答への影響を調べていきたい。

6. まとめ

基本周波数の周波数変動のある純音や複合音を聴取した際の発声実験を実施し、聴取する音の構成が発声の不随意応答にどのように影響するかを調査した。その結果、複合音条件では、基本周波数の周波数変動に対する補償応答が見られたが、純音条件では見られなかった。また、高次の高調波のみ含むミッシングファンダメンタル音はその他の条件より、応答が小さかった。この結果から、瞬間的な基本周波数の周波数変動に対する応答には低次の高調波が影響していることが示唆された。また、同じ種類の刺激音に対して、ピッチマッチング実験を行なった。試行ごとのばらつきはあったが、発声応答との関係性は不明瞭であった。

謝辞 本研究の一部は、JSPS 科研費 18K10708, 21K19794, 21H03468, 21H03759, 22K03193 の支援を受けた。

参考文献

- [1] Kawahara, H., Matsui, T., Yatabe, K., Sakakibara, K.-I., Tsuzaki, M., Morise, M. and Irino, T., *Proc. Interspeech 2021*, pp. 3206–3210 (2021).
- [2] Elman, J. L., *The Journal of the Acoustical Society of America*, Vol. 70, No. 1, pp. 45–50 (1981).
- [3] Hain, T. C., Burnett, T. A., Kiran, S., Larson, C. R., Singh, S. and Kenney, M. K., *Experimental Brain Research*, Vol. 130, No. 2, pp. 133–141 (2000).
- [4] Zarate, J. M., Wood, S. and Zatorre, R. J., *Neuropsychologia*, Vol. 48, No. 2, pp. 607–618 (2010).
- [5] Sapir, S., McClean, M. D. and Luschei, E. S., *The Journal of the Acoustical Society of America*, Vol. 73, No. 3, pp. 1070–1073 (1983).
- [6] Schroeder, M. R., *The Journal of the Acoustical Society of America*, Vol. 66, No. 2, pp. 497–500 (1979).
- [7] Kawahara, H., *Acoustical Society of Japan*, Vol. 15, No. 3, pp. 201–202 (1994).
- [8] Kawahara, H., Kato, H. and Williams, J., *Proceeding of Fourth International Conference on Spoken Language Processing. ICSLP'96*, Vol. 1, IEEE, pp. 287–290 (1996).
- [9] Farina, A., *Audio engineering society convention 108*, Audio Engineering Society (2000).
- [10] Stan, G., Embrechts, J. and Archambeau, D., *Journal of the Audio engineering society*, Vol. 50, No. 4, pp. 249–262 (2002).
- [11] Kawahara, H. and Yatabe, K., *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, pp. 306–310 (2021).
- [12] 河原英紀, 矢田部浩平, 情報処理学会研究報告, Vol. 2020-NL246, No. 32, p. 1–8 (2020).
- [13] Kawahara, H., Sakakibara, K., Mizumachi, M., Morise, M. and Banno, H., *2020 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, IEEE, pp. 174–183 (2020).
- [14] Sivasankar, M., Bauer, J. J., Babu, T. and Larson, C. R., *The Journal of the Acoustical Society of America*, Vol. 117, No. 2, pp. 850–857 (2005).