

逆シャドーイング法を用いた瞬時了解度アノテーションと その高精度化に関する分析的検討

箱田 峻¹ 朱 伝博¹ 齋藤 大輔¹ 峯松 信明¹ 中西 のりこ²

概要：コロナウィルスの影響により諸外国との往来は激減したが、オンラインツールの発達により、依然として英語を用いた国際コミュニケーション能力の習得は重要である。英語発音の指導では、「母語話者らしい発音」の習得を目指すのではなく、「十分相手に伝わる発音」を目指すことが多い。この場合、発話の「伝わりやすさ」すなわち「了解度 (intelligibility)」を、任意の発話に対して定量的にスコア化 (ラベル付け) する必要がある。筆者らはこれまで、学習者音声をも母語話者 (相当) にシャドーさせ、そのシャドー音声の崩れを分析する逆シャドーイング法を検討してきた。逆シャドーイング法では、シャドー音声の他に (学習者が読み上げ時に参照した) 原稿を見ながら行うスクリプト・シャドー音声も別途収録し、両者の posteriorgram DTW のスコアを、自動計測された了解度スコアとして用いている。本研究では、シャドー音声の崩れの様子を考慮し、また DTW 計算において、スクリプト・シャドー音声が常に規範音声として使われることを鑑み、非対称な局所パスを導入し、精度向上を検討した。更に、学習者音声やシャドー音声に対する音声認識結果 (自動書き起こし) を使うことの是非についても検討した。

キーワード：英語教育、瞬時了解度、逆シャドーイング、スクリプト・シャドーイング、DTW

1. はじめに

外国語学習の発音指導において、母語話者らしい発音を学習者に求めるのはコストが高く、また、母語訛りはその学習者のお国柄 (民族的な identity) を表すため、それを否定することは教育的に不適切である。即ち、相手に十分伝わる発話能力の獲得が目的となる [1-5]。ある L2 音声に対して、その伝わりやすさ (了解度, intelligibility) をスコア化する場合、二つの言語的多様性を考慮する必要がある。一つ目は、発話者、つまり学習者の多様性である。様々な母語の学習者が英語を学んでおり、その結果、多様な母語訛りが存在する [6]。もう一つは、聴取者の多様性である。ある言語を母語とする英語利用者は、どのような英語発音が聞き取りやすく / 聞き取り難いのか、を考えると、了解度は聞き手の言語背景に大きく依存する [7, 8]。即ち了解度とは、話者と聴取者、両方に依存することとなり、両者に依存して活用できる簡便な計測方法が求められている。

従来、外国語教育や応用言語学の分野では、L2 音声の了解度を計測する手法として、音声を聴取者 (多くは母語話

者) に書き起こさせ、その書き起こしと L2 音声の原稿 (スクリプト) との単語単位の一貫度を算出している [9-13]。一般的に単語を単位とすることが多いが、音節や音素など細かい単位を採択する場合は、原稿や書き起こしを、音節や音素を単位として表記しなおして比較すれば良い。聴取者の多様性を考慮するならば、上記の作業を様々な言語背景の聴取者を対象として実施することになる。

しかしながら、この手法には大きな問題が存在する。まず書き起こしは一般に、提示音声の数倍の時間を要する。了解度計測は複数回聴取は許されず、また、聴取者の短期記憶の容量にも制限があるため [14, 15]、短文しか利用できない [10-13, 16]。更に、書き起こし中に発話内容を深く推測する余地を与えることとなる。以上の問題は、従来方法が聴取プロセスを「聴取の後」に観測していることがそもその原因であり [11]、聴取プロセスを「聴取中」のみ、観測することで解決される。

筆者らは、発音訓練として学習者に広く導入されているシャドーイングを、母語話者 (相当) に導入し、彼らが L2 音声を聴取している様子をモニタリングする手法を提案した (逆シャドーイング法) [17-20]。従来法では L2 音声を書き起こし「その書き起こしがどのくらい崩れたのか」を計測するが、逆シャドーイング法では、L2 音声をシャドーし「そのシャドー音声のどのくらい崩れたのか」を計測す

¹ 東京大学
The University of Tokyo, Bunkyo, Tokyo, Japan
² 神戸学院大学
Kobe Gakuin University, Kobe, Hyogo, Japan

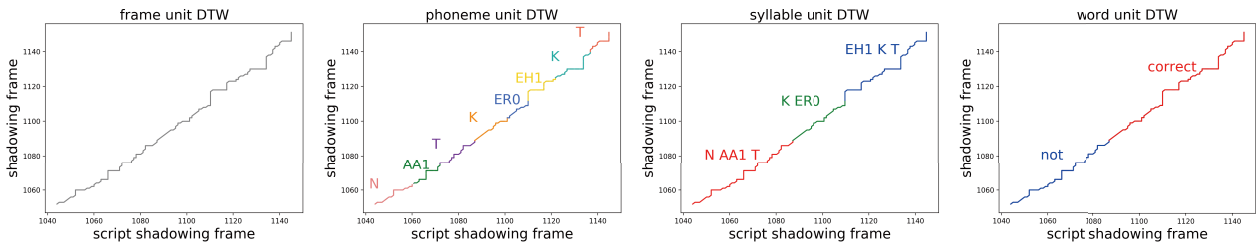
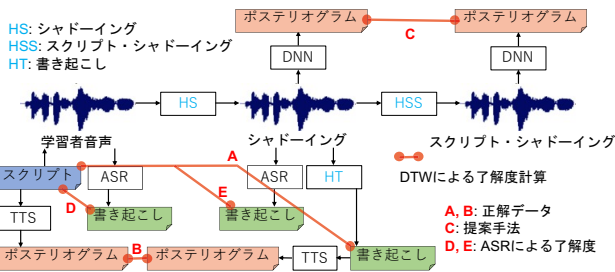


図 2 4つの言語単位による学習者音声“not correct”のDTW



HS, HSS, HT の H は人間 (human) による作業という意図である。
図 1 本実験全体の概略図

る。この場合、推測の余地が入らず、また、30 秒程度であればシャドーし続けることは容易である。即ち、聴取しているその場での了解度、瞬時了解度をスコア化できる。本研究では、逆シャドーイング法において了解度を算出する際のアルゴリズムを工夫し、より高精度な了解度のアノテーションができるか検討した。

なお、本報告の実験の全体像を図 1 に示す。以降の説明は適宜、この図を参照しながら行う。

2. 逆シャドーイング法

2.1 シャドーイングと逆シャドーイング

シャドーイングとは、モデル音声を聴きながら即座に復唱する訓練法のことであり、外国語学習において広く導入されている。逆シャドーイングとは、シャドーする/される関係性を逆にしたもので、L2 音声を母語話者 (相当) の聴取者がシャドーする。この場合、外国語訛りまでを真似ることは求めず、学習者が意図した単語列を自身の発音で復唱することが求められる。逆シャドー音声の崩れの定量化は、シャドーの後に、(学習者が参照した) 原稿・スクリプトを見ながらのシャドーである、スクリプト・シャドーも課し、二種類のシャドー音声の比較を通して定量化する [20]。スクリプト・シャドーは内容既知であり、一番流暢なシャドー音声と解釈でき、それを基準としてシャドー音声の崩れを計測する。両音声を posteriorgram-DTW 法によって比較すれば、平均 DTW スコア (調音崩れの平均に相当) が瞬時明瞭度の指標となる (図 1C)。

2.2 複数粒度による逆シャドーイング分析

Posteriorgram DTW で利用する局所距離尺度としてはバタチャリヤ距離を用い、また、音声の posteriorgram 化には、WSJ-KALDI [21] のレシピを用いたフロントエンド

を用いている。平均 DTW スコアはフレーム、音素、音節、単語の 4 つの単位で求めることができる。フレームは音響的分析の単位、音素は発音の最小単位、音節はリズム知覚、聞こえに関する単位、単語は音声コミュニケーションの最小単位である。DTW 処理はフレームを単位として局所距離を累積するが、構成された DTW パスを音素、音節、単語で区分 (forced alignment) し、各単位でのサブスコアを求め、これらのサブスコアの (全パスに対する) 平均値を計算する。各単位での DTW 平均スコアを、fDTW (frame)、pDTW (phoneme)、sDTW (syllable)、wDTW (word) と以降、呼ぶこととする。図 2 に、学習者発声 “not correct” に対するシャドー音声、スクリプト・シャドー音声から得られた DTW パスと、4 種類のスコア計算の様子を示す。言語単位の明示的な導入は、その単位の継続長の違いを正規化する効果を有する。例えば母音は長く、子音は短い、fDTW は継続長の偏りを正規化していないため、母音のシャドー崩れの影響が出やすい。pDTW は音素継続長の偏りは回避できるが、例えば、長い名詞と短い名詞の差は正規化できていない。wDTW は単語毎にサブスコアを計算するので、単語単位での継続長の正規化が行われている。了解度と言っても、どの粒度で了解度を検討したいのかは教育指針に依存する。教育指針に柔軟に対応できるよう、複数粒度での分析を行なっている。

3. 書き起こしによる瞬時了解度の算出

逆シャドーを課すことで、瞬時的聴解の崩れを定量的に予測することが可能であることを説明した。本節では、得られた自動スコアに対する ground truth である、手動スコアについて説明する。従来研究では、学習者音声を書き起こし、これと学習者が参照した原稿とを比較していたが、本研究では、母語話者 (相当) によるシャドー音声の手動書き起こしと、学習者の原稿とを比較する (図 1A)。即ち、自動スコアである「シャドー音声とスクリプト・シャドー音声の DTW スコア」に対応する手動スコアとして、「シャドー音声の書き起こしと学習者の原稿との差異」を使う。書き起こし、原稿、両方とも単語書き起こしであるが、これを発音辞書を参照して音素書き起こし化し、また、音節分かち書きすることで、音節書き起こし化できる。即ち手動スコアも、単語、音節、音素の粒度のスコアを準備できる。以上は二種類のテキスト比較に基づくスコア化である

が、このテキストを音声合成により音声化し、posteriorgram DTW を通して比較することもできる (図 1B)。

3.1 記号的な書き起こしとその比較

シャドー書き起こしと学習者が参照した原稿は、文字列を対象とした DTW で比較する (図 1A)。シャドー音声中に検出された置換誤り (S)、削除誤り (D)、挿入誤り (I) の数を用いて計算される正解率 %Accuracy (Acc) $= (N - D - S - I) / N$ (N は単語総数) を、書き起こしに基づく瞬時了解度の手動スコアとする。以上を単語、音節、音素を単位として行い、wAcc, sAcc, pAcc を計算する。

pAcc と sAcc では、wAcc よりも、シャドワーの聴解精度をより詳細に評価できる。例えば、“conception” と “consumption” を単語単位で比較する場合、両者は異なるだけで判定されるが、音節単位で評価した場合、中間の /sep/ と /samp/ が異なっているだけで、それ以外は合致する。音素単位の評価すると、/e/ が /a/ と /m/ に変化しただけである。より小さい単位での評価により、より細かい比較が可能だが、後述の確率的な書き起こし (posteriorgram) を用いれば、より細かい比較が可能となる。

3.2 確率的な書き起こしとその比較

上記の/e/が/a/と/m/に変化した場合では、/e/と/a/との違いの方が、/m/との違いよりも小さい。これは、/e/と/a/がともに母音である一方、/m/は子音であるためであるが、pAcc による記号的な評価では、音素間の音声学の差異が全て等価に扱われてしまう。そこで、シャドー書き起こしと、学習者の原稿を音声合成器 (Amazon Polly [22]) により一旦音声化し、posteriorgram 化して比較する。二種類のテキストの差異の定量化を、posteriorgram を通して行うことで、より詳細な差異の定量化が可能になる (図 1B)。なお、posteriorgram は音素事後確率ベクトルの時系列であるので、確率的な音素書き起こしと解釈でき、逆に通常の音素書き起こしは、one-hot ベクトルだけで構成された posteriorgram と解釈できる。posteriorgram DTW を書き起こしに適用したものを S-DTW (Synthesis-based DTW) とする。S-DTW も、複数の粒度で算出できる。

4. 非対称 DTW パスの利用

筆者らの先行研究では、Posteriorgram DTW で用いる局所パスとして図 2-1 に示す、最も一般的な、対称的な局所パスを用いていた。この場合、比較する二発声が等価に扱われる。しかし逆シャドーイング法の場合、比較対象はシャドー音声とスクリプト・シャドー音声であり、発話の崩れは常に前者にあり、後者は常に模範発声として扱われる。また、シャドー音声中に散見される (無音や繰り返し) の挿入や、発声の間延び [23] は、発話崩れとは考えずにスコア計算から除外すべき発話現象である。これらを考慮し

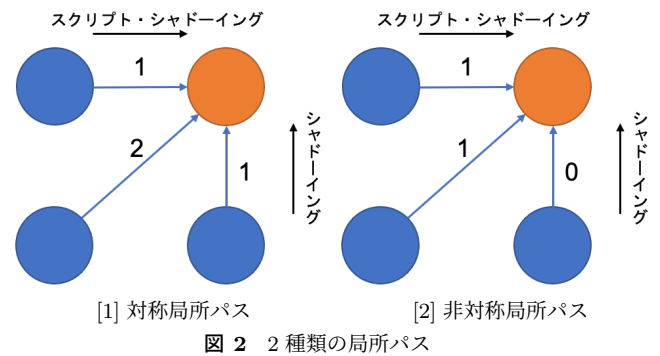


図 2 2 種類の局所パス

て、図 2-2 に示す非対称局所パスも検討した。これを用いると、局所距離の累積は常にスクリプト・シャドー音声の時間軸に沿って行われ、累積回数は当該音声長と等しくなる。即ち、スクリプト・シャドー音声の各時刻を等価に扱える。なお、このままだと母音長、子音長の偏りの影響を直接的に受けることになるので、非対称パスによる DTW 結果も、音素、音節、単語を単位として再集計する。

5. 実験

5.1 実験条件

了解度を計測する学習者音声として、30 人の日本人大学生に英語で作文させ、それを音読させて録音した。内容は学術研究から社会問題に至るまで多岐にわたる。音声長は 30~50 秒で、文法ミスや発音間違いが一部含まれている。学習者は Versant English Test [24] を受けており、そのスコアから、幅広い習熟度を一様に網羅できる 12 名を選定した。それらに米国人、英国人の母語話者の英語音声を加え、合計 14 音声を提示刺激として用いた。なお、聴取により、音声長が約 30 秒となるよう編集した。

5.2 実験手順

14 種類の音声をシャドーするシャドワーとして、英語母語話者 (N1, N2)、英語上級者であるが非母語話者 (NN1–NN3)、英語上級者の日本人 (J1, J2) の 3 グループのシャドワーが本実験に参加した。最初の二グループは日本語の知識はない。なお、NN1–NN3 の母語は、ベトナム語が一人、中国語が二人である。シャドワーは各々の音声をシャドーし、その後スクリプト・シャドーした。

シャドー音声の録音はウェブ上で行った。練習として、英検 2 級のリスニング問題用 3 音声と日本人英語 1 音声を用意した。練習は、希望に応じて複数回行わせた。実験刺激である 14 音声の方は、シャドー、スクリプトシャドー、共に一度ずつ行わせた。その後、シャドーイング音声全てを書き起こし、書き起こしによる了解度を計算した (図 1A)。音声合成による音声化、posteriorgram 化することによる了解度も計算した (図 1B)。

了解度のサンプル数を増やすため、読み上げ原稿を文ごとに区切り、posteriorgram を分割し、それぞれ独立に了解

表 1 手動スコア (A, B) と自動スコア (C, D, E) の相関

学習者音声全体に対して集計した場合								
Corr.	J1	J2	N1	N2	NN1	NN2	NN3	mean
D-A	0.481	0.501	0.734	0.626	0.629	0.728	0.743	0.635
E-A	0.874	0.690	0.889	0.854	0.904	0.910	0.847	0.853
C-A*	0.814	0.920	0.886	0.847	0.920	0.923	0.914	0.889
C-B*	0.890	0.932	0.961	0.835	0.961	0.951	0.976	0.929
学習者音声中文毎に集計した場合								
Corr.	J1	J2	N1	N2	NN1	NN2	NN3	mean
D-A	0.334	0.275	0.520	0.425	0.376	0.398	0.499	0.404
E-A	0.785	0.583	0.869	0.776	0.791	0.829	0.795	0.776
C-A*	0.699	0.722	0.736	0.771	0.787	0.822	0.825	0.766
C-B*	0.760	0.805	0.824	0.770	0.806	0.860	0.873	0.814

* 相関は負となるが、符号は除いている。

度を算出した。また、単語数の少ない文では、text-based DTW のスコアリングの粒度が極端に粗くなるため、8 単語以上の文のみ分析対象とし、全 64 文のうち、それ未満の 12 文は除外した。最終的に 14 のシャドー音声から、学習者発話全体に対するスコアが 14、文毎のスコアが 52、得られることになる (図 1C)。図 1A, B に示す手動スコアも同様、発話全体のスコア、文毎のスコアを用意した。

5.3 実験結果と考察

5.3.1 ASR による書き起こしとの比較

本実験の結果を示す前に、シャドー音声を人手で書き起こす代わりに、シャドー音声 (および学習者音声) を音声認識して得られる (自動) 書き起こしを用いた場合の結果を示す。シャドー音声の認識結果が十分機能すれば (図 1E), スクリプト・シャドー音声は不要となる。また、学習者音声の認識結果が了解度を示す指標として十分機能すれば (図 1D), 逆シャドーそのものが不要となる。Amazon の米英語用音声認識器 [25] を用い、サブスコア計算の単位は単語であり、DTW の局所パスは対称型を用いている。

表 1 に手動スコアである A, B と、自動スコアである C, D, E との相関を、各学習者発話全体に対して求めた場合と、文毎に求めた場合の結果を示す。D-A の相関は低く、音声認識結果を、人間の即時聴解を模擬するために使うこと適切ではないだろう。一方で E-A の相関は高く、聴取者のシャドー音声に対する音声認識結果を使うことは一定の効果がある。L2 音声全体を用いた場合は平均値では、C-A が E-A より高いが、文毎に集計すると後者の方が高くなった。シャドワー別にみると、J2, NN3 の相関は C-A の方が高い。彼らの英語音声を見ると、母語訛りが比較的強く、音声認識精度が十分に高くなかったことが考えられる。一方、シャドー音声とスクリプト・シャドー音声の DTW 比較は、当該話者間で発話を比較しており、外国語訛りには頑健に対応できている。国際語である英語は、世界諸英語 (World Englishes) と呼ばれるように、ある英語音声の了解度を議論する場合、様々な母語を持つ英語利用者に対して調査する必要がある。そのような状況を考えると、外国語訛りへの頑健性は大きな意味を持つ。

シャドー音声に対する wDTW はバタチャリヤ距離の平

表 2 各種手動スコア (A) と各種自動スコア (C) の相関
シャドワー間平均, 対称局所パス使用

	A	C	mean (発話毎)	mean (文毎)
1	pAcc	fDTW	0.841	0.726
2	pAcc	pDTW	0.904	0.813
3	sAcc	fDTW	0.844	0.751
4	sAcc	sDTW	0.890	0.816
5	wAcc	fDTW	0.848	0.726
6	wAcc	wDTW	0.889	0.766

相関は負となるが、符号は除いている。

表 3 各種手動スコア (B) と各種自動スコア (C) の相関
シャドワー間平均, 対称局所パス使用

	B	C	mean (発話毎)	mean (文毎)
1	S-fDTW	fDTW	0.875	0.743
2	S-pDTW	pDTW	0.922(+0.018)	0.845(+0.032)
3	S-sDTW	sDTW	0.915(+0.025)	0.827(+0.011)
4	S-wDTW	wDTW	0.929(+0.040)	0.814(+0.048)

括弧内は表 2 からの増分

表 4 各種手動スコア (A) と各種自動スコア (C) の相関
シャドワー間平均, 非対称局所パス使用

	A	C	mean (発話毎)	mean (文毎)
1	pAcc	fDTW	0.897(+0.056)	0.795(+0.069)
2	pAcc	pDTW	0.904(+0.000)	0.820(+0.007)
3	sAcc	fDTW	0.893(+0.049)	0.809(+0.058)
4	sAcc	sDTW	0.896(+0.006)	0.830(+0.014)
5	wAcc	fDTW	0.863(+0.015)	0.763(+0.037)
6	wAcc	wDTW	0.893(+0.004)	0.780(+0.014)

括弧内は表 2 からの増分

表 5 各種手動スコア (B) と各種自動スコア (C) の相関
シャドワー間平均, 非対称局所パス使用

	B	C	mean (発話毎)	mean (文毎)
1	S-fDTW	fDTW	0.922(+0.047)	0.802(+0.059)
2	S-pDTW	pDTW	0.924(+0.002)	0.853(+0.008)
3	S-sDTW	sDTW	0.923(+0.008)	0.842(+0.015)
4	S-wDTW	wDTW	0.935(+0.006)	0.831(+0.017)

括弧内は表 3 からの増分

均値として計算されるが、書き起こしに対する wAcc は、単語の合致/非合致を 1/0 で計算しており、スコアリングの粒度に乖離がある。後者を posterigram 化してスコア計算すると、この乖離は解消される (図 1B)。表 1 中、最高相関値は、多くの場合で C-B で得られている。

5.3.2 複数粒度による了解度アノテーション

表 2 に、対称局所パス使用時の手動スコア (各種 Acc, A) と自動スコア (各種 DTW, C) との相関を示す。3 種類の粒度の手動スコアに対して、フレーム単位及び当該言語単位の DTW スコアとの相関を各シャドワー毎に算出し、その平均値を示している。また、約 30 秒の発話毎の集計と、文毎に集計を別々に示している。まず、表 1 同様、分析対象の音声長が長いほど自動スコアと手動スコアの相関は高くなっている。また、各種粒度 (言語単位) に対して、フレーム単位の平均 DTW スコアを直接使うよりも、当該言語単位でサブスコアを計算し、その平均値を最終的なスコアとした方が、いずれの場合も相関が高い。L2 音声の了解度を考える場合、どの単位に基づいて検討するのかは教育方針に依存するが、いずれの方針に対しても、より適切な自動スコアが計算できている。なお、単語を単位と

して L2 音声 を 文 単 位 で 集 計 し た 場 合 の 相 関 が 低 い が、こ れ は、単 語 単 位 の text-based DTW を 行 う と、単 語 の 合 致 / 非 合 致 を 1/0 で 計 算 し て お り、文 中 の 単 語 数 が 少 な い 場 合、ス コ ア リ ン グ の 粒 度 が 粗 く な る こ と が 原 因 で あ る。

こ の 問 題 は、書 き 起 こ し と 原 稿 と の DTW を テ キ ス ト で 行 う の で は な く、音 声 化・posteriorgram 化 し て 行 う (S-DTW, **B**) こ と で 回 避 で き る と 予 想 さ れ る。各 種 単 位 の S-DTW ス コ ア を 用 い た 場 合 の 相 関 を 表 3 に 示 す。括 弧 内 の 数 値 は 表 2 か ら の 増 分 で あ る。い ず れ の 単 位 で も 相 関 は 増 加 し て お り そ の 効 果 が 窺 わ れ る が、単 語 を 単 位 と し て 文 毎 に 集 計 し た 場 合 の 増 分 が 一 番 高 い。

5.3.3 非対称 DTW パスによる了解度アノテーション

表 4, 5 に **C** に お け る DTW に 非 対 称 局 所 パ ス を 導 入 し て、前 節 と 同 じ 分 析 を 行 っ た 結 果 を 示 す。い ず れ の 場 合 も、発 話 毎、文 毎 で の 比 較 と も に、非 対 称 局 所 パ ス の 導 入 に よ っ て、手 動 ス コ ア と の 相 関 は 高 く な っ て お り、シャ ドー 音 声 の 崩 れ の 特 性 を 適 切 に 捉 え ら れ て い る と い え る。ま た、相 関 の 増 分 は fDTW で の 増 分 が 顕 著 に 高 く、そ れ 以 外 は 凡 そ、 $wDTW \approx sDTW > pDTW$ と な っ て い る。非 対 称 局 所 パ ス に よ る fDTW は 図 2 の 左 図 に お け て、縦 方 向 に 伸 び る 局 所 パ ス の ス コ ア を 無 視 し、局 所 ス コ ア の 累 積 は 規 範 発 声 で あ る ス ク リ プ ト・シャ ドー の 各 音 声 フ レー ム 毎 に 行 う こ と と な る。そ の 結 果、fDTW は、手 動 ス コ ア と し て ど の 言 語 単 位 を 使 っ た 場 合 で も、よ り 類 似 し た 自 動 ス コ ア を 呈 ず る よ う に な っ た が、当 該 単 位 を 使 っ た DTW ス コ ア に は 至 っ て い な い。ま た、 $wDTW \approx sDTW > pDTW$ の 傾 向 で あ る が、サ ブ ス コ ア を 算 出 す る 区 間 が 長 い ほ ど シャ ドー イ ン グ が 間 延 び た 部 分 の 影 響 が 大 き く、非 対 称 局 所 パ ス に よ る そ の 影 響 の 打 ち 消 し の 効 果 が 大 き い こ と に よ る。

な お、対 称 局 所 パ ス を 用 い た 表 1 で は、文 毎 で の 比 較 の 場 合、 $wDTW$ を 用 い た 平 均 相 関 (**C-A**) が、シャ ドー 音 声 の 音 声 認 識 結 果 を 用 い た 平 均 相 関 (**E-A**) を 下 回 っ て い た が、表 4 に お け て 非 対 称 局 所 パ ス を 導 入 す る こ と で、 $wDTW$ の 平 均 相 関 が 上 回 る 結 果 と な っ た。最 終 的 に テ キ ス ト 比 較 に お け て も posterior-gram DTW を 導 入 す る こ と で (表 5)、音 素、音 節、単 語 い ず れ の 言 語 単 位 に 対 し て も、発 話 全 体 の 場 合 は 0.92 以 上 の 相 関 を 示 し、文 単 位 で も 0.83 以 上 の 相 関 を 示 す こ と が で き た。シャ ドー 音 声 の 手 動 書 き 起 こ し を 用 い な く て も、シャ ドー 音 声 と ス ク リ プ ト シャ ドー 音 声 を 確 率 的 に 自 動 音 素 書 き 起 こ し す る こ と で、高 い 精 度 で 瞬 時 了 解 度 の ア ノ テー シ ョ ン は 取 得 で き る。

6. まとめと今後の課題

本 報 告 で は、逆 シャ ドー イ ン グ 法 に よ る 瞬 時 了 解 度 の ア ノ テー シ ョ ン を、複 数 の 言 語 単 位 に 基 づ く 粒 度 で 定 量 化 し、さ ら に 非 対 称 な 局 所 パ ス を 用 い た DTW に よ り ア ノ テー シ ョ ン の 高 精 度 化 を 分 析 的 に 検 討 し た。そ の 結 果、DTW を 教 育 方 針 に 応 じ た 粒 度 で、非 対 称 局 所 パ ス を 用 い て 行 う

こ と で、手 動 書 き 起 こ し に よ る 正 解 ラ ベ ル と よ り 高 い 相 関 を 持 つ 了 解 度 ア ノ テー シ ョ ン を 実 現 で き た。同 時 に L2 音 声 や シャ ドー イ ン グ 音 声 に 音 声 認 識 器 を 適 用 し、自 動 で 書 き 起 こ す こ と に よ っ て シャ ドー イ ン グ や ス ク リ プ ト・シャ ドー イ ン グ を 代 替 で き な い か 検 討 し た。音 声 認 識 器 に よ る ア ノ テー シ ョ ン は、シャ ドー イ ン グ 音 声 を 自 動 で 書 き 起 こ し、学 習 者 の 原 稿 と 比 較 す る こ と で、手 動 に よ る 了 解 度 と あ る 程 度 の 相 関 を 示 し た が、シャ ドー 自 身 の 訛 り の 影 響 を 受 け、精 度 が あ ま り 上 が ら な い も の も 見 ら れ た。

今 回 は ア ノ テー シ ョ ン を 複 数 粒 度 で 行 い、局 所 パ ス を 変 え る と い う ア プ ロ ー チ で 高 精 度 化 を 図 っ た が、精 度 の よ り 一 層 高 い ア ノ テー シ ョ ン を 行 う に は、い く つ か の 手 法 が 考 え ら れ る。一 つ は、ア ノ テー シ ョ ン 粒 度 の 階 層 化 で あ る。今 回、複 数 粒 度 に よ る ア ノ テー シ ョ ン と し て、言 語 単 位 に 応 じ て サ ブ ス コ ア を 求 め た が、そ の サ ブ ス コ ア は そ の 区 間 に お け て fDTW を す る こ と に よ っ て 求 め ら れ る。例 え ば $wDTW$ に お け る サ ブ ス コ ア の 計 算 を $pDTW$ の 結 果 を 使 う と い っ た よ う に、あ る 言 語 単 位 に お け る サ ブ ス コ ア の 計 算 を、そ れ よ り 細 か い 言 語 単 位 を 粒 度 と し た DTW で 改 め て 行 う こ と で、サ ブ ス コ ア 内 で も 継 続 長 の 正 規 化 が な さ れ、よ り ア ノ テー シ ョ ン の 高 精 度 化 が 期 待 で き る。二 つ 目 と し て、posteriorgram 以 外 の 特 徴 量 (特 に 韻 律) の 使 用 が 考 え ら れ る。逆 シャ ドー イ ン グ 法 は、シャ ドー の 即 時 聴 解 の 様 子 を シャ ドー 音 声 と ス ク リ プ ト・シャ ドー 音 声 の 差 異 を 通 し て 観 測 し て い る が、両 者 の 差 異 は、調 音 構 造 の 崩 れ (posteriorgram DTW に 相 当) 以 外 に も、韻 律 の 崩 れ と し て も 観 測 さ れ る。こ れ ら を 併 用 す る こ と で、瞬 時 了 解 度 ア ノ テー シ ョ ン の 更 な る 高 精 度 化 を 検 討 し た い。

参考文献

- [1] Munro, M. J. and Derwing, T. M.: Foreign accent, comprehensibility, and intelligibility in the speech of second language learners, *Language learning*, Vol. 45, No. 1, pp. 73-97 (1995).
- [2] Murphy, J. M.: Intelligible, comprehensible, non-native models in ESL/EFL pronunciation teaching, *System*, Vol. 42, pp. 258-269 (2014).
- [3] Derwing, T. M. and Munro, M. J.: *Pronunciation fundamentals: evidence-based perspectives for L2 teaching and research*, John Benjamins (2015).
- [4] Liontas, J. I., Association, T. I. et al.: *The TESOL Encyclopedia of English Language Teaching*, John Wiley & Sons (2018).
- [5] Levis, J.: Revisiting the intelligibility and nativeness principles, *Journal of Second Language Pronunciation*, Vol. 6, No. 3, pp. 310-328 (2020).
- [6] Smith, B.: *Learner English: A teacher's guide to interference and other problems*, Ernst Klett Sprachen (2001).
- [7] Zhu, C., Lin, Z., Minematsu, N. and Nakanishi, N.: Analyses on instantaneous perception of Japanese English by listeners with various language profiles, *Proc. The Phonetic Society of Japan General Meeting*, pp. 26-31 (2020).

- [8] Zhu, C., Minematsu, N. and Nakanishi, N.: Objective and semi-automatic measurement of smoothness of instantaneous understanding of L2 English speech (to appear), *Proc. Pronunciation in Second Language Learning and Teaching Conference* (2020).
- [9] Derwing, T. M. and Munro, M. J.: Accent, intelligibility, and comprehensibility: Evidence from four L1s, *Studies in second language acquisition*, pp. 1–16 (1997).
- [10] Bernstein, J.: Objective measurement of intelligibility, *Proc. ICPHS*, pp. 1581–1584 (2003).
- [11] Minematsu, N., Guo, C. and Hirose, K.: CART-based factor analysis of intelligibility reduction in Japanese English, *Proc. EUROSPEECH*, pp. 2069–2072 (2003).
- [12] Minematsu, N., Okabe, K., Ogaki, K. and Hirose, K.: Measurement of objective intelligibility of Japanese accented English using ERJ database, *Proc. INTERSPEECH*, pp. 1481–1484 (2011).
- [13] Kang, O., Thomson, R. I. and Moran, M.: Empirical approaches to measuring the intelligibility of different varieties of English in predicting listener comprehension, *Language Learning*, Vol. 68, No. 1, pp. 115–146 (2018).
- [14] Alison, M., Rebecca, A., Catherine, S. and Paula, W.: Exploring the Relationship Between Modified Output and Working Memory Capacity, *Language Learning*, Vol. 60, No. 3, pp. 501–533 (2010).
- [15] Munro, M. J. and Derwing, T. M.: Foreign accent, comprehensibility and intelligibility, redux, *Journal of Second Language Pronunciation*, Vol. 6, No. 3, pp. 283–309 (2020).
- [16] Field, J.: Intelligibility and the listener: the role of lexical stress, *TESOL quarterly*, Vol. 39, No. 3, pp. 399–423 (2005).
- [17] Inoue, Y., Kabashima, S., Saito, D., Minematsu, N., Kanamura, K. and Yamauchi, Y.: A study of objective measurement of comprehensibility through native speakers shadowing of learners’ utterances, *Proc. INTERSPEECH*, pp. 1651–1655 (2018).
- [18] Trisitchoke, T., Ando, S., Inoue, Y., Saito, D. and Minematsu, N.: Influence of content variations on smoothness of native speakers’ reverse shadowing, *Proc. ICPHS* (2019).
- [19] Ando, S., Lin, Z., Trisitchoke, T., Inoue, Y., Yoshizawa, F., Saito, D. and Minematsu, N.: A Large Collection of Sentences Read Aloud by Vietnamese Learners of Japanese and Native Speaker’s Reverse Shadowings, *Proc. O-COCOSDA*, pp. 1–6 (2019).
- [20] Lin, Z., Takashima, R., Saito, D., Minematsu, N. and Nakanishi, N.: Shadowability annotation with fine granularity on L2 utterances and its improvement with native listeners’ script-shadowing, *Proc. INTERSPEECH*, pp. 3865–3869 (2020).
- [21] Povey, D., Ghoshal, A., Boulianne, G., Glembek, L. B., Goel, N., Hannemann, M., Motlíček, P., Y., Q., P., S., Silovský, J., Stemmer, G. and Veselý, K.: The KALDI speech recognition toolkit, *Proc. ASRU* (2011).
- [22] Amazon: *Amazon Polly*. <https://aws.amazon.com/polly/>.
- [23] Shi, S. and Minematsu, N.: A corpus-based analysis of shadowing speech: case of L2 English by Japanese learners, *Proc. ISAPh*, No. 3, pp. 34–37 (2016).
- [24] VERSANT: *Versant English Test*. <https://www.pearson.com/english/versant/tests.html>.
- [25] Amazon: *Amazon Transcribe*. <https://aws.amazon.com/transcribe/>.